



US 20170211142A1

(19) **United States**

(12) **Patent Application Publication**
Smargon et al.

(10) **Pub. No.: US 2017/0211142 A1**

(43) **Pub. Date: Jul. 27, 2017**

(54) **NOVEL CRISPR ENZYMES AND SYSTEMS**

filed on Feb. 17, 2016, provisional application No. 62/376,382, filed on Aug. 17, 2016, provisional application No. 62/376,367, filed on Aug. 17, 2016.

(71) Applicants: **Aaron Andrew Smargon**, Cambridge, MA (US); **Fenq Zhang**, Cambridge, MA (US)

Publication Classification

(72) Inventors: **Aaron Andrew Smargon**, Cambridge, MA (US); **Fenq Zhang**, Cambridge, MA (US)

(51) **Int. Cl.**
C12Q 1/68 (2006.01)
G06F 19/20 (2006.01)
G06F 19/22 (2006.01)

(73) Assignees: **THE BROAD INSTITUTE, INC.**, Cambridge, MA (US);
MASSACHUSETTS INSTITUTE OF TECHNOLOGY, Cambridge, MA (US)

(52) **U.S. Cl.**
CPC *C12Q 1/6876* (2013.01); *G06F 19/22* (2013.01); *G06F 19/20* (2013.01); *C12Q 2600/156* (2013.01)

(21) Appl. No.: **15/331,792**

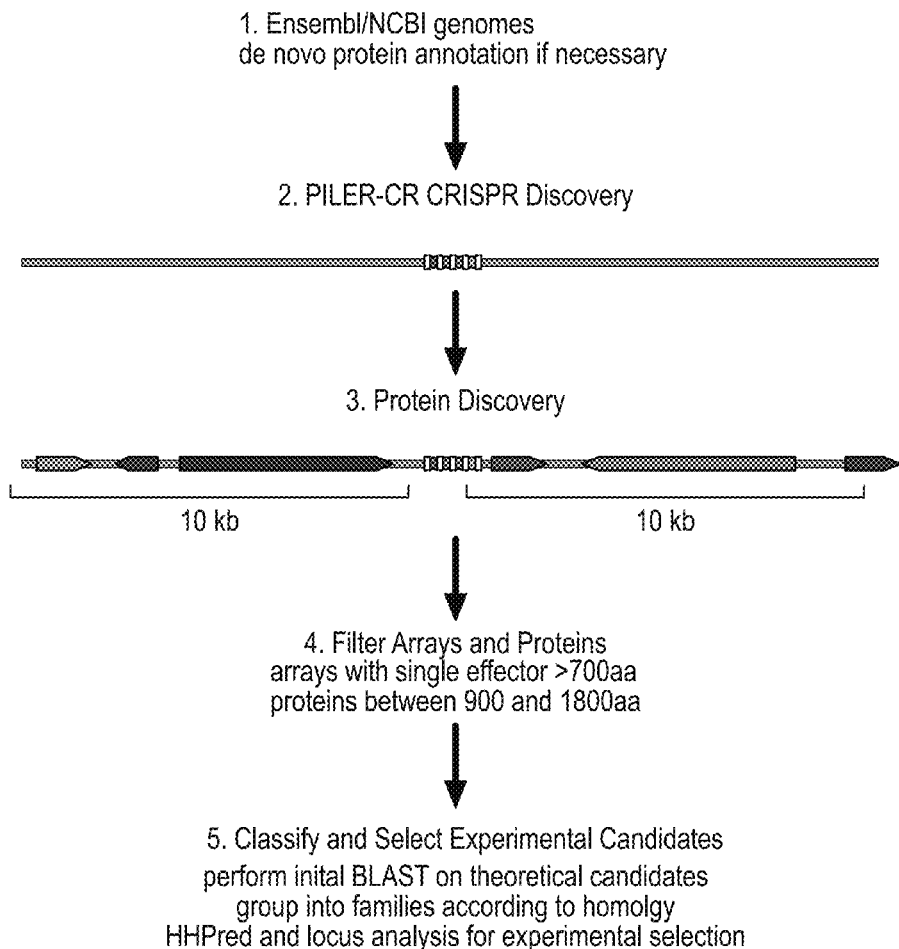
(57) **ABSTRACT**

(22) Filed: **Oct. 21, 2016**

The invention provides for systems, methods, and compositions for targeting nucleic acids. In particular, the invention provides non-naturally occurring or engineered DNA or RNA-targeting systems comprising a novel DNA or RNA-targeting CRISPR effector protein and at least one targeting nucleic acid component like a guide RNA.

Related U.S. Application Data

(60) Provisional application No. 62/245,270, filed on Oct. 22, 2015, provisional application No. 62/296,548,



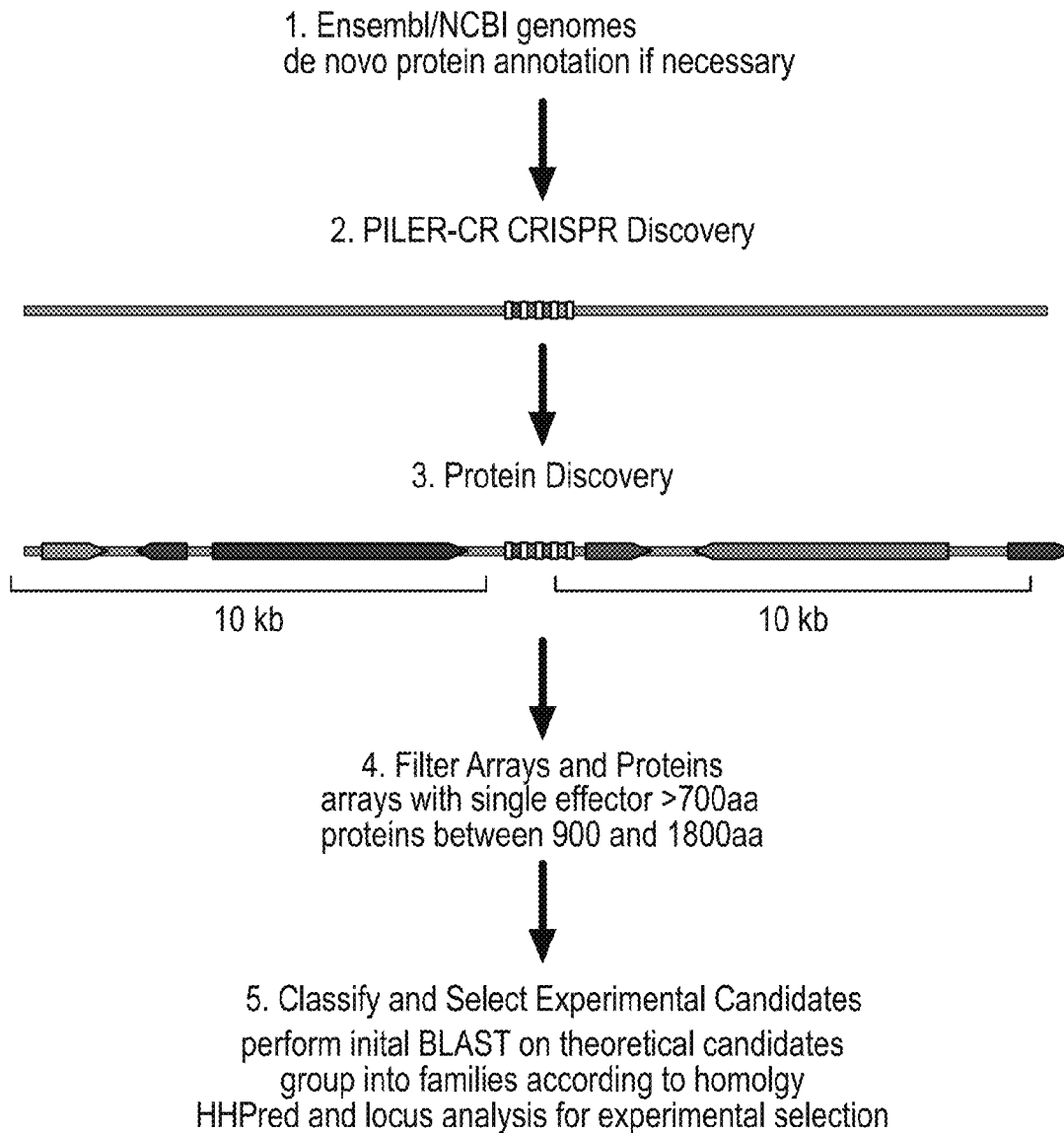


FIG. 1

23,001 full, annotated bacterial and archaeal genomes from Ensembl Genomes Release 27

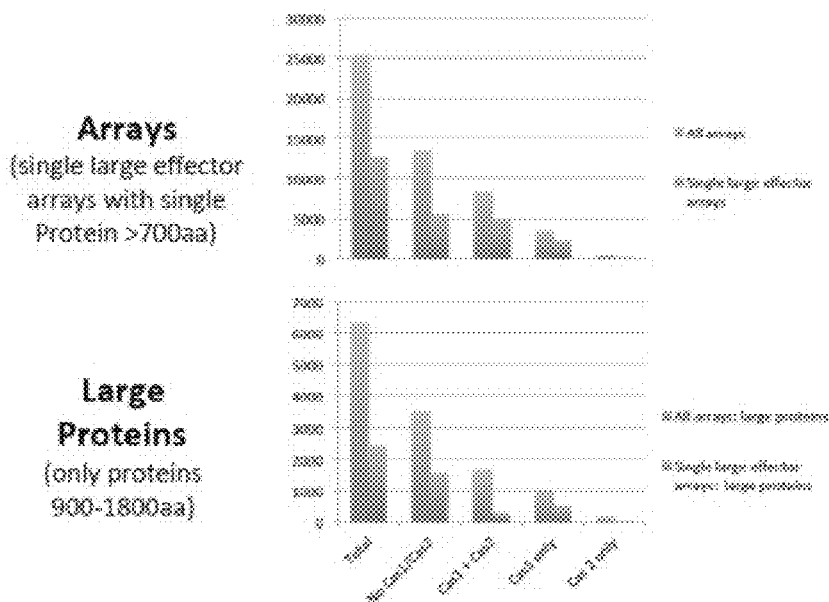


FIG. 2

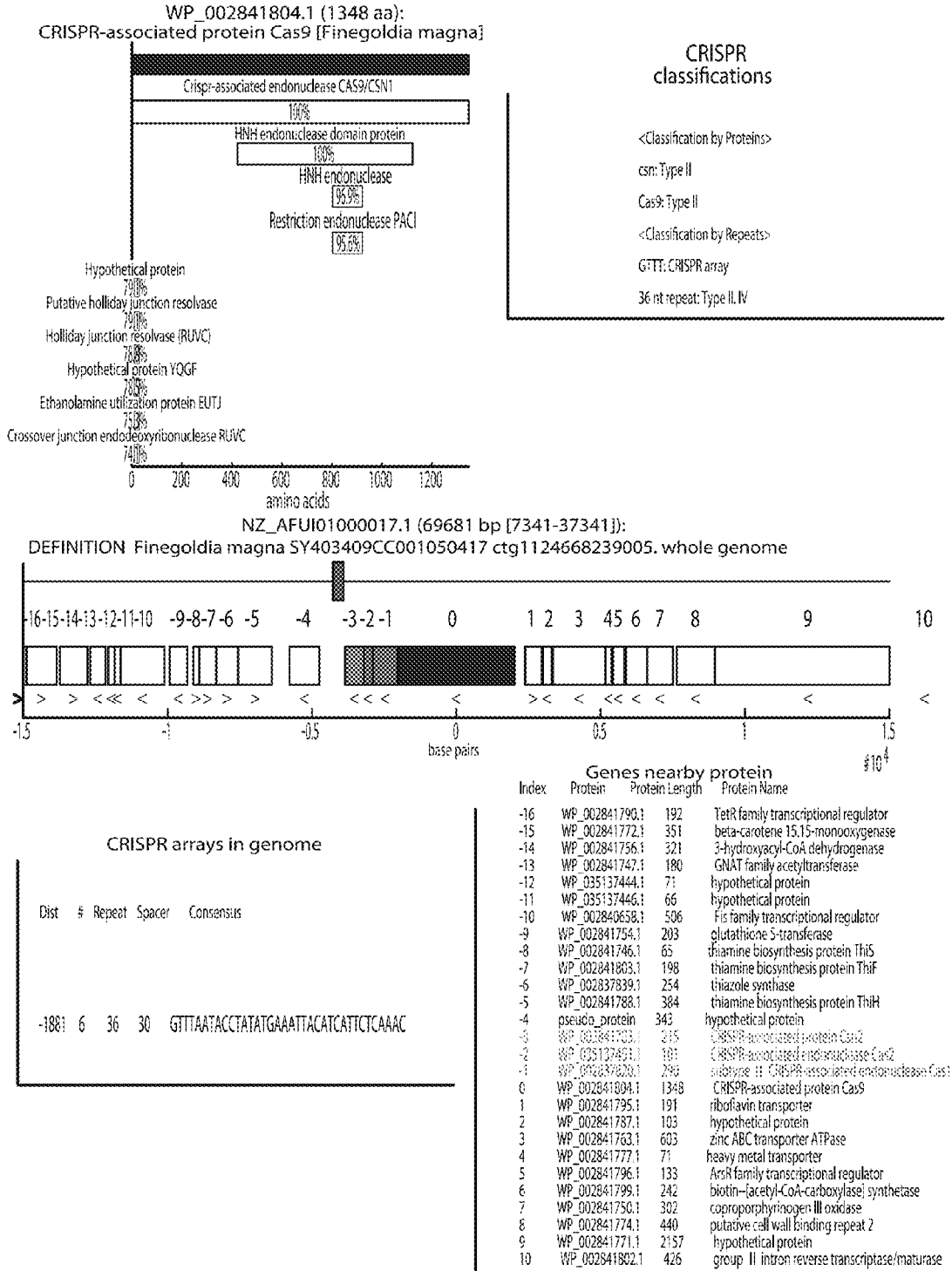
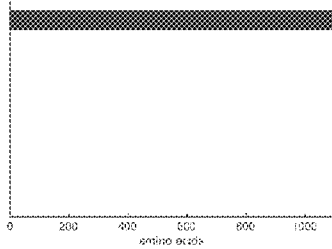


FIG. 3

ADZ11800.1 (1096 aa):
hypothetical protein RIA_0639 [Riemerella anatispestifer RA-GD]

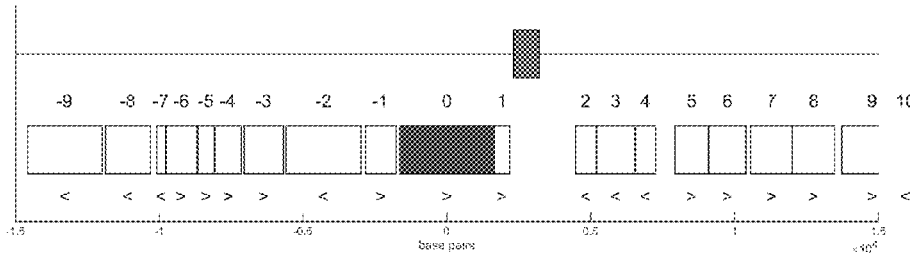


CRISPR
classifications

<Classification by Repeats>

98 of repeat Type II-IV

CP002562.1 (2166384 bp [606551.5-636551.5]):
DEFINITION Riemerella anatispestifer RA-GD, complete genome.



CRISPR arrays in genome

Dist	# Repeat	Spacer	Consensus	
400351	8	82	26	GTTFGCGAWCTATGTAAGAGATAGTAAAGATTGAAAGCAATTCAGAGATT
800902	9	47	20	GTTCGGAAATTTGACAAAGATAGTAAAGATTGAAAGCAATTCAGAGATT
912	14	26	38	GTTCGGAAATTTGACAAAGATAGTAAAGATTGAAAGCAATTCAGAGATT

Genes nearby protein

Index	Protein	Protein Length	Protein Name
-9	ADZ11797.1	993	Glucanase hydrolase, family 5, conserved site
-8	ADZ11792.1	519	Ribonucleosase G and E
-7	ADZ11793.1	97	Bacteroid phage-associated DNA-binding protein
-6	ADZ11794.1	993	ADP-specific DNA glycosylase
-5	ADZ11795.1	187	GelD
-4	ADZ11796.1	206	Phis
-3	ADZ11797.1	495	peptide:glycyl isomerase
-2	ADZ11798.1	897	Alanine-RNA synthetase
-1	ADZ11799.1	969	Glycosylase
0	ADZ11800.1	1096	hypothetical protein
1	ADZ11801.1	178	hypothetical protein
2	ADZ11802.1	293	hypothetical protein
3	ADZ11803.1	442	C-type lectin
4	ADZ11804.1	226	hypothetical protein
5	ADZ11805.1	987	Acetylglucosyl 1-3-glucosylase
6	ADZ11806.1	429	Acetyl-CoA lyase
7	ADZ11807.1	482	Aspartylase:phagey-RNA synthetase
8	ADZ11808.1	489	DNA-dependent RNA polymerase specialized sigma
9	ADZ11809.1	709	TonB-dependent outer membrane receptor
10	ADZ11810.1	70	hypothetical protein

FIG. 4D

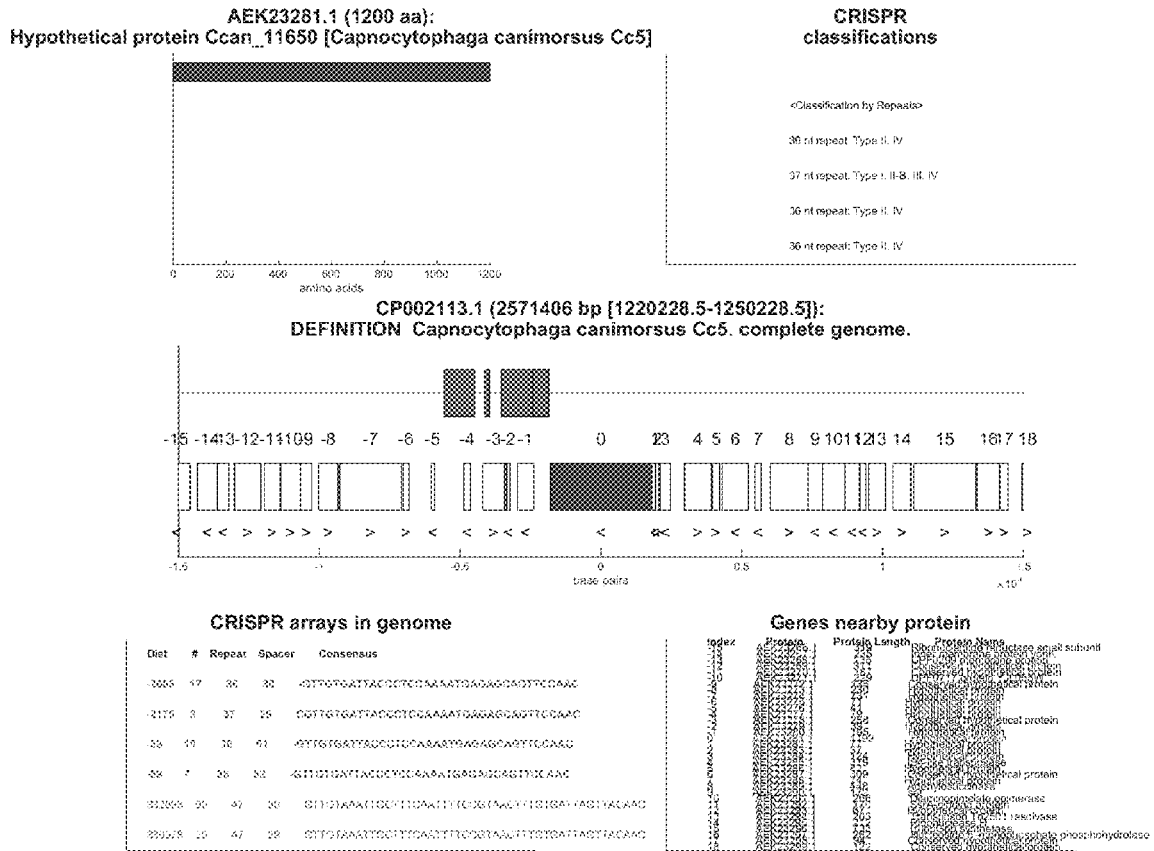
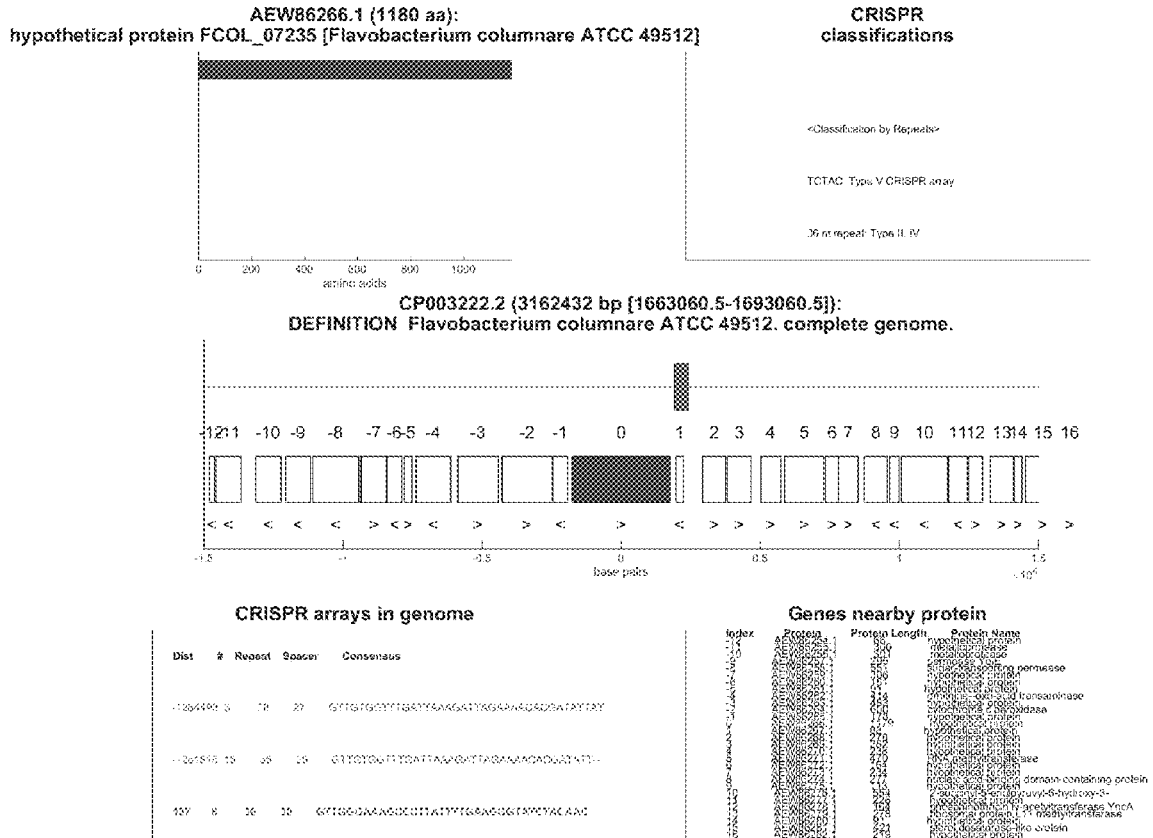


FIG. 4E



CRISPR arrays in genome

Dist	# Repeat	Spacer	Consensus
1284489	3	27	GTTGTGCTTTGATTAAGAAATYGGAAAGAGGAGATATCTAT
1361518	15	25	GTTGTGCTTTGATTAAGAAATYGGAAAGAGGAGATATCTAT
137	8	10	GTTGTGCTTTGATTAAGAAATYGGAAAGAGGAGATATCTAT

Genes nearby protein

Index	Protein	Protein Length	Protein Name
1	permease
2	consensinase
3	idase
4	ase
5	domain-containing protein
6	curv-8-hydroxylase
7	arylsulfonamide YacA
8	ribitol transferase
9	pro. protein

FIG. 4F

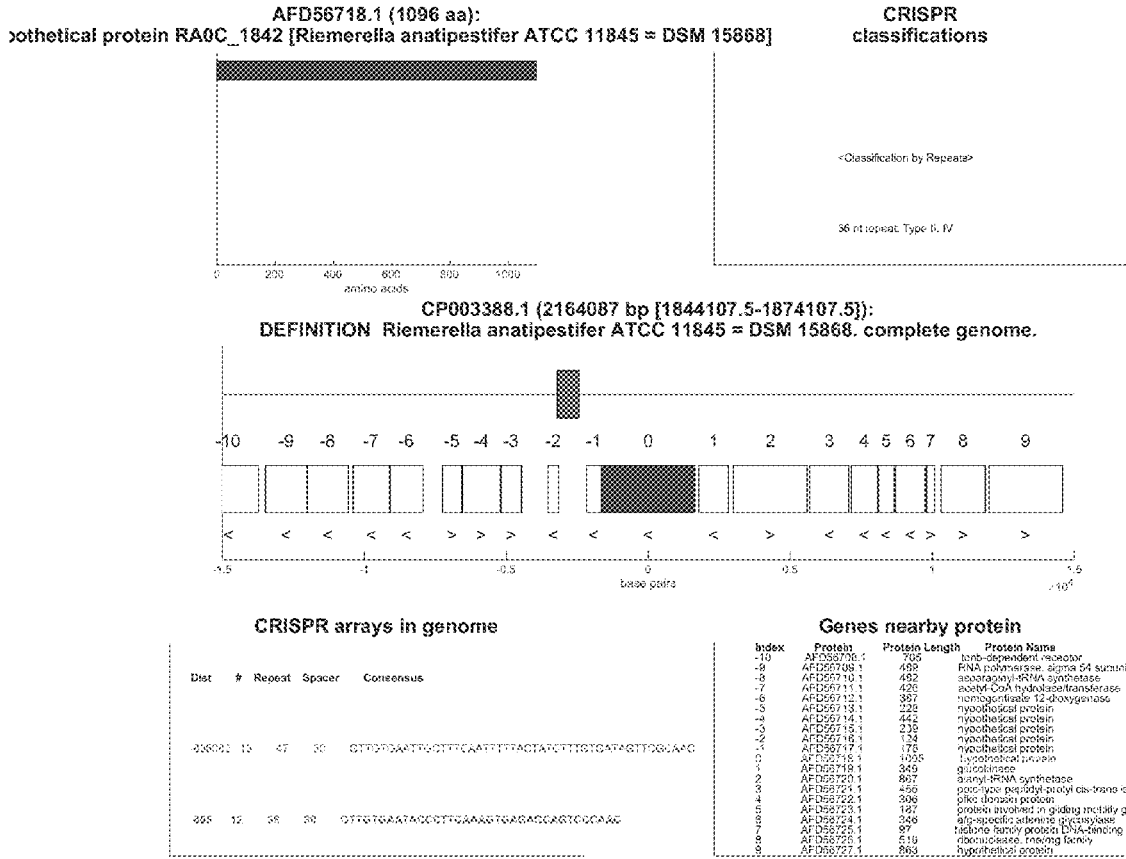


FIG. 4G

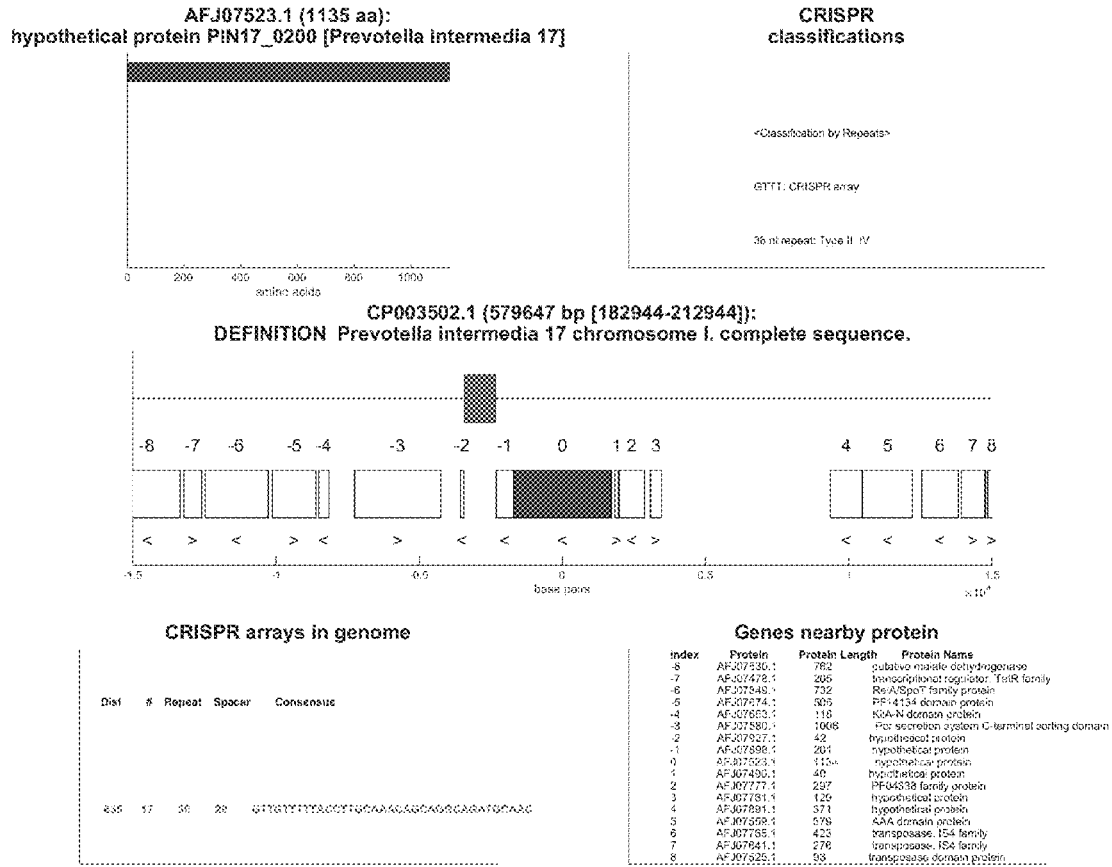


FIG. 4H

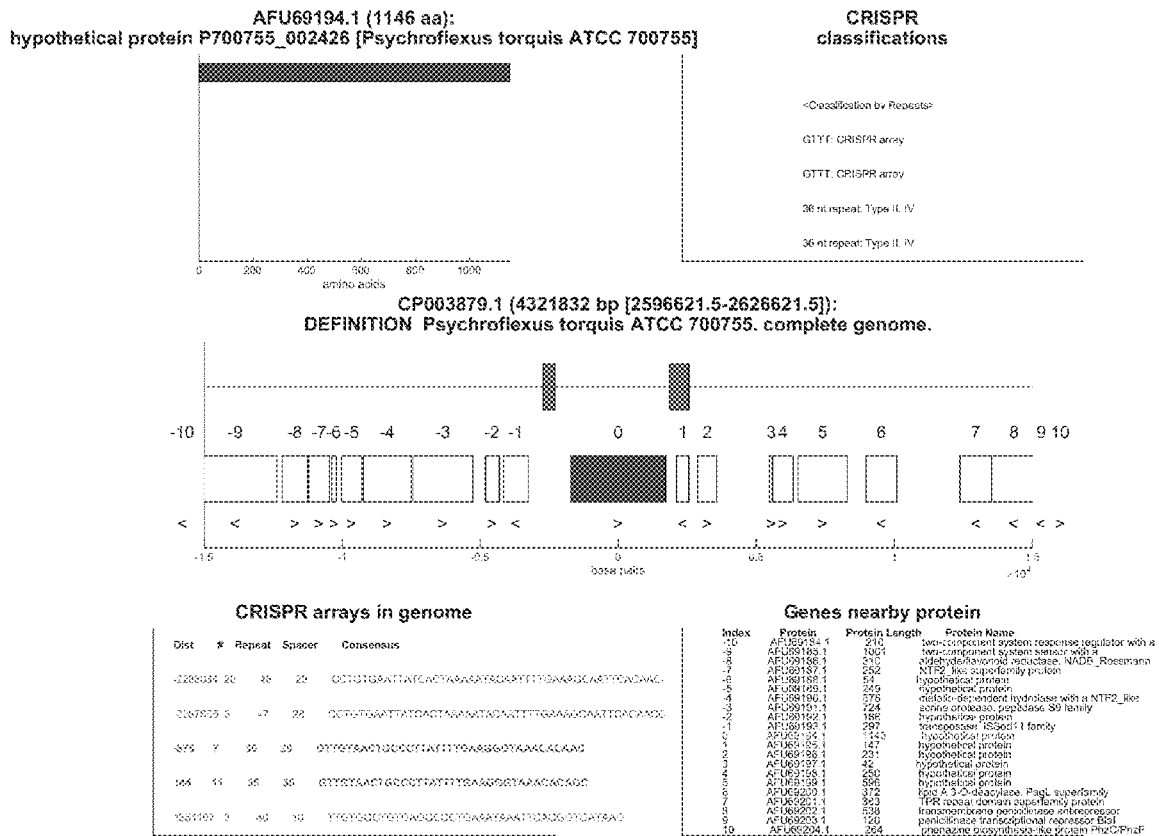


FIG. 41

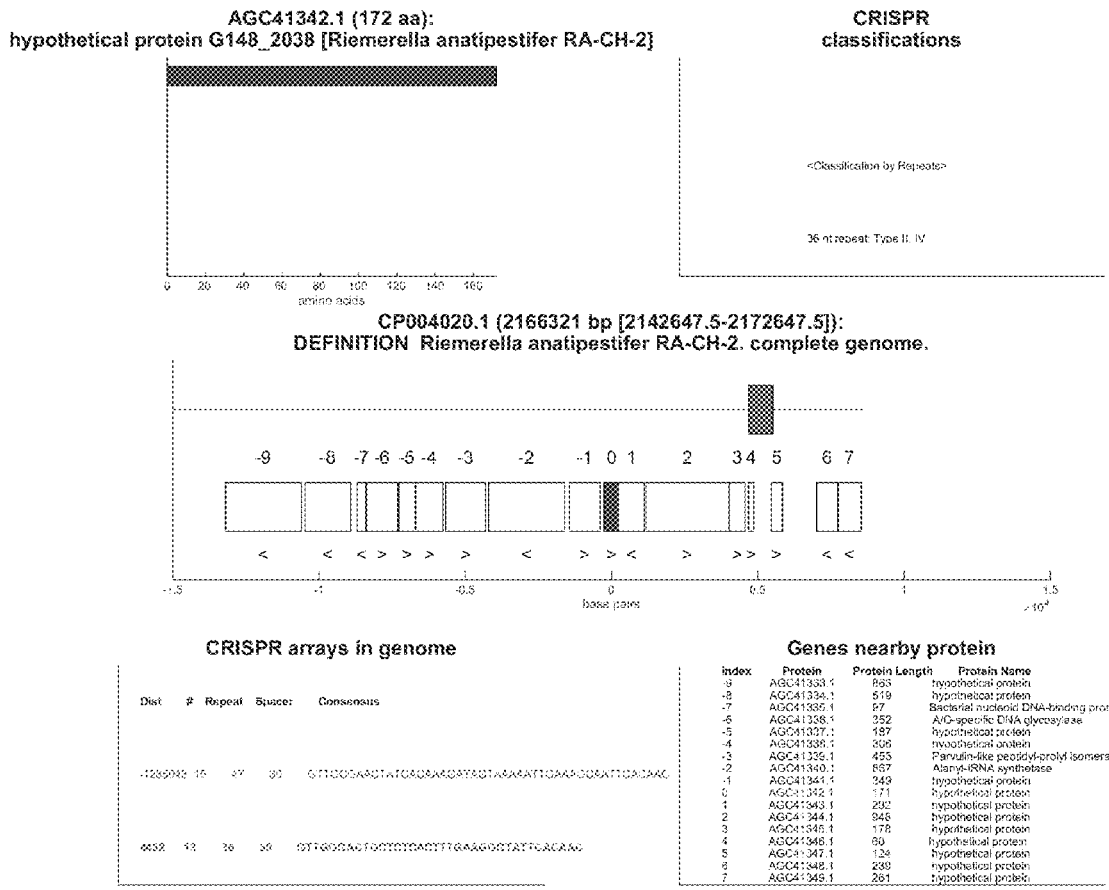


FIG. 4J

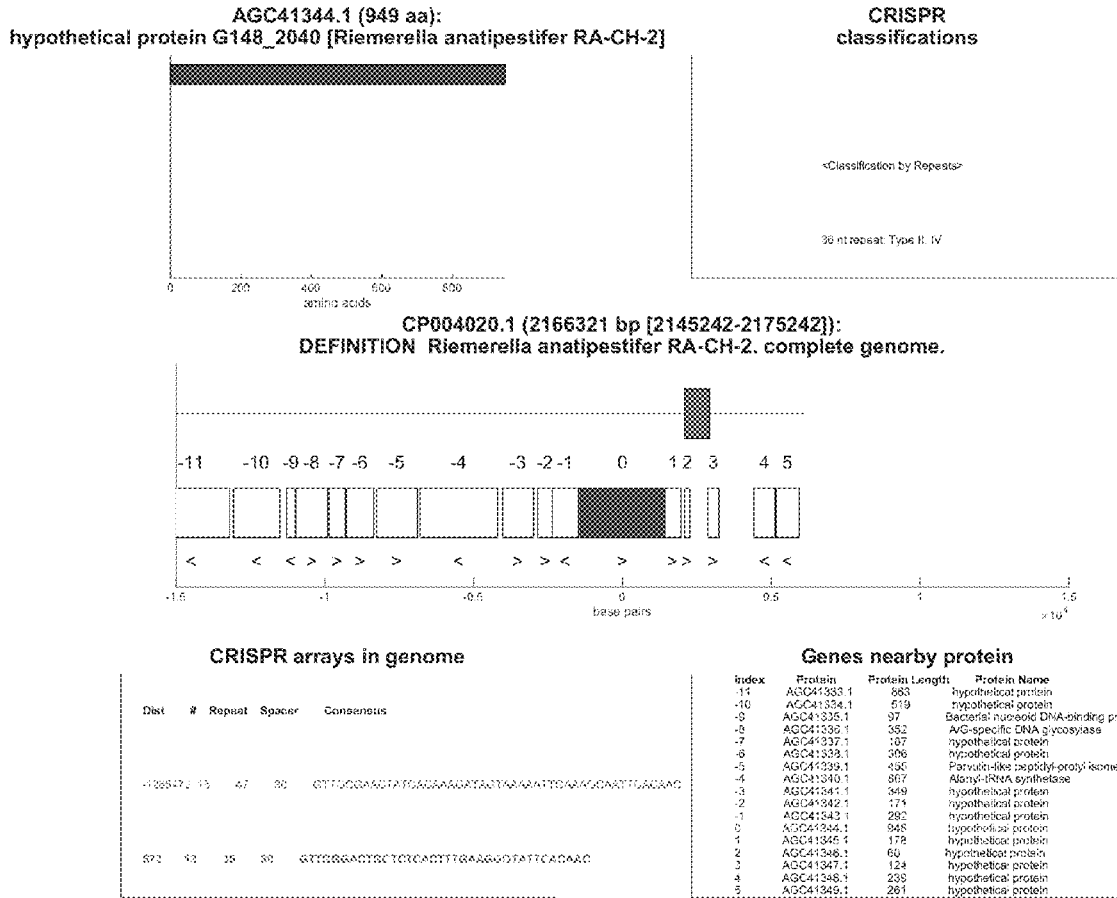
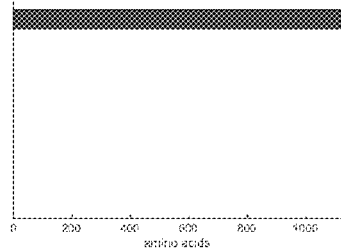


FIG. 4K

AJ36392.1 (1120 aa):
hypothetical protein EG14_10345 [Porphyromonas gingivalis]



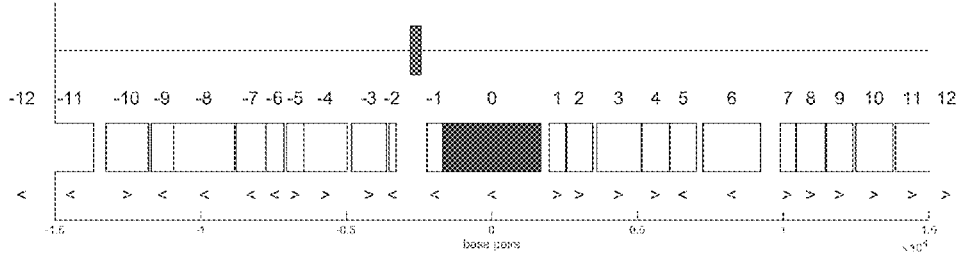
CRISPR
classifications

<Classification by Repeats

TCTAC, Type V CRISPR array

35 nt repeat, Type II, IV

CP007756.1 (2441780 bp [2298856.5-2326856.5]):
DEFINITION Porphyromonas gingivalis strain HG56 genome.



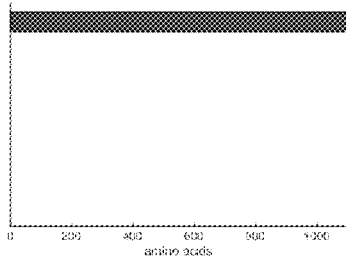
CRISPR arrays in genome

Dist	#	Repeat	Spacer	Consensus
2593290	99	10	35	GTTTTAAATTCCTGTATGSDTGGAAATGAAAT
582440	11	20	20	GTTTGGGAAATAGCTGTACTTCAAGATGTTGGAGGAGG
690	6	36	30	GTTCGCGTACGCTTCGAAATAGAGCGTACGATCGAAG

Genes nearby protein

Index	Protein	Protein Length	Protein Name
1	AJ36392.1	1120	hypothetical protein
2	AJ36393.1	255	cellulase
3	AJ36394.1	255	cellulase
4	AJ36395.1	255	cellulase
5	AJ36396.1	255	cellulase
6	AJ36397.1	255	cellulase
7	AJ36398.1	255	cellulase
8	AJ36399.1	255	cellulase
9	AJ36400.1	255	cellulase
10	AJ36401.1	255	cellulase
11	AJ36402.1	255	cellulase
12	AJ36403.1	255	cellulase
13	AJ36404.1	255	cellulase
14	AJ36405.1	255	cellulase
15	AJ36406.1	255	cellulase
16	AJ36407.1	255	cellulase
17	AJ36408.1	255	cellulase
18	AJ36409.1	255	cellulase
19	AJ36410.1	255	cellulase
20	AJ36411.1	255	cellulase
21	AJ36412.1	255	cellulase
22	AJ36413.1	255	cellulase
23	AJ36414.1	255	cellulase
24	AJ36415.1	255	cellulase
25	AJ36416.1	255	cellulase
26	AJ36417.1	255	cellulase
27	AJ36418.1	255	cellulase
28	AJ36419.1	255	cellulase
29	AJ36420.1	255	cellulase
30	AJ36421.1	255	cellulase
31	AJ36422.1	255	cellulase
32	AJ36423.1	255	cellulase
33	AJ36424.1	255	cellulase
34	AJ36425.1	255	cellulase
35	AJ36426.1	255	cellulase
36	AJ36427.1	255	cellulase
37	AJ36428.1	255	cellulase
38	AJ36429.1	255	cellulase
39	AJ36430.1	255	cellulase
40	AJ36431.1	255	cellulase
41	AJ36432.1	255	cellulase
42	AJ36433.1	255	cellulase
43	AJ36434.1	255	cellulase
44	AJ36435.1	255	cellulase
45	AJ36436.1	255	cellulase
46	AJ36437.1	255	cellulase
47	AJ36438.1	255	cellulase
48	AJ36439.1	255	cellulase
49	AJ36440.1	255	cellulase
50	AJ36441.1	255	cellulase
51	AJ36442.1	255	cellulase
52	AJ36443.1	255	cellulase
53	AJ36444.1	255	cellulase
54	AJ36445.1	255	cellulase
55	AJ36446.1	255	cellulase
56	AJ36447.1	255	cellulase
57	AJ36448.1	255	cellulase
58	AJ36449.1	255	cellulase
59	AJ36450.1	255	cellulase
60	AJ36451.1	255	cellulase
61	AJ36452.1	255	cellulase
62	AJ36453.1	255	cellulase
63	AJ36454.1	255	cellulase
64	AJ36455.1	255	cellulase
65	AJ36456.1	255	cellulase
66	AJ36457.1	255	cellulase
67	AJ36458.1	255	cellulase
68	AJ36459.1	255	cellulase
69	AJ36460.1	255	cellulase
70	AJ36461.1	255	cellulase
71	AJ36462.1	255	cellulase
72	AJ36463.1	255	cellulase
73	AJ36464.1	255	cellulase
74	AJ36465.1	255	cellulase
75	AJ36466.1	255	cellulase
76	AJ36467.1	255	cellulase
77	AJ36468.1	255	cellulase
78	AJ36469.1	255	cellulase
79	AJ36470.1	255	cellulase
80	AJ36471.1	255	cellulase
81	AJ36472.1	255	cellulase
82	AJ36473.1	255	cellulase
83	AJ36474.1	255	cellulase
84	AJ36475.1	255	cellulase
85	AJ36476.1	255	cellulase
86	AJ36477.1	255	cellulase
87	AJ36478.1	255	cellulase
88	AJ36479.1	255	cellulase
89	AJ36480.1	255	cellulase
90	AJ36481.1	255	cellulase
91	AJ36482.1	255	cellulase
92	AJ36483.1	255	cellulase
93	AJ36484.1	255	cellulase
94	AJ36485.1	255	cellulase
95	AJ36486.1	255	cellulase
96	AJ36487.1	255	cellulase
97	AJ36488.1	255	cellulase
98	AJ36489.1	255	cellulase
99	AJ36490.1	255	cellulase
100	AJ36491.1	255	cellulase
101	AJ36492.1	255	cellulase
102	AJ36493.1	255	cellulase
103	AJ36494.1	255	cellulase
104	AJ36495.1	255	cellulase
105	AJ36496.1	255	cellulase
106	AJ36497.1	255	cellulase
107	AJ36498.1	255	cellulase
108	AJ36499.1	255	cellulase
109	AJ36500.1	255	cellulase
110	AJ36501.1	255	cellulase
111	AJ36502.1	255	cellulase
112	AJ36503.1	255	cellulase
113	AJ36504.1	255	cellulase
114	AJ36505.1	255	cellulase
115	AJ36506.1	255	cellulase
116	AJ36507.1	255	cellulase
117	AJ36508.1	255	cellulase
118	AJ36509.1	255	cellulase
119	AJ36510.1	255	cellulase
120	AJ36511.1	255	cellulase
121	AJ36512.1	255	cellulase
122	AJ36513.1	255	cellulase
123	AJ36514.1	255	cellulase
124	AJ36515.1	255	cellulase
125	AJ36516.1	255	cellulase
126	AJ36517.1	255	cellulase
127	AJ36518.1	255	cellulase
128	AJ36519.1	255	cellulase
129	AJ36520.1	255	cellulase
130	AJ36521.1	255	cellulase
131	AJ36522.1	255	cellulase
132	AJ36523.1	255	cellulase
133	AJ36524.1	255	cellulase
134	AJ36525.1	255	cellulase
135	AJ36526.1	255	cellulase
136	AJ36527.1	255	cellulase
137	AJ36528.1	255	cellulase
138	AJ36529.1	255	cellulase
139	AJ36530.1	255	cellulase
140	AJ36531.1	255	cellulase
141	AJ36532.1	255	cellulase
142	AJ36533.1	255	cellulase
143	AJ36534.1	255	cellulase
144	AJ36535.1	255	cellulase
145	AJ36536.1	255	cellulase
146	AJ36537.1	255	cellulase
147	AJ36538.1	255	cellulase
148	AJ36539.1	255	cellulase
149	AJ36540.1	255	cellulase
150	AJ36541.1	255	cellulase
151	AJ36542.1	255	cellulase
152	AJ36543.1	255	cellulase
153	AJ36544.1	255	cellulase
154	AJ36545.1	255	cellulase
155	AJ36546.1	255	cellulase
156	AJ36547.1	255	cellulase
157	AJ36548.1	255	cellulase
158	AJ36549.1	255	cellulase
159	AJ36550.1	255	cellulase
160	AJ36551.1	255	cellulase
161	AJ36552.1	255	cellulase
162	AJ36553.1	255	cellulase
163	AJ36554.1	255	cellulase
164	AJ36555.1	255	cellulase
165	AJ36556.1	255	cellulase
166	AJ36557.1	255	cellulase
167	AJ36558.1	255	cellulase
168	AJ36559.1	255	cellulase
169	AJ36560.1	255	cellulase
170	AJ36561.1	255	cellulase
171	AJ36562.1	255	cellulase
172	AJ36563.1	255	cellulase
173	AJ36564.1	255	cellulase
174	AJ36565.1	255	cellulase
175	AJ36566.1	255	cellulase
176	AJ36567.1	255	cellulase
177	AJ36568.1	255	cellulase
178	AJ36569.1	255	cellulase
179	AJ36570.1	255	cellulase
180	AJ36571.1	255	cellulase
181	AJ36572.1	255	cellulase
182	AJ36573.1	255	cellulase
183	AJ36574.1	255	cellulase
184	AJ36575.1	255	cellulase
185	AJ36576.1	255	cellulase
186	AJ36577.1	255	cellulase
187	AJ36578.1	255	cellulase
188	AJ36579.1	255	cellulase
189	AJ36580.1	255	cellulase
190	AJ36581.1	255	cellulase
191	AJ36582.1	255	cellulase
192	AJ36583.1	255	cellulase
193	AJ36584.1	255	cellulase
194	AJ36585.1	255	cellulase
195	AJ36586.1	255	cellulase
196	AJ36587.1	255	cellulase
197	AJ36588.1	255	cellulase
198	AJ36589.1	255	cellulase
199	AJ36590.1	255	cellulase
200	AJ36591.1	255	cellulase
201	AJ36592.1	255	cellulase
202	AJ36593.1	255	cellulase
203	AJ36594.1	255	cellulase
204	AJ36595.1	255	cellulase
205	AJ36596.1	255	cellulase
206	AJ36597.1	255	cellulase
207	AJ36598.1	255	cellulase
208	AJ36599.1	255	cellulase
209	AJ36600.1	255	cellulase
210	AJ36601.1	255	cellulase
211	AJ36602.1	255	cellulase
212	AJ36603.1	255	cellulase
213	AJ36604.1	255	cellulase
214	AJ36605.1	255	cellulase
215	AJ36606.1	255	cellulase
216	AJ36607.1	255	cellulase
217	AJ36608.1	255	cellulase
218	AJ36609.1	255	cellulase
219	AJ36610.1	255	cellulase
220	AJ36611.1	255	cellulase
221	AJ36612.1	255	cellulase
222	AJ36613.1	255	cellulase
223	AJ36614.1	255	cellulase
224	AJ36615.1	255	cellulase
225	AJ36616.1	255	cellulase
226	AJ36617.1	255	cellulase
227	AJ36618.1	255	cellulase
228	AJ36619.1	255	cellulase
229	AJ36620.1	255	cellulase
230	AJ36621.1	255	cellulase
231	AJ36622.1	255	cellulase
232	AJ36623.1	255	cellulase
233	AJ36624.1	255	cellulase
234	AJ36625.1	255	cellulase
235	AJ36626.1	255	cellulase
236	AJ36627.1	255	cellulase
237	AJ36628.1	255	cellulase
238	AJ36629.1	255	cellulase
239	AJ36630.1	255	cellulase
240	AJ36631.1	255	cellulase
241	AJ36632.1	255	cellulase
242	AJ36633.1	255	cellulase
243	AJ36634.1	255	cellulase
244	AJ36635.1	255	cellulase
245	AJ36636.1	255	cellulase
246	AJ36637.1	255	cellulase
247	AJ36638.1	255	cellulase
248	AJ36639.1	255	cellulase
249	AJ36640.1	255	cellulase
250	AJ36641.1	255	cellulase
251	AJ36642.1	255	cellulase
252	AJ36643.1	255	cellulase
253	AJ36644.1	255	cellulase
254	AJ36645.1	255	cellulase
255	AJ36646.1	255	cellulase
256	AJ36647.1	255	cellulase
257	AJ36648.1	255	cellulase
258	AJ36649.1	255	cellulase
259	AJ36650.1	255	cellulase
260	AJ36651.1	255	cellulase
261	AJ36652.1	255	cellulase
262	AJ36653.1	255	cellulase
263	AJ36654.1	255	cellulase
264	AJ36655.1	255	cellulase
265	AJ36656.1	255	cellulase
266	AJ36657.1	255	cellulase
267	AJ36658.1	255	cellulase
268	AJ36659.1	255	cellulase
269	AJ36660.1	255	cellulase
270	AJ36661.1	255	cellulase
271	AJ36662.1	255	

AKP69887.1 (1096 aa):
hypothetical protein CG08_1741 [Riemerella anatispestifer]

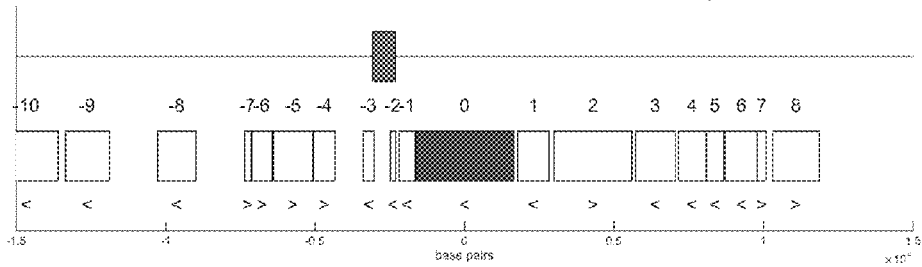


CRISPR
classifications

<Classification by Repeats>

36 nt repeat Type II, IV

CP007503.1 (2152519 bp [1673683.5-1703683.5]):
DEFINITION Riemerella anatispestifer strain 17, complete genome.



CRISPR arrays in genome

Dist	# Repeat	Spacer	Consensus
658465	47	80	GTTCGCAATTCGGTTCGAATTTTACATATGTTTGTGATAGTTTGGCAAC
872	32	50	CTTGTGCAATGAGCCCTTCAGAAATGAGGAGCCAGTGGTCAAC

Genes nearby protein

Index	Protein	Protein Length	Protein Name
-10	AKP69877.1	705	ionb-dependent receptor
-8	AKP69870.1	439	RNA polymerase, sigma 34 subunit, rpoN/σ ³⁴
-6	AKP69878.1	428	osafyl-CeA hydrolase/esterase
-7	AKP69860.1	70	hypothetical protein
-6	AKP69851.1	228	hypothetical protein
-5	AKP69852.1	442	hypothetical protein
-4	AKP69853.1	239	hypothetical protein
-3	AKP69854.1	124	hypothetical protein
-2	AKP69845.1	80	hypothetical protein
-1	AKP69855.1	178	hypothetical protein
0	AKP69887.1	1096	hyp. bacterial protei
1	AKP69862.1	343	glucokinase
2	AKP69889.1	867	alanine-tRNA synthetase
3	AKP69890.1	445	pac-type peptidyl-prolyl cis-trans isomerase
4	AKP69891.1	308	pho domain-containing protein
5	AKP69892.1	127	gliding motility protein GldD
6	AKP69893.1	371	MG domain; DNA gyrasease
7	AKP69894.1	97	nistone family protein DNA-binding protein
8	AKP69895.1	519	ribonuclease, omeiga family

FIG. 4N

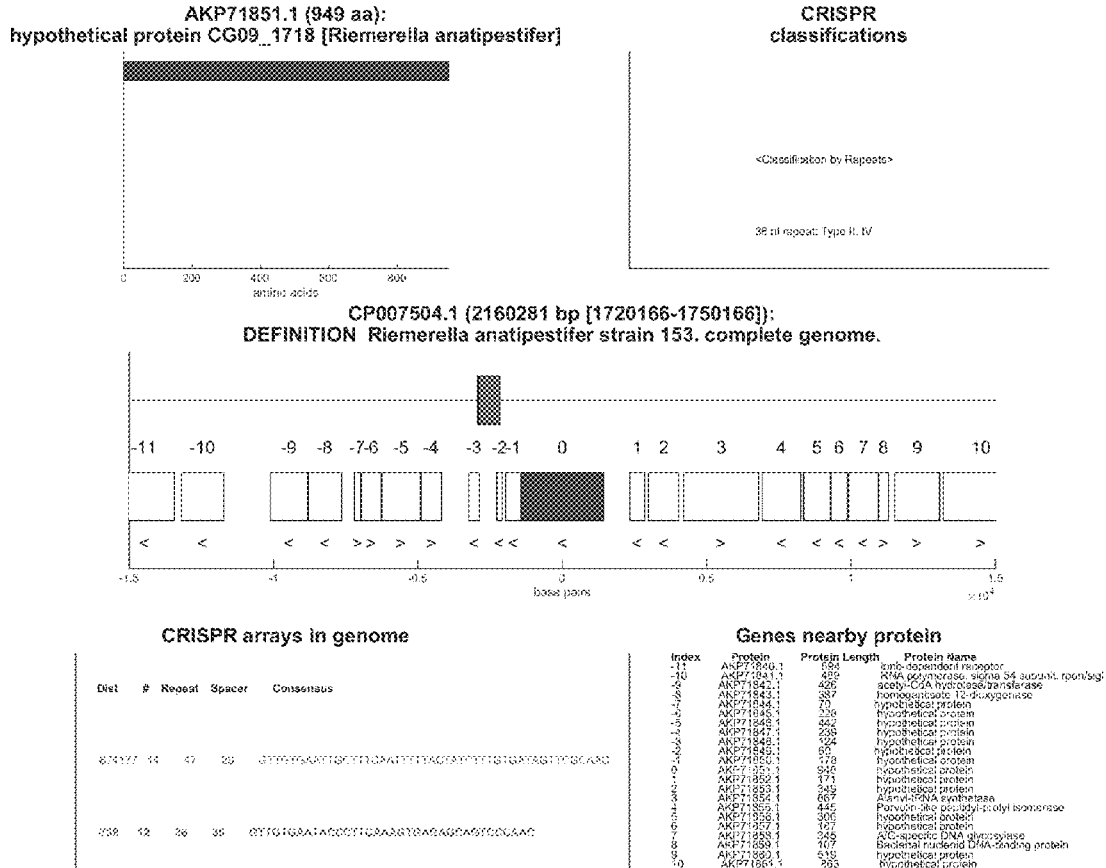


FIG. 40

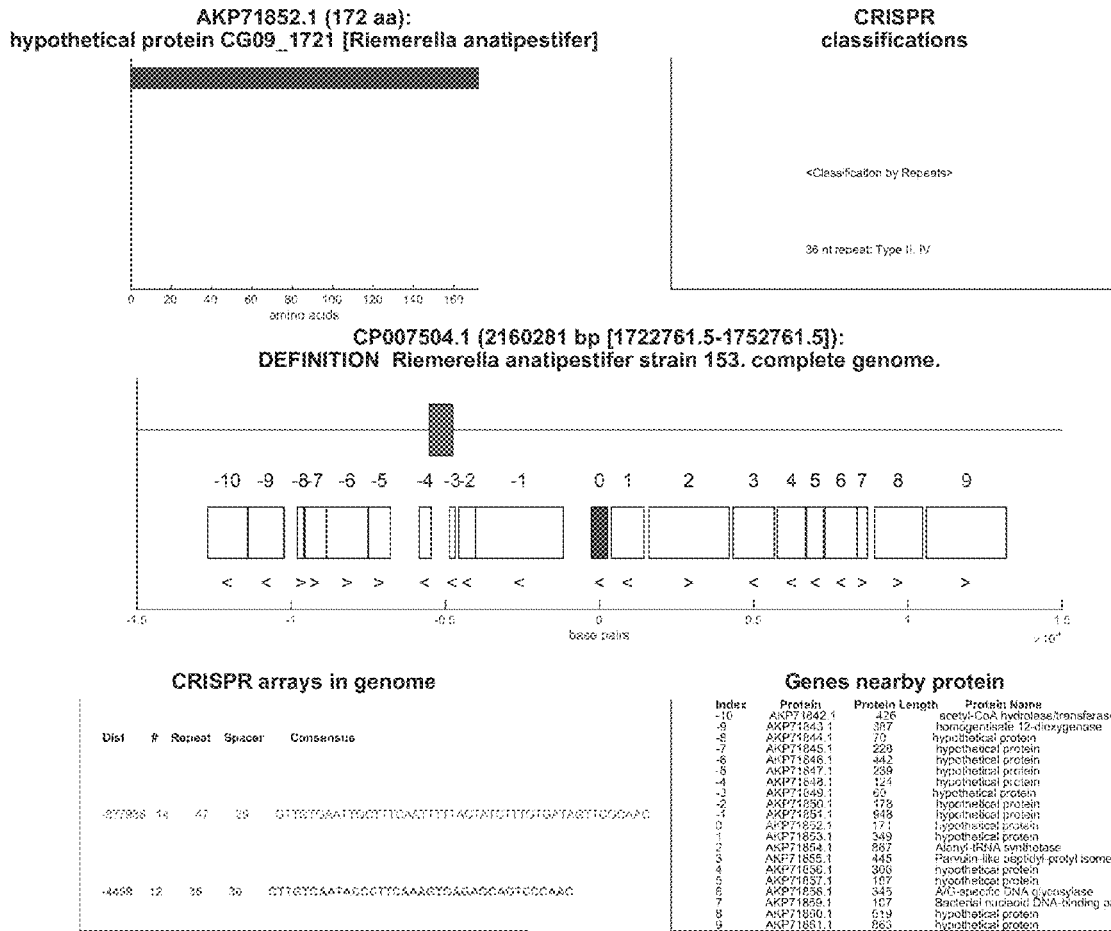


FIG. 4P

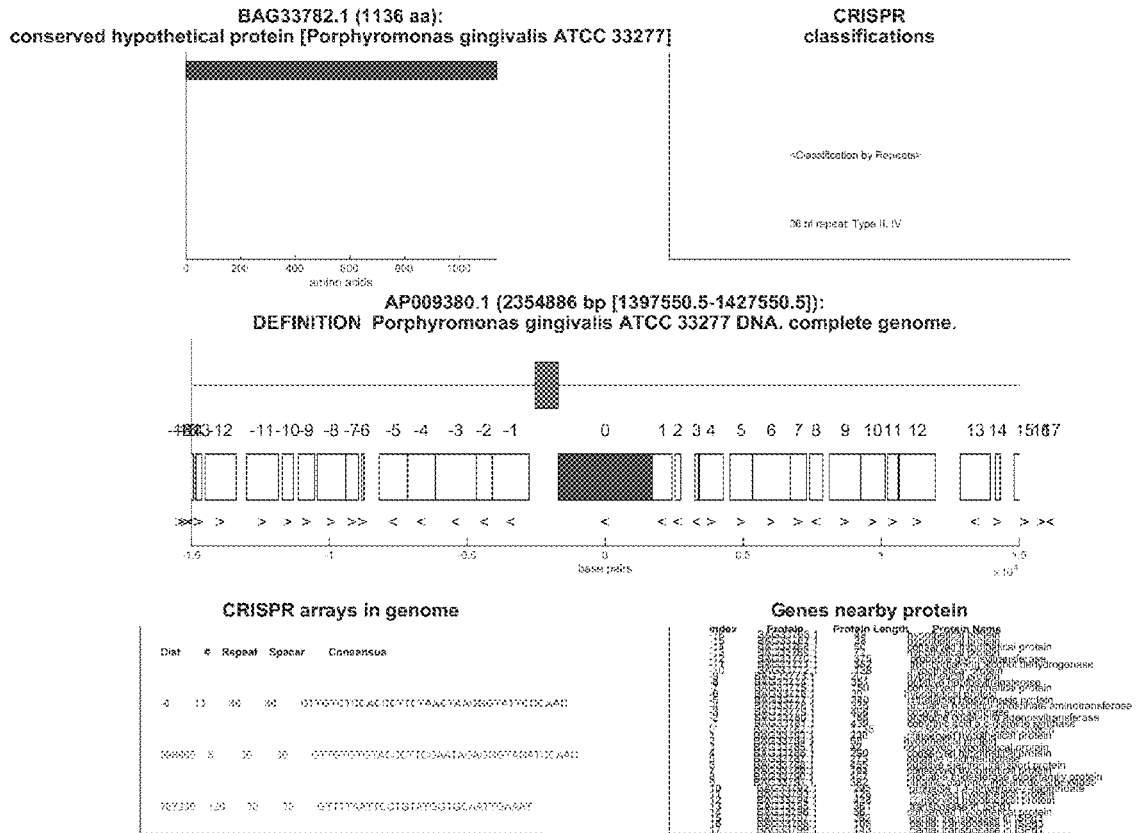


FIG. 4T

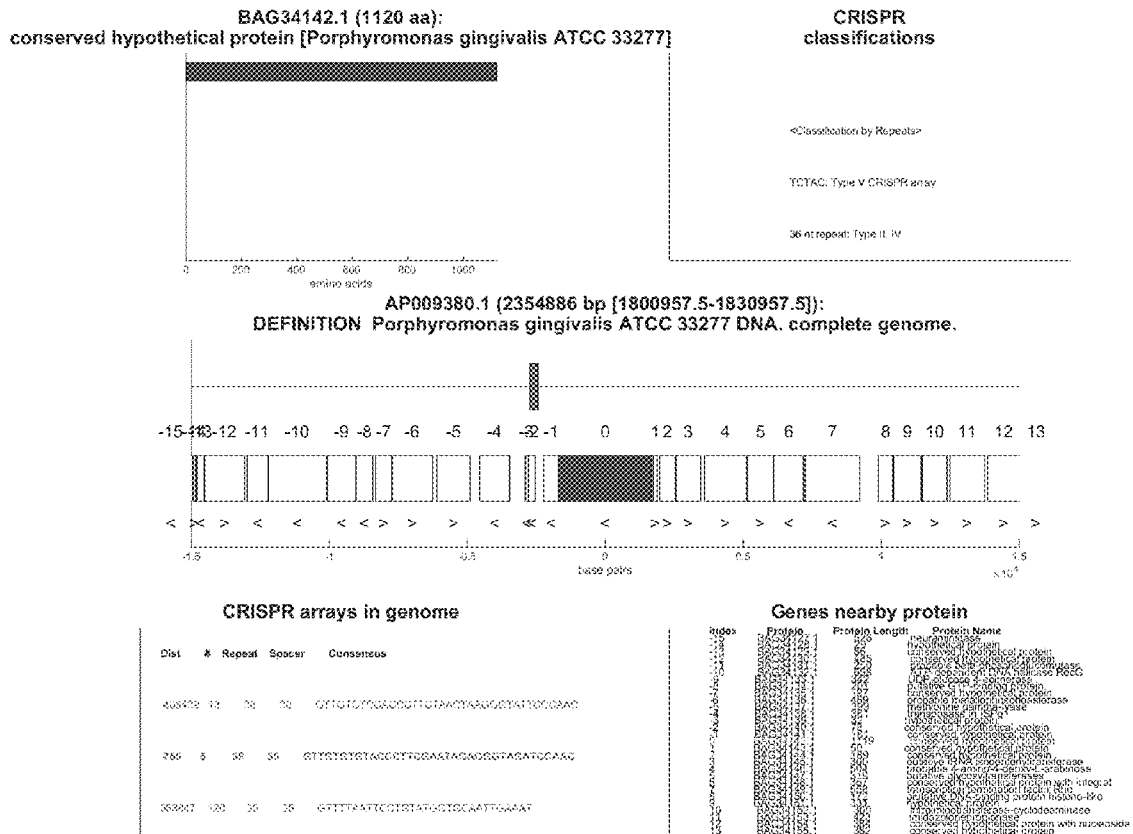
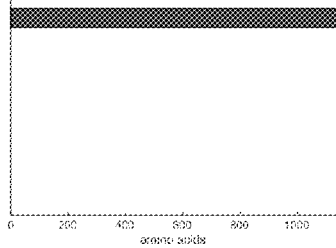


FIG. 4U

BAR96998.1 (1133 aa):
hypothetical protein PI172_2270 [Prevotella intermedia]



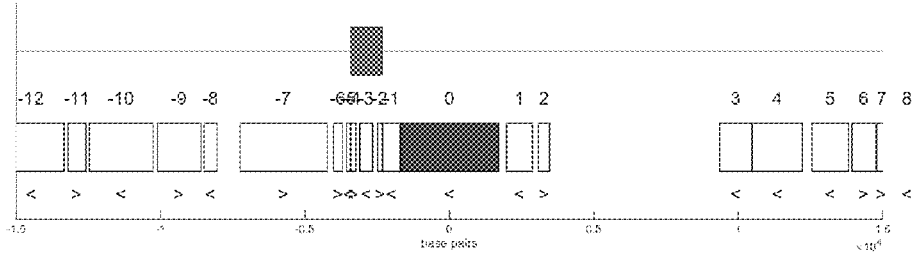
CRISPR
classifications

<Classification by Repeats>

GTTC CRISPR array

36 nt repeat Type II, IV

AP014926.1 (606227 bp [328581-358581]):
DEFINITION Prevotella intermedia DNA, chromosome 1, complete genome, strain:



CRISPR arrays in genome

Dist	# Repeat	Spacer	Consensus	
610	17	36	25	GTGTTTTCAGCTTTCAGACAGCAAGGCGAATGCAATG

Genes nearby protein

Index	Protein	Protein Length	Protein Name
-12	BAR96986.1	762	NAD(P)-dependent methylase
-11	BAR96987.1	205	transcriptional regulator, TetR family
-10	BAR96988.1	732	GTP diphosphokinase, (putative) synthetase 1
-9	BAR96989.1	905	nicotinamide adenine dinucleotide synthetase
-8	BAR96990.1	133	hypothetical protein
-7	BAR96991.1	1006	hypothetical protein
-6	BAR96992.1	102	hypothetical protein
-5	BAR96993.1	42	hypothetical protein
-4	BAR96994.1	56	hypothetical protein
-3	BAR96995.1	147	hypothetical protein
-2	BAR96996.1	53	hypothetical protein
-1	BAR96997.1	201	hypothetical protein
0	BAR96998.1	1133	hypothetical protein
1	BAR96999.1	287	hypothetical protein
2	BAR97000.1	122	hypothetical protein
3	BAR97001.1	271	hypothetical protein
4	BAR97002.1	575	hypothetical protein
5	BAR97003.1	425	mobile element protein
6	BAR97004.1	276	mobile element protein
7	BAR97005.1	108	transposase
8	BAR97006.1	108	IS1257, transposase

FIG. 4W

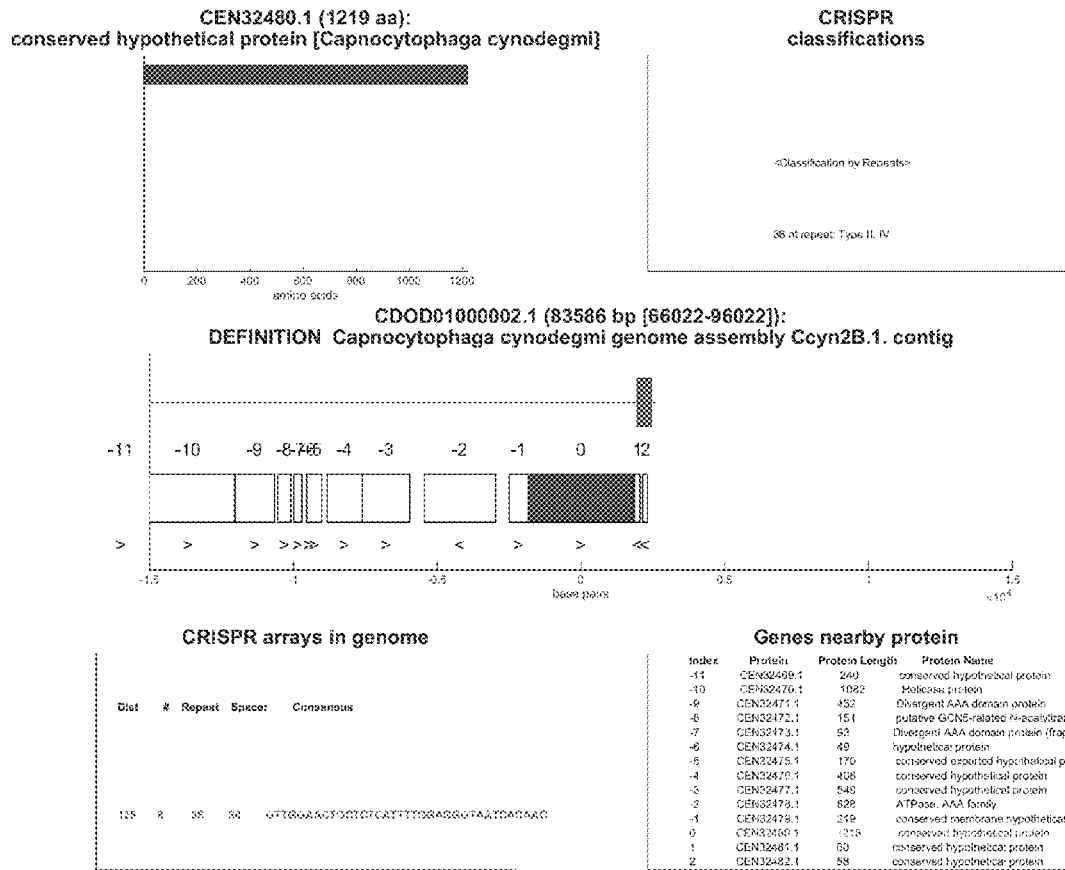


FIG. 4X

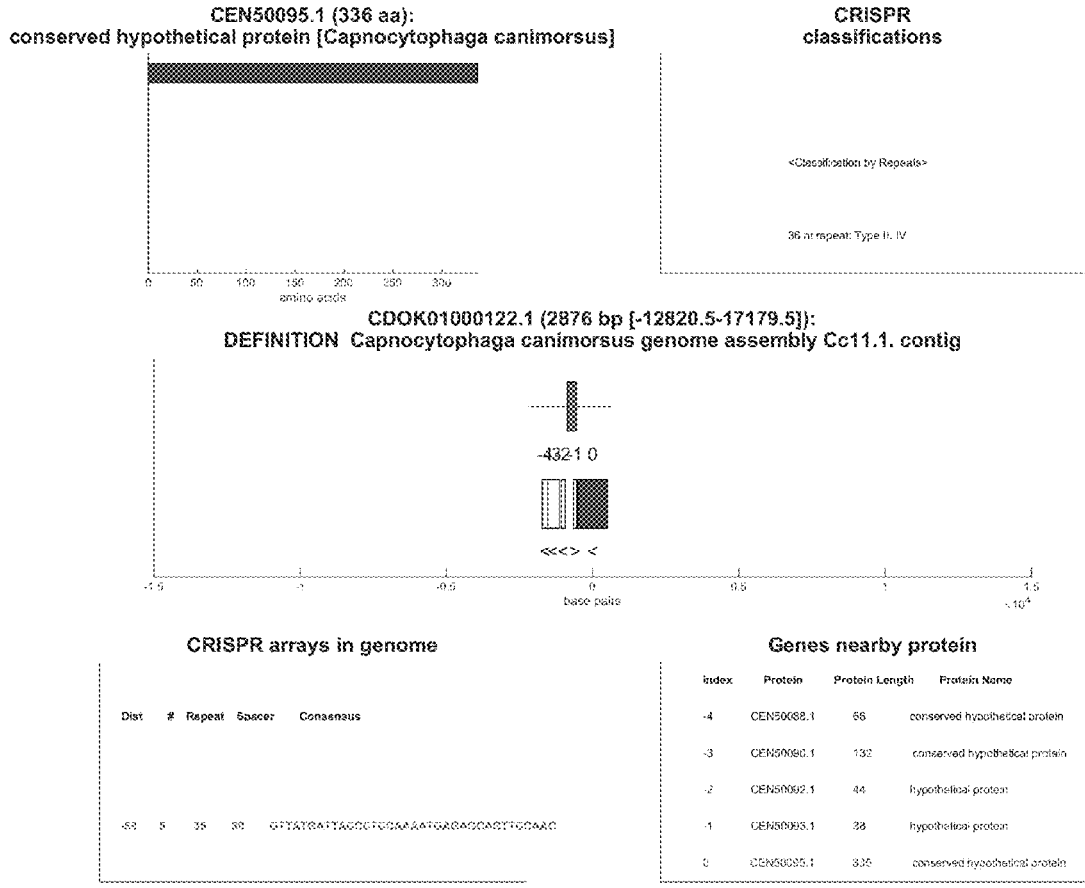


FIG. 4Y

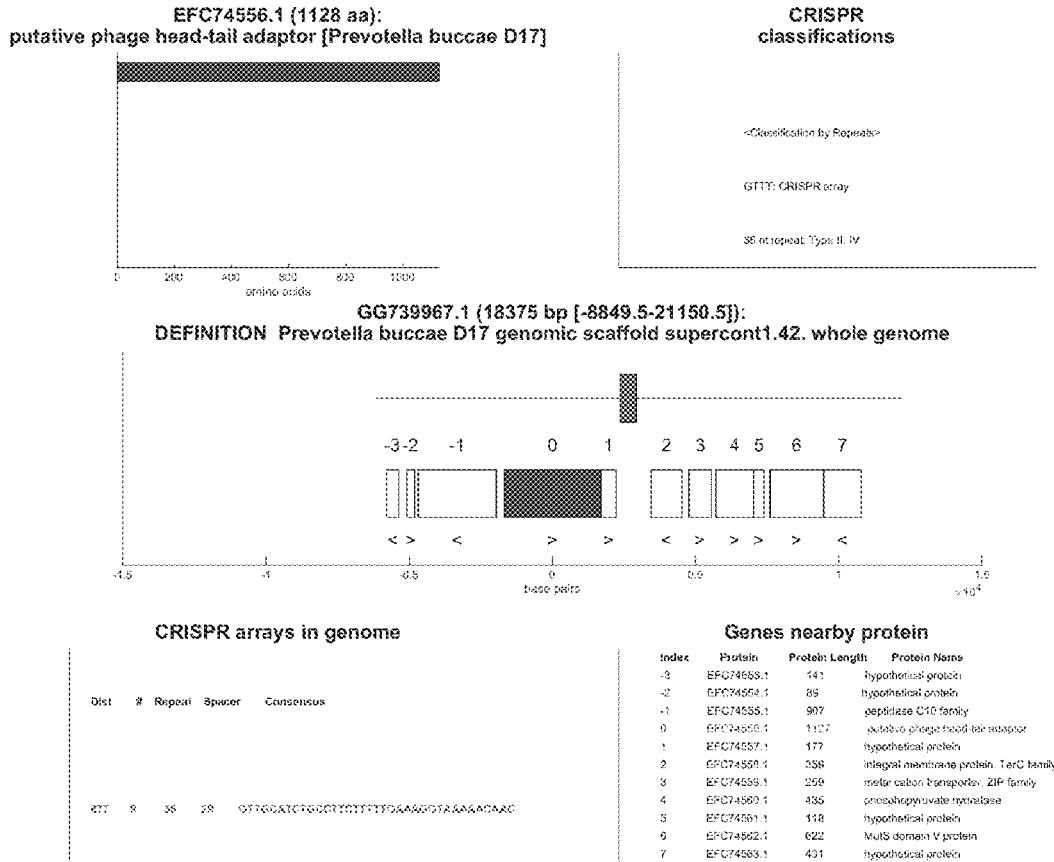
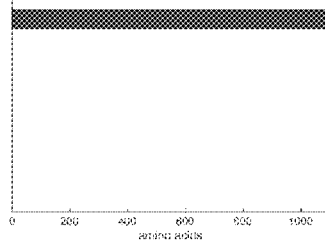


FIG. 4Z

EFT35596.1 (1096 aa):
hypothetical protein RAYM_05191 [Riemerella anatipestifer RA-YM]

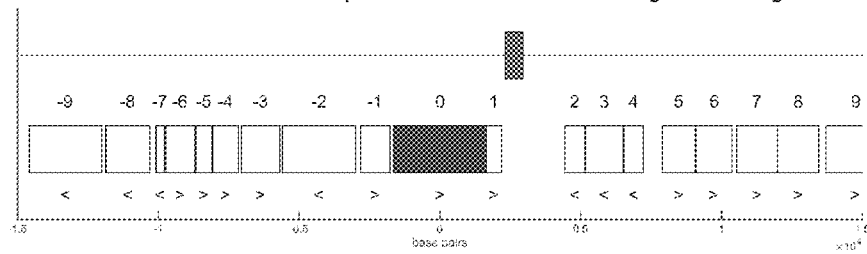


CRISPR
classifications

<Classification by Repeats>

30 rd repeat, Type II, IV

AENH0100026.1 (219800 bp [188684.5-219884.5]):
DEFINITION Riemerella anatipestifer RA-YM scaffold6, whole genome shotgun



CRISPR arrays in genome

Dist	#	Repeat	Spacer	Consensus	
	602	16	36	36	GTTRGGAGTGCCTTCATTTTAAJGGGTATTACAAAC

Genes nearby protein

Index	Protein	Protein Length	Protein Name
-9	EFT35597.1	636	hypothetical protein
-8	EFT35596.1	1096	hypothetical protein
-7	EFT35595.1	67	ribonuclease G
-6	EFT35594.1	348	RNA-binding protein HU
-5	EFT35593.1	127	ATC-specific adenosine glycosylase
-4	EFT35592.1	127	SDS
-3	EFT35591.1	308	ribonuclease
-2	EFT35590.1	455	PpC-type peptidyl-prolyl isomerase
-1	EFT35589.1	287	glycyl-tRNA synthetase
0	EFT35588.1	345	glucosylase
1	EFT35587.1	1015	hypothetical protein
2	EFT35586.1	173	hypothetical protein
3	EFT35585.1	239	hypothetical protein
4	EFT35584.1	442	rasc_1tag family protein
5	EFT35583.1	228	hypothetical protein
6	EFT35582.1	387	teicoplaninase 1,2-dioxygenase
7	EFT35581.1	428	4-hydroxybutyrate coenzyme A transferase
8	EFT35580.1	452	asparaginyl-tRNA synthetase
9	EFT35579.1	489	NWA polymerase sigma-54 factor rpoN
10	EFT35578.1	705	TanR-dependent outer membrane receptor

FIG. 4AA

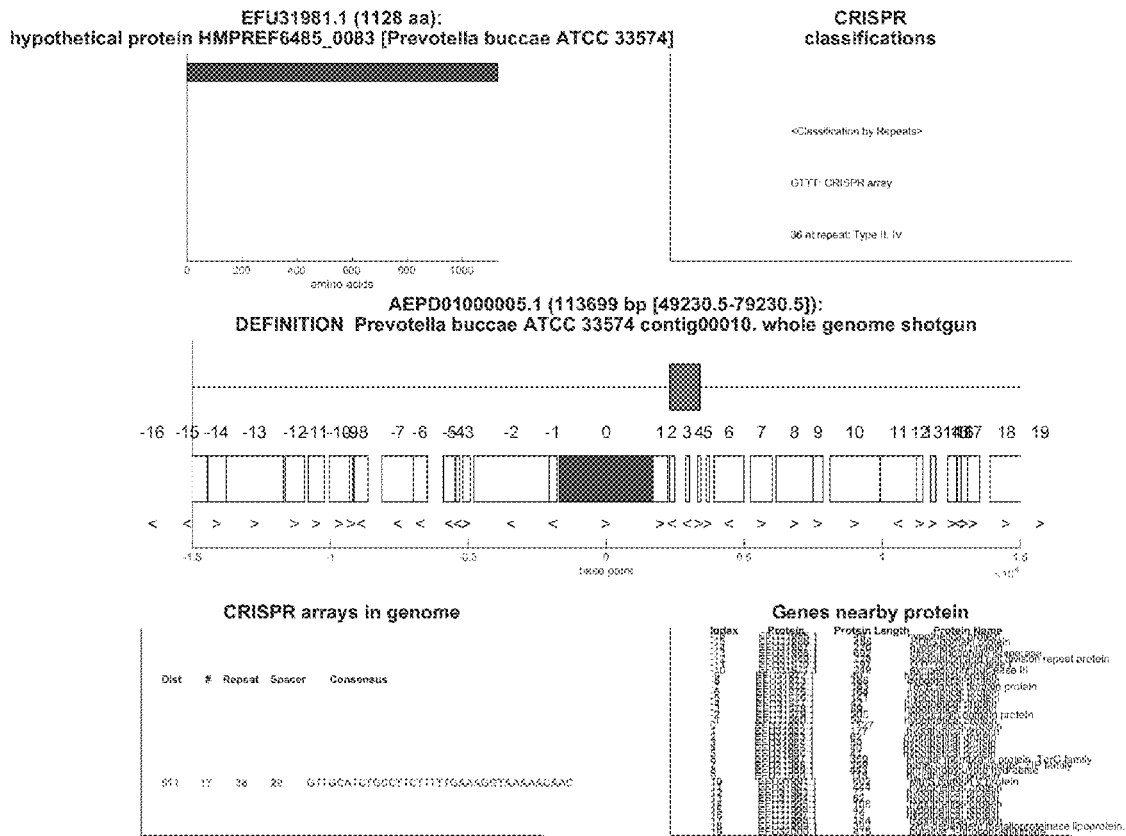


FIG. 4BB

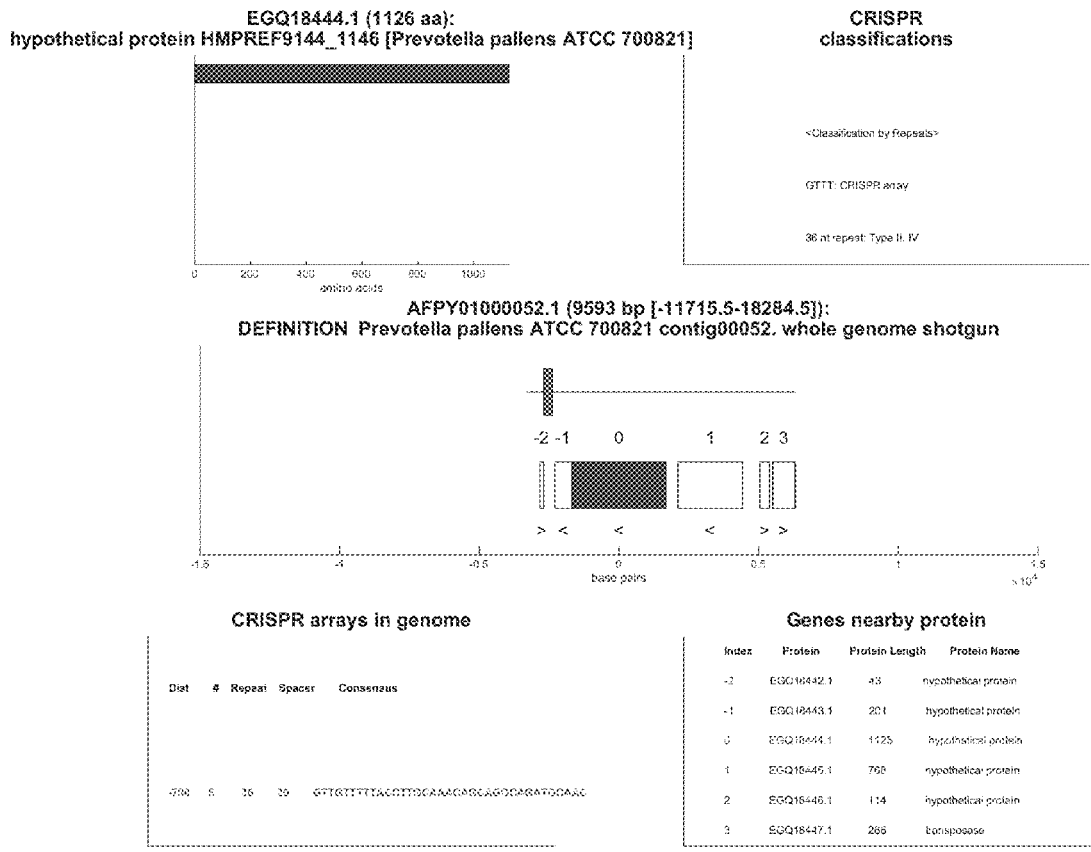


FIG. 4CC

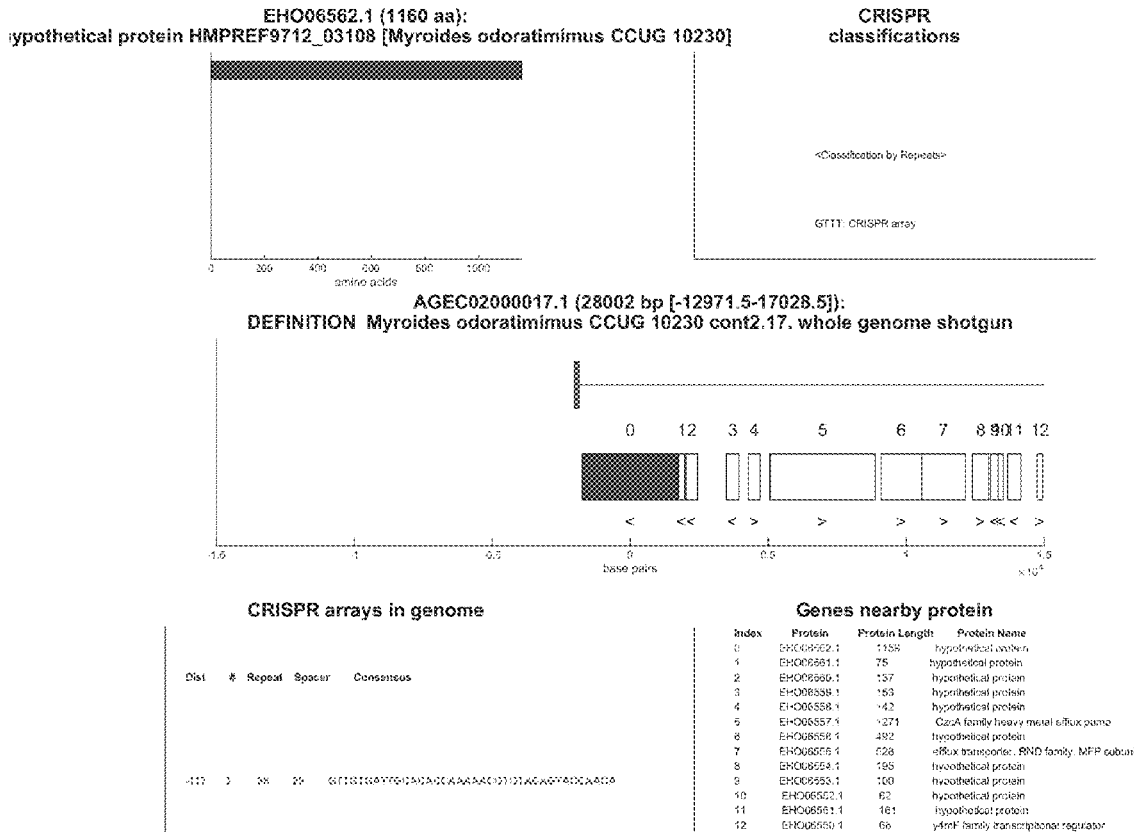


FIG. 4DD

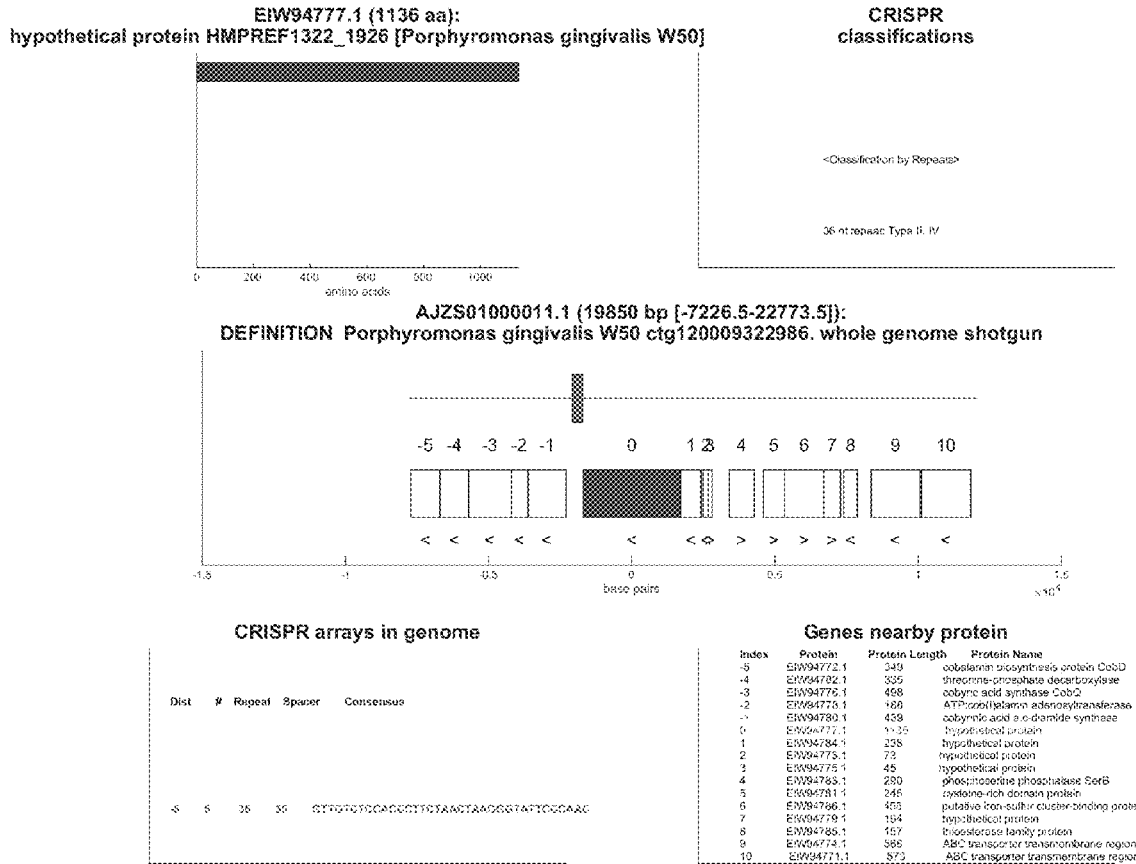


FIG. 4FF

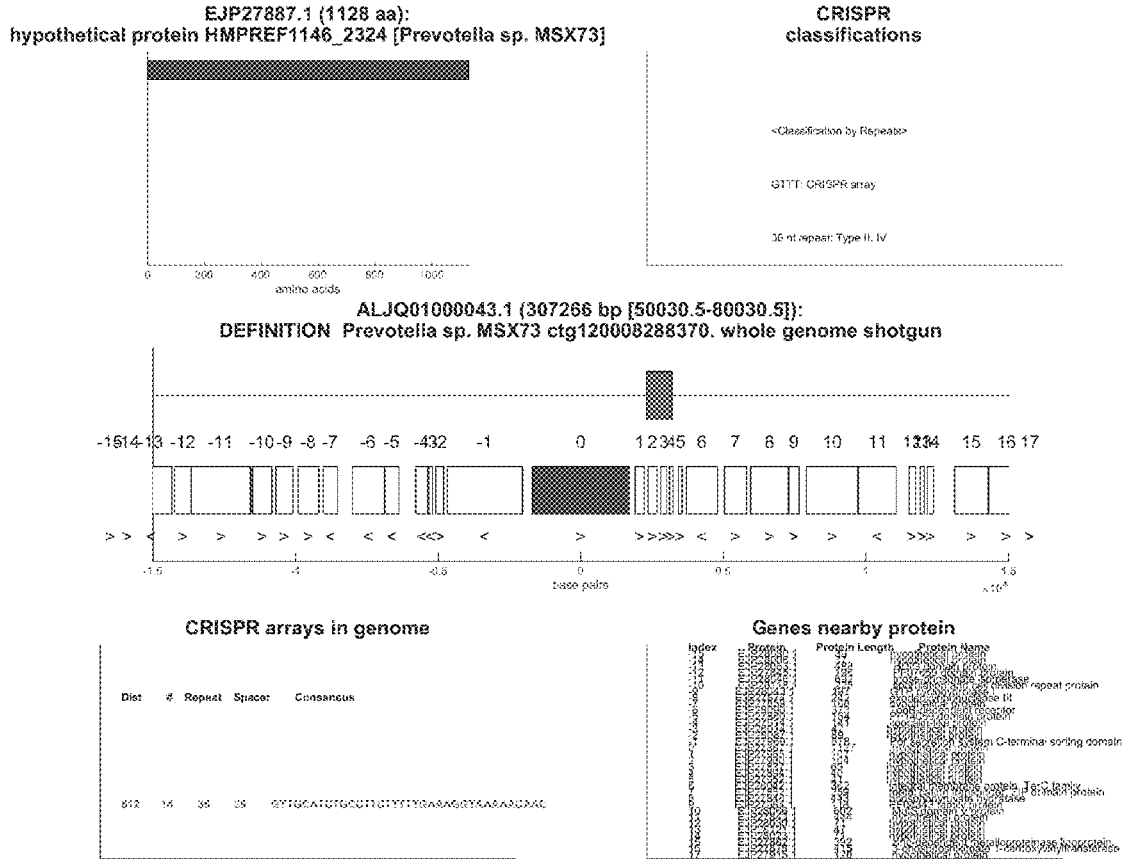


FIG. 4GG

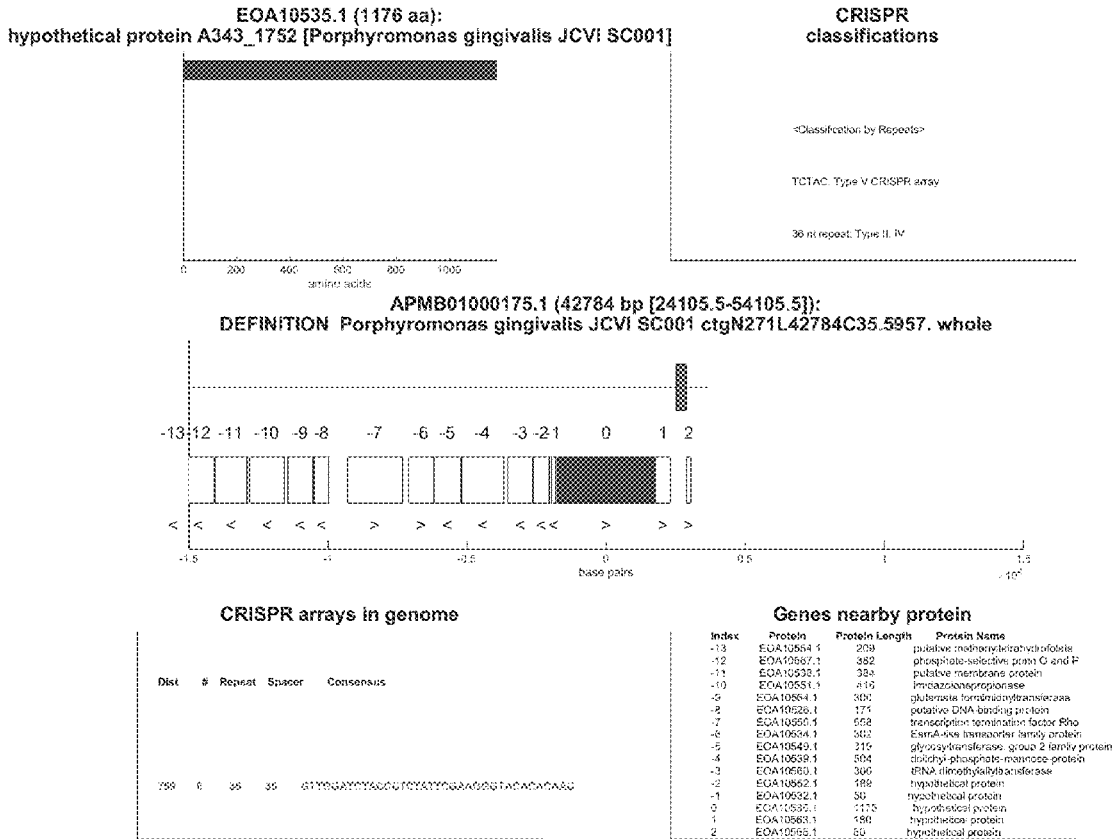


FIG. 4II

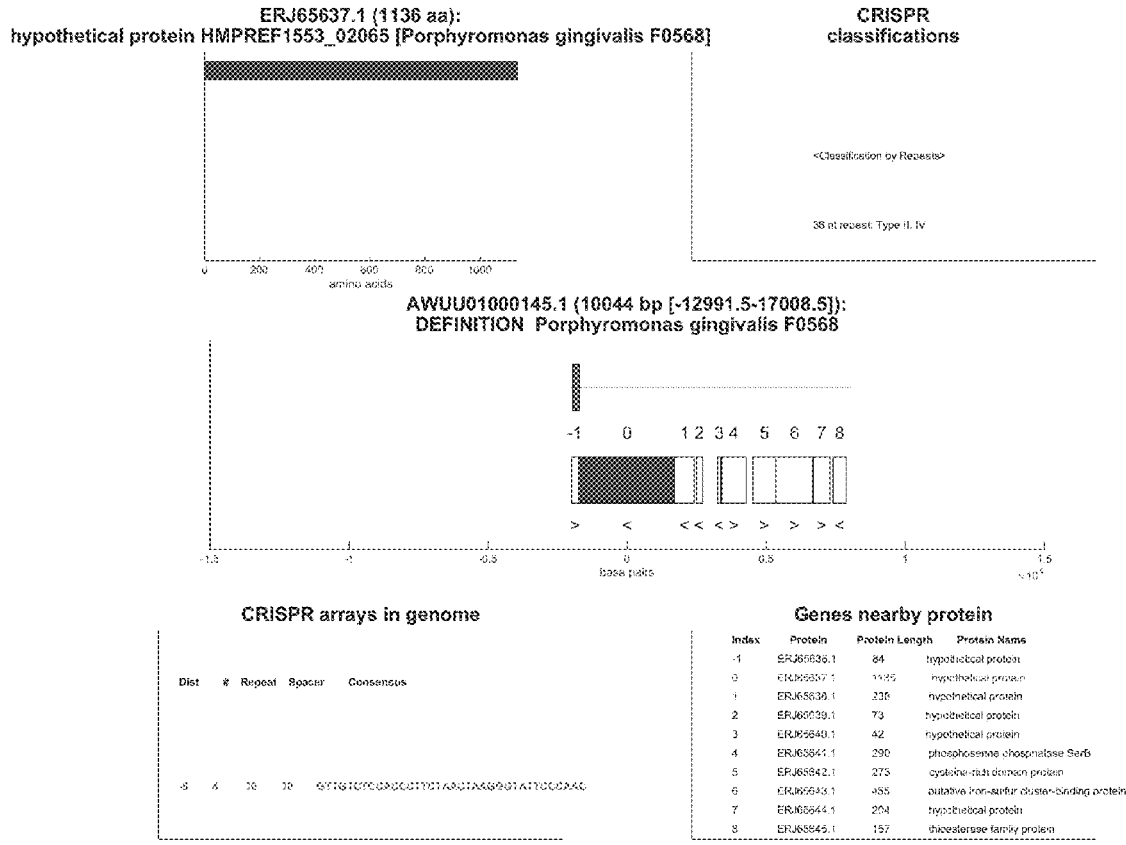


FIG. 4KK

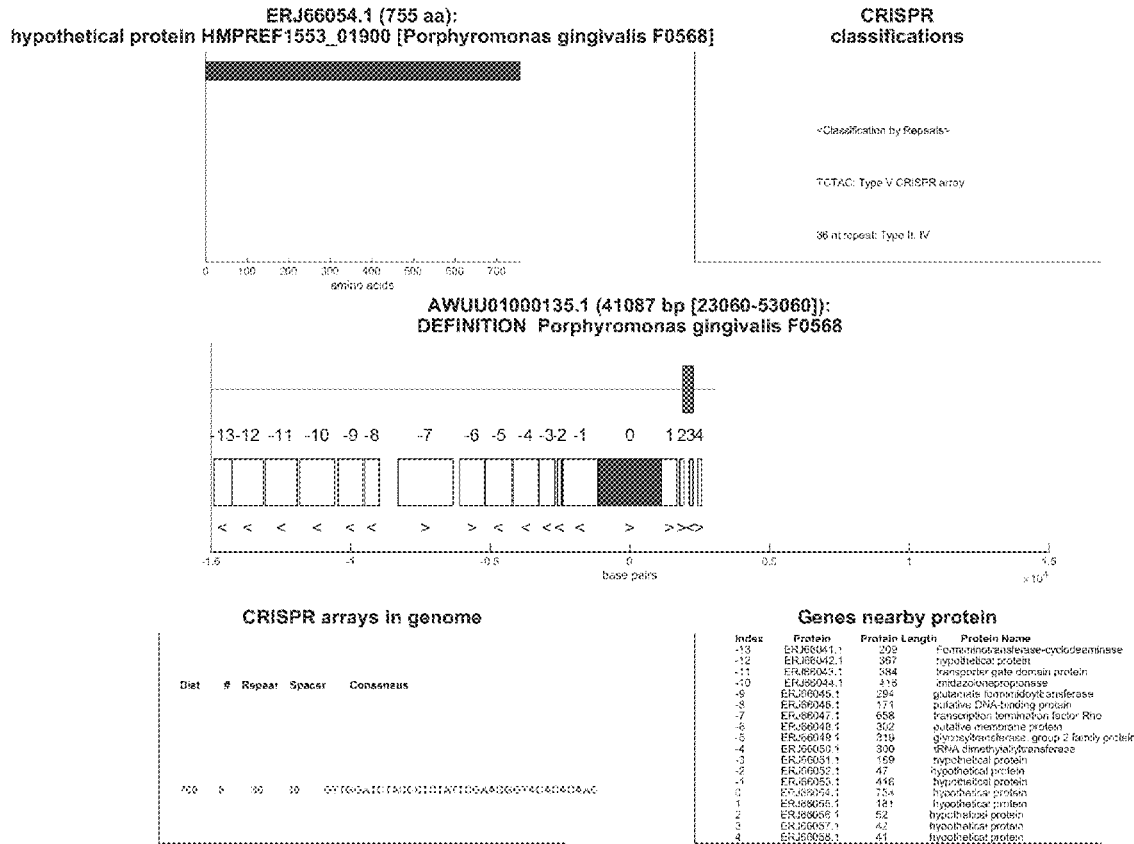


FIG. 4LL

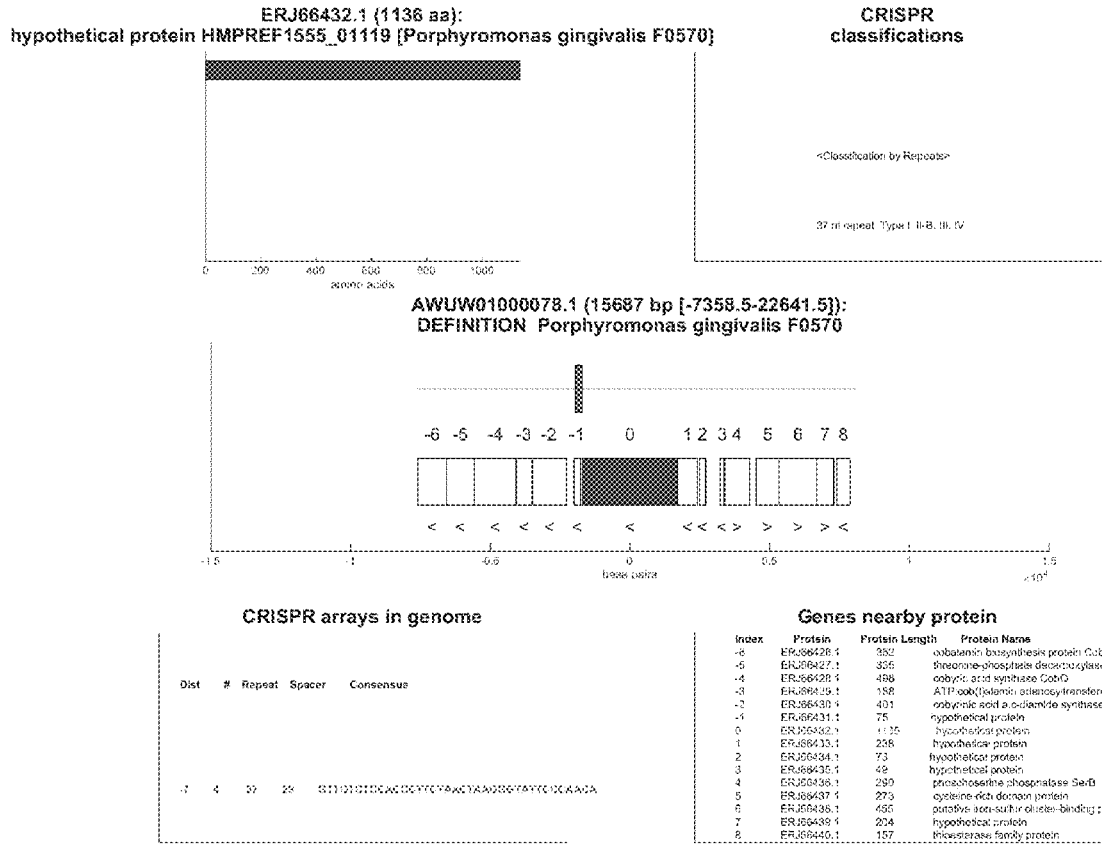


FIG. 4MM

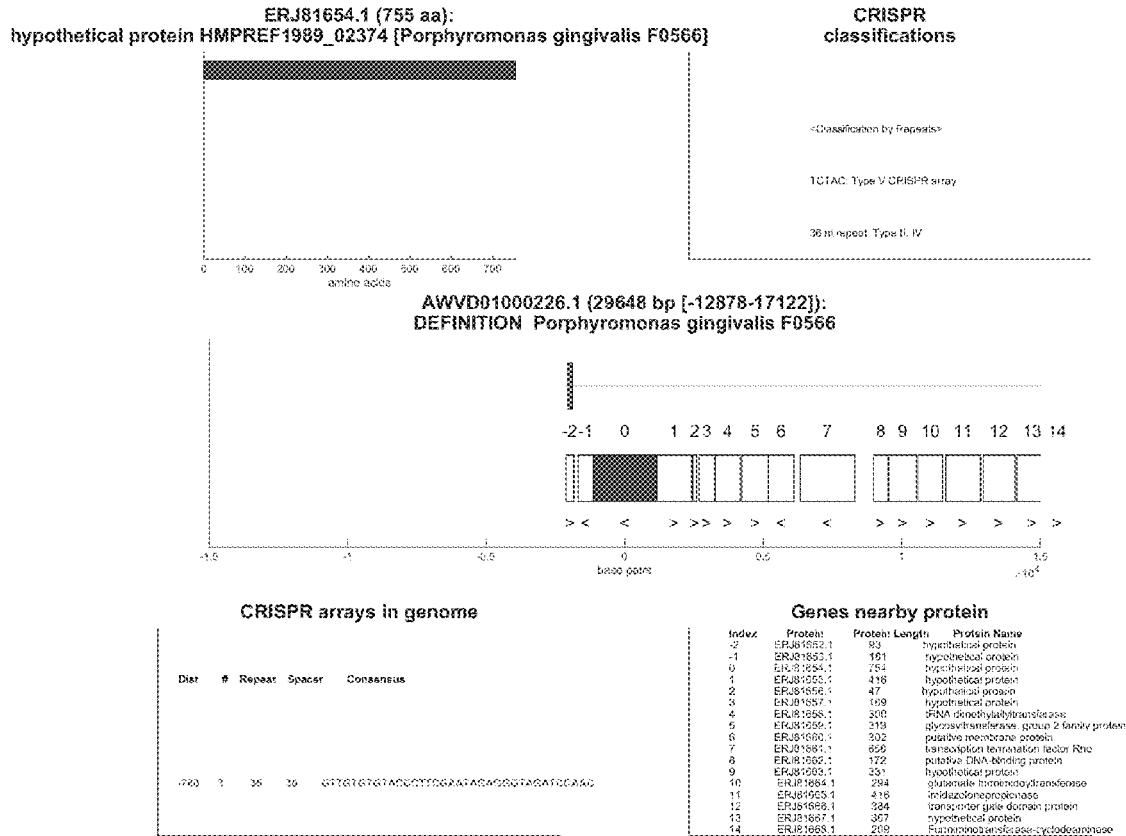


FIG. 4NN

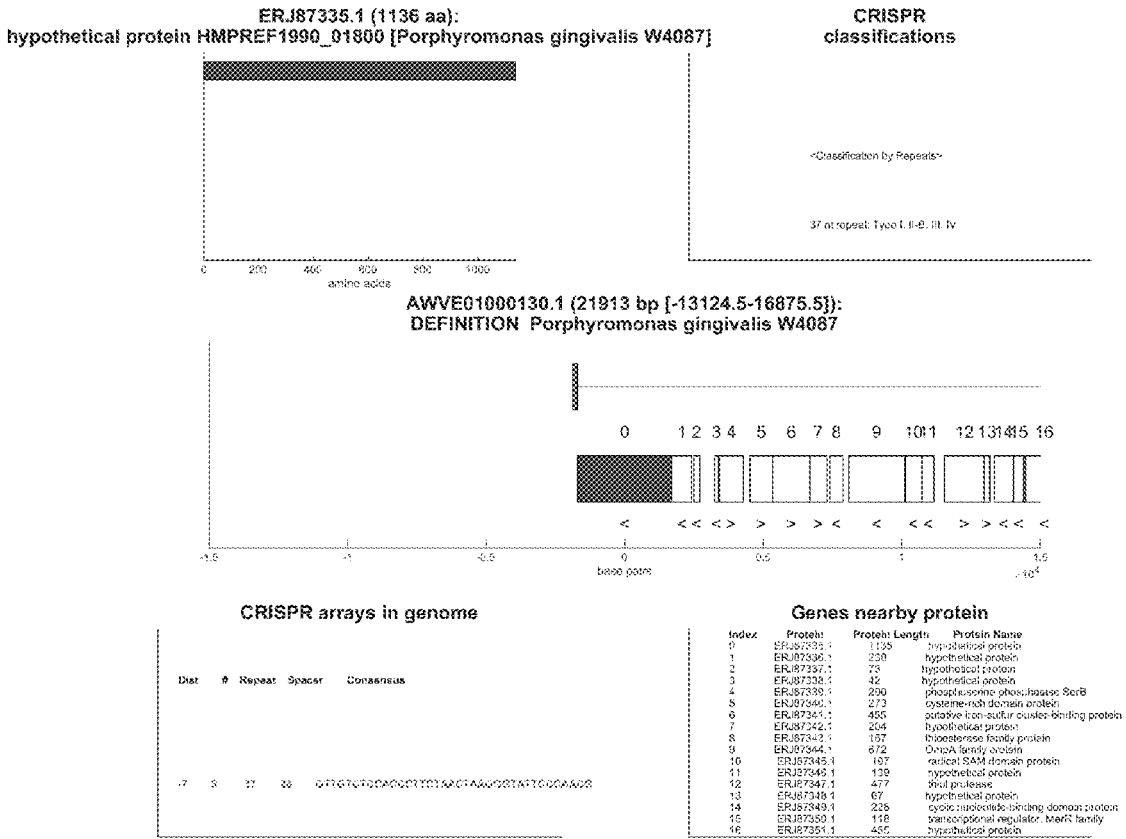
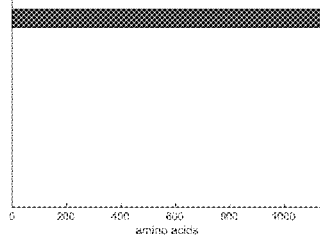


FIG. 400

ERJ98772.1 (1140 aa):
hypothetical protein HMPREF1218_0639 [Prevotella pleuritidis F0068]



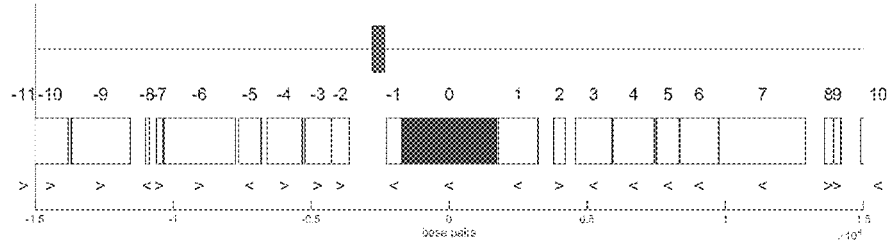
CRISPR
classifications

<Classification by Repeat>

GTTT: CRISPR array

DE in repeat: Type II, IV

AWET01000045.1 (167660 bp [24484.5-54484.5]):
DEFINITION Prevotella pleuritidis F0068 contig00004. whole genome shotgun



CRISPR arrays in genome

Dist	#	Repeat	Spacer	Consensus
-442	7	80	80	GTGTTTSTGCGCTGTGAATAGGAAGGCAAGTGCAGC

Genes nearby protein

Index	Protein	Protein Length	Protein Name
1	ER_088118.1	215	putative lipoprotein
2	ER_088119.1	413	putative lipoprotein, PFD1636 family
3	ER_088120.1	713	PP1C-type 1 PHASE domain protein
4	ER_088121.1	46	hypothetical protein
5	ER_088122.1	77	hypothetical protein
6	ER_088123.1	46	hypothetical protein
7	ER_088124.1	684	hypothetical protein
8	ER_088125.1	384	hypothetical protein
9	ER_088126.1	513	Myosin XXIV A1/Pass family protein
10	ER_088127.1	519	4-oxo-3-oxoacyl-CoA synthetase
11	ER_088128.1	519	putative membrane protein
12	ER_088129.1	186	hypothetical protein
13	ER_088130.1	119	hypothetical protein
14	ER_088131.1	119	hypothetical protein
15	ER_088132.1	140	transceptor, ABC family
16	ER_088133.1	140	transposase, DDE domain protein
17	ER_088134.1	438	ribonop
18	ER_088135.1	493	hypothetical protein
19	ER_088136.1	263	putative lipoprotein
20	ER_088137.1	470	sporobactin protein, SusD-like family
21	ER_088138.1	1544	Toxin-antitoxin (membrane protein), SusCRegA
22	ER_088139.1	182	zinc-finger protein 12
23	ER_088140.1	87	ribosomal protein L27
24	ER_088141.1	414	DENR/DENR3 helicase

FIG. 4PP

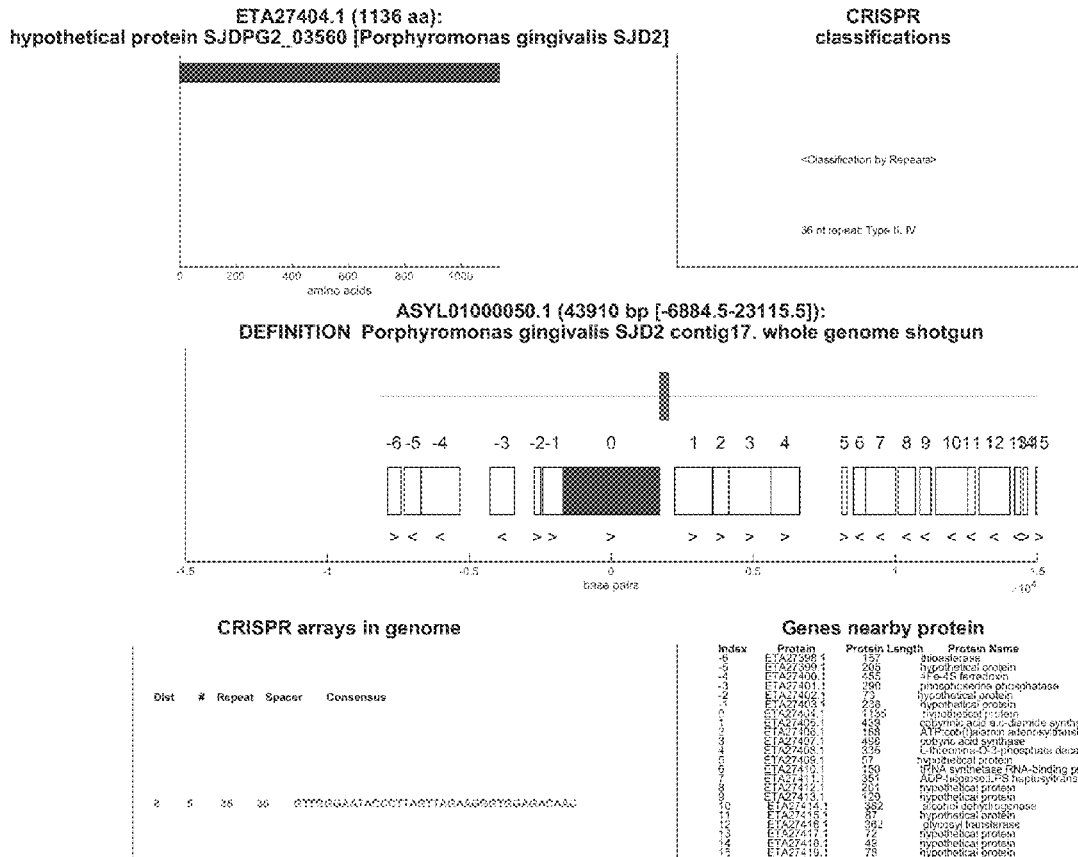


FIG. 4QQ

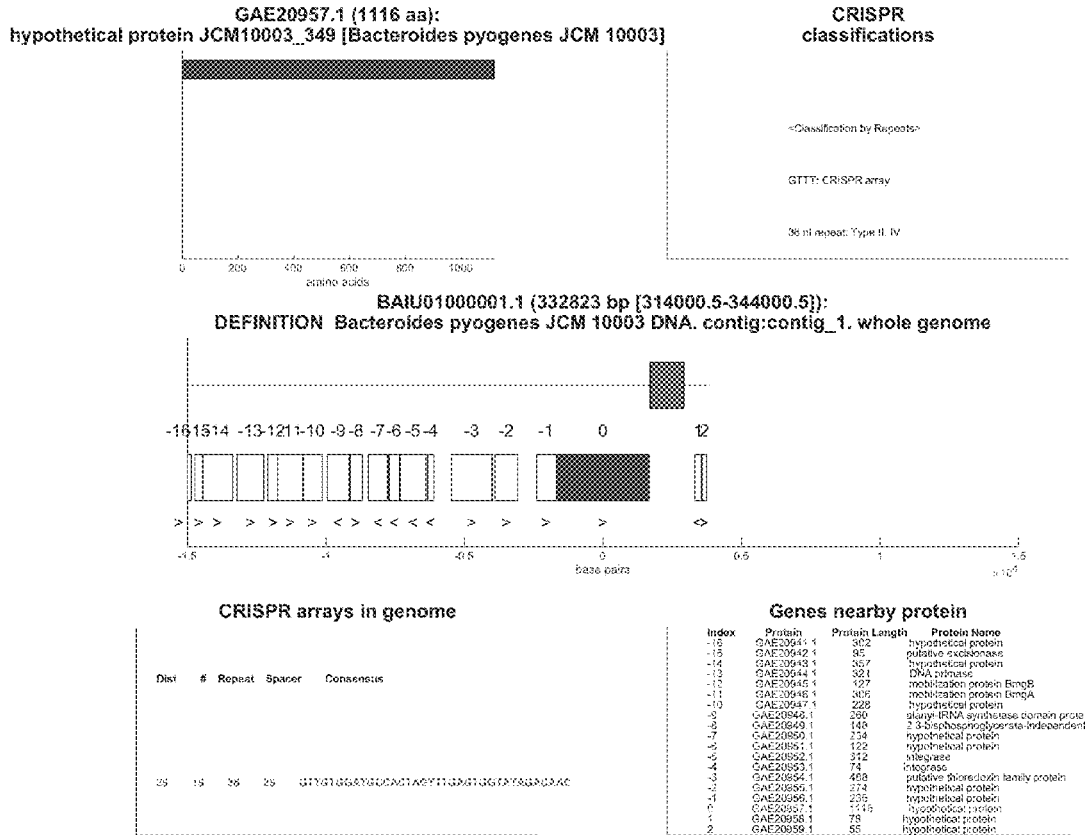
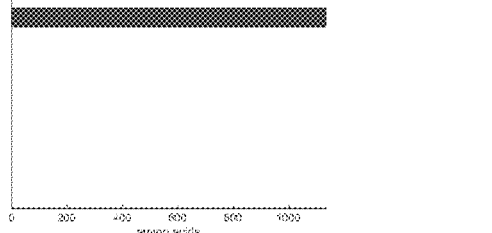


FIG. 4RR

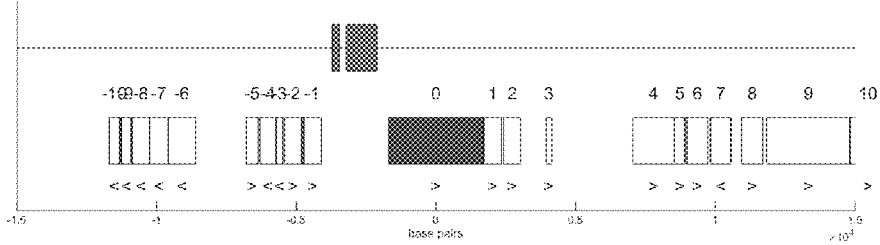
KGE88582.1 (1132 aa):
hypothetical protein IX84_07840 [Phaeodactylibacter xiamenensis]



CRISPR
classifications

<Classification by Repeats>
GTTT: CRISPR array
36 nt repeat: Type II: IV
36 nt repeat: Type II: IV

JPOS01000018.1 (442685 bp [169162.5-199162.5]):
DEFINITION Phaeodactylibacter xiamenensis strain KD52 Scaffold1_18. whole



CRISPR arrays in genome

Dist	#	Repeat	Spacer	Consensus
477	4	34	31	GTCTGAGTTCATCTGCTCAGAGCGGAGGCGCTATCGGAGC
412	10	30	30	GTCGTAAGTTCATCTGCTCAGAGCGGAGGCGCTATCGGAGC

Genes nearby protein

Index	Protein	Protein Length	Protein Name
-10	KGE88372.1	226	hypothetical protein
-9	KGE88371.1	112	hypothetical protein
-8	KGE88374.1	208	16S rRNA methyltransferase
-7	KGE88375.1	228	hypothetical protein
-6	KGE88376.1	357	family 2 glyoxyl transferase
-5	KGE88377.1	137	hypothetical protein
-4	KGE88378.1	192	hypothetical protein
-3	KGE88379.1	85	hypothetical protein
-2	KGE88380.1	207	hypothetical protein
-1	KGE88381.1	211	hypothetical protein
0	KGE88382.1	1151	hypothetical protein
1	KGE88383.1	211	hypothetical protein
2	KGE88384.1	207	hypothetical protein
3	KGE88385.1	86	hypothetical protein
4	KGE88386.1	459	hypothetical protein
5	KGE88387.1	123	hypothetical protein
6	KGE88388.1	251	hypothetical protein
7	KGE88389.1	246	hypothetical protein
8	KGE88390.1	264	hypothetical protein
9	KGE88391.1	694	OSCAR/DCAM box helicase
10	KGE88392.1	438	hypothetical protein

FIG. 4SS

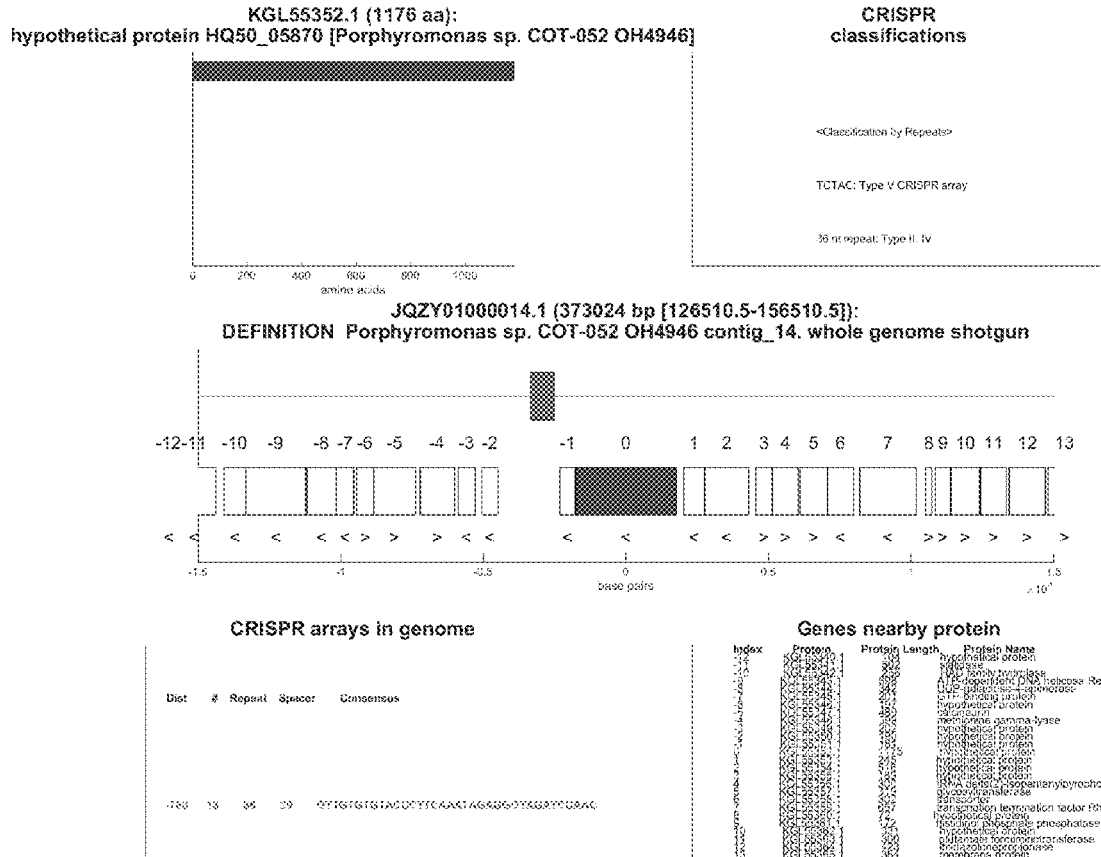


FIG. 4UU

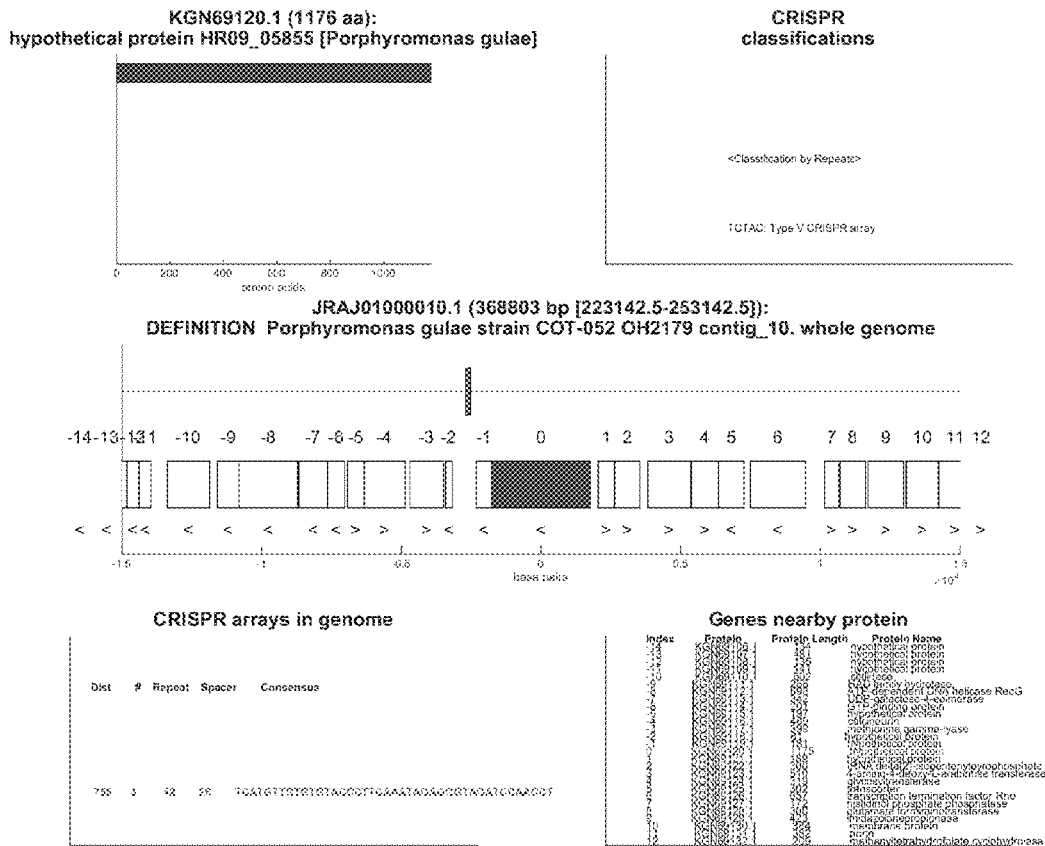


FIG. 4W

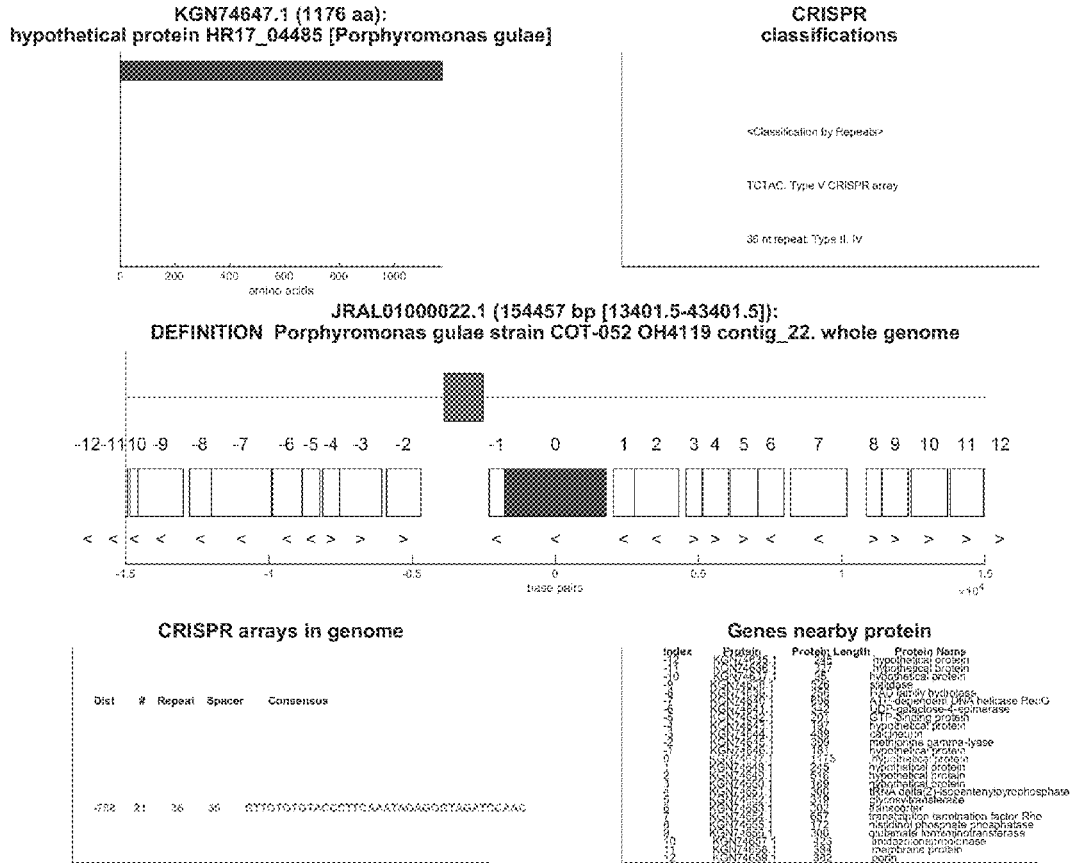


FIG. 4WW

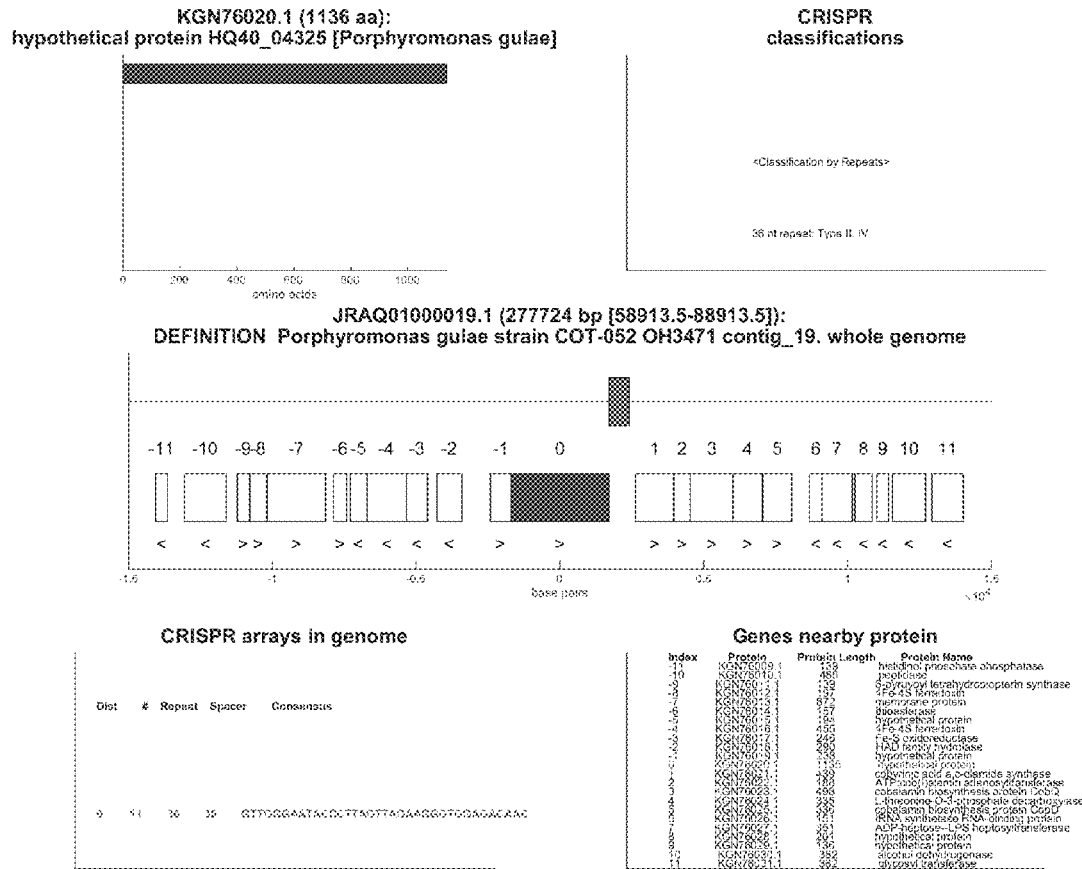
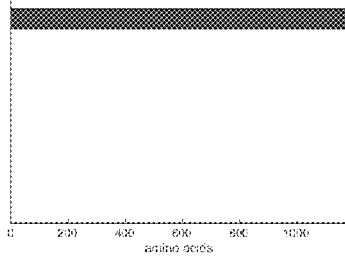


FIG. 4XX

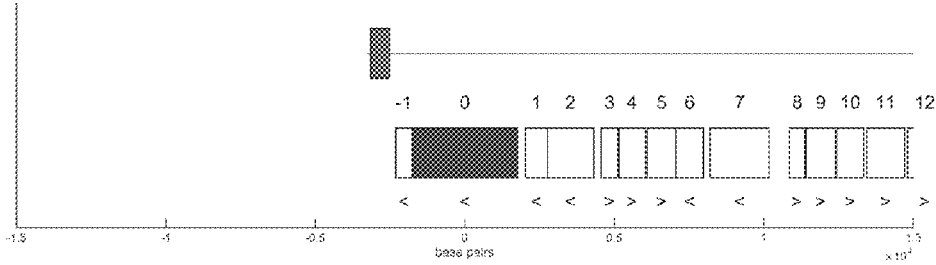
KG87312.1 (1175 aa):
hypothetical protein HQ46_09365 [Porphyromonas gulae]



CRISPR
classifications

<Classification by Repeats>
TCTAC: Type V CRISPR array
SS nt repeat: Type II, IV

JRFD01000046.1 (75493 bp [-11751-18249]):
DEFINITION Porphyromonas gulae strain COT-052 OH2857 contig_46, whole genome



CRISPR arrays in genome

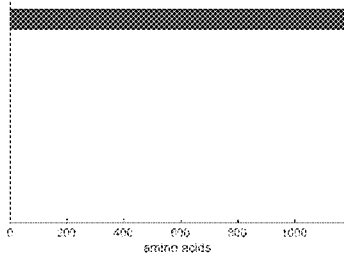
Dist	#	Repeat	Spacer	Consensus
268	10	25	23	GTGTTGTGTATCCCTTCMAATAGAGGGTADATGCAAG

Genes nearby protein

Index	Protein	Protein Length	Protein Name
-1	KG87311.1	181	hypothetical protein
0	KG87312.1	1174	hypothetical protein
1	KG87313.1	245	hypothetical protein
2	KG87314.1	516	hypothetical protein
3	KG87315.1	189	hypothetical protein
4	KG87316.1	300	tRNA delta(2)-accepton/lysochlophosphate
5	KG87317.1	319	glycosyltransferase
6	KG87318.1	302	transporter
7	KG87319.1	657	transcription termination factor Rho
8	KG87320.1	172	hexokinol phosphate phosphotase
9	KG87321.1	331	hypothetical protein
10	KG87322.1	300	glutamate formyltransferase
11	KG87323.1	423	nitidazolecarbamoylase
12	KG87324.1	384	metalloprotease

FIG. 4ZZ

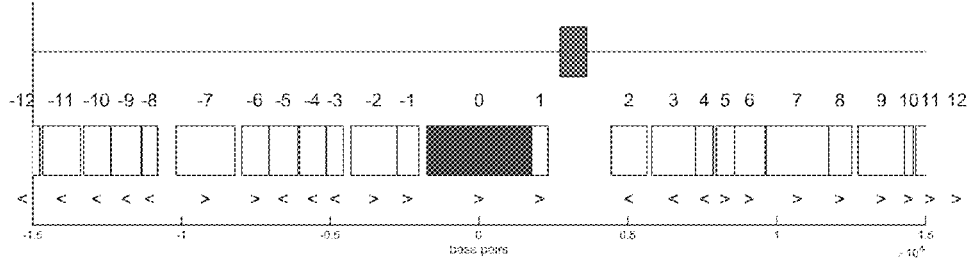
KGO05347.1 (1176 aa):
hypothetical protein HR16_00525 [Porphyromonas gulae]



CRISPR
classifications

<Classification by Repeats>
TCTAC: Type V CRISPR array
36 nt repeat: Type II, IV

JRAF01000001.1 (160244 bp [117096.5-147096.5]):
DEFINITION Porphyromonas gulae strain COT-052 OH3498 contig_1, whole genome



CRISPR arrays in genome

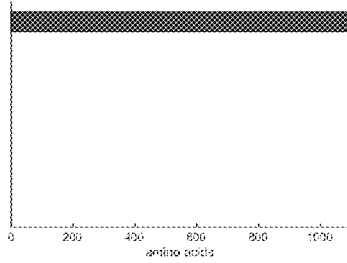
Dist	#	Repeat	Spacer	Consensus
898	14	36	30	GTTGGGATCTGCGCTGTATTTGAGGGGTACACAGAC

Genes nearby protein

Index	Protein	Protein Length	Protein Name
1	membrane protein
2	nitrogenase prokaryotic
3	nitrogenase prokaryotic
4	hypothetical protein
5	nitrogenase prokaryotic
6	transcription termination factor Rho
7	transposase
8	nitrogenase prokaryotic
9	RNA GTPase 2-oxoglutarate/ATP-dependent
10	hypothetical protein
11	hypothetical protein
12	hypothetical protein
13	hypothetical protein
14	hypothetical protein
15	hypothetical protein
16	hypothetical protein
17	hypothetical protein
18	hypothetical protein
19	hypothetical protein
20	hypothetical protein
21	hypothetical protein
22	hypothetical protein
23	hypothetical protein
24	hypothetical protein
25	hypothetical protein
26	hypothetical protein
27	hypothetical protein
28	hypothetical protein
29	hypothetical protein
30	hypothetical protein
31	hypothetical protein
32	hypothetical protein
33	hypothetical protein
34	hypothetical protein
35	hypothetical protein
36	hypothetical protein
37	hypothetical protein
38	hypothetical protein
39	hypothetical protein
40	hypothetical protein
41	hypothetical protein
42	hypothetical protein
43	hypothetical protein
44	hypothetical protein
45	hypothetical protein
46	hypothetical protein
47	hypothetical protein
48	hypothetical protein
49	hypothetical protein
50	hypothetical protein
51	hypothetical protein
52	hypothetical protein
53	hypothetical protein
54	hypothetical protein
55	hypothetical protein
56	hypothetical protein
57	hypothetical protein
58	hypothetical protein
59	hypothetical protein
60	hypothetical protein
61	hypothetical protein
62	hypothetical protein
63	hypothetical protein
64	hypothetical protein
65	hypothetical protein
66	hypothetical protein
67	hypothetical protein
68	hypothetical protein
69	hypothetical protein
70	hypothetical protein
71	hypothetical protein
72	hypothetical protein
73	hypothetical protein
74	hypothetical protein
75	hypothetical protein
76	hypothetical protein
77	hypothetical protein
78	hypothetical protein
79	hypothetical protein
80	hypothetical protein
81	hypothetical protein
82	hypothetical protein
83	hypothetical protein
84	hypothetical protein
85	hypothetical protein
86	hypothetical protein
87	hypothetical protein
88	hypothetical protein
89	hypothetical protein
90	hypothetical protein
91	hypothetical protein
92	hypothetical protein
93	hypothetical protein
94	hypothetical protein
95	hypothetical protein
96	hypothetical protein
97	hypothetical protein
98	hypothetical protein
99	hypothetical protein
100	hypothetical protein

FIG. 4BBB

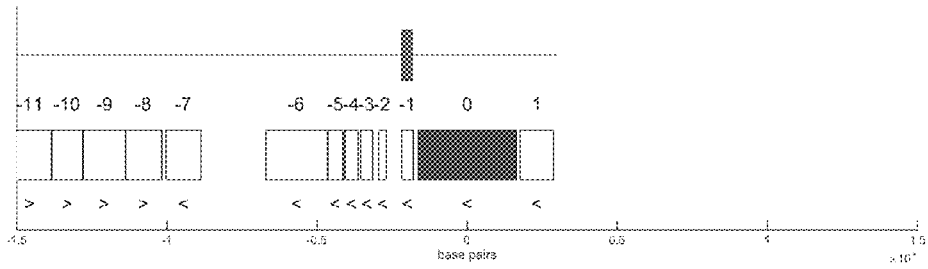
KIP58950.1 (1091 aa):
hypothetical protein ST43_06385 [Prevotella sp. P5-60]



CRISPR
classifications

<Classification by Repeat>
GTTT: CRISPR array
37 nt repeat: Type I, II-B, III, IV
20 nt repeat: Type V

JXQJ0100080.1 (32870 bp [14894-44894]):
DEFINITION Prevotella sp. P5-60 config80, whole genome shotgun sequence.



CRISPR arrays in genome

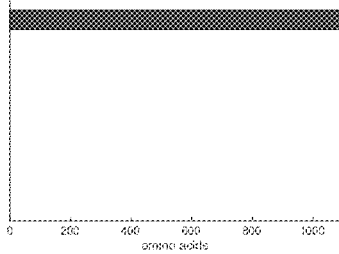
Dist	#	Repeat	Spacer	Consensus
-13632	22	37	26	CTTGGGAGAGCATATGGTTCAGTAAATATAGGATTTAGAGAG
-191	0	30	48	GTTTGAATAGTCCCTTAAAGA

Genes nearby protein

Index	Protein	Protein Length	Protein Name
-11	KIP58942.1	492	hydrogenase
-10	KIP58943.1	343	beta synthase
-9	KIP58952.1	470	thiamine biosynthesis protein TRH
-8	KIP58944.1	400	GTP-binding protein
-7	KIP58953.1	366	hypothetical protein
-6	KIP58945.1	684	hypothetical protein
-5	KIP58946.1	168	hypothetical protein
-4	KIP58947.1	146	hypothetical protein
-3	KIP58954.1	131	hypothetical protein
-2	KIP58948.1	84	hypothetical protein
-1	KIP58949.1	127	hypothetical protein
0	KIP58950.1	1091	hypothetical protein
1	KIP58951.1	370	hypothetical protein

FIG. 4CCC

KIP62088.1 (1091 aa):
hypothetical protein ST45_06300 [Prevotella sp. P5-125]



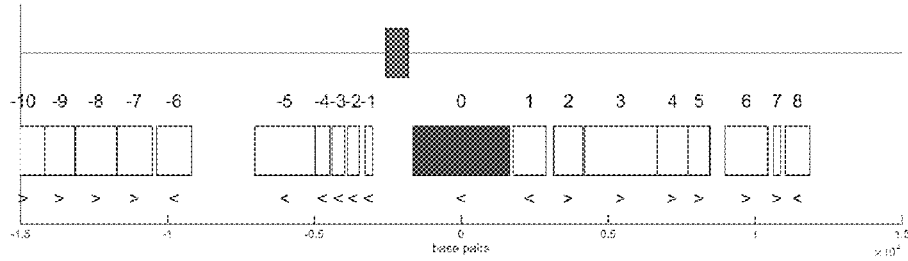
CRISPR
classifications

<Classification by Repeats>

GTTT: CRISPR array

88 nt repeat Type II IV

JXQL01000055.1 (74092 bp [33881-63881]):
DEFINITION Prevotella sp. P5-125 contig55. whole genome shotgun sequence.



CRISPR arrays in genome

Dist	#	Repeat	Spacer	Consensus
	1	26	28	GTTCCTAATGATGCGCCGCTCCAAAGACCGGAGCCCTTCACACAC

Genes nearby protein

Index	Protein	Protein Length	Protein Name
-10	KIP62081.1	480	hydrogenase
-9	KIP62082.1	343	bovin synthase
-2	KIP62103.1	470	threonine biosynthesis protein ThiB
-7	KIP62093.1	400	GTP-binding protein
-6	KIP62104.1	597	hypothetical protein
-5	KIP62084.1	684	hypothetical protein
-4	KIP62085.1	188	hypothetical protein
-3	KIP62086.1	148	hypothetical protein
-2	KIP62105.1	431	hypothetical protein
-1	KIP62087.1	84	hypothetical protein
0	KIP62088.1	1090	hypothetical protein
1	KIP62089.1	370	hypothetical protein
2	KIP62090.1	358	hypothetical protein
3	KIP62091.1	223	hypothetical protein
4	KIP62092.1	346	hypothetical protein
5	KIP62093.1	245	hypothetical protein
6	KIP62094.1	484	ATCase
7	KIP62095.1	78	hypothetical protein
8	KIP62096.1	276	phospholipase

FIG. 4DDD

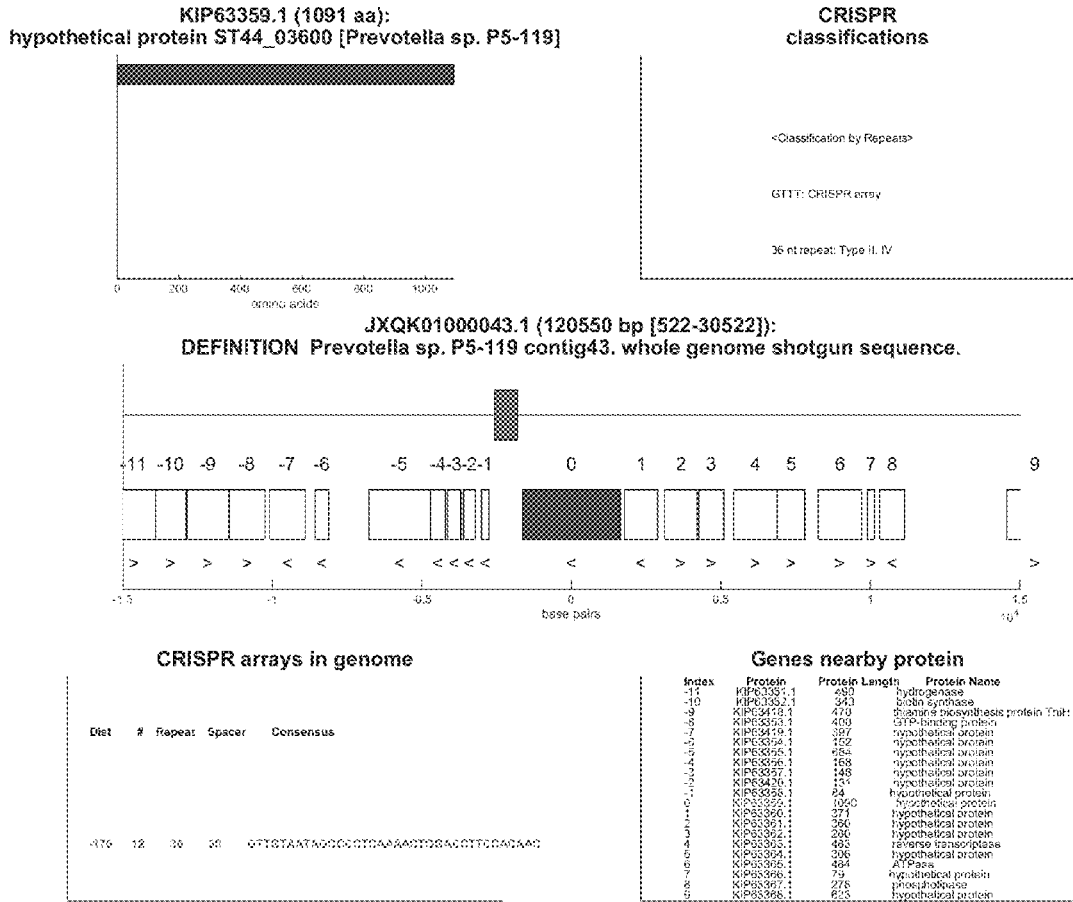
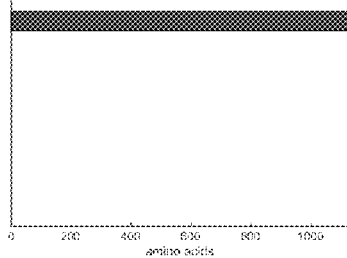


FIG. 4EEE

KJJ86756.1 (1126 aa):
 hypothetical protein M573_117042 [Prevotella intermedia ZT]

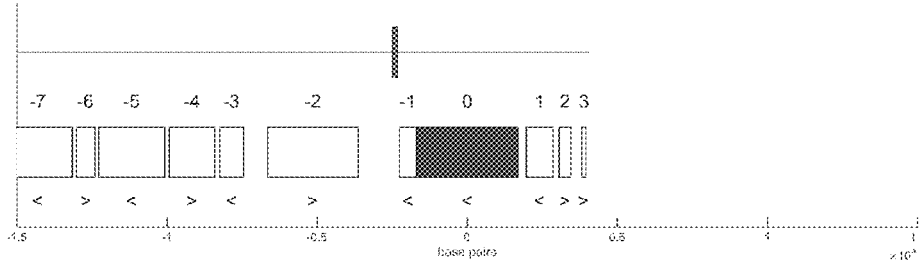


CRISPR classifications

<Classification by Repeats>

GTTT: CRISPR array

ATMK01000017.1 (58607 bp [39579.5-69579.5]):
 DEFINITION Prevotella intermedia ZT contig00017. whole genome shotgun



CRISPR arrays in genome

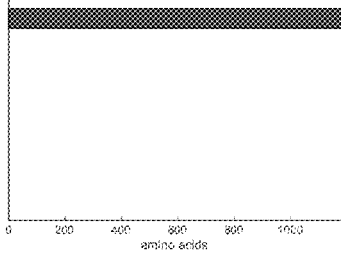
Dist	#	Repeat	Spacer	Consensus
-34	5	42	23	AAABAATGTTGTTTAACTTGGAAAPGACAGSBCAGATGCAAC

Genes nearby protein

Index	Protein	Protein Length	Protein Name
-7	KJJ86749.1	762	putative maleate dehydrogenase
-6	KJJ86750.1	205	TetR family transcriptional regulator
-5	KJJ86751.1	732	MFS transporter
-4	KJJ86752.1	905	hypothetical protein
-3	KJJ86753.1	295	DNA-binding protein
-2	KJJ86754.1	1006	hypothetical protein
-1	KJJ86755.1	180	hypothetical protein
0	KJJ86756.1	1126	hypothetical protein
1	KJJ86757.1	297	hypothetical protein
2	KJJ86758.1	128	hypothetical protein
3	KJJ86759.1	48	hypothetical protein

FIG. 4FFF

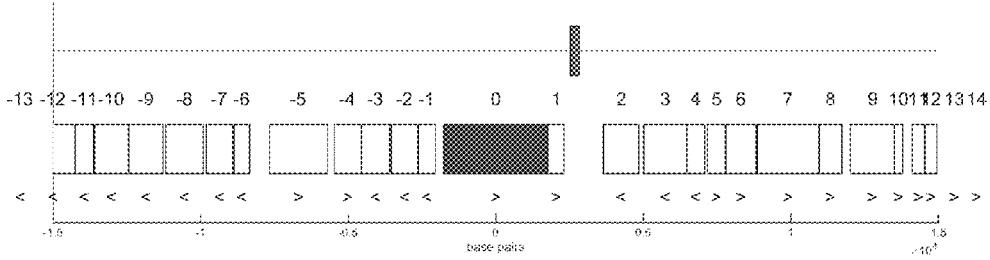
KKC50278.1 (1176 aa):
 hypothetical protein HR10_10685 [Porphyromonas gulae]



CRISPR
 classifications

<Classification by Repeats>
 TCTAC: Type V CRISPR array
 88 of repeats: Type II, IV

JRAE01000093.1 (161241 bp [113199.5-143199.5]):
 DEFINITION Porphyromonas gulae strain COT-052_OH2199 contig_93. whole genome



CRISPR arrays in genome

Dist	#	Repeat	Spacer	Consensus
729	6	36	20	CTTGGATGCTAGCCGCGTAAATTGGAAAGGNTAGGAGACAGC

Genes nearby protein

Index	Protein	Protein Length	Protein Name
8415	YP_009848246.1	414	unannotated protein
8416	YP_009848247.1	407	isoflavin glucosyl-transferase
8417	YP_009848248.1	129	non-symbiotic hemolysin
8418	YP_009848249.1	184	cyclohydrolase
8419	YP_009848250.1	204	unannotated protein
8420	YP_009848251.1	184	unannotated protein
8421	YP_009848252.1	184	unannotated protein
8422	YP_009848253.1	184	unannotated protein
8423	YP_009848254.1	184	unannotated protein
8424	YP_009848255.1	184	unannotated protein
8425	YP_009848256.1	184	unannotated protein
8426	YP_009848257.1	184	unannotated protein
8427	YP_009848258.1	184	unannotated protein
8428	YP_009848259.1	184	unannotated protein
8429	YP_009848260.1	184	unannotated protein
8430	YP_009848261.1	184	unannotated protein
8431	YP_009848262.1	184	unannotated protein
8432	YP_009848263.1	184	unannotated protein
8433	YP_009848264.1	184	unannotated protein
8434	YP_009848265.1	184	unannotated protein
8435	YP_009848266.1	184	unannotated protein
8436	YP_009848267.1	184	unannotated protein
8437	YP_009848268.1	184	unannotated protein
8438	YP_009848269.1	184	unannotated protein
8439	YP_009848270.1	184	unannotated protein
8440	YP_009848271.1	184	unannotated protein
8441	YP_009848272.1	184	unannotated protein
8442	YP_009848273.1	184	unannotated protein
8443	YP_009848274.1	184	unannotated protein
8444	YP_009848275.1	184	unannotated protein
8445	YP_009848276.1	184	unannotated protein
8446	YP_009848277.1	184	unannotated protein
8447	YP_009848278.1	184	unannotated protein
8448	YP_009848279.1	184	unannotated protein
8449	YP_009848280.1	184	unannotated protein
8450	YP_009848281.1	184	unannotated protein
8451	YP_009848282.1	184	unannotated protein
8452	YP_009848283.1	184	unannotated protein
8453	YP_009848284.1	184	unannotated protein
8454	YP_009848285.1	184	unannotated protein
8455	YP_009848286.1	184	unannotated protein
8456	YP_009848287.1	184	unannotated protein
8457	YP_009848288.1	184	unannotated protein
8458	YP_009848289.1	184	unannotated protein
8459	YP_009848290.1	184	unannotated protein
8460	YP_009848291.1	184	unannotated protein
8461	YP_009848292.1	184	unannotated protein
8462	YP_009848293.1	184	unannotated protein
8463	YP_009848294.1	184	unannotated protein
8464	YP_009848295.1	184	unannotated protein
8465	YP_009848296.1	184	unannotated protein
8466	YP_009848297.1	184	unannotated protein
8467	YP_009848298.1	184	unannotated protein
8468	YP_009848299.1	184	unannotated protein
8469	YP_009848300.1	184	unannotated protein
8470	YP_009848301.1	184	unannotated protein
8471	YP_009848302.1	184	unannotated protein
8472	YP_009848303.1	184	unannotated protein
8473	YP_009848304.1	184	unannotated protein
8474	YP_009848305.1	184	unannotated protein
8475	YP_009848306.1	184	unannotated protein
8476	YP_009848307.1	184	unannotated protein
8477	YP_009848308.1	184	unannotated protein
8478	YP_009848309.1	184	unannotated protein
8479	YP_009848310.1	184	unannotated protein
8480	YP_009848311.1	184	unannotated protein
8481	YP_009848312.1	184	unannotated protein
8482	YP_009848313.1	184	unannotated protein
8483	YP_009848314.1	184	unannotated protein
8484	YP_009848315.1	184	unannotated protein
8485	YP_009848316.1	184	unannotated protein
8486	YP_009848317.1	184	unannotated protein
8487	YP_009848318.1	184	unannotated protein
8488	YP_009848319.1	184	unannotated protein
8489	YP_009848320.1	184	unannotated protein
8490	YP_009848321.1	184	unannotated protein
8491	YP_009848322.1	184	unannotated protein
8492	YP_009848323.1	184	unannotated protein
8493	YP_009848324.1	184	unannotated protein
8494	YP_009848325.1	184	unannotated protein
8495	YP_009848326.1	184	unannotated protein
8496	YP_009848327.1	184	unannotated protein
8497	YP_009848328.1	184	unannotated protein
8498	YP_009848329.1	184	unannotated protein
8499	YP_009848330.1	184	unannotated protein

FIG. 4GGG

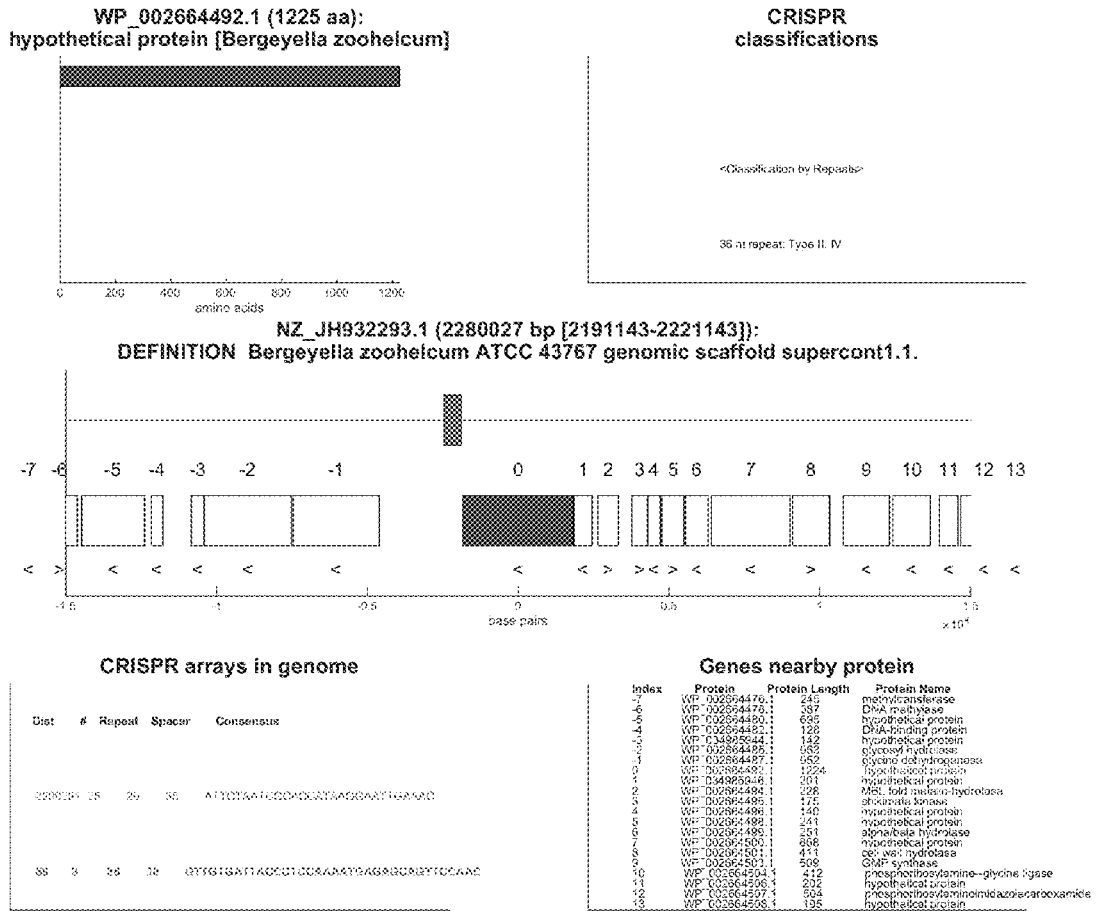


FIG. 4HHH

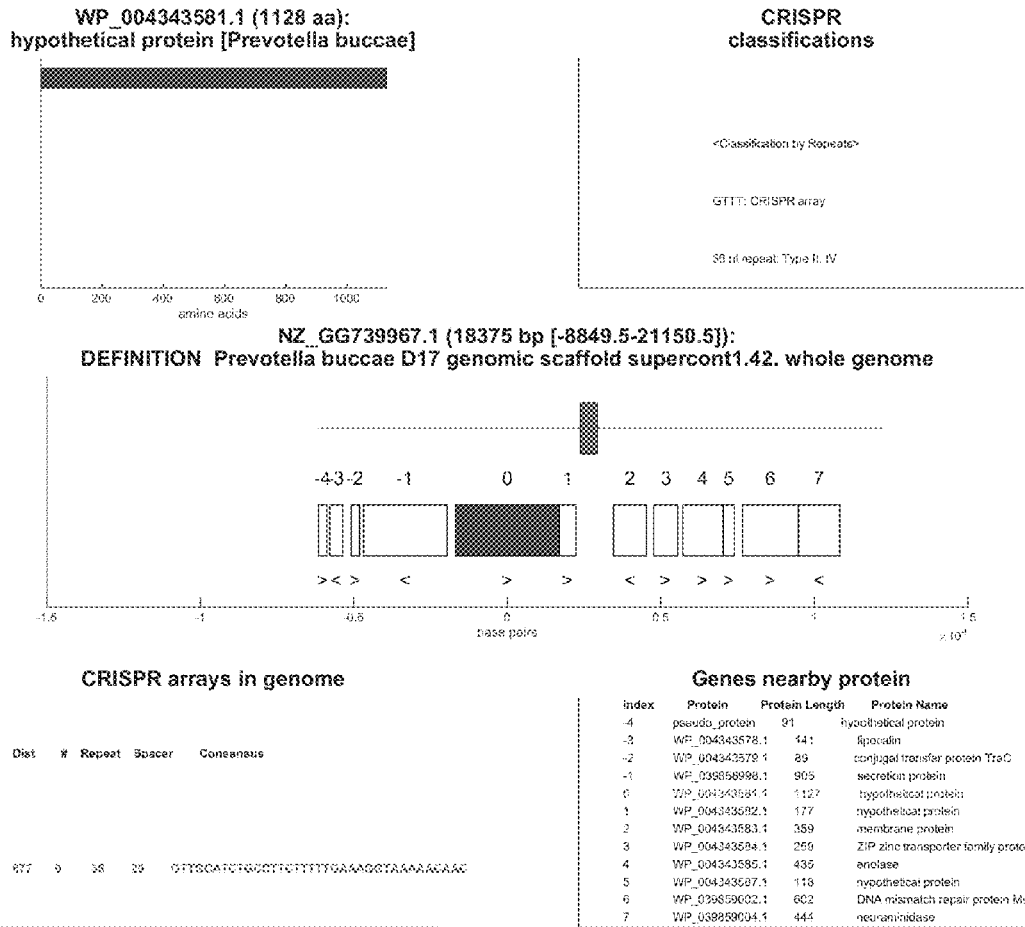


FIG. 4III

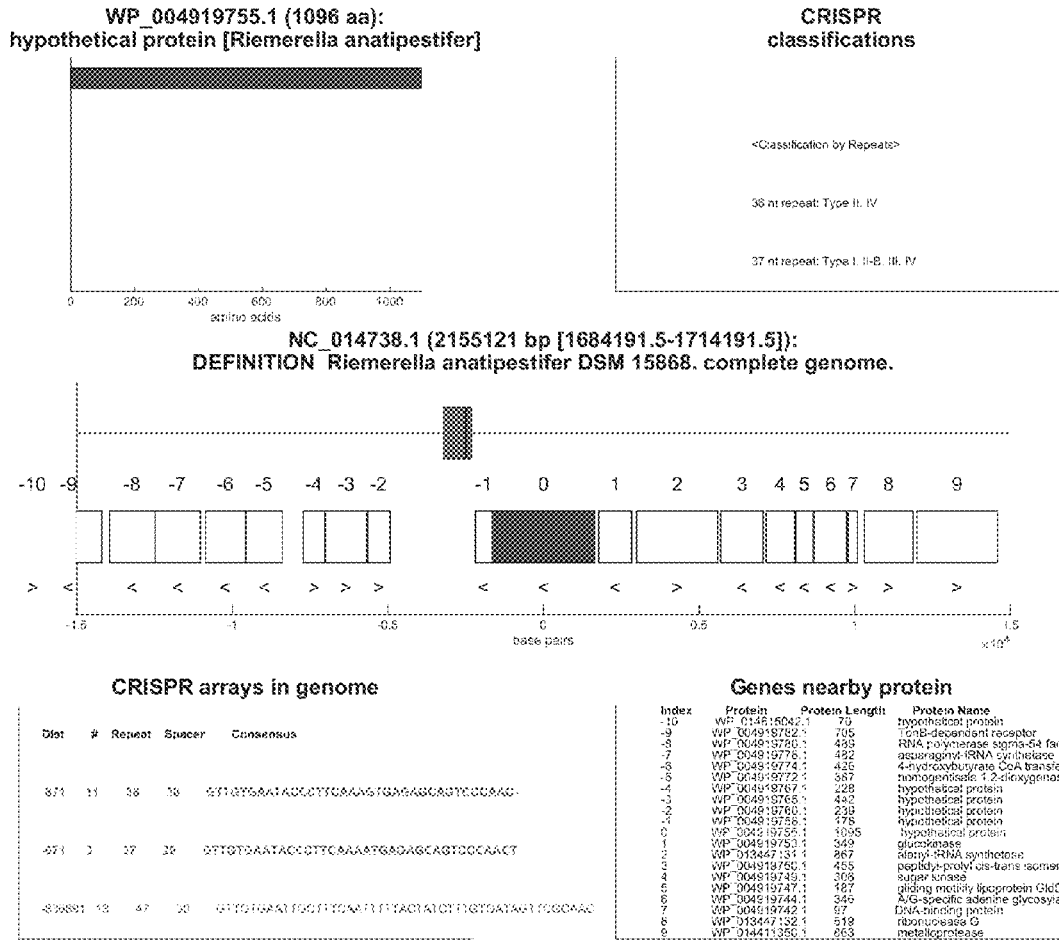


FIG. 4KKK

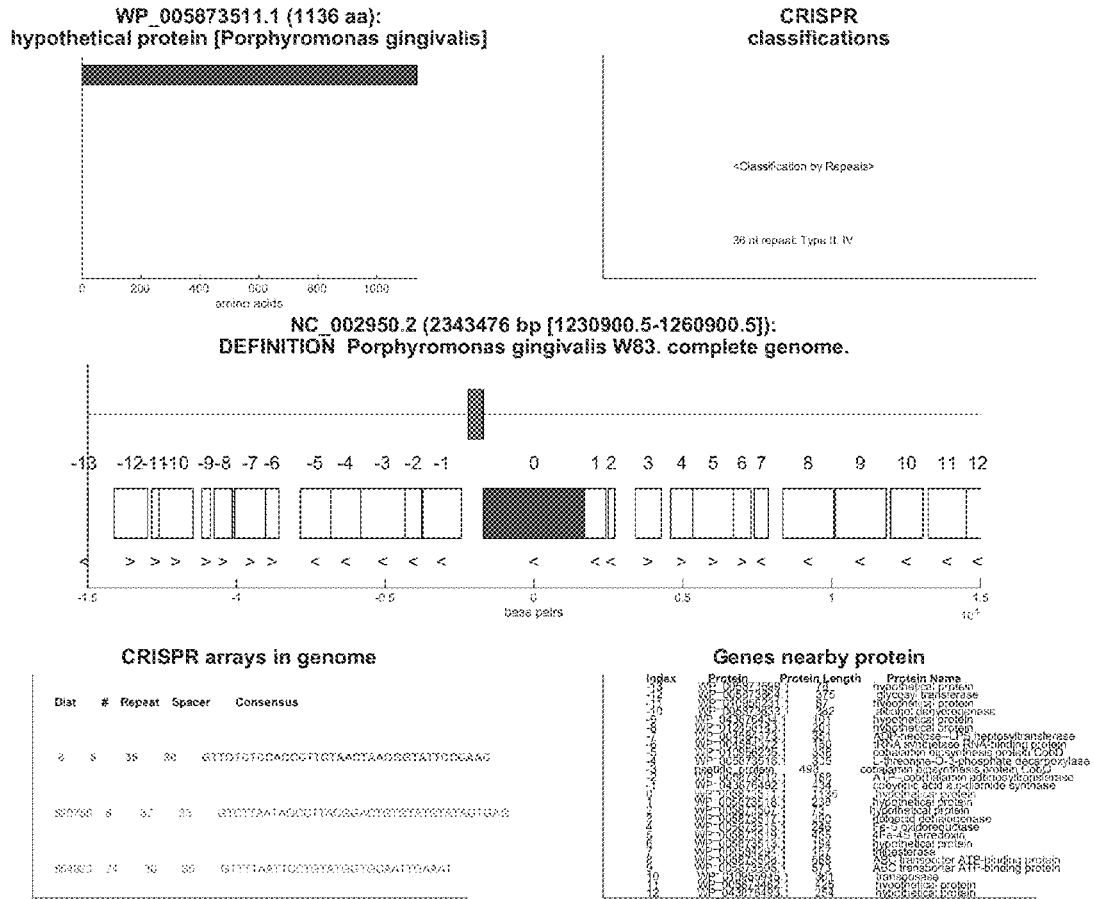
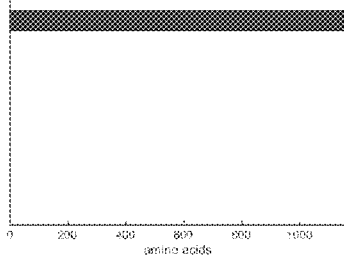


FIG. 4LLL

WP_006259957.1 (1160 aa):
hypothetical protein [Myroides odoratimimus]

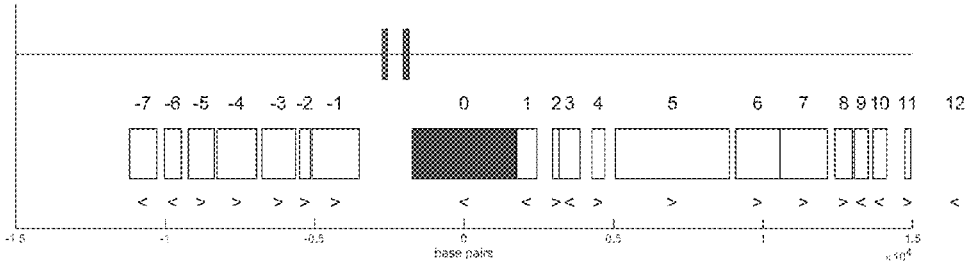


CRISPR
classifications

<Classification by Repeats>

GTTT: CRISPR array

NZ_KE161016.1 (1024274 bp [190937.5-220937.5]):
DEFINITION Myroides odoratimimus CCUG 10230 genomic scaffold supercont2.2.



CRISPR arrays in genome

Dist	#	Repeat	Spacer	Consensus
226	3	30	27	TTTTCGATATTCGATGAGCAAGCAATACGTTTACAGTACCAAGCT
117	1	12	25	CTTGTGATTCGATGAGCAAGCAATACGTTTACAGTACCAAGCT
322525	44	30	25	ATTTCGATACCTATATCGTCCAGTTAAAGC
55511	30	30	25	ATTTCGATACCTATATCGTCCAGTTAAAGC

Genes nearby protein

Index	Protein	Protein Length	Protein Name
-7	WP_006259954.1	206	hypothetical protein
-6	WP_006259953.1	186	DNA-directed RNA polymerase sigma-70 factor
-5	WP_006259952.1	205	AraC family transcriptional regulator
-4	WP_006259951.1	419	mercaptopyruvate sulfurtransferase
-3	WP_006259950.1	372	hemolysin D
-2	WP_006259949.1	124	four-helix bundle protein
-1	WP_006259948.1	526	transporter
0	WP_006259947.1	1160	hypothetical protein
1	pseudo protein	227	
2	WP_030851756.1	52	hypothetical protein
3	pseudo protein	235	
4	WP_006259944.1	142	hypothetical protein
5	WP_006259943.1	1171	carbon transport
6	WP_006259942.1	452	transporter
7	WP_006259941.1	525	cation transporter
8	WP_006259940.1	195	hypothetical protein
9	pseudo protein	182	
10	WP_030851849.1	135	hypothetical protein
11	WP_006259945.1	65	transcriptional regulator
12	WP_006259944.1	660	regular methyl protein Methyl

FIG. 4NNN

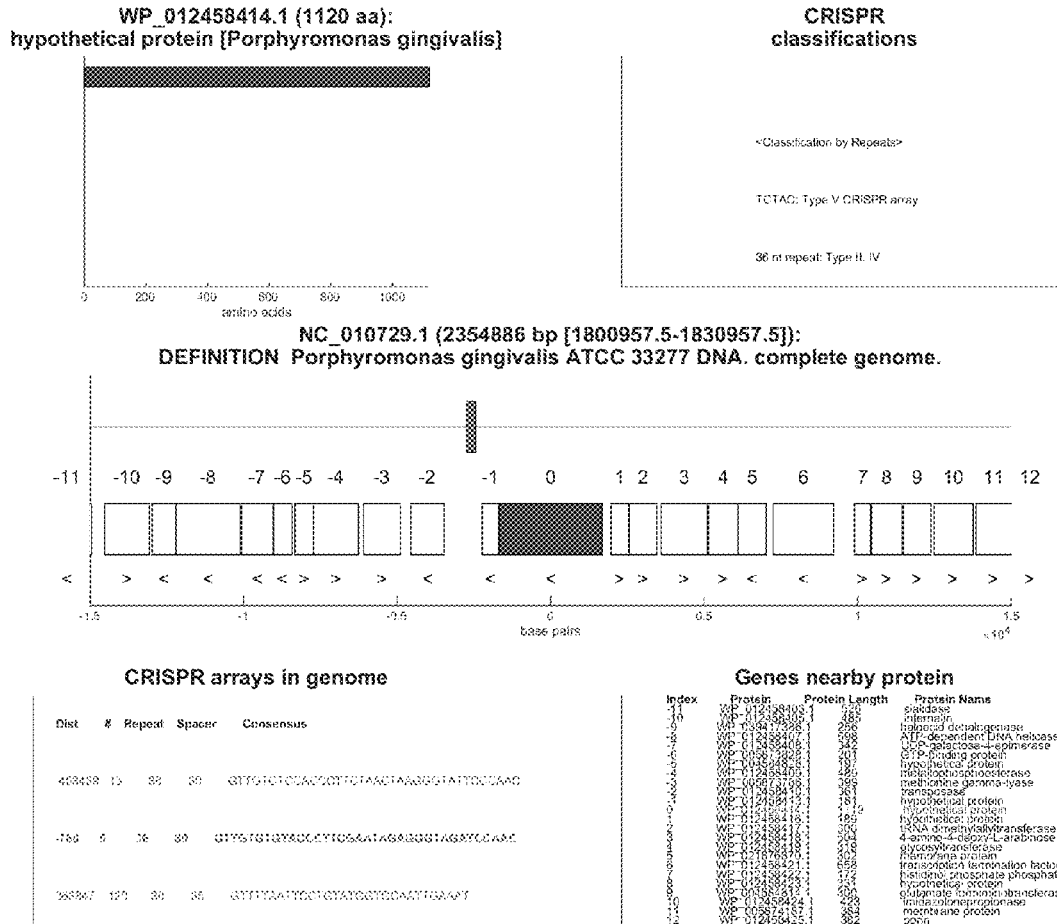


FIG. 4QQQ

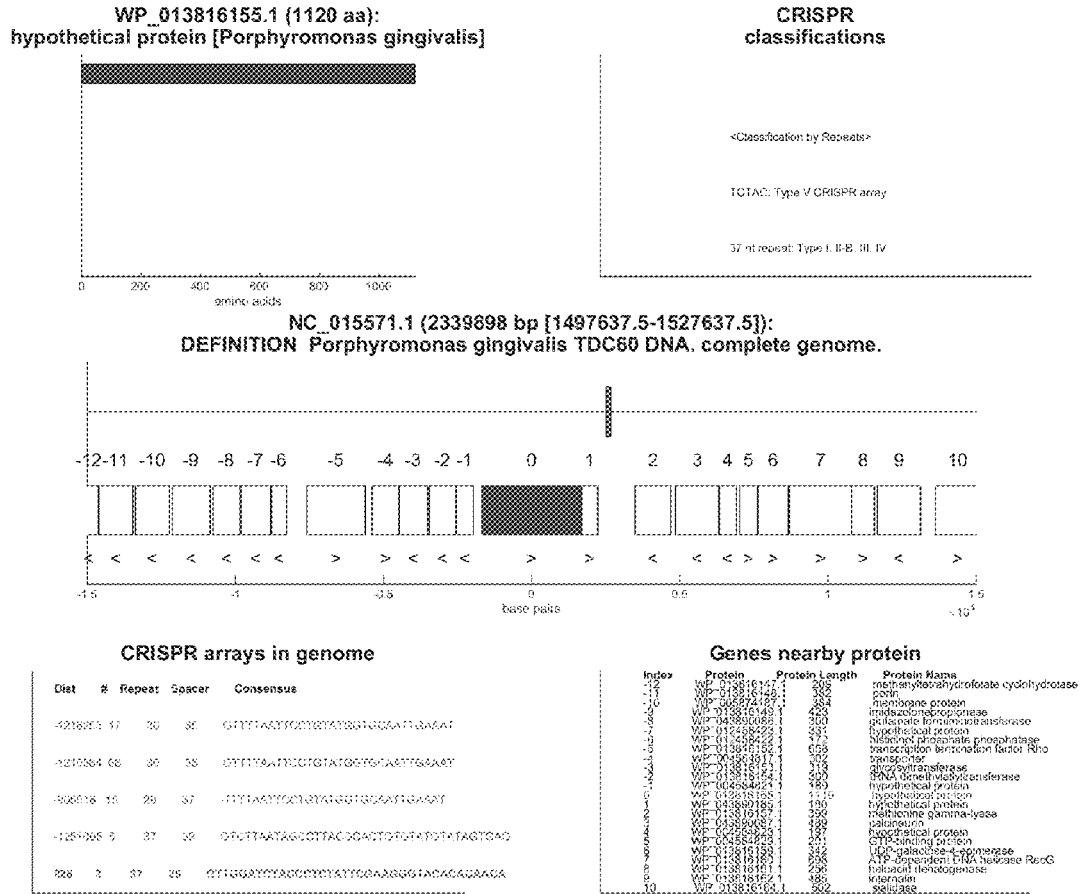
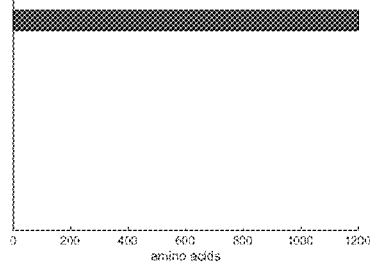


FIG. 4RRR

WP_013997271.1 (1200 aa):
hypothetical protein [Capnocytophaga canimorsus]

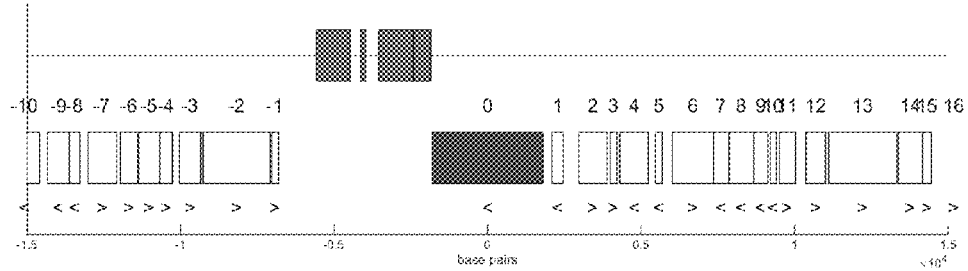


CRISPR
classifications

<Classification by Repeats>

- 36 nt repeat: Type II, IV
- 37 nt repeat: Type I, II-B, III, IV
- 36 nt repeat: Type II, IV
- 36 nt repeat: Type II, IV

NC_015846.1 (2571406 bp [1220228.5-1250228.5]):
DEFINITION Capnocytophaga canimorsus Cc5. complete genome.



CRISPR arrays in genome

Dist	#	Repeat	Spacer	Consensus
3583	17	36	36	-GTTTGTATTACCGTCCAAATGACAGCACTTTCGAAC
3178	3	37	36	CGTTGTGATYACGCTCCAAATGACAGCACTTTCGAAC
38	10	36	61	GTTTGTGATYACGCTCCAAATGACAGCACTTTCGAAC
49	7	36	62	GTTTGTGATYACGCTCCAAATGACAGCACTTTCGAAC
135895	66	47	30	GTTTGTGATYACGCTCCAAATGACAGCACTTTCGAAC
392172	20	47	19	GTTTGTGATYACGCTCCAAATGACAGCACTTTCGAAC

Genes nearby protein

Index	Protein ID	Protein Length	Protein Name
1	WP_013997271.1	1200	hypothetical protein
2	WP_013997272.1	230	adenosine deaminase
3	WP_013997273.1	111	membrane protein
4	WP_013997274.1	236	membrane protein
5	WP_013997275.1	236	membrane protein
6	WP_013997276.1	236	membrane protein
7	WP_013997277.1	236	membrane protein
8	WP_013997278.1	236	membrane protein
9	WP_013997279.1	236	membrane protein
10	WP_013997280.1	236	membrane protein
11	WP_013997281.1	236	membrane protein
12	WP_013997282.1	236	membrane protein
13	WP_013997283.1	236	membrane protein
14	WP_013997284.1	236	membrane protein
15	WP_013997285.1	236	membrane protein
16	WP_013997286.1	236	membrane protein
17	WP_013997287.1	236	membrane protein
18	WP_013997288.1	236	membrane protein
19	WP_013997289.1	236	membrane protein
20	WP_013997290.1	236	membrane protein
21	WP_013997291.1	236	membrane protein
22	WP_013997292.1	236	membrane protein
23	WP_013997293.1	236	membrane protein
24	WP_013997294.1	236	membrane protein
25	WP_013997295.1	236	membrane protein
26	WP_013997296.1	236	membrane protein
27	WP_013997297.1	236	membrane protein
28	WP_013997298.1	236	membrane protein
29	WP_013997299.1	236	membrane protein
30	WP_013997300.1	236	membrane protein
31	WP_013997301.1	236	membrane protein
32	WP_013997302.1	236	membrane protein
33	WP_013997303.1	236	membrane protein
34	WP_013997304.1	236	membrane protein
35	WP_013997305.1	236	membrane protein
36	WP_013997306.1	236	membrane protein
37	WP_013997307.1	236	membrane protein
38	WP_013997308.1	236	membrane protein
39	WP_013997309.1	236	membrane protein
40	WP_013997310.1	236	membrane protein
41	WP_013997311.1	236	membrane protein
42	WP_013997312.1	236	membrane protein
43	WP_013997313.1	236	membrane protein
44	WP_013997314.1	236	membrane protein
45	WP_013997315.1	236	membrane protein
46	WP_013997316.1	236	membrane protein
47	WP_013997317.1	236	membrane protein
48	WP_013997318.1	236	membrane protein
49	WP_013997319.1	236	membrane protein
50	WP_013997320.1	236	membrane protein
51	WP_013997321.1	236	membrane protein
52	WP_013997322.1	236	membrane protein
53	WP_013997323.1	236	membrane protein
54	WP_013997324.1	236	membrane protein
55	WP_013997325.1	236	membrane protein
56	WP_013997326.1	236	membrane protein
57	WP_013997327.1	236	membrane protein
58	WP_013997328.1	236	membrane protein
59	WP_013997329.1	236	membrane protein
60	WP_013997330.1	236	membrane protein
61	WP_013997331.1	236	membrane protein
62	WP_013997332.1	236	membrane protein
63	WP_013997333.1	236	membrane protein
64	WP_013997334.1	236	membrane protein
65	WP_013997335.1	236	membrane protein
66	WP_013997336.1	236	membrane protein
67	WP_013997337.1	236	membrane protein
68	WP_013997338.1	236	membrane protein
69	WP_013997339.1	236	membrane protein
70	WP_013997340.1	236	membrane protein
71	WP_013997341.1	236	membrane protein
72	WP_013997342.1	236	membrane protein
73	WP_013997343.1	236	membrane protein
74	WP_013997344.1	236	membrane protein
75	WP_013997345.1	236	membrane protein
76	WP_013997346.1	236	membrane protein
77	WP_013997347.1	236	membrane protein
78	WP_013997348.1	236	membrane protein
79	WP_013997349.1	236	membrane protein
80	WP_013997350.1	236	membrane protein
81	WP_013997351.1	236	membrane protein
82	WP_013997352.1	236	membrane protein
83	WP_013997353.1	236	membrane protein
84	WP_013997354.1	236	membrane protein
85	WP_013997355.1	236	membrane protein
86	WP_013997356.1	236	membrane protein
87	WP_013997357.1	236	membrane protein
88	WP_013997358.1	236	membrane protein
89	WP_013997359.1	236	membrane protein
90	WP_013997360.1	236	membrane protein
91	WP_013997361.1	236	membrane protein
92	WP_013997362.1	236	membrane protein
93	WP_013997363.1	236	membrane protein
94	WP_013997364.1	236	membrane protein
95	WP_013997365.1	236	membrane protein
96	WP_013997366.1	236	membrane protein
97	WP_013997367.1	236	membrane protein
98	WP_013997368.1	236	membrane protein
99	WP_013997369.1	236	membrane protein
100	WP_013997370.1	236	membrane protein

FIG. 4SSS

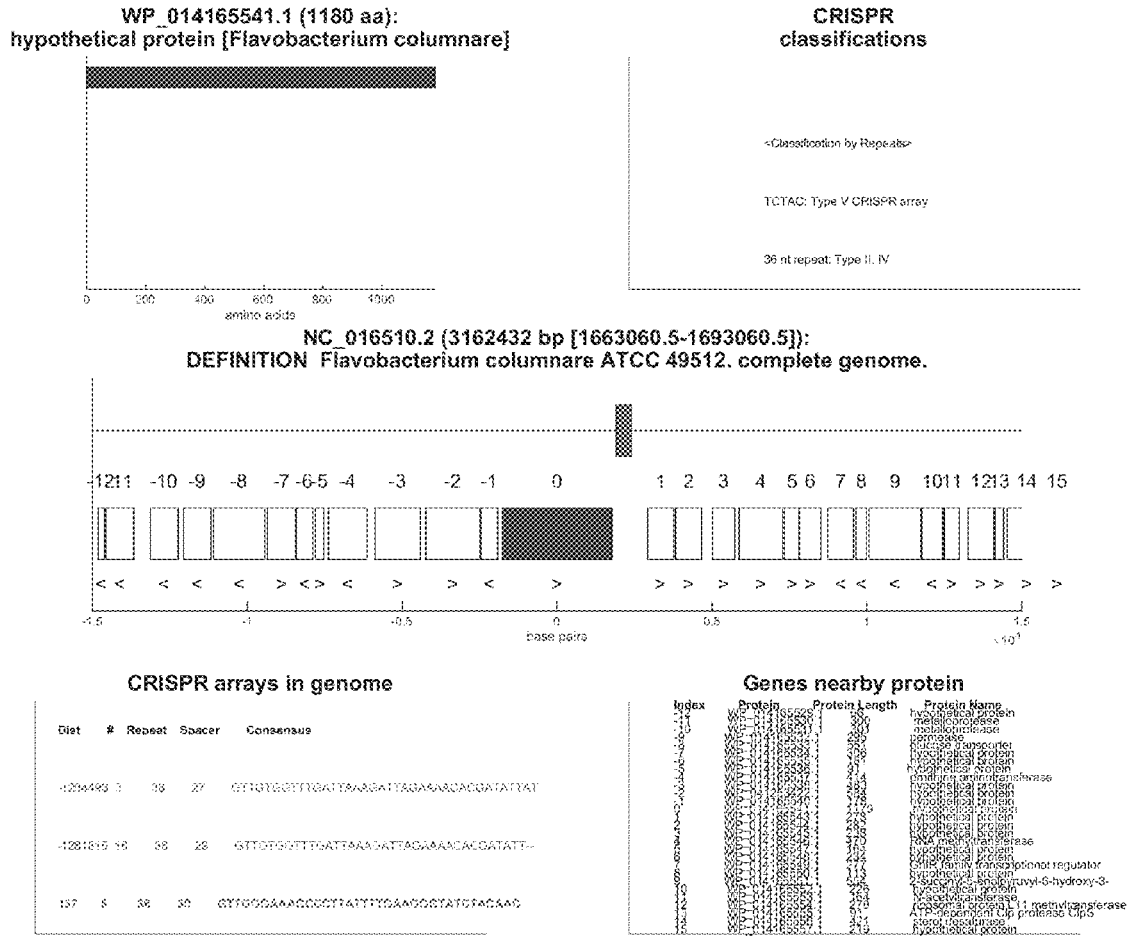
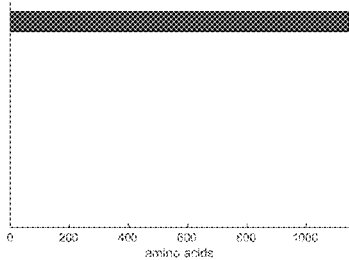


FIG. 4TTT

WP_015024765.1 (1146 aa):
hypothetical protein [Psychroflexus torquus]



CRISPR
classifications

<Classification by Repeats>

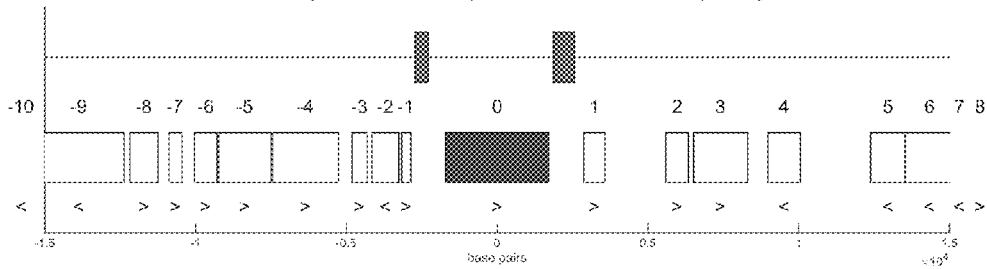
GTTT: CRISPR array

GTTT: CRISPR array

36 nt repeat: Type II, IV

36 nt repeat: Type II, IV

NC_018721.1 (4321832 bp [2596621.5-2626621.5]):
DEFINITION Psychroflexus torquus ATCC 700755, complete genome.



CRISPR arrays in genome

Dist	# Repeat	Spacer	Consensus	
-2281034	20	40	29	CGTGTGATTATLACTAAATATAGATTCTTGAAGCAATTGACAGC
2287685	3	67	28	CGTGTGATTATLACTAAATATAGATTCTTGAAGCAATTGACAGC
678	7	36	28	GTGTGATCTGCGCTTATTTGAGGGSFARAGAGAGC
144	15	36	38	GTTTGAATGCGCTTATTTGAGGGSFARAGAGAGC
1851107	3	40	10	TTGTGATCTGCGCTTATTTGAGGGSFARAGAGAGC

Genes nearby protein

Index	Protein	Protein Length	Protein Name
-10	WP_015024750.1	210	LuxR family transcriptional regulator
-6	WP_015024757.1	1001	two-component system sensor with a
-6	WP_015024756.1	210	diacylglycerol 4-reductase
-7	WP_015024755.1	143	ketosteroid isomerase
-8	WP_015024751.1	249	hypothetical protein
-5	WP_015024752.1	573	aminosynthase
-4	WP_015024753.1	724	peptidase 56
-3	WP_015024754.1	186	hypothetical protein
-2	WP_015024117.1	297	transposase
-1	WP_041758328.1	98	hypothetical protein
0	WP_015024765.1	1146	hypothetical protein
1	WP_015024767.1	231	hypothetical protein
2	WP_015024768.1	232	hypothetical protein
3	WP_015024770.1	596	hypothetical protein
4	WP_041758788.1	321	60S A 3-O-deacetylase
5	WP_015024772.1	382	hypothetical protein
6	WP_015024773.1	496	penicillinase antirepressor BlaR-like protein
7	WP_015024774.1	120	transcriptional regulator
8	WP_015024775.1	294	phenazine biosynthesis-like protein PhzC/PhzF

FIG. 4UUU

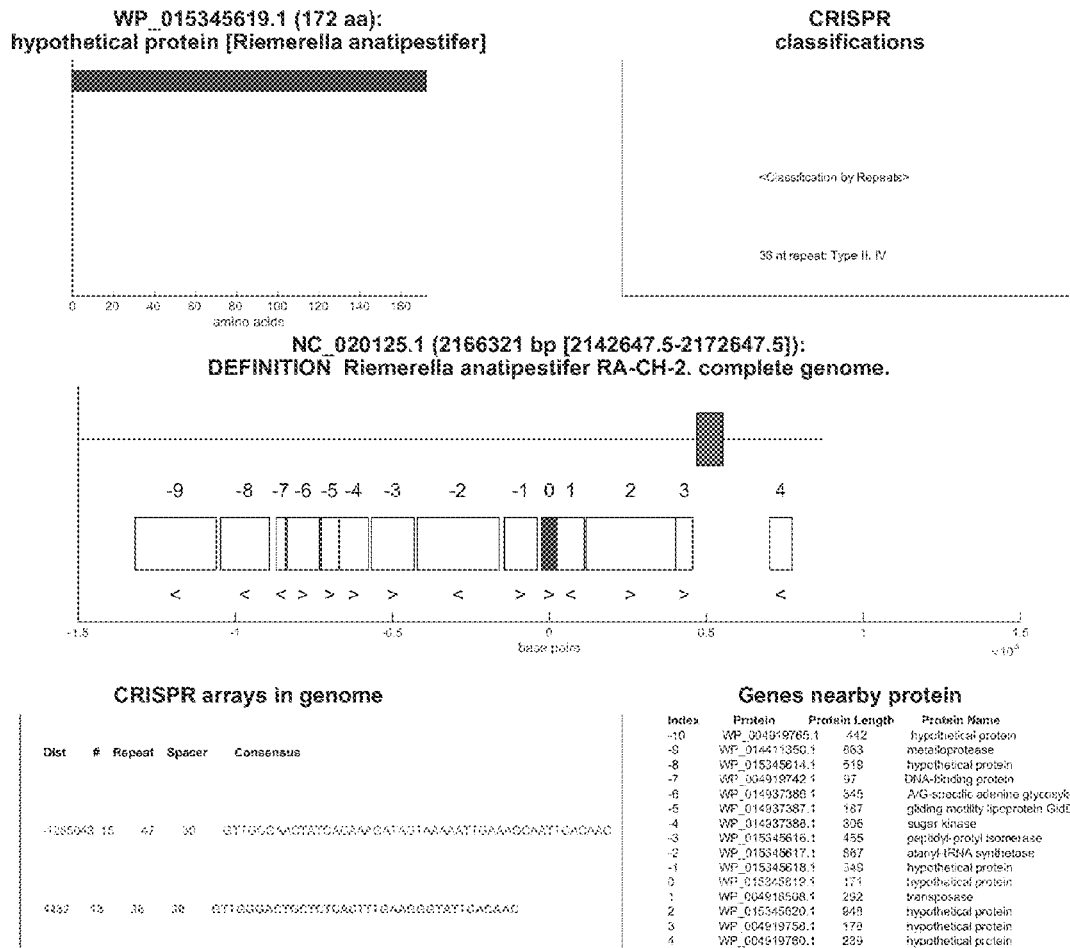


FIG. 4VVV

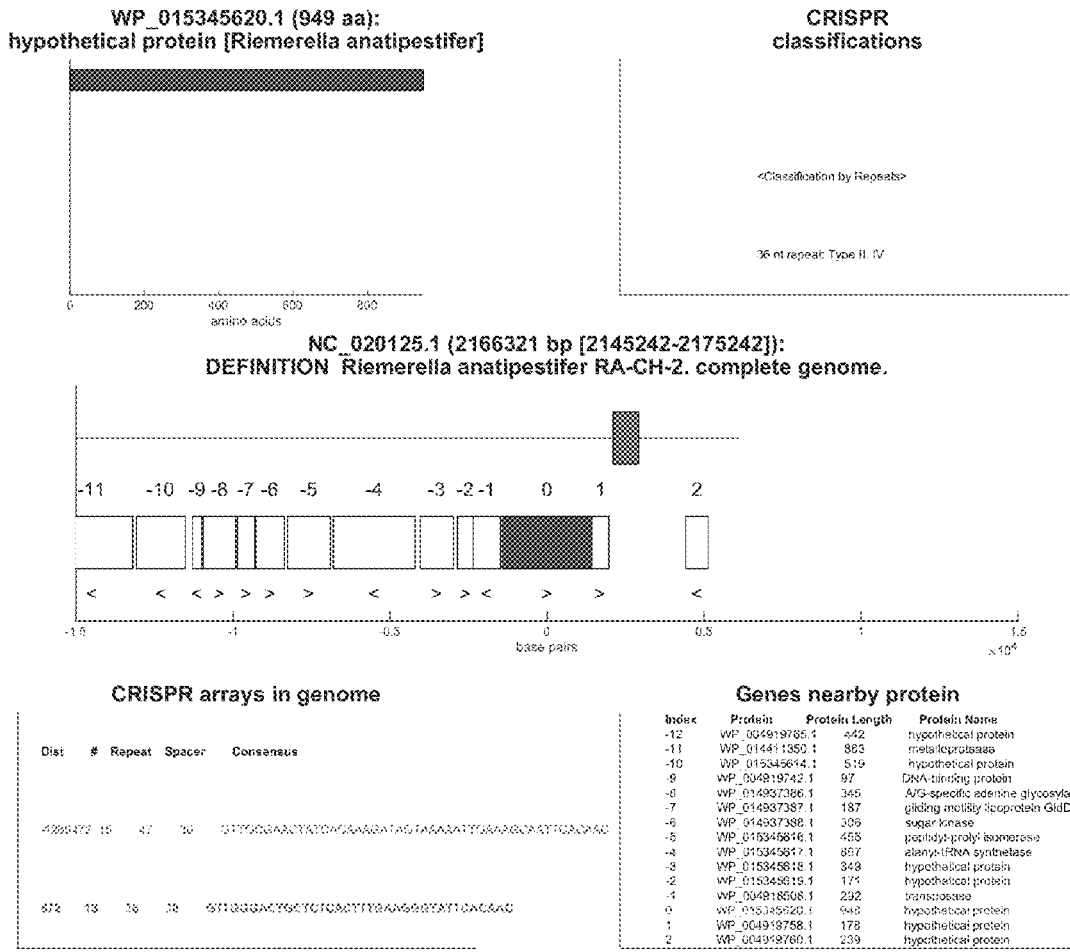
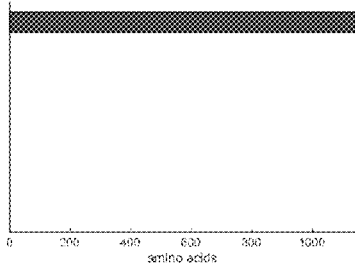


FIG. 4WWW

WP_021584635.1 (1140 aa):
hypothetical protein [Prevotella pleuritidis]



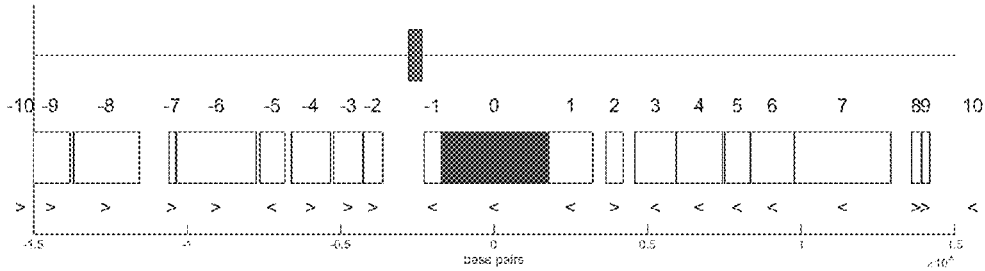
CRISPR
classifications

<Classification by Repeats>

GTTC CRISPR array

36 nt repeat Type II, IV

NZ_AWET01000045.1 (167660 bp [24484.5-54484.5]):
DEFINITION Prevotella pleuritidis F0068 contig00004. whole genome shotgun



CRISPR arrays in genome

Dist	# Repeat	Spacer	Consensus	
-842	7	56	28	GTCTTTTGTGCTTTTCAAAATAGAAAGCCACATCCAAC

Genes nearby protein

Index	Protein	Protein Length	Protein Name
-10	WP_021584621.1	213	putative lipoprotein
-9	WP_021584725.1	423	hypothetical protein
-8	WP_021584650.1	713	PP1C-type PPIase domain protein
-7	WP_039079912.1	98	hypothetical protein
-6	WP_021584703.1	557	hypothetical protein
-5	WP_021584673.1	208	hypothetical protein
-4	WP_021584715.1	423	MgSA, AAA+ ATPase
-3	WP_021584768.1	132	L-glutamate:gamma-butyrate dihydrogenase
-2	WP_021584678.1	212	membrane protein
-1	WP_021584703.1	486	hypothetical protein
0	WP_021584635.1	1140	CRISPR array
1	WP_021584697.1	478	transposon SRS family
2	WP_021584703.1	186	serine protease
3	WP_02951766.1	347	hypothetical protein
4	WP_021584632.1	499	hypothetical protein
5	WP_02931766.1	278	hypothetical protein
6	WP_021584740.1	470	signal-binding protein, SusD-like family
7	WP_021584674.1	1043	SusC/RagA family TonB-linked outer membrane
8	WP_021584727.1	135	SUS histone-like protein L21
9	WP_021584669.1	87	SUS histone-like protein L27
10	WP_021584635.1	367	hypothetical protein

FIG. 4XXX

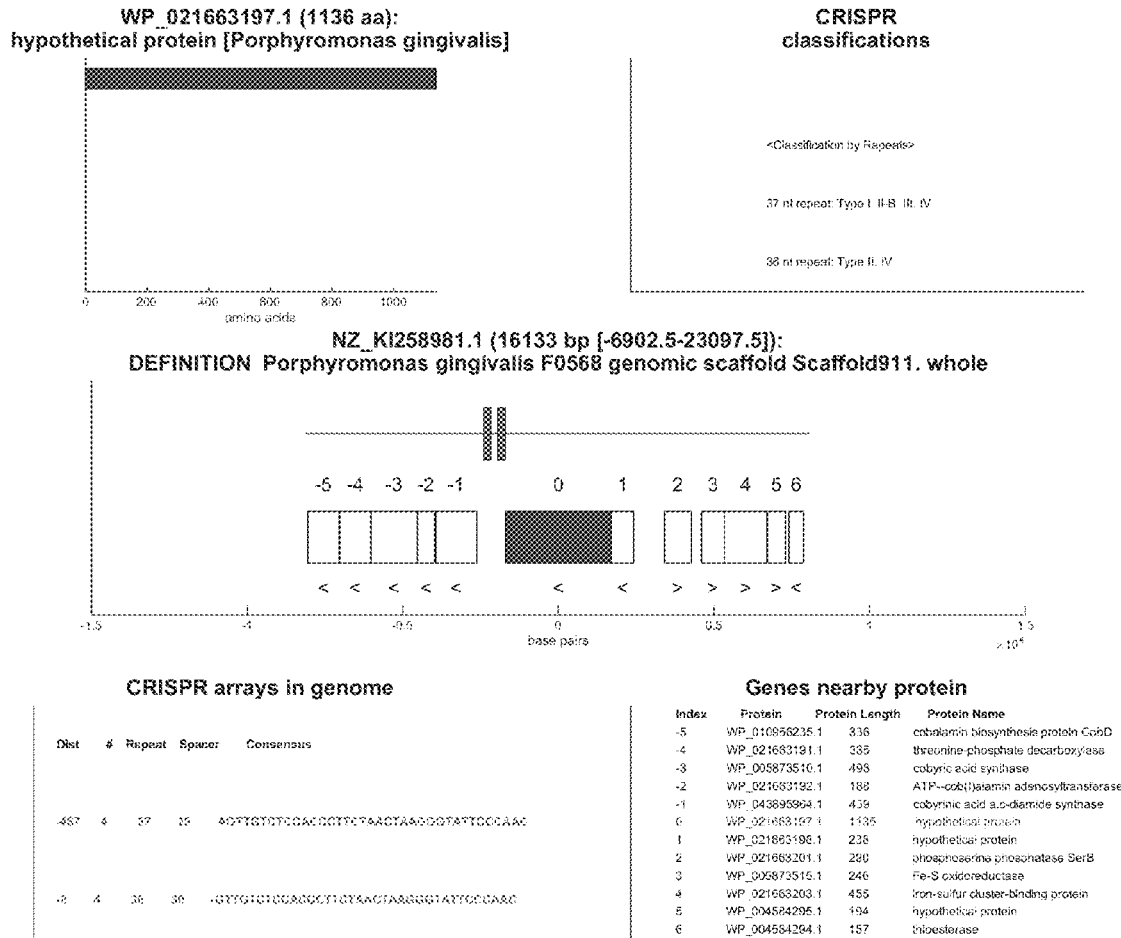
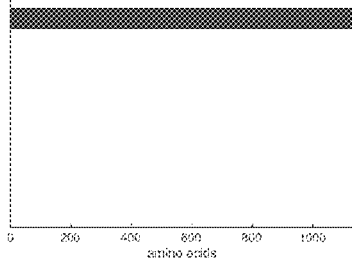


FIG. 4YYY

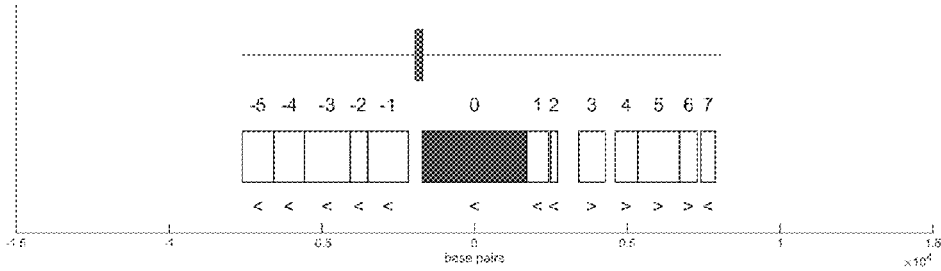
WP_021665475.1 (1136 aa):
hypothetical protein [Porphyromonas gingivalis]



CRISPR
classifications



NZ_KI259168.1 (15667 bp [-7358.5-22641.5]):
DEFINITION Porphyromonas gingivalis F0570 genomic scaffold Scaffold500. whole



CRISPR arrays in genome

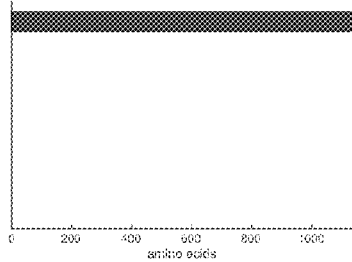
Dist	#	Repeat	Spacer	Consensus
0	4	31	29	GTTTCCTCCACCCCTTCTAAGTAAAGSUGTATTCGCAAGCA

Genes nearby protein

Index	Protein	Protein Length	Protein Name
-6	WP_043894287.1	358	cellulolactin biosynthesis protein CobD
-4	WP_021665471.1	355	trifunctional protein
-3	WP_021665472.1	496	aspartate aminotransferase
-2	WP_012458145.1	586	ATP-cob(I)alamin adenosyltransferase
-1	WP_043894296.1	439	aspartate aminotransferase
0	WP_021665475.1	1135	hypothetical protein
1	WP_021665476.1	238	hypothetical protein
2	WP_021665477.1	71	hypothetical protein
3	WP_021665479.1	280	phosphoserine phosphatase SerB
4	WP_043894295.1	246	Fe-S oxidoreductase
5	WP_021665481.1	455	iron-sulfur cluster-binding protein
6	WP_043894174.1	154	hypothetical protein
7	WP_004594294.1	457	thioesterase

FIG. 4ZZZ

WP_021877657.1 (1136 aa):
hypothetical protein [Porphyromonas gingivalis]



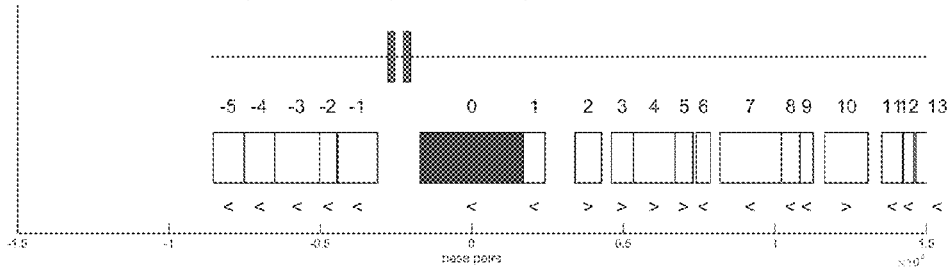
CRISPR
classifications

<Classification by Repeats>

36 nt repeat: Type II, IV

36 nt repeat: Type II, IV

NZ_KI259960.1 (64964 bp [-6417.5-23582.5]):
DEFINITION Porphyromonas gingivalis F0185 genomic scaffold Scaffold698, whole



CRISPR arrays in genome

Dist	#	Repeat	Spacer	Consensus
848	8	28	29	CTTGTGTCGCAAGCTTCGTAAGTAAAGGGTATTCCGAAAG
328	8	28	29	CTTGTGTCGCAAGCTTCGTAAGTAAAGGGTATTCCGAAAG

Genes nearby protein

Index	Protein	Protein Length	Protein Name
-3	WP_043894625.1	356	cobalamin biosynthesis protein CobD
-4	WP_021877656.1	335	treonine-phosphate decarboxylase
-3	WP_021877655.1	498	oxo-oxo acid synthase OxaO
-2	WP_021863192.1	180	ATP-guaf (histidin adenocyclotransferase
-1	WP_043894170.1	438	coenzyme A acid α -methylase synthase
0	WP_021877657.1	1136	hypothetical protein
1	WP_021877658.1	209	hydroxymethylase
2	WP_021877659.1	260	phosphoserine phosphatase SerB
3	WP_006873515.1	245	Fe-S oxidoreductase
4	WP_021877660.1	455	iron-sulfur cluster-binding protein
5	WP_004584286.1	184	hydroxymethylase
6	WP_021877661.1	167	hydroxymethylase
7	WP_021862134.1	572	OmpA family protein
8	WP_021877662.1	157	radical S-adenosyl domain protein
9	WP_006873450.1	159	6-pyruvoyl tetrahydropteridine synthase
10	WP_021864734.1	477	cysteine protease
11	WP_021877664.1	229	cysteine protease-like (ubiquitin domain) protein
12	WP_013815792.1	118	transcriptional regulator
13	WP_043894626.1	471	lipase

FIG. 4AAAA

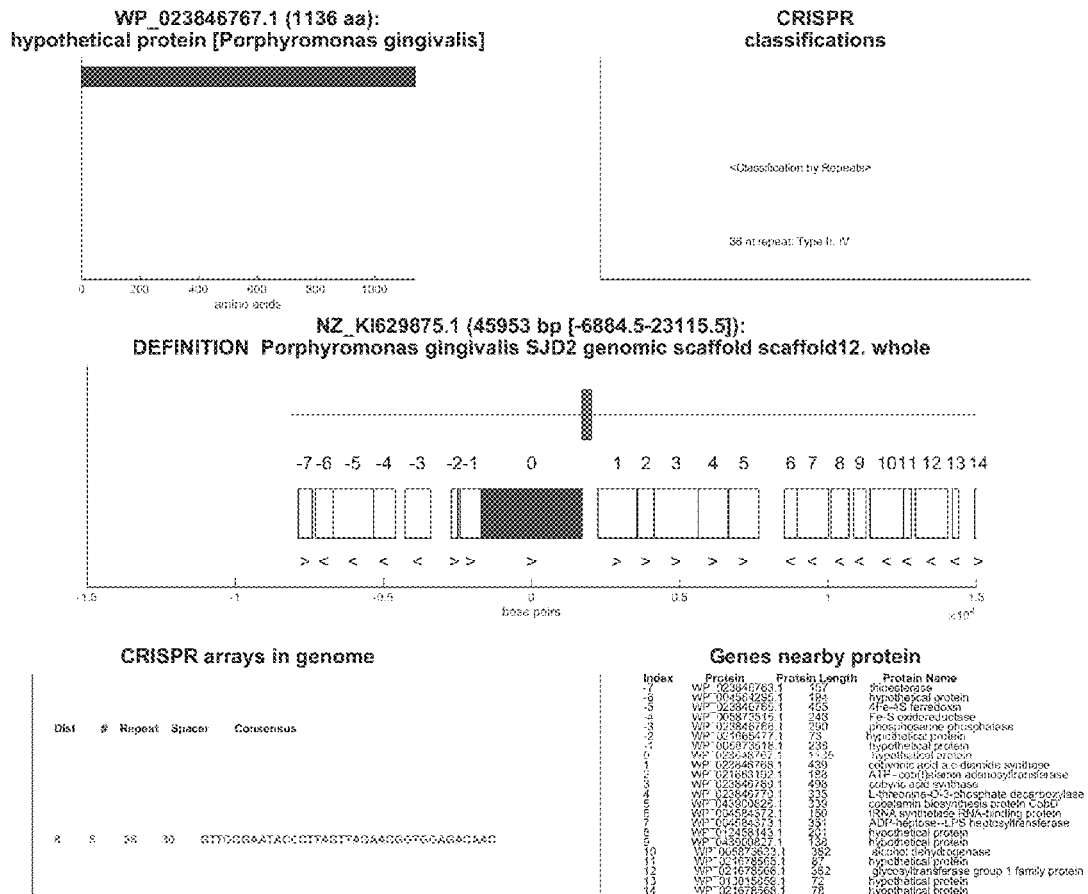


FIG. 4CCCC

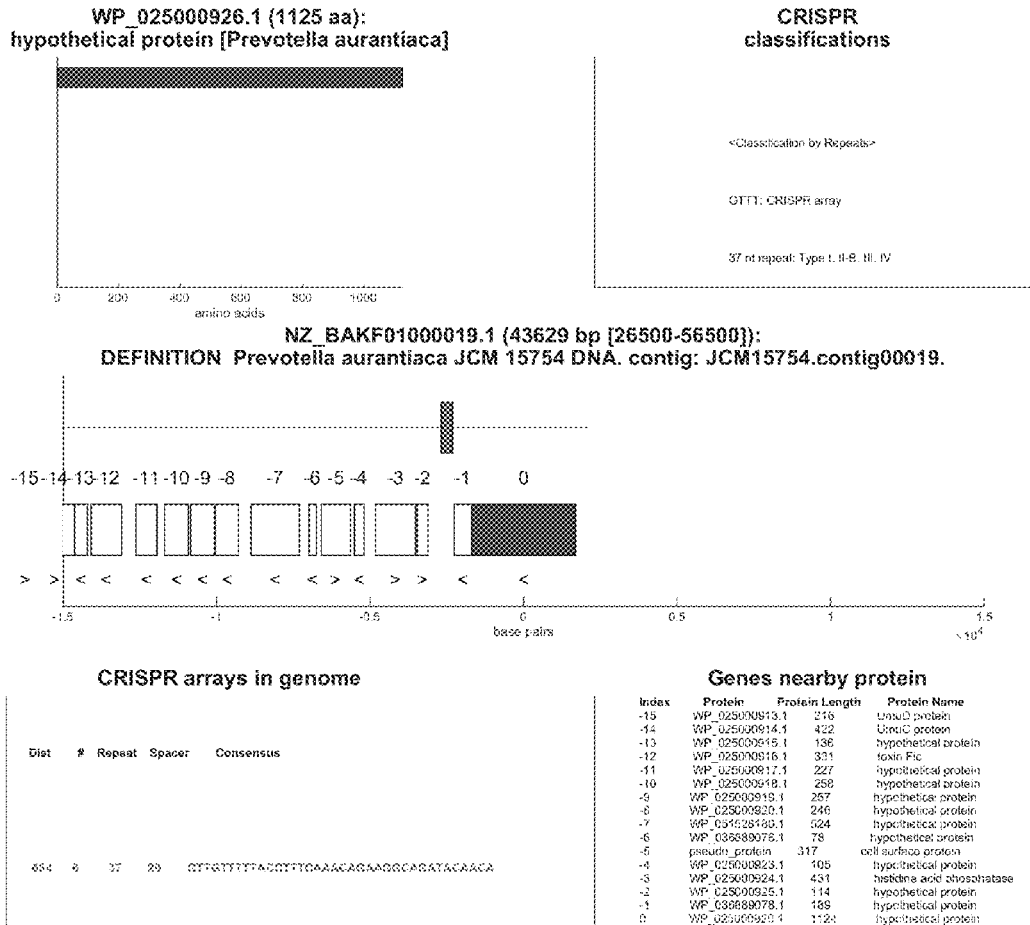
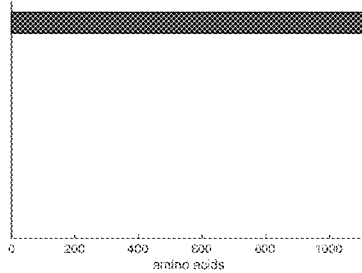


FIG. 4DDDD

WP_034542281.1 (1116 aa):
hypothetical protein [Bacteroides pyogenes]



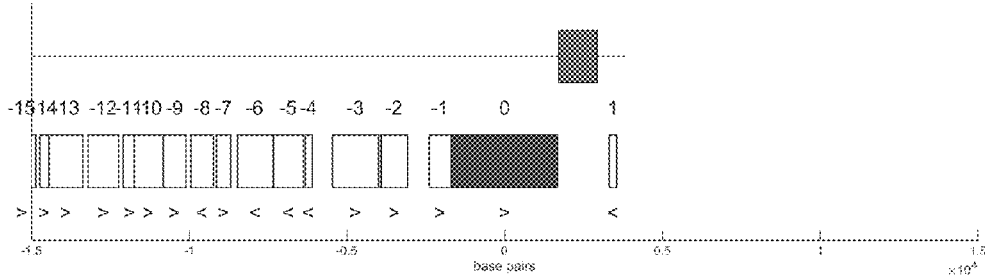
CRISPR
classifications

<Classification by Repeats>

GTTT CRISPR array

SS of repeat: Type II, IV

NZ_BAIU01000001.1 (332823 bp [314000.5-344000.5]):
DEFINITION Bacteroides pyogenes JCM 10003 DNA, contig:contig_1, whole genome



CRISPR arrays in genome

Dist	# Repeat	Spacer	Consensus
28	18	38	GTTTGTGGATGCGCACTAGTTTGGAGTGGATATGAGAGG

Genes nearby protein

Index	Protein	Protein Length	Protein Name
-15	WP_034542265.1	352	hypothetical protein
-14	WP_034543371.1	95	transcriptional regulator
-13	WP_034542267.1	357	hypothetical protein
-12	WP_034542267.1	321	DNA primase
-11	WP_034542268.1	127	nucleation protein
-10	WP_034542270.1	308	nucleation protein
-9	pseudo_protein	239	hypothetical protein
-8	WP_034542234.1	238	hypothetical protein
-7	WP_034542272.1	149	MarK family transcriptional regulator
-6	pseudo_protein	373	
-5	WP_034542273.1	312	integrase
-4	WP_034542276.1	74	hypothetical protein
-3	WP_034542277.1	482	histone family protein
-2	WP_034542277.1	278	hypothetical protein
-1	WP_034542278.1	285	hypothetical protein
0	WP_034542281.1	1116	hypothetical protein
1	WP_021678849.1	78	metal-binding protein

FIG. 4EEEE

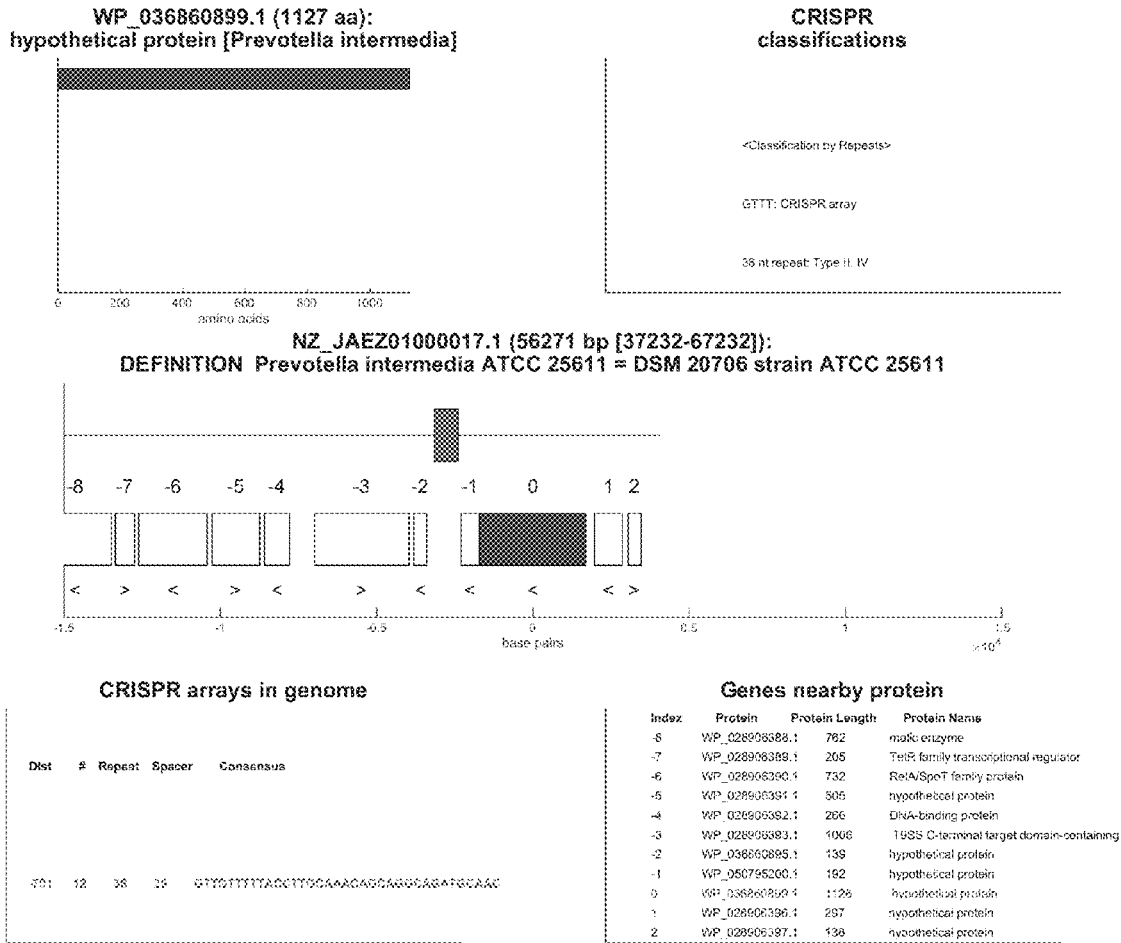


FIG. 4FFFF

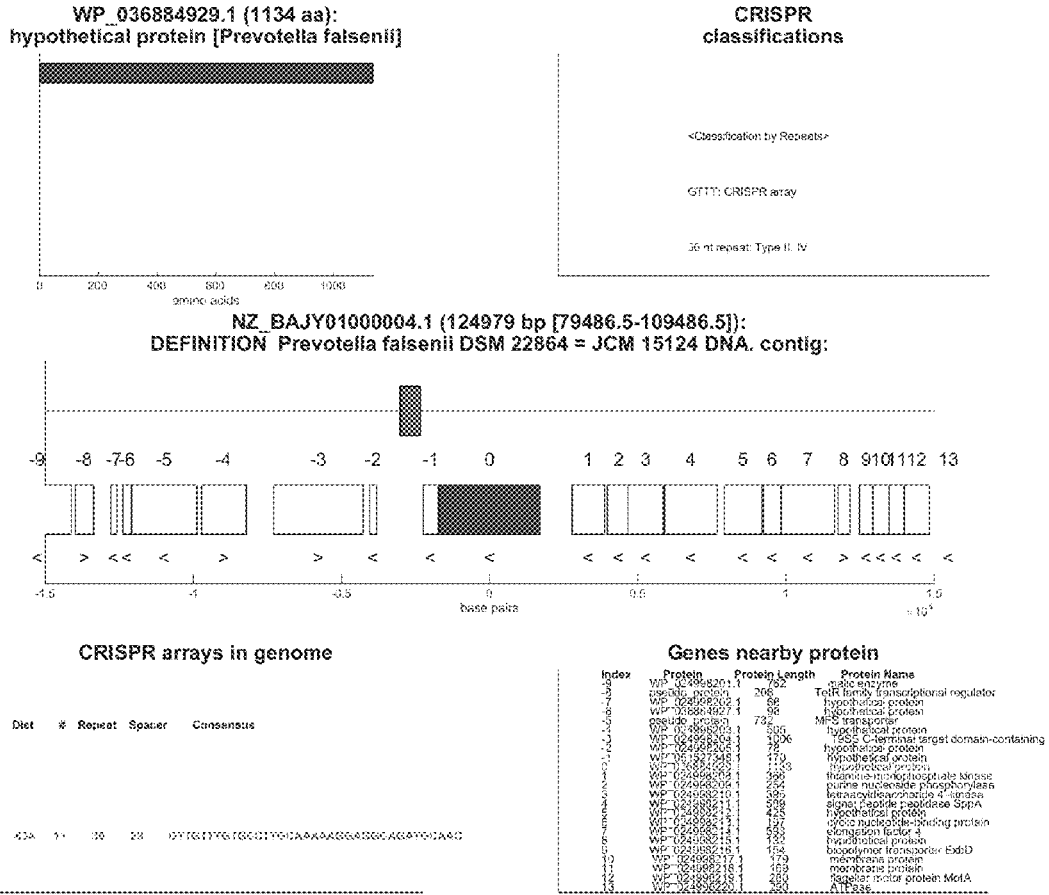


FIG. 4GGGG

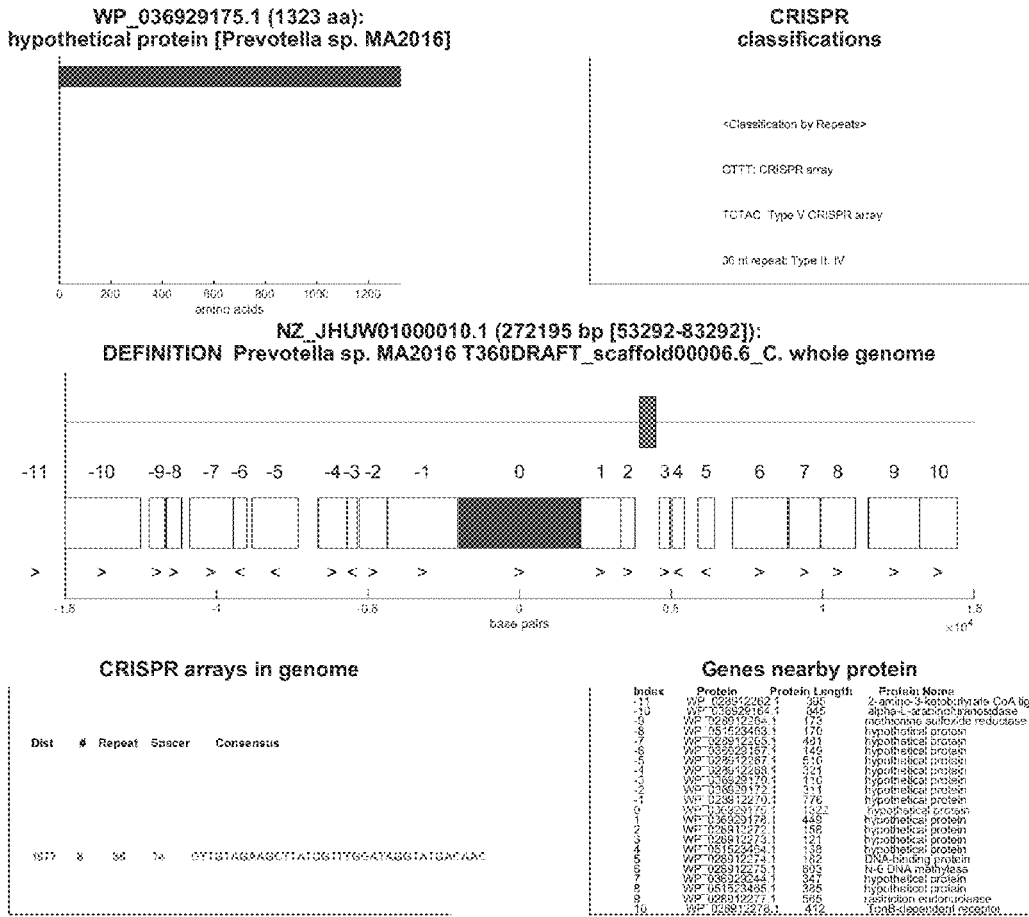


FIG. 4HHHH

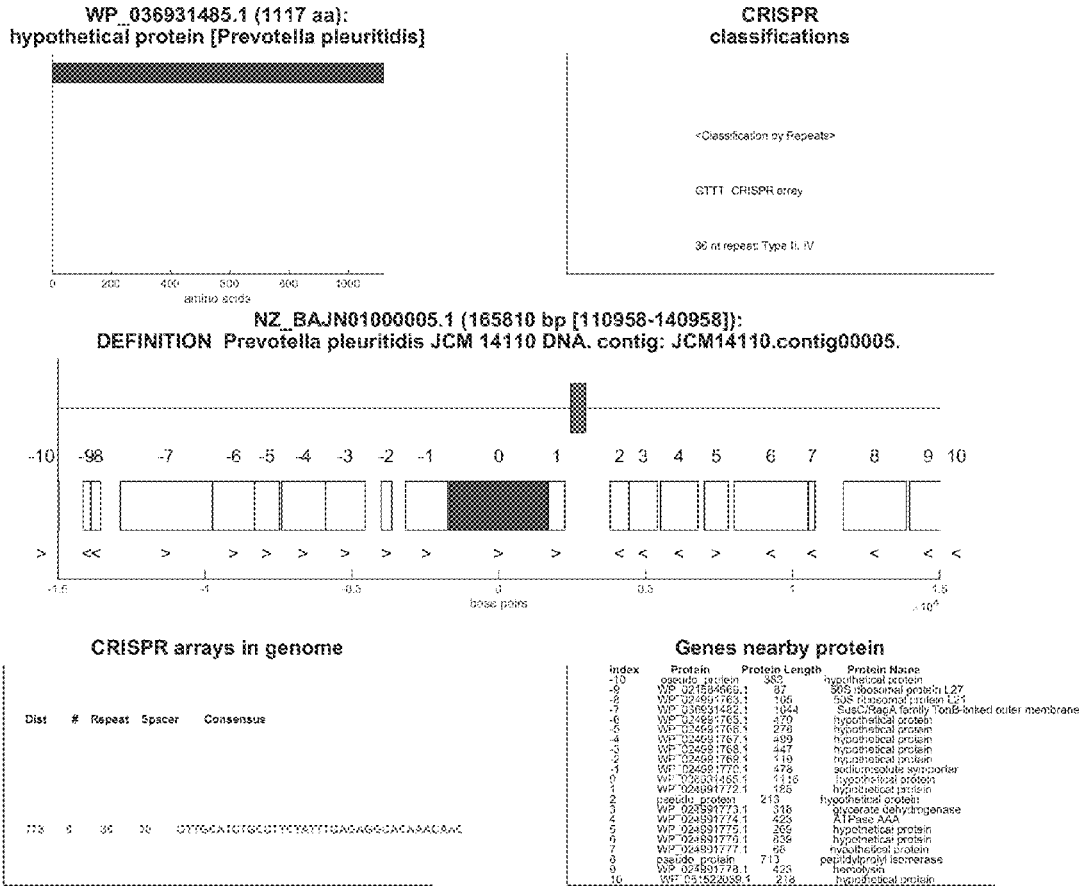
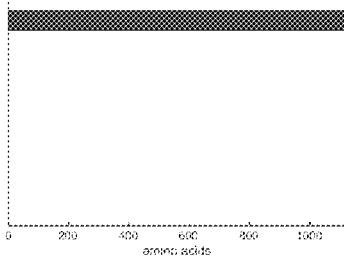


FIG. 4III

WP_039417390.1 (1120 aa):
hypothetical protein [Porphyromonas gingivalis]



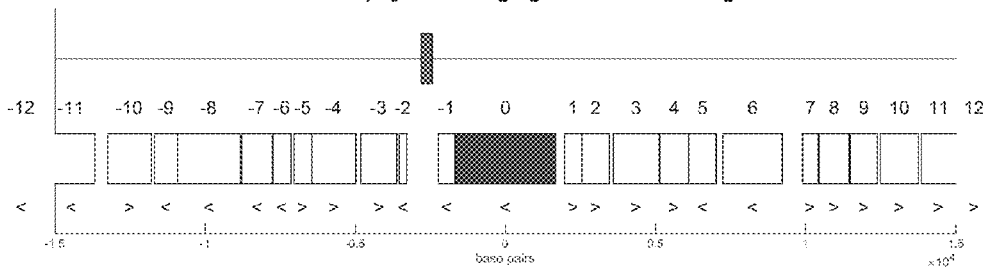
CRISPR
classifications

<Classification by Repeats>

TCTAC Type V CRISPR array

36 nt repeat Type II, IV

NZ_CP007756.1 (2441780 bp [2298856.5-2326856.5]):
DEFINITION Porphyromonas gingivalis strain HG66 genome.



CRISPR arrays in genome

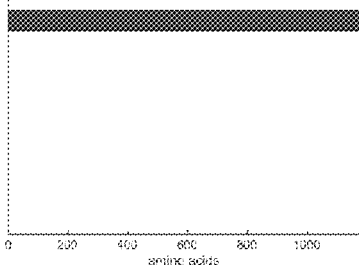
Dist	#	Repeat	Spacer	Consensus
-807026	36	30	73	GTCTTAAATTCGTTATGATGATGGAAATTGAAAT
-882448	15	30	39	GTTGCGAATACCTTAGTTAAAGCGGTGGAGACAC
-780	6	30	30	GTTCGATGATACGCTTCGAAATGAGCGGTGATGACAC

Genes nearby protein

Index	Protein	Protein Length	Protein Name
10	WP_039417390.1	1120	hypothetical protein
9	WP_039417389.1	255	hypothetical protein
8	WP_039417388.1	255	hypothetical protein
7	WP_039417387.1	255	hypothetical protein
6	WP_039417386.1	255	hypothetical protein
5	WP_039417385.1	255	hypothetical protein
4	WP_039417384.1	255	hypothetical protein
3	WP_039417383.1	255	hypothetical protein
2	WP_039417382.1	255	hypothetical protein
1	WP_039417381.1	255	hypothetical protein
0	WP_039417380.1	255	hypothetical protein
-1	WP_039417379.1	255	hypothetical protein
-2	WP_039417378.1	255	hypothetical protein
-3	WP_039417377.1	255	hypothetical protein
-4	WP_039417376.1	255	hypothetical protein
-5	WP_039417375.1	255	hypothetical protein
-6	WP_039417374.1	255	hypothetical protein
-7	WP_039417373.1	255	hypothetical protein
-8	WP_039417372.1	255	hypothetical protein
-9	WP_039417371.1	255	hypothetical protein
-10	WP_039417370.1	255	hypothetical protein
-11	WP_039417369.1	255	hypothetical protein
-12	WP_039417368.1	255	hypothetical protein

FIG. 4JJJ

WP_039418912.1 (1176 aa):
hypothetical protein [Porphyromonas gulae]



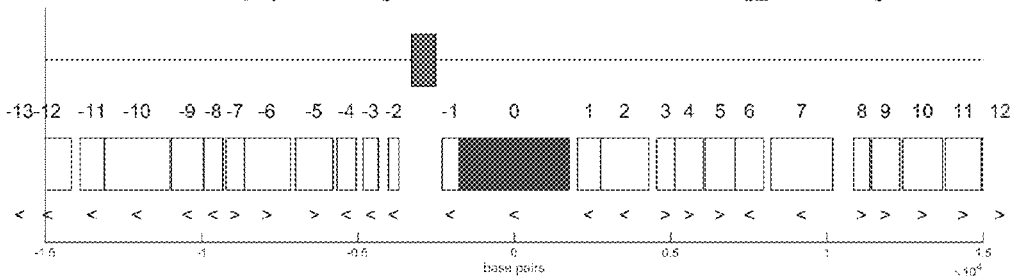
CRISPR
classifications

<Classification by Repeat>

TCTAC, Type V CRISPR array

26 nt repeat: Type II IV

NZ_JRAT01000012.1 (653056 bp [567073.5-597073.5]):
DEFINITION Porphyromonas gulae strain COT-052 OH3856 config_12, whole genome



CRISPR arrays in genome

Dist	# Repeat	Spacer	Consensus	
999097	18	17	18	GTCTTAATAGCTTACGACTGTGTATGTATAGTAGG
1088730	40	39	39	GTCTTAATAGCTTACGACTGTGTATGTATAGTAGG
1760	12	35	35	GTCTTAATAGCTTACGACTGTGTATGTATAGTAGG

Genes nearby protein

Index	Protein	Protein Length	Protein Name
1	hypothetical protein	1176	hypothetical protein
2	hypothetical protein	1176	hypothetical protein
3	hypothetical protein	1176	hypothetical protein
4	hypothetical protein	1176	hypothetical protein
5	hypothetical protein	1176	hypothetical protein
6	hypothetical protein	1176	hypothetical protein
7	hypothetical protein	1176	hypothetical protein
8	hypothetical protein	1176	hypothetical protein
9	hypothetical protein	1176	hypothetical protein
10	hypothetical protein	1176	hypothetical protein
11	hypothetical protein	1176	hypothetical protein
12	hypothetical protein	1176	hypothetical protein
13	hypothetical protein	1176	hypothetical protein
14	hypothetical protein	1176	hypothetical protein
15	hypothetical protein	1176	hypothetical protein
16	hypothetical protein	1176	hypothetical protein
17	hypothetical protein	1176	hypothetical protein
18	hypothetical protein	1176	hypothetical protein
19	hypothetical protein	1176	hypothetical protein
20	hypothetical protein	1176	hypothetical protein
21	hypothetical protein	1176	hypothetical protein
22	hypothetical protein	1176	hypothetical protein
23	hypothetical protein	1176	hypothetical protein
24	hypothetical protein	1176	hypothetical protein
25	hypothetical protein	1176	hypothetical protein
26	hypothetical protein	1176	hypothetical protein
27	hypothetical protein	1176	hypothetical protein
28	hypothetical protein	1176	hypothetical protein
29	hypothetical protein	1176	hypothetical protein
30	hypothetical protein	1176	hypothetical protein
31	hypothetical protein	1176	hypothetical protein
32	hypothetical protein	1176	hypothetical protein
33	hypothetical protein	1176	hypothetical protein
34	hypothetical protein	1176	hypothetical protein
35	hypothetical protein	1176	hypothetical protein
36	hypothetical protein	1176	hypothetical protein
37	hypothetical protein	1176	hypothetical protein
38	hypothetical protein	1176	hypothetical protein
39	hypothetical protein	1176	hypothetical protein
40	hypothetical protein	1176	hypothetical protein
41	hypothetical protein	1176	hypothetical protein
42	hypothetical protein	1176	hypothetical protein
43	hypothetical protein	1176	hypothetical protein
44	hypothetical protein	1176	hypothetical protein
45	hypothetical protein	1176	hypothetical protein
46	hypothetical protein	1176	hypothetical protein
47	hypothetical protein	1176	hypothetical protein
48	hypothetical protein	1176	hypothetical protein
49	hypothetical protein	1176	hypothetical protein
50	hypothetical protein	1176	hypothetical protein
51	hypothetical protein	1176	hypothetical protein
52	hypothetical protein	1176	hypothetical protein
53	hypothetical protein	1176	hypothetical protein
54	hypothetical protein	1176	hypothetical protein
55	hypothetical protein	1176	hypothetical protein
56	hypothetical protein	1176	hypothetical protein
57	hypothetical protein	1176	hypothetical protein
58	hypothetical protein	1176	hypothetical protein
59	hypothetical protein	1176	hypothetical protein
60	hypothetical protein	1176	hypothetical protein
61	hypothetical protein	1176	hypothetical protein
62	hypothetical protein	1176	hypothetical protein
63	hypothetical protein	1176	hypothetical protein
64	hypothetical protein	1176	hypothetical protein
65	hypothetical protein	1176	hypothetical protein
66	hypothetical protein	1176	hypothetical protein
67	hypothetical protein	1176	hypothetical protein
68	hypothetical protein	1176	hypothetical protein
69	hypothetical protein	1176	hypothetical protein
70	hypothetical protein	1176	hypothetical protein
71	hypothetical protein	1176	hypothetical protein
72	hypothetical protein	1176	hypothetical protein
73	hypothetical protein	1176	hypothetical protein
74	hypothetical protein	1176	hypothetical protein
75	hypothetical protein	1176	hypothetical protein
76	hypothetical protein	1176	hypothetical protein
77	hypothetical protein	1176	hypothetical protein
78	hypothetical protein	1176	hypothetical protein
79	hypothetical protein	1176	hypothetical protein
80	hypothetical protein	1176	hypothetical protein
81	hypothetical protein	1176	hypothetical protein
82	hypothetical protein	1176	hypothetical protein
83	hypothetical protein	1176	hypothetical protein
84	hypothetical protein	1176	hypothetical protein
85	hypothetical protein	1176	hypothetical protein
86	hypothetical protein	1176	hypothetical protein
87	hypothetical protein	1176	hypothetical protein
88	hypothetical protein	1176	hypothetical protein
89	hypothetical protein	1176	hypothetical protein
90	hypothetical protein	1176	hypothetical protein
91	hypothetical protein	1176	hypothetical protein
92	hypothetical protein	1176	hypothetical protein
93	hypothetical protein	1176	hypothetical protein
94	hypothetical protein	1176	hypothetical protein
95	hypothetical protein	1176	hypothetical protein
96	hypothetical protein	1176	hypothetical protein
97	hypothetical protein	1176	hypothetical protein
98	hypothetical protein	1176	hypothetical protein
99	hypothetical protein	1176	hypothetical protein
100	hypothetical protein	1176	hypothetical protein

FIG. 4K K K K

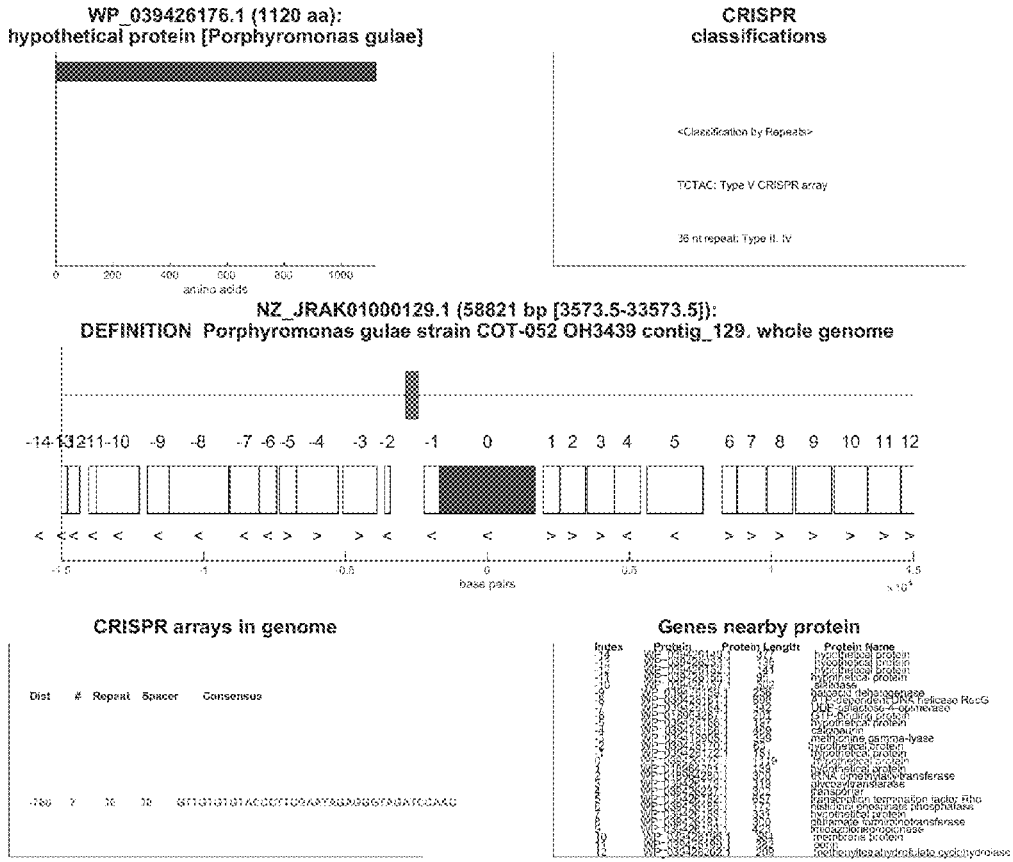


FIG. 4MMMM

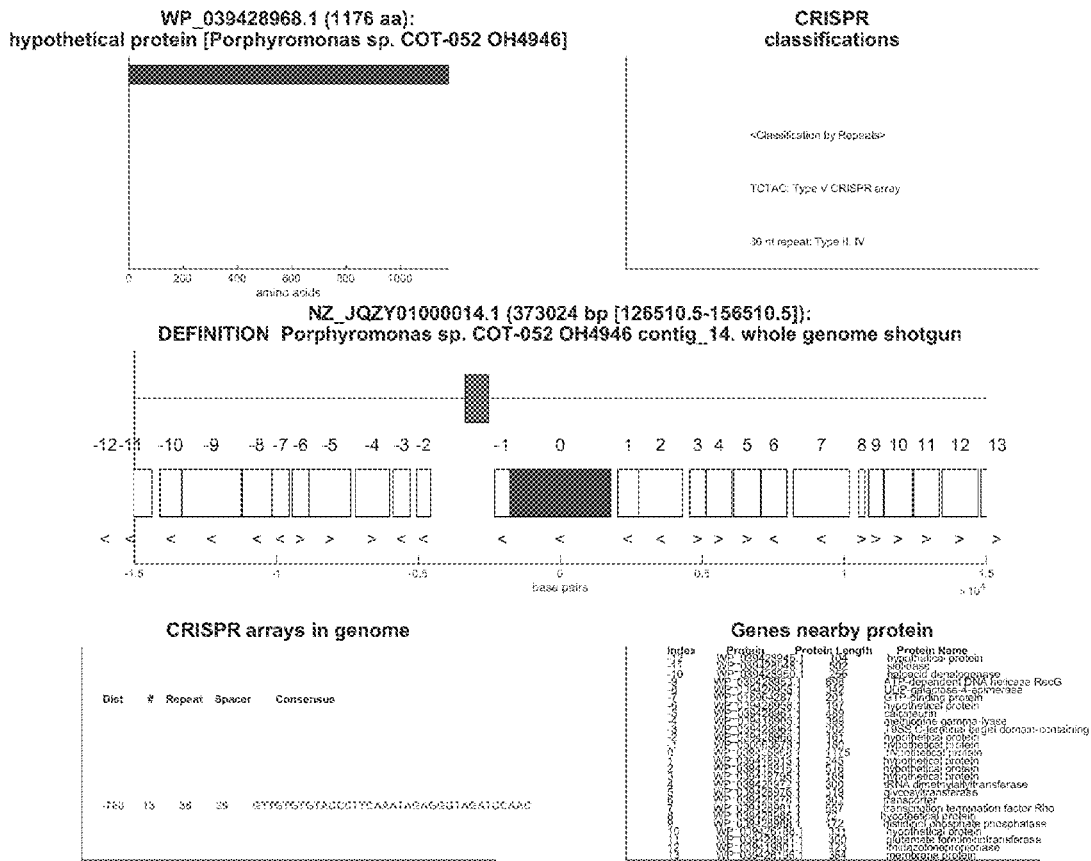


FIG. 4NNNN

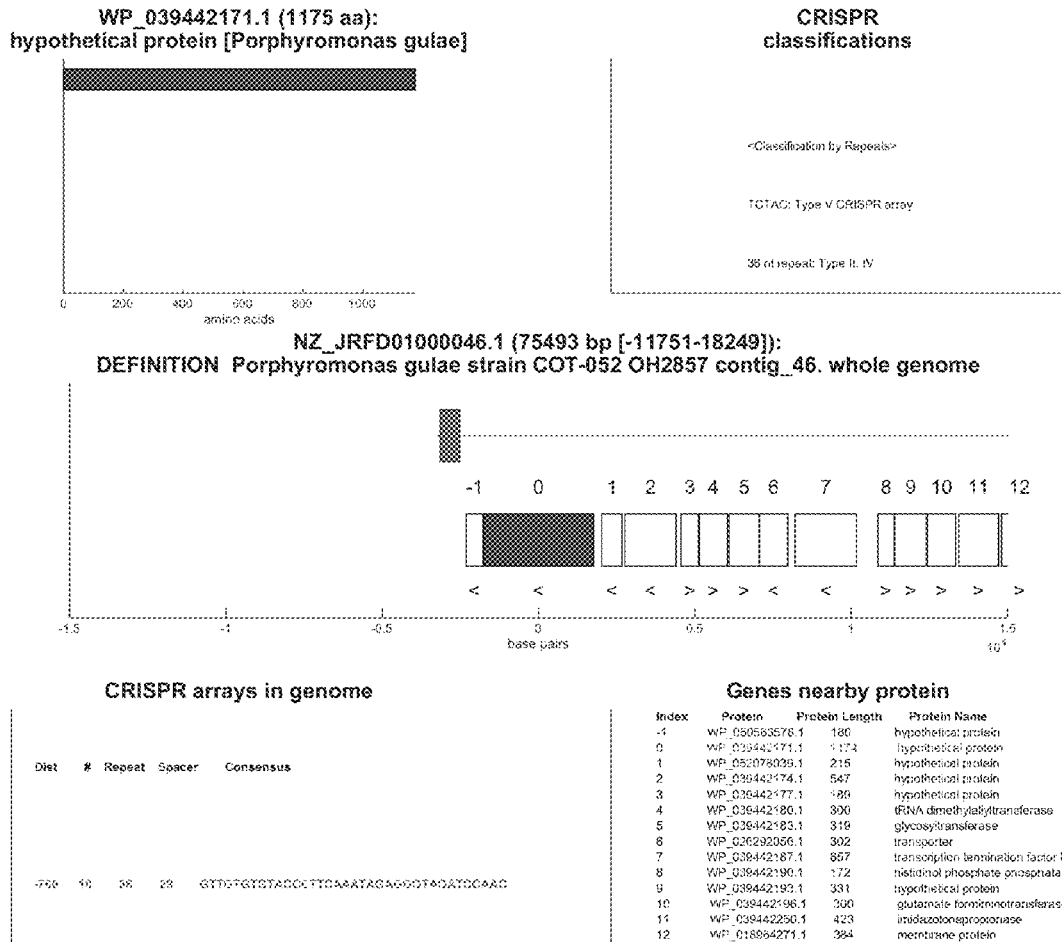


FIG. 4QQQQ

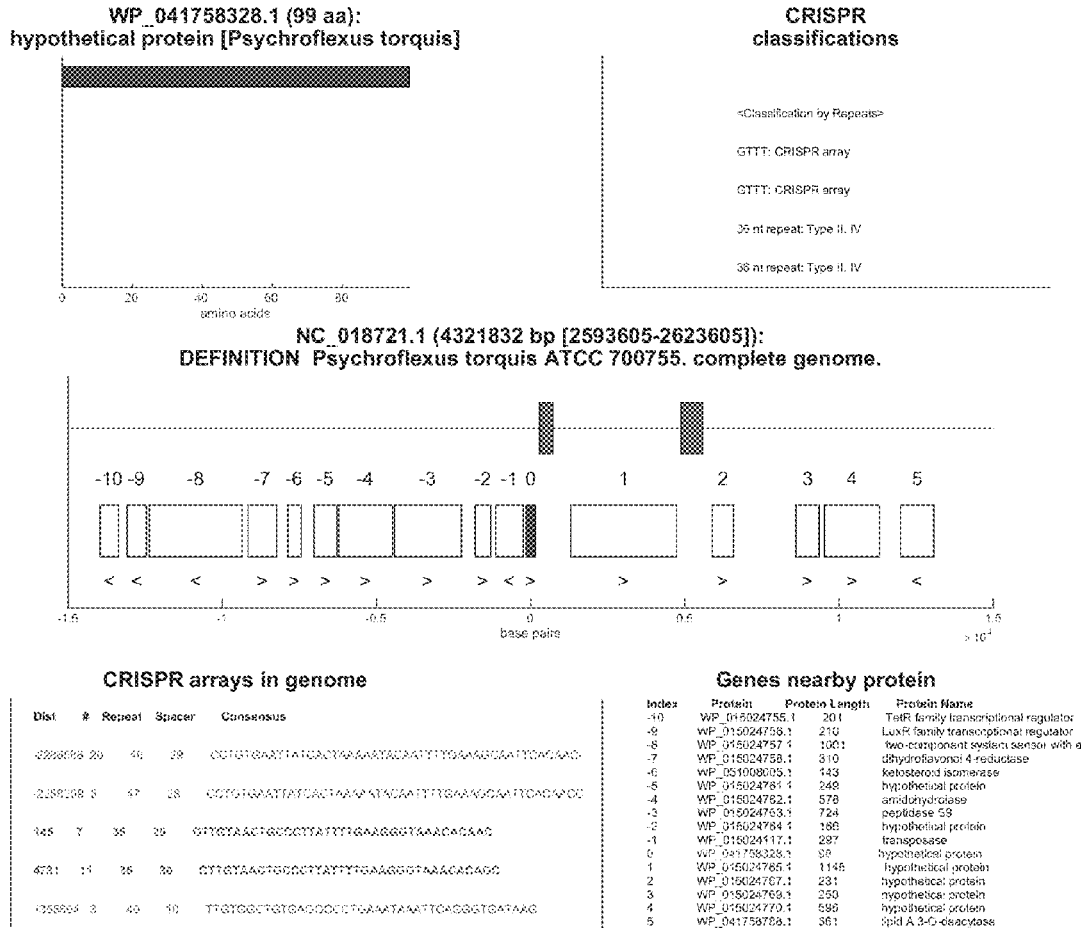


FIG. 4SSSS

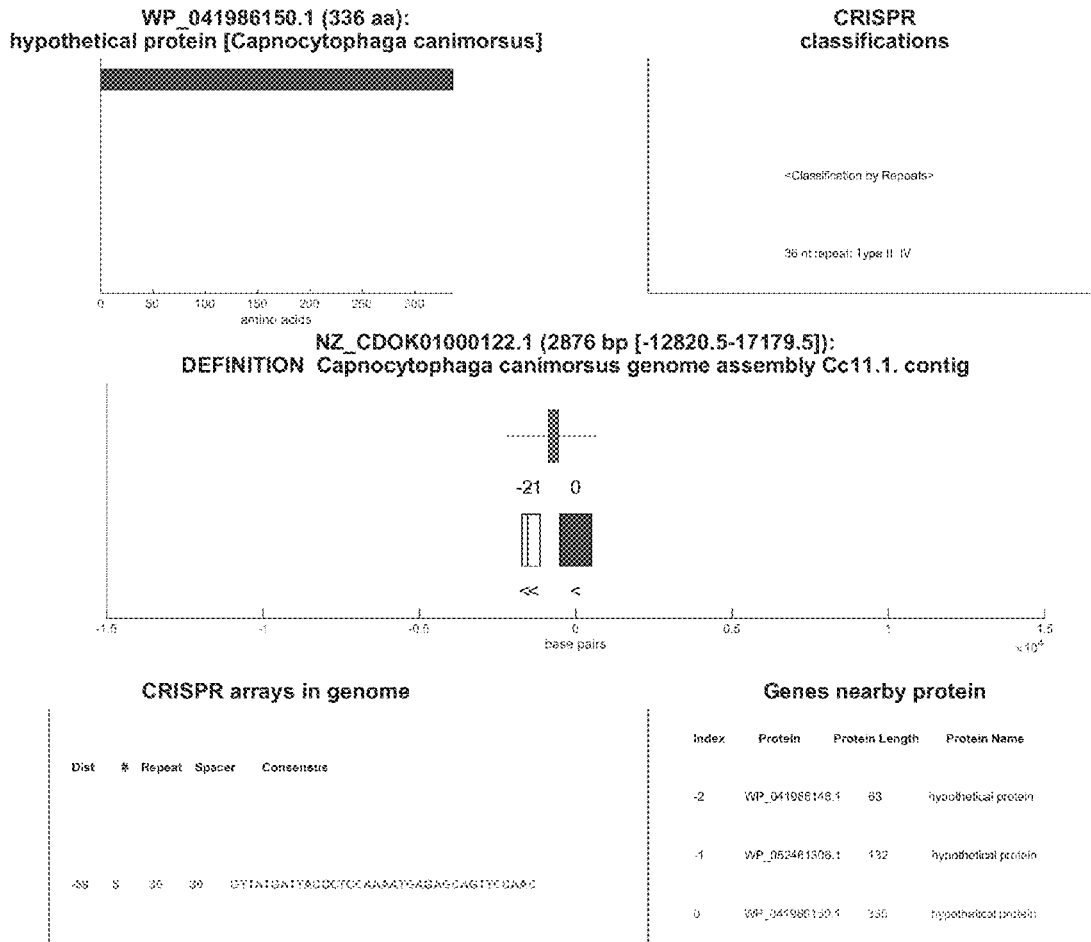
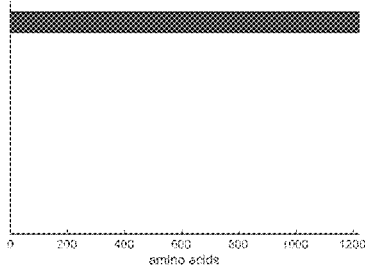


FIG. 4TTTT

WP_041989581.1 (1219 aa):
hypothetical protein [Capnocytophaga cynodegmi]

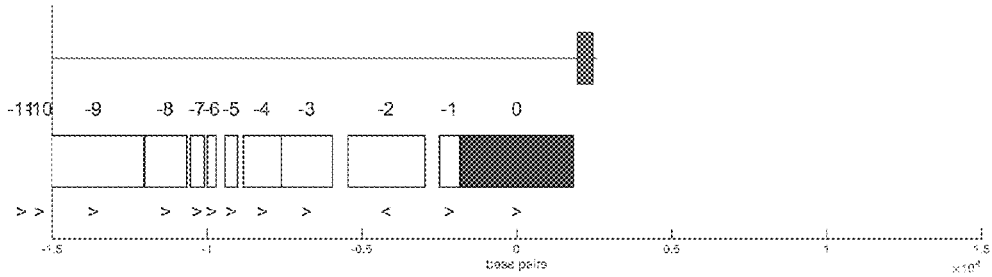


CRISPR
classifications

<Classification by Repetition>

36 nt repeat Type II, IV

NZ_CDOD01008002.1 (83586 bp [66022-96022]):
DEFINITION Capnocytophaga cynodegmi genome assembly Ccyn2B.1. contig



CRISPR arrays in genome

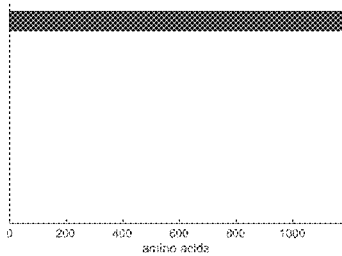
Dist	#	Repeat	Spacer	Consensus
925	2	26	26	GTTCGAAAGCTGCTGTGATTTTCGAGGCTTAAFLACAAAC

Genes nearby protein

Index	Protein	Protein Length	Protein Name
-11	WP_041989551.1	240	hypothetical protein
-10	WP_018360123.1	67	transcriptional regulator
-9	WP_041989558.1	1062	DNA/RNA helicase
-8	WP_041989561.1	452	transcriptional regulator
-7	WP_041989586.1	151	GCN5 family acetyltransferase
-6	WP_052456862.1	95	hypothetical protein
-5	WP_041989610.1	133	polyketide cyclase
-4	WP_041989570.1	406	hypothetical protein
-3	WP_041989572.1	549	hypothetical protein
-2	WP_041989575.1	828	ATPase AAA
-1	WP_041989576.1	219	hypothetical protein
0	WP_041989581.1	1219	hypothetical protein

FIG. 4UUUU

WP_043894148.1 (1176 aa):
hypothetical protein [*Porphyromonas gingivalis*]



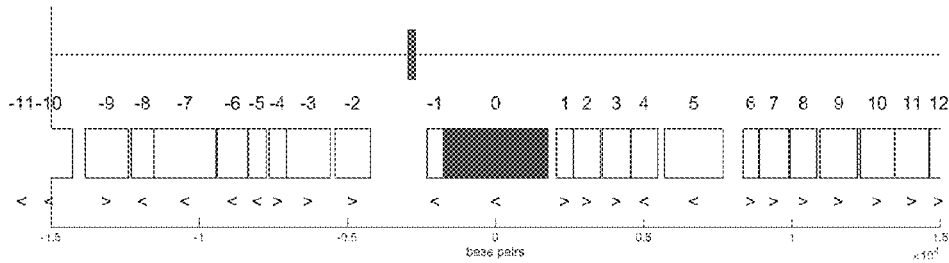
CRISPR
classifications

<Classification by Repeats>

TCTAC: Type V CRISPR array

37 nt repeat: Type I, II-B, III, IV

NZ_KI259080.1 (104120 bp [25184.5-55184.5]):
DEFINITION *Porphyromonas gingivalis* F0569 genomic scaffold Scaffold677. whole



CRISPR arrays in genome

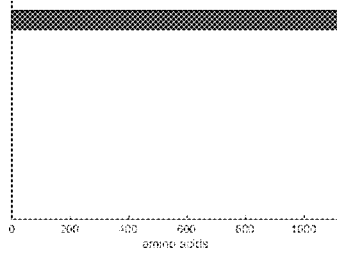
Dist	# Repeat	Spacer	Consensus	
954	4	37	29	CTTGTGTTGTTTACGTTTCTTGGATATAGAGCGGTCAGATCCAGCA

Genes nearby protein

Index	Protein	Protein Length	Protein Name
-11	WP_043894148.1	98	hypothetical protein
-10	WP_043894148.1	98	hypothetical protein
-9	WP_043894148.1	98	hypothetical protein
-8	WP_043894148.1	98	hypothetical protein
-7	WP_043894148.1	98	hypothetical protein
-6	WP_043894148.1	98	hypothetical protein
-5	WP_043894148.1	98	hypothetical protein
-4	WP_043894148.1	98	hypothetical protein
-3	WP_043894148.1	98	hypothetical protein
-2	WP_043894148.1	98	hypothetical protein
-1	WP_043894148.1	98	hypothetical protein
0	WP_043894148.1	98	hypothetical protein
1	WP_043894148.1	98	hypothetical protein
2	WP_043894148.1	98	hypothetical protein
3	WP_043894148.1	98	hypothetical protein
4	WP_043894148.1	98	hypothetical protein
5	WP_043894148.1	98	hypothetical protein
6	WP_043894148.1	98	hypothetical protein
7	WP_043894148.1	98	hypothetical protein
8	WP_043894148.1	98	hypothetical protein
9	WP_043894148.1	98	hypothetical protein
10	WP_043894148.1	98	hypothetical protein
11	WP_043894148.1	98	hypothetical protein
12	WP_043894148.1	98	hypothetical protein

FIG. 4WWWW

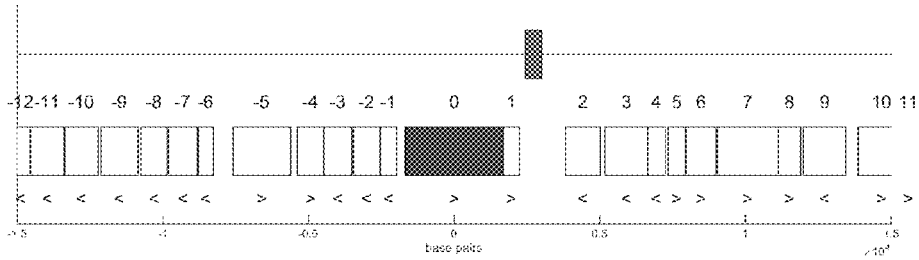
WP_043894372.1 (1120 aa):
 hypothetical protein [Porphyromonas gingivalis]



CRISPR
 classifications

<Classification by Repeat>
 TCTAG: Type V CRISPR array
 38 nt repeat: Type II, IV

NZ_KI259218.1 (340581 bp [60214.5-90214.5]):
 DEFINITION Porphyromonas gingivalis F0570 genomic scaffold Scaffold965. whole



CRISPR arrays in genome

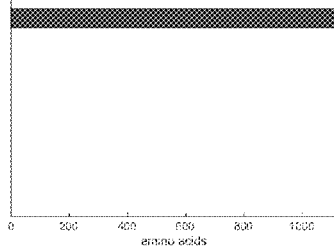
Dist	#	Repeat	Spacer	Consensus
598	9	38	39	GTTGGGATYAGCTCTCTATATTCGAAAGGGTACACGACAC

Genes nearby protein

Index	Protein	Protein Length	Protein Name
114	WP_043894372.1	1120	hypothetical protein
115	WP_043894372.1	1120	hypothetical protein
116	WP_043894372.1	1120	hypothetical protein
117	WP_043894372.1	1120	hypothetical protein
118	WP_043894372.1	1120	hypothetical protein
119	WP_043894372.1	1120	hypothetical protein
120	WP_043894372.1	1120	hypothetical protein
121	WP_043894372.1	1120	hypothetical protein
122	WP_043894372.1	1120	hypothetical protein
123	WP_043894372.1	1120	hypothetical protein
124	WP_043894372.1	1120	hypothetical protein
125	WP_043894372.1	1120	hypothetical protein
126	WP_043894372.1	1120	hypothetical protein
127	WP_043894372.1	1120	hypothetical protein
128	WP_043894372.1	1120	hypothetical protein
129	WP_043894372.1	1120	hypothetical protein
130	WP_043894372.1	1120	hypothetical protein
131	WP_043894372.1	1120	hypothetical protein
132	WP_043894372.1	1120	hypothetical protein
133	WP_043894372.1	1120	hypothetical protein
134	WP_043894372.1	1120	hypothetical protein
135	WP_043894372.1	1120	hypothetical protein
136	WP_043894372.1	1120	hypothetical protein
137	WP_043894372.1	1120	hypothetical protein
138	WP_043894372.1	1120	hypothetical protein
139	WP_043894372.1	1120	hypothetical protein
140	WP_043894372.1	1120	hypothetical protein

FIG. 4XXXX

WP_043894751.1 (1120 aa):
hypothetical protein [Porphyromonas gingivalis]



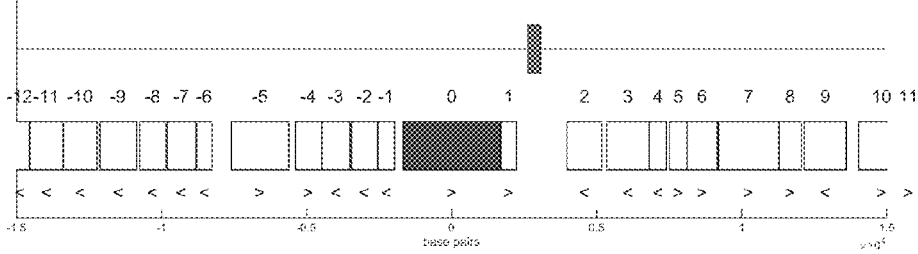
CRISPR
classifications

<Classification by Repeats>

TCTAG: Type V CRISPR array

36 nt repeat: Type II, IV

NZ_KI259867.1 (172891 bp [48682.5-78682.5]):
DEFINITION Porphyromonas gingivalis F0185 genomic scaffold Scaffold1008. whole



CRISPR arrays in genome

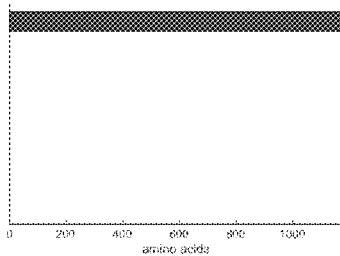
Dist	#	Repeat	Spacer	Consensus
547	7	36	10	CTTGGATCTAGCGTGTATTGSAAGGGTACAGAAAC

Genes nearby protein

Index	Protein	Protein Length	Protein Name
-12	WP_021870664.1	409	formate dehydrogenase-cytochrome b558
-11	WP_043894751.1	1120	hypothetical protein
-10	WP_021870664.1	409	formate dehydrogenase-cytochrome b558
-9	WP_021870664.1	409	formate dehydrogenase-cytochrome b558
-8	WP_043894751.1	1120	hypothetical protein
-7	WP_021870664.1	409	formate dehydrogenase-cytochrome b558
-6	WP_021870664.1	409	formate dehydrogenase-cytochrome b558
-5	WP_021870664.1	409	formate dehydrogenase-cytochrome b558
-4	WP_021870664.1	409	formate dehydrogenase-cytochrome b558
-3	WP_021870664.1	409	formate dehydrogenase-cytochrome b558
-2	WP_021870664.1	409	formate dehydrogenase-cytochrome b558
-1	WP_021870664.1	409	formate dehydrogenase-cytochrome b558
0	WP_043894751.1	1120	hypothetical protein
1	WP_021870664.1	409	formate dehydrogenase-cytochrome b558
2	WP_021870664.1	409	formate dehydrogenase-cytochrome b558
3	WP_021870664.1	409	formate dehydrogenase-cytochrome b558
4	WP_021870664.1	409	formate dehydrogenase-cytochrome b558
5	WP_021870664.1	409	formate dehydrogenase-cytochrome b558
6	WP_021870664.1	409	formate dehydrogenase-cytochrome b558
7	WP_021870664.1	409	formate dehydrogenase-cytochrome b558
8	WP_021870664.1	409	formate dehydrogenase-cytochrome b558
9	WP_021870664.1	409	formate dehydrogenase-cytochrome b558
10	WP_021870664.1	409	formate dehydrogenase-cytochrome b558
11	WP_043894751.1	1120	hypothetical protein

FIG. 4ZZZZ

WP_043895240.1 (1176 aa):
hypothetical protein [Porphyromonas gingivalis]



CRISPR
classifications

<Classification by Repeat>

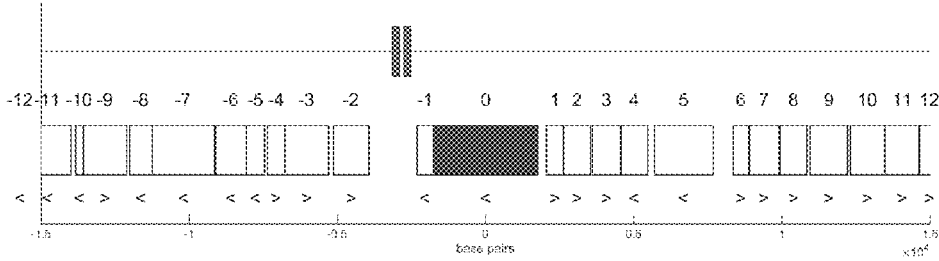
TCTAC: Type V CRISPR array

TCTAC: Type V CRISPR array

36 nt repeat: Type II, IV

36 nt repeat: Type II, IV

NZ_KI260014.1 (193603 bp [151708.5-181708.5]):
DEFINITION Porphyromonas gingivalis F0566 genomic scaffold Scaffold1339. whole



CRISPR arrays in genome

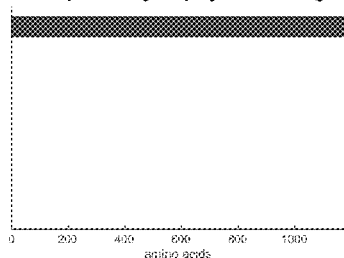
Dist	#	Repeat	Spacer	Consensus
1152	4	28	28	GTTTGTCTTCTAGCCCTTCTGSAATAGAGGCTGATGCAAGC
1260	4	28	36	GTTTGTCTTCTAGCCCTTCTGSAATAGAGGCTGATGCAAGC

Genes nearby protein

Index	Protein	Protein Length	Protein Name
1	hypothetical protein	360	hypothetical protein
2	hypothetical protein	360	hypothetical protein
3	hypothetical protein	360	hypothetical protein
4	hypothetical protein	360	hypothetical protein
5	hypothetical protein	360	hypothetical protein
6	hypothetical protein	360	hypothetical protein
7	hypothetical protein	360	hypothetical protein
8	hypothetical protein	360	hypothetical protein
9	hypothetical protein	360	hypothetical protein
10	hypothetical protein	360	hypothetical protein
11	hypothetical protein	360	hypothetical protein
12	hypothetical protein	360	hypothetical protein
13	hypothetical protein	360	hypothetical protein
14	hypothetical protein	360	hypothetical protein
15	hypothetical protein	360	hypothetical protein
16	hypothetical protein	360	hypothetical protein
17	hypothetical protein	360	hypothetical protein
18	hypothetical protein	360	hypothetical protein
19	hypothetical protein	360	hypothetical protein
20	hypothetical protein	360	hypothetical protein
21	hypothetical protein	360	hypothetical protein
22	hypothetical protein	360	hypothetical protein
23	hypothetical protein	360	hypothetical protein
24	hypothetical protein	360	hypothetical protein
25	hypothetical protein	360	hypothetical protein
26	hypothetical protein	360	hypothetical protein
27	hypothetical protein	360	hypothetical protein
28	hypothetical protein	360	hypothetical protein
29	hypothetical protein	360	hypothetical protein
30	hypothetical protein	360	hypothetical protein
31	hypothetical protein	360	hypothetical protein
32	hypothetical protein	360	hypothetical protein
33	hypothetical protein	360	hypothetical protein
34	hypothetical protein	360	hypothetical protein
35	hypothetical protein	360	hypothetical protein
36	hypothetical protein	360	hypothetical protein
37	hypothetical protein	360	hypothetical protein
38	hypothetical protein	360	hypothetical protein
39	hypothetical protein	360	hypothetical protein
40	hypothetical protein	360	hypothetical protein
41	hypothetical protein	360	hypothetical protein
42	hypothetical protein	360	hypothetical protein
43	hypothetical protein	360	hypothetical protein
44	hypothetical protein	360	hypothetical protein
45	hypothetical protein	360	hypothetical protein
46	hypothetical protein	360	hypothetical protein
47	hypothetical protein	360	hypothetical protein
48	hypothetical protein	360	hypothetical protein
49	hypothetical protein	360	hypothetical protein
50	hypothetical protein	360	hypothetical protein
51	hypothetical protein	360	hypothetical protein
52	hypothetical protein	360	hypothetical protein
53	hypothetical protein	360	hypothetical protein
54	hypothetical protein	360	hypothetical protein
55	hypothetical protein	360	hypothetical protein
56	hypothetical protein	360	hypothetical protein
57	hypothetical protein	360	hypothetical protein
58	hypothetical protein	360	hypothetical protein
59	hypothetical protein	360	hypothetical protein
60	hypothetical protein	360	hypothetical protein
61	hypothetical protein	360	hypothetical protein
62	hypothetical protein	360	hypothetical protein
63	hypothetical protein	360	hypothetical protein
64	hypothetical protein	360	hypothetical protein
65	hypothetical protein	360	hypothetical protein
66	hypothetical protein	360	hypothetical protein
67	hypothetical protein	360	hypothetical protein
68	hypothetical protein	360	hypothetical protein
69	hypothetical protein	360	hypothetical protein
70	hypothetical protein	360	hypothetical protein
71	hypothetical protein	360	hypothetical protein
72	hypothetical protein	360	hypothetical protein
73	hypothetical protein	360	hypothetical protein
74	hypothetical protein	360	hypothetical protein
75	hypothetical protein	360	hypothetical protein
76	hypothetical protein	360	hypothetical protein
77	hypothetical protein	360	hypothetical protein
78	hypothetical protein	360	hypothetical protein
79	hypothetical protein	360	hypothetical protein
80	hypothetical protein	360	hypothetical protein
81	hypothetical protein	360	hypothetical protein
82	hypothetical protein	360	hypothetical protein
83	hypothetical protein	360	hypothetical protein
84	hypothetical protein	360	hypothetical protein
85	hypothetical protein	360	hypothetical protein
86	hypothetical protein	360	hypothetical protein
87	hypothetical protein	360	hypothetical protein
88	hypothetical protein	360	hypothetical protein
89	hypothetical protein	360	hypothetical protein
90	hypothetical protein	360	hypothetical protein
91	hypothetical protein	360	hypothetical protein
92	hypothetical protein	360	hypothetical protein
93	hypothetical protein	360	hypothetical protein
94	hypothetical protein	360	hypothetical protein
95	hypothetical protein	360	hypothetical protein
96	hypothetical protein	360	hypothetical protein
97	hypothetical protein	360	hypothetical protein
98	hypothetical protein	360	hypothetical protein
99	hypothetical protein	360	hypothetical protein
100	hypothetical protein	360	hypothetical protein

FIG. 4AAAAA

WP_043895933.1 (1176 aa):
hypothetical protein [Porphyromonas gingivalis]



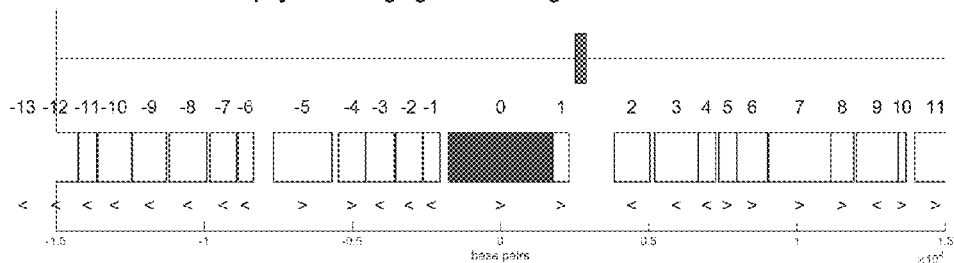
CRISPR
classifications

<Classification by Repeats>

TCTAC: Type V CRISPR array

86 nt repeat: Type II, IV

NZ_KI258974.1 (168213 bp [22428.5-52428.5]):
DEFINITION Porphyromonas gingivalis F0568 genomic scaffold Scaffold847. whole



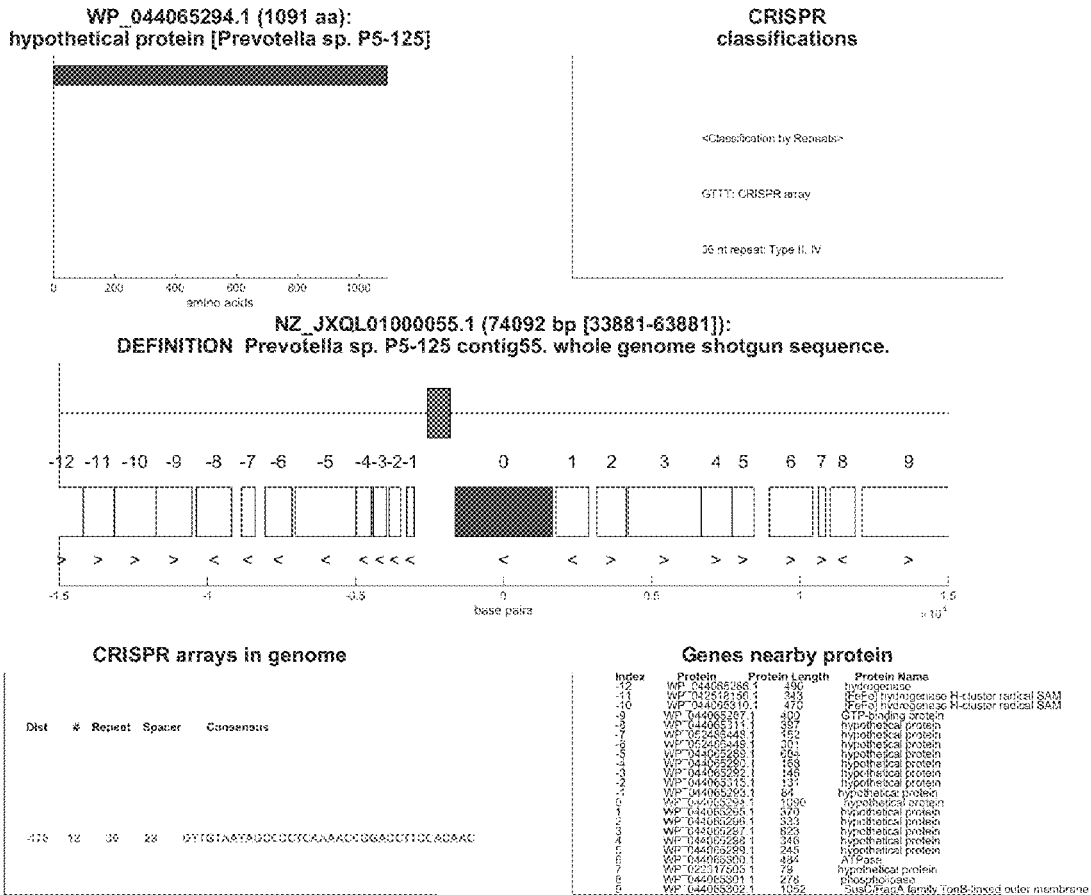
CRISPR arrays in genome

Dist	# Repeat	Spacer	Consensus	
700	2	56	10	GTTCGATTTACGGTCTATTGGAAGGCTTACACACAGC

Genes nearby protein

Index	Protein	Protein Length	Protein Name
1	WP_043895933.1	1176	hypothetical protein
2	WP_043895933.2	1176	hypothetical protein
3	WP_043895933.3	1176	hypothetical protein
4	WP_043895933.4	1176	hypothetical protein
5	WP_043895933.5	1176	hypothetical protein
6	WP_043895933.6	1176	hypothetical protein
7	WP_043895933.7	1176	hypothetical protein
8	WP_043895933.8	1176	hypothetical protein
9	WP_043895933.9	1176	hypothetical protein
10	WP_043895933.10	1176	hypothetical protein
11	WP_043895933.11	1176	hypothetical protein
12	WP_043895933.12	1176	hypothetical protein
13	WP_043895933.13	1176	hypothetical protein
14	WP_043895933.14	1176	hypothetical protein
15	WP_043895933.15	1176	hypothetical protein
16	WP_043895933.16	1176	hypothetical protein
17	WP_043895933.17	1176	hypothetical protein
18	WP_043895933.18	1176	hypothetical protein
19	WP_043895933.19	1176	hypothetical protein
20	WP_043895933.20	1176	hypothetical protein
21	WP_043895933.21	1176	hypothetical protein
22	WP_043895933.22	1176	hypothetical protein
23	WP_043895933.23	1176	hypothetical protein
24	WP_043895933.24	1176	hypothetical protein
25	WP_043895933.25	1176	hypothetical protein
26	WP_043895933.26	1176	hypothetical protein
27	WP_043895933.27	1176	hypothetical protein
28	WP_043895933.28	1176	hypothetical protein
29	WP_043895933.29	1176	hypothetical protein
30	WP_043895933.30	1176	hypothetical protein
31	WP_043895933.31	1176	hypothetical protein
32	WP_043895933.32	1176	hypothetical protein
33	WP_043895933.33	1176	hypothetical protein
34	WP_043895933.34	1176	hypothetical protein
35	WP_043895933.35	1176	hypothetical protein
36	WP_043895933.36	1176	hypothetical protein
37	WP_043895933.37	1176	hypothetical protein
38	WP_043895933.38	1176	hypothetical protein
39	WP_043895933.39	1176	hypothetical protein
40	WP_043895933.40	1176	hypothetical protein
41	WP_043895933.41	1176	hypothetical protein
42	WP_043895933.42	1176	hypothetical protein
43	WP_043895933.43	1176	hypothetical protein
44	WP_043895933.44	1176	hypothetical protein
45	WP_043895933.45	1176	hypothetical protein
46	WP_043895933.46	1176	hypothetical protein
47	WP_043895933.47	1176	hypothetical protein
48	WP_043895933.48	1176	hypothetical protein
49	WP_043895933.49	1176	hypothetical protein
50	WP_043895933.50	1176	hypothetical protein
51	WP_043895933.51	1176	hypothetical protein
52	WP_043895933.52	1176	hypothetical protein
53	WP_043895933.53	1176	hypothetical protein
54	WP_043895933.54	1176	hypothetical protein
55	WP_043895933.55	1176	hypothetical protein
56	WP_043895933.56	1176	hypothetical protein
57	WP_043895933.57	1176	hypothetical protein
58	WP_043895933.58	1176	hypothetical protein
59	WP_043895933.59	1176	hypothetical protein
60	WP_043895933.60	1176	hypothetical protein
61	WP_043895933.61	1176	hypothetical protein
62	WP_043895933.62	1176	hypothetical protein
63	WP_043895933.63	1176	hypothetical protein
64	WP_043895933.64	1176	hypothetical protein
65	WP_043895933.65	1176	hypothetical protein
66	WP_043895933.66	1176	hypothetical protein
67	WP_043895933.67	1176	hypothetical protein
68	WP_043895933.68	1176	hypothetical protein
69	WP_043895933.69	1176	hypothetical protein
70	WP_043895933.70	1176	hypothetical protein
71	WP_043895933.71	1176	hypothetical protein
72	WP_043895933.72	1176	hypothetical protein
73	WP_043895933.73	1176	hypothetical protein
74	WP_043895933.74	1176	hypothetical protein
75	WP_043895933.75	1176	hypothetical protein
76	WP_043895933.76	1176	hypothetical protein
77	WP_043895933.77	1176	hypothetical protein
78	WP_043895933.78	1176	hypothetical protein
79	WP_043895933.79	1176	hypothetical protein
80	WP_043895933.80	1176	hypothetical protein
81	WP_043895933.81	1176	hypothetical protein
82	WP_043895933.82	1176	hypothetical protein
83	WP_043895933.83	1176	hypothetical protein
84	WP_043895933.84	1176	hypothetical protein
85	WP_043895933.85	1176	hypothetical protein
86	WP_043895933.86	1176	hypothetical protein
87	WP_043895933.87	1176	hypothetical protein
88	WP_043895933.88	1176	hypothetical protein
89	WP_043895933.89	1176	hypothetical protein
90	WP_043895933.90	1176	hypothetical protein
91	WP_043895933.91	1176	hypothetical protein
92	WP_043895933.92	1176	hypothetical protein
93	WP_043895933.93	1176	hypothetical protein
94	WP_043895933.94	1176	hypothetical protein
95	WP_043895933.95	1176	hypothetical protein
96	WP_043895933.96	1176	hypothetical protein
97	WP_043895933.97	1176	hypothetical protein
98	WP_043895933.98	1176	hypothetical protein
99	WP_043895933.99	1176	hypothetical protein
100	WP_043895933.100	1176	hypothetical protein

FIG. 4BBBBB



CRISPR arrays in genome

Dist	# Repeat	Spacer	Consensus
-12	12	28	GTTGTAATYAGCCCTCCAAACCGGAGCCGTCCAGAAC

Genes nearby protein

Index	Protein	Protein Length	Protein Name
12	WP_044065285.1	436	hydrogenase
-11	WP_044065291.1	343	ProF1 hydrogenase H-cluster related SAM
-10	WP_044065293.1	473	ProF2 hydrogenase H-cluster related SAM
-9	WP_044065267.1	300	GTP-binding protein
-8	WP_044065314.1	287	hypothetical protein
-7	WP_044065343.1	182	hypothetical protein
-6	WP_044065345.1	201	hypothetical protein
-5	WP_044065298.1	266	hypothetical protein
-4	WP_044065282.1	158	hypothetical protein
-3	WP_044065284.1	146	hypothetical protein
-2	WP_044065315.1	137	hypothetical protein
-1	WP_044065293.1	61	hypothetical protein
0	WP_044065288.1	1396	hypothetical protein
1	WP_044065295.1	370	hypothetical protein
2	WP_044065296.1	333	hypothetical protein
3	WP_044065297.1	323	hypothetical protein
4	WP_044065298.1	346	hypothetical protein
5	WP_044065299.1	285	hypothetical protein
6	WP_044065303.1	434	ATPase
7	WP_02217505.1	73	hypothetical protein
8	WP_044065311.1	278	phosphatase
9	WP_044065302.1	1026	SusC/FraA family TonB-dependent outer membrane

FIG. 4CCCC

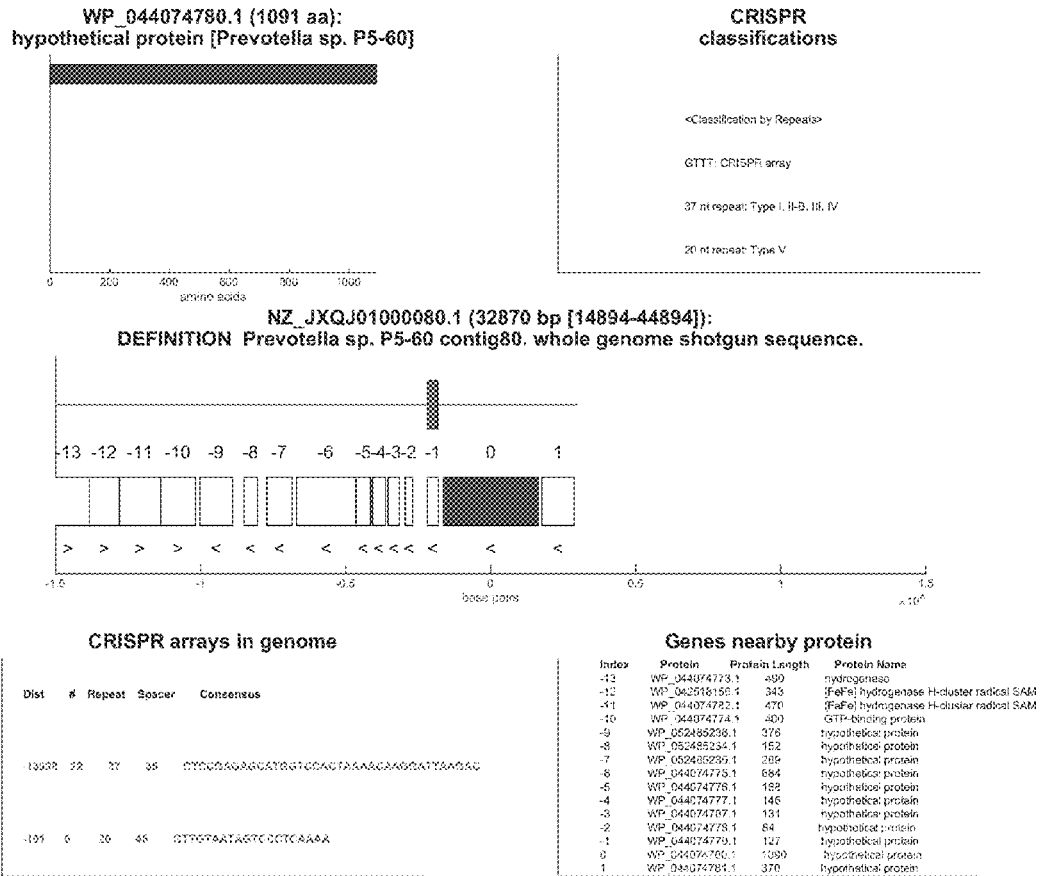


FIG. 4DDDDD

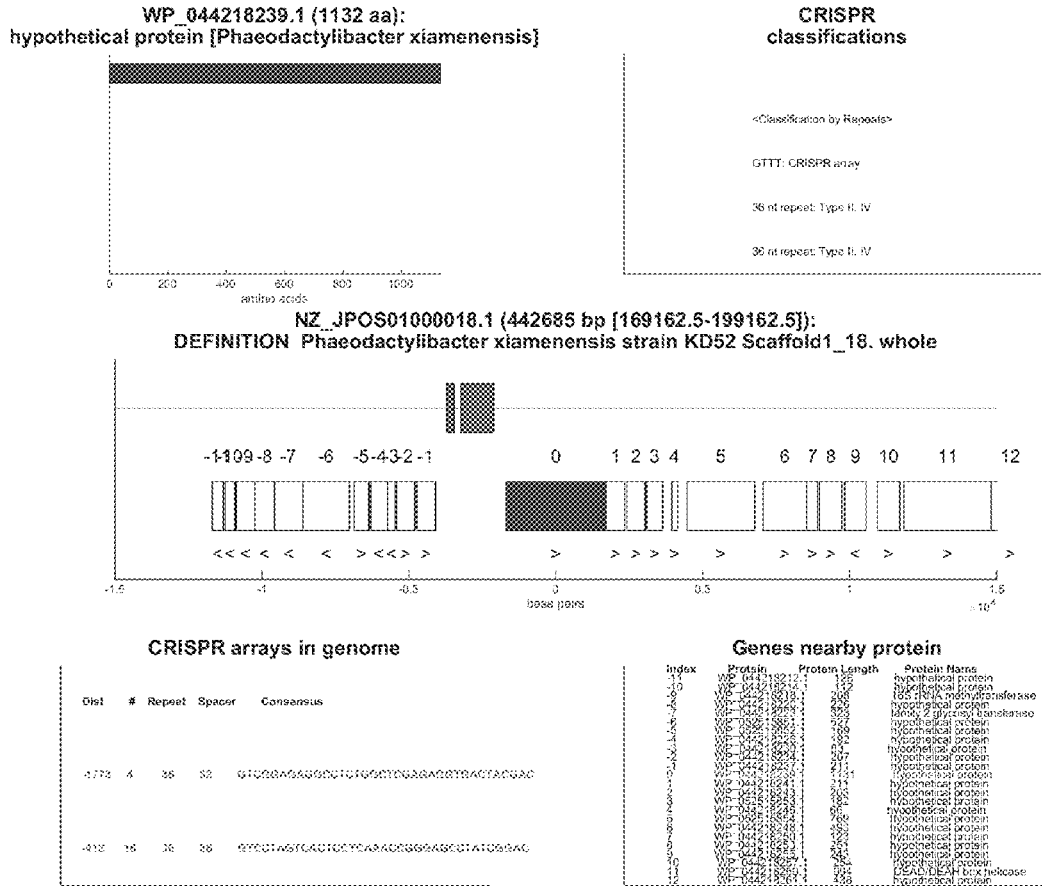


FIG. 4EEEEEE

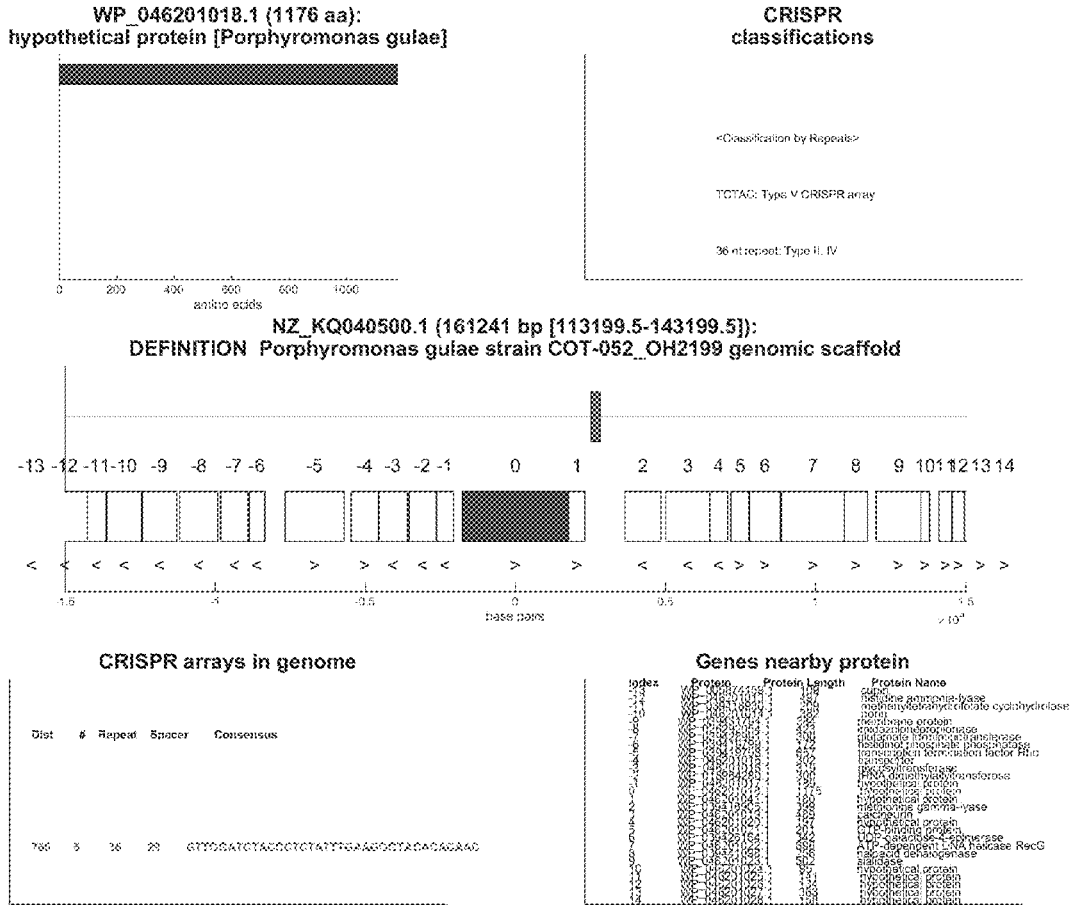


FIG. 4FFFFF

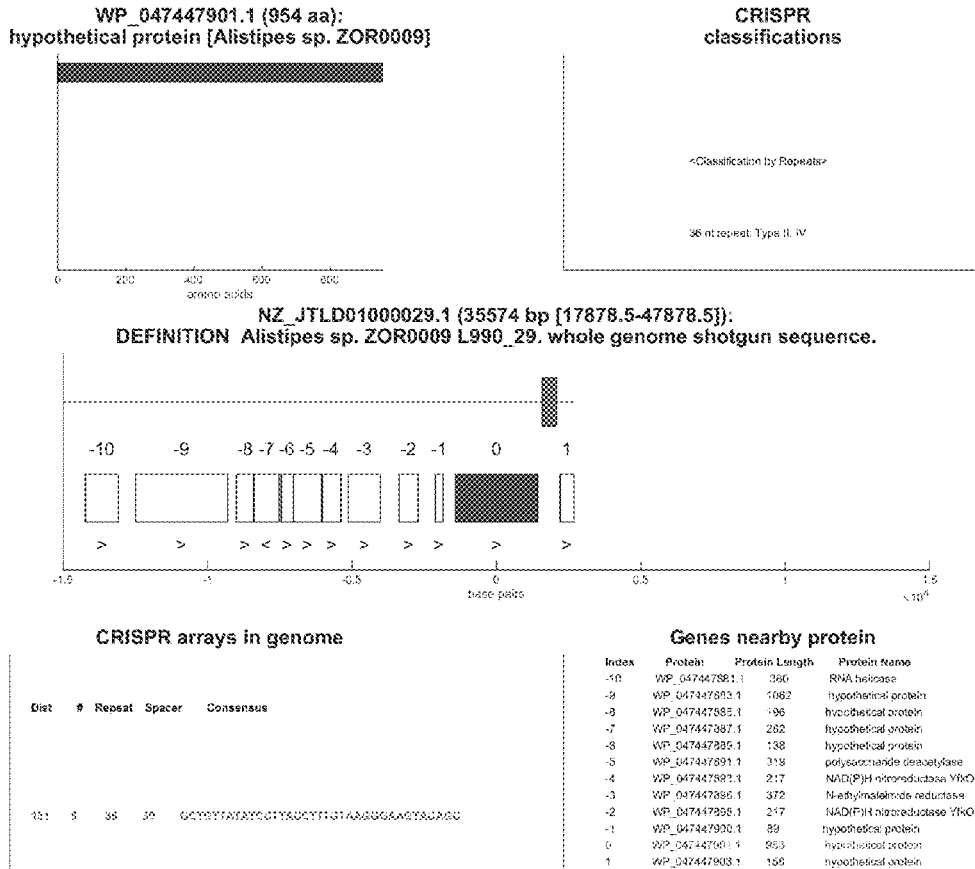
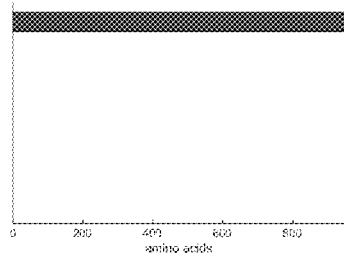


FIG. 4GGGGG

WP_049354263.1 (949 aa):
hypothetical protein [Riemerella anatipestifer]

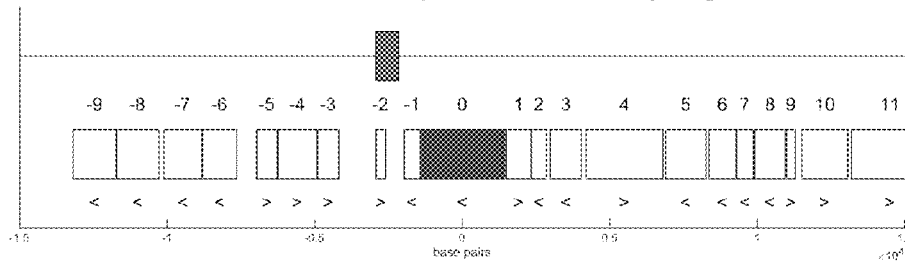


CRISPR
classifications

<Classification by Repeats>

38 nt repeat: Type II, IV

NZ_CP007504.1 (2160281 bp [1720166-1750166]):
DEFINITION Riemerella anatipestifer strain 153. complete genome.



CRISPR arrays in genome

Dist	#	Repeat	Spacer	Consensus
674177	14	47	28	GTTCGTGAATTCGTTTCGAATTTTACGATCTTTTGTGATAGTTTCGCAAC
736	12	30	35	GTTCGTGAATTCGTTTCGAATTCGATGAGGAGGATTCGCAAC

Genes nearby protein

Index	Protein	Protein Length	Protein Name
0	WP_049354263.1	949	RNA polymerase sigma-G4 factor
1	WP_049354264.1	482	serine/threonine phosphatase
2	WP_049354265.1	426	4-hydroxybutyrate CoA transferase
3	WP_049354266.1	287	hydroxyethylglutathione transferase
4	WP_049354267.1	228	hypothetical protein
5	WP_049354268.1	422	hypothetical protein
6	WP_049354269.1	299	hypothetical protein
7	WP_049354270.1	135	hypothetical protein
8	WP_049354271.1	178	hypothetical protein
9	WP_049354272.1	348	hypothetical protein
10	WP_049354273.1	205	hypothetical protein
11	WP_049354274.1	171	hypothetical protein
12	WP_049354275.1	148	hypothetical protein
13	WP_049354276.1	867	adenine deaminase
14	WP_049354277.1	455	phosphotransferase
15	WP_049354278.1	900	sugar kinase
16	WP_049354279.1	187	gliding mobility protein GldD
17	WP_049354280.1	325	ATG-specific adenine glycosylase
18	WP_049354281.1	97	DNA-binding protein
19	WP_049354282.1	618	hypothetical protein
20	WP_049354283.1	683	metalloprotease

FIG. 4HHHHH

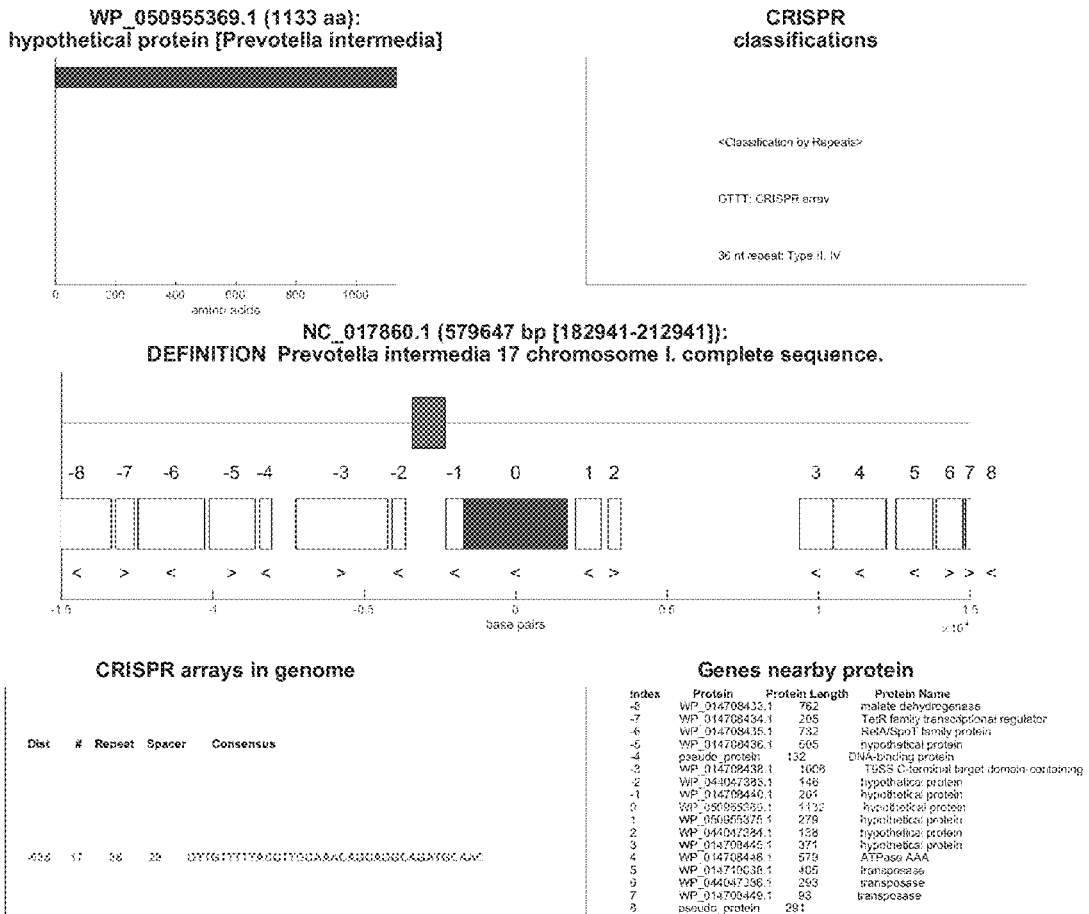


FIG. 4IIIII

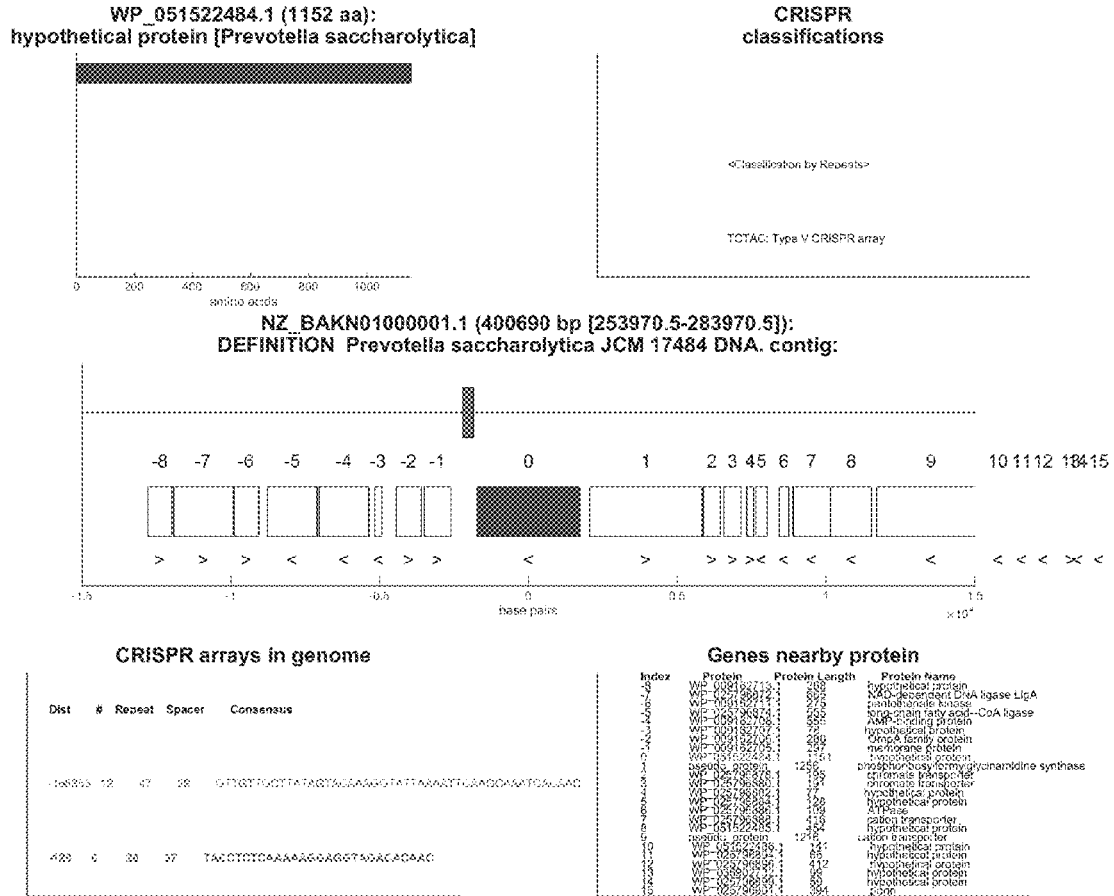


FIG. 4JJJJJ

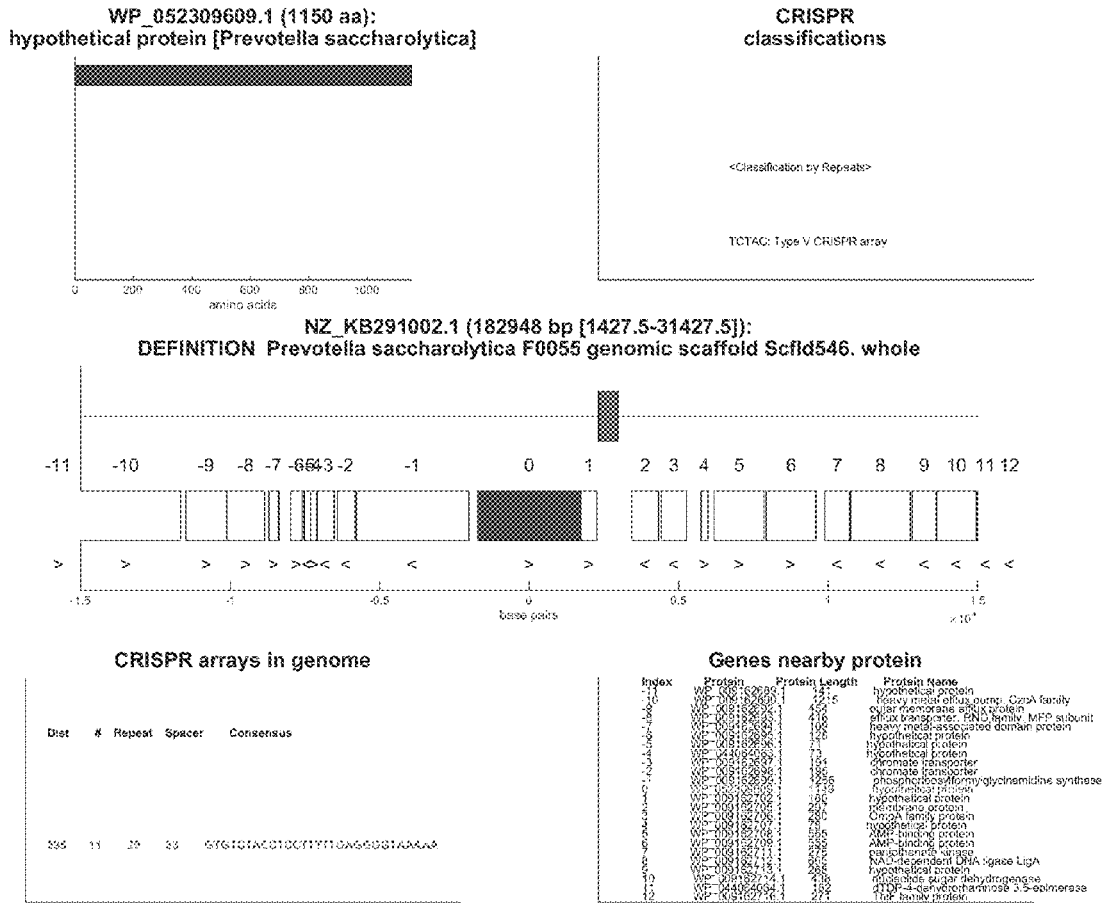


FIG. 4KKKKK

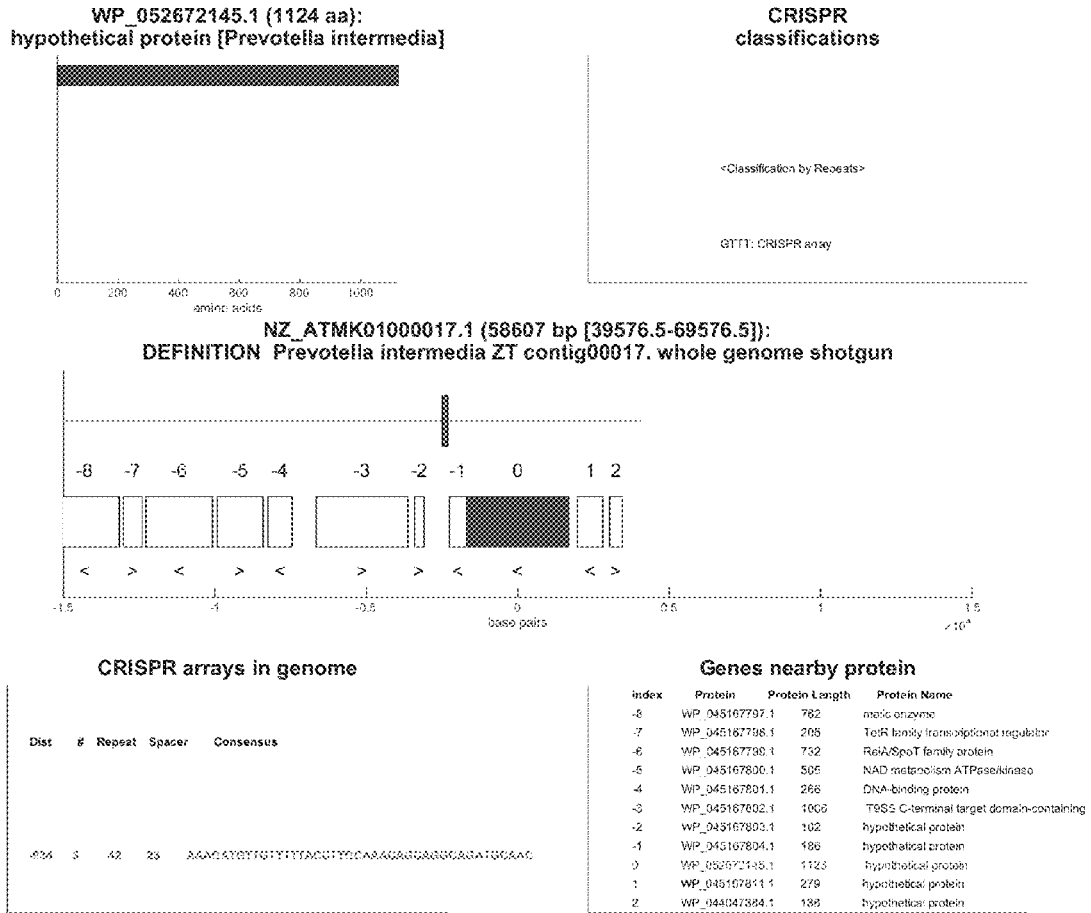
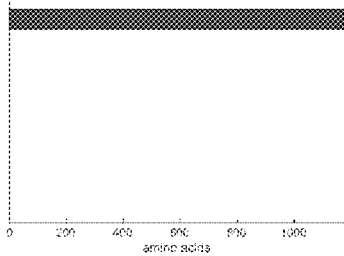


FIG. 4LLLLL

WP_052912312.1 (1176 aa):
hypothetical protein [Porphyromonas gingivalis]



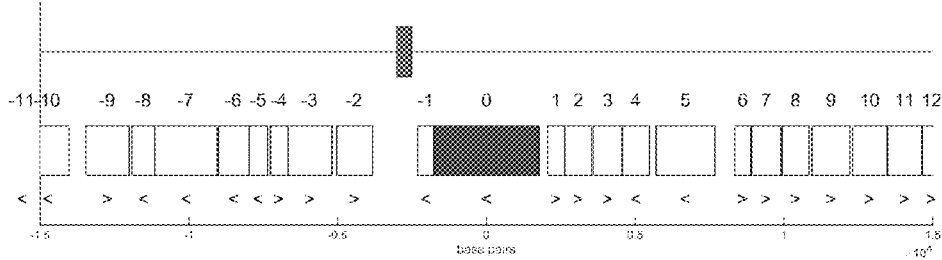
CRISPR
classifications

<Classification by Repeats>

TCTAC: Type V CRISPR array

36 nt repeat: Type II, IV

NZ_CP011995.1 (2367029 bp [1772994.5-1802994.5]):
DEFINITION Porphyromonas gingivalis strain A7436, complete genome.



CRISPR arrays in genome

Dist	#	Repeat	Spacer	Consensus
306982	0	37	38	CTTTAAATAGGCTTACGGACTCTGTATGTATGATCAG
301573	4	38	34	AGTGTAAATAGGCTTACGGACTCTGTATGTATGATGATGAT
699816	0	36	31	TTCGGATATGCTTAAATAGGCTTACGGACTCTGTATGATGATGATGAT
758	8	38	39	GTATGATGATGATGATGATGATGATGATGATGATGATGATGATGAT
325894	7	36	36	GTTTAAATGCTTATGATGATGATGATGATGATGATGATGATGATGAT

Genes nearby protein

Index	Protein	Protein Length	Protein Name
101	WP_052912312.1	1176	hypothetical protein
102	WP_052912313.1	38	glycosylase
103	WP_052912314.1	284	liponase
104	WP_052912315.1	229	lipoteichoic acid hydrolase
105	WP_052912316.1	232	ATP-dependent DNA helicase RecQ
106	WP_052912317.1	232	GDP dehydrogenase 4-nucleotidase
107	WP_052912318.1	217	GTP-binding protein
108	WP_052912319.1	417	hypothetical protein
109	WP_052912320.1	448	ribonuclease H1c phosphatase
110	WP_052912321.1	289	hypothetical protein
111	WP_052912322.1	289	hypothetical protein
112	WP_052912323.1	338	hypothetical protein
113	WP_052912324.1	270	RNA dependent DNA polymerase
114	WP_052912325.1	270	glyoxalase
115	WP_052912326.1	258	transposase
116	WP_052912327.1	472	transcription termination factor Rho
117	WP_052912328.1	472	adenosine 3' phosphatase phosphatase
118	WP_052912329.1	370	hypothetical protein
119	WP_052912330.1	370	glutamate formyltransferase
120	WP_052912331.1	366	thiazolone synthase
121	WP_052912332.1	366	melanocyte protein
122	WP_052912333.1	268	poly(ADP-ribose) polymerase 1
123	WP_052912334.1	268	ribitol 5-phosphate 4-epimerase

FIG. 4MMMMM

Protein Accession	Nucleotide Accession	Species	Strain
AFJ07523.1	CP003502.1	Prevotella intermedia 17	17
EOA10535.1	APMB01000175.1	Porphyromonas gingivalis JCVI SC001	JCVI SC001
WP_002664492.1	NZ_JH932293.1	Bergeyella zoohelcum ATCC 43767	ATCC 43767
EKB54193.1	AGYA01000037.1	Bergeyella zoohelcum ATCC 43767	ATCC 43767
WP_004343581.1	NZ_GG739967.1	Prevotella buccae D17	D17
EFC74556.1	GG739967.1	Prevotella buccae D17	D17
WP_004343973.1	NZ_GL586311.1	Prevotella buccae ATCC 33574	ATCC 33574
EFU31981.1	AEPD01000005.1	Prevotella buccae ATCC 33574	ATCC 33574
WP_004913755.1	NC_014738.1	Riemerella anatipestifer ATCC 11845 = DSM 15868	DSM 15868
AKP69887.1	CP007503.1	Riemerella anatipestifer	17
ADQ82708.1	CP002346.1	Riemerella anatipestifer ATCC 11845 = DSM 15868	DSM 15868
AFD56718.1	CP003388.1	Riemerella anatipestifer ATCC 11845 = DSM 15868	ATCC 11845
ADZ11800.1	CP002562.1	Riemerella anatipestifer RA-GD	RA-GD
EFT35596.1	AENH01000026.1	Riemerella anatipestifer RA-YM	RA-YM
AKQ40303.1	CP007204.1	Riemerella anatipestifer Yb2	Yb2
WP_005873511.1	NC_002950.2	Porphyromonas gingivalis W83	W83
AKV64040.1	CP011995.1	Porphyromonas gingivalis	A7436
EIW94777.1	AJZS01000011.1	Porphyromonas gingivalis W50	W50
AAQ66265.1	AE015924.1	Porphyromonas gingivalis W83	W83
EIW93190.1	AJZS01000051.1	Porphyromonas gingivalis W50	W50
WP_006044833.1	NZ_GL982513.1	Prevotella pallens ATCC 700821	ATCC 700821
EGQ18444.1	AFPY01000052.1	Prevotella pallens ATCC 700821	ATCC 700821
WP_006259957.1	NZ_KE161016.1	Myroides odoratimimus CCUG 10230	CCUG 10230
EHO06562.1	AGEC02000017.1	Myroides odoratimimus CCUG 10230	CCUG 10230
WP_007412163.1	NZ_ALIQ01000043.1	Prevotella sp. MSX73	MSX73
EJP27887.1	ALIQ01000043.1	Prevotella sp. MSX73	MSX73
WP_012458151.1	NC_010729.1	Porphyromonas gingivalis ATCC 33277	
AIJ35608.1	CP007756.1	Porphyromonas gingivalis	HG66
BAG33782.1	AP009380.1	Porphyromonas gingivalis ATCC 33277	ATCC 33277
WP_012458414	NC_010729.1	Porphyromonas gingivalis ATCC 33277	

FIG. 5A

Protein Accession	Nucleotide Accession	Species	Strain
.1			
BAG34142.1	AP009380.1	Porphyromonas gingivalis ATCC 33277	ATCC 33277
WP_013816155.1	NC_015571.1	Porphyromonas gingivalis TDC60	
BAK25606.1	AP012203.1	Porphyromonas gingivalis TDC60	TDC60
WP_013997271.1	NC_015846.1	Capnocytophaga canimorsus Cc5	Cc5
AEK23281.1	CP002113.1	Capnocytophaga canimorsus Cc5	Cc5
WP_014165541.1	NC_016510.2	Flavobacterium columnare ATCC 49512	ATCC 49512
AEW86266.1	CP003222.2	Flavobacterium columnare ATCC 49512	ATCC 49512
WP_015024765.1	NC_018721.1	Psychroflexus torquis ATCC 700755	ATCC 700755
AFU69194.1	CP003879.1	Psychroflexus torquis ATCC 700755	ATCC 700755
WP_015345619.1	NC_020125.1	Riemerella anatipestifer RA-CH-2	RA-CH-2
AKP71852.1	CP007504.1	Riemerella anatipestifer	153
AGC41342.1	CP004020.1	Riemerella anatipestifer RA-CH-2	RA-CH-2
WP_015345620.1	NC_020125.1	Riemerella anatipestifer RA-CH-2	RA-CH-2
AGC41344.1	CP004020.1	Riemerella anatipestifer RA-CH-2	RA-CH-2
ERJ64231.1	AWUW01000137.1	Porphyromonas gingivalis F0570	F0570
ERJ66054.1	AWUU01000135.1	Porphyromonas gingivalis F0568	F0568
ERJ81654.1	AWVDC01000226.1	Porphyromonas gingivalis F0566	F0566
WP_021584635.1	NZ_AWET01000045.1	Prevotella pleuritidis F0068	F0068
ERJ98772.1	AWET01000045.1	Prevotella pleuritidis F0068	F0068
WP_021663197.1	NZ_KI258981.1	Porphyromonas gingivalis F0568	F0568
ERJ65637.1	AWUU01000145.1	Porphyromonas gingivalis F0568	F0568
WP_021665475.1	NZ_KI259158.1	Porphyromonas gingivalis F0570	F0570
ERJ66432.1	AWUW01000078.1	Porphyromonas gingivalis F0570	F0570
WP_021677657.1	NZ_KI259960.1	Porphyromonas gingivalis F0185	F0185
WP_021680012.1	NZ_KI260263.1	Porphyromonas gingivalis W4087	W4087
ERJ87335.1	AWVE01000130.1	Porphyromonas gingivalis W4087	W4087
WP_023846767.1	NZ_KI629875.1	Porphyromonas gingivalis SJD2	SJD2
ETA27404.1	ASYL01000050.1	Porphyromonas gingivalis SJD2	SJD2

FIG. 5B

Protein Accession	Nucleotide Accession	Species	Strain
WP_025000926.1	NZ_BAKF01000019.1	Prevotella aurantiaca JCM 15754	JCM 15754
WP_034542281.1	NZ_BAIU01000001.1	Bacteroides pyogenes JCM 10003	
GAE20957.1	BAIU01000001.1	Bacteroides pyogenes JCM 10003	JCM 10003
WP_036860899.1	NZ_JAEZ01000017.1	Prevotella intermedia ATCC 25611 = DSM 20706	ATCC 25611
WP_036884929.1	NZ_BAJY01000004.1	Prevotella falsenii DSM 22864 = JCM 15124	JCM 15124
WP_036929175.1	NZ_JHUW01000010.1	Prevotella sp. MA2016	MA2016
WP_036931485.1	NZ_BAJN01000005.1	Prevotella pleuritidis JCM 14110	JCM 14110
WP_039417390.1	NZ_CP007756.1	Porphyromonas gingivalis	HG66
AU36392.1	CP007756.1	Porphyromonas gingivalis	HG66
WP_039418912.1	NZ_JRAT01000012.1	Porphyromonas gulae	COT-052 OH3856
KGL48767.1	JRAT01000012.1	Porphyromonas gulae	COT-052 OH3856
WP_039419792.1	NZ_JRAI01000002.1	Porphyromonas gulae	COT-052 OH1451
KGN88274.1	JRAI01000002.1	Porphyromonas gulae	COT-052 OH1451
WP_039426176.1	NZ_JRAK01000129.1	Porphyromonas gulae	COT-052 OH3439
KGN85385.1	JRAK01000129.1	Porphyromonas gulae	COT-052 OH3439
WP_039428968.1	NZ_JQZY01000014.1	Porphyromonas sp. COT-052 OH4946	COT-052 OH4946
KGL55352.1	JQZY01000014.1	Porphyromonas sp. COT-052 OH4946	COT-052 OH4946
WP_039431778.1	NZ_JRAJ01000010.1	Porphyromonas gulae	COT-052 OH2179
KGN69120.1	JRAJ01000010.1	Porphyromonas gulae	COT-052 OH2179
WP_039434803.1	NZ_JRAL01000022.1	Porphyromonas gulae	COT-052 OH4119
KGN74647.1	JRAL01000022.1	Porphyromonas gulae	COT-052 OH4119
KGO05347.1	JRAF01000001.1	Porphyromonas gulae	COT-052 OH3498
WP_039442171.1	NZ_JRFD01000046.1	Porphyromonas gulae	COT-052 OH2857

FIG. 5C

Protein Accession	Nucleotide Accession	Species	Strain
KGN87312.1	JRFDD1000046.1	Porphyromonas gulae	COT-052 OH2857
WP_039445055.1	NZ_JRAQ01000019.1	Porphyromonas gulae	COT-052 OH3471
KGN76020.1	JRAQ01000019.1	Porphyromonas gulae	COT-052 OH3471
WP_041758328.1	NC_018721.1	Psychroflexus torquis ATCC 700755	ATCC 700755
WP_041986150.1	NZ_CDOK0100012.2.1	Capnocytophaga canimorsus	Cc11
CEN50095.1	CDOK01000122.1	Capnocytophaga canimorsus	Cc11
WP_041989581.1	NZ_CDOD0100000.2.1	Capnocytophaga cynodegmi	Ccyn28
CEN32480.1	CDOD01000002.1	Capnocytophaga cynodegmi	Ccyn28
WP_042518169.1	NZ_JXQK01000043.1	Prevotella sp. P5-119	P5-119
KIP63359.1	JXQK01000043.1	Prevotella sp. P5-119	P5-119
WP_043894148.1	NZ_KI259080.1	Porphyromonas gingivalis F0569	F0569
WP_043894372.1	NZ_KI259218.1	Porphyromonas gingivalis F0570	F0570
WP_043894537.1	NZ_KI260229.1	Porphyromonas gingivalis W4087	W4087
WP_043894751.1	NZ_KI259867.1	Porphyromonas gingivalis F0185	F0185
WP_043895240.1	NZ_KI260014.1	Porphyromonas gingivalis F0566	F0566
WP_043895933.1	NZ_KI258974.1	Porphyromonas gingivalis F0568	F0568
WP_044065294.1	NZ_JXQL01000055.1	Prevotella sp. P5-125	P5-125
KIP62088.1	JXQL01000055.1	Prevotella sp. P5-125	P5-125
WP_044074780.1	NZ_JXQJ01000080.1	Prevotella sp. P5-60	P5-60
KIP58950.1	JXQJ01000080.1	Prevotella sp. P5-60	P5-60
WP_044218239.1	NZ_JPOS01000018.1	Phaeodactylibacter xiamenensis	KD52
KGE88582.1	JPOS01000018.1	Phaeodactylibacter xiamenensis	KD52
KJJ86756.1	ATMK01000017.1	Prevotella intermedia ZT	ZT
WP_046201018.1	NZ_KQ040500.1	Porphyromonas gulae	COT-052_OH2199
KKC50278.1	JRAE01000093.1	Porphyromonas gulae	COT-052_OH2199

FIG. 5D

Protein Accession	Nucleotide Accession	Species	Strain
WP_047447901.1	NZ_JTLD01000029.1	Alistipes sp. ZOR0009	ZOR0009
WP_049354263.1	NZ_CP007504.1	Riemerella anatipestifer	
AKP71851.1	CP007504.1	Riemerella anatipestifer	153
WP_050955369.1	NC_017860.1	Prevotella intermedia 17	17
BAR96998.1	AP014926.1	Prevotella intermedia	17-2
WP_051322523.1	NZ_KE993153.1	Bacteroides pyogenes F0041	F0041
WP_051522484.1	NZ_BAKN01000001.1	Prevotella saccharolytica JCM 17484	JCM 17484
WP_052309609.1	NZ_KB291002.1	Prevotella saccharolytica F0055	F0055
WP_052672145.1	NZ_ATMK0100001.7.1	Prevotella intermedia ZT	ZT
WP_052912312.1	NZ_CP011995.1	Porphyromonas gingivalis	A7436
AKV64762.1	CP011995.1	Porphyromonas gingivalis	A7436

FIG. 5E

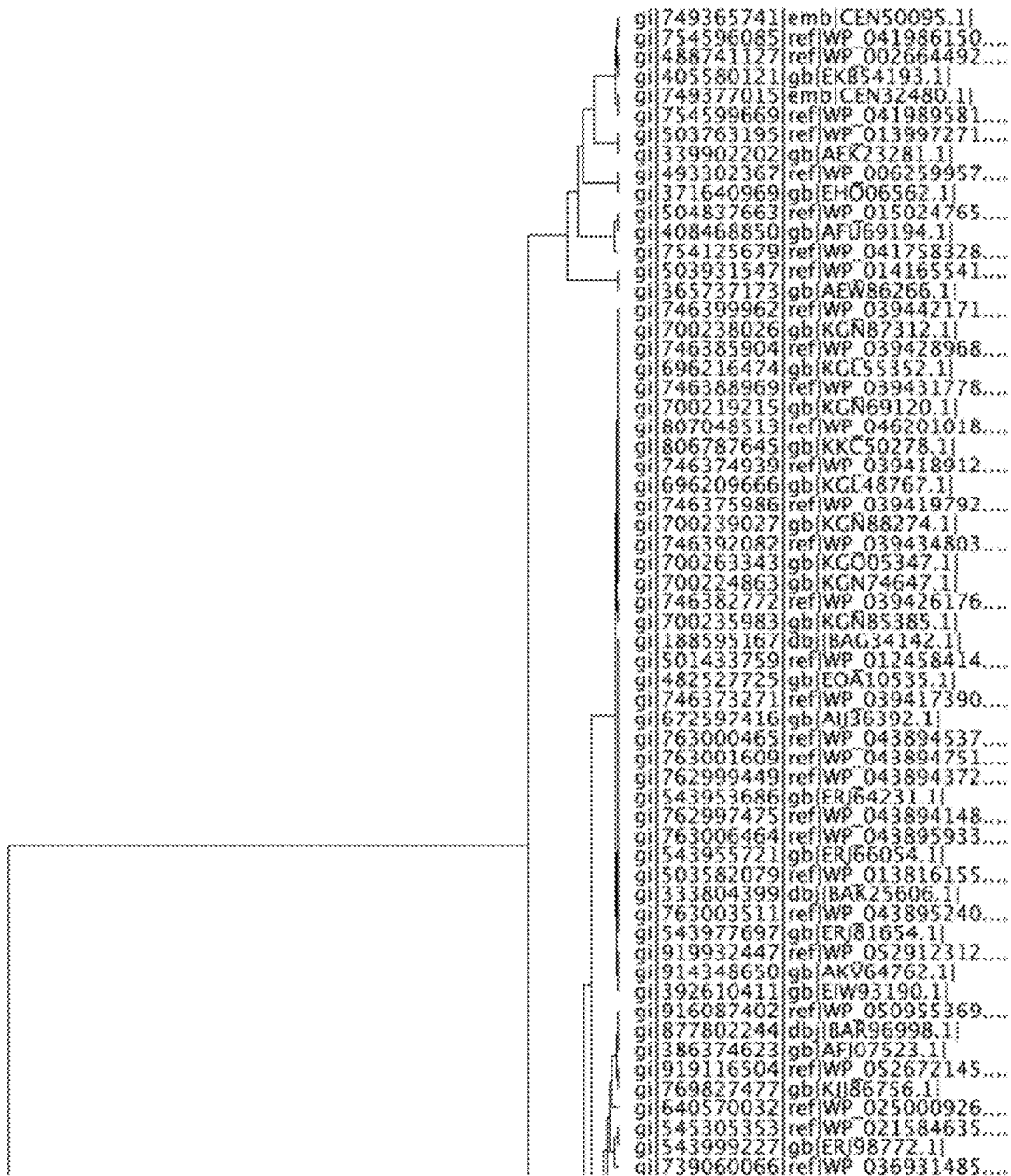
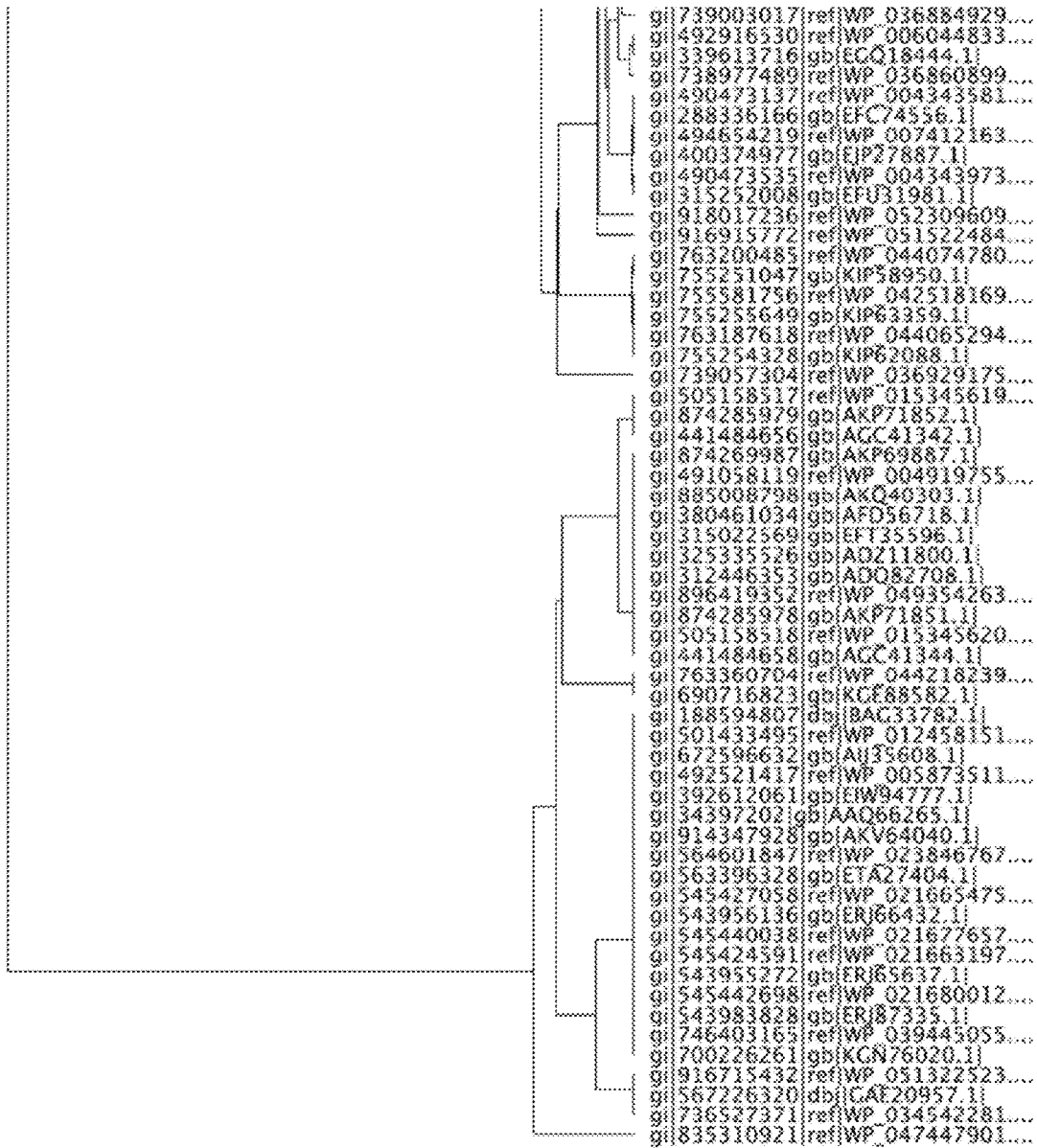


FIG. 6A



2.0

FIG. 6B

#	Prot	Nuc	Species	Strain
1	WP_002664492.1	NZ_JH932293.1	Bergeyella zoohelcum ATCC 43767	ATCC 43767
2	WP_036860899.1	NZ_JAEZ01000017.1	Prevotella intermedia ATCC 25611 = DSM 20706	ATCC 25611
3	WP_004343973.1	NZ_GLS86311.1	Prevotella buccae ATCC 33574	ATCC 33574
4	WP_012458151.1	NC_010729.1	Porphyromonas gingivalis ATCC 33277	ATCC 33277
5	WP_034542281.1	NZ_BAIU01000001.1	Bacteroides pyogenes JCM 10003	
6	WP_047447901.1	NZ_JTLD01000029.1	Alistipes sp. ZOR0009	ZOR0009
7	WP_036929175.1	NZ_JHUW01000010.1	Prevotella sp. MA2016	MA2016
8	WP_004919755.1	NC_014738.1	Riemerella anatipestifer ATCC 11845 = DSM 15868	DSM 15868
9	WP_025000926.1	NZ_BAKF01000019.1	Prevotella aurantiaca JCM 15754	JCM 15754
10	WP_051522484.1	NZ_BAKN01000001.1	Prevotella saccharolytica JCM 17484	JCM 17484
11	WP_006259957.1	NZ_KE161016.1	Myroides odoratimimus CCUG 10230	CCUG 10230
12	WP_014165541.1	NC_016510.2	Flavobacterium columnare ATCC 49512	ATCC 49512
13	WP_051322523.1	NZ_KE993153.1	Bacteroides pyogenes	F0041
14	WP_013997271.1	NC_015846.1	Capnocytophaga canimorsus Cc5	Cc5
15	WP_041989581.1	NZ_CDOD01000002.1	Capnocytophaga cynodegmi	Ccyn2B
16	WP_047431796.1	NZ_KN549099.1	Chryseobacterium sp.	YR477
17	WP_044218239.1	NZ_JPOS01000018.1	Phaendactylibacter xiamenensis	KD52
18	WP_043894148.1	NZ_KI259080.1	Porphyromonas gingivalis F0569	F0569
19	WP_039443055.1	NZ_JRAQ01000019.1	Porphyromonas gulae	COT-052 OH3471
20	WP_039434803.1	NZ_JRAL01000022.1	Porphyromonas gulae	COT-052 OH4119
21	WP_039419792.1	NZ_JRAI01000002.1	Porphyromonas gulae	COT-052 OH1451
22	WP_036884929.1	NZ_BAIY01000004.1	Prevotella falsenii DSM 22864 = JCM 15124	JCM 15124
23	WP_050955369.1	NC_017860.1	Prevotella intermedia 17	17
24	WP_052672145.1	NZ_ATMK01000017.1	Prevotella intermedia ZT	ZT
25	WP_006044833.1	NZ_GL982513.1	Prevotella pallens ATCC 700821	ATCC 700821
26	WP_021584635.1	NZ_AWET01000045.1	Prevotella pleuritidis F0068	F0068
27	WP_036931485.1	NZ_BAIND1000005.1	Prevotella pleuritidis JCM 14110	JCM 14110
28	WP_052309609.1	NZ_KB291002.1	Prevotella saccharolytica F0055	F0055
29	WP_042518169.1	NZ_JXQK01000043.1	Prevotella sp. P5-119	P5-119
30	WP_015024765.1	NC_018721.1	Psychroflexus torquis ATCC 700755	ATCC 700755

FIG. 7

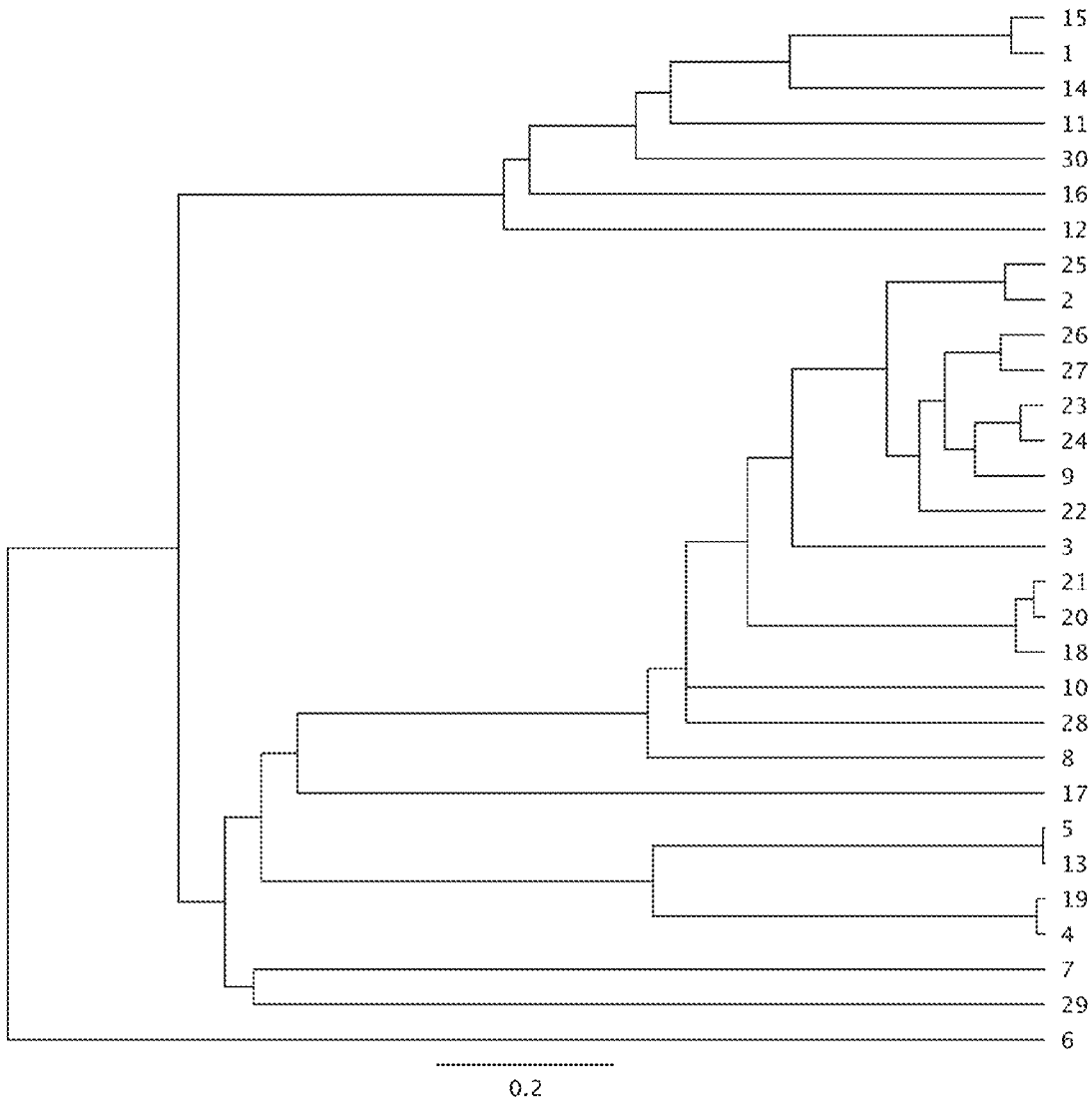


FIG. 8

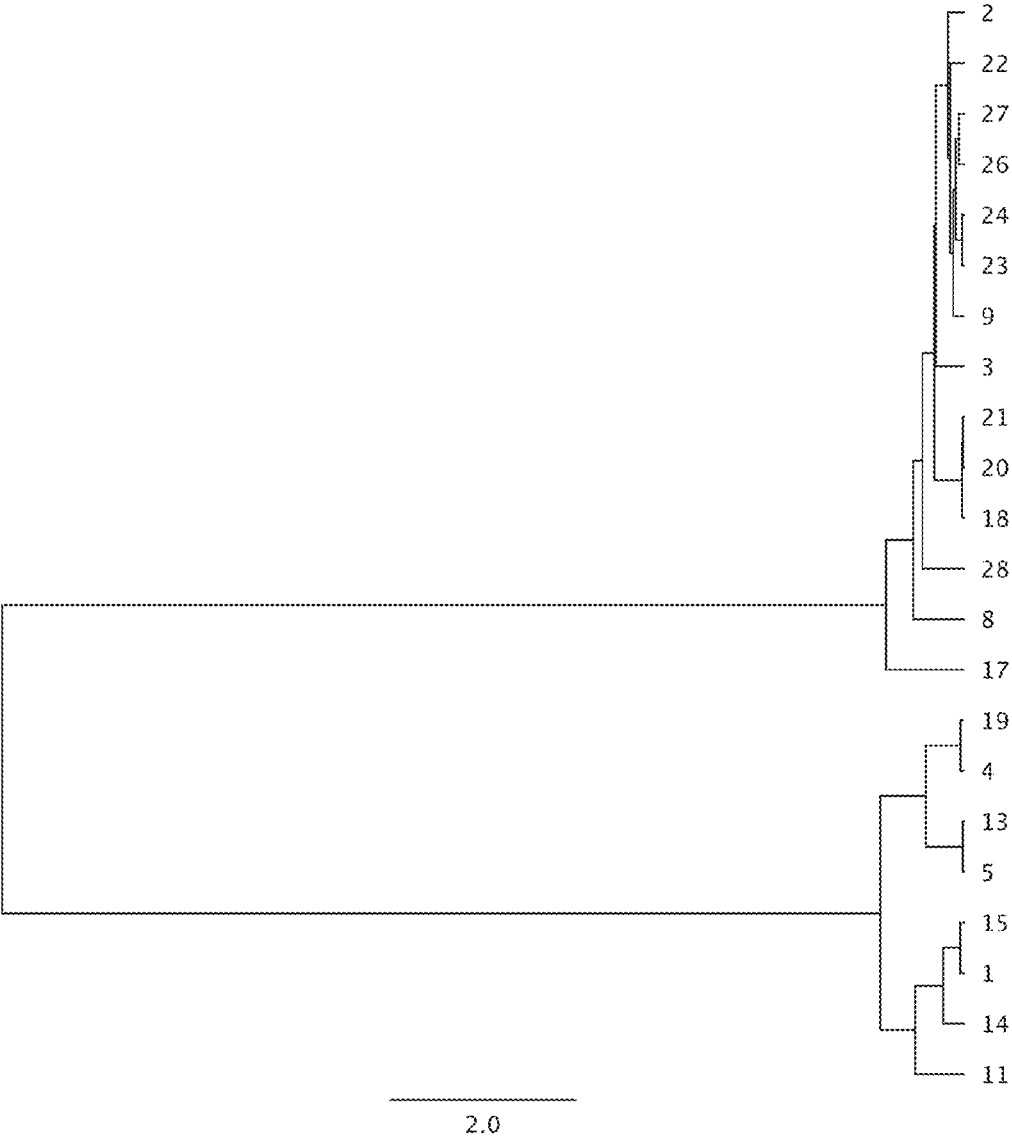


FIG. 9

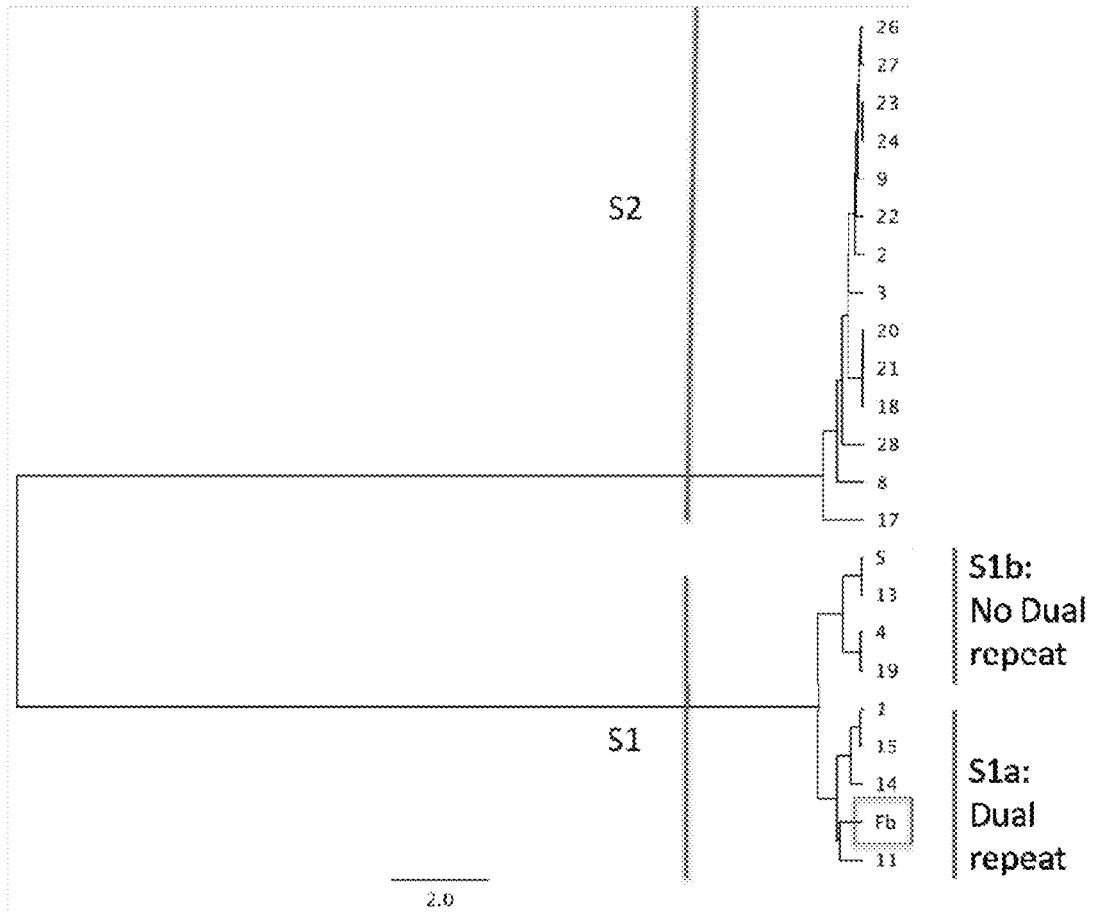
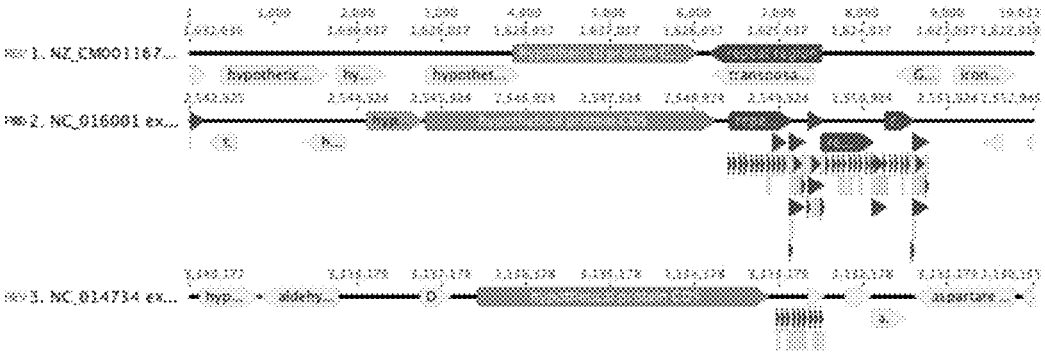


FIG. 10

FIG. 11



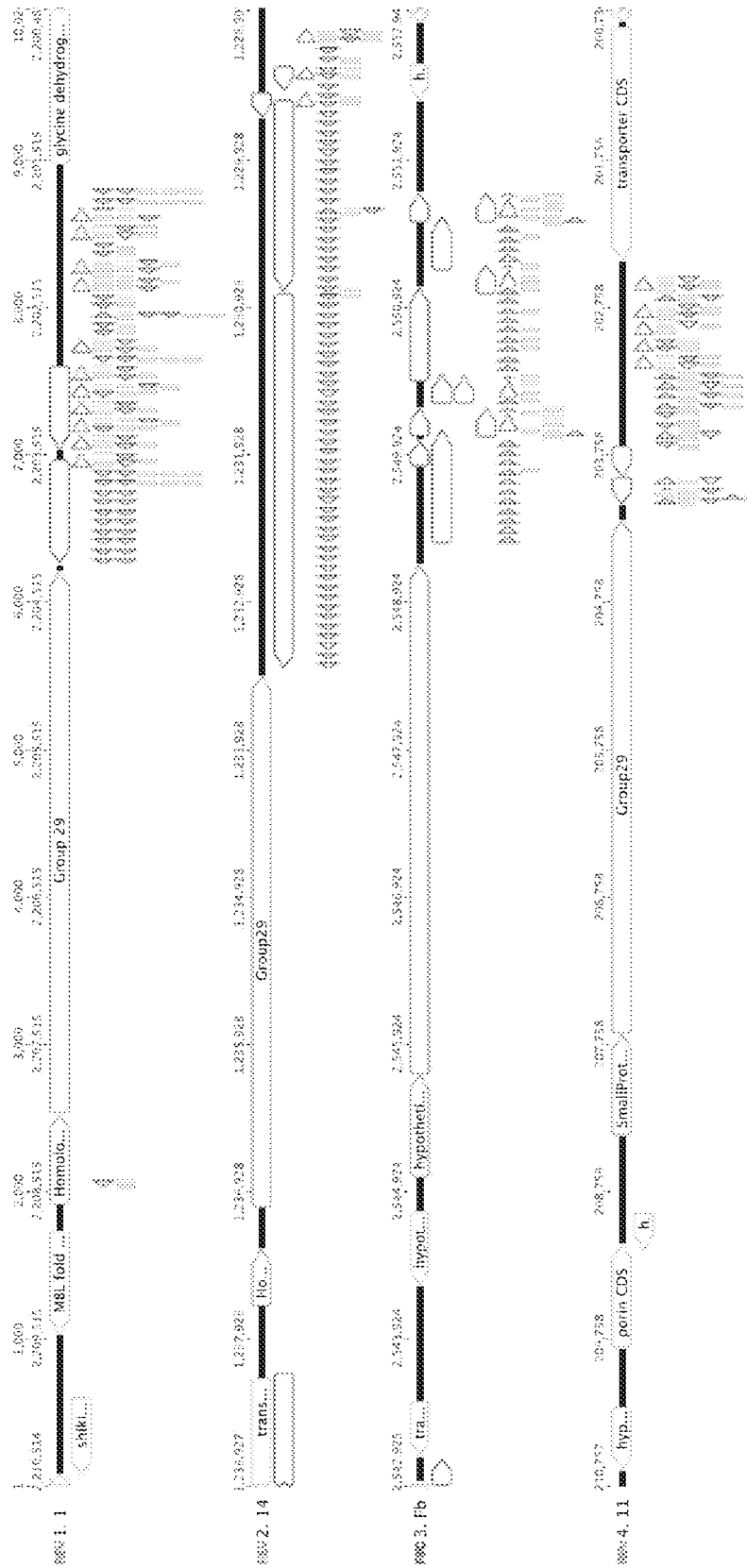


FIG. 12

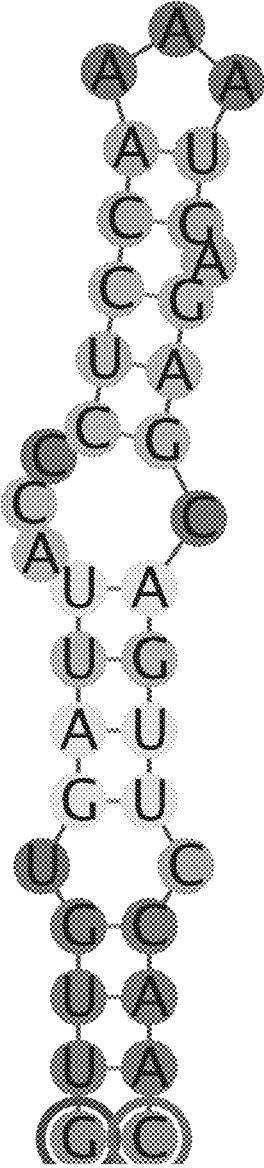


FIG. 13A

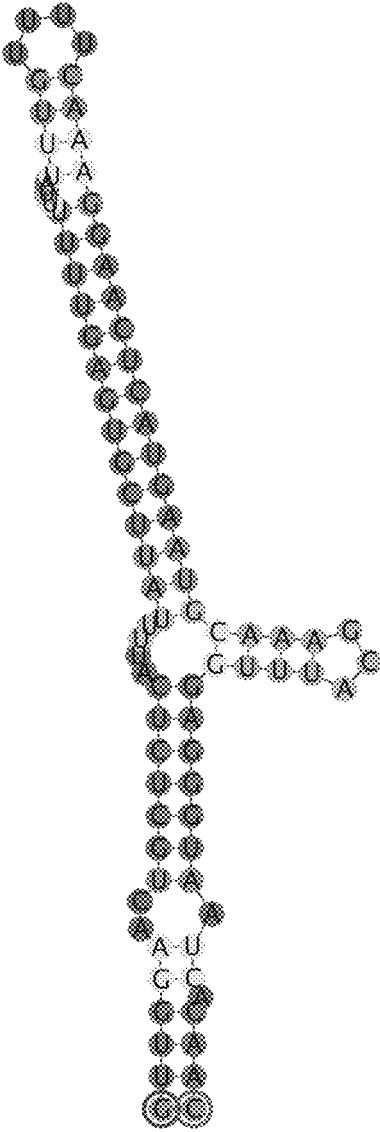


FIG. 13B

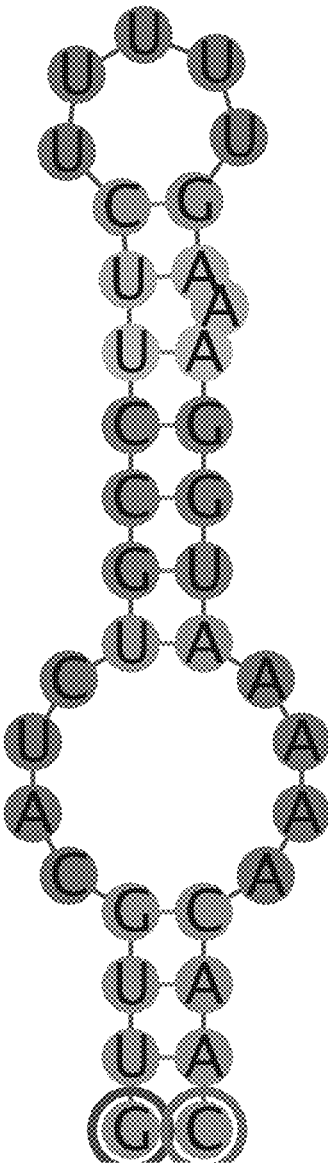


FIG. 13D

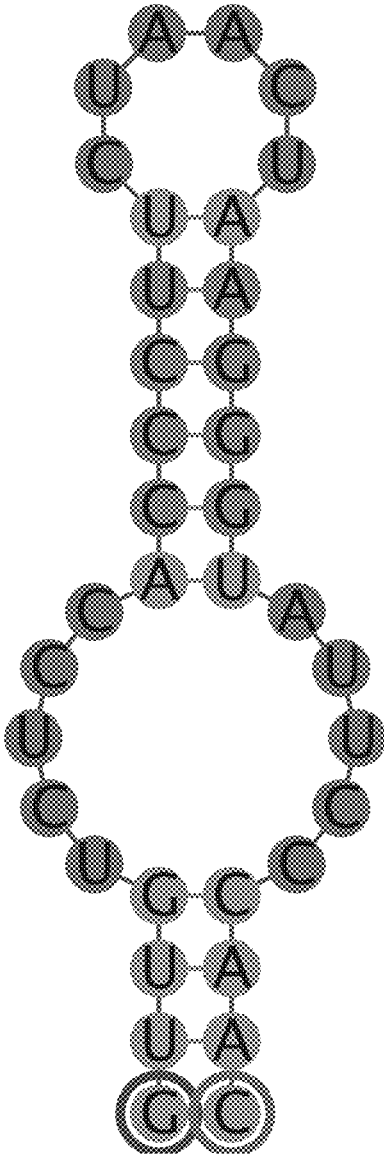


FIG. 13E

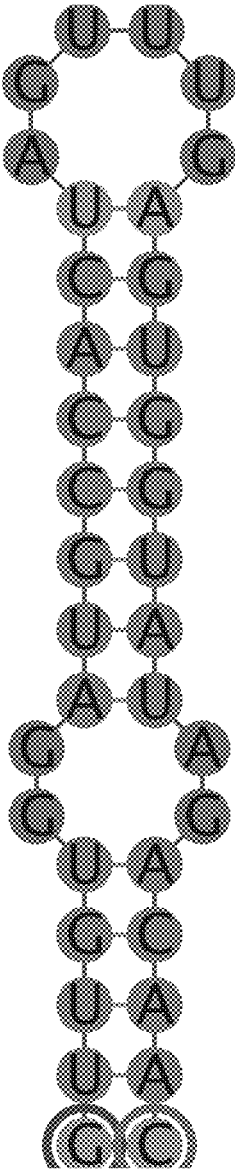


FIG. 13F

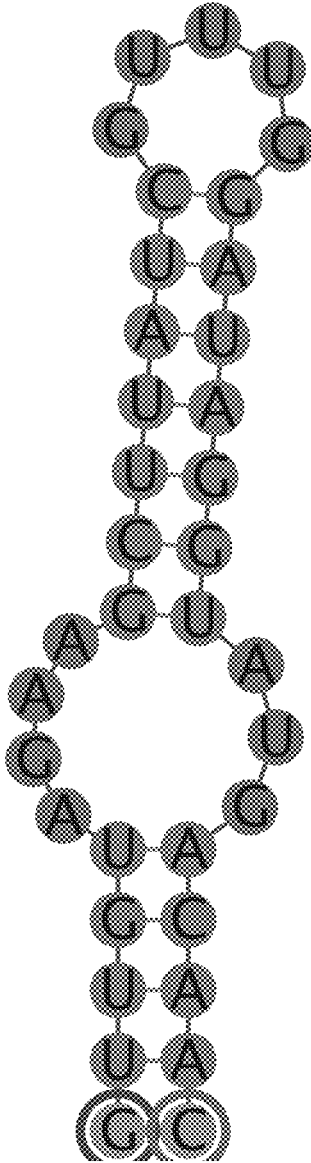


FIG. 13H

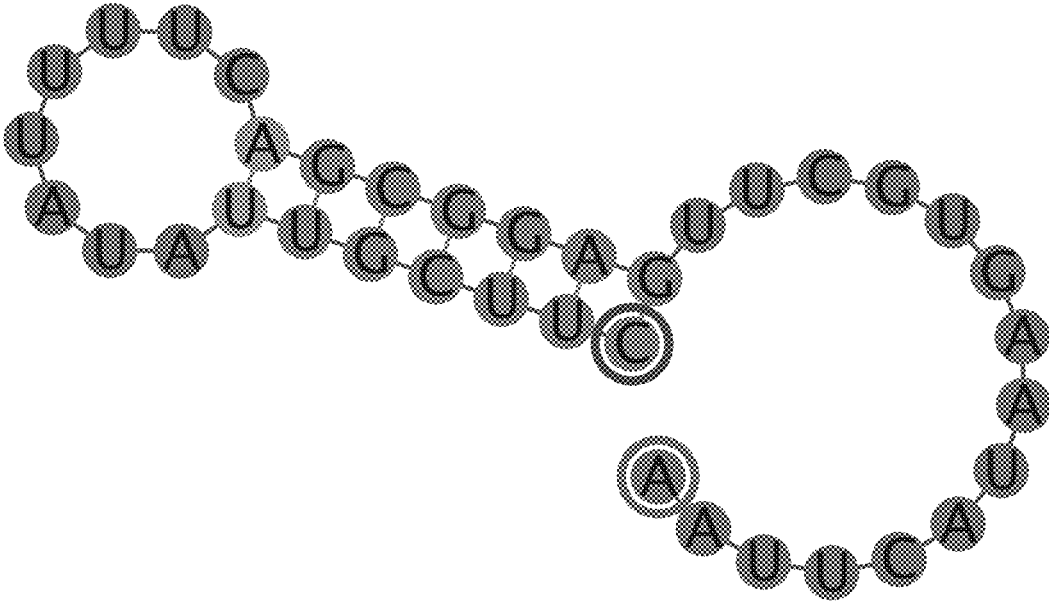


FIG. 13I

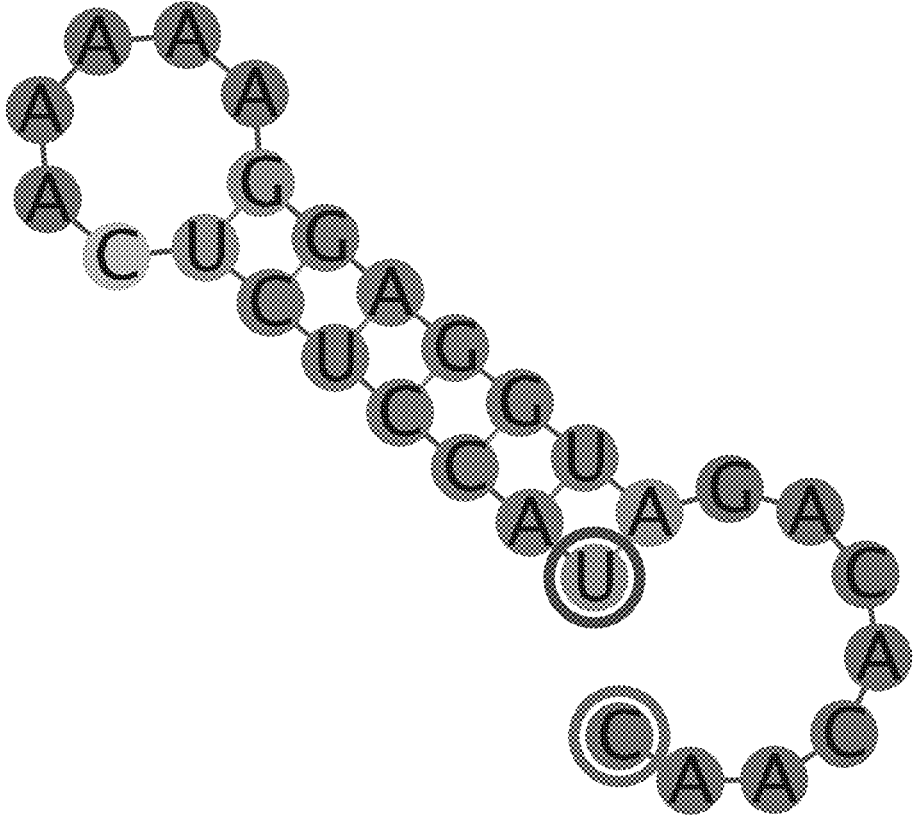


FIG. 13K

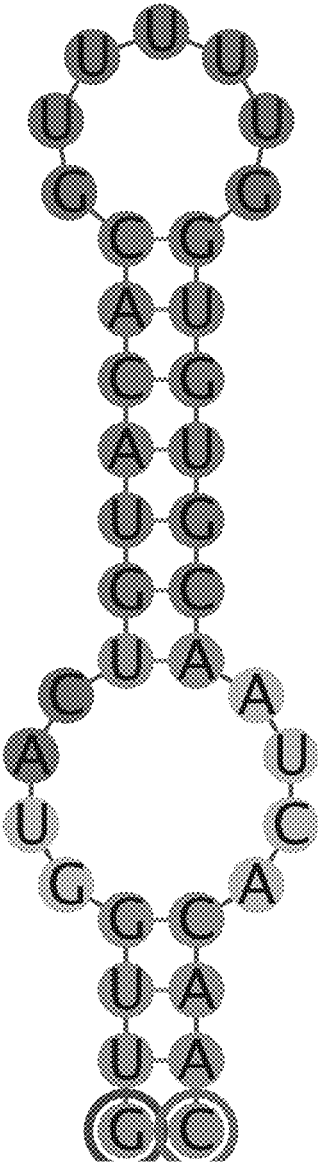


FIG. 13L

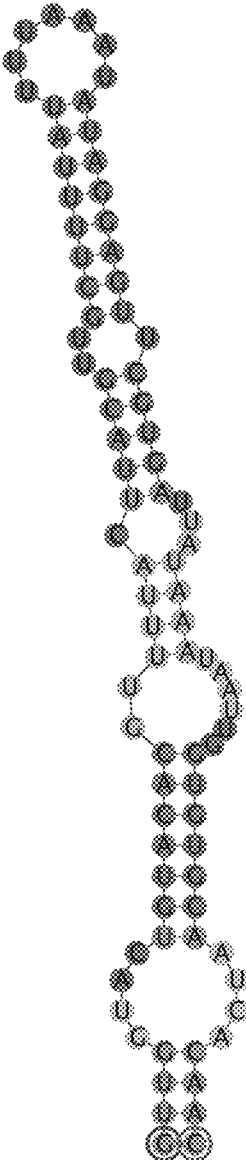


FIG. 13M

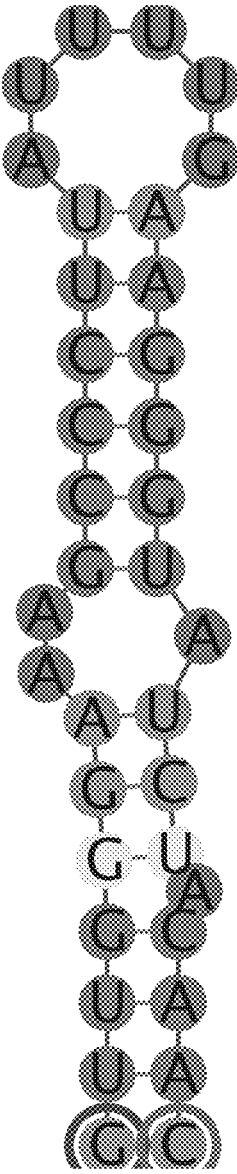


FIG. 13N

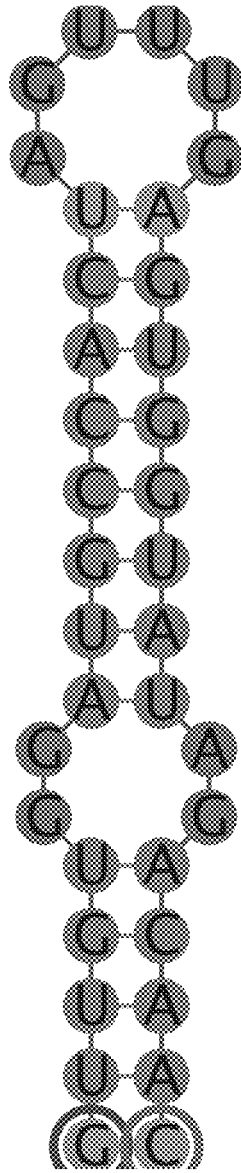


FIG. 130

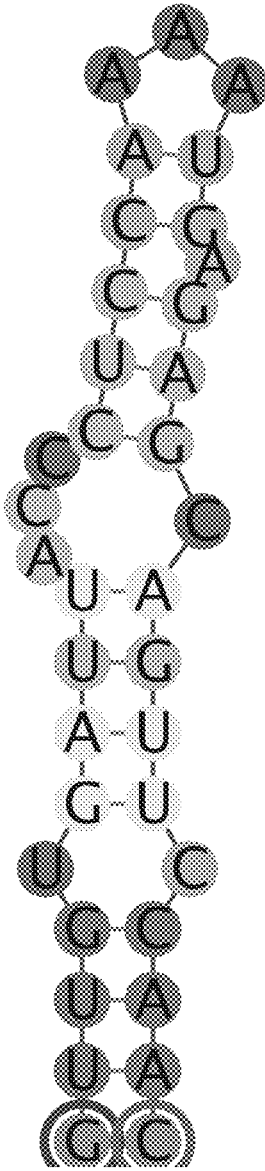


FIG. 13P

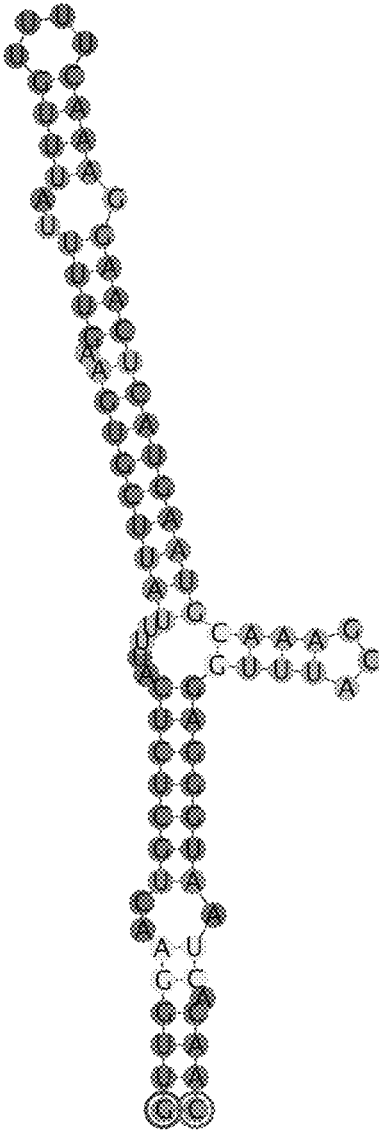


FIG. 13Q

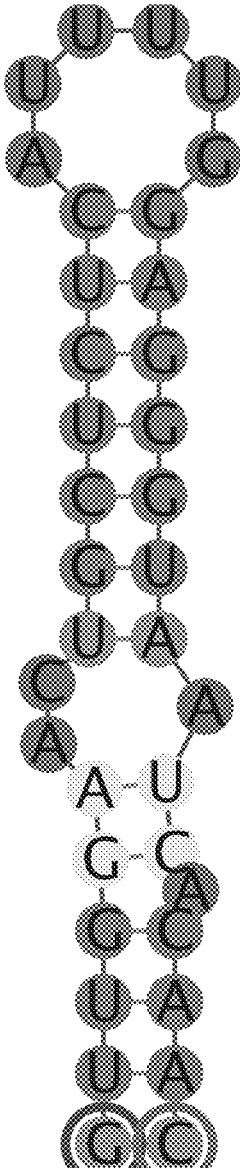


FIG. 13R

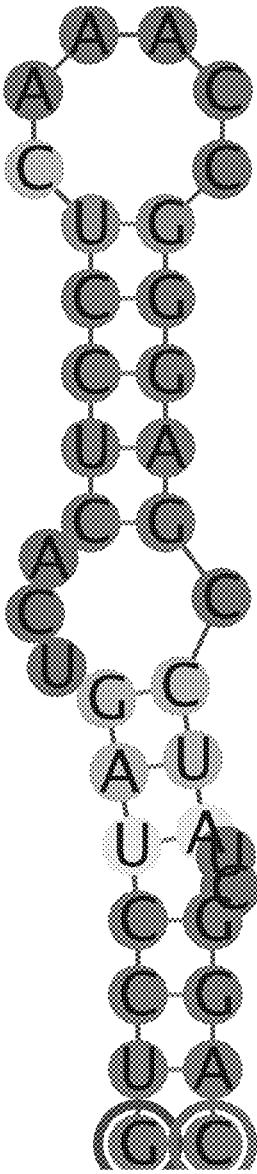


FIG. 13S

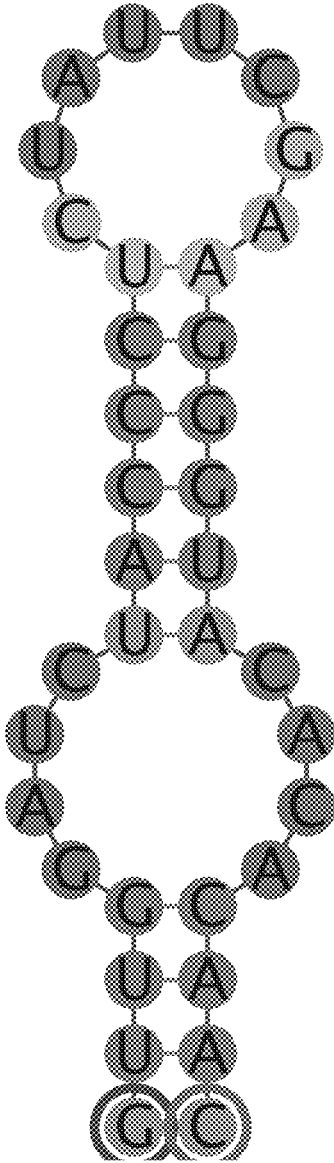


FIG. 13T

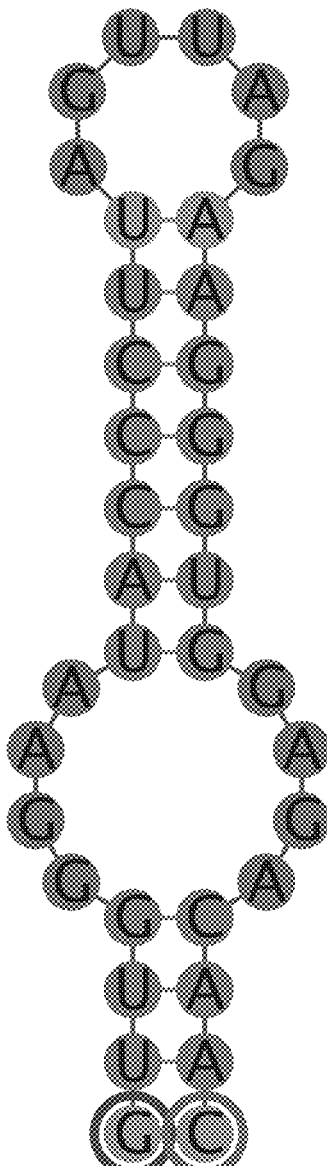


FIG. 13U

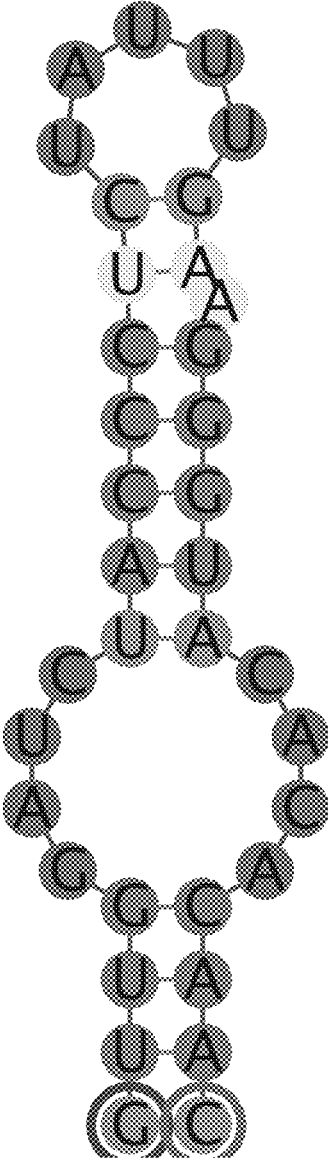


FIG. 13V

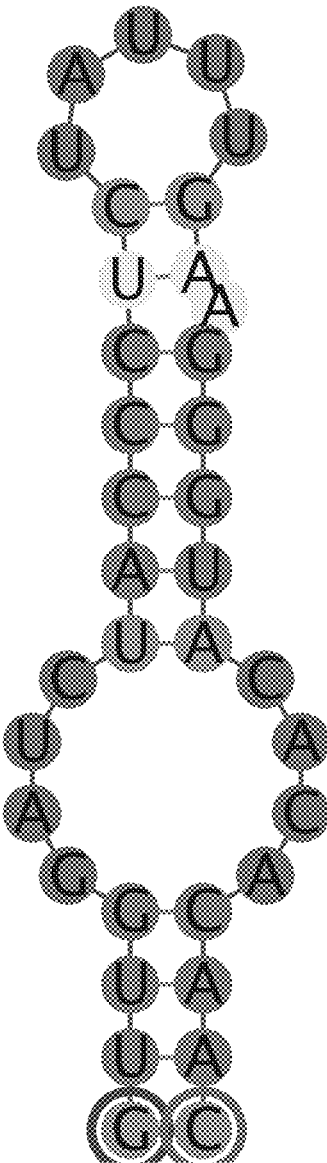


FIG. 13W

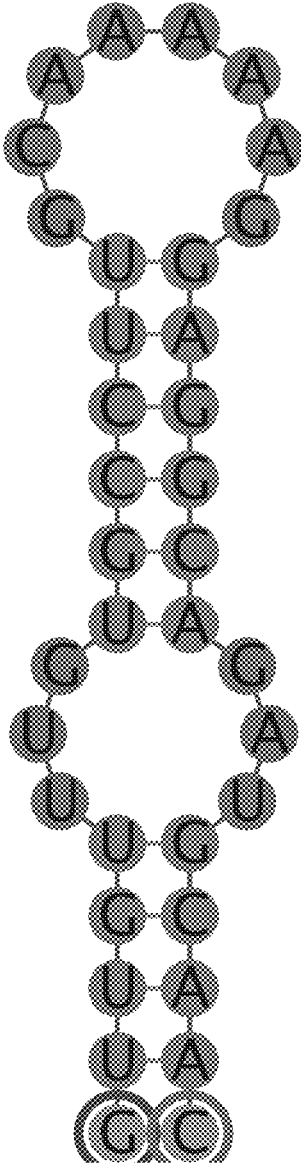


FIG. 13X

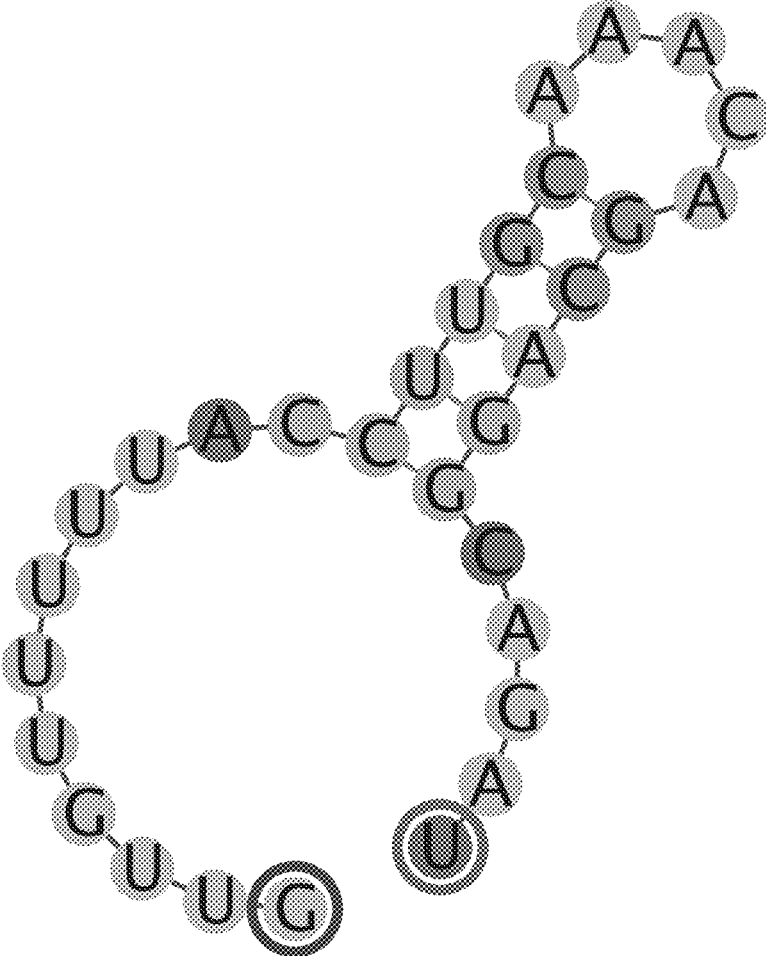


FIG. 13Y

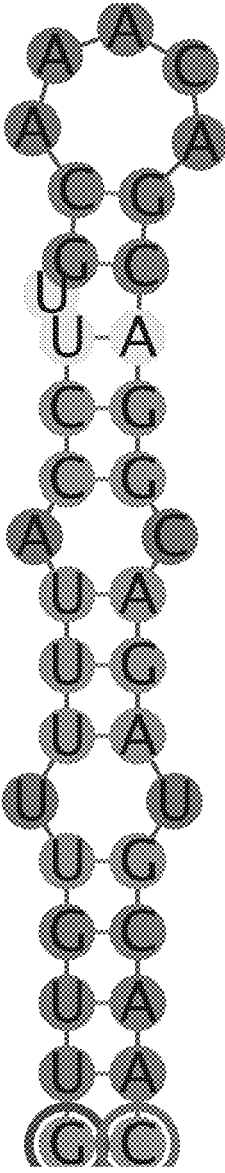


FIG. 13Z

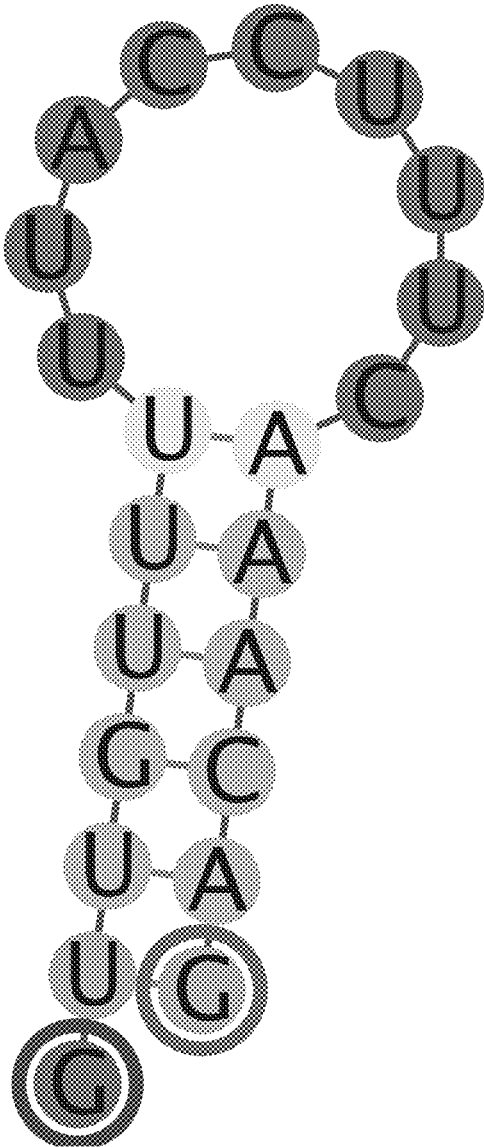


FIG. 13AA

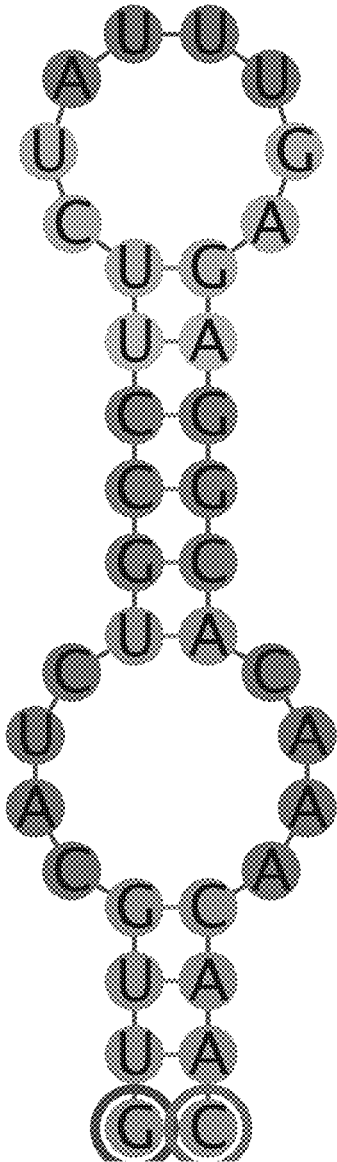


FIG. 13BB

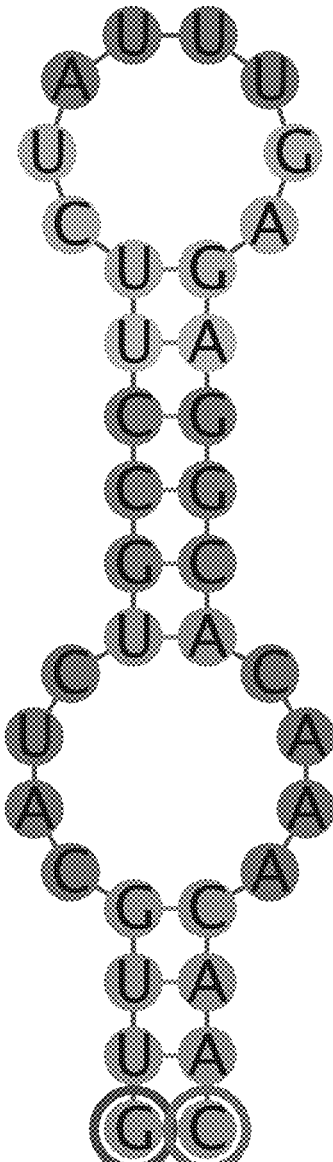


FIG. 13CC

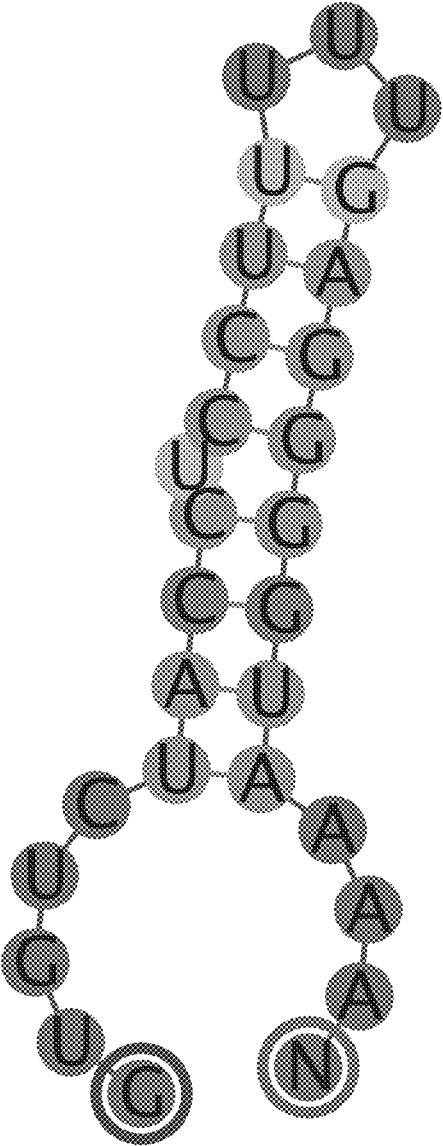


FIG. 13DD

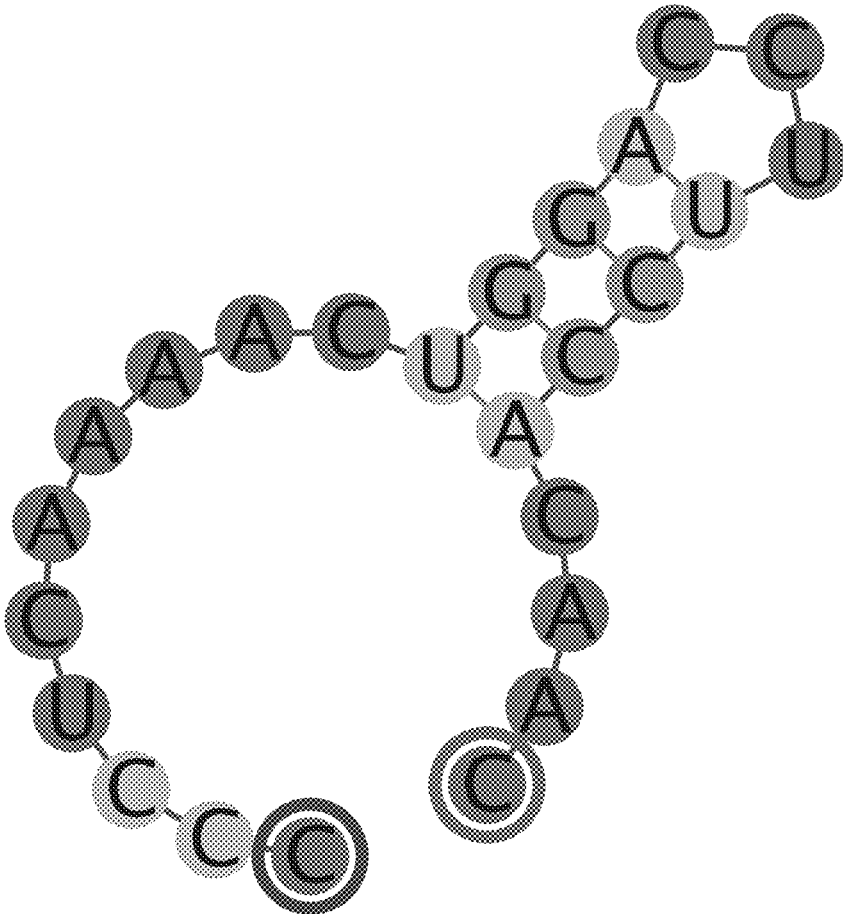


FIG. 13EE

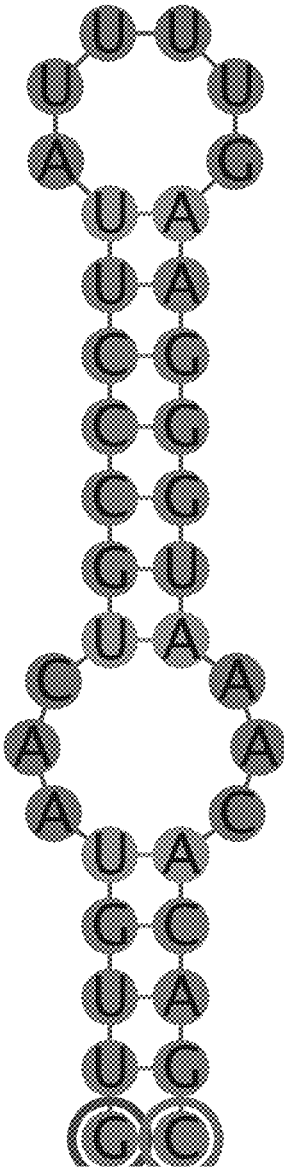


FIG. 13FF

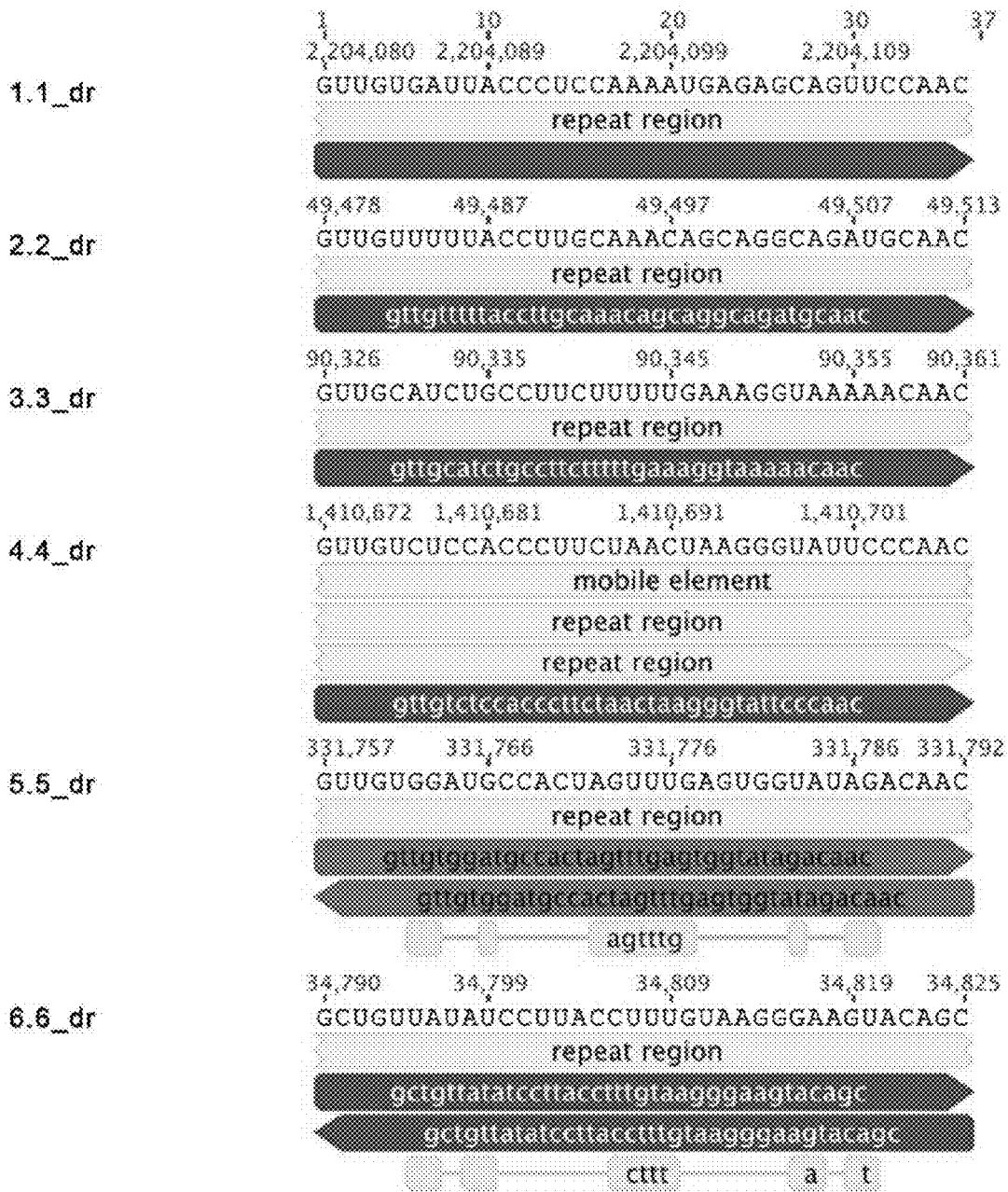


FIG. 14A

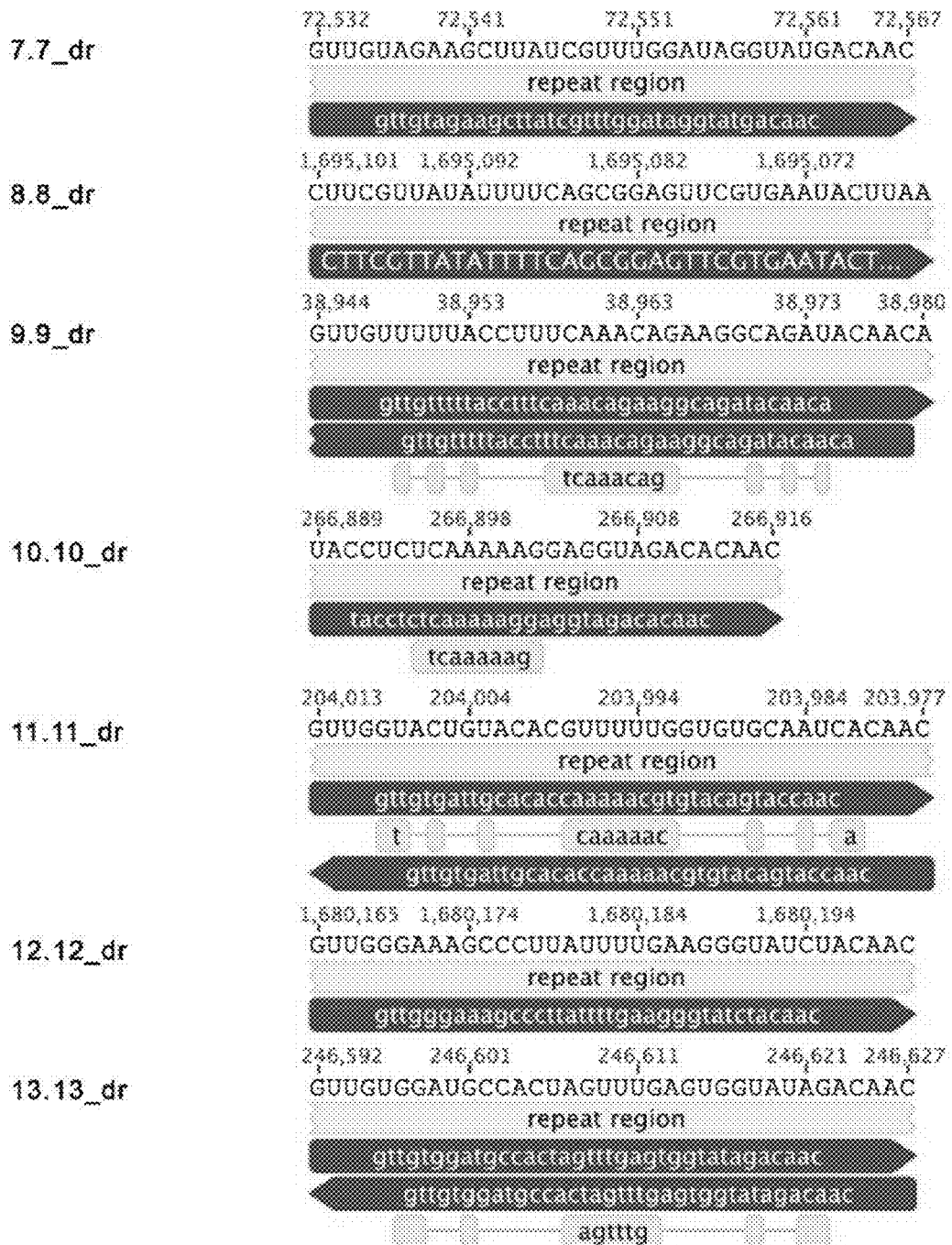


FIG. 14B

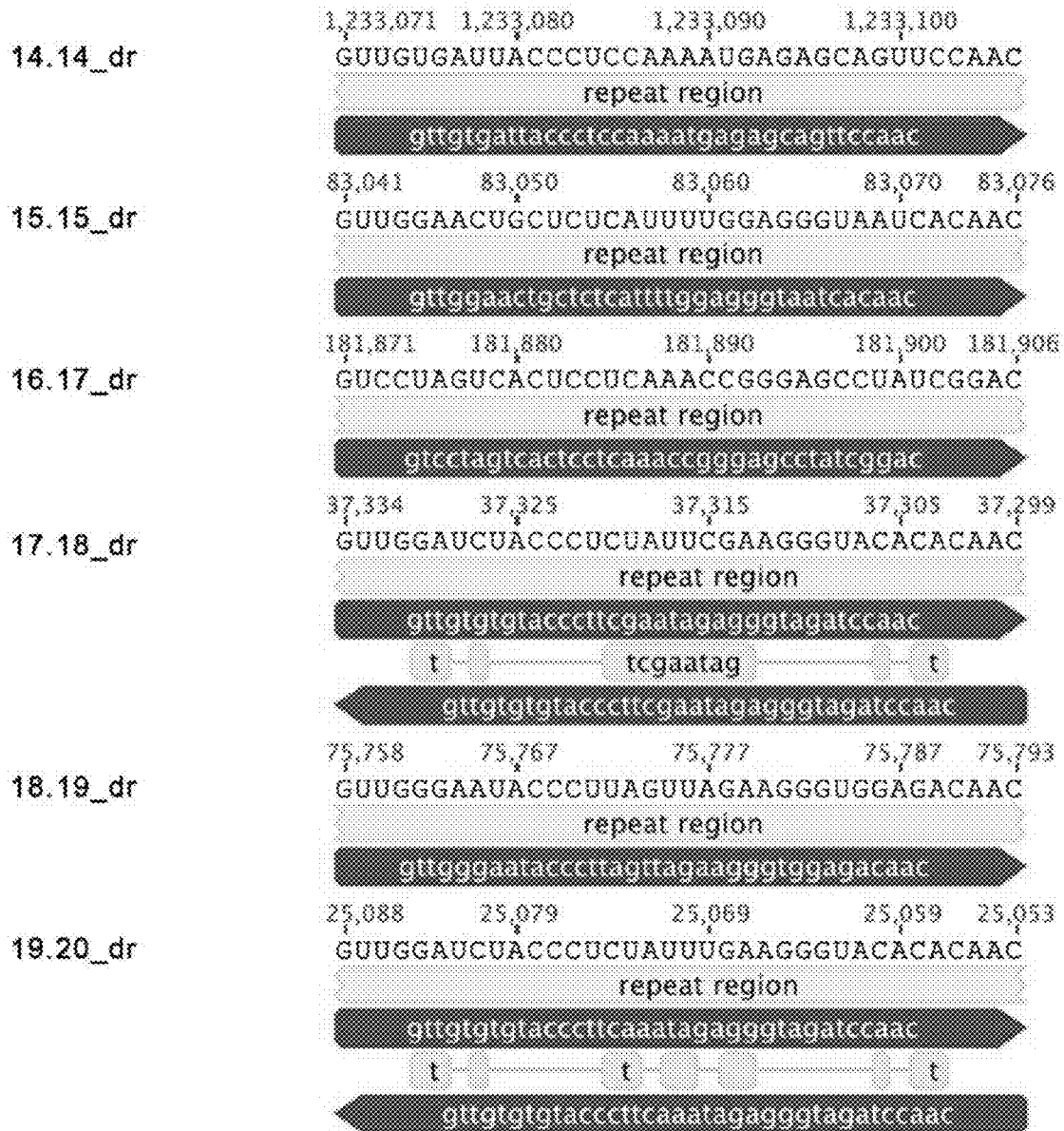


FIG. 14C

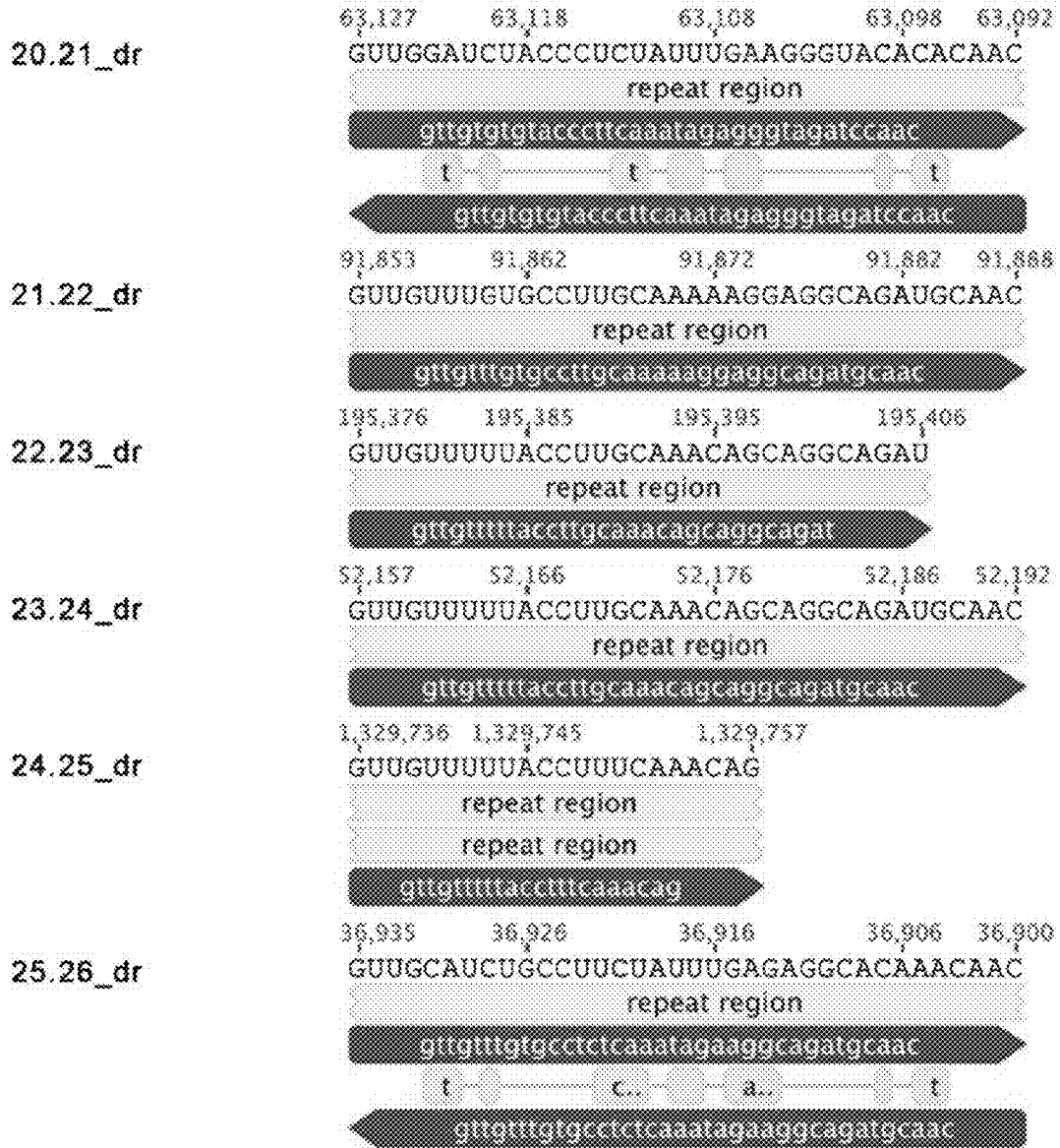


FIG. 14D

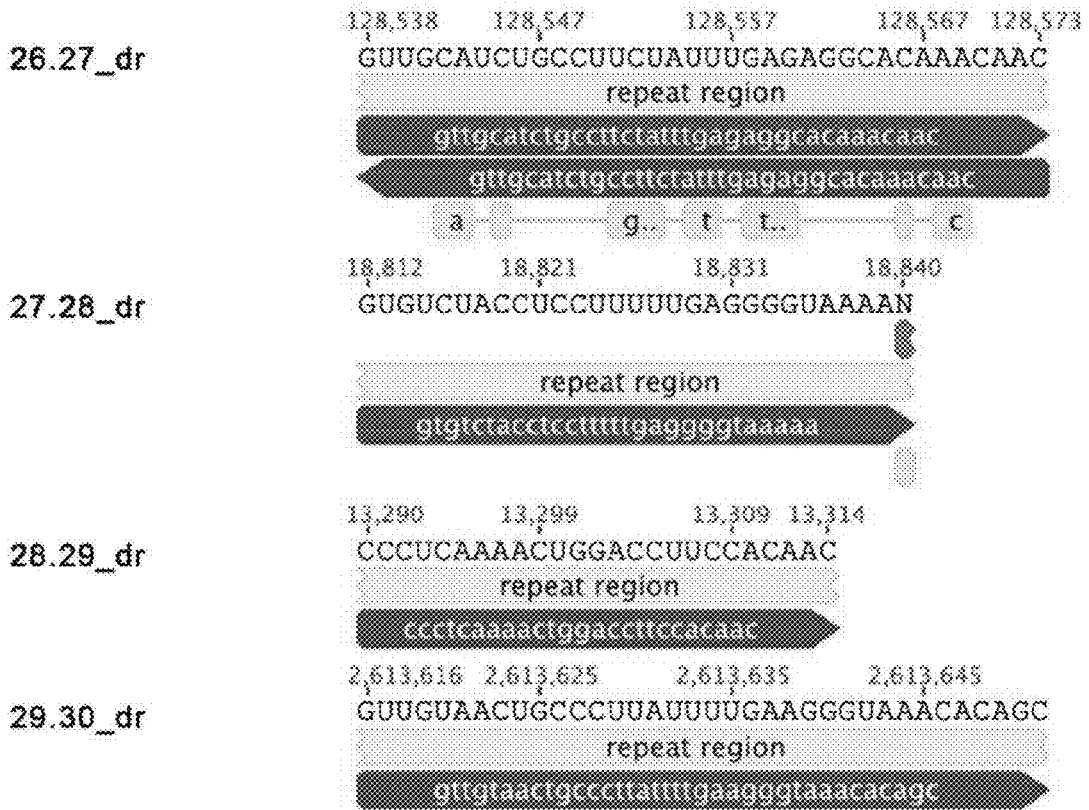


FIG. 14E

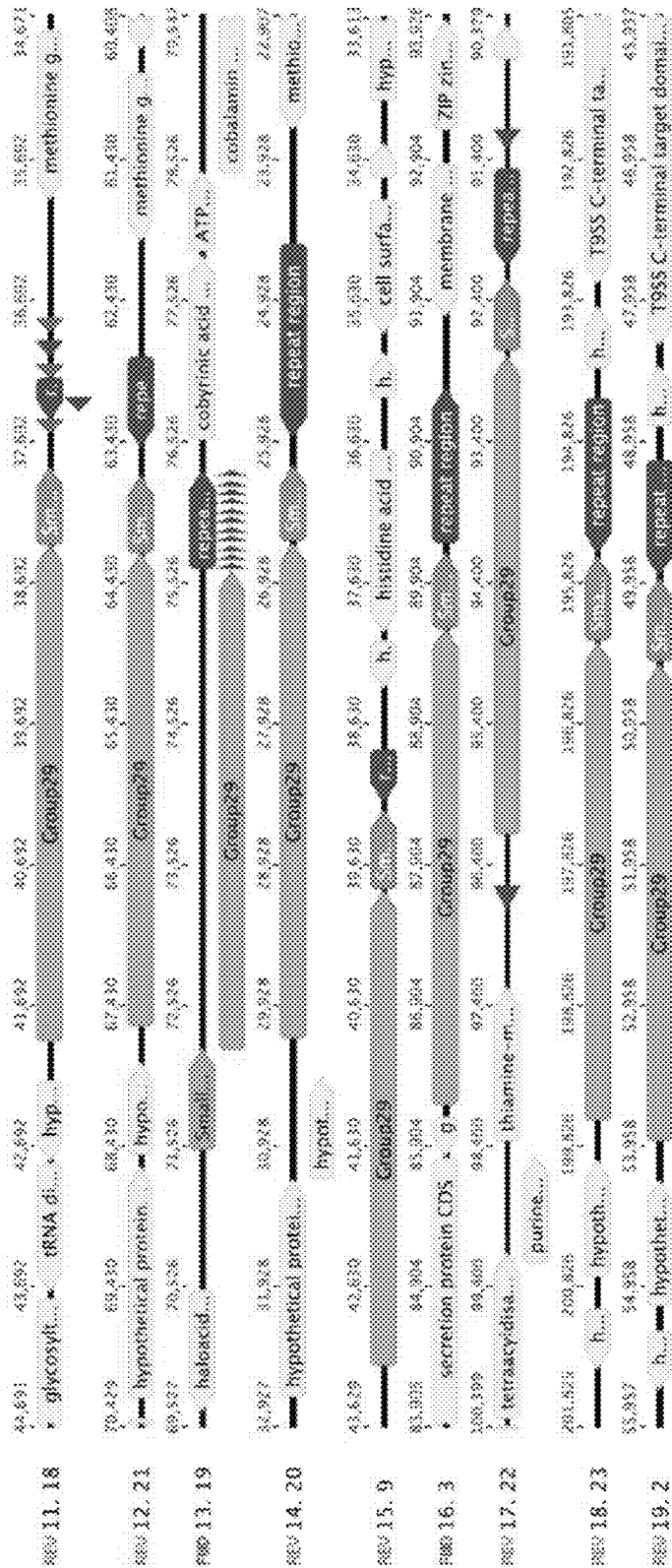


FIG. 16B

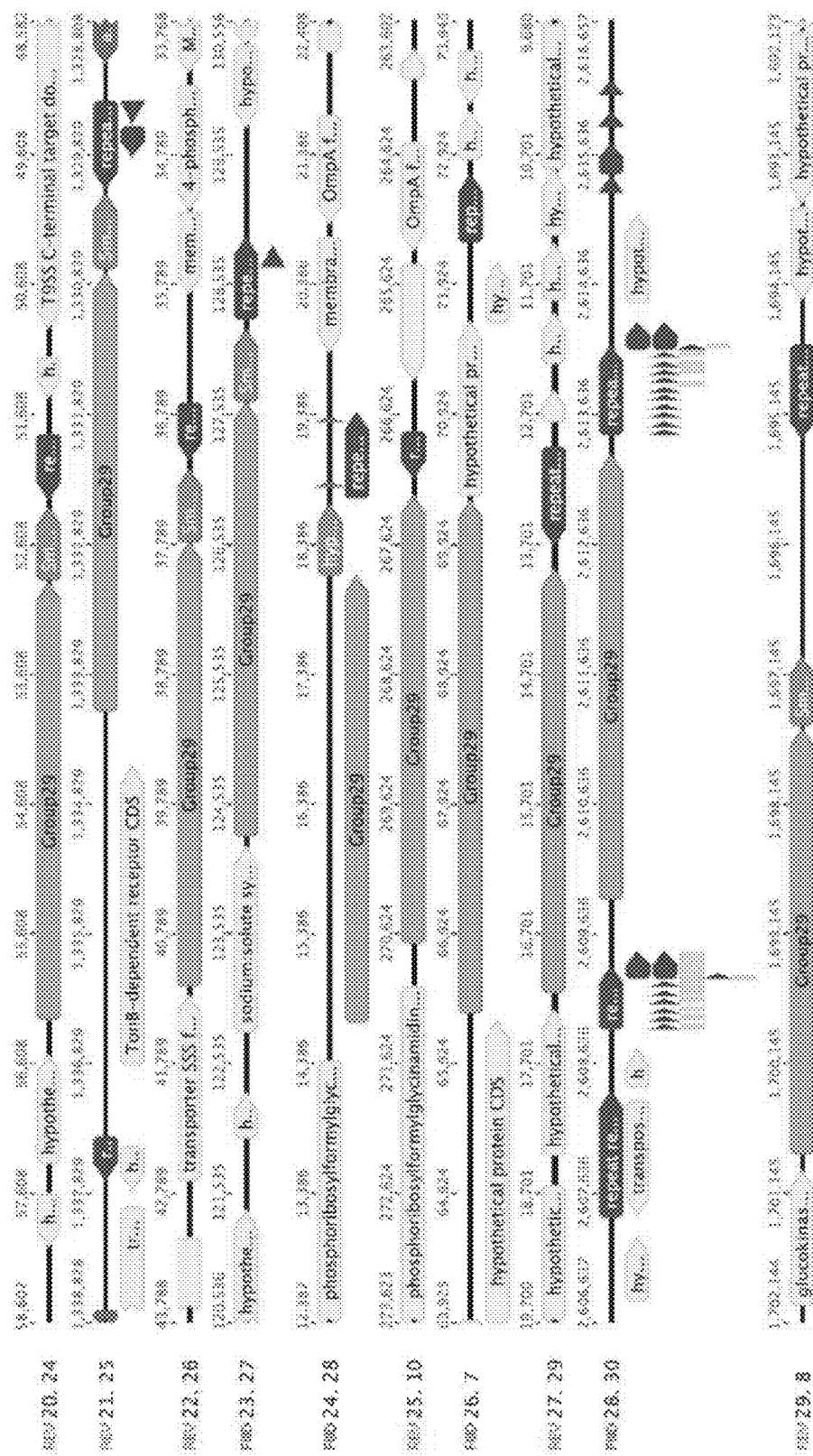


FIG. 16C

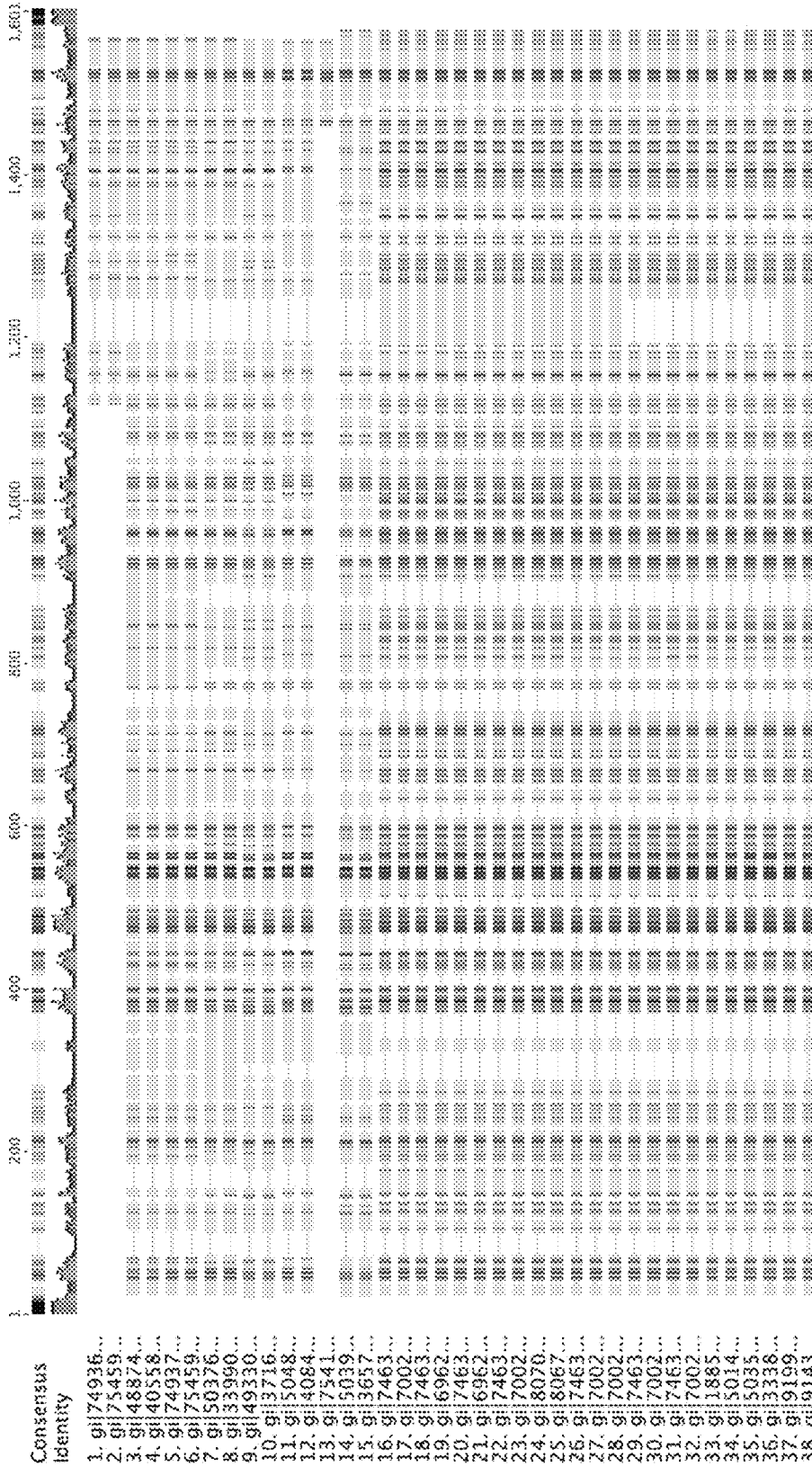


FIG. 17A

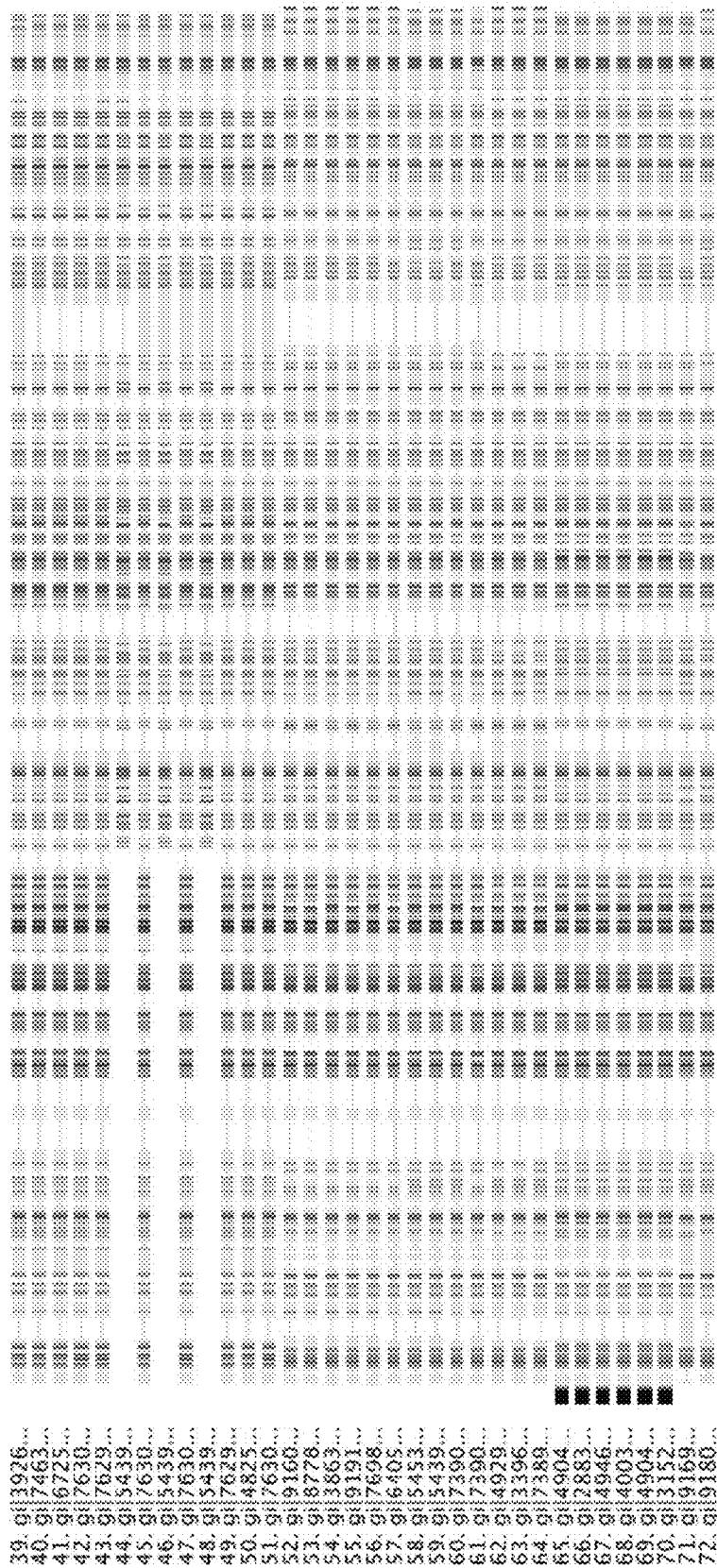
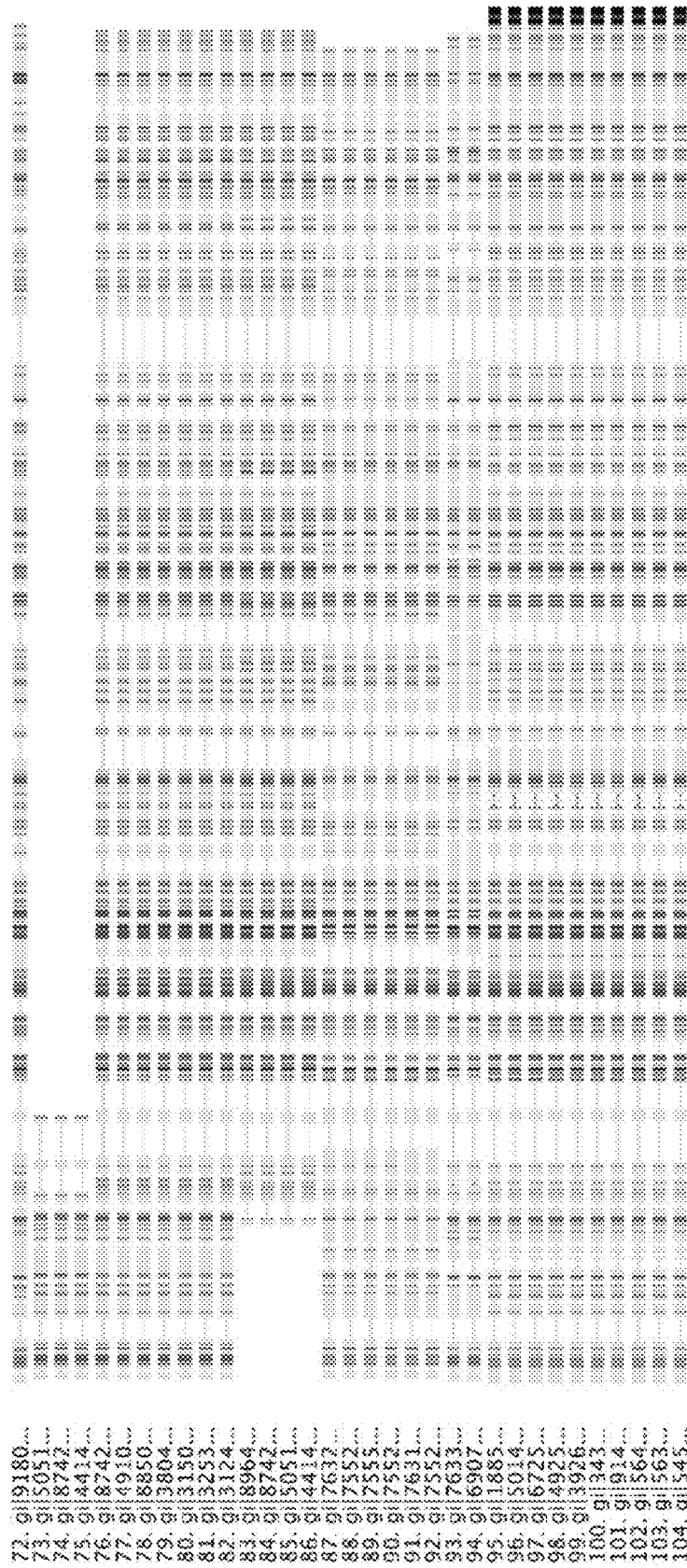


FIG. 17B



72. 9|19180...
73. 9|15051...
74. 9|18742...
75. 9|14414...
76. 9|18742...
77. 9|14910...
78. 9|18850...
79. 9|13804...
80. 9|13150...
81. 9|12533...
82. 9|13124...
83. 9|18964...
84. 9|18742...
85. 9|15051...
86. 9|14414...
87. 9|17632...
88. 9|17552...
89. 9|17552...
90. 9|17552...
91. 9|17631...
92. 9|17552...
93. 9|17633...
94. 9|19007...
95. 9|11885...
96. 9|15014...
97. 9|16725...
98. 9|14925...
99. 9|19266...
100. 9|1343...
101. 9|1814...
102. 9|1564...
103. 9|1563...
104. 9|1545...

Fig. 17C

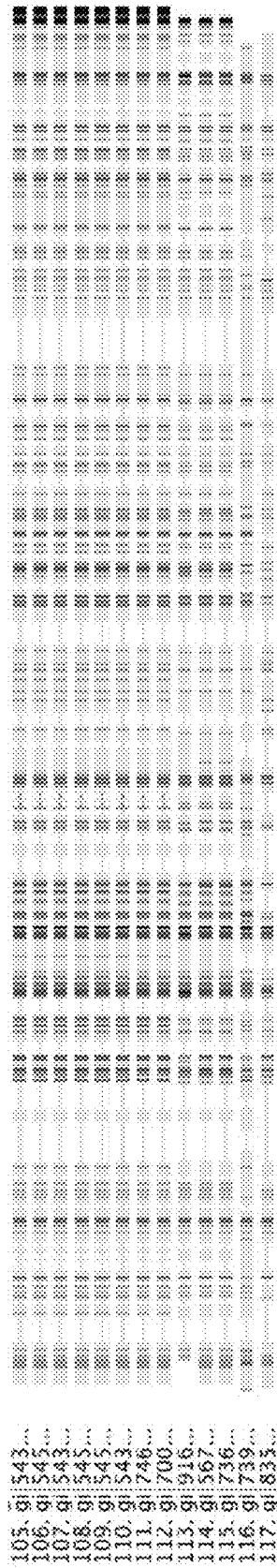
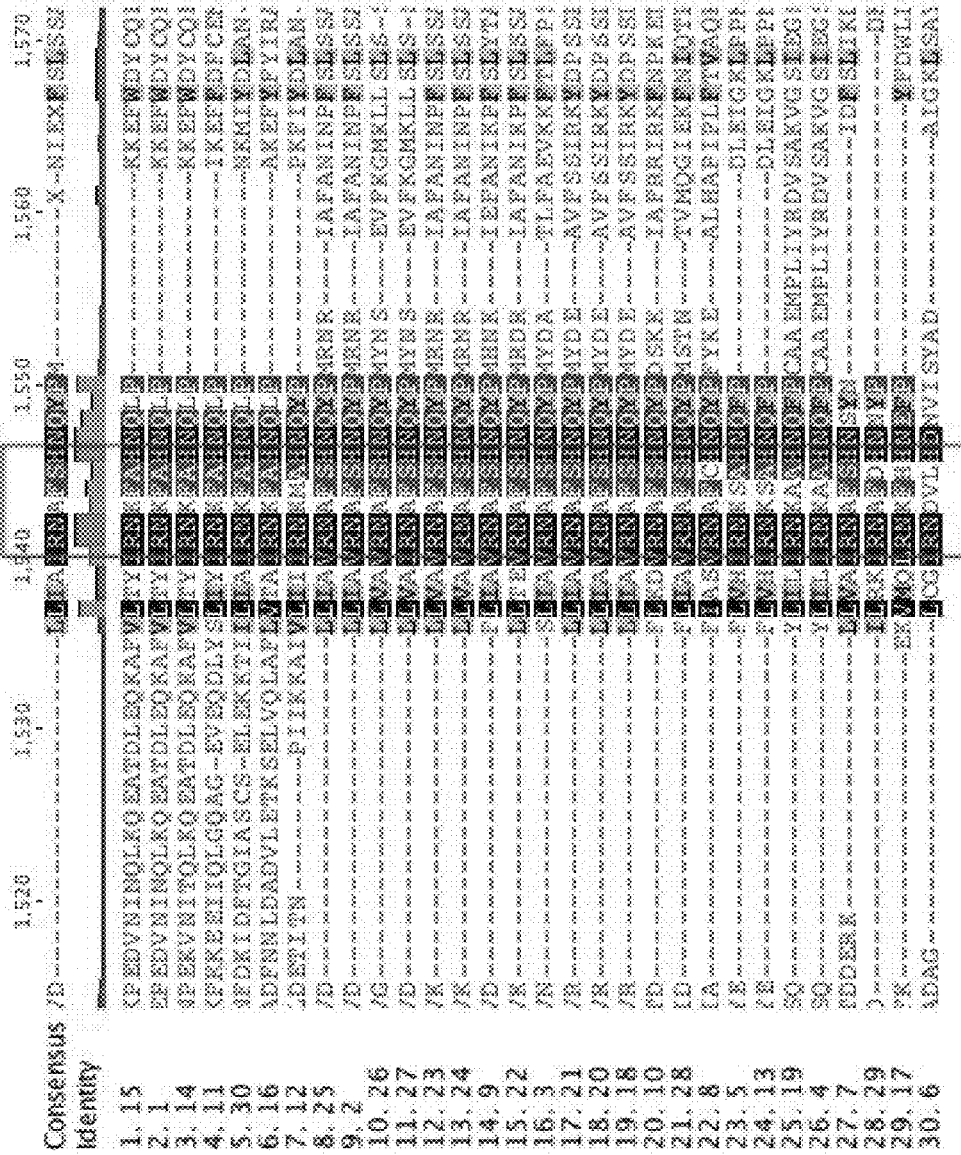


FIG. 17D

Consensus Identity	170	180	190	200	210
1. 15	EEVPIKERREY	ENFRGIT	EDN	EEVEHITD	EEFGVLEEM
2. 1	EEVPIKERREY	ENFRGIT	EDN	EEVEHITD	EEFGVLEEM
3. 14	LEEDAFKRVIC	GVKRL	EDN	EEPIKED	EEFVPLVY
4. 11	IRRRRTIDLL	SRIRQF	RRHISLWR	EEEDIVIE	EEVLDPLN
5. 30	EEENNENRWVE	ENFILLAA	EDN	EEPIFEDN	EEPLLKLL
6. 16	EEENAKT	EEEDTLPES	EDN	EEHLPFDAS	EEELPILL
7. 12	EEENAKT	EEEDTLPES	EDN	EEKPIINP	EEKIYDFL
8. 25	EEQAEQ	EESEK	EEFTV	EEKSE	EEKNP
9. 2	EEQAEQ	EESEK	EEFTV	EEKSE	EEKNP
10. 26	EEKAEQ	EESEK	EEFTV	EEKSE	EEKNP
11. 27	EEKAEQ	EESEK	EEFTV	EEKSE	EEKNP
12. 23	EEKAEQ	EESEK	EEFTV	EEKSE	EEKNP
13. 24	EEKAEQ	EESEK	EEFTV	EEKSE	EEKNP
14. 9	EEKAEQ	EESEK	EEFTV	EEKSE	EEKNP
15. 22	EEKAEQ	EESEK	EEFTV	EEKSE	EEKNP
16. 3	EEKAEQ	EESEK	EEFTV	EEKSE	EEKNP
17. 21	EEKAEQ	EESEK	EEFTV	EEKSE	EEKNP
18. 20	EEKAEQ	EESEK	EEFTV	EEKSE	EEKNP
19. 18	EEKAEQ	EESEK	EEFTV	EEKSE	EEKNP
20. 10	EEKAEQ	EESEK	EEFTV	EEKSE	EEKNP
21. 28	EEKAEQ	EESEK	EEFTV	EEKSE	EEKNP
22. 8	EEKAEQ	EESEK	EEFTV	EEKSE	EEKNP
23. 5	EEKAEQ	EESEK	EEFTV	EEKSE	EEKNP
24. 13	EEKAEQ	EESEK	EEFTV	EEKSE	EEKNP
25. 19	EEKAEQ	EESEK	EEFTV	EEKSE	EEKNP
26. 4	EEKAEQ	EESEK	EEFTV	EEKSE	EEKNP
27. 7	EEKAEQ	EESEK	EEFTV	EEKSE	EEKNP
28. 29	EEKAEQ	EESEK	EEFTV	EEKSE	EEKNP
29. 17	EEKAEQ	EESEK	EEFTV	EEKSE	EEKNP
30. 6	EEKAEQ	EESEK	EEFTV	EEKSE	EEKNP

HEPN motif R-X₍₄₋₆₎H where X is any amino acid

FIG. 18A



HEPN motif R-X₍₄₋₆₎-H where X is any amino acid

FIG. 18B

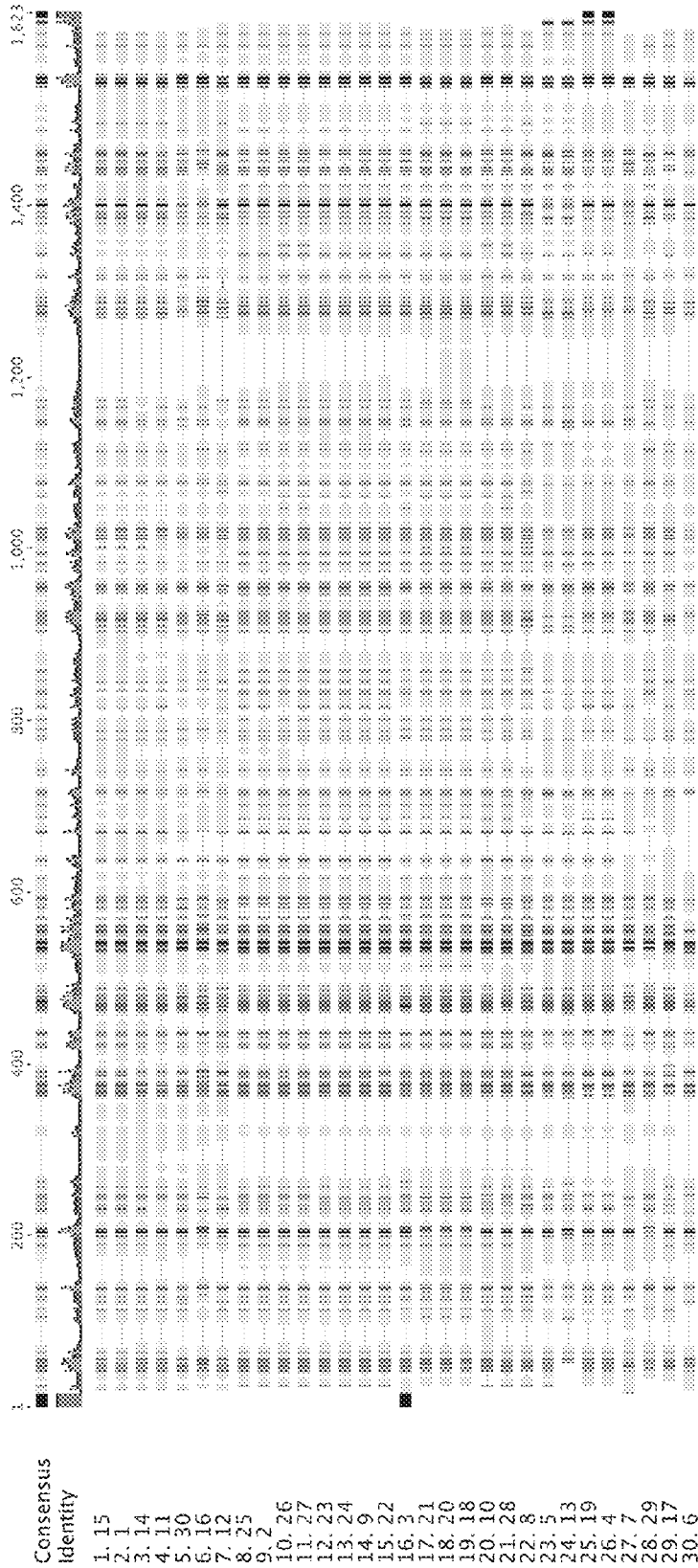


FIG. 19

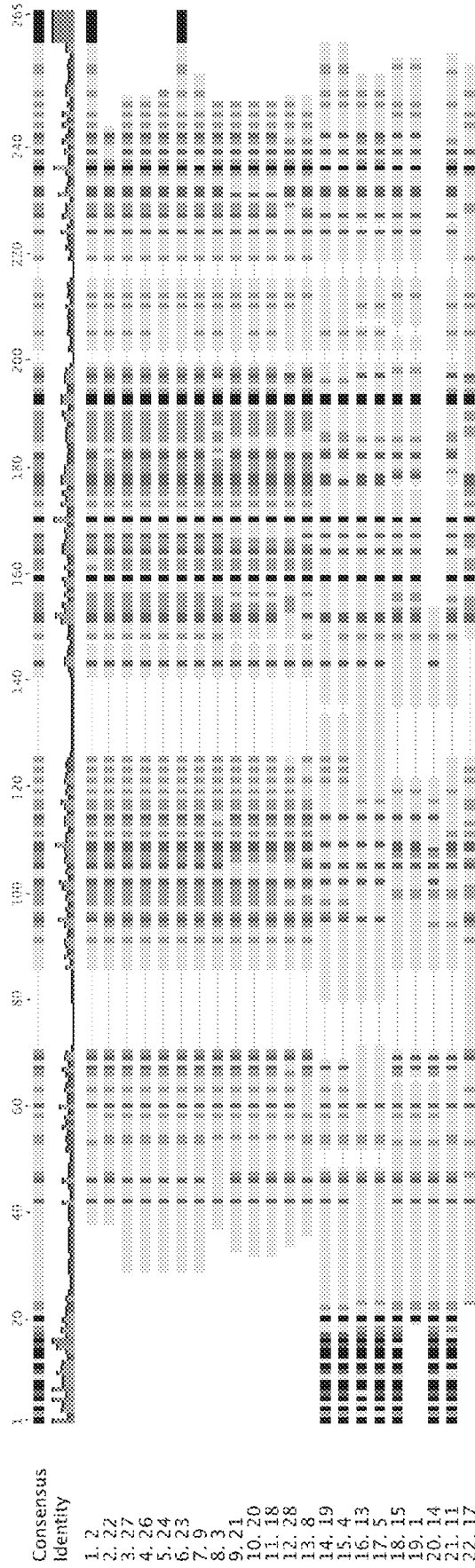


FIG. 20

RNA-seq for small transcripts reveals expression of the CRISPR array adjacent to the group 29 gene in *Bergeyella Zoohelcum*

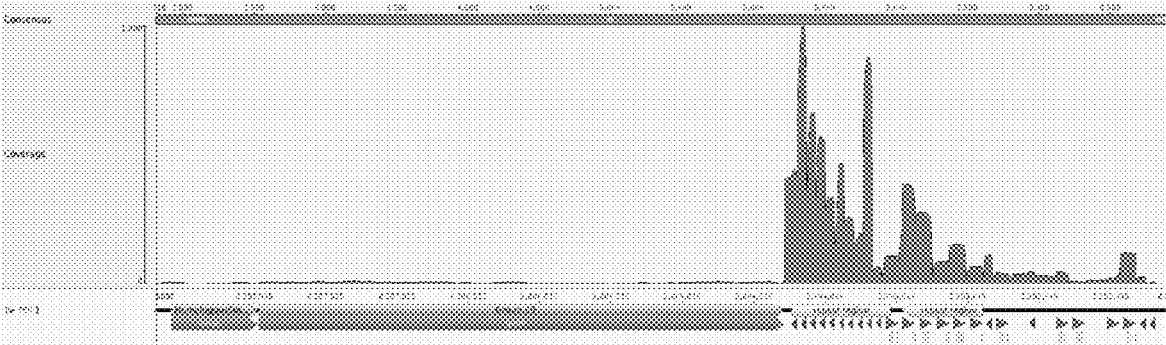


FIG. 21

The CRISPR array next to *B. zoohelcum* group 29 protein contains two DRs, both of which are processed into putative crRNAs

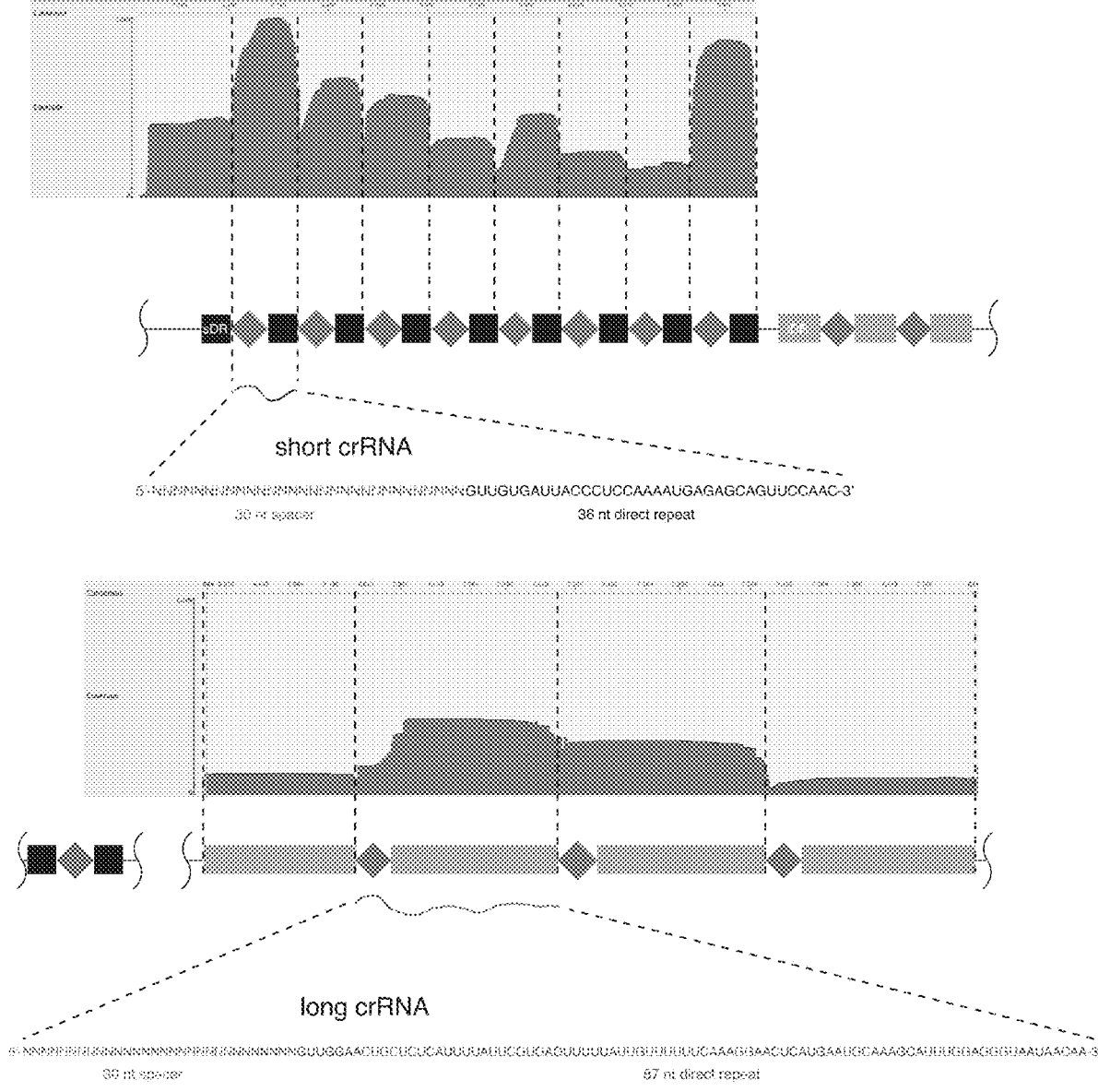


FIG. 22

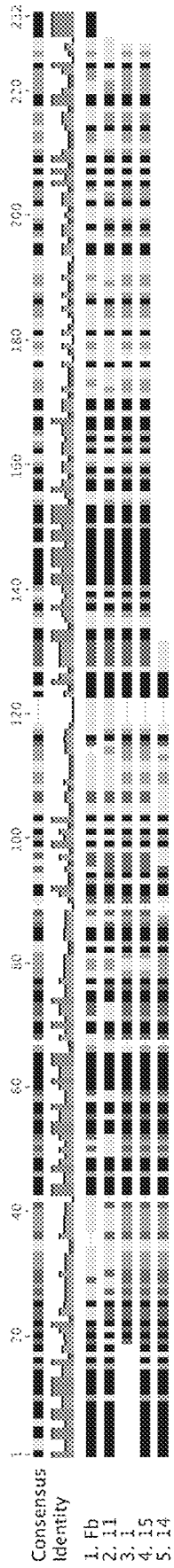


FIG. 23

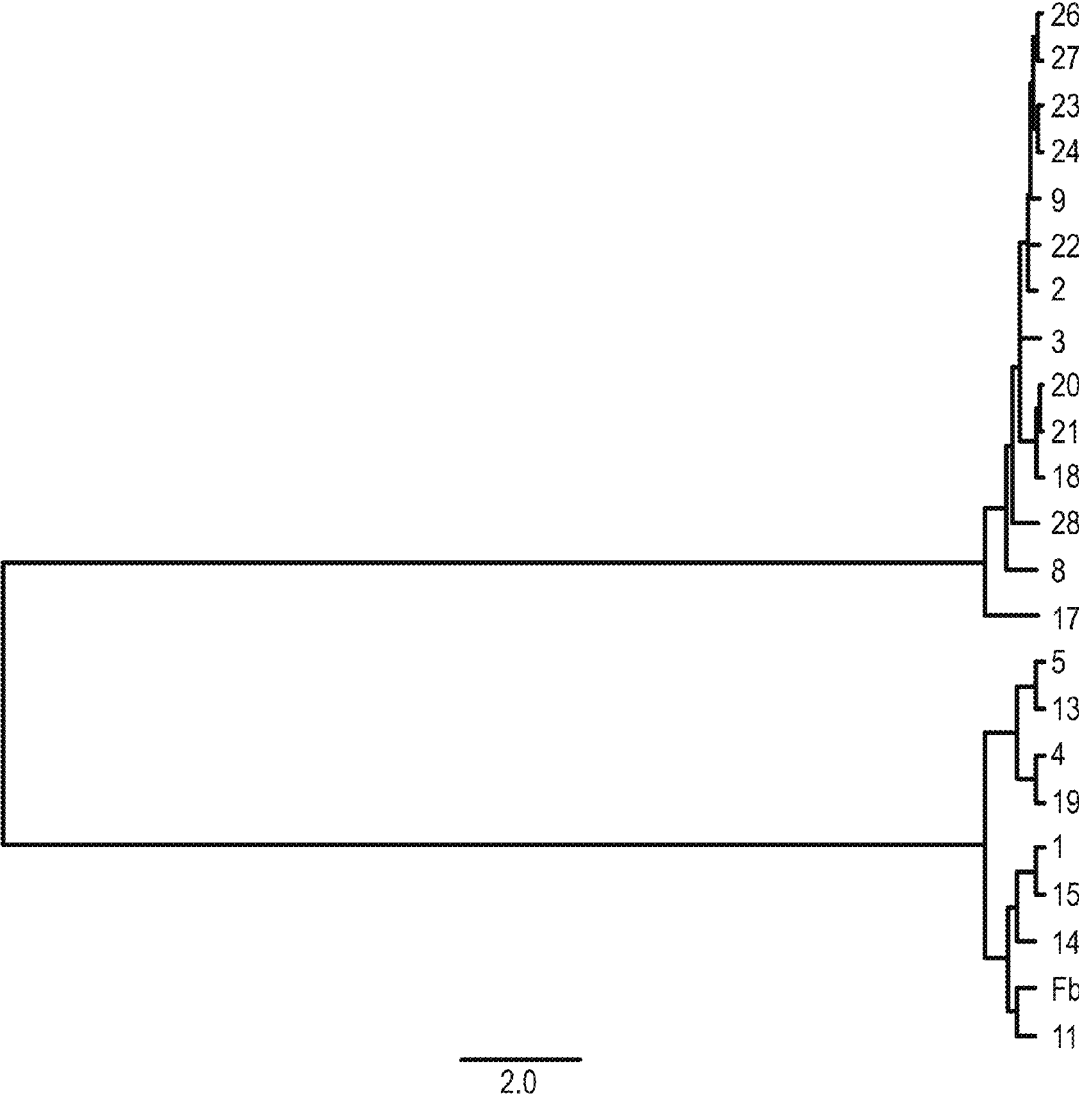


FIG. 24

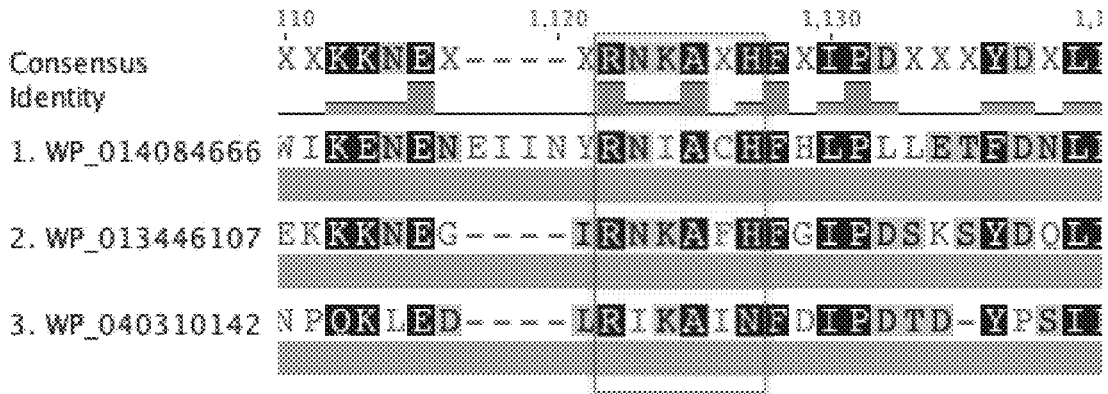
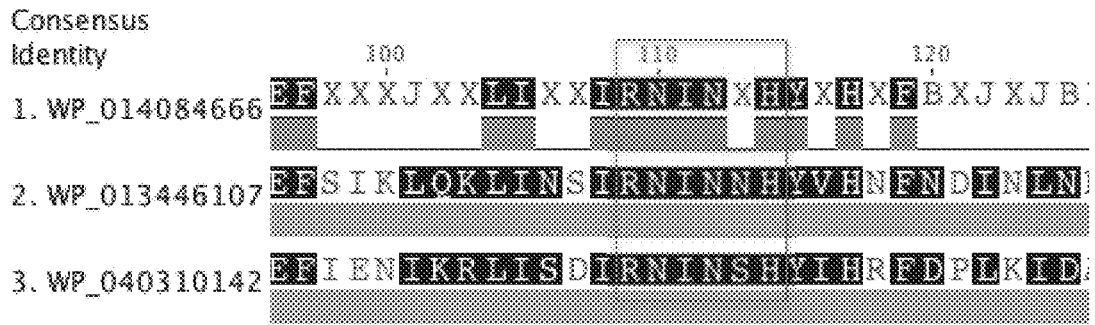


FIG. 25

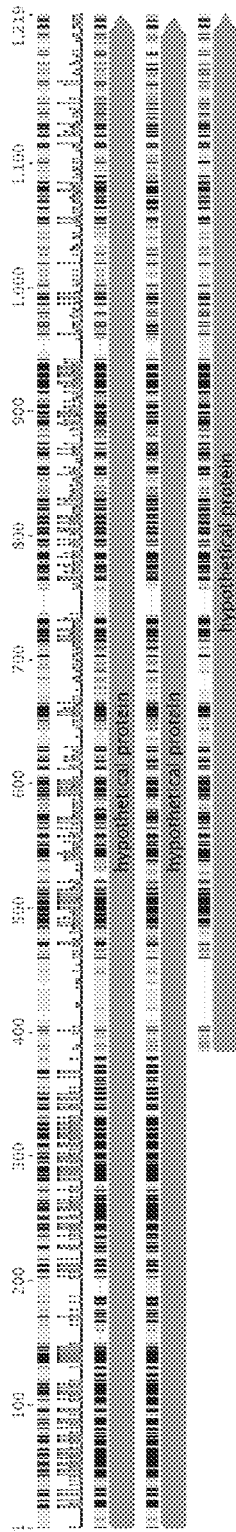


FIG. 27

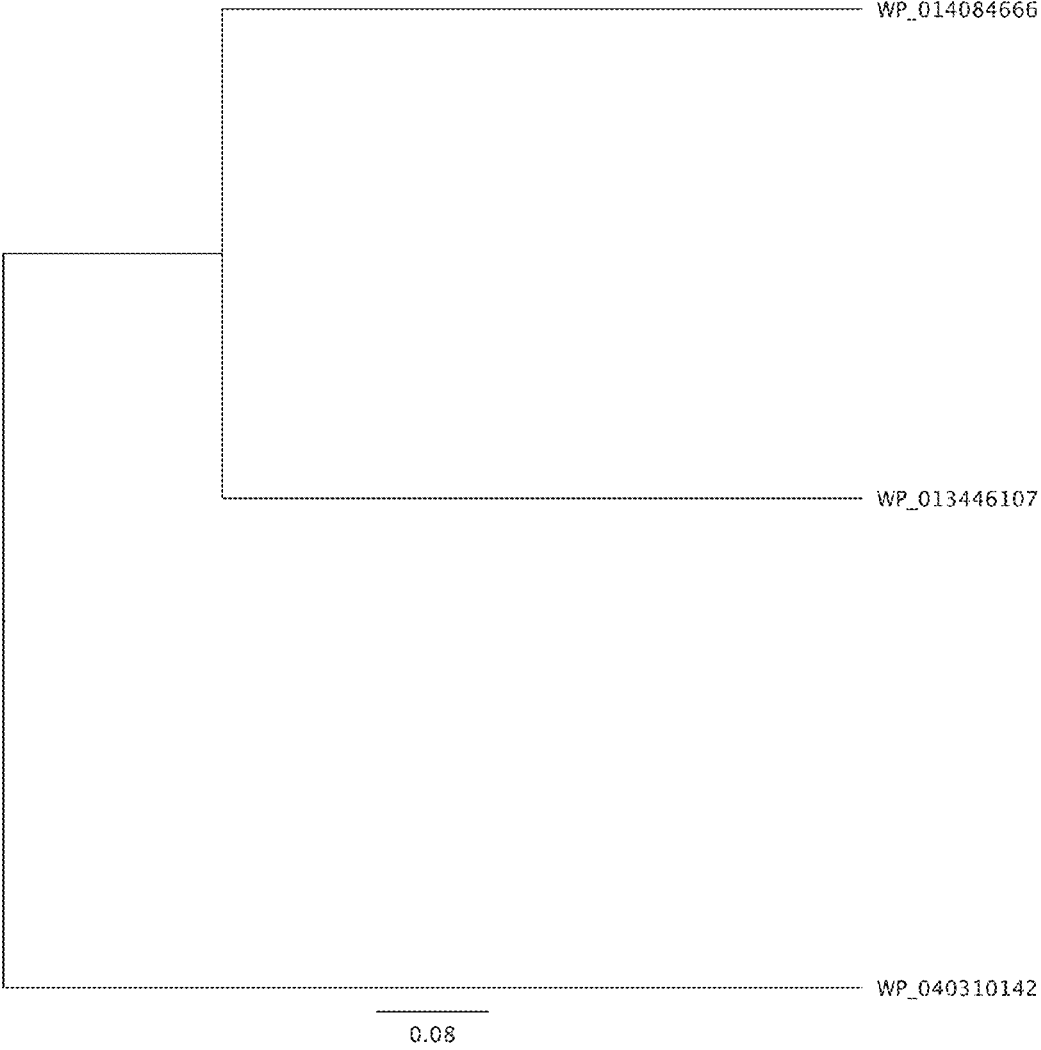


FIG. 28

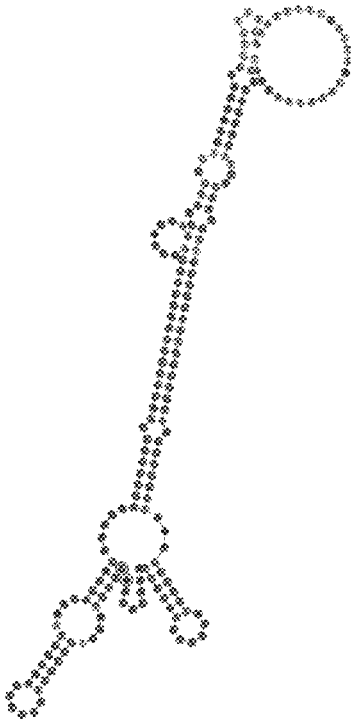


FIG. 29

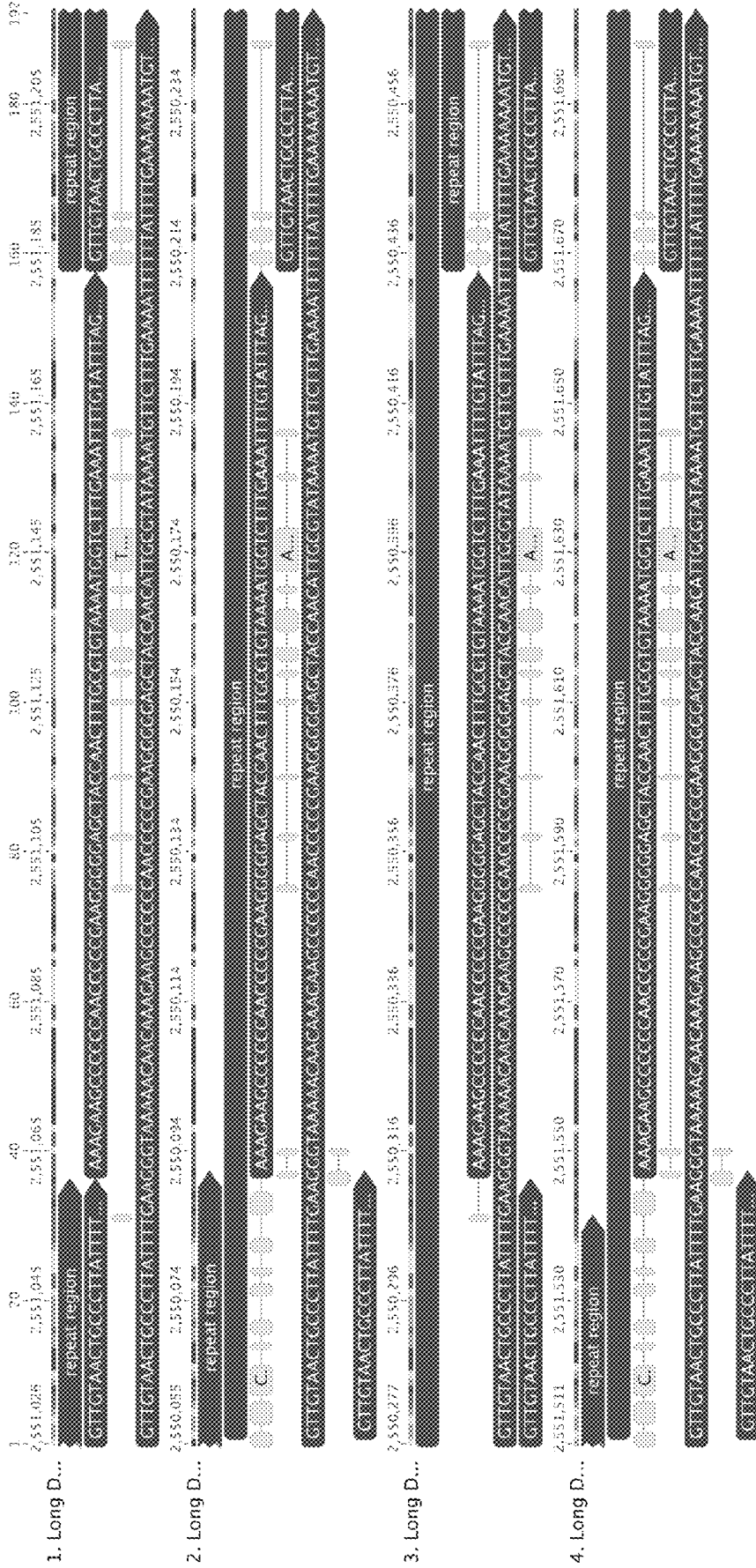


FIG. 30

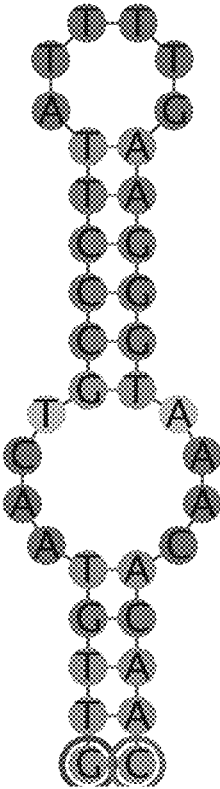


FIG. 31

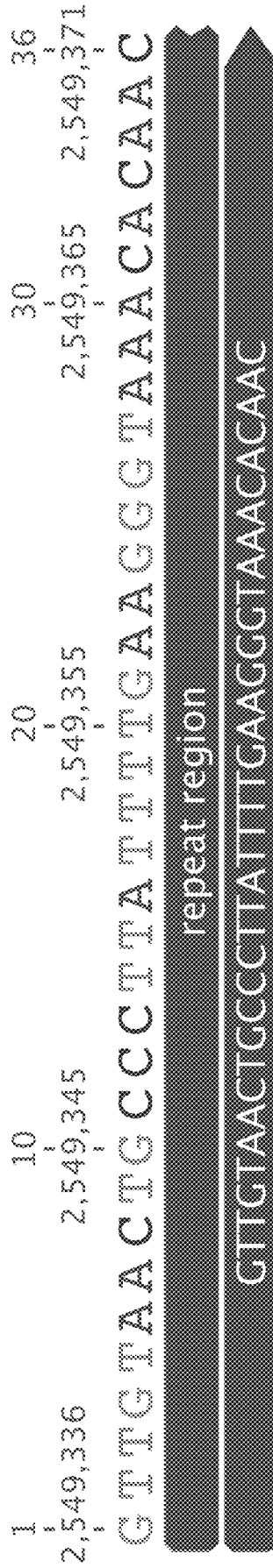


FIG. 32

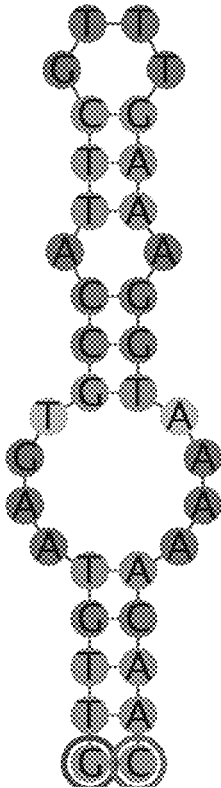


FIG. 33

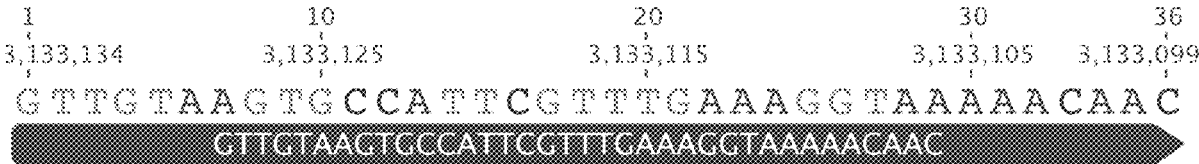


FIG. 34

FIG. 35A

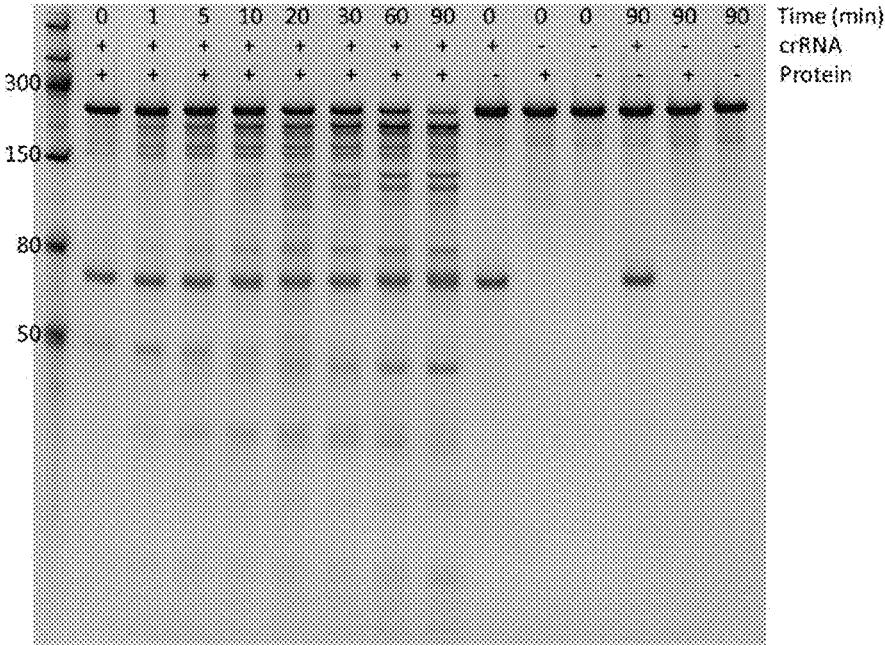


FIG. 35B

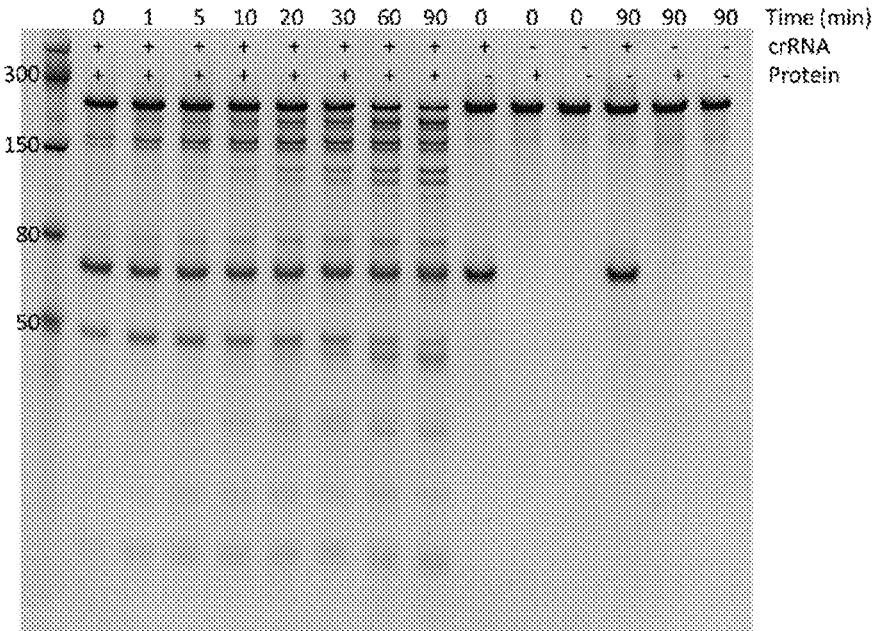


FIG. 36A

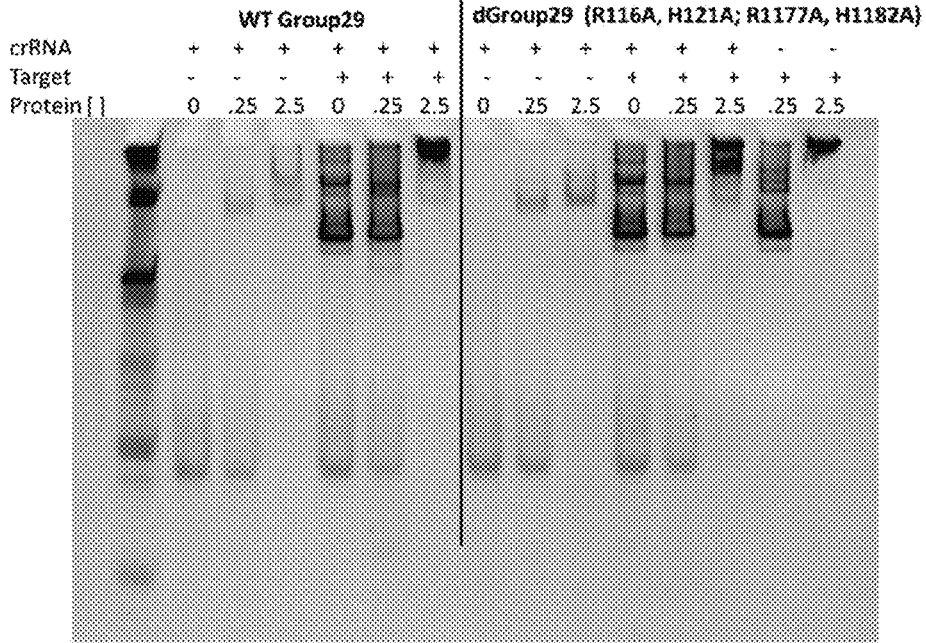


FIG. 36B

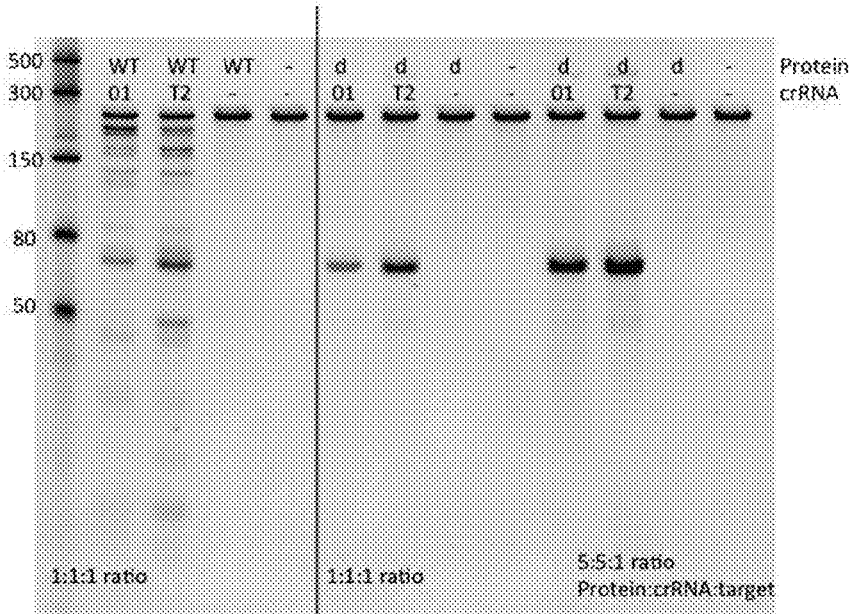


FIG. 36C

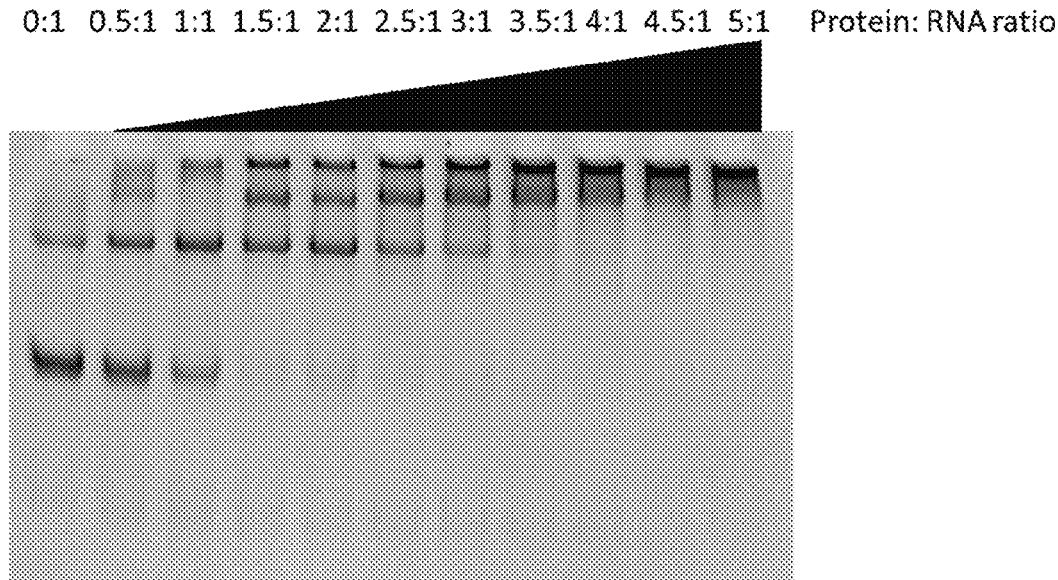


FIG. 36D

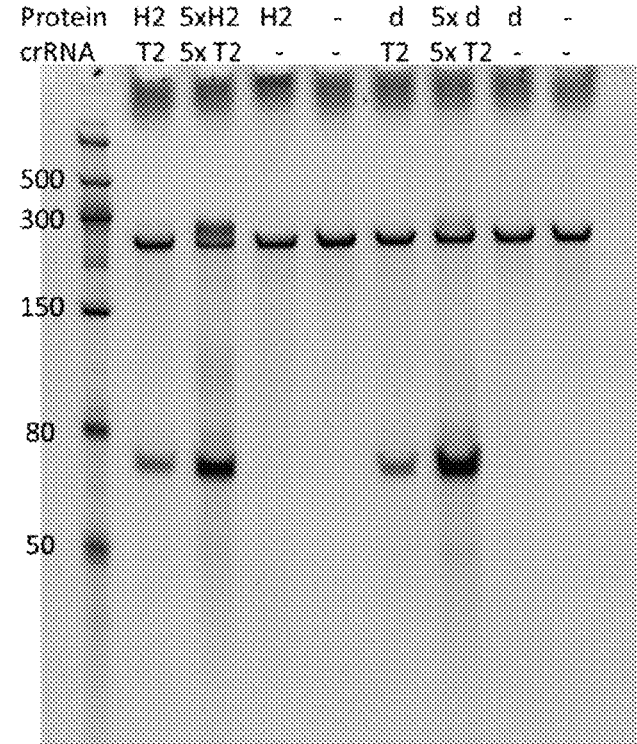


FIG. 37

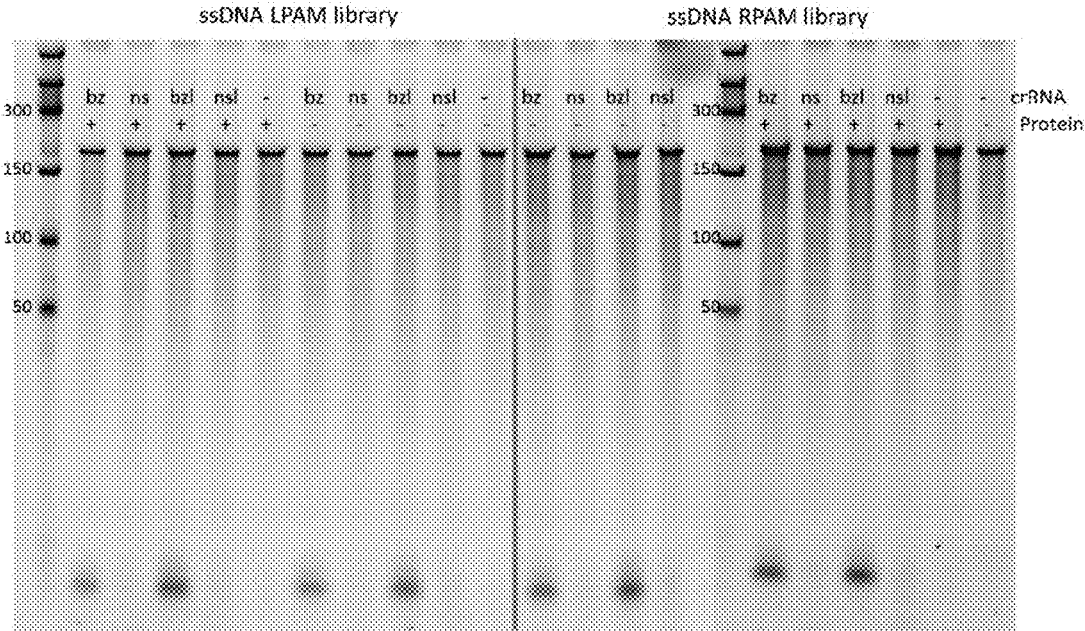


FIG. 38A

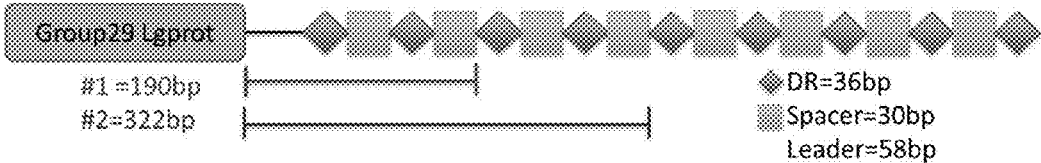


FIG. 38B

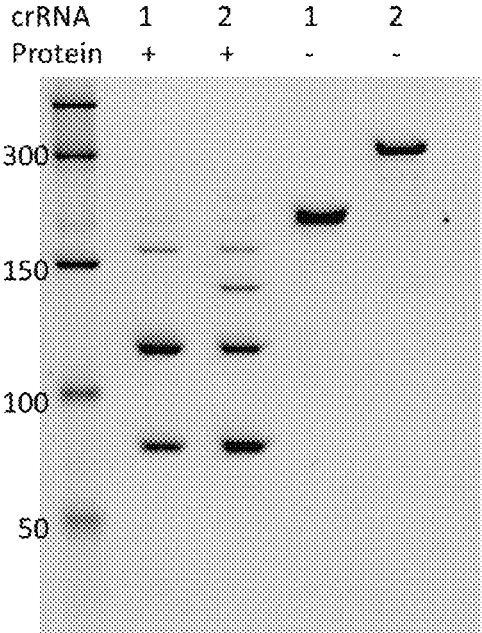


FIG. 39

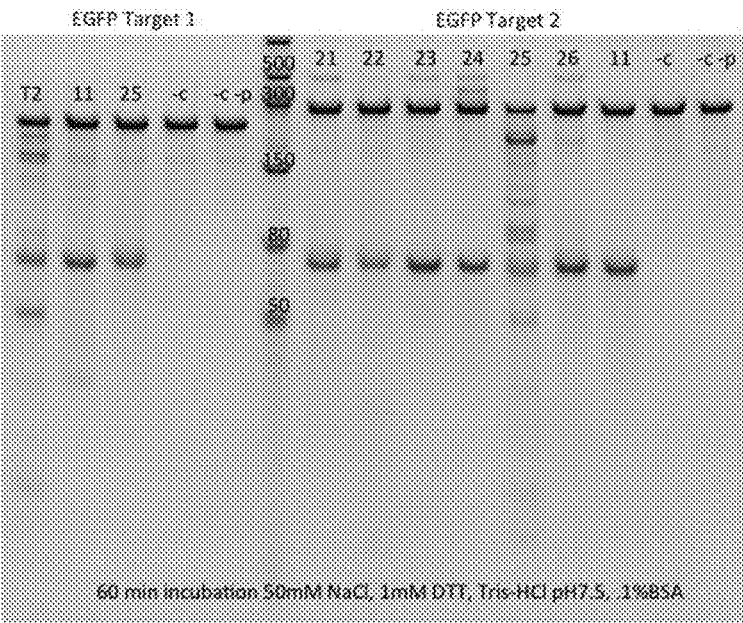
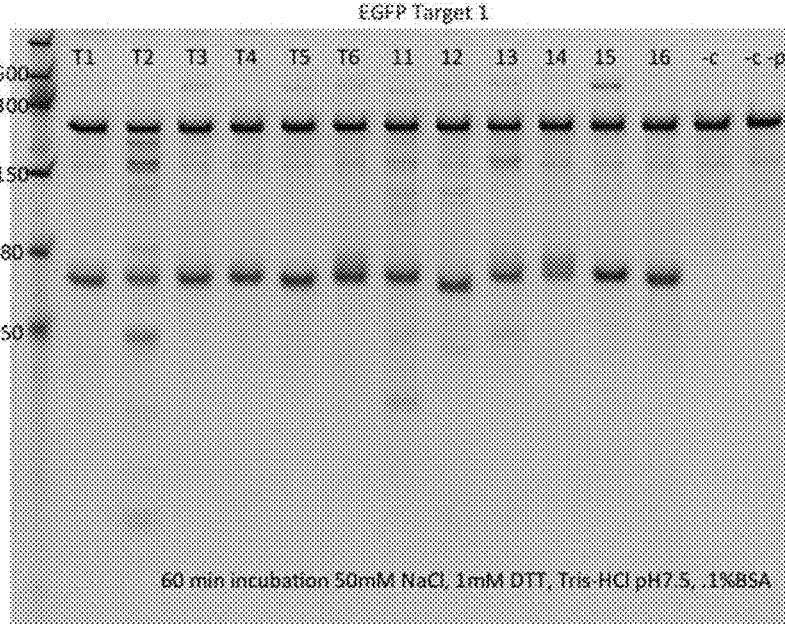


FIG. 40A

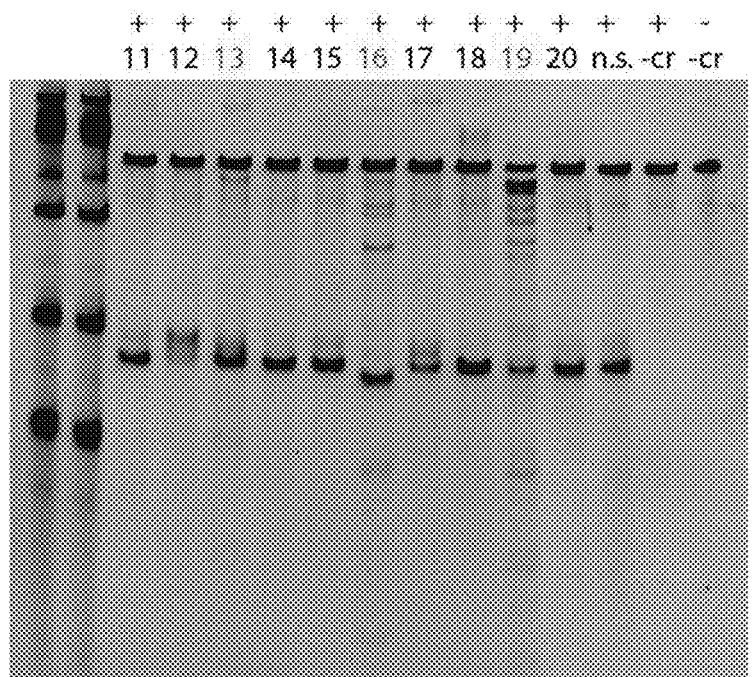
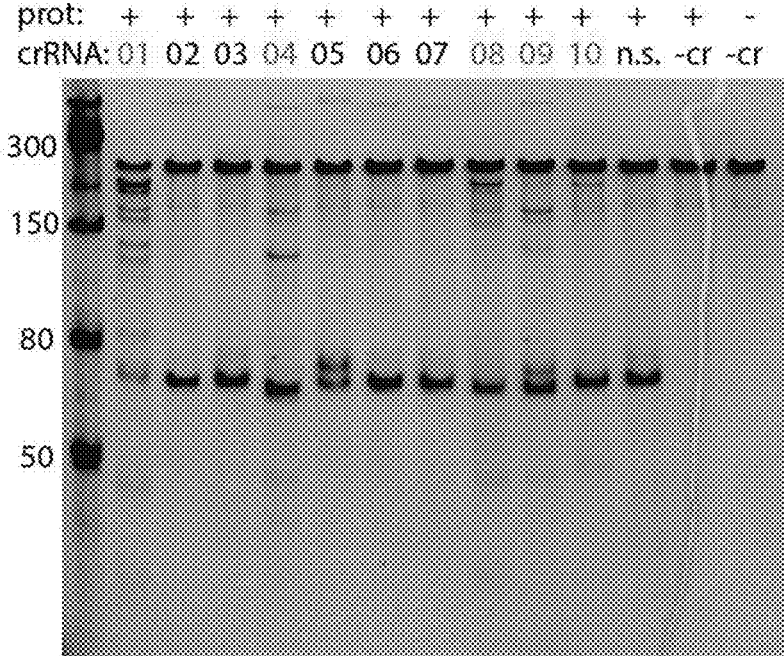
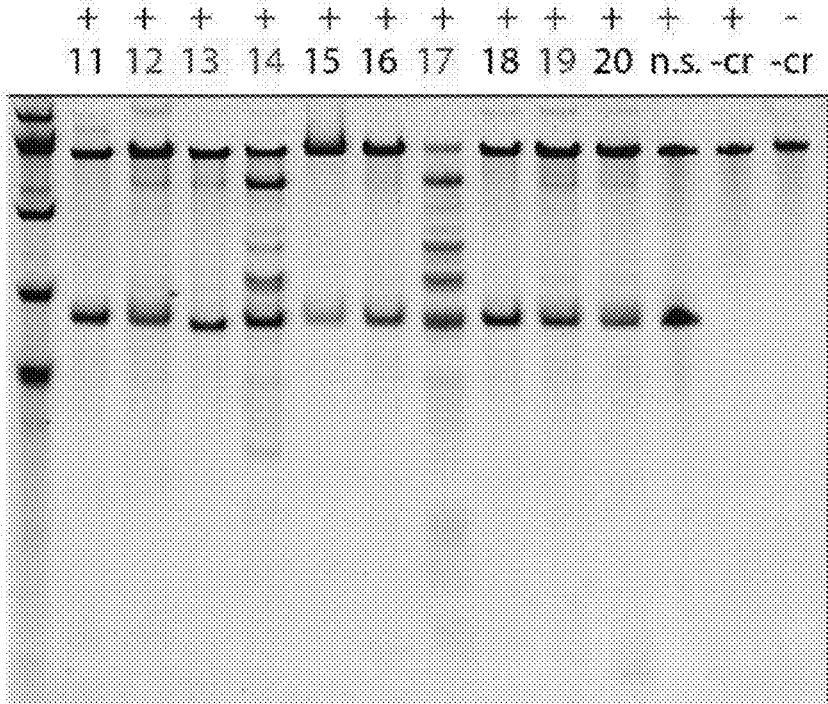
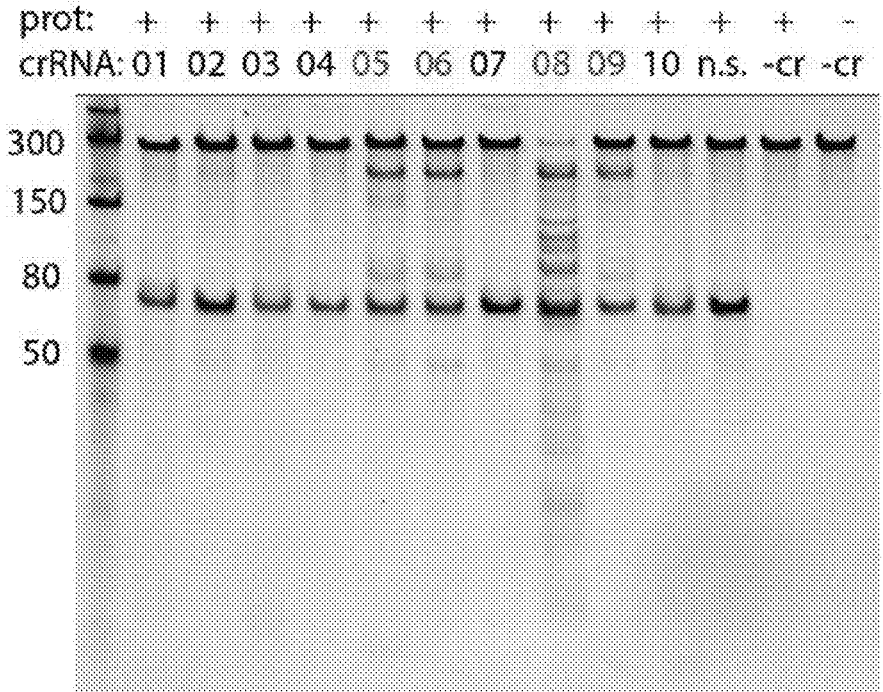


FIG. 40B



EGFP Target 1

2 10 20 30 40 50 60 70
GAUGGGGGUUCUGCUGGUAGUGGUCGGCGAGCUGCACCGCCGUCUCCGAUGUUGUGCCGGAUCUUG
80 90 100 110 120 130 140
AAGUUCACCUUGAUGCCGUUCUUCUGCCUUGUCGGCCAUAGADAGACGDUUGDGGCCUGDUAGUUGUACU
150 160 170 180 190 200 210
CCAGCUUGUGCCCCAGGAUGUUGCCGUCUCCUUGAAGUCGAUGCCCUUCAGCUCGGAUGCCGGUUCACCAG
220 230 240
GGUGUCGGCCUCGAACUUCACCUCCGGCGCGI

EGFP Target 2

1 30 60 90 120 150 180 210 240 270 300
CGAGCUGGACGGCGACGUAAAACGGCCACAAGUUCAGCGUGUCCGGCGAGGGCCGAGGGCCGAUGCCACCUAC
30 40 50 60 70 80 90 100 110 120 130 140
GGCAAGCUGACCCUGAAGUUCaucugcaccaccgscagcugccccgugccccuggccccaccucggugacca
150 160 170 180 190 200 210 220 230 240 250 260
CCUGACCUACGGCGUGCAGUGCUUCAGCCGCUACCCCGACCACAUAGAAGCAGCAGCAGUUCUUCAGUC
270 280 290 300
CGCCAUSCCCSAAGSCUACGUCCAGGAGCCACCAUCUUCUUCUUAAGGACCACGGCAACUACAAGACCCCGC
310 320 330
GCCGAGGUGAAGUU]

FIG. 41

FIG. 42A

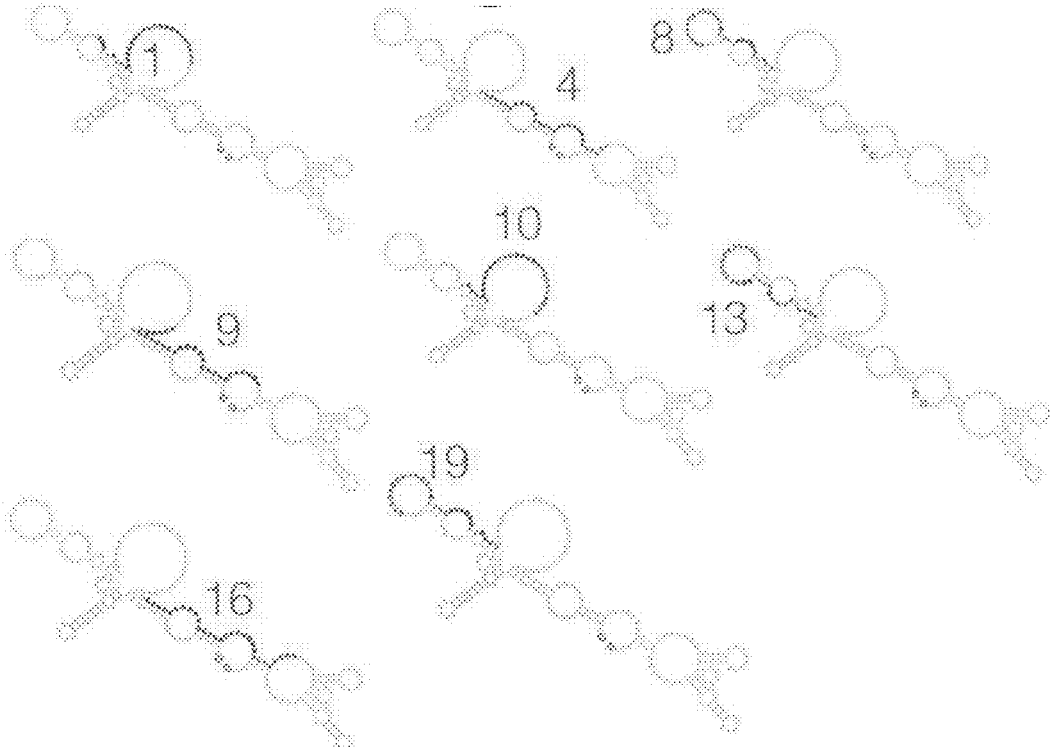


FIG. 42B

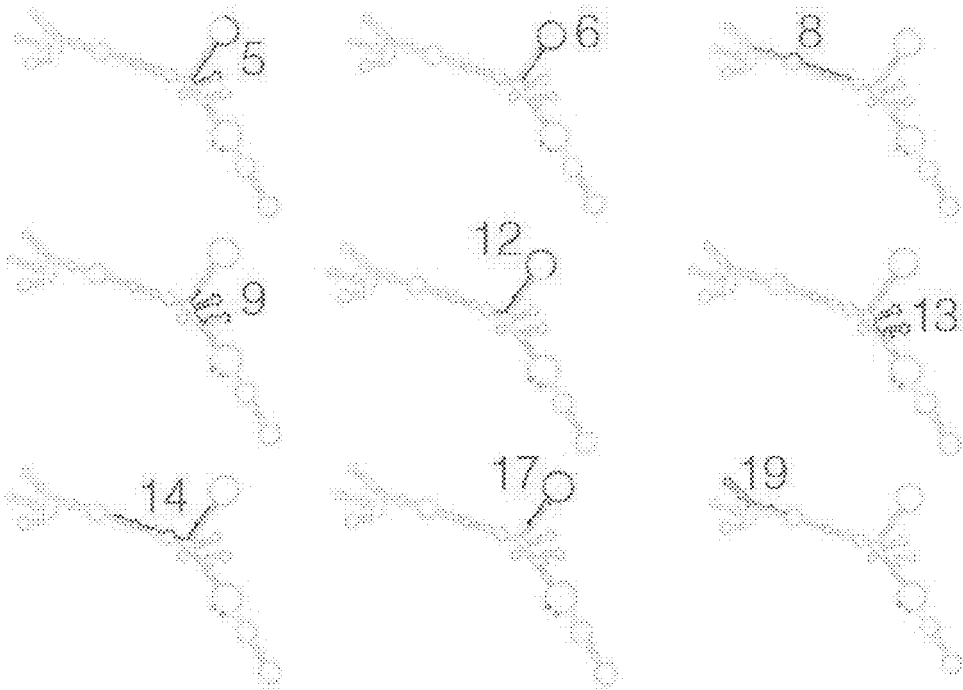


FIG. 44

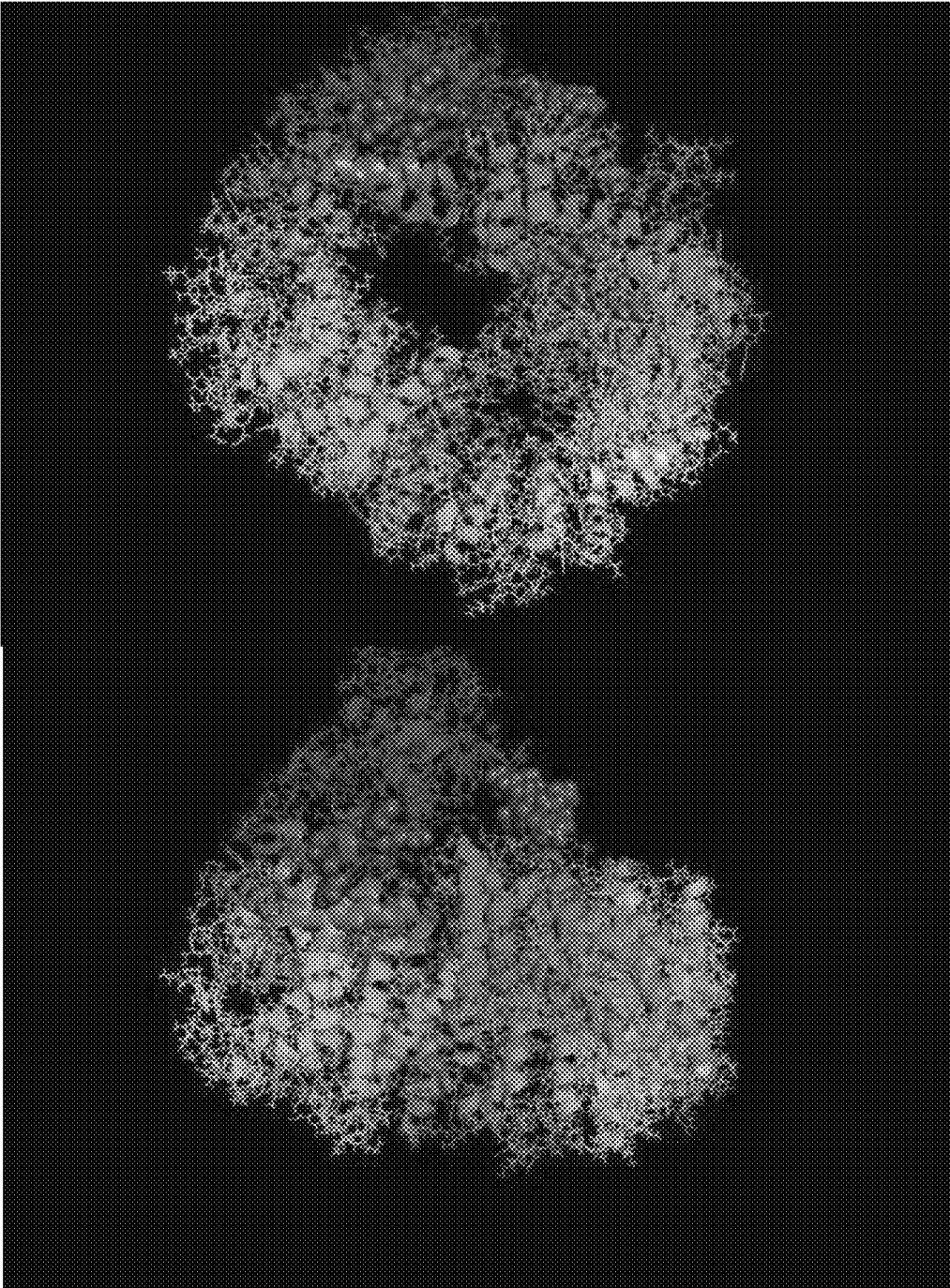


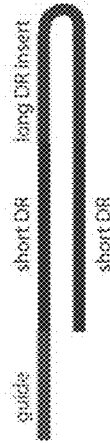
Fig. 45A

Different crRNAs

1) short crRNA (native)



2) long crRNA (native)



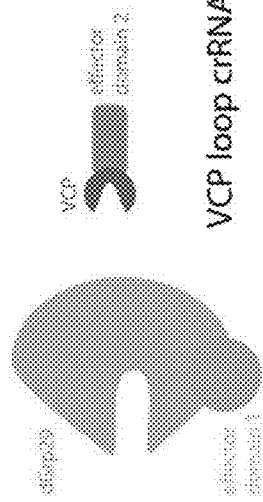
3) VCP loop crRNA (engineered)



Fig. 45B

Grp29 SAM system

Grp29 SAM system



Grp29 SAM system bound to RNA

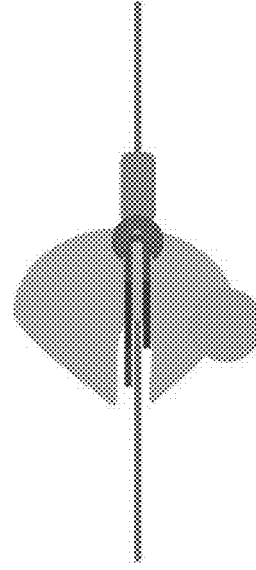


Fig. 46A

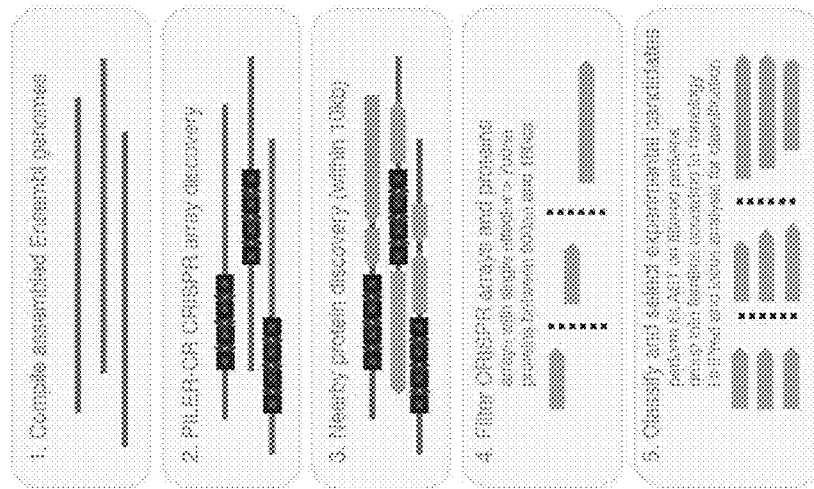


Fig. 46B

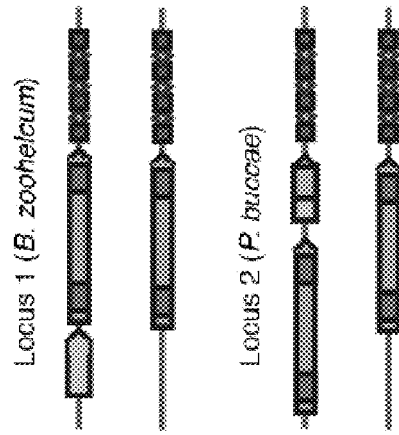


Fig. 46C

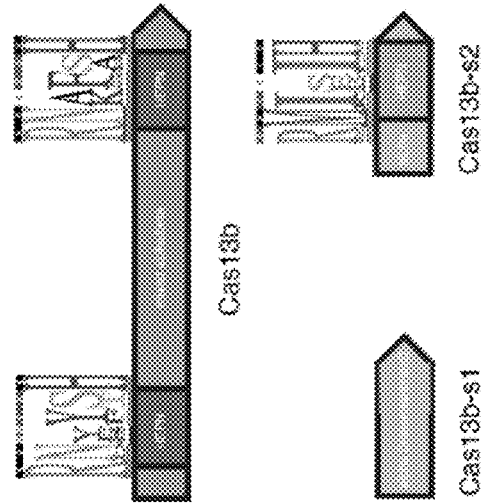


Fig. 47

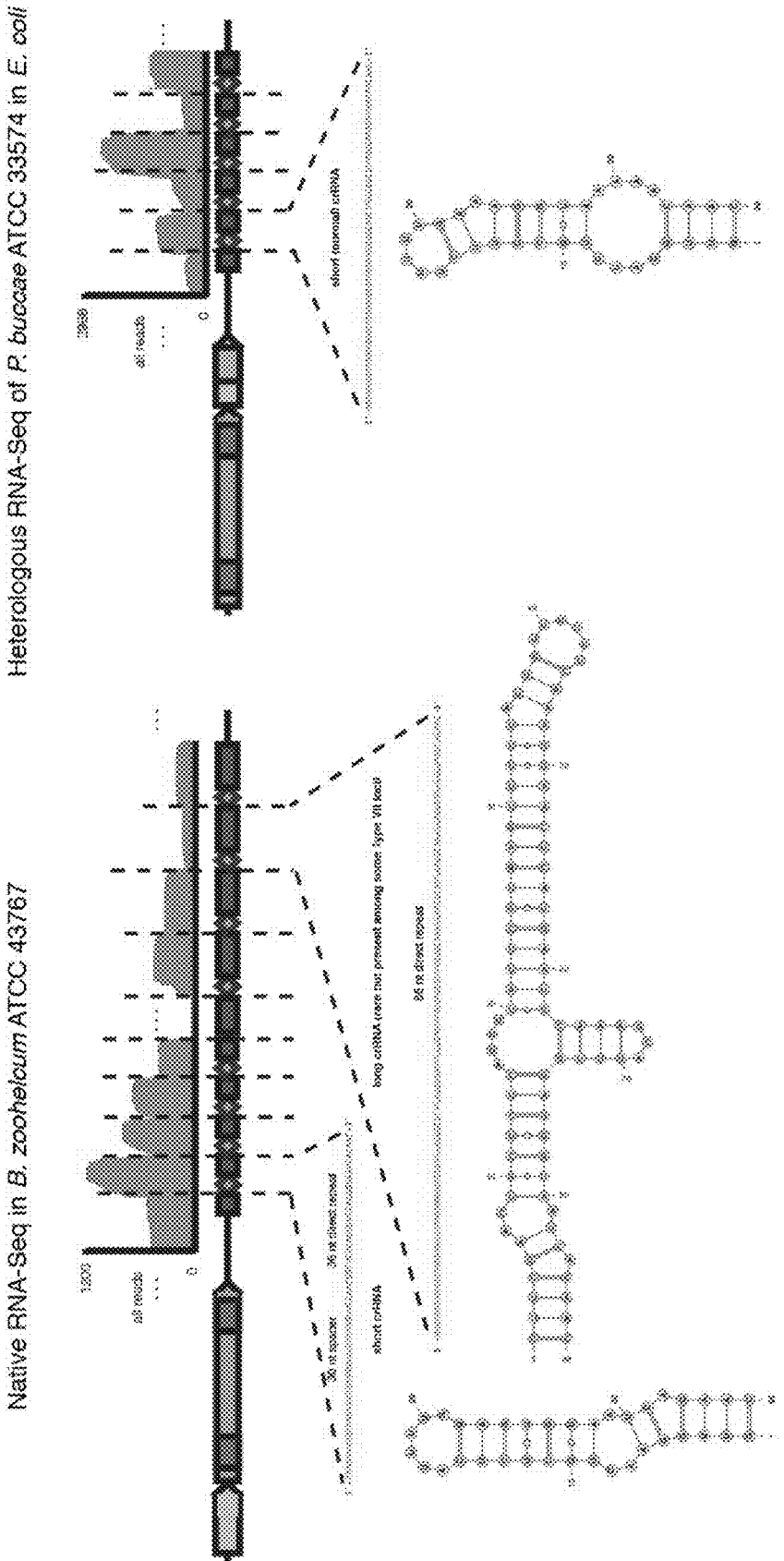


Fig. 48

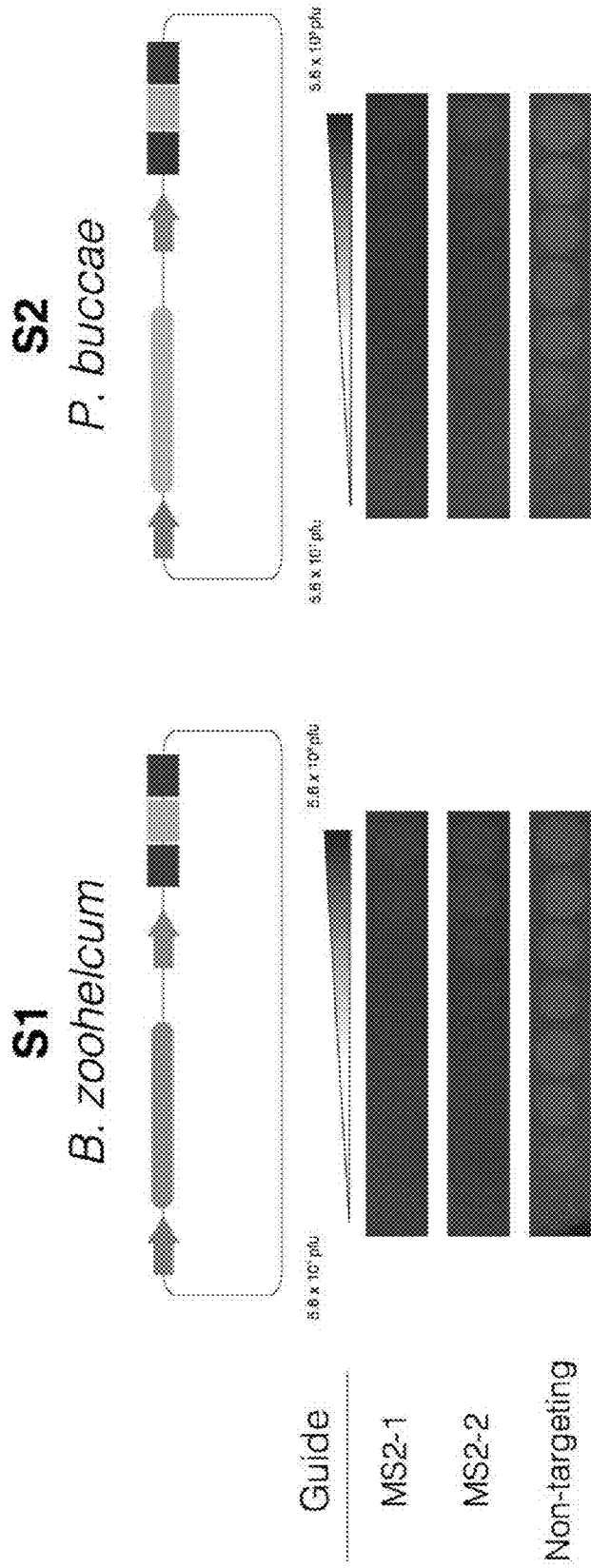


Fig. 49

targeted RNA

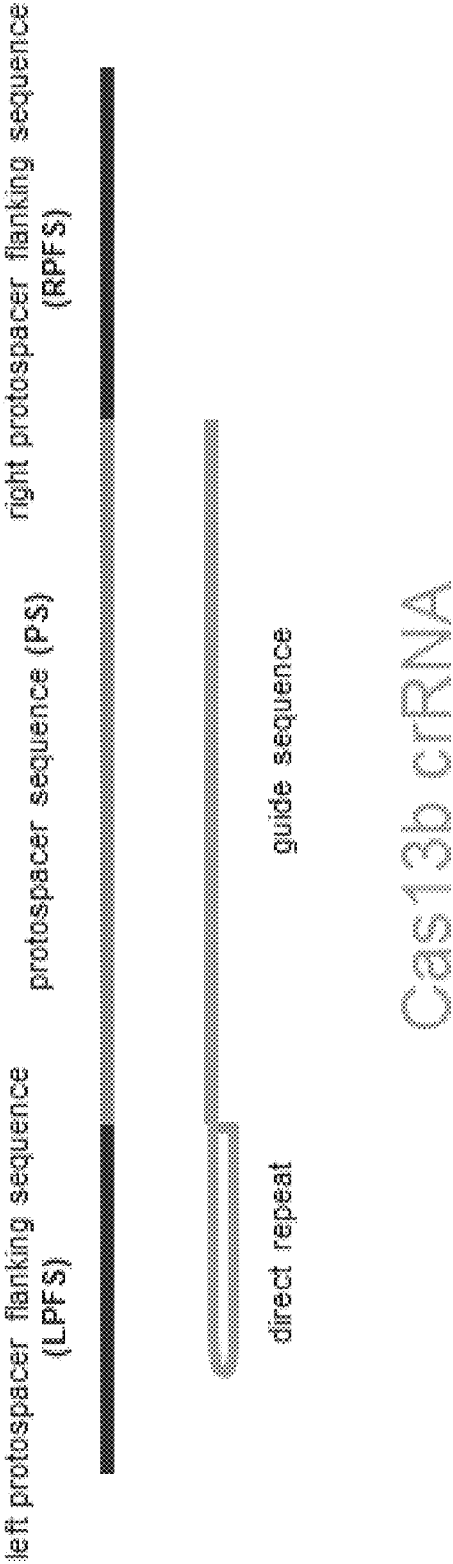
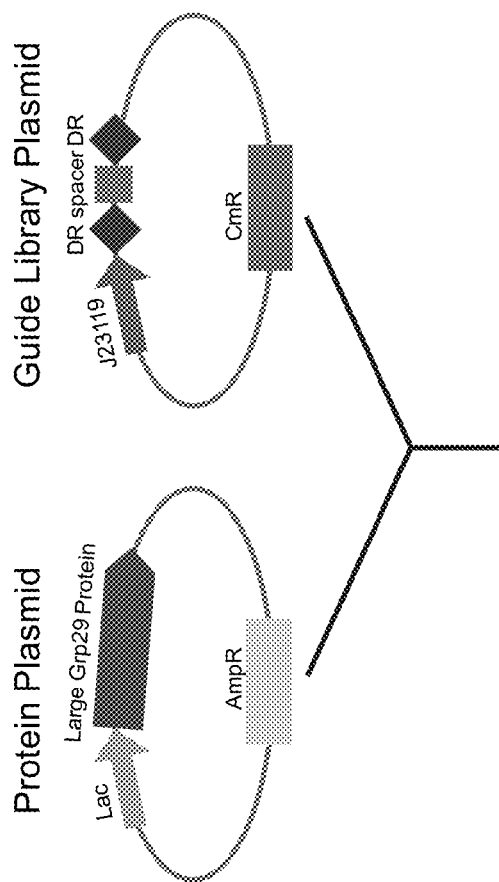


Fig. 50



Co-electroporate at 50/50 ratio into electrocompetent MegaX DH10B T1R line (>3e10 cfu/ug)

Outgrow for 1 hour at 37C, then incubate overnight on Chloramphenicol & Carbenicillin agar plates

Collect grown colonies 11 hours post-transformation, maxiprep, NGS library prep, and load on NextSeq

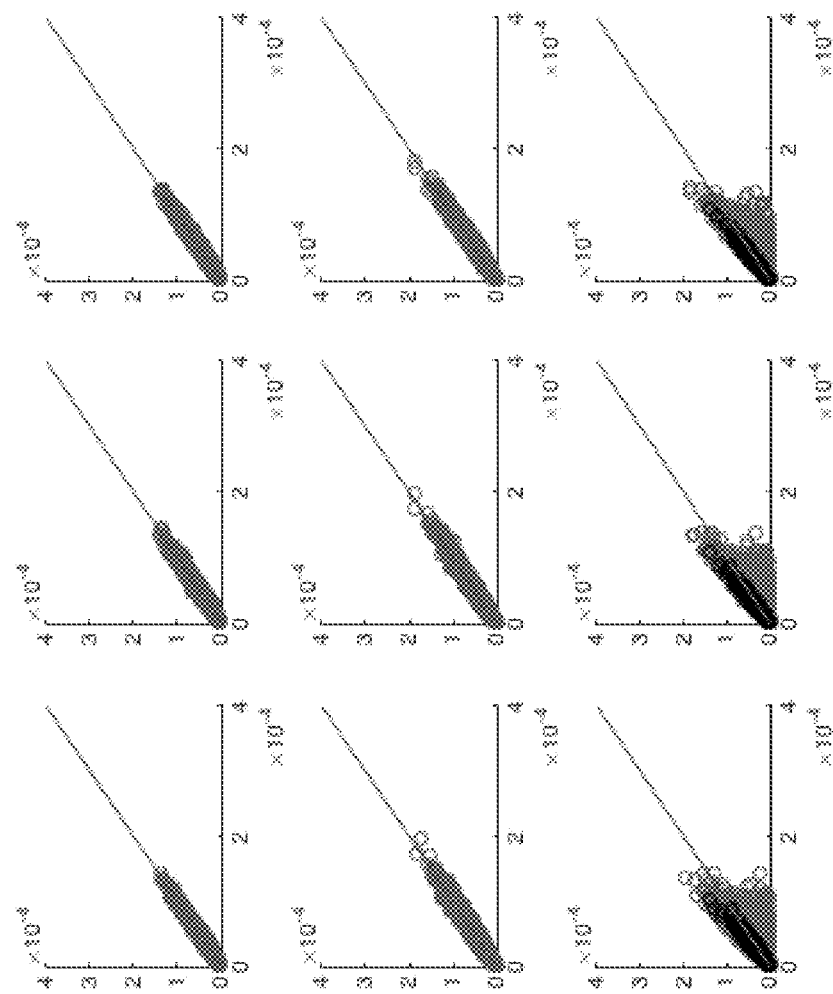


Fig. 51

3 controls (-BZ)
vs. each other
(1v2, 1v3, 2v3)

3 experimentals
vs. each other
(1v2, 1v3, 2v3)

each experimental
vs. its control

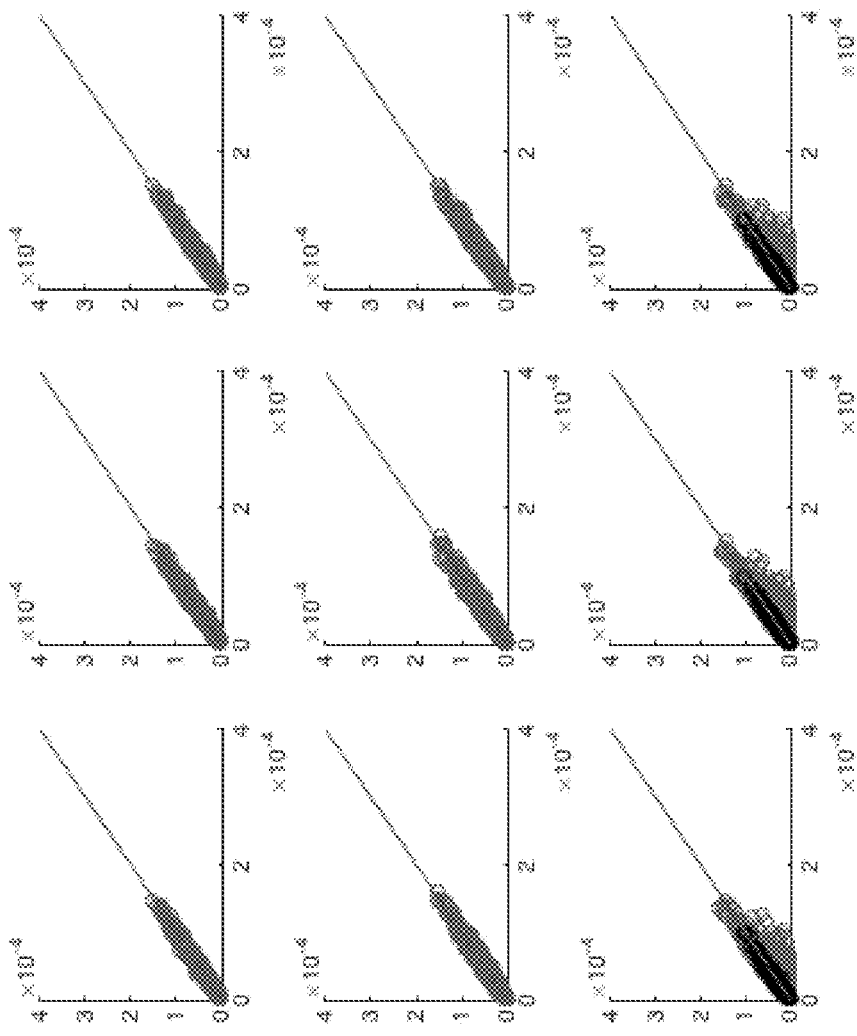


Fig. 52

3 controls (-PB)
vs. each other
(1v2, 1v3, 2v3)

3 experimentals
vs. each other
(1v2, 1v3, 2v3)

each experimental
vs. its control

Fig. 53
BZ Screen

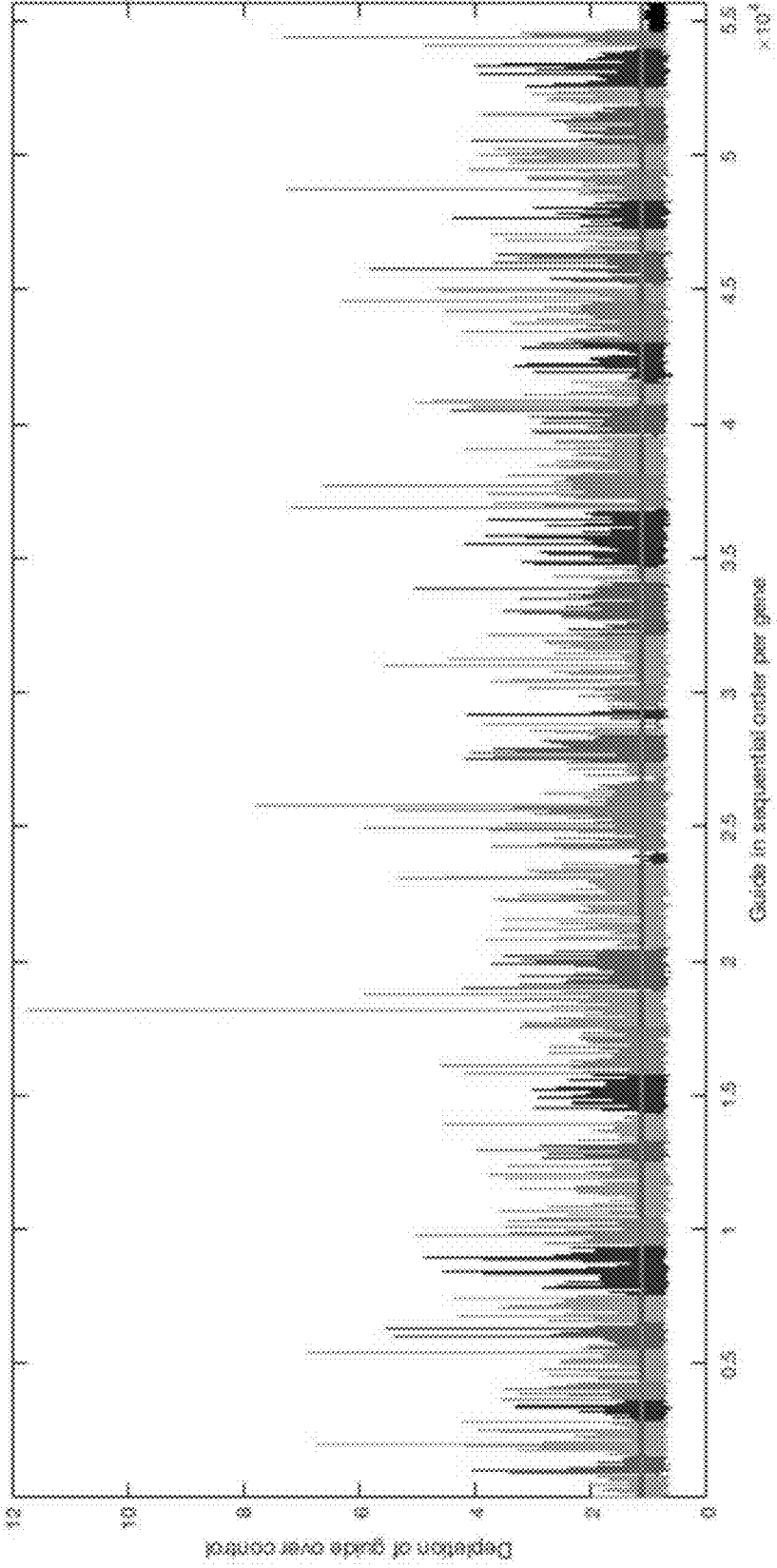


Fig. 54
PB Screen

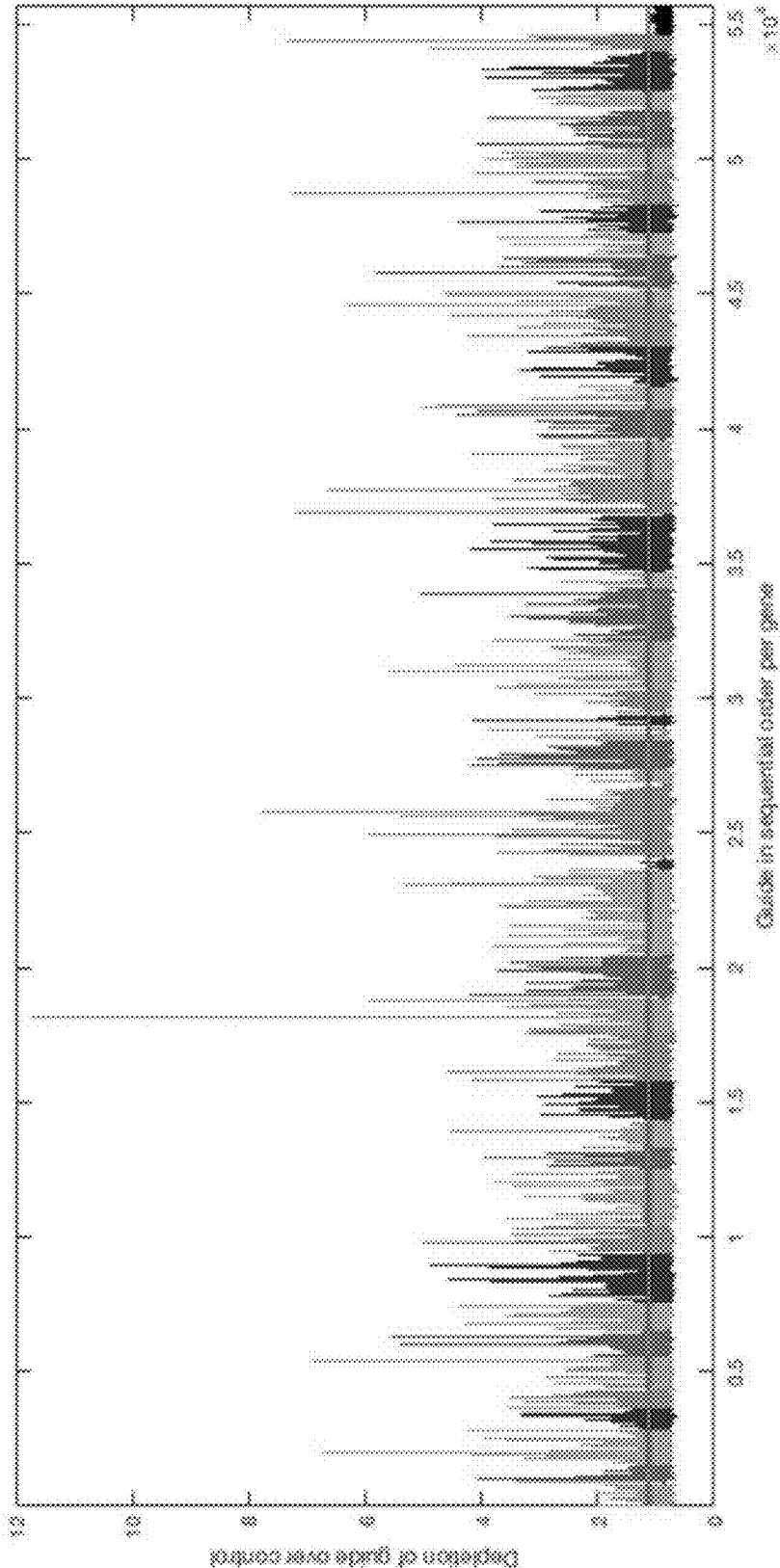


Fig. 55

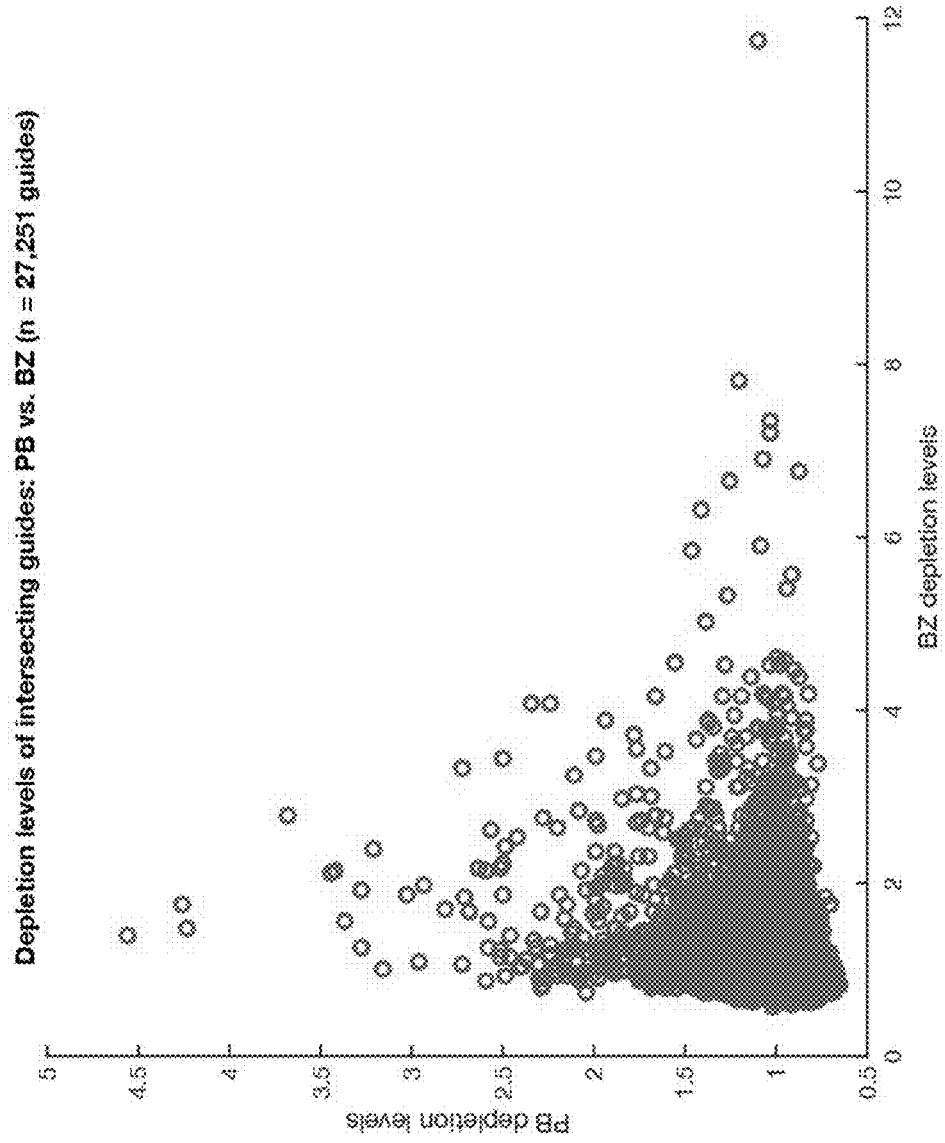


Fig. 56
BZ Aggregate depletions

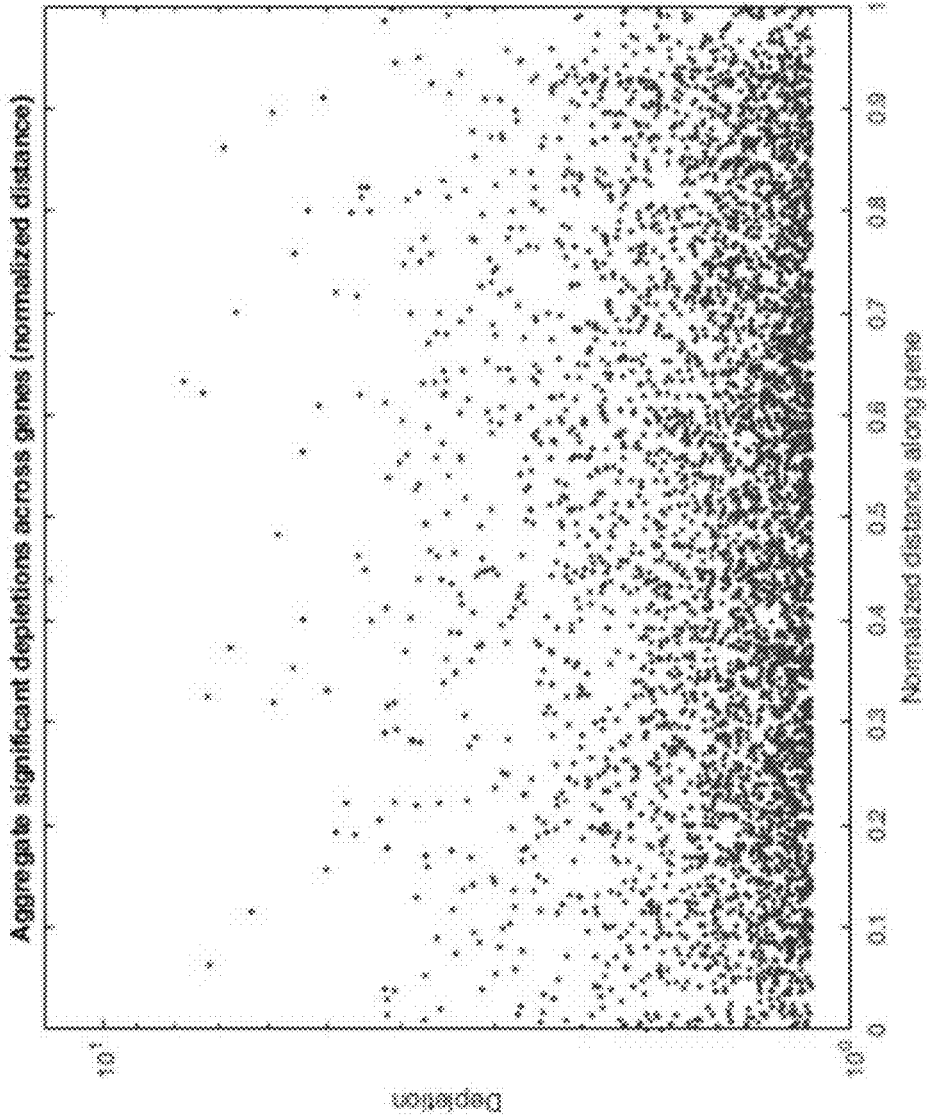


Fig. 57
PB Aggregate Depletions

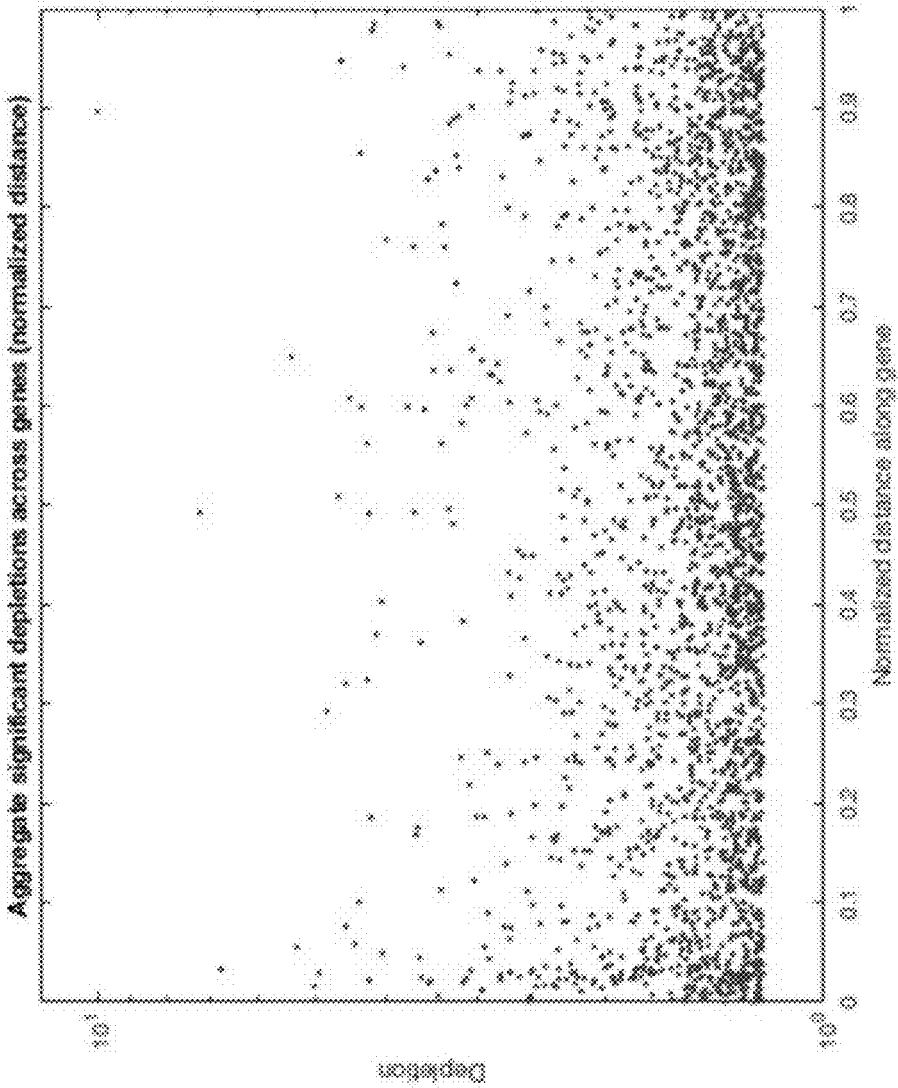


Fig. 58

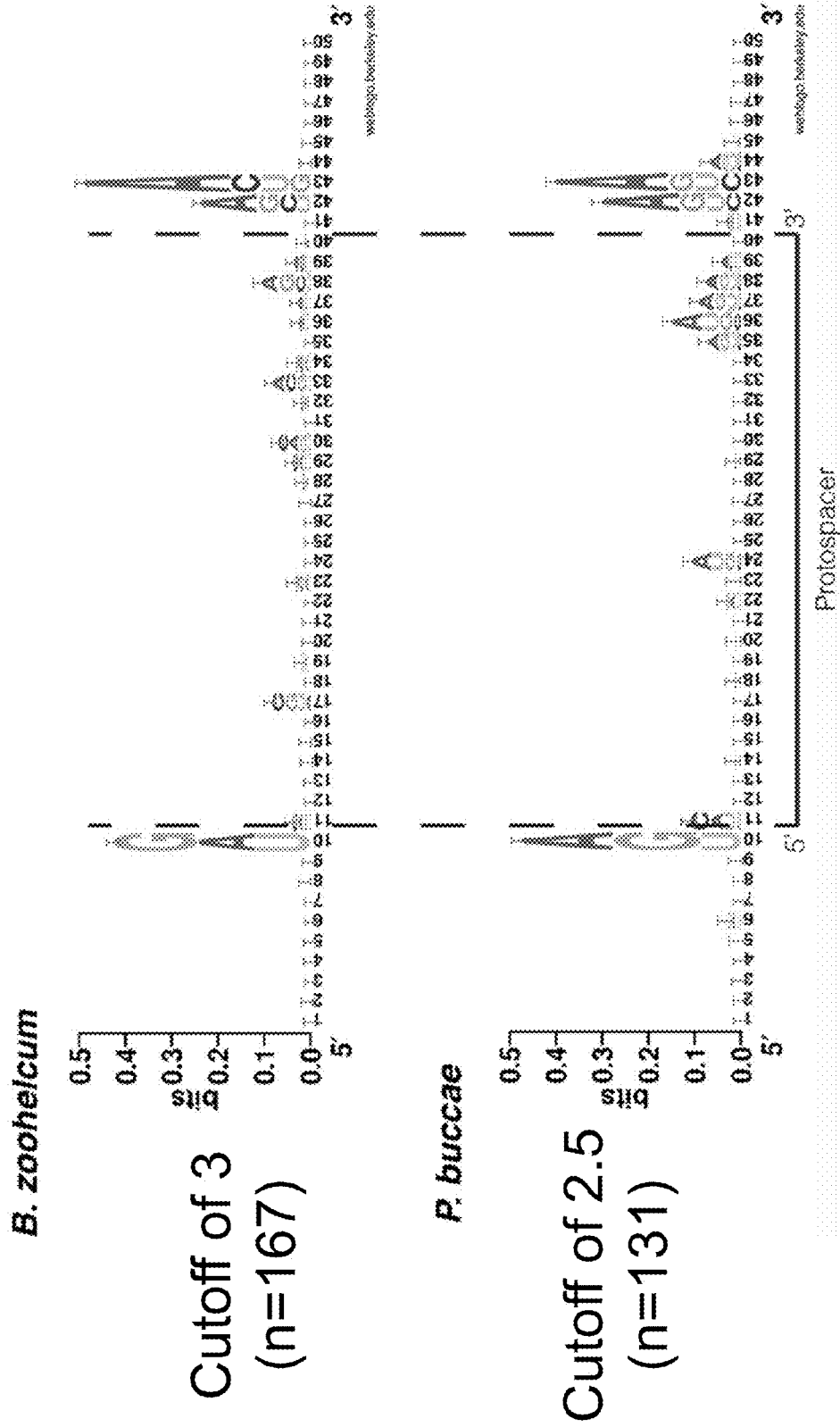
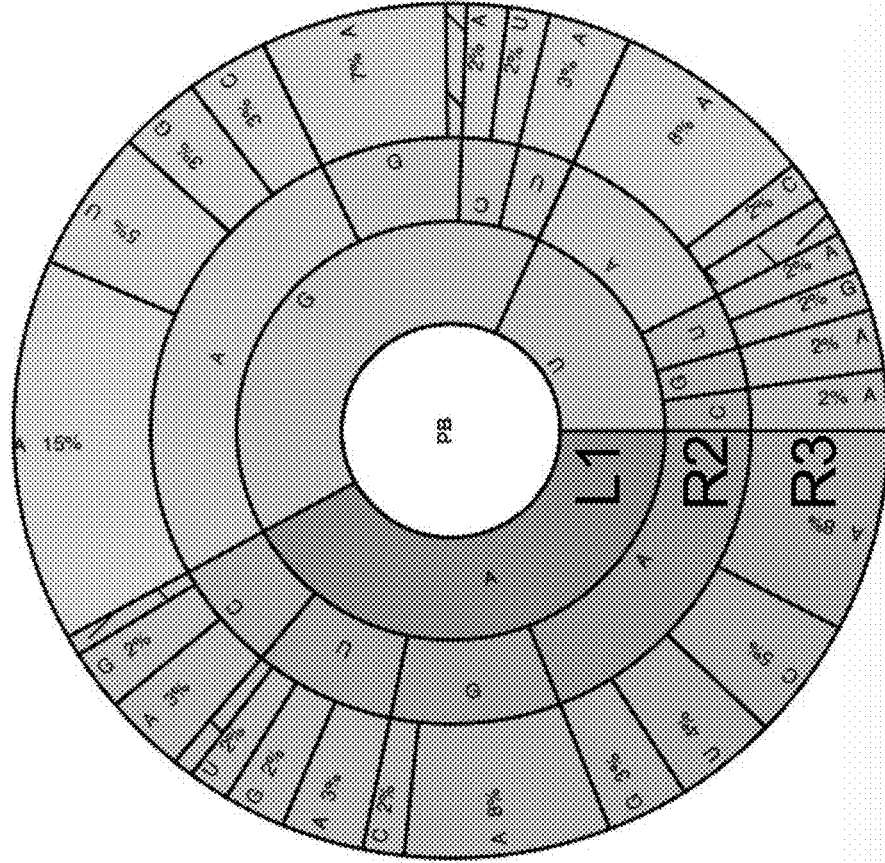
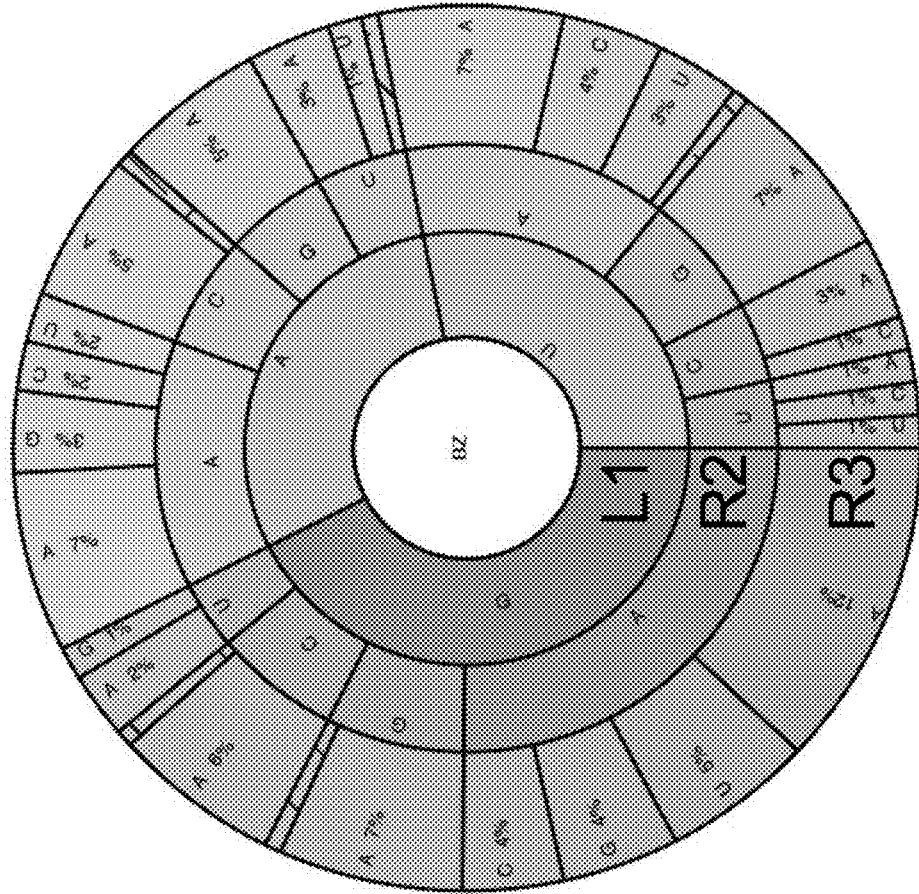


Fig. 59
 BZ PFS Wheel



Cutoff of 3
 (n=167)

Fig. 60
PB PFS Wheel



Cutoff of 2.5
(n=131)

Fig. 61

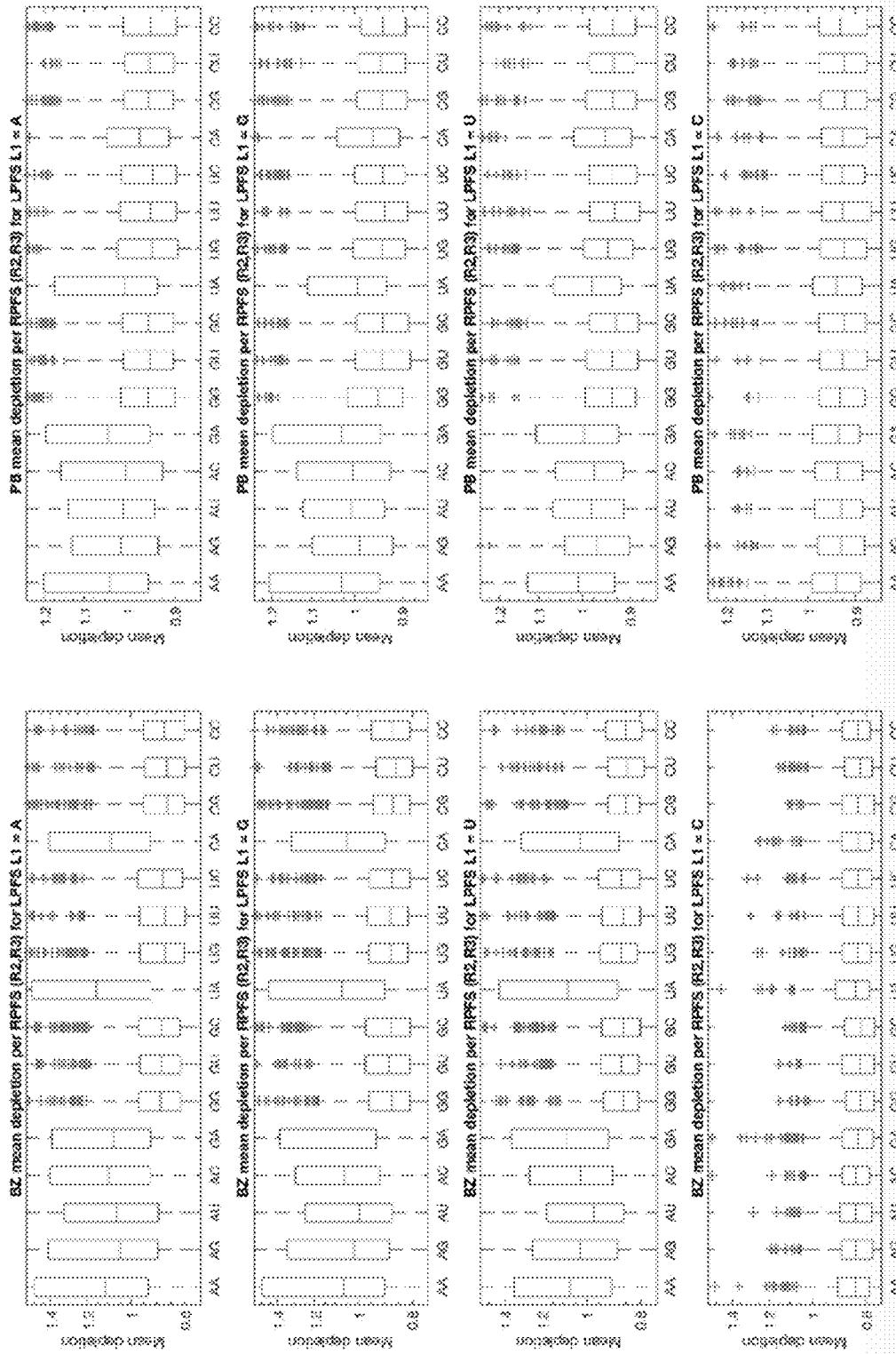


Fig. 62

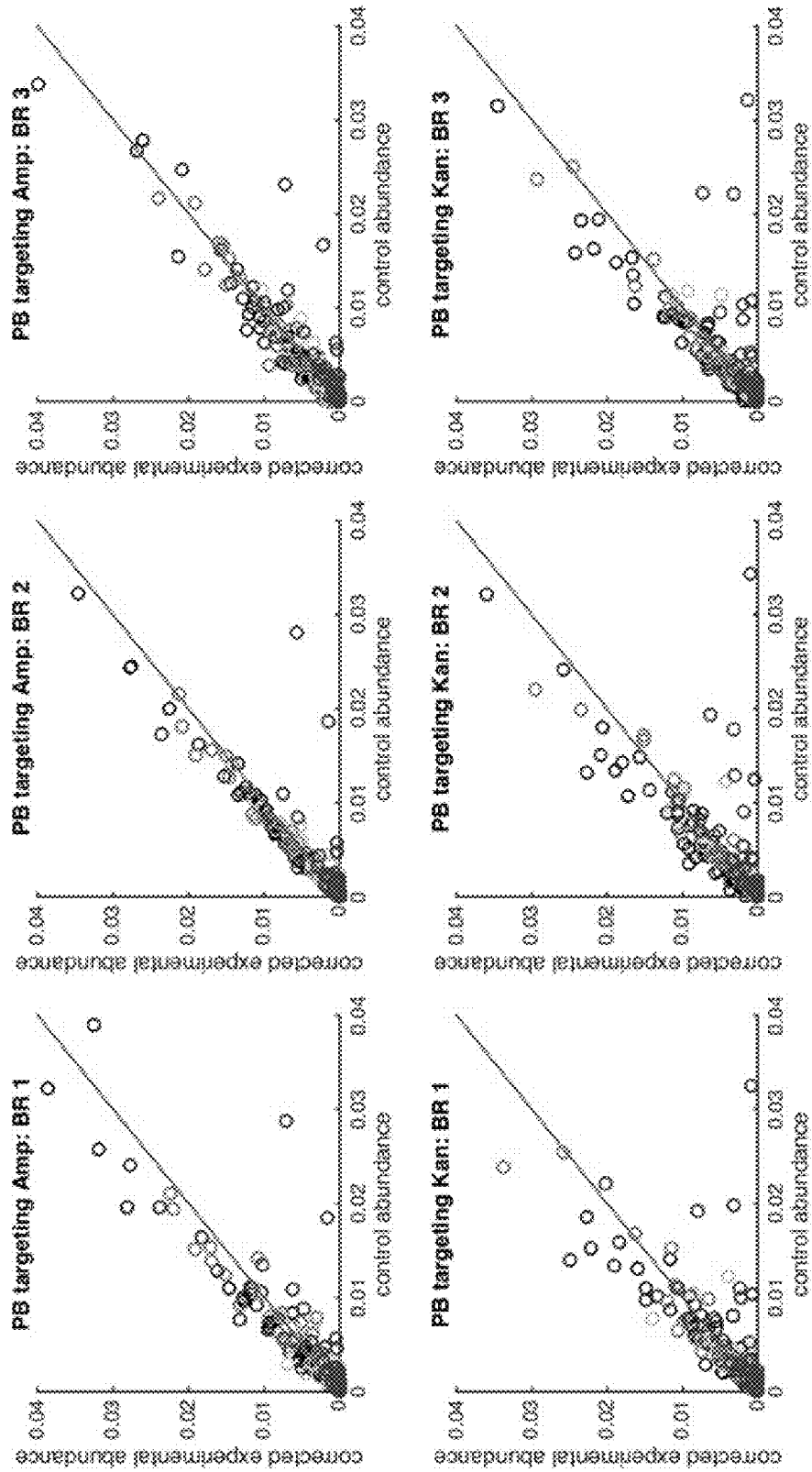


Fig. 63

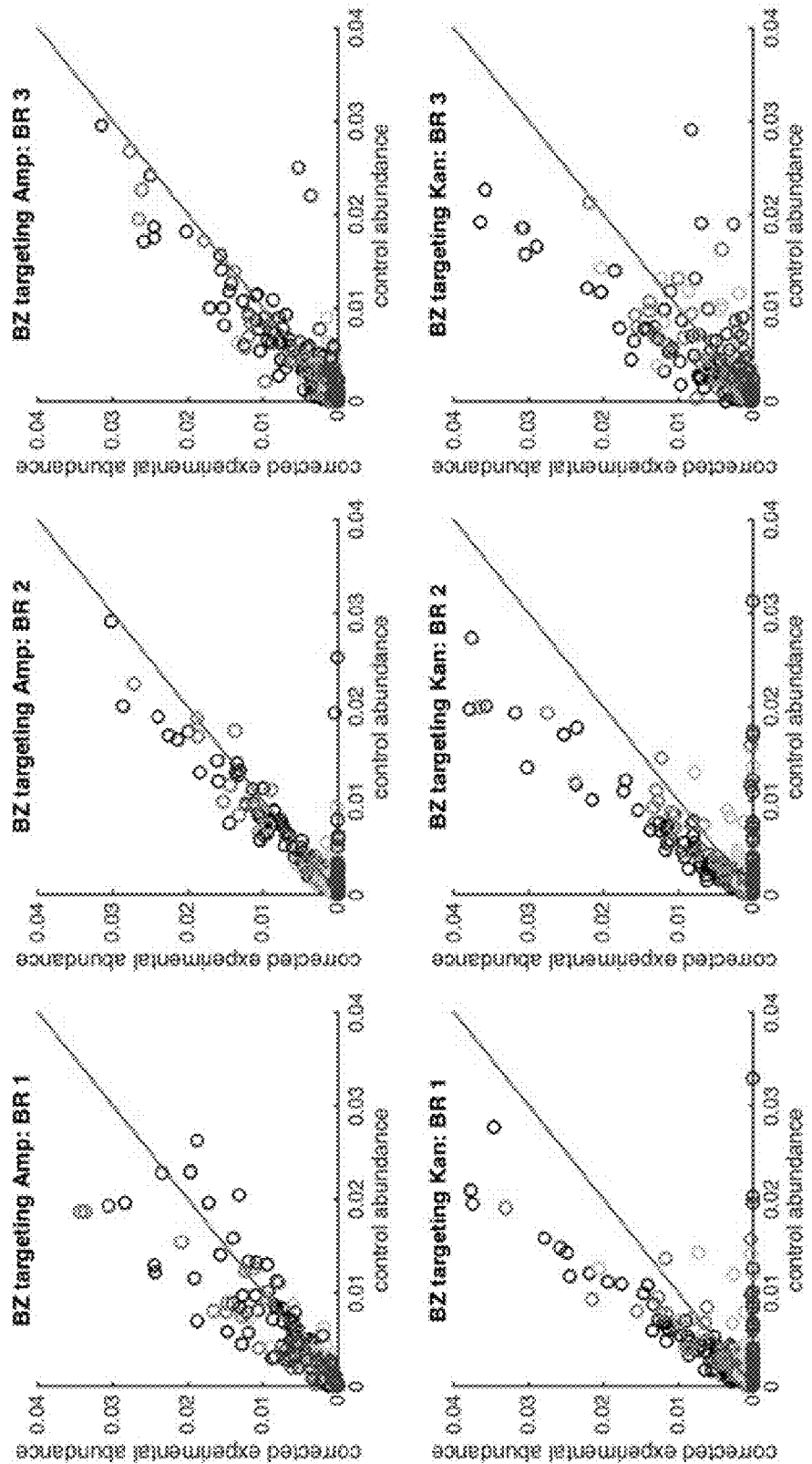


Fig. 64A

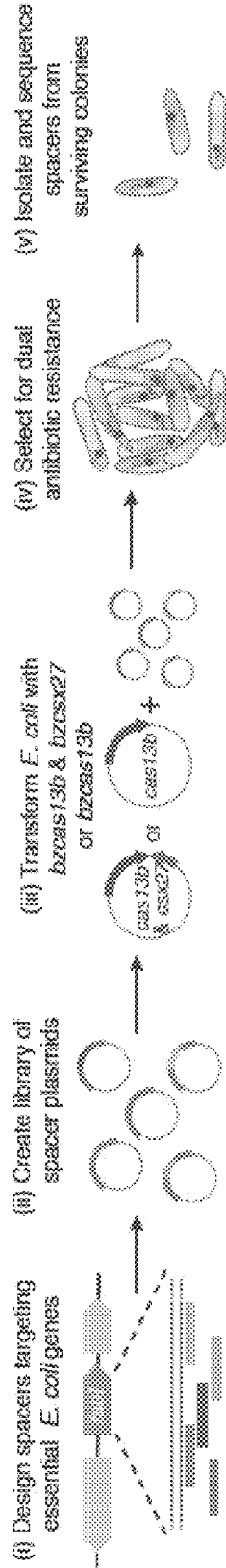


Fig. 64B

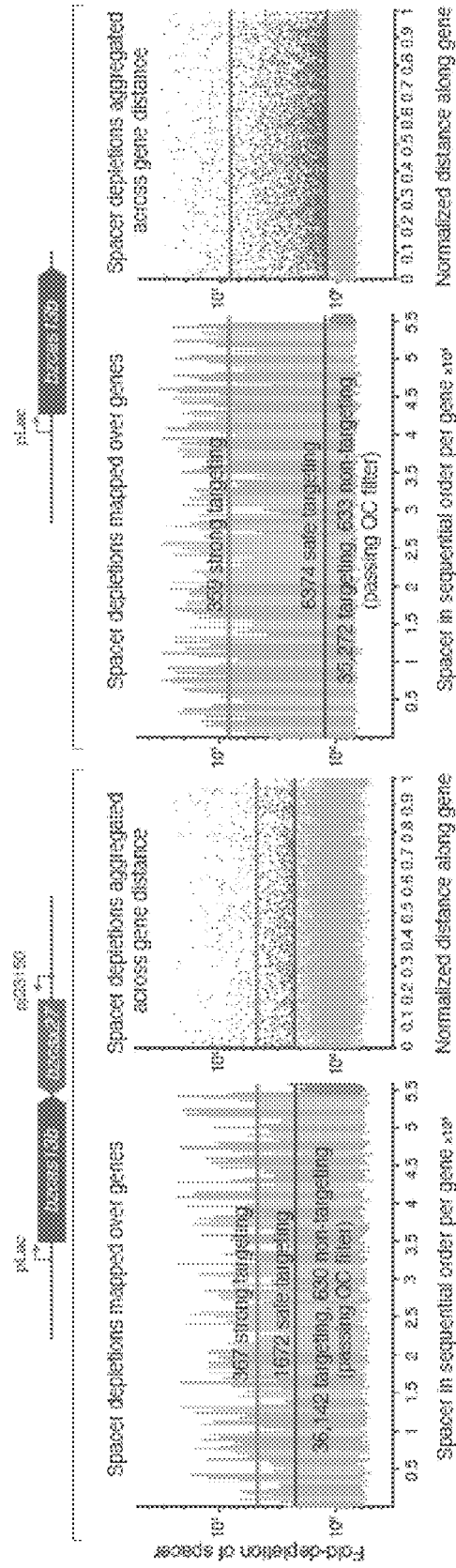


Fig. 64C

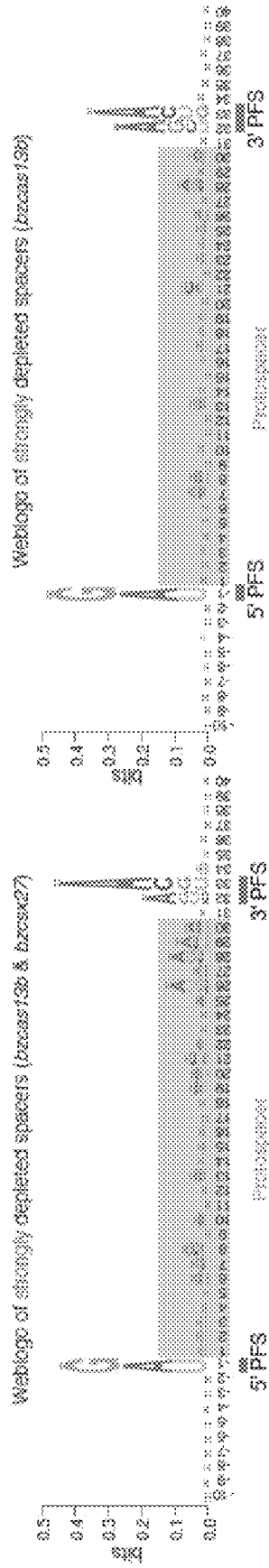


Fig. 64D

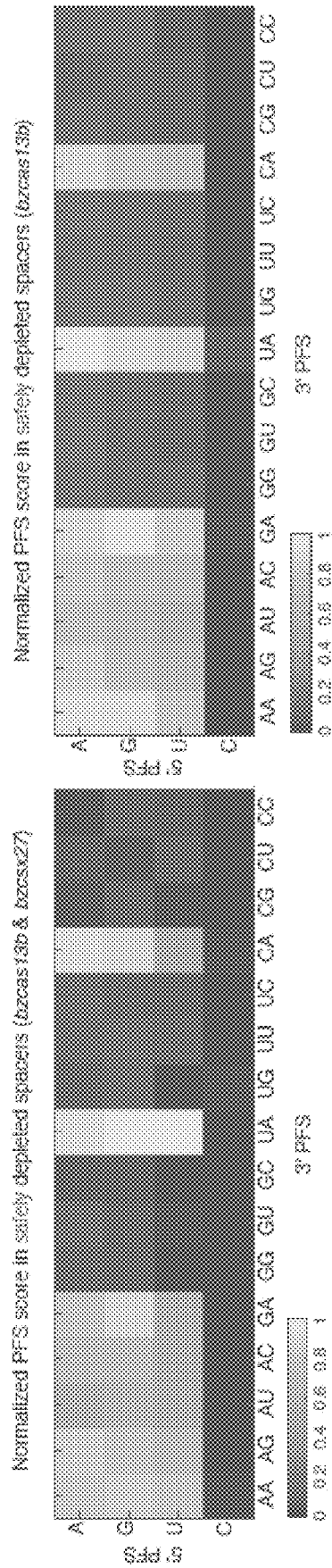


Fig. 64E

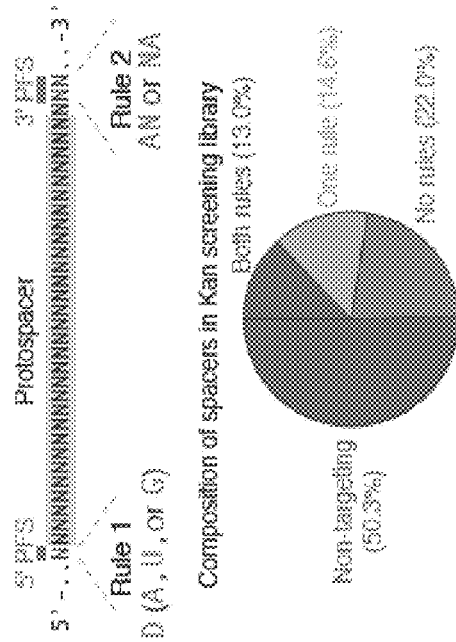


Fig. 64F

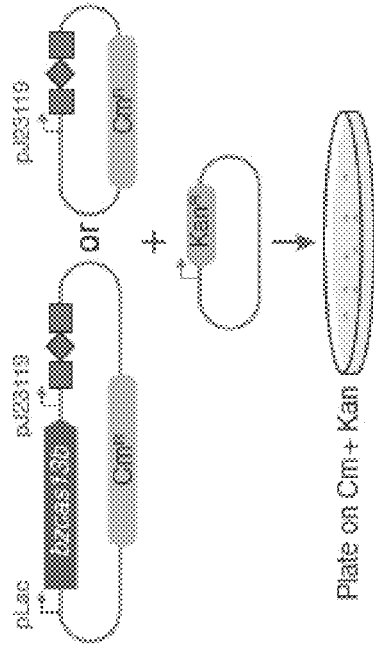


Fig. 64G

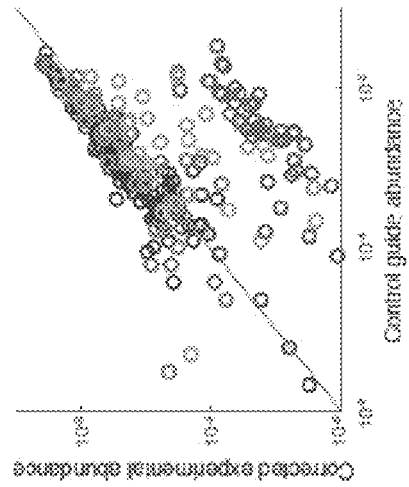


Fig. 65D

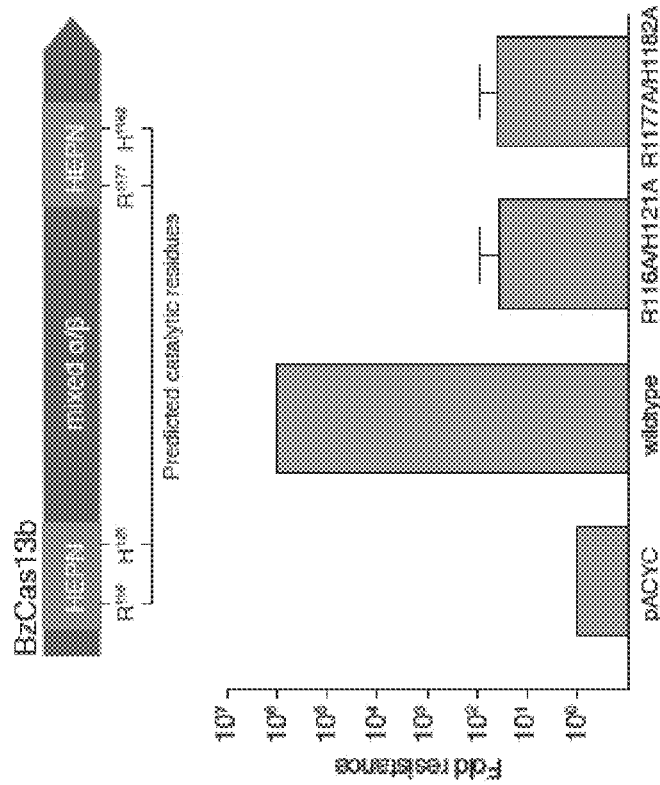


Fig. 65C

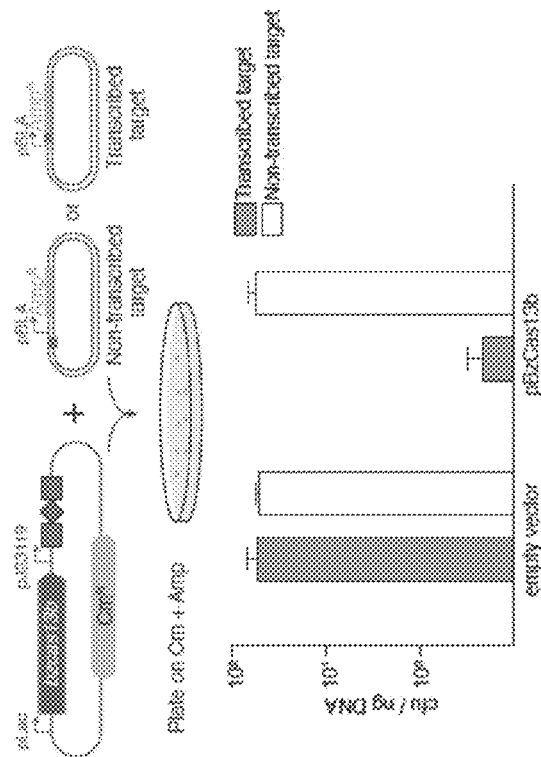


Fig. 66A

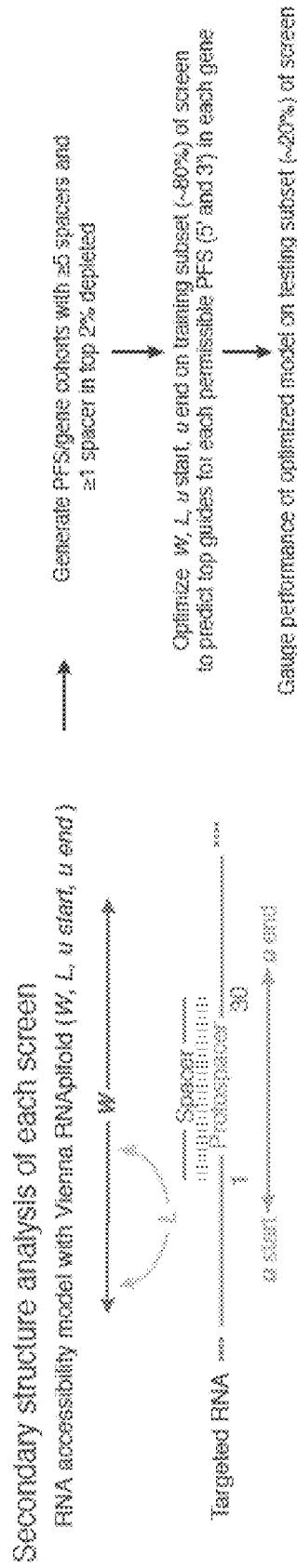
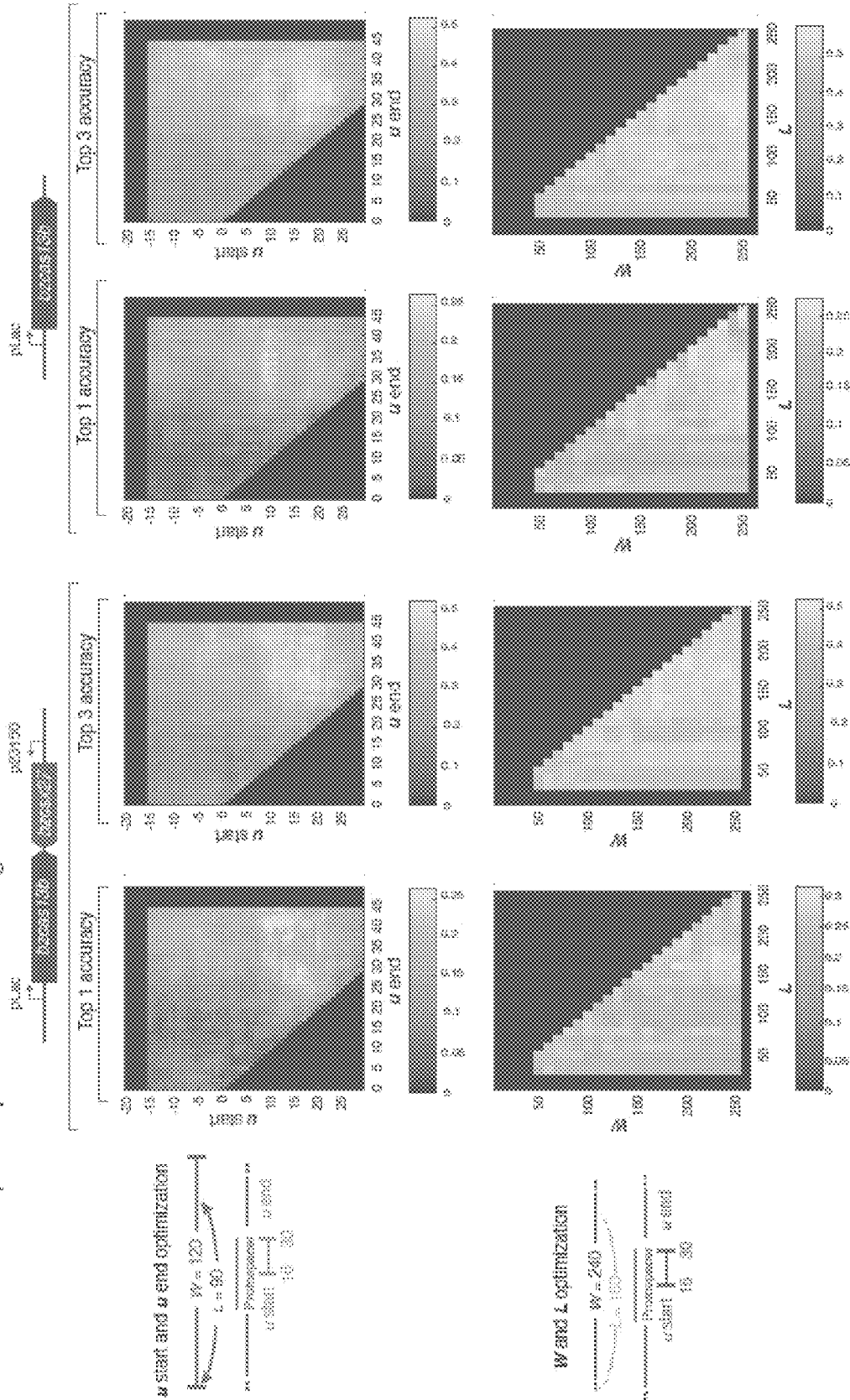


Fig. 66B

B. zooheicum spacer optimization on training set



W and L optimization
 W = 120
 L = 80
 Placeholder
 gradient 10 and 10 80

W and L optimization
 W = 240
 L = 160
 Placeholder
 gradient 10 and 10 80

Fig. 66C

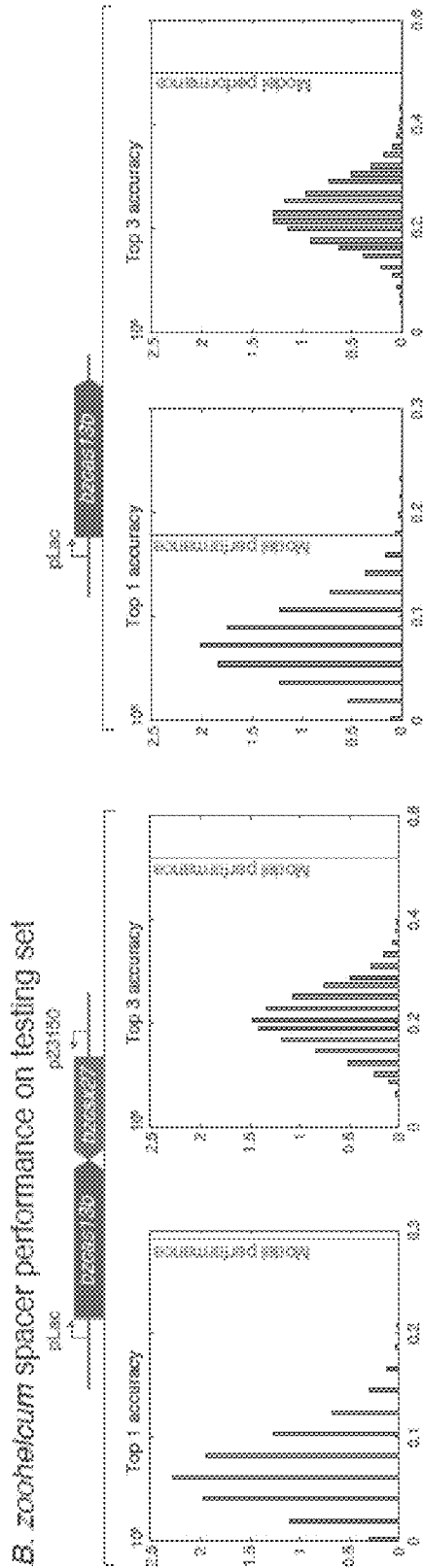


Fig. 66D

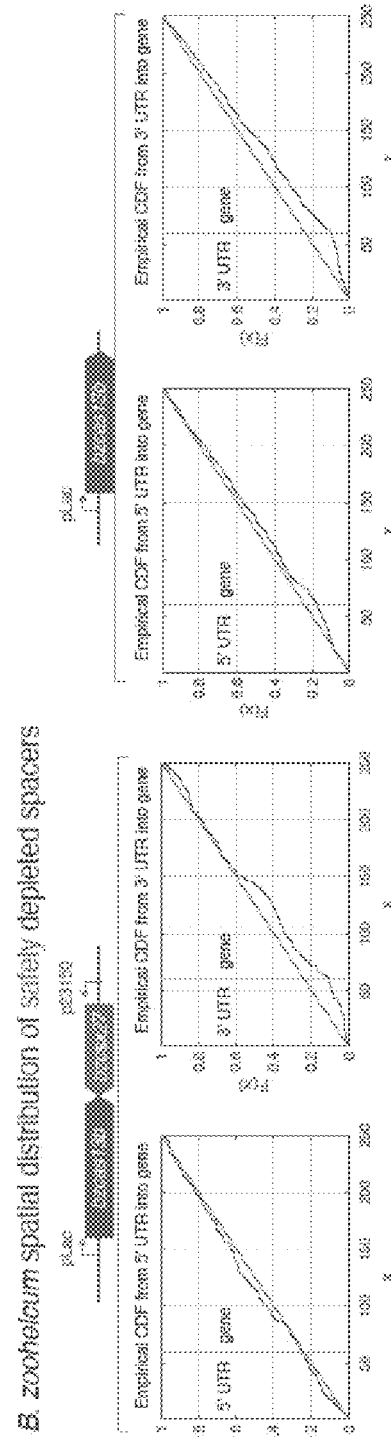


Fig. 67A

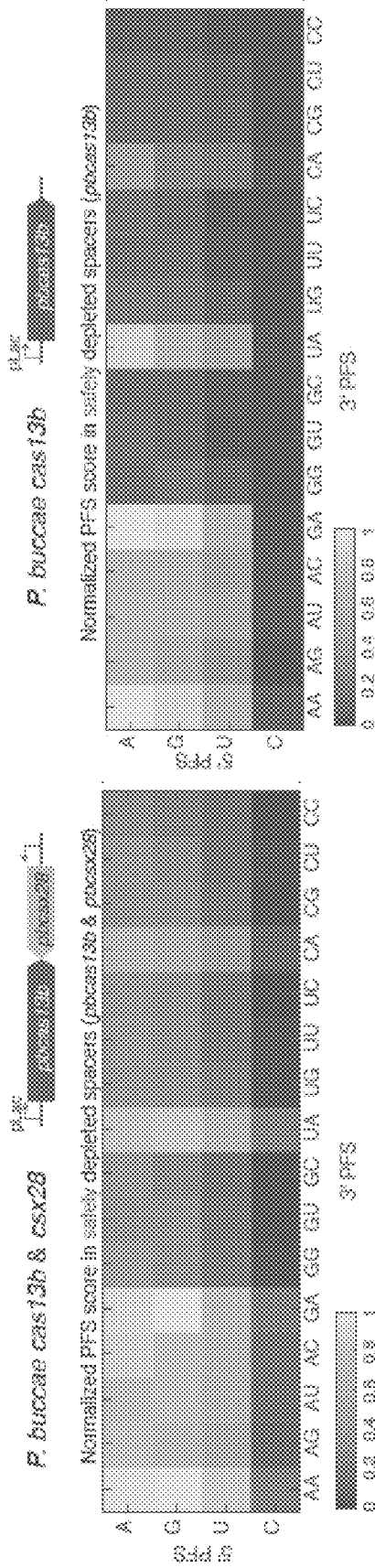


Fig. 67B

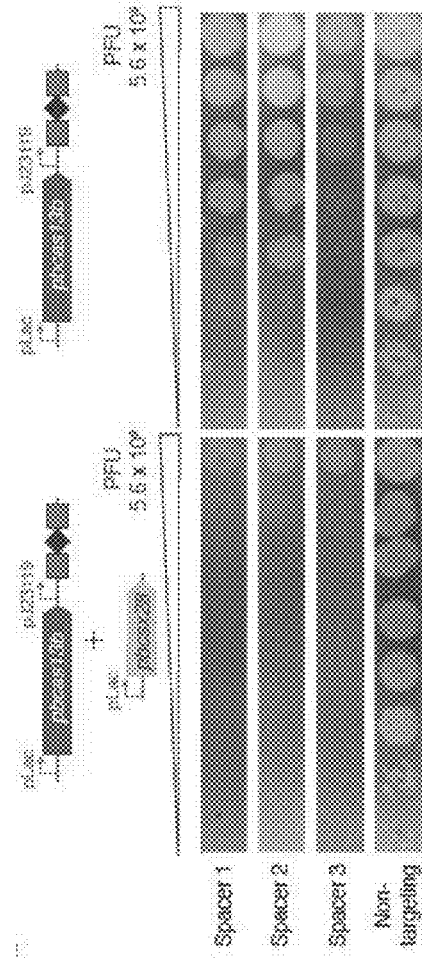


Fig. 68

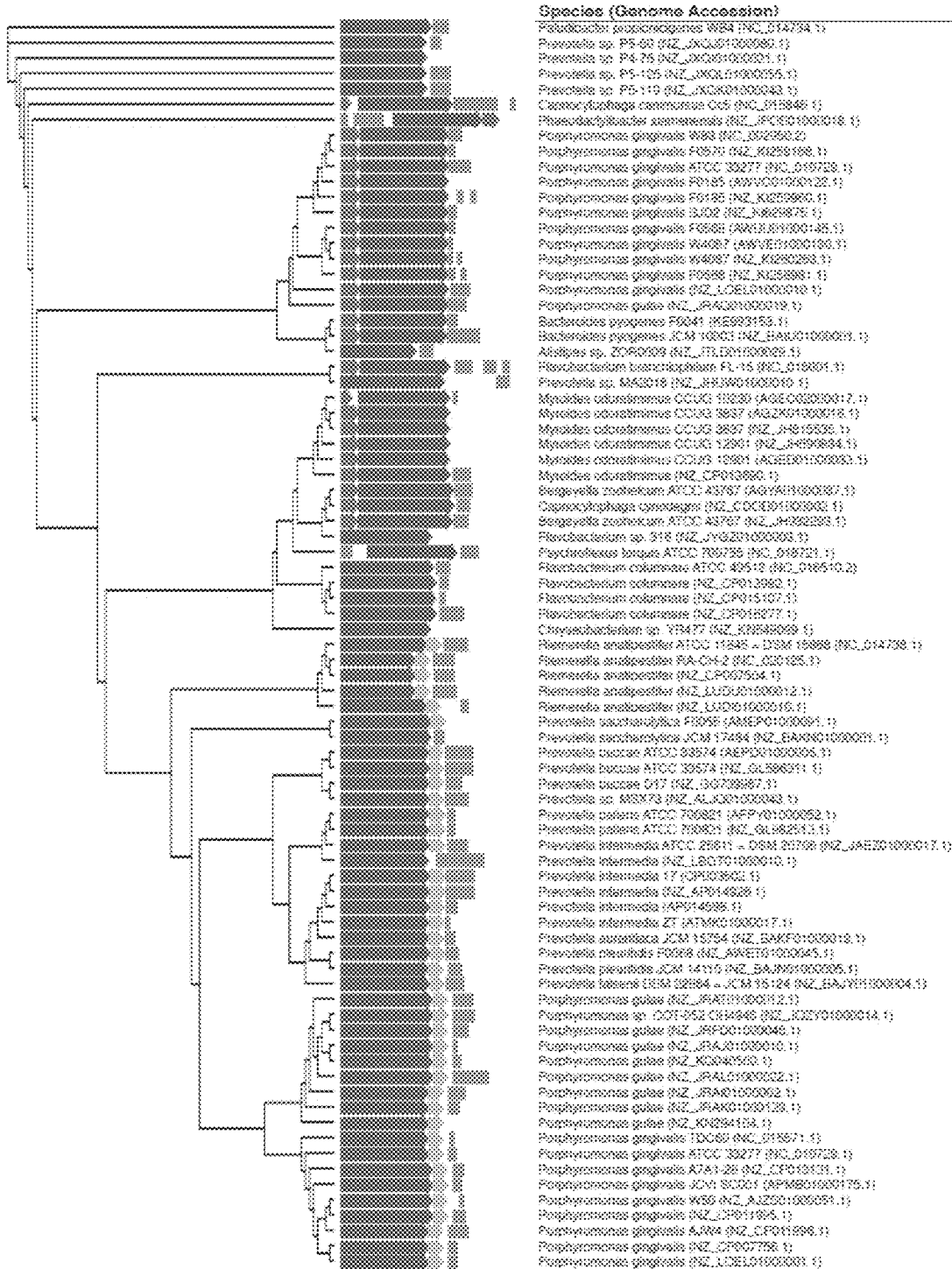


Fig. 69A
 CRISPR Class 2 Type VI-B Locus 1 direct repeat RNA folds - 36nt*

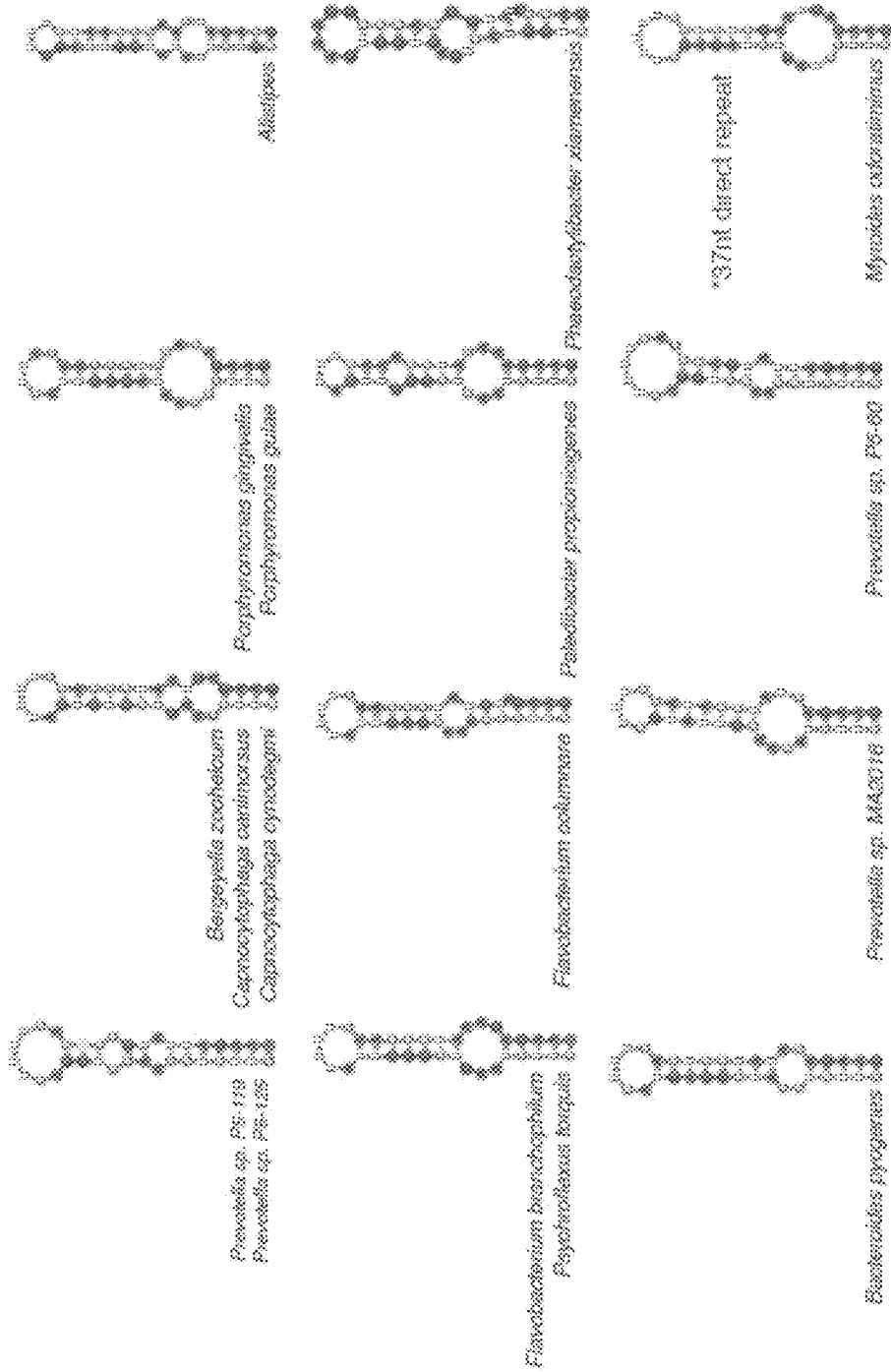


Fig. 69B

CRISPR Class 2 Type VI-B Locus 2 direct repeat RNA folds - 36nt

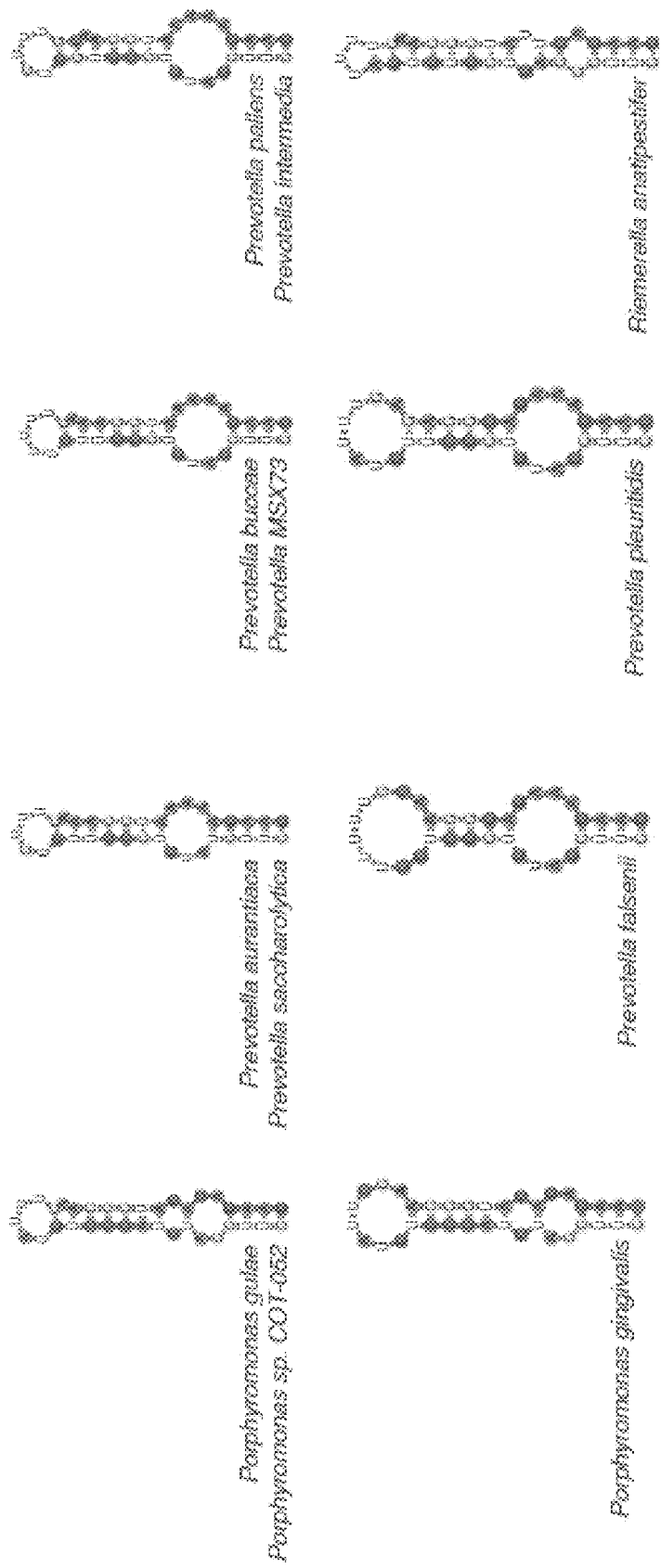


Fig. 70A

Conservation of 36nt direct repeats (24 unique)

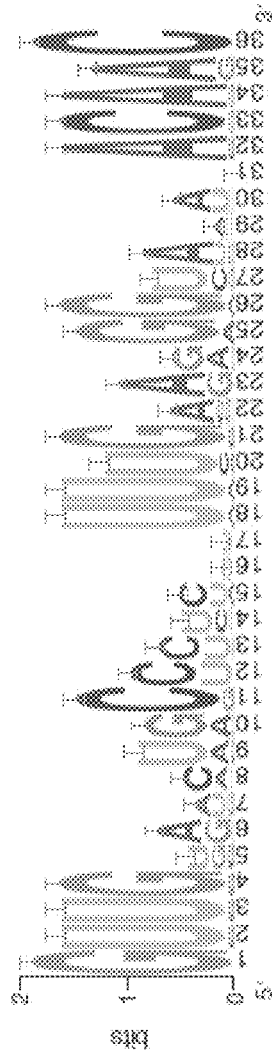


Fig. 70B

Flanking regions of mapped protospacers (114 unique)

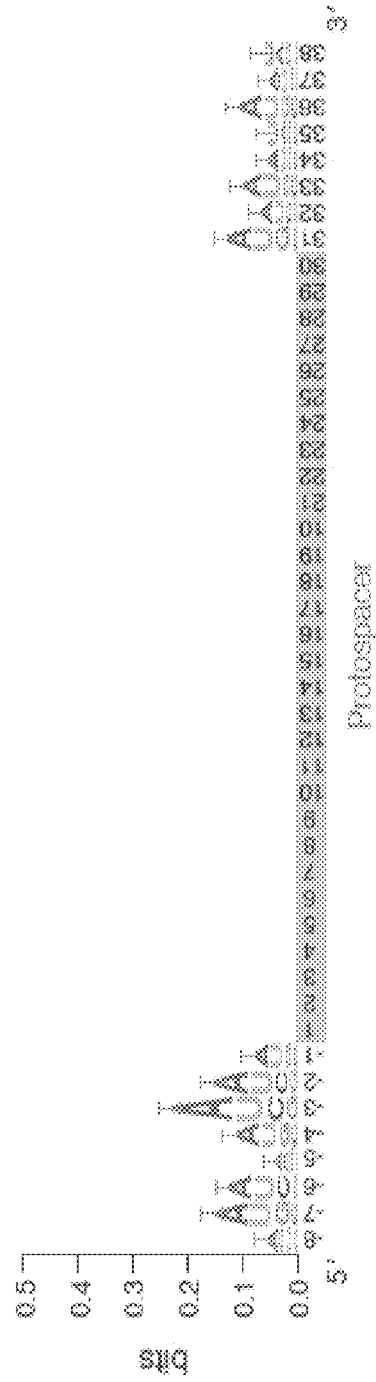


Fig. 71A

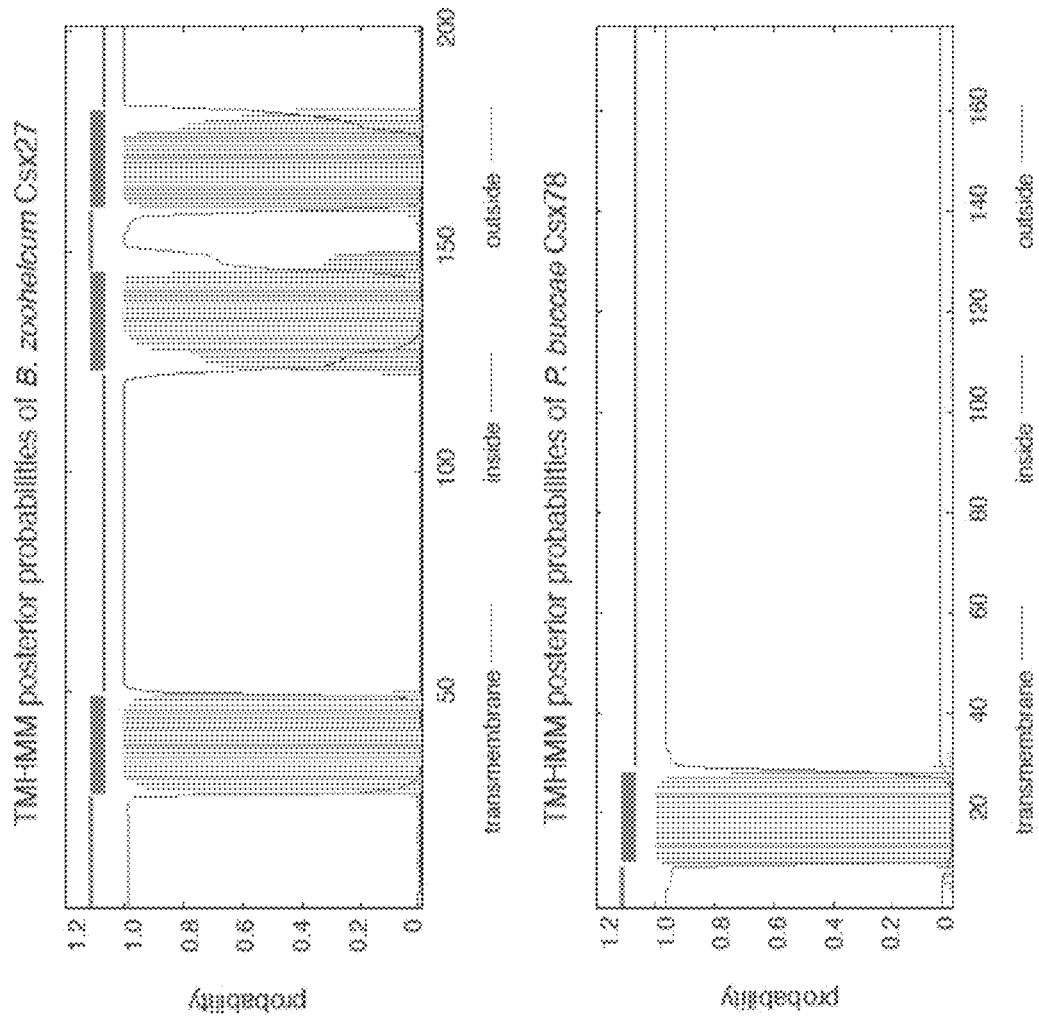


Fig. 71B

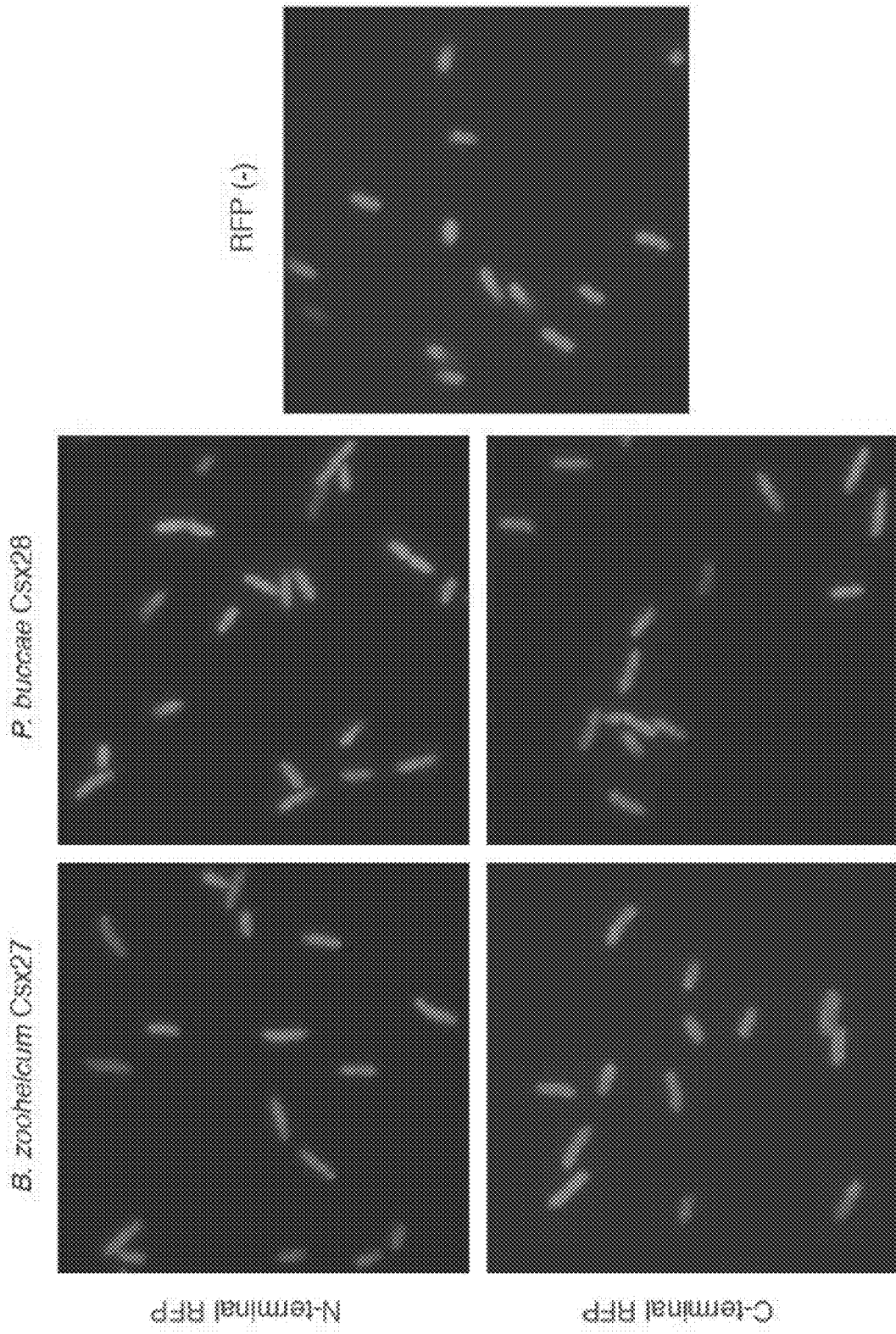


Fig. 73

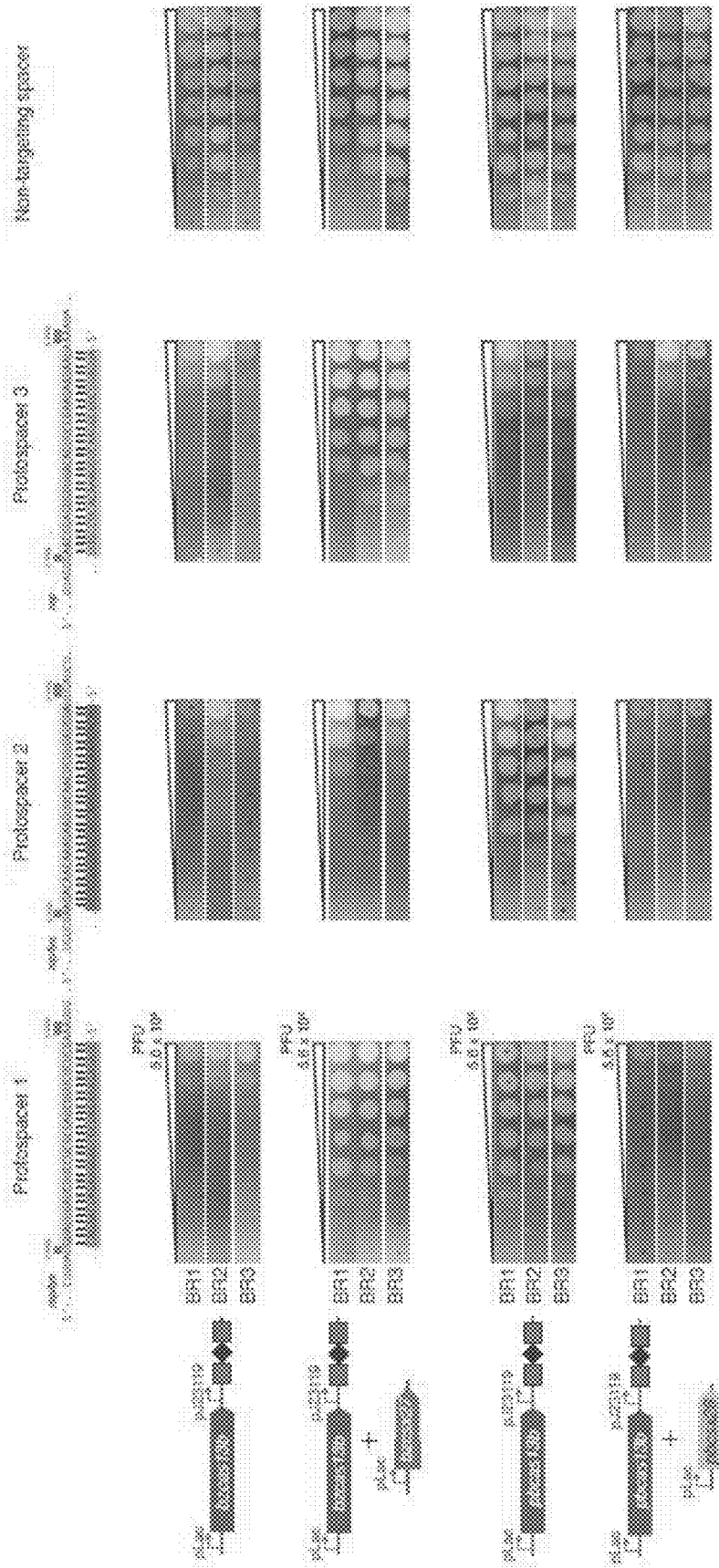


Fig. 75A

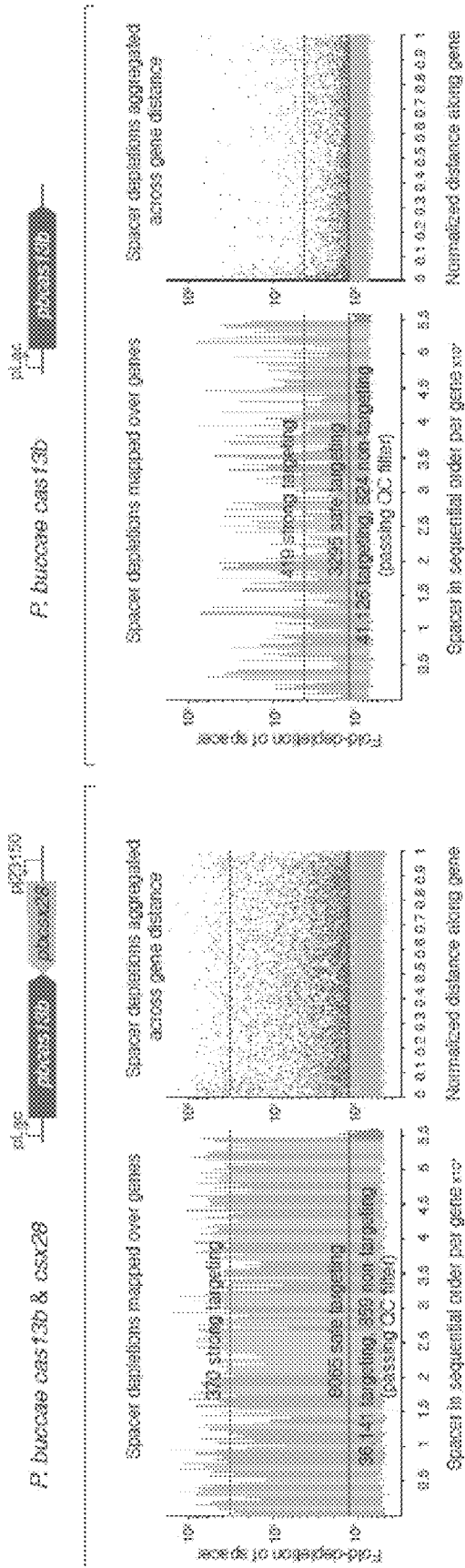


Fig. 75B

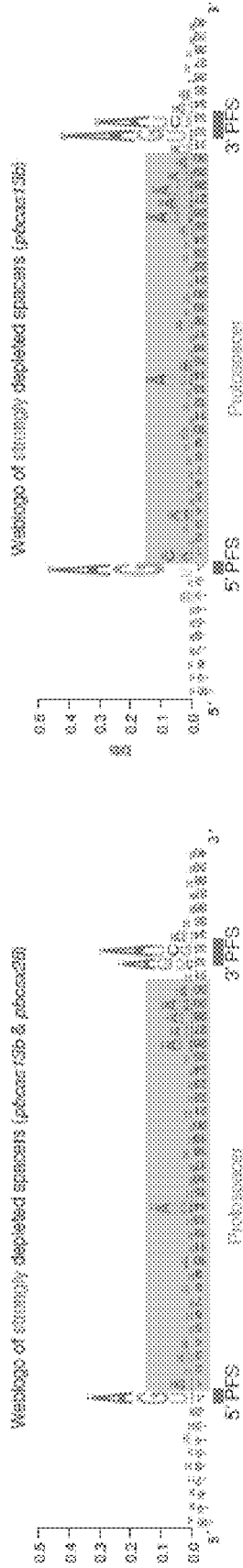
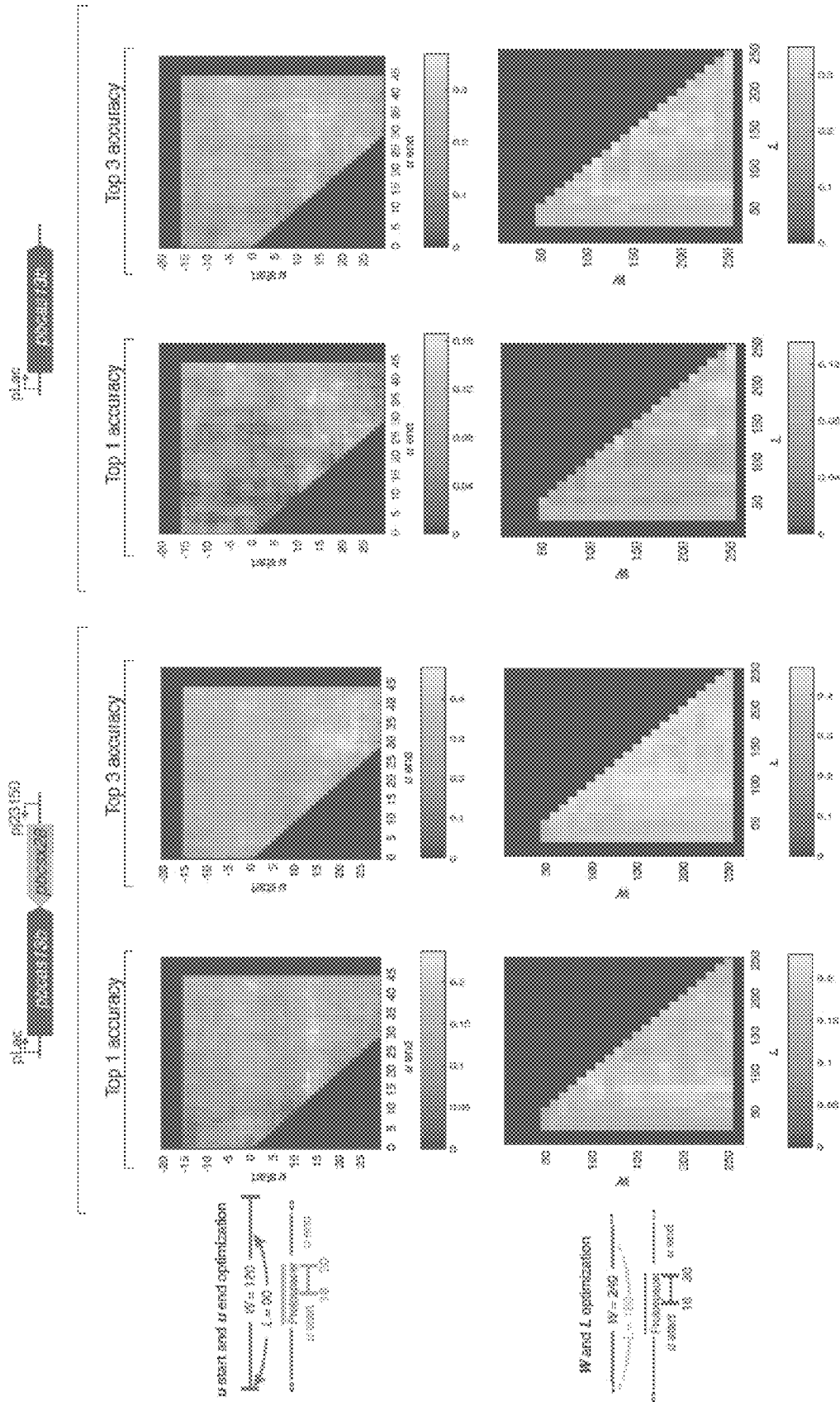


Fig. 76A

P. buccae spacer optimization on training set



u start and u end optimization
 $W = 100$
 $L = 90$
 u start: 0 5 10 15 20 25 30 35 40 45
 u end: 0 5 10 15 20 25 30 35 40 45

W and L optimization
 $W = 245$
 $L = 100$
 u start: 0 5 10 15 20 25 30 35 40 45
 u end: 0 5 10 15 20 25 30 35 40 45

P. buccae spacer performance on testing set

Fig. 76B

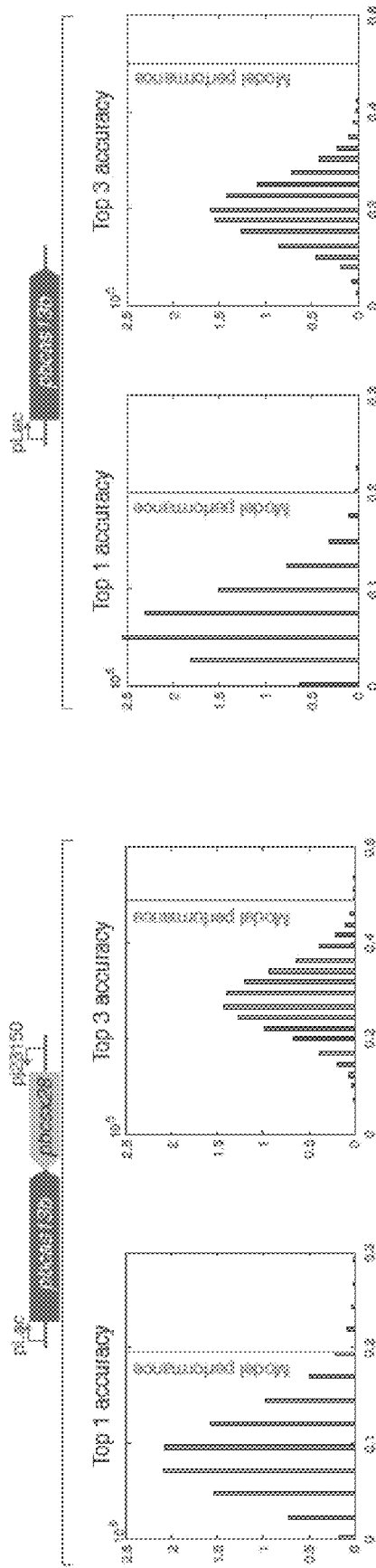


Fig. 76 C
P. buccae spatial distribution of safely depleted guides

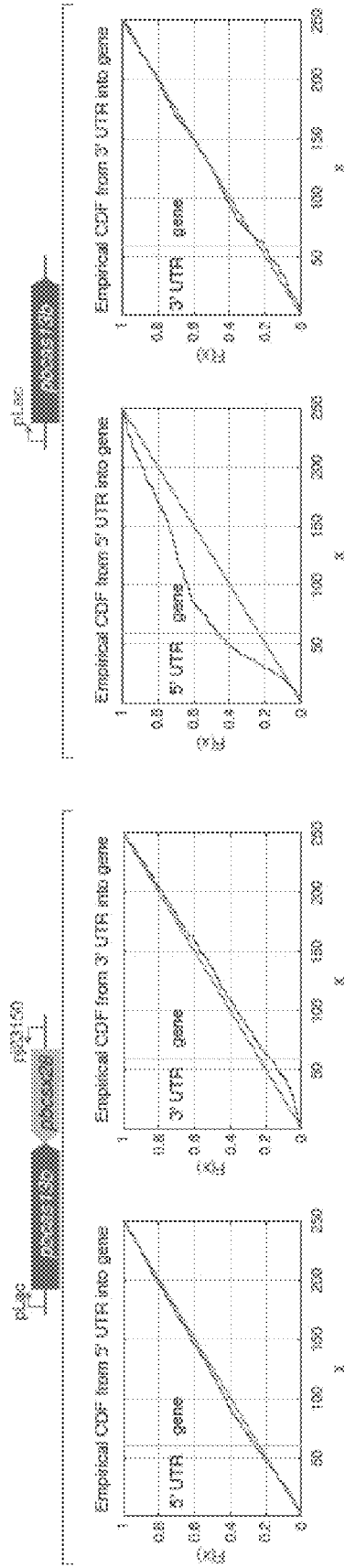
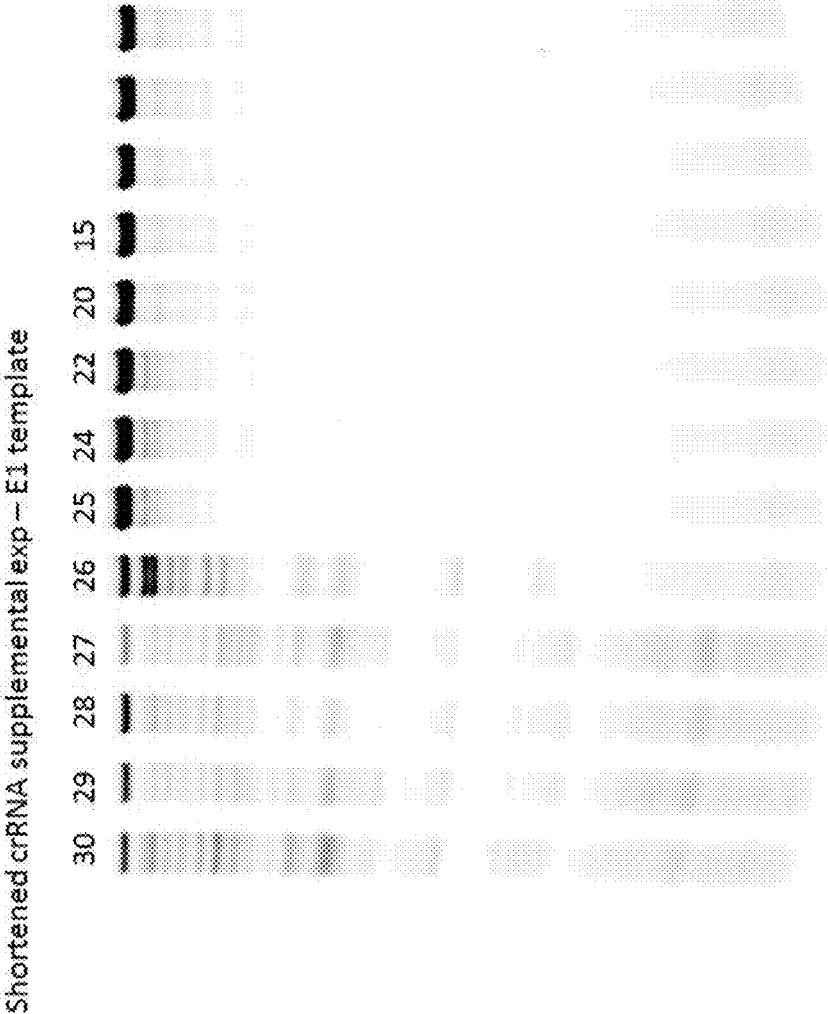


Fig. 77



EGFP 1 target, Body labeling
60 min incubation 50mM NaCl, 1mM DTT, Tris-HCl pH7.5, .1%BSA

Fig 78

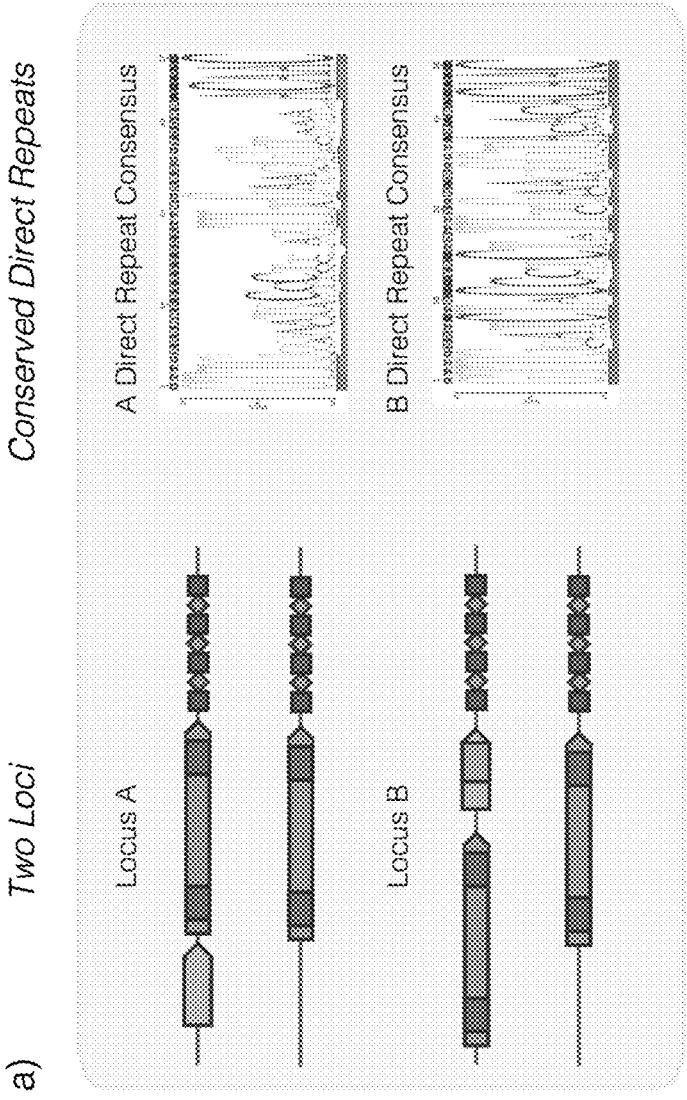


Fig. 79

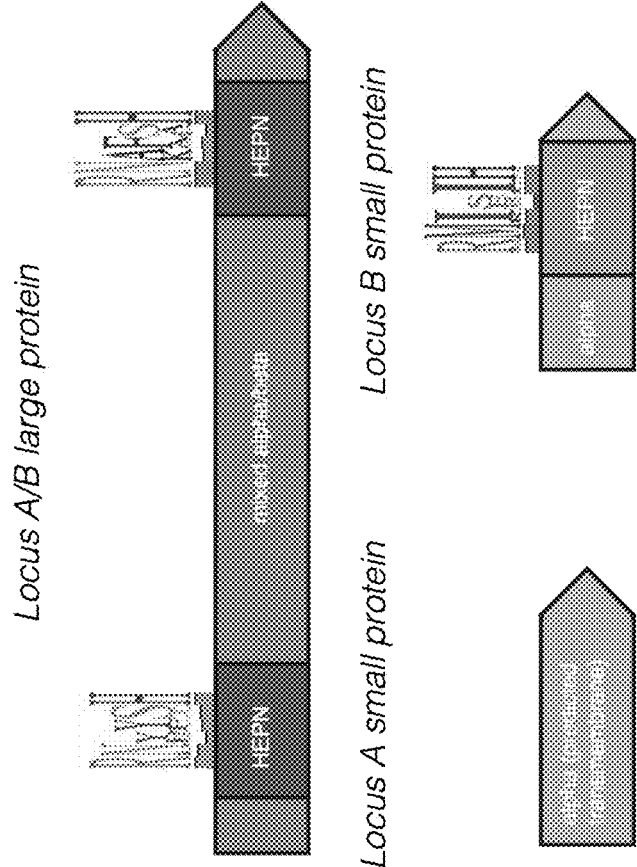


Fig. 80

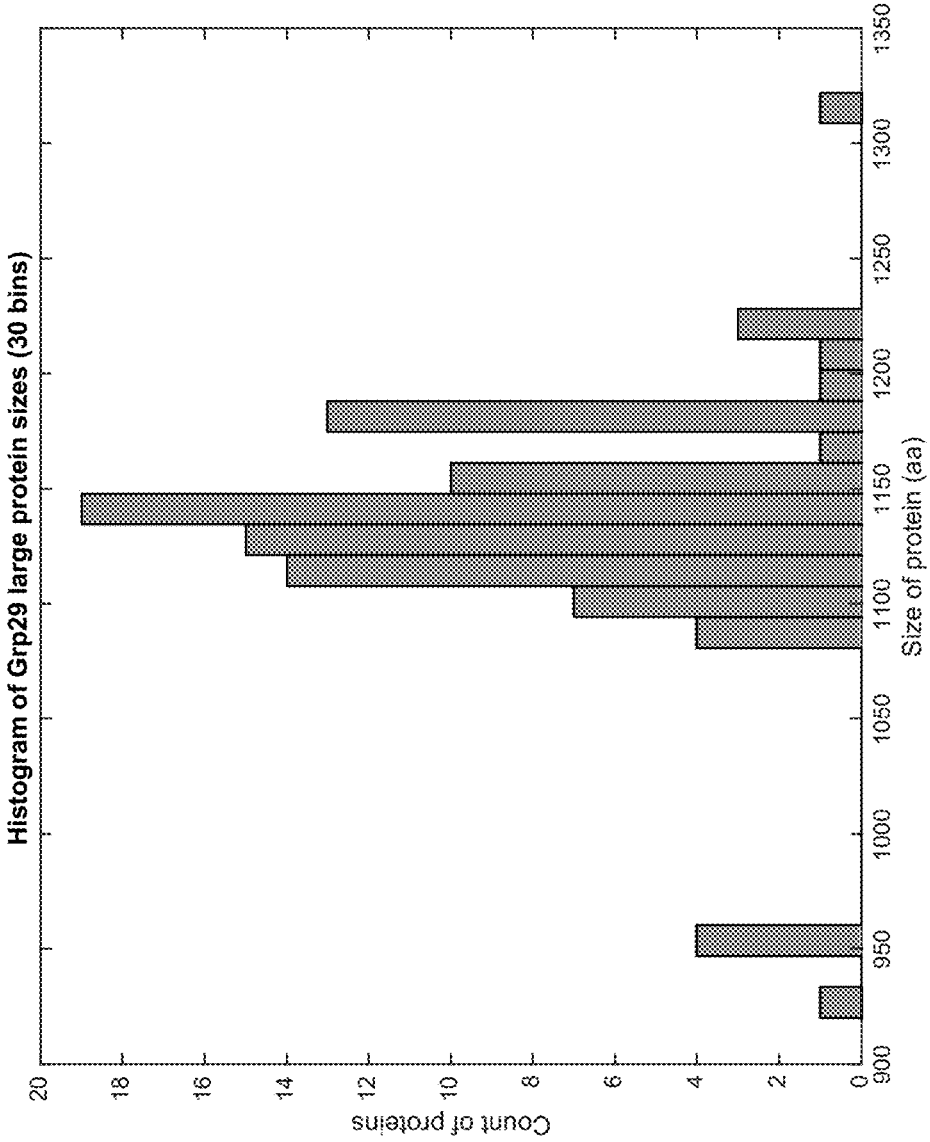
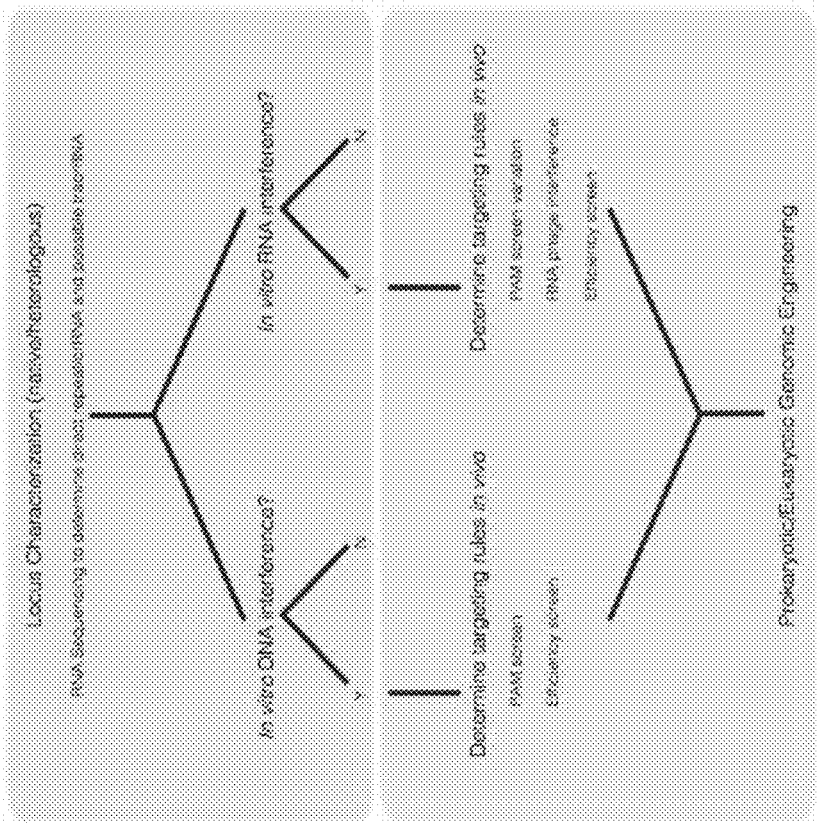


Fig. 82



NOVEL CRISPR ENZYMES AND SYSTEMS

RELATED APPLICATIONS AND INCORPORATION BY REFERENCE

[0001] This application claims priority to U.S. Provisional 62/245,270, filed on Oct. 22, 2015, U.S. Provisional 62/296,548, filed on Feb. 17, 2016, U.S. Provisional 62/376,382, filed on Aug. 17, 2016, and U.S. Provisional 62/376,367, filed Aug. 17, 2016.

STATEMENT AS TO FEDERALLY SPONSORED RESEARCH

[0002] This invention was made with government support under grant numbers MH100706, DK097768 and MH110049 awarded by the National Institutes of Health. The government has certain rights in the invention.

[0003] All documents cited or referenced in herein cited documents, together with any manufacturer's instructions, descriptions, product specifications, and product sheets for any products mentioned herein or in any document incorporated by reference herein, are hereby incorporated herein by reference, and may be employed in the practice of the invention. More specifically, all referenced documents are incorporated by reference to the same extent as if each individual document was specifically and individually indicated to be incorporated by reference.

FIELD OF THE INVENTION

[0004] The present invention generally relates to systems, methods and compositions used for the control of gene expression involving sequence targeting, such as perturbation of gene transcripts or nucleic acid editing, that may use vector systems related to Clustered Regularly Interspaced Short Palindromic Repeats (CRISPR) and components thereof.

BACKGROUND OF THE INVENTION

[0005] Recent advances in genome sequencing techniques and analysis methods have significantly accelerated the ability to catalog and map genetic factors associated with a diverse range of biological functions and diseases. Precise genome targeting technologies are needed to enable systematic reverse engineering of causal genetic variations by allowing selective perturbation of individual genetic elements, as well as to advance synthetic biology, biotechnological, and medical applications. Although genome-editing techniques such as designer zinc fingers, transcription activator-like effectors (TALEs), or homing meganucleases are available for producing targeted genome perturbations, there remains a need for new genome engineering technologies that employ novel strategies and molecular mechanisms and are affordable, easy to set up, scalable, and amenable to targeting multiple positions within the eukaryotic genome. This would provide a major resource for new applications in genome engineering and biotechnology.

[0006] The CRISPR-CRISPR associated (Cas) systems of bacterial and archaeal adaptive immunity are some such systems that show extreme diversity of protein composition and genomic loci architecture. The CRISPR-Cas system loci has more than 50 gene families and there is no strictly universal genes indicating fast evolution and extreme diversity of loci architecture. So far, adopting a multi-pronged approach, there is comprehensive cas gene identification of

about 395 profiles for 93 Cas proteins. Classification includes signature gene profiles plus signatures of locus architecture. A new classification of CRISPR-Cas systems is proposed in which these systems are broadly divided into two classes, Class 1 with multisubunit effector complexes and Class 2 with single-subunit effector modules exemplified by the Cas9 protein (FIGS. 1A and 1B). Novel effector proteins associated with Class 2 CRISPR-Cas systems may be developed as powerful genome engineering tools and the prediction of putative novel effector proteins and their engineering and optimization is important.

[0007] However, no researcher has undertaken a Cas protein-agnostic approach to identify novel single-effector systems. Doing so requires an unbiased and comprehensive bioinformatic analysis of all prokaryotic genomes and an accompanying sophisticated annotation methodology, and would potentially exhaust the possibility space for RNA-programmable CRISPR single effectors. The present application adopts such an approach and utilizes the ever-growing number of publicly accessible bacterial genomes and changing CRISPR rules to identify novel effector proteins that have expanded genome engineering capabilities.

[0008] Citation or identification of any document in this application is not an admission that such document is available as prior art to the present invention.

SUMMARY OF THE INVENTION

[0009] There exists a pressing need for alternative and robust systems and techniques for targeting nucleic acids or polynucleotides (e.g. DNA or RNA or any hybrid or derivative thereof) with a wide array of applications. This invention addresses this need and provides related advantages. Adding the novel DNA or RNA-targeting systems of the present application to the repertoire of genomic and epigenomic targeting technologies may transform the study and perturbation or editing of specific target sites through direct detection, analysis and manipulation. To utilize the DNA or RNA-targeting systems of the present application effectively for genomic or epigenomic targeting without deleterious effects, it is critical to understand aspects of engineering and optimization of these DNA or RNA targeting tools.

[0010] The invention provides computational methods for identifying and classifying new single effector protein systems that have applications in editing and modulating nucleic acids, expression products or genomes. In a preferred embodiment of the invention, the single effector protein does not have a proximate Cas protein (e.g. Cas1 or Cas2) in its corresponding locus. In an embodiment of the invention, the single effector protein comprises a Class 2 Type VI-B effector protein. Class 2 Type VI-B effector proteins include two subgroups, Type VI-B1 and Type VI-B2, which are also referred to as Group 29 proteins and Group 30 proteins, and include members which are RNA-programmable nucleases, RNA-interfering and may be involved in bacterial adoptive immunity against RNA phages.

[0011] Group 29 and group 30 systems comprise a large single effector (approximately 1100 amino acids in length), termed Cas13b, and one or none of two small putative accessory proteins (approximately 200 amino acids in length) nearby a CRISPR array. Based on the nearby small protein, the system is bifurcated into two Loci A and B. No additional proteins out to 25 kilobase pairs upstream or downstream from the array are conserved across species

with each locus. With minor exceptions, the CRISPR array comprises direct repeat sequences 36 nucleotides in length and spacer sequences 30 nucleotides in length. The direct repeat is generally well conserved, especially at the ends, with a GTTG/GUUG at the 5' end reverse complementary to a CAAC at the 3' end. This conservation suggests strong base pairing for an RNA loop structure that potentially interacts with the protein(s) in the locus. A motif search complementary to the direct repeats revealed no candidate tracrRNAs nearby the arrays, possibly indicative of a single crRNA like that found in the Cpf1 locus.

[0012] In embodiments of the invention, a Type VI-B system (e.g., a Group 29 or group 30 system) comprises a Cas13b effector protein and optionally a small accessory protein encoded upstream or downstream of the Cas13b effector protein. In certain embodiments, the small accessory protein enhances the Cas13b effector's ability to target RNA.

[0013] In certain embodiments of the invention, a Group 29 or group 30 system comprises a Cas13b effector protein and optionally a small accessory protein encoded upstream or downstream of the Cas13b effector protein. In certain embodiments, the small accessory protein represses the Cas13b effector's ability to target RNA.

[0014] The invention provides a non-naturally occurring or engineered composition comprising i) a Type VI-B CRISPR-Cas effector protein, and ii) a Type VI-B CRISPR-Cas crRNA, wherein the crRNA comprises a) a guide sequence that is capable of hybridizing to a target RNA sequence, and b) a direct repeat sequence. The Type VI-B CRISPR-Cas effector protein forms a complex with the crRNA, and the guide sequence directs sequence-specific binding of the complex to the target RNA sequence, whereby there is formed a CRISPR complex comprising the Type VI-B CRISPR-Cas effector protein complexed with the guide sequence that is hybridized to the target RNA sequence. The complex that is formed when the guide sequence hybridizes to the target RNA sequence includes interaction (recognition) of the protospacer flanking sequence (PFS).

[0015] In some embodiments, a non-naturally occurring or engineered composition of the invention may comprise a Type VI-B CRISPR-Cas accessory protein that enhances Type VI-B CRISPR-Cas effector protein activity. In certain embodiments, the accessory protein that enhances Type VI-B CRISPR-Cas effector protein activity is a csx28 protein. In such embodiments, the Type VI-B CRISPR-Cas effector protein and the Type VI-B CRISPR-Cas accessory protein may be from the same source or from a different source.

[0016] In some embodiment, a non-naturally occurring or engineered composition of the invention comprises a Type VI-B CRISPR-Cas accessory protein that represses Type VI-B CRISPR-Cas effector protein activity. In certain such embodiment, the accessory protein that represses Type VI-B CRISPR-Cas effector protein activity is a csx27 protein. In such embodiments, the Type VI-B CRISPR-Cas effector protein and the Type VI-B CRISPR-Cas accessory protein may be from the same source or from a different source. In certain embodiments of the invention, the Type VI-B CRISPR-Cas effector protein is from Table 1. In certain embodiments, the Type VI-B CRISPR-Cas accessory protein is from Table 1.

[0017] In some embodiments, a non-naturally occurring or engineered composition of the invention comprises two or more Type VI-B CRISPR-Cas crRNAs.

[0018] In some embodiments, a non-naturally occurring or engineered composition of the invention comprises a guide sequence that hybridizes to a target RNA sequence in a prokaryotic cell. In some embodiments, a non-naturally occurring or engineered composition of the invention comprises a guide sequence that hybridizes to a target RNA sequence in a eukaryotic cell. In some embodiment, the Type VI-B CRISPR-Cas effector protein comprises one or more nuclear localization signals (NLSs).

[0019] In some embodiment of the non-naturally occurring or engineered composition of the invention, the Type VI-B CRISPR-Cas effector protein is associated with one or more functional domains. The association can be by direct linkage of the effector protein to the functional domain, or by association with the crRNA. In a non-limiting example, the crRNA comprises an added or inserted sequence that can be associated with a functional domain of interest, including, for example, an aptamer or a nucleotide that binds to a nucleic acid binding adapter protein.

[0020] In certain non-limiting embodiments, a non-naturally occurring or engineered composition of the invention comprises a functional domain cleaves the target RNA sequence. In certain non-limiting embodiments, the non-naturally occurring or engineered composition of the invention comprises a functional domain that modifies transcription or translation of the target RNA sequence.

[0021] In some embodiment of the composition of the invention, the Type VI-B CRISPR-Cas effector protein is associated with one or more functional domains; and the effector protein contains one or more mutations within an HEPN domain, whereby the complex can deliver an epigenetic modifier or a transcriptional or translational activation or repression signal.

[0022] The invention also provides a Type VI-B CRISPR-Cas vector system, which comprises one or more vectors comprising a first regulatory element operably linked to a nucleotide sequence encoding the Type VI-B CRISPR-Cas effector protein, and a second regulatory element operably linked to a nucleotide sequence encoding the Type VI-B CRISPR-Cas crRNA. In certain embodiments, the vector system of the invention further comprises a regulatory element operably linked to a nucleotide sequence of a Type VI-B CRISPR-Cas accessory protein. When appropriate, the nucleotide sequence encoding the Type VI-B CRISPR-Cas effector protein and/or the nucleotide sequence encoding the Type VI-B CRISPR-Cas accessory protein are codon optimized for expression in a host cell, such as, for example, a eukaryotic cell.

[0023] In some embodiments, the vector system of the invention comprises a single vector. In other embodiments, the vector system comprises two or more vectors. In certain embodiments, the single vector, or the two or more vectors, comprise one or more viral vectors. Non-limiting examples of viral vectors include retroviral, lentiviral, adenoviral, adeno-associated or herpes simplex viral vectors.

[0024] Also provided is a delivery system configured to deliver a Type VI-B CRISPR-Cas effector protein and one or more nucleic acid components of a non-naturally occurring or engineered composition comprising i) a Type VI-B CRISPR-Cas effector protein, and ii) a Type VI-B CRISPR-Cas crRNA, wherein the crRNA comprises a) a guide

sequence that hybridizes to a target RNA sequence in a cell, and b) a direct repeat sequence, wherein the Type VI-B CRISPR-Cas system effector protein forms a complex with the crRNA, and wherein the guide sequence directs sequence-specific binding to the target RNA sequence, to form a CRISPR complex comprising the Type VI-B CRISPR-Cas effector protein complexed with the guide sequence that is hybridized to the target RNA sequence. Optionally, the delivery system delivers a Type VI-B CRISPR-Cas accessory protein. The system may be designed to deliver the components separately, for example at different times or under different conditions, or to deliver the components together.

[0025] In various embodiments, the delivery system comprises one or more vectors or one or more polynucleotide molecules, or combinations of one or more vectors and polynucleotide molecules, or combinations of one or more vectors, one or more polynucleotide molecules, Type VI-B CRISPR-Cas effector proteins, and Type VI-B CRISPR-Cas accessory proteins.

[0026] In certain embodiments, the delivery system of the invention comprises a delivery vehicle comprising liposome(s), particle(s), exosome(s), microvesicle(s), a gene-gun or one or more viral vector(s).

[0027] In some embodiments, the non-naturally occurring or engineered composition of the invention is for use in a therapeutic method of treatment. In some embodiments, the non-naturally occurring or engineered vector system of the invention is for use in a therapeutic method of treatment. In some embodiment, the non-naturally occurring or engineered delivery system of the invention is for use in a therapeutic method of treatment.

[0028] The invention provides a method of modifying expression of a target gene of interest, the method comprising contacting a target RNA with one or more non-naturally occurring or engineered compositions comprising i) a Type VI-B CRISPR-Cas effector protein, and ii) a Type VI-B CRISPR-Cas crRNA, wherein the crRNA comprises a) a guide sequence that hybridizes to a target RNA sequence in a cell, and b) a direct repeat sequence, wherein the Type VI-B CRISPR-Cas system effector protein forms a complex with the crRNA, wherein the guide sequence directs sequence-specific binding to the target RNA sequence in a cell, whereby there is formed a CRISPR complex comprising the Type VI-B CRISPR-Cas effector protein complexed with the guide sequence that is hybridized to the target RNA sequence, whereby expression of the target locus of interest is modified.

[0029] In certain embodiments, the method of modifying expression of a target gene of interest further comprises contacting the the target RNA with a Type VI-B CRISPR-Cas accessory protein that enhances Type VI-B CRISPR-Cas effector protein activity, such as a csx28 protein, non-limiting examples of which are provided.

[0030] In certain embodiments, the method of modifying expression of a target gene of interest further comprises contacting the the target RNA with a Type VI-B CRISPR-Cas accessory protein that represses Type VI-B CRISPR-Cas effector protein activity, such as a csx27 protein, non-limiting examples of which are provided.

[0031] In some embodiments, modifying expression of a target gene of interest comprises cleaving the target RNA. In some embodiments, modifying expression of a target gene of interest comprises increasing or decreasing expression of

the target RNA. In certain embodiments, modifying expression of a target gene causes increased or decreased expression of non-target components of a cell. In certain embodiments, the cell is a prokaryotic cell. In other embodiments, the target cell is a eukaryotic cell, non-limiting examples of which include plant cells and tissues, animal cells ad tissues, and human cells and tissues.

[0032] The invention provides a cell, tissue or organism comprising a modified target of interest, wherein the target of interest has been modified according to any of the method disclosed herein. In some embodiment, modification of the target of interest in a cell results in: a cell comprising altered expression of at least one gene product; a cell comprising altered expression of at least one gene product, wherein the expression of the at least one gene product is increased; or a cell comprising altered expression of at least one gene product, wherein the expression of the at least one gene product is decreased. Non-limiting examples include mammalian cell, human cell, and plant cells. The invention also provides a cell line of a cell disclosed herein or a cell modified by any of the methods disclosed herein, or progeny thereof.

[0033] The invention provides a multicellular organism comprising one or more cells disclosed herein or one or more cells modified according to any of the methods disclosed herein. The invention provides a plant or animal model comprising one or more cells disclosed herein or one or more cells modified according to any of the methods disclosed herein. The invention provides a gene product from a cell or the cell line or the organism or the plant or animal model disclosed herein. In some embodiments, the amount of gene product expressed is greater than or less than the amount of gene product from a cell that does not have altered expression.

[0034] The invention provides an isolated Type VI-B CRISPR-Cas effector protein. In non-limiting embodiments, the isolated Type VI-B CRISPR-Cas effector protein is from a micro-organism selected from *Bergeyella zoohelcum* or *Prevotella buccae*. The invention provides an isolated Type VI-B CRISPR-Cas accessory protein. In non-limiting embodiments, the isolated Type VI-B CRISPR-Cas accessory protein is from a micro-organism selected from *Bergeyella zoohelcum* or *Prevotella buccae*.

[0035] The invention provides an isolated nucleic acid encoding the Type VI-B CRISPR-Cas effector protein, as well as an isolated nucleic acid encoding the Type VI-B CRISPR-Cas accessory protein. In a non-limiting example, the isolated nucleic acid encodes a Type VI-B CRISPR-Cas effector protein from a micro-organism selected from *Bergeyella zoohelcum* or *Prevotella buccae*. The isolated nucleic acid can be DNA or RNA. Also provided is an isolated cell comprising the the nucleic acid encoding the Type VI-B CRISPR-Cas effector protein or accessory protein.

[0036] The invention provides a method of identifying the requirements of a suitable guide sequence for an RNA-targeting protein, said method comprising:

[0037] (a) selecting a set of essential genes within an organism

[0038] (b) designing a library of targeting guide sequences capable of hybridizing to regions the coding regions of these genes as well as 5' and 3' UTRs of these genes

[0039] (c) generating randomized guide sequences that do not hybridize to any region within the genome of said organism as control guides

[0040] (d) preparing a plasmid comprising the RNA-targeting protein and a first resistance gene and a guide plasmid library comprising said library of targeting guides and said control guides and a second resistance gene,

[0041] (e) co-introducing said plasmids into a host cell

[0042] (f) introducing said host cells on a selective medium for said first and second resistance genes

[0043] (g) sequencing essential genes of growing host cells

[0044] (h) determining significance of depletion of cells transformed with targeting guides by comparing depletion of cells with control guides; and

[0045] (i) determining based on the depleted guide sequences the requirements of a suitable guide sequence.

[0046] In one aspect of such method, determining the PFS sequence for suitable guide sequence of the RNA-targeting protein is by comparison of sequences targeted by guides in depleted cells.

[0047] In one aspect of such method, the method further comprises comparing the guide abundance for the different conditions in different replicate experiments.

[0048] In one aspect of such method, the control guides are selected in that they are determined to show limited deviation in guide depletion in replicate experiments.

[0049] In one aspect of such method, the significance of depletion is determined as (a) a depletion which is more than the most depleted control guide; or (b) a depletion which is more than the average depletion plus two times the standard deviation for the control guides.

[0050] In one aspect of such method, the host cell is a bacterial host cell.

[0051] In one aspect of such method, the step of co-introducing the plasmids is by electroporation and the host cell is an electro-competent host cell.

[0052] The invention provides a method of modifying sequences associated with or at a target locus of interest, the method comprising delivering to said locus a non-naturally occurring or engineered composition comprising a Group 29 or Group 30 effector protein and one or more nucleic acid components, wherein the effector protein forms a complex with the one or more nucleic acid components and upon binding of the said complex to the locus of interest the effector protein induces the modification of the sequences associated with or at the target locus of interest. In a preferred embodiment, the modification is the introduction of a strand break. In a preferred embodiment, the sequences associated with or at the target locus of interest comprises RNA or consists of RNA.

[0053] The invention provides a method of modifying sequences associated with or at a target locus of interest, the method comprising delivering to said locus a non-naturally occurring or engineered composition comprising a Cas13b effector protein, optionally a small accessory protein, and one or more nucleic acid components, wherein the effector protein forms a complex with the one or more nucleic acid components and upon binding of the said complex to the locus of interest the effector protein induces the modification of the sequences associated with or at the target locus of interest. In a preferred embodiment, the modification is the introduction of a strand break. In a preferred embodiment,

the sequences associated with or at the target locus of interest comprises RNA or consists of RNA.

[0054] It will be appreciated that the terms Cas enzyme, CRISPR enzyme, CRISPR protein, Cas protein, and CRISPR Cas are generally used interchangeably and at all points of reference herein refer by analogy to novel CRISPR effector proteins further described in this application, unless otherwise apparent, such as by specific reference to Cas9. The CRISPR effector proteins described herein are preferably Group 29 and Group 30 effector proteins.

[0055] The invention provides a method of modifying sequences associated with or at a target locus of interest, the method comprising delivering to said sequences associated with or at the locus a non-naturally occurring or engineered composition comprising a Group 29 loci effector protein and one or more nucleic acid components, wherein the Group 29 effector protein forms a complex with the one or more nucleic acid components and upon binding of the said complex to the locus of interest the effector protein induces the modification of sequences associated with or at the target locus of interest. In a preferred embodiment, the modification is the introduction of a strand break. In a preferred embodiment the Group 29 or Group 30 effector protein forms a complex with one nucleic acid component; advantageously an engineered or non-naturally occurring nucleic acid component. The induction of modification of sequences associated with or at the target locus of interest can be Group 29 or Group 30 effector protein-nucleic acid guided. In a preferred embodiment the one nucleic acid component is a CRISPR RNA (crRNA). In a preferred embodiment the one nucleic acid component is a mature crRNA or guide RNA, wherein the mature crRNA or guide RNA comprises a spacer sequence (or guide sequence) and a direct repeat (DR) sequence or derivatives thereof. In a preferred embodiment the spacer sequence or the derivative thereof comprises a seed sequence, wherein the seed sequence is critical for recognition and/or hybridization to the sequence at the target locus. In a preferred embodiment of the invention the crRNA is a short crRNA that may be associated with a short DR sequence. In another embodiment of the invention the crRNA is a long crRNA that may be associated with a long DR sequence (or dual DR). Aspects of the invention relate to Group 29 or Group 30 effector protein complexes having one or more non-naturally occurring or engineered or modified or optimized nucleic acid components. In a preferred embodiment the nucleic acid component comprises RNA. In a preferred embodiment the nucleic acid component of the complex may comprise a guide sequence linked to a direct repeat sequence, wherein the direct repeat sequence comprises one or more stem loops or optimized secondary structures. In preferred embodiments of the invention, the direct repeat may be a short DR or a long DR (dual DR). In a preferred embodiment the direct repeat may be modified to comprise one or more protein-binding RNA aptamers. In a preferred embodiment, one or more aptamers may be included such as part of optimized secondary structure. Such aptamers may be capable of binding a bacteriophage coat protein. The bacteriophage coat protein may be selected from the group comprising Q β , F2, GA, fr, JP501, MS2, M12, R17, BZ13, JP34, JP500, KU1, M11, MX1, TW18, VK, SP, FI, ID2, NL95, TW19, AP205, ϕ Cb5, ϕ Cb8r, ϕ Cb12r, ϕ Cb23r, 7s and PRR1. In a preferred embodiment the bacteriophage coat protein is MS2. The invention also

provides for the nucleic acid component of the complex being 30 or more, 40 or more or 50 or more nucleotides in length.

[0056] The invention provides methods of genome editing or modifying sequences associated with or at a target locus of interest wherein the method comprises introducing a Group 29 or Group 30 effector protein or Group 29 or Group 30 effector protein complex into any desired cell type, prokaryotic or eukaryotic cell, whereby the Group 29 or Group 30 effector protein or Group 29 or Group 30 effector protein complex effectively functions to interfere with RNA in the eukaryotic or prokaryotic cell. In preferred embodiments, the cell is a eukaryotic cell and the RNA is transcribed from a mammalian genome or is present in a mammalian cell. In preferred methods of RNA editing or genome editing in human cells, the Group 29 or Group 30 effector proteins may include but are not limited to the specific species of Group 29 or Group 30 effector proteins disclosed herein.

[0057] The invention also provides a method of modifying a target locus of interest, the method comprising delivering to said locus a non-naturally occurring or engineered composition comprising a Group 29 or Group 30 effector protein and one or more nucleic acid components, wherein the Group 29 or Group 30 effector protein forms a complex with the one or more nucleic acid components and upon binding of the said complex to the locus of interest the effector protein induces the modification of the target locus of interest. In a preferred embodiment, the modification is the introduction of a strand break.

[0058] In such methods the target locus of interest may be comprised within a RNA molecule. In such methods the target locus of interest may be comprised in a RNA molecule in vitro.

[0059] In such methods the target locus of interest may be comprised in a RNA molecule within a cell. The cell may be a prokaryotic cell or a eukaryotic cell. The cell may be a mammalian cell. The modification introduced to the cell by the present invention may be such that the cell and progeny of the cell are altered for improved production of biologic products such as an antibody, starch, alcohol or other desired cellular output. The modification introduced to the cell by the present invention may be such that the cell and progeny of the cell include an alteration that changes the biologic product produced.

[0060] The mammalian cell may be a non-human mammal, e.g., primate, bovine, ovine, porcine, canine, rodent, Leporidae such as monkey, cow, sheep, pig, dog, rabbit, rat or mouse cell. The cell may be a non-mammalian eukaryotic cell such as poultry bird (e.g., chicken), vertebrate fish (e.g., salmon) or shellfish (e.g., oyster, clam, lobster, shrimp) cell. The cell may also be a plant cell. The plant cell may be of a monocot or dicot or of a crop or grain plant such as cassava, corn, sorghum, soybean, wheat, oat or rice. The plant cell may also be of an algae, tree or production plant, fruit or vegetable (e.g., trees such as citrus trees, e.g., orange, grapefruit or lemon trees; peach or nectarine trees; apple or pear trees; nut trees such as almond or walnut or pistachio trees; nightshade plants; plants of the genus *Brassica*; plants of the genus *Lactuca*; plants of the genus *Spinacia*; plants of the genus *Capsicum*; cotton, tobacco, asparagus, carrot, cabbage, broccoli, cauliflower, tomato, eggplant, pepper, lettuce, spinach, strawberry, blueberry, raspberry, blackberry, grape, coffee, cocoa, etc).

[0061] The invention provides a method of modifying a target locus of interest, the method comprising delivering to said locus a non-naturally occurring or engineered composition comprising a Group 29 or Group 30 effector protein and one or more nucleic acid components, wherein the effector protein forms a complex with the one or more nucleic acid components and upon binding of the said complex to the locus of interest the effector protein induces the modification of the target locus of interest. In a preferred embodiment, the modification is the introduction of a strand break.

[0062] In such methods the target locus of interest may be comprised within an RNA molecule. In a preferred embodiment, the target locus of interest comprises or consists of RNA.

[0063] The invention also provides a method of modifying a target locus of interest, the method comprising delivering to said locus a non-naturally occurring or engineered composition comprising a Group 29 or Group 30 effector protein and one or more nucleic acid components, wherein the Group 29 or Group 30 effector protein forms a complex with the one or more nucleic acid components and upon binding of the said complex to the locus of interest the effector protein induces the modification of the target locus of interest. In a preferred embodiment, the modification is the introduction of a strand break.

[0064] Preferably, in such methods the target locus of interest may be comprised in a RNA molecule in vitro. Also preferably, in such methods the target locus of interest may be comprised in a RNA molecule within a cell. The cell may be a prokaryotic cell or a eukaryotic cell. The cell may be a mammalian cell. The cell may be a rodent cell. The cell may be a mouse cell.

[0065] In any of the described methods the target locus of interest may be a genomic or epigenomic locus of interest. In any of the described methods the complex may be delivered with multiple guides for multiplexed use. In any of the described methods more than one protein(s) may be used.

[0066] In further aspects of the invention the nucleic acid components may comprise a CRISPR RNA (crRNA) sequence and/or a trans-activating crRNA (tracr RNA) sequence. In certain embodiments, cleavage such as biochemical or in vitro cleavage or cleavage in cells, can result without a transactivating crRNA (tracr RNA) sequence. In other embodiments, cleavage such as biochemical or in vitro cleavage or cleavage in cells, can result with a transactivating crRNA (tracr RNA) sequence.

[0067] In further aspects of the invention the nucleic acid components may comprise a CRISPR RNA (crRNA) sequence and not comprise any trans-activating crRNA (tracr RNA) sequence. Without limitation, the Applicants hypothesize that in such instances, the pre-crRNA may comprise secondary structure that is sufficient for processing to yield the mature crRNA as well as crRNA loading onto the effector protein. By means of example and not limitation, such secondary structure may comprise, consist essentially of or consist of one or more stem loop within the pre-crRNA, more particularly within the direct repeat.

[0068] In certain embodiments, where the effector protein is a Group 29 or Group 30 effector protein, the nucleic acid components may comprise a CRISPR RNA (crRNA) sequence and may not comprise any trans-activating crRNA (tracr RNA) sequence.

[0069] In any of the described methods the effector protein and nucleic acid components may be provided via one or more polynucleotide molecules encoding the protein and/or nucleic acid component(s), and wherein the one or more polynucleotide molecules are operably configured to express the protein and/or the nucleic acid component(s). The one or more polynucleotide molecules may comprise one or more regulatory elements operably configured to express the protein and/or the nucleic acid component(s). The one or more polynucleotide molecules may be comprised within one or more vectors. In any of the described methods the target locus of interest may be a genomic, epigenomic, or transcriptomic locus of interest. In any of the described methods the complex may be delivered with multiple guides for multiplexed use. In any of the described methods more than one protein(s) may be used.

[0070] In any of the described methods the strand break may be a single strand break or a double strand break. In preferred embodiments the double strand break may refer to the breakage of two sections of RNA, such as the two sections of RNA formed when a single strand RNA molecule has folded onto itself or putative double helices that are formed with an RNA molecule which contains self-complementary sequences allows parts of the RNA to fold and pair with itself.

[0071] Regulatory elements may comprise inducible promoters. Polynucleotides and/or vector systems may comprise inducible systems.

[0072] In any of the described methods the one or more polynucleotide molecules may be comprised in a delivery system, or the one or more vectors may be comprised in a delivery system.

[0073] In any of the described methods the non-naturally occurring or engineered composition may be delivered via liposomes, particles including particles, exosomes, microvesicles, a gene-gun or one or more viral vectors.

[0074] The invention also provides a non-naturally occurring or engineered composition which is a composition having the characteristics as discussed herein or defined in any of the herein described methods.

[0075] In certain embodiments, the invention thus provides a non-naturally occurring or engineered composition, such as particularly a composition capable of or configured to modify a target locus of interest, said composition comprising a Group 29 or Group 30 effector protein and one or more nucleic acid components, wherein the effector protein forms a complex with the one or more nucleic acid components and upon binding of the said complex to the locus of interest the effector protein induces the modification of the target locus of interest. In certain embodiments, the effector protein may be a Group 29 effector protein.

[0076] The invention also provides in a further aspect a non-naturally occurring or engineered composition, such as particularly a composition capable of or configured to modify a target locus of interest, said composition comprising: (a) a guide RNA molecule (or a combination of guide RNA molecules, e.g., a first guide RNA molecule and a second guide RNA molecule) or a nucleic acid encoding the guide RNA molecule (or one or more nucleic acids encoding the combination of guide RNA molecules); (b) a Group 29 or Group 30 effector protein. In certain embodiments, the effector protein may be a Group 29 effector protein.

[0077] The invention also provides in a further aspect a non-naturally occurring or engineered composition compris-

ing: (I.) one or more CRISPR-Cas system polynucleotide sequences comprising (a) a guide sequence capable of hybridizing to a target sequence in a polynucleotide locus, (b) a tracr mate sequence, and (c) a tracrRNA sequence, and (II.) a second polynucleotide sequence encoding a Group 29 or Group 30 effector protein, wherein when transcribed, the tracr mate sequence hybridizes to the tracrRNA sequence and the guide sequence directs sequence-specific binding of a CRISPR complex to the target sequence, and wherein the CRISPR complex comprises the Group 29 or Group 30 effector protein complexed with (1) the guide sequence that is hybridized to the target sequence, and (2) the tracr mate sequence that is hybridized to the tracrRNA sequence. In certain embodiments, the effector protein may be a Group 29 effector protein.

[0078] In certain embodiments, a tracrRNA may not be required. Hence, the invention also provides in certain embodiments a non-naturally occurring or engineered composition comprising: (I.) one or more CRISPR-Cas system polynucleotide sequences comprising (a) a guide sequence capable of hybridizing to a target sequence in a polynucleotide locus, and (b) a direct repeat sequence, and (II.) a second polynucleotide sequence encoding a Group 29 or Group 30 effector protein, wherein when transcribed, the guide sequence directs sequence-specific binding of a CRISPR complex to the target sequence, and wherein the CRISPR complex comprises the Group 29 or Group 30 effector protein complexed with (1) the guide sequence that is hybridized to the target sequence, and (2) the direct repeat sequence. Preferably, the effector protein may be a Group 29 effector protein. Without limitation, the Applicants hypothesize that in such instances, the direct repeat sequence may comprise secondary structure that is sufficient for crRNA loading onto the effector protein. By means of example and not limitation, such secondary structure may comprise, consist essentially of or consist of a stem loop (such as one or more stem loops) within the direct repeat.

[0079] The invention also provides a vector system comprising one or more vectors, the one or more vectors comprising one or more polynucleotide molecules encoding components of a non-naturally occurring or engineered composition which is a composition having the characteristics as defined in any of the herein described methods.

[0080] The invention also provides a delivery system comprising one or more vectors or one or more polynucleotide molecules, the one or more vectors or polynucleotide molecules comprising one or more polynucleotide molecules encoding components of a non-naturally occurring or engineered composition which is a composition having the characteristics discussed herein or as defined in any of the herein described methods.

[0081] The invention also provides a non-naturally occurring or engineered composition, or one or more polynucleotides encoding components of said composition, or vector or delivery systems comprising one or more polynucleotides encoding components of said composition for use in a therapeutic method of treatment. The therapeutic method of treatment may comprise gene or genome editing, or gene therapy.

[0082] The invention also encompasses computational methods and algorithms to predict new Cas protein agnostic single effector protein CRISPR systems and identify the components therein.

[0083] The invention also provides for methods and compositions wherein one or more amino acid residues of the effector protein may be modified e.g., an engineered or non-naturally-occurring Group 20 or Group 30 effector protein. In an embodiment, the modification may comprise mutation of one or more amino acid residues of the effector protein. The one or more mutations may be in one or more catalytically active domains of the effector protein. The effector protein may have reduced or abolished nuclease activity compared with an effector protein lacking said one or more mutations. The effector protein may not direct cleavage of one RNA strand at the target locus of interest. In a preferred embodiment, the one or more mutations may comprise two mutations. In a preferred embodiment the one or more amino acid residues are modified in a Group 29 or Group 30 effector protein, e.g., an engineered or non-naturally-occurring Group 29 or Group 30 effector protein. In certain embodiments of the invention the effector protein comprises one or more HEPN domains. In a preferred embodiment, the effector protein comprises two HEPN domains. In another preferred embodiment, the effector protein comprises one HEPN domain at the C-terminus and another HEPN domain at the N-terminus of the protein. In certain embodiments, the one or more mutations or the two or more mutations may be in a catalytically active domain of the effector protein comprising a HEPN domain, or a catalytically active domain which is homologous to a HEPN domain. In certain embodiments, the effector protein comprises one or more of the following mutations: R116A, H121A, R1177A, H1182A (wherein amino acid positions correspond to amino acid positions of Group 29 protein originating from *Bergeyella zoohelcum* ATCC 43767), such as R116A, H121A, R1177A, and H1182A; R116A, H121A, and R1177A; R116A, H121A, and H1182A; R116A, R1177A, and H1182A; H121A, R1177A, and H1182A; R116A and H121A; R116A and R1177A; R116A and H1182A; H121A and R1177A; H121A and H1182A; R1177A and H1182A; R116A; H121A; R1177A; H1182A. The skilled person will understand that corresponding amino acid positions in different Group 29 or Group 30 proteins may be mutated to the same effect. In certain embodiments, one or more of mutations R116A, H121A, R1177A, H1182A abolish catalytic activity of the protein completely or partially (e.g. altered cleavage rate, altered specificity, etc.), such as R116A, H121A, R1177A, and H1182A; R116A, H121A, and R1177A; R116A, H121A, and H1182A; R116A, R1177A, and H1182A; H121A, R1177A, and H1182A; R116A and H121A; R116A and R1177A; R116A and H1182A; H121A and R1177A; H121A and H1182A; R1177A and H1182A; R116A; H121A; R1177A; H1182A. In certain embodiments, the effector protein as described herein is a “dead” effector protein, such as a dead Group 29 or dead Group 30 effector protein (i.e. dGroup 29 or dGroup 30). In certain embodiments, the effector protein has one or more mutations in HEPN domain 1. In certain embodiments, the effector protein has one or more mutations in HEPN domain 2. In certain embodiments, the effector protein has one or more mutations in HEPN domain 1 and HEPN domain 2. The effector protein may comprise one or more heterologous functional domains. The one or more heterologous functional domains may comprise one or more nuclear localization signal (NLS) domains. The one or more heterologous functional domains may comprise at least two or more NLS domains. The one or more NLS domain(s) may

be positioned at or near or in proximity to a terminus of the effector protein (e.g., Group 29 or Group 30 effector protein) and if two or more NLSs, each of the two may be positioned at or near or in proximity to a terminus of the effector protein (e.g., Group 29 or Group 30 effector protein). The one or more heterologous functional domains may comprise one or more transcriptional activation domains. In a preferred embodiment the transcriptional activation domain may comprise VP64. The one or more heterologous functional domains may comprise one or more transcriptional repression domains. In a preferred embodiment the transcriptional repression domain comprises a KRAB domain or a SID domain (e.g. SID4X). The one or more heterologous functional domains may comprise one or more nuclease domains. In a preferred embodiment a nuclease domain comprises FokI.

[0084] The invention also provides for the one or more heterologous functional domains to have one or more of the following activities: methylase activity, demethylase activity, transcription activation activity, transcription repression activity, transcription release factor activity, histone modification activity, nuclease activity, single-strand RNA cleavage activity, double-strand RNA cleavage activity, single-strand DNA cleavage activity, double-strand DNA cleavage activity and nucleic acid binding activity. At least one or more heterologous functional domains may be at or near the amino-terminus of the effector protein and/or wherein at least one or more heterologous functional domains is at or near the carboxy-terminus of the effector protein. The one or more heterologous functional domains may be fused to the effector protein. The one or more heterologous functional domains may be tethered to the effector protein. The one or more heterologous functional domains may be linked to the effector protein by a linker moiety.

[0085] The invention also provides for the effector protein comprising an effector protein from an organism from a genus comprising *Bergeyella*, *Prevotella*, *Porphyromonas*, *Bacteroides*, *Alistipes*, *Riemerella*, *Myroides*, *Flavobacterium*, *Capnocytophaga*, *Chryseobacterium*, *Phaeodactylibacter*, *Paludibacter* or *Psychroflexus*. In a preferred embodiment the Group 29 or Group 30 effector protein may be from a group which includes but is not limited to the following organisms: *Bergeyella zoohelcum*, *Prevotella intermedia*, *Prevotella buccae*, *Porphyromonas gingivalis*, *Bacteroides pyogenes*, *Alistipes* sp., *Prevotella* sp. MA2016, *Riemerella anatipestifer*, *Prevotella aurantiaca*, *Prevotella saccharolytica*, *Myroides odoratimimus*, *Flavobacterium columnare*, *Flavobacterium branchiophilum*, *Bacteroides pyogenes*, *Bacteroides coprosuis*, *Capnocytophaga canimorsus*, *Capnocytophaga cynodegmi*, *Chryseobacterium* sp., *Paludibacter propionicigenes*, *Phaeodactylibacter xiamenensis*, *Porphyromonas gingivalis*, *Porphyromonas gulae*, *Prevotella falsenii*, *Prevotella intermedia*, *Prevotella pallens*, *Prevotella pleuritidis*, *Prevotella saccharolytica*, *Prevotella* sp. P5-119 and *Psychroflexus torquis*.

[0086] The effector protein may comprise a chimeric effector protein comprising a first fragment from a first effector protein ortholog and a second fragment from a second effector protein ortholog, and wherein the first and second effector protein orthologs are different. At least one of the first and second effector protein orthologs may comprise an effector protein from an organism comprising *Bergeyella*, *Prevotella*, *Porphyromonas*, *Bacteroides*, *Alis-*

tipes, *Riemerella*, *Myroides*, *Flavobacterium*, *Capnocytophaga*, *Chryseobacterium*, *Paludibacter*, *Phaeodactylibacter* or *Psychroflexus*.

[0087] In certain embodiments, the effector protein, particularly a Group 29 or Group 30 effector protein effector protein may be at least 700 amino acids long. In preferred embodiments, the effector protein may be about 1100 to about 1500 amino acids long, e.g., about 1100 to about 1200 amino acids long, or about 1200 to about 1300 amino acids long, or about 1300 to about 1400 amino acids long, or about 1400 to about 1500 amino acids long, e.g., about 900, about 1000, about 1100, about 1200, about 1300, about 1400, about 1500, about 1600, about 1700, or about 1800 amino acids long.

[0088] In certain embodiments, the Group 29 or Group 30 effector proteins as intended herein may be associated with a locus comprising short CRISPR repeats between 30 and 40 bp long, more typically between 34 and 38 bp long, even more typically between 36 and 37 bp long, e.g., 30, 31, 32, 33, 34, 35, 36, 37, 38, 39, or 40 bp long. In certain embodiments the CRISPR repeats are long or dual repeats between 80 and 350 bp long such as between 80 and 200 bp long, even more typically between 86 and 88 bp long, e.g., 80, 81, 82, 83, 84, 85, 86, 87, 88, 89, or 90 bp long.

[0089] In certain embodiments, a protospacer adjacent motif (PAM) or PAM-like motif directs binding of the effector protein (e.g. a Group 29 or Group 30 effector protein) complex as disclosed herein to the target locus of interest. In some embodiments, the PAM may be a 5' PAM (i.e., located upstream of the 5' end of the protospacer). In other embodiments, the PAM may be a 3' PAM (i.e., located downstream of the 5' end of the protospacer). In other embodiments, both a 5' PAM and a 3' PAM are required. In certain embodiments of the invention, a PAM or PAM-like motif may not be required for directing binding of the effector protein (e.g. a Group 29 or Group 30 effector protein). In certain embodiments, a 5' PAM is D (i.e. A, G, or U). In certain embodiments, a 5' PAM is D for Group 29 effectors. In certain embodiments, a 5' PAM is D for Group 29 effectors originating from *Bergeyella zoohelcum* (such as *Bergeyella zoohelcum* ATCC 43767). In certain embodiments of the invention, cleavage at repeat sequences may generate crRNAs (e.g. short or long crRNAs) containing a full spacer sequence flanked by a short nucleotide (e.g. 5, 6, 7, 8, 9, or 10 nt or longer if it is a dual repeat) repeat sequence at the 5' end (this may be referred to as a crRNA "tag") and the rest of the repeat at the 3' end. In certain embodiments, targeting by the effector proteins described herein may require the lack of homology between the crRNA tag and the target 5' flanking sequence. This requirement may be similar to that described further in Samai et al. "Co-transcriptional DNA and RNA Cleavage during Type III CRISPR-Cas Immunity" Cell 161, 1164-1174, May 21, 2015, where the requirement is thought to distinguish between bona fide targets on invading nucleic acids from the CRISPR array itself, and where the presence of repeat sequences will lead to full homology with the crRNA tag and prevent autoimmunity.

[0090] In certain embodiments, the Group 29 or Group 30 effector proteins engineered and can comprise one or more mutations that reduce or eliminate nuclease activity, thereby reducing or eliminating RNA interfering activity. Mutations can also be made at neighboring residues, e.g., at amino acids near those that participate in the nuclease activity. In

some embodiments, one or more putative catalytic nuclease domains are inactivated and the effector protein complex lacks cleavage activity and functions as an RNA binding complex. In a preferred embodiment, the resulting RNA binding complex may be linked with one or more functional domains as described herein.

[0091] In certain embodiments, the one or more functional domains are controllable, i.e. inducible.

[0092] In certain embodiments of the invention, the guide RNA or mature crRNA comprises, consists essentially of, or consists of a direct repeat sequence and a guide sequence or spacer sequence. In certain embodiments, the guide RNA or mature crRNA comprises, consists essentially of, or consists of a direct repeat sequence linked to a guide sequence or spacer sequence. In preferred embodiments of the invention, the mature crRNA comprises a stem loop or an optimized stem loop structure or an optimized secondary structure. In preferred embodiments the mature crRNA comprises a stem loop or an optimized stem loop structure in the direct repeat sequence, wherein the stem loop or optimized stem loop structure is important for cleavage activity. In certain embodiments, the mature crRNA preferably comprises a single stem loop. In certain embodiments, the direct repeat sequence preferably comprises a single stem loop. In certain embodiments, the cleavage activity of the effector protein complex is modified by introducing mutations that affect the stem loop RNA duplex structure. In preferred embodiments, mutations which maintain the RNA duplex of the stem loop may be introduced, whereby the cleavage activity of the effector protein complex is maintained. In other preferred embodiments, mutations which disrupt the RNA duplex structure of the stem loop may be introduced, whereby the cleavage activity of the effector protein complex is completely abolished.

[0093] The invention also provides for the nucleotide sequence encoding the effector protein being codon optimized for expression in a eukaryote or eukaryotic cell in any of the herein described methods or compositions. In an embodiment of the invention, the codon optimized effector protein is any Group 29 or Group 30 effector protein discussed herein and is codon optimized for operability in a eukaryotic cell or organism, e.g., such cell or organism as elsewhere herein mentioned, for instance, without limitation, a yeast cell, or a mammalian cell or organism, including a mouse cell, a rat cell, and a human cell or non-human eukaryote organism, e.g., plant.

[0094] In certain embodiments of the invention, at least one nuclear localization signal (NLS) is attached to the nucleic acid sequences encoding the Group 29 or Group 30 Group 29 or Group 30 effector proteins. In preferred embodiments at least one or more C-terminal or N-terminal NLSs are attached (and hence nucleic acid molecule(s) coding for the Group 29 or Group 30 effector protein can include coding for NLS(s) so that the expressed product has the NLS(s) attached or connected). In a preferred embodiment a C-terminal NLS is attached for optimal expression and nuclear targeting in eukaryotic cells, preferably human cells. The invention also encompasses methods for delivering multiple nucleic acid components, wherein each nucleic acid component is specific for a different target locus of interest thereby modifying multiple target loci of interest. The nucleic acid component of the complex may comprise one or more protein-binding RNA aptamers. The one or more aptamers may be capable of binding a bacteriophage

coat protein. The bacteriophage coat protein may be selected from the group comprising Q β , F2, GA, fr, JP501, MS2, M12, R17, BZ13, JP34, JP500, KU1, M11, MX1, TW18, VK, SP, FI, ID2, NL95, TW19, AP205, ϕ Cb5, ϕ Cb8r, ϕ Cb12r, ϕ Cb23r, 7s and PRR1. In a preferred embodiment the bacteriophage coat protein is MS2. The invention also provides for the nucleic acid component of the complex being 30 or more, 40 or more or 50 or more nucleotides in length.

[0095] In a further aspect, the invention provides a eukaryotic cell comprising a modified target locus of interest, wherein the target locus of interest has been modified according to in any of the herein described methods. A further aspect provides a cell line of said cell. Another aspect provides a multicellular organism comprising one or more said cells.

[0096] In certain embodiments, the modification of the target locus of interest may result in: the eukaryotic cell comprising altered expression of at least one gene product; the eukaryotic cell comprising altered expression of at least one gene product, wherein the expression of the at least one gene product is increased; the eukaryotic cell comprising altered expression of at least one gene product, wherein the expression of the at least one gene product is decreased; or the eukaryotic cell comprising an edited genome.

[0097] In certain embodiments, the eukaryotic cell may be a mammalian cell or a human cell.

[0098] In further embodiments, the non-naturally occurring or engineered compositions, the vector systems, or the delivery systems as described in the present specification may be used for: site-specific gene knockout; site-specific genome editing; RNA sequence-specific interference; or multiplexed genome engineering.

[0099] Also provided is a gene product from the cell, the cell line, or the organism as described herein. In certain embodiments, the amount of gene product expressed may be greater than or less than the amount of gene product from a cell that does not have altered expression or edited genome. In certain embodiments, the gene product may be altered in comparison with the gene product from a cell that does not have altered expression or edited genome.

[0100] In another aspect, the invention provides a method for identifying novel nucleic acid modifying effectors, comprising: identifying putative nucleic acid modifying loci from a set of nucleic acid sequences encoding the putative nucleic acid modifying enzyme loci that are within a defined distance from a conserved genomic element of the loci, that comprise at least one protein above a defined size limit, or both; grouping the identified putative nucleic acid modifying loci into subsets comprising homologous proteins; identifying a final set of candidate nucleic acid modifying loci by selecting nucleic acid modifying loci from one or more subsets based on one or more of the following; subsets comprising loci with putative effector proteins with low domain homology matches to known protein domains relative to loci in other subsets, subsets comprising putative proteins with minimal distances to the conserved genomic element relative to loci in other subsets, subsets with loci comprising large effector proteins having a same orientations as putative adjacent accessory proteins relative to large effector proteins in other subsets, subset comprising putative effector proteins with lower existing nucleic acid modifying classifications relative to other loci, subsets comprising loci

with a lower proximity to known nucleic acid modifying loci relative to other subsets, and total number of candidate loci in each subset.

[0101] In one embodiment, the set of nucleic acid sequences is obtained from a genomic or metagenomic database, such as a genomic or metagenomic database comprising prokaryotic genomic or metagenomic sequences.

[0102] In one embodiment, the defined distance from the conserved genomic element is between 1 kb and 25 kb.

[0103] In one embodiment, the conserved genomic element comprises a repetitive element, such as a CRISPR array. In a specific embodiment, the defined distance from the conserved genomic element is within 10 kb of the CRISPR array.

[0104] In one embodiment, the defined size limit of a protein comprised within the putative nucleic acid modifying locus is greater than 200 amino acids, or more particularly, the defined size limit is greater than 700 amino acids.

[0105] In one embodiment, the conserved genomic elements are identified using a repeat or pattern finding analysis of the set of nucleic acids, such as PILER-CR.

[0106] In one embodiment, the grouping step of the method described herein is based, at least in part, on results of a domain homology search or an HHpred protein domain homology search.

[0107] In one embodiment, the defined threshold is a BLAST nearest-neighbor cut-off value of 0 to 1e-7.

[0108] In one embodiment, the method described herein further comprises a filtering step that includes only loci with putative proteins between 900 and 1800 amino acids.

[0109] In one embodiment, the method described herein further comprises experimental validation of the nucleic acid modifying function of the candidate nucleic acid modifying effectors comprising generating a set of nucleic acid constructs encoding the nucleic acid modifying effectors and performing one or more biochemical validation assays, such as through the use of PAM validation in bacterial colonies, in vitro cleavage assays, the Surveyor method, experiments in mammalian cells, PFS validation, or a combination thereof.

[0110] In one embodiment, the method described herein further comprises preparing a non-naturally occurring or engineered composition comprising one or more proteins from the identified nucleic acid modifying loci.

[0111] In one embodiment, the identified loci comprise a Class 2 CRISPR effector, or the identified loci lack Cas1 or Cas2, or the identified loci comprise a single effector.

[0112] In one embodiment, the single large effector protein is greater than 900, or greater than 1100 amino acids in length, or comprises at least one HEPN domain.

[0113] In one embodiment, the at least one HEPN domain is near a N- or C-terminus of the effector protein, or is located in an interior position of the effector protein.

[0114] In one embodiment, the single large effector protein comprises a HEPN domain at the N- and C-terminus and two HEPN domains internal to the protein.

[0115] In one embodiment, the identified loci further comprise one or two small putative accessory proteins within 2 kb to 10 kb of the CRISPR array.

[0116] In one embodiment, the small accessory protein comprises multiple predicted transmembrane domains, or comprises four predicted transmembrane domains, or comprises at least one HEPN domain.

[0117] In one embodiment, the small accessory protein comprises at least one HEPN domain and at least one transmembrane domain.

[0118] In one embodiment, the loci comprise no additional proteins out to 25 kb from the CRISPR array.

[0119] In one embodiment, the CRISPR array comprises direct repeat sequences comprising about 36 nucleotides in length. In a specific embodiment, the direct repeat comprises a GTTG/GUUG at the 5' end that is reverse complementary to a CAAC at the 3' end.

[0120] In one embodiment, the CRISPR array comprises spacer sequences comprising about 30 nucleotides in length.

[0121] In one embodiment, the identified loci lack a small accessory protein.

[0122] Accordingly, it is an object of the invention not to encompass within the invention any previously known product, process of making the product, or method of using the product such that Applicants reserve the right and hereby disclose a disclaimer of any previously known product, process, or method. It is further noted that the invention does not intend to encompass within the scope of the invention any product, process, or making of the product or method of using the product, which does not meet the written description and enablement requirements of the USPTO (35 U.S.C. §112, first paragraph) or the EPO (Article 83 of the EPC), such that Applicants reserve the right and hereby disclose a disclaimer of any previously described product, process of making the product, or method of using the product. It may be advantageous in the practice of the invention to be in compliance with Art. 53(c) EPC and Rule 28(b) and (c) EPC. Nothing herein is to be construed as a promise.

[0123] It is noted that in this disclosure and particularly in the claims and/or paragraphs, terms such as “comprises”, “comprising”, “comprising” and the like can have the meaning attributed to it in U.S. Patent law; e.g., they can mean “includes”, “included”, “including”, and the like; and that terms such as “consisting essentially of” and “consists essentially of” have the meaning ascribed to them in U.S. Patent law, e.g., they allow for elements not explicitly recited, but exclude elements that are found in the prior art or that affect a basic or novel characteristic of the invention.

[0124] These and other embodiments are disclosed or are obvious from and encompassed by, the following Detailed Description.

BRIEF DESCRIPTION OF THE DRAWINGS

[0125] The novel features of the invention are set forth with particularity in the appended claims. A better understanding of the features and advantages of the present invention will be obtained by reference to the following detailed description that sets forth illustrative embodiments, in which the principles of the invention are utilized, and the accompanying drawings of which:

[0126] FIG. 1 shows the bioinformatic pipeline for new CRISPR single effector protein discovery.

[0127] FIG. 2 provides preliminary counts of discovered arrays (all arrays and arrays with single large effectors >700 aa) and large proteins (900-1800 aa) in arrays.

[0128] FIG. 3 provides a protein discovery visualization tool with Cas9 example. (Top left) Select protein accession number, length, and name above top 10 HHpred homology hits and their respective probabilities. (Top right) Tentative CRISPR classifications based on existing literature. (Middle) Genome accession number, length, range, and

name above graph of genomic locus containing CRISPR array in red and proteins in blue (selected protein), magenta (annotated Cas1, Cas 2), orange (other annotated CRISPR associated proteins), and white (all other proteins) with protein orientations depicted. Arrows below the graph dictate predicted direction of transcription (5' to 3'). Numbers above graph represent kilobase distance from selected protein with negative numbers indicating “downstream” and positive numbers indicating “upstream.” (Bottom left) CRISPR arrays in genome with distance from select protein (negative if downstream of protein), number of repeats, length of repeats, length of spacers, and consensus repeat sequence. (Bottom right) Key of genes nearby select proteins with index, accession number, length in amino acids, and name.

[0129] FIGS. 4A-4MMMM show output of bioinformatic pipeline that discovered Group 29 proteins (refer to FIG. 3 legend for explanation of files). Proteins were further discovered by NCBI BLAST and then grouped by homology with a nearest neighbor E-value minimum cutoff of 1e-7 among all proteins. Most proteins have one small accessory protein.

[0130] FIGS. 5A-5E show a list of all currently known Group 29 proteins by protein accession, genome accession, species, and strain of bacteria.

[0131] FIGS. 6A and 6B show a phylogenetic tree of Group 29 proteins from FIGS. 5A-5E.

[0132] FIG. 7 shows a select subset of Group 29 chosen for uniqueness, derived from the listing in FIG. 5.

[0133] FIG. 8 shows a phylogenetic tree of select subset of unique Group 29 proteins.

[0134] FIG. 9 shows a phylogenetic tree of accessory small proteins to select subset of unique Group 29 proteins.

[0135] FIG. 10 shows a phylogenetic tree of small proteins indicating that *Flavobacterium branchiophilum* (Fb) is similar to other small proteins associated with loci that have dual direct repeats.

[0136] FIG. 11 depicts the Fb and *Paludibacter propionici*-genes (Pp) gene loci.

[0137] FIG. 12 shows a comparison of gene loci indicating that small S1a proteins are associated with CRISPR loci with dual repeats.

[0138] FIGS. 13A-13FF show predicted RNA folding of direct repeats of subset of Group 29 proteins selected for uniqueness. Includes putative dual direct repeats for proteins 1, 11, and 14, corresponding to Group 29 proteins with accessory small protein S1a.

[0139] FIGS. 14A-14E show select unique Group 29 direct repeat sequences.

[0140] FIG. 15 shows select unique Group 29 dual direct repeat sequences.

[0141] FIGS. 16A-16C show select unique Group 29 extracted genomic loci, depicting Group 29 protein, repeat region with CRISPR array, and at times accessory small protein.

[0142] FIGS. 17A-17D show the amino acid alignment of all Group 29 proteins.

[0143] FIGS. 18A and 18B show the HEPN domains indicated at N terminus and C terminus of all Group 29 proteins.

[0144] FIG. 19 shows the HEPN domains indicated at N terminus and C terminus of select unique Group 29 proteins (reference numbering in FIG. 7).

[0145] FIG. 20 shows the amino acid alignment of accessory small proteins to select subset of unique Group 29 proteins.

[0146] FIG. 21 shows the analysis of RNA-Seq data for *Bergeyella zoohelcum*, with direction of transcription from left to right.

[0147] FIG. 22 shows a comparison of short DRs and long DRs of RNA transcripts, from which we discovered short crRNAs and long crRNAs (with dual direct repeats.)

[0148] FIG. 23 shows Fb short protein homology.

[0149] FIG. 24 shows the Fb short protein tree.

[0150] FIG. 25 shows the HEPN domains in Group 30 proteins.

[0151] FIG. 26 shows the loci organization of Group 30 proteins.

[0152] FIG. 27 shows the protein alignment of Group 30 proteins.

[0153] FIG. 28 shows the Group 30 protein tree.

[0154] FIG. 29 shows the fold structure of Fb long DR.

[0155] FIG. 30 shows the Fb long DR.

[0156] FIG. 31 shows the fold structure of Fb short DR.

[0157] FIG. 32 shows the Fb short DR.

[0158] FIG. 33 shows the fold structure of Pp short DR.

[0159] FIG. 34 shows the Pp short DR.

[0160] FIGS. 35A and 35B show results of cleavage assay of an EGFP RNA target by Group 29 effector (from *Bergeyella zoohelcum* ATCC 43767) over the indicated time points. FIG. 35A: assay with E1-N2-01; guide. FIG. 35B: assay with E1-T2 guide.

[0161] FIGS. 36A-36D. FIG. 36A shows results of binding assay of Group 29 effector (from *Bergeyella zoohelcum* ATCC 43767) to EGFP target via EMSA at indicated effector protein concentrations. Left panel: wild type Group 29 effector protein. Right panel: mutated Group 29 effector protein (R116A; H121A; R1177A; H1182A). FIG. 36B shows results of cleavage assay of an EGFP RNA target by Group 29 effector (from *Bergeyella zoohelcum* ATCC 43767). Wild type ("WT") and dead ("d") Group 29 enzymes were evaluated with the indicated crRNAs. FIG. 36C shows results of binding assay of mutated Group 29 effector (from *Bergeyella zoohelcum* ATCC 43767) to EGFP target via EMSA at indicated effector protein/RNA ratios. (mutations: R1177A; H1182A). EGFP 1 target, 60 min incubation 50 mM NaCl, 1 mM DTT, Tris-HCl pH7.5, 0.1% BSA. FIG. 36D shows results of cleavage assay of an EGFP RNA target by mutated Group 29 effector (from *Bergeyella zoohelcum* ATCC 43767). Mutations: R1177A, H1182A. EGFP 1 target, 60 min incubation 50 mM NaCl, 1 mM DTT, Tris-HCl pH7.5, 0.1% BSA.

[0162] FIG. 37 shows results of cleavage assay of a single stranded DNA LPAM library (left panel) and single stranded DNA RPAM library (right panel) of Group 29 effector protein (from *Bergeyella zoohelcum* ATCC 43767). bz=targeted crRNA. ns=non-specific crRNA. bzI=targeted crRNA with long DR. nsl=non-specific crRNA with long DR.

[0163] FIGS. 38A and 38B show that Group 29 effector proteins process their own crRNA. FIG. 38A: schematic representation of Group 29 locus (from *Bergeyella zoohelcum* ATCC 43767) and targets used. FIG. 38B: results of cleavage assay of the indicated targets.

[0164] FIG. 39 shows results of cleavage assays of EGFP RNA target 1 (top) and both target 1 and target 2 (bottom) by Group 29 effector (from *Bergeyella zoohelcum* ATCC 43767) with the different guides for the indicated targets.

[0165] FIGS. 40A and 40B show results of cleavage assays of an EGFP RNA target by Group 29 effector (from *Bergeyella zoohelcum* ATCC 43767) with the different guides for the indicated targets.

[0166] FIG. 41 shows EGFP Target sequences (artificial sequences).

[0167] FIGS. 42A and 42B depict computer-simulated secondary structures of EGFP Targets, to determine the effect of secondary structure on cleavage.

[0168] FIGS. 43A and 43B depict the loci of 3 known group 30 effector proteins; FIG. 43B depicts three types of group 30 effector proteins.

[0169] FIG. 44 provides *Bergeyella zoohelcum* group 29 protein structure using PyMOL (a structure software). Blue is the N-Terminus and Red is the C-terminus. Purple dots are the HEPN domains. C-terminus is shown to be free and amenable to linking of functional domains.

[0170] FIGS. 45A and 45B provides crRNAs. The short and long crRNAs are natively processed by BzGrp29. The VCP loop crRNA can be engineered by inserting a viral coat protein (VCP, e.g., MS2, PP7) loop into the direct repeat. The simplest way is to replace the long DR insert in the long crRNA with the VCP. One could also shorten or extend the scaffolding up to the VCP loop, and add additional VCP loops to recruit more domains to the dGRP29. FIG. 45B provides a Grp29 SAM system. The Grp29 SAM system consists of dGrp29 fused to an effector domain(s) (for example at the C terminus), a VCP fused to additional effector domain(s), and the VCP loop crRNA. The three components synergistically act to effect the substrate RNA.

[0171] FIGS. 46A-46C show the discovery of a novel Class 2 CRISPR-Cas system, type VI-B. FIG. 46A: Steps 1 through 5 of a bioinformatic pipeline to discover novel putative Class 2 CRISPR loci without Cas1 or Cas2. FIG. 46B: Abbreviated phylogenetic tree of type VI-B loci. Loci with cas13b (blue) and csx27 (brown) comprise type VI-B1; loci with cas13b (blue) and csx28 (gold) comprise VI-B2. FIG. 46C: HEPN domains in Cas13b and Csx28 from multiple sequence alignments. Top, two HEPN sequences identified via multiple sequence alignment (BLOSUM62) of putative non-redundant Cas13b proteins. Bottom, divergent HEPN sequence identified via multiple sequence alignment (BLOSUM62) of putative non-redundant Csx28 (also referred to as Cas13b-s2) proteins.

[0172] FIG. 47 shows the RNA-sequencing of VI-B CRISPR locus. Left panel: RNA-Sequencing of native VI-B1 locus from *Bergeyella zoohelcum* ATCC 43767. Right panel: RNA-Sequencing of heterologously expressed VI-B2 locus from *P. buccae* ATCC 33574 in *E. coli*.

[0173] FIG. 48 depicts inhibition of MS2 plaque development by *B. zoohelcum* (S1, left panel) or *P. buccae* (S2, right panel) targeted to MS2.

[0174] FIG. 49 depicts the structure of a Cas13b crRNA and shows the arrangement flanking a protospacer of the left protospacer flanking sequence (LPFS) and right protospacer flanking sequence (RPFS).

[0175] FIG. 50 depicts the experimental design used to create and identify GRP29 effector-dependent depletions of spacer libraries.

[0176] FIG. 51 depicts guide abundance cross-correlations for BZ between spacer libraries in the absence of BZ effector protein (top), among identical experiments (middle), and between experiment and control (bottom).

[0177] FIG. 52 depicts guide abundance cross-correlations for PB between spacer libraries in the absence of PB effector protein (top), among identical experiments (middle), and between experiment and control (bottom).

[0178] FIG. 53 depicts targeting guide depletions with BZ in alternating colors by gene, non-targeting in black. A red horizontal line indicates the safely depleted line.

[0179] FIG. 54 depicts targeting guide depletions with PB in alternating colors by gene, non-targeting in black. A red horizontal line indicates the safely depleted line.

[0180] FIG. 55 depicts depletion levels of intersecting guides common to both BZ and PB.

[0181] FIG. 56 depicts the aggregate depletion normalized across each gene for BZ.

[0182] FIG. 57 depicts the aggregate depletion normalized across each gene for PB.

[0183] FIG. 58 depicts sequence logos for left protospacer flanking sequence (LPFS) and right protospacer flanking sequence (RPFS) for BZ (top) and PB (bottom).

[0184] FIG. 59 shows a protospacer flanking sequence (PFS) wheel illustrating nucleotide frequencies for L1, R2, and R3 for BZ.

[0185] FIG. 60 shows a protospacer flanking sequence (PFS) wheel illustrating nucleotide frequencies for L1, R2, and R3 for PB.

[0186] FIG. 61 depicts mean depletions for different combination of nucleotides at L1, R2, and R3. In each panel, L1 is fixed. Left side panels show mean depletions with BZ. Right side panels show depletions with PB.

[0187] FIG. 62 shows PB targeting of ampicillin and kanamycin resistance genes.

[0188] FIG. 63 shows BZ targeting of ampicillin and kanamycin resistance genes.

[0189] FIGS. 64A-64G show heterologous expression of *B. zoohelcum* VI-B locus mediates knockdown of *E. coli* essential genes. FIG. 64A: Design of *E. coli* essential gene screen to determine targeting rules of nucleic acid interference. FIG. 64B: Manhattan plots of mean spacer depletions mapped over 45 genes and aggregated across normalized gene distance for full *B. zoohelcum* VI-B1 locus (left) and *cas13b* (right), with non-targeting spacers in gray, safely depleted spacers ($>5\sigma$ above mean depletion of non-targeting spacers) above blue line, and strongly depleted spacers (top 1% depleted) above red line. We performed three independent bioreplicates of each condition of the screen and quality-control-filtered the mean depletions of spacers by both a maximum coefficient of variation of 0.2 and a minimum abundance in the control screen of $\frac{1}{3}N$, where $N=55,700$. This step reduced the total number of spacers included in the analysis to $N' \sim 35,000-40,000$, depending on the screen. FIG. 64C: Weblogo of sequence motifs of strongly depleted *B. zoohelcum* spacers. FIG. 64D: Normalized PFS score matrix, where each score is the ratio of number of safely depleted *B. zoohelcum* spacers to total number of spacers for a given PFS, scaled so that maximum PFS score is 1. FIG. 64E: Spacers targeting kanamycin to validate PFS targeting rules of 5' PFS (D) and 3' PFS (N followed by AN or NA). FIG. 64F: Results of kanamycin validation screen for *B. zoohelcum cas13b* in *E. coli*. FIG. 64G: Spacer abundances versus control for individual *B. zoohelcum* spacers, with abundances colored by type of spacer.

[0190] FIGS. 65A-65D depict that Cas13b interferes with RNA, but not DNA, through HEPN domains. FIG. 65A:

Protospacer design for MS2 phage plaque drop assay to test RNA interference. FIG. 65B: Plaque drop assay for full *B. zoohelcum* VI-B1 locus (left) and *bzcas13b* (right) FIG. 65C: DNA interference assay schematic and results. A target sequence is placed in frame at the start of the transcribed bla gene that confers ampicillin resistance or in a non-transcribed region of the same target plasmid. Target plasmids were co-transformed with *bzcas13b* plasmid or empty vectors conferring chloramphenicol resistance and plated on double selection antibiotic plates. FIG. 65D: Quantification of MS2 phage plaque drop assay with *B. zoohelcum Cas13b* HEPN mutants.

[0191] FIGS. 66A-66D show efficient RNA targeting by Cas13b correlated with local RNA accessibility. FIG. 66A: Methodology of secondary structure-mediated spacer efficiency analysis of *E. coli* essential gene screen data with Vienna RNAfold. FIG. 66B: Optimization of top 1 accuracy (computationally predicted best spacer matches the experimentally top depleted spacer) and top 3 accuracy (computationally predicted top spacer falls in top 3 depleted spacers) on randomly selected *B. zoohelcum* training dataset using RNAfold, first with u start and u end, and then with W and L. FIG. 66C: Performance of optimized RNAfold model on randomly selected *B. zoohelcum* testing dataset (48 cohorts for full *B. zoohelcum* VI-B1 locus, 56 cohorts for *bzcas13b*) against 10^6 Monte Carlo simulations: empirical P-values from left to right of $3e-6$, $1e-6$, $8.7e-3$, $6e-6$. FIG. 66D: Empirical cumulative distribution function of safely depleted *B. zoohelcum* spacers over all genes from 5' UTR into gene and from 3' UTR into gene. Yellow line separates UTR and gene, red line is theoretical cumulative distribution function of uniformly distributed spacers, and blue line is empirical cumulative distribution of safely depleted *B. zoohelcum* spacers.

[0192] FIGS. 67A-67E demonstrate Class 2 type VI-B systems differentially modulated by Csx27 and Csx28. FIG. 67A: Normalized PFS matrix, for *P. buccae* VI-B2 locus. FIG. 67B: MS2 Plaque drop assay for full *P. buccae* VI-B2 locus (left) and *pbcas13b* (right). FIG. 67C: Spacer depletions of *bzcas13b* with and without *bzcsx27* (brown), as compared to *pbcas13b* with and without *pbcxs28* (gold). FIG. 67D: A sequence, structure, and spatial RNA targeting model for Cas13b. (On our website is a downloadable *bzcas13b* target design script with instructions for researchers, as well as a protocol for performing the *E. coli* essential gene screen on another CRISPR system.) FIG. 67E: A bimodal functional model for the VI-B system, with Csx27 repressing and Csx28 enhancing Cas13b-mediated RNA interference.

[0193] FIG. 68 shows a phylogenetic tree of Cas13b bifurcated into two type VI-B CRISPR loci. A phylogenetic tree (alignment generated by BLOSUM62) of non-redundant Cas13b effectors, with the full type VI-B locus depicted in every instance. Accession numbers for genome, Cas13b (blue), and Csx27 (brown)/Csx28 (gold) are included, as well as number of nearby spacers detected by PILER-CR, the presence of other CRISPR-Cas elements in the genome, and the size of Cas13b are found in Table 1.

[0194] FIGS. 69A and 69B reveal that the predicted secondary structure of type VI-B direct repeats is well-conserved. FIG. 69A: Predicted secondary structure folds of structurally unique CRISPR class 2 type VI-B1 direct

repeats (Vienna RNAfold). FIG. 69B: Predicted secondary structure folds of structurally unique CRISPR Class 2 type VI-B2 direct repeats.

[0195] FIGS. 70A and 70B show the well-conserved Type VI-B direct repeats; predicted protospacer flanking sequences are inconclusive. FIG. 70A: Weblogo of all unique VI-B direct repeat sequences of length 36 nt, taken as the same transcriptional orientation as Cas13b. FIG. 70B: Weblogo of all unique VI-B protospacer flanking sequences from CRISPRtarget mapping of protospacers to phage and plasmid databases.

[0196] FIGS. 71A and 71B show Predicted transmembrane domains of Csx27 and Csx28 not validated experimentally. FIG. 71A: Transmembrane domain prediction in Csx27 of *B. zoohelcum* and Csx28 of *P. buccae* using TMHMM v2. FIG. 71B: N- and C-terminally fused RFP imaging of Csx27 of *B. zoohelcum* and Csx28 of *P. buccae*.

[0197] FIGS. 72A and 72B show agreement of second bioreplicate of kanamycin validation screen for Cas13b from *B. zoohelcum* to the first bioreplicate. FIG. 72A: Spacers targeting kanamycin to validate PFS targeting rules of 5' PFS (D) and 3'PFS (N followed by AN or NA).

[0198] FIG. 72B: Spacer abundances versus control for individual *B. zoohelcum* spacers, with abundances colored by type of spacer.

[0199] FIG. 73 demonstrates bioreplicates of MS2 phage plaque drop assay. Plaque drop assay with bioreplicates for *B. zoohelcum* VI-B1 locus and cas13b, and for *P. buccae* VI-B2 locus and cas13b.

[0200] FIG. 74 demonstrates bioreplicates of HEPN mutant plaque drop assay. Plaque drop assay with bioreplicates for *B. zoohelcum* Cas13b HEPN mutants (R116A/H121A and R1177A/H1182A) versus wildtype Cas13b.

[0201] FIGS. 75A and 75B show *E. coli* essential gene screen of *P. buccae* VI-B2 CRISPR locus. FIG. 75A: Manhattan plots of spacer depletions mapped over 45 genes and aggregated across normalized gene distance for full *P. buccae* VI-B2 locus (left) and cas13b (right), with non-targeting spacers in gray, safely depleted ($>5\sigma$ above mean depletion of non-targeting spacers) spacers above blue line, and strongly depleted (top 1% depleted) spacers above red line.

[0202] FIG. 75B: Sequence weblogs of strongly depleted *P. buccae* spacers, revealing double-sided PFS (protospacer flanking sequence).

[0203] FIGS. 76A-76C show computational analysis of secondary structure and spatial rules for *P. buccae* cas13b RNA targeting. FIG. 76A: Optimization of top 1 accuracy (computationally predicted spacer is top depleted) and top 3 accuracy (computationally predicted spacer falls in top 3 depleted) on randomly selected *P. buccae* training dataset using RNAplfold, first with u start and u end, and then with W and L. FIG. 76B: Performance of optimized RNAplfold model on randomly selected *P. buccae* testing dataset (41 cohorts for full *P. buccae* VI-B2 locus, 40 cohorts for pbcas13b) against 10^6 Monte Carlo simulations: empirical P-values from left to right of 3.3e-2, 2.7e-3, 3.9e-3, 1.5e-5. FIG. 76C: Empirical cumulative distribution function of safely depleted *P. buccae* spacers over all genes from 5'UTR into gene and from 3' UTR into gene. Yellow line separates UTR and gene, red line is theoretical cumulative distribution function of uniformly distributed spacers, and blue line is empirical cumulative distribution of safely depleted *P. buccae* spacers.

[0204] FIG. 77 shows cleavage activity of target sequence by Cas13b with crRNAs having different spacer sequence lengths, ranging from 30 nucleotides or less. As shown, the spacer length of 27-30 nucleotides works for optimal cutting. At a spacer length of 26-22 nucleotides, the cutting is reduced, and at a length of 20 nucleotides or less, there was no cutting activity of the target sequence.

[0205] FIG. 78 is a schematic of an example Type VI-B locus.

[0206] FIG. 79 is a schematic of an example large effector protein from a Type VI-B locus.

[0207] FIG. 80 is a graph showing the size distribution of putative large effector proteins identified in candidate Type VI-B locus.

[0208] FIG. 81 is an alignment of putative large effector proteins from candidate Type VI-B locus.

[0209] FIG. 82 is a schematic workflow from putative effector discovery to potential genomic engineering applications.

[0210] The figures herein are for illustrative purposes only and are not necessarily drawn to scale.

DETAILED DESCRIPTION OF THE INVENTION

[0211] In general, the CRISPR-Cas or CRISPR system refers collectively to transcripts and other elements involved in the expression of or directing the activity of CRISPR-associated (“Cas”) genes, including sequences encoding a Cas gene, a tracr (trans-activating CRISPR) sequence (e.g. tracrRNA or an active partial tracrRNA), a tracr-mate sequence (encompassing a “direct repeat” and a tracrRNA-processed partial direct repeat in the context of an endogenous CRISPR system), a guide sequence (also referred to as a “spacer” in the context of an endogenous CRISPR system), or “RNA(s)” as that term is herein used (e.g., RNA(s) to guide Cas, such as Cas9, e.g. CRISPR RNA and transactivating (tracr) RNA or a single guide RNA (sgRNA) (chimeric RNA)) or other sequences and transcripts from a CRISPR locus. In general, a CRISPR system is characterized by elements that promote the formation of a CRISPR complex at the site of a target sequence (also referred to as a protospacer in the context of an endogenous CRISPR system).

[0212] When the CRISPR protein is a Class 2 Type VI-B effector, a tracrRNA is not required. In an engineered system of the invention, the direct repeat may encompass naturally-occurring sequences or nonnaturally-occurring sequences. The direct repeat of the invention is not limited to naturally occurring lengths and sequences. A direct repeat can be 36 nt in length, but a longer or shorter direct repeat can vary. For example, a direct repeat can be 30 nt or longer, such as 30-100 nt or longer. For example, a direct repeat can be 30 nt, 40 nt, 50 nt, 60 nt, 70 nt, 80 nt, 90 nt, 100 nt, or longer in length. In some embodiments, a direct repeat of the invention can include synthetic nucleotide sequences inserted between the 5' and 3' ends of naturally occurring direct repeat. In certain embodiments, the inserted sequence may be self-complementary, for example, 20%, 30%, 40%, 50%, 60%, 70%, 80%, 90%, or 100% self complementary. Furthermore, a direct repeat of the invention may include insertions of nucleotides such as an aptamer or sequences that bind to an adapter protein (for association with functional domains). In certain embodiments, one end of a direct

repeat containing such an insertion is roughly the first half of a short DR and the other end is roughly the second half of the short DR.

[0213] In the context of formation of a CRISPR complex, “target sequence” refers to a sequence to which a guide sequence is designed to have complementarity, where hybridization between a target sequence and a guide sequence promotes the formation of a CRISPR complex. A target sequence may comprise any polynucleotide, such as DNA or RNA polynucleotides. In some embodiments, a target sequence is located in the nucleus or cytoplasm of a cell. In some embodiments, direct repeats may be identified in silico by searching for repetitive motifs that fulfill any or all of the following criteria: 1. found in a 2 Kb window of genomic sequence flanking the type II CRISPR locus; 2. span from 20 to 50 bp; and 3. interspaced by 20 to 50 bp. In some embodiments, 2 of these criteria may be used, for instance 1 and 2, 2 and 3, or 1 and 3. In some embodiments, all 3 criteria may be used.

[0214] In embodiments of the invention the terms guide sequence and guide RNA, i.e. RNA capable of guiding Type VI-B effector proteins, such as Cas13b and Group 29 or Group 30 proteins to a target locus, are used interchangeably as in foregoing cited documents such as WO 2014/093622 (PCT/US2013/074667). In general, a guide sequence (or spacer sequence) is any polynucleotide sequence having sufficient complementarity with a target polynucleotide sequence to hybridize with the target sequence and direct sequence-specific binding of a CRISPR complex to the target sequence. In some embodiments, the degree of complementarity between a guide sequence and its corresponding target sequence, when optimally aligned using a suitable alignment algorithm, is about or more than about 50%, 60%, 75%, 80%, 85%, 90%, 95%, 97.5%, 99%, or more. Optimal alignment may be determined with the use of any suitable algorithm for aligning sequences, non-limiting example of which include the Smith-Waterman algorithm, the Needleman-Wunsch algorithm, algorithms based on the Burrows-Wheeler Transform (e.g. the Burrows Wheeler Aligner), ClustalW, Clustal X, BLAT, Novoalign (Novocraft Technologies; available at www.novocraft.com), ELAND (Illumina, San Diego, Calif.), SOAP (available at soap.genomics.org.cn), and Maq (available at maq.sourceforge.net). In some embodiments, a guide sequence (or spacer sequence) is about or more than about 5, 10, 11, 12, 13, 14,

15, 16, 17, 18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30, 35, 40, 45, 50, 75, or more nucleotides in length. In some embodiments, a guide sequence is less than about 75, 50, 45, 40, 35, 30, 25, 20, 15, 12, or fewer nucleotides in length. Preferably the guide sequence is 10-40 nucleotides long, such as 20-30 or 20-40 nucleotides long or longer, such as 30 nucleotides long or about 30 nucleotides long. In certain embodiments, the guide sequence is 10-30 nucleotides long, such as 20-30 or 20-40 nucleotides long or longer, such as 30 nucleotides long or about 30 nucleotides long for Group 29 or Group 30 effectors. In certain embodiments, the guide sequence is 10-30 nucleotides long, such as 20-30 nucleotides long, such as 30 nucleotides long or about 30 nucleotides long for Group 29 effectors originating from *Bergeyella zoohelcum* (such as *Bergeyella zoohelcum* ATCC 43767). The ability of a guide sequence to direct sequence-specific binding of a CRISPR complex to a target sequence may be assessed by any suitable assay. For example, the components of a CRISPR system sufficient to form a CRISPR complex, including the guide sequence to be tested, may be provided to a host cell having the corresponding target sequence, such as by transfection with vectors encoding the components of the CRISPR sequence, followed by an assessment of preferential cleavage within the target sequence, such as by Surveyor assay as described herein. Similarly, cleavage of a target polynucleotide sequence may be evaluated in a test tube by providing the target sequence, components of a CRISPR complex, including the guide sequence to be tested and a control guide sequence different from the test guide sequence, and comparing binding or rate of cleavage at the target sequence between the test and control guide sequence reactions. Other assays are possible, and will occur to those skilled in the art.

[0215] The instant invention provides Type VI-B CRISPR-Cas effectors, nucleic acids, systems, vectors, and methods of use. Type VI-B CRISPR-Cas effectors and nucleic acids encompass both Group 29 and Group 30. All Group VI-B effectors are distinguishable from VI-A by structure, and also by the location of the HEPN domains). As used herein, the terms Cas13b-s1 accessory protein, Cas13b-s1 protein, Cas13b-s1, Csx27, and Csx27 protein are used interchangeably and the terms Cas13b-s2 accessory protein, Cas13b-s2 protein, Cas13b-S2, Csx28, and Csx28 protein are used interchangeably.

[0216] Group VI-B CRISPR-Cas effectors include the following examples:

TABLE 1

Representative Type VI-B Effectors and Accessory Proteins							
Species (Genome Accession)	Cas13b Accession	Csx27/28 Accession	# Spacers	CRISPR-Cas?	Cas1?	Cas2?	Cas13b size (aa)
<i>Paludibacter propionigenes</i> WB4 (NC_014734.1)	WP_013446107.1	NA	8	N	N	N	1155
<i>Prevotella</i> sp. P5-60 (NZ_IJXQJ01000080.1)	WP_044074780.1	NA	5	Y	?	?	1091
<i>Prevotella</i> sp. P4-76 (NZ_IJXQI01000021.1)	WP_044072147.1	NA	0	?	?	?	1091
<i>Prevotella</i> sp. P5-125 (NZ_IJXQL01000055.1)	WP_044065294.1	NA	11	?	?	?	1091
<i>Prevotella</i> sp. P5-119 (NZ_IJXQK01000043.1)	WP_042518169.1	NA	11	?	?	?	1091
<i>Capnocytophaga canimorsus</i> Cc5 (NC_015846.1)	WP_013997271.1	WP_013997274.1	51	Y	Y	Y	1200

TABLE 1-continued

Representative Type VI-B Effectors and Accessory Proteins							
Species (Genome Accession)	Cas13b Accession	Csx27/28 Accession	# Spacers	CRISPR- Cas?	Cas1?	Cas2?	Cas13b size (aa)
<i>Phaeodactylibacter xiamenensis</i> (NZ_JPOS01000018.1)	WP_044218239.1	WP_044218241.1	19	?	?	?	1132
<i>Porphyromonas gingivalis</i> W83 (NC_002950.2)	WP_005873511.1	WP_005873518.1	7	Y	Y	Y	1136
<i>Porphyromonas gingivalis</i> F0570 (NZ_KI259168.1)	WP_021665475.1	WP_021665476.1	3	?	?	?	1136
<i>Porphyromonas gingivalis</i> ATCC 33277 (NC_010729.1)	WP_012458151.1	WP_012458152.1	12	Y	Y	Y	1136
<i>Porphyromonas gingivalis</i> F0185 (AWVC01000122.1)	ERJ81987.1	ERJ81988.1	0	?	?	?	1136
<i>Porphyromonas gingivalis</i> F0185 (NZ_KI259960.1)	WP_021677657.1	WP_021677658.1	6	?	?	?	1136
<i>Porphyromonas gingivalis</i> SJD2 (NZ_KI629875.1)	WP_023846767.1	WP_005873518.1	4	?	?	?	1136
<i>Porphyromonas gingivalis</i> F0568 (AWUU01000145.1)	ERJ65637.1	ERJ65638.1	3	?	?	?	1136
<i>Porphyromonas gingivalis</i> W4087 (AWVE01000130.1)	ERJ87335.1	ERJ87336.1	2	?	?	?	1136
<i>Porphyromonas gingivalis</i> W4087 (NZ_KI260263.1)	WP_021680012.1	WP_005873518.1	4	?	?	?	1136
<i>Porphyromonas gingivalis</i> F0568 (NZ_KI258981.1)	WP_021663197.1	WP_021663198.1	6	?	?	?	1136
<i>Porphyromonas gingivalis</i> (NZ_LOEL01000010.1)	WP_061156637.1	WP_005873518.1	11	?	?	?	1136
<i>Porphyromonas gulae</i> (NZ_JRAQ01000019.1)	WP_039445055.1	WP_039445052.1	10	?	?	?	1136
<i>Bacteroides pyogenes</i> F0041 (KE993153.1)	ERI81700.1	ERI81699.1	5	?	?	?	1116
<i>Bacteroides pyogenes</i> JCM 10003 (NZ_BAIU0100000 1.1)	WP_034542281.1	WP_034542279.1	18	?	?	?	1116
<i>Alistipes</i> sp. ZOR0009 (NZ_JTLD01000029.1)	WP_047447901.1	NA	7	?	?	?	954
<i>Flavobacterium branchiophilum</i> FL-15 (NC_016001.1)	WP_014084666.1	WP_014084665.1	19	Y	N	Y	1151
<i>Prevotella</i> sp. MA2016 (NZ_JHUW01000010.1)	WP_036929175.1	NA	7	?	?	?	1323
<i>Myroides odoratimimus</i> CCUG 10230 (AGEC02000017.1)	EHO06562.1	EHO06560.1	2	?	?	?	1160
<i>Myroides odoratimimus</i> CCUG 3837 (AGZK01000016.1)	EKB06014.1	EKB06015.1	0	?	?	?	1158
<i>Myroides odoratimimus</i> CCUG 3837 (NZ_JH815535.1)	WP_006265509.1	WP_006265510.1	0	?	?	?	1158
<i>Myroides odoratimimus</i> CCUG 12901 (NZ_JH590834.1)	WP_006261414.1	WP_006261415.1	0	?	?	?	1158
<i>Myroides odoratimimus</i> CCUG 12901 (AGED01000033.1)	EHO08761.1	EHO08762.1	0	?	?	?	1158
<i>Myroides odoratimimus</i> (NZ_CP013690.1)	WP_058700060.1	WP_006261415.1	10	Y	Y	Y	1160
<i>Bergeyella zoohelcum</i> ATCC 43767 (AGYA01000037.1)	EKB54193.1	EKB54194.1	9	?	?	?	1225
<i>Capnocytophaga cynodegmi</i> (NZ_CDOD01000002.1)	WP_041989581.1	WP_041989578.1	7	?	?	?	1219
<i>Bergeyella zoohelcum</i> ATCC 43767 (NZ_JH932293.1)	WP_002664492.1	WP_034985946.1	8	Y	Y	Y	1225
<i>Flavobacterium</i> sp. 316 (NZ_JYGZ01000003.1)	WP_045968377.1	NA	0	?	?	?	1156
<i>Psychroflexus torquis</i> ATCC 700755 (NC_018721.1)	WP_015024765.1	NA	16	Y	Y	Y	1146
<i>Flavobacterium columnare</i> ATCC 49512 (NC_016510.2)	WP_014165541.1	NA	7	Y	Y	Y	1180
<i>Flavobacterium columnare</i> (NZ_CP013992.1)	WP_060381855.1	NA	5	Y	Y	Y	1214
<i>Flavobacterium columnare</i> (NZ_CP015107.1)	WP_063744070.1	NA	3	Y	Y	Y	1214
<i>Flavobacterium columnare</i> (NZ_CP016277.1)	WP_065213424.1	NA	14	Y	Y	Y	1215
<i>Chryseobacterium</i> sp. YR477 (NZ_KN549099.1)	WP_047431796.1	NA	0	?	?	?	1146
<i>Riemerella anatipestifer</i> ATCC 11845 = DSM 15868 (NC_014738.1)	WP_004919755.1	WP_004919758.1	12	Y	Y	Y	1096

TABLE 1-continued

Representative Type VI-B Effectors and Accessory Proteins							
Species (Genome Accession)	Cas13b Accession	Csx27/28 Accession	# Spacers	CRISPR- Cas?	Cas1?	Cas2?	Cas13b size (aa)
<i>Riemerella anatipestifer</i> RA-CH-2 (NC_020125.1)	WP_015345620.1	WP_004919758.1	12	Y	Y	Y	949
<i>Riemerella anatipestifer</i> (NZ_CP007504.1)	WP_049354263.1	WP_004919758.1	11	Y	Y	Y	949
<i>Riemerella anatipestifer</i> (NZ_LUDU01000012.1)	WP_061710138.1	WP_061710139.1	13	?	?	?	951
<i>Riemerella anatipestifer</i> (NZ_LUDJ01000010.1)	WP_064970887.1	WP_064970885.1	4	?	?	?	1096
<i>Prevotella saccharolytica</i> F0055 (AMEP01000091.1)	EKY00089.1	EKY00090.1	0	?	?	?	1151
<i>Prevotella saccharolytica</i> JCM 17484 (NZ_BAKN01000001.1)	WP_051522484.1	NA	5	Y	Y	Y	1152
<i>Prevotella buccae</i> ATCC 33574 (AEPD01000005.1)	EFU31981.1	EFU31982.1	16	?	?	?	1128
<i>Prevotella buccae</i> ATCC 33574 (NZ_GL586311.1)	WP_004343973.1	WP_004343974.1	16	Y	Y	Y	1128
<i>Prevotella buccae</i> D17 (NZ_GG739967.1)	WP_004343581.1	WP_004343582.1	8	?	?	?	1128
<i>Prevotella</i> sp. MSX73 (NZ_ALJQ01000043.1)	WP_007412163.1	WP_036927782.1	13	?	?	?	1128
<i>Prevotella pallens</i> ATCC 700821 (AFPY01000052.1)	EGQ18444.1	EGQ18443.1	4	?	?	?	1126
<i>Prevotella pallens</i> ATCC 700821 (NZ_GL982513.1)	WP_006044833.1	WP_050795200.1	4	?	?	?	1126
<i>Prevotella intermedia</i> ATCC 25611 = DSM 20706 (NZ_JAEZ01000017.1)	WP_036860899.1	WP_050795200.1	11	?	?	?	1127
<i>Prevotella intermedia</i> (NZ_LBGT01000010.1)	WP_061868553.1	NA	27	?	?	?	1121
<i>Prevotella intermedia</i> 17 (CP003502.1)	AFJ07523.1	AFJ07898.1	16	N	N	N	1135
<i>Prevotella intermedia</i> (NZ_AP014926.1)	WP_050955369.1	WP_014708440.1	16	N	N	N	1133
<i>Prevotella intermedia</i> (AP014598.1)	BAU18623.1	BAU18624.1	6	N	N	N	1134
<i>Prevotella intermedia</i> ZT (ATMK01000017.1)	KJJ86756.1	KJJ86755.1	2	?	?	?	1126
<i>Prevotella aurantiaca</i> JCM 15754 (NZ_BAKF01000019.1)	WP_025000926.1	WP_036889078.1	5	?	?	?	1125
<i>Prevotella pleuritidis</i> F0068 (NZ_AWET01000045.1)	WP_021584635.1	WP_021584705.1	6	?	?	?	1140
<i>Prevotella pleuritidis</i> JCM 14110 (NZ_BAJN01000005.1)	WP_036931485.1	WP_024991772.1	7	?	?	?	1117
<i>Prevotella falsenii</i> DSM 22864 = JCM 15124 (NZ_BAJY01000004.1)	WP_036884929.1	WP_051527348.1	10	?	?	?	1134
<i>Porphyromonas gulae</i> (NZ_JRAT01000012.1)	WP_039418912.1	WP_052073447.1	11	Y	Y	Y	1176
<i>Porphyromonas</i> sp. COT-052 OH4946 (NZ_JQZY01000014.1)	WP_039428968.1	WP_050563578.1	12	?	?	?	1176
<i>Porphyromonas gulae</i> (NZ_JRFD01000046.1)	WP_039442171.1	WP_050563578.1	9	?	?	?	1175
<i>Porphyromonas gulae</i> (NZ_JRAJ01000010.1)	WP_039431778.1	WP_046201041.1	2	?	?	?	1176
<i>Porphyromonas gulae</i> (NZ_KQ040500.1)	WP_046201018.1	WP_046201041.1	4	?	?	?	1176
<i>Porphyromonas gulae</i> (NZ_JRAL01000022.1)	WP_039434803.1	WP_039434800.1	20	?	?	?	1176
<i>Porphyromonas gulae</i> (NZ_JRAI01000002.1)	WP_039419792.1	WP_052078041.1	9	?	?	?	1120
<i>Porphyromonas gulae</i> (NZ_JRAK01000129.1)	WP_039426176.1	WP_039426172.1	6	?	?	?	1120
<i>Porphyromonas gulae</i> (NZ_KN294104.1)	WP_039437199.1	WP_052102013.1	0	?	?	?	1120
<i>Porphyromonas gingivalis</i> TDC60 (NC_015571.1)	WP_013816155.1	WP_043890185.1	2	Y	Y	Y	1120
<i>Porphyromonas gingivalis</i> ATCC 33277 (NC_010729.1)	WP_012458414.1	WP_012458413.1	4	Y	Y	Y	1120
<i>Porphyromonas gingivalis</i> A7A1- 28 (NZ_CP013131.1)	WP_058019250.1	WP_043898408.1	6	Y	Y	Y	1176

TABLE 1-continued

Representative Type VI-B Effectors and Accessory Proteins							
Species (Genome Accession)	Cas13b Accession	Csx27/28 Accession	# Spacers	CRISPR- Cas?	Cas1?	Cas2?	Cas13b size (aa)
<i>Porphyromonas gingivalis</i> JCVI SC001 (APMB01000175.1)	EOA10535.1	EOA10563.1	5	?	?	?	1176
<i>Porphyromonas gingivalis</i> W50 (NZ_AJZS01000051.1)	WP_005874195.1	WP_010955981.1	2	?	?	?	1176
<i>Porphyromonas gingivalis</i> (NZ_CP011995.1)	WP_052912312.1	WP_010955981.1	7	Y	Y	Y	1176
<i>Porphyromonas gingivalis</i> AJW4 (NZ_CP011996.1)	WP_053444417.1	WP_043898408.1	11	N	N	N	1120
<i>Porphyromonas gingivalis</i> (NZ_CP007756.1)	WP_039417390.1	WP_021665928.1	5	Y	Y	Y	1120
<i>Porphyromonas gingivalis</i> (NZ_LOEL01000001.1)	WP_061156470.1	WP_021663076.1	5	?	?	?	1120

[0217] In a classic CRISPR-Cas systems, the degree of complementarity between a guide sequence and its corresponding target sequence can be about or more than about 50%, 60%, 75%, 80%, 85%, 90%, 95%, 97.5%, 99%, or 100%; a guide or RNA or sgRNA can be about or more than about 5, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30, 35, 40, 45, 50, 75, or more nucleotides in length; or guide or RNA or sgRNA can be less than about 75, 50, 45, 40, 35, 30, 25, 20, 15, 12, or fewer nucleotides in length; and advantageously tracr RNA is 30 or 50 nucleotides in length. However, an aspect of the invention is to reduce off-target interactions, e.g., reduce the guide interacting with a target sequence having low complementarity. Indeed, in the examples, it is shown that the invention involves mutations that result in the CRISPR-Cas system being able to distinguish between target and off-target sequences that have greater than 80% to about 95% complementarity, e.g., 83%-84% or 88-89% or 94-95% complementarity (for instance, distinguishing between a target having 18 nucleotides from an off-target of 18 nucleotides having 1, 2 or 3 mismatches). Accordingly, in the context of the present invention the degree of complementarity between a guide sequence and its corresponding target sequence is greater than 94.5% or 95% or 95.5% or 96% or 96.5% or 97% or 97.5% or 98% or 98.5% or 99% or 99.5% or 99.9%, or 100%. Off target is less than 100% or 99.9% or 99.5% or 99% or 99% or 98.5% or 98% or 97.5% or 97% or 96.5% or 96% or 95.5% or 95% or 94.5% or 94% or 93% or 92% or 91% or 90% or 89% or 88% or 87% or 86% or 85% or 84% or 83% or 82% or 81% or 80% complementarity between the sequence and the guide, with it advantageous that off target is 100% or 99.9% or 99.5% or 99% or 99% or 98.5% or 98% or 97.5% or 97% or 96.5% or 96% or 95.5% or 95% or 94.5% complementarity between the sequence and the guide.

[0218] In particularly preferred embodiments according to the invention, the guide RNA (capable of guiding Cas to a target locus) may comprise (1) a guide sequence capable of hybridizing to a target locus (a polynucleotide target locus, such as an RNA target locus) in the eukaryotic cell; (2) a direct repeat (DR) sequence which reside in a single RNA, i.e. an sgRNA (arranged in a 5' to 3' orientation). In particular embodiments, the CRISPR/Cas protein is characterized in that it makes use of a guide RNA comprising a guide sequence capable of hybridizing to a genomic target locus in the eukaryotic cell and a direct repeat sequence, and

does not require a tracrRNA. In particular embodiments, where the CRISPR/Cas protein is characterized in that it makes use of a tracrRNA, the guide sequence, tracr mate and tracr sequence may reside in a single RNA, i.e. an sgRNA (arranged in a 5' to 3' orientation), or the tracr RNA may be a different RNA than the RNA containing the guide and tracr mate sequence. In these embodiments, the tracr hybridizes to the tracr mate sequence and directs the CRISPR/Cas complex to the target sequence.

[0219] As used herein, "database" refers to a repository of genomic sequence information. Non-limiting examples of include the NCBI database and the European Ensembl database.

[0220] As used herein, "CRISPR effector" or "effector" refers to an RNA-guided DNA-targeting polypeptide or an RNA-guided RNA-targeting polypeptide possessing enzymatic activity. Non-limiting examples of enzymatic activity include endonuclease, nickase, integrase, or transposase activity.

[0221] As used herein, "nuclease domain" refers to a protein domain capable of cleaving RNA, DNA, or both. "Cleaving" refers to the breaking of the covalent backbone of one or more strands of a target polynucleotide.

[0222] As used herein, the term "CRISPR array" refers to the DNA segment which includes all of the CRISPR repeats and spacers, starting with the first nucleotide of the first CRISPR repeat and ending with the last nucleotide of the last (terminal) CRISPR repeat. Typically, each spacer sequence in a CRISPR array is located between two repeats.

[0223] As used herein, the terms "CRISPR repeat," "direct repeat," "repeat sequence," or "repeat" have the conventional meaning as used in the art, i.e., multiple short direct repeating sequences, which show very little or no sequence variation within a given CRISPR array. As used herein, "CRISPR spacer," "spacer sequence," or "spacer" refer to the non-repetitive sequences that are located between the repeats of a CRISPR array.

[0224] A person of skill in art would use experimental techniques well known in the art to characterize an effector system identified by the method disclosed in the current invention. For example, the components of a CRISPR system sufficient to form a CRISPR complex, including the CRISPR effector and RNA guide sequence to be tested, may be provided to a host cell having a corresponding target sequence, such as by transfection with vectors encoding the components of the CRISPR system, followed by an assess-

ment of preferential cleavage within the target sequence, such as by Surveyor assay. Similarly, cleavage of a target polynucleotide sequence may be evaluated in a test tube by providing the target sequence, components of a CRISPR complex, including the guide sequence to be tested and a control guide sequence different from the test guide sequence, and comparing binding or rate of cleavage at the target sequence between the test and control guide sequence reactions. Other assays are possible, and will occur to those skilled in the art.

[0225] In particular embodiments, the wildtype group29/group30 effector protein has RNA binding and cleaving function.

[0226] In particular embodiments, the group29/group30 effector protein may have DNA cleaving function. In these embodiments, methods may be provided based on the effector proteins provided herein which comprehend inducing one or more mutations in a eukaryotic cell (in vitro, i.e. in an isolated eukaryotic cell) as herein discussed comprising delivering to cell a vector as herein discussed. The mutation (s) can include the introduction, deletion, or substitution of one or more nucleotides at each target sequence of cell(s) via the guide(s) RNA(s) or sgRNA(s). The mutations can include the introduction, deletion, or substitution of 1-75 nucleotides at each target sequence of said cell(s) via the guide(s) RNA(s) or sgRNA(s). The mutations can include the introduction, deletion, or substitution of 1, 5, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30, 35, 40, 45, 50, or 75 nucleotides at each target sequence of said cell(s) via the guide(s) RNA(s) or sgRNA(s). The mutations can include the introduction, deletion, or substitution of 5, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30, 35, 40, 45, 50, or 75 nucleotides at each target sequence of said cell(s) via the guide(s) RNA(s) or sgRNA(s). The mutations can include the introduction, deletion, or substitution of 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30, 35, 40, 45, 50, or 75 nucleotides at each target sequence of said cell(s) via the guide(s) RNA(s) or sgRNA(s). The mutations can include the introduction, deletion, or substitution of 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30, 35, 40, 45, 50, or 75 nucleotides at each target sequence of said cell(s) via the guide(s) RNA(s) or sgRNA(s). The mutations can include the introduction, deletion, or substitution of 40, 45, 50, 75, 100, 200, 300, 400 or 500 nucleotides at each target sequence of said cell(s) via the guide(s) RNA(s) or sgRNA(s).

[0227] For minimization of toxicity and off-target effect, it will be important to control the concentration of Cas mRNA and guide RNA delivered. Optimal concentrations of Cas mRNA and guide RNA can be determined by testing different concentrations in a cellular or non-human eukaryote animal model and using deep sequencing to analyze the extent of modification at potential off-target genomic loci. Alternatively, to minimize the level of toxicity and off-target effect, Cas nickase mRNA (for example *S. pyogenes* Cas9 with the D10A mutation) can be delivered with a pair of guide RNAs targeting a site of interest. Guide sequences and strategies to minimize toxicity and off-target effects can be as in WO 2014/093622 (PCT/US2013/074667); or, via mutation as herein.

[0228] Typically, in the context of an endogenous CRISPR system, formation of a CRISPR complex (comprising a guide sequence hybridized to a target sequence and com-

plexed with one or more Cas proteins) results in cleavage of one or both strands (if applicable) in or near (e.g. within 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 20, 50, or more base pairs from) the target sequence. Without wishing to be bound by theory, the tracr sequence (if applicable or present), which may comprise or consist of all or a portion of a wild-type tracr sequence (e.g. about or more than about 20, 26, 32, 45, 48, 54, 63, 67, 85, or more nucleotides of a wild-type tracr sequence), may also form part of a CRISPR complex, such as by hybridization along at least a portion of the tracr sequence to all or a portion of a tracr mate sequence that is operably linked to the guide sequence.

[0229] The nucleic acid molecule encoding a Cas is advantageously codon optimized Cas. An example of a codon optimized sequence, is in this instance a sequence optimized for expression in a eukaryote, e.g., humans (i.e. being optimized for expression in humans), or for another eukaryote, animal or mammal as herein discussed; see, e.g., SaCas9 human codon optimized sequence in WO 2014/093622 (PCT/US2013/074667). Whilst this is preferred, it will be appreciated that other examples are possible and codon optimization for a host species other than human, or for codon optimization for specific organs is known. In some embodiments, an enzyme coding sequence encoding a Cas is codon optimized for expression in particular cells, such as eukaryotic cells. The eukaryotic cells may be those of or derived from a particular organism, such as a mammal, including but not limited to human, or non-human eukaryote or animal or mammal as herein discussed, e.g., mouse, rat, rabbit, dog, livestock, or non-human mammal or primate. In some embodiments, processes for modifying the germ line genetic identity of human beings and/or processes for modifying the genetic identity of animals which are likely to cause them suffering without any substantial medical benefit to man or animal, and also animals resulting from such processes, may be excluded. In general, codon optimization refers to a process of modifying a nucleic acid sequence for enhanced expression in the host cells of interest by replacing at least one codon (e.g. about or more than about 1, 2, 3, 4, 5, 10, 15, 20, 25, 50, or more codons) of the native sequence with codons that are more frequently or most frequently used in the genes of that host cell while maintaining the native amino acid sequence. Various species exhibit particular bias for certain codons of a particular amino acid. Codon bias (differences in codon usage between organisms) often correlates with the efficiency of translation of messenger RNA (mRNA), which is in turn believed to be dependent on, among other things, the properties of the codons being translated and the availability of particular transfer RNA (tRNA) molecules. The predominance of selected tRNAs in a cell is generally a reflection of the codons used most frequently in peptide synthesis. Accordingly, genes can be tailored for optimal gene expression in a given organism based on codon optimization. Codon usage tables are readily available, for example, at the "Codon Usage Database" available at www.kazusa.or.jp/codon/ and these tables can be adapted in a number of ways. See Nakamura, Y., et al. "Codon usage tabulated from the international DNA sequence databases: status for the year 2000" *Nucl. Acids Res.* 28:292 (2000). Computer algorithms for codon optimizing a particular sequence for expression in a particular host cell are also available, such as Gene Forge (Aptagen; Jacobus, Pa.), are also available. In some embodiments, one or more codons (e.g. 1, 2, 3, 4, 5, 10, 15, 20, 25, 50, or more,

or all codons) in a sequence encoding a Cas correspond to the most frequently used codon for a particular amino acid.

[0230] In certain embodiments, the methods as described herein may comprise providing a Cas transgenic cell in which one or more nucleic acids encoding one or more guide RNAs are provided or introduced operably connected in the cell with a regulatory element comprising a promoter of one or more gene of interest. As used herein, the term “Cas transgenic cell” refers to a cell, such as a eukaryotic cell, in which a Cas gene has been genomically integrated. The nature, type, or origin of the cell are not particularly limiting according to the present invention. Also the way how the Cas transgene is introduced in the cell is may vary and can be any method as is known in the art. In certain embodiments, the Cas transgenic cell is obtained by introducing the Cas transgene in an isolated cell. In certain other embodiments, the Cas transgenic cell is obtained by isolating cells from a Cas transgenic organism. By means of example, and without limitation, the Cas transgenic cell as referred to herein may be derived from a Cas transgenic eukaryote, such as a Cas knock-in eukaryote. Reference is made to WO 2014/093622 (PCT/US13/74667), incorporated herein by reference. Methods of US Patent Publication Nos. 20120017290 and 20110265198 assigned to Sangamo Bio-Sciences, Inc. directed to targeting the Rosa locus may be modified to utilize the CRISPR Cas system of the present invention. Methods of US Patent Publication No. 20130236946 assigned to Collectis directed to targeting the Rosa locus may also be modified to utilize the CRISPR Cas system of the present invention. By means of further example reference is made to Platt et. al. (Cell; 159(2):440-455 (2014)), describing a Cas9 knock-in mouse, which is incorporated herein by reference. The Cas transgene can further comprise a Lox-Stop-polyA-Lox(LSL) cassette thereby rendering Cas expression inducible by Cre recombinase. Alternatively, the Cas transgenic cell may be obtained by introducing the Cas transgene in an isolated cell. Delivery systems for transgenes are well known in the art. By means of example, the Cas transgene may be delivered in for instance eukaryotic cell by means of vector (e.g., AAV, adenovirus, lentivirus) and/or particle and/or particle delivery, as also described herein elsewhere.

[0231] It will be understood by the skilled person that the cell, such as the Cas transgenic cell, as referred to herein may comprise further genomic alterations besides having an integrated Cas gene or the mutations arising from the sequence specific action of Cas when complexed with RNA capable of guiding Cas to a target locus, such as for instance one or more oncogenic mutations, as for instance and without limitation described in Platt et al. (2014), Chen et al., (2014) or Kumar et al., (2009).

[0232] In some embodiments, the Cas sequence is fused to one or more nuclear localization sequences (NLSs), such as about or more than about 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, or more NLSs. In some embodiments, the Cas comprises about or more than about 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, or more NLSs at or near the amino-terminus, about or more than about 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, or more NLSs at or near the carboxy-terminus, or a combination of these (e.g. zero or at least one or more NLS at the amino-terminus and zero or at one or more NLS at the carboxy terminus). When more than one NLS is present, each may be selected independently of the others, such that a single NLS may be present in more than one copy and/or in combination with one or more other

NLSs present in one or more copies. In a preferred embodiment of the invention, the Cas comprises at most 6 NLSs. In some embodiments, an NLS is considered near the N- or C-terminus when the nearest amino acid of the NLS is within about 1, 2, 3, 4, 5, 10, 15, 20, 25, 30, 40, 50, or more amino acids along the polypeptide chain from the N- or C-terminus. Non-limiting examples of NLSs include an NLS sequence derived from: the NLS of the SV40 virus large T-antigen, having the amino acid sequence PKK-KRKV (SEQ ID NO: X); the NLS from nucleoplasmin (e.g. the nucleoplasmin bipartite NLS with the sequence KRPAATKKAGQAKKKK) (SEQ ID NO: X, artificial); the c-myc NLS having the amino acid sequence PAAKRVKLD (SEQ ID NO: X, artificial) or RQRRNELKRSP (SEQ ID NO: X, artificial); the hRNP1 M9 NLS having the sequence NQSSNFGPMKGGNFGGRSSGPYGGGGQY-FAKPRNQGGY (SEQ ID NO: X, artificial); the sequence RMRIZFKNKGKDTAELRRRRVEVSVELRKAKKDEQ-ILKRRNV (SEQ ID NO: X, artificial) of the IBB domain from importin-alpha; the sequences VSRKRPRP (SEQ ID NO: X, artificial) and PPKKARED (SEQ ID NO: X, artificial) of the myoma T protein; the sequence POPKKKPL (SEQ ID NO: X) of human p53; the sequence SALIKKK-KKMAP (SEQ ID NO: X, mouse) of mouse c-ab1 IV; the sequences DRLRR (SEQ ID NO: X) and PKQKKRK (SEQ ID NO: X) of the influenza virus NS1; the sequence RKLK-KKIKKL (SEQ ID NO: X) of the Hepatitis virus delta antigen; the sequence REKKKFLKRR (SEQ ID NO: X) of the mouse Mx1 protein; the sequence KRKGDEVDGVDV-VAKKKSKK (SEQ ID NO: X) of the human poly(ADP-ribose) polymerase; and the sequence RKCLQAGMN-LEARKTKK (SEQ ID NO: X) of the steroid hormone receptors (human) glucocorticoid. In general, the one or more NLSs are of sufficient strength to drive accumulation of the Cas in a detectable amount in the nucleus of a eukaryotic cell. In general, strength of nuclear localization activity may derive from the number of NLSs in the Cas, the particular NLS(s) used, or a combination of these factors. Detection of accumulation in the nucleus may be performed by any suitable technique. For example, a detectable marker may be fused to the Cas, such that location within a cell may be visualized, such as in combination with a means for detecting the location of the nucleus (e.g. a stain specific for the nucleus such as DAPI). Cell nuclei may also be isolated from cells, the contents of which may then be analyzed by any suitable process for detecting protein, such as immunohistochemistry, Western blot, or enzyme activity assay. Accumulation in the nucleus may also be determined indirectly, such as by an assay for the effect of CRISPR complex formation (e.g. assay for DNA cleavage or mutation at the target sequence, or assay for altered gene expression activity affected by CRISPR complex formation and/or Cas enzyme activity), as compared to a control not exposed to the Cas or complex, or exposed to a Cas lacking the one or more NLSs.

[0233] In certain aspects the invention involves vectors, e.g. for delivering or introducing in a cell Cas and/or RNA capable of guiding Cas to a target locus (i.e. guide RNA), but also for propagating these components (e.g. in prokaryotic cells). As used herein, a “vector” is a tool that allows or facilitates the transfer of an entity from one environment to another. It is a replicon, such as a plasmid, phage, or cosmid, into which another DNA segment may be inserted so as to bring about the replication of the inserted segment. Generally, a vector is capable of replication when associated with

the proper control elements. In general, the term “vector” refers to a nucleic acid molecule capable of transporting another nucleic acid to which it has been linked. Vectors include, but are not limited to, nucleic acid molecules that are single-stranded, double-stranded, or partially double-stranded; nucleic acid molecules that comprise one or more free ends, no free ends (e.g. circular); nucleic acid molecules that comprise DNA, RNA, or both; and other varieties of polynucleotides known in the art. One type of vector is a “plasmid,” which refers to a circular double stranded DNA loop into which additional DNA segments can be inserted, such as by standard molecular cloning techniques. Another type of vector is a viral vector, wherein virally-derived DNA or RNA sequences are present in the vector for packaging into a virus (e.g. retroviruses, replication defective retroviruses, adenoviruses, replication defective adenoviruses, and adeno-associated viruses (AAVs)). Viral vectors also include polynucleotides carried by a virus for transfection into a host cell. Certain vectors are capable of autonomous replication in a host cell into which they are introduced (e.g. bacterial vectors having a bacterial origin of replication and episomal mammalian vectors). Other vectors (e.g., non-episomal mammalian vectors) are integrated into the genome of a host cell upon introduction into the host cell, and thereby are replicated along with the host genome. Moreover, certain vectors are capable of directing the expression of genes to which they are operatively-linked. Such vectors are referred to herein as “expression vectors.” Common expression vectors of utility in recombinant DNA techniques are often in the form of plasmids.

[0234] Recombinant expression vectors can comprise a nucleic acid of the invention in a form suitable for expression of the nucleic acid in a host cell, which means that the recombinant expression vectors include one or more regulatory elements, which may be selected on the basis of the host cells to be used for expression, that is operatively-linked to the nucleic acid sequence to be expressed. Within a recombinant expression vector, “operably linked” is intended to mean that the nucleotide sequence of interest is linked to the regulatory element(s) in a manner that allows for expression of the nucleotide sequence (e.g. in an *in vitro* transcription/translation system or in a host cell when the vector is introduced into the host cell). With regards to recombination and cloning methods, mention is made of U.S. patent application Ser. No. 10/815,730, published Sep. 2, 2004 as US 2004-0171156 A1, the contents of which are herein incorporated by reference in their entirety.

[0235] The vector(s) can include the regulatory element(s), e.g., promoter(s). The vector(s) can comprise Cas encoding sequences, and/or a single, but possibly also can comprise at least 3 or 8 or 16 or 32 or 48 or 50 guide RNA(s) (e.g., sgRNAs) encoding sequences, such as 1-2, 1-3, 1-4 1-5, 3-6, 3-7, 3-8, 3-9, 3-10, 3-8, 3-16, 3-30, 3-32, 3-48, 3-50 RNA(s) (e.g., sgRNAs). In a single vector there can be a promoter for each RNA (e.g., sgRNA), advantageously when there are up to about 16 RNA(s) (e.g., sgRNAs); and, when a single vector provides for more than 16 RNA(s) (e.g., sgRNAs), one or more promoter(s) can drive expression of more than one of the RNA(s) (e.g., sgRNAs), e.g., when there are 32 RNA(s) (e.g., sgRNAs), each promoter can drive expression of two RNA(s) (e.g., sgRNAs), and when there are 48 RNA(s) (e.g., sgRNAs), each promoter can drive expression of three RNA(s) (e.g., sgRNAs). By simple arithmetic and well established cloning protocols and

the teachings in this disclosure one skilled in the art can readily practice the invention as to the RNA(s) (e.g., sgRNA(s)) for a suitable exemplary vector such as AAV, and a suitable promoter such as the U6 promoter, e.g., U6-sgRNAs. For example, the packaging limit of AAV is ~4.7 kb. The length of a single U6-sgRNA (plus restriction sites for cloning) is 361 bp. Therefore, the skilled person can readily fit about 12-16, e.g., 13 U6-sgRNA cassettes in a single vector. This can be assembled by any suitable means, such as a golden gate strategy used for TALE assembly (<http://www.genome-engineering.org/talectors/>). The skilled person can also use a tandem guide strategy to increase the number of U6-sgRNAs by approximately 1.5 times, e.g., to increase from 12-16, e.g., 13 to approximately 18-24, e.g., about 19 U6-sgRNAs. Therefore, one skilled in the art can readily reach approximately 18-24, e.g., about 19 promoter-RNAs, e.g., U6-sgRNAs in a single vector, e.g., an AAV vector. A further means for increasing the number of promoters and RNAs, e.g., sgRNA(s) in a vector is to use a single promoter (e.g., U6) to express an array of RNAs, e.g., sgRNAs separated by cleavable sequences. And an even further means for increasing the number of promoter-RNAs, e.g., sgRNAs in a vector, is to express an array of promoter-RNAs, e.g., sgRNAs separated by cleavable sequences in the intron of a coding sequence or gene; and, in this instance it is advantageous to use a polymerase II promoter, which can have increased expression and enable the transcription of long RNA in a tissue specific manner. (see, e.g., <http://nar.oxfordjournals.org/content/34/7/e53>. short, <http://www.nature.com/mt/journal/v16/n9/abs/mt2008144a.html>). In an advantageous embodiment, AAV may package U6 tandem sgRNA targeting up to about 50 genes. Accordingly, from the knowledge in the art and the teachings in this disclosure the skilled person can readily make and use vector(s), e.g., a single vector, expressing multiple RNAs or guides or sgRNAs under the control or operatively or functionally linked to one or more promoters—especially as to the numbers of RNAs or guides or sgRNAs discussed herein, without any undue experimentation.

[0236] The guide RNA(s), e.g., sgRNA(s) encoding sequences and/or Cas encoding sequences, can be functionally or operatively linked to regulatory element(s) and hence the regulatory element(s) drive expression. The promoter(s) can be constitutive promoter(s) and/or conditional promoter(s) and/or inducible promoter(s) and/or tissue specific promoter(s). The promoter can be selected from the group consisting of RNA polymerases, pol I, pol II, pol III, T7, U6, H1, retroviral Rous sarcoma virus (RSV) LTR promoter, the cytomegalovirus (CMV) promoter, the SV40 promoter, the dihydrofolate reductase promoter, the β -actin promoter, the phosphoglycerol kinase (PGK) promoter, and the EF1 α promoter. An advantageous promoter is the promoter is U6.

[0237] In general, the CRISPR-Cas9 system is as used in the foregoing documents, such as WO 2014/093622 (PCT/US2013/074667) and refers collectively to transcripts and other elements involved in the expression of or directing the activity of CRISPR-associated (“Cas”) enzyme, e.g. Cas9, including sequences encoding or delivering a Cas enzyme (DNA and/or RNA-targeting) enzyme, a tracr (trans-activating CRISPR) sequence (e.g., tracrRNA or an active partial tracrRNA), a tracr-mate sequence (encompassing a “direct repeat” and a tracrRNA-processed partial direct repeat in the context of an endogenous CRISPR system), a guide sequence (also referred to as a “spacer” in the context of an

endogenous CRISPR system), or “RNA(s)” as that term is herein used (e.g., RNA(s) to guide Cas9, e.g., CRISPR RNA (crRNA) and trans-activating crRNA (tracrRNA) or a single guide RNA (sgRNA) (chimeric RNA)) or other sequences and transcripts from a CRISPR locus. In general, a CRISPR system is characterized by elements that promote the formation of a CRISPR complex at the site of a target sequence (also referred to as a protospacer in the context of an endogenous CRISPR system). In the context of formation of a CRISPR complex, “target sequence” refers to a sequence to which a guide sequence is designed to target, e.g. have complementarity, where hybridization between a target sequence and a guide sequence promotes the formation of a CRISPR complex. The section of the guide sequence through which complementarity to the target sequence is important for cleavage activity is referred to herein as the seed sequence. A target sequence may comprise any polynucleotide, such as DNA or RNA polynucleotides and is comprised within a target locus of interest. In some embodiments, a target sequence is located in the nucleus or cytoplasm of a cell. The herein described invention encompasses novel effector proteins of Class 2 CRISPR-Cas systems, of which Cas9 is an exemplary effector protein and hence terms used in this application to describe novel effector proteins, may correlate to the terms used to describe the CRISPR-Cas9 system.

[0238] The CRISPR-Cas loci has more than 50 gene families and there is no strictly universal genes. Therefore, no single evolutionary tree is feasible and a multi-pronged approach is needed to identify new families. So far, there is comprehensive cas gene identification of 395 profiles for 93 Cas proteins. Classification includes signature gene profiles plus signatures of locus architecture. Class 1 includes multimeric crRNA-effector complexes (Cascade) and Class 2 includes Single-subunit crRNA-effector complexes (Cas9-like).

[0239] The action of the CRISPR-Cas system is usually divided into three stages: (1) adaptation or spacer integration, (2) processing of the primary transcript of the CRISPR locus (pre-crRNA) and maturation of the crRNA which includes the spacer and variable regions corresponding to 5' and 3' fragments of CRISPR repeats, and (3) DNA or RNA interference. Two proteins, Cas1 and Cas2, that are present in the great majority of the known CRISPR-Cas systems are sufficient for the insertion of spacers into the CRISPR cassettes. These two proteins form a complex that is required for this adaptation process; the endonuclease activity of Cas1 is required for spacer integration whereas Cas2 appears to perform a nonenzymatic function. The Cas1-Cas2 complex represents the highly conserved “information processing” module of CRISPR-Cas that appears to be quasi-autonomous from the rest of the system. (See Annotation and Classification of CRISPR-Cas Systems. Makarova K S, Koonin E V. *Methods Mol Biol.* 2015; 1311:47-75).

[0240] The previously described Class 2 systems, namely Type II and the putative Type V, consisted of only three or four genes in the cas operon, namely the cas1 and cas2 genes comprising the adaptation module (the cas1-cas2 pair of genes are not involved in interference), a single multidomain effector protein that is responsible for interference but also contributes to the pre-crRNA processing and adaptation, and often a fourth gene with uncharacterized functions that is dispensable in at least some Type II systems (and in some cases the fourth gene is cas4 (biochemical or in silico

evidence shows that Cas4 is a PD-(DE)xK superfamily nuclease with three-cysteine C-terminal cluster; possesses 5'-ssDNA exonuclease activity) or csn2, which encodes an inactivated ATPase). In most cases, a CRISPR array and a gene for a distinct RNA species known as tracrRNA, a trans-encoded small CRISPR RNA, are adjacent to Class 2 cas operons. The tracrRNA is partially homologous to the repeats within the respective CRISPR array and is essential for the processing of pre-crRNA that is catalyzed by RNase III, a ubiquitous bacterial enzyme that is not associated with the CRISPR-cas loci.

[0241] Aspects of the invention relate to the identification and engineering of novel effector proteins associated with CRISPR systems that are Cas protein agnostic. In a preferred embodiment, the effector protein comprises a single-subunit effector module. In a further embodiment the effector protein is functional in prokaryotic or eukaryotic cells for in vitro, in vivo or ex vivo applications. An aspect of the invention encompasses computational methods and algorithms to predict new Cas protein agnostic CRISPR systems and identify the components therein.

[0242] In one aspect the identifying all predicted protein coding genes is carried out by comparing the identified genes with Cas protein-specific profiles and annotating them according to NCBI Conserved Domain Database (CDD) which is a protein annotation resource that consists of a collection of well-annotated multiple sequence alignment models for ancient domains and full-length proteins. These are available as position-specific score matrices (PSSMs) for fast identification of conserved domains in protein sequences via RPS-BLAST. CDD content includes NCBI-curated domains, which use 3D-structure information to explicitly define domain boundaries and provide insights into sequence/structure/function relationships, as well as domain models imported from a number of external source databases (Pfam, SMART, COG, PRK, TIGRFAM). In a further aspect, CRISPR arrays were predicted using a PILER-CR program which is a public domain software for finding CRISPR repeats as described in “PILER-CR: fast and accurate identification of CRISPR repeats”, Edgar, R. C., *BMC Bioinformatics*, January 20; 8:18(2007), herein incorporated by reference.

[0243] In a further aspect, the case by case analysis is performed using PSI-BLAST (Position-Specific Iterative Basic Local Alignment Search Tool). PSI-BLAST derives a position-specific scoring matrix (PSSM) or profile from the multiple sequence alignment of sequences detected above a given score threshold using protein-protein BLAST. This PSSM is used to further search the database for new matches, and is updated for subsequent iterations with these newly detected sequences. Thus, PSI-BLAST provides a means of detecting distant relationships between proteins.

[0244] In another aspect, the case by case analysis is performed using HHpred, a method for sequence database searching and structure prediction that is as easy to use as BLAST or PSI-BLAST and that is at the same time much more sensitive in finding remote homologs. In fact, HHpred's sensitivity is competitive with the most powerful servers for structure prediction currently available. HHpred is the first server that is based on the pairwise comparison of profile hidden Markov models (HMMs). Whereas most conventional sequence search methods search sequence databases such as UniProt or the NR, HHpred searches alignment databases, like Pfam or SMART. This greatly

simplifies the list of hits to a number of sequence families instead of a clutter of single sequences. All major publicly available profile and alignment databases are available through HHpred. HHpred accepts a single query sequence or a multiple alignment as input. Within only a few minutes it returns the search results in an easy-to-read format similar to that of PSI-BLAST. Search options include local or global alignment and scoring secondary structure similarity. HHpred can produce pairwise query-template sequence alignments, merged query-template multiple alignments (e.g. for transitive searches), as well as 3D structural models calculated by the MODELLER software from HHpred alignments.

[0245] In one aspect, the present disclosure is directed to a method for identifying novel nucleic acid modifying effectors. In certain example embodiments, the method is directed to identifying novel CRISPR effectors. However, the methods disclosed herein are applicable to other nucleic acid modifying loci that contain conserved genetic elements, such as the conserved repetitive elements found in TALEs and potentially other nucleic acid modifying effectors. Likewise, the nucleic acid modifying activity is not limited to endonuclease activity but any nucleic acid modifying activity including, but not limited to, transposases, recombinases, ligases, glycosylases, topoisomerases, nickases, methylases. It should be further understood that one or more of the steps of the method may be executed on one or more computing devices.

[0246] In certain example embodiments, putative nucleic acid modifying loci may be identified from a set of nucleic acid sequences. The set of nucleic acid sequences may be obtained from a genomic or metagenomic database. The genomic or metagenomic database may comprise only prokaryotic genomic sequences, only eukaryotic genomic sequences, or a combination thereof, depending on the nucleic acid modifying loci to be searched. In certain example embodiments, the method comprises obtaining all available prokaryotic genomic sequence from a genomic or metagenomic database repository. Example database repositories include European Ensemble and NCBI databases. A sub-set of nucleic acid sequences comprising the conserved genomic element may then be selected. Tools for searching and identifying the conserved genomic element will depend on the composition of the conserved genomic element to be identified. In certain example embodiments, PILER-CR CRISPR is used to identify genomic sequences comprising one or more CRISPR arrays. R. C. Edgar. "PILER-CR: fast and accurate identification of CRISPR repeats." *BMC Bioinformatics*. 2007; 8:18.

[0247] The assembled set of nucleic acid sequences are then searched to identify putative nucleic acid modifying loci within either a defined distance of the conserved genomic element, loci comprising at least one protein above a defined size limit, or both. Sequences comprising the conserved genomic element are selected from the set of all nucleic acid sequences. The defined distance from the conserved genomic element may be informed by the architecture of known nucleic acid modifying loci of interest. In certain example embodiments, the defined distance is within 25 kb, 24 kb, 23 kb, 22 kb, 21 kb, 20 kb, 19 kb, 18 kb, 17 kb, 16 kb, 15 kb, 14 kb, 13 kb, 12 kb, 11 kb, 10 kb, 9 kb, 8 kb, 7 kb, 6 kb, 5 kb, 4 kb, 3 kb, 2 kb, or 1 kb of the conserved genomic element. In certain example embodiments, the conserved genomic element is a CRISPR array

and the defined distance is within 10 kb of the CRISPR array. The putative nucleic acid modifying loci may be further screened based on whether the loci encode for one or more protein products of a given size. The define size limit of the may be informed by the known characteristics of similar effectors. In certain example embodiments, loci are selected only if the defined size limit is greater than 100, 150, 200, 225, 250, 275, 300, 325, 350, 375, 400, 425, 450, 475, 400, 525, 550, 575, 600, 625, 650, 675, 700, 725, 750, 775, 800, 825, 850, 875, 900, 925, 950, 975, or 1000 amino acids. In certain example embodiments, the size limit is greater than 700 amino acids.

[0248] The identified putative nucleic acid modifying loci are then grouped into subsets of loci comprising homologous proteins. In certain example embodiments, homologous proteins are identified by conducting a NCBI BLAST homology search against a protein database, such as the NCBI protein database. Cutoff parameters for the NCBI BLAST homology search may vary between 0 and 1e-15. In certain example embodiments, the cut-off parameter is 1e-7. In certain example embodiments HHpred protein domain homology searches (Soding et al. *Nucleic Acids Research*. 2005;33(Web Server Issue):W244-8. doi: 10.1093/nar/gki408) may be conducted on all proteins found in this manner to map putative domains to each putative protein.

[0249] In certain example embodiments, a final set of candidate nucleic acid modifying loci are identified from the grouped subsets. In certain example embodiments, the final set of candidate nucleic acid modifying loci are selected based on low HHpred homology matches to known protein domains. Low homology may refer to predicted protein domains of greater than 5, 10, 15, 20, 25, or 30 residues and greater than 50% match to existing known nucleic acid catalytic domains. In certain other example embodiments, low homology may refer to predicted protein domains of greater than 5, 10, 15, 20, 25, or 30 residues and greater than 50% match to existing known nucleic acid catalytic domains. The final set of candidate nucleic acid modifying loci may be further selected based on identical orientation of large putative protein effectors with respect to putative adjacent accessory proteins. Further, based on the total number of single effector candidates within a group, it is possible to assess whether the group is genetically conserved and thus whether it is likely to be biologically functional. Further criteria for selecting a final candidate set of nucleic acid modifying loci may include subsets comprising putative proteins with minimal existing nucleic acid modifying classifications relative to other loci, subsets comprising large putative proteins with a same orientation as any putative small accessory proteins relative to other loci, and loci comprising large putative proteins with shorter distances to the conserved genomic element relative to other loci. This short distance may be within 3000 kb, 2500 kb, 2000 kb, 1500 kb, or 1000 kb. In certain example embodiments, the final set of candidate nucleic acid modifying loci is selected from the group comprising the largest number of putative effector proteins. In certain example embodiments, any combination of the above criteria may be used to select a final set of candidate nucleic acid modifying loci.

[0250] In certain example embodiments, the final set of candidate nucleic acid modifying loci may be further screened by experimentally validating the nucleic acid modifying function of the candidate nucleic modifying

effectors by performing one or more biochemical validation assays, including those biochemical validation assays further disclosed herein.

[0251] Candidate nucleic acid modifying loci identified using the above method may then serve as a basis for preparing non-naturally occurring or engineered compositions comprising one or more proteins from the identified nucleic acid modifying loci as further disclosed herein.

[0252] In certain example embodiments, a method of identifying a novel CRISPR effector from a genomic database comprises the steps of:

[0253] a) selecting sequences from the database encoding a CRISPR array,

[0254] b) identifying loci located within 10 kb of the CRISPR array comprising Open Reading Frames (ORFs) in the selected sequences from (a),

[0255] c) Selecting from (b) loci comprising ORFs of which only a single ORF encodes a novel CRISPR effector having greater than 700 amino acids and no more than 90% homology to a known CRISPR effector.

[0256] In certain other example embodiments, a method of identifying a novel CRISPR effector from a genomic database comprises the steps of:

[0257] a) selecting sequences from the database encoding a CRISPR array,

[0258] b) identifying loci located within 10 kb of the CRISPR array comprising Open Reading Frames (ORFs) in the selected sequences from (a),

[0259] c) Selecting from (b) loci comprising ORFs of which only a single ORF encodes a novel CRISPR effector having greater than 700 amino acids; and

[0260] comparing the secondary structure of the novel CRISPR effector to the secondary structure of a known CRISPR effector, thereby identifying the novel CRISPR effector.

[0261] In certain other example embodiments, a method of identifying a novel CRISPR effector from a genomic database comprises the steps of:

[0262] a) selecting sequences from the database encoding a CRISPR array,

[0263] b) identifying loci located within 10 kb of the CRISPR array comprising Open Reading Frames (ORFs) in the selected sequences from (a),

[0264] c) Selecting from (b) loci comprising ORFs of which only a single ORF encodes a novel CRISPR effector having greater than 700 amino acids; and

[0265] comparing the secondary structure of the novel CRISPR effector to the secondary structure of a known CRISPR effector, thereby identifying the novel CRISPR effector. The term “nucleic acid-targeting system”, wherein nucleic acid is DNA or RNA, and in some aspects may also refer to DNA-RNA hybrids or derivatives thereof, refers collectively to transcripts and other elements involved in the expression of or directing the activity of DNA or RNA-targeting CRISPR-associated (“Cas”) genes, which may include sequences encoding a DNA or RNA-targeting Cas protein and a DNA or RNA-targeting guide RNA comprising a CRISPR RNA (crRNA) sequence and (in some but not all systems) a trans-activating CRISPR/Cas system RNA (tracrRNA) sequence, or other sequences and transcripts from a DNA or RNA-targeting CRISPR locus. In general, a RNA-targeting system is characterized by elements that promote the formation of a DNA or RNA-targeting complex at the site of a target DNA or RNA sequence. In the context

of formation of a DNA or RNA-targeting complex, “target sequence” refers to a DNA or RNA sequence to which a DNA or RNA-targeting guide RNA is designed to have complementarity, where hybridization between a target sequence and a RNA-targeting guide RNA promotes the formation of a RNA-targeting complex. In some embodiments, a target sequence is located in the nucleus or cytoplasm of a cell.

[0266] In an aspect of the invention, novel RNA targeting systems also referred to as RNA- or RNA-targeting CRISPR systems of the present application are based on identified Group 29 or Group 30 proteins which do not require the generation of customized proteins to target specific RNA sequences but rather a single enzyme can be programmed by a RNA molecule to recognize a specific RNA target, in other words the enzyme can be recruited to a specific RNA target using said RNA molecule.

[0267] The nucleic acids-targeting systems, the vector systems, the vectors and the compositions described herein may be used in various nucleic acids-targeting applications, altering or modifying synthesis of a gene product, such as a protein, nucleic acids cleavage, nucleic acids editing, nucleic acids splicing; trafficking of target nucleic acids, tracing of target nucleic acids, isolation of target nucleic acids, visualization of target nucleic acids, etc.

[0268] In an advantageous embodiment, the present invention encompasses effector proteins identified in loci without a proximate Cas1 or Cas2.

Group 29 Nucleases

[0269] The activity of Group 29 proteins depends on the presence of two HEPN domains. These have been shown to be RNase domains, i.e. nuclease (in particular an endonuclease) cutting RNA. HEPN may also target DNA, or potentially DNA and/or RNA. On the basis that that the HEPN domains of Group 29 proteins are at least capable of binding to and, in their wild-type form, cutting RNA, then it is preferred that the Group 29 effector protein has RNase function. It may also, or alternatively, have DNase function. DNase function, the ability to bind and, potentially cut or nick, DNA is discussed in detail herein.

[0270] Thus, in some embodiments, the effector protein may be a RNA-binding protein, such as a dead-Cas type effector protein, which may be optionally functionalised as described herein for instance with a transcriptional activator or repressor domain, NLS or other functional domain. In some embodiments, the effector protein may be a RNA-binding protein that cleaves a single strand of RNA. If the RNA bound is ssRNA, then the ssRNA is fully cleaved. In some embodiments, the effector protein may be a RNA-binding protein that cleaves a double strand of RNA, for example if it comprises two RNase domains. If the RNA bound is dsRNA, then the dsRNA is fully cleaved. In some embodiments, the effector protein may be a RNA-binding protein that has nickase activity, i.e. it binds dsRNA, but only cleaves one of the RNA strands.

[0271] RNase function in CRISPR systems is known, for example mRNA targeting has been reported for certain type III CRISPR-Cas systems (Hale et al., 2014, *Genes Dev.* vol. 28, 2432-2443; Hale et al., 2009, *Cell*, vol. 139, 945-956; Peng et al., 2015, *Nucleic acids research*, vol. 43, 406-417) and provides significant advantages. A CRISPR-Cas system, composition or method targeting RNA via the present effector proteins is thus provided.

[0272] The target RNA, i.e. the RNA of interest, is the RNA to be targeted by the present invention leading to the recruitment to, and the binding of the effector protein at, the target site of interest on the target RNA. The target RNA may be any suitable form of RNA. This may include, in some embodiments, mRNA. In other embodiments, the target RNA may include tRNA or rRNA.

Interfering RNA (RNAi) and microRNA (miRNA)

[0273] In other embodiments, the target RNA may include interfering RNA, i.e. RNA involved in an RNA interference pathway, such as shRNA, siRNA and so forth. In other embodiments, the target RNA may include microRNA (miRNA). Control over interfering RNA or miRNA may help reduce off-target effects (OTE) seen with those approaches by reducing the longevity of the interfering RNA or miRNA *in vivo* or *in vitro*.

[0274] If the effector protein and suitable guide are selectively expressed (for example spatially or temporally under the control of a suitable promoter, for example a tissue- or cell cycle-specific promoter and/or enhancer) then this could be used to 'protect' the cells or systems (*in vivo* or *in vitro*) from RNAi in those cells. This may be useful in neighbouring tissues or cells where RNAi is not required or for the purposes of comparison of the cells or tissues where the effector protein and suitable guide are and are not expressed (i.e. where the RNAi is not controlled and where it is, respectively). The effector protein may be used to control or bind to molecules comprising or consisting of RNA, such as ribozymes, ribosomes or riboswitches. In embodiments of the invention, the RNA guide can recruit the effector protein to these molecules so that the effector protein is able to bind to them.

Ribosomal RNA (rRNA)

[0275] For example, azalide antibiotics such as azithromycin, are well known. They target and disrupt the 50S ribosomal subunit. The present effector protein, together with a suitable guide RNA to target the 50S ribosomal subunit, may be, in some embodiments, recruited to and bind to the 50S ribosomal subunit. Thus, the present effector protein in concert with a suitable guide directed at a ribosomal (especially the 50s ribosomal subunit) target is provided. Use of this use effector protein in concert with the suitable guide directed at the ribosomal (especially the 50s ribosomal subunit) target may include antibiotic use. In particular, the antibiotic use is analogous to the action of azalide antibiotics, such as azithromycin. In some embodiments, prokaryotic ribosomal subunits, such as the 70S subunit in prokaryotes, the 50S subunit mentioned above, the 30S subunit, as well as the 16S and 5S subunits may be targeted. In other embodiments, eukaryotic ribosomal subunits, such as the 80S subunit in eukaryotes, the 60S subunit, the 40S subunit, as well as the 28S, 18S, 5.8S and 5S subunits may be targeted.

[0276] In some embodiments, the effector protein may be a RNA-binding protein, optionally functionalised, as described herein. In some embodiments, the effector protein may be a RNA-binding protein that cleaves a single strand of RNA. In either case, but particularly where the RNA-binding protein cleaves a single strand of RNA, then ribosomal function may be modulated and, in particular, reduced or destroyed. This may apply to any ribosomal RNA and any ribosomal subunit and the sequences of rRNA are well known.

[0277] Control of ribosomal activity is thus envisaged through use of the present effector protein in concert with a suitable guide to the ribosomal target. This may be through cleavage of, or binding to, the ribosome. In particular, reduction of ribosomal activity is envisaged. This may be useful in assaying ribosomal function *in vivo* or *in vitro*, but also as a means of controlling therapies based on ribosomal activity, *in vivo* or *in vitro*. Furthermore, control (i.e. reduction) of protein synthesis in an *in vivo* or *in vitro* system is envisaged, such control including antibiotic and research and diagnostic use.

Riboswitches

[0278] A riboswitch (also known as an aptozyme) is a regulatory segment of a messenger RNA molecule that binds a small molecule. This typically results in a change in production of the proteins encoded by the mRNA. Thus, control of riboswitch activity is thus envisaged through use of the present effector protein in concert with a suitable guide to the riboswitch target. This may be through cleavage of, or binding to, the riboswitch. In particular, reduction of riboswitch activity is envisaged. This may be useful in assaying riboswitch function *in vivo* or *in vitro*, but also as a means of controlling therapies based on riboswitch activity, *in vivo* or *in vitro*. Furthermore, control (i.e. reduction) of protein synthesis in an *in vivo* or *in vitro* system is envisaged. This control, as for rRNA may include antibiotic and research and diagnostic use.

Ribozymes

[0279] Ribozymes are RNA molecules having catalytic properties, analogous to enzymes (which are of course proteins). As ribozymes, both naturally occurring and engineered, comprise or consist of RNA, they may also be targeted by the present RNA-binding effector protein. In some embodiments, the effector protein may be a RNA-binding protein cleaves the ribozyme to thereby disable it. Control of ribozymal activity is thus envisaged through use of the present effector protein in concert with a suitable guide to the ribozymal target. This may be through cleavage of, or binding to, the ribozyme. In particular, reduction of ribozymal activity is envisaged. This may be useful in assaying ribozymal function *in vivo* or *in vitro*, but also as a means of controlling therapies based on ribozymal activity, *in vivo* or *in vitro*.

Gene Expression, Including RNA Processing

[0280] The effector protein may also be used, together with a suitable guide, to target gene expression, including via control of RNA processing. The control of RNA processing may include RNA processing reactions such as RNA splicing, including alternative splicing, via targeting of RNAPol; viral replication (in particular of satellite viruses, bacteriophages and retroviruses, such as HBV, HBC and HIV and others listed herein) including virioids in plants; and tRNA biosynthesis. The effector protein and suitable guide may also be used to control RNA activation (RNAa). RNAa leads to the promotion of gene expression, so control of gene expression may be achieved that way through disruption or reduction of RNAa and thus less promotion of gene expression.

RNAi Screens

[0281] Identifying gene products whose knockdown is associated with phenotypic changes, biological pathways can be interrogated and the constituent parts identified, via RNAi screens. Control may also be exerted over or during these screens by use of the effector protein and suitable guide to remove or reduce the activity of the RNAi in the screen and thus reinstate the activity of the (previously interfered with) gene product (by removing or reducing the interference/repression).

[0282] Satellite RNAs (satRNAs) and satellite viruses may also be treated.

[0283] Control herein with reference to RNase activity generally means reduction, negative disruption or knock-down or knock out.

In Vivo RNA Applications

Inhibition of Gene Expression

[0284] The target-specific RNases provided herein allow for very specific cutting of a target RNA. The interference at RNA level allows for modulation both spatially and temporally and in a non-invasive way, as the genome is not modified.

[0285] A number of diseases have been demonstrated to be treatable by mRNA targeting. While most of these studies relate to administration of siRNA, it is clear that the RNA targeting effector proteins provided herein can be applied in a similar way.

[0286] Examples of mRNA targets (and corresponding disease treatments) are VEGF, VEGF-R1 and RTP801 (in the treatment of AMD and/or DME), Caspase 2 (in the treatment of Naion) ADRB2 (in the treatment of intraocular pressure), TRPVI (in the treatment of Dry eye syndrome), Syk kinase (in the treatment of asthma), Apo B (in the treatment of hypercholesterolemia), PLK1, KSP and VEGF (in the treatment of solid tumors), Ber-Abl (in the treatment of CML) (Burnett and Rossi Chem Biol. 2012, 19(1): 60-71)). Similarly, RNA targeting has been demonstrated to be effective in the treatment of RNA-virus mediated diseases such as HIV (targeting of HIV Tet and Rev), RSV (targeting of RSV nucleocapsid) and HCV (targeting of miR-122) (Burnett and Rossi Chem Biol. 2012, 19(1): 60-71).

[0287] It is further envisaged that the RNA targeting effector protein of the invention can be used for mutation specific or allele specific knockdown. Guide RNA's can be designed that specifically target a sequence in the transcribed mRNA comprising a mutation or an allele-specific sequence. Such specific knockdown is particularly suitable for therapeutic applications relating to disorders associated with mutated or allele-specific gene products. For example, most cases of familial hypobetalipoproteinemia (FHBL) are caused by mutations in the ApoB gene. This gene encodes two versions of the apolipoprotein B protein: a short version (ApoB-48) and a longer version (ApoB-100). Several ApoB gene mutations that lead to FHBL cause both versions of ApoB to be abnormally short. Specifically targeting and knockdown of mutated ApoB mRNA transcripts with an RNA targeting effector protein of the invention may be beneficial in treatment of FHBL. As another example, Huntington's disease (HD) is caused by an expansion of CAG triplet repeats in the gene coding for the Huntingtin protein, which results in an abnormal protein. Specifically targeting

and knockdown of mutated or allele-specific mRNA transcripts encoding the Huntingtin protein with an RNA targeting effector protein of the invention may be beneficial in treatment of HD.

[0288] It is noted that in this context, and more generally for the various applications as described herein, the use of a split version of the RNA targeting effector protein can be envisaged. Indeed, this may not only allow increased specificity but may also be advantageous for delivery. The Cas13b is split in the sense that the two parts of the Cas13b enzyme substantially comprise a functioning Cas13b. Ideally, the split should always be so that the catalytic domain(s) are unaffected. That Cas13b may function as a nuclease or it may be a dead-Cas13b which is essentially an RNA-binding protein with very little or no catalytic activity, due to typically mutation(s) in its catalytic domains.

[0289] Each half of the split Cas13b may be fused to a dimerization partner. By means of example, and without limitation, employing rapamycin sensitive dimerization domains, allows to generate a chemically inducible split Cas13b for temporal control of Cas13b activity. Cas13b can thus be rendered chemically inducible by being split into two fragments and that rapamycin-sensitive dimerization domains may be used for controlled reassembly of the Cas13b. The two parts of the split Cas13b can be thought of as the N' terminal part and the C' terminal part of the split Cas13b. The fusion is typically at the split point of the Cas13b. In other words, the C' terminal of the N' terminal part of the split Cas13b is fused to one of the dimer halves, whilst the N' terminal of the C' terminal part is fused to the other dimer half.

[0290] The Cas13b does not have to be split in the sense that the break is newly created. The split point is typically designed in silico and cloned into the constructs. Together, the two parts of the split Cas13b, the N' terminal and C' terminal parts, form a full Cas13b, comprising preferably at least 70% or more of the wildtype amino acids (or nucleotides encoding them), preferably at least 80% or more, preferably at least 90% or more, preferably at least 95% or more, and most preferably at least 99% or more of the wildtype amino acids (or nucleotides encoding them). Some trimming may be possible, and mutants are envisaged. Non-functional domains may be removed entirely. What is important is that the two parts may be brought together and that the desired Cas13b function is restored or reconstituted. The dimer may be a homodimer or a heterodimer.

[0291] In certain embodiments, the Cas13b effector as described herein may be used for mutation-specific, or allele-specific targeting, such as for mutation-specific, or allele-specific knockdown.

[0292] The RNA targeting effector protein can moreover be fused to another functional RNase domain, such as a non-specific RNase or Argonaute 2, which acts in synergy to increase the RNase activity or to ensure further degradation of the message.

Modulation of Gene Expression Through Modulation of RNA Function

[0293] Apart from a direct effect on gene expression through cleavage of the mRNA, RNA targeting can also be used to impact specific aspects of the RNA processing within the cell, which may allow a more subtle modulation of gene expression. Generally, modulation can for instance be mediated by interfering with binding of proteins to the

RNA, such as for instance blocking binding of proteins, or recruiting RNA binding proteins. Indeed, modulations can be ensured at different levels such as splicing, transport, localization, translation and turnover of the mRNA. Similarly in the context of therapy, it can be envisaged to address (pathogenic) malfunctioning at each of these levels by using RNA-specific targeting molecules. In these embodiments it is in many cases preferred that the RNA targeting protein is a “dead” Cas13b that has lost the ability to cut the RNA target but maintains its ability to bind thereto, such as the mutated forms of Cas13b described herein.

a) Alternative Splicing

[0294] Many of the human genes express multiple mRNAs as a result of alternative splicing. Different diseases have been shown to be linked to aberrant splicing leading to loss of function or gain of function of the expressed gene. While some of these diseases are caused by mutations that cause splicing defects, a number of these are not. One therapeutic option is to target the splicing mechanism directly. The RNA targeting effector proteins described herein can for instance be used to block or promote splicing, include or exclude exons and influence the expression of specific isoforms and/or stimulate the expression of alternative protein products. Such applications are described in more detail below.

[0295] A RNA targeting effector protein binding to a target RNA can sterically block access of splicing factors to the RNA sequence. The RNA targeting effector protein targeted to a splice site may block splicing at the site, optionally redirecting splicing to an adjacent site. For instance a RNA targeting effector protein binding to the 5' splice site binding can block the recruitment of the U1 component of the spliceosome, favoring the skipping of that exon. Alternatively, a RNA targeting effector protein targeted to a splicing enhancer or silencer can prevent binding of transacting regulatory splicing factors at the target site and effectively block or promote splicing. Exon exclusion can further be achieved by recruitment of ILF2/3 to precursor mRNA near an exon by an RNA targeting effector protein as described herein. As yet another example, a glycine rich domain can be attached for recruitment of hnRNP A1 and exon exclusion (Del Gatto-Konczak et al. *Mol Cell Biol.* 1999 January; 19(1):251-60).

[0296] In certain embodiments, through appropriate selection of gRNA, specific splice variants may be targeted, while other splice variants will not be targeted

[0297] In some cases the RNA targeting effector protein can be used to promote splicing (e.g. where splicing is defective). For instance a RNA targeting effector protein can be associated with an effector capable of stabilizing a splicing regulatory stem-loop in order to further splicing. The RNA targeting effector protein can be linked to a consensus binding site sequence for a specific splicing factor in order to recruit the protein to the target DNA.

[0298] Examples of diseases which have been associated with aberrant splicing include, but are not limited to Paraneoplastic Opsoclonus Myoclonus Ataxia (or POMA), resulting from a loss of Nova proteins which regulate splicing of proteins that function in the synapse, and Cystic Fibrosis, which is caused by defective splicing of a cystic fibrosis transmembrane conductance regulator, resulting in the production of nonfunctional chloride channels. In other diseases aberrant RNA splicing results in gain-of-function.

This is the case for instance in myotonic dystrophy which is caused by a CUG triplet-repeat expansion (from 50 to >1500 repeats) in the 3'UTR of an mRNA, causing splicing defects.

[0299] The RNA targeting effector protein can be used to include an exon by recruiting a splicing factor (such as U1) to a 5' splicing site to promote excision of introns around a desired exon. Such recruitment could be mediated through a fusion with an arginine/serine rich domain, which functions as splicing activator (Gravely B R and Maniatis T, *Mol Cell.* 1998 (5):765-71).

[0300] It is envisaged that the RNA targeting effector protein can be used to block the splicing machinery at a desired locus, resulting in preventing exon recognition and the expression of a different protein product. An example of a disorder that may be treated is Duchenne muscular dystrophy (DMD), which is caused by mutations in the gene encoding for the dystrophin protein. Almost all DMD mutations lead to frameshifts, resulting in impaired dystrophin translation. The RNA targeting effector protein can be paired with splice junctions or exonic splicing enhancers (ESEs) thereby preventing exon recognition, resulting in the translation of a partially functional protein. This converts the lethal Duchenne phenotype into the less severe Becker phenotype.

b) RNA Modification

[0301] RNA editing is a natural process whereby the diversity of gene products of a given sequence is increased by minor modification in the RNA. Typically, the modification involves the conversion of adenosine (A) to inosine (I), resulting in an RNA sequence which is different from that encoded by the genome. RNA modification is generally ensured by the ADAR enzyme, whereby the pre-RNA target forms an imperfect duplex RNA by base-pairing between the exon that contains the adenosine to be edited and an intronic non-coding element. A classic example of A-I editing is the glutamate receptor GluR-B mRNA, whereby the change results in modified conductance properties of the channel (Higuchi M, et al. *Cell.* 1993; 75:1361-70).

[0302] In humans, a heterozygous functional-null mutation in the ADAR1 gene leads to a skin disease, human pigmentary genodermatosis (Miyamura Y, et al. *Am J Hum Genet.* 2003; 73:693-9). It is envisaged that the RNA targeting effector proteins of the present invention can be used to correct malfunctioning RNA modification.

[0303] It is further envisaged that RNA adenosine methylase (N(6)-methyladenosine) can be fused to the RNA targeting effector proteins of the invention and targeted to a transcript of interest. This methylase causes reversible methylation, has regulatory roles and may affect gene expression and cell fate decisions by modulating multiple RNA-related cellular pathways (Fu et al *Nat Rev Genet.* 2014; 15(5):293-306).

c) Polyadenylation

[0304] Polyadenylation of an mRNA is important for nuclear transport, translation efficiency and stability of the mRNA., and all of these, as well as the process of polyadenylation, depend on specific RBPs. Most eukaryotic mRNAs receive a 3' poly(A) tail of about 200 nucleotides after transcription. Polyadenylation involves different RNA-binding protein complexes which stimulate the activity of a poly(A)polymerase (Minvielle-Sebastia L et al. *Curr Opin Cell Biol.* 1999; 11:352-7). It is envisaged that the RNA-

targeting effector proteins provided herein can be used to interfere with or promote the interaction between the RNA-binding proteins and RNA.

[0305] Examples of diseases which have been linked to defective proteins involved in polyadenylation are oculopharyngeal muscular dystrophy (OPMD) (Brais B, et al. *Nat Genet.* 1998; 18:164-7).

d) RNA Export

[0306] After pre-mRNA processing, the mRNA is exported from the nucleus to the cytoplasm. This is ensured by a cellular mechanism which involves the generation of a carrier complex, which is then translocated through the nuclear pore and releases the mRNA in the cytoplasm, with subsequent recycling of the carrier.

[0307] Overexpression of proteins (such as TAP) which play a role in the export of RNA has been found to increase export of transcripts that are otherwise inefficiently exported in *Xenopus* (Katahira J, et al. *EMBO J.* 1999; 18:2593-609).

e) mRNA Localization

[0308] mRNA localization ensures spatially regulated protein production. Localization of transcripts to a specific region of the cell can be ensured by localization elements. In particular embodiments, it is envisaged that the effector proteins described herein can be used to target localization elements to the RNA of interest. The effector proteins can be designed to bind the target transcript and shuttle them to a location in the cell determined by its peptide signal tag. More particularly for instance, a RNA targeting effector protein fused to a nuclear localization signal (NLS) can be used to alter RNA localization.

[0309] Further examples of localization signals include the zipcode binding protein (ZBP1) which ensures localization of β -actin to the cytoplasm in several asymmetric cell types, KDEL retention sequence (localization to endoplasmic reticulum), nuclear export signal (localization to cytoplasm), mitochondrial targeting signal (localization to mitochondria), peroxisomal targeting signal (localization to peroxisome) and m6A marking/YTHDF2 (localization to p-bodies). Other approaches that are envisaged are fusion of the RNA targeting effector protein with proteins of known localization (for instance membrane, synapse).

[0310] Alternatively, the effector protein according to the invention may for instance be used in localization-dependent knockdown. By fusing the effector protein to a appropriate localization signal, the effector is targeted to a particular cellular compartment. Only target RNAs residing in this compartment will effectively be targeted, whereas otherwise identical targets, but residing in a different cellular compartment will not be targeted, such that a localization dependent knockdown can be established.

f) translation

[0311] The RNA targeting effector proteins described herein can be used to enhance or repress translation. It is envisaged that upregulating translation is a very robust way to control cellular circuits. Further, for functional studies a protein translation screen can be favorable over transcriptional upregulation screens, which have the shortcoming that upregulation of transcript does not translate into increased protein production.

[0312] It is envisaged that the RNA targeting effector proteins described herein can be used to bring translation initiation factors, such as EIF4G in the vicinity of the 5' untranslated repeat (5'UTR) of a messenger RNA of interest

to drive translation (as described in De Gregorio et al. *EMBO J.* 1999; 18(17):4865-74 for a non-reprogrammable RNA binding protein). As another example GLD2, a cytoplasmic poly(A) polymerase, can be recruited to the target mRNA by an RNA targeting effector protein. This would allow for directed polyadenylation of the target mRNA thereby stimulating translation.

[0313] Similarly, the RNA targeting effector proteins envisaged herein can be used to block translational repressors of mRNA, such as ZBP1 (Huttelmaier S, et al. *Nature.* 2005; 438:512-5). By binding to translation initiation site of a target RNA, translation can be directly affected.

[0314] In addition, fusing the RNA targeting effector proteins to a protein that stabilizes mRNAs, e.g. by preventing degradation thereof such as RNase inhibitors, it is possible to increase protein production from the transcripts of interest.

[0315] It is envisaged that the RNA targeting effector proteins described herein can be used to repress translation by binding in the 5' UTR regions of a RNA transcript and preventing the ribosome from forming and beginning translation.

[0316] Further, the RNA targeting effector protein can be used to recruit Caf1, a component of the CCR4—NOT deadenylase complex, to the target mRNA, resulting in deadenylation of the target transcript and inhibition of protein translation.

[0317] For instance, the RNA targeting effector protein of the invention can be used to increase or decrease translation of therapeutically relevant proteins. Examples of therapeutic applications wherein the RNA targeting effector protein can be used to downregulate or upregulate translation are in amyotrophic lateral sclerosis (ALS) and cardiovascular disorders. Reduced levels of the glial glutamate transporter EAAT2 have been reported in ALS motor cortex and spinal cord, as well as multiple abnormal EAAT2 mRNA transcripts in ALS brain tissue. Loss of the EAAT2 protein and function thought to be the main cause of excitotoxicity in ALS. Restoration of EAAT2 protein levels and function may provide therapeutic benefit. Hence, the RNA targeting effector protein can be beneficially used to upregulate the expression of EAAT2 protein, e.g. by blocking translational repressors or stabilizing mRNA as described above. Apolipoprotein A1 is the major protein component of high density lipoprotein (HDL) and ApoA1 and HDL are generally considered as atheroprotective. It is envisaged that the RNA targeting effector protein can be beneficially used to upregulate the expression of ApoA1, e.g. by blocking translational repressors or stabilizing mRNA as described above.

g) mRNA Turnover

[0318] Translation is tightly coupled to mRNA turnover and regulated mRNA stability. Specific proteins have been described to be involved in the stability of transcripts (such as the ELAV/Hu proteins in neurons, Keene J D, 1999, *Proc Natl Acad Sci USA.* 96:5-7) and tristetraprolin (TTP). These proteins stabilize target mRNAs by protecting the messages from degradation in the cytoplasm (Peng S S et al., 1988, *EMBO J.* 17:3461-70).

[0319] It can be envisaged that the RNA-targeting effector proteins of the present invention can be used to interfere with or to promote the activity of proteins acting to stabilize mRNA transcripts, such that mRNA turnover is affected. For instance, recruitment of human TTP to the target RNA using the RNA targeting effector protein would allow for

adenylate-uridylylate-rich element (AU-rich element) mediated translational repression and target degradation. AU-rich elements are found in the 3' UTR of many mRNAs that code for proto-oncogenes, nuclear transcription factors, and cytokines and promote RNA stability. As another example, the RNA targeting effector protein can be fused to HuR, another mRNA stabilization protein (Hinman M N and Lou H, *Cell Mol Life Sci* 2008; 65:3168-81), and recruit it to a target transcript to prolong its lifetime or stabilize short-lived mRNA.

[0320] It is further envisaged that the RNA-targeting effector proteins described herein can be used to promote degradation of target transcripts. For instance, m6A methyltransferase can be recruited to the target transcript to localize the transcript to P-bodies leading to degradation of the target.

[0321] As yet another example, an RNA targeting effector protein as described herein can be fused to the non-specific endonuclease domain PiT N-terminus (PIN), to recruit it to a target transcript and allow degradation thereof.

[0322] Patients with paraneoplastic neurological disorder (PND)—associated encephalomyelitis and neuropathy are patients who develop autoantibodies against Hu-proteins in tumors outside of the central nervous system (Szabo A et al. 1991, *Cell*; 67:325-33 which then cross the blood-brain barrier. It can be envisaged that the RNA-targeting effector proteins of the present invention can be used to interfere with the binding of auto-antibodies to mRNA transcripts.

[0323] Patients with dystrophy type 1 (DM1), caused by the expansion of (CUG)_n in the 3' UTR of dystrophin myotonia-protein kinase (DMPK) gene, are characterized by the accumulation of such transcripts in the nucleus. It is envisaged that the RNA targeting effector proteins of the invention fused with an endonuclease targeted to the (CUG)_n repeats could inhibit such accumulation of aberrant transcripts.

h) Interaction with Multi-Functional Proteins

[0324] Some RNA-binding proteins bind to multiple sites on numerous RNAs to function in diverse processes. For instance, the hnRNP A1 protein has been found to bind exonic splicing silencer sequences, antagonizing the splicing factors, associate with telomere ends (thereby stimulating telomere activity) and bind miRNA to facilitate Drosha-mediated processing thereby affecting maturation. It is envisaged that the RNA-binding effector proteins of the present invention can interfere with the binding of RNA-binding proteins at one or more locations.

i) RNA Folding

[0325] RNA adopts a defined structure in order to perform its biological activities. Transitions in conformation among alternative tertiary structures are critical to most RNA-mediated processes. However, RNA folding can be associated with several problems. For instance, RNA may have a tendency to fold into, and be upheld in, improper alternative conformations and/or the correct tertiary structure may not be sufficiently thermodynamically favored over alternative structures. The RNA targeting effector protein, in particular a cleavage-deficient or dead RNA targeting protein, of the invention may be used to direct folding of (m)RNA and/or capture the correct tertiary structure thereof.

Use of RNA-Targeting Effector Protein in Modulating Cellular Status

[0326] In certain embodiments Cas13b in a complex with crRNA is activated upon binding to target RNA and subsequently cleaves any nearby ssRNA targets (i.e. “collateral” or “bystander” effects). Cas13b, once primed by the cognate target, can cleave other (non-complementary) RNA molecules. Such promiscuous RNA cleavage could potentially cause cellular toxicity, or otherwise affect cellular physiology or cell status.

[0327] Accordingly, in certain embodiments, the non-naturally occurring or engineered composition, vector system, or delivery systems as described herein are used for or are for use in induction of cell dormancy. In certain embodiments, the non-naturally occurring or engineered composition, vector system, or delivery systems as described herein are used for or are for use in induction of cell cycle arrest. In certain embodiments, the non-naturally occurring or engineered composition, vector system, or delivery systems as described herein are used for or are for use in reduction of cell growth and/or cell proliferation. In certain embodiments, the non-naturally occurring or engineered composition, vector system, or delivery systems as described herein are used for or are for use in induction of cell energy. In certain embodiments, the non-naturally occurring or engineered composition, vector system, or delivery systems as described herein are used for or are for use in induction of cell apoptosis. In certain embodiments, the non-naturally occurring or engineered composition, vector system, or delivery systems as described herein are used for or are for use in induction of cell necrosis. In certain embodiments, the non-naturally occurring or engineered composition, vector system, or delivery systems as described herein are used for or are for use in induction of cell death. In certain embodiments, the non-naturally occurring or engineered composition, vector system, or delivery systems as described herein are used for or are for use in induction of programmed cell death.

[0328] In certain embodiments, the invention relates to a method for induction of cell dormancy comprising introducing or inducing the non-naturally occurring or engineered composition, vector system, or delivery systems as described herein. In certain embodiments, the invention relates to a method for induction of cell cycle arrest comprising introducing or inducing the non-naturally occurring or engineered composition, vector system, or delivery systems as described herein. In certain embodiments, the invention relates to a method for reduction of cell growth and/or cell proliferation comprising introducing or inducing the non-naturally occurring or engineered composition, vector system, or delivery systems as described herein. In certain embodiments, the invention relates to a method for induction of cell energy comprising introducing or inducing the non-naturally occurring or engineered composition, vector system, or delivery systems as described herein. In certain embodiments, the invention relates to a method for induction of cell apoptosis comprising introducing or inducing the non-naturally occurring or engineered composition, vector system, or delivery systems as described herein. In certain embodiments, the invention relates to a method for induction of cell necrosis comprising introducing or inducing the non-naturally occurring or engineered composition, vector system, or delivery systems as described herein. In certain embodiments, the invention relates to a method for induc-

tion of cell death comprising introducing or inducing the non-naturally occurring or engineered composition, vector system, or delivery systems as described herein. In certain embodiments, the invention relates to a method for induction of programmed cell death comprising introducing or inducing the non-naturally occurring or engineered composition, vector system, or delivery systems as described herein.

[0329] The methods and uses as described herein may be therapeutic or prophylactic and may target particular cells, cell (sub)populations, or cell/tissue types. In particular, the methods and uses as described herein may be therapeutic or prophylactic and may target particular cells, cell (sub)populations, or cell/tissue types expressing one or more target sequences, such as one or more particular target RNA (e.g. ss RNA). Without limitation, target cells may for instance be cancer cells expressing a particular transcript, e.g. neurons of a given class, (immune) cells causing e.g. autoimmunity, or cells infected by a specific (e.g. viral) pathogen, etc.

[0330] Accordingly, in certain embodiments, the invention relates to a method for treating a pathological condition characterized by the presence of undesirable cells (host cells), comprising introducing or inducing the non-naturally occurring or engineered composition, vector system, or delivery systems as described herein. In certain embodiments, the invention relates the use of the non-naturally occurring or engineered composition, vector system, or delivery systems as described herein for treating a pathological condition characterized by the presence of undesirable cells (host cells). In certain embodiments, the invention relates the non-naturally occurring or engineered composition, vector system, or delivery systems as described herein for use in treating a pathological condition characterized by the presence of undesirable cells (host cells). It is to be understood that preferably the CRISPR-Cas system targets a target specific for the undesirable cells. In certain embodiments, the invention relates to the use of the non-naturally occurring or engineered composition, vector system, or delivery systems as described herein for treating, preventing, or alleviating cancer. In certain embodiments, the invention relates to the non-naturally occurring or engineered composition, vector system, or delivery systems as described herein for use in treating, preventing, or alleviating cancer. In certain embodiments, the invention relates to a method for treating, preventing, or alleviating cancer comprising introducing or inducing the non-naturally occurring or engineered composition, vector system, or delivery systems as described herein. It is to be understood that preferably the CRISPR-Cas system targets a target specific for the cancer cells. In certain embodiments, the invention relates to the use of the non-naturally occurring or engineered composition, vector system, or delivery systems as described herein for treating, preventing, or alleviating infection of cells by a pathogen. In certain embodiments, the invention relates to the non-naturally occurring or engineered composition, vector system, or delivery systems as described herein for use in treating, preventing, or alleviating infection of cells by a pathogen. In certain embodiments, the invention relates to a method for treating, preventing, or alleviating infection of cells by a pathogen comprising introducing or inducing the non-naturally occurring or engineered composition, vector system, or delivery systems as described herein. It is to be understood that preferably the CRISPR-Cas system targets a

target specific for the cells infected by the pathogen (e.g. a pathogen derived target). In certain embodiments, the invention relates to the use of the non-naturally occurring or engineered composition, vector system, or delivery systems as described herein for treating, preventing, or alleviating an autoimmune disorder. In certain embodiments, the invention relates to the non-naturally occurring or engineered composition, vector system, or delivery systems as described herein for use in treating, preventing, or alleviating an autoimmune disorder. In certain embodiments, the invention relates to a method for treating, preventing, or alleviating an autoimmune disorder comprising introducing or inducing the non-naturally occurring or engineered composition, vector system, or delivery systems as described herein. It is to be understood that preferably the CRISPR-Cas system targets a target specific for the cells responsible for the autoimmune disorder (e.g. specific immune cells).

Use of RNA-Targeting Effector Protein in RNA Detection

[0331] It is further envisaged that the RNA targeting effector protein can be used in Northern blot assays. Northern blotting involves the use of electrophoresis to separate RNA samples by size. The RNA targeting effector protein can be used to specifically bind and detect the target RNA sequence.

[0332] A RNA targeting effector protein can be fused to a fluorescent protein (such as GFP) and used to track RNA localization in living cells. More particularly, the RNA targeting effector protein can be inactivated in that it no longer cleaves RNA. In particular embodiments, it is envisaged that a split RNA targeting effector protein can be used, whereby the signal is dependent on the binding of both subproteins, in order to ensure a more precise visualization. Alternatively, a split fluorescent protein can be used that is reconstituted when multiple RNA targeting effector protein complexes bind to the target transcript. It is further envisaged that a transcript is targeted at multiple binding sites along the mRNA so the fluorescent signal can amplify the true signal and allow for focal identification. As yet another alternative, the fluorescent protein can be reconstituted from a split intein.

[0333] RNA targeting effector proteins are for instance suitably used to determine the localization of the RNA or specific splice variants, the level of mRNA transcript, up- or down-regulation of transcripts and disease-specific diagnosis. The RNA targeting effector proteins can be used for visualization of RNA in (living) cells using e.g. fluorescent microscopy or flow cytometry, such as fluorescence-activated cell sorting (FACS) which allows for high-throughput screening of cells and recovery of living cells following cell sorting. Further, expression levels of different transcripts can be assessed simultaneously under stress, e.g. inhibition of cancer growth using molecular inhibitors or hypoxic conditions on cells. Another application would be to track localization of transcripts to synaptic connections during a neural stimulus using two photon microscopy.

[0334] In certain embodiments, the components or complexes according to the invention as described herein can be used in multiplexed error-robust fluorescence in situ hybridization (MERFISH; Chen et al. Science; 2015; 348(6233)), such as for instance with (fluorescently) labeled Cas13b effectors.

In Vitro Apex Labeling

[0335] Cellular processes depend on a network of molecular interactions among protein, RNA, and DNA. Accurate detection of protein—DNA and protein—RNA interactions is key to understanding such processes. In vitro proximity labeling technology employs an affinity tag combined with e.g. a photoactivatable probe to label polypeptides and RNAs in the vicinity of a protein or RNA of interest in vitro. After UV irradiation the photoactivatable group reacts with proteins and other molecules that are in close proximity to the tagged molecule, thereby labelling them. Labelled interacting molecules can subsequently be recovered and identified. The RNA targeting effector protein of the invention can for instance be used to target a probe to a selected RNA sequence.

[0336] These applications could also be applied in animal models for in vivo imaging of disease relevant applications or difficult-to culture cell types. Use of RNA-targeting effector protein in RNA origami/in vitro assembly lines—combinatorics RNA origami refers to nanoscale folded structures for creating two-dimensional or three-dimensional structures using RNA as integrated template. The folded structure is encoded in the RNA and the shape of the resulting RNA is thus determined by the synthesized RNA sequence (Geary, et al. 2014. *Science*, 345 (6198). pp. 799-804). The RNA origami may act as scaffold for arranging other components, such as proteins, into complexes. The RNA targeting effector protein of the invention can for instance be used to target proteins of interest to the RNA origami using a suitable guide RNA.

[0337] These applications could also be applied in animal models for in vivo imaging of disease relevant applications or difficult-to culture cell types.

Use of RNA-Targeting Effector Protein in RNA Isolation or Purification, Enrichment or Depletion

[0338] It is further envisaged that the RNA targeting effector protein when complexed to RNA can be used to isolate and/or purify the RNA. The RNA targeting effector protein can for instance be fused to an affinity tag that can be used to isolate and/or purify the RNA-RNA targeting effector protein complex. Such applications are for instance useful in the analysis of gene expression profiles in cells.

[0339] In particular embodiments, it can be envisaged that the RNA targeting effector proteins can be used to target a specific noncoding RNA (ncRNA) thereby blocking its activity, providing a useful functional probe. In certain embodiments, the effector protein as described herein may be used to specifically enrich for a particular RNA (including but not limited to increasing stability, etc.), or alternatively to specifically deplete a particular RNA (such as without limitation for instance particular splice variants, isoforms, etc.).

Interrogation of lincRNA Function and Other Nuclear RNAs

[0340] Current RNA knockdown strategies such as siRNA have the disadvantage that they are mostly limited to targeting cytosolic transcripts since the protein machinery is cytosolic. The advantage of a RNA targeting effector protein of the present invention, an exogenous system that is not essential to cell function, is that it can be used in any compartment in the cell. By fusing a NLS signal to the RNA targeting effector protein, it can be guided to the nucleus, allowing nuclear RNAs to be targeted. It is for instance

envisaged to probe the function of lincRNAs. Long intergenic non-coding RNAs (lincRNAs) are a vastly underexplored area of research. Most lincRNAs have as of yet unknown functions which could be studied using the RNA targeting effector protein of the invention.

Identification of RNA Binding Proteins

[0341] Identifying proteins bound to specific RNAs can be useful for understanding the roles of many RNAs. For instance, many lincRNAs associate with transcriptional and epigenetic regulators to control transcription. Understanding what proteins bind to a given lincRNA can help elucidate the components in a given regulatory pathway. A RNA targeting effector protein of the invention can be designed to recruit a biotin ligase to a specific transcript in order to label locally bound proteins with biotin. The proteins can then be pulled down and analyzed by mass spectrometry to identify them.

Assembly of Complexes on RNA and Substrate Shuttling

[0342] RNA targeting effector proteins of the invention can further be used to assemble complexes on RNA. This can be achieved by functionalizing the RNA targeting effector protein with multiple related proteins (e.g. components of a particular synthesis pathway). Alternatively, multiple RNA targeting effector proteins can be functionalized with such different related proteins and targeted to the same or adjacent target RNA. Useful application of assembling complexes on RNA are for instance facilitating substrate shuttling between proteins.

Synthetic Biology

[0343] The development of biological systems have a wide utility, including in clinical applications. It is envisaged that the programmable RNA targeting effector proteins of the invention can be used fused to split proteins of toxic domains for targeted cell death, for instance using cancer-linked RNA as target transcript. Further, pathways involving protein-protein interaction can be influenced in synthetic biological systems with e.g. fusion complexes with the appropriate effectors such as kinases or other enzymes.

Protein Splicing: Inteins

[0344] Protein splicing is a post-translational process in which an intervening polypeptide, referred to as an intein, catalyzes its own excision from the polypeptides flanking it, referred to as exteins, as well as subsequent ligation of the exteins. The assembly of two or more RNA targeting effector proteins as described herein on a target transcript could be used to direct the release of a split intein (Topilina and Mills *Mob DNA*. 2014 Feb. 4; 5(1):5), thereby allowing for direct computation of the existence of a mRNA transcript and subsequent release of a protein product, such as a metabolic enzyme or a transcription factor (for downstream actuation of transcription pathways). This application may have significant relevance in synthetic biology (see above) or large-scale bioproduction (only produce product under certain conditions).

Inducible, Dosed and Self-Inactivating Systems

[0345] In one embodiment, fusion complexes comprising an RNA targeting effector protein of the invention and an effector component are designed to be inducible, for instance

light inducible or chemically inducible. Such inducibility allows for activation of the effector component at a desired moment in time.

[0346] Light inducibility is for instance achieved by designing a fusion complex wherein CRY2PHR/CIBN pairing is used for fusion. This system is particularly useful for light induction of protein interactions in living cells (Konermann S, et al. *Nature*. 2013; 500:472-476).

[0347] Chemical inducibility is for instance provided for by designing a fusion complex wherein FKBP/FRB (FK506 binding protein/FKBP rapamycin binding) pairing is used for fusion. Using this system rapamycin is required for binding of proteins (Zetsche et al. *Nat Biotechnol*. 2015; 33(2):139-42 describes the use of this system for Cas9).

[0348] Further, when introduced in the cell as DNA, the RNA targeting effector protein of the inventions can be modulated by inducible promoters, such as tetracycline or doxycycline controlled transcriptional activation (Tet-On and Tet-Off expression system), hormone inducible gene expression system such as for instance an ecdysone inducible gene expression system and an arabinose-inducible gene expression system. When delivered as RNA, expression of the RNA targeting effector protein can be modulated via a riboswitch, which can sense a small molecule like tetracycline (as described in Goldfless et al. *Nucleic Acids Res*. 2012; 40(9):e64).

[0349] In one embodiment, the delivery of the RNA targeting effector protein of the invention can be modulated to change the amount of protein or crRNA in the cell, thereby changing the magnitude of the desired effect or any undesired off-target effects.

[0350] In one embodiment, the RNA targeting effector proteins described herein can be designed to be self-inactivating. When delivered to a cell as RNA, either mRNA or as a replication RNA therapeutic (Wroblecka et al *Nat Biotechnol*. 2015 August; 33(8): 839-841), they can self-inactivate expression and subsequent effects by destroying the own RNA, thereby reducing residency and potential undesirable effects.

[0351] For further in vivo applications of RNA targeting effector proteins as described herein, reference is made to Mackay J P et al (*Nat Struct Mol Biol*. 2011 March; 18(3):256-61), Nelles et al (*Bioessays*. 2015 July; 37(7): 732-9) and Abil Z and Zhao H (*Mol Biosyst*. 2015 October; 11(10):2658-65), which are incorporated herein by reference. In particular, the following applications are envisaged in certain embodiments of the invention, preferably in certain embodiments by using catalytically inactive Cas13b: enhancing translation (e.g. Cas13b—translation promotion factor fusions (e.g. eIF4 fusions)); repressing translation (e.g. gRNA targeting ribosome binding sites); exon skipping (e.g. gRNAs targeting splice donor and/or acceptor sites); exon inclusion (e.g. gRNA targeting a particular exon splice donor and/or acceptor site to be included or Cas13b fused to or recruiting spliceosome components (e.g. U1 snRNA)); accessing RNA localization (e.g. Cas13b—marker fusions (e.g. EGFP fusions)); altering RNA localization (e.g. Cas13b—localization signal fusions (e.g. NLS or NES fusions)); RNA degradation (in this case no catalytically inactive Cas13b is to be used if relied on the activity of Cas13b, alternatively and for increased specificity, a split Cas13b may be used); inhibition of non-coding RNA function (e.g. miRNA), such as by degradation or binding of

gRNA to functional sites (possibly titrating out at specific sites by relocalization by Cas13b-signal sequence fusions).

[0352] As described herein before and demonstrated in the Examples, Cas13b function is robust to 5' or 3' extensions of the crRNA and to extension of the crRNA loop. It is therefore envisaged that MS2 loops and other recruitment domains can be added to the crRNA without affecting complex formation and binding to target transcripts. Such modifications to the crRNA for recruitment of various effector domains are applicable in the uses of a RNA targeted effector proteins described above.

[0353] As demonstrated in the Examples, Cas13b, in particular BzCas13b, is capable of mediating resistance to RNA phages. It is therefore envisaged that Cas13b can be used to immunize, e.g. animals, humans and plants, against RNA-only pathogens, including but not limited to Ebola virus and Zika virus.

[0354] In certain embodiments, Cas13b can process (cleave) its own array. This applies to both the wildtype Cas13b protein and the mutated Cas13b protein containing one or more mutated amino acid residues R116, H121, R1177 and H1182, such as one or more of the modifications selected from R116A, H121A, R1177A and H1182A. It is therefore envisaged that multiple crRNAs designed for different target transcripts and/or applications can be delivered as a single pre-crRNA or as a single transcript driven by one promoter. Such method of delivery has the advantages that it is substantially more compact, easier to synthesize and easier to delivery in viral systems. Preferably, amino acid numbering as described herein refers to BzCas13b protein. It will be understood that exact amino acid positions may vary for orthologues of BzCas13b, which can be adequately determined by protein alignment, as is known in the art, and as described herein elsewhere. Aspects of the invention also encompass methods and uses of the compositions and systems described herein in genome engineering, e.g. for altering or manipulating the expression of one or more genes or the one or more gene products, in prokaryotic or eukaryotic cells, in vitro, in vivo or ex vivo.

[0355] In an aspect, the invention provides methods and compositions for modulating, e.g., reducing, expression of a target RNA in cells. In the subject methods, a Cas13b system of the invention is provided that interferes with transcription, stability, and/or translation of an RNA.

[0356] In certain embodiments, an effective amount of Cas13b system is used to cleave RNA or otherwise inhibit RNA expression. In this regard, the system has uses similar to siRNA and shRNA, thus can also be substituted for such methods. The method includes, without limitation, use of a Cas13b system as a substitute for e.g., an interfering ribonucleic acid (such as an siRNA or shRNA) or a transcription template thereof, e.g., a DNA encoding an shRNA. The Cas13b system is introduced into a target cell, e.g., by being administered to a mammal that includes the target cell.

[0357] Advantageously, a Cas13b system of the invention is specific. For example, whereas interfering ribonucleic acid (such as an siRNA or shRNA) polynucleotide systems are plagued by design and stability issues and off-target binding, a Cas13b system of the invention can be designed with high specificity.

Destabilized Cas13b

[0358] In certain embodiments, the effector protein (CRISPR enzyme; Cas13b) according to the invention as

described herein is associated with or fused to a destabilization domain (DD). In some embodiments, the DD is ER50. A corresponding stabilizing ligand for this DD is, in some embodiments, 4HT. As such, in some embodiments, one of the at least one DDs is ER50 and a stabilizing ligand therefor is 4HT or CMP8. In some embodiments, the DD is DHFR50. A corresponding stabilizing ligand for this DD is, in some embodiments, TMP. As such, in some embodiments, one of the at least one DDs is DHFR50 and a stabilizing ligand therefor is TMP. In some embodiments, the DD is ER50. A corresponding stabilizing ligand for this DD is, in some embodiments, CMP8. CMP8 may therefore be an alternative stabilizing ligand to 4HT in the ER50 system. While it may be possible that CMP8 and 4HT can/should be used in a competitive matter, some cell types may be more susceptible to one or the other of these two ligands, and from this disclosure and the knowledge in the art the skilled person can use CMP8 and/or 4HT.

[0359] In some embodiments, one or two DDs may be fused to the N-terminal end of the CRISPR enzyme with one or two DDs fused to the C-terminal of the CRISPR enzyme. In some embodiments, the at least two DDs are associated with the CRISPR enzyme and the DDs are the same DD, i.e. the DDs are homologous. Thus, both (or two or more) of the DDs could be ER50 DDs. This is preferred in some embodiments. Alternatively, both (or two or more) of the DDs could be DHFR50 DDs. This is also preferred in some embodiments. In some embodiments, the at least two DDs are associated with the CRISPR enzyme and the DDs are different DDs, i.e. the DDs are heterologous. Thus, one of the DDs could be ER50 while one or more of the DDs or any other DDs could be DHFR50. Having two or more DDs which are heterologous may be advantageous as it would provide a greater level of degradation control. A tandem fusion of more than one DD at the N or C-term may enhance degradation; and such a tandem fusion can be, for example ER50-ER50-Cas13b or DHFR-DHFR-Cas13b. It is envisaged that high levels of degradation would occur in the absence of either stabilizing ligand, intermediate levels of degradation would occur in the absence of one stabilizing ligand and the presence of the other (or another) stabilizing ligand, while low levels of degradation would occur in the presence of both (or two of more) of the stabilizing ligands. Control may also be imparted by having an N-terminal ER50 DD and a C-terminal DHFR50 DD.

[0360] In some embodiments, the fusion of the CRISPR enzyme with the DD comprises a linker between the DD and the CRISPR enzyme. In some embodiments, the linker is a GlySer linker. In some embodiments, the DD-CRISPR enzyme further comprises at least one Nuclear Export Signal (NES). In some embodiments, the DD-CRISPR enzyme comprises two or more NESs. In some embodiments, the DD-CRISPR enzyme comprises at least one Nuclear Localization Signal (NLS). This may be in addition to an NES. In some embodiments, the CRISPR enzyme comprises or consists essentially of or consists of a localization (nuclear import or export) signal as, or as part of, the linker between the CRISPR enzyme and the DD. HA or Flag tags are also within the ambit of the invention as linkers. Applicants use NLS and/or NES as linker and also use artificial Glycine Serine linkers as short as GS up to (GGGGS)₃.

[0361] Destabilizing domains have general utility to confer instability to a wide range of proteins; see, e.g., Miyazaki, *J Am Chem Soc.* Mar. 7, 2012; 134(9): 3942-

3945, incorporated herein by reference. CMP8 or 4-hydroxytamoxifen can be destabilizing domains. More generally, A temperature-sensitive mutant of mammalian DHFR (DHFRts), a destabilizing residue by the N-end rule, was found to be stable at a permissive temperature but unstable at 37° C. The addition of methotrexate, a high-affinity ligand for mammalian DHFR, to cells expressing DHFRts inhibited degradation of the protein partially. This was an important demonstration that a small molecule ligand can stabilize a protein otherwise targeted for degradation in cells. A rapamycin derivative was used to stabilize an unstable mutant of the FRB domain of mTOR (FRB*) and restore the function of the fused kinase, GSK-3 β .6,7 This system demonstrated that ligand-dependent stability represented an attractive strategy to regulate the function of a specific protein in a complex biological environment. A system to control protein activity can involve the DD becoming functional when the ubiquitin complementation occurs by rapamycin induced dimerization of FK506-binding protein and FKBP12. Mutants of human FKBP12 or eCDHFR protein can be engineered to be metabolically unstable in the absence of their high-affinity ligands, Shield-1 or trimethoprim (TMP), respectively. These mutants are some of the possible destabilizing domains (DDs) useful in the practice of the invention and instability of a DD as a fusion with a CRISPR enzyme confers to the CRISPR protein degradation of the entire fusion protein by the proteasome. Shield-1 and TMP bind to and stabilize the DD in a dose-dependent manner. The estrogen receptor ligand binding domain (ERLBD, residues 305-549 of ERS1) can also be engineered as a destabilizing domain. Since the estrogen receptor signaling pathway is involved in a variety of diseases such as breast cancer, the pathway has been widely studied and numerous agonist and antagonists of estrogen receptor have been developed. Thus, compatible pairs of ERLBD and drugs are known. There are ligands that bind to mutant but not wild-type forms of the ERLBD. By using one of these mutant domains encoding three mutations (L384M, M421G, G521R)₁₂, it is possible to regulate the stability of an ERLBD-derived DD using a ligand that does not perturb endogenous estrogen-sensitive networks. An additional mutation (Y537S) can be introduced to further destabilize the ERLBD and to configure it as a potential DD candidate. This tetra-mutant is an advantageous DD development. The mutant ERLBD can be fused to a CRISPR enzyme and its stability can be regulated or perturbed using a ligand, whereby the CRISPR enzyme has a DD. Another DD can be a 12-kDa (107-amino-acid) tag based on a mutated FKBP protein, stabilized by Shield1 ligand; see, e.g., *Nature Methods* 5, (2008). For instance a DD can be a modified FK506 binding protein 12 (FKBP12) that binds to and is reversibly stabilized by a synthetic, biologically inert small molecule, Shield-1; see, e.g., Banaszynski L A, Chen L C, Maynard-Smith L A, Ooi A G, Wandless T J. A rapid, reversible, and tunable method to regulate protein function in living cells using synthetic small molecules. *Cell.* 2006; 126:995-1004; Banaszynski L A, Sellmyer M A, Contag C H, Wandless T J, Thorne S H. Chemical control of protein stability and function in living mice. *Nat Med.* 2008; 14:1123-1127; Maynard-Smith L A, Chen L C, Banaszynski L A, Ooi A G, Wandless T J. A directed approach for engineering conditional protein stability using biologically silent small molecules. *The Journal of biological chemistry.* 2007; 282: 24866-24872; and Rodriguez, *Chem Biol.* Mar. 23, 2012;

19(3): 391-398—all of which are incorporated herein by reference and may be employed in the practice of the invention in selected a DD to associate with a CRISPR enzyme in the practice of this invention. As can be seen, the knowledge in the art includes a number of DDs, and the DD can be associated with, e.g., fused to, advantageously with a linker, to a CRISPR enzyme, whereby the DD can be stabilized in the presence of a ligand and when there is the absence thereof the DD can become destabilized, whereby the CRISPR enzyme is entirely destabilized, or the DD can be stabilized in the absence of a ligand and when the ligand is present the DD can become destabilized; the DD allows the CRISPR enzyme and hence the CRISPR-Cas complex or system to be regulated or controlled—turned on or off so to speak, to thereby provide means for regulation or control of the system, e.g., in an in vivo or in vitro environment. For instance, when a protein of interest is expressed as a fusion with the DD tag, it is destabilized and rapidly degraded in the cell, e.g., by proteasomes. Thus, absence of stabilizing ligand leads to a D associated Cas being degraded. When a new DD is fused to a protein of interest, its instability is conferred to the protein of interest, resulting in the rapid degradation of the entire fusion protein. Peak activity for Cas is sometimes beneficial to reduce off-target effects. Thus, short bursts of high activity are preferred. The present invention is able to provide such peaks. In some senses the system is inducible. In some other senses, the system repressed in the absence of stabilizing ligand and de-repressed in the presence of stabilizing ligand.

Application of RNA Targeting—CRISPR System to Plants and Yeast

Definitions

[0362] In general, the term “plant” relates to any various photosynthetic, eukaryotic, unicellular or multicellular organism of the kingdom Plantae characteristically growing by cell division, containing chloroplasts, and having cell walls comprised of cellulose. The term plant encompasses monocotyledonous and dicotyledonous plants. Specifically, the plants are intended to comprise without limitation angiosperm and gymnosperm plants such as acacia, alfalfa, amaranth, apple, apricot, artichoke, ash tree, asparagus, avocado, banana, barley, beans, beet, birch, beech, blackberry, blueberry, broccoli, Brussel’s sprouts, cabbage, canola, cantaloupe, carrot, cassava, cauliflower, cedar, a cereal, celery, chestnut, cherry, Chinese cabbage, citrus, clementine, clover, coffee, corn, cotton, cowpea, cucumber, cypress, eggplant, elm, endive, eucalyptus, fennel, figs, fir, geranium, grape, grapefruit, groundnuts, ground cherry, gum hemlock, hickory, kale, kiwifruit, kohlrabi, larch, lettuce, leek, lemon, lime, locust, pine, maidenhair, maize, mango, maple, melon, millet, mushroom, mustard, nuts, oak, oats, oil palm, okra, onion, orange, an ornamental plant or flower or tree, papaya, palm, parsley, parsnip, pea, peach, peanut, pear, peat, pepper, persimmon, pigeon pea, pine, pineapple, plantain, plum, pomegranate, potato, pumpkin, radicchio, radish, rapeseed, raspberry, rice, rye, sorghum, safflower, sallow, soybean, spinach, spruce, squash, strawberry, sugar beet, sugarcane, sunflower, sweet potato, sweet corn, tangerine, tea, tobacco, tomato, trees, triticale, turf grasses, turnips, vine, walnut, watercress, watermelon, wheat, yams, yew, and zucchini. The term plant also encompasses Algae, which are mainly

photoautotrophs unified primarily by their lack of roots, leaves and other organs that characterize higher plants.

[0363] The methods for modulating gene expression using the RNA targeting system as described herein can be used to confer desired traits on essentially any plant. A wide variety of plants and plant cell systems may be engineered for the desired physiological and agronomic characteristics described herein using the nucleic acid constructs of the present disclosure and the various transformation methods mentioned above. In preferred embodiments, target plants and plant cells for engineering include, but are not limited to, those monocotyledonous and dicotyledonous plants, such as crops including grain crops (e.g., wheat, maize, rice, millet, barley), fruit crops (e.g., tomato, apple, pear, strawberry, orange), forage crops (e.g., alfalfa), root vegetable crops (e.g., carrot, potato, sugar beets, yam), leafy vegetable crops (e.g., lettuce, spinach); flowering plants (e.g., *petunia*, rose, *chrysanthemum*), conifers and pine trees (e.g., pine fir, spruce); plants used in phytoremediation heavy metal accumulating plants); oil crops (e.g., sunflower, rape seed) and plants used for experimental purposes (e.g., *Arabidopsis*). Thus, the methods and CRISPR-Cas systems can be used over a broad range of plants, such as for example with dicotyledonous plants belonging to the orders Magnoliales, Illiciales, Laurales, Piperales, Aristochiales, Nymphaeales, Ranunculales, Papeverales, Sarraceniaceae, Trochodendrales, Hamamelidales, Eucomiales, Leitneriales, Myricales, Fagales, Casuarinales, Caryophyllales, Batales, Polygonales, Plumbaginales, Dilleniales, Theales, Malvales, Urticales, Lecythidales, Violales, Salicales, Capparales, Ericales, Diapensales, Ebenales, Primulales, Rosales, Fabales, Podostemales, Haloragales, Myrtales, Cornales, Proteales, San tales, Rafflesiales, Celastrales, Euphorbiales, Rhamnales, Sapindales, Juglandales, Geraniales, Polygalales, Umbellales, Gentianales, Polemoniales, Lami ales, Plantaginiales, Scrophulariales, Campanulales, Rubiales, Dipsacales, and Asterales, the methods and CRISPR-Cas systems can be used with monocotyledonous plants such as those belonging to the orders Alismatales, Hydrocharitales, Najadales, Triuridales, Commelinales, Eriocaulales, Restionales, Poales, Juncales, Cyperales, Typhales, Bromeliales, Zingiberales, Arecales, Cyclanthales, Pandanales, Arales, Lilliales, and Orchid ales, or with plants belonging to Gymnospermae, e.g. those belonging to the orders Pinales, Ginkgoales, Cycadales, Araucariales, Cupressales and Gnetales.

[0364] The RNA targeting CRISPR systems and methods of use described herein can be used over a broad range of plant species, included in the non-limitative list of dicot, monocot or gymnosperm genera hereunder: *Atropa*, *Alseodaphne*, *Anacardium*, *Arachis*, *Beilschmiedia*, *Brassica*, *Carthamus*, *Cocculus*, *Croton*, *Cucumis*, *Citrus*, *Citrullus*, *Capsicum*, *Catharanthus*, *Cocos*, *Coffea*, *Cucurbita*, *Daucus*, *Duguetia*, *Eschscholzia*, *Ficus*, *Fragaria*, *Glauca*, *Glycine*, *Gossypium*, *Helianthus*, *Hevea*, *Hyoscyamus*, *Lactuca*, *Landolphia*, *Linum*, *Litsea*, *Lycopersicon*, *Lupinus*, *Manihot*, *Majorana*, *Malus*, *Medicago*, *Nicotiana*, *Olea*, *Parthenium*, *Papaver*, *Persea*, *Phaseolus*, *Pistacia*, *Pisum*, *Pyrus*, *Prunus*, *Raphanus*, *Ricinus*, *Senecio*, *Sinomenium*, *Stephania*, *Sinapis*, *Solanum*, *Theobroma*, *Trifolium*, *Trigonella*, *Vicia*, *Vinca*, and *Vigna*; and the genera *Allium*, *Andropogon*, *Aragrostis*, *Asparagus*, *Avena*, *Cynodon*, *Elaeis*, *Festuca*, *Festulolium*, *Heterocallis*, *Hordeum*, *Lemna*, *Lolium*, *Musa*, *Oryza*, *Panicum*, *Pannisetum*,

Phleum, Poa, Secale, Sorghum, Triticum, Zea, Abies, Cunninghamia, Ephedra, Picea, Pinus, and Pseudotsuga.

[0365] The RNA targeting CRISPR systems and methods of use can also be used over a broad range of “algae” or “algae cells”; including for example algae selected from several eukaryotic phyla, including the Rhodophyta (red algae), Chlorophyta (green algae), Phaeophyta (brown algae), Bacillariophyta (diatoms), Eustigmatophyta and dinoflagellates as well as the prokaryotic phylum Cyanobacteria (blue-green algae). The term “algae” includes for example algae selected from *Amphora, Anabaena, Anikstrodesmis, Botryococcus, Chaetoceros, Chlamydomonas, Chlorella, Chlorococcum, Cyclotella, Cylinthotheca, Dunaliella, Emiliana, Euglena, Hematococcus, Isochrysis, Monochrysis, Monoraphidium, Nannochloris, Nannochloropsis, Navicula, Nephrochloris, Nephroselmis, Nitzschia, Nodularia, Nostoc, Oochromonas, Oocystis, Oscillatoria, Pavlova, Phaeodactylurn, Playmonas, Pleurochrysis, Porphyra, Pseudoanabaena, Pyramimonas, Stichococcus, Synechococcus, Synechocystis, Tetrasselmis, Thalassiosira, and Trichodesmium.*

[0366] A part of a plant, i.e., a “plant tissue” may be treated according to the methods of the present invention to produce an improved plant. Plant tissue also encompasses plant cells. The term “plant cell” as used herein refers to individual units of a living plant, either in an intact whole plant or in an isolated form grown in in vitro tissue cultures, on media or agar, in suspension in a growth media or buffer or as a part of higher organized unites, such as, for example, plant tissue, a plant organ, or a whole plant.

[0367] A “protoplast” refers to a plant cell that has had its protective cell wall completely or partially removed using, for example, mechanical or enzymatic means resulting in an intact biochemical competent unit of living plant that can reform their cell wall, proliferate and regenerate grow into a whole plant under proper growing conditions.

[0368] The term “transformation” broadly refers to the process by which a plant host is genetically modified by the introduction of DNA by means of *Agrobacteria* or one of a variety of chemical or physical methods. As used herein, the term “plant host” refers to plants, including any cells, tissues, organs, or progeny of the plants. Many suitable plant tissues or plant cells can be transformed and include, but are not limited to, protoplasts, somatic embryos, pollen, leaves, seedlings, stems, calli, stolons, microtubers, and shoots. A plant tissue also refers to any clone of such a plant, seed, progeny, propagule whether generated sexually or asexually, and descendants of any of these, such as cuttings or seed.

[0369] The term “transformed” as used herein, refers to a cell, tissue, organ, or organism into which a foreign DNA molecule, such as a construct, has been introduced. The introduced DNA molecule may be integrated into the genomic DNA of the recipient cell, tissue, organ, or organism such that the introduced DNA molecule is transmitted to the subsequent progeny. In these embodiments, the “transformed” or “transgenic” cell or plant may also include progeny of the cell or plant and progeny produced from a breeding program employing such a transformed plant as a parent in a cross and exhibiting an altered phenotype resulting from the presence of the introduced DNA molecule. Preferably, the transgenic plant is fertile and capable of transmitting the introduced DNA to progeny through sexual reproduction.

[0370] The term “progeny”, such as the progeny of a transgenic plant, is one that is born of, begotten by, or derived from a plant or the transgenic plant. The introduced DNA molecule may also be transiently introduced into the recipient cell such that the introduced DNA molecule is not inherited by subsequent progeny and thus not considered “transgenic”. Accordingly, as used herein, a “non-transgenic” plant or plant cell is a plant which does not contain a foreign DNA stably integrated into its genome.

[0371] The term “plant promoter” as used herein is a promoter capable of initiating transcription in plant cells, whether or not its origin is a plant cell. Exemplary suitable plant promoters include, but are not limited to, those that are obtained from plants, plant viruses, and bacteria such as *Agrobacterium* or *Rhizobium* which comprise genes expressed in plant cells.

[0372] As used herein, a “fungal cell” refers to any type of eukaryotic cell within the kingdom of fungi. Phyla within the kingdom of fungi include Ascomycota, Basidiomycota, Blastocladiomycota, Chytridiomycota, Glomeromycota, Microsporidia, and Neocallimastigomycota. Fungal cells may include yeasts, molds, and filamentous fungi. In some embodiments, the fungal cell is a yeast cell.

[0373] As used herein, the term “yeast cell” refers to any fungal cell within the phyla Ascomycota and Basidiomycota. Yeast cells may include budding yeast cells, fission yeast cells, and mold cells. Without being limited to these organisms, many types of yeast used in laboratory and industrial settings are part of the phylum Ascomycota. In some embodiments, the yeast cell is an *S. cerevisiae*, *Kluyveromyces marxianus*, or *Issatchenkia orientalis* cell. Other yeast cells may include without limitation *Candida* spp. (e.g., *Candida albicans*), *Yarrowia* spp. (e.g., *Yarrowia lipolytica*), *Pichia* spp. (e.g., *Pichia pastoris*), *Kluyveromyces* spp. (e.g., *Kluyveromyces lactis* and *Kluyveromyces marxianus*), *Neurospora* spp. (e.g., *Neurospora crassa*), *Fusarium* spp. (e.g., *Fusarium oxysporum*), and *Issatchenkia* spp. (e.g., *Issatchenkia orientalis*, a.k.a. *Pichia kudriavzevii* and *Candida acidothermophilum*). In some embodiments, the fungal cell is a filamentous fungal cell. As used herein, the term “filamentous fungal cell” refers to any type of fungal cell that grows in filaments, i.e., hyphae or mycelia. Examples of filamentous fungal cells may include without limitation *Aspergillus* spp. (e.g., *Aspergillus niger*), *Trichoderma* spp. (e.g., *Trichoderma reesei*), *Rhizopus* spp. (e.g., *Rhizopus oryzae*), and *Mortierella* spp. (e.g., *Mortierella isabellina*).

[0374] In some embodiments, the fungal cell is an industrial strain. As used herein, “industrial strain” refers to any strain of fungal cell used in or isolated from an industrial process, e.g., production of a product on a commercial or industrial scale. Industrial strain may refer to a fungal species that is typically used in an industrial process, or it may refer to an isolate of a fungal species that may be also used for non-industrial purposes (e.g., laboratory research). Examples of industrial processes may include fermentation (e.g., in production of food or beverage products), distillation, biofuel production, production of a compound, and production of a polypeptide. Examples of industrial strains may include, without limitation, JAY270 and ATCC4124.

[0375] In some embodiments, the fungal cell is a polyploid cell. As used herein, a “polyploid” cell may refer to any cell whose genome is present in more than one copy. A polyploid cell may refer to a type of cell that is naturally

found in a polyploid state, or it may refer to a cell that has been induced to exist in a polyploid state (e.g., through specific regulation, alteration, inactivation, activation, or modification of meiosis, cytokinesis, or DNA replication). A polyploid cell may refer to a cell whose entire genome is polyploid, or it may refer to a cell that is polyploid in a particular genomic locus of interest. Without wishing to be bound to theory, it is thought that the abundance of guide RNA may more often be a rate-limiting component in genome engineering of polyploid cells than in haploid cells, and thus the methods using the Cas13b CRISPR system described herein may take advantage of using a certain fungal cell type.

[0376] In some embodiments, the fungal cell is a diploid cell. As used herein, a “diploid” cell may refer to any cell whose genome is present in two copies. A diploid cell may refer to a type of cell that is naturally found in a diploid state, or it may refer to a cell that has been induced to exist in a diploid state (e.g., through specific regulation, alteration, inactivation, activation, or modification of meiosis, cytokinesis, or DNA replication). For example, the *S. cerevisiae* strain S228C may be maintained in a haploid or diploid state. A diploid cell may refer to a cell whose entire genome is diploid, or it may refer to a cell that is diploid in a particular genomic locus of interest. In some embodiments, the fungal cell is a haploid cell. As used herein, a “haploid” cell may refer to any cell whose genome is present in one copy. A haploid cell may refer to a type of cell that is naturally found in a haploid state, or it may refer to a cell that has been induced to exist in a haploid state (e.g., through specific regulation, alteration, inactivation, activation, or modification of meiosis, cytokinesis, or DNA replication). For example, the *S. cerevisiae* strain S228C may be maintained in a haploid or diploid state. A haploid cell may refer to a cell whose entire genome is haploid, or it may refer to a cell that is haploid in a particular genomic locus of interest.

[0377] As used herein, a “yeast expression vector” refers to a nucleic acid that contains one or more sequences encoding an RNA and/or polypeptide and may further contain any desired elements that control the expression of the nucleic acid(s), as well as any elements that enable the replication and maintenance of the expression vector inside the yeast cell. Many suitable yeast expression vectors and features thereof are known in the art; for example, various vectors and techniques are illustrated in *Yeast Protocols*, 2nd edition, Xiao, W., ed. (Humana Press, New York, 2007) and Buckholz, R. G. and Gleeson, M. A. (1991) *Biotechnology* (NY) 9(11): 1067-72. Yeast vectors may contain, without limitation, a centromeric (CEN) sequence, an autonomous replication sequence (ARS), a promoter, such as an RNA Polymerase III promoter, operably linked to a sequence or gene of interest, a terminator such as an RNA polymerase III terminator, an origin of replication, and a marker gene (e.g., auxotrophic, antibiotic, or other selectable markers). Examples of expression vectors for use in yeast may include plasmids, yeast artificial chromosomes, 2 μ plasmids, yeast integrative plasmids, yeast replicative plasmids, shuttle vectors, and episomal plasmids.

Stable Integration of RNA Targeting CRISPR System Components in the Genome of Plants and Plant Cells

[0378] In particular embodiments, it is envisaged that the polynucleotides encoding the components of the RNA targeting CRISPR system are introduced for stable integration

into the genome of a plant cell. In these embodiments, the design of the transformation vector or the expression system can be adjusted depending on when, where and under what conditions the guide RNA and/or the RNA targeting gene(s) are expressed.

[0379] In particular embodiments, it is envisaged to introduce the components of the RNA targeting CRISPR system stably into the genomic DNA of a plant cell. Additionally or alternatively, it is envisaged to introduce the components of the RNA targeting CRISPR system for stable integration into the DNA of a plant organelle such as, but not limited to a plastid, a mitochondrion or a chloroplast.

[0380] The expression system for stable integration into the genome of a plant cell may contain one or more of the following elements: a promoter element that can be used to express the guide RNA and/or RNA targeting enzyme in a plant cell; a 5' untranslated region to enhance expression; an intron element to further enhance expression in certain cells, such as monocot cells; a multiple-cloning site to provide convenient restriction sites for inserting the one or more guide RNAs and/or the RNA targeting gene sequences and other desired elements; and a 3' untranslated region to provide for efficient termination of the expressed transcript.

[0381] The elements of the expression system may be on one or more expression constructs which are either circular such as a plasmid or transformation vector, or non-circular such as linear double stranded DNA.

[0382] In a particular embodiment, a RNA targeting CRISPR expression system comprises at least:

[0383] (a) a nucleotide sequence encoding a guide RNA (gRNA) that hybridizes with a target sequence in a plant, and wherein the guide RNA comprises a guide sequence and a direct repeat sequence, and

[0384] (b) a nucleotide sequence encoding a RNA targeting protein,

wherein components (a) or (b) are located on the same or on different constructs, and whereby the different nucleotide sequences can be under control of the same or a different regulatory element operable in a plant cell.

[0385] DNA construct(s) containing the components of the RNA targeting CRISPR system, may be introduced into the genome of a plant, plant part, or plant cell by a variety of conventional techniques. The process generally comprises the steps of selecting a suitable host cell or host tissue, introducing the construct(s) into the host cell or host tissue, and regenerating plant cells or plants therefrom.

[0386] In particular embodiments, the DNA construct may be introduced into the plant cell using techniques such as but not limited to electroporation, microinjection, aerosol beam injection of plant cell protoplasts, or the DNA constructs can be introduced directly to plant tissue using biolistic methods, such as DNA particle bombardment (see also Fu et al., *Transgenic Res.* 2000 February; 9(1):11-9). The basis of particle bombardment is the acceleration of particles coated with gene/s of interest, toward cells, resulting in the penetration of the protoplasm by the particles and typically stable integration into the genome. (see e.g. Klein et al, *Nature* (1987), Klein et al, *Bio/Technology* (1992), Casas et al, *Proc. Natl. Acad. Sci. USA* (1993)).

[0387] In particular embodiments, the DNA constructs containing components of the RNA targeting CRISPR system may be introduced into the plant by *Agrobacterium*-mediated transformation. The DNA constructs may be combined with suitable T-DNA flanking regions and introduced

into a conventional *Agrobacterium tumefaciens* host vector. The foreign DNA can be incorporated into the genome of plants by infecting the plants or by incubating plant protoplasts with *Agrobacterium* bacteria, containing one or more Ti (tumor-inducing) plasmids. (see e.g. Fraley et al., (1985), Rogers et al., (1987) and U.S. Pat. No. 5,563,055).

Plant Promoters

[0388] In order to ensure appropriate expression in a plant cell, the components of the Cas13b CRISPR system described herein are typically placed under control of a plant promoter, i.e. a promoter operable in plant cells. The use of different types of promoters is envisaged.

[0389] A constitutive plant promoter is a promoter that is able to express the open reading frame (ORF) that it controls in all or nearly all of the plant tissues during all or nearly all developmental stages of the plant (referred to as “constitutive expression”). One non-limiting example of a constitutive promoter is the cauliflower mosaic virus 35S promoter. The present invention envisages methods for modifying RNA sequences and as such also envisages regulating expression of plant biomolecules. In particular embodiments of the present invention it is thus advantageous to place one or more elements of the RNA targeting CRISPR system under the control of a promoter that can be regulated. “Regulated promoter” refers to promoters that direct gene expression not constitutively, but in a temporally- and/or spatially-regulated manner, and includes tissue-specific, tissue-preferred and inducible promoters. Different promoters may direct the expression of a gene in different tissues or cell types, or at different stages of development, or in response to different environmental conditions. In particular embodiments, one or more of the RNA targeting CRISPR components are expressed under the control of a constitutive promoter, such as the cauliflower mosaic virus 35S promoter issue-preferred promoters can be utilized to target enhanced expression in certain cell types within a particular plant tissue, for instance vascular cells in leaves or roots or in specific cells of the seed. Examples of particular promoters for use in the RNA targeting CRISPR system are found in Kawamata et al., (1997) *Plant Cell Physiol* 38:792-803; Yamamoto et al., (1997) *Plant J* 12:255-65; Hire et al., (1992) *Plant Mol Biol* 20:207-18, Kuster et al., (1995) *Plant Mol Biol* 29:759-72, and Capana et al., (1994) *Plant Mol Biol* 25:681-91.

[0390] Examples of promoters that are inducible and that allow for spatiotemporal control of gene editing or gene expression may use a form of energy. The form of energy may include but is not limited to sound energy, electromagnetic radiation, chemical energy and/or thermal energy. Examples of inducible systems include tetracycline inducible promoters (Tet-On or Tet-Off), small molecule two-hybrid transcription activations systems (FKBP, ABA, etc), or light inducible systems (Phytochrome, LOV domains, or cryptochrome), such as a Light Inducible Transcriptional Effector (LITE) that direct changes in transcriptional activity in a sequence-specific manner. The components of a light inducible system may include a RNA targeting CRISPR enzyme, a light-responsive cytochrome heterodimer (e.g. from *Arabidopsis thaliana*), and a transcriptional activation/repression domain. Further examples of inducible DNA binding proteins and methods for their use are provided in U.S. 61/736,465 and U.S. 61/721,283, which is hereby incorporated by reference in its entirety.

[0391] In particular embodiments, transient or inducible expression can be achieved by using, for example, chemical-regulated promoters, i.e. whereby the application of an exogenous chemical induces gene expression. Modulating of gene expression can also be obtained by a chemical-repressible promoter, where application of the chemical represses gene expression. Chemical-inducible promoters include, but are not limited to, the maize In2-2 promoter, activated by benzene sulfonamide herbicide safeners (De Veylder et al., (1997) *Plant Cell Physiol* 38:568-77), the maize GST promoter (GST-11-27, WO93/01294), activated by hydrophobic electrophilic compounds used as pre-emergent herbicides, and the tobacco PR-1 a promoter (Ono et al., (2004) *Biosci Biotechnol Biochem* 68:803-7) activated by salicylic acid. Promoters which are regulated by antibiotics, such as tetracycline-inducible and tetracycline-repressible promoters (Gatz et al., (1991) *Mol Gen Genet* 227:229-37; U.S. Pat. Nos. 5,814,618 and 5,789,156) can also be used herein.

Translocation to and/or Expression in Specific Plant Organelles

[0392] The expression system may comprise elements for translocation to and/or expression in a specific plant organelle.

Chloroplast Targeting

[0393] In particular embodiments, it is envisaged that the RNA targeting CRISPR system is used to specifically modify expression and/or translation of chloroplast genes or to ensure expression in the chloroplast. For this purpose use is made of chloroplast transformation methods or compartmentalization of the RNA targeting CRISPR components to the chloroplast. For instance, the introduction of genetic modifications in the plastid genome can reduce biosafety issues such as gene flow through pollen.

[0394] Methods of chloroplast transformation are known in the art and include Particle bombardment, PEG treatment, and microinjection. Additionally, methods involving the translocation of transformation cassettes from the nuclear genome to the plastid can be used as described in WO2010061186.

[0395] Alternatively, it is envisaged to target one or more of the RNA targeting CRISPR components to the plant chloroplast. This is achieved by incorporating in the expression construct a sequence encoding a chloroplast transit peptide (CTP) or plastid transit peptide, operably linked to the 5' region of the sequence encoding the RNA targeting protein. The CTP is removed in a processing step during translocation into the chloroplast. Chloroplast targeting of expressed proteins is well known to the skilled artisan (see for instance Protein Transport into Chloroplasts, 2010, Annual Review of Plant Biology, Vol. 61: 157-180). In such embodiments it is also desired to target the one or more guide RNAs to the plant chloroplast. Methods and constructs which can be used for translocating guide RNA into the chloroplast by means of a chloroplast localization sequence are described, for instance, in US 20040142476, incorporated herein by reference. Such variations of constructs can be incorporated into the expression systems of the invention to efficiently translocate the RNA targeting-guide RNA(s).

Introduction of Polynucleotides Encoding the CRISPR—RNA Targeting System in Algal Cells.

[0396] Transgenic algae (or other plants such as rape) may be particularly useful in the production of vegetable oils or biofuels such as alcohols (especially methanol and ethanol) or other products. These may be engineered to express or overexpress high levels of oil or alcohols for use in the oil or biofuel industries.

[0397] U.S. Pat. No. 8,945,839 describes a method for engineering Micro-Algae (*Chlamydomonas reinhardtii* cells) species) using Cas9. Using similar tools, the methods of the RNA targeting CRISPR system described herein can be applied on *Chlamydomonas* species and other algae. In particular embodiments, RNA targeting protein and guide RNA(s) are introduced in algae expressed using a vector that expresses RNA targeting protein under the control of a constitutive promoter such as Hsp70A-Rbc S2 or Beta2-tubulin. Guide RNA is optionally delivered using a vector containing T7 promoter. Alternatively, RNA targeting mRNA and in vitro transcribed guide RNA can be delivered to algal cells. Electroporation protocols are available to the skilled person such as the standard recommended protocol from the GeneArt *Chlamydomonas* Engineering kit.

Introduction of Polynucleotides Encoding RNA Targeting Components in Yeast Cells

[0398] In particular embodiments, the invention relates to the use of the RNA targeting CRISPR system for RNA editing in yeast cells. Methods for transforming yeast cells which can be used to introduce polynucleotides encoding the RNA targeting CRISPR system components are well known to the artisan and are reviewed by Kawai et al., 2010, *Bioeng Bugs*. 2010 November-December; 1(6): 395-403). Non-limiting examples include transformation of yeast cells by lithium acetate treatment (which may further include carrier DNA and PEG treatment), bombardment or by electroporation.

Transient Expression of RNA Targeting CRISPR System Components in Plants and Plant Cell

[0399] In particular embodiments, it is envisaged that the guide RNA and/or RNA targeting gene are transiently expressed in the plant cell. In these embodiments, the RNA targeting CRISPR system can ensure modification of RNA target molecules only when both the guide RNA and the RNA targeting protein is present in a cell, such that gene expression can further be controlled. As the expression of the RNA targeting enzyme is transient, plants regenerated from such plant cells typically contain no foreign DNA. In particular embodiments the RNA targeting enzyme is stably expressed by the plant cell and the guide sequence is transiently expressed.

[0400] In particularly preferred embodiments, the RNA targeting CRISPR system components can be introduced in the plant cells using a plant viral vector (Scholthof et al. 1996, *Annu Rev Phytopathol*. 1996; 34:299-323). In further particular embodiments, said viral vector is a vector from a DNA virus. For example, geminivirus (e.g., cabbage leaf curl virus, bean yellow dwarf virus, wheat dwarf virus, tomato leaf curl virus, maize streak virus, tobacco leaf curl virus, or tomato golden mosaic virus) or nanovirus (e.g., Faba bean necrotic yellow virus). In other particular embodiments, said viral vector is a vector from an RNA virus. For

example, tobnavirus (e.g., tobacco rattle virus, tobacco mosaic virus), potexvirus potato virus X), or hordeivirus (e.g., barley stripe mosaic virus). The replicating genomes of plant viruses are non-integrative vectors, which is of interest in the context of avoiding the production of GMO plants.

[0401] In particular embodiments, the vector used for transient expression of RNA targeting CRISPR constructs is for instance a pEAQ vector, which is tailored for *Agrobacterium*-mediated transient expression (Sainsbury F. et al., *Plant Biotechnol J*. 2009 September; 7(7):682-93) in the protoplast. Precise targeting of genomic locations was demonstrated using a modified Cabbage Leaf Curl virus (CaL-CuV) vector to express gRNAs in stable transgenic plants expressing a CRISPR enzyme (Scientific Reports 5, Article number: 14926 (2015), doi:10.1038/srep14926).

[0402] In particular embodiments, double-stranded DNA fragments encoding the guide RNA and/or the RNA targeting gene can be transiently introduced into the plant cell. In such embodiments, the introduced double-stranded DNA fragments are provided in sufficient quantity to modify RNA molecule(s) in the cell but do not persist after a contemplated period of time has passed or after one or more cell divisions. Methods for direct DNA transfer in plants are known by the skilled artisan (see for instance Davey et al. *Plant Mol Biol*. 1989 September; 13(3):273-85.)

[0403] In other embodiments, an RNA polynucleotide encoding the RNA targeting protein is introduced into the plant cell, which is then translated and processed by the host cell generating the protein in sufficient quantity to modify the RNA molecule(s) cell (in the presence of at least one guide RNA) but which does not persist after a contemplated period of time has passed or after one or more cell divisions. Methods for introducing mRNA to plant protoplasts for transient expression are known by the skilled artisan (see for instance in Gallie, *Plant Cell Reports* (1993), 13; 119-122). Combinations of the different methods described above are also envisaged.

Delivery of RNA Targeting CRISPR Components to the Plant Cell

[0404] In particular embodiments, it is of interest to deliver one or more components of the RNA targeting CRISPR system directly to the plant cell. This is of interest, inter alia, for the generation of non-transgenic plants (see below). In particular embodiments, one or more of the RNA targeting components is prepared outside the plant or plant cell and delivered to the cell. For instance in particular embodiments, the RNA targeting protein is prepared in vitro prior to introduction to the plant cell. RNA targeting protein can be prepared by various methods known by one of skill in the art and include recombinant production. After expression, the RNA targeting protein is isolated, refolded if needed, purified and optionally treated to remove any purification tags, such as a His-tag. Once crude, partially purified, or more completely purified RNA targeting protein is obtained, the protein may be introduced to the plant cell.

[0405] In particular embodiments, the RNA targeting protein is mixed with guide RNA targeting the RNA of interest to form a pre-assembled ribonucleoprotein.

[0406] The individual components or pre-assembled ribonucleoprotein can be introduced into the plant cell via electroporation, by bombardment with RNA targeting-associated gene product coated particles, by chemical transfection or by some other means of transport across a cell

membrane. For instance, transfection of a plant protoplast with a pre-assembled CRISPR ribonucleoprotein has been demonstrated to ensure targeted modification of the plant genome (as described by Woo of al. *Nature Biotechnology*, 2015; DOI: 10.1038/nbt.3389). These methods can be modified to achieve targeted modification of RNA molecules in the plants.

[0407] In particular embodiments, the RNA targeting CRISPR system components are introduced into the plant cells using particles. The components, either as protein or nucleic acid or in a combination thereof, can be uploaded onto or packaged in particles and applied to the plants (such as for instance described in WO 2008042156 and US 20130185823). In particular, embodiments of the invention comprise particles uploaded with or packed with DNA molecule(s) encoding the RNA targeting protein, DNA molecules encoding the guide RNA and/or isolated guide RNA as described in WO2015089419.

[0408] Further means of introducing one or more components of the RNA targeting CRISPR system to the plant cell is by using cell penetrating peptides (CPP). Accordingly, in particular, embodiments the invention comprises compositions comprising a cell penetrating peptide linked to an RNA targeting protein. In particular embodiments of the present invention, an RNA targeting protein and/or guide RNA(s) is coupled to one or more CPPs to effectively transport them inside plant protoplasts (Ramakrishna (2014, *Genome Res.* 2014 June; 24(6):1020-7 for Cas9 in human cells). In other embodiments, the RNA targeting gene and/or guide RNA(s) are encoded by one or more circular or non-circular DNA molecule(s) which are coupled to one or more CPPs for plant protoplast delivery. The plant protoplasts are then regenerated to plant cells and further to plants. CPPs are generally described as short peptides of fewer than 35 amino acids either derived from proteins or from chimeric sequences which are capable of transporting biomolecules across cell membrane in a receptor independent manner. CPP can be cationic peptides, peptides having hydrophobic sequences, amphipatic peptides, peptides having proline-rich and antimicrobial sequence, and chimeric or bipartite peptides (Pooga and Langel 2005). CPPs are able to penetrate biological membranes and as such trigger the movement of various biomolecules across cell membranes into the cytoplasm and to improve their intracellular routing, and hence facilitate interaction of the biomolecule with the target. Examples of CPP include amongst others: Tat, a nuclear transcriptional activator protein required for viral replication by HIV type1, penetratin, Kaposi fibroblast growth factor (FGF) signal peptide sequence, integrin β 3 signal peptide sequence; polyarginine peptide Arg₅ sequence, Guanine rich-molecular transporters, sweet arrow peptide, etc. . . .

Target RNA Envisaged for Plant, Algae or Fungal Applications

[0409] The target RNA, i.e. the RNA of interest, is the RNA to be targeted by the present invention leading to the recruitment to, and the binding of the RNA targeting protein at, the target site of interest on the target RNA. The target RNA may be any suitable form of RNA. This may include, in some embodiments, mRNA. In other embodiments, the target RNA may include transfer RNA (tRNA) or ribosomal RNA (rRNA). In other embodiments the target RNA may include interfering RNA (RNAi), microRNA (miRNA), microswitches, microzymes, satellite RNAs and RNA

viruses. The target RNA may be located in the cytoplasm of the plant cell, or in the cell nucleus or in a plant cell organelle such as a mitochondrion, chloroplast or plastid.

[0410] In particular embodiments, the RNA targeting CRISPR system is used to cleave RNA or otherwise inhibit RNA expression.

Use of RNA Targeting CRISPR System for Modulating Plant Gene Expression Via RNA Modulation

[0411] The RNA targeting protein may also be used, together with a suitable guide RNA, to target gene expression, via control of RNA processing. The control of RNA processing may include RNA processing reactions such as RNA splicing, including alternative splicing; viral replication (in particular of plant viruses, including viroids in plants and tRNA biosynthesis. The RNA targeting protein in combination with a suitable guide RNA may also be used to control RNA activation (RNAa). RNAa leads to the promotion of gene expression, so control of gene expression may be achieved that way through disruption or reduction of RNAa and thus less promotion of gene expression.

[0412] The RNA targeting effector protein of the invention can further be used for antiviral activity in plants, in particular against RNA viruses. The effector protein can be targeted to the viral RNA using a suitable guide RNA selective for a selected viral RNA sequence. In particular, the effector protein may be an active nuclease that cleaves RNA, such as single stranded RNA, provided is therefore the use of an RNA targeting effector protein of the invention as an antiviral agent. Examples of viruses that can be counteracted in this way include, but are not limited to, Tobacco mosaic virus (TMV), Tomato spotted wilt virus (TSW V), Cucumber mosaic virus (CMV), Potato virus Y (PVY), Cauliflower mosaic virus (CaMV) (RT virus), Plum pox virus (PPV), Brome mosaic virus (BMV) and Potato virus X (PVX).

[0413] Examples of modulating RNA expression in plants, algae or fungi, as an alternative of targeted gene modification are described herein further.

[0414] Of particular interest is the regulated control of gene expression through regulated cleavage of mRNA. This can be achieved by placing elements of the RNA targeting under the control of regulated promoters as described herein. Use of the RNA Targeting CRISPR System to Restore the Functionality of tRNA Molecules.

[0415] Pring et al describe RNA editing in plant mitochondria and chloroplasts that alters mRNA sequences to code for different proteins than the DNA. (*Plant Mol. Biol.* (1993) 21 (6): 1163-1170. doi:10.1007/BF00023611). In particular embodiments of the invention, the elements of the RNA targeting CRISPR system specifically targeting mitochondrial and chloroplast mRNA can be introduced in a plant or plant cell to express different proteins in such plant cell organelles mimicking the processes occurring in vivo.

Use of the RNA Targeting CRISPR System as an Alternative to RNA Interference to Inhibit RNA Expression.

[0416] The RNA targeting CRISPR system has uses similar to RNA inhibition or RNA interference, thus can also be substituted for such methods. In particular embodiment, the methods of the present invention include the use of the RNA targeting CRISPR as a substitute for e.g. an interfering ribonucleic acid (such as an siRNA or shRNA or a dsRNA).

Examples of inhibition of RNA expression in plants, algae or fungi as an alternative of targeted gene modification are described herein further.

Use of the RNA Targeting CRISPR System to Control RNA Interference.

[0417] Control over interfering RNA or miRNA may help reduce off-target effects (OTE) seen with those approaches by reducing the longevity of the interfering RNA or miRNA in vivo or in vitro. In particular embodiments, the target RNA may include interfering RNA, i.e. RNA involved in an RNA interference pathway, such as shRNA, siRNA and so forth. In other embodiments, the target RNA may include microRNA (miRNA) or double stranded RNA (dsRNA).

[0418] In other particular embodiments, if the RNA targeting protein and suitable guide RNA(s) are selectively expressed (for example spatially or temporally under the control of a regulated promoter, for example a tissue- or cell cycle-specific promoter and/or enhancer) this can be used to "protect" the cells or systems (in vivo or in vitro) from RNAi in those cells. This may be useful in neighbouring tissues or cells where RNAi is not required or for the purposes of comparison of the cells or tissues where the effector protein and suitable guide are and are not expressed (i.e. where the RNAi is not controlled and where it is, respectively). The RNA targeting protein may be used to control or bind to molecules comprising or consisting of RNA, such as ribozymes, ribosomes or riboswitches. In embodiments of the invention, the guide RNA can recruit the RNA targeting protein to these molecules so that the RNA targeting protein is able to bind to them.

[0419] The RNA targeting CRISPR system of the invention can be applied in areas of in-planta RNAi technologies, without undue experimentation, from this disclosure, including insect pest management, plant disease management and management of herbicide resistance, as well as in plant assay and for other applications (see, for instance Kim et al., in *Pesticide Biochemistry and Physiology* (Impact Factor: 2.01). 01/2015; 120. DOI: 10.1016/j.pestbp.2015.01.002; Sharma et al. in *Academic Journals* (2015), Vol. 12(18) pp2303-2312); Green J. M, in *Pest Management Science*, Vol 70(9), pp 1351-1357), because the present application provides the foundation for informed engineering of the system.

Use of RNA Targeting CRISPR System to Modify Riboswitches and Control Metabolic Regulation in Plants, Algae and Fungi

[0420] Riboswitches (also known as aptozymes) are regulatory segments of messenger RNA that bind small molecules and in turn regulate gene expression. This mechanism allows the cell to sense the intracellular concentration of these small molecules. A particular riboswitch typically regulates its adjacent gene by altering the transcription, the translation or the splicing of this gene. Thus, in particular embodiments of the present invention, control of riboswitch activity is envisaged through the use of the RNA targeting protein in combination with a suitable guide RNA to target the riboswitch. This may be through cleavage of, or binding to, the riboswitch. In particular embodiments, reduction of riboswitch activity is envisaged. Recently, a riboswitch that binds thiamin pyrophosphate (TPP) was characterized and found to regulate thiamin biosynthesis in plants and algae. Furthermore it appears that this element is an essential

regulator of primary metabolism in plants (Bocobza and Aharoni, *Plant J.* 2014 August; 79(4):693-703. doi: 10.1111/tbj.12540. Epub 2014 Jun. 17). TPP riboswitches are also found in certain fungi, such as in *Neurospora crassa*, where it controls alternative splicing to conditionally produce an Upstream Open Reading Frame (uORF), thereby affecting the expression of downstream genes (Cheah M T et al., (2007) *Nature* 447 (7143): 497-500. doi:10.1038/nature05769) The RNA targeting CRISPR system described herein may be used to manipulate the endogenous riboswitch activity in plants, algae or fungi and as such alter the expression of downstream genes controlled by it. In particular embodiments, the RNA targeting CRISPR system may be used in assaying riboswitch function in vivo or in vitro and in studying its relevance for the metabolic network. In particular embodiments the RNA targeting CRISPR system may potentially be used for engineering of riboswitches as metabolite sensors in plants and platforms for gene control.

Use of RNA Targeting CRISPR System in RNAi Screens for Plants, Algae or Fungi

[0421] Identifying gene products whose knockdown is associated with phenotypic changes, biological pathways can be interrogated and the constituent parts identified, via RNAi screens. In particular embodiments of the invention, control may also be exerted over or during these screens by use of the Guide 29 or Guide 30 protein and suitable guide RNA described herein to remove or reduce the activity of the RNAi in the screen and thus reinstate the activity of the (previously interfered with) gene product (by removing or reducing the interference/repression).

Use of RNA Targeting Proteins for Visualization of RNA Molecules In Vivo and In Vitro

[0422] In particular embodiments, the invention provides a nucleic acid binding system. In situ hybridization of RNA with complementary probes is a powerful technique. Typically fluorescent DNA oligonucleotides are used to detect nucleic acids by hybridization. Increased efficiency has been attained by certain modifications, such as locked nucleic acids (LNAs), but there remains a need for efficient and versatile alternatives. As such, labelled elements of the RNA targeting system can be used as an alternative for efficient and adaptable system for in situ hybridization

Further Applications of the RNA Targeting CRISPR System in Plants and Yeasts

Use of RNA Targeting CRISPR System in Biofuel Production

[0423] The term "biofuel" as used herein is an alternative fuel made from plant and plant-derived resources. Renewable biofuels can be extracted from organic matter whose energy has been obtained through a process of carbon fixation or are made through the use or conversion of biomass. This biomass can be used directly for biofuels or can be converted to convenient energy containing substances by thermal conversion, chemical conversion, and biochemical conversion. This biomass conversion can result in fuel in solid, liquid, or gas form. There are two types of biofuels: bioethanol and biodiesel. Bioethanol is mainly produced by the sugar fermentation process of cellulose (starch), which is mostly derived from maize and sugar cane.

Biodiesel on the other hand is mainly produced from oil crops such as rapeseed, palm, and soybean. Biofuels are used mainly for transportation.

Enhancing Plant Properties for Biofuel Production

[0424] In particular embodiments, the methods using the RNA targeting CRISPR system as described herein are used to alter the properties of the cell wall in order to facilitate access by key hydrolysing agents for a more efficient release of sugars for fermentation. In particular embodiments, the biosynthesis of cellulose and/or lignin are modified. Cellulose is the major component of the cell wall. The biosynthesis of cellulose and lignin are co-regulated. By reducing the proportion of lignin in a plant the proportion of cellulose can be increased. In particular embodiments, the methods described herein are used to downregulate lignin biosynthesis in the plant so as to increase fermentable carbohydrates. More particularly, the methods described herein are used to downregulate at least a first lignin biosynthesis gene selected from the group consisting of 4-coumarate 3-hydroxylase (C3H), phenylalanine ammonia-lyase (PAL), cinnamate 4-hydroxylase (C4H), hydroxycinnamoyl transferase (HCT), caffeic acid O-methyltransferase (COMT), caffeoyl CoA 3-O-methyltransferase (CCoAOMT), ferulate 5-hydroxylase (F5H), cinnamyl alcohol dehydrogenase (CAD), cinnamoyl CoA-reductase (CCR), 4-coumarate-CoA ligase (4CL), monolignol-lignin-specific glycosyltransferase, and aldehyde dehydrogenase (ALDH) as disclosed in WO 2008064289 A2.

[0425] In particular embodiments, the methods described herein are used to produce plant mass that produces lower levels of acetic acid during fermentation (see also WO 2010096488).

Modifying Yeast for Biofuel Production

[0426] In particular embodiments, the RNA targeting enzyme provided herein is used for bioethanol production by recombinant micro-organisms. For instance, RNA targeting enzymes can be used to engineer micro-organisms, such as yeast, to generate biofuel or biopolymers from fermentable sugars and optionally to be able to degrade plant-derived lignocellulose derived from agricultural waste as a source of fermentable sugars. More particularly, the invention provides methods whereby the RNA targeting CRISPR complex is used to modify the expression of endogenous genes required for biofuel production and/or to modify endogenous genes which may interfere with the biofuel synthesis. More particularly the methods involve stimulating the expression in a micro-organism such as a yeast of one or more nucleotide sequence encoding enzymes involved in the conversion of pyruvate to ethanol or another product of interest. In particular embodiments the methods ensure the stimulation of expression of one or more enzymes which allows the micro-organism to degrade cellulose, such as a cellulase. In yet further embodiments, the RNA targeting CRISPR complex is used to suppress endogenous metabolic pathways which compete with the biofuel production pathway.

Modifying Algae and Plants for Production of Vegetable Oils or Biofuels

[0427] Transgenic algae or other plants such as rape may be particularly useful in the production of vegetable oils or

biofuels such as alcohols (especially methanol and ethanol), for instance. These may be engineered to express or over-express high levels of oil or alcohols for use in the oil or biofuel industries.

[0428] U.S. Pat. No. 8,945,839 describes a method for engineering Micro-Algae (*Chlamydomonas reinhardtii* cells) species) using Cas9. Using similar tools, the methods of the RNA targeting CRISPR system described herein can be applied on *Chlamydomonas* species and other algae. In particular embodiments, the RNA targeting effector protein and guide RNA are introduced in algae expressed using a vector that expresses the RNA targeting effector protein under the control of a constitutive promoter such as Hsp70A-Rbc S2 or Beta2-tubulin. Guide RNA will be delivered using a vector containing T7 promoter. Alternatively, in vitro transcribed guide RNA can be delivered to algae cells. Electroporation protocol follows standard recommended protocol from the GeneArt *Chlamydomonas* Engineering kit.

Particular Applications of the RNA Targeting Enzymes in Plants

[0429] In particular embodiments, present invention can be used as a therapy for virus removal in plant systems as it is able to cleave viral RNA. Previous studies in human systems have demonstrated the success of utilizing CRISPR in targeting the single strand RNA virus, hepatitis C (A. Price, et al., Proc. Natl. Acad. Sci, 2015). These methods may also be adapted for using the RNA targeting CRISPR system in plants.

Improved Plants

[0430] The present invention also provides plants and yeast cells obtainable and obtained by the methods provided herein. The improved plants obtained by the methods described herein may be useful in food or feed production through the modified expression of genes which, for instance ensure tolerance to plant pests, herbicides, drought, low or high temperatures, excessive water, etc.

[0431] The improved plants obtained by the methods described herein, especially crops and algae may be useful in food or feed production through expression of, for instance, higher protein, carbohydrate, nutrient or vitamin levels than would normally be seen in the wildtype. In this regard, improved plants, especially pulses and tubers are preferred.

[0432] Improved algae or other plants such as rape may be particularly useful in the production of vegetable oils or biofuels such as alcohols (especially methanol and ethanol), for instance. These may be engineered to express or over-express high levels of oil or alcohols for use in the oil or biofuel industries.

[0433] The invention also provides for improved parts of a plant. Plant parts include, but are not limited to, leaves, stems, roots, tubers, seeds, endosperm, ovule, and pollen. Plant parts as envisaged herein may be viable, nonviable, regeneratable, and/or non-regeneratable.

[0434] It is also encompassed herein to provide plant cells and plants generated according to the methods of the invention. Gametes, seeds, embryos, either zygotic or somatic, progeny or hybrids of plants comprising the genetic modification, which are produced by traditional breeding methods, are also included within the scope of the present

invention. Such plants may contain a heterologous or foreign DNA sequence inserted at or instead of a target sequence. Alternatively, such plants may contain only an alteration (mutation, deletion, insertion, substitution) in one or more nucleotides. As such, such plants will only be different from their progenitor plants by the presence of the particular modification.

[0435] In an embodiment of the invention, a Cas13b system is used to engineer pathogen resistant plants, for example by creating resistance against diseases caused by bacteria, fungi or viruses. In certain embodiments, pathogen resistance can be accomplished by engineering crops to produce a Cas13b system that will be ingested by an insect pest, leading to mortality. In an embodiment of the invention, a Cas13b system is used to engineer abiotic stress tolerance. In another embodiment, a Cas13b system is used to engineer drought stress tolerance or salt stress tolerance, or cold or heat stress tolerance. Younis et al. 2014, Int. J. Biol. Sci. 10; 1150 reviewed potential targets of plant breeding methods, all of which are amenable to correction or improvement through use of a Cas13b system described herein. Some non-limiting target crops include *Arabidopsis thaliana*, *Zea mays*, *Oryza sativa* L., *Prunus domestica* L., *Gossypium hirsutum*, *Nicotiana rustica*, *Zea mays*, *Medicago sativa*, *Nicotiana benthamiana* and *Arabidopsis thaliana*.

[0436] In an embodiment of the invention, a Cas13b system is used for management of crop pests. For example, a Cas13b system operable in a crop pest can be expressed from a plant host or transferred directly to the target, for example using a viral vector.

[0437] In an embodiment, the invention provides a method of efficiently producing homozygous organisms from a heterozygous non-human starting organism. In an embodiment, the invention is used in plant breeding. In another embodiment, the invention is used in animal breeding. In such embodiments, a homozygous organism such as a plant or animal is made by preventing or suppressing recombination by interfering with at least one target gene involved in double strand breaks, chromosome pairing and/or strand exchange.

Application of the CAS13B Proteins in Optimized Functional RNA Targeting Systems

[0438] In an aspect the invention provides a system for specific delivery of functional components to the RNA environment. This can be ensured using the CRISPR systems comprising the RNA targeting effector proteins of the present invention which allow specific targeting of different components to RNA. More particularly such components include activators or repressors, such as activators or repressors of RNA translation, degradation, etc. Applications of this system are described elsewhere herein.

[0439] According to one aspect the invention provides non-naturally occurring or engineered composition comprising a guide RNA comprising a guide sequence capable of hybridizing to a target sequence in a genomic locus of interest in a cell, wherein the guide RNA is modified by the insertion of one or more distinct RNA sequence(s) that bind an adaptor protein. In particular embodiments, the RNA sequences may bind to two or more adaptor proteins (e.g. aptamers), and wherein each adaptor protein is associated with one or more functional domains. The guide RNAs of the Cas13b enzymes described herein are shown to be

amenable to modification of the guide sequence. In particular embodiments, the guide RNA is modified by the insertion of distinct RNA sequence(s) 5' of the direct repeat, within the direct repeat, or 3' of the guide sequence. When there is more than one functional domain, the functional domains can be same or different, e.g., two of the same or two different activators or repressors. In an aspect the invention provides a herein-discussed composition, wherein the one or more functional domains are attached to the RNA targeting enzyme so that upon binding to the target RNA the functional domain is in a spatial orientation allowing for the functional domain to function in its attributed function; In an aspect the invention provides a herein-discussed composition, wherein the composition comprises a CRISPR-Cas complex having at least three functional domains, at least one of which is associated with the RNA targeting enzyme and at least two of which are associated with the gRNA.

[0440] Accordingly, In an aspect the invention provides non-naturally occurring or engineered CRISPR-Cas complex composition comprising the guide RNA as herein-discussed and a CRISPR enzyme which is an RNA targeting enzyme, wherein optionally the RNA targeting enzyme comprises at least one mutation, such that the RNA targeting enzyme has no more than 5% of the nuclease activity of the enzyme not having the at least one mutation, and optionally one or more comprising at least one or more nuclear localization sequences. In particular embodiments, the guide RNA is additionally or alternatively modified so as to still ensure binding of the RNA targeting enzyme but to prevent cleavage by the RNA targeting enzyme (as detailed elsewhere herein).

[0441] In particular embodiments, the RNA targeting enzyme is a Cas13b enzyme which has a diminished nuclease activity of at least 97%, or 100% as compared with the Cas13b enzyme not having the at least one mutation. In an aspect the invention provides a herein-discussed composition, wherein the Cas13b enzyme comprises two or more mutations. The mutations may be selected from mutations of one or more of the following amino acid residues: R116, H121, R1177, and H1182, such as for instance one or more of the following mutations: R116A, H121A, R1177A, and H1182A, according to *Bergeyella zoohelcum* Cas13b protein or a corresponding position in an ortholog.

[0442] In particular embodiments, an RNA targeting system is provided as described herein above comprising two or more functional domains. In particular embodiments, the two or more functional domains are heterologous functional domain. In particular embodiments, the system comprises an adaptor protein which is a fusion protein comprising a functional domain, the fusion protein optionally comprising a linker between the adaptor protein and the functional domain. In particular embodiments, the linker includes a GlySer linker. Additionally or alternatively, one or more functional domains are attached to the RNA effector protein by way of a linker, optionally a GlySer linker. In particular embodiments, the one or more functional domains are attached to the RNA targeting enzyme through one or both of the HEPN domains.

[0443] In an aspect the invention provides a herein-discussed composition, wherein the one or more functional domains associated with the adaptor protein or the RNA targeting enzyme is a domain capable of activating or repressing RNA translation. In an aspect the invention provides a herein-discussed composition, wherein at least

one of the one or more functional domains associated with the adaptor protein have one or more activities comprising methylase activity, demethylase activity, transcription activation activity, transcription repression activity, transcription release factor activity, histone modification activity, DNA integration activity RNA cleavage activity, DNA cleavage activity or nucleic acid binding activity, or molecular switch activity or chemical inducibility or light inducibility.

[0444] In an aspect the invention provides a herein-discussed composition comprising an aptamer sequence. In particular embodiments, the aptamer sequence is two or more aptamer sequences specific to the same adaptor protein. In an aspect the invention provides a herein-discussed composition, wherein the aptamer sequence is two or more aptamer sequences specific to different adaptor protein. In an aspect the invention provides a herein-discussed composition, wherein the adaptor protein comprises MS2, PP7, Q β , F2, GA, fr, JP501, M12, R17, BZ13, JP34, JP500, KU1, M11, MX1, TW18, VK, SP, FI, ID2, NL95, TW19, AP205, ϕ Cb5, ϕ Cb8r, ϕ Cb12r, ϕ Cb23r, 7s, PRR1. Accordingly, in particular embodiments, the aptamer is selected from a binding protein specifically binding any one of the adaptor proteins listed above. In an aspect the invention provides a herein-discussed composition, wherein the cell is a eukaryotic cell. In an aspect the invention provides a herein-discussed composition, wherein the eukaryotic cell is a mammalian cell, a plant cell or a yeast cell, whereby the mammalian cell is optionally a mouse cell. In an aspect the invention provides a herein-discussed composition, wherein the mammalian cell is a human cell.

[0445] In an aspect the invention provides a herein above-discussed composition wherein there is more than one gRNA, and the gRNAs target different sequences whereby when the composition is employed, there is multiplexing. In an aspect the invention provides a composition wherein there is more than one gRNA modified by the insertion of distinct RNA sequence(s) that bind to one or more adaptor proteins.

[0446] In an aspect the invention provides a herein-discussed composition wherein one or more adaptor proteins associated with one or more functional domains is present and bound to the distinct RNA sequence(s) inserted into the guide RNA(s).

[0447] In an aspect the invention provides a herein-discussed composition wherein the guide RNA is modified to have at least one non-coding functional loop; e.g., wherein the at least one non-coding functional loop is repressive; for instance, wherein at least one non-coding functional loop comprises Alu.

[0448] In an aspect the invention provides a method for modifying gene expression comprising the administration to a host or expression in a host *in vivo* of one or more of the compositions as herein-discussed.

[0449] In an aspect the invention provides a herein-discussed method comprising the delivery of the composition or nucleic acid molecule(s) coding therefor, wherein said nucleic acid molecule(s) are operatively linked to regulatory sequence(s) and expressed *in vivo*. In an aspect the invention provides a herein-discussed method wherein the expression *in vivo* is via a lentivirus, an adenovirus, or an AAV.

[0450] In an aspect the invention provides a mammalian cell line of cells as herein-discussed, wherein the cell line is, optionally, a human cell line or a mouse cell line. In an

aspect the invention provides a transgenic mammalian model, optionally a mouse, wherein the model has been transformed with a herein-discussed composition or is a progeny of said transformant.

[0451] In an aspect the invention provides a nucleic acid molecule(s) encoding guide RNA or the RNA targeting CRISPR-Cas complex or the composition as herein-discussed. In an aspect the invention provides a vector comprising: a nucleic acid molecule encoding a guide RNA (gRNA) comprising a guide sequence capable of hybridizing to a target sequence in a genomic locus of interest in a cell, wherein the direct repeat of the gRNA is modified by the insertion of distinct RNA sequence(s) that bind(s) to two or more adaptor proteins, and wherein each adaptor protein is associated with one or more functional domains; or, wherein the gRNA is modified to have at least one non-coding functional loop. In an aspect the invention provides vector(s) comprising nucleic acid molecule(s) encoding: non-naturally occurring or engineered CRISPR-Cas complex composition comprising the gRNA herein-discussed, and an RNA targeting enzyme, wherein optionally the RNA targeting enzyme comprises at least one mutation, such that the RNA targeting enzyme has no more than 5% of the nuclease activity of the RNA targeting enzyme not having the at least one mutation, and optionally one or more comprising at least one or more nuclear localization sequences. In an aspect a vector can further comprise regulatory element(s) operable in a eukaryotic cell operably linked to the nucleic acid molecule encoding the guide RNA (gRNA) and/or the nucleic acid molecule encoding the RNA targeting enzyme and/or the optional nuclear localization sequence(s).

[0452] In one aspect, the invention provides a kit comprising one or more of the components described herein-above. In some embodiments, the kit comprises a vector system as described above and instructions for using the kit.

[0453] In an aspect the invention provides a method of screening for gain of function (GOF) or loss of function (LOF) or for screening non-coding RNAs or potential regulatory regions (e.g. enhancers, repressors) comprising the cell line of as herein-discussed or cells of the model herein-discussed containing or expressing the RNA targeting enzyme and introducing a composition as herein-discussed into cells of the cell line or model, whereby the gRNA includes either an activator or a repressor, and monitoring for GOF or LOF respectively as to those cells as to which the introduced gRNA includes an activator or as to those cells as to which the introduced gRNA includes a repressor.

[0454] In an aspect the invention provides a library of non-naturally occurring or engineered compositions, each comprising a RNA targeting CRISPR guide RNA (gRNA) comprising a guide sequence capable of hybridizing to a target RNA sequence of interest in a cell, an RNA targeting enzyme, wherein the RNA targeting enzyme comprises at least one mutation, such that the RNA targeting enzyme has no more than 5% of the nuclease activity of the RNA targeting enzyme not having the at least one mutation, wherein the gRNA is modified by the insertion of distinct RNA sequence(s) that bind to one or more adaptor proteins, and wherein the adaptor protein is associated with one or more functional domains, wherein the composition comprises one or more or two or more adaptor proteins, wherein the each protein is associated with one or more functional domains, and wherein the gRNAs comprise a genome wide library comprising a plurality of RNA targeting guide RNAs

(gRNAs). In an aspect the invention provides a library as herein-discussed, wherein the RNA targeting RNA targeting enzyme has a diminished nuclease activity of at least 97%, or 100% as compare with the RNA targeting enzyme not having the at least one mutation. In an aspect the invention provides a library as herein-discussed, wherein the adaptor protein is a fusion protein comprising the functional domain. In an aspect the invention provides a library as herein discussed, wherein the gRNA is not modified by the insertion of distinct RNA sequence(s) that bind to the one or two or more adaptor proteins. In an aspect the invention provides a library as herein discussed, wherein the one or two or more functional domains are associated with the RNA targeting enzyme. In an aspect the invention provides a library as herein discussed, wherein the cell population of cells is a population of eukaryotic cells. In an aspect the invention provides a library as herein discussed, wherein the eukaryotic cell is a mammalian cell, a plant cell or a yeast cell. In an aspect the invention provides a library as herein discussed, wherein the mammalian cell is a human cell. In an aspect the invention provides a library as herein discussed, wherein the population of cells is a population of embryonic stem (ES) cells.

[0455] In an aspect the invention provides a library as herein discussed, wherein the targeting is of about 100 or more RNA sequences. In an aspect the invention provides a library as herein discussed, wherein the targeting is of about 1000 or more RNA sequences. In an aspect the invention provides a library as herein discussed, wherein the targeting is of about 20,000 or more sequences. In an aspect the invention provides a library as herein discussed, wherein the targeting is of the entire transcriptome. In an aspect the invention provides a library as herein discussed, wherein the targeting is of a panel of target sequences focused on a relevant or desirable pathway. In an aspect the invention provides a library as herein discussed, wherein the pathway is an immune pathway. In an aspect the invention provides a library as herein discussed, wherein the pathway is a cell division pathway.

[0456] In one aspect, the invention provides a method of generating a model eukaryotic cell comprising a gene with modified expression. In some embodiments, a disease gene is any gene associated an increase in the risk of having or developing a disease. In some embodiments, the method comprises (a) introducing one or more vectors encoding the components of the system described herein above into a eukaryotic cell, and (b) allowing a CRISPR complex to bind to a target polynucleotide so as to modify expression of a gene, thereby generating a model eukaryotic cell comprising modified gene expression.

[0457] The structural information provided herein allows for interrogation of guide RNA interaction with the target RNA and the RNA targeting enzyme permitting engineering or alteration of guide RNA structure to optimize functionality of the entire RNA targeting CRISPR-Cas system. For example, the guide RNA may be extended, without colliding with the RNA targeting protein by the insertion of adaptor proteins that can bind to RNA. These adaptor proteins can further recruit effector proteins or fusions which comprise one or more functional domains.

[0458] An aspect of the invention is that the above elements are comprised in a single composition or comprised in

individual compositions. These compositions may advantageously be applied to a host to elicit a functional effect on the genomic level.

[0459] The skilled person will understand that modifications to the guide RNA which allow for binding of the adapter+functional domain but not proper positioning of the adapter+functional domain (e.g. due to steric hindrance within the three dimensional structure of the CRISPR complex) are modifications which are not intended. The one or more modified guide RNA may be modified, by introduction of a distinct RNA sequence(s) 5' of the direct repeat, within the direct repeat, or 3' of the guide sequence.

[0460] The modified guide RNA, the inactivated RNA targeting enzyme (with or without functional domains), and the binding protein with one or more functional domains, may each individually be comprised in a composition and administered to a host individually or collectively. Alternatively, these components may be provided in a single composition for administration to a host. Administration to a host may be performed via viral vectors known to the skilled person or described herein for delivery to a host (e.g. lentiviral vector, adenoviral vector, AAV vector). As explained herein, use of different selection markers (e.g. for lentiviral gRNA selection) and concentration of gRNA (e.g. dependent on whether multiple gRNAs are used) may be advantageous for eliciting an improved effect.

[0461] Using the provided compositions, the person skilled in the art can advantageously and specifically target single or multiple loci with the same or different functional domains to elicit one or more genomic events. The compositions may be applied in a wide variety of methods for screening in libraries in cells and functional modeling in vivo (e.g. gene activation of lincRNA and identification of function; gain-of-function modeling; loss-of-function modeling; the use the compositions of the invention to establish cell lines and transgenic animals for optimization and screening purposes).

[0462] The current invention comprehends the use of the compositions of the current invention to establish and utilize conditional or inducible CRISPR RNA targeting events. (See, e.g., Platt et al., *Cell* (2014), <http://dx.doi.org/10.1016/j.cell.2014.09.014>, or PCT patent publications cited herein, such as WO 2014/093622 (PCT/US2013/074667), which are not believed prior to the present invention or application). For example, the target cell comprises RNA targeting CRISPR enzyme conditionally or inducibly (e.g. in the form of Cre dependent constructs) and/or the adaptor protein conditionally or inducibly and, on expression of a vector introduced into the target cell, the vector expresses that which induces or gives rise to the condition of s RNA targeting enzyme expression and/or adaptor expression in the target cell. By applying the teaching and compositions of the current invention with the known method of creating a CRISPR complex, inducible gene expression affected by functional domains are also an aspect of the current invention. Alternatively, the adaptor protein may be provided as a conditional or inducible element with a conditional or inducible s RNA targeting enzyme to provide an effective model for screening purposes, which advantageously only requires minimal design and administration of specific gRNAs for a broad number of applications.

Guide RNA According to the Invention Comprising a Dead Guide Sequence

[0463] In one aspect, the invention provides guide sequences which are modified in a manner which allows for formation of the CRISPR complex and successful binding to the target, while at the same time, not allowing for successful nuclease activity (i.e. without nuclease activity/without indel activity). For matters of explanation such modified guide sequences are referred to as “dead guides” or “dead guide sequences”. These dead guides or dead guide sequences can be thought of as catalytically inactive or conformationally inactive with regard to nuclease activity. Indeed, dead guide sequences may not sufficiently engage in productive base pairing with respect to the ability to promote catalytic activity or to distinguish on-target and off-target binding activity. Briefly, the assay involves synthesizing a CRISPR target RNA and guide RNAs comprising mismatches with the target RNA, combining these with the RNA targeting enzyme and analyzing cleavage based on gels based on the presence of bands generated by cleavage products, and quantifying cleavage based upon relative band intensities.

[0464] Hence, in a related aspect, the invention provides a non-naturally occurring or engineered composition RNA targeting CRISPR-Cas system comprising a functional RNA targeting as described herein, and guide RNA (gRNA) wherein the gRNA comprises a dead guide sequence whereby the gRNA is capable of hybridizing to a target sequence such that the RNA targeting CRISPR-Cas system is directed to a genomic locus of interest in a cell without detectable RNA cleavage activity of a non-mutant RNA targeting enzyme of the system. It is to be understood that any of the gRNAs according to the invention as described herein elsewhere may be used as dead gRNAs/gRNAs comprising a dead guide sequence as described herein below. Any of the methods, products, compositions and uses as described herein elsewhere is equally applicable with the dead gRNAs/gRNAs comprising a dead guide sequence as further detailed below. By means of further guidance, the following particular aspects and embodiments are provided.

[0465] The ability of a dead guide sequence to direct sequence-specific binding of a CRISPR complex to an RNA target sequence may be assessed by any suitable assay. For example, the components of a CRISPR system sufficient to form a CRISPR complex, including the dead guide sequence to be tested, may be provided to a host cell having the corresponding target sequence, such as by transfection with vectors encoding the components of the CRISPR sequence, followed by an assessment of preferential cleavage within the target sequence. For instance, cleavage of a target RNA polynucleotide sequence may be evaluated in a test tube by providing the target sequence, components of a CRISPR complex, including the dead guide sequence to be tested and a control guide sequence different from the test dead guide sequence, and comparing binding or rate of cleavage at the target sequence between the test and control guide sequence reactions. Other assays are possible, and will occur to those skilled in the art. A dead guide sequence may be selected to target any target sequence. In some embodiments, the target sequence is a sequence within a genome of a cell.

[0466] As explained further herein, several structural parameters allow for a proper framework to arrive at such dead guides. Dead guide sequences are typically shorter than respective guide sequences which result in active RNA

cleavage. In particular embodiments, dead guides are 5%, 10%, 20%, 30%, 40%, 50%, shorter than respective guides directed to the same.

[0467] As explained below and known in the art, one aspect of gRNA—RNA targeting specificity is the direct repeat sequence, which is to be appropriately linked to such guides. In particular, this implies that the direct repeat sequences are designed dependent on the origin of the RNA targeting enzyme. Thus, structural data available for validated dead guide sequences may be used for designing Cas13b specific equivalents. Structural similarity between, e.g., the orthologous nuclease domains HEPN of two or more Cas13b effector proteins may be used to transfer design equivalent dead guides. Thus, the dead guide herein may be appropriately modified in length and sequence to reflect such Cas13b specific equivalents, allowing for formation of the CRISPR complex and successful binding to the target RNA, while at the same time, not allowing for successful nuclease activity.

[0468] The use of dead guides in the context herein as well as the state of the art provides a surprising and unexpected platform for network biology and/or systems biology in both *in vitro*, *ex vivo*, and *in vivo* applications, allowing for multiplex gene targeting, and in particular bidirectional multiplex gene targeting. Prior to the use of dead guides, addressing multiple targets has been challenging and in some cases not possible. With the use of dead guides, multiple targets, and thus multiple activities, may be addressed, for example, in the same cell, in the same animal, or in the same patient. Such multiplexing may occur at the same time or staggered for a desired timeframe.

[0469] For example, the dead guides allow to use gRNA as a means for gene targeting, without the consequence of nuclease activity, while at the same time providing directed means for activation or repression. Guide RNA comprising a dead guide may be modified to further include elements in a manner which allow for activation or repression of gene activity, in particular protein adaptors (e.g. aptamers) as described herein elsewhere allowing for functional placement of gene effectors (e.g. activators or repressors of gene activity). One example is the incorporation of aptamers, as explained herein and in the state of the art. By engineering the gRNA comprising a dead guide to incorporate protein-interacting aptamers (Konermann et al., “Genome-scale transcription activation by an engineered CRISPR-Cas9 complex,” doi:10.1038/nature14136, incorporated herein by reference), one may assemble multiple distinct effector domains. Such may be modeled after natural processes.

[0470] Thus, one aspect is a gRNA of the invention which comprises a dead guide, wherein the gRNA further comprises modifications which provide for gene activation or repression, as described herein. The dead gRNA may comprise one or more aptamers. The aptamers may be specific to gene effectors, gene activators or gene repressors. Alternatively, the aptamers may be specific to a protein which in turn is specific to and recruits/binds a specific gene effector, gene activator or gene repressor. If there are multiple sites for activator or repressor recruitment, it is preferred that the sites are specific to either activators or repressors. If there are multiple sites for activator or repressor binding, the sites may be specific to the same activators or same repressors. The sites may also be specific to different activators or different repressors. The effectors, activators, repressors may be present in the form of fusion proteins.

[0471] In an aspect, the invention provides a method of selecting a dead guide RNA targeting sequence for directing a functionalized CRISPR system to a gene locus in an organism, which comprises: a) locating one or more CRISPR motifs in the gene locus; b) analyzing the 20 nt sequence downstream of each CRISPR motif by: i) determining the GC content of the sequence; and ii) determining whether there are off-target matches of the first 15 nt of the sequence in the genome of the organism; c) selecting the sequence for use in a guide RNA if the GC content of the sequence is 70% or less and no off-target matches are identified. In an embodiment, the sequence is selected if the GC content is 50% or less. In an embodiment, the sequence is selected if the GC content is 40% or less. In an embodiment, the sequence is selected if the GC content is 30% or less. In an embodiment, two or more sequences are analyzed and the sequence having the lowest GC content is selected. In an embodiment, off-target matches are determined in regulatory sequences of the organism. In an embodiment, the gene locus is a regulatory region. An aspect provides a dead guide RNA comprising the targeting sequence selected according to the aforementioned methods.

[0472] In an aspect, the invention provides a dead guide RNA for targeting a functionalized CRISPR system to a gene locus in an organism. In an embodiment of the invention, the dead guide RNA comprises a targeting sequence wherein the CG content of the target sequence is 70% or less, and the first 15 nt of the targeting sequence does not match an off-target sequence downstream from a CRISPR motif in the regulatory sequence of another gene locus in the organism. In certain embodiments, the GC content of the targeting sequence 60% or less, 55% or less, 50% or less, 45% or less, 40% or less, 35% or less or 30% or less. In certain embodiments, the GC content of the targeting sequence is from 70% to 60% or from 60% to 50% or from 50% to 40% or from 40% to 30%. In an embodiment, the targeting sequence has the lowest CG content among potential targeting sequences of the locus.

[0473] In an embodiment of the invention, the first 15 nt of the dead guide match the target sequence. In another embodiment, first 14 nt of the dead guide match the target sequence. In another embodiment, the first 13 nt of the dead guide match the target sequence. In another embodiment first 12 nt of the dead guide match the target sequence. In another embodiment, first 11 nt of the dead guide match the target sequence. In another embodiment, the first 10 nt of the dead guide match the target sequence. In an embodiment of the invention the first 15 nt of the dead guide does not match an off-target sequence downstream from a CRISPR motif in the regulatory region of another gene locus. In other embodiments, the first 14 nt, or the first 13 nt of the dead guide, or the first 12 nt of the guide, or the first 11 nt of the dead guide, or the first 10 nt of the dead guide, does not match an off-target sequence downstream from a CRISPR motif in the regulatory region of another gene locus. In other embodiments, the first 15 nt, or 14 nt, or 13 nt, or 12 nt, or 11 nt of the dead guide do not match an off-target sequence downstream from a CRISPR motif in the genome.

[0474] In certain embodiments, the dead guide RNA includes additional nucleotides at the 3'-end that do not match the target sequence. Thus, a dead guide RNA that includes the first 20-28 nt, downstream of a CRISPR motif can be extended in length at the 3' end.

General Provisions

[0475] In an aspect, the invention provides a nucleic acid binding system. In situ hybridization of RNA with complementary probes is a powerful technique. Typically fluorescent DNA oligonucleotides are used to detect nucleic acids by hybridization. Increased efficiency has been attained by certain modifications, such as locked nucleic acids (LNAs), but there remains a need for efficient and versatile alternatives. The invention provides an efficient and adaptable system for in situ hybridization.

[0476] In embodiments of the invention the terms guide sequence and guide RNA are used interchangeably as in foregoing cited documents such as WO 2014/093622 (PCT/US2013/074667). In general, a guide sequence is any polynucleotide sequence having sufficient complementarity with a target polynucleotide sequence to hybridize with the target sequence and direct sequence-specific binding of a CRISPR complex to the target sequence. In some embodiments, the degree of complementarity between a guide sequence and its corresponding target sequence, when optimally aligned using a suitable alignment algorithm, is about or more than about 50%, 60%, 75%, 80%, 85%, 90%, 95%, 97.5%, 99%, or more. Optimal alignment may be determined with the use of any suitable algorithm for aligning sequences, non-limiting example of which include the Smith-Waterman algorithm, the Needleman-Wunsch algorithm, algorithms based on the Burrows-Wheeler Transform (e.g., the Burrows Wheeler Aligner), ClustalW, Clustal X, BLAT, Novoalign (Novocraft Technologies; available at www.novocraft.com), ELAND (Illumina, San Diego, Calif.), SOAP (available at soap.genomics.org.cn), and Maq (available at maq.sourceforge.net). In some embodiments, a guide sequence is about or more than about 5, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30, 35, 40, 45, 50, 75, or more nucleotides in length. In some embodiments, a guide sequence is less than about 75, 50, 45, 40, 35, 30, 25, 20, 15, 12, or fewer nucleotides in length. Preferably the guide sequence is 10-30 nucleotides long, such as 30 nucleotides long. The ability of a guide sequence to direct sequence-specific binding of a CRISPR complex to a target sequence may be assessed by any suitable assay. For example, the components of a CRISPR system sufficient to form a CRISPR complex, including the guide sequence to be tested, may be provided to a host cell having the corresponding target sequence, such as by transfection with vectors encoding the components of the CRISPR sequence, followed by an assessment of preferential cleavage within the target sequence, such as by Surveyor assay as described herein. Similarly, cleavage of a target polynucleotide sequence may be evaluated in a test tube by providing the target sequence, components of a CRISPR complex, including the guide sequence to be tested and a control guide sequence different from the test guide sequence, and comparing binding or rate of cleavage at the target sequence between the test and control guide sequence reactions. Other assays are possible, and will occur to those skilled in the art. A guide sequence may be selected to target any target sequence. In some embodiments, the target sequence is a sequence within a genome of a cell. Exemplary target sequences include those that are unique in the target genome.

[0477] As used herein, "homology" refers to the degree of similarity between sequences of amino acids or nucleic acids. Sequence similarity searches can identify "homolo-

gous” proteins or genes by detecting statistically significant similarity that reflects common ancestry. BLAST, FASTA, SSEARCH, and other commonly used similarity searching programs produce accurate statistical estimates that can be used to reliably infer homology. A common rule of thumb is that two sequences are homologous if they are more than 30% identical over their entire lengths (much higher identities are seen by chance in short alignments).

[0478] By “homologous sequence” is meant, a nucleotide sequence or amino acid sequence that is shared by one or more polynucleotide or amino acid sequences, such as genes, gene transcripts and/or non-coding polynucleotides or proteins. Homologous sequences can also include conserved sequence regions shared by more than one polynucleotide or amino acid sequence. The invention provides for homologous sequences having less than 90% homology (e.g., 89%, 88%, 87%, 86%, 85%, 84%, 83%, 82%, 81%, 80%, 70%, 60%, 50%, 40%, 30%, 20% or less etc.).

[0479] Sequence identity may be determined by sequence comparison and alignment algorithms known in the art. To determine the percent identity of two nucleic acid sequences (or of two amino acid sequences), the sequences are aligned for optimal comparison purposes (e.g., gaps can be introduced in the first sequence or second sequence for optimal alignment). The nucleotides (or amino acid residues) at corresponding nucleotide (or amino acid) positions are then compared. When a position in the first sequence is occupied by the same residue as the corresponding position in the second sequence, then the molecules are identical at that position. The percent identity between the two sequences is a function of the number of identical positions shared by the sequences (i.e., % homology=# of identical positions/total # of positions×100), optionally penalizing the score for the number of gaps introduced and/or length of gaps introduced.

[0480] The comparison of sequences and determination of percent identity between two sequences can be accomplished using a mathematical algorithm. In one embodiment, the alignment generated over a certain portion of the sequence aligned having sufficient identity but not over portions having low degree of identity (i.e., a local alignment). A preferred, non-limiting example of a local alignment algorithm utilized for the comparison of sequences is the algorithm of Karlin and Altschul (1990) Proc. Natl. Acad. Sci. USA 87:2264-68, modified as in Karlin and Altschul (1993) Proc. Natl. Acad. Sci. USA 90:5873-77. Such an algorithm is incorporated into the BLAST programs (version 2.0) of Altschul, et al. (1990) J. Mol. Biol. 215: 403-10.

[0481] The alignment is optimized by introducing appropriate gaps and percent identity is determined over the length of the aligned sequences (i.e., a gapped alignment). To obtain gapped alignments for comparison purposes, Gapped BLAST can be utilized as described in Altschul et al., (1997) Nucleic Acids Res. 25(17):3389-3402. The alignment can also be optimized by introducing appropriate gaps and percent identity is determined over the entire length of the sequences aligned (i.e., a global alignment). A preferred, non-limiting example of a mathematical algorithm utilized for the global comparison of sequences is the algorithm of Myers and Miller, CABIOS (1989). Such an algorithm is incorporated into the ALIGN program (version 2.0) which is part of the GCG sequence alignment software package. When utilizing the ALIGN program for comparing amino

acid sequences, a PAM120 weight residue table, a gap length penalty of 12, and a gap penalty of 4 can be used.

[0482] By “conserved sequence region” is meant, a nucleotide or amino acid sequence of one or more regions in a polynucleotide or protein that does not vary significantly between generations or from one biological system, subject, or organism to another biological system, subject, or organism.

[0483] In general, and throughout this specification, the term “vector” refers to a nucleic acid molecule capable of transporting another nucleic acid to which it has been linked. Vectors include, but are not limited to, nucleic acid molecules that are single-stranded, double-stranded, or partially double-stranded; nucleic acid molecules that comprise one or more free ends, no free ends (e.g., circular); nucleic acid molecules that comprise DNA, RNA, or both; and other varieties of polynucleotides known in the art. One type of vector is a “plasmid,” which refers to a circular double stranded DNA loop into which additional DNA segments can be inserted, such as by standard molecular cloning techniques. Another type of vector is a viral vector, wherein virally-derived DNA or RNA sequences are present in the vector for packaging into a virus (e.g., retroviruses, replication defective retroviruses, adenoviruses, replication defective adenoviruses, and adeno-associated viruses). Viral vectors also include polynucleotides carried by a virus for transfection into a host cell. Certain vectors are capable of autonomous replication in a host cell into which they are introduced (e.g., bacterial vectors having a bacterial origin of replication and episomal mammalian vectors). Other vectors (e.g., non-episomal mammalian vectors) are integrated into the genome of a host cell upon introduction into the host cell, and thereby are replicated along with the host genome. Moreover, certain vectors are capable of directing the expression of genes to which they are operatively-linked. Such vectors are referred to herein as “expression vectors.” Vectors for and that result in expression in a eukaryotic cell can be referred to herein as “eukaryotic expression vectors.” Common expression vectors of utility in recombinant DNA techniques are often in the form of plasmids.

[0484] Recombinant expression vectors can comprise a nucleic acid of the invention in a form suitable for expression of the nucleic acid in a host cell, which means that the recombinant expression vectors include one or more regulatory elements, which may be selected on the basis of the host cells to be used for expression, that is operatively-linked to the nucleic acid sequence to be expressed. Within a recombinant expression vector, “operably linked” is intended to mean that the nucleotide sequence of interest is linked to the regulatory element(s) in a manner that allows for expression of the nucleotide sequence (e.g., in an in vitro transcription/translation system or in a host cell when the vector is introduced into the host cell).

[0485] The term “regulatory element” is intended to include promoters, enhancers, internal ribosomal entry sites (IRES), and other expression control elements (e.g., transcription termination signals, such as polyadenylation signals and poly-U sequences). Such regulatory elements are described, for example, in Goeddel, GENE EXPRESSION TECHNOLOGY: METHODS IN ENZYMOLOGY 185, Academic Press, San Diego, Calif. (1990). Regulatory elements include those that direct constitutive expression of a nucleotide sequence in many types of host cell and those that direct expression of the nucleotide sequence only in certain

host cells (e.g., tissue-specific regulatory sequences). A tissue-specific promoter may direct expression primarily in a desired tissue of interest, such as muscle, neuron, bone, skin, blood, specific organs (e.g., liver, pancreas), or particular cell types (e.g., lymphocytes). Regulatory elements may also direct expression in a temporal-dependent manner, such as in a cell-cycle dependent or developmental stage-dependent manner, which may or may not also be tissue or cell-type specific. In some embodiments, a vector comprises one or more pol III promoter (e.g., 1, 2, 3, 4, 5, or more pol III promoters), one or more pol II promoters (e.g., 1, 2, 3, 4, 5, or more pol II promoters), one or more pol I promoters (e.g., 1, 2, 3, 4, 5, or more pol I promoters), or combinations thereof. Examples of pol III promoters include, but are not limited to, U6 and H1 promoters. Examples of pol II promoters include, but are not limited to, the retroviral Rous sarcoma virus (RSV) LTR promoter (optionally with the RSV enhancer), the cytomegalovirus (CMV) promoter (optionally with the CMV enhancer) [see, e.g., Boshart et al, *Cell*, 41:521-530 (1985)], the SV40 promoter, the dihydrofolate reductase promoter, the β -actin promoter, the phosphoglycerol kinase (PGK) promoter, and the EF1 α promoter. Also encompassed by the term "regulatory element" are enhancer elements, such as WPRE; CMV enhancers; the R-U5' segment in LTR of HTLV-I (*Mol. Cell. Biol.*, Vol. 8(1), p. 466-472, 1988); SV40 enhancer; and the intron sequence between exons 2 and 3 of rabbit β -globin (*Proc. Natl. Acad. Sci. USA.*, Vol. 78(3), p. 1527-31, 1981). It will be appreciated by those skilled in the art that the design of the expression vector can depend on such factors as the choice of the host cell to be transformed, the level of expression desired, etc. A vector can be introduced into host cells to thereby produce transcripts, proteins, or peptides, including fusion proteins or peptides, encoded by nucleic acids as described herein (e.g., clustered regularly interspersed short palindromic repeats (CRISPR) transcripts, proteins, enzymes, mutant forms thereof, fusion proteins thereof, etc.).

[0486] Advantageous vectors include lentiviruses and adeno-associated viruses, and types of such vectors can also be selected for targeting particular types of cells.

[0487] As used herein, the term "crRNA" or "guide RNA" or "single guide RNA" or "sgRNA" or "one or more nucleic acid components" of a Group 29 or Group 30 effector protein loci comprises any polynucleotide sequence having sufficient complementarity with a target nucleic acid sequence to hybridize with the target nucleic acid sequence and direct sequence-specific binding of a nucleic acid-targeting complex to the target nucleic acid sequence. In some embodiments, the degree of complementarity, when optimally aligned using a suitable alignment algorithm, is about or more than about 50%, 60%, 75%, 80%, 85%, 90%, 95%, 97.5%, 99%, or more. Optimal alignment may be determined with the use of any suitable algorithm for aligning sequences, non-limiting example of which include the Smith-Waterman algorithm, the Needleman-Wunsch algorithm, algorithms based on the Burrows-Wheeler Transform (e.g., the Burrows Wheeler Aligner), ClustalW, Clustal X, BLAT, Novoalign (Novocraft Technologies; available at www.novocraft.com), ELAND (Illumina, San Diego, Calif.), SOAP (available at soap.genomics.org.cn), and Maq (available at maq.sourceforge.net). The ability of a guide sequence (within a nucleic acid-targeting guide RNA) to direct sequence-specific binding of a nucleic acid-targeting

complex to a target nucleic acid sequence may be assessed by any suitable assay. For example, the components of a nucleic acid-targeting CRISPR system sufficient to form a nucleic acid-targeting complex, including the guide sequence to be tested, may be provided to a host cell having the corresponding target nucleic acid sequence, such as by transfection with vectors encoding the components of the nucleic acid-targeting complex, followed by an assessment of preferential targeting (e.g., cleavage) within the target nucleic acid sequence, such as by Surveyor assay as described herein. Similarly, cleavage of a target nucleic acid sequence may be evaluated in a test tube by providing the target nucleic acid sequence, components of a nucleic acid-targeting complex, including the guide sequence to be tested and a control guide sequence different from the test guide sequence, and comparing binding or rate of cleavage at the target sequence between the test and control guide sequence reactions. Other assays are possible, and will occur to those skilled in the art. A guide sequence, and hence a nucleic acid-targeting guide RNA may be selected to target any target nucleic acid sequence. The target sequence may be DNA. The target sequence may be any RNA sequence. In some embodiments, the target sequence may be a sequence within a RNA molecule selected from the group consisting of messenger RNA (mRNA), pre-mRNA, ribosomal RNA (rRNA), transfer RNA (tRNA), micro-RNA (miRNA), small interfering RNA (siRNA), small nuclear RNA (snRNA), small nucleolar RNA (snoRNA), double stranded RNA (dsRNA), non coding RNA (ncRNA), long non-coding RNA (lncRNA), and small cytoplasmic RNA (scrRNA). In some preferred embodiments, the target sequence may be a sequence within a RNA molecule selected from the group consisting of mRNA, pre-mRNA, and rRNA. In some preferred embodiments, the target sequence may be a sequence within a RNA molecule selected from the group consisting of ncRNA, and lncRNA. In some more preferred embodiments, the target sequence may be a sequence within an mRNA molecule or a pre-mRNA molecule.

[0488] In some embodiments, a nucleic acid-targeting guide RNA is selected to reduce the degree secondary structure within the RNA-targeting guide RNA. In some embodiments, about or less than about 75%, 50%, 40%, 30%, 25%, 20%, 15%, 10%, 5%, 1%, or fewer of the nucleotides of the nucleic acid-targeting guide RNA participate in self-complementary base pairing when optimally folded. Optimal folding may be determined by any suitable polynucleotide folding algorithm. Some programs are based on calculating the minimal Gibbs free energy. An example of one such algorithm is mFold, as described by Zuker and Stiegler (*Nucleic Acids Res.* 9 (1981), 133-148). Another example folding algorithm is the online webserver RNAfold, developed at Institute for Theoretical Chemistry at the University of Vienna, using the centroid structure prediction algorithm (see e.g., A. R. Gruber et al., 2008, *Cell* 106(1): 23-24; and PA Carr and GM Church, 2009, *Nature Biotechnology* 27(12): 1151-62).

[0489] In certain embodiments, a guide RNA or crRNA may comprise, consist essentially of, or consist of a direct repeat (DR) sequence and a guide sequence or spacer sequence. In certain embodiments, the guide RNA or crRNA may comprise, consist essentially of, or consist of a direct repeat sequence fused or linked to a guide sequence or spacer sequence. In certain embodiments, the direct repeat sequence may be located upstream (i.e., 5') from the guide

sequence or spacer sequence. In other embodiments, the direct repeat sequence may be located downstream (i.e., 3') from the guide sequence or spacer sequence. In other embodiments, multiple DRs (such as dual DRs) may be present.

[0490] In certain embodiments, the crRNA comprises a stem loop, preferably a single stem loop. In certain embodiments, the direct repeat sequence forms a stem loop, preferably a single stem loop.

[0491] In certain embodiments, the spacer length of the guide RNA is from 15 to 35 nt. In certain embodiments, the spacer length of the guide RNA is at least 15 nucleotides. In certain embodiments, the spacer length is from 15 to 17 nt, e.g., 15, 16, or 17 nt, from 17 to 20 nt, e.g., 17, 18, 19, or 20 nt, from 20 to 24 nt, e.g., 20, 21, 22, 23, or 24 nt, from 23 to 25 nt, e.g., 23, 24, or 25 nt, from 24 to 27 nt, e.g., 24, 25, 26, or 27 nt, from 27-30 nt, e.g., 27, 28, 29, or 30 nt, from 30-35 nt, e.g., 30, 31, 32, 33, 34, or 35 nt, or 35 nt or longer.

[0492] The “tracrRNA” sequence or analogous terms includes any polynucleotide sequence that has sufficient complementarity with a crRNA sequence to hybridize. In general, degree of complementarity is with reference to the optimal alignment of the sca sequence and tracr sequence, along the length of the shorter of the two sequences. Optimal alignment may be determined by any suitable alignment algorithm, and may further account for secondary structures, such as self-complementarity within either the sca sequence or tracr sequence. In some embodiments, the degree of complementarity between the tracr sequence and sca sequence along the length of the shorter of the two when optimally aligned is about or more than about 25%, 30%, 40%, 50%, 60%, 70%, 80%, 90%, 95%, 97.5%, 99%, or higher. In certain embodiments, the tracrRNA may not be required. Indeed, as demonstrated herein, the group 29 effector protein from *Bergeyella zoohelcum* and orthologs thereof do not require a tracrRNA to ensure cleavage of an RNA target.

[0493] In further detail, the assay is as follows for a RNA target, provided that a PAM sequence is required to direct recognition. Two *E. coli* strains are used in this assay. One carries a plasmid that encodes the endogenous effector protein locus from the bacterial strain. The other strain carries an empty plasmid (e.g. pACYC184, control strain). All possible 7 or 8 bp PAM sequences are presented on an antibiotic resistance plasmid (pUC19 with ampicillin resistance gene). The PAM is located next to the sequence of proto-spacer 1 (the RNA target to the first spacer in the endogenous effector protein locus). Two PAM libraries were cloned. One has a 8 random bp 5' of the proto-spacer (e.g. total of 65536 different PAM sequences=complexity). The other library has 7 random bp 3' of the proto-spacer (e.g. total complexity is 16384 different PAMs). Both libraries were cloned to have in average 500 plasmids per possible PAM. Test strain and control strain were transformed with 5'PAM and 3'PAM library in separate transformations and transformed cells were plated separately on ampicillin plates. Recognition and subsequent cutting/interference with the plasmid renders a cell vulnerable to ampicillin and prevents growth. Approximately 12 h after transformation, all colonies formed by the test and control strains were harvested and plasmid RNA was isolated. Plasmid RNA was used as template for PCR amplification and subsequent deep sequencing. Representation of all PAMs in the untransformed libraries showed the expected representation of

PAMs in transformed cells. Representation of all PAMs found in control strains showed the actual representation. Representation of all PAMs in test strain showed which PAMs are not recognized by the enzyme and comparison to the control strain allows extracting the sequence of the depleted PAM. In particular embodiments, the cleavage, such as the RNA cleavage is not PAM dependent. Indeed, as demonstrated herein, for the *Bergeyella zoohelcum* effector protein and its orthologs, RNA target cleavage appears to be PAM independent.

[0494] For minimization of toxicity and off-target effect, it will be important to control the concentration of nucleic acid-targeting guide RNA delivered. Optimal concentrations of nucleic acid—targeting guide RNA can be determined by testing different concentrations in a cellular or non-human eukaryote animal model and using deep sequencing the analyze the extent of modification at potential off-target genomic loci. The concentration that gives the highest level of on-target modification while minimizing the level of off-target modification should be chosen for in vivo delivery. The nucleic acid-targeting system is derived advantageously from a Type VI-B CRISPR system. In some embodiments, one or more elements of a nucleic acid-targeting system is derived from a particular organism comprising an endogenous RNA-targeting system. In particular embodiments, the Cas protein CRISPR system is a Group 29 or Group 30 effector protein system.

[0495] The terms “orthologue” (also referred to as “ortholog” herein) and “homologue” (also referred to as “homolog” herein) are well known in the art. By means of further guidance, a “homologue” of a protein as used herein is a protein of the same species which performs the same or a similar function as the protein it is a homologue of. Homologous proteins may but need not be structurally related, or are only partially structurally related. An “orthologue” of a protein as used herein is a protein of a different species which performs the same or a similar function as the protein it is an orthologue of. Orthologous proteins may but need not be structurally related, or are only partially structurally related. In particular embodiments, the homologue or orthologue of a Group 29 or Group 30 protein as referred to herein has a sequence homology or identity of at least 80%, more preferably at least 85%, even more preferably at least 90%, such as for instance at least 95% with said Group 29 or Group 30 effector protein. In a preferred embodiment, the Group 29 or Group 30 effector protein may be an ortholog of an organism of a genus which includes but is not limited to *Bergeyella*, *Prevotella*, *Porphyromonas*, *Bacteroides*, *Alistipes*, *Riemerella*, *Myroides*, *Flavobacterium*, *Capnocytophaga*, *Chryseobacterium*, *Phaeodactylibacter*, *Paludibacter* or *Psychroflexus*. Some methods of identifying orthologs of CRISPR system enzymes may involve identifying tracr sequences in genomes of interest. Identification of tracr sequences may relate to the following steps: Search for the direct repeats or tracr mate sequences in a database to identify a CRISPR region comprising a CRISPR enzyme. Search for homologous sequences in the CRISPR region flanking the CRISPR enzyme in both the sense and antisense directions. Look for transcriptional terminators and secondary structures. Identify any sequence that is not a direct repeat or a tracr mate sequence but has more than 50% identity to the direct repeat or tracr mate sequence as a

potential tracr sequence. Take the potential tracr sequence and analyze for transcriptional terminator sequences associated therewith.

[0496] It will be appreciated that any of the functionalities described herein may be engineered into CRISPR enzymes from other orthologs, including chimeric enzymes comprising fragments from multiple orthologs. Examples of such orthologs are described elsewhere herein. Thus, chimeric enzymes may comprise fragments of CRISPR enzyme orthologs of an organism which includes but is not limited to *Bergeyella*, *Prevotella*, *Porphyromonas*, *Bacteroides*, *Alistipes*, *Riemerella*, *Myroides*, *Flavobacterium*, *Capnocytophaga*, *Chryseobacterium*, *Phaeodactylibacter*, *Paludibacter* or *Psychroflexus*. A chimeric enzyme can comprise a first fragment and a second fragment, and the fragments can be of CRISPR enzyme orthologs of organisms of genres herein mentioned or of species herein mentioned; advantageously the fragments are from CRISPR enzyme orthologs of different species.

[0497] In embodiments, the Type VI-B RNA-targeting effector proteins referred to herein also encompasses a functional variant of the effector protein or a homologue or an orthologue thereof. A “functional variant” of a protein as used herein refers to a variant of such protein which retains at least partially the activity of that protein. Functional variants may include mutants (which may be insertion, deletion, or replacement mutants), including polymorphs, etc. Also included within functional variants are fusion products of such protein with another, usually unrelated, nucleic acid, protein, polypeptide or peptide. Functional variants may be naturally occurring or may be man-made. In an embodiment, nucleic acid molecule(s) encoding the Group 29 or Group 30 RNA-targeting effector proteins, or an ortholog or homolog thereof, may be codon-optimized for expression in an eukaryotic cell. A eukaryote can be as herein discussed. Nucleic acid molecule(s) can be engineered or non-naturally occurring.

[0498] In an embodiment, the Type VI-B RNA-targeting (e.g. Group 29 or Group 30) effector protein or an ortholog or homolog thereof, may comprise one or more mutations. The mutations may be artificially introduced mutations and may include but are not limited to one or more mutations in a catalytic domain. Examples of catalytic domains with reference to an effector enzyme may include but are not limited to RuvC I, RuvC II, RuvC III, HNH domains, and HEPN domains. In particular embodiments, the one or more mutations are introduced into one or more of the HEPN domains.

[0499] In an embodiment, the Group 29 or Group 30 protein or an ortholog or homolog thereof, may be used as a generic nucleic acid binding protein with fusion to or being operably linked to a functional domain. Exemplary functional domains may include but are not limited to translational initiator, translational activator, translational repressor, nucleases, in particular ribonucleases, a spliceosome, beads, a light inducible/controllable domain or a chemically inducible/controllable domain.

[0500] In some embodiments, the unmodified nucleic acid-targeting effector protein may have cleavage activity. In some embodiments, the RNA-targeting effector protein may direct cleavage of one or both nucleic acid (DNA or RNA) strands at the location of or near a target sequence, such as within the target sequence and/or within the complement of the target sequence or at sequences associated with the target

sequence. In some embodiments, the nucleic acid-targeting CRISPR protein may direct cleavage of one or both DNA or RNA strands within about 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 15, 20, 25, 50, 100, 200, 500, or more base pairs from the first or last nucleotide of a target sequence. In some embodiments, the nucleic acid-targeting CRISPR protein may direct more than one cleavage (such as one, two three, four, five, or more cleavages) of one or both DNA or RNA strands within the target sequence and/or within the complement of the target sequence or at sequences associated with the target sequence and/or within about 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 15, 20, 25, 50, 100, 200, 500, or more base pairs from the first or last nucleotide of a target sequence. In some embodiments, the cleavage may be blunt, i.e., generating blunt ends. In some embodiments, the cleavage may be staggered, i.e., generating sticky ends. In some embodiments, the cleavage may be a staggered cut with a 5' overhang, e.g., a 5' overhang of 1 to 5 nucleotides. In some embodiments, the cleavage may be a staggered cut with a 3' overhang, e.g., a 3' overhang of 1 to 5 nucleotides. In some embodiments, a vector encodes a nucleic acid-targeting Cas protein that may be mutated with respect to a corresponding wild-type enzyme such that the mutated nucleic acid-targeting Cas protein lacks the ability to cleave one or both DNA or RNA strands of a target polynucleotide containing a target sequence. As a further example, two or more catalytic domains of Cas (RuvC I, RuvC II, and RuvC III or the HNH domain, or HEPN domain) may be mutated to produce a mutated Cas substantially lacking all RNA cleavage activity. As described herein, corresponding catalytic domains of the effector protein may also be mutated to produce a mutated Group 29 or Group 30 effector protein lacking all nucleic acid cleavage activity or having substantially reduced nucleic acid cleavage activity. In some embodiments, a nucleic acid-targeting effector protein may be considered to substantially lack all RNA cleavage activity when the RNA cleavage activity of the mutated enzyme is about no more than 25%, 10%, 5%, 1%, 0.1%, 0.01%, or less of the nucleic acid cleavage activity of the non-mutated form of the enzyme; an example can be when the nucleic acid cleavage activity of the mutated form is nil or negligible as compared with the non-mutated form. An effector protein may be identified with reference to the general class of enzymes that share homology to the biggest nuclease with multiple nuclease domains from previously described enzyme systems. By derived, Applicants mean that the derived enzyme is largely based, in the sense of having a high degree of sequence homology with, a wildtype enzyme, but that it has been mutated (modified) in some way as known in the art or as described herein.

[0501] Again, it will be appreciated that the terms Cas and CRISPR enzyme and CRISPR protein and Cas protein are generally used interchangeably and at all points of reference herein refer by analogy to novel CRISPR effector proteins further described in this application, unless otherwise apparent, such as by specific reference to Cas9. As mentioned above, many of the residue numberings used herein refer to the effector protein from Group 29 or Group 30 CRISPR locus. However, it will be appreciated that this invention includes many more effector proteins from other species of microbes. In certain embodiments, Cas may be constitutively present or inducibly present or conditionally present or administered or delivered. Cas optimization may be used to enhance function or to develop new functions, one can

generate chimeric Cas proteins. And Cas may be used as a generic nucleic acid binding protein.

[0502] Typically, in the context of an endogenous nucleic acid-targeting system, formation of a nucleic acid-targeting complex (comprising a guide RNA hybridized to a target sequence and complexed with one or more nucleic acid-targeting effector proteins) results in cleavage of one or both DNA or RNA strands in or near (e.g., within 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 20, 50, or more base pairs from) the target sequence. As used herein the term “sequence(s) associated with a target locus of interest” refers to sequences near the vicinity of the target sequence (e.g. within 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 20, 50, or more base pairs from the target sequence, wherein the target sequence is comprised within a target locus of interest).

[0503] An example of a codon optimized sequence, is in this instance a sequence optimized for expression in a eukaryote, e.g., humans (i.e. being optimized for expression in humans), or for another eukaryote, animal or mammal as herein discussed; see, e.g., SaCas9 human codon optimized sequence in WO 2014/093622 (PCT/US2013/074667) as an example of a codon optimized sequence (from knowledge in the art and this disclosure, codon optimizing coding nucleic acid molecule(s), especially as to effector protein (e.g., Group 29 or Group 30) is within the ambit of the skilled artisan). Whilst this is preferred, it will be appreciated that other examples are possible and codon optimization for a host species other than human, or for codon optimization for specific organs is known. In some embodiments, an enzyme coding sequence encoding a DNA/RNA-targeting Cas protein is codon optimized for expression in particular cells, such as eukaryotic cells. The eukaryotic cells may be those of or derived from a particular organism, such as a mammal, including but not limited to human, or non-human eukaryote or animal or mammal as herein discussed, e.g., mouse, rat, rabbit, dog, livestock, or non-human mammal or primate. In some embodiments, processes for modifying the germ line genetic identity of human beings and/or processes for modifying the genetic identity of animals which are likely to cause them suffering without any substantial medical benefit to man or animal, and also animals resulting from such processes, may be excluded. In general, codon optimization refers to a process of modifying a nucleic acid sequence for enhanced expression in the host cells of interest by replacing at least one codon (e.g., about or more than about 1, 2, 3, 4, 5, 10, 15, 20, 25, 50, or more codons) of the native sequence with codons that are more frequently or most frequently used in the genes of that host cell while maintaining the native amino acid sequence. Various species exhibit particular bias for certain codons of a particular amino acid. Codon bias (differences in codon usage between organisms) often correlates with the efficiency of translation of messenger RNA (mRNA), which is in turn believed to be dependent on, among other things, the properties of the codons being translated and the availability of particular transfer RNA (tRNA) molecules. The predominance of selected tRNAs in a cell is generally a reflection of the codons used most frequently in peptide synthesis. Accordingly, genes can be tailored for optimal gene expression in a given organism based on codon optimization. Codon usage tables are readily available, for example, at the “Codon Usage Database” available at www.kazusa.or.jp/codon/ and these tables can be adapted in a number of ways. See Nakamura, Y., et al. “Codon usage tabulated from the international DNA

sequence databases: status for the year 2000” *Nucl. Acids Res.* 28:292 (2000). Computer algorithms for codon optimizing a particular sequence for expression in a particular host cell are also available, such as Gene Forge (Aptagen; Jacobus, Pa.), are also available. In some embodiments, one or more codons (e.g., 1, 2, 3, 4, 5, 10, 15, 20, 25, 50, or more, or all codons) in a sequence encoding a DNA/RNA-targeting Cas protein corresponds to the most frequently used codon for a particular amino acid.

[0504] In some embodiments, a vector encodes a nucleic acid-targeting effector protein such as the Group 29 or Group 30 effectors as described herein, or an ortholog or homolog thereof comprising one or more nuclear localization sequences (NLSs), such as about or more than about 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, or more NLSs. In some embodiments, the RNA-targeting effector protein comprises about or more than about 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, or more NLSs at or near the amino-terminus, about or more than about 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, or more NLSs at or near the carboxy-terminus, or a combination of these (e.g., zero or at least one or more NLS at the amino-terminus and zero or at one or more NLS at the carboxy terminus). When more than one NLS is present, each may be selected independently of the others, such that a single NLS may be present in more than one copy and/or in combination with one or more other NLSs present in one or more copies. In some embodiments, an NLS is considered near the N- or C-terminus when the nearest amino acid of the NLS is within about 1, 2, 3, 4, 5, 10, 15, 20, 25, 30, 40, 50, or more amino acids along the polypeptide chain from the N- or C-terminus. Non-limiting examples of NLSs include an NLS sequence derived from: the NLS of the SV40 virus large T-antigen, having the amino acid sequence PKKKRKKV; the NLS from nucleoplasmin (e.g., the nucleoplasmin bipartite NLS with the sequence KRPAATKKAGQAKKKK); the c-myc NLS having the amino acid sequence PAAKRVKLD or RQRRNELKRSP; the hRNPA1 M9 NLS having the sequence NQSSNFGPM-KGGNFGGRSSGPGYGGGGQYFAKPRNQGGY; the sequence RMRIZFKNKGKDTAELRRRRVEVSVEL-RKAKKDEQILKRRNV of the IBB domain from importin-alpha; the sequences VSRKRPRP and PPKKARED of the myoma T protein; the sequence POPKKKPL of human p53; the sequence SALIKKKKKMAP of mouse c-ab1 IV; the sequences DRLRR and PKQKKRK of the influenza virus NS1; the sequence RKLKKKIKKL of the Hepatitis virus delta antigen; the sequence REKKKFLKRR of the mouse Mx1 protein; the sequence KRKGDEVDGVDEVAK-KKSKK of the human poly(ADP-ribose) polymerase; and the sequence RKCLQAGMNLEARKTKK of the steroid hormone receptors (human) glucocorticoid. In general, the one or more NLSs are of sufficient strength to drive accumulation of the DNA/RNA-targeting Cas protein in a detectable amount in the nucleus of a eukaryotic cell. In general, strength of nuclear localization activity may derive from the number of NLSs in the nucleic acid-targeting effector protein, the particular NLS(s) used, or a combination of these factors. Detection of accumulation in the nucleus may be performed by any suitable technique. For example, a detectable marker may be fused to the nucleic acid-targeting protein, such that location within a cell may be visualized, such as in combination with a means for detecting the location of the nucleus (e.g., a stain specific for the nucleus such as DAPI). Cell nuclei may also be isolated from cells, the contents of which may then be analyzed by any suitable

process for detecting protein, such as immunohistochemistry, Western blot, or enzyme activity assay. Accumulation in the nucleus may also be determined indirectly, such as by an assay for the effect of nucleic acid-targeting complex formation (e.g., assay for DNA or RNA cleavage or mutation at the target sequence, or assay for altered gene expression activity affected by DNA or RNA-targeting complex formation and/or DNA or RNA-targeting Cas protein activity), as compared to a control not exposed to the nucleic acid-targeting Cas protein or nucleic acid-targeting complex, or exposed to a nucleic acid-targeting Cas protein lacking the one or more NLSs. In preferred embodiments of the herein described Group 29 or Group 30 effector protein complexes and systems the codon optimized Cpf1 effector proteins comprise an NLS attached to the C-terminal of the protein.

[0505] In some embodiments, one or more vectors driving expression of one or more elements of a nucleic acid-targeting system are introduced into a host cell such that expression of the elements of the nucleic acid-targeting system direct formation of a nucleic acid-targeting complex at one or more target sites. For example, a nucleic acid-targeting effector enzyme and a nucleic acid-targeting guide RNA could each be operably linked to separate regulatory elements on separate vectors. RNA(s) of the nucleic acid-targeting system can be delivered to a transgenic nucleic acid-targeting effector protein animal or mammal, e.g., an animal or mammal that constitutively or inducibly or conditionally expresses nucleic acid-targeting effector protein; or an animal or mammal that is otherwise expressing nucleic acid-targeting effector protein or has cells containing nucleic acid-targeting effector protein, such as by way of prior administration thereto of a vector or vectors that code for and express in vivo nucleic acid-targeting effector protein. Alternatively, two or more of the elements expressed from the same or different regulatory elements, may be combined in a single vector, with one or more additional vectors providing any components of the nucleic acid-targeting system not included in the first vector. nucleic acid-targeting system elements that are combined in a single vector may be arranged in any suitable orientation, such as one element located 5' with respect to ("upstream" of) or 3' with respect to ("downstream" of) a second element. The coding sequence of one element may be located on the same or opposite strand of the coding sequence of a second element, and oriented in the same or opposite direction. In some embodiments, a single promoter drives expression of a transcript encoding a nucleic acid-targeting effector protein and the nucleic acid-targeting guide RNA, embedded within one or more intron sequences (e.g., each in a different intron, two or more in at least one intron, or all in a single intron). In some embodiments, the nucleic acid-targeting effector protein and the nucleic acid-targeting guide RNA may be operably linked to and expressed from the same promoter. Delivery vehicles, vectors, particles, formulations and components thereof for expression of one or more elements of a nucleic acid-targeting system are as used in the foregoing documents, such as WO 2014/093622 (PCT/US2013/074667). In some embodiments, a vector comprises one or more insertion sites, such as a restriction endonuclease recognition sequence (also referred to as a "cloning site"). In some embodiments, one or more insertion sites (e.g., about or more than about 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, or more insertion sites) are located upstream and/or downstream of one or more sequence elements of one or more

vectors. In some embodiments, a vector comprises an insertion site upstream of a tracr mate sequence, and optionally downstream of a regulatory element operably linked to the tracr mate sequence, such that following insertion of a guide sequence into the insertion site and upon expression the guide sequence directs sequence-specific binding of a nucleic acid-targeting complex to a target sequence in a eukaryotic cell. In some embodiments, a vector comprises two or more insertion sites, so as to allow insertion of a guide sequence at each site. In such an arrangement, the two or more guide sequences may comprise two or more copies of a single guide sequence, two or more different guide sequences, or combinations of these. When multiple different guide sequences are used, a single expression construct may be used to target nucleic acid-targeting activity to multiple different, corresponding target sequences within a cell. For example, a single vector may comprise about or more than about 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 15, 20, or more guide sequences. In some embodiments, about or more than about 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, or more such guide-sequence-containing vectors may be provided, and optionally delivered to a cell. In some embodiments, a vector comprises a regulatory element operably linked to an enzyme-coding sequence encoding a nucleic acid-targeting effector protein. Nucleic acid-targeting effector protein or nucleic acid-targeting guide RNA or RNA(s) can be delivered separately; and advantageously at least one of these is delivered via a particle or particle complex. Nucleic acid-targeting effector protein mRNA can be delivered prior to the nucleic acid-targeting guide RNA to give time for nucleic acid-targeting effector protein to be expressed. Nucleic acid-targeting effector protein mRNA might be administered 1-12 hours (preferably around 2-6 hours) prior to the administration of nucleic acid-targeting guide RNA. Alternatively, nucleic acid-targeting effector protein mRNA and nucleic acid-targeting guide RNA can be administered together. Advantageously, a second booster dose of guide RNA can be administered 1-12 hours (preferably around 2-6 hours) after the initial administration of nucleic acid-targeting effector protein mRNA+guide RNA. Additional administrations of nucleic acid-targeting effector protein mRNA and/or guide RNA might be useful to achieve the most efficient levels of genome modification.

[0506] In one aspect, the invention provides methods for using one or more elements of a nucleic acid-targeting system. The nucleic acid-targeting complex of the invention provides an effective means for modifying a target DNA or RNA single or double stranded, linear or super-coiled). The nucleic acid-targeting complex of the invention has a wide variety of utility including modifying (e.g., deleting, inserting, translocating, inactivating, activating) a target DNA or RNA in a multiplicity of cell types. As such the nucleic acid-targeting complex of the invention has a broad spectrum of applications in, e.g., gene therapy, drug screening, disease diagnosis, and prognosis. An exemplary nucleic acid-targeting complex comprises a DNA or RNA-targeting effector protein complexed with a guide RNA hybridized to a target sequence within the target locus of interest.

[0507] In other embodiments, this invention provides a method of modifying expression of a RNA in a eukaryotic cell. The method comprises increasing or decreasing expression of a target polynucleotide by using a nucleic acid-targeting complex that binds to the RNA (e.g., mRNA or pre-mRNA). In some methods, a target RNA can be inac-

tivated to effect the modification of the expression in a cell. For example, upon the binding of a RNA-targeting complex to a target sequence in a cell, the target RNA is inactivated such that the sequence is not translated, the coded protein is not produced, or the sequence does not function as the wild-type sequence does. For example, a protein or microRNA coding sequence may be inactivated such that the protein or microRNA or pre-microRNA transcript is not produced. The target RNA of a RNA-targeting complex can be any RNA endogenous or exogenous to the eukaryotic cell. For example, the target RNA can be a RNA residing in the nucleus of the eukaryotic cell. The target RNA can be a sequence (e.g., mRNA or pre-mRNA) coding a gene product (e.g., a protein) or a non-coding sequence (e.g., ncRNA, lncRNA, tRNA, or rRNA). Examples of target RNA include a sequence associated with a signaling biochemical pathway, e.g., a signaling biochemical pathway-associated RNA. Examples of target RNA include a disease associated RNA. A "disease-associated" RNA refers to any RNA which is yielding translation products at an abnormal level or in an abnormal form in cells derived from a disease-affected tissues compared with tissues or cells of a non disease control. It may be a RNA transcribed from a gene that becomes expressed at an abnormally high level; it may be a RNA transcribed from a gene that becomes expressed at an abnormally low level, where the altered expression correlates with the occurrence and/or progression of the disease. A disease-associated RNA also refers to a RNA transcribed from a gene possessing mutation(s) or genetic variation that is directly responsible or is in linkage disequilibrium with a gene(s) that is responsible for the etiology of a disease. The translated products may be known or unknown, and may be at a normal or abnormal level. The target RNA of a RNA-targeting complex can be any RNA endogenous or exogenous to the eukaryotic cell. For example, the target RNA can be a RNA residing in the nucleus of the eukaryotic cell. The target RNA can be a sequence (e.g., mRNA or pre-mRNA) coding a gene product (e.g., a protein) or a non-coding sequence (e.g., ncRNA, lncRNA, tRNA, or rRNA).

[0508] In some embodiments, the method may comprise allowing a nucleic acid-targeting complex to bind to the target DNA or RNA to effect cleavage of said target DNA or RNA thereby modifying the target DNA or RNA, wherein the nucleic acid-targeting complex comprises a nucleic acid-targeting effector protein complexed with a guide RNA hybridized to a target sequence within said target DNA or RNA. In one aspect, the invention provides a method of modifying expression of DNA or RNA in a eukaryotic cell. In some embodiments, the method comprises allowing a nucleic acid-targeting complex to bind to the DNA or RNA such that said binding results in increased or decreased expression of said DNA or RNA; wherein the nucleic acid-targeting complex comprises a nucleic acid-targeting effector protein complexed with a guide RNA. Similar considerations and conditions apply as above for methods of modifying a target DNA or RNA. In fact, these sampling, culturing and re-introduction options apply across the aspects of the present invention. In one aspect, the invention provides for methods of modifying a target DNA or RNA in a eukaryotic cell, which may be in vivo, ex vivo or in vitro. In some embodiments, the method comprises sampling a cell or population of cells from a human or non-human animal, and modifying the cell or cells. Culturing may occur at any stage ex vivo. The cell or cells may even be re-

introduced into the non-human animal or plant. For re-introduced cells it is particularly preferred that the cells are stem cells.

[0509] Indeed, in any aspect of the invention, the nucleic acid-targeting complex may comprise a nucleic acid-targeting effector protein complexed with a guide RNA hybridized to a target sequence.

[0510] The invention relates to the engineering and optimization of systems, methods and compositions used for the control of gene expression involving DNA or RNA sequence targeting, that relate to the nucleic acid-targeting system and components thereof. An advantage of the present methods is that the CRISPR system minimizes or avoids off-target binding and its resulting side effects. This is achieved using systems arranged to have a high degree of sequence specificity for the target DNA or RNA.

[0511] In relation to a nucleic acid-targeting complex or system preferably, the tracr sequence has one or more hairpins and is 30 or more nucleotides in length, 40 or more nucleotides in length, or 50 or more nucleotides in length; the crRNA sequence is between 10 to 30 nucleotides in length, the nucleic acid-targeting effector protein is a Group 29 or Group 30 effector protein.

[0512] The use of two different aptamers (each associated with a distinct nucleic acid-targeting guide RNAs) allows an activator-adaptor protein fusion and a repressor-adaptor protein fusion to be used, with different nucleic acid-targeting guide RNAs, to activate expression of one DNA or RNA, whilst repressing another. They, along with their different guide RNAs can be administered together, or substantially together, in a multiplexed approach. A large number of such modified nucleic acid-targeting guide RNAs can be used all at the same time, for example 10 or 20 or 30 and so forth, whilst only one (or at least a minimal number) of effector protein molecules need to be delivered, as a comparatively small number of effector protein molecules can be used with a large number modified guides. The adaptor protein may be associated (preferably linked or fused to) one or more activators or one or more repressors. For example, the adaptor protein may be associated with a first activator and a second activator. The first and second activators may be the same, but they are preferably different activators. Three or more or even four or more activators (or repressors) may be used, but package size may limit the number being higher than 5 different functional domains. Linkers are preferably used, over a direct fusion to the adaptor protein, where two or more functional domains are associated with the adaptor protein. Suitable linkers might include the GlySer linker.

[0513] It is also envisaged that the nucleic acid-targeting effector protein-guide RNA complex as a whole may be associated with two or more functional domains. For example, there may be two or more functional domains associated with the nucleic acid-targeting effector protein, or there may be two or more functional domains associated with the guide RNA (via one or more adaptor proteins), or there may be one or more functional domains associated with the nucleic acid-targeting effector protein and one or more functional domains associated with the guide RNA (via one or more adaptor proteins).

[0514] The fusion between the adaptor protein and the activator or repressor may include a linker. For example, GlySer linkers GGGGS can be used. They can be used in repeats of 3 ((GGGGS)₃) or 6, 9 or even 12 or more, to provide suitable lengths, as required. Linkers can be used

between the guide RNAs and the functional domain (activator or repressor), or between the nucleic acid-targeting effector protein and the functional domain (activator or repressor). The linker the user to engineer appropriate amounts of “mechanical flexibility”.

[0515] The invention comprehends a nucleic acid-targeting complex comprising a nucleic acid-targeting effector protein and a guide RNA, wherein the nucleic acid-targeting effector protein comprises at least one mutation, such that the nucleic acid-targeting Cas protein has no more than 5% of the activity of the nucleic acid-targeting Cas protein not having the at least one mutation and, optionally, at least one or more nuclear localization sequences; the guide RNA comprises a guide sequence capable of hybridizing to a target sequence in a RNA of interest in a cell; and wherein: the nucleic acid-targeting effector protein is associated with two or more functional domains; or at least one loop of the guide RNA is modified by the insertion of distinct RNA sequence(s) that bind to one or more adaptor proteins, and wherein the adaptor protein is associated with two or more functional domains; or the nucleic acid-targeting effector protein is associated with one or more functional domains and at least one loop of the guide RNA is modified by the insertion of distinct RNA sequence(s) that bind to one or more adaptor proteins, and wherein the adaptor protein is associated with one or more functional domains.

Group 29 or Group 30 Effector Protein Complexes can Deliver Functional Effectors

[0516] Unlike CRISPR-Cas-mediated gene knockout, which permanently eliminates expression by mutating the gene at the DNA level, CRISPR-Cas knockdown allows for temporary reduction of gene expression through the use of artificial transcription factors. Mutating key residues in one or both DNA cleavage domains of the Group 29 or Group 30 protein results in the generation of a catalytically inactive Group 29 or Group 30 protein. A catalytically inactive Group 29 or Group 30 complexes with a guide RNA and localizes to the DNA sequence specified by that guide RNA's targeting domain, however, it does not cleave the target DNA. Fusion of the inactive Group 29 or Group 30 protein to an effector domain, e.g., a transcription repression domain, enables recruitment of the effector to any DNA site specified by the guide RNA. In certain embodiments, Group 29 or Group 30 may be fused to a transcriptional repression domain and recruited to the promoter region of a gene. Especially for gene repression, it is contemplated herein that blocking the binding site of an endogenous transcription factor would aid in downregulating gene expression. In another embodiment, an inactive Group 29 or Group 30 can be fused to a chromatin modifying protein. Altering chromatin status can result in decreased expression of the target gene.

Optimized Functional RNA Targeting Systems

[0517] In an aspect the invention thus provides a system for specific delivery of functional components to the RNA environment. This can be ensured using the CRISPR systems comprising the RNA targeting effector proteins of the present invention which allow specific targeting of different components to RNA. More particularly such components include activators or repressors, such as activators or repressors of RNA translation, degradation, etc. According to one

aspect the invention provides non-naturally occurring or engineered composition comprising a guide RNA comprising a guide sequence capable of hybridizing to a target sequence in a genomic locus of interest in a cell, wherein the guide RNA is modified by the insertion of one or more distinct RNA sequence(s) that bind an adaptor protein. In particular embodiments, the RNA sequences may bind to two or more adaptor proteins (e.g. aptamers), and wherein each adaptor protein is associated with one or more functional domains. The guide RNAs of at least the group29 enzymes described herein are shown to be amenable to modification of the guide sequence. In particular embodiments, the guide RNA is modified by the insertion of distinct RNA sequence(s) 5' of the direct repeat, within the direct repeat, or 3' of the guide sequence. When there is more than one functional domain, the functional domains can be same or different, e.g., two of the same or two different activators or repressors. In an aspect the invention provides a herein-discussed composition, wherein the one or more functional domains are attached to the RNA targeting enzyme so that upon binding to the target RNA the functional domain is in a spatial orientation allowing for the functional domain to function in its attributed function; In an aspect the invention provides a herein-discussed composition, wherein the composition comprises a CRISPR-Cas complex having at least three functional domains, at least one of which is associated with the RNA targeting enzyme and at least two of which are associated with the gRNA.

[0518] In particular embodiments, the effector protein is a group29 effector protein from *Bergeyella zoohelcum* ATCC 43767 and the functional domain is linked to the guide RNA through a linker inserted after the minimal direct repeat. Indeed, it has been established that at least some of the group29 effector proteins comprise both long (about 87 nt) and short (about 36 nt) direct repeat sequences. Accordingly this implies that a non-functional sequence can be inserted 3' of the short direct repeat sequence while maintaining a functional guide.

[0519] Accordingly, In an aspect the invention provides non-naturally occurring or engineered CRISPR-Cas complex composition comprising the guide RNA as herein-discussed and a CRISPR enzyme which is an RNA targeting enzyme, wherein optionally the RNA targeting enzyme comprises at least one mutation, such that the RNA targeting enzyme has no more than 5% of the nuclease activity of the enzyme not having the at least one mutation, and optionally one or more comprising at least one or more nuclear localization sequences. In particular embodiments, the guide RNA is additionally or alternatively modified so as to still ensure binding of the RNA targeting enzyme but to prevent cleavage by the RNA targeting enzyme (as detailed elsewhere herein).

[0520] In particular embodiments, the RNA targeting enzyme is a group29 enzyme which has a diminished nuclease activity of at least 97%, or 100% as compared with the group29 enzyme not having the at least one mutation. In an aspect the invention provides a herein-discussed composition, wherein the group29 enzyme comprises two or more mutations. The mutations may be selected from R116A, H121A, R1177A, and H1182A, according to *Bergeyella zoohelcum* ATCC 43767 group29 protein or a corresponding position in an ortholog.

[0521] In particular embodiments, an RNA targeting system is provided as described herein above comprising two or

more functional domains. In particular embodiments, the two or more functional domains are heterologous functional domain. In particular embodiments, the system comprises an adaptor protein which is a fusion protein comprising a functional domain, the fusion protein optionally comprising a linker between the adaptor protein and the functional domain. In particular embodiments, the linker includes a GlySer linker. Additionally or alternatively, one or more functional domains are attached to the RNA effector protein by way of a linker, optionally a GlySer linker. In particular embodiments, the one or more functional domains are attached to the RNA targeting enzyme through one or both of the HEPN domains. In particular embodiments, the one or more functional domains are attached to the C-terminal end of the RNA effector protein, such as the C-terminal end of the group29 effector protein from *Bergeyella zoohelcum* ATCC 43767.

[0522] In an aspect the invention provides a herein-discussed composition, wherein the one or more functional domains associated with the adaptor protein or the RNA targeting enzyme is a domain capable of activating or repressing RNA translation. In an aspect the invention provides a herein-discussed composition, wherein at least one of the one or more functional domains associated with the adaptor protein have one or more activities comprising methylase activity, demethylase activity, transcription activation activity, transcription repression activity, transcription release factor activity, histone modification activity, DNA integration activity RNA cleavage activity, DNA cleavage activity or nucleic acid binding activity, or molecular switch activity or chemical inducibility or light inducibility. Envisaged applications are detailed herein below.

[0523] In an aspect the invention provides a herein-discussed composition comprising an aptamer sequence. In particular embodiments, the aptamer sequence is two or more aptamer sequences specific to the same adaptor protein. In an aspect the invention provides a herein-discussed composition, wherein the aptamer sequence is two or more aptamer sequences specific to different adaptor protein. In an aspect the invention provides a herein-discussed composition, wherein the adaptor protein comprises MS2, PP7, Q β , F2, GA, fr, JP501, M12, R17, BZ13, JP34, JP500, KU1, M11, MX1, TW18, VK, SP, FI, ID2, NL95, TW19, AP205, ϕ Cb5, ϕ Cb8r, ϕ Cb12r, ϕ Cb23r, 7s, PRR1. Accordingly, in particular embodiments, the aptamer is selected from a binding protein specifically binding any one of the adaptor proteins listed above. In an aspect the invention provides a herein-discussed composition, wherein the cell is a eukaryotic cell. In an aspect the invention provides a herein-discussed composition, wherein the eukaryotic cell is a mammalian cell, a plant cell or a yeast cell, whereby the mammalian cell is optionally a mouse cell. In an aspect the invention provides a herein-discussed composition, wherein the mammalian cell is a human cell.

[0524] In an aspect the invention provides a herein above-discussed composition wherein there is more than one gRNA, and the gRNAs target different sequences whereby when the composition is employed, there is multiplexing. In an aspect the invention provides a composition wherein there is more than one gRNA modified by the insertion of distinct RNA sequence(s) that bind to one or more adaptor proteins.

[0525] In an aspect the invention provides a herein-discussed composition wherein one or more adaptor proteins

associated with one or more functional domains is present and bound to the distinct RNA sequence(s) inserted into the guide RNA(s).

[0526] In an aspect the invention provides a herein-discussed composition wherein the guide RNA is modified to have at least one non-coding functional loop; e.g., wherein the at least one non-coding functional loop is repressive; for instance, wherein at least one non-coding functional loop comprises Alu.

[0527] In an aspect the invention provides a method for modifying gene expression comprising the administration to a host or expression in a host in vivo of one or more of the compositions as herein-discussed.

[0528] In an aspect the invention provides a herein-discussed method comprising the delivery of the composition or nucleic acid molecule(s) coding therefor, wherein said nucleic acid molecule(s) are operatively linked to regulatory sequence(s) and expressed in vivo. In an aspect the invention provides a herein-discussed method wherein the expression in vivo is via a lentivirus, an adenovirus, or an AAV.

[0529] In an aspect the invention provides a mammalian cell line of cells as herein-discussed, wherein the cell line is, optionally, a human cell line or a mouse cell line. In an aspect the invention provides a transgenic mammalian model, optionally a mouse, wherein the model has been transformed with a herein-discussed composition or is a progeny of said transformant.

[0530] In an aspect the invention provides a nucleic acid molecule(s) encoding guide RNA or the RNA targeting CRISPR-Cas complex or the composition as herein-discussed. In an aspect the invention provides a vector comprising: a nucleic acid molecule encoding a guide RNA (gRNA) comprising a guide sequence capable of hybridizing to a target sequence in a genomic locus of interest in a cell, wherein the direct repeat of the gRNA is modified by the insertion of distinct RNA sequence(s) that bind(s) to two or more adaptor proteins, and wherein each adaptor protein is associated with one or more functional domains; or, wherein the gRNA is modified to have at least one non-coding functional loop. In an aspect the invention provides vector(s) comprising nucleic acid molecule(s) encoding: non-naturally occurring or engineered CRISPR-Cas complex composition comprising the gRNA herein-discussed, and an RNA targeting enzyme, wherein optionally the RNA targeting enzyme comprises at least one mutation, such that the RNA targeting enzyme has no more than 5% of the nuclease activity of the RNA targeting enzyme not having the at least one mutation, and optionally one or more comprising at least one or more nuclear localization sequences. In an aspect a vector can further comprise regulatory element(s) operable in a eukaryotic cell operably linked to the nucleic acid molecule encoding the guide RNA (gRNA) and/or the nucleic acid molecule encoding the RNA targeting enzyme and/or the optional nuclear localization sequence(s).

[0531] In one aspect, the invention provides a kit comprising one or more of the components described herein-above. In some embodiments, the kit comprises a vector system as described above and instructions for using the kit.

[0532] In an aspect the invention provides a method of screening for gain of function (GOF) or loss of function (LOF) or for screen non-coding RNAs or potential regulatory regions (e.g. enhancers, repressors) comprising the cell line of as herein-discussed or cells of the model herein-discussed containing or expressing the RNA targeting

enzyme and introducing a composition as herein-discussed into cells of the cell line or model, whereby the gRNA includes either an activator or a repressor, and monitoring for GOF or LOF respectively as to those cells as to which the introduced gRNA includes an activator or as to those cells as to which the introduced gRNA includes a repressor.

[0533] In an aspect the invention provides a library of non-naturally occurring or engineered compositions, each comprising a RNA targeting CRISPR guide RNA (gRNA) comprising a guide sequence capable of hybridizing to a target RNA sequence of interest in a cell, an RNA targeting enzyme, wherein the RNA targeting enzyme comprises at least one mutation, such that the RNA targeting enzyme has no more than 5% of the nuclease activity of the RNA targeting enzyme not having the at least one mutation, wherein the gRNA is modified by the insertion of distinct RNA sequence(s) that bind to one or more adaptor proteins, and wherein the adaptor protein is associated with one or more functional domains, wherein the composition comprises one or more or two or more adaptor proteins, wherein the each protein is associated with one or more functional domains, and wherein the gRNAs comprise a genome wide library comprising a plurality of RNA targeting guide RNAs (gRNAs). In an aspect the invention provides a library as herein-discussed, wherein the RNA targeting RNA targeting enzyme has a diminished nuclease activity of at least 97%, or 100% as compare with the RNA targeting enzyme not having the at least one mutation. In an aspect the invention provides a library as herein-discussed, wherein the adaptor protein is a fusion protein comprising the functional domain. In an aspect the invention provides a library as herein discussed, wherein the gRNA is not modified by the insertion of distinct RNA sequence(s) that bind to the one or two or more adaptor proteins. In an aspect the invention provides a library as herein discussed, wherein the one or two or more functional domains are associated with the RNA targeting enzyme. In an aspect the invention provides a library as herein discussed, wherein the cell population of cells is a population of eukaryotic cells. In an aspect the invention provides a library as herein discussed, wherein the eukaryotic cell is a mammalian cell, a plant cell or a yeast cell. In an aspect the invention provides a library as herein discussed, wherein the mammalian cell is a human cell. In an aspect the invention provides a library as herein discussed, wherein the population of cells is a population of embryonic stem (ES) cells.

[0534] In an aspect the invention provides a library as herein discussed, wherein the targeting is of about 100 or more RNA sequences. In an aspect the invention provides a library as herein discussed, wherein the targeting is of about 1000 or more RNA sequences. In an aspect the invention provides a library as herein discussed, wherein the targeting is of about 20,000 or more sequences. In an aspect the invention provides a library as herein discussed, wherein the targeting is of the entire transcriptome. In an aspect the invention provides a library as herein discussed, wherein the targeting is of a panel of target sequences focused on a relevant or desirable pathway. In an aspect the invention provides a library as herein discussed, wherein the pathway is an immune pathway. In an aspect the invention provides a library as herein discussed, wherein the pathway is a cell division pathway.

[0535] In one aspect, the invention provides a method of generating a model eukaryotic cell comprising a gene with

modified expression. In some embodiments, a disease gene is any gene associated an increase in the risk of having or developing a disease. In some embodiments, the method comprises (a) introducing one or more vectors encoding the components of the system described herein above into a eukaryotic cell, and (b) allowing a CRISPR complex to bind to a target polynucleotide so as to modify expression of a gene, thereby generating a model eukaryotic cell comprising modified gene expression.

[0536] The structural information provided herein allows for interrogation of guide RNA interaction with the target RNA and the RNA targeting enzyme permitting engineering or alteration of guide RNA structure to optimize functionality of the entire RNA targeting CRISPR-Cas system. For example, the guide RNA may be extended, without colliding with the RNA targeting protein by the insertion of adaptor proteins that can bind to RNA. These adaptor proteins can further recruit effector proteins or fusions which comprise one or more functional domains.

[0537] An aspect of the invention is that the above elements are comprised in a single composition or comprised in individual compositions. These compositions may advantageously be applied to a host to elicit a functional effect on the genomic level.

[0538] The skilled person will understand that modifications to the guide RNA which allow for binding of the adapter+functional domain but not proper positioning of the adapter+functional domain (e.g. due to steric hindrance within the three dimensional structure of the CRISPR complex) are modifications which are not intended. The one or more modified guide RNA may be modified, by introduction of a distinct RNA sequence(s) 5' of the direct repeat, within the direct repeat, or 3' of the guide sequence.

[0539] The modified guide RNA, the inactivated RNA targeting enzyme (with or without functional domains), and the binding protein with one or more functional domains, may each individually be comprised in a composition and administered to a host individually or collectively. Alternatively, these components may be provided in a single composition for administration to a host. Administration to a host may be performed via viral vectors known to the skilled person or described herein for delivery to a host (e.g. lentiviral vector, adenoviral vector, AAV vector). As explained herein, use of different selection markers (e.g. for lentiviral gRNA selection) and concentration of gRNA (e.g. dependent on whether multiple gRNAs are used) may be advantageous for eliciting an improved effect.

[0540] Using the provided compositions, the person skilled in the art can advantageously and specifically target single or multiple loci with the same or different functional domains to elicit one or more genomic events. The compositions may be applied in a wide variety of methods for screening in libraries in cells and functional modeling in vivo (e.g. gene activation of lincRNA and identification of function; gain-of-function modeling; loss-of-function modeling; the use the compositions of the invention to establish cell lines and transgenic animals for optimization and screening purposes).

[0541] The current invention comprehends the use of the compositions of the current invention to establish and utilize conditional or inducible CRISPR RNA targeting events. (See, e.g., Platt et al., *Cell* (2014), <http://dx.doi.org/10.1016/j.cell.2014.09.014>, or PCT patent publications cited herein, such as WO 2014/093622 (PCT/US2013/074667), which

are not believed prior to the present invention or application). For example, the target cell comprises RNA targeting CRISPR enzyme conditionally or inducibly (e.g. in the form of Cre dependent constructs) and/or the adaptor protein conditionally or inducibly and, on expression of a vector introduced into the target cell, the vector expresses that which induces or gives rise to the condition of s RNA targeting enzyme expression and/or adaptor expression in the target cell. By applying the teaching and compositions of the current invention with the known method of creating a CRISPR complex, inducible gene expression affected by functional domains are also an aspect of the current invention. Alternatively, the adaptor protein may be provided as a conditional or inducible element with a conditional or inducible s RNA targeting enzyme to provide an effective model for screening purposes, which advantageously only requires minimal design and administration of specific gRNAs for a broad number of applications.

[0542] In an embodiment, a guide RNA molecule can be targeted to a known transcription response elements (e.g., promoters, enhancers, etc.), a known upstream activating sequences, and/or sequences of unknown or known function that are suspected of being able to control expression of the target DNA.

[0543] In some methods, a target polynucleotide can be inactivated to effect the modification of the expression in a cell. For example, upon the binding of a CRISPR complex to a target sequence in a cell, the target polynucleotide is inactivated such that the sequence is not transcribed, the coded protein is not produced, or the sequence does not function as the wild-type sequence does. For example, a protein or microRNA coding sequence may be inactivated such that the protein is not produced.

[0544] Further examples of applications of this system are described elsewhere herein.

Guide RNA According to the Invention Comprising a Dead Guide Sequence

[0545] In one aspect, the invention provides guide sequences which are modified in a manner which allows for formation of the CRISPR complex and successful binding to the target, while at the same time, not allowing for successful nuclease activity (i.e. without nuclease activity/without indel activity). For matters of explanation such modified guide sequences are referred to as “dead guides” or “dead guide sequences”. These dead guides or dead guide sequences can be thought of as catalytically inactive or conformationally inactive with regard to nuclease activity. Indeed, dead guide sequences may not sufficiently engage in productive base pairing with respect to the ability to promote catalytic activity or to distinguish on-target and off-target binding activity. Briefly, the assay involves synthesizing a CRISPR target RNA and guide RNAs comprising mismatches with the target RNA, combining these with the RNA targeting enzyme and analyzing cleavage based on gels based on the presence of bands generated by cleavage products, and quantifying cleavage based upon relative band intensities.

[0546] Hence, in a related aspect, the invention provides a non-naturally occurring or engineered composition RNA targeting CRISPR-Cas system comprising a functional RNA targeting as described herein, and guide RNA (gRNA) wherein the gRNA comprises a dead guide sequence whereby the gRNA is capable of hybridizing to a target

sequence such that the RNA targeting CRISPR-Cas system is directed to a genomic locus of interest in a cell without detectable RNA cleavage activity of a non-mutant RNA targeting enzyme of the system. It is to be understood that any of the gRNAs according to the invention as described herein elsewhere may be used as dead gRNAs/gRNAs comprising a dead guide sequence as described herein below. Any of the methods, products, compositions and uses as described herein elsewhere is equally applicable with the dead gRNAs/gRNAs comprising a dead guide sequence as further detailed below. By means of further guidance, the following particular aspects and embodiments are provided.

[0547] The ability of a dead guide sequence to direct sequence-specific binding of a CRISPR complex to an RNA target sequence may be assessed by any suitable assay. For example, the components of a CRISPR system sufficient to form a CRISPR complex, including the dead guide sequence to be tested, may be provided to a host cell having the corresponding target sequence, such as by transfection with vectors encoding the components of the CRISPR sequence, followed by an assessment of preferential cleavage within the target sequence. For instance, cleavage of a target RNA polynucleotide sequence may be evaluated in a test tube by providing the target sequence, components of a CRISPR complex, including the dead guide sequence to be tested and a control guide sequence different from the test dead guide sequence, and comparing binding or rate of cleavage at the target sequence between the test and control guide sequence reactions. Other assays are possible, and will occur to those skilled in the art. A dead guide sequence may be selected to target any target sequence. In some embodiments, the target sequence is a sequence within a genome of a cell.

[0548] As explained further herein, several structural parameters allow for a proper framework to arrive at such dead guides. Dead guide sequences are typically shorter than respective guide sequences which result in active RNA cleavage. In particular embodiments, dead guides are 5%, 10%, 20%, 30%, 40%, 50%, shorter than respective guides directed to the same.

[0549] As explained below and known in the art, one aspect of gRNA—RNA targeting specificity is the direct repeat sequence, which is to be appropriately linked to such guides. In particular, this implies that the direct repeat sequences are designed dependent on the origin of the RNA targeting enzyme. Thus, structural data available for validated dead guide sequences may be used for designing group29 specific equivalents. Structural similarity between, e.g., the orthologous nuclease domains HEPN of two or more group29 effector proteins may be used to transfer design equivalent dead guides. Thus, the dead guide herein may be appropriately modified in length and sequence to reflect such group29 specific equivalents, allowing for formation of the CRISPR complex and successful binding to the target RNA, while at the same time, not allowing for successful nuclease activity.

[0550] The use of dead guides in the context herein as well as the state of the art provides a surprising and unexpected platform for network biology and/or systems biology in both in vitro, ex vivo, and in vivo applications, allowing for multiplex gene targeting, and in particular bidirectional multiplex gene targeting. Prior to the use of dead guides, addressing multiple targets has been challenging and in some cases not possible. With the use of dead guides, multiple targets, and thus multiple activities, may be

addressed, for example, in the same cell, in the same animal, or in the same patient. Such multiplexing may occur at the same time or staggered for a desired timeframe.

[0551] For example, the dead guides allow to use gRNA as a means for gene targeting, without the consequence of nuclease activity, while at the same time providing directed means for activation or repression. Guide RNA comprising a dead guide may be modified to further include elements in a manner which allow for activation or repression of gene activity, in particular protein adaptors (e.g. aptamers) as described herein elsewhere allowing for functional placement of gene effectors (e.g. activators or repressors of gene activity). One example is the incorporation of aptamers, as explained herein and in the state of the art. By engineering the gRNA comprising a dead guide to incorporate protein-interacting aptamers (Koneremann et al., "Genome-scale transcription activation by an engineered CRISPR-Cas9 complex," doi:10.1038/nature14136, incorporated herein by reference), one may assemble multiple distinct effector domains. Such may be modeled after natural processes.

[0552] Thus, one aspect is a gRNA of the invention which comprises a dead guide, wherein the gRNA further comprises modifications which provide for gene activation or repression, as described herein. The dead gRNA may comprise one or more aptamers. The aptamers may be specific to gene effectors, gene activators or gene repressors. Alternatively, the aptamers may be specific to a protein which in turn is specific to and recruits/binds a specific gene effector, gene activator or gene repressor. If there are multiple sites for activator or repressor recruitment, it is preferred that the sites are specific to either activators or repressors. If there are multiple sites for activator or repressor binding, the sites may be specific to the same activators or same repressors. The sites may also be specific to different activators or different repressors. The effectors, activators, repressors may be present in the form of fusion proteins.

[0553] In an aspect, the invention provides a method of selecting a dead guide RNA targeting sequence for directing a functionalized CRISPR system to a gene locus in an organism, which comprises: a) locating one or more CRISPR motifs in the gene locus; b) analyzing the 20 nt sequence downstream of each CRISPR motif by: i) determining the GC content of the sequence; and ii) determining whether there are off-target matches of the first 15 nt of the sequence in the genome of the organism; c) selecting the sequence for use in a guide RNA if the GC content of the sequence is 70% or less and no off-target matches are identified. In an embodiment, the sequence is selected if the GC content is 50% or less. In an embodiment, the sequence is selected if the GC content is 40% or less. In an embodiment, the sequence is selected if the GC content is 30% or less. In an embodiment, two or more sequences are analyzed and the sequence having the lowest GC content is selected. In an embodiment, off-target matches are determined in regulatory sequences of the organism. In an embodiment, the gene locus is a regulatory region. An aspect provides a dead guide RNA comprising the targeting sequence selected according to the aforementioned methods.

[0554] In an aspect, the invention provides a dead guide RNA for targeting a functionalized CRISPR system to a gene locus in an organism. In an embodiment of the invention, the dead guide RNA comprises a targeting sequence wherein the CG content of the target sequence is 70% or less, and the first 15 nt of the targeting sequence does not

match an off-target sequence downstream from a CRISPR motif in the regulatory sequence of another gene locus in the organism. In certain embodiments, the GC content of the targeting sequence 60% or less, 55% or less, 50% or less, 45% or less, 40% or less, 35% or less or 30% or less. In certain embodiments, the GC content of the targeting sequence is from 70% to 60% or from 60% to 50% or from 50% to 40% or from 40% to 30%. In an embodiment, the targeting sequence has the lowest CG content among potential targeting sequences of the locus.

[0555] In an embodiment of the invention, the first 15 nt of the dead guide match the target sequence. In another embodiment, first 14 nt of the dead guide match the target sequence. In another embodiment, the first 13 nt of the dead guide match the target sequence. In another embodiment first 12 nt of the dead guide match the target sequence. In another embodiment, first 11 nt of the dead guide match the target sequence. In another embodiment, the first 10 nt of the dead guide match the target sequence. In an embodiment of the invention the first 15 nt of the dead guide does not match an off-target sequence downstream from a CRISPR motif in the regulatory region of another gene locus. In other embodiments, the first 14 nt, or the first 13 nt of the dead guide, or the first 12 nt of the guide, or the first 11 nt of the dead guide, or the first 10 nt of the dead guide, does not match an off-target sequence downstream from a CRISPR motif in the regulatory region of another gene locus. In other embodiments, the first 15 nt, or 14 nt, or 13 nt, or 12 nt, or 11 nt of the dead guide do not match an off-target sequence downstream from a CRISPR motif in the genome.

[0556] In certain embodiments, the dead guide RNA includes additional nucleotides at the 3'-end that do not match the target sequence. Thus, a dead guide RNA that includes the first 36 nt, downstream of a CRISPR motif can be extended in length at the 3' end.

Delivery of the Group 29 or Group 30 Effector Protein Complex or Components Thereof

[0557] Through this disclosure and the knowledge in the art, TALEs, CRISPR-Cas systems, or components thereof or nucleic acid molecules thereof (including, for instance HDR template) or nucleic acid molecules encoding or providing components thereof may be delivered by a delivery system herein described both generally and in detail.

[0558] Vector delivery, e.g., plasmid, viral delivery: The CRISPR enzyme, and/or any of the present RNAs, for instance a guide RNA, can be delivered using any suitable vector, e.g., plasmid or viral vectors, such as adeno associated virus (AAV), lentivirus, adenovirus or other viral vector types, or combinations thereof. Effector proteins and one or more guide RNAs can be packaged into one or more vectors, e.g., plasmid or viral vectors. In some embodiments, the vector, e.g., plasmid or viral vector is delivered to the tissue of interest by, for example, an intramuscular injection, while other times the delivery is via intravenous, transdermal, intranasal, oral, mucosal, or other delivery methods. Such delivery may be either via a single dose, or multiple doses. One skilled in the art understands that the actual dosage to be delivered herein may vary greatly depending upon a variety of factors, such as the vector choice, the target cell, organism, or tissue, the general condition of the subject to be treated, the degree of transformation/modification sought, the administration route, the administration mode, the type of transformation/modification sought, etc.

[0559] Such a dosage may further contain, for example, a carrier (water, saline, ethanol, glycerol, lactose, sucrose, calcium phosphate, gelatin, dextran, agar, pectin, peanut oil, sesame oil, etc.), a diluent, a pharmaceutically-acceptable carrier (e.g., phosphate-buffered saline), a pharmaceutically-acceptable excipient, and/or other compounds known in the art. The dosage may further contain one or more pharmaceutically acceptable salts such as, for example, a mineral acid salt such as a hydrochloride, a hydrobromide, a phosphate, a sulfate, etc.; and the salts of organic acids such as acetates, propionates, malonates, benzoates, etc. Additionally, auxiliary substances, such as wetting or emulsifying agents, pH buffering substances, gels or gelling materials, flavorings, colorants, microspheres, polymers, suspension agents, etc. may also be present herein. In addition, one or more other conventional pharmaceutical ingredients, such as preservatives, humectants, suspending agents, surfactants, antioxidants, anticaking agents, fillers, chelating agents, coating agents, chemical stabilizers, etc. may also be present, especially if the dosage form is a reconstitutable form. Suitable exemplary ingredients include microcrystalline cellulose, carboxymethylcellulose sodium, polysorbate 80, phenylethyl alcohol, chlorobutanol, potassium sorbate, sorbic acid, sulfur dioxide, propyl gallate, the parabens, ethyl vanillin, glycerin, phenol, parachlorophenol, gelatin, albumin and a combination thereof. A thorough discussion of pharmaceutically acceptable excipients is available in REMINGTON'S PHARMACEUTICAL SCIENCES (Mack Pub. Co., N.J. 1991) which is incorporated by reference herein.

[0560] In an embodiment herein the delivery is via an adenovirus, which may be at a single booster dose containing at least 1×10^5 particles (also referred to as particle units, pu) of adenoviral vector. In an embodiment herein, the dose preferably is at least about 1×10^6 particles (for example, about 1×10^6 - 1×10^{12} particles), more preferably at least about 1×10^7 particles, more preferably at least about 1×10^8 particles (e.g., about 1×10^8 - 1×10^{11} particles or about 1×10^8 - 1×10^{12} particles), and most preferably at least about 1×10^9 particles (e.g., about 1×10^9 - 1×10^{10} particles or about 1×10^9 - 1×10^{12} particles), or even at least about 1×10^{10} particles (e.g., about 1×10^{10} - 1×10^{12} particles) of the adenoviral vector. Alternatively, the dose comprises no more than about 1×10^{14} particles, preferably no more than about 1×10^{13} particles, even more preferably no more than about 1×10^{12} particles, even more preferably no more than about 1×10^{11} particles, and most preferably no more than about 1×10^{10} particles (e.g., no more than about 1×10^9 articles). Thus, the dose may contain a single dose of adenoviral vector with, for example, about 1×10^6 particle units (pu), about 2×10^6 pu, about 4×10^6 pu, about 1×10^7 pu, about 2×10^7 pu, about 4×10^7 pu, about 1×10^8 pu, about 2×10^8 pu, about 4×10^8 pu, about 1×10^9 pu, about 2×10^9 pu, about 4×10^9 pu, about 1×10^{10} pu, about 2×10^{10} pu, about 4×10^{10} pu, about 1×10^{11} pu, about 2×10^{11} pu, about 4×10^{11} pu, about 1×10^{12} pu, about 2×10^{12} pu, or about 4×10^{12} pu of adenoviral vector. See, for example, the adenoviral vectors in U.S. Pat. No. 8,454,972 B2 to Nabel, et. al., granted on Jun. 4, 2013; incorporated by reference herein, and the dosages at col 29, lines 36-58 thereof. In an embodiment herein, the adenovirus is delivered via multiple doses.

[0561] In an embodiment herein, the delivery is via an AAV. A therapeutically effective dosage for in vivo delivery of the AAV to a human is believed to be in the range of from about 20 to about 50 ml of saline solution containing from

about 1×10^{10} to about 1×10^{10} functional AAV/ml solution. The dosage may be adjusted to balance the therapeutic benefit against any side effects. In an embodiment herein, the AAV dose is generally in the range of concentrations of from about 1×10^5 to 1×10^{50} genomes AAV, from about 1×10^8 to 1×10^{20} genomes AAV, from about 1×10^{10} to about 1×10^{16} genomes, or about 1×10^{11} to about 1×10^{16} genomes AAV. A human dosage may be about 1×10^{13} genomes AAV. Such concentrations may be delivered in from about 0.001 ml to about 100 ml, about 0.05 to about 50 ml, or about 10 to about 25 ml of a carrier solution. Other effective dosages can be readily established by one of ordinary skill in the art through routine trials establishing dose response curves. See, for example, U.S. Pat. No. 8,404,658 B2 to Hajjar, et al., granted on Mar. 26, 2013, at col. 27, lines 45-60.

[0562] In an embodiment herein the delivery is via a plasmid. In such plasmid compositions, the dosage should be a sufficient amount of plasmid to elicit a response. For instance, suitable quantities of plasmid DNA in plasmid compositions can be from about 0.1 to about 2 mg, or from about 1 μ g to about 10 μ g per 70 kg individual. Plasmids of the invention will generally comprise (i) a promoter; (ii) a sequence encoding an nucleic acid-targeting CRISPR enzyme, operably linked to said promoter; (iii) a selectable marker; (iv) an origin of replication; and (v) a transcription terminator downstream of and operably linked to (ii). The plasmid can also encode the RNA components of a CRISPR complex, but one or more of these may instead be encoded on a different vector.

[0563] The doses herein are based on an average 70 kg individual. The frequency of administration is within the ambit of the medical or veterinary practitioner (e.g., physician, veterinarian), or scientist skilled in the art. It is also noted that mice used in experiments are typically about 20 g and from mice experiments one can scale up to a 70 kg individual.

[0564] In some embodiments the RNA molecules of the invention are delivered in liposome or lipofectin formulations and the like and can be prepared by methods well known to those skilled in the art. Such methods are described, for example, in U.S. Pat. Nos. 5,593,972, 5,589,466, and 5,580,859, which are herein incorporated by reference. Delivery systems aimed specifically at the enhanced and improved delivery of siRNA into mammalian cells have been developed, (see, for example, Shen et al FEBS Let. 2003, 539:111-114; Xia et al., Nat. Biotech. 2002, 20:1006-1010; Reich et al., Mol. Vision. 2003, 9: 210-216; Sorensen et al., J. Mol. Biol. 2003, 327: 761-766; Lewis et al., Nat. Gen. 2002, 32: 107-108 and Simeoni et al., NAR 2003, 31, 11: 2717-2724) and may be applied to the present invention. siRNA has recently been successfully used for inhibition of gene expression in primates (see for example, Tolentino et al., Retina 24(4):660 which may also be applied to the present invention.

[0565] Indeed, RNA delivery is a useful method of in vivo delivery. It is possible to deliver nucleic acid-targeting Cas protein and guide RNA (and, for instance, HR repair template) into cells using liposomes or particles. Thus delivery of the nucleic acid-targeting Cas protein/CRISPR enzyme, such as a Cas9 and/or delivery of the guide RNAs of the invention may be in RNA form and via microvesicles, liposomes or particles. For example, Cas mRNA and guide RNA can be packaged into liposomal particles for delivery in vivo. Liposomal transfection reagents such as lipo-

fectamine from Life Technologies and other reagents on the market can effectively deliver RNA molecules into the liver.

[0566] Means of delivery of RNA also preferred include delivery of RNA via particles (Cho, S., Goldberg, M., Son, S., Xu, Q., Yang, F., Mei, Y., Bogatyrev, S., Langer, R. and Anderson, D., Lipid-like nanoparticles for small interfering RNA delivery to endothelial cells, *Advanced Functional Materials*, 19: 3112-3118, 2010) or exosomes (Schroeder, A., Levins, C., Cortez, C., Langer, R., and Anderson, D., Lipid-based nanotherapeutics for siRNA delivery, *Journal of Internal Medicine*, 267: 9-21, 2010, PMID: 20059641). Indeed, exosomes have been shown to be particularly useful in delivery siRNA, a system with some parallels to the RNA-targeting system. For instance, El-Andaloussi S, et al. ("Exosome-mediated delivery of siRNA in vitro and in vivo." *Nat Protoc.* 2012 December; 7(12):2112-26. doi: 10.1038/nprot.2012.131. Epub 2012 Nov. 15.) describe how exosomes are promising tools for drug delivery across different biological barriers and can be harnessed for delivery of siRNA in vitro and in vivo. Their approach is to generate targeted exosomes through transfection of an expression vector, comprising an exosomal protein fused with a peptide ligand. The exosomes are then purified and characterized from transfected cell supernatant, then RNA is loaded into the exosomes. Delivery or administration according to the invention can be performed with exosomes, in particular but not limited to the brain. Vitamin E (α -tocopherol) may be conjugated with nucleic acid-targeting Cas protein and delivered to the brain along with high density lipoprotein (HDL), for example in a similar manner as was done by Uno et al. (*HUMAN GENE THERAPY* 22:711-719 (June 2011)) for delivering short-interfering RNA (siRNA) to the brain. Mice were infused via Osmotic minipumps (model 1007D; Alzet, Cupertino, Calif.) filled with phosphate-buffered saline (PBS) or free TocsiBACE or TocsiBACE/HDL and connected with Brain Infusion Kit 3 (Alzet). A brain-infusion cannula was placed about 0.5 mm posterior to the bregma at midline for infusion into the dorsal third ventricle. Uno et al. found that as little as 3 nmol of Toc-siRNA with HDL could induce a target reduction in comparable degree by the same ICV infusion method. A similar dosage of nucleic acid-targeting effector protein conjugated to α -tocopherol and co-administered with HDL targeted to the brain may be contemplated for humans in the present invention, for example, about 3 nmol to about 3 μ mol of nucleic acid-targeting effector protein targeted to the brain may be contemplated. Zou et al. (*HUMAN GENE THERAPY* 22:465-475 (April 2011)) describes a method of lentiviral-mediated delivery of short-hairpin RNAs targeting PKC γ for in vivo gene silencing in the spinal cord of rats. Zou et al. administered about 10 μ l of a recombinant lentivirus having a titer of 1×10^9 transducing units (TU)/ml by an intrathecal catheter. A similar dosage of nucleic acid-targeting effector protein expressed in a lentiviral vector targeted to the brain may be contemplated for humans in the present invention, for example, about 10-50 ml of nucleic acid-targeting effector protein targeted to the brain in a lentivirus having a titer of 1×10^9 transducing units (TU)/ml may be contemplated.

[0567] In terms of local delivery to the brain, this can be achieved in various ways. For instance, material can be delivered intrastrially e.g., by injection. Injection can be performed stereotactically via a craniotomy.

Packaging and Promoters Generally

[0568] Ways to package nucleic acid-targeting effector protein (such as Group 29 or Group 30 proteins) coding nucleic acid molecules, e.g., DNA, into vectors, e.g., viral vectors, to mediate genome modification in vivo include:

[0569] To achieve RNA cleavage (and hence knockdown):

[0570] Single virus vector:

[0571] Vector containing two or more expression cassettes:

[0572] Promoter-nucleic acid-targeting effector protein coding nucleic acid molecule-terminator

[0573] Promoter-guide RNA1-terminator

[0574] Promoter-guide RNA (N)-terminator (up to size limit of vector)

[0575] Double virus vector:

[0576] Vector 1 containing one expression cassette for driving the expression of nucleic acid-targeting effector protein (such as a Group 29 or Group 30)

[0577] Promoter-nucleic acid-targeting effector protein coding nucleic acid molecule-terminator

[0578] Vector 2 containing one more expression cassettes for driving the expression of one or more guide-eRNAs

[0579] Promoter-guide RNA1-terminator

[0580] Promoter-guide RNA1 (N)-terminator (up to size limit of vector)

[0581] The promoter used to drive nucleic acid-targeting effector protein coding nucleic acid molecule expression can include:

[0582] AAV ITR can serve as a promoter: this is advantageous for eliminating the need for an additional promoter element (which can take up space in the vector). The additional space freed up can be used to drive the expression of additional elements (gRNA, etc.). Also, ITR activity is relatively weaker, so can be used to reduce potential toxicity due to over expression of nucleic acid-targeting effector protein.

[0583] For ubiquitous expression, can use promoters: CMV, CAG, CBh, PGK, SV40, Ferritin heavy or light chains, etc.

[0584] For brain or other CNS expression, can use promoters: Synapsin1 for all neurons, CaMKIIalpha for excitatory neurons, GAD67 or GAD65 or VGAT for GABAergic neurons, etc.

[0585] For liver expression, can use Albumin promoter.

[0586] For lung expression, can use SP-B.

[0587] For endothelial cells, can use ICAM.

[0588] For hematopoietic cells can use IFNbeta or CD45.

[0589] For Osteoblasts can use OG-2.

[0590] The promoter used to drive guide RNA can include:

[0591] Pol III promoters such as U6 or H1

[0592] Use of Pol II promoter and intronic cassettes to express guide RNA

Adeno Associated Virus (AAV)

[0593] nucleic acid-targeting effector protein (such as a Group 29 or Group 30 effector protein) and one or more guide RNA can be delivered using adeno associated virus (AAV), lentivirus, adenovirus or other plasmid or viral vector types, in particular, using formulations and doses from, for example, U.S. Pat. No. 8,454,972 (formulations, doses for adenovirus), U.S. Pat. No. 8,404,658 (formula-

tions, doses for AAV) and U.S. Pat. No. 5,846,946 (formulations, doses for DNA plasmids) and from clinical trials and publications regarding the clinical trials involving lentivirus, AAV and adenovirus. For examples, for AAV, the route of administration, formulation and dose can be as in U.S. Pat. No. 8,454,972 and as in clinical trials involving AAV. For

AAV2, AAV5 or any combination thereof for targeting brain or neuronal cells; and one can select AAV4 for targeting cardiac tissue. AAV8 is useful for delivery to the liver. The herein promoters and vectors are preferred individually. A tabulation of certain AAV serotypes as to these cells (see Grimm, D. et al, J. Virol. 82: 5887-5911 (2008)) is as follows:

Cell Line	AAV-1	AAV-2	AAV-3	AAV-4	AAV-5	AAV-6	AAV-8	AAV-9
Huh-7	13	100	2.5	0.0	0.1	10	0.7	0.0
HEK293	25	100	2.5	0.1	0.1	5	0.7	0.1
HeLa	3	100	2.0	0.1	6.7	1	0.2	0.1
HepG2	3	100	16.7	0.3	1.7	5	0.3	ND
Hep1A	20	100	0.2	1.0	0.1	1	0.2	0.0
911	17	100	11	0.2	0.1	17	0.1	ND
CHO	100	100	14	1.4	333	50	10	1.0
COS	33	100	33	3.3	5.0	14	2.0	0.5
MeWo	10	100	20	0.3	6.7	10	1.0	0.2
NIH3T3	10	100	2.9	2.9	0.3	10	0.3	ND
A549	14	100	20	ND	0.5	10	0.5	0.1
HT1180	20	100	10	0.1	0.3	33	0.5	0.1
Monocytes	1111	100	ND	ND	125	1429	ND	ND
Immature DC	2500	100	ND	ND	222	2857	ND	ND
Mature DC	2222	100	ND	ND	333	3333	ND	ND

Adenovirus, the route of administration, formulation and dose can be as in U.S. Pat. No. 8,404,658 and as in clinical trials involving adenovirus. For plasmid delivery, the route of administration, formulation and dose can be as in U.S. Pat. No. 5,846,946 and as in clinical studies involving plasmids. Doses may be based on or extrapolated to an average 70 kg individual (e.g., a male adult human), and can be adjusted for patients, subjects, mammals of different weight and species. Frequency of administration is within the ambit of the medical or veterinary practitioner (e.g., physician, veterinarian), depending on usual factors including the age, sex, general health, other conditions of the patient or subject and the particular condition or symptoms being addressed. The viral vectors can be injected into the tissue of interest. For cell-type specific genome modification, the expression of nucleic acid-targeting effector protein (such as a Group 29 or Group 30 effector protein) can be driven by a cell-type specific promoter. For example, liver-specific expression might use the Albumin promoter and neuron-specific expression (e.g., for targeting CNS disorders) might use the Synapsin I promoter.

[0594] In terms of in vivo delivery, AAV is advantageous over other viral vectors for a couple of reasons:

[0595] Low toxicity (this may be due to the purification method not requiring ultra centrifugation of cell particles that can activate the immune response) and

[0596] Low probability of causing insertional mutagenesis because it doesn't integrate into the host genome.

[0597] AAV has a packaging limit of 4.5 or 4.75 Kb. This means that nucleic acid-targeting effector protein (such as a Group 29 or Group 30 effector protein) as well as a promoter and transcription terminator have to be all fit into the same viral vector. Therefore embodiments of the invention include utilizing homologs of nucleic acid-targeting effector protein (such as a Group 29 or Group 30 effector protein) that are shorter.

[0598] As to AAV, the AAV can be AAV1, AAV2, AAV5 or any combination thereof. One can select the AAV of the AAV with regard to the cells to be targeted; e.g., one can select AAV serotypes 1, 2, 5 or a hybrid capsid AAV1,

Lentivirus

[0599] Lentiviruses are complex retroviruses that have the ability to infect and express their genes in both mitotic and post-mitotic cells. The most commonly known lentivirus is the human immunodeficiency virus (HIV), which uses the envelope glycoproteins of other viruses to target a broad range of cell types.

[0600] Lentiviruses may be prepared as follows. After cloning pCasES10 (which contains a lentiviral transfer plasmid backbone), HEK293FT at low passage (p=5) were seeded in a T-75 flask to 50% confluence the day before transfection in DMEM with 10% fetal bovine serum and without antibiotics. After 20 hours, media was changed to OptiMEM (serum-free) media and transfection was done 4 hours later. Cells were transfected with 10 µg of lentiviral transfer plasmid (pCasES10) and the following packaging plasmids: 5 µg of pMD2.G (VSV-g pseudotype), and 7.5 µg of psPAX2 (gag/pol/rev/tat). Transfection was done in 4 mL OptiMEM with a cationic lipid delivery agent (50 µL Lipofectamine 2000 and 100 µL Plus reagent). After 6 hours, the media was changed to antibiotic-free DMEM with 10% fetal bovine serum. These methods use serum during cell culture, but serum-free methods are preferred.

[0601] Lentivirus may be purified as follows. Viral supernatants were harvested after 48 hours. Supernatants were first cleared of debris and filtered through a 0.45 µm low protein binding (PVDF) filter. They were then spun in a ultracentrifuge for 2 hours at 24,000 rpm. Viral pellets were resuspended in 50 µl of DMEM overnight at 4 C. They were then aliquotted and immediately frozen at -80° C.

[0602] In another embodiment, minimal non-primate lentiviral vectors based on the equine infectious anemia virus (EIAV) are also contemplated, especially for ocular gene therapy (see, e.g., Balagaan, J Gene Med 2006; 8: 275-285). In another embodiment, RetinoStat®, an equine infectious anemia virus-based lentiviral gene therapy vector that expresses angiostatic proteins endostatin and angiostatin that is delivered via a subretinal injection for the treatment of the web form of age-related macular degeneration is also contemplated (see, e.g., Binley et al., HUMAN GENE

THERAPY 23:980-991 (September 2012)) and this vector may be modified for the nucleic acid-targeting system of the present invention.

[0603] In another embodiment, self-inactivating lentiviral vectors with an siRNA targeting a common exon shared by HIV tat/rev, a nucleolar-localizing TAR decoy, and an anti-CCR5-specific hammerhead ribozyme (see, e.g., DiGiusto et al. (2010) *Sci Transl Med* 2:36ra43) may be used/and or adapted to the nucleic acid-targeting system of the present invention. A minimum of 2.5×10^6 CD34+ cells per kilogram patient weight may be collected and prestimulated for 16 to 20 hours in X-VIVO 15 medium (Lonza) containing 2 μ mol-glutamine, stem cell factor (100 ng/ml), Flt-3 ligand (Flt-3L) (100 ng/ml), and thrombopoietin (10 ng/ml) (Cell-Genix) at a density of 2×10^6 cells/ml. Prestimulated cells may be transduced with lentiviral at a multiplicity of infection of 5 for 16 to 24 hours in 75-cm² tissue culture flasks coated with fibronectin (25 mg/cm²) (RetroNectin, Takara Bio Inc.).

[0604] Lentiviral vectors have been disclosed as in the treatment for Parkinson's Disease, see, e.g., US Patent Publication No. 20120295960 and U.S. Pat. Nos. 7,303,910 and 7,351,585. Lentiviral vectors have also been disclosed for the treatment of ocular diseases, see e.g., US Patent Publication Nos. 20060281180, 20090007284, US20110117189; US20090017543; US20070054961, US20100317109. Lentiviral vectors have also been disclosed for delivery to the brain, see, e.g., US Patent Publication Nos. US20110293571; US20110293571, US20040013648, US20070025970, US20090111106 and U.S. Pat. No. 7,259,015.

RNA Delivery

[0605] RNA delivery: The nucleic acid-targeting Cas protein, and/or guide RNA, can also be delivered in the form of RNA. nucleic acid-targeting Cas protein (such as a Group 29 or Group 30 effector protein) mRNA can be generated using in vitro transcription. For example, nucleic acid-targeting effector protein (such as a Group 29 or Group 30 effector protein) mRNA can be synthesized using a PCR cassette containing the following elements: T7_promoter-kozak sequence (GCCACC)-effector protein-3' UTR (synthetic sequence) from beta globin-polyA tail (a string of 120 or more adenines). The cassette can be used for transcription by T7 polymerase. Guide RNAs can also be transcribed using in vitro transcription from a cassette containing T7_promoter-GG-guide RNA sequence.

[0606] To enhance expression and reduce possible toxicity, the nucleic acid-targeting effector protein-coding sequence and/or the guide RNA can be modified to include one or more modified nucleoside e.g., using pseudo-U or 5-Methyl-C.

[0607] mRNA delivery methods are especially promising for liver delivery currently.

[0608] Much clinical work on RNA delivery has focused on RNAi or antisense, but these systems can be adapted for delivery of RNA for implementing the present invention. References below to RNAi etc. should be read accordingly. Particle Delivery Systems and/or Formulations:

[0609] Several types of particle delivery systems and/or formulations are known to be useful in a diverse spectrum of biomedical applications. In general, a particle is defined as a small object that behaves as a whole unit with respect to its transport and properties. Particles are further classified

according to diameter. Coarse particles cover a range between 2,500 and 10,000 nanometers. Fine particles are sized between 100 and 2,500 nanometers. Ultrafine particles, or particles, are generally between 1 and 100 nanometers in size. The basis of the 100-nm limit is the fact that novel properties that differentiate particles from the bulk material typically develop at a critical length scale of under 100 nm.

[0610] As used herein, a particle delivery system/formulation is defined as any biological delivery system/formulation which includes a particle in accordance with the present invention. A particle in accordance with the present invention is any entity having a greatest dimension (e.g. diameter) of less than 100 microns (μ m). In some embodiments, inventive particles have a greatest dimension of less than 10 μ m. In some embodiments, inventive particles have a greatest dimension of less than 2000 nanometers (nm). In some embodiments, inventive particles have a greatest dimension of less than 1000 nanometers (nm). In some embodiments, inventive particles have a greatest dimension of less than 900 nm, 800 nm, 700 nm, 600 nm, 500 nm, 400 nm, 300 nm, 200 nm, or 100 nm. Typically, inventive particles have a greatest dimension (e.g., diameter) of 500 nm or less. In some embodiments, inventive particles have a greatest dimension (e.g., diameter) of 250 nm or less. In some embodiments, inventive particles have a greatest dimension (e.g., diameter) of 200 nm or less. In some embodiments, inventive particles have a greatest dimension (e.g., diameter) of 150 nm or less. In some embodiments, inventive particles have a greatest dimension (e.g., diameter) of 100 nm or less. Smaller particles, e.g., having a greatest dimension of 50 nm or less are used in some embodiments of the invention. In some embodiments, inventive particles have a greatest dimension ranging between 25 nm and 200 nm.

[0611] Particle characterization (including e.g., characterizing morphology, dimension, etc.) is done using a variety of different techniques. Common techniques are electron microscopy (TEM, SEM), atomic force microscopy (AFM), dynamic light scattering (DLS), X-ray photoelectron spectroscopy (XPS), powder X-ray diffraction (XRD), Fourier transform infrared spectroscopy (FTIR), matrix-assisted laser desorption/ionization time-of-flight mass spectrometry (MALDI-TOF), ultraviolet-visible spectroscopy, dual polarisation interferometry and nuclear magnetic resonance (NMR). Characterization (dimension measurements) may be made as to native particles (i.e., preloading) or after loading of the cargo (herein cargo refers to e.g., one or more components of CRISPR-Cas system e.g., CRISPR enzyme or mRNA or guide RNA, or any combination thereof, and may include additional carriers and/or excipients) to provide particles of an optimal size for delivery for any in vitro, ex vivo and/or in vivo application of the present invention. In certain preferred embodiments, particle dimension (e.g., diameter) characterization is based on measurements using dynamic laser scattering (DLS). Mention is made of U.S. Pat. No. 8,709,843; U.S. Pat. No. 6,007,845; U.S. Pat. No. 5,855,913; U.S. Pat. No. 5,985,309; U.S. Pat. No. 5,543,158; and the publication by James E. Dahlman and Carmen Barnes et al. *Nature Nanotechnology* (2014) published online 11 May 2014, doi:10.1038/nano.2014.84, concerning particles, methods of making and using them and measurements thereof.

[0612] Particles delivery systems within the scope of the present invention may be provided in any form, including

but not limited to solid, semi-solid, emulsion, or colloidal particles. As such any of the delivery systems described herein, including but not limited to, e.g., lipid-based systems, liposomes, micelles, microvesicles, exosomes, or gene gun may be provided as particle delivery systems within the scope of the present invention.

Particles

[0613] CRISPR enzyme mRNA and guide RNA may be delivered simultaneously using particles or lipid envelopes; for instance, CRISPR enzyme and RNA of the invention, e.g., as a complex, can be delivered via a particle as in Dahlman et al., WO2015089419 A2 and documents cited therein, such as 7C1 (see, e.g., James E. Dahlman and Carmen Barnes et al. *Nature Nanotechnology* (2014) published online 11 May 2014, doi:10.1038/nnano.2014.84), e.g., delivery particle comprising lipid or lipidoid and hydrophilic polymer, e.g., cationic lipid and hydrophilic polymer, for instance wherein the cationic lipid comprises 1,2-dioleoyl-3-trimethylammonium-propane (DOTAP) or 1,2-ditetradecanoyl-sn-glycero-3-phosphocholine (DMPC) and/or wherein the hydrophilic polymer comprises ethylene glycol or polyethylene glycol (PEG); and/or wherein the particle further comprises cholesterol (e.g., particle from formulation 1=DOTAP 100, DMPC 0, PEG 0, Cholesterol 0; formulation number 2=DOTAP 90, DMPC 0, PEG 10, Cholesterol 0; formulation number 3=DOTAP 90, DMPC 0, PEG 5, Cholesterol 5), wherein particles are formed using an efficient, multistep process wherein first, effector protein and RNA are mixed together, e.g., at a 1:1 molar ratio, e.g., at room temperature, e.g., for 30 minutes, e.g., in sterile, nuclease free 1×PBS; and separately, DOTAP, DMPC, PEG, and cholesterol as applicable for the formulation are dissolved in alcohol, e.g., 100% ethanol; and, the two solutions are mixed together to form particles containing the complexes).

[0614] Nucleic acid-targeting effector proteins (such as a Group 29 or Group 30 effector protein) mRNA and guide RNA may be delivered simultaneously using particles or lipid envelopes.

[0615] For example, Su X, Fricke J, Kavanagh D G, Irvine D J (“In vitro and in vivo mRNA delivery using lipid-enveloped pH-responsive polymer nanoparticles” *Mol Pharm.* 2011 Jun. 6; 8(3):774-87. doi: 10.1021/mp100390w. Epub 2011 Apr. 1) describes biodegradable core-shell structured particles with a poly(β -amino ester) (PBAE) core enveloped by a phospholipid bilayer shell. These were developed for in vivo mRNA delivery. The pH-responsive PBAE component was chosen to promote endosome disruption, while the lipid surface layer was selected to minimize toxicity of the polycation core. Such are, therefore, preferred for delivering RNA of the present invention.

[0616] In one embodiment, particles based on self-assembling bioadhesive polymers are contemplated, which may be applied to oral delivery of peptides, intravenous delivery of peptides and nasal delivery of peptides, all to the brain. Other embodiments, such as oral absorption and ocular delivery of hydrophobic drugs are also contemplated. The molecular envelope technology involves an engineered polymer envelope which is protected and delivered to the site of the disease (see, e.g., Mazza, M. et al. *ACS Nano*, 2013. 7(2): 1016-1026; Siew, A., et al. *Mol Pharm*, 2012. 9(1):14-28; Lalatsa, A., et al. *J Contr Rel*, 2012. 161(2):523-36; Lalatsa, A., et al., *Mol Pharm*, 2012. 9(6):1665-80;

Lalatsa, A., et al. *Mol Pharm*, 2012. 9(6):1764-74; Garrett, N. L., et al. *J Biophotonics*, 2012. 5(5-6):458-68; Garrett, N. L., et al. *J Raman Spect*, 2012. 43(5):681-688; Ahmad, S., et al. *J Royal Soc Interface* 2010. 7:S423-33; Uchegbu, I. F. *Expert Opin Drug Deliv*, 2006. 3(5):629-40; Qu, X., et al. *Biomacromolecules*, 2006. 7(12):3452-9 and Uchegbu, I. F., et al. *Int J Pharm*, 2001. 224:185-199). Doses of about 5 mg/kg are contemplated, with single or multiple doses, depending on the target tissue.

[0617] In one embodiment, particles that can deliver RNA to a cancer cell to stop tumor growth developed by Dan Anderson’s lab at MIT may be used/and or adapted to the nucleic acid-targeting system of the present invention. In particular, the Anderson lab developed fully automated, combinatorial systems for the synthesis, purification, characterization, and formulation of new biomaterials and nanoformulations. See, e.g., Alabi et al., *Proc Natl Acad Sci USA*. 2013 Aug. 6; 110(32):12881-6; Zhang et al., *Adv Mater*. 2013 Sep. 6; 25(33):4641-5; Jiang et al., *Nano Lett*. 2013 Mar. 13; 13(3):1059-64; Karagiannis et al., *ACS Nano*. 2012 Oct. 23; 6(10):8484-7; Whitehead et al., *ACS Nano*. 2012 Aug. 28; 6(8):6922-9 and Lee et al., *Nat Nanotechnol*. 2012 Jun. 3; 7(6):389-93.

[0618] US patent application 20110293703 relates to lipidoid compounds are also particularly useful in the administration of polynucleotides, which may be applied to deliver the nucleic acid-targeting system of the present invention. In one aspect, the aminoalcohol lipidoid compounds are combined with an agent to be delivered to a cell or a subject to form microparticles, particles, liposomes, or micelles. The agent to be delivered by the particles, liposomes, or micelles may be in the form of a gas, liquid, or solid, and the agent may be a polynucleotide, protein, peptide, or small molecule. The aminoalcohol lipidoid compounds may be combined with other aminoalcohol lipidoid compounds, polymers (synthetic or natural), surfactants, cholesterol, carbohydrates, proteins, lipids, etc. to form the particles. These particles may then optionally be combined with a pharmaceutical excipient to form a pharmaceutical composition.

[0619] US Patent Publication No. 20110293703 also provides methods of preparing the aminoalcohol lipidoid compounds. One or more equivalents of an amine are allowed to react with one or more equivalents of an epoxide-terminated compound under suitable conditions to form an aminoalcohol lipidoid compound of the present invention. In certain embodiments, all the amino groups of the amine are fully reacted with the epoxide-terminated compound to form tertiary amines. In other embodiments, all the amino groups of the amine are not fully reacted with the epoxide-terminated compound to form tertiary amines thereby resulting in primary or secondary amines in the aminoalcohol lipidoid compound. These primary or secondary amines are left as is or may be reacted with another electrophile such as a different epoxide-terminated compound. As will be appreciated by one skilled in the art, reacting an amine with less than excess of epoxide-terminated compound will result in a plurality of different aminoalcohol lipidoid compounds with various numbers of tails. Certain amines may be fully functionalized with two epoxide-derived compound tails while other molecules will not be completely functionalized with epoxide-derived compound tails. For example, a diamine or polyamine may include one, two, three, or four epoxide-derived compound tails off the various amino moi-

eties of the molecule resulting in primary, secondary, and tertiary amines. In certain embodiments, all the amino groups are not fully functionalized. In certain embodiments, two of the same types of epoxide-terminated compounds are used. In other embodiments, two or more different epoxide-terminated compounds are used. The synthesis of the aminoalcohol lipidoid compounds is performed with or without solvent, and the synthesis may be performed at higher temperatures ranging from 30-100° C., preferably at approximately 50-90° C. The prepared aminoalcohol lipidoid compounds may be optionally purified. For example, the mixture of aminoalcohol lipidoid compounds may be purified to yield an aminoalcohol lipidoid compound with a particular number of epoxide-derived compound tails. Or the mixture may be purified to yield a particular stereo- or regioisomer. The aminoalcohol lipidoid compounds may also be alkylated using an alkyl halide (e.g., methyl iodide) or other alkylating agent, and/or they may be acylated.

[0620] US Patent Publication No. 20110293703 also provides libraries of aminoalcohol lipidoid compounds prepared by the inventive methods. These aminoalcohol lipidoid compounds may be prepared and/or screened using high-throughput techniques involving liquid handlers, robots, microtiter plates, computers, etc. In certain embodiments, the aminoalcohol lipidoid compounds are screened for their ability to transfect polynucleotides or other agents (e.g., proteins, peptides, small molecules) into the cell.

[0621] US Patent Publication No. 20130302401 relates to a class of poly(beta-amino alcohols) (PBAs) has been prepared using combinatorial polymerization. The inventive PBAs may be used in biotechnology and biomedical applications as coatings (such as coatings of films or multilayer films for medical devices or implants), additives, materials, excipients, non-biofouling agents, micropatterning agents, and cellular encapsulation agents. When used as surface coatings, these PBAs elicited different levels of inflammation, both in vitro and in vivo, depending on their chemical structures. The large chemical diversity of this class of materials allowed us to identify polymer coatings that inhibit macrophage activation in vitro. Furthermore, these coatings reduce the recruitment of inflammatory cells, and reduce fibrosis, following the subcutaneous implantation of carboxylated polystyrene microparticles. These polymers may be used to form polyelectrolyte complex capsules for cell encapsulation. The invention may also have many other biological applications such as antimicrobial coatings, DNA or siRNA delivery, and stem cell tissue engineering. The teachings of US Patent Publication No. 20130302401 may be applied to the nucleic acid-targeting system of the present invention.

[0622] In another embodiment, lipid particles (LNPs) are contemplated. An antitranshyretin small interfering RNA has been encapsulated in lipid particles and delivered to humans (see, e.g., Coelho et al., *N Engl J Med* 2013; 369:819-29), and such a system may be adapted and applied to the nucleic acid-targeting system of the present invention. Doses of about 0.01 to about 1 mg per kg of body weight administered intravenously are contemplated. Medications to reduce the risk of infusion-related reactions are contemplated, such as dexamethasone, acetaminophen, diphenhydramine or cetirizine, and ranitidine are contemplated. Multiple doses of about 0.3 mg per kilogram every 4 weeks for five doses are also contemplated.

[0623] LNPs have been shown to be highly effective in delivering siRNAs to the liver (see, e.g., Tabernero et al., *Cancer Discovery*, April 2013, Vol. 3, No. 4, pages 363-470) and are therefore contemplated for delivering RNA encoding nucleic acid-targeting effector protein to the liver. A dosage of about four doses of 6 mg/kg of the LNP every two weeks may be contemplated. Tabernero et al. demonstrated that tumor regression was observed after the first 2 cycles of LNPs dosed at 0.7 mg/kg, and by the end of 6 cycles the patient had achieved a partial response with complete regression of the lymph node metastasis and substantial shrinkage of the liver tumors. A complete response was obtained after 40 doses in this patient, who has remained in remission and completed treatment after receiving doses over 26 months. Two patients with RCC and extrahepatic sites of disease including kidney, lung, and lymph nodes that were progressing following prior therapy with VEGF pathway inhibitors had stable disease at all sites for approximately 8 to 12 months, and a patient with PNET and liver metastases continued on the extension study for 18 months (36 doses) with stable disease.

[0624] However, the charge of the LNP must be taken into consideration. As cationic lipids combined with negatively charged lipids to induce nonbilayer structures that facilitate intracellular delivery. Because charged LNPs are rapidly cleared from circulation following intravenous injection, ionizable cationic lipids with pKa values below 7 were developed (see, e.g., Rosin et al, *Molecular Therapy*, vol. 19, no. 12, pages 1286-2200, December 2011). Negatively charged polymers such as RNA may be loaded into LNPs at low pH values (e.g., pH 4) where the ionizable lipids display a positive charge. However, at physiological pH values, the LNPs exhibit a low surface charge compatible with longer circulation times. Four species of ionizable cationic lipids have been focused upon, namely 1,2-dilinoeyl-3-dimethylammonium-propane (DLinDAP), 1,2-dilinoeyloxy-3-N,N-dimethylaminopropane (DLinDMA), 1,2-dilinoeyloxyketo-N,N-dimethyl-3-aminopropane (DLinKDMA), and 1,2-dilinoeyl-4-(2-dimethylaminoethyl)-[1,3]-dioxolane (DLinKC2-DMA). It has been shown that LNP siRNA systems containing these lipids exhibit remarkably different gene silencing properties in hepatocytes in vivo, with potencies varying according to the series DLinKC2-DMA>DLinKDMA>DLinDMA>>DLinDAP employing a Factor VII gene silencing model (see, e.g., Rosin et al, *Molecular Therapy*, vol. 19, no. 12, pages 1286-2200, December 2011). A dosage of 1 µg/ml of LNP or CRISPR-Cas RNA in or associated with the LNP may be contemplated, especially for a formulation containing DLinKC2-DMA.

[0625] Preparation of LNPs and CRISPR-Cas encapsulation may be used/and or adapted from Rosin et al, *Molecular Therapy*, vol. 19, no. 12, pages 1286-2200, December 2011). The cationic lipids 1,2-dilinoeyl-3-dimethylammonium-propane (DLinDAP), 1,2-dilinoeyloxy-3-N,N-dimethylaminopropane (DLinDMA), 1,2-dilinoeyloxyketo-N,N-dimethyl-3-aminopropane (DLinK-DMA), 1,2-dilinoeyl-4-(2-dimethylaminoethyl)-[1,3]-dioxolane (DLinKC2-DMA), (3-o-[2"-(methoxypolyethyleneglycol 2000) succinoyl]-1,2-dimyristoyl-sn-glycol (PEG-S-DMG), and R-3-[ω-methoxy-poly(ethylene glycol)2000) carbamoyl]-1,2-dimyristoylpropyl-3-amine (PEG-C-DOMG) may be provided by Tekmira Pharmaceuticals (Vancouver, Canada) or synthesized. Cholesterol may be purchased from Sigma

(St Louis, Mo.). The specific nucleic acid-targeting complex (CRISPR-Cas) RNA may be encapsulated in LNPs containing DLinDAP, DLinDMA, DLinK-DMA, and DLinKC2-DMA (cationic lipid:DSPC:CHOL: PEGS-DMG or PEG-C-DOMG at 40:10:40:10 molar ratios). When required, 0.2% SP-DiOC18 (Invitrogen, Burlington, Canada) may be incorporated to assess cellular uptake, intracellular delivery, and biodistribution. Encapsulation may be performed by dissolving lipid mixtures comprised of cationic lipid:DSPC: cholesterol:PEG-c-DOMG (40:10:40:10 molar ratio) in ethanol to a final lipid concentration of 10 mmol/l. This ethanol solution of lipid may be added drop-wise to 50 mmol/l citrate, pH 4.0 to form multilamellar vesicles to produce a final concentration of 30% ethanol vol/vol. Large unilamellar vesicles may be formed following extrusion of multilamellar vesicles through two stacked 80 nm Nuclepore polycarbonate filters using the Extruder (Northern Lipids, Vancouver, Canada). Encapsulation may be achieved by adding RNA dissolved at 2 mg/ml in 50 mmol/l citrate, pH 4.0 containing 30% ethanol vol/vol drop-wise to extruded preformed large unilamellar vesicles and incubation at 31° C. for 30 minutes with constant mixing to a final RNA/lipid weight ratio of 0.06/1 wt/wt. Removal of ethanol and neutralization of formulation buffer were performed by dialysis against phosphate-buffered saline (PBS), pH 7.4 for 16 hours using Spectra/Por 2 regenerated cellulose dialysis membranes. Particle size distribution may be determined by dynamic light scattering using a NICOMP 370 particle sizer, the vesicle/intensity modes, and Gaussian fitting (Nicomp Particle Sizing, Santa Barbara, Calif.). The particle size for all three LNP systems may be ~70 nm in diameter. RNA encapsulation efficiency may be determined by removal of free RNA using VivaPureD MiniH columns (Sartorius Stedim Biotech) from samples collected before and after dialysis. The encapsulated RNA may be extracted from the eluted particles and quantified at 260 nm. RNA to lipid ratio was determined by measurement of cholesterol content in vesicles using the Cholesterol E enzymatic assay from Wako Chemicals USA (Richmond, Va.). In conjunction with the herein discussion of LNPs and PEG lipids, PEGylated liposomes or LNPs are likewise suitable for delivery of a nucleic acid-targeting system or components thereof

[0626] Preparation of large LNPs may be used/and or adapted from Rosin et al, *Molecular Therapy*, vol. 19, no. 12, pages 1286-2200, December 2011. A lipid premix solution (20.4 mg/ml total lipid concentration) may be prepared in ethanol containing DLinKC2-DMA, DSPC, and cholesterol at 50:10:38.5 molar ratios. Sodium acetate may be added to the lipid premix at a molar ratio of 0.75:1 (sodium acetate: DLinKC2-DMA). The lipids may be subsequently hydrated by combining the mixture with 1.85 volumes of citrate buffer (10 mmol/l, pH 3.0) with vigorous stirring, resulting in spontaneous liposome formation in aqueous buffer containing 35% ethanol. The liposome solution may be incubated at 37° C. to allow for time-dependent increase in particle size. Aliquots may be removed at various times during incubation to investigate changes in liposome size by dynamic light scattering (Zetasizer Nano Z S, Malvern Instruments, Worcestershire, UK). Once the desired particle size is achieved, an aqueous PEG lipid solution (stock=10 mg/ml PEG-DMG in 35% (vol/vol) ethanol) may be added to the liposome mixture to yield a final PEG molar concentration of 3.5% of total lipid. Upon addition of PEG-lipids, the liposomes should their size, effectively quenching fur-

ther growth. RNA may then be added to the empty liposomes at a RNA to total lipid ratio of approximately 1:10 (wt:wt), followed by incubation for 30 minutes at 37° C. to form loaded LNPs. The mixture may be subsequently dialyzed overnight in PBS and filtered with a 0.45- μ m syringe filter.

[0627] Spherical Nucleic Acid (SNATM) constructs and other particles (particularly gold particles) are also contemplated as a means to delivery nucleic acid-targeting system to intended targets. Significant data show that AuraSense Therapeutics' Spherical Nucleic Acid (SNATM) constructs, based upon nucleic acid-functionalized gold particles, are useful.

[0628] Literature that may be employed in conjunction with herein teachings include: Cutler et al., *J. Am. Chem. Soc.* 2011 133:9254-9257, Hao et al., *Small.* 2011 7:3158-3162, Zhang et al., *ACS Nano.* 2011 5:6962-6970, Cutler et al., *J. Am. Chem. Soc.* 2012 134:1376-1391, Young et al., *Nano Lett.* 2012 12:3867-71, Zheng et al., *Proc. Natl. Acad. Sci. USA.* 2012 109:11975-80, Mirkin, *Nanomedicine* 2012 7:635-638 Zhang et al., *J. Am. Chem. Soc.* 2012 134:16488-1691, Weintraub, *Nature* 2013 495:S14-S16, Choi et al., *Proc. Natl. Acad. Sci. USA.* 2013 110(19):7625-7630, Jensen et al., *Sci. Transl. Med.* 5, 209ra152 (2013) and Mirkin, et al., *Small*, 10:186-192.

[0629] Self-assembling particles with RNA may be constructed with polyethyleneimine (PEI) that is PEGylated with an Arg-Gly-Asp (RGD) peptide ligand attached at the distal end of the polyethylene glycol (PEG). This system has been used, for example, as a means to target tumor neovasculature expressing integrins and deliver siRNA inhibiting vascular endothelial growth factor receptor-2 (VEGF R2) expression and thereby achieve tumor angiogenesis (see, e.g., Schiffelers et al., *Nucleic Acids Research*, 2004, Vol. 32, No. 19). Nanoplexes may be prepared by mixing equal volumes of aqueous solutions of cationic polymer and nucleic acid to give a net molar excess of ionizable nitrogen (polymer) to phosphate (nucleic acid) over the range of 2 to 6. The electrostatic interactions between cationic polymers and nucleic acid resulted in the formation of polyplexes with average particle size distribution of about 100 nm, hence referred to here as nanoplexes. A dosage of about 100 to 200 mg of nucleic acid-targeting complex RNA is envisioned for delivery in the self-assembling particles of Schiffelers et al.

[0630] The nanoplexes of Bartlett et al. (*PNAS*, Sep. 25, 2007, vol. 104, no. 39) may also be applied to the present invention. The nanoplexes of Bartlett et al. are prepared by mixing equal volumes of aqueous solutions of cationic polymer and nucleic acid to give a net molar excess of ionizable nitrogen (polymer) to phosphate (nucleic acid) over the range of 2 to 6. The electrostatic interactions between cationic polymers and nucleic acid resulted in the formation of polyplexes with average particle size distribution of about 100 nm, hence referred to here as nanoplexes. The DOTA-siRNA of Bartlett et al. was synthesized as follows:

1,4,7,10-tetraazacyclododecane-1,4,7,10-tetraacetic acid mono(N-hydroxysuccinimide ester) (DOTA-NHSester) was ordered from Macrocyclics (Dallas, Tex.). The amine modified RNA sense strand with a 100-fold molar excess of DOTA-NHS-ester in carbonate buffer (pH 9) was added to a microcentrifuge tube. The contents were reacted by stirring for 4 h at room temperature. The DOTA-RNAsense conjugate was ethanol-precipitated, resuspended in water, and annealed to the unmodified antisense strand to

yield DOTA-siRNA. All liquids were pretreated with Chelex-100 (Bio-Rad, Hercules, Calif.) to remove trace metal contaminants. Tf-targeted and nontargeted siRNA particles may be formed by using cyclodextrin-containing polycations. Typically, particles were formed in water at a charge ratio of 3 (+/-) and an siRNA concentration of 0.5 g/liter. One percent of the adamantane-PEG molecules on the surface of the targeted particles were modified with Tf (adamantane-PEG-Tf). The particles were suspended in a 5% (wt/vol) glucose carrier solution for injection.

[0631] Davis et al. (Nature, Vol 464, 15 Apr. 2010) conducts a RNA clinical trial that uses a targeted particle-delivery system (clinical trial registration number NCT00689065). Patients with solid cancers refractory to standard-of-care therapies are administered doses of targeted particles on days 1, 3, 8 and 10 of a 21-day cycle by a 30-min intravenous infusion. The particles comprise, consist essentially of, or consist of a synthetic delivery system containing: (1) a linear, cyclodextrin-based polymer (CDP), (2) a human transferrin protein (TF) targeting ligand displayed on the exterior of the particle to engage TF receptors (TFR) on the surface of the cancer cells, (3) a hydrophilic polymer (polyethylene glycol (PEG) used to promote particle stability in biological fluids), and (4) siRNA designed to reduce the expression of the RRM2 (sequence used in the clinic was previously denoted siR2B+5). The TFR has long been known to be upregulated in malignant cells, and RRM2 is an established anti-cancer target. These particles (clinical version denoted as CALAA-01) have been shown to be well tolerated in multi-dosing studies in non-human primates. Although a single patient with chronic myeloid leukaemia has been administered siRNA by liposomal delivery, Davis et al.'s clinical trial is the initial human trial to systemically deliver siRNA with a targeted delivery system and to treat patients with solid cancer. To ascertain whether the targeted delivery system can provide effective delivery of functional siRNA to human tumours, Davis et al. investigated biopsies from three patients from three different dosing cohorts; patients A, B and C, all of whom had metastatic melanoma and received CALAA-01 doses of 18, 24 and 30 mg m⁻² siRNA, respectively. Similar doses may also be contemplated for the nucleic acid-targeting system of the present invention. The delivery of the invention may be achieved with particles containing a linear, cyclodextrin-based polymer (CDP), a human transferrin protein (TF) targeting ligand displayed on the exterior of the particle to engage TF receptors (TFR) on the surface of the cancer cells and/or a hydrophilic polymer (for example, polyethylene glycol (PEG) used to promote particle stability in biological fluids).

[0632] In terms of this invention, it is preferred to have one or more components of nucleic acid-targeting complex, e.g., nucleic acid-targeting effector protein or mRNA, or guide RNA delivered using particles or lipid envelopes. Other delivery systems or vectors may be used in conjunction with the particle aspects of the invention.

[0633] In general, a "nanoparticle" refers to any particle having a diameter of less than 1000 nm. In certain preferred embodiments, nanoparticles of the invention have a greatest dimension (e.g., diameter) of 500 nm or less. In other preferred embodiments, nanoparticles of the invention have a greatest dimension ranging between 25 nm and 200 nm. In other preferred embodiments, particles of the invention have a greatest dimension of 100 nm or less. In other preferred

embodiments, nanoparticles of the invention have a greatest dimension ranging between 35 nm and 60 nm.

[0634] Particles encompassed in the present invention may be provided in different forms, e.g., as solid particles (e.g., metal such as silver, gold, iron, titanium), non-metal, lipid-based solids, polymers), suspensions of particles, or combinations thereof. Metal, dielectric, and semiconductor particles may be prepared, as well as hybrid structures (e.g., core-shell particles). Particles made of semiconducting material may also be labeled quantum dots if they are small enough (typically sub 10 nm) that quantization of electronic energy levels occurs. Such nanoscale particles are used in biomedical applications as drug carriers or imaging agents and may be adapted for similar purposes in the present invention.

[0635] Semi-solid and soft particles have been manufactured, and are within the scope of the present invention. A prototype particle of semi-solid nature is the liposome. Various types of liposome particles are currently used clinically as delivery systems for anticancer drugs and vaccines. Particles with one half hydrophilic and the other half hydrophobic are termed Janus particles and are particularly effective for stabilizing emulsions. They can self-assemble at water/oil interfaces and act as solid surfactants.

[0636] U.S. Pat. No. 8,709,843, incorporated herein by reference, provides a drug delivery system for targeted delivery of therapeutic agent-containing particles to tissues, cells, and intracellular compartments. The invention provides targeted particles comprising polymer conjugated to a surfactant, hydrophilic polymer or lipid.

[0637] U.S. Pat. No. 6,007,845, incorporated herein by reference, provides particles which have a core of a multi-block copolymer formed by covalently linking a multifunctional compound with one or more hydrophobic polymers and one or more hydrophilic polymers, and contain a biologically active material.

[0638] U.S. Pat. No. 5,855,913, incorporated herein by reference, provides a particulate composition having aerodynamically light particles having a tap density of less than 0.4 g/cm³ with a mean diameter of between 5 µm and 30 µm, incorporating a surfactant on the surface thereof for drug delivery to the pulmonary system.

[0639] U.S. Pat. No. 5,985,309, incorporated herein by reference, provides particles incorporating a surfactant and/or a hydrophilic or hydrophobic complex of a positively or negatively charged therapeutic or diagnostic agent and a charged molecule of opposite charge for delivery to the pulmonary system.

[0640] U.S. Pat. No. 5,543,158, incorporated herein by reference, provides biodegradable injectable particles having a biodegradable solid core containing a biologically active material and poly(alkylene glycol) moieties on the surface.

[0641] WO2012135025 (also published as US20120251560), incorporated herein by reference, describes conjugated polyethyleneimine (PEI) polymers and conjugated aza-macrocycles (collectively referred to as "conjugated lipomer" or "lipomers"). In certain embodiments, it can be envisioned that such methods and materials of herein-cited documents, e.g., conjugated lipomers can be used in the context of the nucleic acid-targeting system to achieve in vitro, ex vivo and in vivo genomic perturbations to modify gene expression, including modulation of protein expression.

[0642] In one embodiment, the particle may be epoxide-modified lipid-polymer, advantageously 7C1 (see, e.g., James E. Dahlman and Carmen Barnes et al. *Nature Nanotechnology* (2014) published online 11 May 2014, doi:10.1038/nnano.2014.84). C71 was synthesized by reacting C15 epoxide-terminated lipids with PEI600 at a 14:1 molar ratio, and was formulated with C14PEG2000 to produce particles (diameter between 35 and 60 nm) that were stable in PBS solution for at least 40 days.

[0643] An epoxide-modified lipid-polymer may be utilized to deliver the nucleic acid-targeting system of the present invention to pulmonary, cardiovascular or renal cells, however, one of skill in the art may adapt the system to deliver to other target organs. Dosage ranging from about 0.05 to about 0.6 mg/kg are envisioned. Dosages over several days or weeks are also envisioned, with a total dosage of about 2 mg/kg.

Exosomes

[0644] Exosomes are endogenous nano-vesicles that transport RNAs and proteins, and which can deliver RNA to the brain and other target organs. To reduce immunogenicity, Alvarez-Erviti et al. (2011, *Nat Biotechnol* 29: 341) used self-derived dendritic cells for exosome production. Targeting to the brain was achieved by engineering the dendritic cells to express Lamp2b, an exosomal membrane protein, fused to the neuron-specific RVG peptide. Purified exosomes were loaded with exogenous RNA by electroporation. Intravenously injected RVG-targeted exosomes delivered GAPDH siRNA specifically to neurons, microglia, oligodendrocytes in the brain, resulting in a specific gene knockdown. Pre-exposure to RVG exosomes did not attenuate knockdown, and non-specific uptake in other tissues was not observed. The therapeutic potential of exosome-mediated siRNA delivery was demonstrated by the strong mRNA (60%) and protein (62%) knockdown of BACE1, a therapeutic target in Alzheimer's disease.

[0645] To obtain a pool of immunologically inert exosomes, Alvarez-Erviti et al. harvested bone marrow from inbred C57BL/6 mice with a homogenous major histocompatibility complex (MHC) haplotype. As immature dendritic cells produce large quantities of exosomes devoid of T-cell activators such as MHC-II and CD86, Alvarez-Erviti et al. selected for dendritic cells with granulocyte/macrophage-colony stimulating factor (GM-CSF) for 7 d. Exosomes were purified from the culture supernatant the following day using well-established ultracentrifugation protocols. The exosomes produced were physically homogenous, with a size distribution peaking at 80 nm in diameter as determined by particle tracking analysis (NTA) and electron microscopy. Alvarez-Erviti et al. obtained 6-12 μg of exosomes (measured based on protein concentration) per 10^6 cells.

[0646] Next, Alvarez-Erviti et al. investigated the possibility of loading modified exosomes with exogenous cargoes using electroporation protocols adapted for nanoscale applications. As electroporation for membrane particles at the nanometer scale is not well-characterized, nonspecific Cy5-labeled RNA was used for the empirical optimization of the electroporation protocol. The amount of encapsulated RNA was assayed after ultracentrifugation and lysis of exosomes. Electroporation at 400 V and 125 μF resulted in the greatest retention of RNA and was used for all subsequent experiments.

[0647] Alvarez-Erviti et al. administered 150 μg of each BACE1 siRNA encapsulated in 150 μg of RVG exosomes to normal C57BL/6 mice and compared the knockdown efficiency to four controls: untreated mice, mice injected with RVG exosomes only, mice injected with BACE1 siRNA complexed to an in vivo cationic liposome reagent and mice injected with BACE1 siRNA complexed to RVG-9R, the RVG peptide conjugated to 9 D-arginines that electrostatically binds to the siRNA. Cortical tissue samples were analyzed 3 d after administration and a significant protein knockdown (45%, $P < 0.05$, versus 62%, $P < 0.01$) in both siRNA-RVG-9R-treated and siRNARVG exosome-treated mice was observed, resulting from a significant decrease in BACE1 mRNA levels (66% [+ or -] 15%, $P < 0.001$ and 61% [+ or -] 13% respectively, $P < 0.01$). Moreover, Applicants demonstrated a significant decrease (55%, $P < 0.05$) in the total [beta]-amyloid 1-42 levels, a main component of the amyloid plaques in Alzheimer's pathology, in the RVG-exosome-treated animals. The decrease observed was greater than the β -amyloid 1-40 decrease demonstrated in normal mice after intraventricular injection of BACE1 inhibitors. Alvarez-Erviti et al. carried out 5'-rapid amplification of cDNA ends (RACE) on BACE1 cleavage product, which provided evidence of RNAi-mediated knockdown by the siRNA.

[0648] Finally, Alvarez-Erviti et al. investigated whether RNA-RVG exosomes induced immune responses in vivo by assessing IL-6, IP-10, TNF α and IFN- α serum concentrations. Following exosome treatment, nonsignificant changes in all cytokines were registered similar to siRNA-transfection reagent treatment in contrast to siRNA-RVG-9R, which potently stimulated IL-6 secretion, confirming the immunologically inert profile of the exosome treatment. Given that exosomes encapsulate only 20% of siRNA, delivery with RVG-exosome appears to be more efficient than RVG-9R delivery as comparable mRNA knockdown and greater protein knockdown was achieved with fivefold less siRNA without the corresponding level of immune stimulation. This experiment demonstrated the therapeutic potential of RVG-exosome technology, which is potentially suited for long-term silencing of genes related to neurodegenerative diseases. The exosome delivery system of Alvarez-Erviti et al. may be applied to deliver the nucleic acid-targeting system of the present invention to therapeutic targets, especially neurodegenerative diseases. A dosage of about 100 to 1000 mg of nucleic acid-targeting system encapsulated in about 100 to 1000 mg of RVG exosomes may be contemplated for the present invention.

[0649] El-Andaloussi et al. (*Nature Protocols* 7, 2112-2126(2012)) discloses how exosomes derived from cultured cells can be harnessed for delivery of RNA in vitro and in vivo. This protocol first describes the generation of targeted exosomes through transfection of an expression vector, comprising an exosomal protein fused with a peptide ligand. Next, El-Andaloussi et al. explain how to purify and characterize exosomes from transfected cell supernatant. Next, El-Andaloussi et al. detail crucial steps for loading RNA into exosomes. Finally, El-Andaloussi et al. outline how to use exosomes to efficiently deliver RNA in vitro and in vivo in mouse brain. Examples of anticipated results in which exosome-mediated RNA delivery is evaluated by functional assays and imaging are also provided. The entire protocol takes ~3 weeks. Delivery or administration according to the invention may be performed using exosomes produced from

self-derived dendritic cells. From the herein teachings, this can be employed in the practice of the invention

[0650] In another embodiment, the plasma exosomes of Wahlgren et al. (Nucleic Acids Research, 2012, Vol. 40, No. 17 e130) are contemplated. Exosomes are nano-sized vesicles (30-90 nm in size) produced by many cell types, including dendritic cells (DC), B cells, T cells, mast cells, epithelial cells and tumor cells. These vesicles are formed by inward budding of late endosomes and are then released to the extracellular environment upon fusion with the plasma membrane. Because exosomes naturally carry RNA between cells, this property may be useful in gene therapy, and from this disclosure can be employed in the practice of the instant invention.

[0651] Exosomes from plasma can be prepared by centrifugation of buffy coat at 900 g for 20 min to isolate the plasma followed by harvesting cell supernatants, centrifuging at 300 g for 10 min to eliminate cells and at 16 500 g for 30 min followed by filtration through a 0.22 mm filter. Exosomes are pelleted by ultracentrifugation at 120 000 g for 70 min. Chemical transfection of siRNA into exosomes is carried out according to the manufacturer's instructions in RNAi Human/Mouse Starter Kit (Quiagen, Hilden, Germany). siRNA is added to 100 ml PBS at a final concentration of 2 mmol/ml. After adding HiPerFect transfection reagent, the mixture is incubated for 10 min at RT. In order to remove the excess of micelles, the exosomes are re-isolated using aldehyde/sulfate latex beads. The chemical transfection of nucleic acid-targeting system into exosomes may be conducted similarly to siRNA. The exosomes may be co-cultured with monocytes and lymphocytes isolated from the peripheral blood of healthy donors. Therefore, it may be contemplated that exosomes containing nucleic acid-targeting system may be introduced to monocytes and lymphocytes of and autologously reintroduced into a human. Accordingly, delivery or administration according to the invention may be performed using plasma exosomes.

Liposomes

[0652] Delivery or administration according to the invention can be performed with liposomes. Liposomes are spherical vesicle structures composed of a uni- or multilamellar lipid bilayer surrounding internal aqueous compartments and a relatively impermeable outer lipophilic phospholipid bilayer. Liposomes have gained considerable attention as drug delivery carriers because they are biocompatible, nontoxic, can deliver both hydrophilic and lipophilic drug molecules, protect their cargo from degradation by plasma enzymes, and transport their load across biological membranes and the blood brain barrier (BBB) (see, e.g., Spuch and Navarro, Journal of Drug Delivery, vol. 2011, Article ID 469679, 12 pages, 2011. doi:10.1155/2011/469679 for review).

[0653] Liposomes can be made from several different types of lipids; however, phospholipids are most commonly used to generate liposomes as drug carriers. Although liposome formation is spontaneous when a lipid film is mixed with an aqueous solution, it can also be expedited by applying force in the form of shaking by using a homogenizer, sonicator, or an extrusion apparatus (see, e.g., Spuch and Navarro, Journal of Drug Delivery, vol. 2011, Article ID 469679, 12 pages, 2011. doi:10.1155/2011/469679 for review).

[0654] Several other additives may be added to liposomes in order to modify their structure and properties. For instance, either cholesterol or sphingomyelin may be added to the liposomal mixture in order to help stabilize the liposomal structure and to prevent the leakage of the liposomal inner cargo. Further, liposomes are prepared from hydrogenated egg phosphatidylcholine or egg phosphatidylcholine, cholesterol, and dicetyl phosphate, and their mean vesicle sizes were adjusted to about 50 and 100 nm. (see, e.g., Spuch and Navarro, Journal of Drug Delivery, vol. 2011, Article ID 469679, 12 pages, 2011. doi:10.1155/2011/469679 for review).

[0655] A liposome formulation may be mainly comprised of natural phospholipids and lipids such as 1,2-distearoyl-sn-glycero-3-phosphatidyl choline (DSPC), sphingomyelin, egg phosphatidylcholines and monosialoganglioside. Since this formulation is made up of phospholipids only, liposomal formulations have encountered many challenges, one of the ones being the instability in plasma. Several attempts to overcome these challenges have been made, specifically in the manipulation of the lipid membrane. One of these attempts focused on the manipulation of cholesterol. Addition of cholesterol to conventional formulations reduces rapid release of the encapsulated bioactive compound into the plasma or 1,2-dioleoyl-sn-glycero-3-phosphoethanolamine (DOPE) increases the stability (see, e.g., Spuch and Navarro, Journal of Drug Delivery, vol. 2011, Article ID 469679, 12 pages, 2011. doi:10.1155/2011/469679 for review).

[0656] In a particularly advantageous embodiment, Trojan Horse liposomes (also known as Molecular Trojan Horses) are desirable and protocols may be found at <http://cshprotocols.cshlp.org/content/2010/4/pdb.prot5407.long>. These particles allow delivery of a transgene to the entire brain after an intravascular injection. Without being bound by limitation, it is believed that neutral lipid particles with specific antibodies conjugated to surface allow crossing of the blood brain barrier via endocytosis. Applicant postulates utilizing Trojan Horse Liposomes to deliver the CRISPR family of nucleases to the brain via an intravascular injection, which would allow whole brain transgenic animals without the need for embryonic manipulation. About 1-5 g of DNA or RNA may be contemplated for in vivo administration in liposomes.

[0657] In another embodiment, the nucleic acid-targeting system or components thereof may be administered in liposomes, such as a stable nucleic-acid-lipid particle (SNALP) (see, e.g., Morrissey et al., Nature Biotechnology, Vol. 23, No. 8, August 2005). Daily intravenous injections of about 1, 3 or 5 mg/kg/day of a specific nucleic acid-targeting system targeted in a SNALP are contemplated. The daily treatment may be over about three days and then weekly for about five weeks. In another embodiment, a specific nucleic acid-targeting system encapsulated SNALP administered by intravenous injection to at doses of about 1 or 2.5 mg/kg are also contemplated (see, e.g., Zimmerman et al., Nature Letters, Vol. 441, 4 May 2006). The SNALP formulation may contain the lipids 3-N-[(w-methoxypoly(ethylene glycol) 2000) carbamoyl]-1,2-dimyristyloxy-propylamine (PEG-C-DMA), 1,2-dilinoleyloxy-N,N-dimethyl-3-aminopropane (DLinDMA), 1,2-distearoyl-sn-glycero-3-phosphocholine (DSPC) and cholesterol, in a 2:40:10:48 molar percent ratio (see, e.g., Zimmerman et al., Nature Letters, Vol. 441, 4 May 2006).

[0658] In another embodiment, stable nucleic-acid-lipid particles (SNALPs) have proven to be effective delivery molecules to highly vascularized HepG2-derived liver tumors but not in poorly vascularized HCT-116 derived liver tumors (see, e.g., Li, *Gene Therapy* (2012) 19, 775-780). The SNALP liposomes may be prepared by formulating D-Lin-DMA and PEG-C-DMA with distearoylphosphatidylcholine (DSPC), Cholesterol and siRNA using a 25:1 lipid/siRNA ratio and a 48/40/10/2 molar ratio of Cholesterol/D-Lin-DMA/DSPC/PEG-C-DMA. The resulted SNALP liposomes are about 80-100 nm in size.

[0659] In yet another embodiment, a SNALP may comprise synthetic cholesterol (Sigma-Aldrich, St Louis, Mo., USA), dipalmitoylphosphatidylcholine (Avanti Polar Lipids, Alabaster, Ala., USA), 3-N-[(w-methoxy poly(ethylene glycol)2000)carbamoyl]-1,2-dimyrestyloxypropylamine, and cationic 1,2-dilinoleoyloxy-3-N,N-dimethylaminopropane (see, e.g., Geisbert et al., *Lancet* 2010; 375: 1896-905). A dosage of about 2 mg/kg total nucleic acid-targeting system dose administered as, for example, a bolus intravenous infusion may be contemplated.

[0660] In yet another embodiment, a SNALP may comprise synthetic cholesterol (Sigma-Aldrich), 1,2-distearoyl-sn-glycero-3-phosphocholine (DSPC; Avanti Polar Lipids Inc.), PEG-cDMA, and 1,2-dilinoleoyloxy-3-(N,N-dimethyl)aminopropane (DLinDMA) (see, e.g., Judge, *J. Clin. Invest.* 119:661-673 (2009)). Formulations used for in vivo studies may comprise a final lipid/RNA mass ratio of about 9:1.

[0661] The safety profile of RNAi nanomedicines has been reviewed by Barros and Gollob of Alnylam Pharmaceuticals (see, e.g., *Advanced Drug Delivery Reviews* 64 (2012) 1730-1737). The stable nucleic acid lipid particle (SNALP) is comprised of four different lipids—an ionizable lipid (DLinDMA) that is cationic at low pH, a neutral helper lipid, cholesterol, and a diffusible polyethylene glycol (PEG)-lipid. The particle is approximately 80 nm in diameter and is charge-neutral at physiologic pH. During formulation, the ionizable lipid serves to condense lipid with the anionic RNA during particle formation. When positively charged under increasingly acidic endosomal conditions, the ionizable lipid also mediates the fusion of SNALP with the endosomal membrane enabling release of RNA into the cytoplasm. The PEG-lipid stabilizes the particle and reduces aggregation during formulation, and subsequently provides a neutral hydrophilic exterior that improves pharmacokinetic properties.

[0662] To date, two clinical programs have been initiated using SNALP formulations with RNA. Tekmira Pharmaceuticals recently completed a phase I single-dose study of SNALP-ApoB in adult volunteers with elevated LDL cholesterol. ApoB is predominantly expressed in the liver and jejunum and is essential for the assembly and secretion of VLDL and LDL. Seventeen subjects received a single dose of SNALP-ApoB (dose escalation across 7 dose levels). There was no evidence of liver toxicity (anticipated as the potential dose-limiting toxicity based on preclinical studies). One (of two) subjects at the highest dose experienced flu-like symptoms consistent with immune system stimulation, and the decision was made to conclude the trial.

[0663] Alnylam Pharmaceuticals has similarly advanced ALN-TTR01, which employs the SNALP technology described above and targets hepatocyte production of both mutant and wild-type TTR to treat TTR amyloidosis (ATTR). Three ATTR syndromes have been described:

familial amyloidotic polyneuropathy (FAP) and familial amyloidotic cardiomyopathy (FAC)—both caused by autosomal dominant mutations in TTR; and senile systemic amyloidosis (SSA) cause by wildtype TTR. A placebo-controlled, single dose-escalation phase I trial of ALN-TTR01 was recently completed in patients with ATTR. ALN-TTR01 was administered as a 15-minute IV infusion to 31 patients (23 with study drug and 8 with placebo) within a dose range of 0.01 to 1.0 mg/kg (based on siRNA). Treatment was well tolerated with no significant increases in liver function tests. Infusion-related reactions were noted in 3 of 23 patients at ≥ 0.4 mg/kg; all responded to slowing of the infusion rate and all continued on study. Minimal and transient elevations of serum cytokines IL-6, IP-10 and IL-1ra were noted in two patients at the highest dose of 1 mg/kg (as anticipated from preclinical and NHP studies). Lowering of serum TTR, the expected pharmacodynamics effect of ALN-TTR01, was observed at 1 mg/kg.

[0664] In yet another embodiment, a SNALP may be made by solubilizing a cationic lipid, DSPC, cholesterol and PEG-lipid e.g., in ethanol, e.g., at a molar ratio of 40:10:40:10, respectively (see, Semple et al., *Nature Nanotechnology*, Volume 28 Number 2 Feb. 2010, pp. 172-177). The lipid mixture was added to an aqueous buffer (50 mM citrate, pH 4) with mixing to a final ethanol and lipid concentration of 30% (vol/vol) and 6.1 mg/ml, respectively, and allowed to equilibrate at 22° C. for 2 min before extrusion. The hydrated lipids were extruded through two stacked 80 nm pore-sized filters (Nuclepore) at 22° C. using a Lipex Extruder (Northern Lipids) until a vesicle diameter of 70-90 nm, as determined by dynamic light scattering analysis, was obtained. This generally required 1-3 passes. The siRNA (solubilized in a 50 mM citrate, pH 4 aqueous solution containing 30% ethanol) was added to the pre-equilibrated (35° C.) vesicles at a rate of ~5 ml/min with mixing. After a final target siRNA/lipid ratio of 0.06 (wt/wt) was reached, the mixture was incubated for a further 30 min at 35° C. to allow vesicle reorganization and encapsulation of the siRNA. The ethanol was then removed and the external buffer replaced with PBS (155 mM NaCl, 3 mM Na₂HPO₄, 1 mM KH₂PO₄, pH 7.5) by either dialysis or tangential flow diafiltration. siRNA were encapsulated in SNALP using a controlled step-wise dilution method process. The lipid constituents of KC2-SNALP were DLin-KC2-DMA (cationic lipid), dipalmitoylphosphatidylcholine (DPPC; Avanti Polar Lipids), synthetic cholesterol (Sigma) and PEG-C-DMA used at a molar ratio of 57.1:7.1:34.3:1.4. Upon formation of the loaded particles, SNALP were dialyzed against PBS and filter sterilized through a 0.2 μ m filter before use. Mean particle sizes were 75-85 nm and 90-95% of the siRNA was encapsulated within the lipid particles. The final siRNA/lipid ratio in formulations used for in vivo testing was ~0.15 (wt/wt). LNP-siRNA systems containing Factor VII siRNA were diluted to the appropriate concentrations in sterile PBS immediately before use and the formulations were administered intravenously through the lateral tail vein in a total volume of 10 ml/kg. This method and these delivery systems may be extrapolated to the nucleic acid-targeting system of the present invention.

[0665] Other Lipids

[0666] Other cationic lipids, such as amino lipid 2,2-dilinoleyl-4-dimethylaminoethyl-[1,3]-dioxolane (DLin-KC2-DMA) may be utilized to encapsulate nucleic acid-targeting system or components thereof or nucleic acid

molecule(s) coding therefor e.g., similar to siRNA (see, e.g., Jayaraman, *Angew. Chem. Int. Ed.* 2012, 51, 8529-8533), and hence may be employed in the practice of the invention. A preformed vesicle with the following lipid composition may be contemplated: amino lipid, distearoylphosphatidylcholine (DSPC), cholesterol and (R)-2,3-bis(octadecyloxy)propyl-1-(methoxy poly(ethylene glycol)2000)propylcarbamate (PEG-lipid) in the molar ratio 40/10/40/10, respectively, and a FVII siRNA/total lipid ratio of approximately 0.05 (w/w). To ensure a narrow particle size distribution in the range of 70-90 nm and a low polydispersity index of 0.11 ± 0.04 ($n=56$), the particles may be extruded up to three times through 80 nm membranes prior to adding the guide RNA. Particles containing the highly potent amino lipid 16 may be used, in which the molar ratio of the four lipid components 16, DSPC, cholesterol and PEG-lipid (50/10/38.5/1.5) which may be further optimized to enhance in vivo activity.

[0667] Michael S D Kormann et al. ("Expression of therapeutic proteins after delivery of chemically modified mRNA in mice: *Nature Biotechnology*, Volume: 29, Pages: 154-157 (2011)) describes the use of lipid envelopes to deliver RNA. Use of lipid envelopes is also preferred in the present invention.

[0668] In another embodiment, lipids may be formulated with the nucleic acid-targeting system of the present invention or component(s) thereof or nucleic acid molecule(s) coding therefor to form lipid particles (LNPs). Lipids include, but are not limited to, DLin-KC2-DMA4, C12-200 and colipids distearylphosphatidyl choline, cholesterol, and PEG-DMG may be formulated with RNA-targeting system instead of siRNA (see, e.g., Novobrantseva, *Molecular Therapy-Nucleic Acids* (2012) 1, e4; doi:10.1038/mtna.2011.3) using a spontaneous vesicle formation procedure. The component molar ratio may be about 50/10/38.5/1.5 (DLin-KC2-DMA or C12-200/distearylphosphatidyl choline/cholesterol/PEG-DMG). The final lipid:siRNA weight ratio may be ~12:1 and 9:1 in the case of DLin-KC2-DMA and C12-200 lipid particles (LNPs), respectively. The formulations may have mean particle diameters of ~80 nm with >90% entrapment efficiency. A 3 mg/kg dose may be contemplated.

[0669] Tekmira has a portfolio of approximately 95 patent families, in the U.S. and abroad, that are directed to various aspects of LNPs and LNP formulations (see, e.g., U.S. Pat. Nos. 7,982,027; 7,799,565; 8,058,069; 8,283,333; 7,901,708; 7,745,651; 7,803,397; 8,101,741; 8,188,263; 7,915,399; 8,236,943 and 7,838,658 and European Pat. Nos 1766035; 1519714; 1781593 and 1664316), all of which may be used and/or adapted to the present invention.

[0670] The nucleic acid-targeting system or components thereof or nucleic acid molecule(s) coding therefor may be delivered encapsulated in PLGA Microspheres such as that further described in US published applications 20130252281 and 20130245107 and 20130244279 (assigned to Moderna Therapeutics) which relate to aspects of formulation of compositions comprising modified nucleic acid molecules which may encode a protein, a protein precursor, or a partially or fully processed form of the protein or a protein precursor. The formulation may have a molar ratio 50:10:38.5:1.5-3.0 (cationic lipid:fusogenic lipid:cholesterol:PEG lipid). The PEG lipid may be selected from, but is not limited to PEG-c-DOMG, PEG-DMG. The fusogenic lipid may be DSPC. See also, Schrum et al.,

Delivery and Formulation of Engineered Nucleic Acids, US published application 20120251618.

[0671] Nanomerics' technology addresses bioavailability challenges for a broad range of therapeutics, including low molecular weight hydrophobic drugs, peptides, and nucleic acid based therapeutics (plasmid, siRNA, miRNA). Specific administration routes for which the technology has demonstrated clear advantages include the oral route, transport across the blood-brain-barrier, delivery to solid tumours, as well as to the eye. See, e.g., Mazza et al., 2013, *ACS Nano*. 2013 Feb. 26; 7(2):1016-26; Uchegbu and Siew, 2013, *J Pharm Sci.* 102(2):305-10 and Lalatsa et al., 2012, *J Control Release*. 2012 Jul. 20; 161(2):523-36.

[0672] US Patent Publication No. 20050019923 describes cationic dendrimers for delivering bioactive molecules, such as polynucleotide molecules, peptides and polypeptides and/or pharmaceutical agents, to a mammalian body. The dendrimers are suitable for targeting the delivery of the bioactive molecules to, for example, the liver, spleen, lung, kidney or heart (or even the brain). Dendrimers are synthetic 3-dimensional macromolecules that are prepared in a step-wise fashion from simple branched monomer units, the nature and functionality of which can be easily controlled and varied. Dendrimers are synthesized from the repeated addition of building blocks to a multifunctional core (divergent approach to synthesis), or towards a multifunctional core (convergent approach to synthesis) and each addition of a 3-dimensional shell of building blocks leads to the formation of a higher generation of the dendrimers. Polypropyleneimine dendrimers start from a diaminobutane core to which is added twice the number of amino groups by a double Michael addition of acrylonitrile to the primary amines followed by the hydrogenation of the nitriles. This results in a doubling of the amino groups. Polypropyleneimine dendrimers contain 100% protonable nitrogens and up to 64 terminal amino groups (generation 5, DAB 64). Protonable groups are usually amine groups which are able to accept protons at neutral pH. The use of dendrimers as gene delivery agents has largely focused on the use of the polyamidoamine. and phosphorous containing compounds with a mixture of amine/amide or N—P(O₂)S as the conjugating units respectively with no work being reported on the use of the lower generation polypropyleneimine dendrimers for gene delivery. Polypropyleneimine dendrimers have also been studied as pH sensitive controlled release systems for drug delivery and for their encapsulation of guest molecules when chemically modified by peripheral amino acid groups. The cytotoxicity and interaction of polypropyleneimine dendrimers with DNA as well as the transfection efficacy of DAB 64 has also been studied.

[0673] US Patent Publication No. 20050019923 is based upon the observation that, contrary to earlier reports, cationic dendrimers, such as polypropyleneimine dendrimers, display suitable properties, such as specific targeting and low toxicity, for use in the targeted delivery of bioactive molecules, such as genetic material. In addition, derivatives of the cationic dendrimer also display suitable properties for the targeted delivery of bioactive molecules. See also, Bioactive Polymers, US published application 20080267903, which discloses "Various polymers, including cationic polyamine polymers and dendrimeric polymers, are shown to possess anti-proliferative activity, and may therefore be useful for treatment of disorders characterised by undesirable cellular proliferation such as neoplasms and tumours,

inflammatory disorders (including autoimmune disorders), psoriasis and atherosclerosis. The polymers may be used alone as active agents, or as delivery vehicles for other therapeutic agents, such as drug molecules or nucleic acids for gene therapy. In such cases, the polymers' own intrinsic anti-tumour activity may complement the activity of the agent to be delivered." The disclosures of these patent publications may be employed in conjunction with herein teachings for delivery of nucleic acid-targeting system(s) or component(s) thereof or nucleic acid molecule(s) coding therefor.

[0674] Supercharged Proteins

[0675] Supercharged proteins are a class of engineered or naturally occurring proteins with unusually high positive or negative net theoretical charge and may be employed in delivery of nucleic acid-targeting system(s) or component(s) thereof or nucleic acid molecule(s) coding therefor. Both supernegatively and superpositively charged proteins exhibit a remarkable ability to withstand thermally or chemically induced aggregation. Superpositively charged proteins are also able to penetrate mammalian cells. Associating cargo with these proteins, such as plasmid DNA, RNA, or other proteins, can enable the functional delivery of these macromolecules into mammalian cells both in vitro and in vivo. David Liu's lab reported the creation and characterization of supercharged proteins in 2007 (Lawrence et al., 2007, *Journal of the American Chemical Society* 129, 10110-10112).

[0676] The nonviral delivery of RNA and plasmid DNA into mammalian cells are valuable both for research and therapeutic applications (Akinc et al., 2010, *Nat. Biotech.* 26, 561-569). Purified +36 GFP protein (or other superpositively charged protein) is mixed with RNAs in the appropriate serum-free media and allowed to complex prior addition to cells. Inclusion of serum at this stage inhibits formation of the supercharged protein-RNA complexes and reduces the effectiveness of the treatment. The following protocol has been found to be effective for a variety of cell lines (McNaughton et al., 2009, *Proc. Natl. Acad. Sci. USA* 106, 6111-6116). However, pilot experiments varying the dose of protein and RNA should be performed to optimize the procedure for specific cell lines.

[0677] (1) One day before treatment, plate 1×10^5 cells per well in a 48-well plate.

[0678] (2) On the day of treatment, dilute purified +36 GFP protein in serumfree media to a final concentration 200 nM. Add RNA to a final concentration of 50 nM. Vortex to mix and incubate at room temperature for 10 min.

[0679] (3) During incubation, aspirate media from cells and wash once with PBS.

[0680] (4) Following incubation of +36 GFP and RNA, add the protein-RNA complexes to cells.

[0681] (5) Incubate cells with complexes at 37° C. for 4 h.

[0682] (6) Following incubation, aspirate the media and wash three times with 20 U/mL heparin PBS. Incubate cells with serum-containing media for a further 48 h or longer depending upon the assay for activity.

[0683] (7) Analyze cells by immunoblot, qPCR, phenotypic assay, or other appropriate method.

[0684] David Liu's lab has further found +36 GFP to be an effective plasmid delivery reagent in a range of cells. As plasmid DNA is a larger cargo than siRNA, proportionately more +36 GFP protein is required to effectively complex plasmids. For effective plasmid delivery Applicants have

developed a variant of +36 GFP bearing a C-terminal HA2 peptide tag, a known endosome-disrupting peptide derived from the influenza virus hemagglutinin protein. The following protocol has been effective in a variety of cells, but as above it is advised that plasmid DNA and supercharged protein doses be optimized for specific cell lines and delivery applications.

[0685] (1) One day before treatment, plate 1×10^5 per well in a 48-well plate.

[0686] (2) On the day of treatment, dilute purified β 36 GFP protein in serumfree media to a final concentration 2 mM. Add 1 mg of plasmid DNA. Vortex to mix and incubate at room temperature for 10 min.

[0687] (3) During incubation, aspirate media from cells and wash once with PBS.

[0688] (4) Following incubation of β 36 GFP and plasmid DNA, gently add the protein-DNA complexes to cells.

[0689] (5) Incubate cells with complexes at 37 C for 4 h.

[0690] (6) Following incubation, aspirate the media and wash with PBS. Incubate cells in serum-containing media and incubate for a further 24-48 h.

[0691] (7) Analyze plasmid delivery (e.g., by plasmid-driven gene expression) as appropriate.

[0692] See also, e.g., McNaughton et al., *Proc. Natl. Acad. Sci. USA* 106, 6111-6116 (2009); Cronican et al., *ACS Chemical Biology* 5, 747-752 (2010); Cronican et al., *Chemistry & Biology* 18, 833-838 (2011); Thompson et al., *Methods in Enzymology* 503, 293-319 (2012); Thompson, D. B., et al., *Chemistry & Biology* 19 (7), 831-843 (2012). The methods of the super charged proteins may be used and/or adapted for delivery of the nucleic acid-targeting system of the present invention. These systems of Dr. Lui and documents herein in conjunction with herein teachings can be employed in the delivery of nucleic acid-targeting system(s) or component(s) thereof or nucleic acid molecule (s) coding therefor.

Cell Penetrating Peptides (CPPs)

[0693] In yet another embodiment, cell penetrating peptides (CPPs) are contemplated for the delivery of the CRISPR Cas system. CPPs are short peptides that facilitate cellular uptake of various molecular cargo (from nanosize particles to small chemical molecules and large fragments of DNA). The term "cargo" as used herein includes but is not limited to the group consisting of therapeutic agents, diagnostic probes, peptides, nucleic acids, antisense oligonucleotides, plasmids, proteins, particles including particles, liposomes, chromophores, small molecules and radioactive materials. In aspects of the invention, the cargo may also comprise any component of the CRISPR Cas system or the entire functional CRISPR Cas system. Aspects of the present invention further provide methods for delivering a desired cargo into a subject comprising: (a) preparing a complex comprising the cell penetrating peptide of the present invention and a desired cargo, and (b) orally, intraarticularly, intraperitoneally, intrathecally, intrarterially, intranasally, intraparenchymally, subcutaneously, intramuscularly, intravenously, dermally, intrarectally, or topically administering the complex to a subject. The cargo is associated with the peptides either through chemical linkage via covalent bonds or through non-covalent interactions.

[0694] The function of the CPPs are to deliver the cargo into cells, a process that commonly occurs through endocytosis with the cargo delivered to the endosomes of living

mammalian cells. Cell-penetrating peptides are of different sizes, amino acid sequences, and charges but all CPPs have one distinct characteristic, which is the ability to translocate the plasma membrane and facilitate the delivery of various molecular cargoes to the cytoplasm or an organelle. CPP translocation may be classified into three main entry mechanisms: direct penetration in the membrane, endocytosis-mediated entry, and translocation through the formation of a transitory structure. CPPs have found numerous applications in medicine as drug delivery agents in the treatment of different diseases including cancer and virus inhibitors, as well as contrast agents for cell labeling. Examples of the latter include acting as a carrier for GFP, Mill contrast agents, or quantum dots. CPPs hold great potential as in vitro and in vivo delivery vectors for use in research and medicine. CPPs typically have an amino acid composition that either contains a high relative abundance of positively charged amino acids such as lysine or arginine or has sequences that contain an alternating pattern of polar/charged amino acids and non-polar, hydrophobic amino acids. These two types of structures are referred to as polycationic or amphiphathic, respectively. A third class of CPPs are the hydrophobic peptides, containing only apolar residues, with low net charge or have hydrophobic amino acid groups that are crucial for cellular uptake. One of the initial CPPs discovered was the trans-activating transcriptional activator (Tat) from Human Immunodeficiency Virus 1 (HIV-1) which was found to be efficiently taken up from the surrounding media by numerous cell types in culture. Since then, the number of known CPPs has expanded considerably and small molecule synthetic analogues with more effective protein transduction properties have been generated. CPPs include but are not limited to Penetratin, Tat (48-60), Transportan, and (R-Ahx-R4) (Ahx=aminohexanoyl).

[0695] U.S. Pat. No. 8,372,951, provides a CPP derived from eosinophil cationic protein (ECP) which exhibits highly cell-penetrating efficiency and low toxicity. Aspects of delivering the CPP with its cargo into a vertebrate subject are also provided. Further aspects of CPPs and their delivery are described in U.S. Pat. Nos. 8,575,305; 8,614,194 and 8,044,019. CPPs can be used to deliver the CRISPR-Cas system or components thereof. That CPPs can be employed to deliver the CRISPR-Cas system or components thereof is also provided in the manuscript "Gene disruption by cell-penetrating peptide-mediated delivery of Cas9 protein and guide RNA", by Suresh Ramakrishna, Abu-Bonsrah Kwaku Dad, Jagadish Beloor, et al. *Genome Res.* 2014 Apr. 2. [Epub ahead of print], incorporated by reference in its entirety, wherein it is demonstrated that treatment with CPP-conjugated recombinant Cas9 protein and CPP-complexed guide RNAs lead to endogenous gene disruptions in human cell lines. In the paper the Cas9 protein was conjugated to CPP via a thioether bond, whereas the guide RNA was complexed with CPP, forming condensed, positively charged particles. It was shown that simultaneous and sequential treatment of human cells, including embryonic stem cells, dermal fibroblasts, HEK293T cells, HeLa cells, and embryonic carcinoma cells, with the modified Cas9 and guide RNA led to efficient gene disruptions with reduced off-target mutations relative to plasmid transfections.

Implantable Devices

[0696] In another embodiment, implantable devices are also contemplated for delivery of the nucleic acid-targeting

system or component(s) thereof or nucleic acid molecule(s) coding therefor. For example, US Patent Publication 20110195123 discloses an implantable medical device which elutes a drug locally and in prolonged period is provided, including several types of such a device, the treatment modes of implementation and methods of implantation. The device comprising of polymeric substrate, such as a matrix for example, that is used as the device body, and drugs, and in some cases additional scaffolding materials, such as metals or additional polymers, and materials to enhance visibility and imaging. An implantable delivery device can be advantageous in providing release locally and over a prolonged period, where drug is released directly to the extracellular matrix (ECM) of the diseased area such as tumor, inflammation, degeneration or for symptomatic objectives, or to injured smooth muscle cells, or for prevention. One kind of drug is RNA, as disclosed above, and this system may be used/and or adapted to the nucleic acid-targeting system of the present invention. The modes of implantation in some embodiments are existing implantation procedures that are developed and used today for other treatments, including brachytherapy and needle biopsy. In such cases the dimensions of the new implant described in this invention are similar to the original implant. Typically a few devices are implanted during the same treatment procedure.

[0697] US Patent Publication 20110195123, provides a drug delivery implantable or insertable system, including systems applicable to a cavity such as the abdominal cavity and/or any other type of administration in which the drug delivery system is not anchored or attached, comprising a biostable and/or degradable and/or bioabsorbable polymeric substrate, which may for example optionally be a matrix. It should be noted that the term "insertion" also includes implantation. The drug delivery system is preferably implemented as a "Loder" as described in US Patent Publication 20110195123.

[0698] The polymer or plurality of polymers are biocompatible, incorporating an agent and/or plurality of agents, enabling the release of agent at a controlled rate, wherein the total volume of the polymeric substrate, such as a matrix for example, in some embodiments is optionally and preferably no greater than a maximum volume that permits a therapeutic level of the agent to be reached. As a non-limiting example, such a volume is preferably within the range of 0.1 m³ to 1000 mm³, as required by the volume for the agent load. The Loder may optionally be larger, for example when incorporated with a device whose size is determined by functionality, for example and without limitation, a knee joint, an intra-uterine or cervical ring and the like.

[0699] The drug delivery system (for delivering the composition) is designed in some embodiments to preferably employ degradable polymers, wherein the main release mechanism is bulk erosion; or in some embodiments, non degradable, or slowly degraded polymers are used, wherein the main release mechanism is diffusion rather than bulk erosion, so that the outer part functions as membrane, and its internal part functions as a drug reservoir, which practically is not affected by the surroundings for an extended period (for example from about a week to about a few months). Combinations of different polymers with different release mechanisms may also optionally be used. The concentration gradient at the surface is preferably maintained effectively constant during a significant period of the total drug releas-

ing period, and therefore the diffusion rate is effectively constant (termed “zero mode” diffusion). By the term “constant” it is meant a diffusion rate that is preferably maintained above the lower threshold of therapeutic effectiveness, but which may still optionally feature an initial burst and/or may fluctuate, for example increasing and decreasing to a certain degree. The diffusion rate is preferably so maintained for a prolonged period, and it can be considered constant to a certain level to optimize the therapeutically effective period, for example the effective silencing period.

[0700] The drug delivery system optionally and preferably is designed to shield the nucleotide based therapeutic agent from degradation, whether chemical in nature or due to attack from enzymes and other factors in the body of the subject.

[0701] The drug delivery system of US Patent Publication 20110195123 is optionally associated with sensing and/or activation appliances that are operated at and/or after implantation of the device, by non and/or minimally invasive methods of activation and/or acceleration/deceleration, for example optionally including but not limited to thermal heating and cooling, laser beams, and ultrasonic, including focused ultrasound and/or RF (radiofrequency) methods or devices.

[0702] According to some embodiments of US Patent Publication 20110195123, the site for local delivery may optionally include target sites characterized by high abnormal proliferation of cells, and suppressed apoptosis, including tumors, active and/or chronic inflammation and infection including autoimmune diseases states, degenerating tissue including muscle and nervous tissue, chronic pain, degenerative sites, and location of bone fractures and other wound locations for enhancement of regeneration of tissue, and injured cardiac, smooth and striated muscle.

[0703] The site for implantation of the composition, or target site, preferably features a radius, area and/or volume that is sufficiently small for targeted local delivery. For example, the target site optionally has a diameter in a range of from about 0.1 mm to about 5 cm.

[0704] The location of the target site is preferably selected for maximum therapeutic efficacy. For example, the composition of the drug delivery system (optionally with a device for implantation as described above) is optionally and preferably implanted within or in the proximity of a tumor environment, or the blood supply associated thereof.

[0705] For example the composition (optionally with the device) is optionally implanted within or in the proximity to pancreas, prostate, breast, liver, via the nipple, within the vascular system and so forth.

[0706] The target location is optionally selected from the group comprising, consisting essentially of, or consisting of (as non-limiting examples only, as optionally any site within the body may be suitable for implanting a Loder): 1. brain at degenerative sites like in Parkinson or Alzheimer disease at the basal ganglia, white and gray matter; 2. spine as in the case of amyotrophic lateral sclerosis (ALS); 3. uterine cervix to prevent HPV infection; 4. active and chronic inflammatory joints; 5. dermis as in the case of psoriasis; 6. sympathetic and sensoric nervous sites for analgesic effect; 7. Intra osseous implantation; 8. acute and chronic infection sites; 9. Intra vaginal; 10. Inner ear-auditory system, labyrinth of the inner ear, vestibular system; 11. Intra tracheal; 12. Intra-cardiac; coronary, epicardiac; 13. urinary bladder; 14. biliary system; 15. parenchymal tissue including and not limited to

the kidney, liver, spleen; 16. lymph nodes; 17. salivary glands; 18. dental gums; 19. Intra-articular (into joints); 20. Intra-ocular; 21. Brain tissue; 22. Brain ventricles; 23. Cavities, including abdominal cavity (for example but without limitation, for ovary cancer); 24. Intra esophageal and 25. Intra rectal.

[0707] Optionally insertion of the system (for example a device containing the composition) is associated with injection of material to the ECM at the target site and the vicinity of that site to affect local pH and/or temperature and/or other biological factors affecting the diffusion of the drug and/or drug kinetics in the ECM, of the target site and the vicinity of such a site.

[0708] Optionally, according to some embodiments, the release of said agent could be associated with sensing and/or activation appliances that are operated prior and/or at and/or after insertion, by non and/or minimally invasive and/or else methods of activation and/or acceleration/deceleration, including laser beam, radiation, thermal heating and cooling, and ultrasonic, including focused ultrasound and/or RF (radiofrequency) methods or devices, and chemical activators.

[0709] According to other embodiments of US Patent Publication 20110195123, the drug preferably comprises a RNA, for example for localized cancer cases in breast, pancreas, brain, kidney, bladder, lung, and prostate as described below. Although exemplified with RNAi, many drugs are applicable to be encapsulated in Loder, and can be used in association with this invention, as long as such drugs can be encapsulated with the Loder substrate, such as a matrix for example, and this system may be used and/or adapted to deliver the nucleic acid-targeting system of the present invention.

[0710] As another example of a specific application, neuro and muscular degenerative diseases develop due to abnormal gene expression. Local delivery of RNAs may have therapeutic properties for interfering with such abnormal gene expression. Local delivery of anti apoptotic, anti inflammatory and anti degenerative drugs including small drugs and macromolecules may also optionally be therapeutic. In such cases the Loder is applied for prolonged release at constant rate and/or through a dedicated device that is implanted separately. All of this may be used and/or adapted to the nucleic acid-targeting system of the present invention.

[0711] As yet another example of a specific application, psychiatric and cognitive disorders are treated with gene modifiers. Gene knockdown is a treatment option. Loders locally delivering agents to central nervous system sites are therapeutic options for psychiatric and cognitive disorders including but not limited to psychosis, bi-polar diseases, neurotic disorders and behavioral maladies. The Loders could also deliver locally drugs including small drugs and macromolecules upon implantation at specific brain sites. All of this may be used and/or adapted to the nucleic acid-targeting system of the present invention.

[0712] As another example of a specific application, silencing of innate and/or adaptive immune mediators at local sites enables the prevention of organ transplant rejection. Local delivery of RNAs and immunomodulating reagents with the Loder implanted into the transplanted organ and/or the implanted site renders local immune suppression by repelling immune cells such as CD8 activated

against the transplanted organ. All of this may be used/and or adapted to the nucleic acid-targeting system of the present invention.

[0713] As another example of a specific application, vascular growth factors including VEGFs and angiogenin and others are essential for neovascularization. Local delivery of the factors, peptides, peptidomimetics, or suppressing their repressors is an important therapeutic modality; silencing the repressors and local delivery of the factors, peptides, macromolecules and small drugs stimulating angiogenesis with the Loder is therapeutic for peripheral, systemic and cardiac vascular disease.

[0714] The method of insertion, such as implantation, may optionally already be used for other types of tissue implantation and/or for insertions and/or for sampling tissues, optionally without modifications, or alternatively optionally only with non-major modifications in such methods. Such methods optionally include but are not limited to brachytherapy methods, biopsy, endoscopy with and/or without ultrasound, such as ERCP, stereotactic methods into the brain tissue, Laparoscopy, including implantation with a laparoscope into joints, abdominal organs, the bladder wall and body cavities.

[0715] Implantable device technology herein discussed can be employed with herein teachings and hence by this disclosure and the knowledge in the art, CRISPR-Cas system or components thereof or nucleic acid molecules thereof or encoding or providing components may be delivered via an implantable device.

Patient-Specific Screening Methods

[0716] A nucleic acid-targeting system that targets RNA, e.g., trinucleotide repeats can be used to screen patients or patent samples for the presence of such repeats. The repeats can be the target of the RNA of the nucleic acid-targeting system, and if there is binding thereto by the nucleic acid-targeting system, that binding can be detected, to thereby indicate that such a repeat is present. Thus, a nucleic acid-targeting system can be used to screen patients or patient samples for the presence of the repeat. The patient can then be administered suitable compound(s) to address the condition; or, can be administered a nucleic acid-targeting system to bind to and cause insertion, deletion or mutation and alleviate the condition.

[0717] The invention uses nucleic acids to bind target RNA sequences.

CRISPR Effector Protein mRNA and Guide RNA

[0718] CRISPR effector protein mRNA and guide RNA might also be delivered separately. CRISPR effector protein mRNA can be delivered prior to the guide RNA to give time for CRISPR effector protein to be expressed. CRISPR effector protein mRNA might be administered 1-12 hours (preferably around 2-6 hours) prior to the administration of guide RNA.

[0719] Alternatively, CRISPR effector protein mRNA and guide RNA can be administered together. Advantageously, a second booster dose of guide RNA can be administered 1-12 hours (preferably around 2-6 hours) after the initial administration of CRISPR effector protein mRNA+guide RNA.

[0720] The CRISPR effector protein of the present invention, i.e. a Group 29 or Group 30 effector protein is sometimes referred to herein as a CRISPR Enzyme. It will be appreciated that the effector protein is based on or derived from an enzyme, so the term 'effector protein' certainly

includes 'enzyme' in some embodiments. However, it will also be appreciated that the effector protein may, as required in some embodiments, have DNA or RNA binding, but not necessarily cutting or nicking, activity, including a dead-Cas effector protein function.

[0721] In some embodiments diseases that may be targeted include those concerned with disease-causing splice defects.

[0722] In some embodiments, cellular targets include Hemopoietic Stem/Progenitor Cells (CD34+); Human T cells; and Eye (retinal cells)—for example photoreceptor precursor cells.

[0723] In some embodiments Gene targets include: Human Beta Globin—HBB (for treating Sickle Cell Anemia, including by stimulating gene-conversion (using closely related HBD gene as an endogenous template)); CD3 (T-Cells); and CEP920—retina (eye).

[0724] In some embodiments disease targets also include: cancer; Sickle Cell Anemia (based on a point mutation); HIV; Beta-Thalassemia; and ophthalmic or ocular disease—for example Leber Congenital Amaurosis (LCA)-causing Splice Defect.

[0725] In some embodiments delivery methods include: Cationic Lipid Mediated "direct" delivery of Enzyme-Guide complex (RiboNucleoProtein) and electroporation of plasmid DNA.

[0726] For minimization of toxicity and off-target effect, it will be important to control the concentration of CRISPR effector protein mRNA and guide RNA delivered. Optimal concentrations of CRISPR effector protein mRNA and guide RNA can be determined by testing different concentrations in a cellular or animal model and using deep sequencing to analyze the extent of modification at potential off-target genomic loci. For example, for the guide sequence targeting the sequence 5'-GAGTCCGAGCAGAAGAAGAA-3' (artificial) in the EMX1 gene of the human genome, deep sequencing can be used to assess the level of modification at the following two off-target loci, 1: 5'-GAGTCTAGCAGGAGAAGAA-3' and 2: 5'-GAGTCTAAGCAGAAGAAGAA-3' (both artificial sequences). The concentration that gives the highest level of on-target modification while minimizing the level of off-target modification should be chosen for in vivo delivery.

Inducible Systems

[0727] In some embodiments, a CRISPR effector protein may form a component of an inducible system. The inducible nature of the system would allow for spatiotemporal control of gene editing or gene expression using a form of energy. The form of energy may include but is not limited to electromagnetic radiation, sound energy, chemical energy and thermal energy. Examples of inducible system include tetracycline inducible promoters (Tet-On or Tet-Off), small molecule two-hybrid transcription activations systems (FKBP, ABA, etc), or light inducible systems (Phytochrome, LOV domains, or cryptochrome). In one embodiment, the CRISPR effector protein may be a part of a Light Inducible Transcriptional Effector (LITE) to direct changes in transcriptional activity in a sequence-specific manner. The components of a light may include a CRISPR effector protein, a light-responsive cytochrome heterodimer (e.g. from *Arabidopsis thaliana*), and a transcriptional activation/repression domain. Further examples of inducible DNA binding proteins and methods for their use are provided in U.S. 61/736,

465 and U.S. 61/721,283, and WO 2014018423 A2 which is hereby incorporated by reference in its entirety.

Self-Inactivating Systems

[0728] Once all copies of a gene in the genome of a cell have been edited, continued CRISPR/ a Group 29 or Group 30 effector protein expression in that cell is no longer necessary. Indeed, sustained expression would be undesirable in case of off-target effects at unintended genomic sites, etc. Thus time-limited expression would be useful. Inducible expression offers one approach, but in addition Applicants have engineered a Self-Inactivating CRISPR system that relies on the use of a non-coding guide target sequence within the CRISPR vector itself. Thus, after expression begins, the CRISPR system will lead to its own destruction, but before destruction is complete it will have time to edit the genomic copies of the target gene (which, with a normal point mutation in a diploid cell, requires at most two edits). Simply, the self inactivating CRISPR-Cas system includes additional RNA (i.e., guide RNA) that targets the coding sequence for the CRISPR effector protein itself or that targets one or more non-coding guide target sequences complementary to unique sequences present in one or more of the following:

- (a) within the promoter driving expression of the non-coding RNA elements,
- (b) within the promoter driving expression of the effector protein gene,
- (c) within 100 bp of the ATG translational start codon in the effector protein coding sequence,
- (d) within the inverted terminal repeat (iTR) of a viral delivery vector, e.g., in the AAV genome.

[0729] Furthermore, that RNA can be delivered via a vector, e.g., a separate vector or the same vector that is encoding the CRISPR complex. When provided by a separate vector, the CRISPR RNA that targets Cas expression can be administered sequentially or simultaneously. When administered sequentially, the CRISPR RNA that targets Cas expression is to be delivered after the CRISPR RNA that is intended for e.g. gene editing or gene engineering. This period may be a period of minutes (e.g. 5 minutes, 10 minutes, 20 minutes, 30 minutes, 45 minutes, 60 minutes). This period may be a period of hours (e.g. 2 hours, 4 hours, 6 hours, 8 hours, 12 hours, 24 hours). This period may be a period of days (e.g. 2 days, 3 days, 4 days, 7 days). This period may be a period of weeks (e.g. 2 weeks, 3 weeks, 4 weeks). This period may be a period of months (e.g. 2 months, 4 months, 8 months, 12 months). This period may be a period of years (2 years, 3 years, 4 years). Where the guide RNA targets the sequences encoding expression of the Cas protein, the effector protein becomes impeded and the system becomes self inactivating. In the same manner, CRISPR RNA that targets Cas expression applied via, for example liposome, lipofection, particles, microvesicles as explained herein, may be administered sequentially or simultaneously. Similarly, self-inactivation may be used for inactivation of one or more guide RNA used to target one or more targets.

[0730] In some aspects, a single gRNA is provided that is capable of hybridization to a sequence downstream of a CRISPR effector protein start codon, whereby after a period of time there is a loss of the CRISPR effector protein expression. In some aspects, one or more gRNA(s) are provided that are capable of hybridization to one or more

coding or non-coding regions of the polynucleotide encoding the CRISPR-Cas system, whereby after a period of time there is a inactivation of one or more, or in some cases all, of the CRISPR-Cas system. In some aspects of the system, and not to be limited by theory, the cell may comprise a plurality of CRISPR-Cas complexes, wherein a first subset of CRISPR complexes comprise a first guide RNA capable of targeting a genomic locus or loci to be edited, and a second subset of CRISPR complexes comprise at least one second guide RNA capable of targeting the polynucleotide encoding the CRISPR-Cas system, wherein the first subset of CRISPR-Cas complexes mediate editing of the targeted genomic locus or loci and the second subset of CRISPR complexes eventually inactivate the CRISPR-Cas system, thereby inactivating further CRISPR-Cas expression in the cell.

[0731] Thus the invention provides a CRISPR-Cas system comprising one or more vectors for delivery to a eukaryotic cell, wherein the vector(s) encode(s): (i) a CRISPR effector protein; (ii) a first guide RNA capable of hybridizing to a target sequence in the cell; (iii) a second guide RNA capable of hybridizing to one or more target sequence(s) in the vector which encodes the CRISPR effector protein; thus differing only by the guide sequence, wherein, when expressed within the cell: the first guide RNA directs sequence-specific binding of a first CRISPR complex to the target sequence in the cell; the second guide RNA directs sequence-specific binding of a second CRISPR complex to the target sequence in the vector which encodes the CRISPR effector protein; the CRISPR complexes comprise a CRISPR effector protein bound to a guide RNA, such that a guide RNA can hybridize to its target sequence; and the second CRISPR complex inactivates the CRISPR-Cas system to prevent continued expression of the CRISPR effector protein by the cell.

[0732] The various coding sequences (CRISPR effector protein and guide RNAs) can be included on a single vector or on multiple vectors. For instance, it is possible to encode the effector protein on one vector and the various RNA sequences on another vector, or to encode the effector protein and one guide RNA on one vector, and the remaining guide RNA on another vector, or any other permutation. In general, a system using a total of one or two different vectors is preferred.

[0733] Where multiple vectors are used, it is possible to deliver them in unequal numbers, and ideally with an excess of a vector which encodes the first guide RNA relative to the second guide RNA, thereby assisting in delaying final inactivation of the CRISPR system until genome editing has had a chance to occur.

[0734] The first guide RNA can target any target sequence of interest within a genome, as described elsewhere herein. The second guide RNA targets a sequence within the vector which encodes the CRISPR Group 29 or Group 30 effector protein, and thereby inactivates the effector protein's expression from that vector. Thus the target sequence in the vector must be capable of inactivating expression. Suitable target sequences can be, for instance, near to or within the translational start codon for the Group 29 or Group 30 effector protein coding sequence, in a non-coding sequence in the promoter driving expression of the non-coding RNA elements, within the promoter driving expression of the Group 29 or Group 30 effector protein gene, within 100 bp of the ATG translational start codon in the Cas coding sequence,

and/or within the inverted terminal repeat (iTR) of a viral delivery vector, e.g., in the AAV genome. A double stranded break near this region can induce a frame shift in the Cas coding sequence, causing a loss of protein expression. An alternative target sequence for the “self-inactivating” guide RNA would aim to edit/inactivate regulatory regions/sequences needed for the expression of the CRISPR-system or for the stability of the vector. For instance, if the promoter for the Cas coding sequence is disrupted then transcription can be inhibited or prevented. Similarly, if a vector includes sequences for replication, maintenance or stability then it is possible to target these. For instance, in a AAV vector a useful target sequence is within the iTR. Other useful sequences to target can be promoter sequences, polyadenylation sites, etc.

[0735] Furthermore, if the guide RNAs are expressed in array format, the “self-inactivating” guide RNAs that target both promoters simultaneously will result in the excision of the intervening nucleotides from within the CRISPR-Cas expression construct, effectively leading to its complete inactivation. Similarly, excision of the intervening nucleotides will result where the guide RNAs target both ITRs, or targets two or more other CRISPR-Cas components simultaneously. Self-inactivation as explained herein is applicable, in general, with CRISPR-Cas systems in order to provide regulation of the CRISPR-Cas. For example, self-inactivation as explained herein may be applied to the CRISPR repair of mutations, for example expansion disorders, as explained herein. As a result of this self-inactivation, CRISPR repair is only transiently active.

[0736] Addition of non-targeting nucleotides to the 5' end (e.g. 1-10 nucleotides, preferably 1-5 nucleotides) of the “self-inactivating” guide RNA can be used to delay its processing and/or modify its efficiency as a means of ensuring editing at the targeted genomic locus prior to CRISPR-Cas shutdown.

[0737] In one aspect of the self-inactivating AAV-CRISPR-Cas system, plasmids that co-express one or more guide RNA targeting genomic sequences of interest (e.g. 1-2, 1-5, 1-10, 1-15, 1-20, 1-30) may be established with “self-inactivating” guide RNAs that target an SpCas9 sequence at or near the engineered ATG start site (e.g. within 5 nucleotides, within 15 nucleotides, within 30 nucleotides, within 50 nucleotides, within 100 nucleotides). A regulatory sequence in the U6 promoter region can also be targeted with an guide RNA. The U6-driven guide RNAs may be designed in an array format such that multiple guide RNA sequences can be simultaneously released. When first delivered into target tissue/cells (left cell) guide RNAs begin to accumulate while Cas levels rise in the nucleus. Cas complexes with all of the guide RNAs to mediate genome editing and self-inactivation of the CRISPR-Cas plasmids.

[0738] One aspect of a self-inactivating CRISPR-Cas system is expression of singly or in tandem array format from 1 up to 4 or more different guide sequences; e.g. up to about 20 or about 30 guides sequences. Each individual self inactivating guide sequence may target a different target. Such may be processed from, e.g. one chimeric pol3 transcript. Pol3 promoters such as U6 or H1 promoters may be used. Pol2 promoters such as those mentioned throughout herein. Inverted terminal repeat (iTR) sequences may flank the Pol3 promoter-guide RNA(s)-Pol2 promoter-Cas.

[0739] One aspect of a tandem array transcript is that one or more guide(s) edit the one or more target(s) while one or

more self inactivating guides inactivate the CRISPR-Cas system. Thus, for example, the described CRISPR-Cas system for repairing expansion disorders may be directly combined with the self-inactivating CRISPR-Cas system described herein. Such a system may, for example, have two guides directed to the target region for repair as well as at least a third guide directed to self-inactivation of the CRISPR-Cas. Reference is made to Application Ser. No. PCT/US2014/069897, entitled “Compositions And Methods Of Use Of Crispr-Cas Systems In Nucleotide Repeat Disorders,” published Dec. 12, 2014 as WO/2015/089351.

Kits

[0740] In one aspect, the invention provides kits containing any one or more of the elements disclosed in the above methods and compositions. In some embodiments, the kit comprises a vector system as taught herein or one or more of the components of the CRISPR/Cas system as taught herein, such as gRNAs and/or effector protein or effector protein encoding mRNA, and instructions for using the kit. Elements may be provide individually or in combinations, and may be provided in any suitable container, such as a vial, a bottle, or a tube. In some embodiments, the kit includes instructions in one or more languages, for example in more than one language. The instructions may be specific to the applications and methods described herein.

[0741] In some embodiments, a kit comprises one or more reagents for use in a process utilizing one or more of the elements described herein. Reagents may be provided in any suitable container. For example, a kit may provide one or more reaction or storage buffers. Reagents may be provided in a form that is usable in a particular assay, or in a form that requires addition of one or more other components before use (e.g., in concentrate or lyophilized form). A buffer can be any buffer, including but not limited to a sodium carbonate buffer, a sodium bicarbonate buffer, a borate buffer, a Tris buffer, a MOPS buffer, a HEPES buffer, and combinations thereof. In some embodiments, the buffer is alkaline. In some embodiments, the buffer has a pH from about 7 to about 10. In some embodiments, the kit comprises one or more oligonucleotides corresponding to a guide sequence for insertion into a vector so as to operably link the guide sequence and a regulatory element. In some embodiments, the kit comprises one or more of the vectors and/or one or more of the polynucleotides described herein. The kit may advantageously allows to provide all elements of the systems of the invention.

[0742] In one aspect, the invention provides methods for using one or more elements of a CRISPR system. The CRISPR complex of the invention provides an effective means for modifying a target polynucleotide. The CRISPR complex of the invention has a wide variety of utility including modifying (e.g., deleting, inserting, translocating, inactivating, activating) a target polynucleotide in a multiplicity of cell types. As such the CRISPR complex of the invention has a broad spectrum of applications in, e.g., gene therapy, drug screening, disease diagnosis, and prognosis. An exemplary CRISPR complex comprises a CRISPR effector protein complexed with a guide sequence hybridized to a target sequence within the target polynucleotide. In certain embodiments, a direct repeat sequence is linked to the guide sequence.

[0743] In one embodiment, this invention provides a method of cleaving a target polynucleotide. The method

comprises modifying a target polynucleotide using a CRISPR complex that binds to the target polynucleotide and effect cleavage of said target polynucleotide. Typically, the CRISPR complex of the invention, when introduced into a cell, creates a break (e.g., a single or a double strand break) in the target sequence. For example, the method can be used to cleave a disease gene or gene product in a cell.

[0744] Where desired, a donor polynucleotide can be DNA, e.g., a DNA plasmid, a bacterial artificial chromosome (BAC), a yeast artificial chromosome (YAC), a viral vector, a linear piece of DNA, a PCR fragment, a naked nucleic acid, or a nucleic acid complexed with a delivery vehicle such as a liposome or poloxamer.

[0745] In other embodiments, this invention provides a method of modifying expression of a polynucleotide in a eukaryotic cell. The method comprises increasing or decreasing expression of a target polynucleotide by using a CRISPR complex that binds to the polynucleotide.

[0746] In some methods, a target polynucleotide can be inactivated to effect the modification of the expression in a cell. For example, upon the binding of a CRISPR complex to a target sequence in a cell, the target polynucleotide is inactivated such that the sequence is not transcribed, the coded protein is not produced, or the sequence does not function as the wild-type sequence does. For example, a protein or microRNA coding sequence may be inactivated such that the protein is not produced.

[0747] In some methods, a control sequence can be inactivated such that it no longer functions as a control sequence. As used herein, "control sequence" refers to any nucleic acid sequence that effects the transcription, translation, or accessibility of a nucleic acid sequence. Examples of a control sequence include, a promoter, a transcription terminator, and an enhancer are control sequences. The inactivated target sequence may include a deletion mutation (i.e., deletion of one or more nucleotides), an insertion mutation (i.e., insertion of one or more nucleotides), or a nonsense mutation (i.e., substitution of a single nucleotide for another nucleotide such that a stop codon is introduced). In some methods, the inactivation of a target sequence results in "knockout" of the target sequence.

Exemplary Methods of Using of CRISPR Cas System

[0748] The invention provides a non-naturally occurring or engineered composition, or one or more polynucleotides encoding components of said composition, or vector or delivery systems comprising one or more polynucleotides encoding components of said composition for use in a modifying a target cell in vivo, ex vivo or in vitro and, may be conducted in a manner alters the cell such that once modified the progeny or cell line of the CRISPR modified cell retains the altered phenotype. The modified cells and progeny may be part of a multi-cellular organism such as a plant or animal with ex vivo or in vivo application of CRISPR system to desired cell types. The CRISPR invention may be a therapeutic method of treatment. The therapeutic method of treatment may comprise gene or genome editing, or gene therapy, and/or gene knockdown.

Modifying a Target with CRISPR Cas System or Complex (e.g., Group 29 or Group 30 Effector Protein-RNA Complex)

[0749] In one aspect, the invention provides for methods of modifying a target polynucleotide in a eukaryotic cell, which may be in vivo, ex vivo or in vitro. In some embodi-

ments, the method comprises sampling a cell or population of cells from a human or non-human animal, and modifying the cell or cells. Culturing may occur at any stage ex vivo. The cell or cells may even be re-introduced into the non-human animal or plant. For re-introduced cells it is particularly preferred that the cells are stem cells.

[0750] In some embodiments, the method comprises allowing a CRISPR complex to bind to the target polynucleotide to effect cleavage of said target polynucleotide thereby modifying the target polynucleotide, wherein the CRISPR complex comprises a CRISPR effector protein complexed with a guide sequence hybridized or hybridizable to a target sequence within said target polynucleotide.

[0751] In one aspect, the invention provides a method of modifying expression of a polynucleotide in a eukaryotic cell. In some embodiments, the method comprises allowing a CRISPR complex to bind to the polynucleotide such that said binding results in increased or decreased expression of said polynucleotide; wherein the CRISPR complex comprises a CRISPR effector protein complexed with a guide sequence hybridized or hybridizable to a target sequence within said polynucleotide. Similar considerations and conditions apply as above for methods of modifying a target polynucleotide. In fact, these sampling, culturing and re-introduction options apply across the aspects of the present invention.

[0752] Indeed, in any aspect of the invention, the CRISPR complex may comprise a CRISPR effector protein complexed with a guide sequence hybridized or hybridizable to a target sequence. Similar considerations and conditions apply as above for methods of modifying a target polynucleotide.

[0753] Thus in any of the non-naturally-occurring CRISPR effector proteins described herein comprise at least one modification and whereby the effector protein has certain improved capabilities. In particular, any of the effector proteins are capable of forming a CRISPR complex with a guide RNA. When such a complex forms, the guide RNA is capable of binding to a target polynucleotide sequence and the effector protein is capable of modifying a target locus. In addition, the effector protein in the CRISPR complex has reduced capability of modifying one or more off-target loci as compared to an unmodified enzyme/effector protein.

[0754] In addition, the modified CRISPR enzymes described herein encompass enzymes whereby in the CRISPR complex the effector protein has increased capability of modifying the one or more target loci as compared to an unmodified enzyme/effector protein. Such function may be provided separate to or provided in combination with the above-described function of reduced capability of modifying one or more off-target loci. Any such effector proteins may be provided with any of the further modifications to the CRISPR effector protein as described herein, such as in combination with any activity provided by one or more associated heterologous functional domains, any further mutations to reduce nuclease activity and the like.

[0755] In advantageous embodiments of the invention, the modified CRISPR effector protein is provided with reduced capability of modifying one or more off-target loci as compared to an unmodified enzyme/effector protein and increased capability of modifying the one or more target loci as compared to an unmodified enzyme/effector protein. In combination with further modifications to the effector protein, significantly enhanced specificity may be achieved. For

example, combination of such advantageous embodiments with one or more additional mutations is provided wherein the one or more additional mutations are in one or more catalytically active domains. Such further catalytic mutations may confer nickase functionality as described in detail elsewhere herein. In such effector proteins, enhanced specificity may be achieved due to an improved specificity in terms of effector protein activity.

[0756] Modifications to reduce off-target effects and/or enhance on-target effects as described above may be made to amino acid residues located in a positively-charged region/groove situated between the RuvC-III and HNH domains. It will be appreciated that any of the functional effects described above may be achieved by modification of amino acids within the aforementioned groove but also by modification of amino acids adjacent to or outside of that groove.

[0757] Additional functionalities which may be engineered into modified CRISPR effector proteins as described herein include the following. 1. Modified CRISPR effector proteins that disrupt DNA: or RNA:protein interactions without affecting protein tertiary or secondary structure. This includes residues that contact any part of the RNA: DNA or RNA:RNA duplex. 2. modified CRISPR effector proteins that weaken intra-protein interactions holding the effector protein in conformation essential for nuclease cutting in response to DNA/RNA binding (on or off target). For example: a modification that mildly inhibits, but still allows, the nuclease conformation of the HNH domain (positioned at the scissile phosphate). 3. Modified CRISPR effector proteins that strengthen intra-protein interactions holding the effector protein in a conformation inhibiting nuclease activity in response to DNA/RNA binding (on or off targets). For example: a modification that stabilizes the HNH domain in a conformation away from the scissile phosphate. Any such additional functional enhancement may be provided in combination with any other modification to the CRISPR effector protein as described in detail elsewhere herein.

[0758] Any of the herein described improved functionalities may be made to any CRISPR effector protein Group 29 or Group 30 effector protein. However, it will be appreciated that any of the functionalities described herein may be engineered into Group 29 or Group 30 Group 29 or Group 30 effector proteins from other orthologs, including chimeric effector proteins comprising fragments from multiple orthologs.

[0759] The invention uses nucleic acids to bind target DNA or RNA sequences. This is advantageous as nucleic acids are much easier and cheaper to produce than proteins, and the specificity can be varied according to the length of the stretch where homology is sought. Complex 3-D positioning of multiple fingers, for example is not required. The terms “polynucleotide”, “nucleotide”, “nucleotide sequence”, “nucleic acid” and “oligonucleotide” are used interchangeably. They refer to a polymeric form of nucleotides of any length, either deoxyribonucleotides or ribonucleotides, or analogs thereof. Polynucleotides may have any three dimensional structure, and may perform any function, known or unknown. The following are non-limiting examples of polynucleotides: coding or non-coding regions of a gene or gene fragment, loci (locus) defined from linkage analysis, exons, introns, messenger RNA (mRNA), transfer RNA, ribosomal RNA, short interfering RNA (siRNA), short-hairpin RNA (shRNA), micro-RNA

(miRNA), ribozymes, cDNA, recombinant polynucleotides, branched polynucleotides, plasmids, vectors, isolated DNA of any sequence, isolated RNA of any sequence, nucleic acid probes, and primers. The term also encompasses nucleic-acid-like structures with synthetic backbones, see, e.g., Eckstein, 1991; Baserga et al., 1992; Milligan, 1993; WO 97/03211; WO 96/39154; Mata, 1997; Strauss-Soukup, 1997; and Samstag, 1996. A polynucleotide may comprise one or more modified nucleotides, such as methylated nucleotides and nucleotide analogs. If present, modifications to the nucleotide structure may be imparted before or after assembly of the polymer. The sequence of nucleotides may be interrupted by non-nucleotide components. A polynucleotide may be further modified after polymerization, such as by conjugation with a labeling component. As used herein the term “wild type” is a term of the art understood by skilled persons and means the typical form of an organism, strain, gene or characteristic as it occurs in nature as distinguished from mutant or variant forms. A “wild type” can be a base line. As used herein the term “variant” should be taken to mean the exhibition of qualities that have a pattern that deviates from what occurs in nature. The terms “non-naturally occurring” or “engineered” are used interchangeably and indicate the involvement of the hand of man. The terms, when referring to nucleic acid molecules or polypeptides mean that the nucleic acid molecule or the polypeptide is at least substantially free from at least one other component with which they are naturally associated in nature and as found in nature. “Complementarity” refers to the ability of a nucleic acid to form hydrogen bond(s) with another nucleic acid sequence by either traditional Watson-Crick base pairing or other non-traditional types. A percent complementarity indicates the percentage of residues in a nucleic acid molecule which can form hydrogen bonds (e.g., Watson-Crick base pairing) with a second nucleic acid sequence (e.g., 5, 6, 7, 8, 9, 10 out of 10 being 50%, 60%, 70%, 80%, 90%, and 100% complementary). “Perfectly complementary” means that all the contiguous residues of a nucleic acid sequence will hydrogen bond with the same number of contiguous residues in a second nucleic acid sequence. “Substantially complementary” as used herein refers to a degree of complementarity that is at least 60%, 65%, 70%, 75%, 80%, 85%, 90%, 95%, 97%, 98%, 99%, or 100% over a region of 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22, 23, 24, 25, 30, 35, 40, 45, 50, or more nucleotides, or refers to two nucleic acids that hybridize under stringent conditions. As used herein, “stringent conditions” for hybridization refer to conditions under which a nucleic acid having complementarity to a target sequence predominantly hybridizes with the target sequence, and substantially does not hybridize to non-target sequences. Stringent conditions are generally sequence-dependent, and vary depending on a number of factors. In general, the longer the sequence, the higher the temperature at which the sequence specifically hybridizes to its target sequence. Non-limiting examples of stringent conditions are described in detail in Tijssen (1993), *Laboratory Techniques In Biochemistry And Molecular Biology-Hybridization With Nucleic Acid Probes Part I, Second Chapter “Overview of principles of hybridization and the strategy of nucleic acid probe assay”*, Elsevier, N.Y. Where reference is made to a polynucleotide sequence, then complementary or partially complementary sequences are also envisaged. These are preferably capable of hybridizing to the reference sequence

under highly stringent conditions. Generally, in order to maximize the hybridization rate, relatively low-stringency hybridization conditions are selected: about 20 to 25° C. lower than the thermal melting point (T_m). The T_m is the temperature at which 50% of specific target sequence hybridizes to a perfectly complementary probe in solution at a defined ionic strength and pH. Generally, in order to require at least about 85% nucleotide complementarity of hybridized sequences, highly stringent washing conditions are selected to be about 5 to 15° C. lower than the T_m . In order to require at least about 70% nucleotide complementarity of hybridized sequences, moderately-stringent washing conditions are selected to be about 15 to 30° C. lower than the T_m . Highly permissive (very low stringency) washing conditions may be as low as 50° C. below the T_m , allowing a high level of mis-matching between hybridized sequences. Those skilled in the art will recognize that other physical and chemical parameters in the hybridization and wash stages can also be altered to affect the outcome of a detectable hybridization signal from a specific level of homology between target and probe sequences. Preferred highly stringent conditions comprise incubation in 50% formamide, 5×SSC, and 1% SDS at 42° C., or incubation in 5×SSC and 1% SDS at 65° C., with wash in 0.2×SSC and 0.1% SDS at 65° C. “Hybridization” refers to a reaction in which one or more polynucleotides react to form a complex that is stabilized via hydrogen bonding between the bases of the nucleotide residues. The hydrogen bonding may occur by Watson Crick base pairing, Hoogsteen binding, or in any other sequence specific manner. The complex may comprise two strands forming a duplex structure, three or more strands forming a multi stranded complex, a single self-hybridizing strand, or any combination of these. A hybridization reaction may constitute a step in a more extensive process, such as the initiation of PCR, or the cleavage of a polynucleotide by an enzyme. A sequence capable of hybridizing with a given sequence is referred to as the “complement” of the given sequence. As used herein, the term “genomic locus” or “locus” (plural loci) is the specific location of a gene or DNA sequence on a chromosome. A “gene” refers to stretches of DNA or RNA that encode a polypeptide or an RNA chain that has functional role to play in an organism and hence is the molecular unit of heredity in living organisms. For the purpose of this invention it may be considered that genes include regions which regulate the production of the gene product, whether or not such regulatory sequences are adjacent to coding and/or transcribed sequences. Accordingly, a gene includes, but is not necessarily limited to, promoter sequences, terminators, translational regulatory sequences such as ribosome binding sites and internal ribosome entry sites, enhancers, silencers, insulators, boundary elements, replication origins, matrix attachment sites and locus control regions. As used herein, “expression of a genomic locus” or “gene expression” is the process by which information from a gene is used in the synthesis of a functional gene product. The products of gene expression are often proteins, but in non-protein coding genes such as rRNA genes or tRNA genes, the product is functional RNA. The process of gene expression is used by all known life-eukaryotes (including multicellular organisms), prokaryotes (bacteria and archaea) and viruses to generate functional products to survive. As used herein “expression” of a gene or nucleic acid encompasses not only cellular gene expression, but also the transcription and translation of

nucleic acid(s) in cloning systems and in any other context. As used herein, “expression” also refers to the process by which a polynucleotide is transcribed from a DNA template (such as into and mRNA or other RNA transcript) and/or the process by which a transcribed mRNA is subsequently translated into peptides, polypeptides, or proteins. Transcripts and encoded polypeptides may be collectively referred to as “gene product.” If the polynucleotide is derived from genomic DNA, expression may include splicing of the mRNA in a eukaryotic cell. The terms “polypeptide”, “peptide” and “protein” are used interchangeably herein to refer to polymers of amino acids of any length. The polymer may be linear or branched, it may comprise modified amino acids, and it may be interrupted by non-amino acids. The terms also encompass an amino acid polymer that has been modified; for example, disulfide bond formation, glycosylation, lipidation, acetylation, phosphorylation, or any other manipulation, such as conjugation with a labeling component. As used herein the term “amino acid” includes natural and/or unnatural or synthetic amino acids, including glycine and both the D or L optical isomers, and amino acid analogs and peptidomimetics. As used herein, the term “domain” or “protein domain” refers to a part of a protein sequence that may exist and function independently of the rest of the protein chain. As described in aspects of the invention, sequence identity is related to sequence homology. Homology comparisons may be conducted by eye, or more usually, with the aid of readily available sequence comparison programs. These commercially available computer programs may calculate percent (%) homology between two or more sequences and may also calculate the sequence identity shared by two or more amino acid or nucleic acid sequences.

[0760] In aspects of the invention the term “guide RNA”, refers to the polynucleotide sequence comprising one or more of a putative or identified tracr sequence and a putative or identified crRNA sequence or guide sequence. In particular embodiments, the “guide RNA” comprises a putative or identified crRNA sequence or guide sequence. In further embodiments, the guide RNA does not comprise a putative or identified tracr sequence.

[0761] As used herein the term “wild type” is a term of the art understood by skilled persons and means the typical form of an organism, strain, gene or characteristic as it occurs in nature as distinguished from mutant or variant forms. A “wild type” can be a base line.

[0762] As used herein the term “variant” should be taken to mean the exhibition of qualities that have a pattern that deviates from what occurs in nature.

[0763] The terms “non-naturally occurring” or “engineered” are used interchangeably and indicate the involvement of the hand of man. The terms, when referring to nucleic acid molecules or polypeptides mean that the nucleic acid molecule or the polypeptide is at least substantially free from at least one other component with which they are naturally associated in nature and as found in nature. In all aspects and embodiments, whether they include these terms or not, it will be understood that, preferably, the may be optional and thus preferably included or not preferably not included. Furthermore, the terms “non-naturally occurring” and “engineered” may be used interchangeably and so can therefore be used alone or in combination and one or other may replace mention of both together. In particular, “engi-

neered” is preferred in place of “non-naturally occurring” or “non-naturally occurring and/or engineered.”

[0764] Sequence homologies may be generated by any of a number of computer programs known in the art, for example BLAST or FASTA, etc. A suitable computer program for carrying out such an alignment is the GCG Wisconsin Bestfit package (University of Wisconsin, U.S.A.; Devereux et al., 1984, *Nucleic Acids Research* 12:387). Examples of other software than may perform sequence comparisons include, but are not limited to, the BLAST package (see Ausubel et al., 1999 *ibid*—Chapter 18), FASTA (Atschul et al., 1990, *J. Mol. Biol.*, 403-410) and the GENWORKS suite of comparison tools. Both BLAST and FASTA are available for offline and online searching (see Ausubel et al., 1999 *ibid*, pages 7-58 to 7-60). However it is preferred to use the GCG Bestfit program. Percentage (%) sequence homology may be calculated over contiguous sequences, i.e., one sequence is aligned with the other sequence and each amino acid or nucleotide in one sequence is directly compared with the corresponding amino acid or nucleotide in the other sequence, one residue at a time. This is called an “ungapped” alignment. Typically, such ungapped alignments are performed only over a relatively short number of residues. Although this is a very simple and consistent method, it fails to take into consideration that, for example, in an otherwise identical pair of sequences, one insertion or deletion may cause the following amino acid residues to be put out of alignment, thus potentially resulting in a large reduction in % homology when a global alignment is performed. Consequently, most sequence comparison methods are designed to produce optimal alignments that take into consideration possible insertions and deletions without unduly penalizing the overall homology or identity score. This is achieved by inserting “gaps” in the sequence alignment to try to maximize local homology or identity. However, these more complex methods assign “gap penalties” to each gap that occurs in the alignment so that, for the same number of identical amino acids, a sequence alignment with as few gaps as possible—reflecting higher relatedness between the two compared sequences—may achieve a higher score than one with many gaps. “Affinity gap costs” are typically used that charge a relatively high cost for the existence of a gap and a smaller penalty for each subsequent residue in the gap. This is the most commonly used gap scoring system. High gap penalties may, of course, produce optimized alignments with fewer gaps. Most alignment programs allow the gap penalties to be modified. However, it is preferred to use the default values when using such software for sequence comparisons. For example, when using the GCG Wisconsin Bestfit package the default gap penalty for amino acid sequences is -12 for a gap and -4 for each extension. Calculation of maximum % homology therefore first requires the production of an optimal alignment, taking into consideration gap penalties. A suitable computer program for carrying out such an alignment is the GCG Wisconsin Bestfit package (Devereux et al., 1984 *Nuc. Acids Research* 12 p387). Examples of other software than may perform sequence comparisons include, but are not limited to, the BLAST package (see Ausubel et al., 1999 *Short Protocols in Molecular Biology*, 4th Ed.—Chapter 18), FASTA (Altschul et al., 1990 *J. Mol. Biol.* 403-410) and the GENWORKS suite of comparison tools. Both BLAST and FASTA are available for offline and online searching (see Ausubel et al., 1999, *Short Protocols in Molecular Biology*,

pages 7-58 to 7-60). However, for some applications, it is preferred to use the GCG Bestfit program. A new tool, called BLAST 2 Sequences is also available for comparing protein and nucleotide sequences (see *FEMS Microbiol Lett.* 1999 174(2): 247-50; *FEMS Microbiol Lett.* 1999 177(1): 187-8 and the website of the National Center for Biotechnology information at the website of the National Institutes for Health). Although the final % homology may be measured in terms of identity, the alignment process itself is typically not based on an all-or-nothing pair comparison. Instead, a scaled similarity score matrix is generally used that assigns scores to each pair-wise comparison based on chemical similarity or evolutionary distance. An example of such a matrix commonly used is the BLOSUM62 matrix—the default matrix for the BLAST suite of programs. GCG Wisconsin programs generally use either the public default values or a custom symbol comparison table, if supplied (see user manual for further details). For some applications, it is preferred to use the public default values for the GCG package, or in the case of other software, the default matrix, such as BLOSUM62. Alternatively, percentage homologies may be calculated using the multiple alignment feature in DNASIS™ (Hitachi Software), based on an algorithm, analogous to CLUSTAL (Higgins D G & Sharp P M (1988), *Gene* 73(1), 237-244). Once the software has produced an optimal alignment, it is possible to calculate % homology, preferably % sequence identity. The software typically does this as part of the sequence comparison and generates a numerical result. The sequences may also have deletions, insertions or substitutions of amino acid residues which produce a silent change and result in a functionally equivalent substance. Deliberate amino acid substitutions may be made on the basis of similarity in amino acid properties (such as polarity, charge, solubility, hydrophobicity, hydrophilicity, and/or the amphipathic nature of the residues) and it is therefore useful to group amino acids together in functional groups. Amino acids may be grouped together based on the properties of their side chains alone. However, it is more useful to include mutation data as well. The sets of amino acids thus derived are likely to be conserved for structural reasons. These sets may be described in the form of a Venn diagram (Livingstone C. D. and Barton G. J. (1993) “Protein sequence alignments: a strategy for the hierarchical analysis of residue conservation” *Comput. Appl. Biosci.* 9: 745-756) (Taylor W. R. (1986) “The classification of amino acid conservation” *J. Theor. Biol.* 119; 205-218). Conservative substitutions may be made, for example according to the table below which describes a generally accepted Venn diagram grouping of amino acids.

Set		Sub-set	
Hydrophobic	FWYHKMILVAGC	Aromatic	F W Y H
		Aliphatic	I L V
Polar	WYHKREDCSTNQ	Charged	H K R E D
		Positively charged	H K R
		Negatively charged	E D
Small	VCAGSPTND	Tiny	A G S

[0765] The terms “subject,” “individual,” and “patient” are used interchangeably herein to refer to a vertebrate,

preferably a mammal, more preferably a human. Mammals include, but are not limited to, murines, simians, humans, farm animals, sport animals, and pets. Tissues, cells and their progeny of a biological entity obtained in vivo or cultured in vitro are also encompassed.

[0766] The terms “therapeutic agent”, “therapeutic capable agent” or “treatment agent” are used interchangeably and refer to a molecule or compound that confers some beneficial effect upon administration to a subject. The beneficial effect includes enablement of diagnostic determinations; amelioration of a disease, symptom, disorder, or pathological condition; reducing or preventing the onset of a disease, symptom, disorder or condition; and generally counteracting a disease, symptom, disorder or pathological condition.

[0767] As used herein, “treatment” or “treating,” or “palliating” or “ameliorating” are used interchangeably. These terms refer to an approach for obtaining beneficial or desired results including but not limited to a therapeutic benefit and/or a prophylactic benefit. By therapeutic benefit is meant any therapeutically relevant improvement in or effect on one or more diseases, conditions, or symptoms under treatment. For prophylactic benefit, the compositions may be administered to a subject at risk of developing a particular disease, condition, or symptom, or to a subject reporting one or more of the physiological symptoms of a disease, even though the disease, condition, or symptom may not have yet been manifested.

[0768] The term “effective amount” or “therapeutically effective amount” refers to the amount of an agent that is sufficient to effect beneficial or desired results. The therapeutically effective amount may vary depending upon one or more of: the subject and disease condition being treated, the weight and age of the subject, the severity of the disease condition, the manner of administration and the like, which can readily be determined by one of ordinary skill in the art. The term also applies to a dose that will provide an image for detection by any one of the imaging methods described herein. The specific dose may vary depending on one or more of: the particular agent chosen, the dosing regimen to be followed, whether it is administered in combination with other compounds, timing of administration, the tissue to be imaged, and the physical delivery system in which it is carried.

[0769] The practice of the present invention employs, unless otherwise indicated, conventional techniques of immunology, biochemistry, chemistry, molecular biology, microbiology, cell biology, genomics and recombinant DNA, which are within the skill of the art. See Sambrook, Fritsch and Maniatis, *MOLECULAR CLONING: A LABORATORY MANUAL*, 2nd edition (1989); *CURRENT PROTOCOLS IN MOLECULAR BIOLOGY* (F. M. Ausubel, et al. eds., (1987)); the series *METHODS IN ENZYMOLOGY* (Academic Press, Inc.): *PCR 2: A PRACTICAL APPROACH* (M. J. MacPherson, B. D. Hames and G. R. Taylor eds. (1995)), Harlow and Lane, eds. (1988) *ANTIBODIES, A LABORATORY MANUAL*, and *ANIMAL CELL CULTURE* (R. I. Freshney, ed. (1987)).

[0770] Several aspects of the invention relate to vector systems comprising one or more vectors, or vectors as such. Vectors can be designed for expression of CRISPR transcripts (e.g. nucleic acid transcripts, proteins, or enzymes) in prokaryotic or eukaryotic cells. For example, CRISPR transcripts can be expressed in bacterial cells such as *Escheri-*

chia coli, insect cells (using baculovirus expression vectors), yeast cells, or mammalian cells. Suitable host cells are discussed further in Goeddel, *GENE EXPRESSION TECHNOLOGY: METHODS IN ENZYMOLOGY* 185, Academic Press, San Diego, Calif. (1990). Alternatively, the recombinant expression vector can be transcribed and translated in vitro, for example using T7 promoter regulatory sequences and T7 polymerase.

[0771] Embodiments of the invention include sequences (both polynucleotide or polypeptide) which may comprise homologous substitution (substitution and replacement are both used herein to mean the interchange of an existing amino acid residue or nucleotide, with an alternative residue or nucleotide) that may occur i.e., like-for-like substitution in the case of amino acids such as basic for basic, acidic for acidic, polar for polar, etc. Non-homologous substitution may also occur i.e., from one class of residue to another or alternatively involving the inclusion of unnatural amino acids such as ornithine (hereinafter referred to as Z), diamino butyric acid ornithine (hereinafter referred to as B), nor-leucine ornithine (hereinafter referred to as O), pyrrolysine, thienylalanine, naphthylalanine and phenylglycine. Variant amino acid sequences may include suitable spacer groups that may be inserted between any two amino acid residues of the sequence including alkyl groups such as methyl, ethyl or propyl groups in addition to amino acid spacers such as glycine or β -alanine residues. A further form of variation, which involves the presence of one or more amino acid residues in peptoid form, may be well understood by those skilled in the art. For the avoidance of doubt, “the peptoid form” is used to refer to variant amino acid residues wherein the α -carbon substituent group is on the residue’s nitrogen atom rather than the α -carbon. Processes for preparing peptides in the peptoid form are known in the art, for example Simon R J et al., *PNAS* (1992) 89(20), 9367-9371 and Horwell D C, *Trends Biotechnol.* (1995) 13(4), 132-134.

[0772] Homology modelling: Corresponding residues in other Group 29 or Group 30 orthologs can be identified by the methods of Zhang et al., 2012 (*Nature*; 490(7421): 556-60) and Chen et al., 2015 (*PLoS Comput Biol*; 11(5): e1004248)—a computational protein-protein interaction (PPI) method to predict interactions mediated by domain-motif interfaces. PrePPI (Predicting PPI), a structure based PPI prediction method, combines structural evidence with non-structural evidence using a Bayesian statistical framework. The method involves taking a pair a query proteins and using structural alignment to identify structural representatives that correspond to either their experimentally determined structures or homology models. Structural alignment is further used to identify both close and remote structural neighbors by considering global and local geometric relationships. Whenever two neighbors of the structural representatives form a complex reported in the Protein Data Bank, this defines a template for modelling the interaction between the two query proteins. Models of the complex are created by superimposing the representative structures on their corresponding structural neighbor in the template. This approach is further described in Dey et al., 2013 (*Prot Sci*; 22: 359-66).

[0773] For purpose of this invention, amplification means any method employing a primer and a polymerase capable of replicating a target sequence with reasonable fidelity. Amplification may be carried out by natural or recombinant

DNA polymerases such as TaqGold™, T7 DNA polymerase, Klenow fragment of *E. coli* DNA polymerase, and reverse transcriptase. A preferred amplification method is PCR.

[0774] In certain aspects the invention involves vectors. A used herein, a “vector” is a tool that allows or facilitates the transfer of an entity from one environment to another. It is a replicon, such as a plasmid, phage, or cosmid, into which another DNA segment may be inserted so as to bring about the replication of the inserted segment. Generally, a vector is capable of replication when associated with the proper control elements. In general, the term “vector” refers to a nucleic acid molecule capable of transporting another nucleic acid to which it has been linked. Vectors include, but are not limited to, nucleic acid molecules that are single-stranded, double-stranded, or partially double-stranded; nucleic acid molecules that comprise one or more free ends, no free ends (e.g., circular); nucleic acid molecules that comprise DNA, RNA, or both; and other varieties of polynucleotides known in the art. One type of vector is a “plasmid,” which refers to a circular double stranded DNA loop into which additional DNA segments can be inserted, such as by standard molecular cloning techniques. Another type of vector is a viral vector, wherein virally-derived DNA or RNA sequences are present in the vector for packaging into a virus (e.g., retroviruses, replication defective retroviruses, adenoviruses, replication defective adenoviruses, and adeno-associated viruses (AAVs)). Viral vectors also include polynucleotides carried by a virus for transfection into a host cell. Certain vectors are capable of autonomous replication in a host cell into which they are introduced (e.g., bacterial vectors having a bacterial origin of replication and episomal mammalian vectors). Other vectors (e.g., non-episomal mammalian vectors) are integrated into the genome of a host cell upon introduction into the host cell, and thereby are replicated along with the host genome. Moreover, certain vectors are capable of directing the expression of genes to which they are operatively-linked. Such vectors are referred to herein as “expression vectors.” Common expression vectors of utility in recombinant DNA techniques are often in the form of plasmids.

[0775] Recombinant expression vectors can comprise a nucleic acid of the invention in a form suitable for expression of the nucleic acid in a host cell, which means that the recombinant expression vectors include one or more regulatory elements, which may be selected on the basis of the host cells to be used for expression, that is operatively-linked to the nucleic acid sequence to be expressed. Within a recombinant expression vector, “operably linked” is intended to mean that the nucleotide sequence of interest is linked to the regulatory element(s) in a manner that allows for expression of the nucleotide sequence (e.g., in an *in vitro* transcription/translation system or in a host cell when the vector is introduced into the host cell). With regards to recombination and cloning methods, mention is made of U.S. patent application Ser. No. 10/815,730, published Sep. 2, 2004 as US 2004-0171156 A1, the contents of which are herein incorporated by reference in their entirety.

[0776] Aspects of the invention relate to bicistronic vectors for guide RNA and wild type, modified or mutated CRISPR effector proteins/enzymes (e.g. Group 29 or Group 30 effector proteins). Bicistronic expression vectors guide RNA and wild type, modified or mutated CRISPR effector proteins/enzymes (e.g. Group 29 or Group 30 effector proteins) are preferred. In general and particularly in this

embodiment and wild type, modified or mutated CRISPR effector proteins/enzymes (e.g. Group 29 or Group 30 effector proteins) is preferably driven by the CBh promoter. The RNA may preferably be driven by a Pol III promoter, such as a U6 promoter. Ideally the two are combined.

[0777] In some embodiments, a loop in the guide RNA is provided. This may be a stem loop or a tetra loop. The loop is preferably GAAA, but it is not limited to this sequence or indeed to being only 4 bp in length. Indeed, preferred loop forming sequences for use in hairpin structures are four nucleotides in length, and most preferably have the sequence GAAA. However, longer or shorter loop sequences may be used, as may alternative sequences. The sequences preferably include a nucleotide triplet (for example, AAA), and an additional nucleotide (for example C or G). Examples of loop forming sequences include CAAA and AAAG.

[0778] In practicing any of the methods disclosed herein, a suitable vector can be introduced to a cell or an embryo via one or more methods known in the art, including without limitation, microinjection, electroporation, sonoporation, biolistics, calcium phosphate-mediated transfection, cationic transfection, liposome transfection, dendrimer transfection, heat shock transfection, nucleofection transfection, magnetofection, lipofection, impalefection, optical transfection, proprietary agent-enhanced uptake of nucleic acids, and delivery via liposomes, immunoliposomes, virosomes, or artificial virions. In some methods, the vector is introduced into an embryo by microinjection. The vector or vectors may be microinjected into the nucleus or the cytoplasm of the embryo. In some methods, the vector or vectors may be introduced into a cell by nucleofection.

[0779] The term “regulatory element” is intended to include promoters, enhancers, internal ribosomal entry sites (IRES), and other expression control elements (e.g., transcription termination signals, such as polyadenylation signals and poly-U sequences). Such regulatory elements are described, for example, in Goeddel, GENE EXPRESSION TECHNOLOGY: METHODS IN ENZYMOLOGY 185, Academic Press, San Diego, Calif. (1990). Regulatory elements include those that direct constitutive expression of a nucleotide sequence in many types of host cell and those that direct expression of the nucleotide sequence only in certain host cells (e.g., tissue-specific regulatory sequences). A tissue-specific promoter may direct expression primarily in a desired tissue of interest, such as muscle, neuron, bone, skin, blood, specific organs (e.g., liver, pancreas), or particular cell types (e.g., lymphocytes). Regulatory elements may also direct expression in a temporal-dependent manner, such as in a cell-cycle dependent or developmental stage-dependent manner, which may or may not also be tissue or cell-type specific. In some embodiments, a vector comprises one or more pol III promoter (e.g., 1, 2, 3, 4, 5, or more pol III promoters), one or more pol II promoters (e.g., 1, 2, 3, 4, 5, or more pol II promoters), one or more pol I promoters (e.g., 1, 2, 3, 4, 5, or more pol I promoters), or combinations thereof. Examples of pol III promoters include, but are not limited to, U6 and H1 promoters. Examples of pol II promoters include, but are not limited to, the retroviral Rous sarcoma virus (RSV) LTR promoter (optionally with the RSV enhancer), the cytomegalovirus (CMV) promoter (optionally with the CMV enhancer) [see, e.g., Boshart et al, Cell, 41:521-530 (1985)], the SV40 promoter, the dihydrofolate reductase promoter, the β -actin promoter, the phosphoglycerol kinase (PGK) promoter, and the EF1 α pro-

motor. Also encompassed by the term “regulatory element” are enhancer elements, such as WPRE; CMV enhancers; the R-U5' segment in LTR of HTLV-I (Mol. Cell. Biol., Vol. 8(1), p. 466-472, 1988); SV40 enhancer; and the intron sequence between exons 2 and 3 of rabbit β -globin (Proc. Natl. Acad. Sci. USA., Vol. 78(3), p. 1527-31, 1981). It will be appreciated by those skilled in the art that the design of the expression vector can depend on such factors as the choice of the host cell to be transformed, the level of expression desired, etc. A vector can be introduced into host cells to thereby produce transcripts, proteins, or peptides, including fusion proteins or peptides, encoded by nucleic acids as described herein (e.g., clustered regularly interspersed short palindromic repeats (CRISPR) transcripts, proteins, enzymes, mutant forms thereof, fusion proteins thereof, etc.). With regards to regulatory sequences, mention is made of U.S. patent application Ser. No. 10/491,026, the contents of which are incorporated by reference herein in their entirety. With regards to promoters, mention is made of PCT publication WO 2011/028929 and U.S. application Ser. No. 12/511,940, the contents of which are incorporated by reference herein in their entirety.

[0780] Vectors can be designed for expression of CRISPR transcripts (e.g., nucleic acid transcripts, proteins, or enzymes) in prokaryotic or eukaryotic cells. For example, CRISPR transcripts can be expressed in bacterial cells such as *Escherichia coli*, insect cells (using baculovirus expression vectors), yeast cells, or mammalian cells. Suitable host cells are discussed further in Goeddel, GENE EXPRESSION TECHNOLOGY: METHODS IN ENZYMOLOGY 185, Academic Press, San Diego, Calif. (1990). Alternatively, the recombinant expression vector can be transcribed and translated in vitro, for example using T7 promoter regulatory sequences and T7 polymerase.

[0781] Vectors may be introduced and propagated in a prokaryote or prokaryotic cell. In some embodiments, a prokaryote is used to amplify copies of a vector to be introduced into a eukaryotic cell or as an intermediate vector in the production of a vector to be introduced into a eukaryotic cell (e.g., amplifying a plasmid as part of a viral vector packaging system). In some embodiments, a prokaryote is used to amplify copies of a vector and express one or more nucleic acids, such as to provide a source of one or more proteins for delivery to a host cell or host organism. Expression of proteins in prokaryotes is most often carried out in *Escherichia coli* with vectors containing constitutive or inducible promoters directing the expression of either fusion or non-fusion proteins. Fusion vectors add a number of amino acids to a protein encoded therein, such as to the amino terminus of the recombinant protein. Such fusion vectors may serve one or more purposes, such as: (i) to increase expression of recombinant protein; (ii) to increase the solubility of the recombinant protein; and (iii) to aid in the purification of the recombinant protein by acting as a ligand in affinity purification. Often, in fusion expression vectors, a proteolytic cleavage site is introduced at the junction of the fusion moiety and the recombinant protein to enable separation of the recombinant protein from the fusion moiety subsequent to purification of the fusion protein. Such enzymes, and their cognate recognition sequences, include Factor Xa, thrombin and enterokinase. Example fusion expression vectors include pGEX (Pharmacia Biotech Inc; Smith and Johnson, 1988. *Gene* 67: 31-40), pMAL (New England Biolabs, Beverly, Mass.) and pRIT5 (Pharmacia,

Piscataway, N.J.) that fuse glutathione S-transferase (GST), maltose E binding protein, or protein A, respectively, to the target recombinant protein.

[0782] Examples of suitable inducible non-fusion *E. coli* expression vectors include pTrc (Amrann et al., (1988) *Gene* 69:301-315) and pET lld (Studier et al., GENE EXPRESSION TECHNOLOGY: METHODS IN ENZYMOLOGY 185, Academic Press, San Diego, Calif. (1990) 60-89).

[0783] In some embodiments, a vector is a yeast expression vector. Examples of vectors for expression in yeast *Saccharomyces cerevisiae* include pYepSec1 (Baldari, et al., 1987. *EMBO J.* 6: 229-234), pMFa (Kuijan and Herskowitz, 1982. *Cell* 30: 933-943), pJRY88 (Schultz et al., 1987. *Gene* 54: 113-123), pYES2 (Invitrogen Corporation, San Diego, Calif.), and picZ (Invitrogen Corp, San Diego, Calif.).

[0784] In some embodiments, a vector drives protein expression in insect cells using baculovirus expression vectors. Baculovirus vectors available for expression of proteins in cultured insect cells (e.g., SF9 cells) include the pAc series (Smith, et al., 1983. *Mol. Cell. Biol.* 3: 2156-2165) and the pVL series (Lucklow and Summers, 1989. *Virology* 170: 31-39).

[0785] In some embodiments, a vector is capable of driving expression of one or more sequences in mammalian cells using a mammalian expression vector. Examples of mammalian expression vectors include pCDM8 (Seed, 1987. *Nature* 329: 840) and pMT2PC (Kaufman, et al., 1987. *EMBO J.* 6: 187-195). When used in mammalian cells, the expression vector's control functions are typically provided by one or more regulatory elements. For example, commonly used promoters are derived from polyoma, adenovirus 2, cytomegalovirus, simian virus 40, and others disclosed herein and known in the art. For other suitable expression systems for both prokaryotic and eukaryotic cells see, e.g., Chapters 16 and 17 of Sambrook, et al., MOLECULAR CLONING: A LABORATORY MANUAL, 2nd ed., Cold Spring Harbor Laboratory, Cold Spring Harbor Laboratory Press, Cold Spring Harbor, N.Y., 1989.

[0786] In some embodiments, the recombinant mammalian expression vector is capable of directing expression of the nucleic acid preferentially in a particular cell type (e.g., tissue-specific regulatory elements are used to express the nucleic acid). Tissue-specific regulatory elements are known in the art. Non-limiting examples of suitable tissue-specific promoters include the albumin promoter (liver-specific; Pinkert, et al., 1987. *Genes Dev.* 1: 268-277), lymphoid-specific promoters (Calame and Eaton, 1988. *Adv. Immunol.* 43: 235-275), in particular promoters of T cell receptors (Winoto and Baltimore, 1989. *EMBO J.* 8: 729-733) and immunoglobulins (Baneiji, et al., 1983. *Cell* 33: 729-740; Queen and Baltimore, 1983. *Cell* 33: 741-748), neuron-specific promoters (e.g., the neurofilament promoter; Byrne and Ruddle, 1989. *Proc. Natl. Acad. Sci. USA* 86: 5473-5477), pancreas-specific promoters (Edlund, et al., 1985. *Science* 230: 912-916), and mammary gland-specific promoters (e.g., milk whey promoter; U.S. Pat. No. 4,873,316 and European Application Publication No. 264,166). Developmentally-regulated promoters are also encompassed, e.g., the murine hox promoters (Kessel and Gruss, 1990. *Science* 249: 374-379) and the α -fetoprotein promoter (Campes and Tilghman, 1989. *Genes Dev.* 3: 537-546). With regards to these prokaryotic and eukaryotic vectors, mention is made of U.S. Pat. No. 6,750,059, the contents of which are incorporated by reference herein in their entirety. Other

embodiments of the invention may relate to the use of viral vectors, with regards to which mention is made of U.S. patent application Ser. No. 13/092,085, the contents of which are incorporated by reference herein in their entirety. Tissue-specific regulatory elements are known in the art and in this regard, mention is made of U.S. Pat. No. 7,776,321, the contents of which are incorporated by reference herein in their entirety.

[0787] In some embodiments, a regulatory element is operably linked to one or more elements of a CRISPR system so as to drive expression of the one or more elements of the CRISPR system. In general, CRISPRs (Clustered Regularly Interspaced Short Palindromic Repeats), also known as SPIDRs (Spacer Interspersed Direct Repeats), constitute a family of DNA loci that are usually specific to a particular bacterial species. The CRISPR locus comprises a distinct class of interspersed short sequence repeats (SSRs) that were recognized in *E. coli* (Ishino et al., *J. Bacteriol.*, 169:5429-5433 [1987]; and Nakata et al., *J. Bacteriol.*, 171:3553-3556 [1989]), and associated genes. Similar interspersed SSRs have been identified in *Haloferax mediterranei*, *Streptococcus pyogenes*, *Anabaena*, and *Mycobacterium tuberculosis* (See, Groenen et al., *Mol. Microbiol.*, 10:1057-1065 [1993]; Hoe et al., *Emerg. Infect. Dis.*, 5:254-263 [1999]; Masepohl et al., *Biochim. Biophys. Acta* 1307: 26-30 [1996]; and Mojica et al., *Mol. Microbiol.*, 17:85-93 [1995]). The CRISPR loci typically differ from other SSRs by the structure of the repeats, which have been termed short regularly spaced repeats (SRSRs) (Janssen et al., *OMICS J. Integ. Biol.*, 6:23-33 [2002]; and Mojica et al., *Mol. Microbiol.*, 36:244-246 [2000]). In general, the repeats are short elements that occur in clusters that are regularly spaced by unique intervening sequences with a substantially constant length (Mojica et al., [2000], supra). Although the repeat sequences are highly conserved between strains, the number of interspersed repeats and the sequences of the spacer regions typically differ from strain to strain (van Embden et al., *J. Bacteriol.*, 182:2393-2401 [2000]). CRISPR loci have been identified in more than 40 prokaryotes (See e.g., Jansen et al., *Mol. Microbiol.*, 43:1565-1575 [2002]; and Mojica et al., [2005]) including, but not limited to *Aeropyrum*, *Pyrobaculum*, *Sulfolobus*, *Archaeoglobus*, *Halocarcula*, *Methanobacterium*, *Methanococcus*, *Methanosarcina*, *Methanopyrus*, *Pyrococcus*, *Picrophilus*, *Thermoplasma*, *Corynebacterium*, *Mycobacterium*, *Streptomyces*, *Aquifex*, *Porphyromonas*, *Chlorobium*, *Thermus*, *Bacillus*, *Listeria*, *Staphylococcus*, *Clostridium*, *Thermoanaerobacter*, *Mycoplasma*, *Fusobacterium*, *Azarcus*, *Chromobacterium*, *Neisseria*, *Nitrosomonas*, *Desulfovibrio*, *Geobacter*, *Myxococcus*, *Campylobacter*, *Wolinella*, *Acinetobacter*, *Envinia*, *Escherichia*, *Legionella*, *Methylococcus*, *Pasteurella*, *Photobacterium*, *Salmonella*, *Xanthomonas*, *Yersinia*, *Treponema*, and *Thermotoga*.

[0788] In general, “nucleic acid-targeting system” as used in the present application refers collectively to transcripts and other elements involved in the expression of or directing the activity of nucleic acid-targeting CRISPR-associated (“Cas”) genes (also referred to herein as an effector protein), including sequences encoding a nucleic acid-targeting Cas (effector) protein and a guide RNA (comprising crRNA sequence and a trans-activating CRISPR/Cas system RNA (tracrRNA) sequence), or other sequences and transcripts from a nucleic acid-targeting CRISPR locus. In some embodiments, one or more elements of a nucleic acid-

targeting system are derived from a Type V/Type VI nucleic acid-targeting CRISPR system. In some embodiments, one or more elements of a nucleic acid-targeting system is derived from a particular organism comprising an endogenous nucleic acid-targeting CRISPR system. In general, a nucleic acid-targeting system is characterized by elements that promote the formation of a nucleic acid-targeting complex at the site of a target sequence. In the context of formation of a nucleic acid-targeting complex, “target sequence” refers to a sequence to which a guide sequence is designed to have complementarity, where hybridization between a target sequence and a guide RNA promotes the formation of a DNA or RNA-targeting complex. Full complementarity is not necessarily required, provided there is sufficient complementarity to cause hybridization and promote formation of a nucleic acid-targeting complex. A target sequence may comprise RNA polynucleotides. In some embodiments, a target sequence is located in the nucleus or cytoplasm of a cell. In some embodiments, the target sequence may be within an organelle of a eukaryotic cell, for example, mitochondrion or chloroplast.

[0789] Typically, in the context of an endogenous nucleic acid-targeting system, formation of a nucleic acid-targeting complex (comprising a guide RNA hybridized to a target sequence and complexed with one or more nucleic acid-targeting effector proteins) results in cleavage of one or both RNA strands in or near (e.g. within 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 20, 50, or more base pairs from) the target sequence. In some embodiments, one or more vectors driving expression of one or more elements of a nucleic acid-targeting system are introduced into a host cell such that expression of the elements of the nucleic acid-targeting system direct formation of a nucleic acid-targeting complex at one or more target sites. For example, a nucleic acid-targeting effector protein and a guide RNA could each be operably linked to separate regulatory elements on separate vectors. Alternatively, two or more of the elements expressed from the same or different regulatory elements, may be combined in a single vector, with one or more additional vectors providing any components of the nucleic acid-targeting system not included in the first vector. nucleic acid-targeting system elements that are combined in a single vector may be arranged in any suitable orientation, such as one element located 5' with respect to (“upstream” of) or 3' with respect to (“downstream” of) a second element. The coding sequence of one element may be located on the same or opposite strand of the coding sequence of a second element, and oriented in the same or opposite direction. In some embodiments, a single promoter drives expression of a transcript encoding a nucleic acid-targeting effector protein and a guide RNA embedded within one or more intron sequences (e.g. each in a different intron, two or more in at least one intron, or all in a single intron). In some embodiments, the nucleic acid-targeting effector protein and guide RNA are operably linked to and expressed from the same promoter.

[0790] In general, a guide sequence is any polynucleotide sequence having sufficient complementarity with a target polynucleotide sequence to hybridize with the target sequence and direct sequence-specific binding of a nucleic acid-targeting complex to the target sequence. In some embodiments, the degree of complementarity between a guide sequence and its corresponding target sequence, when optimally aligned using a suitable alignment algorithm, is

about or more than about 50%, 60%, 75%, 80%, 85%, 90%, 95%, 97.5%, 99%, or more. Optimal alignment may be determined with the use of any suitable algorithm for aligning sequences, non-limiting example of which include the Smith-Waterman algorithm, the Needleman-Wunsch algorithm, algorithms based on the Burrows-Wheeler Transform (e.g. the Burrows Wheeler Aligner), ClustalW, Clustal X, BLAT, Novoalign (Novocraft Technologies, ELAND (Illumina, San Diego, Calif.), SOAP (available at soap.genomics.org.cn), and Maq (available at maq.sourceforge.net). In some embodiments, a guide sequence is about or more than about 5, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30, 35, 40, 45, 50, 75, or more nucleotides in length. In some embodiments, a guide sequence is less than about 75, 50, 45, 40, 35, 30, 25, 20, 15, 12, or fewer nucleotides in length. The ability of a guide sequence to direct sequence-specific binding of a nucleic acid-targeting complex to a target sequence may be assessed by any suitable assay. For example, the components of a nucleic acid-targeting system sufficient to form a nucleic acid-targeting complex, including the guide sequence to be tested, may be provided to a host cell having the corresponding target sequence, such as by transfection with vectors encoding the components of the nucleic acid-targeting CRISPR sequence, followed by an assessment of preferential cleavage within or in the vicinity of the target sequence, such as by Surveyor assay as described herein. Similarly, cleavage of a target polynucleotide sequence (or a sequence in the vicinity thereof) may be evaluated in a test tube by providing the target sequence, components of a nucleic acid-targeting complex, including the guide sequence to be tested and a control guide sequence different from the test guide sequence, and comparing binding or rate of cleavage at or in the vicinity of the target sequence between the test and control guide sequence reactions. Other assays are possible, and will occur to those skilled in the art.

[0791] A guide sequence may be selected to target any target sequence. In some embodiments, the target sequence is a sequence within a gene transcript or mRNA.

[0792] In some embodiments, the target sequence is a sequence within a genome of a cell.

[0793] In some embodiments, a guide sequence is selected to reduce the degree of secondary structure within the guide sequence. Secondary structure may be determined by any suitable polynucleotide folding algorithm. Some programs are based on calculating the minimal Gibbs free energy. An example of one such algorithm is mFold, as described by Zuker and Stiegler (Nucleic Acids Res. 9 (1981), 133-148). Another example folding algorithm is the online webserver RNAfold, developed at Institute for Theoretical Chemistry at the University of Vienna, using the centroid structure prediction algorithm (see e.g. A. R. Gruber et al., 2008, *Cell* 106(1): 23-24; and PA Carr and GM Church, 2009, *Nature Biotechnology* 27(12): 1151-62). Further algorithms may be found in U.S. application Ser. No. _____ (attorney docket 44790.11.2022; Broad Reference BI-2013/004A); incorporated herein by reference.

[0794] In some embodiments, the nucleic acid-targeting effector protein is part of a fusion protein comprising one or more heterologous protein domains (e.g., about or more than about 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, or more domains in addition to the nucleic acid-targeting effector protein). In some embodiments, the CRISPR effector protein/enzyme is part of a fusion protein comprising one or more heterologous

protein domains (e.g. about or more than about 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, or more domains in addition to the CRISPR enzyme). A CRISPR effector protein/enzyme fusion protein may comprise any additional protein sequence, and optionally a linker sequence between any two domains. Examples of protein domains that may be fused to an effector protein include, without limitation, epitope tags, reporter gene sequences, and protein domains having one or more of the following activities: methylase activity, demethylase activity, transcription activation activity, transcription repression activity, transcription release factor activity, histone modification activity, RNA cleavage activity and nucleic acid binding activity. Non-limiting examples of epitope tags include histidine (His) tags, V5 tags, FLAG tags, influenza hemagglutinin (HA) tags, Myc tags, VSV-G tags, and thioredoxin (Trx) tags. Examples of reporter genes include, but are not limited to, glutathione-S-transferase (GST), horseradish peroxidase (HRP), chloramphenicol acetyltransferase (CAT) beta-galactosidase, beta-glucuronidase, luciferase, green fluorescent protein (GFP), HcRed, DsRed, cyan fluorescent protein (CFP), yellow fluorescent protein (YFP), and autofluorescent proteins including blue fluorescent protein (BFP). A nucleic acid-targeting effector protein may be fused to a gene sequence encoding a protein or a fragment of a protein that bind DNA molecules or bind other cellular molecules, including but not limited to maltose binding protein (MBP), S-tag, Lex A DNA binding domain (DBD) fusions, GAL4 DNA binding domain fusions, and herpes simplex virus (HSV) BP16 protein fusions. Additional domains that may form part of a fusion protein comprising a nucleic acid-targeting effector protein are described in US20110059502, incorporated herein by reference. In some embodiments, a tagged nucleic acid-targeting effector protein is used to identify the location of a target sequence.

[0795] In some embodiments, a CRISPR enzyme may form a component of an inducible system. The inducible nature of the system would allow for spatiotemporal control of gene editing or gene expression using a form of energy. The form of energy may include but is not limited to electromagnetic radiation, sound energy, chemical energy and thermal energy. Examples of inducible system include tetracycline inducible promoters (Tet-On or Tet-Off), small molecule two-hybrid transcription activations systems (FKBP, ABA, etc), or light inducible systems (Phytochrome, LOV domains, or cryptochrome). In one embodiment, the CRISPR enzyme may be a part of a Light Inducible Transcriptional Effector (LITE) to direct changes in transcriptional activity in a sequence-specific manner. The components of a light may include a CRISPR enzyme, a light-responsive cytochrome heterodimer (e.g. from *Arabidopsis thaliana*), and a transcriptional activation/repression domain. Further examples of inducible DNA binding proteins and methods for their use are provided in U.S. 61/736,465 and U.S. 61/721,283 and WO 2014/018423 and U.S. Pat. No. 8,889,418, U.S. Pat. No. 8,895,308, US20140186919, US20140242700, US20140273234, US20140335620, WO2014093635, which is hereby incorporated by reference in its entirety.

[0796] In some aspects, the invention provides methods comprising delivering one or more polynucleotides, such as or one or more vectors as described herein, one or more transcripts thereof, and/or one or more proteins transcribed therefrom, to a host cell. In some aspects, the invention further provides cells produced by such methods, and organisms

(such as animals, plants, or fungi) comprising or produced from such cells. In some embodiments, a nucleic acid-targeting effector protein in combination with (and optionally complexed with) a guide RNA is delivered to a cell. Conventional viral and non-viral based gene transfer methods can be used to introduce nucleic acids in mammalian cells or target tissues. Such methods can be used to administer nucleic acids encoding components of a nucleic acid-targeting system to cells in culture, or in a host organism. Non-viral vector delivery systems include DNA plasmids, RNA (e.g. a transcript of a vector described herein), naked nucleic acid, and nucleic acid complexed with a delivery vehicle, such as a liposome. Viral vector delivery systems include DNA and RNA viruses, which have either episomal or integrated genomes after delivery to the cell. For a review of gene therapy procedures, see Anderson, *Science* 256:808-813 (1992); Nabel & Felgner, *TIBTECH* 11:211-217 (1993); Mitani & Caskey, *TIBTECH* 11:162-166 (1993); Dillon, *TIBTECH* 11:167-175 (1993); Miller, *Nature* 357:455-460 (1992); Van Brunt, *Biotechnology* 6(10):1149-1154 (1988); Vigne, *Restorative Neurology and Neuroscience* 8:35-36 (1995); Kremer & Perricaudet, *British Medical Bulletin* 51(1):31-44 (1995); Haddada et al., in *Current Topics in Microbiology and Immunology*, Doerfler and Bohm (eds) (1995); and Yu et al., *Gene Therapy* 1:13-26 (1994).

[0797] Methods of non-viral delivery of nucleic acids include lipofection, nucleofection, microinjection, biolistics, virosomes, liposomes, immunoliposomes, polycation or lipid:nucleic acid conjugates, naked DNA, artificial virions, and agent-enhanced uptake of DNA. Lipofection is described in e.g., U.S. Pat. Nos. 5,049,386, 4,946,787; and 4,897,355 and lipofection reagents are sold commercially (e.g., Transfectam™ and Lipofectin™). Cationic and neutral lipids that are suitable for efficient receptor-recognition lipofection of polynucleotides include those of Felgner, WO 91/17424; WO 91/16024. Delivery can be to cells (e.g. in vitro or ex vivo administration) or target tissues (e.g. in vivo administration).

[0798] The preparation of lipid:nucleic acid complexes, including targeted liposomes such as immunolipid complexes, is well known to one of skill in the art (see, e.g., Crystal, *Science* 270:404-410 (1995); Blaese et al., *Cancer Gene Ther.* 2:291-297 (1995); Behr et al., *Bioconjugate Chem.* 5:382-389 (1994); Remy et al., *Bioconjugate Chem.* 5:647-654 (1994); Gao et al., *Gene Therapy* 2:710-722 (1995); Ahmad et al., *Cancer Res.* 52:4817-4820 (1992); U.S. Pat. Nos. 4,186,183, 4,217,344, 4,235,871, 4,261,975, 4,485,054, 4,501,728, 4,774,085, 4,837,028, and 4,946,787).

[0799] The use of RNA or DNA viral based systems for the delivery of nucleic acids takes advantage of highly evolved processes for targeting a virus to specific cells in the body and trafficking the viral payload to the nucleus. Viral vectors can be administered directly to patients (in vivo) or they can be used to treat cells in vitro, and the modified cells may optionally be administered to patients (ex vivo). Conventional viral based systems could include retroviral, lentivirus, adenoviral, adeno-associated and herpes simplex virus vectors for gene transfer. Integration in the host genome is possible with the retrovirus, lentivirus, and adeno-associated virus gene transfer methods, often resulting in long term expression of the inserted transgene. Additionally, high transduction efficiencies have been observed in many different cell types and target tissues.

[0800] The tropism of a retrovirus can be altered by incorporating foreign envelope proteins, expanding the potential target population of target cells. Lentiviral vectors are retroviral vectors that are able to transduce or infect non-dividing cells and typically produce high viral titers. Selection of a retroviral gene transfer system would therefore depend on the target tissue. Retroviral vectors are comprised of cis-acting long terminal repeats with packaging capacity for up to 6-10 kb of foreign sequence. The minimum cis-acting LTRs are sufficient for replication and packaging of the vectors, which are then used to integrate the therapeutic gene into the target cell to provide permanent transgene expression. Widely used retroviral vectors include those based upon murine leukemia virus (MuLV), gibbon ape leukemia virus (GaLV), Simian Immuno deficiency virus (SIV), human immuno deficiency virus (HIV), and combinations thereof (see, e.g., Buchscher et al., *J. Virol.* 66:2731-2739 (1992); Johann et al., *J. Virol.* 66:1635-1640 (1992); Sommerfelt et al., *Virol.* 176:58-59 (1990); Wilson et al., *J. Virol.* 63:2374-2378 (1989); Miller et al., *J. Virol.* 65:2220-2224 (1991); PCT/US94/05700). In applications where transient expression is preferred, adenoviral based systems may be used. Adenoviral based vectors are capable of very high transduction efficiency in many cell types and do not require cell division. With such vectors, high titer and levels of expression have been obtained. This vector can be produced in large quantities in a relatively simple system. Adeno-associated virus ("AAV") vectors may also be used to transduce cells with target nucleic acids, e.g., in the in vitro production of nucleic acids and peptides, and for in vivo and ex vivo gene therapy procedures (see, e.g., West et al., *Virology* 160:38-47 (1987); U.S. Pat. No. 4,797,368; WO 93/24641; Kotin, *Human Gene Therapy* 5:793-801 (1994); Muzyczka, *J. Clin. Invest.* 94:1351 (1994). Construction of recombinant AAV vectors are described in a number of publications, including U.S. Pat. No. 5,173,414; Tratschin et al., *Mol. Cell. Biol.* 5:3251-3260 (1985); Tratschin, et al., *Mol. Cell. Biol.* 4:2072-2081 (1984); Hermonat & Muzyczka, *PNAS* 81:6466-6470 (1984); and Samulski et al., *J. Virol.* 63:03822-3828 (1989).

Models of Genetic and Epigenetic Conditions

[0801] A method of the invention may be used to create a plant, an animal or cell that may be used to model and/or study genetic or epigenetic conditions of interest, such as a through a model of mutations of interest or a disease model. As used herein, "disease" refers to a disease, disorder, or indication in a subject. For example, a method of the invention may be used to create an animal or cell that comprises a modification in one or more nucleic acid sequences associated with a disease, or a plant, animal or cell in which the expression of one or more nucleic acid sequences associated with a disease are altered. Such a nucleic acid sequence may encode a disease associated protein sequence or may be a disease associated control sequence. Accordingly, it is understood that in embodiments of the invention, a plant, subject, patient, organism or cell can be a non-human subject, patient, organism or cell. Thus, the invention provides a plant, animal or cell, produced by the present methods, or a progeny thereof. The progeny may be a clone of the produced plant or animal, or may result from sexual reproduction by crossing with other individuals of the same species to introgress further desirable traits into their offspring. The cell may be in vivo or ex vivo in the

cases of multicellular organisms, particularly animals or plants. In the instance where the cell is in cultured, a cell line may be established if appropriate culturing conditions are met and preferably if the cell is suitably adapted for this purpose (for instance a stem cell). Bacterial cell lines produced by the invention are also envisaged. Hence, cell lines are also envisaged.

[0802] In some methods, the disease model can be used to study the effects of mutations, or more general altered, such as reduced, expression of genes or gene products on the animal or cell and development and/or progression of the disease using measures commonly used in the study of the disease. Alternatively, such a disease model is useful for studying the effect of a pharmaceutically active compound on the disease.

[0803] In some methods, the disease model can be used to assess the efficacy of a potential gene therapy strategy. That is, a disease-associated gene or polynucleotide can be modified such that the disease development and/or progression is inhibited or reduced. In particular, the method comprises modifying a disease-associated gene or polynucleotide such that an altered protein is produced and, as a result, the animal or cell has an altered response. Accordingly, in some methods, a genetically modified animal may be compared with an animal predisposed to development of the disease such that the effect of the gene therapy event may be assessed.

[0804] In another embodiment, this invention provides a method of developing a biologically active agent that modulates a cell signaling event associated with a disease gene. The method comprises contacting a test compound with a cell comprising one or more vectors that drive expression of one or more of a CRISPR enzyme, and a direct repeat sequence linked to a guide sequence; and detecting a change in a readout that is indicative of a reduction or an augmentation of a cell signaling event associated with, e.g., a mutation in a disease gene contained in the cell.

[0805] A cell model or animal model can be constructed in combination with the method of the invention for screening a cellular function change. Such a model may be used to study the effects of a genome sequence modified by the CRISPR complex of the invention on a cellular function of interest. For example, a cellular function model may be used to study the effect of a modified genome sequence on intracellular signaling or extracellular signaling. Alternatively, a cellular function model may be used to study the effects of a modified genome sequence on sensory perception. In some such models, one or more genome sequences associated with a signaling biochemical pathway in the model are modified.

[0806] Several disease models have been specifically investigated. These include de novo autism risk genes CHD8, KATNAL2, and SCN2A; and the syndromic autism (Angelman Syndrome) gene UBE3A. These genes and resulting autism models are of course preferred, but serve to show the broad applicability of the invention across genes and corresponding models.

[0807] An altered expression of one or more genome sequences associated with a signalling biochemical pathway can be determined by assaying for a difference in the mRNA levels of the corresponding genes between the test model cell and a control cell, when they are contacted with a candidate agent. Alternatively, the differential expression of the sequences associated with a signaling biochemical path-

way is determined by detecting a difference in the level of the encoded polypeptide or gene product.

[0808] To assay for an agent-induced alteration in the level of mRNA transcripts or corresponding polynucleotides, nucleic acid contained in a sample is first extracted according to standard methods in the art. For instance, mRNA can be isolated using various lytic enzymes or chemical solutions according to the procedures set forth in Sambrook et al. (1989), or extracted by nucleic-acid-binding resins following the accompanying instructions provided by the manufacturers. The mRNA contained in the extracted nucleic acid sample is then detected by amplification procedures or conventional hybridization assays (e.g. Northern blot analysis) according to methods widely known in the art or based on the methods exemplified herein.

[0809] For purpose of this invention, amplification means any method employing a primer and a polymerase capable of replicating a target sequence with reasonable fidelity. Amplification may be carried out by natural or recombinant DNA polymerases such as TaqGold™, T7 DNA polymerase, Klenow fragment of *E. coli* DNA polymerase, and reverse transcriptase. A preferred amplification method is PCR. In particular, the isolated RNA can be subjected to a reverse transcription assay that is coupled with a quantitative polymerase chain reaction (RT-PCR) in order to quantify the expression level of a sequence associated with a signaling biochemical pathway.

[0810] Detection of the gene expression level can be conducted in real time in an amplification assay. In one aspect, the amplified products can be directly visualized with fluorescent DNA-binding agents including but not limited to DNA intercalators and DNA groove binders. Because the amount of the intercalators incorporated into the double-stranded DNA molecules is typically proportional to the amount of the amplified DNA products, one can conveniently determine the amount of the amplified products by quantifying the fluorescence of the intercalated dye using conventional optical systems in the art. DNA-binding dye suitable for this application include SYBR green, SYBR blue, DAPI, propidium iodine, Hoeste, SYBR gold, ethidium bromide, acridines, proflavine, acridine orange, acriflavine, fluorocoumanin, ellipticine, daunomycin, chloroquine, distamycin D, chromomycin, homidium, mithramycin, ruthenium polypyridyls, anthramycin, and the like.

[0811] In another aspect, other fluorescent labels such as sequence specific probes can be employed in the amplification reaction to facilitate the detection and quantification of the amplified products. Probe-based quantitative amplification relies on the sequence-specific detection of a desired amplified product. It utilizes fluorescent, target-specific probes (e.g., TaqMan® probes) resulting in increased specificity and sensitivity. Methods for performing probe-based quantitative amplification are well established in the art and are taught in U.S. Pat. No. 5,210,015.

[0812] In yet another aspect, conventional hybridization assays using hybridization probes that share sequence homology with sequences associated with a signaling biochemical pathway can be performed. Typically, probes are allowed to form stable complexes with the sequences associated with a signaling biochemical pathway contained within the biological sample derived from the test subject in a hybridization reaction. It will be appreciated by one of skill in the art that where antisense is used as the probe nucleic acid, the target polynucleotides provided in the sample are

chosen to be complementary to sequences of the antisense nucleic acids. Conversely, where the nucleotide probe is a sense nucleic acid, the target polynucleotide is selected to be complementary to sequences of the sense nucleic acid.

[0813] Hybridization can be performed under conditions of various stringency. Suitable hybridization conditions for the practice of the present invention are such that the recognition interaction between the probe and sequences associated with a signaling biochemical pathway is both sufficiently specific and sufficiently stable. Conditions that increase the stringency of a hybridization reaction are widely known and published in the art. See, for example, (Sambrook, et al., (1989); Nonradioactive In Situ Hybridization Application Manual, Boehringer Mannheim, second edition). The hybridization assay can be formed using probes immobilized on any solid support, including but are not limited to nitrocellulose, glass, silicon, and a variety of gene arrays. A preferred hybridization assay is conducted on high-density gene chips as described in U.S. Pat. No. 5,445,934.

[0814] For a convenient detection of the probe-target complexes formed during the hybridization assay, the nucleotide probes are conjugated to a detectable label. Detectable labels suitable for use in the present invention include any composition detectable by photochemical, biochemical, spectroscopic, immunochemical, electrical, optical or chemical means. A wide variety of appropriate detectable labels are known in the art, which include fluorescent or chemiluminescent labels, radioactive isotope labels, enzymatic or other ligands. In preferred embodiments, one will likely desire to employ a fluorescent label or an enzyme tag, such as digoxigenin, β -galactosidase, urease, alkaline phosphatase or peroxidase, avidin/biotin complex.

[0815] The detection methods used to detect or quantify the hybridization intensity will typically depend upon the label selected above. For example, radiolabels may be detected using photographic film or a phosphorimager. Fluorescent markers may be detected and quantified using a photodetector to detect emitted light. Enzymatic labels are typically detected by providing the enzyme with a substrate and measuring the reaction product produced by the action of the enzyme on the substrate; and finally colorimetric labels are detected by simply visualizing the colored label.

[0816] An agent-induced change in expression of sequences associated with a signaling biochemical pathway can also be determined by examining the corresponding gene products. Determining the protein level typically involves a) contacting the protein contained in a biological sample with an agent that specifically bind to a protein associated with a signaling biochemical pathway; and (b) identifying any agent:protein complex so formed. In one aspect of this embodiment, the agent that specifically binds a protein associated with a signaling biochemical pathway is an antibody, preferably a monoclonal antibody.

[0817] The reaction is performed by contacting the agent with a sample of the proteins associated with a signaling biochemical pathway derived from the test samples under conditions that will allow a complex to form between the agent and the proteins associated with a signaling biochemical pathway. The formation of the complex can be detected directly or indirectly according to standard procedures in the art. In the direct detection method, the agents are supplied with a detectable label and unreacted agents may be removed from the complex; the amount of remaining label

thereby indicating the amount of complex formed. For such method, it is preferable to select labels that remain attached to the agents even during stringent washing conditions. It is preferable that the label does not interfere with the binding reaction. In the alternative, an indirect detection procedure may use an agent that contains a label introduced either chemically or enzymatically. A desirable label generally does not interfere with binding or the stability of the resulting agent:polypeptide complex. However, the label is typically designed to be accessible to an antibody for an effective binding and hence generating a detectable signal.

[0818] A wide variety of labels suitable for detecting protein levels are known in the art. Non-limiting examples include radioisotopes, enzymes, colloidal metals, fluorescent compounds, bioluminescent compounds, and chemiluminescent compounds.

[0819] The amount of agent:polypeptide complexes formed during the binding reaction can be quantified by standard quantitative assays. As illustrated above, the formation of agent:polypeptide complex can be measured directly by the amount of label remained at the site of binding. In an alternative, the protein associated with a signaling biochemical pathway is tested for its ability to compete with a labeled analog for binding sites on the specific agent. In this competitive assay, the amount of label captured is inversely proportional to the amount of protein sequences associated with a signaling biochemical pathway present in a test sample.

[0820] A number of techniques for protein analysis based on the general principles outlined above are available in the art. They include but are not limited to radioimmunoassays, ELISA (enzyme linked immunoradiometric assays), "sandwich" immunoassays, immunoradiometric assays, in situ immunoassays (using e.g., colloidal gold, enzyme or radioisotope labels), western blot analysis, immunoprecipitation assays, immunofluorescent assays, and SDS-PAGE.

[0821] Antibodies that specifically recognize or bind to proteins associated with a signaling biochemical pathway are preferable for conducting the aforementioned protein analyses. Where desired, antibodies that recognize a specific type of post-translational modifications (e.g., signaling biochemical pathway inducible modifications) can be used. Post-translational modifications include but are not limited to glycosylation, lipidation, acetylation, and phosphorylation. These antibodies may be purchased from commercial vendors. For example, anti-phosphotyrosine antibodies that specifically recognize tyrosine-phosphorylated proteins are available from a number of vendors including Invitrogen and Perkin Elmer. Antiphosphotyrosine antibodies are particularly useful in detecting proteins that are differentially phosphorylated on their tyrosine residues in response to an ER stress. Such proteins include but are not limited to eukaryotic translation initiation factor 2 alpha (eIF-2 α). Alternatively, these antibodies can be generated using conventional polyclonal or monoclonal antibody technologies by immunizing a host animal or an antibody-producing cell with a target protein that exhibits the desired post-translational modification.

[0822] In practicing the subject method, it may be desirable to discern the expression pattern of a protein associated with a signaling biochemical pathway in different bodily tissue, in different cell types, and/or in different subcellular structures. These studies can be performed with the use of tissue-specific, cell-specific or subcellular struc-

ture specific antibodies capable of binding to protein markers that are preferentially expressed in certain tissues, cell types, or subcellular structures.

[0823] An altered expression of a gene associated with a signaling biochemical pathway can also be determined by examining a change in activity of the gene product relative to a control cell. The assay for an agent-induced change in the activity of a protein associated with a signaling biochemical pathway will depend on the biological activity and/or the signal transduction pathway that is under investigation. For example, where the protein is a kinase, a change in its ability to phosphorylate the downstream substrate(s) can be determined by a variety of assays known in the art. Representative assays include but are not limited to immunoblotting and immunoprecipitation with antibodies such as anti-phosphotyrosine antibodies that recognize phosphorylated proteins. In addition, kinase activity can be detected by high throughput chemiluminescent assays such as AlphaScreen™ (available from Perkin Elmer) and eTag™ assay (Chan-Hui, et al. (2003) *Clinical Immunology* 111: 162-174).

[0824] Where the protein associated with a signaling biochemical pathway is part of a signaling cascade leading to a fluctuation of intracellular pH condition, pH sensitive molecules such as fluorescent pH dyes can be used as the reporter molecules. In another example where the protein associated with a signaling biochemical pathway is an ion channel, fluctuations in membrane potential and/or intracellular ion concentration can be monitored. A number of commercial kits and high-throughput devices are particularly suited for a rapid and robust screening for modulators of ion channels. Representative instruments include FLIPR™ (Molecular Devices, Inc.) and VIPR (Aurora Biosciences). These instruments are capable of detecting reactions in over 1000 sample wells of a microplate simultaneously, and providing real-time measurement and functional data within a second or even a minisecond.

[0825] In practicing any of the methods disclosed herein, a suitable vector can be introduced to a cell or an embryo via one or more methods known in the art, including without limitation, microinjection, electroporation, sonoporation, biolistics, calcium phosphate-mediated transfection, cationic transfection, liposome transfection, dendrimer transfection, heat shock transfection, nucleofection transfection, magnetofection, lipofection, impalefection, optical transfection, proprietary agent-enhanced uptake of nucleic acids, and delivery via liposomes, immunoliposomes, virosomes, or artificial virions. In some methods, the vector is introduced into an embryo by microinjection. The vector or vectors may be microinjected into the nucleus or the cytoplasm of the embryo. In some methods, the vector or vectors may be introduced into a cell by nucleofection.

[0826] The target polynucleotide of a CRISPR complex can be any polynucleotide endogenous or exogenous to the eukaryotic cell. For example, the target polynucleotide can be a polynucleotide residing in the nucleus of the eukaryotic cell. The target polynucleotide can be a sequence coding or encoded by a gene product (e.g., a protein) or a non-coding (RNA) sequence (e.g., a regulatory polynucleotide or a junk DNA).

[0827] Examples of target polynucleotides include a sequence associated with a signaling biochemical pathway, e.g., a signaling biochemical pathway-associated gene or polynucleotide. Examples of target polynucleotides include

a disease associated gene or polynucleotide. A “disease-associated” gene or polynucleotide refers to any gene or polynucleotide which is yielding transcription or translation products at an abnormal level or in an abnormal form in cells derived from a disease-affected tissues compared with tissues or cells of a non disease control. It may be a gene that becomes expressed at an abnormally high level; it may be a gene that becomes expressed at an abnormally low level, where the altered expression correlates with the occurrence and/or progression of the disease. A disease-associated gene also refers to a gene possessing mutation(s) or genetic variation that is directly responsible or is in linkage disequilibrium with a gene(s) that is responsible for the etiology of a disease. The transcribed or translated products may be known or unknown, and may be at a normal or abnormal level.

[0828] The target polynucleotide of a CRISPR complex can be any polynucleotide endogenous or exogenous to the eukaryotic cell. For example, the target polynucleotide can be a polynucleotide residing in the nucleus of the eukaryotic cell. The target polynucleotide can be a sequence coding a gene product (e.g., a protein) or a non-coding sequence (e.g., a regulatory polynucleotide or a junk DNA). Without wishing to be bound by theory, it is believed that the target sequence may be associated with a PAM (protospacer adjacent motif); that is, a short sequence recognized by the CRISPR complex. The precise sequence and length requirements for the PAM differ depending on the CRISPR enzyme used, but PAMs are typically 2-5 base pair sequences adjacent the protospacer (that is, the target sequence) Examples of PAM sequences are given in the examples section below, and the skilled person will be able to identify further PAM sequences for use with a given CRISPR enzyme.

[0829] The target polynucleotide of a CRISPR complex may include a number of disease associated genes and polynucleotides as well as signaling biochemical pathway-associated genes and polynucleotides as listed in U.S. provisional patent applications 61/736,527 and 61/748,427 both entitled SYSTEMS METHODS AND COMPOSITIONS FOR SEQUENCE MANIPULATION filed on Dec. 12, 2012 and Jan. 2, 2013, respectively, and PCT Application PCT/US 2013/074667, entitled DELIVERY, ENGINEERING AND OPTIMIZATION OF SYSTEMS, METHODS AND COMPOSITIONS FOR SEQUENCE MANIPULATION AND THERAPEUTIC APPLICATIONS, filed Dec. 12, 2013, the contents of all of which are herein incorporated by reference in their entirety.

[0830] Examples of target polynucleotides include a sequence associated with a signaling biochemical pathway, e.g., a signaling biochemical pathway-associated gene or polynucleotide. Examples of target polynucleotides include a disease associated gene or polynucleotide. A “disease-associated” gene or polynucleotide refers to any gene or polynucleotide which is yielding transcription or translation products at an abnormal level or in an abnormal form in cells derived from a disease-affected tissues compared with tissues or cells of a non disease control. It may be a gene that becomes expressed at an abnormally high level; it may be a gene that becomes expressed at an abnormally low level, where the altered expression correlates with the occurrence and/or progression of the disease. A disease-associated gene also refers to a gene possessing mutation(s) or genetic variation that is directly responsible or is in linkage disequi-

librium with a gene(s) that is responsible for the etiology of a disease. The transcribed or translated products may be known or unknown, and may be at a normal or abnormal level.

Genome/Transcriptome Wide Knock-Out or Knock-Down Screening

[0831] The CRISPR effector protein complexes described herein can be used to perform efficient and cost effective functional genomic screens. Such screens can utilize CRISPR effector protein based genome wide libraries. Such screens and libraries can provide for determining the function of genes, cellular pathways genes are involved in, and how any alteration in gene expression can result in a particular biological process. An advantage of the present invention is that the CRISPR system avoids off-target binding and its resulting side effects. This is achieved using systems arranged to have a high degree of sequence specificity for the target DNA/RNA. In preferred embodiments of the invention, the CRISPR effector protein complexes are Group 29 or Group 30 effector protein complexes.

[0832] In embodiments of the invention, a genome wide library may comprise a plurality of Group 29 or Group 30 effector protein guide RNAs, as described herein, comprising guide sequences that are capable of targeting a plurality of target sequences in a plurality of loci in a population of eukaryotic cells. The population of cells may be a population of embryonic stem (ES) cells. The target sequence in the genomic locus may be a non-coding sequence. The non-coding sequence may be an intron, regulatory sequence, splice site, 3' UTR, 5' UTR, or polyadenylation signal. Gene function of one or more gene products may be altered by said targeting. The targeting may result in a knockout of gene function. The targeting of a gene product may comprise more than one guide RNA. A gene product may be targeted by 2, 3, 4, 5, 6, 7, 8, 9, or 10 guide RNAs, for example 3 to 4 per gene. Off-target modifications may be minimized by exploiting the staggered double strand breaks generated by Group 29 or Group 30 effector protein complexes or by utilizing methods analogous to those used in CRISPR-Cas9 systems (See, e.g., DNA targeting specificity of RNA-guided Cas9 nucleases. Hsu, P., Scott, D., Weinstein, J., Ran, F. A., Konermann, S., Agarwala, V., Li, Y., Fine, E., Wu, X., Shalem, O., Cradick, T. J., Marraffini, L. A., Bao, G., & Zhang, F. *Nat Biotechnol* doi:10.1038/nbt.2647 (2013)), incorporated herein by reference. The targeting may be of about 100 or more sequences. The targeting may be of about 1000 or more sequences. The targeting may be of about 20,000 or more sequences. The targeting may be of the entire genome. The targeting may be of a panel of target sequences focused on a relevant or desirable pathway. The pathway may be an immune pathway. The pathway may be a cell division pathway.

[0833] One aspect of the invention comprehends a genome or transcriptome wide library that may comprise a plurality of Group 29 or Group 30 guide RNAs that may comprise guide sequences that are capable of targeting a plurality of target sequences in a plurality of loci, wherein said targeting results in a knockout or knockdown of gene function. This library may potentially comprise guide RNAs that target each and every gene in the genome of an organism.

[0834] In some embodiments of the invention the organism or subject is a eukaryote (including mammal including human) or a non-human eukaryote or a non-human animal or a non-human mammal. In some embodiments, the organism or subject is a non-human animal, and may be an arthropod, for example, an insect, or may be a nematode. In some methods of the invention the organism or subject is a plant. In some methods of the invention the organism or

subject is a mammal or a non-human mammal. A non-human mammal may be for example a rodent (preferably a mouse or a rat), an ungulate, or a primate. In some methods of the invention the organism or subject is algae, including microalgae, or is a fungus.

[0835] The knockout or knockdown of gene function may comprise: introducing into each cell in the population of cells a vector system of one or more vectors comprising an engineered, non-naturally occurring Group 29 or Group 30 effector protein of Group 29 or Group 30 effector protein system comprising I. a Group 29 or Group 30 effector protein, and II. one or more guide RNAs, wherein components I and II may be same or on different vectors of the system, integrating components I and II into each cell, wherein the guide sequence targets a unique gene, or transcript thereof, or other RNA sequence in each cell, wherein the Group 29 or Group 30 effector protein is operably linked to a regulatory element, wherein when transcribed, the guide RNA comprising the guide sequence directs sequence-specific binding of the Group 29 or Group 30 effector protein Group 29 or Group 30 effector protein system to a target sequence in the loci of the unique gene (transcript) or other RNA as described herein elsewhere, inducing cleavage of the loci by the Group 29 or Group 30 effector protein, and confirming different knockout/knockdown mutations in a plurality of unique targets in each cell of the population of cells thereby generating a gene knockout/knockdown cell library. The invention comprehends that the population of cells is a population of eukaryotic cells, and in a preferred embodiment, the population of cells is a population of embryonic stem (ES) cells.

[0836] The one or more vectors may be plasmid vectors. The vector may be a single vector comprising a Group 29 or Group 30 effector protein, a sgRNA, and optionally, a selection marker into target cells. Not being bound by a theory, the ability to simultaneously deliver a Group 29 or Group 30 effector protein and sgRNA through a single vector enables application to any cell type of interest, without the need to first generate cell lines that express the Group 29 or Group 30 effector protein. The regulatory element may be an inducible promoter. The inducible promoter may be a doxycycline inducible promoter. In some methods of the invention the expression of the guide sequence is under the control of the T7 promoter and is driven by the expression of T7 polymerase. The confirming of different knockout/knockdown mutations may be by whole exome/transcriptome sequencing. The knockout/knockdown mutation may be achieved in 100 or more unique genes. The knockout/knockdown mutation may be achieved in 1000 or more unique genes. The knockout mutation may be achieved in 20,000 or more unique genes. The knockout/knockdown mutation may be achieved in the entire genome/transcriptome. The knockout/knockdown of gene function may be achieved in a plurality of unique genes which function in a particular physiological pathway or condition. The pathway or condition may be an immune pathway or condition. The pathway or condition may be a cell division pathway or condition.

[0837] The invention also provides kits that comprise the genome/transcriptome (or other RNAs, as described herein elsewhere) wide libraries mentioned herein. The kit may comprise a single container comprising vectors or plasmids comprising the library of the invention. The kit may also comprise a panel comprising a selection of unique Group 29 or Group 30 effector protein Group 29 or Group 30 effector protein system guide RNAs comprising guide sequences from the library of the invention, wherein the selection is indicative of a particular physiological condition. The invention comprehends that the targeting is of about 100 or more sequences, about 1000 or more sequences or about 20,000 or

more sequences or the entire genome/transcriptome. Furthermore, a panel of target sequences may be focused on a relevant or desirable pathway, such as an immune pathway or cell division.

[0838] In an additional aspect of the invention, the Group 29 or Group 30 effector protein may comprise one or more mutations and may be used as a generic DNA/RNA binding protein with or without fusion to a functional domain. The mutations may be artificially introduced mutations or gain- or loss-of-function mutations. The mutations have been characterized as described herein. In one aspect of the invention, the functional domain may be a transcriptional activation domain, which may be VP64. In other aspects of the invention, the functional domain may be a transcriptional repressor domain, which may be KRAB or SID4X. Other aspects of the invention relate to the mutated Group 29 or Group 30 effector protein being fused to domains which include but are not limited to a transcriptional activator, repressor, a recombinase, a transposase, a histone remodeler, a demethylase, a DNA methyltransferase, a cryptochrome, a light inducible/controllable domain or a chemically inducible/controllable domain. Some methods of the invention can include inducing expression of targeted genes. In one embodiment, inducing expression by targeting a plurality of target sequences in a plurality of genomic loci in a population of eukaryotic cells is by use of a functional domain.

[0839] Useful in the practice of the instant invention utilizing Group 29 or Group 30 effector protein complexes are methods used in CRISPR-Cas9 systems and reference is made to:

[0840] Genome-Scale CRISPR-Cas9 Knockout Screening in Human Cells. Shalem, O., Sanjana, N E., Hartenian, E., Shi, X., Scott, D A., Mikkelsen, T., Heckl, D., Ebert, B L., Root, D E., Doench, J G., Zhang, F. Science December 12, (2013). [Epub ahead of print]; Published in final edited form as: Science. 2014 Jan. 3; 343(6166): 84-87.

[0841] Shalem et al. involves a new way to interrogate gene function on a genome-wide scale. Their studies showed that delivery of a genome-scale CRISPR-Cas9 knockout (GeCKO) library targeted 18,080 genes with 64,751 unique guide sequences enabled both negative and positive selection screening in human cells. First, the authors showed use of the GeCKO library to identify genes essential for cell viability in cancer and pluripotent stem cells. Next, in a melanoma model, the authors screened for genes whose loss is involved in resistance to vemurafenib, a therapeutic that inhibits mutant protein kinase BRAF. Their studies showed that the highest-ranking candidates included previously validated genes NF1 and MED12 as well as novel hits NF2, CUL3, TADA2B, and TADA1. The authors observed a high level of consistency between independent guide RNAs targeting the same gene and a high rate of hit confirmation, and thus demonstrated the promise of genome-scale screening with Cas9.

[0842] Reference is also made to US patent publication number US20140357530; and PCT Patent Publication WO2014093701, hereby incorporated herein by reference.

Functional Alteration and Screening

[0843] In another aspect, the present invention provides for a method of functional evaluation and screening of genes/transcripts (or other RNAs). The use of the CRISPR system of the present invention to precisely deliver functional domains, to activate or repress genes or to alter epigenetic state by precisely altering the methylation site on a specific locus of interest, can be with one or more guide RNAs applied to a single cell or population of cells or with a library applied to genome in a pool of cells ex vivo or in vivo comprising the administration or expression of a library

comprising a plurality of guide RNAs (sgRNAs) and wherein the screening further comprises use of a Group 29 or Group 30 effector protein, wherein the CRISPR complex comprising the Group 29 or Group 30 effector protein is modified to comprise a heterologous functional domain. In an aspect the invention provides a method for screening a genome comprising the administration to a host or expression in a host in vivo of a library. In an aspect the invention provides a method as herein discussed further comprising an activator administered to the host or expressed in the host. In an aspect the invention provides a method as herein discussed wherein the activator is attached to a Group 29 or Group 30 effector protein. In an aspect the invention provides a method as herein discussed wherein the activator is attached to the N terminus or the C terminus of the Group 29 or Group 30 effector protein. In an aspect the invention provides a method as herein discussed wherein the activator is attached to a sgRNA loop. In an aspect the invention provides a method as herein discussed further comprising a repressor administered to the host or expressed in the host. In an aspect the invention provides a method as herein discussed, wherein the screening comprises affecting and detecting gene activation, gene inhibition, or cleavage in the locus.

[0844] In an aspect, the invention provides efficient on-target activity and minimizes off target activity. In an aspect, the invention provides efficient on-target cleavage by Group 29 or Group 30 effector protein and minimizes off-target cleavage by the Group 29 or Group 30 effector protein. In an aspect, the invention provides guide specific binding of Group 29 or Group 30 effector protein at a gene locus without DNA/RNA cleavage. Accordingly, in an aspect, the invention provides target-specific gene regulation. In an aspect, the invention provides guide specific binding of Group 29 or Group 30 effector protein at a gene locus without DNA/RNA cleavage. Accordingly, in an aspect, the invention provides for cleavage at one gene locus and gene regulation at a different gene locus using a single Group 29 or Group 30 effector protein. In an aspect, the invention provides orthogonal activation and/or inhibition and/or cleavage of multiple targets using one or more Group 29 or Group 30 effector protein and/or enzyme.

[0845] In an aspect the invention provides a method as herein discussed, wherein the host is a eukaryotic cell. In an aspect the invention provides a method as herein discussed, wherein the host is a mammalian cell. In an aspect the invention provides a method as herein discussed, wherein the host is a non-human eukaryote. In an aspect the invention provides a method as herein discussed, wherein the non-human eukaryote is a non-human mammal. In an aspect the invention provides a method as herein discussed, wherein the non-human mammal is a mouse. An aspect the invention provides a method as herein discussed comprising the delivery of the Group 29 or Group 30 effector protein complexes or component(s) thereof or nucleic acid molecule (s) coding therefor, wherein said nucleic acid molecule(s) are operatively linked to regulatory sequence(s) and expressed in vivo. In an aspect the invention provides a method as herein discussed wherein the expressing in vivo is via a lentivirus, an adenovirus, or an AAV. In an aspect the invention provides a method as herein discussed wherein the delivery is via a particle, a particle, a lipid or a cell penetrating peptide (CPP).

[0846] In an aspect the invention provides a pair of CRISPR complexes comprising Group 29 or Group 30 effector protein, each comprising a guide RNA (sgRNA) comprising a guide sequence capable of hybridizing to a target sequence in a locus of interest in a cell, wherein at least one loop of each sgRNA is modified by the insertion of distinct RNA sequence(s) that bind to one or more adaptor

proteins, and wherein the adaptor protein is associated with one or more functional domains, wherein each sgRNA of each Group 29 or Group 30 effector protein complex comprises a functional domain having a DNA and/or RNA cleavage activity. In an aspect the invention provides paired Group 29 and/or Group 30 effector protein complexes as herein-discussed, wherein the DNA or RNA cleavage activity is due to a FokI nuclease.

[0847] In an aspect the invention provides a library, method or complex as herein-discussed wherein the sgRNA is modified to have at least one non-coding functional loop, e.g., wherein the at least one non-coding functional loop is repressive; for instance, wherein the at least one non-coding functional loop comprises Alu.

[0848] In one aspect, the invention provides a method for altering or modifying expression of a gene product. The said method may comprise introducing into a cell containing and expressing a DNA/RNA molecule encoding the gene product an engineered, non-naturally occurring CRISPR system comprising a Group 29 or Group 30 effector protein and guide RNA that targets the DNA/RNA molecule, whereby the guide RNA targets the DNA/RNA molecule encoding the gene product and the Group 29 or Group 30 effector protein cleaves the DNA/RNA molecule encoding the gene product, whereby expression of the gene product is altered; and, wherein the Group 29 or Group 30 effector protein and the guide RNA do not naturally occur together. The invention comprehends the guide RNA comprising a guide sequence linked to a direct repeat sequence. The invention further comprehends the Group 29 or Group 30 effector protein being codon optimized for expression in a Eukaryotic cell. In a preferred embodiment the Eukaryotic cell is a mammalian cell and in a more preferred embodiment the mammalian cell is a human cell. In a further embodiment of the invention, the expression of the gene product is decreased.

[0849] In some embodiments, one or more functional domains are associated with the Group 29 or Group 30 effector protein. In some embodiments, one or more functional domains are associated with an adaptor protein, for example as used with the modified guides of Konnerman et al. (Nature 517, 583-588, 29 Jan. 2015). In some embodiments, one or more functional domains are associated with an dead sgRNA (dRNA). In some embodiments, a dRNA complex with active Group 29 or Group 30 effector protein directs gene regulation by a functional domain at on locus while an sgRNA directs DNA/RNA cleavage by the active Group 29 or Group 30 effector protein at another locus, for example as described analogously in CRISPR-Cas9 systems by Dahlman et al., 'Orthogonal gene control with a catalytically active Cas9 nuclease' (in press). In some embodiments, dRNAs are selected to maximize selectivity of regulation for a gene locus of interest compared to off-target regulation. In some embodiments, dRNAs are selected to maximize target gene regulation and minimize target cleavage

[0850] For the purposes of the following discussion, reference to a functional domain could be a functional domain associated with the Group 29 or Group 30 effector protein or a functional domain associated with the adaptor protein.

[0851] In some embodiments, the one or more functional domains is an NLS (Nuclear Localization Sequence) or an NES (Nuclear Export Signal). In some embodiments, the one or more functional domains is a transcriptional activation domain comprises VP64, p65, MyoD1, HSF1, RTA, SET7/9 and a histone acetyltransferase. Other references herein to activation (or activator) domains in respect of those associated with the CRISPR enzyme include any known

transcriptional activation domain and specifically VP64, p65, MyoD1, HSF1, RTA, SET7/9 or a histone acetyltransferase.

[0852] In some embodiments, the one or more functional domains is a transcriptional repressor domain. In some embodiments, the transcriptional repressor domain is a KRAB domain. In some embodiments, the transcriptional repressor domain is a NuE domain, NcoR domain, SID domain or a SID4X domain.

[0853] In some embodiments, the one or more functional domains have one or more activities comprising methylase activity, demethylase activity, transcription activation activity, transcription repression activity, transcription release factor activity, histone modification activity, RNA cleavage activity, DNA cleavage activity, DNA integration activity or nucleic acid binding activity.

[0854] Histone modifying domains are also preferred in some embodiments. Exemplary histone modifying domains are discussed below. Transposase domains, HR (Homologous Recombination) machinery domains, recombinase domains, and/or integrase domains are also preferred as the present functional domains. In some embodiments, DNA integration activity includes HR machinery domains, integrase domains, recombinase domains and/or transposase domains. Histone acetyltransferases are preferred in some embodiments.

[0855] In some embodiments, the DNA cleavage activity is due to a nuclease. In some embodiments, the nuclease comprises a FokI nuclease. See, "Dimeric CRISPR RNA-guided FokI nucleases for highly specific genome editing", Shengdar Q. Tsai, Nicolas Wyvekens, Cyd Khayter, Jennifer A. Foden, Vishal Thapar, Deepak Reyon, Mathew J. Goodwin, Martin J. Aryee, J. Keith Joung Nature Biotechnology 32(6): 569-77 (2014), relates to dimeric RNA-guided FokI Nucleases that recognize extended sequences and can edit endogenous genes with high efficiencies in human cells.

[0856] In some embodiments, the one or more functional domains is attached to the Group 29 or Group 30 effector protein so that upon binding to the sgRNA and target the functional domain is in a spatial orientation allowing for the functional domain to function in its attributed function.

[0857] In some embodiments, the one or more functional domains is attached to the adaptor protein so that upon binding of the Group 29 or Group 30 effector protein to the sgRNA and target, the functional domain is in a spatial orientation allowing for the functional domain to function in its attributed function.

[0858] In an aspect the invention provides a composition as herein discussed wherein the one or more functional domains is attached to the Group 29 or Group 30 effector protein or adaptor protein via a linker, optionally a GlySer linker, as discussed herein.

[0859] Endogenous transcriptional repression is often mediated by chromatin modifying enzymes such as histone methyltransferases (HMTs) and deacetylases (HDACs). Repressive histone effector domains are known and an exemplary list is provided below. In the exemplary table, preference was given to proteins and functional truncations of small size to facilitate efficient viral packaging (for instance via AAV). In general, however, the domains may include HDACs, histone methyltransferases (HMTs), and histone acetyltransferase (HAT) inhibitors, as well as HDAC and HMT recruiting proteins. The functional domain may be or include, in some embodiments, HDAC Effector Domains, HDAC Recruiter Effector Domains, Histone Methyltransferase (HMT) Effector Domains, Histone Methyltransferase (HMT) Recruiter Effector Domains, or Histone Acetyltransferase Inhibitor Effector Domains.

HDAC Effector Domains								
Subtype/ Complex	Name	Substrate (if known)	Modification (if known)	Organism	Full size (aa)	Selected truncation (aa)	Final size (aa)	Catalytic domain
HDAC I	HDAC8	—	—	<i>X. laevis</i>	325	1-325	325	1-272: HDAC
HDAC I	RPD3	—	—	<i>S. cerevisiae</i>	433	19-340	322 (Vannier)	19-331: HDAC
HDAC IV	MesoLo4	—	—	<i>M. luti</i>	300	1-300 (Gregoretta)	300	—
HDAC IV	HDAC11	—	—	<i>H. sapiens</i>	347	1-347 (Gao)	347	14-326: HDAC
HD2	HDT1	—	—	<i>A. thaliana</i>	245	1-211 (Wu)	211	—
SIRT I	SIRT3	H3K9Ac H4K16 Ac H3K56 Ac	—	<i>H. sapiens</i>	399	143-399 (Scher)	257	126-382: SIRT
SIRT I	HST2	—	—	<i>C. albicans</i>	331	1-331 (Hnisz)	331	—
SIRT I	CobB	—	—	<i>E. coli</i>	242	1-242 (Landry)	242	—
SIRT I	HST2	—	—	<i>S. cerevisiae</i>	357	8-298 (Wilson)	291	—
SIRT III	SIRT5	H4K8Ac H4K16 Ac	—	<i>H. sapiens</i>	310	37-310 (Gertz)	274	41-309: SIRT
SIRT III	Sir2A	—	—	<i>P. falciparum</i>	273	1-273 (Zhu)	273	19-273: SIRT
SIRT IV	SIRT6	H3K9Ac H3K56 Ac	—	<i>H. sapiens</i>	355	1-289 (Tennen)	289	35-274: SIRT

[0860] Accordingly, the repressor domains of the present invention may be selected from histone methyltransferases (HMTs), histone deacetylases (HDACs), histone acetyltransferase (HAT) inhibitors, as well as HDAC and HMT recruiting proteins.

[0861] The HDAC domain may be any of those in the table above, namely: HDAC8, RPD3, MesoLo4, HDAC11,

HDT1, SIRT3, HST2, CobB, HST2, SIRT5, Sir2A, or SIRT6.

[0862] In some embodiment, the functional domain may be a HDAC Recruiter Effector Domain. Preferred examples include those in the Table below, namely MeCP2, MBD2b, Sin3a, NcoR, SALL1, RCOR1. NcoR is exemplified in the present Examples and, although preferred, it is envisaged that others in the class will also be useful.

Table of HDAC Recruiter Effector Domains

Subtype/ Complex	Name	Substrate (if known)	Modification (if known)	Organism	Full size (aa)	Selected truncation (aa)	Final size (aa)	Catalytic domain
Sin3a	MeCP2	—	—	<i>R. norvegicus</i>	492	207-492 (Nan)	286	—
Sin3a	MBD2b	—	—	<i>H. sapiens</i>	262	45-262 (Boeke)	218	—
Sin3a	Sin3a	—	—	<i>H. sapiens</i>	1273	524-851 (Laherty)	328	627-829: HDAC1 interaction
NcoR	NcoR	—	—	<i>H. sapiens</i>	2440	420-488 (Zhang)	69	—
NuRD	SALL1	—	—	<i>M. musculus</i>	1322	1-93 (Lauberth)	93	—
CoREST	RCOR1	—	—	<i>H. sapiens</i>	482	81-300 (Gu, Ouyang)	220	—

[0863] In some embodiment, the functional domain may be a Methyltransferase (HMT) Effector Domain. Preferred examples include those in the Table below, namely NUE, vSET, EHMT2/G9A, SUV39H1, dim-5, KYP, SUV39H1,

SET4, SET1, SETD8, and TgSET8. NUE is exemplified in the present Examples and, although preferred, it is envisaged that others in the class will also be useful.

Subtype/Complex	Name	Substrate (if known)	Modification (if known)	Organism	Full size (aa)	Selected truncation (aa)	Final size (aa)	Catalytic domain
SET	NUE	H2B, H3, H4	—	<i>C. trachomatis</i>	219	1-219 (Pennini)	219	—
SET	vSET	—	H3K27me3	<i>M. bursaria</i> <i>chlorella virus</i>	119	1-119 (Mujtaba)	119	4-112: SET2
SUV39 family	EHMT2/G9A	H1.4K2, H3K9, H3K27	H3K9me1/2, HIK25me1	<i>M. musculus</i>	1263	969-1263 (Tachibana)	295	1025-1233: preSET, SET, postSET
SUV39	SUV39H1	—	H3K9me2/3	<i>H. sapiens</i>	412	79-412 (Snowden)	334	172-412: preSET, SET, postSET
Suvar3-9	dim-5	—	H3K9me3	<i>N. crassa</i>	331	1-331 (Rathert)	331	77-331: preSET, SET, postSET
Suvar3-9 (SUVH subfamily)	KYP	—	H3K9me1/2	<i>A. thaliana</i>	624	335-601	267 (Jackson)	—
Suvar3-9 (SUVR subfamily)	SUVR4	H3K9me1	H3K9me2/3	<i>A. thaliana</i>	492	180-492	313 (Thorstensen)	192-462: preSET, SET, postSET
Suvar4-20	SET4	—	H4K20me3	<i>C. elegans</i>	288	1-288 (Vielle)	288	—
SET8	SET1	—	H4K20me1	<i>C. elegans</i>	242	1-242 (Vielle)	242	—
SET8	SETD8	—	H4K20me1	<i>H. sapiens</i>	393	185-393	209 (Couture)	256-382: SET
SET8	TgSET8	—	H4K20me1/2/3	<i>T. gondii</i>	1893	1590-1893 (Sautel)	304	1749-1884: SET

[0864] In some embodiment, the functional domain may be a Histone Methyltransferase (HMT) Recruiter Effector Domain. Preferred examples include those in the Table below, namely Hp1a, PHF 19, and NIPP1.

Subtype/Complex	Name	Substrate (if known)	Modification (if known)	Organism	Full size (aa)	Selected truncation (aa)	Final size (aa)	Catalytic domain
—	Hp1a	—	H3K9me3	<i>M. musculus</i>	191	73-191	119 (Hathaway)	121-179: chromoshadow
—	PHF19	—	H3K27me3	<i>H. sapiens</i>	580	(1-250) + GGSG linker + (500-580)	335 (Ballaré)	163-250: PHD2
—	NIPP1	—	H3K27me3	<i>H. sapiens</i>	351	1-329 (Jin)	329	310-329: EED

[0865] In some embodiment, the functional domain may be Histone Acetyltransferase Inhibitor Effector Domain. Preferred examples include SET/TAF-1 β listed in the Table below.

target epigenomic sequences can be used to activate or repress promoters, silencer or enhancers.

[0870] Examples of acetyltransferases are known but may include, in some embodiments, histone acetyltransferases. In

Table of Histone Acetyltransferase Inhibitor Effector Domains

Subtype/ Complex	Name	Substrate (if known)	Modification (if known)	Organism	Full size (aa)	Selected truncation (aa)	Final size (aa)	Catalytic domain
—	SET/TAF-1 β	—	—	<i>M. musculus</i>	289	1-289 (Cervoni)	289	—

[0866] It is also preferred to target endogenous (regulatory) control elements (such as enhancers and silencers) in addition to a promoter or promoter-proximal elements. Thus, the invention can also be used to target endogenous control elements (including enhancers and silencers) in addition to targeting of the promoter. These control elements can be located upstream and downstream of the transcriptional start site (TSS), starting from 200 bp from the TSS to 100 kb away. Targeting of known control elements can be used to activate or repress the gene of interest. In some cases, a single control element can influence the transcription of multiple target genes. Targeting of a single control element could therefore be used to control the transcription of multiple genes simultaneously.

[0867] Targeting of putative control elements on the other hand (e.g. by tiling the region of the putative control element as well as 200 bp up to 100 kb around the element) can be used as a means to verify such elements (by measuring the transcription of the gene of interest) or to detect novel control elements (e.g. by tiling 100 kb upstream and downstream of the TSS of the gene of interest). In addition, targeting of putative control elements can be useful in the context of understanding genetic causes of disease. Many mutations and common SNP variants associated with disease phenotypes are located outside coding regions. Targeting of such regions with either the activation or repression systems described herein can be followed by readout of transcription of either a) a set of putative targets (e.g. a set of genes located in closest proximity to the control element) or b) whole-transcriptome readout by e.g. RNAseq or microarray. This would allow for the identification of likely candidate genes involved in the disease phenotype. Such candidate genes could be useful as novel drug targets.

[0868] Histone acetyltransferase (HAT) inhibitors are mentioned herein. However, an alternative in some embodiments is for the one or more functional domains to comprise an acetyltransferase, preferably a histone acetyltransferase. These are useful in the field of epigenomics, for example in methods of interrogating the epigenome. Methods of interrogating the epigenome may include, for example, targeting epigenomic sequences. Targeting epigenomic sequences may include the guide being directed to an epigenomic target sequence. Epigenomic target sequence may include, in some embodiments, include a promoter, silencer or an enhancer sequence.

[0869] Use of a functional domain linked to a Group 29 or Group 30 effector protein as described herein, preferably a dead-Group 29 or Group 30 effector protein, more preferably a dead-Fn Group 29 or Group 30 effector protein, to

some embodiments, the histone acetyltransferase may comprise the catalytic core of the human acetyltransferase p300 (Gerbasch & Reddy, Nature Biotech 6th April 2015).

[0871] In some preferred embodiments, the functional domain is linked to a dead-Group 29 or Group 30 effector protein to target and activate epigenomic sequences such as promoters or enhancers. One or more guides directed to such promoters or enhancers may also be provided to direct the binding of the CRISPR enzyme to such promoters or enhancers.

[0872] The term “associated with” is used here in relation to the association of the functional domain to the Group 29 or Group 30 effector protein or the adaptor protein. It is used in respect of how one molecule ‘associates’ with respect to another, for example between an adaptor protein and a functional domain, or between the Group 29 or Group 30 effector protein and a functional domain. In the case of such protein-protein interactions, this association may be viewed in terms of recognition in the way an antibody recognizes an epitope. Alternatively, one protein may be associated with another protein via a fusion of the two, for instance one subunit being fused to another subunit. Fusion typically occurs by addition of the amino acid sequence of one to that of the other, for instance via splicing together of the nucleotide sequences that encode each protein or subunit. Alternatively, this may essentially be viewed as binding between two molecules or direct linkage, such as a fusion protein. In any event, the fusion protein may include a linker between the two subunits of interest (i.e. between the enzyme and the functional domain or between the adaptor protein and the functional domain). Thus, in some embodiments, the Group 29 or Group 30 effector protein or adaptor protein is associated with a functional domain by binding thereto. In other embodiments, the Group 29 or Group 30 effector protein or adaptor protein is associated with a functional domain because the two are fused together, optionally via an intermediate linker.

Functional Screen

[0873] In one aspect, the present invention provides for a method of screening for functional elements associated with a change in a phenotype. The library may be introduced into a population of cells that are adapted to contain a Group 29 or Group 30 effector protein. The cells may be sorted into at least two groups based on the phenotype. The phenotype may be expression of a gene, cell growth, or cell viability. The relative representation of the guide RNAs present in each group are determined, whereby genomic sites associated with the change in phenotype are determined by the

representation of guide RNAs present in each group. The change in phenotype may be a change in expression of a gene of interest. The gene of interest may be upregulated, downregulated, or knocked out/down. The cells may be sorted into a high expression group and a low expression group. The population of cells may include a reporter construct that is used to determine the phenotype. The reporter construct may include a detectable marker. Cells may be sorted by use of the detectable marker.

[0874] In another aspect, the present invention provides for a method of screening for sites associated with resistance to a chemical compound. The chemical compound may be a drug or pesticide. The library may be introduced into a population of cells that are adapted to contain a Group 29 or Group 30 effector protein, wherein each cell of the population contains no more than one guide RNA; the population of cells are treated with the chemical compound; and the representation of guide RNAs are determined after treatment with the chemical compound at a later time point as compared to an early time point, whereby sites associated with resistance to the chemical compound are determined by enrichment of guide RNAs. Representation of sgRNAs may be determined by deep sequencing methods.

[0875] Useful in the practice of the instant invention utilizing Group 29 or Group 30 effector protein complexes are methods used in CRISPR-Cas9 systems and reference is made to the article entitled BCL11A enhancer dissection by Cas9-mediated in situ saturating mutagenesis. Canver, M. C., Smith, E. C., Sher, F., Pinello, L., Sanjana, N. E., Shalem, O., Chen, D. D., Schupp, P. G., Vinjamur, D. S., Garcia, S. P., Luc, S., Kurita, R., Nakamura, Y., Fujiwara, Y., Maeda, T., Yuan, G., Zhang, F., Orkin, S. H., & Bauer, D. E. DOI:10.1038/nature15521, published online Sep. 16, 2015, the article is herein incorporated by reference and discussed briefly below:

[0876] Canver et al. involves novel pooled CRISPR-Cas9 guide RNA libraries to perform in situ saturating mutagenesis of the human and mouse BCL11A erythroid enhancers previously identified as an enhancer associated with fetal hemoglobin (HbF) level and whose mouse ortholog is necessary for erythroid BCL11A expression. This approach revealed critical minimal features and discrete vulnerabilities of these enhancers. Through editing of primary human progenitors and mouse transgenesis, the authors validated the BCL11A erythroid enhancer as a target for HbF reinduction. The authors generated a detailed enhancer map that informs therapeutic genome editing.

Method of Using Group 29 or Group 30 Systems to Modify a Cell or Organism

[0877] The invention in some embodiments comprehends a method of modifying a cell or organism. The cell may be a prokaryotic cell or a eukaryotic cell. The cell may be a mammalian cell. The mammalian cell may be a non-human primate, bovine, porcine, rodent or mouse cell. The cell may be a non-mammalian eukaryotic cell such as poultry, fish or shrimp. The cell may also be a plant cell. The plant cell may be of a crop plant such as cassava, corn, sorghum, wheat, or rice. The plant cell may also be of an algae, tree or vegetable. The modification introduced to the cell by the present invention may be such that the cell and progeny of the cell are altered for improved production of biologic products such as an antibody, starch, alcohol or other desired cellular output. The modification introduced to the cell by the

present invention may be such that the cell and progeny of the cell include an alteration that changes the biologic product produced.

[0878] The system may comprise one or more different vectors. In an aspect of the invention, the effector protein is codon optimized for expression of the desired cell type, preferentially a eukaryotic cell, preferably a mammalian cell or a human cell.

[0879] Packaging cells are typically used to form virus particles that are capable of infecting a host cell. Such cells include 293 cells, which package adenovirus, and ψ 2 cells or PA317 cells, which package retrovirus. Viral vectors used in gene therapy are usually generated by producing a cell line that packages a nucleic acid vector into a viral particle. The vectors typically contain the minimal viral sequences required for packaging and subsequent integration into a host, other viral sequences being replaced by an expression cassette for the polynucleotide(s) to be expressed. The missing viral functions are typically supplied in trans by the packaging cell line. For example, AAV vectors used in gene therapy typically only possess ITR sequences from the AAV genome which are required for packaging and integration into the host genome. Viral DNA is packaged in a cell line, which contains a helper plasmid encoding the other AAV genes, namely rep and cap, but lacking ITR sequences. The cell line may also be infected with adenovirus as a helper. The helper virus promotes replication of the AAV vector and expression of AAV genes from the helper plasmid. The helper plasmid is not packaged in significant amounts due to a lack of ITR sequences. Contamination with adenovirus can be reduced by, e.g., heat treatment to which adenovirus is more sensitive than AAV. Additional methods for the delivery of nucleic acids to cells are known to those skilled in the art. See, for example, US20030087817, incorporated herein by reference.

[0880] In some embodiments, a host cell is transiently or non-transiently transfected with one or more vectors described herein. In some embodiments, a cell is transfected as it naturally occurs in a subject. In some embodiments, a cell that is transfected is taken from a subject. In some embodiments, the cell is derived from cells taken from a subject, such as a cell line. A wide variety of cell lines for tissue culture are known in the art. Examples of cell lines include, but are not limited to, C8161, CCRF-CEM, MOLT, mIMCD-3, NHDF, HeLa-S3, Huh1, Huh4, Huh7, HUVEC, HASMC, HEK_n, HEK_a, MiaPaCell, Panc1, PC-3, TF1, CTLL-2, C1R, Rat6, CV1, RPTE, A10, T24, J82, A375, ARH-77, Calu1, SW480, SW620, SKOV3, SK-UT, CaCo2, P388D1, SEM-K2, WEHI-231, HB56, TIB55, Jurkat, J45.01, LRMB, Bcl-1, BC-3, IC21, DLD2, Raw264.7, NRK, NRK-52E, MRC5, MEF, Hep G2, HeLa B, HeLa T4, COS, COS-1, COS-6, COS-M6A, BS-C-1 monkey kidney epithelial, BALB/3T3 mouse embryo fibroblast, 3T3 Swiss, 3T3-L1, 132-d5 human fetal fibroblasts; 10.1 mouse fibroblasts, 293-T, 3T3, 721, 9L, A2780, A2780ADR, A2780cis, A172, A20, A253, A431, A-549, ALC, B16, B35, BCP-1 cells, BEAS-2B, bEnd.3, BHK-21, BR 293, BxPC3, C3H-10T1/2, C6/36, Cal-27, CHO, CHO-7, CHO-IR, CHO-K1, CHO-K2, CHO-T, CHO Dhfr $-/-$, COR-L23, COR-L23/CPR, COR-L23/5010, COR-L23/R23, COS-7, COV-434, CML T1, CMT, CT26, D17, DH82, DU145, DuCaP, EL4, EM2, EM3, EMT6/AR1, EMT6/AR10.0, FM3, H1299, H69, HB54, HB55, HCA2, HEK-293, HeLa, Hepa1c1c7, HL-60, HMEC, HT-29, Jurkat, JY cells, K562 cells, Ku812, KCL22,

KG1, KYO1, LNCap, Ma-Mel 1-48, MC-38, MCF-7, MCF-10A, MDA-MB-231, MDA-MB-468, MDA-MB-435, MDCK II, MDCK II, MOR/0.2R, MONO-MAC 6, MTD-1A, MyEnd, NCI-H69/CPR, NCI-H69/LX10, NCI-H69/LX20, NCI-H69/LX4, NIH-3T3, NALM-1, NW-145, OPCN/OPCT cell lines, Peer, PNT-1A/PNT 2, RenCa, RIN-5F, RMA/RMAS, Saos-2 cells, Sf-9, SkBr3, T2, T-47D, T84, THP1 cell line, U373, U87, U937, VCaP, Vero cells, WM39, WT-49, X63, YAC-1, YAR, and transgenic varieties thereof. Cell lines are available from a variety of sources known to those with skill in the art (see, e.g., the American Type Culture Collection (ATCC) (Manassas, Va.)). In some embodiments, a cell transfected with one or more vectors described herein is used to establish a new cell line comprising one or more vector-derived sequences. In some embodiments, a cell transiently transfected with the components of a nucleic acid-targeting system as described herein (such as by transient transfection of one or more vectors, or transfection with RNA), and modified through the activity of a nucleic acid-targeting complex, is used to establish a new cell line comprising cells containing the modification but lacking any other exogenous sequence. In some embodiments, cells transiently or non-transiently transfected with one or more vectors described herein, or cell lines derived from such cells are used in assessing one or more test compounds.

[0881] In some embodiments, one or more vectors described herein are used to produce a non-human transgenic animal or transgenic plant. In some embodiments, the transgenic animal is a mammal, such as a mouse, rat, or rabbit. In certain embodiments, the organism or subject is a plant. In certain embodiments, the organism or subject or plant is algae. Methods for producing transgenic plants and animals are known in the art, and generally begin with a method of cell transfection, such as described herein.

[0882] In one aspect, the invention provides for methods of modifying a target polynucleotide in a eukaryotic cell. In some embodiments, the method comprises allowing a nucleic acid-targeting complex to bind to the target polynucleotide to effect cleavage of said target polynucleotide thereby modifying the target polynucleotide, wherein the nucleic acid-targeting complex comprises a nucleic acid-targeting effector protein complexed with a guide RNA hybridized to a target sequence within said target polynucleotide.

[0883] In one aspect, the invention provides a method of modifying expression of a polynucleotide in a eukaryotic cell. In some embodiments, the method comprises allowing a nucleic acid-targeting complex to bind to the polynucleotide such that said binding results in increased or decreased expression of said polynucleotide; wherein the nucleic acid-targeting complex comprises a nucleic acid-targeting effector protein complexed with a guide RNA hybridized to a target sequence within said polynucleotide.

Group 29 or Group 30 Effector Protein Complexes can be Used in Plants

[0884] The invention in some embodiments comprehends a method of modifying an cell or organism. The cell may be a prokaryotic cell or a eukaryotic cell. The cell may be a mammalian cell. The mammalian cell may be a non-human primate, bovine, porcine, rodent or mouse cell. The cell may be a non-mammalian eukaryotic cell such as poultry, fish or shrimp. The cell may also be a plant cell. The plant cell may

be of a crop plant such as cassava, corn, sorghum, wheat, or rice. The plant cell may also be of an algae, tree or vegetable. The modification introduced to the cell by the present invention may be such that the cell and progeny of the cell are altered for improved production of biologic products such as an antibody, starch, alcohol or other desired cellular output. The modification introduced to the cell by the present invention may be such that the cell and progeny of the cell include an alteration that changes the biologic product produced.

[0885] The system may comprise one or more different vectors. In an aspect of the invention, the effector protein is codon optimized for expression the desired cell type, preferentially a eukaryotic cell, preferably a mammalian cell or a human cell.

[0886] CRISPR effector protein system(s) (e.g., single or multiplexed) can be used in conjunction with recent advances in crop genomics. Such CRISPR system(s) can be used to perform efficient and cost effective plant gene or genome or transcriptome interrogation or editing or manipulation—for instance, for rapid investigation and/or selection and/or interrogations and/or comparison and/or manipulations and/or transformation of plant genes or genomes; e.g., to create, identify, develop, optimize, or confer trait(s) or characteristic(s) to plant(s) or to transform a plant genome. There can accordingly be improved production of plants, new plants with new combinations of traits or characteristics or new plants with enhanced traits. Such CRISPR system(s) can be used with regard to plants in Site-Directed Integration (SDI) or Gene Editing (GE) or any Near Reverse Breeding (NRB) or Reverse Breeding (RB) techniques. With respect to use of the CRISPR systems in plants, mention is made of the University of Arizona web site “CRISPR-PLANT” (<http://www.genome.arizona.edu/crispr/>) (supported by Penn State and AGI). Embodiments of the invention can be used in genome editing in plants or where RNAi or similar genome editing techniques have been used previously; see, e.g., Nekrasov, “Plant genome editing made easy: targeted mutagenesis in model and crop plants using the CRISPR/Cas system,” *Plant Methods* 2013, 9:39 (doi:10.1186/1746-4811-9-39); Brooks, “Efficient gene editing in tomato in the first generation using the CRISPR/Cas9 system,” *Plant Physiology* September 2014 pp 114.247577; Shan, “Targeted genome modification of crop plants using a CRISPR-Cas system,” *Nature Biotechnology* 31, 686-688 (2013); Feng, “Efficient genome editing in plants using a CRISPR/Cas system,” *Cell Research* (2013) 23:1229-1232. doi:10.1038/cr.2013.114; published online 20 Aug. 2013; Xie, “RNA-guided genome editing in plants using a CRISPR-Cas system,” *Mol Plant*. 2013 November; 6(6):1975-83. doi: 10.1093/mp/sst119. Epub 2013 Aug. 17; Xu, “Gene targeting using the *Agrobacterium tumefaciens*-mediated CRISPR-Cas system in rice,” *Rice* 2014, 7:5 (2014), Zhou et al., “Exploiting SNPs for biallelic CRISPR mutations in the outcrossing woody perennial *Populus* reveals 4-coumarate: CoA ligase specificity and Redundancy,” *New Phytologist* (2015) (Forum) 1-4 (available online only at www.newphytologist.com); Caliando et al, “Targeted DNA degradation using a CRISPR device stably carried in the host genome, *NATURE COMMUNICATIONS* 6:6989, DOI: 10.1038/ncomms7989, www.nature.com/naturecommunications DOI: 10.1038/ncomms7989; U.S. Pat. No. 6,603,061—*Agrobacterium*-Mediated Plant Transformation Method; U.S. Pat. No. 7,868,149—Plant Genome

Sequences and Uses Thereof and US 2009/0100536—Transgenic Plants with Enhanced Agronomic Traits, all the contents and disclosure of each of which are herein incorporated by reference in their entirety. In the practice of the invention, the contents and disclosure of Morrell et al “Crop genomics: advances and applications,” Nat Rev Genet. 2011 Dec. 29; 13(2):85-96; each of which is incorporated by reference herein including as to how herein embodiments may be used as to plants. Accordingly, reference herein to animal cells may also apply, mutatis mutandis, to plant cells unless otherwise apparent; and, the enzymes herein having reduced off-target effects and systems employing such enzymes can be used in plant applications, including those mentioned herein.

[0887] Sugano et al. (Plant Cell Physiol. 2014 March; 55(3):475-81. doi: 10.1093/pcp/pcu014. Epub 2014 Jan. 18) reports the application of CRISPR/Cas9 to targeted mutagenesis in the liverwort *Marchantia polymorpha* L., which has emerged as a model species for studying land plant evolution. The U6 promoter of *M. polymorpha* was identified and cloned to express the gRNA. The target sequence of the gRNA was designed to disrupt the gene encoding auxin response factor 1 (ARF1) in *M. polymorpha*. Using *Agrobacterium*-mediated transformation, Sugano et al. isolated stable mutants in the gametophyte generation of *M. polymorpha*. CRISPR/Cas9-based site-directed mutagenesis in vivo was achieved using either the Cauliflower mosaic virus 35S or *M. polymorpha* EF1 α promoter to express Cas9. Isolated mutant individuals showing an auxin-resistant phenotype were not chimeric. Moreover, stable mutants were produced by asexual reproduction of T1 plants. Multiple arf1 alleles were easily established using CRISPR/Cas9-based targeted mutagenesis. The methods of Sugano et al. may be applied to the CRISPR Cas system of the present invention.

[0888] Kabadi et al. (Nucleic Acids Res. 2014 Oct. 29; 42(19):e147. doi: 10.1093/nar/gku749. Epub 2014 Aug. 13) developed a single lentiviral system to express a Cas9 variant, a reporter gene and up to four sgRNAs from independent RNA polymerase III promoters that are incorporated into the vector by a convenient Golden Gate cloning method. Each sgRNA was efficiently expressed and can mediate multiplex gene editing and sustained transcriptional activation in immortalized and primary human cells. The methods of Kabadi et al. may be applied to the CRISPR Cas system of the present invention.

[0889] Xing et al. (BMC Plant Biology 2014, 14:327) developed a CRISPR/Cas9 binary vector set based on the pGreen or pCambia backbone, as well as a gRNA. This toolkit requires no restriction enzymes besides BsaI to generate final constructs harboring maize-codon optimized Cas9 and one or more gRNAs with high efficiency in as little as one cloning step. The toolkit was validated using maize protoplasts, transgenic maize lines, and transgenic *Arabidopsis* lines and was shown to exhibit high efficiency and specificity. More importantly, using this toolkit, targeted mutations of three *Arabidopsis* genes were detected in transgenic seedlings of the T1 generation. Moreover, the multiple-gene mutations could be inherited by the next generation. (guide RNA) module vector set, as a toolkit for multiplex genome editing in plants. The toolbox of Xing et al. may be applied to the CRISPR Cas system of the present invention.

[0890] Protocols for targeted plant genome editing via CRISPR/Cas9 are also available in volume 1284 of the series Methods in Molecular Biology pp 239-255 10 Feb. 2015. A detailed procedure to design, construct, and evaluate dual gRNAs for plant codon optimized Cas9 (pcoCas9) mediated genome editing using *Arabidopsis thaliana* and *Nicotiana benthamiana* protoplasts s model cellular systems are described. Strategies to apply the CRISPR/Cas9 system to generating targeted genome modifications in whole plants are also discussed. The protocols described in the chapter may be applied to the CRISPR Cas system of the present invention.

[0891] Ma et al. (Mol Plant. 2015 Aug. 3; 8(8):1274-84. doi: 10.1016/j.molp.2015.04.007) reports robust CRISPR/Cas9 vector system, utilizing a plant codon optimized Cas9 gene, for convenient and high-efficiency multiplex genome editing in monocot and dicot plants. Ma et al. designed PCR-based procedures to rapidly generate multiple sgRNA expression cassettes, which can be assembled into the binary CRISPR/Cas9 vectors in one round of cloning by Golden Gate ligation or Gibson Assembly. With this system, Ma et al. edited 46 target sites in rice with an average 85.4% rate of mutation, mostly in biallelic and homozygous status. Ma et al. provide examples of loss-of-function gene mutations in T0 rice and T1 *Arabidopsis* plants by simultaneous targeting of multiple (up to eight) members of a gene family, multiple genes in a biosynthetic pathway, or multiple sites in a single gene. The methods of Ma et al. may be applied to the CRISPR Cas system of the present invention.

[0892] Lowder et al. (Plant Physiol. 2015 Aug. 21. pii: pp. 00636.2015) also developed a CRISPR/Cas9 toolbox enables multiplex genome editing and transcriptional regulation of expressed, silenced or non-coding genes in plants. This toolbox provides researchers with a protocol and reagents to quickly and efficiently assemble functional CRISPR/Cas9 T-DNA constructs for monocots and dicots using Golden Gate and Gateway cloning methods. It comes with a full suite of capabilities, including multiplexed gene editing and transcriptional activation or repression of plant endogenous genes. T-DNA based transformation technology is fundamental to modern plant biotechnology, genetics, molecular biology and physiology. As such, Applicants developed a method for the assembly of Cas9 (WT, nickase or dCas9) and gRNA(s) into a T-DNA destination-vector of interest. The assembly method is based on both Golden Gate assembly and MultiSite Gateway recombination. Three modules are required for assembly. The first module is a Cas9 entry vector, which contains promoterless Cas9 or its derivative genes flanked by attL1 and attR5 sites. The second module is a gRNA entry vector which contains entry gRNA expression cassettes flanked by attL5 and attL2 sites. The third module includes attR1-attR2-containing destination T-DNA vectors that provide promoters of choice for Cas9 expression. The toolbox of Lowder et al. may be applied to the CRISPR Cas system of the present invention.

[0893] Organisms such as yeast and microalgae are widely used for synthetic biology. Stovicek et al. (Metab. Eng. Comm., 2015; 2:13 describes genome editing of industrial yeast, for example, *Saccharomyces cerevisiae*, to efficiently produce robust strains for industrial production. Stovicek used a CRISPR-Cas9 system codon-optimized for yeast to simultaneously disrupt both alleles of an endogenous gene and knock in a heterologous gene. Cas9 and gRNA were expressed from genomic or episomal 2 μ -based vector loca-

tions. The authors also showed that gene disruption efficiency could be improved by optimization of the levels of Cas9 and gRNA expression. Hlavová et al. (Biotechnol. Adv. 2015) discusses development of species or strains of microalgae using techniques such as CRISPR to target nuclear and chloroplast genes for insertional mutagenesis and screening.

[0894] Petersen (“Towards precisely glycol engineered plants,” Plant Biotech Denmark Annual meeting 2015, Copenhagen, Denmark) developed a method of using CRISPR/Cas9 to engineer genome changes in *Arabidopsis*, for example to glyco engineer *Arabidopsis* for production of proteins and products having desired posttranslational modifications. Hebelstrup et al. (Front Plant Sci. 2015 Apr. 23; 6:247) outlines in planta starch bioengineering, providing crops that express starch modifying enzymes and directly produce products that normally are made by industrial chemical and/or physical treatments of starches. The methods of Petersen and Hebelstrup may be applied to the effector protein system of the present invention.

[0895] Kurthe et al., J Virol. 2012 June; 86(11):6002-9. doi: 10.1128/JVI.00436-12. Epub 2012 Mar. 21) developed an RNA virus-based vector for the introduction of desired traits into grapevine without heritable modifications to the genome. The vector provided the ability to regulate expression of endogenous genes by virus-induced gene silencing. The systems and proteins of the instant invention can be used to silence genes and proteins without heritable modification to the genome.

[0896] In an embodiment, the plant may be a legume. The present invention may utilize the herein disclosed CRISPR-Cas system for exploring and modifying, for example, without limitation, soybeans, peas, and peanuts. Curtin et al. provides a toolbox for legume function genomics. (See Curtin et al., “A genome engineering toolbox for legume Functional genomics,” International Plant and Animal Genome Conference XXII 2014). Curtin used the genetic transformation of CRISPR to knock-out/down single copy and duplicated legume genes both in hairy root and whole plant systems. Some of the target genes were chosen in order to explore and optimize the features of knock-out/down systems (e.g., phytoene desaturase), while others were identified by soybean homology to *Arabidopsis* Dicer-like genes or by genome-wide association studies of nodulation in Medicago.

[0897] Peanut allergies and allergies to legumes generally are a real and serious health concern. The present invention can be used to identify and then edit or silence genes encoding allergenic proteins of such legumes. Without limitation as to such genes and proteins, Nicolaou et al. identifies allergenic proteins in peanuts, soybeans, lentils, peas, lupin, green beans, and mung beans. See, Nicolaou et al., Current Opinion in Allergy and Clinical Immunology 2011; 11(3): 222-228).

[0898] In an advantageous embodiment, the plant may be a tree. The present invention may also utilize the herein disclosed CRISPR Cas system for herbaceous systems (see, e.g., Belhaj et al., Plant Methods 9: 39 and Harrison et al., Genes & Development 28: 1859-1872). In a particularly advantageous embodiment, the CRISPR Cas system of the present invention may target single nucleotide polymorphisms (SNPs) in trees (see, e.g., Zhou et al., New Phytologist, Volume 208, Issue 2, pages 298-301, October 2015). In the Zhou et al. study, the authors applied a CRISPR Cas

system in the woody perennial *Populus* using the 4-coumarate:CoA ligase (4CL) gene family as a case study and achieved 100% mutational efficiency for two 4CL genes targeted, with every transformant examined carrying biallelic modifications. In the Zhou et al., study, the CRISPR/Cas9 system was highly sensitive to single nucleotide polymorphisms (SNPs), as cleavage for a third 4CL gene was abolished due to SNPs in the target sequence.

[0899] The methods of Zhou et al. (New Phytologist, Volume 208, Issue 2, pages 298-301, October 2015) may be applied to the present invention as follows. Two 4CL genes, 4CL1 and 4CL2, associated with lignin and flavonoid biosynthesis, respectively are targeted for CRISPR/Cas9 editing. The *Populus tremula* × *alba* clone 717-1B4 routinely used for transformation is divergent from the genome-sequenced *Populus trichocarpa*. Therefore, the 4CL1 and 4CL2 gRNAs designed from the reference genome are interrogated with in-house 717 RNA-Seq data to ensure the absence of SNPs which could limit Cas efficiency. A third gRNA designed for 4CL5, a genome duplicate of 4CL1, is also included. The corresponding 717 sequence harbors one SNP in each allele near/within the PAM, both of which are expected to abolish targeting by the 4CL5-gRNA. All three gRNA target sites are located within the first exon. For 717 transformation, the gRNA is expressed from the Medicago U6.6 promoter, along with a human codon-optimized Cas under control of the CaMV 35S promoter in a binary vector. Transformation with the Cas-only vector can serve as a control. Randomly selected 4CL1 and 4CL2 lines are subjected to amplicon-sequencing. The data is then processed and biallelic mutations are confirmed in all cases. These methods may be applied to the effector protein system of the present invention.

[0900] In plants, pathogens are often host-specific. For example, *Fusarium oxysporum* f. sp. *lycopersici* causes tomato wilt but attacks only tomato, and *F. oxysporum* f. *dianthii Puccinia graminis* f. sp. *tritici* attacks only wheat. Plants have existing and induced defenses to resist most pathogens. Mutations and recombination events across plant generations lead to genetic variability that gives rise to susceptibility, especially as pathogens reproduce with more frequency than plants. In plants there can be non-host resistance, e.g., the host and pathogen are incompatible. There can also be Horizontal Resistance, e.g., partial resistance against all races of a pathogen, typically controlled by many genes and Vertical Resistance, e.g., complete resistance to some races of a pathogen but not to other races, typically controlled by a few genes. In a Gene-for-Gene level, plants and pathogens evolve together, and the genetic changes in one balance changes in other. Accordingly, using Natural Variability, breeders combine most useful genes for Yield, Quality, Uniformity, Hardiness, Resistance. The sources of resistance genes include native or foreign Varieties, Heirloom Varieties, Wild Plant Relatives, and Induced Mutations, e.g., treating plant material with mutagenic agents. Using the present invention, plant breeders are provided with a new tool to induce mutations. Accordingly, one skilled in the art can analyze the genome of sources of resistance genes, and in Varieties having desired characteristics or traits employ the present invention to induce the rise of resistance genes, with more precision than previous mutagenic agents and hence accelerate and improve plant breeding programs.

[0901] Aside from the plants otherwise discussed herein and above, engineered plants modified by the effector protein and suitable guide, and progeny thereof, as provided. These may include disease or drought resistant crops, such as wheat, barley, rice, soybean or corn; plants modified to remove or reduce the ability to self-pollinate (but which can instead, optionally, hybridise instead); and allergenic foods such as peanuts and nuts where the immunogenic proteins have been disabled, destroyed or disrupted by targeting via a effector protein and suitable guide. Any aspect of using classical CRISPR-Cas systems may be adapted to use in CRISPR systems that are Cas protein agnostic, e.g. Group 29 and Group 30 effector protein systems.

Therapeutic Treatment

[0902] The system of the invention can be applied in areas of former RNA cutting technologies, without undue experimentation, from this disclosure, including therapeutic, assay and other applications, because the present application provides the foundation for informed engineering of the system. The present invention provides for therapeutic treatment of

a disease caused by overexpression of RNA, toxic RNA and/or mutated RNA (such as, for example, splicing defects or truncations). Expression of the toxic RNA may be associated with formation of nuclear inclusions and late-onset degenerative changes in brain, heart or skeletal muscle. In the best studied example, myotonic dystrophy, it appears that the main pathogenic effect of the toxic RNA is to sequester binding proteins and compromise the regulation of alternative splicing (Hum. Mol. Genet. (2006) 15 (suppl 2): R162-R169). Myotonic dystrophy [dystrophia myotonica (DM)] is of particular interest to geneticists because it produces an extremely wide range of clinical features. A partial listing would include muscle wasting, cataracts, insulin resistance, testicular atrophy, slowing of cardiac conduction, cutaneous tumors and effects on cognition. The classical form of DM, which is now called DM type 1 (DM1), is caused by an expansion of CTG repeats in the 3'-untranslated region (UTR) of DMPK, a gene encoding a cytosolic protein kinase.

[0903] The below table presents a list of exons shown to have misregulated alternative splicing in DM1 skeletal muscle, heart or brain.

Tissue/gene	Target	Reference
Skeletal muscle		
ALP	ex 5a, 5b	Lin X., et al. Failure of MBNL1-dependent postnatal splicing transitions in myotonic dystrophy. Hum. Mol. Genet 2006; 15: 2087-2097
CAPN3	ex 16	Lin X., et al. Failure of MBNL1-dependent postnatal splicing transitions in myotonic dystrophy. Hum. Mol. Genet 2006; 15: 2087-2097
CLCN1	int 2, ex 7a, 8a	Mankodi A., et al. Expanded CUG repeats trigger aberrant splicing of CC-1 chloride channel pre-mRNA and hyperexcitability of skeletal muscle in myotonic dystrophy. Mol. Cell 2002; 10: 35-44 Charlet-B N., et al. Loss of the muscle-specific chloride channel in type 1 myotonic dystrophy due to misregulated alternative splicing. Mol. Cell 2002; 10: 45-53
FHOS	ex 11a	Lin X., et al. Failure of MBNL1-dependent postnatal splicing transitions in myotonic dystrophy. Hum. Mol. Genet 2006; 15: 2087-2097
GFAT1	ex 10	Lin X., et al. Failure of MBNL1-dependent postnatal splicing transitions in myotonic dystrophy. Hum. Mol. Genet 2006; 15: 2087-2097
IR	ex 11	Savkur R. S., et al. Aberrant regulation of insulin receptor alternative splicing is associated with insulin resistance in myotonic dystrophy. Nat. Genet. 2001; 29: 40-47
MBNL1	ex 7	Lin X., et al. Failure of MBNL1-dependent postnatal splicing transitions in myotonic dystrophy. Hum. Mol. Genet 2006; 15: 2087-2097
MBNL2	ex 7	Lin X., et al. Failure of MBNL1-dependent postnatal splicing transitions in myotonic dystrophy. Hum. Mol. Genet 2006; 15: 2087-2097
MTMR1	ex 2.1, 2.2	Buj-Bello A., et al. Muscle-specific alternative splicing of myotubularin-related 1 gene is impaired in DM1 muscle cells. Hum. Mol. Genet. 2002; 11: 2297-2307
NRAP	ex 12	Lin X., et al. Failure of MBNL1-dependent postnatal splicing transitions in myotonic dystrophy. Hum. Mol. Genet 2006; 15: 2087-2097
RYR1	ex 70	Kimura T., et al. Altered mRNA splicing of the skeletal muscle ryanodine receptor and sarcoplasmic/endoplasmic reticulum Ca ²⁺ -ATPase in myotonic dystrophy type 1. Hum. Mol. Genet. 2005; 14: 2189-2200
SERCA1	ex 22	Kimura T., et al. Altered mRNA splicing of the skeletal muscle ryanodine receptor and sarcoplasmic/endoplasmic reticulum Ca ²⁺ -ATPase in myotonic dystrophy type 1. Hum. Mol. Genet. 2005; 14: 2189-2200 Lin X., et al. Failure of MBNL1-dependent postnatal splicing transitions in myotonic dystrophy. Hum. Mol. Genet 2006; 15: 2087-2097
z-Titin	ex Zr4, Zr5	Lin X., et al. Failure of MBNL1-dependent postnatal splicing transitions in myotonic dystrophy. Hum. Mol. Genet 2006; 15: 2087-2097
m-Titin	M-line ex5	Lin X., et al. Failure of MBNL1-dependent postnatal splicing transitions in myotonic dystrophy. Hum. Mol. Genet 2006; 15: 2087-2097
TNNT3	fetal ex	Kanadia R. N. et al. A muscleblind knockout model for myotonic dystrophy. Science 2003; 302: 1978-1980
ZASP	ex 11	Lin X., et al. Failure of MBNL1-dependent postnatal splicing transitions in myotonic dystrophy. Hum. Mol. Genet 2006; 15: 2087-2097

-continued

Tissue/gene	Target	Reference
Heart		
TNNT2	ex 5	Philips A. V., et al. Disruption of splicing regulated by a CUG-binding protein in myotonic dystrophy. <i>Science</i> 1998; 280: 737-741
ZASP	ex 11	Mankodi A., et al. Nuclear RNA foci in the heart in myotonic dystrophy. <i>Circ. Res.</i> 2005; 97: 1152-1155
m-Titin	M-line ex 5	Mankodi A., et al. Nuclear RNA foci in the heart in myotonic dystrophy. <i>Circ. Res.</i> 2005; 97: 1152-1155
KCNAB1	ex 2	Mankodi A., et al. Nuclear RNA foci in the heart in myotonic dystrophy. <i>Circ. Res.</i> 2005; 97: 1152-1155
ALP	ex 5	(Mankodi A., et al. Nuclear RNA foci in the heart in myotonic dystrophy. <i>Circ. Res.</i> 2005; 97: 1152-1155
Brain		
TAU	ex 2, ex 10	Sergeant N., et al. Dysregulation of human brain microtubule-associated tau mRNA maturation in myotonic dystrophy type 1. <i>Hum. Mol. Genet.</i> 2001; 10: 2143-2155 Jiang H., et al. Myotonic dystrophy type 1 associated with nuclear foci of mutant RNA, sequestration of muscleblind proteins, and deregulated alternative splicing in neurons. <i>Hum. Mol. Genet.</i> 2004; 13: 3079-3088
APP	ex 7	Jiang H., et al. Myotonic dystrophy type 1 associated with nuclear foci of mutant RNA, sequestration of muscleblind proteins, and deregulated alternative splicing in neurons. <i>Hum. Mol. Genet.</i> 2004; 13: 3079-3088
NMDAR1	ex 5	Jiang H., et al. Myotonic dystrophy type 1 associated with nuclear foci of mutant RNA, sequestration of muscleblind proteins, and deregulated alternative splicing in neurons. <i>Hum. Mol. Genet.</i> 2004; 13: 3079-3088

[0904] The enzymes of the present invention may target overexpressed RNA or toxic RNA, such as for example, the DMPK gene or any of the misregulated alternative splicing in DM1 skeletal muscle, heart or brain in, for example, the above table.

[0905] The enzymes of the present invention may also target trans-acting mutations affecting RNA-dependent functions that cause disease (summarized in *Cell*. 2009 Feb. 20; 136(4): 777-793) as indicated in the below table.

DISEASE	GENE/MUTATION	FUNCTION
Prader Willi syndrome	SNORD116	ribosome biogenesis
Spinal muscular atrophy (SMA)	SMN2	splicing
Dyskeratosis congenita (X-linked)	DKC1	telomerase/translation
Dyskeratosis congenita (autosomal dominant)	TERC	telomerase
Dyskeratosis congenita (autosomal dominant)	TERT	telomerase
Diamond-Blackfan anemia	RPS19, RPS24	ribosome biogenesis
Shwachman-Diamond syndrome	SBDS	ribosome biogenesis
Treacher-Collins syndrome	TCOF1	ribosome biogenesis
Prostate cancer	SNHG5	ribosome biogenesis
Myotonic dystrophy, type 1 (DM1)	DMPK (RNA gain-of-function)	protein kinase
Myotonic dystrophy type 2 (DM2)	ZNF9 (RNA gain-of-function)	RNA binding
Spinocerebellar ataxia 8 (SCA8)	ATXN8/ATXN8OS (RNA gain-of-function)	unknown/noncoding RNA
Huntington's disease-like 2 (HDL2)	JPH3 (RNA gain-of-function)	ion channel function
Fragile X-associated tremor ataxia syndrome (FXTAS)	FMR1 (RNA gain-of-function)	translation/mRNA localization
Fragile X syndrome	FMR1	translation/mRNA localization
X-linked mental retardation	UPF3B	translation/nonsense mediated decay
Oculopharyngeal muscular dystrophy (OPMD)	PABPN1	3' end formation
Human pigmentary genodermatosis	DSRAD	editing
Retinitis pigmentosa	PRPF31	splicing
Retinitis pigmentosa	PRPF8	splicing
Retinitis pigmentosa	HPRP3	splicing
Retinitis pigmentosa	PAP1	splicing

-continued

DISEASE	GENE/MUTATION	FUNCTION
Cartilage-hair hypoplasia (recessive)	RMRP	splicing
Autism	7q22-q33 locus breakpoint	noncoding RNA
Beckwith-Wiedemann syndrome (BWS)	H19	noncoding RNA
Charcot-Marie-Tooth (CMT) Disease	GRS	translation
Charcot-Marie-Tooth (CMT) Disease	YRS	translation
Amyotrophic lateral sclerosis (ALS)	TARDBP	splicing, transcription
Leukoencephalopathy with vanishing white matter	EIF2B1	translation
Wolcott-Rallison syndrome	EIF2AK3	translation (protease)
Mitochondrial myopathy and sideroblastic anemia (MLASA)	PUS1	translation
Encephalomyopathy and hypertrophic cardiomyopathy	TSFM	translation (mitochondrial)
Hereditary spastic paraplegia	SPG7	ribosome biogenesis
Leukoencephalopathy	DARS2	translation (mitochondrial)
Susceptibility to diabetes mellitus	LARS2	translation (mitochondrial)
Deafness	MTRNR1	ribosome biogenesis (mitochondrial)
MELAS syndrome, deafness	MTRNR2	ribosome biogenesis (mitochondrial)
Cancer	SFRS1	splicing, translation, export
Cancer	RBM5	splicing
Multiple disorders	mitochondrial tRNA mutations	translation (mitochondrial)
Cancer	miR-17-92 cluster	RNA interference
Cancer	miR-372/miR-373	RNA interference

[0906] The enzyme of the present invention may also be used in the treatment of various tauopathies, including primary and secondary tauopathies, such as primary age-related tauopathy (PART)/Neurofibrillary tangle-predominant senile dementia, with NFTs similar to AD, but without plaques, dementia pugilistica (chronic traumatic encephalopathy), progressive supranuclear palsy, corticobasal degeneration, frontotemporal dementia and parkinsonism linked to chromosome 17, lytico-Bodig disease (Parkinson-dementia complex of Guam), ganglioglioma and gangliocytoma, meningioangiomas, postencephalitic parkinsonism, subacute sclerosing panencephalitis, as well as lead encephalopathy, tuberous sclerosis, Hallervorden-Spatz disease, and lipofuscinosis, alzheimers disease. The enzymes of the present invention may also target mutations disrupting the cis-acting splicing code cause splicing defects and disease (summarized in Cell. 2009 Feb. 20; 136(4): 777-793). The motor neuron degenerative disease SMA results from deletion of the SMN1 gene. The remaining SMN2 gene has a C->T substitution in exon 7 that inactivates an exonic splicing enhancer (ESE), and creates an exonic splicing silencer (ESS), leading to exon 7 skipping and a truncated protein (SMN Δ 7). A T->A substitution in exon 31 of the dystrophin gene simultaneously creates a premature termination codon (STOP) and an ESS, leading to exon 31 skipping. This mutation causes a mild form of DMD because the mRNA lacking exon 31 produces a partially functional protein. Mutations within and downstream of exon 10 of the MAPT gene encoding the tau protein affect splicing regulatory elements and disrupt the normal 1:1 ratio of mRNAs including or excluding exon 10. This results in a perturbed balance between tau proteins containing either four or three microtubule-binding domains (4R-tau and 3R-tau, respectively), causing the neuropathological disorder FTDP-17.

The example shown is the N279K mutation which enhances an ESE function promoting exon 10 inclusion and shifting the balance toward increased 4R-tau. Polymorphic (UG)m (U)n tracts within the 3' splice site of the CFTR gene exon 9 influence the extent of exon 9 inclusion and the level of full-length functional protein, modifying the severity of cystic fibrosis (CF) caused by a mutation elsewhere in the CFTR gene.

[0907] The innate immune system detects viral infection primarily by recognizing viral nucleic acids inside an infected cell, referred to as DNA or RNA sensing. In vitro RNA sensing assays can be used to detect specific RNA substrates. The RNA targeting effector protein can for instance be used for RNA-based sensing in living cells. Examples of applications are diagnostics by sensing of, for examples, disease-specific RNAs.

[0908] The RNA targeting effector protein of the invention can further be used for antiviral activity, in particular against RNA viruses. The effector protein can be targeted to the viral RNA using a suitable guide RNA selective for a selected viral RNA sequence. In particular, the effector protein may be an active nuclease that cleaves RNA, such as single stranded RNA. provided is therefore the use of an RNA targeting effector protein of the invention as an antiviral agent.

[0909] Therapeutic dosages of the enzyme system of the present invention to target RNA the above-referenced RNAs are contemplated to be about 0.1 to about 2 mg/kg the dosages may be administered sequentially with a monitored response, and repeated dosages if necessary, up to about 7 to 10 doses per patient. Advantageously, samples are collected from each patient during the treatment regimen to ascertain the effectiveness of treatment. For example, RNA samples may be isolated and quantified to determine if expression is

reduced or ameliorated. Such a diagnostic is within the purview of one of skill in the art.

Transcript Detection Methods

[0910] The effector proteins and systems of the invention are useful for specific detection of RNAs in a cell or other sample. In the presence of an RNA target of interest, guide-dependent Cas13b nuclease activity may be accompanied by non-specific RNase activity against collateral targets. To take advantage of the RNase activity, all that is needed is a reporter substrate that can be detectably cleaved. For example, a reporter molecule can comprise RNA, tagged with a fluorescent reporter molecule (fluor) on one end and a quencher on the other. In the absence of Cas13b RNase activity, the physical proximity of the quencher dampens fluorescence from the fluor to low levels. When Cas13b target specific cleavage is activated by the presence of an RNA target-of-interest and suitable guide RNA, the RNA-containing reporter molecule is non-specifically cleaved and the fluor and quencher are spatially separated. This causes the fluor to emit a detectable signal when excited by light of the appropriate wavelength.

[0911] In one exemplary assay method, Cas13b effector, target-of-interest-specific guide RNA, and reporter molecule are added to a cellular sample. An increase in fluorescence indicates the presence of the RNA target-of-interest. In another exemplary method, a detection array is provided. Each location of the array is provided with Cas13b effector, reporter molecule, and a target-of-interest-specific guide RNA. Depending on the assay to be performed, the target-of-interest-specific guide RNAs at each location of the array can be the same, different, or a combination thereof. Different target-of-interest-specific guide RNAs might be provided, for example when it is desired to test for one or more targets in a single source sample. The same target-of-interest-specific guide RNA might be provided at each location, for example when it is desired to test multiple samples for the same target.

[0912] With respect to general information on CRISPR-Cas Systems, components thereof, and delivery of such components, including methods, materials, delivery vehicles, vectors, particles, AAV, and making and using thereof, including as to amounts and formulations, all useful in the practice of the instant invention, reference is made to: U.S. Pat. Nos. 8,999,641, 8,993,233, 8,945,839, 8,932,814, 8,906,616, 8,895,308, 8,889,418, 8,889,356, 8,871,445, 8,865,406, 8,795,965, 8,771,945 and 8,697,359; US Patent Publications US 2014-0310830 (U.S. application Ser. No. 14/105,031), US 2014-0287938 A1 (U.S. application Ser. No. 14/213,991), US 2014-0273234 A1 (U.S. application Ser. No. 14/293,674), US2014-0273232 A1 (U.S. application Ser. No. 14/290,575), US 2014-0273231 (U.S. application Ser. No. 14/259,420), US 2014-0256046 A1 (U.S. application Ser. No. 14/226,274), US 2014-0248702 A1 (U.S. application Ser. No. 14/258,458), US 2014-0242700 A1 (U.S. application Ser. No. 14/222,930), US 2014-0242699 A1 (U.S. application Ser. No. 14/183,512), US 2014-0242664 A1 (U.S. application Ser. No. 14/104,990), US 2014-0234972 A1 (U.S. application Ser. No. 14/183,471), US 2014-0227787 A1 (U.S. application Ser. No. 14/256,912), US 2014-0189896 A1 (U.S. application Ser. No. 14/105,035), US 2014-0186958 (U.S. application Ser. No. 14/105,017), US 2014-0186919 A1 (U.S. application Ser. No. 14/104,977), US 2014-0186843 A1 (U.S. applica-

tion Ser. No. 14/104,900), US 2014-0179770 A1 (U.S. application Ser. No. 14/104,837) and US 2014-0179006 A1 (U.S. application Ser. No. 14/183,486), US 2014-0170753 (U.S. application Ser. No. 14/183,429); European Patents EP 2 784 162 B1 and EP 2 771 468 B1; European Patent Applications EP 2 771 468 (EP13818570.7), EP 2 764 103 (EP13824232.6), and EP 2 784 162 (EP14170383.5); and PCT Patent Publications PCT Patent Publications WO 2014/093661 (PCT/US2013/074743), WO 2014/093694 (PCT/US2013/074790), WO 2014/093595 (PCT/US2013/074611), WO 2014/093718 (PCT/US2013/074825), WO 2014/093709 (PCT/US2013/074812), WO 2014/093622 (PCT/US2013/074667), WO 2014/093635 (PCT/US2013/074691), WO 2014/093655 (PCT/US2013/074736), WO 2014/093712 (PCT/US2013/074819), WO 2014/093701 (PCT/US2013/074800), WO 2014/018423 (PCT/US2013/051418), WO 2014/204723 (PCT/US2014/041790), WO 2014/204724 (PCT/US2014/041800), WO 2014/204725 (PCT/US2014/041803), WO 2014/204726 (PCT/US2014/041804), WO 2014/204727 (PCT/US2014/041806), WO 2014/204728 (PCT/US2014/041808), WO 2014/204729 (PCT/US2014/041809). Reference is also made to U.S. provisional patent applications 61/758,468; 61/802,174; 61/806,375; 61/814,263; 61/819,803 and 61/828,130, filed on Jan. 30, 2013; Mar. 15, 2013; Mar. 28, 2013; Apr. 20, 2013; May 6, 2013 and May 28, 2013 respectively. Reference is also made to U.S. provisional patent application 61/836,123, filed on Jun. 17, 2013. Reference is additionally made to U.S. provisional patent applications 61/835,931, 61/835,936, 61/836,127, 61/836,101, 61/836,080 and 61/835,973, each filed Jun. 17, 2013. Further reference is made to U.S. provisional patent applications 61/862,468 and 61/862,355 filed on Aug. 5, 2013; 61/871,301 filed on Aug. 28, 2013; 61/960,777 filed on Sep. 25, 2013 and 61/961,980 filed on Oct. 28, 2013. Reference is yet further made to: PCT Patent applications Nos: PCT/US2014/041803, PCT/US2014/041800, PCT/US2014/041809, PCT/US2014/041804 and PCT/US2014/041806, each filed Jun. 10, 2014; PCT/US2014/041808 filed Jun. 11, 2014; and PCT/US2014/62558 filed Oct. 28, 2014, and U.S. Provisional patent applications Ser. Nos. 61/915,150, 61/915,301, 61/915,267 and 61/915,260, each filed Dec. 12, 2013; 61/757,972 and 61/768,959, filed on Jan. 29, 2013 and Feb. 25, 2013; 61/835,936, 61/836,127, 61/836,101, 61/836,080, 61/835,973, and 61/835,931, filed Jun. 17, 2013; 62/010,888 and 62/010,879, both filed Jun. 11, 2014; 62/010,329 and 62/010,441, each filed Jun. 10, 2014; 61/939,228 and 61/939,242, each filed Feb. 12, 2014; 61/980,012, filed Apr. 15, 2014; 62/038,358, filed Aug. 17, 2014; 62/054,490, 62/055,484, 62/055,460 and 62/055,487, each filed Sep. 25, 2014; and 62/069,243, filed Oct. 27, 2014. Reference is also made to U.S. provisional patent applications Nos. 62/055,484, 62/055,460, and 62/055,487, filed Sep. 25, 2014; U.S. provisional patent application 61/980,012, filed Apr. 15, 2014; and U.S. provisional patent application 61/939,242 filed Feb. 12, 2014. Reference is made to PCT application designating, inter alia, the United States, application No. PCT/US14/41806, filed Jun. 10, 2014. Reference is made to U.S. provisional patent application 61/930,214 filed on Jan. 22, 2014. Reference is made to U.S. provisional patent applications 61/915,251; 61/915,260 and 61/915,267, each filed on Dec. 12, 2013. Reference is made to US provisional patent application U.S. Ser. No. 61/980,012 filed Apr. 15, 2014. Reference is made to PCT application designating,

inter alia, the United States, application No. PCT/US14/41806, filed Jun. 10, 2014. Reference is made to U.S. provisional patent application 61/930,214 filed on Jan. 22, 2014. Reference is made to U.S. provisional patent applications 61/915,251; 61/915,260 and 61/915,267, each filed on Dec. 12, 2013.

[0913] Mention is also made of U.S. application 62/091,455, filed, 12 Dec. 2014, PROTECTED GUIDE RNAs (PGRNAS); U.S. application 62/096,708, 24 Dec. 2014, PROTECTED GUIDE RNAs (PGRNAS); U.S. application 62/091,462, 12 Dec. 2014, DEAD GUIDES FOR CRISPR TRANSCRIPTION FACTORS; U.S. application 62/096,324, 23 Dec. 2014, DEAD GUIDES FOR CRISPR TRANSCRIPTION FACTORS; U.S. application 62/091,456, 12 Dec. 2014, ESCORTED AND FUNCTIONALIZED GUIDES FOR CRISPR-CAS SYSTEMS; U.S. application 62/091,461, 12 Dec. 2014, DELIVERY, USE AND THERAPEUTIC APPLICATIONS OF THE CRISPR-CAS SYSTEMS AND COMPOSITIONS FOR GENOME EDITING AS TO HEMATOPOETIC STEM CELLS (HSCs); U.S. application 62/094,903, 19 Dec. 2014, UNBIASED IDENTIFICATION OF DOUBLE-STRAND BREAKS AND GENOMIC REARRANGEMENT BY GENOME-WISE INSERT CAPTURE SEQUENCING; U.S. application 62/096,761, 24 Dec. 2014, ENGINEERING OF SYSTEMS, METHODS AND OPTIMIZED ENZYME AND GUIDE SCAFFOLDS FOR SEQUENCE MANIPULATION; U.S. application 62/098,059, 30 Dec. 2014, RNA-TARGETING SYSTEM; U.S. application 62/096,656, 24 Dec. 2014, CRISPR HAVING OR ASSOCIATED WITH DESTABILIZATION DOMAINS; U.S. application 62/096,697, 24 Dec. 2014, CRISPR HAVING OR ASSOCIATED WITH AAV; U.S. application 62/098,158, 30 Dec. 2014, ENGINEERED CRISPR COMPLEX INSERTIONAL TARGETING SYSTEMS; U.S. application 62/151,052, 22 Apr. 2015, CELLULAR TARGETING FOR EXTRACELLULAR EXOSOMAL REPORTING; U.S. application 62/054,490, 24 Sep. 2014, DELIVERY, USE AND THERAPEUTIC APPLICATIONS OF THE CRISPR-CAS SYSTEMS AND COMPOSITIONS FOR TARGETING DISORDERS AND DISEASES USING PARTICLE DELIVERY COMPONENTS; U.S. application 62/055,484, 25 Sep. 2014, SYSTEMS, METHODS AND COMPOSITIONS FOR SEQUENCE MANIPULATION WITH OPTIMIZED FUNCTIONAL CRISPR-CAS SYSTEMS; U.S. application 62/087,537, 4 Dec. 2014, SYSTEMS, METHODS AND COMPOSITIONS FOR SEQUENCE MANIPULATION WITH OPTIMIZED FUNCTIONAL CRISPR-CAS SYSTEMS; U.S. application 62/054,651, 24 Sep. 2014, DELIVERY, USE AND THERAPEUTIC APPLICATIONS OF THE CRISPR-CAS SYSTEMS AND COMPOSITIONS FOR MODELING COMPETITION OF MULTIPLE CANCER MUTATIONS IN VIVO; U.S. application 62/067,886, 23 Oct. 2014, DELIVERY, USE AND THERAPEUTIC APPLICATIONS OF THE CRISPR-CAS SYSTEMS AND COMPOSITIONS FOR MODELING COMPETITION OF MULTIPLE CANCER MUTATIONS IN VIVO; U.S. application 62/054,675, 24 Sep. 2014, DELIVERY, USE AND THERAPEUTIC APPLICATIONS OF THE CRISPR-CAS SYSTEMS AND COMPOSITIONS IN NEURONAL CELLS/TISSUES; U.S. application 62/054,528, 24 Sep. 2014, DELIVERY, USE AND THERAPEUTIC APPLICATIONS OF THE CRISPR-CAS SYSTEMS AND COMPOSITIONS IN IMMUNE DISEASES OR DISORDERS; U.S.

application 62/055,454, 25 Sep. 2014, DELIVERY, USE AND THERAPEUTIC APPLICATIONS OF THE CRISPR-CAS SYSTEMS AND COMPOSITIONS FOR TARGETING DISORDERS AND DISEASES USING CELL PENETRATION PEPTIDES (CPP); U.S. application 62/055,460, 25 Sep. 2014, MULTIFUNCTIONAL-CRISPR COMPLEXES AND/OR OPTIMIZED ENZYME LINKED FUNCTIONAL-CRISPR COMPLEXES; U.S. application 62/087,475, 4 Dec. 2014, FUNCTIONAL SCREENING WITH OPTIMIZED FUNCTIONAL CRISPR-CAS SYSTEMS; U.S. application 62/055,487, 25 Sep. 2014, FUNCTIONAL SCREENING WITH OPTIMIZED FUNCTIONAL CRISPR-CAS SYSTEMS; U.S. application 62/087,546, 4 Dec. 2014, MULTIFUNCTIONAL CRISPR COMPLEXES AND/OR OPTIMIZED ENZYME LINKED FUNCTIONAL-CRISPR COMPLEXES; and U.S. application 62/098,285, 30 Dec. 2014, CRISPR MEDIATED IN VIVO MODELING AND GENETIC SCREENING OF TUMOR GROWTH AND METASTASIS.

[0914] Each of these patents, patent publications, and applications, and all documents cited therein or during their prosecution ("appin cited documents") and all documents cited or referenced in the appin cited documents, together with any instructions, descriptions, product specifications, and product sheets for any products mentioned therein or in any document therein and incorporated by reference herein, are hereby incorporated herein by reference, and may be employed in the practice of the invention. All documents (e.g., these patents, patent publications and applications and the appin cited documents) are incorporated herein by reference to the same extent as if each individual document was specifically and individually indicated to be incorporated by reference.

[0915] Also with respect to general information on CRISPR-Cas Systems, mention is made of the following (also hereby incorporated herein by reference):

[0916] Multiplex genome engineering using CRISPR/Cas systems. Cong, L., Ran, F. A., Cox, D., Lin, S., Barretto, R., Habib, N., Hsu, P. D., Wu, X., Jiang, W., Marraffini, L. A., & Zhang, F. *Science* February 15; 339(6121):819-23 (2013);

[0917] RNA-guided editing of bacterial genomes using CRISPR-Cas systems. Jiang W., Bikard D., Cox D., Zhang F, Marraffini L A. *Nat Biotechnol* March; 31(3): 233-9 (2013);

[0918] One-Step Generation of Mice Carrying Mutations in Multiple Genes by CRISPR/Cas-Mediated Genome Engineering. Wang H., Yang H., Shivalila C S., Dawlaty M M., Cheng A W., Zhang F., Jaenisch R. *Cell* May 9; 153(4):910-8 (2013);

[0919] Optical control of mammalian endogenous transcription and epigenetic states. Konermann S, Brigham M D, Trevino A E, Hsu P D, Heidenreich M, Cong L, Platt R J, Scott D A, Church G M, Zhang F. *Nature*. August 22; 500(7463):472-6. doi: 10.1038/Nature12466. Epub 2013 Aug. 23 (2013);

[0920] Double Nicking by RNA-Guided CRISPR Cas9 for Enhanced Genome Editing Specificity. Ran, F A., Hsu, P D., Lin, C Y., Gootenberg, J S., Konermann, S., Trevino, A E., Scott, D A., Inoue, A., Matoba, S., Zhang, Y., & Zhang, F. *Cell* August 28. pii: S0092-8674(13)01015-5 (2013-A);

[0921] DNA targeting specificity of RNA-guided Cas9 nucleases. Hsu, P., Scott, D., Weinstein, J., Ran, F A.,

- Konermann, S., Agarwala, V., Li, Y., Fine, E., Wu, X., Shalem, O., Cradick, T J., Marraffini, L A., Bao, G., & Zhang, F. *Nat Biotechnol* doi:10.1038/nbt.2647 (2013);
- [0922] Genome engineering using the CRISPR-Cas9 system. Ran, F A., Hsu, P D., Wright, J., Agarwala, V., Scott, D A., Zhang, F. *Nature Protocols* November; 8(11):2281-308 (2013-B);
- [0923] Genome-Scale CRISPR-Cas9 Knockout Screening in Human Cells. Shalem, O., Sanjana, N E., Hartenian, E., Shi, X., Scott, D A., Mikkelsen, T., Heckl, D., Ebert, B L., Root, D E., Doench, J G., Zhang, F. *Science* December 12. (2013). [Epub ahead of print];
- [0924] Crystal structure of cas9 in complex with guide RNA and target DNA. Nishimasu, H., Ran, F A., Hsu, P D., Konermann, S., Shehata, S I., Dohmae, N., Ishitani, R., Zhang, F., Nureki, O. *Cell* February 27, 156(5):935-49 (2014);
- [0925] Genome-wide binding of the CRISPR endonuclease Cas9 in mammalian cells. Wu X., Scott D A., Kriz A J., Chiu A C., Hsu P D., Dadon D B., Cheng A W., Trevino A E., Konermann S., Chen S., Jaenisch R., Zhang F., Sharp P A. *Nat Biotechnol.* April 20. doi: 10.1038/nbt.2889 (2014);
- [0926] CRISPR-Cas9 Knockin Mice for Genome Editing and Cancer Modeling. Platt R J, Chen S, Zhou Y, Yim M J, Swiech L, Kempton H R, Dahlman J E, Parnas O, Eisenhaure T M, Jovanovic M, Graham D B, Jhunjhunwala S, Heidenreich M, Xavier R J, Langer R, Anderson D G, Hacohen N, Regev A, Feng G, Sharp P A, Zhang F. *Cell* 159(2): 440-455 DOI: 10.1016/j.cell.2014.09.014 (2014);
- [0927] Development and Applications of CRISPR-Cas9 for Genome Engineering, Hsu P D, Lander E S, Zhang F., *Cell.* June 5; 157(6):1262-78 (2014).
- [0928] Genetic screens in human cells using the CRISPR/Cas9 system, Wang T, Wei J J, Sabatini D M, Lander E S., *Science.* January 3; 343(6166): 80-84. doi:10.1126/science.1246981 (2014);
- [0929] Rational design of highly active sgRNAs for CRISPR-Cas9-mediated gene inactivation, Doench J G, Hartenian E, Graham D B, Tothova Z, Hegde M, Smith I, Sullender M, Ebert B L, Xavier R J, Root D E., (published online 3 Sep. 2014) *Nat Biotechnol.* December; 32(12): 1262-7 (2014);
- [0930] In vivo interrogation of gene function in the mammalian brain using CRISPR-Cas9, Swiech L, Heidenreich M, Banerjee A, Habib N, Li Y, Trombetta J, Sur M, Zhang F., (published online 19 Oct. 2014) *Nat Biotechnol.* January; 33(1):102-6 (2015);
- [0931] Genome-scale transcriptional activation by an engineered CRISPR-Cas9 complex, Konermann S, Brigham M D, Trevino A E, Joung J, Abudayyeh O O, Barcena C, Hsu P D, Habib N, Gootenberg J S, Nishimasu H, Nureki O, Zhang F., *Nature.* January 29; 517(7536): 583-8 (2015).
- [0932] A split-Cas9 architecture for inducible genome editing and transcription modulation, Zetsche B, Volz S E, Zhang F., (published online 2 Feb. 2015) *Nat Biotechnol.* February; 33(2):139-42 (2015);
- [0933] Genome-wide CRISPR Screen in a Mouse Model of Tumor Growth and Metastasis, Chen S, Sanjana N E, Zheng K, Shalem O, Lee K, Shi X, Scott D A, Song J, Pan J Q, Weissleder R, Lee H, Zhang F, Sharp P A. *Cell* 160, 1246-1260, Mar. 12, 2015 (multiplex screen in mouse), and
- [0934] In vivo genome editing using *Staphylococcus aureus* Cas9, Ran F A, Cong L, Yan W X, Scott D A, Gootenberg J S, Kriz A J, Zetsche B, Shalem O, Wu X, Makarova K S, Koonin E V, Sharp P A, Zhang F., (published online 1 Apr. 2015), *Nature.* April 9; 520 (7546):186-91 (2015).
- [0935] Shalem et al., "High-throughput functional genomics using CRISPR-Cas9," *Nature Reviews Genetics* 16, 299-311 (May 2015).
- [0936] Xu et al., "Sequence determinants of improved CRISPR sgRNA design," *Genome Research* 25, 1147-1157 (August 2015).
- [0937] Parnas et al., "A Genome-wide CRISPR Screen in Primary Immune Cells to Dissect Regulatory Networks," *Cell* 162, 675-686 (Jul. 30, 2015).
- [0938] Ramanan et al., CRISPR/Cas9 cleavage of viral DNA efficiently suppresses hepatitis B virus," *Scientific Reports* 5:10833. doi: 10.1038/srep10833 (Jun. 2, 2015)
- [0939] Nishimasu et al., Crystal Structure of *Staphylococcus aureus* Cas9," *Cell* 162, 1113-1126 (Aug. 27, 2015)
- [0940] Zetsche et al., "Cpf1 Is a Single RNA-Guided Endonuclease of a Class 2 CRISPR-Cas System," *Cell* 163, 1-13 (Oct. 22, 2015)
- [0941] Shmakov et al., "Discovery and Functional Characterization of Diverse Class 2 CRISPR-Cas Systems," *Molecular Cell* 60, 1-13 (Available online Oct. 22, 2015)
- each of which is incorporated herein by reference, may be considered in the practice of the instant invention, and discussed briefly below:
- [0942] Cong et al. engineered type II CRISPR-Cas systems for use in eukaryotic cells based on both *Streptococcus thermophilus* Cas9 and also *Streptococcus pyogenes* Cas9 and demonstrated that Cas9 nucleases can be directed by short RNAs to induce precise cleavage of DNA in human and mouse cells. Their study further showed that Cas9 as converted into a nicking enzyme can be used to facilitate homology-directed repair in eukaryotic cells with minimal mutagenic activity. Additionally, their study demonstrated that multiple guide sequences can be encoded into a single CRISPR array to enable simultaneous editing of several at endogenous genomic loci sites within the mammalian genome, demonstrating easy programmability and wide applicability of the RNA-guided nuclease technology. This ability to use RNA to program sequence specific DNA cleavage in cells defined a new class of genome engineering tools. These studies further showed that other CRISPR loci are likely to be transplantable into mammalian cells and can also mediate mammalian genome cleavage. Importantly, it can be envisaged that several aspects of the CRISPR-Cas system can be further improved to increase its efficiency and versatility.
- [0943] Jiang et al. used the clustered, regularly interspaced, short palindromic repeats (CRISPR)—associated Cas9 endonuclease complexed with dual-RNAs to introduce precise mutations in the genomes of *Streptococcus pneumoniae* and *Escherichia coli*. The approach relied on dual-RNA:Cas9-directed cleavage at the targeted genomic site to kill unmutated cells and circumvents the need for selectable markers or counter-selection systems. The study reported reprogramming dual-RNA:Cas9

- specificity by changing the sequence of short CRISPR RNA (crRNA) to make single- and multinucleotide changes carried on editing templates. The study showed that simultaneous use of two crRNAs enabled multiplex mutagenesis. Furthermore, when the approach was used in combination with recombineering, in *S. pneumoniae*, nearly 100% of cells that were recovered using the described approach contained the desired mutation, and in *E. coli*, 65% that were recovered contained the mutation.
- [0944] Wang et al. (2013) used the CRISPR/Cas system for the one-step generation of mice carrying mutations in multiple genes which were traditionally generated in multiple steps by sequential recombination in embryonic stem cells and/or time-consuming intercrossing of mice with a single mutation. The CRISPR/Cas system will greatly accelerate the in vivo study of functionally redundant genes and of epistatic gene interactions.
- [0945] Konermann et al. (2013) addressed the need in the art for versatile and robust technologies that enable optical and chemical modulation of DNA-binding domains based CRISPR Cas9 enzyme and also Transcriptional Activator Like Effectors
- [0946] Ran et al. (2013-A) described an approach that combined a Cas9 nickase mutant with paired guide RNAs to introduce targeted double-strand breaks. This addresses the issue of the Cas9 nuclease from the microbial CRISPR-Cas system being targeted to specific genomic loci by a guide sequence, which can tolerate certain mismatches to the DNA target and thereby promote undesired off-target mutagenesis. Because individual nicks in the genome are repaired with high fidelity, simultaneous nicking via appropriately offset guide RNAs is required for double-stranded breaks and extends the number of specifically recognized bases for target cleavage. The authors demonstrated that using paired nicking can reduce off-target activity by 50- to 1,500-fold in cell lines and to facilitate gene knockout in mouse zygotes without sacrificing on-target cleavage efficiency. This versatile strategy enables a wide variety of genome editing applications that require high specificity.
- [0947] Hsu et al. (2013) characterized SpCas9 targeting specificity in human cells to inform the selection of target sites and avoid off-target effects. The study evaluated >700 guide RNA variants and SpCas9-induced indel mutation levels at >100 predicted genomic off-target loci in 293T and 293FT cells. The authors that SpCas9 tolerates mismatches between guide RNA and target DNA at different positions in a sequence-dependent manner, sensitive to the number, position and distribution of mismatches. The authors further showed that SpCas9-mediated cleavage is unaffected by DNA methylation and that the dosage of SpCas9 and sgRNA can be titrated to minimize off-target modification. Additionally, to facilitate mammalian genome engineering applications, the authors reported providing a web-based software tool to guide the selection and validation of target sequences as well as off-target analyses.
- [0948] Ran et al. (2013-B) described a set of tools for Cas9-mediated genome editing via non-homologous end joining (NHEJ) or homology-directed repair (HDR) in mammalian cells, as well as generation of modified cell lines for downstream functional studies. To minimize off-target cleavage, the authors further described a double-nicking strategy using the Cas9 nickase mutant with paired guide RNAs. The protocol provided by the authors experimentally derived guidelines for the selection of target sites, evaluation of cleavage efficiency and analysis of off-target activity. The studies showed that beginning with target design, gene modifications can be achieved within as little as 1-2 weeks, and modified clonal cell lines can be derived within 2-3 weeks.
- [0949] Shalem et al. described a new way to interrogate gene function on a genome-wide scale. Their studies showed that delivery of a genome-scale CRISPR-Cas9 knockout (GeCKO) library targeted 18,080 genes with 64,751 unique guide sequences enabled both negative and positive selection screening in human cells. First, the authors showed use of the GeCKO library to identify genes essential for cell viability in cancer and pluripotent stem cells. Next, in a melanoma model, the authors screened for genes whose loss is involved in resistance to vemurafenib, a therapeutic that inhibits mutant protein kinase BRAF. Their studies showed that the highest-ranking candidates included previously validated genes NF1 and MED12 as well as novel hits NF2, CUL3, TADA2B, and TADA1. The authors observed a high level of consistency between independent guide RNAs targeting the same gene and a high rate of hit confirmation, and thus demonstrated the promise of genome-scale screening with Cas9.
- [0950] Nishimasu et al. reported the crystal structure of *Streptococcus pyogenes* Cas9 in complex with sgRNA and its target DNA at 2.5 Å resolution. The structure revealed a bilobed architecture composed of target recognition and nuclease lobes, accommodating the sgRNA: DNA heteroduplex in a positively charged groove at their interface. Whereas the recognition lobe is essential for binding sgRNA and DNA, the nuclease lobe contains the HNH and RuvC nuclease domains, which are properly positioned for cleavage of the complementary and non-complementary strands of the target DNA, respectively. The nuclease lobe also contains a carboxyl-terminal domain responsible for the interaction with the protospacer adjacent motif (PAM). This high-resolution structure and accompanying functional analyses have revealed the molecular mechanism of RNA-guided DNA targeting by Cas9, thus paving the way for the rational design of new, versatile genome-editing technologies.
- [0951] Wu et al. mapped genome-wide binding sites of a catalytically inactive Cas9 (dCas9) from *Streptococcus pyogenes* loaded with single guide RNAs (sgRNAs) in mouse embryonic stem cells (mESCs). The authors showed that each of the four sgRNAs tested targets dCas9 to between tens and thousands of genomic sites, frequently characterized by a 5-nucleotide seed region in the sgRNA and an NGG protospacer adjacent motif (PAM). Chromatin inaccessibility decreases dCas9 binding to other sites with matching seed sequences; thus 70% of off-target sites are associated with genes. The authors showed that targeted sequencing of 295 dCas9 binding sites in mESCs transfected with catalytically active Cas9 identified only one site mutated above background levels. The authors proposed a two-state model for Cas9 binding and cleavage, in which a seed match triggers binding but extensive pairing with target DNA is required for cleavage.
- [0952] Platt et al. established a Cre-dependent Cas9 knockin mouse. The authors demonstrated in vivo as well

- as ex vivo genome editing using adeno-associated virus (AAV)-, lentivirus-, or particle-mediated delivery of guide RNA in neurons, immune cells, and endothelial cells.
- [0953] Hsu et al. (2014) is a review article that discusses generally CRISPR-Cas9 history from yogurt to genome editing, including genetic screening of cells.
- [0954] Wang et al. (2014) relates to a pooled, loss-of-function genetic screening approach suitable for both positive and negative selection that uses a genome-scale lentiviral single guide RNA (sgRNA) library.
- [0955] Doench et al. created a pool of sgRNAs, tiling across all possible target sites of a panel of six endogenous mouse and three endogenous human genes and quantitatively assessed their ability to produce null alleles of their target gene by antibody staining and flow cytometry. The authors showed that optimization of the PAM improved activity and also provided an on-line tool for designing sgRNAs.
- [0956] Swiech et al. demonstrate that AAV-mediated SpCas9 genome editing can enable reverse genetic studies of gene function in the brain.
- [0957] Konermann et al. (2015) discusses the ability to attach multiple effector domains, e.g., transcriptional activator, functional and epigenomic regulators at appropriate positions on the guide such as stem or tetraloop with and without linkers.
- [0958] Zetsche et al. demonstrates that the Cas9 enzyme can be split into two and hence the assembly of Cas9 for activation can be controlled.
- [0959] Chen et al. relates to multiplex screening by demonstrating that a genome-wide in vivo CRISPR-Cas9 screen in mice reveals genes regulating lung metastasis.
- [0960] Ran et al. (2015) relates to SaCas9 and its ability to edit genomes and demonstrates that one cannot extrapolate from biochemical assays. Shalem et al. (2015) described ways in which catalytically inactive Cas9 (dCas9) fusions are used to synthetically repress (CRISPRi) or activate (CRISPRa) expression, showing advances using Cas9 for genome-scale screens, including arrayed and pooled screens, knockout approaches that inactivate genomic loci and strategies that modulate transcriptional activity. End Edits
- [0961] Shalem et al. (2015) described ways in which catalytically inactive Cas9 (dCas9) fusions are used to synthetically repress (CRISPRi) or activate (CRISPRa) expression, showing advances using Cas9 for genome-scale screens, including arrayed and pooled screens, knockout approaches that inactivate genomic loci and strategies that modulate transcriptional activity.
- [0962] Xu et al. (2015) assessed the DNA sequence features that contribute to single guide RNA (sgRNA) efficiency in CRISPR-based screens. The authors explored efficiency of CRISPR/Cas9 knockout and nucleotide preference at the cleavage site. The authors also found that the sequence preference for CRISPRi/a is substantially different from that for CRISPR/Cas9 knockout.
- [0963] Parnas et al. (2015) introduced genome-wide pooled CRISPR-Cas9 libraries into dendritic cells (DCs) to identify genes that control the induction of tumor necrosis factor (Tnf) by bacterial lipopolysaccharide (LPS). Known regulators of Tlr4 signaling and previously unknown candidates were identified and classified into three functional modules with distinct effects on the canonical responses to LPS.
- [0964] Ramanan et al (2015) demonstrated cleavage of viral episomal DNA (cccDNA) in infected cells. The HBV genome exists in the nuclei of infected hepatocytes as a 3.2 kb double-stranded episomal DNA species called covalently closed circular DNA (cccDNA), which is a key component in the HBV life cycle whose replication is not inhibited by current therapies. The authors showed that sgRNAs specifically targeting highly conserved regions of HBV robustly suppresses viral replication and depleted cccDNA.
- [0965] Nishimasu et al. (2015) reported the crystal structures of SaCas9 in complex with a single guide RNA (sgRNA) and its double-stranded DNA targets, containing the 5'-TTGAAT-3' PAM and the 5'-TTGGGT-3' PAM. A structural comparison of SaCas9 with SpCas9 highlighted both structural conservation and divergence, explaining their distinct PAM specificities and orthologous sgRNA recognition.
- [0966] Zetsche et al. (2015) reported the characterization of Cpf1, a putative class 2 CRISPR effector. It was demonstrated that Cpf1 mediates robust DNA interference with features distinct from Cas9. Identifying this mechanism of interference broadens our understanding of CRISPR-Cas systems and advances their genome editing applications.
- [0967] Shmakov et al. (2015) reported the characterization of three distinct Class 2 CRISPR-Cas systems. The effectors of two of the identified systems, C2c1 and C2c3, contain RuvC like endonuclease domains distantly related to Cpf1. The third system, C2c2, contains an effector with two predicted HEPN RNase domains.
- [0968] Also, "Dimeric CRISPR RNA-guided FokI nucleases for highly specific genome editing", Shengdar Q. Tsai, Nicolas Wyvekens, Cyd Khayter, Jennifer A. Foden, Vishal Thapar, Deepak Reyon, Mathew J. Goodwin, Martin J. Aryee, J. Keith Joung Nature Biotechnology 32(6): 569-77 (2014), relates to dimeric RNA-guided FokI Nucleases that recognize extended sequences and can edit endogenous genes with high efficiencies in human cells.
- [0969] In addition, mention is made of PCT application PCT/US14/70057, Attorney Reference 47627.99.2060 and BI-2013/107 entitled "DELIVERY, USE AND THERAPEUTIC APPLICATIONS OF THE CRISPR-CAS SYSTEMS AND COMPOSITIONS FOR TARGETING DISORDERS AND DISEASES USING PARTICLE DELIVERY COMPONENTS (claiming priority from one or more or all of US provisional patent applications: 62/054,490, filed Sep. 24, 2014; 62/010,441, filed Jun. 10, 2014; and 61/915,118, 61/915,215 and 61/915,148, each filed on Dec. 12, 2013) ("the Particle Delivery PCT"), incorporated herein by reference, with respect to a method of preparing an sgRNA-and-Cas9 protein containing particle comprising admixing a mixture comprising an sgRNA and Cas9 protein (and optionally HDR template) with a mixture comprising or consisting essentially of or consisting of surfactant, phospholipid, biodegradable polymer, lipoprotein and alcohol; and particles from such a process. For example, wherein Cas9 protein and sgRNA were mixed together at a suitable, e.g., 3:1 to 1:3 or 2:1 to 1:2 or 1:1 molar ratio, at a suitable temperature, e.g., 15-30 C, e.g., 20-25 C, e.g., room temperature, for a suitable time, e.g., 15-45, such as 30 minutes, advantageously in sterile, nuclease free buffer, e.g., 1xPBS. Separately, particle components such as or comprising: a surfactant, e.g., cationic lipid, e.g., 1,2-dioleoyl-3-trimeth-

ylammonium-propane (DOTAP); phospholipid, e.g., dimyristoylphosphatidylcholine (DMPC); biodegradable polymer, such as an ethylene-glycol polymer or PEG, and a lipoprotein, such as a low-density lipoprotein, e.g., cholesterol were dissolved in an alcohol, advantageously a C₁₋₆ alkyl alcohol, such as methanol, ethanol, isopropanol, e.g., 100% ethanol. The two solutions were mixed together to form particles containing the Cas9-sgRNA complexes. Accordingly, sgRNA may be pre-complexed with the Cas9 protein, before formulating the entire complex in a particle. Formulations may be made with a different molar ratio of different components known to promote delivery of nucleic acids into cells (e.g. 1,2-dioleoyl-3-trimethylammonium-propane (DOTAP), 1,2-ditetradecanoyl-sn-glycero-3-phosphocholine (DMPC), polyethylene glycol (PEG), and cholesterol) For example DOTAP:DMPC:PEG:Cholesterol Molar Ratios may be DOTAP 100, DMPC 0, PEG 0, Cholesterol 0; or DOTAP 90, DMPC 0, PEG 10, Cholesterol 0; or DOTAP 90, DMPC 0, PEG 5, Cholesterol 5. DOTAP 100, DMPC 0, PEG 0, Cholesterol 0. That application accordingly comprehends admixing sgRNA, Cas9 protein and components that form a particle; as well as particles from such admixing. Aspects of the instant invention can involve particles; for example, particles using a process analogous to that of the Particle Delivery PCT, e.g., by admixing a mixture comprising sgRNA and/or Cas9 as in the instant invention and components that form a particle, e.g., as in the Particle Delivery PCT, to form a particle and particles from such admixing (or, of course, other particles involving sgRNA and/or Cas9 as in the instant invention).

[0970] The present invention will be further illustrated in the following Examples which are given for illustration purposes only and are not intended to limit the invention in any way.

Examples

Example 1: An Expanded Search for RNA Programmable Genomic Engineering Effectors

[0971] Cas9 is an RNA-programmable endonuclease associated with CRISPRs found in prokaryotes that has transformed biology. Recently additional unique candidates for RNA-programmable single effectors have been discovered, and in one case validated (e.g. Cpf1), based on proximity to CRISPR-associated Cas1 and Cas2 proteins. The work described in the present application expands this search with the aim of exhausting the possibility space for RNA-programmable CRISPR single effectors in a full and up-to-date genomic dataset.

[0972] Background: The discovery of the CRISPR/Cas9 endonuclease system found near CRISPRs (clustered regularly interspaced short palindromic repeats) of prokaryotes has radically transformed biology. CRISPR systems composed of CRISPR arrays and CRISPR-associated proteins confer adaptive immunity to bacteria and archaea against invading bacteriophages (Makarova et al., 2011; 2015). By reengineering CRISPR systems, in the span of a few years it has become possible to modify precisely the genome of an organism on a system level (Ran et al., 2015).

[0973] CRISPR arrays contain spacer sequences acquired from bacteriophage DNA separated by repeat sequences for array processing. The CRISPR-associated Cas9 protein operates by recruiting a crRNA from the CRISPR array and a tracrRNA nearby the array to cleave DNA and induce a

double strand break at a specific programmable sequence, encoded by an N-base pair target sequence and short adjacent M-base pair protospacer adjacent motif (PAM), with N, M, and the sequence of the PAM varying among species (Chylinski et al., 2013; Fonfara et al., 2014).

[0974] Unlike multi-effector CRISPR systems, Type II CRISPR systems rely on the single effector Cas9 (Chylinski et al., 2014). Thus, precise genome engineering can be achieved with a minimal two-component system: the Cas9 protein and a modified sgRNA (unimolecular combination of tracrRNA and crRNA).

[0975] Several new single-effector CRISPR systems with effectors distinct from Cas9 have recently been discovered through a bioinformatic search of single large effectors proximate to Cas1 and Cas2, two small proteins implicated in CRISPR spacer acquisition (Shmakov et al., submitted 2015). In addition, Cpf1, a protein previously identified but whose function had not yet been elucidated, has been found to operate analogously to Cas9, but without the use of a tracrRNA (Zetsche et al., submitted 2015). It is hypothesized that there may be yet more single-effector systems to discover, with DNA-targeting, RNA-targeting, or perhaps even DNA integrating or transposing capabilities inherited from a common evolutionary origin (Koonin & Krupovic, 2015).

[0976] Prior to the present invention, a CRISPR-associated protein-agnostic approach to identify novel single-effector systems was not available. As described in detail below, the present invention provides a bioinformatic pipeline for discovery of novel Class 2 CRISPR systems excluding Cas1 and Cas2. Beyond uncovering new biology, the invention provides means to identify an RNA-guided single effector for translation into a useful genomic engineering tool through the following steps:

[0977] 1) computationally classifying new single-effector systems; and

[0978] 2) experimentally following up to expand genome engineering capabilities.

[0979] Computational Procedure: A biocomputational pipeline was designed and implemented to mine through metagenomics datasets for novel Class 2 CRISPR single effectors, akin to the revolutionary genome-engineering tools Cas9 and Cpf1. Whereas previous search strategies relied on protein proximity to CRISPR-associated spacer acquisition proteins Cas1 and Cas2, this approach was seeded on CRISPR arrays alone. What resulted was the discovery of a potential new Class 2 CRISPR system, with two subtypes as characterized by distinct putative accessory proteins.

[0980] The implementation of the computational procedure has already yielded a few novel potential single-effector systems. FIG. 1 depicts the basic workflow. First, all prokaryotic genomes are downloaded (the European Ensembl and American NCBI databases contain the largest publicly accessible repositories). Second, the PILER-CR CRISPR array discovery algorithm is run on all genomes with default settings to store all possible CRISPR arrays in data (Edgar, 2007). Third, nearby proteins within 10 kilobase pairs of discovered arrays are stored in data. Fourth, in a single-effector filter, all systems with no nearby protein greater than 700 amino acids or more than one nearby protein greater than 700 amino acids are filtered out. Additionally, only systems with a large single effector between 900 and 1800 amino acids—informed by functionality and compatibility considerations—are chosen. Fifth, these pro-

teins are subjected to an NCBI BLAST homology search against the entire NCBI protein database to identify homologous proteins, and HHpred protein domain homology searches are conducted on all homologous proteins found (Soding et al., 2005). From all data and the literature, genomic loci are visualized. The candidate list was limited to families in which the putative effector was near a CRISPR array in at least 50% of the members to focus on potentially functional CRISPR loci.

[0981] FIG. 2 shows overall preliminary data from an input of Ensembl Genomes Release 27, with the number of arrays and large proteins discovered graphed according to category (all arrays vs. single large effector arrays; arrays with varying combinations of Cas1 and Cas2 proximate). The pertinent value for the proposed research is the number of large proteins in single large effector arrays with neither Cas1 nor Cas 2 proximate, namely 1,531.

[0982] These large single effectors in systems lacking Cas1 and Cas2 were subjected to an NCBI BLAST homology search with an E-value cutoff of $1e-7$ against the entire NCBI protein database to identify homologous proteins (Camacho et al., 2009). HHpred protein domain homology searches were conducted on all proteins found in this manner (Soding et al., 2005).

[0983] From the data obtained and annotated, genomic loci are visualized, as in FIG. 3 for the case of Cas9 of the pathogen *Finnegoldia magna*. In order to select experimental candidates, it is relevant to consider proteins with low HHpred homology matches, minimal existing CRISPR classifications, short distance of protein from CRISPR array in genomic locus, identical orientation of protein with respect to adjacent proteins, consistent nature of CRISPR arrays in the genome, and presence of few annotated CRISPR proteins nearby. Based on the number of protein candidates within a family determined by NCBI BLAST homology, one can assess whether said family is evolutionarily conserved, and thus whether it is likely to be biologically functional. Applicants looked for conserved domains, even if previously unannotated, among proteins within a family.

Experimental Procedure: Discovering New Class 2 loci—Group 29

[0984] When subjected to NCBI BLAST homology searches, the 1,531 candidates grew to roughly 7,200 proteins. These proteins and their respective loci were grouped into 266 folders according to a nearest-neighbor E-value cutoff of $1e-7$. Many of these groups were false positives, for example transposons and methylases, the majority in the group of which was not correlated with any CRISPR arrays. Several of the groups were in fact Class 1 CRISPR loci that had passed the large single effector-filter. Quite a few groups were bona fide Class 2 CRISPR loci—including Cas9, Cpf1, and C2c2—some in the group of which had no Cas1 or Cas2 proximate. One group in particular, sorted into folder 29, possessed a diverse set of novel proteins all lacking Cas1 or Cas2 proximate. Referred to in this specification as either “Group 29” or “Grp29”, this system embodying two loci will now be described in bio computational detail.

[0985] The Group 29 system consists of a large single effector (~1100 amino acids in length) and one or none-of two small putative accessory proteins (~200 amino acids in length) nearby a CRISPR array. Based on the nearby small protein, the system is bifurcated into two Loci A and B, depicted in FIGS. 78 and 79. No additional proteins out to

25 kilobase pairs upstream or downstream from the array are conserved across species with each locus. With minor exceptions, the CRISPR array consists of direct repeat sequences 36 nucleotides in length and spacer sequences 30 nucleotides in length. The direct repeat is generally well conserved, especially at the ends, with a GTTG/GUUG at the 5' end reverse complementary to a CAAC at the 3' end. This conservation suggests strong base pairing for an RNA loop structure that potentially interacts with the protein(s) in the locus. A motif search complementary to the direct repeats revealed no candidate tracrRNAs nearby the arrays, possibly indicative of a single crRNA like that found in the Cpf1 locus.

[0986] A phage genome search of Grp29 spacer sequences using CRISPR Target revealed no sequence matches above 90% (Biswas et al., 2013). The result is unsurprising, given how few phage genomes have been sequenced, in addition to the relatively high mutation rates phages endure to evade bacterial adaptive immunity. From the unreliability of matches, a target-flanking PAM sequence could not be predicted.

[0987] Each of the proteins possesses interesting structure. The large protein contains many conserved residues forming putative domains interspersed within a mixed alpha helix/beta sheet structure. Two HEPN domains (each of the amino acid motif RXXXXH) are located near the N and C termini of the protein, respectively. Given the predicted RNase activity of HEPN domains, two of which are also present in the interior of the C2c2 protein, it is predicted that Grp29 may be an RNA-targeting effector, though rigorous experimental follow-up outside the scope of this thesis is required to validate this assumption (Anantharaman et al., 2013; Shmakov et al., 2015).

[0988] The small proteins fall into two distinct categories. In Locus A the small protein is predicted to contain four putative transmembrane domains throughout, whereas in Locus B it is predicted to contain a single putative transmembrane domain near the N-terminus and a HEPN domain (amino acid motif RXXXXH) near the C-terminus. Each of the transmembrane domains was recognized with high probability by TMHMM, a membrane protein topology prediction model implementing a hidden Markov model (Krogh et al., 2001). Transmembrane domains typically contain hydrophobic residues, and are therefore likely to localize to the cellular membrane. (It is worth noting that often signaling peptides at the N-terminus of a protein may be mistaken for a transmembrane domain, as may be the case with the small protein of Locus B.). Characterizing Group 29 loci

[0989] A total of 111 large proteins from Group 29 loci were found through a series of exhaustive NCBI BLAST homology searches. Multiple sequence alignment of the 111 revealed that only 94 were complete annotated proteins, the remainder typically falling at the end of a sequencing fragment. The size distribution of the 94 is shown in FIG. 80. Most large proteins fall tightly in the size range of 1100-1150 amino acids, suggestive of low divergence either due to relatively short genetic time span or strong selective pressure to maintain similar form.

[0990] Of these 94 large proteins, only 77 were found to be non-redundant protein entries on NCBI. The BLOSUM62 algorithm was used to align large proteins and generate a phylogenetic tree, depicted along with species and locus information in FIG. 81.

[0991] The species represented are found in Bacterioidetes, a broad phylum of Gram-negative, rod-shaped bacteria. Notably, the gut microbes *Porphyromonas gingivalis* and *Porphyromonas gulae* possess both loci. There are a few exceptions to the rules outlined in the previous section: i) Locus A has quite a few more species without a small protein than does Locus B, ii) CRISPR arrays as recognized by PILERCR are not present in all loci, and iii) the recently discovered chemoheterotrophic *Phaeodactylibacter xiame-nensis* has both an outlier locus orientation and direct repeat sequence:

5' - GTCCTAGTCACTCCTCAAACCGGGAGCCTATCGGAC - 3' .

[0992] Most bacterial strains do not have fully annotated genomes, and so it cannot be concluded whether there exist other CRISPR systems in their genomes. (This is indicated by a question mark in the last three columns of the FIG. 81.) A few strains of species with the Group 29 locus, *Prevotella intermedia* and *Porphyromonas gingivalis*, possess no Cas1 or Cas2 in their entire genomes. Given that CRISPR systems are suspected to propagate through horizontal gene transfer and not evolution, and that a CRISPR system does not explicitly require adaptive Cas1/Cas2 nearby (it can recruit them in trans, possibly the case for all Group 29 loci) or even in the entire genome (CRISPR systems can continue to function without adaptation), the results do not argue strongly in favor of self-adaptation in the Group 29 system. Translating Computational Work into Experimental Workflow to Yield Novel Effectors

[0993] FIG. 82 depicts a minimalist experimental workflow from Group 29 locus discovery to potential genomic engineering applications. The workflow begins with locus characterization. From biocomputational work, the bacterial strain containing the locus of interest is known. One can therefore order the strain and its genomic DNA, and perform RNA sequencing, either on the native strain or on the locus heterologously transplanted into a well-established model organism such as *Escherichia coli*. Heterologous expression entails the integration of the locus (protein and CRISPR array) into a plasmid transformed under drug selection into the model organism. In the case where native regulatory elements of the locus are insufficient to express the protein and/or CRISPR array, synthetic gene promoter(s) must be inserted to drive expression. From the RNA sequencing data, one can gather whether the protein sufficiently processes the CRISPR array into crRNA for RNA-programmable nucleic acid targeting. RNA sequencing can be performed in a number of ways, all of which convert bacterially harvested RNA into DNA libraries to be sequenced on Next Generation Sequencing machines (Adiconis et al., 2013; Mardis, 2013).

[0994] If the protein does process its CRISPR array, RNA sequencing will indicate the nature of the mature crRNA, both in length and orientation. This information can be used to design a two-component system consisting of the protein and mature crRNA to perform in vitro interference experiments. In these experiments, protein purified from bacteria and crRNA synthesized by in vitro transcription are incubated with a target nucleic acid (DNA or RNA). The product of the incubation is run on an agarose gel, and derivation from the control condition (typically without protein and/or crRNA) indicates cleavage of the target.

[0995] It is expected that the HEPN domains that Grp29 may interfere with RNA, therefore both DNA and RNA targets must be tested. If either shows interference, the next stage is to determine the targeting rules in vivo. The PAM screen is designed to reveal any cleavage-mediating sequence motifs flanking the target sequence reverse complementary to the spacer sequence, and may be adapted to discover RNA targeting motifs (Zetsche et al., 2015). Efficiency screens may also be performed, especially if the targeting rules are more complex than a PAM (Doench et al., 2014). For RNA targeting in particular, RNA phage interference experiments ought to be conducted to show that RNA interference alone, and not transcriptionally dependent RNA interference, is responsible (Tamulaitis et al., 2014).

[0996] For potential endonucleases, Applicants ordered the bacterial species with the proteins of interest and cloned their genomic loci into plasmids for transformation into competent *E. coli*. Applicants then performed RNA sequencing on the bacteria and the transformed *E. coli* to determine transcribed RNAs in the vicinity of the CRISPR array. Following this initial characterization, Applicants screened for PAMs by cloning a library of randomly synthesized nucleotides flanking spacer sequences and co-transforming with the plasmids containing genomic loci. Through Next Generation Sequencing and observing depletion of certain randomized nucleotides, Applicants predicted PAM sequences. After PAM validation in bacterial colonies, Applicants synthesized predicted crRNAs (and tracrRNAs and sgRNAs if possible), purified protein of interest from cell lysate, and performed in vitro cleavage assays, observing cutting through the Surveyor method. Applicants followed up with experiments in mammalian cells to assess for eukaryotic genomic engineering potential.

[0997] Applicants identified Group 29 proteins to be a set of RNA-programmable, RNA-interfering nucleases possibly involved in bacterial adaptive immunity against RNA phages. This is based on the predicted presence of known RNA-catalytic HEPN domains near their N and C termini, and RNA-Seq data indicating processing of an adjacent CRISPR array into putative crRNAs. The Group 30 proteins identified are a couple proteins with little local alignment to Group 29 but which congregate with one of the accessory small proteins from Group 29.

[0998] From Group 29 and 30 proteins, Applicants identified two novel, genetically diverse putative Class 2 CRISPR-Cas systems (105 total genomic loci; 81 loci with a non-redundant single effector) represented in Gram-positive bacteria. All these loci encode a large (~1100 aa) candidate effector protein and, in about 80% of the loci, an additional small (~200 aa) protein (FIG. 46b, FIG. 68). The putative effector proteins contain two predicted HEPN domains at their N- and C-termini (Anantharaman et al., 2013), similar to the domain architecture of the large effector in the Class 2 type VI-A system (Cas13a). Beyond the catalytic signatures of the HEPN domains, however, there was no significant sequence similarity between the new predicted effector and Cas13a. Accordingly, we classified these systems as subtype VI-B and denoted its putative effector protein Cas13b.

[0999] The presence of one or the other putative accessory protein correlates with the two distinct branches in the phylogenetic tree of Cas13b, indicative of the existence of two variant systems, which we denoted VI-B1 (accessory protein Csx27) and VI-B2 (accessory protein Csx28). The

type VI-B2 systems we identified almost invariably contain *csx28*, whereas *csx27* is less consistently present in VI-B1 loci. The *Csx27* and *Csx28* protein sequences show no significant similarity to any previously identified CRISPR-associated proteins. Both *Csx27* and *Csx28* were predicted to contain one or more transmembrane segments (FIG. 71A). However, RFP-tagged *Csx27* of *Bergeyella zoohelcum* and *Csx28* *Prevotella buccae* at either the N- or C-terminus did not noticeably localize to the membrane (FIG. 71B), suggesting that both are soluble proteins despite the presence of hydrophobic α -helicase domains. In addition, analysis of the multiple alignment of *Csx28* proteins indicated the presence of a highly diverged HEPN domain (FIG. 46C).

[1000] Unlike their divergent putative accessory proteins, both type VI-B systems show distinct, conserved features in the CRISPR arrays. The direct repeats in the CRISPR arrays are conserved in size, sequence, and structure (FIGS. 69 & 70A), with a length of 36 nt, a poly-U stretch in the open loop region, and complementary sequences 5'-GUUG and CAAC-3' at the ends of the repeat predicted to yield a defined secondary structure mediated by intramolecular base-pairing. Although a CRISPRTarget probe of protospacers revealed no consensus flanking region sequence (FIG. 70B) (Biswas et al., 2013), a well-conserved protospacer length of 30 nt combined with the well-conserved direct repeat sequence and length suggests the possibility of similar nucleic acid targeting rules among different VI-B loci.

[1001] RNA sequencing of the CRISPR array in native *Bergeyella zoohelcum* (type VI-B1) showed processing of the pre-crRNA into a 66 nt mature crRNA, with the full 30 nt 5' spacer followed by the full 36 nt 3' direct repeat (FIG. 47A). A longer 118 nt crRNA, distal to the nominal crRNAs in the CRISPR array and consisting of 5' and 3' truncated direct repeat sequences with an intervening repeat sequence, was also processed, a phenomenon predicted to occur in other VI-B loci, suggesting it has some functional role.

[1002] To validate the interference activity of the VI-B system and to determine the targeting rules for the VI-B1 locus from *B. zoohelcum*, Applicants developed an *E. coli* essential gene screen (FIG. 64A). For this negative selection screen, Applicants generated a library of 54,600 unique spacers tiled with single-nucleotide resolution over the coding region of 45 operon-exclusive essential genes (Gerdes et al., 2003; Baba et al., 2006), plus 60 nt into the 5' and 3' UTRs, plus 1100 randomly generated non-targeting spacers. Applicants then transformed this library with plasmids carrying *bzcas13b* (*cas13b* gene from *B. zoohelcum*) and *bzcsx27*, just *bzcas13b*, or a control empty vector. After quality-control-filtering of all screen spacers, Applicants found statistically significant depletion of targeting spacers over non-targeting spacers, indicating that *Cas13b* alone or with *Csx27* can achieve nucleic acid interference (FIG. 64B).

[1003] To assess *Cas13b* targeting rules, Applicants established two depletion levels: strongly depleted (top 1% of depleted spacers) and safely depleted (spacers depleted 5 σ above the mean depletion of the filtered non-targeting spacers). From spacers passing the strongly depleted cutoff, Applicants generated weblogs qualitatively identifying a double-sided protospacer flanking sequence (PFS) (FIG. 64C) (Crooks et al., 2004). Given that each base in a sequence weblog is independent, Applicants developed a more quantitative, base-dependent PFS score defined as the

ratio of the number of safely depleted spacers to the number of all spacers, normalized across all PFS scores (FIG. 64D).

[1004] The normalized PFS scores revealed a 5' PFS of D (A, U, or G) and 3' PFS of N followed by AN or NA, consistent for both *Cas13b* and *Csx27*, as well as for *Cas13b* alone. To validate this sequence-targeting rule, Applicants performed an orthogonal depletion screen with *Cas13b* alone, targeting the Kanamycin resistance gene (FIG. 64E). Four classes of spacers were created: non-targeting, targeting with both 5' and 3' PFS rules, targeting with just the 5' or 3' PFS rule, and targeting with neither rule. Consistent with findings from the *E. coli* essential gene screen, the 5' and 3' PFS spacers resulted in the highest Kanamycin sensitivity (FIG. 64G, FIG. 72).

[1005] Based on the computationally predicted HEPN domains, which function as RNases in other CRISPR-Cas systems such as VI-A, Applicants predicted that *Cas13b* interferes with RNA. To test this prediction, Applicants assayed interference against the lytic, single-stranded RNA bacteriophage MS2, whose life cycle contains no DNA intermediates. Applicants performed an MS2 drop plaque assay at serial dilutions of phage for both *bzcas13b* with *bzcsx27* and *bzcas13b* alone with three spacers targeting the MS2 genome, two at the *lys*-rep interface and one in rep, and one non-targeting spacer (FIG. 65A). Upon completion of the assay, substantial reductions in plaque formation for all targeting spacers compared to the non-targeting spacer were observed, consistent with sequence-specific RNA targeting by VI-B1 systems. (FIG. 65B, FIG. 73). Notably, the presence of *bzcsx27* appeared to weaken RNA interference by *bzcas13b* for all three spacers. Applicants confirmed the importance of the predicted enzymatic domains for interference by mutating the conserved, catalytic arginines and histidines in the two HEPN domains to alanine (R116A/H121A and R1177A/H1182A), leading to ~4.5 orders of magnitude decrease in resistance to MS2 phage (FIG. 65D, FIG. 74).

[1006] To rule out DNA interference as an additional mode of action, Applicants modified an existing plasmid interference assay with a protospacer placed either in-frame at the 5' end of the *bla* ampicillin-resistance gene transcript (transcribed target) or upstream of the *bla* gene promoter (non-transcribed target). Bacteria co-transformed with *bzcas13b* and spacer, as well as the non-transcribed target plasmid, survived at a comparable rate to co-transformation of the same target with the empty vector on dual antibiotic selection. For bacteria co-transformed with the transcribed target, the colony forming unit rate under dual antibiotic selection was reduced by ~2 orders of magnitude in the presence of *bzcas13b*, suggesting that *Cas13b* targets RNA exclusively (FIG. 65C).

[1007] Given that RNA targets likely contain a variety of secondary structures and the results of the *E. coli* essential gene screen suggested that additional targeting rules beyond the PFS exist (only ~18% of spacers were safely depleted for *bzcas13b*; with sufficient gene essentiality the expected value would be ~33% from the PFS rule alone), Applicants sought to determine how RNA accessibility impacts targeting. Using the Vienna RNAplfold model (Bernhart et al., 2006), which has been successfully used to predict RNAi efficiency (Tafer et al., 2008), Applicants trained and tested the RNA accessibility model for spacer efficiency on the screen data, and found that RNA accessibility matters most in the protospacer region most distal to the direct repeat of

the crRNA (FIG. 66B, 66C). Finally, Applicants calculated the empirical cumulative distribution functions of safely depleted spacers from 5' UTR into gene from 3' UTR into gene (FIG. 66D), finding marginal underrepresentation in spacers in the 3' UTR versus a uniform distribution, but otherwise no significant spatial dependence in spacer efficiency.

[1008] To determine if these RNA targeting rules generalize across type VI-B loci, Applicants next characterized the type VI-B2 locus from *P. buccae*. RNA sequencing of the CRISPR array revealed identical processing to that of *B. zoohelcum*, excluding the long crRNA (FIG. 47B). Performing the *E. coli* essential gene screen with pbcas13b and pbcx28 or just pbcas13 revealed a similar PFS matrix to that of *B. zoohelcum*, with certain PFS's disfavored (FIG. 67A, FIG. 75). As with bzcxs27, the presence of pbcx28 did not appreciably alter the PFS. Applicants also repeated the secondary structure analysis with pbcas13b, and a comparable RNAPfold model applied (FIGS. 76A and B), but strikingly the safely depleted spacers for pbcas13b without pbcx28 were highly biased to the 5' UTR of genes, suggestive of inhibited or more spatially localized RNase activity in the absence of pbcx28 (FIG. 76C). Applicants confirmed the presumed reduced activity of pbcas13b relative to its full locus using the MS2 phage plaque drop assay, which showed that pbcx28 enhances MS2 phage interference by up to four orders of magnitude (FIG. 67B, FIG. 73). The divergent ability of csx27 to repress and csx28 to enhance cas13b activity generalizes across thousands of spacers in

the *E. coli* essential gene screen (FIG. 67C), highlighting a unique, bimodal property of type VI-B within Class 2 CRISPR-Cas systems.

[1009] Previous studies have shown that Cas13a exhibits non-specific RNase activity after target cleavage. Cas13b may possess similar collateral activity, in which case repression of Cas13b by Csx27 in VI-B1 systems could be an important regulatory mechanism (Makarova et al., 2009; Hayes et al., 2011). In the case of type VI-B2 systems, Csx28 might enhance Cas13b activity to overcome the numerous transcripts of invading bacteriophages or even to promote programmed cell death (Makarova et al., 2012), but this remains to be determined (FIG. 67E).

[1010] The novel, highly efficient, RNA-targeting Class 2 type VI-B CRISPR system resembles the previously described subtype VI-A in general functional scheme and effector protein architecture, but differs substantially in the specifics, in particular, by the apparent absence of local adaptation machinery and the presence of accessory proteins that modulate interference. The novel *E. coli* essential gene screen designed by Applicants provides a fast, unbiased, and reproducible method for elucidating the complex nucleic acid targeting rules of CRISPR-Cas systems and enabled the development of a multifaceted model for Cas13b RNA targeting dependent on sequence, structure, and space (FIG. 67D). An advanced understanding of the targeting rules of Cas13b and other RNA-targeting CRISPR systems raises the prospect of a new suite of precise and robust in vivo RNA manipulation tools for studying a wide range of biological processes.

Amino acid sequences of all Group 29 proteins:

```
>gi|34397202|gb|AAQ66265.1| hypothetical protein PG_1164
[Porphyromonas gingivalis W83]
MNIVPASENKQSRIVEDDPQYFGLYLNLARENLIEVESHVRIKFGKKKLN EESLKQSLLCDHLLSVDRW
TKVYGHARRYL PFLHYFDPDSQIEKDHDSTKGVDPDSAQRLLIRELYSLDDFLRND FSHNRLDGTFEHL
VSPDISSFITGTYSLACGRAQSRFADFFKPD DFLAKNRKEQLI SVADGKCELI VSGLAFFICLFLDREQ
ASGMLSRIRGFKRIDENWARAVHETFCDLCTRHPDRLESSNIKEALLDMLNELNRCPRILYDMLPEEE
RAQFLPALDENSMNLSSENSLN EESRLLWDGSSDWAEALTKRIRHQDRFPYLM LRFIEEMDLLKGRFRV
DLGEIELDSYSKVGVRNGEYDRTITDHALAFGKLSDFQNEEEVSRMISGEASYPVRFSLFAPRYAIDYDK
IGYCHTSDPVYPKSKTGEKRALSNPQSMGFI SVHNLKLLMLLCEGSGSRMQSDFLRKANRILDETAE
GKLOFSALFP EPMHRFIPQP NPKSKDRREKAETLEKYKQEI KGRKDKLNSQLLSAFDMNQRLPSRLLD
EWMNIRPASHSVKLRITYVKQLNEDCRLRLRKRFRKDGDKARAIPLVGE MATFLSQDI VRMIISEETKCLI
TSAYYNEQMSLAQYAGEENRRQFRAI VAEHLHLLDPSGHPFLSATMETAHRYTEDFYKCYLEKKEWLA
KTFYRPEQDENTKRRISVFFVPDGEARKLLPLLIRRRMKEQNDLQDWIRNKQAHPIDLPSHLFDSKIMEL
LKVKDGGKKNWNAFQDWSTKYPDGMQPFYGLRRELNIHGKSVSYI PSDGKFFADCYTHLMEKTVQDKKR
ELRTAGKVPVPPDLAADIKRSFHRAVNEREFMLRLVQEDDRMLMLAINKMMTDREEDI LPLGKNIDSILDE
ENQFSLAVHAKVLEKEGEGGDNLSLVPATIEIKSKRKDW SKYIRYRYDRRVPGLMSHPHKAALDEVK
ILLGEYDRCRIKIFDWAFALEGAIMSRRDLKPYLHES SRREGKSGEHS TLVKMLVEKKGCLTPDESQYLI
LIRNKAAHNQFPCAAEMPLIYRDVSAKVGSI EGSSAKDLPEGSSLVDSLWKKYEMIRKILPLILD PENRF
FGKLLNMSQP INDL
```

```
>gi|188594807|dbj|BAG33782.1| conserved hypothetical protein
[Porphyromonas gingivalis ATCC 33277]
MNIVPASENKQSRIVEDDPQYFGLYLNLARENLIEVESHVRIKFGKKKLN EESLKQSLLCDHLLSVDRW
TKVYGHARRYL PFLHYFDPDSQIEKDHDSTKGVDPDSAQRLLIRELYSLDDFLRND FSHNRLDGTFEHL
VSPDISSFITGTYSLACGRAQSRFAVFFKPD DFLAKNRKEQLI SVADGKCELI VSGFAFFICLFLDREQ
ASGMLSRIRGFKRIDENWARAVHETFCDLCTRHPDRLESSNIKEALLDMLNELNRCPRILYDMLPEEE
RAQFLPALDENSMNLSSENSLN EESRLLWDGSSDWAEALTKRIRHQDRFPYLM LRFIEEMDLLKGRFRV
DLGEIELDSYSKVGVRNGEYDRTITDHALAFGKLSDFQNEEEVSRMISGEASYPVRFSLFAPRYAIDYDK
IGYCHTSDPVYPKSKTGEKRALSNPQSMGFI SVHDLRKLMLLCEGSGSRMQSDFLRKANRILDETAE
GKLOFSALFP EPMHRFIPQP NPKSKDRREKAETLEKYKQEI KGRKDKLNSQLLSAFDMNQRLPSRLLD
EWMNIRPASHSVKLRITYVKQLNEDCRLRLRKRFRKDGDKARAIPLVGE MATFLSQDI VRMIISEETKCLI
TSAYYNEQMSLAQYAGEENRRQFRAI VAEHLHLLDPSGHPFLSATMETAHRYTEGPFYKCYLEKKEWLA
KIFYRPEQDENTKRRISVFFVPDGEARKLLPLLIRRRMKEQNDLQDWIRNKQAHPIDLPSHLFDSKIMEL
LKVKDGGKKNWNAFQDWSTKYPDGMQPFYGLRRELNIHGKSVSYI PSDGKFFADCYTHLMEKTVQDKKR
ELRTAGKVPVPPDLAADIKRSFHRAVNEREFMLRLVQEDDRMLMLAINKMMTDREEDI LPLGKNIDSILDE
```

-continued

ENQFSLAVHAKVLEKEGEGGDNLSLVPATIEIKSKRKDWSKYIRYRYDRRVPGLMSHFPEHKAILDEVK
 ILLGEYDRCRIKIPDWAFALEGAIMSDRDLKPYLHESSSREGKSGEHSTLVKMLVEKKGCLTPDESQYLI
 LIRNKAHNQFPCAAEMPIYRDVSAKVGSI EGGSAKDLPEGSSLVDSLWKYEMIIRKILPILDPENRF
 FGKLLNNMSQPINDL

>gi|188595167|dbj|BAG34142.1| conserved hypothetical protein

[*Porphyromonas gingivalis* ATCC 33277]

MTEQNERPYNGTYTLEDKHFWAAPFNLARHNAYITLAHIDRQLAYSKADITNDEDILFFPKGQWKNLNDND
 LERKARLSLILKHFSLFLEGAAYGKKLPESSQSGNKSSKKKELIKKEKEELQANALSLDNLKSLIFDFLQ
 KLKDFRNYSSHYPHPESELPDFDGNMLQRLYNVFDVSVQRVKRDHEHNDKVDPHRHFNHLVRKGGKDRY
 GNNDNPPFKHHPVDREKVTAGLLFFVSLPLEKRDALWQKIRGPKGGTETYQOMTNEVFCR.SRISLP
 KLKLESLRIDDWMLDMLNELVRCPKSLYDRLEEDRARFRVPVILSDEDDIDGTEEDPFKNI LVRHQD
 RPPYFALRYFDLKKVFTSLRPHIDLGTYHFAIYKKNIGEOPEDRHLTRNLVYGFGRIQDFAEEHRPEEWKR
 LVRDLDFPETGDKPYITQTPHYHI EKGI GLRFVPEGQHLWSPPEVGA TRTGRSKYAQDKRLTAEAFSL
 VHELPMMFYFLLREKYSDASAEVQGRIRKRVIEDVYAVYDAFARGEINTRDELDACLADKGI RRGHL
 PRQMGIGLSQEHKDMEEKVRKQLQEMI VTDHRLDMLDRQTDRKIRIGRKNAGLPKSGVIADWLVRDMMR
 PQPVAKDTSGKPLNNSKANSTERYMLQRALALFGGEKERLIPYFRQMNLI GGNPHPLHETRWESHINI
 LSFYRSYLKARKAFIQS IGRSDRVENHRFL LKPEKPTDRQTLVAGWKGEFHLPRGI FTEAVRDCL IEMGL
 DEVGSKYEVGFMAKAVPLYFERACKDRVQPFYDYPFVNGNSLKPCKGRFLSKERAEWESGKERFRLAK
 LKKEILEAKHFPYLDKFSWQKFERELRLVKNQDIITWMI CRDLMEENKVEGLDIGTLYLKD IRTDVQEQG
 NLNVLNRYKPMRLPVVVYRADSRGHVHKEQAPLATVYIEERDTKLLKQGNFKSFKVDRRLNGLFSFVDTG
 ALAMEQYPI SKLRVEYELAKYQARVCAFEQTLLEESLLTRYPHLPDKNFRKMLESWSDP LDKWPD LH
 GNVRLLI AVRNAPFSHNQYPMYDEAVFSSIRKYDPSDDAI EERMGLNIAHRLSEEVKQAKEMAERI IQA

>gi|288336166|gb|EFC74556.1| putative phage head-tail adaptor

[*Prevotella buccae* D17]

MQKQDKLFPVDRKKNIAIFAPKYITIMENQEKPEPIYELTDKHFWAAPFNLARHNVTYITINHINRRELEIA
 ELKDDGYMMDIKGSWNEQAKKLDKVKRLRDLIMKHFPFLEAAAYEITNSKSPNNKEQREKEQSEALSNN
 LKNVLFIFLEKLVLRNYSYKYSSESPKPIFETSLKNMYKVF DANVRLVKRDYMHENIDMQRDPFH
 LNRKKQVGRINKNIDSPNFHYHFADEKGNMTIAGLLFPVSLFLDKKDAIWMQKLLKGFKDGRLNREQMIN
 EVFCRSRISLPKLEENVTQKDWMLDMLNELVRCPKSLYERLREKDRSFVFPDIFSDDYDAEEEPFK
 NTLVRHQDRFPYFVLRYPDLNEIPEQLRFQIDLGTYPHSIYNKRIGDEEVRHLLTHLYGFARIQDFAAQ
 NQPEVWRKLVKLDLDFEASQEPYIPKTAAPHYHLENEKIGIKFCSIHNNLFP SLKTEKTCNGRSKFNLGIQ
 FTAEAFSLVHHELLPMFYLLIKDYSRKESADKVEGIRKETSNIYAIYDAFANGEINSIADLCRLQK
 TNILQGHLPKQMSILEGRQDMKEEAERKIGEMIDDIQRRLDLLCKQINQKIRIGKRNAGLLKSGKIDAT
 WLVNDMMRFQPVQKQNNIPINNSKANSTERYMLQRALALFGSENFRKAYFNQMNLVGNDNPHPLAET
 QWEHQTNISFYRNYLEARKKYLKGLKPNWKQYQHFLILKVKQTNRNTLV TGWKNFNLPRGIPTQPIR
 EWFEKHNNSKRIDYQILSDFDRVFAKAIPLYFAEEYKDNVQPFYDYPFNI GNKLPKQKQGLDKERVE
 LWQKNKELFKNYPSKKKIDLAYLDFLSWKKFERELRLIKNQD IVTWLWMLFELFNMATVEGLKIGIHLR
 DIDINTANESNNILNRIMPMLVKIYETDNKGNILKERPLATFYI BETETKVLKQGNFKVLAKDRRLN
 GLLSFAETTDIDLEKPNITKLSVDHELKIKYQTRISIFEMTLGLEKLIKINKYPTLPTDSFRNMLERWLQC
 KANRPELKNYVNSLI AVRNAPFSHNQYPMYDATALFAEVKFTLFPSSVDTKKI ELNIA PQLLLEIVGKAIKEI
 EKSENKN

>gi|312446353|gb|ADQ82708.1| hypothetical protein Riean_1551

[*Riemerella anatipestifer* ATCC 11845 = DSM 15868]

MEKPLLPNVYTLKHKFFWGAFLNIARHNAFITICHINEQLGLKTPSNDDKIVDVCETWNNILNNDHDL
 KKSQTLTELILKHFPLTAMCYHPPKKEGKKGHQKEQQKEKESAEQSAEALNPSKLEALEILVNQLHS
 LRNYSHYKHKKPAEAKDIFKHLKAFDASLRMVKEDYKAHFI VNLIRDFAHNLNRKGNKQDNDFNRYR
 FEKDGFFTESGLLFFINLFLDKRDAYWMLKVKVSGFKASHKQREKMIT EVFCRSRILLPKLRLESRYDHNQ
 MLLDMLSELSRCPKLLYKLS EENKHFQVEADGFLDEIEEENQPFKDTLIRHQDRFPYFALRYLDLNE
 FKSIRFQVDLGTYYHYCIYDKKIGDEQEKRLHTRILLSFGRLQDFTIENRPOEWKALTKDLDYKETSNOFF
 ISKTI PHYHITDNKIGFRLGTSKELYPSEIKDGANRI AKYPYNSGFVAHAFI SVHELLPLMFYQHLIGK
 SEDLLKETVRIHQRIYKDFEERINTIEDLEKANQGRPLGAFPKQMLGGLQNKQPDLSKAKIKIEKLI
 AETKLLSHRLNLIKSSPKLGRREKLIKTVLADWLVKDFMRFPQVAYDAQNP I KSSKANSTEFWFIR
 RALALYGGKRNLEGYFKINLIGNINPHPLNKFNWKACRNLVDFYQQYLEQREKFL EAIKNQWPWEYQ
 YCLLLKIPKENRKNLVKGEWQGGISLPRGLFTEAIRETLESDLMMLSKPIRKEIKKHGRVGFISRAITLYF
 KEKYQDKHQSFYNLSYKLEAKAPLLKREEHYEWQNKQPSPTESQRLLELHSDRWKDYLLYKRWQHLEK
 KLRLYRNQDVMWMLTLELTKNHPKELNLYHQLKLENLAVNVQEADAKLNLNQLTLPMLVLPVKVYPATA
 FGEVQYHKTPIRTVYIREEHTKALKMGNFKALVKDRRLNGLFSFIKEENDTQKHPISQLRLRRELEIYQS
 LRVDFAKETLSLEKLLNKHISLSSLENEFRALLEEWKKEYAASSMIDEHIAFIASVRNAPFCHNQYFPY
 KEALHAPIPLFTVAQPI TEEDKGLGIAEALLKVLREYCEIVKSQI

>gi|315022569|gb|EFT35596.1| hypothetical protein RAYM_05191

[*Riemerella anatipestifer* RA-YM]

MEKPLLPNVYTLKHKFFWGAFLNIARHNAFITICHINEQLGLKTPSNDDKIVDVCETWNNILNNDHDL
 KKSQTLTELILKHFPLTAMCYHPPKKEGKKGHQKEQQKEKESAEQSAEALNPSKLEALEILVNQLHS
 LRNYSHYKHKKPAEAKDIFKHLKAFDASLRMVKEDYKAHFI VNLIRDFAHNLNRKGNKQDNDFNRYR
 FEKDGFFTESGLLFFINLFLDKRDAYWMLKVKVSGFKASHKQREKMIT EVFCRSRILLPKLRLESRYDHNQ
 MLLDMLSELSRCPKLLYKLS EENKHFQVEADGFLDEIEEENQPFKDTLIRHQDRFPYFALRYLDLNE
 FKSIRFQVDLGTYYHYCIYDKKIGDEQEKRLHTRILLSFGRLQDFTIENRPOEWKALTKDLDYKETSNOFF
 ISKTI PHYHITDNKIGFRLGTSKELYPSEIKDGANRI AKYPYNSGFVAHAFI SVHELLPLMFYQHLIGK
 SEDLLKETVRIHQRIYKDFEERINTIEDLEKANQGRPLGAFPKQMLGGLQNKQPDLSKAKIKIEKLI
 AETKLLSHRLNLIKSSPKLGRREKLIKTVLADWLVKDFMRFPQVAYDAQNP I KSSKANSTEFWFIR
 RALALYGGKRNLEGYFKINLIGNINPHPLNKFNWKACRNLVDFYQQYLEQREKFL EAIKNQWPWEYQ
 YCLLLKIPKENRKNLVKGEWQGGISLPRGLFTEAIRETLESDLMMLSKPIRKEIKKHGRVGFISRAITLYF
 KEKYQDKHQSFYNLSYKLEAKAPLLKREEHYEWQNKQPSPTESQRLLELHSDRWKDYLLYKRWQHLEK
 KLRLYRNQDVMWMLTLELTKNHPKELNLYHQLKLENLAVNVQEADAKLNLNQLTLPMLVLPVKVYPATA

- continued

FGEVQYHKTPIRTVYIREEHTKALKMGNFKALVKDRRLNGLFSPIKEENDTQKHPISQLRRLRRELEIYQS
LRVDAFKETLSLEEKLLNKHISLSSLENEFRALLEEWKKEYAASSMVIDEHIAFIASVRNAPFCHNQYPPY
KEALHAPIPLFTVAQPIITEBKDGLGIAEALLKVLREYCEIVKSQI

>gi|315252008|gb|EFU31981.1| hypothetical protein HMPREF6485_0083
[Prevotella buccae ATCC 33574]

MQKQDKLFDVRKKNIAIFAPPKYITIMENKEKPEPIYYELTDKHFWAFLNLRHNHYTTINHNRRLEIA
ELKDDGYMMGIKGSWNEQAKKLDKKVLRDLIMKHFPFLAAYEMINSKSPNNKEQREKEQSEALSNN
LKNVLFIFLEKLVLRNYSYHYKSEESPKPIFETSLKNMYKVFANDVRLVKRDMYMHENIDMQDRPTH
LNRKKQVGRINKNIDSPNFHYHFADKEGNMTIAGLLFFVSLFLDKKDAIWMQKCLKGPKDGRNLRQMIN
EVFCRSRISLPLKLENVTKDWMQDMLNELVRCPSLYERLREKDRSEFKVPPDIFSDDYNAEEEPFK
NTLVRHQDRFPYFVLRYPDLNEIFEQLRFDIDLGTYPHSIYNKRI GDEDEVRHLTHLYGFARIQDFAPQ
NQPEWRKLVKDLDFHETSQEPYISKTAPHYHLENEKIGIKFCSAHNNLFPSLQTDKICKGRSKFNLGIQ
PTAEAFSLVHELLPMMFYLLLIKDYSRKESADKVEGIRKEISNIYAIYDAFANNEINSIADLTRLQN
TNILQGHLPKQMISLKGQRQDMGKEAERKIGEMDDIQRRLDLCKQINQKIRIGKRNAGLLKSGKIAD
WLVNDMMRFQPVQKQDNNIPINNSKANSTEYRMLQALALFGSENFRLKAYFNQMNLVGNDNPHFLAET
QWEHQTNILSFYRNYLEARKKYLKGLKPKQNWQYQHFLILKVKQKTRNRNLVTVGWKNSFNLPKRGIFTQPIR
EWFKEHNNISKRIYDQILSFDRVGFVAKAIPLYFAEEYKDNVQPFYDYPFNI GNRLLPKKRFPLDKKERVE
LWQKNKELFKNYPSEKKKIDLAYLDFLWKKFERELRLIKNQDIVTWLMPKELFNMATVEGLKIGEIHLR
DIDINTANEESNNILNRMIPKMLPVKIIYETDNKGNILKERPLATFYIETEETKVLKQGNFKALVKDRRLN
GLFSFAETTDLMLLEHPI SKLSDVLELKYQTRISIFEMTLGLEKLLIDKYSTLPTDSFRNMLERWLQC
KANRPELKNYVNSLIVARNAPSHNQYPMYDATALFAEVKFKTLPFSPVDTTKIELNIAPQLEIVGKAIKEI
EKSENKN

>gi|325335526|gb|ADZ11800.1| hypothetical protein RIA_0639 [Riemerella
anatipestifer RA-GD]

MEKPLLPNVYTLKHKFFWGAFLNIARHNAFITICHINEQLGLKTPSNDDKIVDVCETWNNILNNDHDL
KKSQTLTELILKHFPLTAMCYHPPKKEGKKKGHQKEQQKEKESQAQSQAELNPSKLEIALEILVNQLHS
LRNYSHYKHKPPDAEKDI FKHLKAPDASLRMVKEDYKAHPIVNLIRDFALNRKGNKQDNDPDRNRYR
FEKDGFFTESGLLFFINFLDKRDAYWMLKKVSGFKASHKQREKMITVFCRSRILLPKLRLESRYDHNQ
MLDMLSELSRCPKLLYEKLS EENKHFQVEADGFLDEIEEONPPFKDTLIRHQDRFPYFALRYLDLNE
PKSIRFQVDLGTYPHYCIYDKKIGDEQEKRLTRILLSFGRLOQDFTIENRPOEWKALTKDLDYKETSNOFF
ISKTIPIHYHITDNKIGFRGLT SKELYPSEIKDGANRI AKYPYNSGFVAHAFISVHELLPLMFYQHLIGK
SEDLKKTVRHRIQRIYKDFEERINTIEDLEKANQGRPLGAPFKQMLGLLQNKQPDLSKAKIKIEKLI
AETKLLSHRLNLIKSSPKLGGKREKLIKTVGLADWLVDPMRFPQVAYDAQNPISKANSTEFWFIR
RALALYGGKNRLEGYFKQINLIGNINPHPLNKFNWACRNLVDFYQOYLEQREKFLAIAKQWPWEYQ
YCLLLKIPKENRKNLVKGEWQGGISLPRGLPTEAIRETLESDMLSKPIRKEIKKHGRVGFISRAITLYF
KEKYQDKHQSFYNLSYKLEAKAPLLKREHHEYEYQONKQSPTEQRLELHSDRWKDYLLYKRWHLEK
KLRLYRQDVMWMLMTLELTKNHPKELNLYHQLKLENLAVNVQEADAKLPLNQLPMLVLPVKVYPATA
FGEVQYHKTPIRTVYIREEHTKALKMGNFKALVKDRRLNGLFSPIKEENDTQKHPISQLRRLRRELEIYQS
LRVDAFKETLSLEEKLLNKHISLSSLENEFRALLEEWKKEYAASSMVIDEHIAFIASVRNAPFCHNQYPPY
KEALHAPIPLFTVAQPIITEBKDGLGIAEALLKVLREYCEIVKSQI

>gi|333804399|dbj|BAK25606.1| hypothetical protein PGTDC60_1457
[Porphyromonas gingivalis TDC60]

MTEQNEKPYNGTYTLEDKHFWAAPFNLRHNAFITLAHIDRQLAYSKADITNDEDILFFKQGWKNDND
LERKARLSLILKHFSPLEGAAYGKLPESQSSGNKSSKNKELIKKEEELQANALSLDNLKSLIFDFLQ
KLDKFRNYSYHRPSESLEPLFDGNMLQRLYVNFVDSVQVRKDRHEHNDKVDPHRFHNLVVKGGKDRY
GMNDNPPFKHFFVDREGIVTEAGLLFFVSLFLEKRDIAIWMQKIRGFKGGTETYQOMINEVFCRSRISLP
KLLKLESLRIDDWMLLDMLNELVRCPSLYDRLEEDRARFRVVPVILSDEEDIDGAEEDPPKNI LVRHQD
RFPYFALRYFDLKKVFTSLRFQIDLGTYPHAIYKKNIGEQPEDRHLTRNLYGFGRIQDFAEHRPEEWR
LVRDLDFPETGDKPYITQTPHYHIEKGI GLRFVPEGLHWPSPVEGATRTRGSKYAQDKRFTAEAFS
AHELPMMPMFYFLREKYS EASAEVRQGRIRKRVIEDVYAVYDAFARDEINTRELDACLADKGRIRRH
PRQMI GILSQEHKMEEKIRKQLQEMMADTDHRLDMLDRQTDKIRIRGRKNAGLPKSGVIADWLVRDMMR
FQPVAKDTSKGPLNNSKANSTEYRMLQALALFGGKERLIPYFRQMLNIGNNPHPLHETRWESHINI
LSFYRSYLKARKAFQSIGRSRDRVENHRFLLKEPKTDRQTLVAGWKGFHLPKRGIFTAEVRDCLIEML
DEVGSYKEVGFMAKAVPLYFERACKDWQPFYNYPFNVGNLKKPKGRFSLKEKRAEWEWSEKGERFLAK
LKKEILEAKEHPYLDKFSWQKFERELRLVKNQDIITWMI CGDLMEENKVEGLDIGTLYLKDVRTDVBQGG
SLNVLNRVKPMLRPVVYRADSRRGHVHKEQAPLATVYIEERDTKLLKQGNPKSVKDRRLNGLFSFVDTG
ALAMEQYPI SKLRVEYELAKYQARVCAFEQTFLELESLLTRCPHLDPKNFRKMLESWSDPLLDKWPDLH
RKVRLLIIVARNAPSHNQYPMYDEAVFSSIRKYDPSFPDAI BERMGLNIAHRLSEEVKQAKETVERIIQA

>gi|339613716|gb|EGQ18444.1| hypothetical protein HMPREF9144_1146
[Prevotella pallens ATCC 700821]

MKEEEKGKTPVVSTYNKDDKHFWAAPFNLRHNAFITVNHINKILGEGEINRDGYENTLEKSWEIKDIN
KKDRLSKLI I KHFPFLVITVYQRNSADTTKQKEEKQAEASLES LKKSFFVFIYKLRDLRNHYSHYKHSK
SLERPKEEDLQEKMYNIFDASIQLVKEDYKHNTDIKTEEDFKHLDRKQGFYKSFADNEGNITESGLLFF
VSLFLEKDAIWMQKLLYKGLKPKQNWQYQHFLILKVKQKTRNRNLVTVGWKNSFNLPKRGIFTQPIR
LYERLREEDRKKFRVPIEIAEDDYDAEQEPFKNALVRHQDRFPYFALRYFDYNEIFTNLRFDIDLGTYPH
SIYKQIGDYKESHHLTHKLYGFERIQEFTKQNRPDWRKFKVIFNSFETSKEPYIPEITPHYHLENQKI
GTRFRNDNDKIWPSLKINS EKNKESKYKLDKSFQAEAFSLVHELLPMMFYLLLIKENTENTNDIETTKK
ENKNDKQEKHKIEEIEENKIPEIYALYDAFANGKINSIDKLEEYCKGKIDIEIGHLPKQMIAILKSEHKDM
ATEAKRQEQEMLADVQKLSLESDNQINEEINENVERKNSLSKSGEIASWLVNDMMRFQPVQKDNENGNPLNN
SKANSTEQMQLQRSALYKKEKPTRYFRQVNLIESSNPHPLNNTWEKCNLILSFYRSYLEAKKNFLP
SLKPEDWEKNQYFLMLKEPKINCEITLVQGWKNGNLPRGIFTPIRKNFMERKNIITVAELKRVGLVAVK
IPLPFS EYKDSVQPFYNYLFWNININKPDEKNFLNCEERRELLRKKKDEFKMTDKEKEENPSYLEFQS
WNKFERELRLVRNQDIVTVLMLCMEGLFNKKKIKELNVEKIYLNKININTIKKEKNTBEKNGEKIIEKNN
ILNRMIPMLRPIKVYGRNENFSKNKKKIRRNIFFTVYIEBKGTKLLKQGNFKALERDRRLGGLFSVFKTH

- continued

SKAESKSNITISKSRVEYELGEYQKARIETIKDMLALBETLIDKYNLSLIDNFHNMLIGWLKLDKDEPKAS
 PQNDVLLIIVARNAFSHQYPMRNRIFAANINPFLSSANTSEEEKGLGIANQLKDKIHKHTIEKIIEIEKPI
 IETKE

>gi|339902202|gb|AEK23281.1| Hypothetical protein Ccan_11650
 [Capnocytophaga canimorsus Cc5]

MKNIQRLGKGNFSPFKKEDKFPYGGFLNLANNNIEDFFKEITITREGIVITDENKKPKETPGEKILNEIF
 KKDISIIVYEKWNIPADYFPPTKYLSLYLEEMQFKNRVIFCRDVMKELLKTVLALRNFYTHYDHEPIKI
 EDRVFFYFLDKVLLDVSILVQNKYLTDKIKEFLNQHIGEELKELCKQRKDYLVGKGRIDKSEIINGIY
 NNAFDFIKCKREKQDDKENHNSVEKILCNKEPQNKQKSSATVWELCSKSSSKYTEKSPFNRENDDKHCLE
 VPISQKGI VFLLSFFLNKGETIYALTSNIKGFKAKITKEEPVTYDKNSIRYMATHRMFSLAYKGLKRKIR
 TSEINYNEDGQASSTYEKETLMLQMLDELNKVDPVVYQNLSEVQKTFIEDWNEYLKENNGDVGTMEEEQ
 VIHPIVRKRYEDKFNYPFAIRFLDEFAQFPILRFQVHLGNYLDCDKRTKQICDITTEEVKKKITVFGRLSE
 LENKKAIFLNEREEIKGWVFPNPSYDFPKENISVNYKDFPIVGSILDREKQPVSNKIGIRVKIADQLQR
 EIDKAIKPKKLRNPKRANQDEKQKERLVNEIVSINSNEQGEVVFVIGQPTAYLSMNDIHSVLYEFLIN
 KISGEALETKIVEKIEIETQKQIIGKDATTIKLPYTNANSNSINREKLLRDLEQEQQILKILLEEQQQE
 KDKDKKSKRKHLELPSKGVAVWLANDIKREMPKAFKEQWRGYHSHLLQKYLAYYEQSKEELNLLPK
 EVFKHFPFKLKGYPQQYLNQFYIDYLRRLSYVNEELLNIQNFKNQKDALKATEKECFKFRKQNYIIN
 PINIQIQSILVYPIFLKRGFLDEKPIIMIDREKFKENKDELADWFMHYKNYKEDNYQKFFAYPLEKVEEK
 EKFKRNKQINQKKNVYTLMMVEYIIQKIFGDKFVEENPLVVLKGIQSKAERQNNTHAATQERNLNG
 ILNQPKDKIQGKIIVKGVKLDIGNFRKYEIDQRVNTFLDYEPKKEWMLPNDWKEKEKQGLPNNV
 IDRQISKYETVRSKILLKDVQLEKIIISDEIKEEHRHDLKQKGYNYFKYIINLGLLRQLKNENVENYKVF
 KLNINPEKVNITQLKQEATDLQKAFVLTIRNKFQAHNQLPKKEFPWYCYQEKYGI EKEKTYAEYFAEVF
 KREKALIK

>gi|365737173|gb|AEW86266.1| hypothetical protein FCOL_07235
 [Flavobacterium columnare ATCC 49512]

MSSKNESYNKQKTFNHYKQEDKYFPGGLNADDNLRQVGKEFKTRINFNHNNELASVPKDYFNKEKSV
 AKREHALNLLSNYPVFLERIQKHINHNFEQTRIEIFELLDDTIKKLRDYTHHYHKPITINPKIYDFLDDT
 LLDVLIITIKKKVKNDTSRELLKEKLRPELTQLKNQKREELIKKGGKLLLEENLNAVFNHCLRPFLLEENK
 TDDKQNKTVSLRKYRKSKEETSITLIQSGLVFLMSFHLRKEFQVFTSGLEGFKAKVNTIKKEEISLN
 KNNIYMI THWSYSYNNPKLKHRIKTDQGVSTLEQNNITHSLINTNIKEALLTQIVDYLKVPNEIYET
 LSEKQKQFEEDINEYMRNPNEDSTPSSIVSHKVIKRYENKFNYPAMRFLDEYAEPLTRFMVNFPGD
 YIKDRQKILIESIQDPSERIIEKKEIHLPEKLSLVTEYKKNVYLKETSNIIDLSRFPNPNPSYVMANNIP
 FYIDRSNNLDEYLNQKKAQSNKRNLIPEKYKQSKDAI IAMLQKEIGVLDQQRSTIGLLSCNEL
 PSMLEYEVLKDIKGAELNKIAQKIREQYQSIIRDPILDSPQKDNIPPTLIKTIINTDSSVIFENQPIDIPR
 LKNAIQKELTLIQEKLNVKEHEIEVDNYNRNKNYKFKNQPKNVDDKQLQKRVVYFNRNIRQEANWLA
 SDLIHFMKNKSLWKGVMHNEIQSPLAFPEDKKNDICIALLETVPFLKEDCILTGLKKNLFLKHGNFDYK
 EYKLEKEDFLNTESTFLENGLIGLPPKILKKEKLSKRFKYIFIVFQKRFIIEKEEKKNNLYADAINLSR
 GIFDEKPIIMPFPKPNPDEFASWVASYQYNNYQSFYELTDPDIVERDKKKYKLNRAINKVKIQDYLLK
 MVDTLYQDLFNPDKLSDFYVSKAEREKIKADAKAYQKRNDSLWNVKVIHLSLQNNRITANPKLKDIG
 KYKRALQDEKIAILLTYDDRTWYALQKPEKENENDYKELHYTALNMLQYEVKVRSKLELQVQLEKQ
 ILEEYTDPLSTQIHADPFEREGNPNFKKYLAKHSILENEDDLKLEKVEAMRELDETITNPIKKAIVLI
 IIRNMAHNQYPPKFIYDLANRFVPPKBEYFATYFNRVFETITKELWENKKEKDKTQV

>gi|371640969|gb|EH006562.1| hypothetical protein HMPREF9712_03108
 [Myroides odoratimimus CCUG 10230]

MKDILTIDTTEKQNRFSHKIADKYFPGYFNLASNNIYEVFEVKNRNIFGKLAKRDNGLNKYI IHVF
 KDELISDFEKRAVAFASYFPILETVDKKSIKERNRTIDLILSQRIRQFREMILSLVTAVDQLRNFYTHY
 HSDIVIENKVLDFLNSSFVSTALHVKDKYLKTDKIKEFLKETIAAELDILIEAYKKKQIEKKNTRPKAN
 KREDILNAINYNAFWSFINDDKDKDKETVVAAGADAYFEKNHHSNDPDPALNISSEKGIYVLLSFFLIN
 KEMDSLKANLIGFGKGVDRSNGSIKYMATORISYFHTYRGLKQKIRTSSEGVKEILLMQMIDELSKVPN
 VYQHLSTIQQNSFIEDWNEYKYDVEDDVEITDLSRVIHPIVRKRYEDRFNYPFAIRFLDEFFDFPIILRFQ
 VHLGDIVVHRRIKQLGKVESDRIIEKVIIVFARLKDINSAKASYFHSLEEQDKBELDNKWLIFPNPSYDF
 PKEHTLQHQGEQKAGKIGIYVKLRTQYKKALEEARSLNPKERSATKASYDIITQIIEANDNVK
 EKPLVFIGQPIAYLSMNDIHSMLFSLIDNAELKKTPEEVAKLIDQIGKQINEILSKDIDTKILKKYK
 NDLKETDIDKTRDLARDKEIEKLELQKQRAADDYNYISSTKFNIDKSRKRLHLLFNAEKGKIGVWLAN
 DIKRFMPKESKWKGYQHTLQKLFAYFDTSKSDLELILSMVMVKDYPIELIDLVKKSRILVDFLNKY
 LEARLEYIENVI TRVKNISGTPQPKIVRKECFPLKKSNTYVTVSLDKQVERILSMPLFIERGFMDDKPI
 LEGKSYKQKKEKPADWVHYKENSNYQNFYDIEVYEITTEDKREKAKVIKKIKQQKNDVFLMMVNYML
 BEVLKSSNDRLSLNELYQTKERIVNKQVAKDTQERNKNIWNKVVDLQLCDGLVHIDNVKLDIGNFR
 KYENDSRVKEPLTYQSDIVWSAYLSNEVDSNKLVIYIERQLDNYESIRSKELLKEVQETIECSVYQVANKE
 SLKQSGNENFKQYVLLQGLLPIGMDVREMLILSTVVKFKKEEIIQLGQAGEVEQDLYSLIYIRNKFQAHNQL
 PIKEFFDFCENNYRSISDNEYAEBYMEIFRSIKEKYAN

>gi|380461034|gb|AFD56718.1| hypothetical protein RA0C_1842
 [Riemerella anatipestifer ATCC 11845 = DSM 15868]

MEKPLPHVYTLKHKFFWGAFLNIARHNAFITICHINEQLGLKTPSNDDKIVDVCETWNNILNNDHDL
 KKSQLELILKHFPLTAMCYHPKKEGKGGHQQEKESEASQAEALNPSKLIIEALEILVNQLHS
 LRNYSHYKHKKPAEKDIFKHYKAFDASLRMVKEDYKAHFI VNLIRDFAHNRKKGKQNDPDNFRY
 FEKGGFTESGLLPIINFLDKRDAYWMLKVKVSGFKASHQREKMIIEVFCRSRILPLKRLRESRYDHNQ
 MLLDMLSELSCRCPKLLYEKLS EENKHFQVEADGFLDEIEEQNPFKDTLIRHQDRFPYFALRYLDLNE
 FKSIRFQVDLGTNYHYCIYDKKIGDEQEKRLHTRILLSFGLQDFTEINRPQEWKALTKDLDYKETSNOFF
 ISKTI PHYITDNKIGFRLGTSKELYP SLEIKDGNRI AKYYPNSGFVAHAFI SVHELLPLMFYQHLIG
 SEDLLKETVRHQRIRYKDFEERINTIEDLEKANQGRPLGAFPKQMLGGLQNKQPDLESEKAKIKIEKLI
 AETKLLSHRLNIKLSKPKLGRREKLIKTVGLADWLVKDFMRFPVAYDAQNP IKSKANSTEFWFIR
 RALALYGGKRNLEGYFQINLIGNINPHPLNFKFNWACRNLDVDFYQQYLEQREKFLKAIKNQWPEPYQ
 YCLLLKPIKRNKLNKVGWEQGGISLPRGLFTEARETLESDLLMSKPIRKEIKHGRVGFISRAITLYF

- continued

KEKYQDKHQSFYNLSYKLEAKAPLLKREEHYEYQONKQSPTEQRLELHTSDRWKDYLLYKRWQHLEK
 KLRLYRNQDVMLWMLTLELTKNHPKELNLYHQLEKLENLAVNVQEADAKLNPLNQTLPMLVLPVKVYPATA
 FGEVQYHKTPIRTVYIREEHTKALMGNFKALVKDRRLNGLFSFIKEENDTQKHPISQLRLRRELEIYQS
 LRVDAPFKETLSLEKLLNKHISLSSLENEFRALLEEWKKEYAASSMIDEHIAFIASVRNAPFCHNQYFPY
 KEALHAPIPLFTVAQPITEEKDGLGIAEALLKVLREYCEIVKSQI

>gi|386374623|gb|AFJ07523.1| hypothetical protein PIN17_0200

[*Prevotella intermedia* 17]

MKMEDDKKIKESTNMLDNKHPWAAPLNLARHNVIYTVNHINKVLELKNKQDQDIIDNDQDILAIKTHWE
 KVNGDLNTERLRELMTHKFPFLETAIYIKNKEDKEEVKQEKQAQSFDSLKHCFLFLFEKLQEARNY
 SHYKSESTKEPMLKELLLKMYNI FDDNIQLVIKDYQHNDINPDEDFKHLDRTEEEFNYYFTINKKGN
 ITASGLLPFVSLFLKEDDAIWMQKLRGFKDNRESKKNMTHVEVFCRSRMLLPKLRLESTQTQDWILLDML
 NELIRCPKSLYERLQGEYRKKENVPDFSADEYDAEQEPFKNILVRHQDRFPYFALRYFDYNEIFTNLRP
 QIDLGTYHFSYKLLIGGQKEDRHLTHKLYGFERIQEFKQNRIDEWKAIVKDFDIYETSEEPYISETAP
 HYHLENQKIGIRFRNDNDEIWPSLKINGENNEKRYKLDKQYQAEAFLSVHELPMMPYLLLLKKEEPPNN
 DKKNASIVEGFIKREIRDIYKLYDAFANGEINNIDDELEKYCEDKGIKPKRHLPKQMVAILYDEHKDMABEA
 KRKQKEMVKDKKEYLKKDAQYNIINNCVQPPYGFYFNVGNHHPKDEKDFLPSSEERKKLWGDKKYKPKGY
 KAKVSKKLLIDKEKEEYSYLFQSWNKFERELRLVRNODIVTWLLCTELIDKLVGEGLNVEELKCLRK
 DIDIDTAKQEKNNILNRVMPQLPVIVYEIDDSHNIKVRDPLHIVYI BETKIKLLKQGNFKALVKDRRLN
 GLFSFVDTSSSETLEKSNPISKSLVEYELGEYQNAIETIKDMLLEETLIEKYKTLPIIDNPSDMLNGWLE
 GKDEADKARFQNDVLLVAVRNAPSHNQYPMRNRIFANINPFLSSADTSEKKLDIANQLKDKIHKII
 KRITIEIEKPIETKE

>gi|392610411|gb|EIW93190.1| hypothetical protein HMPREF1322_2050

[*Porphyromonas gingivalis* W50]

MTEQNEKPYNGTYYTLEDKHPWAAPFNLARHNAYITLAHIDRQLAYSKADITNDEDILFPFKQWKNLND
 LERKARLRSLILKHFSLFLEGAAYGKLPESQSSGNKSSKKEELIKKEEELQANALSNDLKSILDFPLQ
 KLKDFRNYSHYRHPSESELPLDFGNMLQRLYNVFDVSVQRVKRDHEHNDKVDPHRHFNHLVRKGGKDKY
 GNNDNPPFKHHPVDREKVTAEAGLLFFVSLFLEKRDIAWMQKIRGPKGTEAYQMTNEVFCRSRISLP
 KLLKLESLRIDDWMLDMLNELVRCPKSLYDRLEEDRARFRVVDILSDEDDIDGTEEDPFKNIIVRHQD
 RPPYFALRYFDLKKVFTSLRPHIDLGTYHFAIYKKNIGEQPEDRHLTRNLVYGFGRIQDFAEEHRPEEWR
 LVRLDLYPETGDKPYITGTTPHYHIEKGIKGRFVPEGLWSPPEVGAIRTRGRSKYAQDKRFTAEAFSL
 VHELMPMFFYFLREKYS EEAESA EKVGRIKRVIEDVYAVYDARFDEINTRELDACLADKGIIRRGHL
 PRQMIILSQAHEKDMEEKVRKLEMIADTDHRLDMLDRQTDRIKIRGRKAGLPKSGVIAWDLVRDMMR
 FQPVAKDTSGKPLNNSKANSTERYMLQALALFGGEKERLIPYFRQMNILGNNPFPFLHETRWESHINI
 LSFYRSYLKARKAFIQSITGRSDRENRHFLKKEPKTDRQTLVAGWKSSEFHLPRGIFTEAVRDCLIMGY
 DEVGYSKEVGFMAKAVPLYFASRACKDRVQPFYDYPFNVGNLSLKKGRFLSKEKRAEEWESGKERFRDLE
 AWSHSAARRIEDAFVGI EYASWENKKEIEQLLQDLSLWETPESKLKVKADKINIAKLLKIELEAKEHPYH
 DFKSWQFERELRLVKNQDIIITWMMCRDLMEENKVEGLDIGTLYLKD IRTDVQEQGSLNVLNHVKPMRLP
 VVVYRADSRRGHVHKEEAPLATVYIEERDTKLLKQGNFKSFVKDRRLNGLFSFVDTGALAMEQYPI SKLRV
 EYELAKYQTARVCAFEQTLLEESLLTRYPHLPDES FREMLESSWSDLLDKWPDLOREVRLLIARVNAFS
 HNQYPMYDETISSIRKYDPSLDAIEERMGLNIAHRLSEEVKLAKEMVERIIQA

>gi|392612061|gb|EIW94777.1| hypothetical protein HMPREF1322_1926

[*Porphyromonas gingivalis* W50]

MNIVPASENKGQSRIVEDDPQYFGLYLNLARENLEEVESHVRIKFGKKNLNEESLQSLDCDHLLSVDRW
 TKVYGHSSRRYLPFLHYFPDSDSQIEKDHDSTKGVDPDSAQRLIRELYSLLDLFRNDFSHNRLDGTFEHL
 VSPDISSTITGTYSLACGRAQSRFADFPPKDDFVLAKNRKEQLISVADGKECLIVSGLAFFICLPLDREQ
 ASGMLSRIRGFKRIDENWARAVHETFCDLICIRHPDRLESNNI KEALLDMLNELNRCPRILYDMLPEEB
 RAQFLPALDENSMNLS ENSLNLESRLWGDGSDWAEALTKIRRHQDRFPYMLLRFIEEMDLLKGI RFRV
 DLGETELDSYKVKGRNGEYDRTITDHALAFGKLSDFQNEEEVSRMISGEASYPVRFSLFAPRYAIYDNK
 IGYCHTSDPVYPKSKTGEKRLSNPQSMGFI SVHNLRKLLMELCEGFSFRMQSDFLRKANRILDETAE
 GKLFQFALFPFMRHRPIPPQNPKSKDRREKAETTLKQYKQETKGRKDKLNSQLLSAFDMNQRLPSRLLD
 EWMNIRPASHSVKLRITYVQKLNEDCRLRLRKRFRKDGDKARAIPLVGMATFSLSDIVRMI ISEETKLI
 TSAYYNEMQRS LAQYAGEENRRQFRAIVAEHLHLLDPSSGHPFLSATMETAHRYTEDFYKCYLEKKREWLA
 KTFYRPEQDENTKRRISVFPVDPGEARKLLPLLI RRRMKEQNDLQDWIRNKQAHPIDLPSHLFDSKIMEL
 LKVKDGGKKNWNAEAFKDWSTKYVPGMQPFYGLRRELNIHGKSVSYIPSDGKFKADCYTHLMKTVQDKKR
 ELRTAGKVPVPLAADIKRSFHRAVNEREFMLRLVQEDDRMLMAINKMMDREEDILPGLKNIDSLDE
 ENQFSLAVHAKVLEKEGEGDNLSSLVLPATIEIKSKRKDWISKYIRYR YDRRVPLGMSHFPEHKA ILDEVK
 ILLGEYDRCRKIIFDWAFALEGAIMSDRDLKPYLHESRSREGKSGEHS TLVKMLVEKKGCLTPDESQYLI
 LIRNKAHNQFPFAAEMPLIYRDVSAKVGSI EGSSAKDLPEGSSLVDSLWKKYEMIIRKILPILDPENRF
 FGKLLNNMSQPINDL

>gi|400374977|gb|EJP27887.1| hypothetical protein HMPREF1146_2324

[*Prevotella* sp. MSX73]

MQKQDKLFPVDRKKNIAIFAPPKYITIMENQEKPEPIYELTDKHPWAAPLNLARHNVIYTVNHINRRELEIA
 ELKDDGYMMGIKGSWNEQAKKLDKVRRLRDLIMKHPFLLEAAAYEITNSKSPNNKEQREKQSEALS LNN
 LKNVLFIFLEKLVLRNYYSHYKSEESPKPIFETSLKNMYKVF DANVRLVKRDMHHEINDMQDRFTH
 LNRKKQVGRINKNIDSPNFHYHFADEKGNMTIAGLLFFVSLFLDKKDAIWMQKLLKGFKDRGNLREQMIN
 EVFCRSRISLPKLELVNQTQKDWMLDMLNELVRCPKSLYERLREKDRSFVFPDFISDDYDAEEEPFK
 NTLVRHQDRFPYFVLYRFDLNEIFQLRFQIDLGTYHFSIYKRIIGDDEVRHLTHLYGFAIQDFAPQ
 NQPEEWRKLVKDLDFHETSQEPYISKTAPHYHLENEKIGIKFCSTHNNLFP SLKREKTCNGRSKFNLGIQ
 FTAEAFSLVHELPMMPYLLLLIKDYSRKESADKVEGIRKETSNIYAIYDAFANNEINSIADLIRLQK
 TNILQGHLPKQMSITLEGRQKDMKEAERKIGEMIDDIQRRLDLLCKQINQKIRIGKRNAGLLKSGKTD
 WLVSDDMMRFPQVQKDIINAPINNSKANSTERYMLQALALFGSESSRLKAYFRQMNVLGNANPHPLAET

- continued

QWEHQTNILSFYRNYLEARKKYLKGLKPKQNWKQYQHFLILKVKQKTNRNTLVGWKNSFNLPRGI FTQPIR
 EWFEKHNSKRIYDQILSFDREVGFVAKAIPLYFAEEYKDNVQPPFYDYPFNI GNKLPKQKQGLDKKERVE
 LWQKNKELFKNYPS EKNKTDLAYLDFLSWKKFERELRLIKNQDIVTWLWFKELFKITIVVEGLKIGIHLR
 DIDINTANEESSNINLRIMPMPKLPVKIYETDNKGNILKERPLATFYI BETETKVLKQGNFVKLAKDRRLN
 GLLSFAETTDIDLEKNPITKLSVDYELIKYQTRRISIFEMTLGLEKKLIDKYSTLPTDSFRNMLERWLQC
 KANRPELKNYVNSLJAVRNAPSHNQYPMYDATALFAEVKFTLFPSSVDTKKI ELNIAQLLEIVGKAIKEI
 EKSENKN

>gi|405580121|gb|EKB54193.1| hypothetical protein HMPREF9699_02005

[*Bergeyella zoohelcum* ATCC 43767]

MENKTSLGNNIYNNPFKPKQDKSYFAGYFNAAMENTSVPRELKGRKLGKEYISENFDAIPKENISLVEY
 ERYVKLLSDYFPMARLLDKKEVPIKERKENPKKFKGI IKAVRDLRNFYTHKEHGEVITDEIFGVLDDEM
 LKSTVLIVKVKKVIDKIKELKKSIEKQLDILCQKKLEYLRDTARKIEEKRRNQREGERKELVAPPKYS
 DKRDDLIAAIYNDADFVYIDKKDSLKESKAKYNIKSDPQQEEDGLKPI SKNGVVFLLSLFLIKQEIH
 AFKSKIAGFKATVIDEATVSEATVSHGKNSICFMATHEIFSHLAYKKLKRVRTAEINYGAEENAEQLSV
 YAKETLMMQMLDELSEKVPDVVYQNSLSEVQKTFI EDWNEYLKENNGDVGTMEEBQVIHPVIRKRYEDKFN
 YFAIRFLDEFAQFPILRFQVHLGNVLDHRSRKPENLISDRRIKEKITVFGRLSELEHKKALPIKNTETNED
 REHYWEIFPNPNYDFPKENISVNDKDFPIAGSILDREKQPVAGKIGIKVLLNQYQVSEVDKAVKAHQK
 QKASKPSIQNIIEEIVPINESNPKEAIVEGGQPTAYLSMNDIHSILYEFDFKWEKKEKLEKKGEKELR
 KEIGKELEKKIVGKIQAQIQOIQIIDKIDINAKILKPYQDGNSTAI DKEKLIKDLKQEQNI LQKLDQVRE
 KEYNDFIAYQDKNRINKVRDRNRHKQYLKDNLRKRYPEAPARKEVLYREKGVAVWLANDIKREMPDF
 KNEWKGEQHSLLQKSLAYEQQEELKNLLPEKVPQHLFPKLGQYQKYLQFYTCYLDKRLIYISGLV
 QQAENFSENKVFKEVNECPFKLKKQNYTHKELDARVQSILGYPILFERGFMDKPTI I KGIKFKGNEA
 LFADWFRYKQNFQTPYDTEHYPLVELEKQADRKRKTKIYQKKNQVDFLLMAKHIKSVFQKQSDID
 QPSLEDLYQSREERLGNQERARQTGERNTNYIWNKTVDLKLCDGKITVENVKLNVDGFIKYEYDQVQA
 FLKYEENI EQAFLIKESKEEENYVYVEREIEQYKVRREELKKEVHIEEYILEKVKDEILKKGDNQ
 NFKYIYLNGLKQLKNEVDVESYKVFNLNTEPEDVNIQLKQEAATDLEQKAFVLTIRNKAHNQLPKKEF
 WDYQEQYKGIIEKKEYAEPFAEVFKKEKEALIK

>gi|408468850|gb|AFU69194.1| hypothetical protein P700755_002426

[*Psychroflexus torquis* ATCC 007555]

MESIIGLGLSFPYKTAADKHYPGSLNVLNNAVFAEFKERISYKAKDENISLIEKHFIDNMSIVDY
 EKKISILNGYLPIDDFLDELNNLNTRVKNFKNFII LAEAI EKLRDYYTHFYHDPITFEDNKEP LLEL
 LDEVLLKTLIDVKKYKTKDKIKEIKLDSLREEMDLVIRKTDREKKKINPKIQHTDSSQIKNSIFND
 AFQGLLYEDKGNKKTQVSHRAKTRLNPKD IHKQBERDFEIPLSSTGLVFLMSLFLSKKEIEDPKSNIKG
 FKGKVVKD ENHNSLKYMATHRVYSILAFKGLKRYIKTDIFSKETLMMQIMDELSKVPDPCVYQNLSETKQK
 DFI EDWNEYFKDNEENTENLENSRVVHPVIRKRYEDKFNYPFAIRFLDEBFANFKILKQVFMGYIHDQRT
 KTI GTTNI TTBERTVKEKINVPKSKMDNLKHFPSQLSDDENTDWEFFPNPSYNFLQADNSPANNIP I
 YLELKNQOIIIEKDAIKAEVNTQNRNPNKPSKRDLLNKILKTYEDFHQGDPTAILSLNEIPALHLPLV
 KPNNTGQQIENIRIKIEKQPKAINHPSKNNKGIKPSLADINVRVNAIKLKDLEAELMDLNKHHIAF
 KENQKASSNYDKLKEHQFTPKNKRPELRKVVYKSEKGEATWLANDIKREMPKDFIKWKQCSHELQ
 RKLAFYDRHTKQDIKELLSGC EFDHSLLDINAYFQKDNFEDFFSKYLENRIETLEGLVKLKHDFKNEPTP
 LKGVFKNCFKFLKRNQNYVTESEPIIKRILAKPTFLPRGVDERPTMKKGNPLKDKNEFAEWFVEYLEN
 KDYQKFNABEYRMRDADFKNNAVKKQKLDKDFYILQMVNLLKEVPGKDEMNQLSLEFQTRQERLKLQ
 GIACKQMNKETSSENTRNQTYIWNKDVVPSFFNGKVTIDKVKLKNYKGYKRYERDERVKIFIGYEVDE
 KWMMYLPHNWKDRYSVKPINVIDLQIQEYEBIRSHELLEIQNLEQYIYDHI TDKNILLQDGNPNFMKYV
 LNGLLIGIKQVNI PDIIVLQKMINFDKIDDFIGIASCELEKIIILIAIRNKAHNQLPNKMIYDLANEF
 LKIEKNETYANYLLKVLKMI SDLA

>gi|441484656|gb|AGC41342.1| hypothetical protein G148_2038

[*Riemerella anatipestifer* RA-CH-2]

MEKPLLPNVYTLKHKFFWGAFLNIARHNAFITICHINEQLGLKTPSNDDKIVDVCETWNNILNNDHLL
 KKSQLELILKHFFLTAMCVHPKKEGKKGHQKEQQEKESEASQAEALNPSKLEIALEILVNLQHS
 LRNYSHYKHKKPAEKDIFNSLCILNNTDF

>gi|441484658|gb|AGC41344.1| hypothetical protein G148_2040

[*Riemerella anatipestifer* RA-CH-2]

MFFSFHNAQRVIFKHLKAFDASLRMVKEDYKAHFIVNLRDFAHLNRKGNKQDNPDFNRYRFEKDGFF
 TESGLLFPINFLDKRDAYWMLKVKVSGPKASHKQREKMI TEVFCRSRILLPKRLRESYDHNQMLLDMLS
 ELSRCPKLLYELKSEENKHFQVEADGFLDEIEEBQNPFKDTRLIRHQDRFPYFALRYLDLNEFSKIRFQ
 VDLGTYHYCIYDKKIGDEQEKRLTRILLSFGRLQDFTIENRQEWKALTKDLDYKETSNOFFISKTIPH
 YHITDNKIGFRLGTSKELYPSLEIKDGANRIKYPYNSGFVAHAFISVHELLPLMFYQHLIGKSEDLLEK
 TVRHITQRIYKDFEERINTIEDLEKANQGRPLGAFPKQMLGLLQNKQPDLESEKAKIKIEKLI AETKLLS
 HRLNIKLSKSPKLGKRRKELIKTVGLADWLKDFMRFPQVAYDAQNP I KSSKANSTEFWFIRRALALYV
 GEKRL EGYFFQINLIGNINPHPLNKNFNWACRNLVDFYQYLEQRKFLAIAKHQWPPEYQYCLLLKV
 PKENRKNLVKGEWQGGISLPRGLFTEAIRETSLKDLILSKPIRKEIKKHGRVGFISRAITLYFKKQYQK
 HQSFYNLSYKLEAKAPLLKKEHEHYEQNQKPSPTESQRLELHTSDRWKDYLLYKRWQHLEKLLRLYRN
 QDIMLWMTLELTKNHFKELNLYHQLKLENLAVNVQEADAKLNPLNQLTLPMLVVKVYPTTAFGEVYH
 ETPIRTVYIREQTKALMGNFKALVKDRRLNGLFSFIKEENDTQKHPISQLRLRRELEIYQSLRVDFAK
 EILSLEKLLNKHASLSLENEFRILLEEWKKYAASMVIDKHIAFIASVRNAPFNHQPYPFYKETLHAP
 ILLFTVAQPI TEEDGLGIAEALLKVLREYCEIVKSI

>gi|482527725|gb|E0A10535.1| hypothetical protein A343_1752

[*Porphyromonas gingivalis* JCVI SC001]

MTEQNEKPYNGTYTLEDKHFWAFAFNLARHNAYITLTHIDRQLAYSKADITNDEDILFFKQGWKNLND
 LERKARLRLSILKHFSPLEGAAYGKLFESQSSGNKSSKKKELIKKEKEELQANALSLDNLKSLIFDFLQ
 KLDKFRNYSHYRHPSEBELPLFDGNMLQRLYVFDVSVQRVRRHEHNDKVDPHRFNHLVRKGGKDKR
 GNNNDNPFKHFVDREKVTAEAGLFFVSLFLEKRDALWQKIRGFKGGTETYQOMTNEVFCRSRISLP

- continued

KLKLESRLIDDWMLDMLNELVRCPSLYDRLEEDRARFRVVDILSDEDDIDGTEDDPFKNI LVRHQD
 RFPYFALRYFDLKKVFTSLRPHIDLGTYHFAIYKKNIGEOPEDRHLTRNLYGFGRIQDFAEHRPEEWR
 LVRDLDFETGDKPYITQTPHYHIEKGI GLRFVPEGQLWPSPEVGA TRTGRSKYAQDKRFTAEPFLS
 VHELMPMMFYFLREKEYSEASAERVQGRIRKVI EDVYAVYDAFARGEIDILDRDLACLADKGI RRGHL
 PRQMIALLSQEHKDMEEKVRKLLQEMIA DTDHRLDMLDRQTD RKRIRI GRKNAGLPKSGVIADWLVRDMMR
 FQPVAKDTSGKPLNNSKANST EYRMLQRALALFGGKERLI PYFRQMNLI GGNPHPPFLHETRWESHINI
 LSFYRSYLKARKAPLQSIGRSDRVENHRFLLKPKTDRQTLVAGWKGFEHLPRGIFTEAVRDC LIEMGL
 DEVGSYKEVGFMAKAVPLYFERACKDRVQPPYDYPFNVGNLSLKKGRFLSKERAEWESGKERFRDLE
 AWSHSAARRIEDAFAGIENASRENKKIEQLLQDLSLWETPESKLVKADKINIAKLKKEILEAKEHPYL
 DFKSWQKFERELRLVKNQDII TWMMCRDLMEENKVEGLDITGLTYLKD IRTDVHEQGS LNVLNRVKPMLRP
 VVVYRADSRGHVHKQAPLATVYIEERDTKLLKQGNFKSFVKDRRLNGLFSFVDTGALAMEQYPI SKLRV
 EYELAKYQ TARVCAFEQTL ELESLLTRYPHLPDKNFRKMLESWSDP LLDKWPDLHGNVRLLI AVRNAFS
 HNQYPMYDETFLSSIRKYDPSPDIAIEERMGLNIAHRLSEEVKQAKEMVERIIQA

>gi|488741127|ref|WP_002664492.1| hypothetical protein [*Bergeyella*
zoohelcum]

MENKTSLGNNIYYNPFKPKDQKSYFAGYFNAAMENTS VVRELGKRLKKEYISENFPDAIFKENISLVEY
 ERYVKLLSDYFPMARLLDKKEVPIKERKENFKKNFKGIKAVRDLRNFYTHKEHGEVEITDEIFGVLD
 LKSTVLIVKKKVVKIDKIKEIKKSLKQDLILCQKKLEYLRDTARKIEEKRRNRERGEKELVAPFKYS
 DKRDLIAAIIYNDADFVYIDKKKSLKESKAKYNIKSDPQOEEGDLKIPISKNGVVFLSLFLIKQEI
 AFKSIAGFKATVIDEATVSEATVSHGKNSICFMATHEIFSHLAYKLLKRVRTAEINYGAEANAQLSV
 YAKETLMMQMLDEL SKVPDVVYQNLSEVQKTFIEDWNEYLKENNGDVGMT EEEQVIHPVIRKRYEDKFN
 YFAIRFLDEFAQFPIILRFQVHLGNVLDHRSRKENLISDRRIKKEITVFRGLSELEHKKALPKNTETNED
 REHYWEIFPNPNYDPPKENSIVNDKDFPIAGSILDREKQPVAGKIGIKVLLNQYVSEVDKAVKAHQLK
 QRKASKPSIQNIEIEIVINIESNPKEAIVEGGQPTAYLSMNDIHSILYEFDPKWEKKKEKLEKKEKELR
 KEIGKELEKKIVGKIQAQIQIIDKINAKILKPYQDGNSTAIKKEKLIKDLKQEQNI LQKLDKDEQTVRE
 KEYNDPIAYQDKNRINIKVRDNHKQYLKDNLKRKYPEAPARKEVLYYREKGVAVWLANDIKREMPD
 KNEWGEGHSLQKSLAYEYEQCKELKNLLPEKVFQHL PFKLGGYFQKQYLYQFYTCYLDKRL EYISGLV
 QQAENFKSENKVFKKVENECKFKLKKQNYTHKELDARVQSILGYPIFLERGFMDKPTI IKGKIFKGN
 LFADWFRYKYEQNFQTPYDTE NYPLVELEKQADRKRKTKIYQKKNVDFILMAKHIKFSVFKQDSID
 QPSLEDLYQSRERLGNQERARQTGERNTNYIWNKTVDLKLCDGKITVENVKLNVDGFIKYEYDQRVQA
 FLKYEENIEWQAPLIEKSEKEENYVYVEREIEQYKVRREELLKEVHLIEEYILEKVKDEILKKGDNQ
 NFKYIYLNLGLKQLKNEDEVESYKVFNLNTEPEDVNIQLKQEA TDLQKAFVLT YIRNKFAHNQLPKKEF
 WDYCQEKYKIEKKEKYAEYFAEVFKKEKEALIK

>gi|490473137|ref|WP_004343581.1| hypothetical protein [*Prevotella*
buccae]

MQKQDKLFVDRKKNIAIFAPPKYITIMENQEKPEPIYYELTDKHFWAAPLNARHNVTY TINHINRRLEIA
 ELKDDGYMMDIKGSWNEQA KKLKVKRRLRDLIMKHFPFLEAAAYEITNSKSPNNKEQREKQSEALSLNN
 LKNVLFIFLEKLQVLRNYYSHYKYESPSPKPIFETSLLKNMYKVF DANVRLVKRDYMHHENIDMQRDPTH
 LNRKKQVGRINKNIDSPNFHYHFADKEGNMTIAGLLFFVSLFLDKKDAIWMQKLLKGFKDGRLREQMIN
 EVFCSRISLPKLELVNVTQKDNWQDLMLNELVRCPSLYERLREKDRRESFKVPDFISDDYDAEEEPFK
 NTLVRHQDRFPYVFLRYFDLNEIFEQLRFQIDLGTYHFSIYNKRIGDEDEVRLTHHLYGFARIQDFAQQ
 NQPEVWRKLVKDLDFEASQEPYIPKTAHYHLENEKIGIKFCSIHNNLPPSLKTEKTCNGRSKFNLGIQ
 FTAEAFSLVHELLPMMFYLLI KDYSRKESADKVEGIRKEISNIYAIYDAFANGEINSIADLCRLQD
 TNILQGHLPKQMISILEGRQDMKEAERKIGEMIDDIQRRLDLLCKQINQKIRIGKRNAGLLKSGKIAD
 WLVNDMMRFQPVQKQNNIPIINNSKANST EYRMLQRALALFGSENFRILKAYFNQMNLVGNDNPHFLAET
 QWEHQTNILSFYRNYLEARKKYLKGLKPKQNWQYQHFLILKVKQKTRNNTLV TGWKNFNLPRGIFTPQPIR
 EWFKEHNSKRIYDQILSFDRVGFVAKAIPLYFAEEYKDNVQPPYDYPFNIGNNLSKPKQKGLDKKERVE
 LWQKNKELFKNYPSEKKKIDLAYLDFLSWKKFERELRLIKNQDIVTWLMFKELFNMATVEGLKIGEIHLR
 DIDINTANEESNNILNRIMPMLKPVKIIYETDNKGNILKERPLATFYIEETETKVLKQGNFKVLAKDRRLN
 GLLSFAETTDIDLEKNPIITKLSVDHELIKYQTRISIFEMTLGLEKKLINKYPTLPTDSFRNMLERWLQC
 KANRPELKNYVNSLI AVRNAFSHNQYPMYDATLFAEVKKFTLFPSSVDTKKEI ELNIAPQLLEIVGKAIKEI
 EKSENKN

>gi|490473535|ref|WP_004343973.1| hypothetical protein [*Prevotella*
buccae]

MQKQDKLFVDRKKNIAIFAPPKYITIMENQEKPEPIYYELTDKHFWAAPLNARHNVTY TINHINRRLEIA
 ELKDDGYMMDIKGSWNEQA KKLKVKRRLRDLIMKHFPFLEAAAYEIMINSKSPNNKEQREKQSEALSLNN
 LKNVLFIFLEKLQVLRNYYSHYKYESPSPKPIFETSLLKNMYKVF DANVRLVKRDYMHHENIDMQRDPTH
 LNRKKQVGRINKNIDSPNFHYHFADKEGNMTIAGLLFFVSLFLDKKDAIWMQKLLKGFKDGRLREQMIN
 EVFCSRISLPKLELVNVTQKDNWQDLMLNELVRCPSLYERLREKDRRESFKVPDFISDDYDAEEEPFK
 NTLVRHQDRFPYVFLRYFDLNEIFEQLRFQIDLGTYHFSIYNKRIGDEDEVRLTHHLYGFARIQDFAQQ
 NQPEVWRKLVKDLDFEASQEPYIPKTAHYHLENEKIGIKFCSAHNNLPPSLQTDKICGKRSKFNLGIQ
 FTAEAFSLVHELLPMMFYLLI KDYSRKESADKVEGIRKEISNIYAIYDAFANNEINSIADLTRRLQD
 TNILQGHLPKQMISILEGRQDMKEAERKIGEMIDDIQRRLDLLCKQINQKIRIGKRNAGLLKSGKIAD
 WLVNDMMRFQPVQKQNNIPIINNSKANST EYRMLQRALALFGSENFRILKAYFNQMNLVGNDNPHFLAET
 QWEHQTNILSFYRNYLEARKKYLKGLKPKQNWQYQHFLILKVKQKTRNNTLV TGWKNFNLPRGIFTPQPIR
 EWFKEHNSKRIYDQILSFDRVGFVAKAIPLYFAEEYKDNVQPPYDYPFNIGNNLSKPKQKGLDKKERVE
 LWQKNKELFKNYPSEKKKIDLAYLDFLSWKKFERELRLIKNQDIVTWLMFKELFNMATVEGLKIGEIHLR
 DIDINTANEESNNILNRIMPMLKPVKIIYETDNKGNILKERPLATFYIEETETKVLKQGNFKVLAKDRRLN
 GLFSFAETTDLNLBEHPISKLSVDLELIKYQTRISIFEMTLGLEKKLIDKYSTLPTDSFRNMLERWLQC
 KANRPELKNYVNSLI AVRNAFSHNQYPMYDATLFAEVKKFTLFPSSVDTKKEI ELNIAPQLLEIVGKAIKEI
 EKSENKN

- continued

>gi|491058119|ref|WP_004919755.1| hypothetical protein [*Riemerella anatipestifer*]
 MEKPLLPNVYTLKHKFFWGAFLNIARHNAFITICHINEQLGLKTPSNDDKI VDVVCETWNNILNNDHDL
 KKSQTELELILKHPFLTAMCYHPPKKEGKKKGHQKEQQKEKESQAQAEALNPSKLI EALEILVNQLHS
 LFNYYSHYKHKKPAEKDI FKHLKAFDASLRMVKEDYKAHFI VNLIRDFALHNRKGNKQNDPDPNRYR
 FEKDGFFTESGLLFFINLFLDRDAYWMLKKVSGFKASHKQREKMITVEVFCRSRILLPKLRLESYDHNQ
 MLLDMLSELSRCPKLLYEKLESEENKHPQVEADGFLDEIEEQNPFKDTLIRHQDRFPYFALRYLDLNE
 FKSIRFQVDLGTTHYCIYDKKIGDEQEKRLHTRILLSPGRLQDFTEINRPQEWKALTKDLDYKETSNO
 P I SKTIPHYHI TDNKGIFRLGT SKELYPSEIKDGANRI AKYPYNSGFVAHAFI SVHELPLPMFYQHLIGK
 SEDLLKETVRHRIQRIYKDFEERINTIEDLEKANQGRPLGAPPKQMLGLLQNKQPDLESEKAKIKI EKL
 I AETKLLSHRLNI KLKSSPKLGRREKLIK TGVLADWLVDPMRFPQVAYDAQNP I KSKANSTEFWPIR
 RALALYGGKRNREGYFKQINLIGNINPHPLNKNFNWACRNLDVFPYQQYLEQREKFL EAKNQPEWYQ
 YCLLLKIPKENRKNLVKWEQGGISLPRGLFTEAIRETLESDMLSKPIRKEIKKHGRVGFISRAITLYF
 KEKYQDKHQS FYNLSYKLEAKAPLLKREBEHYEWQNKQPSPTESQRLELHSDRWKDYLLYKRWQHLEK
 KLRLYRNQDVMLMTLELTKMHPKELNLYHQLKLENLAVNVQEADAKLNLPLNQTLPMVLPVKVYPATA
 FGEVQYHKTPIRTVYIREEHTKALKMGNFKALVKDRRLNGLFSFIKEENDTQKHPISQLRLRRELEIYQS
 LRVDAFKETLSLEEKLLNKHISLSLENEFRALLEEWKKEYAASSMVIDEHIAFIASVRNAPCHNQYFPY
 KEALHAPIPLFTVAQPI TBEKDG LGIAEALLKVLREYCEIVKSMI

>gi|492521417|ref|WP_005873511.1| hypothetical protein [*Porphyromonas gingivalis*]
 MNIVPASENKGQSRIVEDDPQYFGLYLNLARENLEEVESHVRIKFGKKKLNEESLKQSLCDHLLSVDRW
 TKVYGHSRRYLPFLHYFDPDSQIEKDHDSKTGVDPSAQRLIRELYSLDLRNDPFSHNRLDGITFEHLE
 VSPDISSFITGTYSLACGRAQSRFADFPKDDFVLAKNRKEQLISVADGKECLIVSGLAFFICLFLDREQ
 ASGMLSRIRGPKRIDENWARAVHETFCDL CIRHPHDLRESSNIKEALLDMLNLRCPRI LYDMLPEEE
 RAQFLPALDENSMMNLSSENSLNEESRLLWDGSSDWAALTKRIRHQDRFPYLMRFIEEMDLLKGI RFRV
 DLGEIELDSYKVKVGRNGEYDRTITDHALAPGKLSDFQNEEVS RMI SGEASYPVRFSLFAPRYAIYDNK
 IGCHTSDPVYPKSKTGEKRALSNPQSMGFI SVHNLKLLMELCEGFSRMOQSDFLRKANRI LDETAE
 GKLQFSALFP EMRHRFI PPQNP KSKDRREKAEATLLEKYKQETKGRKDKLNSQLLSAFDMNQRLPSRLLD
 EWMNIRPASHSVKLRITYKQLNEDCRLRLRKRFRKDGDKARAIPLVGMATFLSQD IVRMI I SEETKLI
 TSAYYNEMQRS LAQYAGEENRRQFRAI VAEHLHLLDPSSGHPFLSATMETAHRYTEDFYKCYLEKKREWLA
 KTFYRPEQDENTKRRISVFPVFDGEARLLPLLI RRRMKEQNDLQDWIRNKQAHPIDLPSHLFDSKIMEL
 LKVKDGKKKWNEAFKDWSTKYPDGMQPFYGLRRELNIHGKSVSYIPSDGKFPADCYTHLMKTVQDKKR
 ELRTAGKPVPPDLAADIKRSFHRVNEREFMLRLVQEDDRMLMMAINKMTDREEDILPGLKNIDSILDE
 ENQFSLAVHAKVLEKEGEGGDNLSLVPATIEIKSKRKDWSKYIRYRDRRVPGLMSPHFHKA ILDEVK
 ILLGEYDRCRIKIFDWAFALEGAIMSDRDLKPYLHESSREGKSGEHS TLVKMLVEKKGCLTPDESQYLI
 LIRNKAHNQFPFAEAEMPLIYRDVSAKVGSI EGS SAKDLPEGSSLVDSLWKYEMIRKILPILDPENRF
 FGKLLNMSQPINDL

>gi|492916530|ref|WP_006044833.1| hypothetical protein [*Prevotella pallens*]
 MKEEEKGKTPVSTYNKDDKHFWAFLNLRHNVYITVNHINKILGEGEINRDGYENTLEKSWNEIKDIN
 KKDRLSKLI IKHFPFLEVI TYQRNSADTTKQKEEKQAEASLESLSKSFVFIYKLRDLRNHYSHYKHSK
 SLERPKFEEDLQEKMYNIFDASIQLVKEDYKHNTDIKTEEDFKHLDRKGQPKYSFADNEGNI TESGLLFF
 VSLFLEKKDAIWWQKLEGEFKCSNESYQKMTNEVFCRSRMLLPKLRLOSTQTDWILLDMLNELIRCPKS
 LYERLREEDRKKFRVPIEIADEDYDAEQEPPKNAALVRHQDRFPYFALRYFDYNEIFTNLRFQIDLGT
 YHFSIYKQIGDYKESHHLTHKLYGFERIQEFTKQNRPDWRKFKVIFNSPETSKEPYIPETTPHYHLENQKI
 GIRFRNDNDK I WPSLKINS EKNEKSKYKLDKSFQAEAPLSVHELPLMFFYLLKTEENTDNDNEIETKKK
 ENKNDKQEKHKIEEIEENKITEIYALYDAFANGKINSIDKLEEYCKGKDI EIGHLPKQMI AILKSEHKDM
 ATEARKQEEMLADVQKLESLDNLQIN EEEI ENVERKNS SLKSGEIASWLVNDMMRFQPVQKSDNEGNPLNN
 SKANST EYQMLQRS LALYNKEEKPTRYPRQVNLIESSNPHPLNNT EWEKCNNILS FYRSYLEAKKNPLE
 SLKPEDWEKNQYFLMLKEPKINCETLVQGWKNGFNLPRI FTEPIRKWFMEHRKNI TVAE LKRVGLVAKV
 I PLFPSEBYKDSVQPFYNYL FNVGNINKPDEKNFLNCEERRELLRKKKDEFKMTDKEEENPSYLEPQS
 WNKFERELRLVRNQD IVTWL CMELFNKKIKELNVEKIYLNKININTIKKEKNT E EKNGE EKI I KEKNN
 I LNRIMPMLRPIKVYGRNFSKNKKKIRRNIFFTVYIEBKGT KLLKQGNFKALERDRRLGGLP SFVKTH
 SKAESKSNITSKSRVEYELGEYQKARIEIKDMLAL EETLIDKYNSLDIDNFHNMLIGWLKLDKDEPKAS
 FQNDVDDL IAVRNAPSHNQYPMRNR IAFANINPFLSSANTSEEKGLGIANQLKDKIKHTIEKIEIEKPI
 IETKE

>gi|493302367|ref|WP_006259957.1| hypothetical protein [*Myroides odoratimimus*]
 MKDILTIDTTEKQNRFYSHKIIDKYFFGGYFNLASNNIYEVFEEVNKRNI FPKLAKRDNGNLKNYI IHVF
 KDELSISDFEKRVAIFASYFP ILETVDKKS I KERNRTIDLILSQRIRQFREMLISLVTAVDQLRNFYTHY
 HSDIV IENKVLDFLNS SFVSTALHVKDKYKLTDKI KEFLKETIAEELDILIEAYKKQKIEKKNTRFKAN
 KREDILNAIYNEAFWSFINDKDKDKKETVVAKGADAYFEKNHHSKNDPFDALNISEKGIYVLLSFFLIN
 KEMDSLKANLIGFGKGVDRRESGNSIKYMATQRIYSFHTYRGLKQKIRTS EEGVKEILLMQMIDELSKVPN
 VVYQHLSTIQQNSFI EDWNEYKDYEDDVETD DLSRVIHPVIRKRYEDRFNYFAIRFLDEFFDFPILRFQ
 VHLGDYVHDRIKQLGKVESDRI I KEKVI VPARLKD INSAKASYPHSLEEQDKBELDNKWLFPNPSYDF
 PKEHTLQHGEQKNAGKIGIYVVKLRDTQYKEKAAL EEARSLNPKERSATKASKYDII TQII EANDNVKS
 EKPLVFIGQPIAYLSMNDI HMLFSLLDIDNAELKKTPEEVEAKLIDQIGKQINEILSKDIDTKILKKYK
 ND LKETDIDK I TRDLARDKEIEKILBQKQRADDYNYISSTKFNIDKSRKRKHLFNAEKGKIGVWLAN

- continued

DIKRFMPKESKSKWKGQHTLQKLFAYFDTSKSDLELILSNMVMVKDYPIELIDLVKKSRIIVDFLNKY
LEARLEYIENVI TRVKNISIGTPQPKIVRKECFITLKKSNYTVVSLDKQVERILSMPLFIERGFMDKPKIM
LEBGKLYKQHKKEPFWVHYKENSYNQNFYDIEVYEITTEDKREKAKVKKIKQQQKNDVFLMMVNYML
EEVLKSSNDRSLNELYQTKKEERIVNKQVAKDTQERNKNIWNKVVDLQLCDGLVHIDNVKLLKDIGNFR
KYENDSRVKEPLTYQSDIVWSAYLSNEVDSNKLVIYIERQLDNYESIRSKELLKEVQEI ECSVYNGVANKE
SLKQSGNENFKQYVQLGLLPI GMDVREMLILSTDVKFKKEETIQLGQAGEVEQDLYSLIYIRNKFAHNQL
PIKEFFDFCENNYRSISDNEYAEBYMEIFRSIKEKYAN

>gi|494654219|ref|WP_007412163.1| hypothetical protein [*Prevotella* sp.
MSX73]

MQKQDKLQVDRKKNIAIFAPPKYITIMENQEKPEIYYELTDKHFWAFLNARHNVTYTTINHINRRLEIA
ELKDDGYMMGIKGSWNEQAKKLDKVRRLRDLIMKHFPPLEAAAYEITNSKSPNNKEQREKQSEALSNLN
LKNVLFIFLEKQLVLRNYSYHYKYSESPKPIFETSLKNMYKVFANVRLVKRDMYHHENIDMQRDFTH
LNRKKQVGRINKNIIDSPNFHYHAFADKEGNMTIAGLLFFVSLFLDKKDAIWMQKLLKGFDRGNLREQMIN
EVFCRSRISLPLKLELVQTKDWMQLDMLNLEVRCPKSLYERLREKDRRESFKVPPDIFSDDYADEEPEFK
NTLVRHQDRFPYFVLRFDLNEIFEQLRFQIDLGTYHFSYKRIIGDEDEVRHLTHHLYGFARIQDFAPQ
NQPEWRKLVKDLDFHETSQEPYISKTAPHYHLENEKIGIKFCSIHNNLPPSLKREKTCNGRSKFNLGIO
FTAEAFSLVHLELLNMFYLLLLIKDYSRKESADKVEGIRKEISNIYAIYDAFANNEINSIADLICRLQK
TNILQGHLPKQMISILEGRQKMEKEAERKIGEMIDDIQRRLDLLCKQINQKIRIGKRNAGLLKSGKIAD
WLVSDMMRFQPVQKDIINNAPINNSKANSTYRMLQHALALFGSESRKAYFRQMNLVGNANPHFLAET
QWEHQTNILSFYRNYLEARKKYLKGLKPKQNKQYQHFLILKVKQTKNRNLTVTGWKNSFNLPGRGIFQPIR
EWFEKHNSKRIYDQILSFDRVGVFAKAIPLYFAEYKDNVQPPYDYPFNI GNKLPKQKQFLDKKERVE
LWQKNKELFKNYPSKKNKTDLAYLDFLSWKKFERELRLIKNQDITVTLWMPKELFKITIVIEGLKIGEIHRL
DIDINTANEESNNILNRIMPMLKPVKIIYETDNKGNILKERPLATFYI EETETKVLKQGNFKVLAKDRRLN
GLLSFAETTDIDLEKNPITKLSVDYELIKYQTRISIFEMTLGLEKLLIDKYSTLPTDSFRNMLERWLQC
KANRPELKNYVNSLIAVRNAPSHNQYPMYDATLFAEVKKTFLPPSVDTYKKEIENIAPQLLEIVGKAIKEI
EKSENKN

>gi|501433495|ref|WP_012458151.1| hypothetical protein [*Porphyromonas*
gingivalis]

MNIVPASENKGQSRIVEDDPQYFGLYLNLAARENLI EYESHVRIKFGKKLNEESLKQSLLCDHLLSVDWR
TKVYGHSSRYLPLFLHYFPDSDSI EKDHDSKTGVDPSAQRLIRELYSLLDLFRNDFSHNRLDGI TFEHLE
VSPDISSTITGTYSLACGRAQSRFAVFPKPDVFLAKNRKEQLISVADGKECLIVSGFAFFICLFLDREQ
ASGMLSRIRGFKRIDENWARAVHETFCDLICIRHPDRLESNI KEALLDMLNLRNRCPRILYDMLPEEE
RAQFLPALDENSMNLS ENSLNDEESRLLWDS SDWAEALTKRIRHQDRFPYMLLRFIEEMDLLKGI RFRV
DLGETELDSYKVKVGRNGEYDRTITDHALAFGKLSDFQNEEEVSRMISGEASYPVRFSLFAPRYAIYDNK
IGYCHTSDPVVYKSKTGEKRALNPSQSMGFI SVHDLRKLMLMELCEGFSFRMQSDFLRKANRI LDETAE
GKLQFSALEFPBMRHRFIPQPNKSKDRREKAEATTLEKYQBEIKGRKDLKNSQLLSAPDMQLPSRLLD
EWMNIRPASHSVKLRITYKQLNEDCRLRLRKRFRKDGDKARAI PLVGMATFSLSDIVRMI ISEETKKLI
TSAYYNEMQRS LAQYAGEENRRQFRAI VAE LRLLDPSSGHPPLSATMETAHRYTEGFYKCYLEKKREWLA
KIFYPREQDENTKRRISVFFVFDGEARKLLPLLI RRRMKEQNDLQDWIRNKQAHPIDLPSHLFDKSVKML
LKVKGKKNWNEAFKDWSTKYPDGMQPFYGLRRELNIHGKSVSYPISDGGKPFADCYTHLMKTVRDKKR
ELRTAGKVPVDDLAADIKRSPHRAVNEREFMLRLVQEDDRMLMAINKMMDREEDILPGLKNIDSLDE
ENQFSLAVHAKVLEKEGEGDMSLSLVPATIEIKSKRKDWSKYIRYR YDRRVPLGMSHFPEHKA ILDEVK
ILLGEYDRCKRIKFDWAFALGAIMSDRDLKPYLHESRSREGKSGEHSTLVKMLVEKKGCLTPDESQYLI
LIRNKAHNQPPCAEMPLIYRDVSAKVGSI EGS SAKDLPEGSSLVDSLWKKYEMI IRKILPILDPENRF
FGKLLNMSQPINDL

>gi|501433759|ref|WP_012458414.1| hypothetical protein [*Porphyromonas*
gingivalis]

MTEQNERPYNGTYITLEDKHFWAFFNARHNAYITLAHIDRQLAYS KADITNDEDILFFKGQWKNLDND
LERKARLRSLILKHFSPLEGAAYGKKLFPESQSSGNKSSKKELIKKEEELQANALSLDNLKSI LFDPLQ
KLKDFRNYYSYHRHPESSELPLFDGNMLQRLYVNFVDSVQVRKDRHEHNDKVDPHRFNHLVRKGGKDRY
GNNDNPFKHHFVDREKVTAEAGLLFFVSLFLEKRDIAWMQKIRGFKGGTETYQMTNEVFCRSRISLP
KLKLESRLIDDWMLDMLNLRVCPKSLYDRLEEDRARFRVVPVILSDEDDIDGTEEDPPFKNI LVRHQD
RFPYFALRYFDLKKVFTSLRFHIDLGTYHFAIYKKNIGEOPEDRHLTRNLVGFGRIQDFAEEHRPEEWR
LVRDLDFETGDKPYITQTPHYHIEKGI GLRFVPEGQHLWSPPEVGATRTRGSKYAQDKRFTAEAFSL
VHELMPMMFYFLREKYSDEASAERVQGRIKRVIEDVYAVYDAFARGEINTRDELDACLADKGI RRGHL
PRQIGILSQEHKDMEEKIRKQLQEMIVDTHRLDMLDRQTDRIIRI GRKNAGLPKSGVIADWLVRDMMR
FQPVAKDTSGKPLNNSKANSTYRMLQRALALFGGKERLIPYFRQMNILGNNPHPLHETRWESHINI
LSFYRSYLKARKAFQSIGRSDRVENHRFLLKKEPKTDRQTLVAGWKGFHLPRGIFTEAVRDCLEMLG
DEVGSYKEVGMKAVPLYFERACKDRVQPPYDYPFNVGNLKPCKGRFLSKEKRAEEWESGKERFRLAK
LKEIL EAKEHPYLDKFSWQKFERELRLVKNQDIITWMI CRDLMEENKVEGLDIGTLYLKDIRTVDVQGG
NLNVLNRVKMRLPVVYRADSRGHVHKEQAPLATVYIEERDTKLLKQGNFKSVKDRRLNGLFSVFDTG
ALAMEQYPSIKLRVEYELAQVTARVCFAFEQTLLEESLLTRYPHLPDKNFRKMLLESWSDPLLDKWPDLH
GNVRLLI AVRNAPSHNQYPMYDEAVFSSIRKYDPSDPDAI EERMGLNIAHRLSEEVKQAKEMAERI IQA

>gi|503582079|ref|WP_013816155.1| hypothetical protein [*Porphyromonas*
gingivalis]

MTEQNEKPYNGTYITLEDKHFWAFFNARHNAYITLAHIDRQLAYS KADITNDEDILFFKGQWKNLDND
LERKARLRSLILKHFSPLEGAAYGKKLFPESQSSGNKSSKNKELIKKEEELQANALSLDNLKSI LFDPLQ
KLKDFRNYYSYHRHPESSELPLFDGNMLQRLYVNFVDSVQVRKDRHEHNDKVDPHRFNHLVRKGGKDRY
GNNDNPFKHHFVDREGIVTEAGLLFFVSLFLEKRDIAWMQKIRGFKGGTETYQMTNEVFCRSRISLP
KLKLESRLIDDWMLDMLNLRVCPKSLYDRLEEDRARFRVVPVILSDEDDIDGAEEDPPFKNI LVRHQD
RFPYFALRYFDLKKVFTSLRFQIDLGTYHFAIYKKNIGEOPEDRHLTRNLVGFGRIQDFAEEHRPEEWR
LVRDLDFETGDKPYITQTPHYHIEKGI GLRFVPEGQHLWSPPEVGATRTRGSKYAQDKRFTAEAFSL
AHELMPMMFYFLREKYSDEASAERVQGRIKRVIEDVYAVYDAFARDEINTRDELDACLADKGI RRGHL
PRQIGILSQEHKDMEEKIRKQLQEMMADTHRLDMLDRQTDRIIRI GRKNAGLPKSGVIADWLVRDMMR

- continued

FQPVAKDTSGKPLNNSKANSTEYRMLQRALALFGGKERLIPYFRQMNLIIGNNPHPLHETRWESHINI
LSFYRSYLKARKAPLQSIGRSDRVENHRFLLKEPKTRQTLVAGWKGFHLPRGIFFTEAVRDCLIMGL
DEVGYSYKVEGFMKAVPLYFERACKDWVQPFYNYPFNVGNSLKPCKGRFSLKKEKRAEWEWSGKERFRLAK
LKKEILEAKEHPYLDKFSWQKFERELRLVKNQDIITWIMICGDLMEENKVEGLDIGTLYLKDITRDVQEQG
SLNVLNRVKMRLPVVVYRADSRGHVHKEQAPLATVYI EERDTLLKQGNFKSVKDRRLNGLFSDVDTG
ALAMBQYPI SKLRVYELAKYQARVCAFEQTLLEESLLTRCPHLPDKNFRKMLESWSDPLLDKWPDLH
RKVRLLI AVRNAFNSHQYPMYDEAVFSSIRKYDPSFPDAI EERMGLNIAHRLSBEVKQAKETVERIIQA

>gi|503763195|ref|WP_013997271.1| hypothetical protein [Capnocytophaga
canimorsus]

MKNIQRLGKGNFSPFKKEDKFFYGGFLNLANNNIEDFFKEIITREGIVITDENKPKETPGEKILNEIF
KKDISIVDYEKWVNIADYFPFTKYLSTLYLEEMQFKNRVIFCRDVMKELKKTVEALRNFYTHDHEPIKI
EDRVFYFLDKVLLDVSILVKNKYLTKDKIKEFLNQHIGEELKELCKQRKDYLVGKGRIDKESEIINGIY
NNAFKDFIKCKREKQDDKENHNSVEKILCNKEPQNKKQSSATVWELCSKSSSKYTEKSPNRENDKHCLE
VPI SQKGI VFLSFLNKGGEIYALTSN IKGPKAKITKEEPVTDKNSIRYMATHRMFSPFLAYKGLKRKIR
TSEINYNEDGQASSTYEKETLMLQMLDELNKVPDVVYQNLSEVDVQKTFIEDWNEYLKENNGDVGTMEEBQ
VIHPVIRKRYEDKPNYFAIRFLDFAQFPILRFQVHLGNLYLCKRTRKQICDTTTEREVKKKITVFGRLSE
LENKKAIFLNEREEIKGWEVFPNPSYDFPKENISVNYKDFPIVGSILDRKQPVSNKIGIRVKI ADELQR
EIDKAIKEKLRNPKNRKANQDEKQKERLVNEIVSINSNEQGEPPVFIGQP TAYLSMNDIHSVLYEFLIN
KISGAELETKIVEKETIQIKQIIGKDATTKILKPYTNANSNSINREKLLRDLEQQEQLKILLEEQQQRE
KDKDKKS KRKHEL YPSEKGVAVWLANDIKREMPKAFKEQWRGYHHSLLQKYLAYEYSKEELKNLLPK
EVFKHFPFLKGYFQQYLNQFYIDYLRRLSYVNELLNLTQNFKNDKALKATEKECFPFKQNYIIN
PINIQIQSILVYPIFLKRGFLDEKIMIDREKFKENKDTLADWFMHYKNYKEDNYQKPYAYPLEKVEEK
EKFKRNKQINKQKNDVYLLMVEYIIQIKIPGDKFVEENPLVLKGIPOSKAERQNNTHAATQERNLNG
ILNQPKDIKI QGKI TVGVKLDKIGNFRKYETDQRVNTFLDYEPKKEWMAYLPNDWKEKEKQGGQLPPNNV
IDRQISKYETVRSKILLKDVQELEKII SDEIKEEHRHDLKQGYNYFKYIINGLLRQLKNENVENYKVF
KLNINPEKVNITQLKQEATDLKQKAFVLTYIRNKFAHNQLPKKEFPWYCYQEKYGI EKEKTYAEYFAEVF
KREKALIK

>gi|503931547|ref|WP_014165541.1| hypothetical protein [Flavobacterium
columnare]

MSSKNESYKQKTFPNHYKQEDKYFFGGFLNADDNLRQVGKFKTRINFNHNHNLASVFKDYFNKEKSV
AKREHALNLLSNYFPVLERIQKHINHNFBQTRIEFELLDITIKKLRDYTHHYHKPITINPKIYDFLDDT
LLDVLVITIKKKVKVNDTSRELLKELRPELTQLKNQKREELIKKGGKLLLEENLNAVFNHCLRPFLSEENK
TDDKQNKTVSLRKYRKS KPNETSITLIQSGLVFLMSFFLHRKEFPVPTSGLEGFKAVNTIKKEEISLN
KNNI VYMI THWSYSYNNFKGLKHRIKTDQGVSTLEQNNITHSLINTNIKEALLTQIVDYLSKVPNEIYET
LSEKQKFEFEDINEYMRNENEDSTFSSIVSHKVIKRYENKFNYPAMRFLDEYAEPLTRFPMVNFPGD
YIKDRQKILES IQPDSERIIKKEIHLFEKLSLVTEYKKNVYLKETSNIIDLRFPLFPNPSYVMANNIP
FYIDSRNNLDEYLNQKKAQSQKRNLIPEKYKQSKDAI IAMLQKTEIGVKDLQQRSTIGLSCNEL
PSMLYEVIVKDIKGAEL ENKI AQKIREQYQSIRDFLDSPQKDNIPPTLIKTIINTDSSVIFENQPIDIPR
LKNAIQKELTILQEKLLNVKEHEIEVDNYNRNKNTYKFNQPKNVDDKQLQRKYVYFNRNIRQEANWLA
SDLIHFPMKNSLWKGYMHNELQSFLAFEDKKNDCI ALLETVFNKEDCILT KGLKNLFLKHGNFDYK
EYLKLEDFLNTESTFLENGLIGLPPKILKKELSKRFKYPFIVFQKQRFIIKELEEKKNLYADAINLSR
GIFDEKPMI PFKPNPDEFASVYQYNNYQSFYELTDPDIVERDKKKYKNLRAINKVIQDYLLK
MVDTLYQDLFNQPLDKSLSDFVSKAEREKIKADAKAYQRNDSLWNVKVIHLSLQNNRI TANPKLDKIG
KYKRALQDEKIAIILTYDDRIWYALQKPEKENENDYKELHYTALNMELQYEVKRSKELKQVQEQLEKQ
ILEEYTFPLSTQIHPADFERENPNFKKYLASHILENEDDLKLPKVEAMRELEDETI TNPIKKAIVLI
IRNKMAHNQYPPKFIYDLANRFVPPKKEEYFATYFNRVPEITIKELWENKKEKDKTQV

>gi|504837663|ref|WP_015024765.1| hypothetical protein [Psychroflexus
torquais]

MESIIGLGLSFPNPKYADKHYFGSFLNVLNENLNAVFAEFKERISYKAKDENISSLIEKHFDNMSIVDY
EKKISILNGYLPIDFLDDELENNLNRVKNPKNFII LAEAI EKLDRDYTHFYHDPI TFEDNKEPLLEL
LDEVLLKTI LDVKKYKLTDKIKEILKDSLREEMDLVIRKTDREKKNINPKIQHTDSSQIKNSIFND
AFQGLLYEDKGMNKKTVSHRAKTRLNPKDIHKQEERDFEIPLSTSGLVFLMSLFLSKKEIEDFKSNIKG
FKGKVVKD ENHNSLYKMATHRVYSI LAFKGLKYRIKTDIFSKETLMMQMIDEL SKVPDCVYQNLSETKQK
DFIEDWNEYFKDNEENTENLENSRVVHPVIRKRYEDKFNYPFAIRFLDE FANFKILKQVFMGYIHDQRT
KTIGTTNITERTVKEKINVPGLSKMDNLKHPFSQLSDDENTDWEFFPNPSYNFLTQADNSPANNPI
YLELKNQOIIIEKDAIKAEVNTQNRNPNKPSKRDLLNKILKTYEDFHQGDPTAILSLNEIPALLHLPLV
KPNNTGQQIENIIRIKIEKQFKAINHPSKNNKGIKSLFADINVRVNAIKLKKDLAEALDMLNKKHIAF
KENQKASSNYDKLKEHQFTPKNKRPELRKYVYFYSKEGEEATWLANDIKREMPKDFIKKQKQHSBLO
RKLAFYDRHTKQDIKELLSGCEFDHSLDLINAYPQKDNFEDFFSKYLENRIETLGBVLKHLHDFNEPT
LKGVPKNCFKPLKRNQNYTESPEI I KKRILAKPTFLPRGVFDERPTMKKGNPLKDKNEFAEWFVEYLEN
KDYQKFNABEYRMRDADFKKNAVKKQKLDKDFYILQMVNYLLKEVFGKDEMNQLSELFQTRQERLKLQ
GIACKQMNKETGDSSENTRNQTYYIWNKDVVPSFFNGKVTIDKVKLKNIGKYRYERDERVVKIFIGYEVDE
KMMYLPHNWKDRYSVKPINVIDLQIQEYETRSHELLKEIQNLEQYIYDHTIDKNILLQDGNPNFKMYV
LNGLLIGIKQVNI PDIIVLQKQINFDKIDIFIGIASCSSELEKKIILIAIRNKFAHNQLPNKMIYDLANEF
LKI EKNETYANYLKVLLKMI SDLA

>gi|505158517|ref|WP_015345619.1| hypothetical protein [Riemerella
anatispestifer]

MEKPLLPNVYTLKHKFPWGAFLNIARHNAFITICHINEQLGLKTPSNDDKIVDVCETWNNILNNDHLL
KKSQLELILKHFPFLTAMCYHPPKKEGKKGHQKEQQKEKESAEQAEALNPSKLI EALEILVNQLHS
LRNYSYSHKHKPDAEKDIFNSLCILNNTDF

- continued

>gi|505158518|ref|WP_015345620.1| hypothetical protein [Riemerella anatipestifer]
 MFFSFHNAQRVIFKHLKYKAFDASLRMVKEDYKAHFIVNLI R D FAHLNRKGNKQDNPDFNRYRFEKDGFF
 TESGLLFFINLFLDKRDAYWMLKKVSGFKASHKQREKMITVFCRSRILLPKLRLESRYDHNQMLLDMLLS
 ELSRCPKLLYEKLS EENKHFQVEADGFLDEI EEBQNPFKDTLIRHQDRFPYFALRYLDLNLNESPFSIRFQ
 VDLGTYHYCIYDKKIGDEQEKRHLLTRILLSFGRLQDFTEINRPQEWKALTKDLDYKETSNOFFISKTIPIH
 YHITDNKIGFRLGTSKELGPELSEIKDGANRI AKYPYNSGFVAHAFISVHELLPLMFYQHLIGKSEDLLEKE
 TVRHITQRIYKDFEERINTIEDLEKANQGRPLGAFPKQMLGLLQNKQPDLSSEKAKIKIEKLI AETKLLS
 HRLNLIKLLKSSPKLGRKREKLIKTGVADLWLVKDFMRFPQVAYDAQNP I KSSKANSTEFWFIRRALALYG
 GEKNRLEGYFKQINLIGNINPHPLFNKFNWKA CRNLVDFYQYQLEQREKFL EAIKHQPWEYPYQCLLLKV
 PKENRKNLVKGNWQGGISLPRGLFTEAIRETL SKDLILSKPIRKEIKKHGRVGFISRAITLYFKKQYQDK
 HQSFYNLSYKLEAKAPLLKKEEYHYWQQNKQSPTESSQRLELHTSDRWKDYLLYKRWQHLKELRLRYRN
 QDIMLWMLTLELTKNHFKLENLNLYHQLENLAVNVQEADAKLNPLNQTLPMLVPVKYPTTAFGEVQYH
 ETPIRTVYIREQT KALMGNFKALVKDRRLNGLFSFIKEENDTQKHPISQLRLRRELEIYQSLRVDFAK
 EILSLEEKLLNKHASLSLENEFRILLBEWKKKYAASMV IDKHIAFIASVRNAPFNHNYPFYKETLHAP
 ILLFTVAQPI TEEDGLGIAEALLKVLREYCEIVKSQI

>gi|543953686|gb|ERJ64231.1| hypothetical protein HMPREF1555_01956
 [Porphyromonas gingivalis F0570]
 MRDLDFYFTGDKPYITQTPHYHIEKGGKIGLRFVPEGQHLWPSPEVGATRTGRSKYAQDKRLTAEAFLSV
 HELMPMMFYFLLREKYSDEASAERVQGRIRKVI EDVYAVYDAFARGEIDI LDRDLACLADKGI RRGHLP
 RQMI A ILSQEHKDMEEKVRKLLQEMIADTDHRLDMLDRQTD R K I R IGRKNAGLPKSGVIADWLVRDMMRF
 QPVAKDTSGKPLNNSKANST EYRMLQRALALFGGEEKERLI PYFRQMNLI GGNMHPHFLHETRWESHINIL
 SFYRSYLKARKAFLQSIGRSDRVENHRFLLLEKPKTD R QTLVAGWKS EBFHLPRGIFTEAVRDCL IEMGHD
 BVGSYKEVGFMAKAVPL YFERACKDRVQPFYDYPFNVGNLSKPKKGRFLSKEDRAEEWESGKERFRLAKL
 KKEILEAKEHPYLD FKSQKFERELRLVKNQDIITWMMCRDLMEENKVEGLDIGTLYLKDITDVTYEQGS
 LNVMLNRVKPMLRPVVYRADS RGHVHKEQAPLATVYIEERDTKLLKQGNFKSFVKDRRLNGLFSFVDTGG
 LAMEQYPI SKLRVEYELAKYQARVCAFEQTL ELESLLTRYPHLPDES FREMLESWSDPLLDKWPDLHR
 KVRLLI A VRNAPFNHNYQPMYDEAVFSSIRKYDPS SPDAIEERMGLNIAHRLSEEVKQAKEMAERI IQA

>gi|543955272|gb|ERJ65637.1| hypothetical protein HMPREF1553_02065
 [Porphyromonas gingivalis F0568]
 MNIVPASENKGQSRIVEDDPQYFGLYLNLARENLI E VESHVRIKFGKLLNEESLKQSLLCDHLLSVDRW
 TKVYGHSSRYL PFLHYFPDSDQIEKDHD SKTGVD PDS AQRLIRELYSLLDFLRNDFSHNRLDGI TFEHLE
 VSPDISSFITGTYSLACGRAQSRFADFPKDDFVLAKNRKEQLISVADGKECLIVSGLAFFICLFLDREQ
 ASGMLSRIRGPKRINDENARAVHETFCDL CIRHPHDLRLESSNI KEALLDMLNELNRCPRILYDMLPEEE
 RAQFLPALDENSMNLS ENSNLEESRLLWDGSDWAEALTKRIRHQDRFPYLM LRFIEEMDLLKGI RFRV
 DLGEI ELD SYSKVGRNGEYDRTITDHALAFGLSD FQNEEEVSRMISGEASYPVRFSLFAPRYAIYDNK
 IGYCHTSDPVYKSKTGEKRLSNRSMGFISVHDLRKL LLMELCEGFSFRMQSDFLRKANRI LDETAE
 GKLQFSALFP EMRHRFIPQNP KSKDRREKAETTL EKYQKEIKGRKDKLNSQLLSAFDMNQRQLPSRLLD
 EWMNIRPASHSVKLRTYVYKQLNEDCRLRLKFRKDG DGKARAIPLVGMATFLSQDIVRMI I SEETKKLI
 TSAYNEMQRS LAQYAGEENRQFR AIVAEHLRLDPS SGGHPFLSATMETAHRYTEDFYKCYLEKKREWLA
 KTFYRPEQDENTKRRISVFFVPDGEARKLLPLLI RRRMKEQNDLQDWIRNKQAHPIDLPSHLFDSKIMEL
 LKVKDGGKKNWNAEAFKDWSTKYPDGMQPFYGLRRELNIHGKSVSYPISDGGKFAFCYTHLMKKTVDKRR
 ELRTAGKVPVDDLAADIKRSPHRAVNEREFMLRLVQEDDRMLMAINKMMDREEDILPGLKNIDSLDE
 ENQFSLAVHAKVLEKEGEGDNSLSLV PATIEIKSKRDKWSKYIRYR YDRRVPGLMSHFPEHKA ILDEVK
 ILLGEYDRCRKI KFDWAFALGAINMSDRDLKPYLHESSSREGKSGEHSTLVKMLVEKKGCLTPDESQYLI
 LIRNKAHNQFP CAEMPLIYRDVSAKVGSI EGS SAKDLPEGSSLVDSLWKKYEMIIRKILPILDPENRF
 FGKLLNMSQPIINDL

>gi|543955721|gb|ERJ66054.1| hypothetical protein HMPREF1553_01900
 [Porphyromonas gingivalis F0568]
 MRDLDFCFTGDKPYITQTPHYHIEKGGKIGLRFVPEGQHLWPSPEVGATRTGRSKYAQDKRLTAEAFLSV
 HELMPMMFYFLLREKYSDEASAERVQGRIRKVI EDVYAVYDAFARGEIDI LDRDLACLADKGI RRGHLP
 RQMI A ILSQEHKDMEEKVRKLLQEMIADTDHRLDMLDRQTD R K I R IGRKNAGLPKSGVIADWLVRDMMRF
 QPVAKDTSGKPLNNSKANST EYRMLQRALALFGGEEKERLI PYFRQMNLI GGNMHPHFLHETRWESHINIL
 SFYRSYLKARKAFLQSIGRSDRVENHRFLLLEKPKTD R QTLVAGWKG EBFHLPRGIFTEAVRDCL IEMGLD
 BVGSYKEVGFMAKAVPL YFERACKDRVQPFYDYPFNVGNISKPKKGRFLSKEDRAEEWESGKERFRDLEA
 WSQSTARRIEDAFAGIKYASPGNKKKIBQLLQDLSLWETFESKLVKADKINLAKLKEILEAKEHPYLD
 FKSQKFERELRLVKNQDIITWMMCRDLMEENKVEGLDIGTLYLKDITDVTYEQGS LNVMLNRVKPMLRPV
 VYRADS RGHVHKEQAPLATVYIEERDTKLLKQGNFKSFVKDRRLNGLFSFVDTGGLAMEQYPI SKLRVE
 YELAKYQARVCAFEQTL ELESLLTRCPHLPDKNFRKMLLESWSDPLLDKWPDLHRKVRLLI A VRNAPFNH
 NYQPMYDEAVFSSIRKYDPS SPDAIEERMGLNIAHRLSEEVKQAKEMAERI IQA

>gi|543956136|gb|ERJ66432.1| hypothetical protein HMPREF1555_01119
 [Porphyromonas gingivalis F0570]
 MNIVPASENKGQSRIVEDDPQYFGLYLNLARENLI E VESHVRIKFGKLLNEESLKQSLLCDHLLSVDRW
 TKVYGHSSRYL PFLHYFPDSDQIEKDHD SKTGVD PDS AQRLIRELYSLLDFLRNDFSHNRLDGI TFEHLE
 VSPDISSFITGTYSLACGRAQSRFADFPKDDFVLAKNRKEQLISVADGKECLIVSGLAFFICLFLDREQ
 ASGMLSRIRGPKRINDENARAVHETFCDL CIRHPHDLRLESSNI KEALLDMLNELNRCPRILYDMLPEEE
 RAQFLPALDENSMNLS ENSNLEESRLLWDGSDWAEALTKRIRHQDRFPYLM LRFIEEMDLLKGI RFRV
 DLGEI ELD SYSKVGRNGEYDRTITDHALAFGLSD FQNEEEVSRMISGEASYPVRFSLFAPRYAIYDNK
 IGYCHTSDPVYKSKTGEKRLSNRSMGFISVHDLRKL LLMELCEGFSFRMQSGFLRKANRI LDETAE
 GKLQFSALFP EMRHRFIPQNP KSKDRREKAETTL EKYQKEIKGRKDKLNSQLLSAFDMNQRQLPSRLLD
 EWMNIRPASHSVKLRTYVYKQLNEDCRLRLKFRKDG DGKARAIPLVGMATFLSQDIVRMI I SEETKKLI
 TSAYNEMQRS LAQYAGEENRQFR AIVAEHLRLDPS SGGHPFLSATMETAHRYTEDFYKCYLEKKREWLA
 KTFYRPEQDENTKRRISVFFVPDGEARKLLPLLI RRRMKEQNDLQDWIRNKQAHPIDLPSHLFDSKIMEL

- continued

LKVKDGKKKWNEAFKDWWS TKYPDGMQPFYGLRRELNIHGKSVSYIPSDGKKFADCYTHLMEKTVQDKKR
ELRTAGKVPVPPDLAADI KRSPHRAVNEREFMLRLVQEDDRMLMLAINKMMTDREEDILPGLKNIDS ILDK
ENQFSLAVHAKVLEKEGEGGDNLSLVPATIEIKSKRKDWSKYIRYRDRRVPGLMSHFPEHKA ILDEVK
ILLGEYDRCRIKI PDWAFALGAIMSDRDLKPYLHSSSRGKSGEHS TLVKMLVEKKGC LTPDESQYLI
LIRNKAHNQFP CAEMPLIYRDVSAKVGSI EGS SAKDLPEGSSLVDSLWKYEMI IRKILPILDHENRF
FGKLLNMSQPINDL

>gi|543977697|gb|ERJ81654.1| hypothetical protein HMPREF1989_02374
[Porphyromonas gingivalis F0566]

MRDLDFETGDKPYITQTTPHYHIEKGI GLRFVPEGQHLWPSPEVGATRTGRSKYAQDKRLTAEAFLSV
HELMPPMFYPLLRKEYSEESAERVQGRIRKVI EDVYAVYDAPARGEIDI LDRDLACLADKGI RRGHLR
RQMI A ILSQEHKDMEEKIRKKLQEMMADTDHRLDMLDRQTRKIRIGRKNAGLPKSGVIADWLVRDMMRF
QPVAKDTSGKPLNNSKANSTEYRMLQALALFEGEKERLI PYFRQMNLIGGNNPHPLHETRWESHINIL
SFYRSYLKARKAFLQSIGRSDRVENHRFLLLKEPKTDRQTLVAGWKGFEHLPRGIFTEAVRDCLEMGHD
EVASYKEVGFMAKAVPL YFERACKDRVQPFYDYPFNVGNSLKP KKGRFLSK EKRAE EWESGKERFRDLEA
WSQS AVRR I EDAPAGIENASRDNKKKI EQLLQDLSLWKAFESKLVKRGDKINI AKLKKETLEAKEHPYHD
FKSWQKFERELRLVKNQDI ITWMMCRDLMEENKVEGLD IGTLYLKD IRTDVHEQGS LNVLNVRKPMRLPV
VVYRADSRGHVHKEQAPLATVYI EERDTKLLKQGNFKS FVKDRRLNGLFSFVDTGALAMEQYPI SKLRVE
YELAKYQ TARVCAFEQTLEESL LTRCPHLPDKNFRKMLSESWS DPLLDKWPDLHRKVRLLI IAVRNAFSH
NQYPMYDEAVFSSIRKYDPSFDAIEERMGLNIAHRLSEEVKQAKETVERI IQA

>gi|543983828|gb|ERJ87335.1| hypothetical protein HMPREF1990_01800
[Porphyromonas gingivalis W4087]

MNIVPASENKQSRIVEDDPQYGLYLNLARENLI EESHVRI KFGKKKLNESLQSSLCDHLLSVDRW
TKVYGHRRYLPFLHYFPDPSQIEKDHDSKTGVDPDS AQRLIRELYSLLDPLRDNFSHNRLDGI TFEHLE
VSPDIS SFI TGTYSLACGRAQSRFADF FPKPDDFVLAKNRKQQLISVADGKECLIVSGLAFFI CLFLDREQ
ASGMLSRIRGFKRIDENWARAVHETFCDL CIRHPDRLESSNI KEALLDMLNELNRCPRILYDMLPEEE
RAQFLPALDENSMNLS ENSLNEESRLLWDGSSDWA EALTKRIRHQDRFPYMLLRFI EEMDLLKGI RFRV
DLGETELDSYKVKVGRNGEYDRTITDHALAFGKLSDFQNEEEVSRMISGEASYPVRFSLFAPIYAIYDNK
IGYCHTSDPVYKSKTGEKRALSNRSMGFI SVHDLRKLMLMELCEGFSFRMQSDFLRKANRILDETAE
GKLQFSALFPFEMRHRFIPFQNP KSKDRREKAETTLEKYQKEIKGRKDKLNSQLLSAFDMDQRLPSRLD
EWMNIRPASHSVKLRTYVKQLNEDCRLRLQKFRKDGDKARAI PLVGMATFLSQDIVRMI ISEETKKLI
TSAYNEMQRLAQA YAGEENRHQFRAI VAE LRLLDPSSGHPFLSATMETAHRYTEDFYKCYLEKKRWLA
KTFYRPEQDENTKRRISVFFV PDGEARKLLPLLI RRRMKEQNDLQDWIRNKQAHPIDLPSHLFD SKVMEL
LKVKDGKKKWNEAFKDWWS TKYPDGMQPFYGLRRELNIHGKSVSYIPSDGKKFADCYTHLMEKTVRDKKR
ELRTAGKVPVPPDLAADI KRSPHRAVNEREFMLRLVQEDDRMLMLAINKIMTDREEDILPGLKNIDS ILDK
ENQFSLAVHAKVLEKEGEGGDNLSLVPATIEIKSKRKDWSKYIRYRDRRVPGLMSHFPEHKA ILDEVK
ILLGEYDRCRIKI PDWAFALGAIMSDRDLKPYLHSSSRGKSGEHS TLVKMLVEKKGC LTPDESQYLI
LIRNKAHNQFP CAEPI IYRDVSAKVGSI EGS SAKDLPEGSSLVDSLWKYEMI IRKILPILDPENRF
FGKLLNMSQPINDL

>gi|543999227|gb|ERJ98772.1| hypothetical protein HMPREF1218_0639
[Prevotella pleuritidis F0068]

MENDKRLEESACYTLNDKHFWAAPLNLARHNVIYITVNHINKTLELKNKKNQEI IIDNDQDILAI KTHWAK
VNGDLNKTDRLELMI KHFPPLAAAIYSNNKEDKEEVKEEKQAKAQS FKS LKDC LFLFLFKLQEARNYYS
HYKYSSESSKEPEFEGLLEKMYNTFDASIRLVKEDYQYKNDIDPEKDFKHLERKEDFNLYLPTDKDNKGI
TKNGLLFFVSLFLEKDAIWMQKFRGFKDNRGNKEKMTHEVFCRSRMLLPKIRLESTQTDWI LLDMLN
ELIRCPKSLYERLQAGYREKFKVPFDS IDEDYDABEQEPFRNTRLVRHQDRFPYFALRYFDYNEIFKNLRFQ
IDLGTYPHSIYKLLIGGKKEDRHLTHKLYGFERIQEFTKQNRPKWQAI IKDLDIYETSNERYISETTPH
YHLENQKIGIRFRNDNDI IWPSLKNINGEKNEKSNYLDKPYQAEAFLSVHELLPMMFYLLKMENTDND
KEDNEVGI KKKGNKNNKQEKHKI EEI I ENKI KDI YALYDAPINGEINS IDELAEQREGKDI EIGHLPKQL
IVILKNKSKDMAEKANRQKEMI KD IKKRLALD KQVKGEIEDGGRNIRLLKSGEIRARWLVNDMMRFQPV
QKDNBEGKPLNNSKANST EYQMLQRS LALYNKEEKPTRYFRQVNLISKSNPHFLBEDTKWEECYNLSFYR
NYLKAKIKFLNKLKPEDWKNQYFLMLKEPKTNRKTLVQGWKNGFNLP RGI FTPEI KEWFKRHQNDSEY
KKVEALDRVGLVAKVIPLFFKEEYFKEDAQKEINNCVQPFYSPYVGNIIHKPEEKNFLHCEERRKLDWK
KKDKFKGYKAKEKSKMMDKEKEHRSYLEFQSWNKFERELRLVRNQDILTWLLC IKLIDKLIKIDELNIE
ELQKLRKLDIDITDAKKEKNIILNRVMPMLRPVIVYEIDKSFNIVKDKPLHIVYI EETGIKLLKQGNFKA
LVKDRRLNGLFSFVKISSEAESKSKPI SKLRVEYELGAYQKARIDII KDMLALEKTLIDNDENLPINKFS
DMLKSWLKGKGEANKARLQNDVGLLVAVRNAFSHNQYPMYNSVEVFKGMKLLSLSSDIPEKEGLGIKQK
DKIKETIERIIEIEKEIRN

>gi|545305353|ref|WP_021584635.1| hypothetical protein [Prevotella
pleuritidis]

MENDKRLEESACYTLNDKHFWAAPLNLARHNVIYITVNHINKTLELKNKKNQEI IIDNDQDILAI KTHWAK
VNGDLNKTDRLELMI KHFPPLAAAIYSNNKEDKEEVKEEKQAKAQS FKS LKDC LFLFLFKLQEARNYYS
HYKYSSESSKEPEFEGLLEKMYNTFDASIRLVKEDYQYKNDIDPEKDFKHLERKEDFNLYLPTDKDNKGI
TKNGLLFFVSLFLEKDAIWMQKFRGFKDNRGNKEKMTHEVFCRSRMLLPKIRLESTQTDWI LLDMLN
ELIRCPKSLYERLQAGYREKFKVPFDS IDEDYDABEQEPFRNTRLVRHQDRFPYFALRYFDYNEIFKNLRFQ
IDLGTYPHSIYKLLIGGKKEDRHLTHKLYGFERIQEFTKQNRPKWQAI IKDLDIYETSNERYISETTPH
YHLENQKIGIRFRNDNDI IWPSLKNINGEKNEKSNYLDKPYQAEAFLSVHELLPMMFYLLKMENTDND
KEDNEVGI KKKGNKNNKQEKHKI EEI I ENKI KDI YALYDAPINGEINS IDELAEQREGKDI EIGHLPKQL
IVILKNKSKDMAEKANRQKEMI KD IKKRLALD KQVKGEIEDGGRNIRLLKSGEIRARWLVNDMMRFQPV
QKDNBEGKPLNNSKANST EYQMLQRS LALYNKEEKPTRYFRQVNLISKSNPHFLBEDTKWEECYNLSFYR
NYLKAKIKFLNKLKPEDWKNQYFLMLKEPKTNRKTLVQGWKNGFNLP RGI FTPEI KEWFKRHQNDSEY
KKVEALDRVGLVAKVIPLFFKEEYFKEDAQKEINNCVQPFYSPYVGNIIHKPEEKNFLHCEERRKLDWK

- continued

KKDKFKGYKAKEKSKKMTDKEKEEHRSYLEFPQSWNKFERELRLVRNQDILTWLLCIKLIDKLIKIDELNIE
ELQKLRLLKIDIDITAKKEKNNILNRVMPMLRPVIVYEIDKSFNIVKDKPLHIVYIETETGIKLLKQGNFKA
LVKDRRLNGLFSFVKISSEAESKSKPI SKLRVEYELGAYQKARIDIIKDMLALEKTLIDNDENLPINKFS
DMLKSWLKGKGEANKARLQNDVGLLVAVRNAPSHNQYPMYNSVEVFKGMKLLSLSSDIPEKEGLGIKAKQLK
DKIKETIERIEIEKEIRN

>gi|545424591|ref|WP_021663197.1| hypothetical protein [*Porphyromonas gingivalis*]
MNIIVPASENKGQSRIVEDDPQYFGLYLNLAARENLEIEVESHVRIKFGKKKLNESLQKSQLCDHLLSVDRW
TKVYGHSTRYLPFLHYFPDPSQIEKDHDSTKGTVDPSAQRLIRELYSLDDFLRNDFSHNRDGI TFEHLE
VSPDISSFITGTYSLACGRAQSRFADFPKDDFVLAKNRKEQLISVADGKECLIVSGLAFFICLFLDREQ
ASGMLSRIRGFKRIDENWARAVHETFCDL CIRHPHDLRELSNI KEALLDMLNLRNRCPRILYDMLPEEE
RAQFLPALDENSMNLS ENSLNEESRLLWDGSSDWA EALTKRIHQDRFPYLM LRFIEEMDLLKGI RFRV
DLGEI ELD SYSKVGRNGEYDRTITDHALAFGKLSDFQNEEEVSRMISGEASYPVRFSLFAPRYAI YDNK
IGYCHTSDPVYPKSKTGKRALSNPRSMGFI SVHDLRKL LLMELLCEGFSFRMQSDFLRKANRI LDETAE
GKLQFSALFP EMRHRFIPPQNPKSKDRREKAETTLEKYKQETIKGRKDKLNSQLLSAFDMDQRLPSRLLD
EWMNIRPASHSVKLRITYVKQLNEDCRLRLQKFRKDGDKARAI PLVGEMATFLSQDIVRMI ISEETKKLI
TSAYNEMQRSLAQYAGEENRHQFRAI VAE LRLLDPSSGHPFLSATMETAHRYTEDFYKCYLEKKREWLA
KTFYRPEQDENTKRRISVFFVDPGEARKLLPLLI RRRMKEQNDLQDWIRNKQAHPIDLPSHLFDSKIMEL
LKVKDGGKKWNEAFKDWWS TKYPDGMQPFYGLRRELNIHGKSVSYIPSDGKKFADCYTHLMEKTVQDKKR
ELRTAGKVPVDPDLAADI KRSPHRAVNEREFMLRLVQEDDRMLMAINKMMDREEDILPGLKNIDSI LDE
ENQFSLAVHAKVLEKEGEGDNSLSLVPATIEIKSKRDKWSKYIRYRDRRVPGLMSHFPEHKA ILDEVK
ILLGEYDRCRIKIFDWAFALEGAIMSDRDLKPYLHES SRGKSGEHSTLVKMLVEKKGCLTPDESQYLI
LIRNKAHNQFP CAEMPLIYRDVSAKVGSI EGSSAKDLPEGSSLVDSLWKYEMI IRKILPILDENRF
FGKLLNNMSQPINDL

>gi|545427058|ref|WP_021665475.1| hypothetical protein [*Porphyromonas gingivalis*]
MNIIVPASENKGQSRIVEDDPQYFGLYLNLAARENLEIEVESHVRIKFGKKKLNESLQKSQLCDHLLSVDRW
TKVYGHSTRYLPFLHYFPDPSQIEKDHDSTKGTVDPSAQRLIRELYSLDDFLRNDFSHNRDGI TFEHLE
VSPDISSFITGTYSLACGRAQSRFADFPKDDFVLAKNRKEQLISVADGKECLIVSGLAFFICLFLDREQ
ASGMLSRIRGFKRIDENWARAVHETFCDL CIRHPHDLRELSNI KEALLDMLNLRNRCPRILYDMLPEEE
RAQFLPALDENSMNLS ENSLNEESRLLWDGSSDWA EALTKRIHQDRFPYLM LRFIEEMDLLKGI RFRV
DLGEI ELD SYSKVGRNGEYDRTITDHALAFGKLSDFQNEEEVSRMISGEASYPVRFSLFAPRYAI YDNK
IGYCHTSDPVYPKSKTGKRALSNPRSMGFI SVHDLRKL LLMELLCEGFSFRMQSGFLRKANRI LDETAE
GKLQFSALFP EMRHRFIPPQNPKSKDRREKAETTLEKYKQETIKGRKDKLNSQLLSAFDMDQRLPSRLLD
EWMNIRPASHSVKLRITYVKQLNEDCRLRLRQKFRKDGDKARAI PLVGEMATFLSQDIVRMI ISEETKKLI
TSAYNEMQRSLAQYAGEENRRQFRAI VAE LHL LDPSSGHPFLSATMETAHRYTEDFYKCYLEKKREWLA
KTFYRPEQDENTKRRISVFFVDPGEARKLLPLLI RRRMKEQNDLQDWIRNKQAHPIDLPSHLFDSKIMEL
LKVKDGGKKWNEAFKDWWS TKYPDGMQPFYGLRRELNIHGKSVSYIPSDGKKFADCYTHLMEKTVQDKKR
ELRTAGKVPVDPDLAADI KRSPHRAVNEREFMLRLVQEDDRMLMAINKMMDREEDILPGLKNIDSI LDK
ENQFSLAVHAKVLEKEGEGDNSLSLVPATIEIKSKRDKWSKYIRYRDRRVPGLMSHFPEHKA ILDEVK
ILLGEYDRCRIKIFDWAFALEGAIMSDRDLKPYLHES SRGKSGEHSTLVKMLVEKKGCLTPDESQYLI
LIRNKAHNQFP CAEMPLIYRDVSAKVGSI EGSSAKDLPEGSSLVDSLWKYEMI IRKILPILDHENRF
FGKLLNNMSQPINDL

>gi|545440038|ref|WP_021677657.1| hypothetical protein [*Porphyromonas gingivalis*]
MNIIVPASENKGQSRIVEDDPQYFGLYLNLAARENLEIEVESHVRIKFGKKKLNESLQKSQLCDHLLSVDRW
TKVYGHSTRYLPFLHYFPDPSQIEKDHDSTKGTVDPSAQRLIRELYSLDDFLRNDFSHNRDGI TFEHLE
VSPDISSFITGTYSLACGRAQSRFADFPKDDFVLAKNRKEQLISVADGKECLIVSGLAFFICLFLDREQ
ASGMLSRIRGFKRIDENWARAVHETFCDL CIRHPHDLRELSNI KEALLDMLNLRNRCPRILYDMLPEEE
RAQFLPALDENSMNLS ENSLNEESRLLWDGSSDWA EALTKRIHQDRFPYLM LRFIEEMDLLKGI RFRV
DLGEI ELD SYSKVGRNGEYDRTITDHALAFGKLSDFQNEEEVSRMISGEASYPVRFSLFAPRYAI YDNK
IGYCHTSDPVYPKSKTGKRALSNPRSMGFI SVHDLRKL LLMELLCEGFSFRMQSGFLRKANRI LDETAE
GKLQFSALFP EMRHRFIPPQNPKSKDRREKAETTLEKYKQETIKGRKDKLNSQLLSAFDMDQRLPSRLLD
EWMNIRPASHSVKLRITYVKQLNEDCRLRLRQKFRKDGDKARAI PLVGEMATFLSQDIVRMI ISEETKKLI
TSAYNEMQRSLAQYAGEENRRQFRAI VAE LHL LDPSSGHPFLSATMETAHRYTEDFYKCYLEKKREWLA
KTFYRPEQDENTKRRISVFFVDPGEARKLLPLLI RRRMKEQNDLQDWIRNKQAHPIDLPSHLFDSKIMEL
LKVKDGGKKWNEAFKDWWS TKYPDGMQPFYGLRRELNIHGKSVSYIPSDGKKFADCYTHLMEKTVQDKKR
ELRTAGKVPVDPDLAADI KRSPHRAVNEREFMLRLVQEDDRMLMAINKMMDREEDILPGLKNIDSI LDE
ENQFSLAVHAKVLEKEGEGDNSLSLVPATIEIKSKRDKWSKYIRYRDRRVPGLMSHFPEHKA ILDEVK
ILLGEYDRCRIKIFDWAFALEGAIMSDRDLKPYLHES SRGKSGEHSTLVKMLVEKKGCLTPDESQYLI
LIRNKAHNQFP CAEMPLIYRDVSAKVGSI EGSSAKDLPEGSSLVDSLWKYEMI IRKILPILDHENRF
FGKLLNNMSQPINDL

>gi|545442698|ref|WP_021680012.1| hypothetical protein [*Porphyromonas gingivalis*]
MNIIVPASENKGQSRIVEDDPQYFGLYLNLAARENLEIEVESHVRIKFGKKKLNESLQKSQLCDHLLSVDRW
TKVYGHSTRYLPFLHYFPDPSQIEKDHDSTKGTVDPSAQRLIRELYSLDDFLRNDFSHNRDGI TFEHLE
VSPDISSFITGTYSLACGRAQSRFADFPKDDFVLAKNRKEQLISVADGKECLIVSGLAFFICLFLDREQ
ASGMLSRIRGFKRIDENWARAVHETFCDL CIRHPHDLRELSNI KEALLDMLNLRNRCPRILYDMLPEEE
RAQFLPALDENSMNLS ENSLNEESRLLWDGSSDWA EALTKRIHQDRFPYLM LRFIEEMDLLKGI RFRV
DLGEI ELD SYSKVGRNGEYDRTITDHALAFGKLSDFQNEEEVSRMISGEASYPVRFSLFAPRYAI YDNK
IGYCHTSDPVYPKSKTGKRALSNPRSMGFI SVHDLRKL LLMELLCEGFSFRMQSDFLRKANRI LDETAE
GKLQFSALFP EMRHRFIPPQNPKSKDRREKAETTLEKYKQETIKGRKDKLNSQLLSAFDMDQRLPSRLLD
EWMNIRPASHSVKLRITYVKQLNEDCRLRLQKFRKDGDKARAI PLVGEMATFLSQDIVRMI ISEETKKLI
TSAYNEMQRSLAQYAGEENRHQFRAI VAE LRLLDPSSGHPFLSATMETAHRYTEDFYKCYLEKKREWLA

- continued

KTFYRPEQDENTKRRISVFFVDPGEARKLLPLLI RRRMKEQNDLQDWIRNKQAHPIDLPSHLFDSKVMEL
 LKVKDGKKKWNEAFKDWWS TKYPDGMQPFYGLRRELNIHGKSVSYIPSDGKKFADCYTHLMEKTVRDKKR
 ELRTAGKVPVDPDLAAYIKRSPHRAVNEREFMLRLVQEDDRMLMLAINKIMTDREEDILPGLKNIDSLDK
 ENQFSLAVHAKVLEKEGEGGDNLSLVPATIEIKSKRKDWSKYIRYRDRRVPGLMSHFPEHKA ILDEVK
 ILLGEYDRCRIKIFDWAFALEGAIMSDRDLKPYLHESSSREGKSGEHSTLVKMLVEKKGCLTPDESQYLI
 LIRNKAHNQFPCAAEMPIIYRDVSAKVGSI EGSSAKDLPEGSSLVDSLWKYEMIIRKILPILDPENRF
 FGKLLNNMSQPINDL

>gi|563396328|gb|ETA27404.1| hypothetical protein SJDPG2_03560

[*Porphyromonas gingivalis* SJD2]

MNIVPASENKGQSRIVEDDPQYFGLYLNLAARENLEEVESHVRIKFGKKKLNEESLKQSLLCDHLLSVDRW
 TKVYGHSSRRYLPFLHYFPDPSQIEKDHDSKTGVDPSAQRLIRELYSLDPLRNDFSHNRLDGI TFEHLE
 VSPDISSFITGTYSLACGRAQSRFADFPPKDDFVLAKNRKEQLISVADGKECLIVSGLAFFICLFLDREQ
 ASGMLSRIRGPKRIDENWARAVHETFCDL CIRHPHDLRESNI KEALLDMLNLRNRCPRILYDMLPEEE
 RAQFLPALDENSMNLS ENSLNEESRLLDWGS SDWAEALT KRIRHQDRFPYMLLRFIEEMDLLKGI RFRV
 DLGEIELDSYSKVKGRNGEYDRTITDHALAPGKLSDFQNEEEVSRMISGEASYPVRFSFLFAPRYAIYDNK
 IGYCHTSDPVYPKSKTGEKRLSNRSMGFISVHDLRKLMLMELCEGFSRMOQSDFLRKNRILDETAE
 GKLFQFSALEPPEMRHRFIPQNPKSKDRREKAETTLEKYKQEI KGRKDKLNSQLLSAFDMNQRLPSRLLD
 EWMNIRPASHSVKLRITYVKQLNEDCRLRLRKRFRKDGDKARAIPLVGMATFSLQDIVRMI I SEETKKLI
 TSAYNEMQRLAQAAGEENRRQFRAI VAEHLHLLDPSSGHPFLSATMETAHRYTEDFYKCYLEKKREWLA
 KTFYRPEQDENTKRRISVFFVDPGEARKLLPLLI RRRMKEQNDLQDWIRNKQAHPIDLPSHLFDSKIMEL
 LKVKDGKKKWNEAFKDWWS TKYPDGMQPFYGLRRELNIHGKSVSYIPSDGKKFADCYTHLMEKTVQDKKR
 ELRTAGKVPVDPDLAADIKRSFHRAVNEREFMLRLVQEDDRMLMLAINKMMDREEDILPGLKNIDSLDE
 ENQFSLAVHAKVLEKEGEGGDNLSLVPATIEIKSKRKDWSKYIRYRDRRVPGLMSHFPEHKA ILDEVK
 ILLGEYDRCRIKIFDWAFALEGAIMSDRDLKPYLHESSSREGKSGEHSTLVKMLVEKKGCLTPDESQYLI
 LIRNKAHNQFPCAAEMPIIYRDVSAKVGSI EGSSAKDLPEGSSLVDSLWKYEMIIRKILPILDPENRF
 FGKLLNNMSQPINDL

>gi|564601847|ref|WP_023846767.1| hypothetical protein [*Porphyromonas*

gingivalis]

MNIVPASENKGQSRIVEDDPQYFGLYLNLAARENLEEVESHVRIKFGKKKLNEESLKQSLLCDHLLSVDRW
 TKVYGHSSRRYLPFLHYFPDPSQIEKDHDSKTGVDPSAQRLIRELYSLDPLRNDFSHNRLDGI TFEHLE
 VSPDISSFITGTYSLACGRAQSRFADFPPKDDFVLAKNRKEQLISVADGKECLIVSGLAFFICLFLDREQ
 ASGMLSRIRGPKRIDENWARAVHETFCDL CIRHPHDLRESNI KEALLDMLNLRNRCPRILYDMLPEEE
 RAQFLPALDENSMNLS ENSLNEESRLLDWGS SDWAEALT KRIRHQDRFPYMLLRFIEEMDLLKGI RFRV
 DLGEIELDSYSKVKGRNGEYDRTITDHALAPGKLSDFQNEEEVSRMISGEASYPVRFSFLFAPRYAIYDNK
 IGYCHTSDPVYPKSKTGEKRLSNRSMGFISVHDLRKLMLMELCEGFSRMOQSDFLRKNRILDETAE
 GKLFQFSALEPPEMRHRFIPQNPKSKDRREKAETTLEKYKQEI KGRKDKLNSQLLSAFDMNQRLPSRLLD
 EWMNIRPASHSVKLRITYVKQLNEDCRLRLRKRFRKDGDKARAIPLVGMATFSLQDIVRMI I SEETKKLI
 TSAYNEMQRLAQAAGEENRRQFRAI VAEHLHLLDPSSGHPFLSATMETAHRYTEDFYKCYLEKKREWLA
 KTFYRPEQDENTKRRISVFFVDPGEARKLLPLLI RRRMKEQNDLQDWIRNKQAHPIDLPSHLFDSKIMEL
 LKVKDGKKKWNEAFKDWWS TKYPDGMQPFYGLRRELNIHGKSVSYIPSDGKKFADCYTHLMEKTVQDKKR
 ELRTAGKVPVDPDLAADIKRSFHRAVNEREFMLRLVQEDDRMLMLAINKMMDREEDILPGLKNIDSLDE
 ENQFSLAVHAKVLEKEGEGGDNLSLVPATIEIKSKRKDWSKYIRYRDRRVPGLMSHFPEHKA ILDEVK
 ILLGEYDRCRIKIFDWAFALEGAIMSDRDLKPYLHESSSREGKSGEHSTLVKMLVEKKGCLTPDESQYLI
 LIRNKAHNQFPCAAEMPIIYRDVSAKVGSI EGSSAKDLPEGSSLVDSLWKYEMIIRKILPILDPENRF
 FGKLLNNMSQPINDL

>gi|567226320|dbj|GAE20957.1| hypothetical protein JCM10003_349

[*Bacteroides pyogenes* JCM 10003]

MESIKNQKSI GKILQKDPYFGLYLNLMALLNVRKVENHIRKWLGDVALLPEKSGFHSLLTIDNLS SAKW
 TRFYKSRKFLPFLPEMFDSDKKS YENRRETAECALDIDRQKISSLLKEVYGLQDIRNAPSHYHIDDQSV
 KHTALI ISSEMRFRICRLKDKIKGEPKSGKLNEMPKLYWDKEFPNGMQRFYRCKRREVEVPDKVVEY
 I ILLGEYDRCRIKIFDWAFALEGAIMSDRDLKPYLHESSSREGKSGEHSTLVKMLVEKKGCLTPDESQYLI
 LIRNKAHNQFPCAAEMPIIYRDVSAKVGSI EGSSAKDLPEGSSLVDSLWKYEMIIRKILPILDPENRF
 FGKLLNNMSQPINDL

>gi|640570032|ref|WP_025000926.1| hypothetical protein [*Prevotella*

aurantiaea]

MEDDKKTIGSISYELKDKHFAAFLNLRARHNVIYITINHINKLLEIREIDNDEKVLDIKTLWQKGNKDLNQ
 KARLRELMTKHFPFLTAIYIKNKEDKKEVKQEKQAEASLES LKDCFLPLDKLQEAARNYSHYKSEF
 SKEPEFEGLLEKMYNIPGNNIQLVINDYQHNDINPDEDPKHLDRKQGFYKSFADNEGNI TESGLLFFV
 SLFLEKKAIDAIWQQKLNFGKDNLENKKMTHEVFCRSRILMPKLRLESTQTDWILLDMLNELIRCPKSL
 YERLQDDREKFKVPFPDAEDYNAEQEPPKNTLIRHQDRFPYFVLR YFDYNEIFKNLRFPQIDLGTYHFS
 IYKLIIGGQKEDRHLTHKLYGFERIQEFAKQNRPDWEKAI VKDLDIYETS NKRYISETTPHYHLENQKIG
 IRFRNGNKETWPSLKINDENNEKSKYKLDKQYQAEAFSLVHHELLPMMFYLLLLKKEKPNNDENASIVEG
 FIKREIRNI FKL YDAFANGENINDDLEKYCADKGI PKRHLPKQMVAILYDEHKDMVKEAKRKQKEMVKG
 TKKLLATLEKQTKQEKEDDGRNVKLLKSGEIA RLVNDDMMRFPQVQKDNQEGKPLNNSKANSTEYQMLQRS

- continued

LALYNNEEKPTRYFRQVNLIESNNPHPLKWI KWEECCNLI LIFYYSYLIKKIEPLNKLKPEDWKKQNYFL
 KLKEPKTNRETLVQGWKNGFNLPRI GTFEPIREWFKRHQNNSEKEYEALDRVGLVTKV I PLFFKEEYF
 KDKENEFKEDTQKEINDCVQPFYFNFPYVGN I HKPKEKDFLHREERI ELWDKKKDKFKGYKEKI KSKKLT
 EKDKEEFRSYLEFQSWNKFERELRLVRNQD I VTWLLCKELIDKLIKIDELNI EELKKLRNNIDIDTAKKE
 KNNILNRRVMPMELPVIVYEIDDSHKIVKDKPLHTIYI KEAETKLLKQGNFKALVKDRRLNGLFSFVKINS
 EAESKRNP I SKLRVYELGEYQEAR I E I IQDMLALEEKLIN KYKDL P INKFS EMLNSWLEGKDEADKARF
 QNDVDFLI AVR NAF SHNQYPMHNKI EFANI KPFS LYTANNSEKGLGIANQLKDKI KETTTDKIKKI EKPI
 ETKE

>gi|672596632|gb|AIJ35608.1| hypothetical protein EG14_06045

[*Porphyromonas gingivalis*]

MNIVPASENKQSRIVEDDPQYFGLYLNLARENLI EVEHSVRI KFGKKKLNNEESLKQSLCDHLLSVDRW
 TKVYGHRRYLPFLHYFPDPSQI EKDHDSKTGVDPDSAQRLIREL YSLLDFLRNDFSHNRLDGI TFEHLE
 VSPDISSPITGTYSLACGRAQSRFAVFPKDDFVLAKNRKEQLISVADGKECLIVSGFAFFI CLFLDREQ
 ASGMLSRIRGFKRIDENWARAVHETFCDL CIRHPDRLES SNI KEALLLDMLNELNRCPR I LYDMLPEEE
 RAQFLPALDENSMNLS ENSLDEESRLLWDGSSDWA EALTKRIRHQDRFPYLM LRF I EEMDLLKGRFRV
 DLGEI ELD SYSKVGRNGEYDRT I TDHALAPGKLSDFQNEEEVSRMI SGEASYPVRFSLFAPRYAI YDNK
 IGYCHTSDPVYKSKTGEKRALSNPQSMGFI SVHDLRKL LLMELLCEGFSFRMQSDFLRKANRI LDETAE
 GKLQFSALFPEMRHRFIPQNPKSKDRREKAETTL EKYKQET I KGRKDKLNSQLLSAFDMDQRQLPSRLLD
 EWMNIRPASHVVKLRITYYKQLNEDCRLRLR KFRKDGDKARAI PLVGEMATFLS QD I VRMI I SEETKKLI
 TSAYINEMQORSLAQYAGEENRRQFRAI VAE LRLLDPSGGHPLSATMETAHRYTEGFKYCYL EKKREWLA
 KIFYPPEQDENTKRRISVFFVDPGEARKLLPL I RRRMKEQNDLQDWIRNKQAHPI DLPSHLFDSKVMEL
 LKVDGKKKWNEAFKDWWS TKYPDGMQPFYGLRRELN I HGKSVSY I PSDGKKFADCYTHLMEKTVRDKKR
 ELRTAGKVPVPLAADI KRSPHRAVNEREFMLRLVQEDDRMLMA INKMMTDREEDI LPGLKNI DSILDE
 ENQFSLAVHAKVLEKEGEGDNL S LVPATIE I KSKRDKWSKYIRYRDRRVPGLMSHFPEHKA I LDEVK
 ILLGEYDRCRI KIPDWAFALEGAIMSDRDLKPYLHES SRBGKSGEHSTLVKMLVEKKGC LTPDESQYL I
 LIRNKAAHNQFP CAEMPLIYRDVSAKVGSI EGGSAKDLPEGSSLVDSLWKYEMI IRKI LPILDPENRF
 FGKLLNNMSQPINDL

>gi|672597416|gb|AIJ36392.1| hypothetical protein EG14_10345

[*Porphyromonas gingivalis*]

MTEQNERPYNGTYTLEDKHPWAAFNLARHNAYITLAIH DRQLAYS KADI TNDED I LFFKQGWKNDND
 LERKARLSL I LKHFSLFLEGAAYGK KLFESQSSGNKSSK KKEKEL I KKEKEELQANALSLDNLKSI LFDPLQ
 KLKDFRNYSHYRHPSESELPLFDGNMLQRLYNVFDVSVQRVKRDHEHNDKVDPHRHFNHLVRKGGKDRY
 GNDNPFKHHFVDREGIVTEAGLLFFVSLFLEKRD A IWMQKIRGFKGGTEAYQMTNEVFCR.SRI SLP
 KLLKLESLRIDDWMLLDMLNELVRCPKSLYDRLEEDRARFRVPID I LSEDDIDGTEEDPFKNI LVHRQD
 RFPYFALRYFDLKKVFTSLRPHIDLGT YHFAI YKKNIGEQPEDRHLTRNL YGFGRI QDFAEEHREPEWKR
 LVRLDLYPETGDKPIYITPTPHYI EKGI GLRFVPEGQHLWSPVEVGATRTRGRSKYAQDKRLTAEAFLS
 VHELMPMMFYFLLREKYS EEEVSAEKVQGR I KRVI EDVYAVYDAFARGE I D I LDRDLACLADKGI RRGHL
 PRQMI A ILSQEHKDMEEKVRKQLQEMI ADTDHRLDMLDRQTDRKI RIRGNKAGLPKSGVI ADWLVRDMMR
 PQPVAKDTSGKPLNNSKANS TEYRMLQRALALFGGEKERLI PYFRQMN I GGNPHPLHETRWESHINI
 LSFYRSYLKARKAFI QSI GRSDREENHRFL LKPEKPTDRQTLVAGWKSEPHLPRGI FTEAVRDCLI EMGY
 DEVGYSKEVGFMAKAVPLYFERACKDRVQPFYDYPFN VGNLSLKP KGRFLS KEKRAEEWESGKERFRLAK
 LKKE I LEAKEHPYLDPKSWQKFERELRLVKNQDI I TWMMCRDLMEENKVEGLDIGTLYLKD I RTDVHEQG
 SLNVLNRYKPMRLPVVVYRADSRGHVHKEQAPLATVYI EERDTKLLKQGNFKS FVKDRRLNGLFSFVDTG
 ALAMEQYPI SKLRVYELAKYQTA R VCAFEQTLLEESLLTRYPHLPDKNFRKMLESWS DPLDKWPDHL
 RKVRLLI AVR NAF SHNQYPMYDEAVFSSIRKYDPS SDAI EERMGLNIAHRLS EEVKQAKEMAERI IQV

>gi|690716823|gb|KGE88582.1| hypothetical protein IX84_07840

[*Phaeodactylibacter xiamenensis*]

MINTPKRRI LHRHPSYFAGFLNIARHNAFMIMEHLS TKYDMEKNI LDEAQLPNAKLFGLCKKRYGKPDV
 TEGVSRDLRRYFPFLNYPFLHLEKQNAEQAAATYD INPEDI EFLKGFRLLNQMRNNYSHYI SNIDYG
 KFDKLPVQDI YEAA I PRLLDGKHTKRPDV FESKHIRHLESNNSEYRPRSLANS PDHENTVAFVTCLEFL
 RKYAFPFLSR LDCFRS INDAABGDPL I RKA SHECYTMFCRRLPQPKLES S D I LDMVNLGRCP SALYNL
 LSEEDQARFHI KREETI GFEEEDPEELEQE I VLKRHSRFPYFALRYFDDTEAFQTLRFVDVYLGRWRIK P
 VYKKRI YGQERDRVLIQSI RTPI RLSRLLPI YENVKHDAVRQNEEDGKLVNPDV I SQPHKSWIQ I TESDDR
 AFLSDR I EHFSPHYNFGDQVI GLKF INPDRYAAI QNVFPKLPGEKKDKDAKLVNETADAI I STHEIRSL
 FLYHYLSKKPI SAGDERR I QVDTETFI KQYIDT I KLF FED I KSGELOPI ADPPNYQKNEPLPYVRGDKE
 KTQEBRAQYRE RQKEI KERRKELNI LLQNR YGLS IQYI PSRLREYLLGYKVPYEKLALQKLRQRKEVK
 KRIKDI EKMRTPRVGEQATWLAEDI VFLIPPKMHI PERKTI KHPQKLNNDQFRINQSSLAYFSVNKKA I K
 KFFQKETGIGLSNRETSHPFLYRIDVGRCRGILDFY IGYLYKMKDWLDDAI KVDNRKHGKKEAKKYEYK
 LPSS IQHKTPLBLDY I RLPVYLPRGLFKKAI VKALAAHADPQVEPEEDNVI FCDLQDLDDGDTQDFYNWQR
 YRSALTEKEIDNQLVLAHPYAEQ I LGTI K TLEGKQKNNKLGNAKQKI KDELIDLKRAKRRLLDREQYL
 RAVQAEDRALWMI QERQKQAEHEE I AFDQLDLKNI TKI LTESIDARLRI PDTKVDI TDKLP LRRYGD L
 RRVAKDRRLVNLAS YHYVAGLSEIPYDLVKKLE EYDRRRVAF FEHVYQFEKEVYDRYAAELRNENPKGE
 STYF SHWEYVAVAKHSADTHFNELFK EKVMLRNKPHHNEFPYFDWLLPEVEKASAAALYADRVFDVAEG
 YYQMRKMLRQ

>gi|696209666|gb|KGL48767.1| hypothetical protein HQ49_06245

[*Porphyromonas gulae*]

MTEQSERPYNGTYTLEDKHPWAAFNLARHNAYITLTHI DRQLAYS KADI TNQDVL SFKALWKNLDND
 LERKSRLRSL I LKHFSLFLEGAAYGK KLFESQSSGNKSSK KKEKEL I KKEKEELQANALSLDNLKSI LFDPLQ
 KLKDFRNYSHYRHSSESELPLFDGNMLQRLYNVFDVSVQRVKRDHEHNDKVDPHRHFNHLVRKGGKDRY
 GNDNPSFKHHFVDSGEMITTEAGLLFFVSLFLEKRD A IWMQKIRGFKGGTEAYQMTNEVFCR.SRI SLP
 KLLKLESLRMDWMLLDMLNELVRCPKPLYDRLEEDRACFRVPVD I LPEDDIDGGEDPFKNI LVHRQD
 RFPYFALRYFDLKKVFTSLRPHIDLGT YHFAI YKKMIGEQPEDRHLTRNL YGFGRI QDFAEEHREPEWKR
 LVRLDLYPETGDKPIYISQTPHYI EKGI GLRFMPEGQHLWSPVEVGI TRTRGRSKYAQDKRLTAEAFLS
 VHELMPMMFYFLLREKYS EEEVSAEKVQGR I KRVI EDVYAI YDAFARDEINTLKELDACLADKGI RRGHL

- continued

PKQMIATLSQEHKNMEEKVRKKLQEMIADTDHRLDMLDRQTDKIRIRGRKNAGLPKSGVIADWLVRDMMR
 FQPVAKDASGKPLNNSKANSTFYRMLQRALALFGGEEKERLIPYFRQMNLIIGGNNPHPLHETRWESHINI
 LSFYRSYLARKAFLEIRGRSDRMENRPFLLKKEPKTDRQTLVAGWKEFHLPRGI FTEAVRDCLIEMGY
 DEVGSYREVGFMKAVPLYFERACEDRVQPFYDS PFNVGNSLKPCKGRFSLKEERAEWEWGERKFRDLE
 AWSHSAARRIEDAFAGIEYASPGNKKKIEQLLRDLSLWEAFESKLKVRADKINLAKLKEILEAQEHPYH
 DFKSWQKFERELRLVKNQDII TWMMCRDLMEENKVEGLDIGTLYLKD IRPNVQEQGSLNVLNHRVKPMRLP
 VVVYRADS RGHVHKEEAPLATVYIEERDTKLLKQGNFKSFVKDRRLNGLFSFVDTGGLAMEQYPI SKLRV
 EYELAKYQATARVCVPELTLRLEESLLTRYPHLPDESFRKMLESWSDPLLAKWPELHGKVRLLIAVRNAPS
 HNQYPMYDEAVFSSIRKYDPSPDIAEERMGLNIAHRLSEEVKQAKETVERIIQA

>gi|696216474|gb|KGL55352.1| hypothetical protein HQ50_05870
 [Porphyromonas sp. COT-052 OH4946]

MTEQSERPYNGTYTLEDKHFWAFLNLARHNAYITLTHIDRQLAYSKADI TNDQDVLSPFKALWKNFND
 LERKSRLRSLILKHFSPLEGAAYGKCLFESKSGNKSSKNKELIKKEKEELQANALSLDNLKSILDFDLQ
 KLKDFRNYYSYHRHSESELELDFDGNMLQRLYNVFDVSVQVRKRDHEHNDKVDPHRHFNHLVRKGGKDRY
 GHNDNPSFKHHFVDSGEMVTEAGLLFFVSLFLEKRDIAWMQKIRGFKGGTETYQOMTNEVFCRSRISLP
 KLKLESRLIDDWMLLDMLNELVRCPKPLYDRLEDDRACFRVPVD ILPDEDDIDGGGEDPPFKNILVRHQD
 RFPYFALRYFDLKKVFTSLRPHIDLGTYHFAIYKMMIGEOPEDRHLTRNLYGFGRIQDFAEEHREPEWKR
 LVRLDLYFETGDKPYISQTPHYHI EKGKILGRFVPEGQHLWSPPEVGI TRTGRSKYAQDKRLTAEAFSL
 VHELMPMMFYFLLREKYS EEVSAEKVQGRIKRVI EDVYAIYDAPARDEINTLKELDACLADKGI RRGHL
 PKQMIATLSQERKDMEEKVRKKLQEMIADTDHRLDMLDRQTDKIRIRGRKNAGLPKSGVIADWLVRDMMR
 FQPVAKDTSGKPLNNSKANSTFYRMLQRALALFGGEEKERLIPYFRQMNLIIGGNNPHPLHETRWESHINI
 LSFYRSYLARKAFLEIRGRSDRVENCPPFLKKEPKTDRQTLVAGWKEFHLPRGI FTEAVRDCLIEMGY
 DEVGSYREVGFMKAVPLYFERACEDRVQPFYDS PFNVGNSLKPCKGRFSLKEDRAEWEWGERKFRDLE
 AWSHSAARRIKDAFAGIEYASPGNKKKIEQLLRDLSLWEAFESKLKVRADKINLAKLKEILEAQEHPYH
 DFKSWQKFERELRLVKNQDII TWMMCRDLMEENKVEGLDIGTLYLKD IRPNVQEQGSLNVLNHRVKPMRLP
 VVVYRADS RGHVHKEEAPLATVYIEERDTKLLKQGNFKSFVKDRRLNGLFSFVDTGGLAMEQYPI SKLRV
 EYELAKYQATARVCVPELTLRLEESLLSRYPHLPDESFRMLESWSDPLLAKWPELHGKVRLLIAVRNAPS
 HNQYPMYDEAVFSSIRKYDPSPDIAEERMGLNIAHRLSEEVKQAKETVERIIQA

>gi|700219215|gb|KGN69120.1| hypothetical protein HR09_05855
 [Porphyromonas gulae]

MTEQSERPYNGTYTLEDKHFWAFLNLARHNAYITLTHIDRQLAYSKADI TNDQDVLSPFKALWKNFND
 LERKSRLRSLILKHFSPLEGAAYGKCLFESKSGNKSSKNKELIKKEKEELQANALSLDNLKSILDFDLQ
 KLKDFRNYYSYHRHSESELELDFDGNMLQRLYNVFDVSVQVRKRDHEHNDKVDPHRHFNHLVRKGGKDRY
 GHNDNPSFKHHFVDSGEMVTEAGLLFFVSLFLEKRDIAWMQKIRGFKGGTETYQOMTNEVFCRSRISLP
 KLKLESRLIDDWMLLDMLNELVRCPKPLYDRLEDDRACFRVPVD ILPDEDDIDGGGEDPPFKNILVRHQD
 RFPYFALRYFDLKKVFTSLRPHIDLGTYHFAIYKMMIGEOPEDRHLTRNLYGFGRIQDFAEEHREPEWKR
 LVRLDLYFETGDKPYISQTPHYHI EKGKILGRFVPEGQHLWSPPEVGI TRTGRSKYAQDKRLTAEAFSL
 VHELMPMMFYFLLREKYS EEVSAEKVQGRIKRVI EDVYAIYDAPARDEINTLKELDACLADKGI RRGHL
 PKQMIATLSQEHKDMEEKIRKLLQEMIADTDHRLDMLDRQTDKIRIRGRKNAGLPKSGVIADWLVRDMMR
 FQPVAKDTSGKPLNNSKANSTFYRMLQRALALFGGEEKERLIPYFRQMNLIIGGNNPHPLHETRWESHINI
 LSFYRSYLARKAFLEIRGRSDRMENRPFLLKKEPKTDRQTLVAGWKEFHLPRGI FTEAVRDCLIEMGY
 DEVGSYREVGFMKAVPLYFERACEDRVQPFYDS PFNVGNSLKPCKGRFSLKEERAEWEWGERKFRDLE
 AWSHSAARRIEDAFAGIEYASPGNKKKIEQLLRDLSLWEAFESKLKVRADKINLAKLKEILEAQEHPYH
 DFKSWQKFERELRLVKNQDII TWMMCRDLMEENKVEGLDIGTLYLKD IRPNVQEQGSLNVLNHRVKPMRLP
 VVVYRADS RGHVHKEEAPLATVYIEERDTKLLKQGNFKSFVKDRRLNGLFSFVDTGGLAMEQYPI SKLRV
 EYELAKYQATARVCVPELTLRLEESLLTRYPHLPDESFRKMLESWSDPLLAKWPELHGKVRLLIAVRNAPS
 HNQYPMYDEAVFSSIRKYDPSPDIAEERMGLNIAHRLSEEVKQAKETVERIIQA

>gi|700224863|gb|KGN74647.1| hypothetical protein HR17_04485
 [Porphyromonas gulae]

MTEQSERPYNGTYTLEDKHFWAFLNLARHNAYITLTHIDRQLAYSKADI TNDQDVLSPFKALWKNFND
 LERKSRLRSLILKHFSPLEGAAYGKCLFESKSGNKSSKNKELIKKEKEELQANALSLDNLKSILDFDLQ
 KLKDFRNYYSYHRHSGSSELELDFDGNMLQRLYNVFDVSVQVRKIDHEHNDEVDPHYHFNHLVRKGGKDRY
 GHNDNPSFKHHFVDSGEMVTEAGLLFFVSLFLEKRDIAWMQKIRGFKGGTETYQOMTNEVFCRSRISLP
 KLKLESRLMDDWMLLDMLNELVRCPKPLYDRLEDDRACFRVPVD ILPDEDDIDGGGEDPPFKNILVRHQD
 RFPYFALRYFDLKKVFTSLRPHIDLGTYHFAIYKMMIGEOPEDRHLTRNLYGFGRIQDFAEEHREPEWKR
 LVRLDLYFETGDKPYISQTPHYHI EKGKILGRFVPEGQHLWSPPEVGI TRTGRSKYAQDKRLTAEAFSL
 VHELMPMMFYFLLREKYS EEVSAERVQGRIKRVI EDVYAVYDAPARDEINTRELDACLADKGI RRGHL
 PRQMIATLSQEHKDMEEKIRKLLQEMMADTDHRLDMLDRQTDKIRIRGRKNAGLPKSGVIADWLVRDMMR
 FQPVAKDASGKPLNNSKANSTFYRMLQRALALFGGEEKERLIPYFRQMNLIIGGNNPHPLHETRWESHINI
 LSFYRSYLARKAFLEIRGRSDRVENRPFLLKKEPKTDRQTLVAGWKEFHLPRGI FTEAVRDCLIEMGH
 DEVASYKEVGFMKAVPLYFERACEDRVQPFYDS PFNVGNSLKPCKGRFSLKEERAEWEWGERKFRDLE
 AWSYSAARRIEDAFAGIEYASPGNKKKIEQLLRDLSLWEAFESKLKVRADKINLAKLKEILEAQEHPYH
 DFKSWQKFERELRLVKNQDII TWMMCRDLMEENKVEGLDIGTLYLKD IRPNVQEQGSLNVLNHRVKPMRLP
 VVVYRADS RGHVHKEEAPLATVYIEERDTKLLKQGNFKSFVKDRRLNGLFSFVDTGGLAMEQYPI SKLRV
 EYELAKYQATARVCVPELTLRLEESLLTRYPHLPDESFRKMLESWSDPLLAKWPELHGKVRLLIAVRNAPS
 HNQYPMYDEAVFSSIRKYDPSPDIAEERMGLNIAHRLSEEVKQAKETVERIIQA

>gi|700226261|gb|KGN76020.1| hypothetical protein HQ40_04325
 [Porphyromonas gulae]

MNIVPATENKQSRIVEDDPQYFGLYLNLARENLEEVESHVRIKFGKKNLNEESLKQSLLCDHLLSIDRW
 TKVYGHRRYLPFLHCFDPDSGIEKDHDSKTGVDPDSAQRLIRELYSLDPLRNDPFSHNRLDGTI FFEHLK
 VSPDISSTITGAYTFACERAQSRFADFFKPDFFLAKNRKQQLISVADGKECLIVSGFAFFICLFLDREQ
 ASGMLSRIRGFKRIDENWARAVHETFCDLICIRHPHDLRESNIEKALLDMLNELNRCRPIIYDMLPEE
 RAQFLPALDENSMNLS ENLSMESRLLWDGSDWAEALTKRIRHQDRFPYLMFLRPIIEMDNLKGI RFRV
 DLGEIELDSYKVKVGRNGBYDRTITDHALAFGLKSDFQNEEVSRMISGEASYVPRFSLFAPRYAIYDNK

- continued

I GYCHTSDPVYPKSKTGEKRALSNPQSMGFI SVHDLRKL LLMELLCEGFSFRMQSDFLRKANRI LDETAE
 GKLOFSALFPEMRHRFI PPQNPKSKDRREKAETTLEKYKQETIKGRKDKLNSQLLSAFDMNQRLPSRLLD
 EWMNIRPASHSVKLRITYVKQLNEDCRLRLRKRFRKDGDKARAI PLVGMATFLSQDIVRMI I SEETKKLI
 TSAYYNEMQRLSQAQYAGEENRRQFRAI VAEHLHLLDPSSGHPFLSATMETAHRYTEDFYKCYLEKKREWLA
 KTFYRPEQDENTKRRISVFFVDPDGEARLKLPLLI RRRMKEQNDLQDWIRNKQAHPIDLPSHLFDSKIMEL
 LKVKDGGKKNWNAFKDWWSTKYPDGMQPFYGLRRELNIHGKSVSYIPSDGKPFADCYTHLMKKTVRDKKR
 ELRTAGKVPVDPDLAAYIKRSPHRAVNEREFMLRLVQEDDRMLMAINKMMDREEDILPGLKNIDSI LDE
 ENQFSLAVHAKVLEKEGGGNSLSLVPATIEIKSKRKDWSKYIRYRDRRVGMLSHFPEHKA ILDEVK
 ILLGGEYDRCKIKIFWAFALGAIMSDRDLKPYLHESSSREGKSGEHSTLVKMLVEKKGCLTPDESQYLI
 LIRNKAHNQFPFAAEMPLIYRDVSAKVGSI EGS SAKDLPEGSSLVDSLWKYEMI IRKILPILDHENRF
 FGKLLNNMSQPINDL

>gi|700235983|gb|KGN85385.1| hypothetical protein HR15_09830

[*Porphyromonas gulae*]

MTEQSERPYNGTYYTLEDKHFWAFLNLARHNAYITLTHIDRQLAYSKADI TNDQDVLSPFKALWKNFDND
 LERKSRLRSLILKHFSPLEGAAYGKKLPESSKSGNKSSKNKELIKKEKEELQANALSLDNLKSI LFDPLQ
 KLDKFRNYYSHYRHSSELPDFGNMLQRLYNVFDVSVQVRKRDHEHNDKVDPHYHFNHLVRKGGKDRY
 GHNDNPSFKHHPVDSGEMVTEAGLLFFVSLFLEKRD AIWMQKKIRGFKGGI GPYEQMTNEVFCRSRISLP
 KLKLESLRIDDWMLDMLNELVRCPKPLYDRLEKDRACFRVPVD ILPDEDDIDGGGEDPFPKNI LVRHQD
 RFPYFALRYFDLKKVFTSLRFHIDLGTYHFAIYKMKIGEOPEDRHLTRNLVGFGRIQDFAEEHRPEEWR
 LVRDLDYFETGDKPYISQTPHYHI EKGKILGRFVPEGQHLWSPPEVGI TRTGRSKYAQDKRLTAEAFSL
 VHELMPMMFYFLLREKYSEEVSAEKVQGR I KRVIKDVYAIYDAFARDEINTLKELDACSADKGI RRGHL
 PKMIGILSQEHKNMEEKVRKLEMIADTDHRLDMLDRQTDKIRI GRKNAGLPKSGVIADWLVRDMMR
 FQPVAKDTSGKPLNNSKANST EYRMLQRALALFGGKERLIPYFRQMNLI GGNPHPLD ETRWESHINI
 LSFYRSYLARKAFLEIRGSDRVENRPFLLKKEPKNDRQTLVAGWKESEFHLPRGIFTEAVRDCL IEMGY
 DEVGSYKEVGFMAKAVPLYFERACKDRVQPPYDS PFNVGNLSLKPCKGRFLS KEKRAEEWESGKERFRLAK
 LKKEILEAQEHPYHDFKSWQKFERELRLVKNQDI ITWMMCRDLMEENKVEGLDIGTLYLKD IRTDVHEQG
 SLNVLNVRKPMRLPVVYRADSRGHVHKEAPLATVYI EERDKLLKQGNFKS FVKDRRLNGLF SFVDTG
 GLAMEQYPI SKLRVEYELAKYQARVCAFEQTLLEESLLTRYPHLPDENPREMLESWSDP LLAQWPDH
 GKVRLLI AVRNAPSHNQYPMYDEAVFSSIRKYDPS SDAI EERMGLNIAHRLSEEVKQAKETVERI IQA

>gi|700238026|gb|KGN87312.1| hypothetical protein HQ46_09365

[*Porphyromonas gulae*]

MTEQSERPYNGTYYTLEDKHFWAFLNLARHNAYITLTHIDRQLAYSKADI TNDQDVLSPFKALWKNLDND
 LERKSRLRSLILKHFSPLEGAAYGKKLPESSKSGNKSSKNKELIKKEKEELQANALSLDNLKSI LFDPLQ
 KLDKFRNYYSHYRHSSELPDFGNMLQRLYNVFDVSVQVRKRDHEHNDKVDPHYHFNHLVRKGGKDRY
 GHNDNPSFKHHPVDSGEMVTEAGLLFFVSLFLEKRD AIWMQKKIRGFKGGI GPYEQMTNEVFCRSRISLP
 KLKLESLRIDDWMLDMLNELVRCPKPLYDRLEKDRACFRVPVD ILPDEDDIDGGGEDPFPKNI LVRHQD
 RFPYFALRYFDLKKVFTSLRFHIDLGTYHFAIYKMKIGEOPEDRHLTRNLVGFGRIQDFAEEHRPEEWR
 LVRDLDYLETGDKPYISQTPHYHI EKGKILGRFVPEGQHLWSPPEVGI TRTGRSKYAQDKRLTAEAFSL
 VHELMPMMFYFLLREKYSEEVSAEKVQGR I KRVIEDVYAIYDAFARDEINTLKELDTCLADKGRIRGHL
 PKMIGILSQEHKNMEEKVRKLEMIADTDHRLDMLDRQTDKIRI GRKNAGLPKSGVIADWLVRDMMR
 FQPVAKDTSGKPLNNSKANST EYRMLQRALALFGGKERLIPYFRQMNLI GGNPHPLD ETRWESHINI
 LSFYRSYLARKAFLEIRGSDRVENRPFLLKKEPKNDRQTLVAGWKESEFHLPRGIFTEAVRDCL IEMGY
 DEVGSYREVGFMKAVPLYFERACKDRVQPPYDS PFNVGNLSLKPCKGRFLS KEDRAEEWESGKERFRDLE
 AWSSAARRIKDAFAGI EYASPGNKKIEQLRDL SLWEAFESKLKVRADKINLAKLKEILEAQEHPYH
 DFKSWQKFERELRLVKNQDI ITWMMCRDLMEENKVEGLDIGTLYLKD IRTDVHEQGSLNVLNVRKPMRLP
 VVVYRADSRGHVHKEAPLATVYI EERNIKLLKQGNFKS FVKDRRLNGLF SFVDTGGLAMEQYPI SKLRVE
 YELAKYQARVCFELTLRLEESLLSRYPHLPDES PREMLESWSDP LLAQWPELHGKVRLLI AVRNAPSH
 NQYPMYDEAVFSSIRKYDPS SDAI EERMGLNIAHRLSEEVKQAKETVERI IQA

>gi|700239027|gb|KGN88274.1| hypothetical protein HR08_00310

[*Porphyromonas gulae*]

MTEQSERPYNGTYYTLEDKHFWAFLNLARHNAYITLTHIDRQLAYSKADI TNDQDVLSPFKALWKNLDND
 LERKSRLRSLILKHFSPLEGAAYGKKLPESSKSGNKSSKNKELIKKEKEELQANALSLDNLKSI LFDPLQ
 KLDKFRNYYSHYRHSSELPDFGNMLQRLYNVFDVSVQVRKRDHEHNDKVDPHRFNHLVRKGGKDRY
 GHNDNPSFKHHPVDSGEMVTEAGLLFFVSLFLEKRD AIWMQKKIRGFKGGT ETYQMTNEVFCRSRISLP
 KLKLESLRIDDWMLDMLNELVRCPKPLYDRLEKDRARFRVPVD ILPDEDDIDGGGEDPFPKNI LVRHQD
 RFPYFALRYFDLKKVFTSLRFHIDLGTYHFAIYKMKIGEOPEDRHLTRNLVGFGRIQDFAEEHRPEEWR
 LVRDLDYFETGDKPYISQTPHYHI EKGKILGRFVPEGQHLWSPPEVGI TRTGRSKYAQDKRLTAEAFSL
 VHELMPMMFYFLLREKYSEEVSAEKVQGR I KRVIEDVYAIYDAFARDEINTRDELDACLADKGI RRGHL
 PKMIGILSQEHKNMEEKVRKLEMIADTDHRLDMLDRQTDKIRI GRKNAGLPKSGVIADWLVRDMMR
 FQPVAKDTSGKPLNNSKANST EYRMLQRALALFGGKERLIPYFRQMNLI GGNPHPLD ETRWESHINI
 LSFYRSYLARKAFLEIRGSDRVENRPFLLKKEPKNDRQTLVAGWKESEFHLPRGIFTEAVRDCL IEMGY
 DEVGSYKEVGFMAKAVPLYFERACKDRVQPPYDS PFNVGNLSLKPCKGRFLS KEKRAEEWESGKERFRLAK
 LKKEILEAQEHPYHDFKSWQKFERELRLVKNQDI ITWMMCRDLMEENKVEGLDIGTLYLKD IRTDVHEQG
 SLNVLNVRKPMRLPVVYRADSRGHVHKEAPLATVYI EERDKLLKQGNFKS FVKDRRLNGLF SFVDTG
 GLAMEQYPI SKLRVEYELAKYQARVCFELTLRLEESLLSRYPHLPDES PREMLESWSDP LLAQWPELH
 GKVRLLI AVRNAPSHNQYPMYDEAVFSSIRKYDPS SDAI EERMGLNIAHRLSEEVKQAKETVERI IQA

>gi|700263343|gb|KG005347.1| hypothetical protein HR16_00525

[*Porphyromonas gulae*]

MTEQSERPYNGTYYTLEDKHFWAFLNLARHNAYITLTHIDRQLAYSKADI TNDQDVLSPFKALWKNFDND
 LERKSRLRSLILKHFSPLEGAAYGKKLPESSKSGNKSSKNKELIKKEKEELQANALSLDNLKSI LFDPLQ
 KLDKFRNYYSHYRHSSELPDFGNMLQRLYNVFDVSVQVRKRDHEHNDKVDPHYHFNHLVRKGGKDRY
 GHNDNPSFKHHPVDSGEMVTEAGLLFFVSLFLEKRD AIWMQKKIRGFKGGT ETYQMTNEVFCRSRISLP
 KLKLESLRMDWMLDMLNELVRCPKPLYDRLEKDRACFRVPVD ILPDEDDIDGGGEDPFPKNI LVRHQD
 RFPYFALRYFDLKKVFTSLRFHIDLGTYHFAIYKMKIGEOPEDRHLTRNLVGFGRIQDFAEEHRPEEWR

- continued

LVRDLDFYFETGDKPYYISQTSPHYHIEKGGKILGRFMPGEGHLLWSPPEVGI TRTGRSKYAQDKRLTAEAFSL
 VHELMPMFYFLLREKYS EESVAERVOGR I KRVI EDVYAVYDAFARDEINTRDELDACLADKGI RRGHL
 PRQMIAILSQEHKDMEEKIRKQLQEMMADTDHRLDMLDRQTDRI RIGRKNAGLPKSGVIADWLVRDMMR
 PFPVAKDASGKPLNNSKANSTERYMLQRALALFGGKERLIPYFRQMNLI GGNPHPLHETRWESHINI
 LSFYRSYLRRARAFLEIRIGRSDRVENRPFLLKKEPKTRQTLVAGWKGEPHLP RGI FTEAVRDCLIEMGH
 DEVASYKEVGFMAKAVPLYFERACEDRVQPFYDS PFNVGNSLKP KGRFLSKEERAE EWERGKERFRDLE
 AWSYSAARRI EDAPAGI EYASPGNNKKIEQLLRDLSLWEAFESKLKVRADRINLAKLKEILEAQEHPYH
 DFKSWQKFERELRLVKNQDII TWMMCRDLMEENKVEGLDITGLYLKDIRPNVQEGS LNVLNRVKPMRLP
 VVVYRADS RGHVHKEEAPLATVYI EERDTKLLKQGNFKSFKVDRRLNGLFSFVDTGGLAMEQYPI SKLRV
 EYELAKYQATARVCFELTLRLEESLLTRYPHLPDES FREMLESWSDPLLAKWPELTHGKVRLLI AVRNAPS
 HNQYPMYDEAVFSSIRKYDPS PDIAEERMGLNIAHRLSEEVKQAKETVERIIQA

>gi|736527371|ref|WP_034542281.1| hypothetical protein [*Bacteroides pyogenes*]

MESIKNQSKSIGKILQKDPYFGLYLNMLLNVRKVENHIRKWLGDVALLPEKSGFHSLLTIDNLS SAKW
 TRFYYSRKFPLPFLMEFSDSKSYENRRETAECCLDTIDRQKISSLLKEVYGLQDIRNAPFSHYHIDDQSV
 KHTALIISEMHRF IENYSPALQKTRARFTGVFVETDFLQAEKGDNKKFPAIGGNEGIKLNKDNALIFL
 ICLFLDREBAFKFLSRAFGFKSTKEKGFLAVRETFCALCCRQPHERLLSVNPREALMDMLNELNRC PD I
 LFEMLDEKQKSFPLLLGEEQAHILENSLNDELCEAIDDPFEMIASLSKRVRYKNRFPYMLLRYIEEKN
 LLPFIRFRIDLGCLELASYPKMGEEENYERSVIDHAMAFGRLLTDHNEADAVLQQI IKGITDEVRFSLYA
 PRYAIYNNKIGFVRTSGSDKI SFPILKKKGEGHCVAYTLQNIKSGFPI SIVDLRKLILLFLDKDKAKN
 IVSGLLEQCEKHWDLS ENLFDARTLQKQEPVPLIRYTLPRSKGGKLVSSKLADKQEKYSEFERRKE
 KLTEILSEKDPDLSQIPRRMIDEWLNVLPTSRKLLKGYVETLKLDCRERLRFVEKREKGEHPLPRIGE
 MAIDLAKDIIRMVIDQGVQRI TSAYYSEIQRCLAQYAGDNRRLHLSI IRELRKLDTKNGHPFLGKVL R
 PGLGHTKLYQRVFEKKEWLEATFYPAASPKRVPRFVNPPIGKQKELPLIIRNLMKERPEWRDQKRN
 SHPIDLPSQLFENEICRLLKDKIKGEPGKLNEMPKLYWDKEFPNGMQRFYRCKRREVEVDFKVVYEY
 SEEGNYKYYEALIDEVVRQKISSSKEKSLQVEDLILSVRRVFKRAINEKEYQLRLLCEDRLLFMAV
 RDLYDWEAQLDLDKIDNMLGEPVSVSQVILEGGQPDVIAKAECKLDVSKLMRYCYDGRVKGMLMPYFA
 NHEATQEQVEMELRHYEDHRRRPNWVFALEKSVLNKNEKLRFRYEESSGGCEHRRCIDALRKASLVSEEE
 YEFLVHIRNKS AHNQFPDLEIGKLPNNVISGFCECIWSKYKAIICRIIPFIDPERRFFGKLEEQ

>gi|738977489|ref|WP_036860899.1| hypothetical protein [*Prevotella intermedia*]

MEDDKTTIDSIRYELKDKHFWAFLNLRHNHYITVNHINKILEEGEINRDGYETTLKNIWNEIKDINKK
 DRLSKLI IKHPPFLAAAYRLNPIDTTKQKEEKQAEQSLERLSPFVFIYKLRDLRNHYSHYKHSKSL
 ERPKFEGLLEKMYNI FNASIRLVKEDYQYKNDINPDEDKHLDRTEEEFNYYFTKDNENI TESGLLFF
 VSLFLEKKDAIWMQKLRGFKDNRENKKMTNEVFCRSRMLLPKLRLOSTQTDWI LLDMLNELIRCPKS
 LYERLREEDREKFRVPIEIADEYDAEQEPKNI LVRHQDRFPYFALRYFDYNEIFTNLRFDIDLGTYHF
 SIYKQIGDYKESHHLTHKLYGFERIQEFTKQNRDPDEWRKFKVIFNSPETSKEPIPETTPHYHLENQKI
 GIRFRNDNDKIWPSLKINS EKNEKSKYKLDKSFQAEAFSLVHELPPMFYLLKTEENTDNDNEIETKKK
 ENKNDKQEKHKI EEI ENKI TEIYALYDTFANGEIKS IDELEEEYCKGKDI EIGHLPKQMIAILKDEHKVM
 ATEAERKQEEMLVDVQKSLSDLNQINEEINENVRKNSSLKSGKIASWLVNDMMRFQPVQDNEGKPLNN
 SKANSTEQQLQRTLAFPGSEHERLAPYFKQTKLIESSNPHPLKDI EWKCNLISFYRSYLEAKKNFL
 ESLKPEDWEEKQYFLKLEKPKTKPKTLVQGWKNGFNLP RGI FTEPIRKFWMKHRENI TVALKRVGLVAK
 VIPLFFSEYKDSVQPFYNYHFNVGNINPKDEKNFLNCEERRELLRKKKDEFFKMTDKEKENPSYLEFK
 SWNKFERELRLVRNQDITVWLLCMELFNKKIKELNVEKIYLNKININTIKKEKNT EEKNGEENIKEN
 NILNRIMPMLRPIKVYGRNENSKNKKKIRRNIFFTVYIEEKGTKLLKQGNFKALDRRLGGLFSFVKT
 PSKAESKNTISKLRVEYELGEGYQKARIEIKDMLALEKTLIDKYNSLDIDNFNMLIDWLELKGEPDKA
 SFQNDVLLI AVRNAPS HNQYPMRRI AFANINPFLSSANTS ECKGLGIANQLDKDIHKTI EKIIEIEK
 PIETKE

>gi|739003017|ref|WP_036884929.1| hypothetical protein [*Prevotella falsenii*]

MKNDNNTKSIDYTLGDKHFWAFLNLRHNHYITVNHINKVLELKNKKDQEI IIDNDQDILAIKTLWKG
 VDTDINKDRLELIMKHFPFLAAATYQSSINNTKQKEEBQAKAQS FESLKDCLFLFLKLEKREARNYIS
 HYKSKSLEEPKLEKLENNYNI FIDINVQLVIKDYEHNKDINPEEDFKHLGRAEGEFNYFTRNKKGNI
 TESGLLFFVSLFLEKKDAIWAQTKIKGFKDNRENKQKMTHEVFCRSRMLLPKLRLESTQTDWI LLDMLN
 ELIRCPKSLYKRLQGEKREKFRVFPDPADEYDAEQEPKNI LVRHQDRFPYFALRYFDYNEIFTNLRFO
 IDLGTYHFSIYKQIGDKKEDRHLTHKLYGFERIQEFAKENRPDEWKALVKDLDFEESNBPYISETTPH
 YHLENQKIGIRNKNKKKKTIWPSLETKITVNERSKYNLGSFKAEAFSLVHELPPMFYLLLNKKEPN
 NGKINASKVEGIEIKKIRDIYKLYGAFANEINNEEELKEYCEGKDIAIRHLPKQMIAILKNEYKDMAKK
 AEDKQKMIKDIKKRALALDQVKGVEEDGGRNIPKLSGRIASWLVNDMMRFQPVQRDRDGYPLNNSKA
 NSTEQQLQRTLALFGSERERLAPYFRQMNLI GKNPHPLKDTKWKHEHNNILSFYRSYLEAKKNFLGSL
 KPEDWKNQYFLKLEKPKTNRETQVQWKNGNLP RGI FTEPIREWFIRHONSESEYKVKDFDRIGLVA
 KVIPLFPKEDYQKIEDYVQPFYGYPFNVGNIHNSQEGTFLNKKEREELWKGNTKPKDYKTKENKKEKT
 NKDKFKKIDEEKEEFRSYLDFQSWKFERELRLVRNQDITVWLLCMELIDKLIKIDELNIEELQKRLKLD
 IDIDTAKKEKNMILNRIIMPMLPVIYVETDSSNNI IKDKPLHTIYIEKAETKLLKQGNFKALVKDRRLNG
 LFSFVETSSEABLKSPIKSLVYELGEGYQARARVEIKDMLRLEETLIGNDEKLPINPKMQLDKWLEH
 KKETDDIDLKNDVKLLTEVRNAPS HNQYPMRDRI AFANIKPFLSSANTS NEELGIAKLLDKDIKETID
 RIEIEEQATKR

>gi|739057304|ref|WP_036929175.1| hypothetical protein [*Prevotella sp. MA2016*]

MSKECKQRQEKRRLOKANFSSISLIGKHVPGAYFNMARINFKVTINYILP IAGVRGNYS ENQINKMLHA
 LFLIQAGRNEELTTEQKQWKKLRLNPEQQTKFQKLFKHPFVLP GPMADVADHKAYLNKKSSTVQTEDE
 TFAMLKGVSLADCLDIICLMADLIT ECRNFYTHKDPYKPSQLADQYLHQEMI AKKLDKVVVASRRILKD
 REGLSVNEVEFLIGDHLHGVLDKEFGNAKVKDGKVMKTFVEYDDFYFKISGKRLVNGYITDKDDKPV
 NVNTMLPALSDFLGYFCVFLSKPYAKLFDIEVRLFEYSPDDKENMIMSEMSIYRIRTPRLHKHIDSH

- continued

DSKATLAMDIFGELRRCPMELYNLLDKNAGQPPFHDEVKHPNSHIPDVS KRLRYDDRFPTLALRYIDETE
LFKRIRFQLQLGSPRYKFDKENCIDGRVVRRIQKEINGYGRMQEVDKRMWKWGLIQKREERSVKLE
HEELYINLDQFLEDADSTPYVTDRRPAYNIHANRIGLYWEDSQNPQKYKVFDENGMYIPELVVIEDKKA
PIKMPAPRCALSVYDLPAMLFYEYLRQDNEFFSABEQVIIEYEDDYRKKFFKAVAEGKLPKPKPKFRD
FLKKEYPKLRMADI PKKLQLPLCSHGLCYNNKPEVYERLDRLLQHLLEERELHIQNRLEHYQKDRDMI
NKDNQYGGKSFSDVRHGALARYLAQSMMEWQPTKLDKKEKGDHKLTLGLNYNVLTAYLATYGHQPVEEFG
TPRTLEQVLINAHLIGGNSPHFPINKVLALGNRNIIEELYLHLEELKHRSRIQSLSNSPDKALSALP
FIHHRMRYHERTS EEMMALAARYTTIQLPDGLFTPYILEILQKHYTESSDLQNALSDQVVPVKNLPTCNA
AYLITLFYQTVLKDNAQPFYLSDKTYTRNKDGEKAESEFSAKRAYELFVSVLNNKDIFFPEMIPLFLTSD
EQERLSAKLLDGDGNPVPEVGEKGPATDSQGNITWKRIIYSEVDDYAEKLTDRDMKISFKGEWEKLP
WKQDKI IKRRDETRQMRDELQRMPIYIRDIKDNERTLRRYKTDQDMVFLLAEMFTNIISEQSSEFNW
KQMRLSKVCNEAFLRQTLFRVPVTVGETTIYVEQENMSLKNYGEFYRFLTDDRMLSLNNIVETLKPNE
NGDLVIRHIDLMSLAAVDQYRSTIFMLIQSINENLITINNAVLDDPADGFWVREDLPKRNFFASLLELI
NQLNNVELTDDRKLLVAIRNAFNSHNSYINDFSLIKDVKHLPEVAKGILQHLQSMGLVEITK

>gi|739060066|ref|WP_036931485.1| hypothetical protein [*Prevotella pleuritidis*]
MENDKRLEESTCYTLNDKHFWAAPFLNLRHNAYITINHINKLLEIRQIDNDEKVLDIKALWQKVKDINQ
KARLRELMIKHPFFLEAAIYSNNKEDKEEVKEEKQAKAQSFKSLKDCFLFLEKLQEARNYSHYKSSSES
SKEPEFEGLLEKMYNTFGVSIKLVKEDYQYKNDIDPEKDPKHLERKEDFNLYFTDKDNKGIITKNGLLF
FVSLFLEKDAIWMQKLRGFKDNRGNKEKMTHEVFCRSRMLLPKIRLESTQTDQWILLDMLNELIRCPK
SLYERLQGAYREKFKVPFDSIDEDYDAEQEPFRNTLVRHQDRFPYFALRYFDYNEIFKNLRFQIDLGTYH
FSIYKLLIDGNKEDRHLTHKLYGPERIQEFAKQKRPNEWQALVKDLDIYETSNEQYISETTPHYHLENQK
IGIRFKNKDKIWPVLENGKNEKSKYNLDKSFAEAFSLIHELLPMMFYDLKKEEPNDEKNAISIV
EGFIKKEIKRMYAIDAFANEEINSKEGLEEYCKNKGQERHLPKQMIATLINKSKNMAEKAKRQKEMI
KDIKRLALIDKQVKEIEDGGNRIRLLKSGEIRARWLNDMMRQSQVKDKKQKPLNNSKANSTYEQMLQ
RSLALYNKEQKPIPIYIQLVNLKSSNPHFLEETKWECCNNILSFYRSYLEAKNFLESLEKPEDWKKNOY
FLMLKEPKTRKTLVQGWKNGFNLRGIFTEPIKEWFKRHQNDSEYKKEVLEALDRVGLVAKVIPLFFKKEE
YFKEDAQKEIINNCVQPFYSFPYVGNIIHKPEEKNFLHCEERRKLDWKKDKFKGYKAKEKSKMTDKEKE
EHRSLYEFQSNWFKERERLVRNQDITVWLLCTELIDKLLIDELNIEELQKLRKIDITDTAKKEKNIL
NRIMPQLPVIYVEIDKSFNIVKDKPLHTIYIEEIGTKLLKQGNFKALVKDRRLNGLFSFVKISSEAESK
SKPISKLRVEYELGAYQKARIDIKMDLAEKTLIDNDENLPIKFSMDLKSRLWKGKGEANKARLQNDVD
LLVAIRNAFNSHNQYPMYNSVFPKMKLSSLSSDIPEKEGLGIKQLKDKIKETIERIIEIEKEIRN

>gi|746373271|ref|WP_039417390.1| hypothetical protein [*Porphyromonas gingivalis*]
MTEQNERPYNGTYTLEDKHFWAAPFLNLRHNAYITLAHIDRQLAYSKADITNDEDILFFKQGWKNLDND
LERKARLSLILKHFSPLEGAAYGKLLPESQSSGNKSSKKEKELIKKEEELQANALSLDNLKSLIFDFLQ
KLKDFRNYSHYRHPSSELPDFDGNMLQRLYVNFVDSVQVRKRDHEHNDKVDPHRFNHLVRKGGKDRY
GMNDNPFKHHFVDREGIVTEAGLLFFVSLFLEKRDALWMQKIRGFKGGTEAYQOMTNEVFCRSRISLP
KLKLESRLIDWMLDMLNELVRCPKSLYDRLREEDRARFRVPIDILSDEDDIDGTEEDPKNILVRHQD
RFPYFALRYFDLKKVFTSLRPHIDLGTYHFAIYKKNIGEQPEDRHLTRNLYGFGRIQDFAEEHRPEEWR
LVRDLDFPETGDKPYITQTPHYHIEKGIKGRFVPEGQHLWSPPEVGAIRTRGRSKYAQDKRLTAEAFSL
VHELMPMFYYFLLEKYS EEVSAEKVQGRIKRVIDEYAVYDAFARGEIDILDRDLACLADKGRIRRHGL
PRQMIAILSQEHKMEKVRKLLQEMIADTDHRLDMLDRQTDKIRIRGRKNAGLPKSGVIADWLVRDMMR
PQPVAKDTSGKPLNNSKANSTYRMLQRALALFGGEKERLIPYFRQMLIGGNNPHFLHETRWESHINI
LSFYRSYLKARKAFLOSIGRSDRRENHRFLLEKPKTDQRQTLVAGWKS EPHLPRGIFTEAVRDCLEIMGY
DEVGSYKEVGFMAKAVPLYFERACKDRVQPFYDYPFNVGNLSLKKKGRFLSKEKRAEWEWESGKERFLAK
LKKEILEAKEHPYLDKFSWQKFERELRLVKNQDITWMMCRDLMEENKVEGLDIGTLYLKDIRTVDVHEQG
SLNVLRVPMRLPVVYRDRSGRHHVHKEQAPLATVYIEERDTKLLKQGNFKSFVKDRRLNGLFSFVDTG
ALAMEQYPI SKLRVEYELAKYQATARVCAFEQTLLELESLLTRYPHLPDKNFKMLLESWSDPLLDKWPDLH
RKVRLLIAVRNAPNSHNQYPMYDEAVFSSIRKYDPSSPDAI EERMGLNIAHRLSEEVKQAKEMAERI IQV

>gi|746374939|ref|WP_039418912.1| hypothetical protein [*Porphyromonas gulae*]
MTEQSERPYNGTYTLEDKHFWAAPFLNLRHNAYITLTHIDRQLAYSKADITNDQDVSFKALWKNLDND
LERKSRLSLILKHFSPLEGAAYGKLLPESKSSGNKSSKKEKELIKKEEELQANALSLDNLKSLIFDFLQ
KLKDFRNYSHYRHSGSSELPDFDGNMLQRLYVNFVDSVQVRKRDHEHNDKVDPHRFNHLVRKGGKDRY
GHNDNPSFKHHFVDSSEGWTEAGLLFFVSLFLEKRDALWMQKIRGFKGGTEAYQOMTNEVFCRSRISLP
KLKLESRLMDDWMLDMLNELVRCPKPLYDRLREDDRACFRVPVDILPDEDDIDGGEDDPKNI LVRHQD
RFPYFALRYFDLKKVFTSLRPHIDLGTYHFAIYKKNIGEQPEDRHLTRNLYGFGRIQDFAEEHRPEEWR
LVRDLDFPETGDKPYISQTSPHYHIEKGIKGRFVPEGQHLWSPPEVGI TRTGRSKYAQDKRLTAEAFSL
VHELMPMFYYFLLEKYS EEVSAEKVQGRIKRVIDEYAVYDAFARGEINTLKELDACLADKGRIRRHGL
PKQMIAILSQEHKMEKVRKLLQEMIADTDHRLDMLDRQTDKIRIRGRKNAGLPKSGVIADWLVRDMMR
PQPVAKDSGKPLNNSKANSTYRMLQRALALFGGEKERLIPYFRQMLIGGNNPHFLHETRWESHINI
LSFYRSYLKARKAFLEIRIGRSDRMENRPFLLLEKPKTDQRQTLVAGWKS EPHLPRGIFTEAVRDCLEIMGY
DEVGSYREVGFMKAVPLYFERACKEDRVQPFYDYPFNVGNLSLKKKGRFLSKEEAAEWEWESGKERFRDLE
AWSHSAARRIEDAFAGI EYASPGNKKIEQLLRDLSLWEAPESKLVKVDKINLAKLKEILEAQEHPHYH
DFKSWQKFERELRLVKNQDITWMMCRDLMEENKVEGLDIGTLYLKDIRTNVQEQGSLNVLNVHVKPMRLP
VVVYRDRSGRHHVHKEEAPLATVYIEERDTKLLKQGNFKSFVKDRRLNGLFSFVDTGGLAMEQYPI SKLRV
EYELAKYQATARVCAFEQTLLELESLLTRYPHLPDKNFKMLLESWSDPLLDKWPDLHAKWPELHGEVRLLI
AVRNAPNSHNQYPMYDEAVFSSIRKYDPSSPDAI EERMGLNIAHRLSEEVKQAKETVERIIQA

>gi|746375986|ref|WP_039419792.1| hypothetical protein [*Porphyromonas gulae*]
MTEQSERPYNGTYTLEDKHFWAAPFLNLRHNAYITLTHIDRQLAYSKADITNDQDVSFKALWKNLDND
LERKSRLSLILKHFSPLEGAAYGKLLPESKSSGNKSSKKEKELIKKEEELQANALSLDNLKSLIFDFLQ
KLKDFRNYSHYRHSGSSELPDFDGNMLQRLYVNFVDSVQVRKRDHEHNDKVDPHRFNHLVRKGGKDRY

- continued

GHNNDNPSFKHHFVDEGEMVTEAGLLFFVSLFLEKRDATWQKIRGFKGGTETYQQMTNEVFCRSRISLP
 KLKLESRLIDDWMLLDMLNELVRCPKPLYDRLEKDRARFRVVD ILPDEDDIDGGGEDPPKNI LVRHQD
 RFPYFALRYFDLKKVFTSLRPHIDLGTYHFAIYKVI GEQPEDRHLTRNL YGFGRIQDFAEEHRPEEWKR
 LVRDLDFETGDKPYISQTPHYHI EKGKILGRFVPEGQHLWPSPEVGI TRTGRSKYAQDKRLTAEAFSL
 VHELMPMMFYFLLREKYS EEVSAEKVQGR I KRVI EDVYAIYDAFARDEINTRDELDACLADKGI RRGHL
 PKQMI GILSQEHKNMEEKVRKQLQEMI ADTDHRLDMLDRQTDRI RIGRKNAGLPKSGVIADWLVRDMMR
 FQPVAKDTSGKPLNNSKANST EYRMLQRALALFGGKERLIPYFRQMNLI GGNPHPLDETRWESHINI
 LSFYRSYLARKAFLEIRGRSDRVENRPFLLKEPKTDRTLVAGWKSEFHLPRGI FTEAVRDCLIMGY
 DEVGSYKEVGFMAKAVPLYFERACKDRVQPPYDS PFNVGNSLKP KGRFLS KEKRAEEWESGKERFRLAK
 LKKEILEAQEHPYHDFKSWQKFERELRLVKNQDI ITWMMCRDLMEENKVEGLDIGTLYLKD IRPNVQEQG
 SLNVLNRVKPMRLPVVYRADSRGHVHKEEAPLATVYI EERDTKLLKQGNFKS FVKDRRLNGLFSFVDTG
 GLAMEQYPI SKLRVEYELAKYQ TARVCVFELTLRLEESLLSRYPHLPDES FREMLESWSDPLLAKWPELH
 GKVRLLI AVRNAF SHNQYPMYDEAVFSSIRKYDPS SPDAI EERMGLNIAHRLSEEVKQAKETVERIIQA

>gi|746382772|ref|WP_039426176.1| hypothetical protein [*Porphyromonas
 gulae*]
 MTEQSERPYNGTYTLEDKHFWAFLNLARHNAYITLTHIDRQLAYS KADI TNDQDVLSPFKALWKNFDND
 LERKSRLRSLILKHFSPLEGAAYGKKL PESKSSGNKSSKNKELIKKEKEELQANALSLDNLKSI LDFDPLQ
 KLKDFRNYSHYRHSSESLP LFDGNMLQRLYNVFDVSVQRVKRDHEHNDKVDPHYHFNHLVRKGGKDRY
 GHNNDNPSFKHHFVDEGEMVTEAGLLFFVSLFLEKRDATWQKIRGFKGGI GPYEQMTNEVFCRSRISLP
 KLKLESRLIDDWMLLDMLNELVRCPKPLYDRLEKDRACFRVVD ILPDEDDIDGGGEDPPKNI LVRHQD
 RFPYFALRYFDLKKVFTSLRPHIDLGTYHFAIYKVI GEQPEDRHLTRNL YGFGRIQDFAEEHRPEEWKR
 LVRDLDFETGDKPYISQTPHYHI EKGKILGRFMPGEGQHLWPSPEVGI TRTGRSKYAQDKRLTAEAFSL
 VHELMPMMFYFLLREKYS EEVSAEKVQGR I KRVI EDVYAIYDAFARDEINTLKELDACSDKGI RRGHL
 PKQMI GILSQEHKNMEEKVRKQLQEMI ADTDHRLDMLDRQTDRI RIGRKNAGLPKSGVIADWLVRDMMR
 FQPVAKDTSGKPLNNSKANST EYRMLQRALALFGGKERLIPYFRQMNLI GGNPHPLDETRWESHINI
 LSFYRSYLARKAFLEIRGRSDRVENRPFLLKEPKTDRTLVAGWKSEFHLPRGI FTEAVRDCLIMGY
 DEVGSYKEVGFMAKAVPLYFERACKDRVQPPYDS PFNVGNSLKP KGRFLS KEKRAEEWESGKERFRLAK
 LKKEILEAQEHPYHDFKSWQKFERELRLVKNQDI ITWMMCRDLMEENKVEGLDIGTLYLKD IRPNVQEQG
 SLNVLNRVKPMRLPVVYRADSRGHVHKEEAPLATVYI EERDTKLLKQGNFKS FVKDRRLNGLFSFVDTG
 GLAMEQYPI SKLRVEYELAKYQ TARVCAFEQTLLEESLLTRYPHLPDEN FREMLESWSDPLLKWPDLH
 GKVRLLI AVRNAF SHNQYPMYDEAVFSSIRKYDPS SPDAI EERMGLNIAHRLSEEVKQAKETVERIIQA

>gi|746385904|ref|WP_039428968.1| hypothetical protein [*Porphyromonas
 sp. COT-052 0H4946*]
 MTEQSERPYNGTYTLEDKHFWAFLNLARHNAYITLTHIDRQLAYS KADI TNDQDVLSPFKALWKNFDND
 LERKSRLRSLILKHFSPLEGAAYGKKL PESKSSGNKSSKNKELIKKEKEELQANALSLDNLKSI LDFDPLQ
 KLKDFRNYSHYRHSSESLP LFDGNMLQRLYNVFDVSVQRVKRDHEHNDKVDPHRHFNHLVRKGGKDRY
 GHNNDNPSFKHHFVDEGEMVTEAGLLFFVSLFLEKRDATWQKIRGFKGGTETYQQMTNEVFCRSRISLP
 KLKLESRLIDDWMLLDMLNELVRCPKPLYDRLEKDRACFRVVD ILPDEDDIDGGGEDPPKNI LVRHQD
 RFPYFALRYFDLKKVFTSLRPHIDLGTYHFAIYKVI GEQPEDRHLTRNL YGFGRIQDFAEEHRPEEWKR
 LVRDLDFETGDKPYISQTPHYHI EKGKILGRFVPEGQHLWPSPEVGI TRTGRSKYAQDKRLTAEAFSL
 VHELMPMMFYFLLREKYS EEVSAEKVQGR I KRVI EDVYAIYDAFARDEINTLKELDACSDKGI RRGHL
 PKQMI GILSQEHKNMEEKVRKQLQEMI ADTDHRLDMLDRQTDRI RIGRKNAGLPKSGVIADWLVRDMMR
 FQPVAKDTSGKPLNNSKANST EYRMLQRALALFGGKERLIPYFRQMNLI GGNPHPLHETRWESHINI
 LSFYRSYLARKAFLEIRGRSDRVENRPFLLKEPKTDRTLVAGWKSEFHLPRGI FTEAVRDCLIMGY
 DEVGSYREVGFMAKAVPLYFERACEDRVQPPYDS PFNVGNSLKP KGRFLS KEDRAEEWERGKERFRDLE
 AWSHSAARRIKDAFAGI EYASPGNKKIEQLRDL SLWEAFES KLKVRADKINLAKLKEILEAQEHPYH
 DFKSWQKFERELRLVKNQDI ITWMMCRDLMEENKVEGLDIGTLYLKD IRPNVQEQGSLNVLNRVKPMRLP
 VVYRADSRGHVHKEEAPLATVYI EERDTKLLKQGNFKS FVKDRRLNGLFSFVDTGGLAMEQYPI SKLRV
 EYELAKYQ TARVCVFELTLRLEESLLSRYPHLPDES FREMLESWSDPLLAKWPELHGKVRLLI AVRNAF
 HNQYPMYDEAVFSSIRKYDPS SPDAI EERMGLNIAHRLSEEVKQAKETVERIIQA

>gi|746388969|ref|WP_039431778.1| hypothetical protein [*Porphyromonas
 gulae*]
 MTEQSERPYNGTYTLEDKHFWAFLNLARHNAYITLTHIDRQLAYS KADI TNDQDVLSPFKALWKNFDND
 LERKSRLRSLILKHFSPLEGAAYGKKL PESKSSGNKSSKNKELIKKEKEELQANALSLDNLKSI LDFDPLQ
 KLKDFRNYSHYRHSSESLP LFDGNMLQRLYNVFDVSVQRVKRDHEHNDKVDPHRHFNHLVRKGGKDRY
 GHNNDNPSFKHHFVDEGEMVTEAGLLFFVSLFLEKRDATWQKIRGFKGGTETYQQMTNEVFCRSRISLP
 KLKLESRLIDDWMLLDMLNELVRCPKPLYDRLEKDRACFRVVD ILPDEDDIDGGGEDPPKNI LVRHQD
 RFPYFALRYFDLKKVFTSLRPHIDLGTYHFAIYKVI GEQPEDRHLTRNL YGFGRIQDFAEEHRPEEWKR
 LVRDLDFETGDKPYISQTPHYHI EKGKILGRFMPGEGQHLWPSPEVGI TRTGRSKYAQDKRLTAEAFSL
 VHELMPMMFYFLLREKYS EEVSAEKVQGR I KRVI EDVYAIYDAFARDEINTLKELDACSDKGI RRGHL
 PKQMI GILSQEHKNMEEKVRKQLQEMI ADTDHRLDMLDRQTDRI RIGRKNAGLPKSGVIADWLVRDMMR
 FQPVAKDTSGKPLNNSKANST EYRMLQRALALFGGKERLIPYFRQMNLI GGNPHPLHETRWESHINI
 LSFYRSYLARKAFLEIRGRSDRVENRPFLLKEPKTDRTLVAGWKSEFHLPRGI FTEAVRDCLIMGY
 DEVGSYREVGFMAKAVPLYFERACEDRVQPPYDS PFNVGNSLKP KGRFLS KEKRAEEWERGKERFRDLE
 AWSHSAARRIKDAFAGI EYASPGNKKIEQLRDL SLWEAFES KLKVRADKINLAKLKEILEAQEHPYH
 DFKSWQKFERELRLVKNQDI ITWMMCRDLMEENKVEGLDIGTLYLKD IRPNVQEQGSLNVLNRVKPMRLP
 VVYRADSRGHVHKEEAPLATVYI EERDTKLLKQGNFKS FVKDRRLNGLFSFVDTGGLAMEQYPI SKLRV
 EYELAKYQ TARVCVFELTLRLEESLLSRYPHLPDES FREMLESWSDPLLAKWPELHGKVRLLI AVRNAF
 HNQYPMYDEAVFSSIRKYDPS SPDAI EERMGLNIAHRLSEEVKQAKETVERIIQA

>gi|746392082|ref|WP_039434803.1| hypothetical protein [*Porphyromonas
 gulae*]
 MTEQSERPYNGTYTLEDKHFWAFLNLARHNAYITLTHIDRQLAYS KADI TNDQDVLSPFKALWKNFDND
 LERKSRLRSLILKHFSPLEGAAYGKKL PESKSSGNKSSKNKELIKKEKEELQANALSLDNLKSI LDFDPLQ
 KLKDFRNYSHYRHSSESLP LFDGNMLQRLYNVFDVSVQRVKIDHEHNDKVDPHYHFNHLVRKGGKDRY

- continued

GHNNDNPSFKHHFVDEGEMVTEAGLFFVSLFLEKRDATWQKKIRGFKGGTETYQQMTNEVFCRSRISLP
 KLKLESLRMDWMLLDMLNELVRCPKPLYDRLEDDRACFRVVPVDILPDEDDIDGGGEDPPFKNILVRHQD
 RFPYFALRYFDLKKVFTSLRFHIDLGTYHFAYKKMIGEQQEDRHLTRNLYGFGRIQDFAEEHRPEEWR
 LVRDLDYFETGDKPYISQTSPHYHIKGGKILGRFVPEGQHLWSPPEVGI TRTGRSKYAQDKRLTAEAFSL
 VHELMPMMFYFLLREKYS EEVSAERVQGR I KRVI EDVYAVYDAFARDEINTRDELDACLADKGIIRRHGL
 PRQMIAILSQEHKDMEEKIRKQLQEMMADTDHRLDMLDRQTDKIRIRGRKNAGLPKSGVIADWLVRDMMR
 FQPVAKDAAGKPLNNSKANSTERYMLQRALALFGGKERLIPYFRQMNLI GGNPHPPFLHETRWESHINI
 LSFYRSYLARKAFLEIRGRSDRVENRPFLLKPKTDRTQTLVAGWKGFHLPRGI FTEAVRDCLIEMGH
 DEVASYKEVGFMAKAVPLYFERACEDRVQPFYDS PFNVGNSLKPCKGRFLSKEERAEWEERGERFRDLE
 AWSYSAARRIEDAFAGI EYASPGNKKIEQLLRDLSSLWEAFESKLKVRADRINLAKLKEILEAQEHPYH
 DFKSWQKFERELRLVKNQDII TWMMCRDLMEENKVEGLDITGLYKDIRPNVQEQGSLNVLNRVKPMRLP
 VVVYRADS RGHVHKEEAPLATVYI EERDTKLLKQGNFKSFKVDRRLNGLFSFVDTGGLAMEQYPI SKLRV
 YELAKYQATARVCFELTLRLEESLLTRYPHLPDES FREMLESWSDDLAKWPELHGKVRLLI AVRNAFS
 HNQYPMYDEAVFSSIRKYDPSPPDAIEERMGLNIAHRLSEEVKQAKETVERIIQA

>gi|746399962|ref|WP_039442171.1| hypothetical protein [*Porphyromonas*
gulae]
 MTEQSERPYNGTYTLEDKHFWAFLNLRHNAYITLTHIDRQLAYSKADITNDQDVLSPKALWKNLNDND
 LERKSRRLSLILKHFSPLEGAAYGKFLPESKSGNKSSKNKELIKKEKEELQANALS LNDLKSILFDLQ
 KLKDFRNYSHYRHSSELEPLDFGNMLQRLYNVFDVSVQVRKRDHEHNDKVDPHYHFNHLVRKGGKDRY
 GHNNDNPSFKHHFVDEGEMVTEAGLFFVSLFLEKRDATWQKKIRGFKGGI GPYEQMTNEVFCRSRISLP
 KLKLESLRIDWMLLDMLNELVRCPKPLYDRLEKDRACFRVVPVDILPDEDDIDGGGEDPPFKNILVRHQD
 RFPYFALRYFDLKKVFTSLRFHIDLGTYHFAYKKMIGEQQEDRHLTRNLYGFGRIQDFAEEHRPEEWR
 LVRDLDYLETGDKPYISQTTPHYHIKGGKILGRFVPEGQHLWSPPEVGI TRTGRSKYAQDKRLTAEAFSL
 VHELMPMMFYFLLREKYS EEVSAEKVQGR I KRVI EDVYAVYDAFARDEINTLKELDTCLADKGIIRRHGL
 PKQMITILSQERKDMKEKIRKQLQEMIADTDHRLDMLDRQTDKIRIRGRKNAGLPKSGVIADWLVRDMMR
 FQPVAKDAAGKPLNNSKANSTERYMLQRALALFGGKERLIPYFRQMNLI GGNPHPPFLHETRWESHINI
 LSFYRSYLARKAFLEIRGRSDRVENRPFLLKPKTDRTQTLVAGWKGFHLPRGI FTEAVRDCLIEMGY
 DEVGSYREVGFMAKAVPLYFERACEDRVQPFYDS PFNVGNSLKPCKGRFLSKEDRAEWEERGERFRDLE
 AWSYSAARRIEDAFAGI EYASPGNKKIEQLLRDLSSLWEAFESKLKVRADKINLAKLKEILEAQEHPYH
 DFKSWQKFERELRLVKNQDII TWMMCRDLMEENKVEGLDITGLYKDIRPNVQEQGSLNVLNRVKPMRLP
 VVVYRADS RGHVHKEEAPLATVYI EERNIKLLKQGNFKSFKVDRRLNGLFSFVDTGGLAMEQYPI SKLRV
 YELAKYQATARVCFELTLRLEESLLSRYPHLPDES FREMLESWSDDLAKWPELHGKVRLLI AVRNAFSH
 NQYPMYDEAVFSSIRKYDPSPPDAIEERMGLNIAHRLSEEVKQAKETVERIIQA

>gi|746403165|ref|WP_039445055.1| hypothetical protein [*Porphyromonas*
gulae]
 MNIVPATENKQSRIVEDDPQYFGLYLNLARENLEIEVESHVRIKFGKKKLNEESLKQSLLCDHLLSIDRW
 TKVYGHSSRYLPLFHCDFPDSGI EKDHDSKTGVDPSAQRLIRELYSLDLRNDFSHNRDLGIFTEHLK
 VSPDISSFITGAYTACERASRFADFPKDDFLAKNRKEQLISVADGKECLIVSGFAFFICLFLDREQ
 ASGMLSRIRGPKRIDENWAVHETFCDL CIRHPHDLRESNI KEALLDMLNELNRCPRILYDMLPEEE
 RAQFLPALDENSMNLS ENSLNEESRLLDWGS SDWAEALTKIRHQDRFPYMLLRFIEEMDLLKGI RFRV
 DLGEIELDSYSKVKGRNGEYDRTITDHALAFGLKSDPFQNEEEVSRMISGEASYPVRFSFLFAPYAIYDNK
 IGYCHTSDPVYPKSTGEBKRLSNPQSMGFI SVHDLRKLMLMELLCEGSFRMQSDFLRKNRILDETAE
 GKLFQSFALPEMRHRFIPQNPKSKDRREKAETTLKYKQETKGRKDKLNSQLSAPFDMNQQLPSRLLD
 EWMNIRPASHSVKLRTYVQNLNEDCRLRLRFRKRDGDKARAIPLVGMATFLSQDIVRMIISSETKKLI
 TSAYNEMQSLAQYAGEENRRQPRAI VAEHLHLLDSSGHPFLSATMETAHRYTEDFYKCYLEKREWLA
 KTFYRPEQDENTKRIISVFFVDPGEARKLLPLLIRRRMKEQNDLQDWIRNKQAHPI DLPSHLFDSKIMEL
 LKVKDGKKNWNEAFKDWWS TKYPDGMQPFYGLRRELNIHGKSVSYIPSDGKFPADCYTHLMEKTVRDKKR
 ELRTAGKVPVDPDLAAYIKRSPHRAVNEREFMLRLVQEDDRMLMLAINKMMTDREEDILPGLKNIDSILDE
 ENQFSLAVHAKVLEKEGEGGDNLSLVPATIEIKSKRKDWSKYIRYRDRRVPGLMSHPPEHKAILEVVK
 ILLGEYDRCRIKIFDWAFALEGAIMSDRDLKPYLHESSREGKSGEHS TLVKMLVEKKGCLTPDESQYLI
 LIRNKAHNQFPFAEEMPLIYRDVSAKVSIEGSSAKDLPEGSSLVDSLWKKYEMIIRKILPLDHENRF
 FGKLLNMSQPINDL

>gi|749365741|emb|CEN50095.1| conserved hypothetical protein
 [*Capnocytophaga canimorsus*]
 MDEKPTIIKGIKFGNEALFAHFRYYKEYQNFQTFYDTENYPLVKLGKKQADRKRKTKIYQKKNVFT
 LLMAKHIFKSVFKQDSIDRFSLEDLYQSRERLGNQERARQTGERNTNYIWNKTVDLKLCDGKI TVENVK
 LKNVGDIFI KYEVDQRVQAFITYEENI EWQAFLIKESKEENYVYVVERIEIQYEVKREELLKEVHLIEE
 YILEKVKDKKELIKKGDGNQNFYIILNGLLKQVKEKEKEDVESYKVFNLNI KPEDVDINQLKQKAI DLBQ
 KAFVLTIRNKAHNQFPFAEEMPLIYRDVSAKVSIEGSSAKDLPEGSSLVDSLWKKYEMIIRKILPLDHENRF
 FGKLLNMSQPINDL

>gi|749377015|emb|CEN32480.1| conserved hypothetical protein
 [*Capnocytophaga cynodegmi*]
 MENKTSLGNNIYYNPFKQDKSYFAGYLNAAENIDSVRELGKRLKGGKEYISENFDAIFKENISLVEY
 ERYVKLSDYFPMARLLDKKEVPIKERENFKNFRGI IKAVRDLRNFYTHKEHGEVEITDEIFGVLDLDM
 LKSTVLIVKKKIKTKDKIEKLSIEKQLDILCQKLELYLKD TARKIEBKRRNRQERGEKLVPRFYS
 DRRDDLAAIYNDAFDVYIDKKKDSLKESSTKYNTESYPOQEEGDLKIPISKNGVVFLLSFLSKQEVH
 AFKSKIAGFKATVIDEATVSHRKNISCFMATHEIPSHLAYKLLKRVRTAEINYSAEANBQLS IYAKET
 LMMQMLDELSKVPDVVYQNSLEDVQKTFI EDWNEYLKENNGDVGTEEBEQV IHPVIRKRYEDKPNYFAIR
 FLDEFAQFPILRFQVHLGNYLHDSRPKEHLISDRRI KEKIVFGRLSLEHKKALFINKTETNEDRKHYW
 EVFPNPNYDPPKENISVNDKDFPIAGS ILDRKQPIAGKIGIKVNLNQQYISEVDKAVKAHQKQRNKNK
 PSTIQNIIEEIVPINGSNPKIEVEGGQPTAYLSMNDIHSILYEPFDKWEKKEKLEKKEKELRKEIGKE
 LEEKIVGKIQTQIQI I DKNINAKILKPYQDDDS TAIDKEKLIKDLKQEQKILQKLNKQETAREKEYQEC
 IAYQEEBRSKIKRSDKSRQKYLNRQLKRKYPEVPIRKEILYEQEKGVAVWLANDIKREMPDTPKNEWKGE
 QHSLQKSLAAYEQCKEGLKNLLPQKVFPHLPFELGGHFQKQYLYQPYTRYLDKRLLEHIIAVRQAENF
 KNENKVFVKVENECPKFLKQNYTHKGLDAQS VLGYPILFERGFMDKPTIIKGIKFGNESLFTDWF

- continued

RYYKEYQNFQTFYD TENYPLVELEKKQADRKRRETKIYQQKNDVFI LLMAKHI P KSVFKQDSIDRFSLED
LYQSREERLENEQKAKQTGERNTNYIWNKTVDLNLCDGKVIENVKLNKVNFIKYEYDQRVQTFPKYEE
NIKWAFLIKESKEEENYPYI VERIEQYKVRREELKEVHLIEEYILEKVKDKEILKKGDNQNFKYI
LNGLLQKLNEDVESYKVFNLNI KPEDVNIINQLKQEA TDLBQKAFVLT YIRNKFAHNQLPKKFEWDYQCQ
KYGKIEKETYAEYFAEVFKREKEALMK

>gi|754125679|ref|WP_041758328.1| hypothetical protein [*Psychroflexus
torquais*]
MYVLNGLLIGIKQINIEDFIVLKQINFPDKIDFKGIASYSVLEKKTIIILIAIRNKFAHNQLPNKIINDLA
NEFVKKKEKNETYANYLLKVLKMMISDLA

>gi|754596085|ref|WP_041986150.1| hypothetical protein [*Capnocytophaga
canimorsus*]
MDEKPTIIKGIKFGNEALFAHWFRYYKEYQNFQTFYD TENYPLVKLGKQADRKRKTKIYQQKNDVFT
LLMAKHI P KSVFKQDSIDRFSLEDLYQSREERLGNQERARQTGERNTNYIWNKTVDLNLCDGKI TVENVK
LKNVGDPIKYEYDQRVQAFLEYEENIEWQAFLIKESKEEENYPVVERIEQYKVRREELKEVHLIEE
YILEKVKDKEILKKGDNQNFKYI LNGLLQKLVKKEKEDVESYKVFNLNI KPEDVDINQLKQKAIDLEQ
KAFVLT YIRNKFAHNQLPKKFEWDYQCQEBGKIAKGKTYAEYFVEVFKREKEALMK

>gi|754599669|ref|WP_041989581.1| hypothetical protein [*Capnocytophaga
cynodegmi*]
MENKTSLGNNIYYPFKPQDKSYFAGYLNAAMENIDSVFRELKRLKKGKEYISENFPDAIFKENISLVEY
ERYVKLLSDYFPMARLLDKKEVPIKERKENPKKNFRGI IKA VRDLRN FYTHKEHGEVEITDEIFGVLD
L KSTVLIVKKKIKTKDKIKELKKSIEKQLDILCQKLEYLKD TARKIEBKRNRQRERGEKLLVPRFEYS
DRRDDLIAAIYNDADFVYIDKKKDSLKESSTKYNTESYPQOEGDLKIPISKNGVVVFLSLFLSKQEVH
AFKSIAGFKATVIDEATVSHRKN SICPMATHEIFSHLAYKLLKRKVRTAEINYSAEANAQLSIYAKET
LMMQMLDRLSKVPDVVYQNLSEDOVQKTFIEDWNEYLKENNGDVGTMEEEQV IHPVIRKRYEDKFN YFAIR
FLDEFAQFPPI LRFPVHLGNVYLDHRPKEHLISDRRIKEKIVTFGRLSLEHKKALFINKNETNDRKH YW
EVFPNPNYDFPKENISVNDKDFPIAGS ILDREKQPIAGKIGIKVNLNQQYISEVDKAVKAHLKQRNKN
PSIQNIIEEIVPINGSNPKIEIVEGGQPTAYLSMNDIHSILYEFDDKWEKKKLEKKEKELRKEIGKE
LEEKIVGKIQTQIQIIDKINAKILKPYQDDSDTAIDKELIKDLKQEQKILQKLNKNEQTAREKEYQEC
IAYQESRKIKRSDKSRQKYLRLNQLKRYPEVPIRKEILYQEKGVAVVWLANDIKREMPDFKNKWEKGE
QHSLLQKSLAYYEQCKEELKLLPQKQVFKHLPFELGGHFQKQYLYQPYTR YLDRLEHISGLVQQAENF
KNENKVPKVENECFKFLKQNYTHKGLDAQS VLGYPFLERGFMDKPTIIKGIKFGNESLFTDWF
RYYKEYQNFQTFYD TENYPLVELEKKQADRKRRETKIYQQKNDVFI LLMAKHI P KSVFKQDSIDRFSLED
LYQSREERLENEQKAKQTGERNTNYIWNKTVDLNLCDGKVIENVKLNKVNFIKYEYDQRVQTFPKYEE
NIKWAFLIKESKEEENYPYI VERIEQYKVRREELKEVHLIEEYILEKVKDKEILKKGDNQNFKYI
LNGLLQKLNEDVESYKVFNLNI KPEDVNIINQLKQEA TDLBQKAFVLT YIRNKFAHNQLPKKFEWDYQCQ
KYGKIEKETYAEYFAEVFKREKEALMK

>gi|755251047|gb|KIP58950.1| hypothetical protein ST43_06385
[*Prevotella* sp. P5-60]
MNI PALVENQKKYFGTYSVMAMLNAQT VLDHIQKVADIEG EQNENNENLWPHVMSHL YNAKNGYDKQPE
KTMFIIERLQSYFPLKMAENQREYSNGKYQNRVEVNSNDI FEVLKRAFVGLKMYRDLINHYKTYEEK
LIDGCEFLTSTEQPLSGMINSKYTYVALRNIKERYGYKAEDLAFIQDNRYKPTKDAYGKRKSQVNI GSFLS
LQDYNQDTIKKLLHLSGVGIALLICFLDKQYINFLSRLPIFSSYNAQSEBERRIIRSFGINSIKQPKDR
IHSEKSNKSVAMDMLNEVKRCPDELFTLSAEKQSRFR IISDDHNEVLMKRSDDRFPVLLQYIDYGLKLF
DHIRFHVNMGLR YLLKADKICIDGQTRVRIEQPLNGFGRLEEVETMRKQENGTFGNSGIRIRDFENMK
RDDANPANYPIVETTYTHYLENNKVMEMFISDEENPTLPLVIEDDRYVVKTI PSCRMSTLEIPAMAFHM
FLFGSEKTEKLIIDVHDYKRLRFQAMQKEEVTAENIASFGIAESDLPQKIMDLISGNAHGKD VDAFIRLT
VDDMLIDTERRIKRFDKDRKRSADNKMGRGFKQISIGKLADFLAKDIVLFPQPSVNDGENKITGLNYR
IMQSAIAVYDSGDDYEAQQQKLMFEKARLIGKITEPHFPYKVFVRSIPANAVDFYERYLIERKFFYL
GLSNEIKKGNRVDPFPIRRDQNKWKIPAMKTLGR IYSEDLVPELPRQMFNDNEIKSHLKS LPMQEGIDFNN
ANVTYLIAEYMKRVLNDDFQTFYQWKRNYR YMDMLRGEYDRKGS LQHCFTSIEERGLWKERASRTERYR
KLASNKIRSNRQMRNASSEIETILDKRLSNCRNEYQKSEKIRRYRVQDALLFLAKKILTELADFDGE
RFKLEIMPDAEKGLSIEIMPMSFIFEKGGKIYITSGGMKLNKYGDPFVLASDKRIGNLLELVGSDIVS
KEDIMEEFKKYDQCRPEISSIVFNLEKWAFTYPEL PARVDRKEKVD FWSILDLVLSNNKIDINNEQSYILR
KIRNAFDHNNYPDKGIVEIKALPEIAMS IKKAFGEYAIMK

>gi|755254328|gb|KIP62088.1| hypothetical protein ST45_06380
[*Prevotella* sp. P5-125]
MNI PALVENQKKYFGTYSVMAMLNAQT VLDHIQKVADIEG EQNENNENLWPHVMSHL YNAKNGYDKQPE
KTMFIIERLQSYFPLKMAENQREYSNGKYQNRVEVNSNDI FEVLKRAFVGLKMYRDLINHYKTYEEK
LNDGCEFLTSTEQPLSGMINNYTYVALRNMNERYGYKTEDLAFIQDKRFKFKV DAYGKKKSQVNI GSFLS
LQDYNQDTIKKLLHLSGVGIALLICFLDKQYINFLSRLPIFSSYNAQSEBERRIIRSFGINSIKLPKDR
IHSEKSNKSVAMDMLNEVKRCPDELFTLSAEKQSRFR IISDDHNEVLMKRSDDRFPVLLQYIDYGLKLF
DHIRFHVNMGLR YLLKADKICIDGQTRVRIEQPLNGFGRLEEAETMRKQENGTFGNSGIRIRDFENMK
RDDANPANYPIVETTYTHYLENNKVMEMFINDKEDSAPLLPVI EDDR YVVKTI PSCRMSTLEIPAMAFHM
FLFGSKTEKLIIDVHDYKRLRFQAMQKEEVTAENIASFGIAESDLPQKIMDLISGNAHGKD VDAFIRLT
VDDMLIDTERRIKRFDKDRKRSADNKMGRGFKQISIGKLADFLAKDIVLFPQPSVNDGENKITGLNYR
IMQSAIAVYDSGDDYEAQQQKLMFEKARLIGKITEPHFPYKVFVRSIPANAVDFYERYLIERKFFYL
GLSNEIKKGNRVDPFPIRRDQNKWKIPAMKTLGR IYSEDLVPELPRQMFNDNEIKSHLKS LPMQEGIDFNN
ANVTYLIAEYMKRVLNDDFQTFYQWKRNYR YMDMLRGEYDRKGS LQHCFTSIEERGLWKERASRTERYR
KQASNKIRSNRQMRNASSEIETILDKRLSNCRNEYQKSEKIRRYRVQDALLFLAKKILTELADFDGE
RFKLEIMPDAEKGLSIEIMPMSFIFEKGGKIYITSGGMKLNKYGDPFVLASDKRIGNLLELVGSDIVS
KEDIMEEFKKYDQCRPEISSIVFNLEKWAFTYPEL SARVDREKVD FKSILKILLNNKINKEQSDILR
KIRNAFDHNNYPDKGIVEIKALPEIAMS IKKAFGEYAIMK

- continued

>gi|755255649|gb|KIP63359.1| hypothetical protein ST44_03600
 [Prevotella sp. P5-119]
 MNI PALVENQKKYFGTYSVMAMLNAQT VLDHIQKVADI EGEQNNENNENLWFPVMSHLYNAKNGYDKQPE
 KTMPIIERLQSYFPPFLKMAENQREYSNGKYQNRVEVNSNDI FEVLKRAFVGLKMYRDLINHYKTYEEK
 LIDGCEFLTSTEQPLSGMI SKYYTVALRNI KERYGYKTEDLAFIQDNIKKI TKDAYGKRKSQVNI GFPLS
 LQDYNQDTQKLLHL SGVGIALLI CLFLDKQYINI FLSRLPIFSSYNAQSEERRII IRSFGINSI KLPKDR
 IHSEKSNKSVAMDMLNEVKRCPDELFTTLSAEKQSRFR IISDDHNEVLMKRSTDRFVPLLLQYIDYGKLF
 DHIRFHVNMGLRYYLLKADKI CIDGQTRVVI EQPLNGFGRLEEAETMRKQENGTFGNSGIRIRDFENVK
 RDDANPANYPYIVDTYTHYIL ENNKVEMFI SDKGSSAPLLPLIEDDRYVVKTI PSCRMTLEIPAMAFHM
 FLFGSKKTEKLI DVVHNRYKRLFQAMQKEEVTAENIASFGIAESDLPQKILDLISGNAHGKDVDAFIRLT
 VDDMLIDTERRI KRFKDDRKS IRSADNKMGKRGFQKQISIGKLADFLAKDIVLFPQSVNDGENKI TGLNYR
 IMQSAIAVYDSGDDYEAQKQFLMPEKARLIGKGI TEPHPFLYKVFARSIPANAVDFYERYLIERKFFYL
 GLCNEIKRGNRVDPFI RRDQNKWKIPAMKTLGR IYSEDLPVELPRQMFNDNEIKSHLKS LPQMEGIDFNN
 ANVTYLIAEYMKRVLNDDFQTFYQWKRNHYMDMLKGEYDRKGS LQHCFTSVEEREGLWKERASRTERYR
 KLASNKIRSNRQMRNASSEETILDKRLSNCRNEYQKSEKVI RRYRVDALLFLLAKKI LTELADFDGE
 RPKLKEIMPDAEKGLSEIMPSPFI FEKGGKYYTITSEGMKLNKYGDFVVLASDKRIGNLLELVGSDIVS
 KEDIMEEPNKYDQCRPEISSIVFNLEKWAFTYPELSARVDREKVDPKSILKILLNNKNINKEQSDILR
 KIRNAFDHNNYPDKGIVEIKALPEIAMS IKKAFGEYAIMK

>gi|755581756|ref|WP_042518169.1| hypothetical protein [Prevotella sp.
 P5-119]
 MNI PALVENQKKYFGTYSVMAMLNAQT VLDHIQKVADI EGEQNNENNENLWFPVMSHLYNAKNGYDKQPE
 KTMPIIERLQSYFPPFLKMAENQREYSNGKYQNRVEVNSNDI FEVLKRAFVGLKMYRDLINHYKTYEEK
 LIDGCEFLTSTEQPLSGMI SKYYTVALRNI KERYGYKTEDLAFIQDNIKKI TKDAYGKRKSQVNI GFPLS
 LQDYNQDTQKLLHL SGVGIALLI CLFLDKQYINI FLSRLPIFSSYNAQSEERRII IRSFGINSI KLPKDR
 IHSEKSNKSVAMDMLNEVKRCPDELFTTLSAEKQSRFR IISDDHNEVLMKRSTDRFVPLLLQYIDYGKLF
 DHIRFHVNMGLRYYLLKADKI CIDGQTRVVI EQPLNGFGRLEEAETMRKQENGTFGNSGIRIRDFENVK
 RDDANPANYPYIVDTYTHYIL ENNKVEMFI SDKGSSAPLLPLIEDDRYVVKTI PSCRMTLEIPAMAFHM
 FLFGSKKTEKLI DVVHNRYKRLFQAMQKEEVTAENIASFGIAESDLPQKILDLISGNAHGKDVDAFIRLT
 VDDMLIDTERRI KRFKDDRKS IRSADNKMGKRGFQKQISIGKLADFLAKDIVLFPQSVNDGENKI TGLNYR
 IMQSAIAVYDSGDDYEAQKQFLMPEKARLIGKGI TEPHPFLYKVFARSIPANAVDFYERYLIERKFFYL
 GLCNEIKRGNRVDPFI RRDQNKWKIPAMKTLGR IYSEDLPVELPRQMFNDNEIKSHLKS LPQMEGIDFNN
 ANVTYLIAEYMKRVLNDDFQTFYQWKRNHYMDMLKGEYDRKGS LQHCFTSVEEREGLWKERASRTERYR
 KLASNKIRSNRQMRNASSEETILDKRLSNCRNEYQKSEKVI RRYRVDALLFLLAKKI LTELADFDGE
 RPKLKEIMPDAEKGLSEIMPSPFI FEKGGKYYTITSEGMKLNKYGDFVVLASDKRIGNLLELVGSDIVS
 KEDIMEEPNKYDQCRPEISSIVFNLEKWAFTYPELSARVDREKVDPKSILKILLNNKNINKEQSDILR
 KIRNAFDHNNYPDKGIVEIKALPEIAMS IKKAFGEYAIMK

>gi|762997475|ref|WP_043894148.1| hypothetical protein [Porphyromonas
 gingivalis]
 MTEQNEKPYNGTYTLEDKHFWAFFNLARHNAYITLTHIDRQLAYS KADITNDEIDLFFKGQWKNLNDND
 LERKARLRS LILKHFSLFLEGAAYGKKL FESQSSGNKSSKKKELIKKEKEELQANALS LDNLKSI LDFDLQ
 KLKDFRNYSHYRHPSESELPFDGNMLQRLYVNFVDSVQVRKRDHEHNDKVDPHRHFNHLVRKGGKDR
 GNNDNPFKHFVDRREGVTEAGLLFFVSLFLEKRD AIWMQKIRGFGKGTETYQMTNEVFCRSRISLP
 KLKLESLRIDDWMLDMLNELVRCPKSLYDRLEEDRACFRVVDILSDEDDIDGAEEDPFKNI LVRHQD
 RFPYFALRYFDLKKVFTSLRPHIDLGTYHFAIYKKNIGEQPEDRHLTRNLYGFGRIQDFAEHRPEEWR
 LVRLDLCFETGDKPYITQTPHYHI EKGKILGRFVPEGQHLWSPPEVGATRTRGSKYAQDKRLTAEAFSL
 VHELMPMFYFYLREKYSDEASAERVQGR IKRVEDVYAVYDAFARGEIDILDRDLACLADKGI RRGHL
 PRQI IAILSQEHKDMEEKVRKQLQEMIADTDHRLDMLDRQTRDKIRI GRKNAGLPKSGVIADWLVRDMMR
 FQPVAKDTSGKPLNNSKANSTEYRMLQRALALFGGEKERLIPYFRQMNIGGNPHFPFHETRWESHINI
 LSFYRSYLKARKAFLQSIGRSDRVENHRFLLKEPKTDRQTLVAGWKSEFHLPRGIFTEAVRDCLIEMGH
 DEVASYKEVGFMAKAVPLYFERACKDRVQPFYDYPFNVGNSLKPCKGRFSLKEKRAEWEESGKERFRDLE
 AWHSAAARRIEDAFVGI EYASWENKKIEQLQDLSLWETPESKLKVKADKINIAKLKKEILEAKEHPYH
 DFKSWQKFERELRLVKNQDIITWMMCRDLMEENKVEGLDITGLYLKDIRTDVHEQGS LNVLNVRKPMRLP
 VVVYRADS RGHVHKEEAPLATVYIEERDTKLLKQGNFKSFVKDRRLNGLFSFVDTGGLAMEQYPI SKLRV
 EYELAKYQTARVCAFEQTLLEESLLTRYPHLPDKNFRKMLSESDP LLDKWPDLHRVRLLI AVRNAPS
 HNQYPMYDEAVFSSIRKYDPSLDAIEERMGLNIAHRLSEEVKQAKETVERIIQA

>gi|762999449|ref|WP_043894372.1| hypothetical protein [Porphyromonas
 gingivalis]
 MTEQNERPYNGTYTLEDKHFWAFFNLARHNAYITLTHIDRQLAYS KADITNDEIDLFFKGQWKNLNDND
 LERKARLRS LILKHFSLFLEGAAYGKKL FESQSSGNKSSKKKELIKKEKEELQANALS LDNLKSI LDFDLQ
 KLKDFRNYSHYRHPSESELPFDGNMLQRLYVNFVDSVQVRKRDHEHNDKVDPHRHFNHLVRKGGKDR
 GNNDNPFKHFVDRREGVTEAGLLFFVSLFLEKRD AIWMQKIRGFGKGTETYQMTNEVFCRSRISLP
 KLKLESLRIDDWMLDMLNELVRCPKSLYDRLEEDRACFRVVDILSDEDDIDGAEEDPFKNI LVRHQD
 RFPYFALRYFDLKKVFTSLRPHIDLGTYHFAIYKKNIGEQPEDRHLTRNLYGFGRIQDFAEHRPEEWR
 LVRLDLCFETGDKPYITQTPHYHI EKGKILGRFVPEGQHLWSPPEVGATRTRGSKYAQDKRLTAEAFSL
 VHELMPMFYFYLREKYSDEASAERVQGR IKRVEDVYAVYDAFARGEIDILDRDLACLADKGI RRGHL
 PRQMI IAILSQEHKDMEEKVRKQLQEMIADTDHRLDMLDRQTRDKIRI GRKNAGLPKSGVIADWLVRDMMR
 FQPVAKDTSGKPLNNSKANSTEYRMLQRALALFGGEKERLIPYFRQMNIGGNPHFPFHETRWESHINI
 LSFYRSYLKARKAFLQSIGRSDRVENHRFLLKEPKTDRQTLVAGWKSEFHLPRGIFTEAVRDCLIEMGH
 DEVGYSKEVGFMAKAVPLYFERACKDRVQPFYDYPFNVGNSLKPCKGRFSLKEDRAEWEESGKERFR
 LKKEILEAKEHPYHDFKSWQKFERELRLVKNQDIITWMMCRDLMEENKVEGLDITGLYLKDIRTDVYEQG
 SLNVLNVRKPMRLPVVYRADS RGHVHKEEAPLATVYIEERDTKLLKQGNFKSFVKDRRLNGLFSFVDTG
 GLAMEQYPI SKLRV EYELAKYQTARVCAFEQTLLEESLLTRYPHLPDESFRMLESWSDP LLDKWPDLH
 RKVRLLI AVRNAPS HNQYPMYDEAVFSSIRKYDPSLDAIEERMGLNIAHRLSEEVKQAKEMAEIRIIQA

- continued

>gi|763000465|ref|WP_043894537.1| hypothetical protein [*Porphyromonas gingivalis*]
MTEQNERPYNGTYTLEDKHFWAFLNLARHNAYITLTHIDRQLAYSKADITNDEDILFFKGQWKNLDND
LERKARLRSLILKHFSPFLGGAAYGKCLFESQSSGNKSSKKKELIKKEKEELQANALSLDNLKSILDFDLQ
KLKDFRNYSHYRHPESSELPDFDGNMLQRLYNVFDVSVQVRKRDHEHNDKVDPHRHFNLVRKGGKDRY
GNNDNPFKHHFVDREKVTTEAGLLFFVSLFLEKRDIAWMQKIRGFKGGTETYQOMTNEVFCR.SRISLP
KLKLESLRTDDWMLLDMLNELVRCPSLYDRLEEDRARFRVPIDILSEDDTDGAEEDPPKNTLVRHQD
RFPYFALRYFDLKKVFTSLRPHIDLGTYHFAIYKKNIGEOPEDRHLTRNLYGFGRIQDFAEEHRPEEWKR
LVRDLDFETGDKPYITQTPHYHIKGGKILGRFVPEGQHLWSPPEVGATRTRGRSKYAQDKRRLTAEAFSL
VHELMPMMFYFLLREKYSDEASAERVQGRIKRVIDVYAVYDAFARGEIDTLDRDLACLADKGIIRRHGL
PRQMIAILSQEHKDMEEKVRKLLQEMIDTDHRLDMLDRQTDKIRIRGRKNAGLPKSGVIADWLVRDMMR
FQPVAKDTSGKPLNNSKANSTEYRMLQRALALFGGEKERLTPYFRQMNLTGGNNPFPFLHETRWESHTNI
LSFYRSYLKARKAFLOSIGRSDRENNHRFLLKPKTDRTLVAGWKGEPHLPRIIFTEAVRDCLIMGL
DEVGSYKEVGFMAKAVPLYFERACKDRVQPPYNYPFNVGNILKPKKGRFLSKEKRAEWEESGKERFRDLE
AWSQSTARREDAFAGIKYASPGNKKIEQLLQDLSLWETPESKLVKADKINLAKLKEILEAKEHPYL
DFKSWQKFERELRLVKNQDIITWMMCRDLMEENKVEGLDTGTLYLKDIRTDVQEQGSLNVLNHRVPMRLP
VVVYRADSRRGHVHKEEAPLATVYIEERDTKLLKQGNFKSFVKDRRLNGLFSFVDTGALAMEQYPIKSLRV
EYELAKYQATARVCAFEQTLLEESLLTRYPHLPDKNFRKMLLESWSDDLKWPDLHRKVRLLIAVRNAPS
HNQYPMYDEAVFSSIRKYDPSPPDAIEERMGLNIAHRLSEEVKQAKETVERIIQA

>gi|763001609|ref|WP_043894751.1| hypothetical protein [*Porphyromonas gingivalis*]
MTEQNERPYNGTYTLEDKHFWAFLNLARHNAYITLAHIDRQLAYSKADITNDEDILFFKGQWKNLDND
LERKARLRSLILKHFSPFLGGAAYGKCLFESQSSGNKSSKKKELIKKEKEELQANALSLDNLKSILDFDLQ
KLKDFRNYSHYRHPESSELPDFDGNMLQRLYNVFDVSVQVRKRDHEHNDKVDPHRHFNLVRKGGKDRY
GNNDNPFKHHFVDREKVTTEAGLLFFVSLFLEKRDIAWMQKIRGFKGGTETYQOMTNEVFCR.SRISLP
KLKLESLRTDDWMLLDMLNELVRCPSLYDRLEEDRARFRVPIDILSEDDTDGTEEDPPKNTLVRHQD
RFPYFALRYFDLKKVFTSLRPHIDLGTYHFAIYKKNIGEOPEDRHLTRNLYGFGRIQDFAEEHRPEEWKR
LVRDLDFETGDKPYITQTPHYHIKGGKILGRFVPEGQHLWSPPEVGATRTRGRSKYAQDKRFTAEAFSL
VHELMPMMFYFLLRENYSDASAERVQGRIKRVIDVYAVYDAFARGEIDTLDRDLACLADKGIIRRHGL
PRQMIAILSQEHKDMEEKVRKLLQEMIDTDHRLDMLDRQTDKIRIRGRKNAGLPKSGVIADWLVRDMMR
FQPVAKDTSGKPLNNSKANSTEYRMLQRALALFGGEKERLTPYFRQMNLTGGNNPFPFLHETRWESHTNI
LSFYRSYLKARKAFLOSIGRSDRVENHRFLLKPKTDRTLVAGWKGEPHLPRIIFTEAARDCLIMGH
DEVASYKEVGFMAKAVPLYFERACKDRVQPPYDYPFNVGNLSLKKKGRFLSKEKRAEWEESGKERFRDLK
LKEILEAKEHPYLDKFSWQKFERELRLVKNQDIITWMMCRDLMEENKVEGLDTGTLYLKDIRTDVQEQG
SLNVLNHRVPMRLPVVVYRADSRRGHVHKEQAPLATVYIEERDTKLLKQGNFKSFVKDRRLNGLFSFVDTG
ALAMEQYPIKSLRVEYELAKYQATARVCAFEQTLLEESLLTRYPHLPDKSFRMLESWSDDLKWPDLH
RKVRLLIAVRNAPSNNQYPMYDEAVFSSIRKYDPSLDAIEERMGLNIAHRLSEEVKQAKEMVERIIQV

>gi|763003511|ref|WP_043895240.1| hypothetical protein [*Porphyromonas gingivalis*]
MTEQNEKPYNGTYTLEDKHFWAFLNLARHNAYITLTHIDRQLAYSKADITNDEDILFFKGQWKNLDND
LERKARLRSLILKHFSPFLGGAAYGKCLFESQSSGNKSSKKKELIKKEKEELQANALSLDNLKSILDFDLQ
KLKDFRNYSHYRHPESSELPDFDGNMLQRLYNVFDVSVQVRKRDHEHNDKVDTHRHFNLVRKGGKDRY
GNNDNPFKHHFVDREKVTTEAGLLFFVSLFLEKRDIAWMQKIRGFKGGTETYQOMTNEVFCR.SRISLP
KLKLESLRTDDWMLLDMLNELVRCPSLYDRLEEDRARFRVPIDILSEDDTDGAEEDPPKNTLVRHQD
RFPYFALRYFDLKKVFTSLRPHIDLGTYHFAIYKKNIGEOPEDRHLTRNLYGFGRIQDFAEEHRPEEWKR
LVRDLDFETGDKPYITQTPHYHIKGGKILGRFVPEGQHLWSPPEVGATRTRGRSKYAQDKRRLTAEAFSL
VHELMPMMFYFLLREKYSDEASAERVQGRIKRVIDVYAVYDAFARGEIDTLDRDLACLADKGIIRRHGL
PRQMIAILSQEHKDMEEKVRKLLQEMIDTDHRLDMLDRQTDKIRIRGRKNAGLPKSGVIADWLVRDMMR
FQPVAKDTSGKPLNNSKANSTEYRMLQRALALFGGEKERLTPYFRQMNLTGGNNPFPFLHETRWESHTNI
LSFYRSYLKARKAFLOSIGRSDRVENHRFLLKPKTDRTLVAGWKGEPHLPRIIFTEAVRDCLIMGH
DEVASYKEVGFMAKAVPLYFERACKDRVQPPYDYPFNVGNLSLKKKGRFLSKEKRAEWEESGKERFRDLE
AWSQSAVRRIDAFAGIENASRENKKEIEQLLQDLSLWKAPEKLVKRGDKINIAKLLKEILEAKEHPYH
DFKSWQKFERELRLVKNQDIITWMMCRDLMEENKVEGLDTGTLYLKDIRTDVHEQGS LNVLNHRVPMRLP
VVVYRADSRRGHVHKEQAPLATVYIEERDTKLLKQGNFKSFVKDRRLNGLFSFVDTGALAMEQYPIKSLRV
EYELAKYQATARVCAFEQTLLEESLLTRCPHLPDKNFRKMLLESWSDDLKWPDLHRKVRLLIAVRNAPS
HNQYPMYDEAVFSSIRKYDPSPPDAIEERMGLNIAHRLSEEVKQAKETVERIIQA

>gi|763006464|ref|WP_043895933.1| hypothetical protein [*Porphyromonas gingivalis*]
MTEQNERPYNGTYTLEDKHFWAFLNLARHNAYITLTHIDRQLAYSKADITNDEDILFFKGQWKNLDND
LERKARLRSLILKHFSPFLGGAAYGKCLFESQSSGNKSSKKKELIKKEKEELQANALSLDNLKSILDFDLQ
KLKDFRNYSHYRHPESSELPDFDGNMLQRLYNVFDVSVQVRKRDHEHNDKVDPHRHFNLVRKGGKDKY
GNNDNPFKHHFVDREKVTTEAGLLFFVSLFLEKRDIAWMQKIRGFKGGTETYQOMTNEVFCR.SRISLP
KLKLESLRIDWMLLDMLNELVRCPSLYDRLEEDRACFRVPVDILSEDDIDGAEEDPPKNI LVRHQD
RFPYFALRYFDLKKVFTSLRPHIDLGTYHFAIYKKNIGEOPEDRHLTRNLYGFGRIQDFAEEHRPEEWKR
LVRDLDFETGDKPYITQTPHYHIKGGKILGRFVPEGQHLWSPPEVGATRTRGRSKYAQDKRRLTAEAFSL
VHELMPMMFYFLLREKYSDEASAERVQGRIKRVIDVYAVYDAFARGEIDILDRDLACLADKGIIRRHGL
PRQI IAILSQEHKDMEEKVRKLLQEMIDTDHRLDMLDRQTDKIRIRDRKNAGLPKSGVIADWLVRDMMR
FQPVAKDTSGKPLNNSKANSTEYRMLQRALALFGGEKERLIPYFRQMNLTGGNNPFPFLHETRWESHINI
LSFYRSYLKARKAFLOSIGRSDRVENHRFLLKPKTDRTLVAGWKGEPHLPRIIFTEAVRDCLIMGL
DEVGSYKEVGFMAKAVPLYFERACKDRVQPPYNYPFNVGNILKPKKGRFLSKEKRAEWEESGKERFRDLE

- continued

AWSQSTARRIEDAFAGIKYASPGNKKKIEQLLQDLSLWETPESKLVKADKINLAKLKEILEAKEHPYL
DFKSWQKFERELRLVKNQDIIITWMMCRDLMEENKVEGLDITGLYLDKIWTDVHEQSGSLNVLNRVKPMRLP
VVVYRADS RGHVHKQAPLATVYIERDTRKLLKQGNFKSFVKDRRLNGLFSFVDTGGLAMEQYPI SKLRV
EYELAKYQATARVCAPEQTELEESLTRCPHLDPKNFRKMLESWSDPLLDKWPDLHRKVRLLIAVRNAFS
HNQYPMYDEAVFSSIRKYDPSPDIAEERMGLNIAHRLSEEVKQAKEMAERIIQA

>gi|763187618|ref|WP_044065294.1| hypothetical protein [Prevotella sp.
P5-125]

MNIPALVENQKKYFGTYSVMAMLNAQTVDLHIQKVADI EGEQNENNENLWPHVMSHLYNAKNGYDKQPE
KTMFIIERLQSYFFFLKMAENQREYSNGKYQNRVEVNSNDI FEVLKRAFGLKMYRDLINHYKTYEEK
LNDGCEFLTSTEQPLSGMINNYTVALRNMNERYGKTEDLAFIQDKRFKFKDAYGKKKSQVNI GPFPLS
LQDYNQDTPKLLHL SGVGIALLICLFLDKQYINIFLSRLPIFSSYNAQSEERRIIIRSFINGINSIKLPKDR
IHSEKSNKSVAMDMLNEVKRCPDELFTTLSAEKQSRFRISDDHNEVLMKRSSDRFVPLLLQYIDYKGLF
DHIRFHVNMGLRLLKADKICIDGQTRVRIEQPLNGFGRLEEAETMRKQENGTFGNSGIRIRDFENMK
RDDANPANYPYIVDYTHYIL ENNKVEMFINDKEDSAPLLPVI EDDR YVVKTI PSCRMSTLEIPAMAFHM
FLFGSKKTEKLI DVHNRKRFQAMQKKEEVTAEINASFGEAEDLPQKILDLISGNAHGKDVDAFIRLT
VDDMLIDTERRIKRFDKDRKIRSADNKMGRGFKQISIGKLADFLAKDIVLFPQPSVNDGENKITGLNRY
IMQSAIAVYDSGDDYEAQKQFKLMEKARLIGKGI TEPHPFLYKVFARSIPANAVEFYERYLIERKFYLT
GLSNEIKKGNRVDPFIRRDQNKWKIPAMKTLGRISYEDLPVELPRQMFNDNEIKSHLKSLPQMEGIDFNN
ANVTYLIAEYMKRVLDDDFQTFYQWNRNRYMDMLRGEYDRKGLQHCFTSVEEREGLWKERASRTERYR
KQASNKIRSNRQMRNASSEETILDKRLSNRNEYQKSEKVI RRYRVDALLFLLAKKITELADFDGE
RPFKLEIMPDAEKILSEIMPMSFIFEKGGKYYTITSEGMKLNKYGDFVVLASDKRIGNLLELVGSDIVS
KEDIMEEFNKYDQCRPEISSIVFNLEKWAFTYPELSARVDREKVDPKSILKILLNNKININKEQSDILR
KIRNAFDHNNYPDKGVVEIKALPEIAMS IKKAFGEYAIMK

>gi|763200485|ref|WP_044074780.1| hypothetical protein [Prevotella sp.
P5-60]

MNIPALVENQKKYFGTYSVMAMLNAQTVDLHIQKVADI EGEQNENNENLWPHVMSHLYNAKNGYDKQPE
KTMFIIERLQSYFFFLKMAENQREYSNGKYQNRVEVNSNDI FEVLKRAFGLKMYRDLINHYKTYEEK
LIDGCEFLTSTEQPFSGMISKYYTVALRNIKERYGYKAEDLAFIQDNRYKPKTDAYGKRKSQVNI GSFPLS
LQDYNQDTIKLLHL SGVGIALLICLFLDKQYINIFLSRLPIFSSYNAQSEERRIIIRSFINGINSIKQPKDR
IHSEKSNKSVAMDMLNEVKRCPDELFTTLSAEKQSRFRISDDHNEVLMKRSSDRFVPLLLQYIDYKGLF
DHIRFHVNMGLRLLKADKICIDGQTRVRIEQPLNGFGRLEEVETMRKQENGTFGNSGIRIRDFENMK
RDDANPANYPYIVDYTHYIL ENNKVEMFISDEENPTLPLPVI EDDR YVVKTI PSCRMSTLEIPAMAFHM
FLFGSEKTEKLI DVHNRKRFQAMQKKEEVTAEINASFGEAEDLPQKIMDLISGNAHGKDVDAFIRLT
VDDMLIDTERRIKRFDKDRKIRSADNKMGRGFKQISIGKLADFLAKDIVLFPQPSVNDGENKITGLNRY
IMQSAIAVYDSGDDYEAQKQFKLMEKARLIGKGI TEPHPFLYKVFVRSIPANAVDFYERYLIERKFYLI
GLSNEIKKGNRVDPFIRRDQNKWKIPAMKTLGRISYEDLPVELPRQMFNDNEIKSHLKSLPQMEGIDFNN
ANVTYLIAEYMKRVLDDDFQTFYQWNRNRYMDMLRGEYDRKGLQHCFTSVEEREGLWKERASRTERYR
KLASNKIRSNRQMRNASSEETILDKRLSNRNEYQKSEKIRRYRVDALLFLLAKKITELADFDGE
RPFKLEIMPDAEKILSEIMPMSFIFEKGGKYYTITSEGMKLNKYGDFVVLASDKRIGNLLELVGSDIVS
KEDIMEEFNKYDQCRPEISSIVFNLEKWAFTYPEL PARVDRKEKVD PWSILDLVLSNNKIDINNEQSYILR
KIRNAFDHNNYPDKGVVEIKALPEIAMS IKKAFGEYAIMK

>gi|763360704|ref|WP_044218239.1| hypothetical protein

[*Phaeodactylibacter xiamenensis*]

MINTPKRRILHRHPSYFGAFLNIRHNAFMIMEHLSTKYDMDKINILDEAQLPNAKLFGLCKKRYGKPDV
TEGVSRDLRRYFPFLNYPFLHLEKQQAQAATYDINPEDI EFLKGFRLNQMRRNYSHYISNIDYG
KFDKLPVQDIYEAIFRLLDRGKHTKRFVDFESKHIRHLESNNSEYRPRSLANSPDHENTVAFVTCFLFE
RKYAFPFLSRDLCPFRSINDAEGDPLIRKASHECYTMPCCRLPQPKLESSDILLDMVNLGRCP SALYNL
LSEEDQARFHIKREBITGFEEDPDELEQEI VLRKHSRDFYPYFALRYPDDTEAFQTLRFDVYLGWRIRK
VYKRIYQERDRVLIQSIRTFIRLSRLPIYENVKHDAVRQNEEDGKLVNPDVISQPHKSWIQIESDR
APLSDRIEHSFYFNGQVIGLKF INPDYAAIQNVFPKLPGEKDKDAKLVNETADAIISTHEIRSL
FLYHYLSKKPISAGDERRFIQVDTEFTIKQYIDTIKLF FEDIKSGELQIADPPNYQKNEPLPYVGRDKE
KTQBERAQRERQKEIKERRKELNILLQNRVGLSIQYIPSRLELYLLQYKVPYEKLLALQKLAQRKEVK
KRIKDI EKMRTPRVGEQATWLAEDI VFLIPPKMHI PERKTIKHPQKLNNDQFRIMQSSLA YFVSNKKAIK
KFFQKETGIGLSNRETSHPFYRIVDVGRCRGLDFYIGYLYKMDWLDLDAIKKVDNRKHGKKEAKKYEKY
LPSSIQHKTPELDDYIRLVPVYLRPLGFKKAI VLAALAHADPQVEPEEDNVI FCLDQLLDGDTQDFYNWQR
YRSALTEKEIDNQLVLAHPYAEQILGTIKTELGKQKNNKLGKAKQKIKDELIDLKRAKRRLLDREQYL
RAVQAEDRALWLMIQERQKQAEHEEIAFDQLDLKNITKILTESIDARLRI PDTKVDITDCLPLRRYGD
RRVAKDRRLVNLASYHVAGLSEIPYDLVKKLELEYDRRRVAFHEVYQFEKEVYDRYAEALRNENPKGE
STYF SHWEYVAVAKHSADTHENELPFKEKVMQLRNLKPHHNEFPYFDWLLPEVEKASAAALYADRVFDVAEG
YYQMRKMLRQ

>gi|769827477|gb|KJ86756.1| hypothetical protein M573_117042

[*Prevotella intermedia* ZT]

MKMEDDKKITES TNLMDNKHFAAFNLARHNVIYITVNHNKVLLELKNKQDIIIDNDQDILAIKTHWE
KVGDLNKTERRLEMLTKHPPLETAIYIKNKEDKEEVKQEKQAEASLESKDCFLFLEKQEARNY
SHYKSESTKEPMLBEGLEKMYNI FDDNIQLVIKDYQHNDINPDEDFKHLDRKGQFKYFADNEGNI
ESGLFFVSLFLEKDAIWMQKLI GFKDNRESKMMTHEVFCRRRMLLPKRLLESTQTQDWILLDMLNE
LIRCPKSLYERLQGEYRKENVYFDSADEYDAEQEPFKNILVRHQDRFPYFALRYFDYNEITFNLRQI
DLGTYHFSIYKLLIGGQKEDRHLTHKLYGFERIQEFQKQNRPEWKAIVKLDLIIYETSNERYIS ETPHY
HLENQKIGIRFRNGNKEIWPSLKINGENNEKSKYKLDKPYQAEAPLSVHELPPMFFYLLLLKKEEPNDK
KNASIVEGFIKREIRDMYKLVDFAFANGEINNI GDLEKYCEDKGI PKRHLPKQMVAILYDEPKDMVKEAKR
KQKEMVKDIKLLLATLEKQTEEEIEDGGRNIRLLKSGEITARWLNDMMRFPQVQKDNENGNPLNNSKANS
EYQMLQRSALYNKEEKPTRYQVNLINSSNPHFPLKWI KWEECNILSFYRNYLIKKEIFLNKLPED
WEKNQYFLKLEKPKTNRETLFGQWKNGFNLPRGIFTEPIREWFKRHQND SKYEKVEALKRVGLVTKVIP
LFFKBEYFKEDAQKEINNCVQPFYSPYVNGNIHKPDEKDFLPSERKKLWGDKKDKFKGYKAKVKSKLL

- continued

TDKEKEEYRSYLEFQSWNKFERELRLVRNQDITVWLLCTELIDMKKVEGLNVEELQKLRKLDIDTDTAKQ
 EKNNILNRIIMPQLPVIYVEIIDDSSHNI VKDRPLHIVYI EETKI KLLKQGNFKALVKDRRLNGLFSFVDT S
 SKAELKDKPI SKSVVEYELGEYQNARI ETI KDMLLLEKTLIKKYEKLPI DNFSMDLNGWLEGGKDESKAR
 FQNDVKLLVAVRNAFSHNQYPMRNR IAFANINPFLSSADISEEKKLIDANQLKDKIHKI IKKI IETIEKP
 IETKE

>gi|806787645|gb|KKC50278.1| hypothetical protein HR1010685

[*Porphyromonas gulae*]

MTEQSERPYNGTYTLEDKHFWAFLNLRHNAVITLTHIDRQLAYSKADI TNDQDVLSPFKALWKNPND
 LERKSRLRSLILKHFSPLEGAAYGKFLPESKSSGNKSSKNKELIKKEKEELQANALSLDNLKSLIFDFLQ
 KLKDFRNYSHYRHSSESELPDFGNMLQRLYNVFDVSVQVRKRDHEHNDKVDPHRHFNHLVRKGGKDRY
 GHNDNPSFKHHPVDESEGMVTEAGLFFVSLFLEKRDALWQKIRGFKGGTETYQOMTNEVFCRSRISLP
 KLKLESLRIDDWMLDMLNELVRCPKPLYDRLEKDRARFRVPVDILPDEDDIDGGGEDPFFKNI LVRHQD
 RFPYFALRYFDLKKVFTSLRPHIDLGTYHFAIYKMI GEQPEDRHLTRNLYGFGRIQDFAEEHRPEEWR
 LVRDLDFETGDKPIYSQTPHYHI EKGI GLRFMPEGQHLWSPPEVGI TRTGRSKYAQDKRLTAEAFLS
 VHELMPMFYFLLREKYS EEVSAEKVQGR I KRVI EDVYAIYDAPARDEINTLKELDACLADKGI RRGHL
 PKQMIAILSQEHKDMEEKIRKLLQEMIADTDHRLDMLDRQTDRIIRGRKNAGLPKSGVIADWLVRDMMR
 FQPVAKDTSGKPLNNSKANSTYRMLQRALALFGGEEKRLIPYFRQMNILGNNPHPLHETRWESHINI
 LSFYRSYLARKAFLEIRGRSDRMENRPFLLKPEKPTDRQTLVAGWKSEFHLPRGIFTEAVRDCLIEMGY
 DEVGSYREVGFMKAVPLYFERACEDRVQPFYDS PFNVGNSLKKKGRFLSKEERAEEWERGKERFRDLE
 AWSSAARRIEDAFAGI EYASPGNKKIEQLRDLDSLWEAFESKLKVRADKINLAKLKEILEAQEHPYH
 DFKSWQKFERELRLVKNQDII TWMMCRDLMEENKVEGLDIGTLYLKD IRPNVQEQGSLNVLNRVKPMLRP
 VVVYRADSRGHVHKEEAPLATVYIEERDTKLLKQGNFKSFVKDRRLNGLFSFVDTGGLAMEQYPI SKLRV
 EYELAKYQTA RVCFELTLRL EESLLTRYPHLPDESFRKMLESWSDPLLAKWPELHGKVRLLIAVRNAPS
 HNQYPMYDEAVFSSIRKYDPSSPDAIEERMGLNIAHRLSEEVKQAKETVERIIQV

>gi|807048513|ref|WP_046201018.1| hypothetical protein [*Porphyromonas*

gulae]

MTEQSERPYNGTYTLEDKHFWAFLNLRHNAVITLTHIDRQLAYSKADI TNDQDVLSPFKALWKNPND
 LERKSRLRSLILKHFSPLEGAAYGKFLPESKSSGNKSSKNKELIKKEKEELQANALSLDNLKSLIFDFLQ
 KLKDFRNYSHYRHSSESELPDFGNMLQRLYNVFDVSVQVRKRDHEHNDKVDPHRHFNHLVRKGGKDRY
 GHNDNPSFKHHPVDESEGMVTEAGLFFVSLFLEKRDALWQKIRGFKGGTETYQOMTNEVFCRSRISLP
 KLKLESLRIDDWMLDMLNELVRCPKPLYDRLEKDRARFRVPVDILPDEDDIDGGGEDPFFKNI LVRHQD
 RFPYFALRYFDLKKVFTSLRPHIDLGTYHFAIYKMI GEQPEDRHLTRNLYGFGRIQDFAEEHRPEEWR
 LVRDLDFETGDKPIYSQTPHYHI EKGI GLRFMPEGQHLWSPPEVGI TRTGRSKYAQDKRLTAEAFLS
 VHELMPMFYFLLREKYS EEVSAEKVQGR I KRVI EDVYAIYDAPARDEINTLKELDACLADKGI RRGHL
 PKQMIAILSQEHKDMEEKIRKLLQEMIADTDHRLDMLDRQTDRIIRGRKNAGLPKSGVIADWLVRDMMR
 FQPVAKDTSGKPLNNSKANSTYRMLQRALALFGGEEKRLIPYFRQMNILGNNPHPLHETRWESHINI
 LSFYRSYLARKAFLEIRGRSDRMENRPFLLKPEKPTDRQTLVAGWKSEFHLPRGIFTEAVRDCLIEMGY
 DEVGSYREVGFMKAVPLYFERACEDRVQPFYDS PFNVGNSLKKKGRFLSKEERAEEWERGKERFRDLE
 AWSSAARRIEDAFAGI EYASPGNKKIEQLRDLDSLWEAFESKLKVRADKINLAKLKEILEAQEHPYH
 DFKSWQKFERELRLVKNQDII TWMMCRDLMEENKVEGLDIGTLYLKD IRPNVQEQGSLNVLNRVKPMLRP
 VVVYRADSRGHVHKEEAPLATVYIEERDTKLLKQGNFKSFVKDRRLNGLFSFVDTGGLAMEQYPI SKLRV
 EYELAKYQTA RVCFELTLRL EESLLTRYPHLPDESFRKMLESWSDPLLAKWPELHGKVRLLIAVRNAPS
 HNQYPMYDEAVFSSIRKYDPSSPDAIEERMGLNIAHRLSEEVKQAKETVERIIQV

>gi|835310921|ref|WP_047447901.1| hypothetical protein [*Alistipes* sp.

ZOR0009]

MSNEIGAPREHQFAYAPGNEKQEEATFATYFNLALS NVGEMMFGEVESNPDKIEKSLDTPPAILRQIAS
 FIWLSKEDHPDKAYSTEEVKVIVIDLVRRLCFYRNYFSHC FYLDTPYFSDLEIETGKLYNPHHF
 ITNRLFRYSLPEITLFRWNEGERKYEILRDGLIFPCCLFLKRGQAERFLNELRFPFKRIDEEGRIKRTI FT
 KYCTRESHKHIGIEEQDLIFQDIIGDLNRVPKCDGVVDLSKENERYIKNRETSNESDENKARVRLIR
 EKDKFPYLMRYIVDFGLPCITFPQNDSYTKEGRGQPHYQDAVAQEERCYNFVVRNNGVYYSYMPQAQ
 NVVRISELQGTISVEELRNMVYASINGKDVNKSVEQYLYHLHLYEKILITISGQIKRGRVDEYRPLD
 DKLLLRPANSGEELRRELRLKLPKRVCDLLSNRFDCEGSAVEKRLKATILRHEQLLSQNPALHIDKI
 KSVIDYLYLFFSDDEKFRQQPIEKAHRLGKDEEFQMYHYLVGDYD SHPLALWKELEASGRKPKEMRKLTS
 ATSLHGLYMLCKGIVBWRKQLMSIGKGTAKVEAIDRVGLKLYDKLKEYTPEQLEREVLVVMHGYAA
 AATPKPKAAQAI PSKLTLELRFYSLGKREMSFAAFIRQDKKAQKLWLRNFYIVENIKTLQKRAAADAAC
 KKLNLVGEVERVHINDKVLVLAQRVRELLNVGSKCAVILDNPERQQKLADVYEVQNAWLSIRFDLDD
 FILLTHVNLNLRKAYNLI PRKHI LAFKEYLDNRVQKLCCECRNVRKEDLCICCSPRYSNLT SWLKENH
 SSSI ERAATMMLLDVERKLLSFLDERRKAIIEYKGFIPFSALVKECRADAGLCGRINDVLHNDVVIS
 YADAIKLSAYFPKEASAVAVYIRRTKEVRBQRREELMANSSQ

>gi|874269987|gb|AKP69887.1| hypothetical protein CG08_1741

[*Riemerella anatipestifer*]

MEKPLLPNVYTLKHKFFWGAFLNLRHNAVITLTHIDRQLAYSKADI TNDQDVLSPFKALWKNPND
 KKSQTLTELILKHFPLTAMCYHPPKKEGKKGHQKEQQKEKESAEQSAEALNPSKLEALEILVNQLHS
 LRNYSHYKHKKPDADAKDI PKHLKAFDASLRMVKEDYKAHPIVNLIRDFALHNRKGNKQNDPNRYR
 FEKDGFTESGLLFFINLFLDKRDAYWMLKKVSGFKASHKQREKMITVEVFCRSRILLPKRLRESYDHNQ
 MLLDMLSELSRCPKLLYEKLS EENKHFQVEADGFLDEIEEENQPFKDTLIRHQDRFPYFALRYLDLNE
 FKSIRFQVDLGTYHYCIYDKIIGDEQEKRLHTRILLSPGRLOQDFTIENRQEWKALTKDLDYKETSNOFF
 ISKTI PHYHI TDNKIGFRLGTSKELYPSEIKDGANRI AKYFYNSSGFVAHAFAI SVHELPLPMFYQHLIGK
 SEDLLKETVRHIQR IYKDFEERINTI EDLEKANQGRPLGAPFKQMLGQLLQNKQPDLSKAKIKIEKLI
 AETKLLSHRLNLIKSSPKLGGKREKLIKTVGLADWLKDFMRPQPVAYDAQNP I KSKANSTEFWFIR
 RALALYGGKRNLEGYFKQINLIGM INPHPLNKNFNWACRNLVDFYQOYLEQRKFL EAIKNQWPWEYQ
 YCLLLKIPKENRKNLVKGEWQGGISLPRGLFTEAIRETLESDMLSKPIRKEIKHGRVGFISRAITLYF
 KEKYQDKHQS FYNLSYKLEAKAPLLKREBEHYEYQQNKQPSPTESQRLLEHTSDRWKDYLLYKRWQHLK
 KLRLYRNQDVMMLWMLTLELTKNHFKELNLYHQLKLENLAVNVQEADAKLNLPLNQLTPMVLVVKVYPATA

- continued

FGEVQYHKTPIRTVYIREEHTKALKMGNFKALVKDRRLNGLFSFIKEENDTQKHPISQLRRLRRELEIYQS
LRVDAFKETLSLEEKLLNKHI SLSSLENEFRALLEEWKKEYAASSMVIDEHIAFIASVRNAPFCHNQYPFY
KEALHAPIPLFTVAQPI TEEDGLGIAEALLKVLREYCEIVKSI

>gi|874285978|gb|AKP71851.1| hypothetical protein CG09_1718
[Riemerella anatipestifer]

MFFSFHNAQRVI FKHLYKAFDASLRMVKEDYKAHFIVNLI R DFAHLNRKGNKQDNPD FNRYRFEKDGFF
TESGLLFFINLFLDKRDAYWMLKVKVSGFKASHKQREKMI TEVFCRSRILLPKLRLESRYDHNQMLLDMLS
ELSRCPKLLYKELKS EENKHFQVEADGFLDEI EEEQNPFKDTLIRHQDRFPYFALRYLDLNLNESFKSIRFQ
VDLGTYYHYCIYDKKIGDEQEKRHLLTRILLSFGRLQDFTEINRPQEWKALTKDLDYKETSNOFFI SKTIPH
YHITDNKI GFR LGT SKELYPSLEIKDGANRI AKYPYNSGFVAHAFI SVHELPLPMFYQHLIGKSEDLLE
TVRHITQRIYKDFEERINTIEDLEKANQGRPLGAFPKQMLGGLQNKQPDLESEKAKIKIEKLI AETKLLS
HRLNLIKLSKSSPKLGKRRKELIKTGV LADWL VKDFMRFPVAYDAQNP I KSSKANSTEFWFIRRALALY
GGEKNRLEGYFKQINLIGNINPHPLNKFNWKACRNLDVDFYQYLEQREKFL EAIKNQWPWEYQYCLLLKI
PKENRNKLVKQWEGGIGSLPRGLFTEAIRETLEDLMLSKPIRKEIKKHGRVGFISRAITLYFKKYQDK
HQSFYNLSYKLEAKAPLLKREEHYEWQONKQSPTE SQRLLELHTSDRWKDYLLYKRWQHLEKLLRLYRN
QDVMLWMLTLELTKNHFKELNLNYHQLEKLENLAVNVQEADAKLNPLNQTLPMLVPVKVYPATAFGEVQYH
KTPIRTVYIREEHTKALKMGNFKALVKDRRLNGLFSFIKEENDTQKHPISQLRRLRRELEIYQSLRVDFAK
EITLSLEEKLLNKHI SLSSLENEFRALLEEWKKEYAASSMVIDEHIAFIASVRNAPFCHNQYPFYKEALHAP
IPLFTVAQPI TEEDGLGIAEALLKVLREYCEIVKSI

>gi|874285979|gb|AKP71852.1| hypothetical protein CG09_1721
[Riemerella anatipestifer]

MEKPLLPHVYTLKHKFFWGAFLNIARHNAFITICHINEQLGLKTPSNDDKI VDVVETWNNILNNDHDL
KKSQLELILKHFPLTAMCYHPPKKEGKKGHQKEQQKESEASQSAEALNPSKLI EALEILVNQLHS
LRNYSHYKHKKPAEAKDI FNSLCLNNTDF

>gi|877802244|dbj|BAR96998.1| hypothetical protein PI17_22270
[Prevotella intermedia]

MEDDKKIKESTNMLDNKHFWAFLNLRHNVYITVNHINKVLELKNKQDII IDNDQDILA I KTHWEKV
NGDLNKTERLRELMTHKFFLETAIYIKNKEDKEVQKQAKAQSFDLSKHCLEFLFLEKLEQEARNYSH
YKSESTKEPMLLEKLLKMYNI PDDNI QLVIKDYQHNKDINPDEDDKHLDRTEEEFNYYPTINKKGNIT
ASGLLFFVSLFLEKKAIDWQKLRGFKDNRESKMKMTEVFCRSRMLLPKLRLESTQTDWVILLDMLNE
LIRCPKSLYERLQGEYRKKENVPFSDAEDYDAEQEPFKNII LVRHQDRFPYFALRYFDYNEIFTNLRFQI
DLGTYHFSIYKLLIGQKEDRHLTHKLYGFERIQEFKQNRIDEWKAIVKDFDIYETSEEPYISETAPHY
HLENQKIGIRFRNDNDEI WPSLKINGENNEKRKYKLDKQYQAEAFLSVHELPLPMFYLLKKEEPNNDK
KNASIVEGFIKREIRDIYKLYDAPANGENINDDLEKYCEDKGI PKRHLPKQMVAI LYDEHKDMAEAEAKR
KQKEMVKDIKLLLA TLEKQTGEI EDGGRNIRLLKSGEIA RLVNDMMRFQPVQKDNENPLNNSKANST
EYQMLQRSALALYNKKEKPTRYFRQVNLINSNPHPLKWI KWEECNNILSFRYSYL I KKI EFLNKLKPED
WEKNQYFLKLEKPTNRET LVQGWKNGFNLRGIFTEP IREWFKRHQNDSEYKQVETLDRVGLVTKVIP
LFFPKEDSKDEEYI LKKAQAQKELNINCVQPFYGFYVNGNI HKPDEKDFLPS EERKLLWGDKKYFKGYKA
KVKSKLIDKEEYRSYLFQSWNKFERELRLVRNQD I VTWLLCTELIDKLKVEGLNVEELKLLRLKDI
DIDTAKQEKNNILNRVMPMQLPVIVYEIDDSHNI VKDRPLHIVYI EETKIKLLKQGNFKALVKDRRLNGL
FSFVDTSSETLKSNI P I SKLSVEYELGEYQNARI ETIKDMLLEETLIEKYKTLPIDNFDMLNGWLEBK
DEADKARFQNDVKLLVAVRNAPFHNQYPMRNR IAFANINPFSLSADTSEEKLDIANQLDKDKIHKI I KR
IIEIEKPIETKE

>gi|885008798|gb|AKQ40303.1| hypothetical protein AS87_08290
[Riemerella anatipestifer Yb2]

MEKPLLPHVYTLKHKFFWGAFLNIARHNAFITICHINEQLGLKTPSNDDKI VDVVETWNNILNNDHDL
KKSQLELILKHFPLTAMCYHPPKKEGKKGHQKEQQKESEASQSAEALNPSKLI EALEILVNQLHS
LRNYSHYKHKKPAEAKDI FKHLYKAFDASLRMVKEDYKAHFIVNLI R DFAHLNRKGNKQDNPD FNRYR
FEKDGFFTESGLLFFINLFLDKRDAYWMLKVKVSGFKASHKQREKMI TEVFCRSRILLPKLRLESRYDHNQ
MLLDMLSELSRCPKLLYKELKS EENKHFQVEADGFLDEI EEEQNPFKDTLIRHQDRFPYFALRYLDLNL
NESFKSIRFQVDLGTYYHYCIYDKKIGDEQEKRHLLTRILLSFGRLQDFTEINRPQEWKALTKDLDYKETS
NOFFI SKTIPH YHITDNKI GFR LGT SKELYPSLEIKDGANRI AKYPYNSGFVAHAFI SVHELPLPMFY
QHLIGKSEDLLEKTVRHITQRIYKDFEERINTIEDLEKANQGRPLGAFPKQMLGGLQNKQPDLESEKAKI
KIEKLI AETKLLSHRLNIKLKSSPKLGKRRKELIKTGV LADWL VKDFMRFPVAYDAQNP I KSSKANST
EFWFIRRALALY GGEKNRLEGYFKQINLIGNINPHPLNKFNWKACRNLDVDFYQYLEQREKFL EAIKNQ
WPWEYQYCLLLKI PKENRNKLVKQWEGGIGSLPRGLFTEAIRETLEDLMLSKPIRKEIKKHGRVGFISRAI
TLYFKKYQDKHQSFYNLSYKLEAKAPLLKREEHYEWQONKQSPTE SQRLLELHTSDRWKDYLLYKRWQH
LEKLLRLYRNQDVMLWMLTLELTKNHFKELNLNYHQLEKLENLAVNVQEADAKLNPLNQTLPMLVPVKVY
PATAFGEVQYHKTPIRTVYIREEHTKALKMGNFKALVKDRRLNGLFSFIKEENDTQKHPISQLRRLRRELEI
YQSLRVDFAK EITLSLEEKLLNKHI SLSSLENEFRALLEEWKKEYAASSMVIDEHIAFIASVRNAPFCHNQ
YPFYKEALHAPIPLFTVAQPI TEEDGLGIAEALLKVLREYCEIVKSI

>gi|896419352|ref|WP_049354263.1| hypothetical protein [Riemerella
anatipestifer]

MFFSFHNAQRVI FKHLYKAFDASLRMVKEDYKAHFIVNLI R DFAHLNRKGNKQDNPD FNRYRFEKDGFF
TESGLLFFINLFLDKRDAYWMLKVKVSGFKASHKQREKMI TEVFCRSRILLPKLRLESRYDHNQMLLDMLS
ELSRCPKLLYKELKS EENKHFQVEADGFLDEI EEEQNPFKDTLIRHQDRFPYFALRYLDLNLNESFKSIRFQ
VDLGTYYHYCIYDKKIGDEQEKRHLLTRILLSFGRLQDFTEINRPQEWKALTKDLDYKETSNOFFI SKTIPH
YHITDNKI GFR LGT SKELYPSLEIKDGANRI AKYPYNSGFVAHAFI SVHELPLPMFYQHLIGKSEDLLE
TVRHITQRIYKDFEERINTIEDLEKANQGRPLGAFPKQMLGGLQNKQPDLESEKAKIKIEKLI AETKLLS
HRLNLIKLSKSSPKLGKRRKELIKTGV LADWL VKDFMRFPVAYDAQNP I KSSKANSTEFWFIRRALALY
GGEKNRLEGYFKQINLIGNINPHPLNKFNWKACRNLDVDFYQYLEQREKFL EAIKNQWPWEYQYCLLLKI
PKENRNKLVKQWEGGIGSLPRGLFTEAIRETLEDLMLSKPIRKEIKKHGRVGFISRAITLYFKKYQDK
HQSFYNLSYKLEAKAPLLKREEHYEWQONKQSPTE SQRLLELHTSDRWKDYLLYKRWQHLEKLLRLYRN

- continued

QDVMLWMLTLELTKNHFKELNLYHQLEKLENLAVNVQEADAKLNPLNQLPMLVLPVKVYPATAPGEVQYH
 KTPIRTVYIREEHTKALKMGKFNKALVKDRRLNGLFSFIKEENDTQKHPI SQLRRLRELEIYQSLRVDFAK
 EILSLEEKLLNKHISLSLENEFRALLEEWKKEYAASSMVIDEHI AFASVRNAPFCHNQYPPYKALHAP
 IPLFTVAQPI TEEDGLGIAEALLKVLREYCEIVKSQI

>gi|914347928|gb|AKV64040.1| hypothetical protein PGA7_00008170
 [*Porphyromonas gingivalis*]
 MNIIVPASENKQQSRIVEDDPOYFGLYLNLAARENLEEVESHVRIKFGKKKLNEESLKQSLLCDHLLSVDRW
 TKVYGHSSRRYLPFLHYFPDPSQIEKDHDSKTGVDPSAQRLIRELYSLDDLRNDFSHNRLDGI TFEHLE
 VSPDISSFITGTYSLACGRAQSRFADFFPKPDDFVLAKNRKEQLISVADGKECLIVSGLAFFICLFLDREQ
 ASGMLSRIRGPKRIDENWARAVHETFCDL CIRHPHDLRESSNI KEALLDMLNENRCPRI LYDMLPEEE
 RAQFLPALDENSMNLS ENSLNEESRLLDWGS SDWAEALTKRIRHQDRFPYMLRFIEEMDLLKGRFRV
 DLGEIELDSYKVKVGRNGEYDRTITDHALAFGKLSDFQNEEEVSRMISGEASYPVRFSFLFAPRYAIYDNK
 IGYCHTSDPVYPKSTGEEKRALSNDPQSMGFI SVHNLRKLLLMELLCEGFSRMSQSDFLRKANRILDETAE
 GKLFQPSALFPENRHRFIPPQNPKSKDRREKAETTLKYEKQEI KGRKDKLNSQLLSAFDMNQRLPSRLD
 EWMNIRPASHSVKLRITYVKQLNEDCRLRLRFRKRDGDGKARAIPLVGE MATFLSQDIVRMI I SEETKKLI
 TSAYNEMQSRSLAQYAGEENRRQFRAI VAELHLLDPSGHPFLSATMETAHRYTEDFYKCYLEKKREWLA
 KTFYRPEQDENTKRRISVFFVDPGEARKLLPLLI RRRMKEQNDLQDWIRNKQAHPIDLPSHLFD SKIMEL
 LKVKDGKKKWEAFKDWSTKYPDGMQPFYGLRRELNIHGKSVSYPISDGKGFADCYTHLMEKTVQDKKR
 ELRTAGKVPVDDLAADIKRSRPHRAVNEREFMLRLVQEDDRMLMAINKMMDREEDILPGLKNIDSILDE
 ENQFSLAVHAKVLEKEGEGGNSLSLVPATIEIKSKRDKWSKYIRYRDRRVPGLMSHPPEHKA ILDEVK
 ILLGEYDRCRIKIPDWAFALEGAIMSDRDLKPYLHSSSRGKSGEHSSTLVKMLVEKKGCLTPDESQYLI
 LIRNKAHNQFP CAEMPLIYRDVSAKVGSI EGS SAKDLPEGSSLVDSLWKKYEMI IRKILPILDPENRF
 FGKLLNNMSQPINDL

>gi|914348650|gb|AKV64762.1| hypothetical protein PGA7_00015700
 [*Porphyromonas gingivalis*]
 MTEQNEKPYNGTYTLEDEKHPWAAPFNLARHNAYITLAHIDRQLAYS KADI TNDEDILFFKQGWKNDND
 LERKARLRSILKHPSPLEGAAYGKLFESQSSGNKSSKKKELIKKEKEELQANALS LDNLKSI LDFDPLQ
 KLDKFRNYSHYRHPSESELPLDFGNMLQRLYVNFVDSVQVRKRDHEHNDKVDPHRHFNLVVRKGGKDKY
 GNNNDNPPFKHHPVDREKVT EAGLLFFVSLFLEKRD AIWMQKIRGFKGGTEAYQQMTNEVFCRSRISLP
 KLKLESRLIDDWMLDMLNELVRCPKLLYDRLEEDRARFRVVDILSEDDIDGTEEDPFKNILVRRHQD
 RFPYFALRYFDLKKVFTSLRPHIDLGTYHFAIYKKNIGEQPEDRHLTRNLVYGFGRIQDFAEHRPEEWR
 LVRLDYPETGDKPYITQTPHYHI EKGKILGRFVPEGQLLWSPPEVGATRTRGSKYAQDKRFTAEAFSL
 VHELPMMPFYYFLLEKYS EEAASA EKVQGR I KRV IEDVYAVYDAFARDEINTRELDACLADKGRIRRGHL
 PRQMIALSLSQEHKMEEKVRKQLQEMIADTDHRLDMLDRQTRKIRIRGRKNAGLPKSGVIA DWLVRDMMR
 PFPVAKDTSKGPLNNSKANSTEYRMLQRALALFGGEKERLIPYFRQMLNIGGNNPHFPFLHETRWESHINI
 LSFYRSYLKARKAFLOSGRSDREENHRFLLEKPKTRDQTLVAGWKS EPHLPRGIFTEAVRDC LIEMGY
 DEVGSYKEVGFMAKAVPLYFERACKDRVQPFYDYPFNVGNLSLKKPKGRFLSKEKRAEWEWESGKERFRDLE
 AWSHAARRIEDAFVGI EYASWENKKIEQLLQDLSLWETPESKLKVKADKINIAKLLKEILEAKEHPYH
 DFKSWQKFERELRLVKNQDII TWMNCRDLMEENKVEGLDIGTLYLKDIRTDVQEQGSLNVLNHVKPMRPL
 VVVYRADS RGHVHKEEAPLATVYIEERDTKLLKQGNFKSFKDRRLNGLFSFVDTGALAMEQYPI SKLRV
 EYELAKYQATARVCAFEQTLLEESLLTRYPHLPDES FREMLESWSDPLLKWPDLQREVRLLI AVRNAPS
 HNQYPMYDETIFSSIRKYDPSLDAIEERMGLNI AHRLESEVKKLAKEMVERIIQA

>gi|916087402|ref|WP_050955369.1| hypothetical protein [*Prevotella
 intermedia*]
 MEDDKKIKESTMNLNKHFWAFLNARHNAYITVNHINKVLELKNKDKQDII IDNDQDILA IKTHTWEKV
 NGDLNKTERLRELMTKHFFLETAIYIKNKEDKEEVKQEKQAKASFDLSLKHCLFLFLEKLEQEARNYYS
 YKYESSTEKPEMLLEKLLKMYNI PDDNIQLV I KDQHNKDINPDEDPKHLDRTEEFNYIPT INKNGNI T
 ASGLLFFVSLFLEKKAIDAWMQKLRGFKDNRESKMKMTEVFCRSRMLLPKRLLESTQTDQWILLDMLNE
 LIRCPKSLYERLQGEYRKKENVPFSDAEDYDAEQEPFKNILVRRHQDRFPYFALRYFDYNEIFTNLRFOI
 DLGTYHFSIYKLLIGGQEDRHLTHKLYGFERIQEFKQNRIDEWKAIVKDFDIYETSEEPYISETAPHY
 HLENQKIGIRFNDNDEI WPSLKINGENNEKRYKLDKQYQAEAFSLVHELPMMPFYYLLKKEEPNNDK
 KNASIVGFIKREIRDIYKLYDAFANGEINNIDLEKYCEDKGI PKRHLPKQMAI LYDEHKDMAEEAKR
 KQKEMVKDIKLLLATLEKQTEI IEDGGRNIRLLKSGE IARWLVNDMMRFQPVQKDNENGPLNNSKANST
 EYQMLQRS LALYNKEEKPTRYFRQVNLINSNPHPLKWI KWEECNIISLFSYRSLIKKIEFLNKLKPED
 WEKNQYFLKLKEPKTNRETLVQGWKNGFNLPRGIFTEP IREWFKRHQNDSEYEVETLDRVGLVTVKVIP
 LFPKEDSKDKEEYLKKAQKEINNCVQPFYGFYVNGNI HKPDEKDFLPS EERKLLWGDKKYKPKGYKA
 KVSKKLLIDKEEYRSYLEFQSWNKFERELRLVRNQDIVTWLLCTELIDKLVKVEGLNVEELKLRLLKDI
 DIDTAKQEKNNILNRVMPMQLPVIVYEIDDSHNI VKDRPLHIVYI EETKIKLLKQGNFKALVKDRRLNGL
 PSFVDTSETBLKSNPISKLSVYEYELGEYQNAIETIKDMLLEETLIEKYTKLPIDNFS DKLNGWLEKG
 DEADKARFQNDVKLLVAVRNAPFHNQYPMRNR IAFANINPFSLSADTSEEKLLDIANQLKDKIHKI IKR
 IIEIEKPIETKE

>gi|916715432|ref|WP_051322523.1| hypothetical protein [*Bacteroides
 pyogenes*]
 MALLNVRKVENHIRKWLGDVALLPEKSGFHSLLTIDNLSAKWIRFYYKSRKFLPFLEMFSDSKSYENR
 RETTECLDTIDRQKISSLLKEVYGLQDIRNAFHYHIDDQSVKHTALIISSEMRHFIENAYSFALQKTR
 ARFTGVFVETDFLQAEKGDNKKVFAIGNEGILKLDNALIFLICLFLDREAEFKLSRATGFKSTKEKG
 FLAVRETFCALCCROPHERLLSVNPREALLMDMLNENRCPDILFEMLDEKDKSFLPLLGEEBQAHILE
 NSLNDELCEALDDPFEMIASLSKRVRYKRRFPYMLRYIEEKNLLPFRFRIDLGLCELASYPKMGEEEN
 NYERSVIDHAMAFGRLLDFHNEAVLQQI IKGITDEVRFSLYAPRYAIYNNKIGFVRI GGSDDKISFPILK
 KGGEGHCVA YTLQNIKSFGPIS IYDLRKL LLSFLDKDKAKNIVSGLLEQCEKHWDLS ENLPDAIRTE
 LQKEFPVPLIRYTLPRSKGGKLVSSKLADKQEKYSEFERREKELTEILSEKDFLSQIPRRMIDEWLNV
 LPTSREKLLKGYVETLKLDCRELRVFEKREKGEHPVPRIGEMAI DLAKDIRMVIDQGVKQRTSAYY
 SEIQRLAQYAGDNNRRHLDISI IRELR LKDTKNGHPFLGKVL RPLGLGHT EKLYQRYFEEKLEWLEATFYP
 AASPKRVRPFRVNP I GQKQELPLIIRNLMKERPEWRDWKQRKNSHPIDLPSQLFENEICRLLDKDIKIGEP

- continued

SGKLGWNEFMFLYWDKEFPNGMQRFYRCRRRVEVFDKVVVEYEESEEGNYKKYEEALIDEVVRQKISSK
 EKSKLQVEDLILSVRRVFKRAINEKEYQLRLLCEDDRLLFMAVRDLYDWKEAQLDLDKIDNMLGEPVSVS
 QVIQLEGGQPDVAI KAECKLKVSKLMRYCYDGRVKGLMPYFANHEATQEQVEMELRHVEDHRRRVFNWV
 FALEKSVLKNKELRRFYBESQGGCEHRRCTDALRKASLVSEEEYEFVLRHNRKSAHNQFPDLEIKGLPPN
 VTSGFCECISWYKAIICRIIPFIDPERRFFGKLEEQ

>gi|916915772|ref|WP_051522484.1| hypothetical protein [*Prevotella
 saccharolytica*]
 MEDKPFWAFFNLARHNVYLI VNHINKLLEKLYDEGKHKEIFEREDI FNISDDVMNDANSNGKKRKL
 IKKIWDLLDIDLTRKYLRELILKHFPFIQPAII GAQTERTTIDKDKRSTSTSNDSLKQGTGEGDINDLL
 SLSNVKSMFRLQLLEQLRNYSHVKHKSAIMPNFDELLNWMRYIFIDSVNKVKEDYSSNSVIDPNT
 SFSHLIYKDEQGIKPCRYPTSKDGSINAFGLLFFVSLFLEKQDSIWMQKKIPGFKKASENYMKMTNEV
 FCRNHILLPKIRLETVYDKWMLDMLNEVVRCP LSLYKRLIPAAQNKFVPEKSSDNANRQEDDNPFSR
 ILVRHQNRFPYFVLRFPDLNEVFTLRFPQINLGCYHFAICKQIGDKKEVHHLIRTLYGF SRLQNTQNT
 RPEEWNTLVKITEPSSGNDGKIVQGVPLPYISYTI PHYQIENEKIGIKIFDGD TAVD TDIWPSVSTEQ
 NKPDKYTLIPGFKADVFLSVHLELLPMMFYQLLCEGLMKIDAGNAVEKVLIDTRNAIFNL YDADFVQEK
 NITITDLENYLDKPLIIGHLPKQIMIDLLKGHQDMLKAVEQKAMLIKD TERRLLKLDKQKQEQTDVAAK
 NIGILLKNGQIADWLVDMMRFQPVKRDKEGNPINC SKANSTEQMLQRAFAYIDS CRLSRYTQLHL
 IHSNSHLFLSRFEYDKQPNLIFAYAYLAKLEFLNELQPQNWASDNYFLLLRAPKNDRQKLAGEGWNG
 FNLPRGLFTEKIKTFWNEHKTIVDISDCDIPKNRVGVARLIPVFFDKKFKDHSQPFYRYDFNVGNVSKP
 TEANVLSKGRRELKFSYQNKFKNNIPAEKTKKEYREYKFNFLWKKFERELRLIKNQDILIWLMCKNLDE
 KIKPKKDI LEPRIVSYIKLDSLQINTSTAGSLNALAKVVPMTLAIHIDSPKPKGAGNNEKENKEFTVY
 IKEEGTKLLKWNFKILLADRRIKGLFSYIEHDDIDLKQHPLIKRRVDELDLYQTCRIDIFQQTLGLEA
 QLLDKYSDLNIDNFYQMLIGWRKKKGI PRNIKEDIDFLKDVNRNAPSHNQYPSKKIAPFRIRKFNPKELI
 LEEEBGLGIATQMYKEVEKVVNRKIRIELF

>gi|918017236|ref|WP_052309609.1| hypothetical protein [*Prevotella
 saccharolytica*]
 MEKENVQSGSHIYEPDKCFWAAFYNLARHNAYLTI AHINSFVNSKKGINDDKVLDI IDDSKFDNDLL
 MGARLNKLLKHFPLKAPLYQLAKRTRKQGGKEQDY EYKGGEDDEVEIQEAIANAFKMANVRKILHAF
 LKQLEDLRNHFSHYNYNSPAKMEVKPDDGFCNKLYYVFDALQMVKDDNRMNPEINMQTDFEHLVRLGR
 NRKI PNIPKYNF INSDGTINNGLLFFVSLFLEKRD AIWMQKKIKGPKGGTENYMRMNEVFCRNRVPI
 KLRLETVDYDNLQMFDMNLVRCPLSLYKRLKQEDQDKFRVPIEFLEDEDNEADNPYQENANSDENPTEE
 TDPLKNILVRHQHRRFPYFVLRFPDLNEVFKQLRFPQINLGCYHFSIYDKTIGERTEKRLHTRILFGFDRLQ
 NFSVKLQPEHWNMVKHLDTEESSDKPYLSDAMPHYQIENEKIGIHLKIDTEKKEVWVSEVEEVSNS
 RNKYKSEKNLTADAFSLTHELPPMFYQLLSSEKTRAAAGDKVQGVLSYRKKIFDIYDDFANGTINS
 MQKLERLAKDNLGRNMPQMLAI LEHQEPDMEQKAKEKLDRLITETKRIKLEDDQFKQKVRIGKRRR
 DLPKVGSTADWLVDMMRFQPAKRNADNIGVPSKANSTERYLLQEAALAFYSAYKDRLEPYRQVNLIGG
 INPHFPLHRVWKKCNHLSFYHLYLEAKEQYLSHLSPADWQKHQHLKVRKDIQNEKKDWKSLVAG
 WKNGFNLPRGLFTEIKTWFS TDADKVQIIDTKLFENRVGLIAKLIPLYDKVYNDKQPPYQYFPNIND
 RYKPEDIRKRP TAAASKLWYDAKRNADNIGVPSKANSTERYLLQEAALAFYSAYKDRLEPYRQVNLIGG
 DLFAQCTVEGVEFADLKLSQL EVDVNVQDNLNVLNNSVSMILPLSVYPSDAQGNVLRNSKPLHIYVYQEN
 NTKLLKQGNFKSLLDKRRNLNGLFSFIAAEGEDLQOHLIKNRLEYELSIYQTMRSVFEQTLQLEKAILT
 RNKTLCGNFNMLNSWSEHRTDKKILQPDIDFLI AVRNFAPSHNQYPMSTNIVMQGIEKFNIQTPKLEEK
 DGLGIASQLAKKIKDAASRLQNIINGGIN

>gi|919116504|ref|WP_052672145.1| hypothetical protein [*Prevotella
 intermedia*]
 MEDDKKITESTNMLDNKHFWAAFNLARHNAYITVNHINKVLELKNKKQDII IDNDQDILAIKTHWEKV
 NGDLNTERLRELMTKHFPLETAIYIKNKEDKEEVKQEKQAEASLESKDCFLFLEKLEARNYYSH
 YKSESTKEPMLLEGLLEKMYNI FDDNIQLVKDYQHNDKINPDEDFKHLDRKGGQKYSFADNEGNITES
 GLLFFVSLFLEKDAIWMQKLGKDNRESKMMTHEVFCRRMLLPKLRLESTQTQDWIILLMLNELI
 RCPKSLYERLQGEYRKKENVPFDSADEDYDABQEPFKNILVRHQDRFPYFALRYFDYNEIPTNLRFPQIDL
 GIYHFSIYKLLGGQKEDRHLTHKLYGFERIQEFAKQNRPEWKALVKDLDIYETSNERYISETTPHYHL
 ENQKIGIRFRNGNKEIWP SLKINGENNEKSKYKLDKPYQAEAFSLVHLELLPMMFYLLKKEEPNNDKKN
 ASIVEGFIKREIRDMYKLYDAPFANGEINNI GDLEKYCEDKGI PKRHLPKQMVAILYDEPKDMVKEAKRKQ
 KEMVKDIKLLATLEKQEQEIEDGGRNIRLLKSGE IARWLVDMMRFQPVQKDNENGNLNNSKANSTEQ
 QMLQRSALALYNKEEKPTRYFRQVNLINSSNPHFLKWI KWEECNMILSFYRNYLIKKIEFLNKLKPEDWE
 KNQYFLKLEKPTNRETLVQGWKNGFNLRGIFTEPIREWFKRHQND SKEYEKVEALKRVGLVTKVIPLF
 FKEEYFKEDAQKEINNCVQPPYSPYVNVGNIHKPDEKDFLSEERKKLWGDKDKFKGKYAKVKSKKLDI
 KEKEEYRSYLEFQSWNKFERELRLVRNQDITVWLCTELIDKMKVEGLNVEELQKLRKLDIDTD TAKQEK
 NNI LNRIMPQLPVIVYIEDSHNVKDRPLHIVYI EETKIKLLKQGNFKALVKDRRLNGLFSFVDTSSK
 AELDKKPI SKSVVEYELGEYQAR IETIKDMLLEKTLIKKYEKLPIDNFSMDMLNGWLEGGKDESDFARFQ
 NDVKLLVAVRNAPSHNQYPMRNI AFANINPFLSSADISEEKKLDIANQLKDKIKHIKIKKIEIEKPIE
 TKE

>gi|919932447|ref|WP_052912312.1| hypothetical protein [*Porphyromonas
 gingivalis*]
 MTEQNEKPYNGTYTLEDKHFWAAFNLARHNAYITLAIH DRQLAYS KADITNDEDILFFKQGWKNLNDND
 LERKARLRSLILKHFSLFLEGAAYGKLFESQSSGNKSSKKELIKKEKEELQANALS LDNLKSI LFDPLQ
 KLDKFRNYYSYHRHPESSELFDGNNLQRLYVNFVDSVQRVQRDEHNDKVDPHRFHNLVLRKGGKDKY
 GNNNDNPFKHHFVDREKVT EAGLLFFVSLFLEKRD AIWMQKKIRGFKGGTEAYQOMTNEVFCRSRISLP
 KKLKLESRLIDWMLDMLNELVRCPLLYDRLEEDRARFRVVPV DILSDEDDIDGTEEDPFPKNI LVRHQ
 RFPYFALRYFDLKKVFTSLRPHI DLGTYHFAIYKNI GEQPEDRHLTRNL YGFGRIQDFAEHRPEEWKR
 LVRDLDFETGDKPIYITQTPHYHI EKGKIGLRFVPEGLWPSPEVGATRTGRSKYAQDKRFTAEAFSL
 VHELPMMFYFLLREKYS EEAESAEKVQGR IKRVI EDVYAVYDARFARDEINTRDELDA CLADKGI RRRGH
 PRQMI ALSQEHKMEKVRKLLQEMI ADTDHRLDMLDRQTDKIRI GRKNAGLPKSGVIADWLVRDMMR
 FQPVAKDTS GKPLNNSKANST EYRMLQRALALFGGKERLIPYFRQMLNIGGNNPHPLHETRWESHINI

- continued

LSFYRSYLKARKAPLQSIGRSDREENHRFLLLKPKTDROTTLVAGWSEFHLPRGIFTEAVRDCLEIMGY
 DEVGSYKEVGPMAKAVPLYFERACKDRVQPPYDYFPNVGNLSLKKPKGRFLSKERAEEWESGKERFRDLE
 AWSHSAARRIEDAFVGI EYASWENKKIEQLQLDLSLWETPESKLVKADKINIAKLKKEILEAKEHPYH
 DFKSWQKFERELRLVKNQDII TWMMCRDLMEENKVEGLDITGLYLKDIRTDVQEQGSLNVLNHNKPMRLP
 VVVYRADSRGHVHKEEAPLATVYIEERDTKLLKQGNFKSFVKDRRLNGLFSFVDTGALAMEQYPI SKLRV
 EYELAKYQTARVCAFEQTLLEESLLTRYPHLPDESFRMLESWSDPLLDKWPDQLQREVRLLI AVRNAFS
 HNQYPMYDETFISSIRKYDPSLDAIEERMGLNIAHRLSEEVKLAKEMVERIIQA

Large protein sequences:

>1 hypothetical protein [*Bergeyella zoohelcum*]
 MENKTSLGNNIYYNPFKQDKSYFAGYFNAAMENTDSVFRLELGRKLGKEYTSENFDDAI FKEN
 ISLVEYERYVLLSDYFPMARLLDKKEVPI KERKENFKKNFKGII KAVRDLRNFYTHKEHEGEVE
 I TDEIFGVLDEMLKSTVTLVKKKKVKTDKTKKELKKSIEKQLDILCQKKLEYLRDRTARKIEEKR
 RNQRERGEKELVAPFKYSDKRDLLIAAIYNDADFVYIDKKDSLKESKAKYNTKSDPQQEGD
 LKIPISKNGVVFLLSFLTKQEI HAFKSKIAGFKATVIDEATVSEATVSHGKNSICFMATHEIF
 SHLAYKLLKRVRTAEINYGAEANAELQSVYAKETLMMQMLDELSKVPDQVYQNLSEDEVQKTFI
 EDWNEYLKENNGDVGTEEBEQVIHPVIRKRYEDKFNYPAIRFLDEFAQFPPTLRFPQVHLGNVLDH
 SRPKENLISDRRIKEKIVTGFRLSELEHKKALFIKNTETNEDREHYWEI FPNPNYDFPKENI SV
 NDKDFPIAGSILDRKQPVAGKIGIKVLLNQQYVSEVDKAVKAHQKQKASKPSIQNIIEEI
 VPINESNPKAEI VFGQPTAYLSMNDIHSILYEFDFKWEKKKLEKKEKELRKEIGKELEKK
 IVGKIQAQIQQIIDKDTNAKILKPYQDGNSTAI DKEKLIKDLKQEQNILQKLKDEQTVREKEYN
 DFIAYQDKNRINKVRDRNRHKQYLLKDNLKRKYPEAPARKEVLYREKGVAVWLANDIKRPMPT
 DFKNEWKEGQSHLLQKSLAYYEQCKEELKNLLPEKVFQHLPEKLGQYQKQYLYQFYTCYLDKR
 LEYISGLVQQAENFKSENKVPKVENECFKFLKKNYTHKELDARVQSILGYPIFLERGFMDK
 PTIIKGGTFKGMELFADWFRYYKEYQNFQTFYD TENYPLVELEKKQADRKRRTKIYQQKNDV
 FTLLMAKHIKFSVFKQDSIDQFSLLEDLYQSREERLGNQERARQTGERNTNYIWNKTVDLKLCDG
 KRTVENVKLKNVGDPIKYEYDQVQAPLKYEENI EWQAFLIKESKEEENYPYVVERIEIQYEK
 RREELLKEVHLIEEYILEKVKDKKELKKGDNQNFKYIYLNGLLQKLKNEDEVESYKVPNLNTEPE
 DVNINQLKQEAETDLBQAFVLTYYIRNKFAHNQLPKKEPFDYCYQEKYGIKEKTYAEYFAEVFK
 KEKEALIK

>2 hypothetical protein [*Prevotella intermedia*]
 MEDDKKTTDSIRYELKDKHFWAAFLNLRHNVIYITVNHINKILEEGEINRDGYETTLKNTWNEI
 KDINKKDLKSLIIKHFPFLEAATYRLNPTDTTKQKEEKQAEQSLERKSPFVFIYKLRDLR
 NHYSHYKHSKSLERPKFEEGLEKMYNIFNASIRLVKEDYQYKNDINPDEDFKHLDRTEEFENY
 YFTKDNENI TESGLLFFVSLFLEKKDAIWMQQLRGPKNRENKMTNEVFCRSRMLLPKLR
 LQSTQTDWI LLDMLNELIRCPKSLYERLREEDREKFRVPIEIADEDYDAEQEPFKNTLVRHQD
 RFPYFALGYFDYNEIFITNLRQIDLGTYHFSIYKQIGDYKESHHLTHKLYGFERIQEFTKQNR
 PDEWRKFKVTFNSFETSKEPIPETTPHYHLENQKIGIRFRNDNDKIWP SLKTNSEKNEKSKYK
 LDKSPQAEAFLSVHELLPMMFYLLKKTENTDNDNEIETKKNENKNDKQEKHKIEEIEENKITE
 IYALYDFPANGIKSIDLEEBEYCKGKIDIEIHLPKQMIAILKDEHKVMATEAERQEEMLVDVQ
 KSLESLDNQINEEINENVERKNSSLSGKIASWLVNDDMRFPQVQKDNKGP LNNSKANSTHEYQL
 LQRTLAFPGSEHERLAPYFKQTKLIESNPHFLKDETEWKCNILSFPYRSYLAEKNFLESKL
 PEDWEKNQYFLKLEKPKTKPKTLVQGWKNGFNLPRGIFTEPIRKNWFMKHRENI TVABLKRVLV
 AKVILPFPSEYKDSVQPPYNYHFNVGNINKPDEKNFLNCEBERLELLRKKKDEPKMTDKEKE
 NPSYLEPKSWNKFERELRLVRNQDITVWLLCMELFNKIKI KELNVEKIYLNKINNTNTTKKENT
 EEKNGEENIKENKNNILNRIMPRLPIKYGRENFSSKNKKKIRRNTFFTVYIEEKGTKLQKQ
 NFKALERDRRLGGLFSFVKTPSKAESKNTISKLRVEYELGEYQKARIEIKDMLALEKTLIDK
 YNSLDTDNFNKMLTDWLELKGEPDKASFNQNDVLLI AVRNAFSHNQYPMRNRIFAANINPFSLS
 SANTSEKGLGIANQLKDKTKHTKIEKIEIEKPIETKE

>3 hypothetical protein [*Prevotella buccae*]
 MQQDKLFPVDRKKNIAIFAPPKYITIMENKEKEPIIYELTDKHFWAAFLNLRHNVIYITNHN
 RRLEIABLKDDGYMMGIIKGSWNEQAKKLDKVRDLRDLIMKHFPFLEAAAYEMTNSKSPNNKEQR
 EKEQSEALSLNNLKNVLPFLFKLQVLRNYSHYKSEESPKPIFETSLKNNMYKVPDANVRLV
 KRDMHHEINIDMQRDFTHLNRKKQVGRTKNIIDS PNFHYHFPADKEGNMTIAGLFFVSLFLDKK
 DAIWMQKLLKGFKDRNLRQMTNEVFCRSRISLPKLLKLENVQTKDWMQLDMLNELVRCPSLY
 ERLREKDRSEFKVPFDIFSDDYNAEEEPFKNTLVRHQDRFPYFVLRYPDLNEIFEQLRFQIDLG
 TYHFSIYNKRI GDEDEVRLHLLHLYGFARIQDFAPQNPQPEWRKLVKDLDFHETSQEPYISKTA
 PHYHLENEKIGIKCSAHNNLFPSSLQTDKTCNGRSKFNLTGQFTAEAFLSVHELLPMMFYLLL
 TKDYSRKESADKVEGIRKEISNIYAYADAFANNEINSIADLTRLRLQNTNLLQGLPKQMSIL
 KGRQKDMGKEAERKIGEMIDDTQRRLDLCQTNQKIRIGKRNAGLLSKGKIADWLVNDDMRFP
 PVQKQNNIPINNSKANSTEYRMLQRALALFGSENFRLKAYFNQMNLVGNDNPHFLAETQWEH
 QNTNLSFYRNYLEARKKYLKGLKPNQWQYQHFLILKVQKTNRNTLVGTWKNFNLPRGIFPTQ
 IREWFEEKHNSKRIYDQILSPDRVGFVAKAIPLYFAEYKDNVQPPYDYFPNI GNRLKPKKRQF
 LDKKERVLEWQKNKELFKNYPSEKKTDLAYLDFLSWKKFERELRLIKNQDITVWLMFKELFNM
 ATVEGLKIGETHLRDITNTANEESNNILNRIIMPMLPKVTKYETDNKGNILKERPLATFYIEET
 ETKVLKQGNFKALVKDRRLNGLFSFAETTDLNLLEHPI SKLSVDLELIKQYQTRISIFEMTLGL
 EKKLIDKYSTLPDTSFRNMLERWLQCKANRPELKNYVNSLIAVRNAFSHNQYPMYDATALFAEVK
 KFTLFPSPVDTTKIELNIAPOQLLEIVGKAIKIEKSENK

>4 hypothetical protein [*Porphyromonas gingivalis*]
 MNTVPASENKGQSRVTEDDPQYFGLYLNLARENLEVEESHVRIKFGKLLNEESLKQSLCDHL
 LSVDRWTKVYGHRSRYLPFLHYFDPSQIEKDHD SKTGVDPSAQRLIRELYSLDLFNRNDFSH
 NRLDGTTFEHLVSPDISFIFTGTYSLACGQAQSRFAVFFKPDFFVLAKNRKEQLISVADGKEC
 LTVSGFAFFICLFLDREQAAGMLSRIRGFKRTDENWARAVHETFCDLICRHPHDLRESNTKEA
 LLLDMLNELNRCRILYDMLPEEERAQFLPALDENSMNLSNSLDEBSRLLDGSSDWAELT

- continued

KRIRHQDRFPYMLRFIEEMDLLKGRFRVDLGEIELDSYSKVKVGRNGEYDRTITDHALAFGKL
SDFQNEEEVSRMISGEASYPVRFSLFAPRYAIYDNKIGYCHTSDPVYPKSKTGEKRALSNPQSM
GFISVHDLRKLLELLMELLCEGFSFRMQSDFLRKANRILDETAEGKLFQFALFPEMRHRFP
KSKDRREKAETTLKQYKQEIKGRKDKLNSQLLSAFDMQDQRLPSRLLEDWMMNIRPASHSVKLR
YVKQLNEDCRLRLRFRKRDGDGKARAIPLVGEMATFSLQDIVRMIISEETKLIITSAYNEMQ
SLAQYAGEENRRQFRAI VAE LRLLDPSSGHPFLSATMETAHRYTEBGFYKCYLEKKREWLA
RPEQDENTKRRISVFFVPDGEARKLLPLLLIRRRMKEQNDLQDWIRNKQAHPIDLPSHLFDSKVM
ELLKVKDGKKKWNFAFKDWSKYPDGMQPPYGLRRELNIHGKSVSYPISDGGKFPADCYTHLME
KTVRDKRELRTAGKPVPPDLAADIKRSHRAVNEREFMLRLVQEDDRMLMAINKMMTDREED
ILPGLKNIDSILDEENQFSLAVHAKVLEKEGEGGDNLSLVPATIEIKSRKRDWSKYIRYRDR
RVPGLMSHFPEHKATLDEVKTLLEGEYDRCRIKIPDWFALLEGAIMSDRDLKPYLHSSSREGKS
GEHSTLVKMLVEKKGCLTPDESQYLILIRNKAHNQFP CAEMPLIYRDVSAKVGSI EGSSAKD
LPEGSSLVDSLWKYEMIRKILPILDPENRFFGKLLNMSQPINDL

>5 hypothetical protein [*Bacteroides pyogenes*]
MESIKNSQKSTGKTLQKDPYFGLYLNMALLNVRKVENHIRKWLGDVALLPEKSGFHSLLTTDN
LSSAKWTRFYKSRKFLPFLMFDFSDKKS YENRRETAECLDTIDRQKISSLLKEVYGLQDIRN
AFSHYHIDDQSVKHTALIISS EMHRFIENAYSFALQKTRARFTGVFVETDPLQAEKGDNKKFF
AIGGNEGIKLKDNALIFLILCLFLDREAEAFKFLSRATGFKSTKEKGF LAVRETFCALCCRPHER
LLSVNPREALMDMLNENLRCPDILFEMLDEKDKSFLPLLGEEEQAHILENSLNDELCEAIDD
PEMIALSLSKRVRYKRPYMLRYIEEKNLPIRFRIDLGLCELASYPKMGGEENYERSVT
DHAMAFGRLLTDFHNEDAVLQQITKGITDEVRFSLYAPRYAIYNNKIGFVRTSGSDKISFPTLKK
KGGEGHCVAITLQNTKSPGFI SIYDLRKLILLSFLDKAKANI VSGLLEQCEKHKDLSENLF
AIRTELQKEFPVPLIRYTLPRSKGGKLVSSKLADKQEKYSEFEFERRKELTEILSEKDFDLSQI
PRRMIDEWLNVLPTREKLLKGYVETLKLDCRERLRVPEKREKEGHEPLPRIGEMATLAKDII
RMVIDQGWKQRIITSAYSEIQRCCLAQYAGDDNRRHLSI IRELRLKDTKNGHPFLGKVLRLPGL
HTEKLYQRYFEKKEWLEATFYPAA SPKRVPRFVNPPTGKQKELPLIIRNLMKERPEWRDVKQR
KNSHPIDLPSQLFENEICRLLKDKIGKEPSGKLNWEMFKLYWDKEFPNGMQRFYRCRKRVEVP
DKVVEYXSEEGGNKYKYLEALIDEVVRQKISSSEKSKLQVEDLTLVRRVFKRAINEKEYQL
RLLCEDDLLFMVAVRDLVDWKEAQLDLDKIDNMLGEPVSVSQVIQLEGGQDAVIAEKCKLDV
SKLMRYCYDGRVKGLMPYFANHEATQEQVEMELRHYEDHRRRVFNWVFALEKSVLKNEKLRRFY
EESQGGCEHRRCIDALRKASLVSEEEYEFVHIRNKS SAHNQFPDLEI GKLPNVTSGFCERICWS
KYKAIICRIIPFIDPERRFFGKLLBQK

>6 hypothetical protein [*Alistipes* sp. ZOR0009]
MSNEIGAPREHQFAYAPGNEKQEEATFATYFNLALSNGVGMFGEVESNPDKIEKSLDTLPPAI
LRQIASFIWLSKEDHPDKAYSTEVEKVIVTDLVRLCFYRNYFSHCYFLDTQYFYSDDELVDTTA
IGEKLPYNFHFIITNRLFRYSLPEITLFRWNEGERKYEILRDGLIFPCLFLKRGQAERFLNEL
RFFKRTDEEGRIKRTIFTKYCTRESHKHIGIEEQDFLIQDII GDLNRPVKCDGVVDLSKENE
RYIKNRETSNESDENKARYRLIREKDKFPYMLRYIVDFGVLPCITFKQNDYSTKEGRGQPHY
QDAVAQEEERCYNFVVRNGVYYSYMPQAQNVVRISELQGTISVEELRNMYASINGKDVNKS
EQYLHYLHLLYKILITISGQTIKEGRVDVEDYRPLLDKLLRPA SNGEELRRELKLLPKRVCD
LLSNRFDSCSEGVSAVEKRLKAILLRHEQLLLSQNPALHIDKIKSVIDYLYLFFSDDEKFRQOQT
EKAHRLKDEBEFQMYHYLVGDYD SHPLALWKELEASGRLLKPEMRKLTSA TSLHGLYMLCLKGT
EWCRLQMSIGKGTAKVEA IADRVGLKLYDKLKEYTPEQLEREVKLVVMHGYAAATPKPKAQA
AIPSKLTELRFYSPFGKREMSFAAFIRQDKKAQKLWLRNFYTVENIKTLQKRQAADAACCKLY
NLVGEVERVHTNDKVLVLAQRYRERLNLVGSKCAVTLDNPERQQKLADVYEVQNAWLSIRFDD
LDFTLTHVNLNLRKAYNLIPRKHILAFKEYLDNRVQKLCCECRNVRKEDLCTCCSPRYSNL
TSLWKENHSESSIEREAATMMLLDVERKLSFLDLDERKAIIEYKGFIPFSALVKECRLADAGL
CGIRNDVLHDNVIISYADAI GKLSAYFPKEASEAVEYIRRTKEVREQRREELMANSSQ

>7 hypothetical protein [*Prevotella* sp. MA2016]
MSKECKQRQEKRRRLQKANPFSILGTGKHVPGAYFN MARTNFVKTINYILP IAGVRGNYS ENQI
NKMLHALFLIQA GRNEELTTEQKQWEKKLRNLNPEQQT KFKLLFKHFPVLGPMMAADVADHKA
YLNKKSTVQTEDETFAMLKGVS LADCLDIICLMADTLTECRNFYTHKDPYNKPSQLADQYLHQEM
IAKLLDKVVVASRRLKDRREGLSVNEVEFLTGIDHLHQEV LKDEFGNAKVKDGKVMKTFVBYDD
FYFKISGKRLVNGYTVTTKDDKPVNVNMTLPA L SDFGLLYFCVLFSLKPYAKLFI DEVRLEFYS
PFDDKENMIMSEMLSIYRIRTPRLHKIDSHDSKATLAMDI FGLRRCRPMELYNLLDKNAGQFFF
HDEVKHPNSHTPPDVSKRLRYDDRFP TLALRYI DE TELFKRIRFQLQLGSFRYKFDKENCIDGR
VRVRRIQKEINGYGRMQEVADKRMWKWGLDIQKREERSVKLEHEBELYINLDQFLBETADSTPYV
TDRRPAYNIHANRIGLYWEDSQNPKQYKVPDENGMYIPELVVTEDEKKAIPKMPAPRCALSVYDL
PAMLFYELRBOQDNEFP S ABQV IIEYEDDYRKFKA VAEGKLPKPKRKEFRDFLKK EYPKLR
MADI PKLQLFLCSHGLCYNNKPE TVYERLDRLTQLHLEERELHI QNRLEHYQKDRDMI GNKDN
QYGGKSFSDVRHGA LARYLAQSMMEWQPTKLDKKEKGHDKLTGLN YNVL TAYLATYGHQPVEE
GFTPRTLQVLI NAHLI GGSNPHFP INKVLALGNRIEELYLHYLEELKHRSRIQSLSNPS
DKALSALPFIHHRMRYHERTSEEMMALAARYTTIQLPDGLFTPYILEILQKH YTENS DLQNAL
SQDVPVKNLPTCNAAYLITLFYQTVLKD NAQFPYLSDKYTRNKDGEKAESFSFKRAYELPSVL
NNKDKTFPFEMIPLFLTSDEIQERLSAKLLDGDGNVPVEVGEKGPATDSQGN TIWKRRYSE
VDDYAEKLTDRDMKISFSGEWKLP RWKQDKI IKRRETRRQMRDELLQRM PRYIRDI KDNERT
LRRYKQDMVLFLLAEKMF TN IIEQSSEFNWKQMRLSKVCNEAFLRQTLTFRVPVTVGETT IY
VEQENMSLKNYGEFYRFLTDDRMSLLNNI VETLKNENGLDVLIRHTDLMSELAAYDQYRSTIF
MLIQS IENLITNNAVLDDPPADGFVWREDFPKRN FASLLELINQLN NVELTDDRKL LVAIR
NAFSHNSYNIDFSLIKDVKHLPEVAKGILQHLQSM LGVEITK

>8 hypothetical protein [*Riemerella anatipestifer*]
MEKPLLPNVYTLKHKFFWGAFLNIARHNAFITICHINEQLGLKTPSNDDKIVDVCETWNNILN
NDHDLKKSQ LTELILKHPFPLTAMCYHFPKKEGKKGKQKQKEKESQAQSQAEBALNPSKLI

- continued

EALEILVNLQLHSRNYSHYKHKKPAEKDIFKHLKAFDASLRMVKEDYKAHFTVNLTRDFAH
 LNRKGNKQDNPDPNRYRFKEDGFFTESGLLFFFTNLFLDKRDAYWMLKKVSGFKASHKQREKMT
 TEVFCRSRILLLPKLRLESRYDHNQMLLDMLSELSRCPKLLYEKLSSEENKHFQVEADGFLDEIE
 EEQNPFKDTLIRHQDRFPYFALRYLDDLNESFKSIRFQVDLGTYHYCIYDKKIGDEQEKRLHTRT
 LLSFGRLODFTEINRPQEWKALTKDLDYKETSNOFPISKTTPHYHITDNKIGFRLGTSKELYPS
 LEIKDGANRIKAKPYNSGFVAHAFISVHELPLMPYQHLTKGSEDLKKTETVRHIQRIYKDFEEE
 RINTIEDLEKANQGRPLGAPFKQMLGLLQNKQPDLSSEKAKIKIEKLIKLAETKLLSHRLNTKLS
 SPKLGKRRREKLIKTGVLDADWLKDFMRPQVAYDAQNQPISKSSKANSTEFWFIRRALALYGEEK
 NRLEGYFKQTNLIGNTPHPFLNKFNWKACRNLDVDFYQYLEQREKFLKAIKNQWPEPYQYCLL
 LKIPKENRNKLVKGWEGGQISLPRGLFTEAIRETLESDMLSKPIRKEIKKHGRVGFISRAITL
 YFKEYQDKHQSFYNSYKLEAKAPLLKREEHYEYQONKQSPTESORLELHTSDRWKDYLLY
 KRWQHLKELKRLRYRNQDVMLWMLTLELTKNHFKELNLNYHQKLENLAVNVQADAKLNPINQ
 LPMVLPVKVYPATAFGEVQYHKTPIRTVYIREEHTKALKMGNFKALVKDRRLNGLFSFIKEEND
 TQKHPISQLRLRRELEIYQSLRVDAPKETLSLEEKLNKHTLSLSENEFRALLEEWKKEYAAS
 SMVTDEHIAFIASVRNAPCHNQPPYKEALHAPIPLFTVAQPTTEKDGGLGIAEALLKVLREYC
 EIVKSIQI

>9 hypothetical protein [*Prevotella aurantiaca*]

MEDDKTTGSISEYELKDKHFWAAPLNLARHNVIITINHINKLLEIREIDNDEKVLDIKTLWQKG
 NKDLNQKARLRELMTKHFPFLETAIYTKNKEDKKEVKQEQKQAEQSLSESKDCFLFLDKLQEA
 RNYSHYKYSFESKPEPEFEGLLKMYNIFGNNIQLVINDYQHNKIDINPDEDFKHLDRKGGQKY
 SPADNEGNITESGLLFFVSLFLEKKDAIWMQKLNKGFKNLENKMKMTEHVFCSRILMPKLRL
 ESTQTQDWILLDMLNELIRCPKSLYERLQGGDDREKFKVPPDPADEDYNAEQEPFKNTLIRHQDR
 FPFVFLRYFDYNEIFKNLRFQIDLGTYHFSYKIKLIGGQKEDRHLTHKLYGFERIQEFKQNR
 DEWKAVKDLDTYETS NKRYISETTPHYHLENQKIGIRFRNGNKEIWPSPKINDENNEKSKYKL
 DKQYQAEAFVSHVHELPLMPFYLLKKEKPNNDENASIVBGFIKREIRNIFKLYDAPANGEIN
 NIDDLKCYCADKGIKPRHLPKQMVAILYDEHKDMVKEAKRQKEMVKDTKLLATLEKQTKQEK
 EDDGRNVKLLKSGEJARWLVDNMMRFPQVQKDNKPLNNSKANSTYQMLQORSLALYNNEEK
 TRYFRQVNLIESNNPFPKPKWTKWEECNILTFYYSYLTCKIEPLNKLKPEDWKNQYFLKKE
 PKTNRETLVQGWKNGFNLRPGIFTEPIREWFKRHQNSKEYEKVEALDRVGLVTKVIFLFPKKE
 YPKDKEENFKEDTQKEINDCVQPPYFNPNVGNVNIHKPKKEDFLHREERIELWDDKKDKFKGYKE
 KIKSKLTKDKKEFRSYLEPQSWNKPERELRLVRNQDITWLLCKELIDKLKIDELNIEELKK
 LRLNNIDTDTAKKEKNNILNRMVMPMELPVTVEIDDSHKIKVKDKPLHTIYIKEAETKLLKQGNF
 KALVKDRRLNGLFSFVKINSEASKRNPISKLRVEYELGEYQEARIEIQDMLALEEKLINLYK
 DLPINKFSEMNSWLEKDEADKARFQNDVDFLI AVRNAFSHNQYPMHNKIEFANIKPFSLYTA
 NNSEKGLGIANQLKDKTKETTDKIKKIEKPIETKE

>10 hypothetical protein [*Prevotella saccharolytica*]

MEDKPFWAAFPNLARHNVIYLVNHINKLLEKLYDEGKHKIEFEREDI FNISDDVMNDANSNG
 KKRKLDIKKIWDLDLTLTRKYQLRELILKHPPFIQPAIIGAQTERTTIDKDKRSTSTSNDSL
 KQTGEGDINDLSSLNPKMFPRLLEQLRNYSHVKHKSATMPNPFDEDLNMMRYIFIDS
 VNKVKEDYSSNSVIDPNTSFHSHIYKDEQGGIKPCRYPFTSKDGSINAFGLLFFVSLFLEKQDS
 IWMQKKIPGFKKASENYMKMTNEVFCRNHILLPKIRLETVYDKDWMLLDMLEVVRCPLSLYKR
 LTPAAQNKFKVPEKSSDNANRQEDDNPFSRILVRHQNRFPYFVLRFPDLNVEVFTLRFQINLGC
 YHFAICFKQIGDKKVEVHHLIRTLYGFSRLQNFQNTNRPPEWNTLVKTTPEPSGNDGKTQGVPL
 PYISYTIHYQIENKIGIKIPDGD TAVDTDIWP SVSTEKQLNPKDKYTLTPGFKADVFLSVHE
 LLPMMFYQLLCEGMLKT DAGNAVEKVLIDTRNAIFNLADAPVQEKINTITDLENYLQDKPIL
 IGHLPKQIMIDLLKHQRDMLKAVEQKAMLIKDTERLKLKDLKQKQETDVAAKNTGTLLKNGQ
 LADWLVDNMMRFPQVVKRDKGPNINCSKANSTYQMLQRAFAYATDSCRLSRYFTQLHLIHS
 NSHLFLSRFEDKQPNLAFYAAYLKAKLEPLNELQPNWASDNYFLLLRAPKNDRQKLABGK
 NGFNLRPLGFTTEKIKTWNEHKTIVDISDCDIFKNRQVQVARLIPVFPDKKFKDHSQPFYRYDF
 NVGNVSKPTEANYLSKGRRELFKSYQNKFNKNI PAEKTKEYREYKFNFLWKKFERELRIKQ
 DILWLMCKNLFDEKIKPKKDILEPRIAVSYIKLDSLQNTNSTAGSLNALAKVVPMTLAIHDS
 PKPKGKAGNNEKENKEFTVYIKEEGTKLLKGNFKILLADRRIKGLFSYIEHDDIDLKQHPLIK
 RRVLDLELDLYQTCRIDIFQQTGLGAEQLLDKYSDLNLTDFYQMLIGWRKKEGIPRNIKEDTDFL
 KDVRNAFSHNQYPSKKAIFRRIKFNPKELILEEEGLGIATQMYKEVEKVVNRIKRILELFD

>11 hypothetical protein [*Myroides odoratimimus*]

MKDIILTDTTEKQNRFPYSHKIDKYPFGGYFNLASNNIYEVPEEVNKRNTFGKLAKRDNGNLKN
 YIIHVFKDELSDIPEKRVAFASYPFLETVDKKSIKERNRTIDLTLISQIRQPREMLISLVT
 AVDQLRNFYTHYHSDIIVENKVLDFLNSFVSTALHVDKYLKTDKTKFELKETIAAELDILI
 BAYKKQTEKKNTRFKANKREDDI LNAYNEAFWFSINDKDKDKKETVVAAGADAYFEKNHHS
 NDPDFALNISEKGIYVLLSFPFLTNKEMDSLKANLTGFKGKVDRESGNSIKYMATQRIYSFHTYR
 GLKQKIRTSSEGVKETLLMQMIDELSKVPNVVYQHLSTTQONSFI EDWNEYKDYEDDVETDDL
 SRVHPVTRKRYEDRPNYFAIRFLDEDFDPPTLRFPQVHLGDYVHRRTKQLGKVESDRILIEK
 TVFARLKDINSKASYPHSLLEEQDKEELDNKWLFPNPSYDFPKHEHTLQHQGEQKNAGKIGIYV
 KLRDTQYKEKAALBEARKSLNPKERSATKASKYDIIITQIIEANDNVKSEKPLVFTGQPIAYLSM
 NDHSMFLSLLTDNAELKKTPEEVEAKLIDQIGKQINEILSKDITDKILKKYKDNLDLKETDITDK
 ITRDLARDKEEIKLILEQKQRADDYNYTSSTKFNIDKSRKRRKHLLFAEKGKIGVWLANDIKR
 FMPKESKSKWGYQHTLQKLFAYPDTSKSDELEILSNMVMVKDYPIELIDLKKSRTLVDLFLN
 KYLEARLEYIENVI TRVNSIGTPQFKTVRKECFPLKKSNTYVSLDKQVERILSMPLFIERG
 FMDDKPTMLEGKSYKQHKKFAWVHYKENSNYQNFYDTEVYEITTEDKREKAKVTKKIKQQQ
 KNDVFTLMMVNYMLEEVLKLSNDRLSLNLQYQTKERIVNKQVAKDQERNKNIYIWNKVVLDLQ
 LCDGLVHIDNVLKDLIGNFRKYENDSRVKEPLTYQSDI VWSAYLSNEVDNKLVIYERQLDNDYE
 SIRSKELLKEVQIECSVYVQVANKESLQSGNENFKQYVQLGGLPIGMDVREMLILSDTVKFK
 KEEIQGLGQAGEVEQDLYSLIYIRNKFAHNQLPIKEFPDFCENNYRSISDNEYAYEYMEIFRS
 I KEKYN

- continued

>12 hypothetical protein [*Flavobacterium columnare*]
 MSSKNESYNKQKTFNHYKQEDKYFPGGFLNNADDNLRQVKGKFKTRINFNHNHNNELASVFKDYF
 NKEKSVAKREHALNLLSNYPVLERIQKHTNHNFBQTRIEIFELLLDITIKKLRDYTHHYHKPIT
 INPKIYDFLDDTLDDVLITIKKKVKNDTSRELLKEKLRPELTQLKNQKREELIKKGGKLEEN
 LENAVFNHCLRPFLLEENKTDKQNKTVSLRKYRKSHPNEETSITLQSGLVFLMSFFLHRKEFQ
 VFTSGLGFKAKVNTIKEEEISLNKNNIVYMI THWSYSYNNFKGLKHRIKTDQGVSTLEQNNTT
 HSLTNTNTKEALLTQIVDYLKVPNEIYETLSEKQKQEFEEIDNEYMRENPENEDSTFSSIVSH
 KVIRKRYENKPNYFAMRFLDEYAELPTLRFMVNFVDYIKDRQKKILESIPQDSERI IKKEIHLF
 EKLSLVTEYKKNVYLKETSNIIDLSRFPFLPNPSYVMANNNI PFYIDRSNNLDEVLNOKKKAQS
 QNKKRNLTFEKYNKQSKDAI IAMLQKEIGVKDLQQRSTIGLLSCNELPMSLYEVIKVDIKGAE
 LENKIAQKIREQYQSIRDFTLDSPOKDNIPPTLIKTINTDSSVTFENQPIDIPRLKNAIQKELT
 LTQEKLLNVKEHEIEVDNYNRNKNTYKPKNQPKIKVDDKQLQKRYVYRNEIRQEANWLASDLI
 HPMKNKSLWKGYMHNELQSFLAFFEDKKNDCIALLETVFNLKEDCILTKGLKNLFLKHGNFIDF
 YKEYLKKLEDPLNTESTPLENGLIGLPPKIKKELSKRFKYYIFIVFQKRQPIIKELBEEKNNLY
 ADAINLSRGI FDEKPTMIPFKKPNPDEFASWVVASVYQYNNYQSFYELTPDIVERKKKKYYKNLR
 AINKVKIQDYLLKMDVTLYQDLFNQPLDKSLSDFFVSKAEREKIKADAKAYQKRNDSLWNV
 IHLSQLNRI TANPKLKI GYKRALQDEKIATLLTYDDRWTYALQKPEKENENDYKELHYTA
 LNMELQYEVKVRSEKELLKQVQELEKQILEEYTDFLSTQIHPADFEREGNPNFKKYLHSAI LENE
 DDLKDLPEKVEAMRELDETITNPIIKKAIIVLIIIRNKMAHNQYPPKFIYDLANRFVPKKEEYF
 ATYFNRFVETITKELWENKEKDKTQV

>13 hypothetical protein [*Bacteroides pyogenes*]
 MALLNVRKVENHIRKWLGDVALLPEKSGFHSLLT TDNLSSAKWTRFYKSRKFLPFLEMFSDSK
 KSYENRRTECLDITIDRQKISLLKVEYVGLQDIRNAFVSHYHIDDQSVKHTALIISSEMHRFI
 ENAYSFALQKTRARFTGVFVETDFLQAEKGDNKKFFAIGGNEGILKDNALIPLICLFLDREE
 AFKFLSRATGPKSTKEKGLAVRETFALCCRQPHERLLSVNPREALMDMLNLRCPDILFE
 MLDEKDKSFLPLLGEEQAHILENSLNDELCEAIDDPFEMIASLSKRVRYKNRFPYMLRYIE
 EKNLLPPIRFRIDLGLCELASYPKMGEBENNYERSVTDHAMAAGRLTDFHNEAVLQQITKGT
 DEVRFSLYAPRYAIYNNKIGPVRTGSGDKISFPTLKKKGGEGHCVAYTLQNTKSGFISIVDLR
 KILLLSFLDKKAKNIVSGLLEQCEKHWDLS ENLFDARTLQKQEFVPLIRYTLPRSKGKGL
 VSSKLADKQEKYSEFERRKEKLEILSEKDFDLSQIPRRMIDEWLNVLPTSRKLLKGYVETL
 KLDRCERLRVFEKREKGRVPPRIGEMATDLAKDIRMVIDQGVKQRI TSAYYSEIQRCLAQY
 AGDDNRRHLSI IRELRKLDTKNGHPFLGKVLRLPGLGHTEKLYQRYFEKKEWLEATFPYASP
 KRVPFRVNPPTGKQKELPLIIRNLMKERPEWRDWRKQKNSHPIDLPSQLFENEICRLKDKIGK
 EPSGKLNEMFKLYWDKEFPNGMQRFYRCKRRVEVPDKVVEYSEEGGNYKYYEALIDEVV
 RQKISSSKEKSKLQVEDLTVSVRRVFKRAINEKEYQLRLLCEDDRLLFMVARDLYDWKEAQLDL
 DKIDNMLGEPVSVSQVILEGGQDAVIAEKCKLDVSKLMRYCYDGRVKGMLMPYFANHEATQE
 QVEMELRHYEDHRRVFNWFALEKSVLKNKELRRFYEESQGGCEHRCIDALRKASVSEEEY
 EPLVHIRNKSAAHQFPDLBIGKLPNVISGFCCEIWSKYKAIICRIIPPIDPERRFPGLLEBQK

>14 hypothetical protein [*Capnocytophaga canimorsus*]
 MKNIQRLGKGNFSPFKKEDKFFYGGFLNLANNNIEDFFKEIITRFGIVITDENKKPKETFGKE
 ILNEIFKKDISIVDYEKVWNI FADYFPFTKYLKSLYLEEMQPKNRVICFRDVMKELLKTVEALRN
 FYTHYDHEPIKIEDRVFYFLDKVLLDVS LTVKNKYLKTDKTKPEFLNQHI GEELKELCKQRKDYL
 VGKGRIDKESI INGIYNAFKDFICKREKQDDKENHNSVEKILCNKEPQNKQKSSATVWEL
 CSKSSSKYTEKSPNRENDKHCLEVPISQKGI VFLSFLNKGIEIYALTSNIKGFKAKITKEEP
 VTYDKNSIRYMATHRMFSFLAYKGLKRRKIRTS EINYNEDEGQASSTYKETLMLQMLDELNKVPD
 VVYQNLSEBDVQKTFIEDWNEYLKENNGDVGTMEEBQVIHPVIRKRYEDKFNFAIRFLDEFAFP
 PTLRFQVHLGNYLCKDKRTKQICDTTTEREVKKITVFGRLSELENKKAIFLNEREEIKGWEVFP
 NPSYDPKKEISVNYKDFP VIGSILDRBKQPVSNKIGIRVKIADLQREIDKAIKELKLRNPKN
 RKANQDEKQKERLVNEIVSTNSNEQGEPPVFIQPTAYLSMNDIHSVLYEFLINKISGEALET
 IVEKIEIQIKQIIGKDATTKIKPYTNANSINREKLLRDLQEQQILKTLLEEQQRKDKK
 DKSKRKHLELVPSEKGVAVWLANDIKRFMPKAPKEQWRGVHSSLLQKYLAYYEQSKEELKNLL
 PKEVFKHFPFKLKGYPQQYLNQFYTDYLRRLSVNELLNINQFNKNDKDALKATEKPKFF
 RKQNYIINPINIQIQSILVYPIFLKRGFLDEKPTMIDREKPKENKDTELADWFMHYKNYKEDNY
 QKPYAYPLEKVEEKEKFKRNKQINKQKKNVYTLMMVEYI IQKIFGDKFVEENPLVLKGI PQSK
 AERQQNNTHAATQERNLNGILNQPKDIKIQGKITVKGVKLKDIGNFRKYEIDQRVNTFLDYEP
 RKEWMAYLPNDWKEKEKQGLPPNNVIDRQISKYETVRSKILLKDVQLEKIIISDEIKEBHRHD
 LKQGYYNFKYIILNGLLRQLKNENVENYKVPKLNINPEKVINI TQLKQBATDLBQKAFVLTYIR
 NKFAHNQLPKKEFDWYDQCKYKGI EKEKTYAEYFAEVFKREKEALIK

>15 hypothetical protein [*Capnocytophaga cynodegmi*]
 MENKTSLGNNIYYPFKPQDKSYFAGYLNAAENIDSVFRELGKRLKKGKEYTS ENFPDAIFKEN
 ISLVEYERYVKKLSDYFPMARLLDKKEVPIKERKENFKKNFRGIIKAVRDLRNFYTHKEHGEVE
 ITDEIFGVLDLMLKSTVLTVKKKIKTDKTKELKKSIEKQLDILCQKLELYLKDARKI BEKR
 RNQREBERGKELVPRFEYSDRRDDLIAIYNDAFDVYIDKDKDSLKESSTKNTYTESYQQEEDG
 LKIPISKNGVFLSFLSKQEVHAFKSKIAGFKATVIDEATVSHRKNISICFMATHEIFSHLAY
 KKLKRVTAETINYSABENABQLS IYAKETLMMQMLDELKVPDVVYQNLSEBDVQKTFIEDWNE
 YLKENNGDVGTMEEBQVIHPVIRKRYEDKFNFAIRFLDEFAQFPPTLRFQVHLGNYLHDSRPKE
 HLISSDRRIKEKITVFGRLSELEHKKALFIKNTETNEDRKHYWEVFPNPNYDPPKENISVNDKDF
 PIAGSILDRKQPTAGKIGIKVNLNQQYI SEVDKAVKAHQKQRNPNKPSIQNIIEEIVPINGS
 NPKIIVFGGQPTAYLSMNDIHSILYEFDFKWEKKEKLEKKEKELRKEIKGEEELIVGKIQ
 TQIQQIIDKIDINAKILKPYQDDSTAI DKEKLIKDLKQEQKILQKLNQBTAREKEYQECIAYQ
 EESRKIKRSDKSRQYLRNQLKRYPEVPTREKIELYQKQKGVAVWLANDIKRFMPPTDFKNEWK
 GEQHSLLQKSLAYYEQCKEELKNLLPQQKVPFKLPELGGHQQKYLQYQFTRYLDKRLHEISG
 LVQQAENFKNENKVPKVENECFKFLKQNYTHKGLDAQASVGLGYPFLFERGFMDKPTI IKG
 KTFKGNESLFTDWFYRYKQYQNFQTFYDTEYPLVLEKQKQADRRETKIYQQKNDVFTLLMA
 KHIKSVFKQDSIDRFSLELDYQSRERLENQEKAKQTGERNTNYIWNKTVDLNLCGDKVTVEN

- continued

VKLKNGVNFIKYEQDQVQTPFLKYEENIKWQAFLIKESKEENYPYIVERETEQYKVRREELL
KEVHLIEEYI LEKVVDKKEI LKKGDNQNFKYI LNLGLKQLKNEVDVESYKVPNLNTKPEDVNIQ
LKQEATDLBQKAFVLTIRNKFAHNQLPKKEFDWYCYQEKYGI EKEKTYAEYFAEVFKREKAL
MK

>16 hypothetical protein [*Chryseobacterium* sp. YR477]

METQTI GHGIAYDHSKI QDKHFFGGFLNLAENNI KAVLKAFSEKFNVDVVKQFADVSLKDNL
PDNDFQKRVSPFLKMYFPVVDV INI PNNRAKFRSDLTTLFKSVDQLRNFYTHYHKKPLDFDASLF
ILLDDI FARTAKEVDRDQMKDDKTRQLLSKSLSEELQKGYELQLERLKELNRLGKKNVNIHQLG
IKNGVLNNAFNHLYKDGESFKTKLYSSALTSFESAENGIETISQSGLLFLLSMFLKRKEIEDL
KNRNKGFPAKVVIVDEGKRVNGLKPMATHWVPSYLCFKGLKSKLSTEFHEETLLIQI IDELSKVP
DELYCAFDFKTRDKFIEDINEYVKEGHQDFLEDAKVIHPVIRKRYENKNFYFAIRFLDFEVPKF
PSLRQVHVGNVYVHRRIKNIDGTTFTETERVVKDRIKVFGRLESEISSYKAQYLVSSVSKDHDETG
WEIFPNPSYVFINNIP IHSVDTSEFKKEIADFKLRRRAQVPELKI RGAEKRRKFEI TQMI GS
KSVLNQEEPI ALLS LNEIPALLYEILINGKEPAEIERI IKDKLNERQDVIKNYNPENWLPASQI
SRRLRSNKGERI INTDKLLQLVTKELLVTEQKLI I SDNREALKQKKEGKYIRKFI FTNSELGR
EAIWLADDI KRMPADVVRKEWKYQHSQQLQSSLAFYNSRPKEALAI LSSWNLDKEKI IWNWEI
LKSFTQNKFFDAFYNEYLKGKRYFAFLSEHIVQYTSNAKNLQKFIKQOMPDLFEKRHYI IED
LQTEKNKI LSKPFI PPRGIFDKKPTFIKGVKVEDSPESFANWYQYGYQKDHQFQKFYDWRDYS
DVFLHGLKAPINNGDRRTLGMEELKERI I IKQDLKIKKI KIQDLFLRLIAENLFOKVPKYSAL
LPLSDFYLTQEERMEKENMAALQNVREEDGKSPNI I KDNFIWSKMIPIYKKGQI IENAVLKDIG
KLNVLSDDDKVTLLSYDDAKPWSKIALENFSGENSYEVIRREKLPKEIQQFSEI LFRSGW
DGINHPAQLEDMRNPKFKMYI VNGILRKSAGLYSQGEDIWFEYNADFNLDADVLETKSELVQL
AFLVTAIRNKFAHNQLPAKEFYFYIRAKYGFADEPSVALVYLNFTKYAINEFKKVI

>17 hypothetical protein [*Phaeodactylibacter xiamenensis*]

MTNTPKRRTLHRHPSYGFAGFLNI ARHNAFMIHEHLSTKYDMEDKNTLDEAQLPNAKLPGLCKKR
YKPKDVTGVSRLDRRYFPFLNPLFLHLEKQNAEQAAATYDINPEDIEFTLKGFFRLLNQMRN
NYSHYI SNTDYGKPKLVPQDIYEAAI FRLLDRGKHKRFDVFEKSKHTRHLESNNSEYRPRSLA
NSPDHENTVAVFTCLFLERKYAFPFLSRLDCFRS TNDAAEGDPLIRKASHECYTMFCRRLPQPK
LESSDILLDMVNELGRCPSALYNLLSEEDQARFHIKREETGFEEDPDEELEQEI VLKRHSDRF
PYFALRYDDTEAFQTLRFDVYLGWRWTKPVYKRI YGQERDVL TQSI RTFTRLSRLLPYEN
VKHDAVRQNEEDGKLVNPDVTSQPHKSWIQIESDDRFLSRI EHFSPHYNFGDQVIGLKFINP
DRYAAIQNVFPKLPGEKKDKDAKLVNETADAI I STHEIRSLFLYHYLSKKPI SAGDERRFIQV
DTETPIKQYIDTIKLFPEEDIKSGELQPIADPNYQKNEPLPYVRGDKEKTEERAQYRERQKEI
KERRKELNLTQLNRVGLSIQYIPSRLEYYLGYKVPYKELALQKLRQKKEVKKRIKDI EKMR
TPRVGEQATWLAEDIVFLIPPKMHTPERKTTKHPQKLNNDQFRIMQSSLAYFSVNKKAIKKFFQ
KETGIGLSNRSTSHFPLRIDVGRCRGILDPTGTGLYKMDWLDADAI KKVDRKHGKKEAKKYE
KYLPSIQHKPTLELDYTRLPVYLPRLGFKKAIKALAAHADQVPEEDNVI FCLDQLLDGDT
QDFYNWQRYRSALTEKETDNQLVLAHPYAEQILGTIKTLEGKQKNNKLGNAKAKIKDELIDL
KRAKRRLLDRBQYLRAVQAEDRALWLMIQERQKQKAEHEEIAFDQLDLKNI TKILTESIDARLR
IPDTKVDITDKLPLRRYGDLRRVAKDRRLVNLASYYHVAGLSEIPYDLVKKLEEYDRRRVAF
EHVYQFEKEVYDRYAAELRNENPKGESTYFSHWEYVAVVAKHSADTHFNELFKEKVMQLRNKPH
HNEFPYFDWLEPEVEKASAAIYADRVFVABEGYYQKMRKLMRQ

>18 hypothetical protein [*Porphyromonas gingivalis*]

MTEQNEKPYNGTYTLEDKHFWAFFNLARHNAYITLTHIDRQLAYSKADI TNDIEDLFFKQW
KNLNDNLERKARLRSILKHF SFLBGAAYGKLFESQS SGNKS KKKELTKKEKEELQANALS
DNLKSI LFDPLQKLDKFRNYSHYRHPESSELPLFDGNMLQRLYNVFDVSVQRVDRDHEHNDK
DPHRHFNHLVRKGGKDRCGNNDNPFKHHFVDREGKVTAEGLLFPVSLFLEKRDIAIWQKKIRG
FKGGTETYQMTNEVFCRSRI SLPKLKLES LRTDDWMLDMLNELVRCPKSLYDRLEEDRACF
RVNVDILSDEDDTDGAEDPFPKNTLVRRHQDRFPYFALRYFDLKKVFTSLRPHIDLGTYHFAIYK
KNI GEQPEDRHLTRNLVYGFGRIQDFAEHRPEEWKRLVRDLDCFETGDKPYITQTPPHYIEKG
KIGLRFVFEQHLWSPVEGATRTGRSKYAQDKRLTAEAFLSVHELMPMPYFLLREKYSDEA
SAERVQGRIKRVIDEYVAVYDAFARGEIDTLDRLDACLADKGI RRGHLPRQIIA ILSQEHKME
EKVRKQLQEMIA DTDHRLMDLDRQTDRIKIRIGRKNAGLPKSGVIADWLVRDMRFPQVAKDTS
KPLNNSKANS TEYRMLQRALALFGGEEKERLTPYFRQMNLTGGNNPHPLHETRWESHNTNLSFY
RSYLKARKAFLQSIGRSDRVENHRFLLLKPKTDRQTLVAGWKGFHLPRIETAVRDCLEIM
GHDEVASYKEVGFMAKAVPLYFERACKDRVQPFYDYPFNVGNLSLKPCKGRPLSKEKRAEWEESG
KERFRDLEAWSHSAARRIEDAFVGI EYASWENKKKI EQLLQDLSLWETFESKLKVKADKINI AK
LKKEILEAKEHPHYDFKSWQFERELRLVKNQDI ITWMMCRDLMEENKVEGLDGTGLYKLDIRT
DVHEGQSLNVLNRVPMKRLPVVVYRADSRGHVHKEEAPLATVYIEERDTKLLKQGNFKSPVKDR
RLNGLFSFVDTGGLAMEQYPI SKLVEYELAKYQTARVCAFEQTLLEESLTRYPHLPDKNFR
KMLSESWDPLLDKWPDLHRKVRLLI AVRNAPSHNQYPMYDEAVFSSIRKYDPSLDAI EERMGL
NIAHRLSEEVQAKETVERIIQA

>19 hypothetical protein [*Porphyromonas gulae*]

MNTVPATENKQQRTRVEDDPQYFGLYLNLARENLEIEVESHVRIKFGKKNLNEESLKQSLCDHL
LSIDRWTKVYGHRSRYLPLFLHCFDPDSGIEKDHD SKTGVDPSAQRLIRELYSLDLFLRNDFSH
NRLDGTTFEHLKVSPISSFI TGAYTFACERAQSRFADFFKPDFFLLAKNRKEQLISVADGKEC
LTVSGFAFFI CLFLDREQASGMLSRIRGFKRITDENWARAVHETFCDLICIRHPHDLRESNTKEA
LLDMLNELNRCPRILYDMLPEEERAQFLPALDENSMNLSSENSLNEESRLWDGSDWAELT
KRIRHQDRFPYMLRFIEEMDLLKGI RFRVLDGEIELDSYKVKVGRNGEYDRTITDHALAFGLM
SDFQNEEVSRMISGEASYPVRFSLFAPRYAIYDNKIGYCHTSDPVYPKSKTGEKRALSNPQSM
GFISVHDLKLLMLLCEGFSRMSQDFLRKANRLIDETAEGKLQFSAFPFEMRHRFPQNP
KSKDRREKAETTLKQYKIKGRKDKLNSQLLSAFDMNQRLPSRLLEDWMMNIRPASHSVKLR
YVKQLNEDCRRLRFRKFDGDKARAIPLVGEMATFLSQDIVRMI ISEBTKLITSAYNEMQR
SLAQYAGEENRRQFRAI VAEHLHLLDPSSGHPFLSATMETAHRYTDEDFYKYLEKKREWLAFTFY

- continued

RPEQDENTKRRI SVFFVPDGEARKLLPLLI RRRMKEQNDLQDWIRNKQAHPIDLPSHLFDSKIM
 ELLKVKDGGKKWNEAFKDWSTKYPDGMQPPYGLRRELNIHGKSVSYIPSDGKKFADCYTHLME
 KTVRDKKRELRTAGKPPDLAAYIKRSFHRVNEREFMLRLVQEDDRMLMAINKMMTDREED
 ILPGLKNIDSILDEENQFSLAVHAKVLEKEGEGGDNLSLVPATIEIKSKRKDWSKYIRYRYDR
 RVPGLMSHFPEHKATLDEVKILLGEYDRCRIKIFDWAFALEGAIMSDRDLKPYLHSSSREGKS
 GEHSTLVKMLVEKKGCLTPDESQYLILIRNKAAHNQFP CAEMPLIYRDVSAKVGSI EGSSAKD
 LPEGSSLVDSLWKYEMIRKILPLLDHENRFPGKLLNMSQPINDL

>20 hypothetical protein [*Porphyromonas gulae*]

MTEQSERPYNGTYTLEDKHFWAAPLNLARHNAYITLTHIDRQLAYSKADI TNDQDVLSFKALW
 KNFDNDLERKSRRLRSLILKHFSLFEGAAYGKLFESKSGNKS SKNKELTKKEKEELQANALSL
 DNLKSLPDLFQKLDKDFRNYSHYRHS GSSLELFDGNMLQRLYNVFDVSVQRVKIDHEHNDEV
 DPHYHFNHLVRGKGDYGHNDNPSFKHHFVDGEGMVTEAGLLPFVSLFLEKRDALWQKKIRG
 PKGGTETYQQMTNEVFCRSRISLPKLLKLES LRMDWMLLDMLNELVRCPKPLYDRLEDDRACF
 RVPVDILPEDEDDTGGGEDPDKNTLVRHQDRFPYFALRYFDLKKVFTSLRPHIDLGTYHFAIYK
 KMIGEOPEDRHLTRNLVGFGRIQDFAEEHRPEEWKRLVRDLDFETGDKPYISQTS PHYHIEKG
 KIGLRFMPGQHLWSPSEVGTTRTGRSKYAQDKRLTAEAFLSVHELMPPMPYFLLREKYSSEV
 SAERVQGR IKRVIEDVYAVYDAFARDEINTRDELDA CLADKGI RRGHLPQMIATLSQEHKDME
 EKIRKKLQEMMADTDHRLDMLDRQTRKIRIGRKNAGLPKSGVIADWLVRDMMRFPQVAKDASG
 KPLNNSKANS TEYRMLQRALALFGGKERLTPYFRQMNLTGGNNPFLHETRWESHNTNISFY
 RSYLRARAFLEIRGRSDRVENRPFLLKPKTRQTLVAGWKGEFHLPRGIFTEAVRDCLLEM
 GHDEVASYKEVGFMAKAVPLYFERACEDRVQPFYDS PFNVGNSLKPKKGRFLSKEERAEWERG
 KERFRDLEAWSYSAARRIEDAFAGIEYASPGNKKKIQLLRDL SLWEAFESKLVKVRADRI NLAK
 LKKEILEAQEHPYHDFKSWQKFERELRLVKNQDIITWMMCRDLMEENKVEGLDTGLTYLKD IRP
 NVQEQGSLNVLNRPVVPVYRADS RGHVHKEEAPLATVYIEERDTKLLKQGNFKS FVKDR
 RLNGLFSVDVTGGLAMEQYPI SKLRVEYELAKYQ TARVCFELTLRLEESLSTRYPHLPDESFR
 EMLESWSDPLLAKWPELHGKVRLLIAVRNAPSHNQYPMYDEAVFSSIRKYDPS SPDAI EERMGL
 NIAHRLSBEVKQAKETVERIIQA

>21 hypothetical protein [*Porphyromonas gulae*]

MTEQSERPYNGTYTLEDKHFWAAPLNLARHNAYITLTHIDRQLAYSKADI TNDQDVLSFKALW
 KNLDNDLERKSRRLRSLILKHFSLFEGAAYGKLFESKSGNKS SKNKELTKKEKEELQANALSL
 DNLKSLPDLFQKLDKDFRNYSHYRHS GSSLELFDGNMLQRLYNVFDVSVQRVKRDHEHNKDV
 DPHRHFHNLVRGKGDYGHNDNPSFKHHFVDGEGMVTEAGLLPFVSLFLEKRDALWQKKIRG
 PKGGTETYQQMTNEVFCRSRISLPKLLKLES LRDDWMLLDMLNELVRCPKPLYDRLEKDRARF
 RVPVDILPEDEDDTGGGEDPDKNTLVRHQDRFPYFALRYFDLKKVFTSLRPHIDLGTYHFAIYK
 KVIGEOPEDRHLTRNLVGFGRIQDFAEEHRPEEWKRLVRDLDFETGDKPYISQTT PHYHIEKG
 KIGLRFVPEGQHLWSPSEVGTTRTGRSKYAQDKRLTAEAFLSVHELMPPMPYFLLREKYSSEV
 SAEKVQGR IKRVIEDVYAIYDAFARDEINTRDELDA CLADKGI RRGHLPKQMI GILSQEHKNME
 EKVRKKLQEM IADTDHRLDMLDRQTRKIRIGRKNAGLPKSGVIADWLVRDMMRFPQVAKDTSG
 KPLNNSKANS TEYRMLQRALALFGGKERLTPYFRQMNLTGGNNPFLDETRWESHNTNISFY
 RSYLRARAFLEIRGRSDRVENRPFLLKPKTRQTLVAGWKSEFHLPRGIFTEAVRDCLLEM
 GYDEVGSYKEVGFMAKAVPLYFERACKDRVQPFYDS PFNVGNSLKPKKGRFLSKEKRAEWEWG
 KERFRLEAKLKEILEAQEHPYHDFKSWQKFERELRLVKNQDIITWMMCRDLMEENKVEGLDTGT
 LYLDKIRPNVQEQGSLNVLNRPVVPVYRADS RGHVHKEEAPLATVYIEERDTKLLKQGN
 FKS FVKDRRLNGLFSVDVTGGLAMEQYPI SKLRVEYELAKYQ TARVCFELTLRLEESLSTRYP
 HLPDES FRELESWSDPLLAKWPELHGKVRLLIAVRNAPSHNQYPMYDEAVFSSIRKYDPS SP
 AIEERMGLNIAHRLSBEVKQAKETVERIIQA

>22 hypothetical protein [*Prevotella falsenii*]

MKNDNNTKSTDTYTLGDKHFWAAPLNLARHNAYITVNHINKVLELKNKKDQEI IDNDQDILAI
 KTLWGKVDTDINKDRLELIMKHFPPLEAATYQQSSNTNTKQEEQAKAQSFESLKDCLFLF
 LEKLEARNYSHYKHSKSLKEEPKLEKLEENMYNI FDTNVQLV IKDYEHNKDINPEEDFKHLG
 RAEGEFNYFTRNKKNITESGLLFFVSLFLEKKDAIWAQTKIKGFKDNRENKQKMTHEVFCRS
 RMLLPKLRLESTQTQDWILLDMLNELIRCPKSLYKRLQGEKREKFRVFPDPAEDYDABEQEPFK
 NTLVRHQDRFPYFALRYFDYNEIFTNLRFQIDLGTYHFSIYKQIGDKKEDRHLTHKLYGFERI
 QEFAKENRPDEWALVKDLDTFEESNEPYISETTPHYHLENQKIGIRNKNK KKKKT IWPSLET
 TIVNERSKYNLGKSFKAFAFLSVHELLPMMFYLLLNKEEPNNGKINASKVEGIEIEKKIRDIYK
 LYGAFANEEINNEELKEYCEGKDAIRHLPKQMIATILKNBYKDMAKKAEDKQKKMIKDTKKRL
 AALDKQVKGEVEDGGRNIKPLKSGRIASWLVNDMMRFPQVQRDRDGYPLNNSKANS TEYQLLQ
 TLALFGSERERLAPYFRQMNLTGKDNPHPLKDTKWEHNNILSFYRSYLEAKNLFGLSKP
 WKINQYFLKLEKPKTRNETLVQGWKNGFNLPRGIFTEP IREWFIRHQNESEYKVKDFDRI GL
 VAKVILPFFKEDYQKEIEDYVQPFYGYPFNVGNIHNSQEGTFLNKKEREELWKGKNTKFKDYKT
 KEKNEKTKNDKFKKKTDEBEKEFRSYLDFQSWKKFERELRLVRNQD IVTWLLCMELIDKLIKID
 ELNIEELQKRLKIDTDITAKKEKNNILNRIMPPELPTVYETDSDNIIKDKPLHTIYIEKEA
 TKLLKQGNFKALVKDRRLNGLFSFVETSSAEELKSKPI SKLVEYELGEYQARVEI IKDMLRL
 BETLIGNDEKLP TNKFRQMLKWLHKKETDDTLKNDVKLLTEVRNAPSHNQYPMRDRIFA
 NIKPFSLSANTSNNEGLGIAKLLKDKTKETIDRI TEIEEQATKR

>23 hypothetical protein [*Prevotella intermedia*]

MEDDKKTKESTNMLDNKHFWAAPLNLARHNAYITVNHINKVLELKNKKDQDI IDNDQDILAIK
 THWEKVNGLNKTERLERLELTKHFPFLETAIYTKNKEDKEEVKQEKQAKAQSFDSLKHCLFLFL
 EKLQEARNYSHYKSESTKEPMLKELLLKMYNIFDDNIQLVIKDYQHNDINPEEDFKHLDR
 TEEEFNYFTTNKKNITASGLLFPVSLFLEKDAIWMQKLRGFKDNRESKKNMTHVFCRSR
 MLLPKLRLESQTQDWILLDMLNELIRCPKSLYERLQGEYRKKFNVPD SADEYDABEQEPFK
 TLVRHQDRFPYFALRYFDYNEIFTNLRFQIDLGTYHFSIYKLLIGGQKEDRHLTHKLYGFERIQ
 EFAKQNRDDEWKAIVKDPDYETSEEPYISETAPHYHLENQKIGIRFRNDNDEIWP SLKTNGEN
 NEKRYKLDKQYQAEAFLSVHELLPMMFYLLLNKKEEPNNDKKNASIVEGFIKREIRDIYKLYD

- continued

DRFPYFALRYFDYNEIFKNLRFQIDLGYHFSIYKLLIGDNKEDRHLTHKLYGFERIQEFAKQK
 RPNEWQALVKDLDIYETSNEQYISETTPHYHLENQKIGIRPKNKDKIWPSETNGKENEKSKY
 NLDKSFQAEAPLSIHELPLMMFYDLLLKKKEFPNDEKNASIVVGF IKKEIKRMYAIYDAFANEE
 INSKGLEBEYCKNKGFQERHLPKQMIAILTNKSNMAEKAKRQKQKEMIKDTPKRLATLQKQVKG
 EIEDGGRRNIRLLKSGEIAARLVNDMMRFQSVQKDKGKPLNNSKANSTEYQMLQRSALALYNKEQ
 KPTPYF IQVNLKLSNPHFPLEETKWECCNNILSFYRSYLEAKKNFLESKPEDWKKQYPLML
 KEPKTRNKTIVQGWKNGFNLPRIETPIKEWFKRHQNDSEYKVEALDRVGLVAKVIPLFFK
 EYFPKEDAQKEINNCVQPFYSPYVNGNIHKPEEKNFLHCEERRKLWDKDKKDFKGYKAKEKSK
 KMTDKKEEHRSYLEFQSWNKFERELRLVRNQDITVWLLCTELIDKLDLNI EELQKRLKLD
 IDTD TAKKEKNILNRI PMPQLPVTVEIDKSFNIVKDKPLHTIYI EETGTLLKQGNFKALVK
 DRRLNGLFSFKTSSEAESKSKPISKLRVEYELGAYQKARIDI IKDMLALEKTLIDNDENLPTN
 KFSMDLKS WLKKGKGEANKARLQNDVLLVAIRNAPSHNQYPMYNSVFKGMKLLSLSSDIPEKE
 GLGI AQKDKKI KETIERI IETEKIERN

>28 hypothetical protein [Prevotella saccharolytica]
 MEKENVQGSIIYEPDKCFWAAFYNLARHNAYLTAHINSFVNSKKGINDDKVLDIIDDWSK
 PDNDLLMGARLNKLLKHFPLKAPLYQLAKRKRKQOQKQDYKQDEDEPEVIQBAIANAF
 KMANVRKTLHAFLEQLDRLNHFSHYNYNSPAKMEVDFDDGFCNKL YVVDAAQVMVQDDNRM
 NPEINMQTDFEHLVRLGRNRKIPNTFKYNTNSDGTINNGLLFFVSLFLEKRDALWMQKKIKG
 PKGGTENYMRMTNEVFCNRNRMVIPKLRLETDYDNHQLMFDMLNELVRCPLSLYKRLKQEDQDKF
 RVP IEFLEEDNEADNPYQENANSDENPTEETDPLKNTLVRHQHFRFPYVFLRYFDLNEVFKQRF
 QINLGCYHFSIYDKTIGERTEKRHLTRTLFGFDRLQNF SVKLQPEHWKNMVKHLDT EESSDKPY
 LSDAMPHYQIENEKIGIHFLKTDTEKKETVWPSLEVEEVSNNRKYKSEKNLTADAPLSTHELL
 PMMFYQLLSSEKTRAAAGDKVQGVLSYRKKIFDIYDDFANGTINSMQKLDERLAKDNLRLG
 NMPQQLALILEHQEPDMEQKAKEKLDRLITETKKRIGKLEDQFKQKVRIGKRRADLPKVGSIAD
 WLVDMMRFQPAKRADNITGVPDSKANSTEYRLLQALAFYSAYKDRLEPYFRQVNLIGGTNPH
 PFLHRVDWKKCNHLLSFYHYDLEAKEQYLSHLSPADWQKHQHFLLLVKVRKDIQNEKDWKKS LV
 AGWKNGNPLPGLFTESIKTWFS TDADKQVITDTKLPENRVGLIAKLIPLYYDKVYNDKQPQPFY
 QYPFNIINDRYKPEDTRKRF TAASSKLWNEKMLYKNAQPDSSDKI EYPOYLDLFLSWKKLERELR
 MLRNQDMMVWLMCKDLFAQCTVEGVFADLKLSQLLEVVDVNVQDNLNVLNNVSSMILPLSVYPSD
 AQGNVLRNSKPLHTVYVQENNTKLLKQGNFKSLLDKRRNLNGLFSFIAAEGEDLQQHPLTKNRLE
 YELSIYQTMRSVFEQTLQLEKAILTRNKTL CGNFNLLNSWSEHRD KKTLPQDIDFLIAVR
 NAFSHNQYPMSTNTVMQGI EKFNITQTPKLEEKDGLGIASQLAKTKDAASRLQNIINGGTN

>29 hypothetical protein [Prevotella sp. P5-119]
 MNIPALVENQKKYFGTYSVMAMLNQAVLDHIQKVADI EGEQENENNENLWFPVMSHLYNAKNG
 YDKQPEKTMFIIERLQSYFPFLKIMAEHQREYSNGKYKQNRVEVNSNDI FEVLKRAPGVLKMYR
 DLTNHYKTYEKKLDIGCEFLTSTEQPLSGMISKYITVALRNTKERYGYKTEDLAFIQDNIKKIT
 KDAYGKRKSQVNTGFFLSLQD YNGDTQKKLHLSGVGIALLICFLDKQYINIFLSRLP I FSSYN
 AQSEERRI IIRSFNGINSIKLPKDR IHS EKSNSKSVAMDMLNEVRCPELFTLSAEKQSRFRI I
 SDDHNEVLMKRS TRDRFVLLQYIDYGLKLFDIRFHVNMGKLRYLKADKICIDGQTRVRIEQ
 PLNGFGRLEEAETMRKQENGTFGNSGIRIRDFENVKRDDANPANYPIVDTYTHYILENNKVEM
 FISDKGSSAPLLPLIEDDRYVVKTI PSCRMSTLEIPAMAFHMLFSGSKTEKLI DVVHNRKYRL
 FQAMQKEEVTAEINISFGIAESDL POKI LDLSGNAHGKDVDAFIRLTVDDMLTDTERRIKRFK
 DDKRSIRSDADNKGKRGPKQISTGKLADFLAKDIVLPQPSVNDGENKITGLNRYIMQSAIAVYD
 SDDYEAQQQFKLMFEKARLIGKGT EPHPLKYVFARSIPANAVDFERYLIERKFFYLTGLCN
 EIKRGNRVDPFIRRDQNKWKTAMKTLGR IYSEDLVPELPRQMFNDNEIKSHLKSLPQMEGIDF
 NNANVTYLIAEYMKRVLNDDFQTFYQWKRNYHYMDMLKGEYDRKGSLOHCFSTSV EEREGLWKER
 ASRTERYRKLASNKIRSNRQMRNASSEIETILDKRLSNCRNEYQKSEKVI RRYRVQDALLFLL
 AKKTLTELADPDGERFKLKEIMPDAEKGLISEIMPMSFTFEKGGKTYITSEGMKLNKYGDFV
 LASDKRIGNLLELVGSDIVSKEDIMEEPNKYDQCRPEISSIVFNLEKWAFTYPELSARVDREE
 KVDFKSLIKI LLLNKNINKEQSDILRKIRNAPDHNNYPDKGIVETKALPEIAMS IKKAFGEYAI
 MK

Small protein sequences

>30 hypothetical protein [Psychroflexus torquus]
 MESIIGLGLSFPYKATADKHYFGSFLNLVNNLNAVFAEFKERISYKAKDENISSLIEKHPIDN
 MSIVDYEKKISILNGLYLPIIDFLDDELENNLNTRVKNFKKNFIILAEAI EKLRDYTHFYHDP
 TFE DNKEP LLELLEDEVLKTLIDVKKYKLTDKTEILKDSLREEMD LLLVIRKTDDELREKKTN
 PKIQHTDSSQIKNSIFNDAFQGLLYEDKGNKKTQVSHRAKTRLNPDKDIHKQEERDFEIPLS
 GLVFLMSLFLSKKEDFKSNIKGFKGVVDENHNLSKYMATHRVYSILAFKGLKYRIKDTDF
 SKETLMMQMI DELSKVPDCVYQNLSETKQKDFIEDWNEYFKDNEENTENLENSRVVHPVIRKY
 EDKPNYFAIRFLDEFANFKILKQVFMGYIHDQRTKTIGTTNIT TERTVKEKINVFGKLSKMD
 NLKHHFFSQLSDDENTDWEFFPNPSYNFLTQADNSPANNIP IYLELKNQOI I KEKDAI KAEVNO
 TQNRNPNKPSKRDLLNKILKTYEDFHQGDPTAILSLNETPALLHLFLVKNPKTGQQIENIRI
 KIEKQFKAINHPSKNNKGI PKSLFADINVRVNAIKLKKDLEAE LDMLNKHIAPKENQKASSNY
 DKLLKEHQFTPKNRPPELRKYVFKSEKGEATWLANDIKRPMKDFKTKWGCQHSSELQKRLA
 FYDRHTQD IKELLSGCEFDHSLDINAYFQKDNFEDFFSKYLENRIETLEGLVKKLHDFKNEP
 TPLKGVFKNCFKPLKRQNYVTESEPEI I KKRILAKFTFLPRGVFDERPTMKGKKNPLKDFEAE
 WFVEYLENKDYQKFNABEYMRDADFKNNAVIKQKLDKDFYTLQMVNLLKEVFGKDEMNLQ
 SELFQTRQERLKLQGI AKKQMNKETGDSSENTRNQT IWNKDVPSFNGKVTIDKVLKNI GK
 YKRYERDERVKTFIGYVEDEKMMYLPHNWKDRYSVKP INVIDLQIEYEEIRSHELLEKIQLN
 BQYIYDHTTDKNI LLQDGNPNFKMYVLNGLLIGIKQVNI PDFI VLKQNTNFKDIDFTGIASCSE
 LEKKTII LIAIRNKFAHNLQPNKMIYDLANEFLKIEKNETYANYLKVLLKMI SDLA

- continued

>1 hypothetical protein [*Bergeyella zoohelcum*]
 MIELINFSAEYVFGKAPLITAGNLSGDLFFAIALFIGLIHIITFWLTKRWYKKLACTPMDK
 MKNDEYITEILPDLKEIYRTNKGISVTENISNGSVFDTAYYYLEAQGKISQAKNFQSIYFLFRN
 IVTSLFVLPVSVIFLLASFFMNDCLSEKIIITIIIGTFVLGGLSSVIAQWFRVMTDRIIFGLY
 YAELTHNKK

>2 MULTISPECIES: hypothetical protein [*Prevotella*]
 MIIALITSLFSSATIIFAGYRTYRNAIKKLREESIENIKRAHYEAILKAHQSIYKLLRFTTDE
 NDDCILVWEQPKDGGEKTYFRQANIRQFFKELTEEVYKNGNGIYLSKEVMSLIFKYRTLHVHKL
 LLAEKNNPDEKLMIHNRVVKEMIEINQSLSIQIRKIDDLQQRNLYHNHKKGKYVWKRKKQNI

>3 hypothetical protein [*Prevotella buccae*]
 MDYMEAKEAFSIICTFIAAYVAYYYAIKQLHQKSVENIEYAKYQAVLQAHKSLYKLLRFTTNT
 ENEDSILWEKTKDGQEAATYFRKENIRKPIKELSKETYNEGCGIFMSKEALSISEYRNIVY
 GFMSAQNNPQETIRITNRESVERMKKIHQNLSIEIRQAINLKKRDLRF

>4 hypothetical protein [*Porphyromonas gingivalis*]
 MEKFSLYDFLAIILPGIAFIVVFRIFSSHLHSLPVDIPLGLESTIVYALICGAVLYVLSFSLV
 KLFPRFLGGLYRHHVADLYQMKALHPIMNDTLNRQAEQWGLGKIYLSSEEFQSEKEKIRMLQS
 DFYDRMWYRLDFRGKLGNAKSFQCYFFFRHSFWGLVLSLILLSYKLLAYIPACDMEDIQWRE
 YSDIAVPIIMLSALFVFLAQWFRKMKVEKMYWTFYISLIEQENSNI

>5 hypothetical protein [*Bacteroides pyogenes*]
 MNKFSLYDFLSILLPGVIFLVAIRVAQPFWRPNTGLYFPQGWEPFLVYSLVIGASLYVLGFSVK
 KNYSVFPRRLGLYEHVITLYHRFETLHPFMNGALNKYAEWNGTKPYCTVEQYDAMDASAQKEI
 EDAQDIFYDHMYRLDCKGKLEGAKAFQSYLCLHSPFLGLLIFGVYLLICILLSYFMDVLLAD
 TWQVAFPLMNLVYLFMRLARWFRQRMVLMKYWAFYESLIE

>8 hypothetical protein [*Riemerella anatipestifer*]
 MYWILEVLKIVPTTIAIVSIAIVSTIVLSRKIRQELKGNIERQKYEDILHAKQMRYLLAYMTD
 QDNPKNLLKWEVPGKQDKIHYINRANAQAPLRELPFLFYGEGCGLFLSEEVTKKFFEYRSIVH
 KLLLAEQNSTEAEPRLKNEEAATRMKVLHKKLSQSIRQCLKIEQRDLKAL

>9 hypothetical protein [*Prevotella aurantiaca*]
 MDCTLISKEVATALSTISSFIPVIAAYLYRYAIKKLHKESIENIERAKYEAAILKAHQSIYK
 LLRFTTDTENDCILLWEQSKGSKKVVYFRQANIRKPIKELTEEYKNGNGIYLSKEVILLVP
 KYRTLHVHKLKLLAEKNNPDEKIMISKQETAKSMIEIHQSLSIQIRNDINLKKRDLCLNSENT

>11 hypothetical protein [*Myroides odoratimimus*]
 MDKLSLYELFSFVIPGGIALHLLNWCNVNLSGTGLFNLSDLSNLSIALVFPALLIGVTLHIITF
 NILLKCGSYRQIYKSVQEKLLDDYIQQVIFPLNQEYFHNKKHEVAANTNNAVPAENLFDYAYY
 YLEVNGKNAQAKNFQSLYFFFRNMPFLGIVSIVILIALVYSTITSVQKDVLSIEVFKIAPFAV
 IIGIAVPVANWLRKKMIITVPGCYADRVRHQTNK

>13 hypothetical protein [*Bacteroides pyogenes*]
 MNKFSYDFLSILLPGVIFLVAIRVAQLFWRPNTGLYFPQGWEPFLVYSLVIGASLYVLGFSVK
 KNYSGFRCLGLYEHVITLYYRFETLHPFMNGALNKYAEWNGTKPYCTVEQYDAMDASAQKEI
 EDAQDIFYDHMYRLDCKGKLEGAKAFQSYLCLHSPFLGLLIFGVYLLICILLSYPMVLLAD
 TWQVAFPLMNLVYLFMRLARWFRQRMVLMKYWAFYESLIE

>14 hypothetical protein [*Capnocytophaga canimorsus*]
 MDRLSIYELLSFVVPGVIMIELNFSAEYVFGKDRITAGNLSGDLFFAIALFIGLIHIITF
 RLKRWYKQLAYKPIINEINNAIKQIFPTLKEYSKIHNSIDELSKLNIPEAVITTT

>15 hypothetical protein [*Capnocytophaga cynodegmi*]
 MERLSIYELLSFVVLGVIMIELMNFSAEYVFGKAPLITAGNLSDSLFFAIALFIGLIHIITF
 RLKRWYKQLACTPMDKMKNDEYITEILPHLKEIYRANKGISATENISNGSVFDTAYYYLEAQ
 GKISQAKNFQSIYFLFRNIVTSLFVLPVSVIFLLASFFMNDCKLSGKIITIVIGTLVIGGISS
 VIAQWFRVMTDRIIFGLYAEALTHHKK

>17 hypothetical protein [*Phaeodactylibacter xiamenensis*]
 MDEGTLINIITGPAKDKPLWFEFIIIGVTTISLAAGYLRAFGAWEKQKRELELQHKQKLEI
 KLQFEKQRYEHELKAAEGVWPLLAYFSLWENDKSVFVKRGDHWYFRQEQGREYILALSENFFNK
 YGVPMGPAPAKENLYHFRGMIYKLLQDSSKNGNDNNEVLLKNQSMVQKDKQAAPNKKSVEQLK
 DEINASLRNMLQKSEIDL

>18 hypothetical protein [*Porphyromonas gingivalis*]
 MMTMTICWTPICVQLIKLSGIFIAAYLAYRYAVHKLKESIENIERCKYQAVLEAHSFYKLLR
 FTTDENADSLVWQKAKGGAKTYFRPACIRGFLSELTEDEFYKNGNGIFLSKEIISRIPEYR
 SIVYGLLSEQRNSDERVVMNKPETAERMSIHQELTQTVREAIALKKRTLN

>19 hypothetical protein [*Porphyromonas gulae*]
 MEKFSLYDFLAIILPGIAFIVVFRVSSHLHSLPVDIPLDLEFTIVYALICGAVLYVLSFSLV
 KLFPRFLGGLYRHHVADLYQMKALHPIMNDTLNRQAEQWGLGRIYLSSEEEYCRSEKEKIRMLQS
 DFYDRMWYRLDFRGKLGNAKSFQCYFFFRHSFLGLVLSLILLSYKLLAYIPACELEDIGWKE
 YSDIAVPIIILSVLFFVFLAQWFRKMKVEKMYWTFYISLIEQENSNI

-continued

>20 hypothetical protein [*Porphyromonas gulae*]
 MMTMTICWTPICVQLIKLSGIFIAAYLAYRYAVRKLKSKESIENIERCKYQAVLEAHRSFYKLLR
 FTDTTENADSLVWQKAKGGGAKTYFRPACIRGFSELTEDEFYKNGNGVFLSKETIISRIFEYR
 SIVYGLLLSERQNSDERIVMKNKPETAERMIRIHQELTQTVREAIALKGRTLNF

>21 hypothetical protein [*Porphyromonas gulae*]
 MTMTICWTPICVQLIELSGIFIAAYLAYRYAVRKLKSKESIENIERCKYQAVLEAHRSFYKLLRF
 TDTTENADSLVWQKAKGGGAKTYFRPACIRGFSELTEDEFYKNGNGVFLSKETIISRIFEYRS
 IVYGLLLSERDSSDERIVMKNKPETAERMIRIHQELTQTVREAIALKGRTLNF

>22 hypothetical protein [*Prevotella falsenii*]
 MIIALIPALINSAAIIVVVYLYFYAVKKLREESFENIERAKYEAILKAHQSVYKLLRFITDTE
 NDDCILVWEQPKGGGKAYYFRQANIRKFIKELTEEYINKNKNGIFLSKKTMSLIFEYRFLVYGL
 LLTAKDNPAEAEIETNKKLAKRMIEIHQSLSIQIRKDRNLQQ

>23 hypothetical protein [*Prevotella intermedia*]
 MDCTLISKEVATALFSTISSLIPIVIAAYLYRYAIKKLRKESFENIERAKYEAILKAHQSIYK
 LLRFITDTENDDCILVWEQPKGGGKTYFRQANIRKFIKELTEEYINKNKNGIFLSKVMPLIF
 EYRSLVYGLLLTAKDKPDETIETKNEKLAKRMIEIHQSLSIQIRKDRNLKQRLDQNLKKGKQW
 CWKRKKQNI

>24 hypothetical protein [*Prevotella intermedia*]
 MDCTLISKEVATALFSTISSLIPIVIAAYLYRYAIKKLRKESFENIERAKYEAVLKAHQSIYK
 LLRFITDTENDDCILVWEQPKGGGKTYFRQANIRKFIKELTEEYINKNKNGIFLSKVMPLIF
 EYRSLVYGLLLTAKDKPDETIETKNEKLAKRMIEIHQSLSIQIRKDRNLKQRLDQFDS

>26 hypothetical protein [*Prevotella pleuritidis*]
 MDCTLISKEVATALFSTLSSLVPIVIAAYLYRYAIKKLHKESFESIECAKYEAILKAHQSIYK
 LLRFITDTENDDCILVWEQPKGGGENVYFRQANIRKFIKELTEEYINKNKNGIFLSKEVMSLIF
 EYRSLVYGLLLTEKDNPKDKI I IKNLKLAKRMIEIHQNL SIKIREAINLQQRDLHFK

>27 hypothetical protein [*Prevotella pleuritidis*]
 MDCTLISKEVATALFSTLSSLVPIVIAAYLYRYAIKKLHKESFESIECAKYEAILKAHQSIYK
 LLRFITDTENDDCILVWEQPKGGGKTYFRQANIRKFIKELTEEYINKNKNGIFLSKEVMSLIF
 KYRTL VHKLLLAKKNPDEKIMIDKRELAKRMIEIHQNL SIKIREAINLQQRDLHFK

>28 hypothetical protein [*Prevotella saccharolytica*]
 MECNDCSQLIITLVSLTAAALSGPLSYFAVRKYRKESQENVERWKYEAIWHAHQSFYKLLRFM
 TDENADSLVFRKQEGSKQPPIFRKDKAREFLAELTEDEFYKNGNGLFLSKEVADSLFKYRHI
 VFGLLQVFSNYDENEVQLKNTQAVENMKKNFQQLSPEIRKSLHLHRRDLFLS

REFERENCES

- [1011] 1. Abil Z, Zhao H. Engineering reprogrammable RNA-binding proteins for study and manipulation of the transcriptome. *Molecular bioSystems*. 2015; 11(10): 2658-65. doi: 10.1039/c5mb00289c. PubMed PMID: 26166256.
- [1012] 2. Abudayyeh O O, Gootenberg J S, Konermann S, Joung J, Slaymaker I M, Cox D B, et al. C2c2 is a single-component programmable RNA-guided RNA-targeting CRISPR effector. *Science*. 2016; 353(6299): aaf5573. doi: 10.1126/science.aaf5573. PubMed PMID: 27256883.
- [1013] 3. Anantharaman V, Makarova K S, Burroughs A M, Koonin E V, Aravind L. Comprehensive analysis of the HEPN superfamily: identification of novel roles in intra-genomic conflicts, defense, pathogenesis and RNA processing. *Biology direct*. 2013; 8:15. doi: 10.1186/1745-6150-8-15. PubMed PMID: 23768067; PubMed Central PMCID: PMC3710099.
- [1014] 4. Bernhart S H, Hofacker I L, Stadler P F. Local RNA base pairing probabilities in large sequences. *Bioinformatics*. 2006; 22(5):614-5. doi: 10.1093/bioinformatics/btk014. PubMed PMID: 16368769.
- [1015] 5. Biswas A, Gagnon J N, Brouns S J, Fineran P C, Brown C M. CRISPRTarget: bioinformatic prediction and analysis of crRNA targets. *RNA biology*. 2013; 10(5): 817-27. doi: 10.4161/rna.24046. PubMed PMID: 23492433; PubMed Central PMCID: PMC3737339.
- [1016] 6. Brouns S J, Jore M M, Lundgren M, Westra E R, Slijkhuys R J, Snijders A P, et al. Small CRISPR RNAs guide antiviral defense in prokaryotes. *Science*. 2008; 321(5891):960-4. doi: 10.1126/science.1159689. PubMed PMID: 18703739.
- [1017] 7. Camacho C, Coulouris G, Avagyan V, Ma N, Papadopoulos J, Bealer K, et al. BLAST+: architecture and applications. *BMC bioinformatics*. 2009; 10:421. doi: 10.1186/1471-2105-10-421. PubMed PMID: 20003500; PubMed Central PMCID: PMC2803857.
- [1018] 8. Crooks G E, Hon G, Chandonia J M, Brenner S E. WebLogo: a sequence logo generator. *Genome research*. 2004; 14(6):1188-90. doi: 10.1101/gr.849004. PubMed PMID: 15173120; PubMed Central PMCID: PMC419797.
- [1019] 9. East-Seletsky A, O'Connell M R, Knight S C, Burstein D, Cate J H, Tjian R, et al. Two distinct RNase activities of CRISPR-C2c2 enable guide-RNA processing and RNA detection. *Nature*. 2016. doi: 10.1038/nature19802. PubMed PMID: 27669025.
- [1020] 10. Edgar R C. PILER-CR: fast and accurate identification of CRISPR repeats. *BMC bioinformatics*. 2007; 8:18. doi: 10.1186/1471-2105-8-18. PubMed PMID: 17239253; PubMed Central PMCID: PMC1790904.
- [1021] 11. Filipovska A, Rackham O. Designer RNA-binding proteins: New tools for manipulating the transcriptome. *RNA biology*. 2011; 8(6):978-83. doi: 10.4161/rna.8.6.17907. PubMed PMID: 21941129.

- [1022] 12. Gerdes S Y, Scholle M D, Campbell J W, Balazsi G, Ravasz E, Daugherty M D, et al. Experimental determination and system level analysis of essential genes in *Escherichia coli* MG1655. *Journal of bacteriology*. 2003; 185(19):5673-84. PubMed PMID: 13129938; PubMed Central PMCID: PMC193955.
- [1023] 13. Hale C R, Cocozaki A, Li H, Terns R M, Terns M P. Target RNA capture and cleavage by the Cmr type III-B CRISPR-Cas effector complex. *Genes & development*. 2014; 28(21):2432-43. doi: 10.1101/gad.250712.114. PubMed PMID: 25367038; PubMed Central PMCID: PMC4215187.
- [1024] 14. Hale C R, Majumdar S, Elmore J, Pfister N, Compton M, Olson S, et al. Essential features and rational design of CRISPR RNAs that function with the Cas RAMP module complex to cleave RNAs. *Molecular cell*. 2012; 45(3):292-302. doi: 10.1016/j.molcel.2011.10.023. PubMed PMID: 22227116; PubMed Central PMCID: PMC3278580.
- [1025] 15. Hale C R, Zhao P, Olson S, Duff M O, Graveley B R, Wells L, et al. RNA-guided RNA cleavage by a CRISPR RNA-Cas protein complex. *Cell*. 2009; 139(5):945-56. doi: 10.1016/j.cell.2009.07.040. PubMed PMID: 19945378; PubMed Central PMCID: PMC2951265.
- [1026] 16. Hayes F, Van Melder L. Toxins-antitoxins: diversity, evolution and function. *Critical reviews in biochemistry and molecular biology*. 2011; 46(5):386-408. doi: 10.3109/10409238.2011.600437. PubMed PMID: 21819231.
- [1027] 17. Jiang W, Samai P, Marraffini L A. Degradation of Phage Transcripts by CRISPR-Associated RNases Enables Type III CRISPR-Cas Immunity. *Cell*. 2016; 164(4):710-21. doi: 10.1016/j.cell.2015.12.053. PubMed PMID: 26853474; PubMed Central PMCID: PMC4752873.
- [1028] 18. Jinek M, Chylinski K, Fonfara I, Hauer M, Doudna J A, Charpentier E. A programmable dual-RNA-guided DNA endonuclease in adaptive bacterial immunity. *Science*. 2012; 337(6096):816-21. doi: 10.1126/science.1225829. PubMed PMID: 22745249.
- [1029] 19. Mackay J P, Font J, Segal D J. The prospects for designer single-stranded RNA-binding proteins. *Nature structural & molecular biology*. 2011; 18(3):256-61. doi: 10.1038/nsmb.2005. PubMed PMID: 21358629.
- [1030] 20. Makarova K S, Anantharaman V, Aravind L, Koonin E V. Live virus-free or die: coupling of antiviral immunity and programmed suicide or dormancy in prokaryotes. *Biology direct*. 2012; 7:40. doi: 10.1186/1745-6150-7-40. PubMed PMID: 23151069; PubMed Central PMCID: PMC3506569.
- [1031] 21. Makarova K S, Haft D H, Barrangou R, Brouns S J, Charpentier E, Horvath P, et al. Evolution and classification of the CRISPR-Cas systems. *Nature reviews Microbiology*. 2011; 9(6):467-77. doi: 10.1038/nrmicro2577. PubMed PMID: 21552286; PubMed Central PMCID: PMC3380444.
- [1032] 22. Makarova K S, Wolf Y I, Alkhnbashi O S, Costa F, Shah S A, Saunders S J, et al. An updated evolutionary classification of CRISPR-Cas systems. *Nature reviews Microbiology*. 2015; 13(11):722-36. doi: 10.1038/nrmicro3569. PubMed PMID: 26411297.
- [1033] 23. Makarova K S, Wolf Y I, Koonin E V. Comprehensive comparative-genomic analysis of type 2 toxin-antitoxin systems and related mobile stress response systems in prokaryotes. *Biology direct*. 2009; 4:19. doi: 10.1186/1745-6150-4-19. PubMed PMID: 19493340; PubMed Central PMCID: PMC2701414.
- [1034] 24. Marraffini L A. CRISPR-Cas immunity in prokaryotes. *Nature*. 2015; 526(7571):55-61. doi: 10.1038/nature15386. PubMed PMID: 26432244.
- [1035] 25. Moller S, Croning M D, Apweiler R. Evaluation of methods for the prediction of membrane spanning regions. *Bioinformatics*. 2001; 17(7):646-53. PubMed PMID: 11448883.
- [1036] 26. Nunez J K, Lee A S, Engelman A, Doudna J A. Integrase-mediated spacer acquisition during CRISPR-Cas adaptive immunity. *Nature*. 2015; 519(7542):193-8. doi: 10.1038/nature14237. PubMed PMID: 25707795; PubMed Central PMCID: PMC4359072.
- [1037] 27. Remmert M, Biegert A, Hauser A, Soding J. HHblits: lightning-fast iterative protein sequence searching by HMM-HMM alignment. *Nature methods*. 2012; 9(2):173-5. doi: 10.1038/nmeth.1818. PubMed PMID: 22198341.
- [1038] 28. Shmakov S, Abudayyeh O O, Makarova K S, Wolf Y I, Gootenberg J S, Semenova E, et al. Discovery and Functional Characterization of Diverse Class 2 CRISPR-Cas Systems. *Molecular cell*. 2015; 60(3):385-97. doi: 10.1016/j.molcel.2015.10.008. PubMed PMID: 26593719; PubMed Central PMCID: PMC4660269.
- [1039] 29. Staals R H, Agari Y, Maki-Yonekura S, Zhu Y, Taylor D W, van Duijn E, et al. Structure and activity of the RNA-targeting Type III-B CRISPR-Cas complex of *Thermus thermophilus*. *Molecular cell*. 2013; 52(1):135-45. doi: 10.1016/j.molcel.2013.09.013. PubMed PMID: 24119403; PubMed Central PMCID: PMC4006948.
- [1040] 30. Staals R H, Zhu Y, Taylor D W, Kornfeld J E, Sharma K, Barendregt A, et al. RNA targeting by the type III-A CRISPR-Cas Csm complex of *Thermus thermophilus*. *Molecular cell*. 2014; 56(4):518-30. doi: 10.1016/j.molcel.2014.10.005. PubMed PMID: 25457165; PubMed Central PMCID: PMC4342149.
- [1041] 31. Tafer H, Ameres S L, Obernosterer G, Gebeshuber C A, Schroeder R, Martinez J, et al. The impact of target site accessibility on the design of effective siRNAs. *Nature biotechnology*. 2008; 26(5):578-83. doi: 10.1038/nbt1404. PubMed PMID: 18438400.
- [1042] 32. Tamulaitis G, Kazlauskienė M, Manakova E, Venclovas C, Nwokeoji A O, Dickman M J, et al. Programmable RNA shredding by the type III-A CRISPR-Cas system of *Streptococcus thermophilus*. *Molecular cell*. 2014; 56(4):506-17. doi: 10.1016/j.molcel.2014.09.027. PubMed PMID: 25458845.
- [1043] 33. van der Oost J, Jore M M, Westra E R, Lundgren M, Brouns S J. CRISPR-based adaptive and heritable immunity in prokaryotes. *Trends in biochemical sciences*. 2009; 34(8):401-7. doi: 10.1016/j.tibs.2009.05.002. PubMed PMID: 19646880.
- [1044] 34. Wright A V, Nunez J K, Doudna J A. Biology and Applications of CRISPR Systems: Harnessing Nature's Toolbox for Genome Engineering. *Cell*. 2016; 164(1-2):29-44. doi: 10.1016/j.cell.2015.12.035. PubMed PMID: 26771484.
- [1045] 35. Yates A, Akanni W, Amode M R, Barrell D, Billis K, Carvalho-Silva D, et al. Ensembl 2016. *Nucleic acids research*. 2016; 44(D1):D710-6. doi: 10.1093/nar/gkv1157. PubMed PMID: 26687719; PubMed Central PMCID: PMC4702834.

- [1046] 36. Yosef I, Goren M G, Qimron U. Proteins and DNA elements essential for the CRISPR adaptation process in *Escherichia coli*. *Nucleic acids research*. 2012; 40(12):5569-76. doi: 10.1093/nar/gks216. PubMed PMID: 22402487; PubMed Central PMCID: PMC3384332.
- [1047] 37. Zhang J, Rouillon C, Kerou M, Reeks J, Brugger K, Graham S, et al. Structure and mechanism of the CMR complex for CRISPR-mediated antiviral immunity. *Molecular cell*. 2012; 45(3):303-13. doi: 10.1016/j.molcel.2011.12.013. PubMed PMID: 22227115; PubMed Central PMCID: PMC3381847.

Example 2: Validation of Group 29 Effector Protein RNA Interference in Bacteria, In Vitro (Biochemically) and in Eukaryotic (Human) Cells

[1048] Mammalian Interference Assays (Protein and crRNA Plasmids): For mammalian interference experiments, Applicants designed two plasmids to co-transfect in HEK 293T cells:

[1049] Protein plasmid containing CMV-driven Group 29 proteins human codon optimized for 11 species (proteins 1-11 from select unique proteins).

[1050] crRNA plasmid containing U6-driven crRNA with repeat and target sequences, repeats varying by species and targets dependent on assay

[1051] Exogenous Luciferase-Based Assay: Applicants conducted the exogenous mammalian interference assay in HEK 293T cells using Promega's Nano-Glo® Dual-Luciferase® Reporter Assay. Luciferase assays offer the highest sensitivity in optical readout of gene transcript activity, and the recently engineered NLuc offers >2 orders of magnitude sensitivity than Fluc and Rluc. Nluc and Fluc for normalization will be expressed on two lentiviral backbone plasmids driven by CMV or another constitutive mammalian promoter (or perhaps an inducible promoter). The plasmids will have antibiotic resistance genes for possible transfection into Nluc/Fluc expressing FMK 293T cell lines. In addition, canonical U2 introns may be included in active Nluc and Rluc regions to allow for crRNA target sequences that span exons across the intron.

[1052] Endogenous Transcript RT-qPCR Assay: Applicants conducted the endogenous mammalian interference assay in HEK 293T cells using crRNA target sequences that span exons across introns in endogenous mammalian genes, for example, CXCR4 and B4GALNT1. Applicants assessed RNA interference through RT-qPCR, from which delta CT values will be used to calculate relative interference across samples and conditions.

[1053] MS2 Phage Interference: To assay potential RNA targeting capabilities of group 29 CRISPR nucleases Applicants designed a series of constructs that may interfere with replication of the MS2 phage. The MS2 phage is a +single-stranded RNA virus that infects *E. coli* as further described in Kannoly, S., Shao, Y. & Wang, I. N. Rethinking the evolution of single-stranded RNA (ssRNA) bacteriophages based on genomic sequences and characterizations of two R-plasmid-dependent ssRNA phages, C-1 and Hgall1. *Journal of bacteriology* 194, 5073-5079 (2012). Its replication does not involve any DNA intermediates. CRISPR systems that specifically act on RNA can be utilized to interfere with MS2 replication_ENREF_2 as further described in Tamulaitis, G., et al. Programmable RNA shredding by the type III-A CRISPR-Cas system of *Streptococcus thermophilus*.

Molecular cell 56, 506-517 (2014). To test whether group 29 proteins can cleave RNA, Applicants designed a series of constructs that may direct group 29 nucleases to interfere with MS2 replication. The assay specifically involves 1) creating *E. coli* carrying group 29 loci on the chloramphenicol resistant plasmid pACYC184, whose elements are expressed under their natural or synthetic promoters and contain crRNA expression cassettes targeting the mat, lys and rep MS2 genes, 2) infecting the *E. coli* from step 1 above with the MS2 phage and 3) determining infection efficiency by standard efficiency of plaquing assays.

[1054] Plasmid Interference Assays: To assay potential RNA targeting capabilities of Group 29 CRISPR nucleases Applicants designed a series of constructs that interfere with gene expression from the pUC19 plasmid. The assay specifically involves 1) creating chemically competent *E. coli* carrying group 29 loci on the chloramphenicol resistant plasmid pACYC184, whose elements are expressed under their natural or synthetic promoters and contain crRNA expression cassettes targeting the ampicillin resistance gene of pUC19, 2) transforming the *E. coli* from step 1 above with the pUC19 plasmid and 3) enumerating the colonies that are resistant to both chloramphenicol and ampicillin through plating on LB agar plates with appropriate antibiotics.

[1055] If Group 29 loci with properly targeted crRNAs display a reduction in the number of chloramphenicol, ampicillin resistant cfu relative to controls, then interference has likely occurred. Furthermore, by designing crRNAs predicted to hybridize to the ampicillin resistance transcript (bla) or non-transcribed parts of pUC19 Applicants expect to assay whether interference is the result of mRNA targeting or DNA cleavage. If Group 29 targets mRNA for destruction, then only crRNAs complementary to the transcript should abrogate ampicillin resistance. If Group 29 targets DNA for destruction, then only crRNAs targeting sequences with the correct PAM (if a PAM is required for cleavage) would be expected to result in a loss of ampicillin resistance for transformed cells.

[1056] Protein Purification: Applicants expressed and purified the Group 29 proteins using standard assays as further described in Rigaut, Guillaume, et al. "A generic protein purification method for protein complex characterization and proteome exploration." *Nature biotechnology* 17.10 (1999): 1030-1032. In brief, Applicants cloned the Group 29 large proteins (as described in the present application) into the pET21 expression vector and fused it to a 6x His tag for Ni-His purification and to Maltose binding protein to aid in solubility. The Group 29 proteins will be expressed in *Escherichia coli* Rosetta 2 (DE3) (Novagen) and purified using Ni-NTA resin (Bio-Rad) using gravity flow columns. Proteins are further purified using size exclusion chromatography.

[1057] Electrophoretic Mobility Shift Assay (EMSA) and Cleavage Assay: Applicants performed an EMSA and a cleavage assay as described by Tamulaitis, Gintautas et al. "Programmable RNA Shredding by the Type III-A CRISPR-Cas System of *Streptococcus thermophilus*." *Molecular Cell* 56, 506-517, 2014 with some minor modifications. For EMSA, Applicants incubated purified Group 29 proteins with fluorescently labeled nucleic acids (including crRNAs and dual repeat crRNAs) and ran them on a native PAGE gel. For the cleavage assay, Applicants incubated Group 29

proteins with crRNAs and the target nucleic acid at 25 C and 37 C and ran the target nucleic acid on a PAGE gel to analyze cleavage.

[1058] RNA Purification for In Vitro Assays: To purify enough RNA for the Group 29 in vitro RNA binding/cleavage assays, Applicants followed a protocol from Rapid purification of RNAs using fast performance liquid chromatography (FPLC) (Kim et al., RNA, 2007). This one-day protocol involves in vitro transcription with T7 polymerase, followed by fast performance liquid chromatography (FPLC) size-exclusion chromatography. Applicants synthesized and purified Group 29 crRNAs (repeats with target sequences) and target RNA. For target RNA, Applicants synthesized both ssRNA complimentary to target sequence and dsRNA with target sequence (complimentary ssRNA annealed to its reverse compliment).

[1059] Protocol for FIG. 35: In vitro cleavage assay. Testing cleavage timecourse with 1:1:1 Target:crRNA:Prot molar ratio of a single target, single guide (EGFP1 target) T2 crRNA and 01 (N2) crRNA

Protein Concentration Variation IVC Protocol

- [1060]** 1. Normalize short crRNA to 36.5 ng/uL
[1061] 2. Normalize each PAM library to 100 ng/uL
[1062] 3. Make 10x buffer+DTT
[1063] 4. Dilute protein in H2O: Need 26 samples=7 uL protein+19 uL H2O
[1064] 5. Make Master Mix

EGFP Target 1-MM

[1065]

RNAse Inhibitor	1 uL	36 uL
10x Buffer w/ DTT	1 uL	36 uL
ddH2O	5 uL x 36	180 uL
E1 (100 ng/uL)	1 uL	36 uL
Total	8 uL	

Layout:

[1066]

Lane	1	2	3	4	5	6	7	8	9	10	11	12	13	14
Time	0 min	1 min	5 min	10 min	20 min	30 min	60 min	90 min	0 min	0 min	0 min	90 min	90 min	90 min
crRNA	+	+	+	+	+	+	+	+	+	-	-	+	-	-
protein	+	+	+	+	+	+	+	+	-	+	-	-	+	-

[1067] 1. MM for time samples for each crRNA (9x for 8 samples). Take 72 uL of MM, add 9 uL of diluted protein and 9 uL. Add 10 uL per well.

[1068] 2. MM for -protein +crRNA. For each crRNA, pipette 17.6 uL MM, 2.2 uL of diluted salt, 2.2 of crRNA. Pipette 10 uL to 90 min timepoint.

[1069] 3. MM for -crRNA (make for all 4 samples). Pipette 35.2 MM, 4.4 H2O, 4.4 protein dil

[1070] 4. MM for -prot -crRNA (make for all 4 samples). Pipette 35.2 MM, 4.4 H2O, 4.4 dil salt.

[1071] 5. Mix components. Incubate at times indicated above.

[1072] 6. For each time point, at the proper time, heat for 5 min at 85 C, then add 1 uL of Proteinase K at RT (15 min)

[1073] 7. Set up gel(s). Heat at 42 C. When 15 min left on incubation, start running gel (to stabilize internal temperature).

[1074] 8. Remove 5 uL of mix, add 5 uL of 2x loading dye.

[1075] 9. Freeze down rest for for sequencing.

[1076] 10. Heat 10 uL at 95 C for 5 minutes.

[1077] 11. Chill on ice and spin down.

[1078] 12. Load 8 uL per lane. Load tightly

[1079] 13. Run gel at 180V for 45 minutes at 42 C.

[1080] 14. Take out gels and wash into DI water

[1081] 15. Image downstairs no later than 1 hr after finishing gel.

[1082] FIG. 35 shows the results of an in vitro cleavage assay of EGFP RNA by Group 29 proteins (purified from *Bergeyella zoohelcum* ATCC 43767). Two different crRNAs were used (FIG. 35A and FIG. 35B respectively). Cleavage reactions were performed for the indicated time points in the presence or absence of crRNA and/or effector protein in a buffer containing 50 mM NaCl, 1 mM DTT, Tris-HCl pH7.5, 0.1% BSA. Components were mixed at 1:1:1 target:protein:crRNA ratio. From FIG. 35 it seems that Group 29 protein makes one or a few initial cuts, and as time goes by, the target is further processed. Cutting appears to be non-random, as cut patterns are reproducible.

[1083] Protocol for FIG. 36A/C: EMSA Gel shift. Testing binding of the WT and dead proteins to the crRNA and crRNA+target using EGFP1 target and T2 crRNA.

Protein Concentration Variation IVC Protocol

- [1084]** 1. Normalize short crRNA to 50 ng/uL (Dil 8 uL crRNA with 8 uL H2O)
[1085] 2. Normalize E1 target to 125 ng/uL (1 uL+7 uL H2O)
[1086] 3. Make 10x buffer+DTT (10 uL DTT+90 10x)
[1087] 4. Make Master Mix

EGFP Target 1

[1088]

RNAse Inhibitor	5 uL	65 uL
10x Buffer w/ DTT	5 uL	65 uL
ddH2O	33 uL x 13	429 uL
crRNA (T2)	1 uL	13 uL
Total	44 uL	

Layout:

[1089]

Lane	1	2	3	4	5	6	7	8	9	10	11	12	13	14
crRNA	+	+	+	+	+	+	+	+	+	+	+	+	-	-
Target	-	-	-	+	+	+	-	-	-	+	+	+	+	+
protein	WT	WT	WT	WT	WT	WT	D	D	D	D	D	D	D	D
	0	.25	2.5	0	.25	2.5	0	.25	2.5	0	.25	2.5	.25	2.5

[1090] 1. Pipette out 44 uL MM**[1091]** 2. Add either 1 uL of dH₂O or E1 target.**[1092]** 3. Make protein dilutions. (WT=1 ug/uL) (D=bug/uL)**[1093]** a. WT0: 5.5 uL 2M NaCl, 5.5 uL H₂O**[1094]** b. WT.25: 0.55 uL prot, 4.95 2M NaCl, 5.5 uL H₂O**[1095]** c. WT2.5: 5.5 uL prot, 5.5 uL H₂O**[1096]** d. D0: 5.5 uL 2M NaCl, 5.5 uL H₂O**[1097]** e. D.25: 8.5 uL NaCl, 8.5 uL H₂O. Take out 1.5 uL. Make dil 0.5 uL prot+4.5 uL 50/50 mix.**[1098]** f. D2.5: 1.5+8.1 2M NaCl, 7.4 uL H₂O**[1099]** 4. Pipette 5 uL per each protein solution into wells.**[1100]** 5. Mix components. Incubate for 30 minutes.**[1101]** 6. Set up gel(s).**[1102]** 7. Remove 10 uL of sample, mix 1 uL of new concentrated loading dye (glycerol and color)**[1103]** 8. Load 10 uL per lane. Load tightly**[1104]** 9. Run gel at 180V for 45 minutes at 42 C.**[1105]** 10. Take out gels and wash into DI water**[1106]** 11. Image downstairs no later than 1 hr after finishing gel.

[1107] FIG. 36A shows the results of a binding assay of the CRISPR complex (or individual components thereof) with an EGFP RNA target. EMSA was performed including one or more of Group 29 protein (purified from *Bergeyella zoohelcum* ATCC 43767) as the indicated concentrations, EGFP target, and crRNA. The binding assay reaction mixture contained 50 ng crRNA, 125 ng target (if present). Addition of 0.25 ug of protein corresponds to ~1:1:1 molar ratio protein:crRNA:target. Addition of 2.5 ug of protein corresponds to ~10:1:1 molar ratio protein:crRNA:target. Binding assays were performed for 30 minutes in a buffer containing 100 mM NaCl, 1 mM DTT, Tris-HCl pH7.5, 0.1% BSA. Binding of wild type Group 29 protein (left

panel of FIG. 36) and mutated Group 29 protein (right panel of FIG. 36) was evaluated. Mutated Group 29 protein contained the following mutations: R116A, H121A, R1177A, and H1182A. From FIG. 36 it seems that mutated Group 29 protein (dGroup 29) binds crRNA, crRNA and its target RNA in complex. Binding is not eliminated by making mutations in the HEPN domain. FIG. 36C shows similar results for HEPN domain 2 mutant Group 29 protein (i.e. having mutations R1177A, H1182A).

[1108] Protocol for FIG. 36B/D: Protocol for testing if dGroup29 cleaves protein Now salt concentration will be 50 mM (2.708 uL of 2M NaCl in a 100 uL rxn)

[1109] 1. Normalize short crRNA to 36.5 ng/uL (5 uL in 8.7 uL ddH₂O).**[1110]** 2. Dilute frozen protein.**[1111]** a. WT: 1.08 uL prot+2.92 H₂O**[1112]** b. D 0.27: 1 uL above mix+4 uL S200**[1113]** c. d 1.35: 1.35 uL prot+4.65 S200 buffer**[1114]** 3. Normalize each PAM library to 100 ng/uL**[1115]** 4. Make 10x Buffer w/DTT (Add 10 uL DTT to 90 uL of 10x buffer without DTT)**[1116]** 5. Make Master Mixes:

EGFP Target 1

[1117]

RNAse Inhibitor	1 uL	14 uL
10x Buffer w/ DTT	1 uL	14 uL
ddH ₂ O	5 uL x 14	70 uL
E1 (100 ng/uL)	1 uL	14 uL
Total	8 uL	

Layout:

[1118]

Gel 1	1	2	3	4	5	6	7	8	9	10	11	12
Tem-plate	+	+	+	+	+	+	+	+	+	+	+	+
Pro-tein	WT	WT	WT	-	d	d	d	-	d	d	d	-
	.27 ug	.27 ug	.27 ug		.27 ug	.27 ug	.27 ug		1.3 ug	1.3 ug	1.3 ug	
	+	+	+		+	+	+		+	+	+	
RNA	01	T2	-	-	01	T2	-	-	01	T2	-	-

- [1119] 1. Aliquot 8 uL MM.
- [1120] 2. Add 1 uL of protein or salt to each well.
- [1121] 3. Add 1 uL of short crRNA (36.5 ng/uL) or ddH2O to each well.
- [1122] 4. Mix components.
- [1123] 5. Incubate for 60 min at 37 C. Then heat inactivate for 5 minutes at 85 C.
- [1124] 6. Prot K treat: 1 uL per sample for 15 min @ RT
- [1125] 7. Run TBE urea gel

[1126] FIG. 36B shows the results of an in vitro cleavage assay of EGFP RNA by Group 29 proteins (purified from *Bergeyella zoohelcum* ATCC 43767), in which wild type and mutated Group 29 proteins were evaluated. Mutated Group 29 effector harbored the following mutations: R116A, H121A, R1177A, and H1182A. two different crRNAs were used as indicated. Different component ratios were tested, as indicated. Assay was evaluated after incubation for 60 minutes; buffer used was: 50 mM NaCl, 1 mM DTT, Tris-HCl pH7.5, 0.1% BSA. From FIG. 36B can be concluded that the mutant Group 29 effector did not cut target RNA, irrespective of component ratios. FIG. 36D shows similar results for HEPN domain 2 mutant Group 29 protein (i.e. having mutations R1177A, H1182A).

[1127] Protocol for FIG. 37: IVC DNA; testing ssDNA cleavage.

Long Vs. Short crRNA IVC Protocol

[1128] Linearize plasmid, Purify, Run on Gel to confirm cut

[1129] Normalize short crRNA to 100 ng/uL, long crRNA to 200 ng/uL

[1130] Make Master Mixes:

	5'PAM MM		3'PAM MM	
10x pH 7.9	2 uL	24 uL	10x pH 7.9	2 uL 24 uL
NaCl	1 uL	12 uL	NaCl	1 uL 12 uL
MgCl2	1 uL x 12	12 uL	MgCl2	1 uL x 12 12 uL
DNA Target (25 ng/uL)	8 uL		DNA Target	8 uL
H2O	2 uL		H2O	2 uL
Total	14 uL		Total	14 uL

Layout:

[1131]

	1	2	3	4	5	6	7	8	9	10
DNA target	+	+	+	+	+	+	+	+	+	+
crRNA	short	short (NS)	long	long (NS)	-	short	short (NS)	long	long (NS)	-
protein	+	+	+	+	+	-	-	-	-	-

[1132] Aliquot 14 uL of MM per well (10 wells of each type of PAM library)

[1133] Dilute protein. Need 10 samples total. Mix 6 uL protein with 18 uL of H2O. Add 2 uL to each well requiring protein

[1134] Add protein buffer to other wells (2 uL) (1:4 dilution of 500 mM 5200 protein solution)

[1135] Add 4 uL of proper crRNA.

[1136] Mix components. Incubate at 37 C for 1 hr. Protein inactivation at 85 C for 5 min.

[1137] Incubate with RNaseA. Dilute 10 uL stock with 50 uL of H2O. Add 1 uL per well.

[1138] Incubate 30 min at 37 C.)

[1139] Incubate with Proteinase K (15 min RT)

[1140] Run on denaturing gel.

[1141] FIG. 37 shows the results of an in vitro cleavage assay of ssDNA by Group 29 proteins (purified from *Bergeyella zoohelcum* ATCC 43767). Binding to a ssDNA LPAM library (left panel of FIG. 37) and a ssDNA RPAM library (right panel of FIG. 37) was evaluated. FIG. 37 indicates that Group 29 proteins do not cut ssDNA, irrespective of direct repeat length.

[1142] Protocol for FIG. 38: IVC; testing cleavage of CRISPR array

[1143] 1. Synthesize crRNAs

[1144] 2. Normalize crRNAs to 100 ng/uL

[1145] 3. Make protein concentrations:

[1146] 1: 1 uL stock

[1147] 2: 0.65 uL stock+0.35 uL 2M NaCl

[1148] 3: Add 0.5 uL stock+2 uL 2M NaCl. Take 1 uL

[1149] 4: 0.65 uL 200 ng/uL+0.35 uL 2M NaCl

[1150] 5/6: 1 uL of 2M

[1151] 4. Make 10x buffer+DTT

[1152] 5. Make Master Mix

Mg MM

[1153]

10x Buffer	2 uL	14 uL
RNase Inh.	1 uL	7 uL
H2O	15 uL x 7	105 uL
Total	18 uL	

Layout:

[1154]

Lane	1	2	3	4	5	6
Time	60 min	60 min	60 min	60 min	60 min	60 min
crRNA	FR1	FR2	FR1	FR2	FR1	FR2
Protein	1 ug (5:1)	.65 ug (5:1)	200 ng (1:1)	130 ng (1:1)	-	-

[1155] 1. Add 1 uL of target crRNA (100 ng per reaction)

[1156] 2. Add protein or Salt buffer to other wells (1 uL).

[1157] 3. Mix components. Incubate for 60 min at 37 C then heat for 5 min at 85 C, then add 1 uL of Proteinase K at RT (15 min)

[1158] 4. Set up gel(s). Heat at 42 C. When 15 min left on incubation, start running gel (to stabilize internal temperature).

- [1159] 5. Remove 7 uL of mix, add 7 uL of 2x loading dye.
- [1160] 6. Heat 10 uL at 95 C for 5 minutes.
- [1161] 7. Chill on ice and spin down.
- [1162] 8. Load 12 uL per lane. Load tightly
- [1163] 9. Run gel at 180V for 45 minutes at 42 C.

[1164] The following LPAM, RPAM and crRNA sequences are artificial sequences.

LPAM dsDNA Sequence*

GGGcgatggtccatcaccttagatgcggttagcattaatcaggcaacggctc
tctagatagNNNNNNNgtaggctctttgtgctctacgccatctccatag
agccctcaaccggagtttgaagcatggcttctaactttactcagttcgtt
ctcgtcgacaatggcggaactggcgacgtgact

RPAM dsDNA Sequence*

GGGcgatggtccatcaccttagatgcggttagcattaatcaggcaacggctc
tctagataggtaggctctttgtgctctacgccatctccNNNNNNNtag
agccctcaaccggagtttgaagcatggcttctaactttactcagttcgtt
ctcgtcgacaatggcggaactggcgacgtgact

LPAM ss/dsDNA Sequence*

TAATACGACTCACTATAGGGcgatggtccatcaccttagatgcggttagcatt
aatcaggcaacggctctctagatagNNNNNNNgtaggctctttgtgctc
taccgccatctccatagagccctcaaccggagtttgaagcatggcttcta
actttactcagttcgttctcgtcgacaatggcggaactggcgacgtgact

RPAM ss/dsDNA Sequence*

TAATACGACTCACTATAGGGcgatggtccatcaccttagatgcggttagcatt
aatcaggcaacggctctctagataggtaggctctttgtgctctacgcc
atctccNNNNNNNtagagccctcaaccggagtttgaagcatggcttcta
actttactcagttcgttctcgtcgacaatggcggaactggcgacgtgact

crRNA DR Sequence of BZ (short)
GUUGGAACUGCUCUCAUUUUGGAGGUAAUACAAC

crRNA DR Sequence of BZ (long)
GUUGGAACUGCUCUCAUUUUAUUCUGAAGUUUUUAUUUGUUUCAA
GGAACUCAUGAAUACAGAGUAUUUGGAGGUAAUACAAC

crRNA Sequence
5'-guide sequence-DR sequence-3'
*Guide sequences reverse complements (in RNA) of highlighted target sequences

[1165] FIG. 38 shows that Group 29 proteins (purified from *Bergeyella zoohelcum* ATCC 43767) process their own crRNA. The Group 29 native locus is schematically represented in FIG. 38A. Two different CRISPR array fragments were evaluated for processing by Group 29 protein (indicated as #1 and #2). FIG. 38B indicates that both fragments are efficiently processed. A mature crRNA is about 66 nt.

[1166] Protocol of FIG. 39: Protocol for repeating EGFP cleavage Now salt concentration will be 50 mM (2.708 uL of 2M NaCl in a 100 uL rxn)

- [1167] 1. Normalize short crRNA to 36.5 ng/uL (5 uL in 8.7 uL ddH2O).
- [1168] 2. Dilute frozen protein. (8.1 uL in 21.9 H2O.) Need 24 (make 30) samples (1 uL per sample), 0.27 ug per rxn. New stock is at 1 ug/uL; 2000 mM NaCl.
- [1169] 3. Normalize each PAM library to 100 ng/uL (10 uL in 90 uL ddH2O)
- [1170] 4. Make 10x Buffer w/DTT (Add 10 uL DTT to 90 uL of 10x buffer without DTT)
- [1171] 5. Make Master Mixes:

EGFP Target 1

[1172]

RNase Inhibitor	1 uL	20 uL
10x Buffer w/ DTT	1 uL	20 uL
ddH2O	5 uL x 20	120 uL
E1 (100 ng/uL)	1 uL	20 uL
Total	8 uL	

EGFP Target 1

[1173]

RNase Inhibitor	1 uL	10 uL
10x Buffer w/ DTT	1 uL	10 uL
ddH2O	5 uL x 10	50 uL
E2 (100 ng/uL)	1 uL	10 uL
Total	8 uL	

Layout:

[1174]

Gel 1	1	2	3	4	5	6	7	8	9	10	11	12	13	14
Tem-plate	1	1	1	1	1	1	1	1	1	1	1	1	1	1
Pro-tein	+	+	+	+	+	+	+	+	+	+	+	+	+	-
RNA	T1	T2	T3	T4	T5	T6	11	12	13	14	15	16	-	-

Gel 2	1	2	3	4	5	6	7	8	9	10	11	12	13	14
Tem-plate	1	1	1	1	1	2	2	2	2	2	2	2	2	2
Pro-tein	+	+	+	+	-	+	+	+	+	+	+	+	+	-
RNA	T2	11	25	-	-	21	22	23	24	25	26	11	-	-

- [1175] 1. Aliquot 8 uL MM.
- [1176] 2. Add 1 uL of protein or ddH2O to each well.
- [1177] 3. Add 1 uL of short crRNA (36.5 ng/uL) or ddH2O to each well.
- [1178] 4. Mix components.
- [1179] 5. Incubate for 60 min at 37 C. Then heat inactivate for 5 minutes at 85 C.
- [1180] 6. Prot K treat: 1 uL per sample for 15 min @ RT
- [1181] 7. Run TBE urea gel

[1182] FIG. 39 shows the results of an in vitro cleavage assay of two different EGFP RNAs by Group 29 proteins (purified from *Bergeyella zoohelcum* ATCC 43767). EGFP RNA template 1 is about 290 nt long. EGFP RNA template 2 is about 240 nt long (see FIG. 42). Reactions were in 50 mM NaCl, pH 7.5 (10 mM Tris-HCl), 0.1% BSA, 1 mM DTT, and conducted for 60 min at 37° C. The molar ratio of protein to crRNA to target RNA was 1:1:1. Different crRNAs were evaluated, as indicated in FIG. 39. From FIG. 39 as well as Table 2, cleaving and non-cleaving guides (artificial sequences) can be distinguished.

TABLE 2

Guides that cut:			
Guide Index	Left Flank	Target	Right Flank
T2	GGCACAAGCT	CGAGTACAACACTACAACAGCCACAACGTCTA	TATCATGGCC
11	ACCTTGGTGA	ACCGCATCGAGCTGAAGGGCATCGACTTCA	AGGAGCACGG
12	ATCGACTTCA	AGGAGGACGGCAACATCCTGGGGCACAAAGC	TGGAGTACAA
13	TGGGGCACAA	GCTGGAGTACAACACTACAACAGCCACAACGT	CTATATCATG
25	ACCACATGAA	GCAGCAGGACTTCTTCAACTCCGCCATGCC	CGAAGGCTAC
26	GTCCGCCATG	CCCGAAGGCTACGTCCAGGAGCGCACCATC	TTCTTCAAGG
Guides that didn't visibly cut:			
Guide Index	Left Flank	Target	Right Flank
T1	GCATCGACTT	CAAGGAGGACGGCAACATCCTGGGGCACAA	GCTGGAGTAC
T3	ACGGCAACAT	CCTGGGGCACAAAGCTGGAGTACAACATAAA	CAGCCACAAC
T4	AAGCTGGAGT	ACAACATAACAGCCACAACCTCTATATCA	TGGCCGACAA
T5	CGAGCTGAAG	GGCATCGACTTCAAGGAGGACGGCAACATC	CTGGGGCACAA
T6	GGGCACAAGC	TGGAGTACAACACTACAACAGCCACAACGTCT	ATATCATGGC
14	AGCCACAACG	TCTATATCATGGCCGACAAGCAGAAGAACG	GCATCAAGGT
15	GCAGAAGAAC	GGCATCAAGGTGAACTTCAAGATCCGCCAC	AACATCGAGG
16	AAGATCCGCC	ACAACATCGAGGACGGCAGCGTGCAGCTCG	CCGACCACTA
21	CACCTACGGC	AAGCTGACCCCTGAAGTTCATCTGCACCACC	GGCAAGCTGC
22	ATCTGCACCA	CCGGCAAGCTGCCCGTGCCTGGCCACCC	TCGTGACCAC
23	TGGCCACCC	TCGTGACCACCTGACCTACGGCGTGCAGT	GCTTCAGCCG
24	TACGGCGTGC	AGTGCTTCAGCCGCTACCCGACCACATGA	AGCAGCACGA

[1183] Protocol of FIG. 40: Protocol for all N² tiling experiments. Now salt concentration will be 50 mM (2.70 μ L of 2M NaCl in a 1004, rxn)

[1184] 6. Normalize short crRNA to 36.5 ng/ μ L (5 μ L in 8.7 μ L ddH₂O).

[1185] 7. Dilute frozen protein. (8.1 μ L in 21.9 H₂O.) Need 24 (make 30) samples (1 μ L per sample), 0.27 μ g per rxn. New stock is at 1 μ g/ μ L; 2000 mM NaCl.

[1186] 8. Normalize each PAM library to 100 ng/ μ L (5 μ L in 45 μ L ddH₂O)

[1187] 9. Make 10x Buffer w/DTT (Add 10 μ L DTT to 90 μ L of 10x buffer without DTT)

[1188] 10. Make Master Mixes:

EGFP Target 1

[1189]

RNase Inhibitor	1 μ L	30 μ L
10x Buffer w/ DTT	1 μ L	30 μ L
ddH ₂ O	5 μ L \times 30	150 μ L
E1 (20 ng/ μ L)	1 μ L	30 μ L
Total	8 μ L	

[1190] 1. Aliquot 8 μ L MM.

[1191] 2. Add 1 μ L of protein or ddH₂O to each well.

[1192] 3. Add 1 μ L of short crRNA (36.5 ng/ μ L) or ddH₂O to each well.

[1193] 4. Mix components.

[1194] 5. Incubate for 60 min at 37 C. Then heat inactivate for 5 minutes at 85 C.

[1195] 6. Prot K treat: 1 μ L per sample for 15 min @ RT

[1196] 7. Run TBE urea gel

[1197] FIG. 40 shows the results of an in vitro cleavage assay of an EGFP RNA (A: Target 1; B: Target 2; see FIG. 41; 100 ng) by Group 29 proteins (purified from *Bergeyella zoohelcum* ATCC 43767). 20 targets chosen randomly while represent N² nucleotides (i.e. AA, AC, AT, AG, . . . , GT, GG) on either side of 30 nucleotide target sequences. Reactions were in 50 mM NaCl, pH 7.5 (10 mM Tris-HCl), 0.1% BSA, 1 mM DTT, and conducted for 60 min at 37° C. The molar ratio of protein to crRNA to target RNA was 1:1:1. Different crRNAs were evaluated, as indicated in FIG. 39. From FIG. 40 as well as Table 3 (Target 1) and Table 4 (Target 2), cleaving and non-cleaving guides (artificial sequences) can be distinguished.

TABLE 3

Guides that cut:		
Guide Index	Left Flank	Target
1	CCTGGGGCCACAAGCTGGAGTACAACCTACAA	CAGCCACAACGTCTATATCATGGCCGACAA
4	ATCGAGCTGAAGGGCATCGACTTCAAGGAG	GACGGCAACATCCTGGGGCCACAAGCTGGAG
8	AATACAACAGCCACACAGTCTATATCATG	GCCGACAAGCAGAAGAACGGCATCAAGGTG
9	TGAAGGCCATCGACTTCAAGGAGGACGGCA	ACATCCTGGGGCACAAGCTGGAGTACAAC
10	ACATCCTGGGGCACAAGCTGGAGTACAAC	ACAACAGCCACAACGTCTATATCATGGCCG
13	ACGTCTATATCATGGCCGACAAGCAGAAGA	ACGGCATCAAGGTGAACCTCAAGATCCGCC
16	GCATCGAGCTGAAGGGCATCGACTTCAAGG	AGGACGGCAACATCCTGGGGCACAAGCTGG
19	CAACTACAACAGCCACAACGTCTATATCAT	GGCCGACAACGAGAAGAACGGCATCAAGGT

Guides that cut:		
Guide Index	Right Flank	Target Position
1	GCAGAAGAACGGCATCAAGGTGAACCTCAA	114
4	TACAACCTACAACAGCCACAACGTCTATATC	73
8	AACTTCAAGTCCGCCACAACATCGAGGAC	136
9	ACAACAGCCACAACGTCTATATCATGGCCG	80
10	ACAAGCAGAAGAACGGCATCAAGGTGAAC	110
13	ACAACATCGAGGACGGCAGCGTGCAGCTCG	152
16	AGTACAACCTACAACAGCCACAACGTCTATA	71
19	GAACTTCAAGTCCGCCACAACATCGAGGA	135

Guides that didn't visibly cut:		
Guide Index	Left Flank	Target
2	GTACAACCTACAACAGCCACAACGTCTATAT	CATGGCCGACAAGCAGAAGAACGGCATCAA
3	GGCAACATCCTGGGGCACAAGCTGGAGTAC	AACTACAACAGCCACAACGTCTATATCATG
5	GGAGTACAACCTACAACAGCCACAACGTCTA	TATCATGGCCGACAAGCAGAAGAACGGCAT
6	CAAGCAGAAGAACGGCATCAAGGTGAACCT	CAAGATCCGCCACAACATCGAGGACGGCAG
7	AAGAACGGCATCAAGGTGAACCTCAAGATC	CGCCACAACATCGAGGACGGCAGCGTGCAG
11	TACAACAGCCACAACGTCTATATCATGGCC	GACAAGCAGAAGAACGGCATCAAGGTGAAC
12	AGCTGGAGTACAACCTACAACAGCCACAACG	TCTATATCATGGCCGACAAGCAGAAGAACG
14	ACGGCAACATCCTGGGGCACAAGCTGGAGT	ACAACCTACAACAGCCACAACGTCTATATCA
15	TATATCATGGCCGACAAGCAGAAGAACGGC	ATCAAGGTGAACCTCAAGATCCGCCACAAC
17	TGGAGTACAACCTACAACAGCCACAACGTCT	ATATCATGGCCGACAAGCAGAAGAACGGCA
18	GCAACATCCTGGGGCACAAGCTGGAGTACA	ACTACAACAGCCACAACGTCTATATCATGG
20	TTCAAGGAGGACGGCAACATCCTGGGGCAC	AAGCTGGAGTACAACCTACAACAGCCACAAC

Guides that didn't visibly cut:		
Guide Index	Right Flank	Target Position
2	GGTGAACCTCAAGATCCGCCACAACATCGA	132
3	GCCGACAAGCAGAAGAACGGCATCAAGGTG	106
5	CAAGGTGAACCTCAAGATCCGCCACAACAT	129
6	CGTGCAGCTCGCCGACCACTACCAGCAGAA	171
7	CTCGCCGACCACTACCAGCAGAACACCCCC	178
11	TTCAAGATCCGCCACAACATCGAGGACGGC	139
12	GCATCAAGGTGAACCTCAAGATCCGCCACA	125
14	TGGCCGACAAGCAGAAGAACGGCATCAAGG	104
15	ATCGAGGACGGCAGCGTGCAGCTCGCCGAC	157
17	TCAAGGTGAACCTCAAGATCCGCCACAACA	128
18	CCGACAAGCAGAAGAACGGCATCAAGGTGA	107
20	GTCTATATCATGGCCGACAAGCAGAAGAAC	94

TABLE 4

Guides that cut:		
Guide Index	Left Flank	Target
5	GACCTACGGCGTGCAGTGCTTCAGCCGCTA	CCCCGACCACATGAAGCAGCAGCACTTCTT
6	GACCAACCTGACCTACGGCGTGCAGTGCTT	CAGCCGCTACCCCGACCACATGAAGCAGCA
8	TTGATCTGCACCCACCGGCAGCTGCCCGTG	CCCTGGCCACCCCTCGTGACCACCTGACC
9	CTTACGCCGCTACCCCGACCACATGAAGCA	GCACGACTTCTTCAAGTCCGCCATGCCCGA
12	GGCCACCCCTCGTGACCACCCTGACCTACG	GCGTGCAGTGCTTCAGCCGCTACCCCGACC
13	CCGCTACCCCGACCACATGAAGCAGCACGA	CTTCTCAAGTCCGCCATGCCCGAAGGCTA
14	CAAGTGCCCCGTGCCCTGGCCACCCTCGT	GACCACCCTGACCTACGGCGTGCAGTGCTT
17	ACCCCTCGTGACCACCCTGACCTACGGCGTG	CAGTCTTCAGCCGCTACCCCGACCACATG
19	ACGGCAAGCTGACCCCTGAAGTTCATCTGCA	CCACCCGCAAGCTGCCCGTGCCCTGGCCCA

Guides that cut:		
Guide Index	Right Flank	Target Position
5	CAAGTCCGCCATGCCCGAAGGCTACGTCCA	174
6	CGACTTCTTCAAGTCCGCCATGCCCGAAGG	165
8	TACGGCGTGCAGTGCTTCAGCCGCTACCCC	118
9	AGGCTACGTCCAGGAGCGCACCATCTTCTT	192
12	ACATGAAGCAGCAGACTTCTTCAAGTCCG	152
13	CGTCCAGGAGCGCACCATCTTCTTCAAGGA	198
14	CAGCCGCTACCCCGACCACATGAAGCAGCA	135
17	AAGCAGCAGACTTCTTCAAGTCCGCCATG	157
19	CCCTCGTGACCACCCTGACCTACGGCGTGC	98

Guides that didn't visibly cut:		
Guide Index	Left Flank	Target
1	GACCCTGAAGTTCATCTGCACCACCGGCAA	GCTGCCCGTGCCCTGGCCACCCTCGTGAC
2	GCAGTGCTTCAGCCGCTACCCCGACCACAT	GAAGCAGCAGCACTTCTTCAAGTCCGCCAT
3	GCCCTGGCCACCCTCGTGACCACCCTGAC	CTACGGCGTGCAGTGCTTCAGCCGCTACCC
4	CACCCTGACCTACGGCGTGCAGTGCTTCAG	CCGCTACCCCGACCACATGAAGCAGCAGCA
7	ACCACCTGACCTACGGCGTGCAGTGCTTC	AGCCGCTACCCCGACCACATGAAGCAGCAGC
10	GCTACCCCGACCACATGAAGCAGCAGCACT	TCTTCAAGTCCGCCATGCCCGAAGGCTACG
11	ATGAAGCAGCAGCACTTCTTCAAGTCCGCC	ATGCCCGAAGGCTACGTCAGGAGCGCACC
15	GCACCACCGGAAGCTGCCCGTGCCCTGGC	CCACCCTCGTGACCACCCTGACCTACGGCG
16	TGCACCACCGGAAGCTGCCCGTGCCCTGG	CCCACCCTCGTGACCACCCTGACCTACGGC
18	ACCCTGACCTACGGCGTGCAGTGCTTCAGC	CGCTACCCCGACCACATGAAGCAGCAGCAGC
20	CCACCTACGGCAAGCTGACCCTGAAGTCA	TCTGCACCACCGGAAGCTGCCCGTGCCCT

Guides that didn't visibly cut:		
Guide Index	Right Flank	Target Position
1	CACCCTGACCTACGGCGTGCAGTGCTTCAG	108
2	GCCCGAAGGCTACGTCCAGGAGCGCACCAT	186
3	CGACCACATGAAGCAGCAGCACTTCTTCAA	147
4	CTTCTTCAAGTCCGCCATGCCCGAAGGCTA	168
7	GACTTCTTCAAGTCCGCCATGCCCGAAGGC	166
10	TCCAGGAGCGCACCATCTTCTTCAAGGACG	200
11	ATCTTCTTCAAGGACGACGGCAACTACAAG	214
15	TGCAGTGCTTCAGCCGCTACCCCGACCACA	125
16	GTGCAGTGCTTCAGCCGCTACCCCGACCAC	124
18	TTCTTCAAGTCCGCCATGCCCGAAGGCTAC	169
20	GGCCACCCTCGTGACCACCCTGACCTACG	92

[1198] From Tables 1, 2, and 3, it appears that if the target sequence is preceded by a C (i.e. "C" directly 5' of the target sequence) no cleavage of the template occurs. From Tables 1 and 2, it appears that if the target sequence is followed by an "A" (i.e. C directly 3' of the target sequence) cleavage of the template is more likely.

[1199] FIG. 42 shows the location of the respective guide sequences leading to target cleavage.

[1200] Cleavage dependence seems not sequence-specific, nor secondary structure dependent, nor position dependent along the RNA targets.

[1201] Similar cleavage intensities for EGFP target 2: #8 and #17, with their first four 5' proximal nucleotides identical and their target positions 39 nucleotides away, give credibility to a 5' proximal cleavage dependence.

[1202] Without wishing to be bound by theory, it may be that perhaps RNA cleavage activity may (at least partially) be different from DNA cleavage activity, it that at small scales DNA secondary structure is much more fixed than RNA secondary structure (DNA nearest neighbor thermodynamic model based on stacking energy much more accurate than that of RNA), which is why the presence of a strong or well conserved PAM sequence is preferred for DNA targeting. For RNA targeting, whose secondary structure is much more variable at the small scales relevant to protein interaction, activity may possibly not be (or not only be) dependent on RNA sequence per se, but rather (or also) RNA secondary structure. The following RNA efficiency screen is employed to test this hypothesis. Infect cells with MS2 (or QBeta, which has a larger genome). For a given species of Group 29 protein, design a crRNA library off of which one can prime for NGS. The crRNA library consists of targets that tile the entire RNA viral genome one base pair at a time, and are transformed with the Group 29 protein into the bacteria prior to viral infection. Perform positive selection screen where most efficient targets for cleavage are overrepresented over control crRNA library when harvested from bacterial colonies that survive viral infection. Some areas of viral genome may be more tolerant to cleavage. Given the single base pair resolution of the screen, this tolerance from cleavage efficiency can be spatially decon-

volved. Further biophysical computational analysis reveals whether it is sequence or secondary structure that determines efficiency of cleavage. Without wishing to be bound by theory, it may be that perhaps RNA cleavage activity may (at least partially) be different from DNA cleavage activity, it that at small scales DNA secondary structure is much more fixed than RNA secondary structure (DNA nearest neighbor thermodynamic model based on stacking energy much more accurate than that of RNA), which is why the presence of a strong or well conserved.

Example 3: Group 30 Proteins

[1203] To learn more about the small proteins, Applicants completed a BLAST protein search. Applicants were interested to see if these small proteins are in other genomes or within other loci. For S2 small proteins, the only other significantly similar proteins were other S2 proteins. For S1, there was one new species *Flavobacterium branchiophilum* (Fb) containing it.

[1204] When the Fb small protein was examined, it was found to be next to a large (1,151 aa) protein next to a CRISPR array. The Fb large protein has no homology to other known proteins associated with a CRISPR array. A BLAST search revealed only three known species with this large protein: Fb, *Paludibacter propionicigenes*, and *Bacteroides coprosuis*. Only a small region of homology was seen—a 34 aa sequence with homology to a L-lysine-6-monooxygenase.

[1205] It seems that the S1a small protein is associated with CRISPR loci that have dual repeat. Applicants will further investigate why there are only three different proteins. This could be a very specific system, a mutated or dead system, or in bacteria types that have little sequencing coverage on NCBI or other databases. Applicants will further investigate the region with homology to L-lysine-6-monooxygenase. Applicants identified 5 proteins in total with significant sequence homology: Fb, Pp, and three Bc protein truncations that were less than 800 aa (732, 680, & 370 aa). Where an organism is not indicated the sequences are artificial.

BLASTP 2.2.32+
Reference: Stephen F. Altschul, Thomas L. Madden, Alejandro A. Schaffer, Jinghui Zhang, Zheng Zhang, Webb Miller, and David J. Lipman (1997), "Gapped BLAST and PSI-BLAST: a new generation of protein database search programs", *Nucleic Acids Res.* 25:3389-3402.

Reference for compositional score matrix adjustment: Stephen F. Altschul, John C. Wootton, E. Michael Gertz, Richa Agarwala, Aleksandr Morgulis, Alejandro A. Schaffer, and Yi-Kuo Yu (2005) "Protein database searches using compositionally adjusted substitution matrices", *FEBS J.* 272:5101-5109.

RID: 2ANBZMGK013

Database: All non-redundant GenBank CDS translations + PDB + SwissProt + PIR + PRF excluding environmental samples from WGS projects

73,055,898 sequences; 26,613,539,447 total letters
Results of PSI-Blast iteration 1
Query = gi|503850672|ref|WP_014084666.1| hypothetical protein [*Flavobacterium branchiophilum*]

-continued

Length = 1150

E-value BETTER than threshold

Score E

Sequences producing significant alignments:

(Bits)

Value		
ref WP_014084666.1 hypothetical protein [Flavobacterium bran...	2271	
0.0		
ref WP_013446107.1 hypothetical protein [Paludibacter propio...	686	
0.0		
ref WP_040310142.1 hypothetical protein [Bacteroides coprosuis]	297	
7e-83		
gb EGJ71546.1 hypothetical protein Bcop_1349 [Bacteroides co...	288	
5e-80		
ref WP_006744704.1 hypothetical protein [Bacteroides coprosuis]	122	
7e-26		

ALIGNMENTS

>ref|WP_014084666.1| hypothetical protein [Flavobacterium branchiophilum]
emb|CCB70204.1| Hypothetical protein PBFL15_2182 [Flavobacterium
branchiophilum
FL-15]

Length = 1150

Score = 2271 bits (5884), Expect = 0.0, Method: Compositional matrix

adjust.

Identities = 1150/1150 (100%), Positives = 1150/1150 (100%), Gaps = 0/1150 (0%)

Query	1	MENLNKILDKENEICISKIENTKGI AAPITEKALDNIKSKQKNDLNKEARLHYFSIGHSF	60
Sbjct	1	MENLNKILDKENEICISKIENTKGI AAPITEKALDNIKSKQKNDLNKEARLHYFSIGHSF	60
Query	61	KQIDIKKVEDYVLI EELKDEKPLKFITLQKDFFTKEFSIKLQKLINSIRNINNHYVHNFN	
120			
Sbjct	61	KQIDIKKVEDYVLI EELKDEKPLKFITLQKDFFTKEFSIKLQKLINSIRNINNHYVHNFN	
120			
Query	121	DINLNKIDS NVEHELKESFELAI IEKYYKVNKKYPLDNEIVLFLKELFIKIDENTALLNYP	
180			
Sbjct	121	DINLNKIDS NVEHELKESFELAI IEKYYKVNKKYPLDNEIVLFLKELFIKIDENTALLNYP	
180			
Query	181	TNLSKDEAIEYILTFTITENKIWNINNEHNI LNI EKGKYLTFEAMLFLITIFLYKNEANH	
240			
Sbjct	181	TNLSKDEAIEYILTFTITENKIWNINNEHNI LNI EKGKYLTFEAMLFLITIFLYKNEANH	
240			
Query	241	LLPKLYDEKNNKSKQELFIFFSKKFTSQDIDAEEGHLIKERDMIQYLNHYPTAWNNDLKL	
300			
Sbjct	241	LLPKLYDEKNNKSKQELFIFFSKKFTSQDIDAEEGHLIKERDMIQYLNHYPTAWNNDLKL	
300			
Query	301	ESENKNKIMTTKLIDSII EFELNSNYPSFATDIQFKKEAKAFLFASNKKRNQTSFSNKS	
360			
Sbjct	301	ESENKNKIMTTKLIDSII EFELNSNYPSFATDIQFKKEAKAFLFASNKKRNQTSFSNKS	
360			
Query	361	NEEIRHNPHIKQYRDEIASALTPISFNVKEDKFKIFVKKHVLEEYFPNSIGYEKFLYND	
420			
Sbjct	361	NEEIRHNPHIKQYRDEIASALTPISFNVKEDKFKIFVKKHVLEEYFPNSIGYEKFLYND	
420			
Query	421	FTEKEKEDEGLKLYSNPKINKLIERIDNHKLVKSHGRNQDREMDFSMRFLAENNYFGKDA	
480			
Sbjct	421	FTEKEKEDEGLKLYSNPKINKLIERIDNHKLVKSHGRNQDREMDFSMRFLAENNYFGKDA	
480			
Query	481	FEKCYKEYDTQEQDEFLOSNENDDVKPHKGKVI TYIKYEEHLKNYSYWDCPFVEENNSM	
540			
		FEKCYKEYDTQEQDEFLOSNENDDVKPHKGKVI TYIKYEEHLKNYSYWDCPFVEENNSM	

-continued

Sbjct	481	FEKCYKEYDTQEQDEFLOSNNENDDVKFHKGKVI TYIKYEEHLKNYSYWDCCPFVEENNSM	
540			
Query	541	SVKISIGSEEKILKIQRNLMIYFLENALYNENVENQGYKLVNNYRELKKDVEESIASLD	
600		SVKISIGSEEKILKIQRNLMIYFLENALYNENVENQGYKLVNNYRELKKDVEESIASLD	
Sbjct	541	SVKISIGSEEKILKIQRNLMIYFLENALYNENVENQGYKLVNNYRELKKDVEESIASLD	
600			
Query	601	LIKSNPDFKSKYKKILPKRLLHNYAPAKQDKAPENAFETLLKKADFEEQYKLLKKAEH	
660		LIKSNPDFKSKYKKILPKRLLHNYAPAKQDKAPENAFETLLKKADFEEQYKLLKKAEH	
Sbjct	601	LIKSNPDFKSKYKKILPKRLLHNYAPAKQDKAPENAFETLLKKADFEEQYKLLKKAEH	
660			
Query	661	EKNKEDEVKRNKGKQFKLHFIRKACQMMYEKEKYNTLKEGNAAFEKKDPVIEKRKNKEHE	
720		EKNKEDEVKRNKGKQFKLHFIRKACQMMYEKEKYNTLKEGNAAFEKKDPVIEKRKNKEHE	
Sbjct	661	EKNKEDEVKRNKGKQFKLHFIRKACQMMYEKEKYNTLKEGNAAFEKKDPVIEKRKNKEHE	
720			
Query	721	FGHHKLNMITREEENDYCKWMFAENGND SYKKYLRDLFSEKHPEDNQEYKNLPESVNL E	
780		FGHHKLNMITREEENDYCKWMFAENGND SYKKYLRDLFSEKHPEDNQEYKNLPESVNL E	
Sbjct	721	FGHHKLNMITREEENDYCKWMFAENGND SYKKYLRDLFSEKHPEDNQEYKNLPESVNL E	
780			
Query	781	AFYAKTKELPKKWIETNKPINNENRYTLENYKNLILQKQVFINVYHFSKYLIDKNLLNSE	
840		AFYAKTKELPKKWIETNKPINNENRYTLENYKNLILQKQVFINVYHFSKYLIDKNLLNSE	
Sbjct	781	AFYAKTKELPKKWIETNKPINNENRYTLENYKNLILQKQVFINVYHFSKYLIDKNLLNSE	
840			
Query	841	NNVIQYKSL ENVEY LISDFYFQSKLSIDQYKTCGKLENKLSKNKLEDCLLYEIA YNYIDK	
900		NNVIQYKSL ENVEY LISDFYFQSKLSIDQYKTCGKLENKLSKNKLEDCLLYEIA YNYIDK	
Sbjct	841	NNVIQYKSL ENVEY LISDFYFQSKLSIDQYKTCGKLENKLSKNKLEDCLLYEIA YNYIDK	
900			
Query	901	KNVHKIDIQKILTSKIILTINDANTPYKISVPPFNKLERYTEMIAIKNQNNLKARFLIDL P	
960		KNVHKIDIQKILTSKIILTINDANTPYKISVPPFNKLERYTEMIAIKNQNNLKARFLIDL P	
Sbjct	901	KNVHKIDIQKILTSKIILTINDANTPYKISVPPFNKLERYTEMIAIKNQNNLKARFLIDL P	
960			
Query	961	LYLSKNKIKK GKDSAGYEII IKN DLEIEDINTINNKIINDSVKPFTEVLM ELEYFILKDK	
1020		LYLSKNKIKK GKDSAGYEII IKN DLEIEDINTINNKIINDSVKPFTEVLM ELEYFILKDK	
Sbjct	961	LYLSKNKIKK GKDSAGYEII IKN DLEIEDINTINNKIINDSVKPFTEVLM ELEYFILKDK	
1020			
Query	1021	CILSKNYIDNSEIPSLKQESK VVIKENENEI INYRN IACHFHLPLETEDNLLL NVEQKF	
1080		CILSKNYIDNSEIPSLKQESK VVIKENENEI INYRN IACHFHLPLETEDNLLL NVEQKF	
Sbjct	1021	CILSKNYIDNSEIPSLKQESK VVIKENENEI INYRN IACHFHLPLETEDNLLL NVEQKF	
1080			
Query	1081	IKEELQNVSTINDLSKPQEYLIL LFIKEKHNNFYLNLENKNESKTIKNDKEVKKNRVLQK	
1140		IKEELQNVSTINDLSKPQEYLIL LFIKEKHNNFYLNLENKNESKTIKNDKEVKKNRVLQK	
Sbjct	1081	IKEELQNVSTINDLSKPQEYLIL LFIKEKHNNFYLNLENKNESKTIKNDKEVKKNRVLQK	
1140			
Query	1141	FINQVILKKK	1150
		FINQVILKKK	
Sbjct	1141	FINQVILKKK	1150

>ref|WP_013446107.1| hypothetical protein [Paludibacter propionicipenes]
 gb|ADQ80738.11| hypothetical protein Palpr 2606 [Paludibacter propionicipenes
 WB4]
 Length = 1154
 Score = 686 bits (1769), Expect = 0.0, Method: Compositional matrix adjust.
 Identities = 465/1202 (39%), Positives = 685/1202 (57%), Gaps = 135/1202
 (11%)

Query	18	KIENTKGIAAPITEKALDNIKSKQKNDLNKEARLHYFSIGHSFKQIDTKKVEDYVLI EEL KIF++KG API EK+ N K +ND+NKE R+HYF++GH+FKQ+DT+ +F+YVL E L	77
-------	----	---	----

-continued

Sbjct	18	KIPDSKGAIAPIAEKSCRNEDIKAQNDVNKEQRIHYFAVGHTFKQLDTENLFEYVLDENL	77
Query	78	KDEKPLKFITLQKDFFTKEFSIKLQKLINSIRNINNHVHNENDINLNKIDSNVFHLKE	
137		+ ++P +FI+LQ+ F KEF +++LI+ IRNIN+HY+H F+ + ++ + +N+ FLKE	
Sbjct	78	RAKRPRTRFISLQQ--FDKEFIENIKRLISDIRNINSHYIHRFDPLKIDAVPTNIIDELKE	
135			
Query	138	SFELAIIEKYK-----VNKKYPLDNEIVLFLKELFIK--DENTALLN-----	
178		SFELA+I+ Y K ++ D ++V FL + F + ++ T++L	
Sbjct	136	SFELAVIQIYLKEKGINYLQFSENPHADQKLVAFPLHDKFLPLDEKKTSMQLQNETPQLKEY	
195			
Query	179	-----YFTNLSKDEAIEYILTFTITENKIWNINNEHNILNIEKGKYLTFEAMLFLITIFL	
233		YF LSK AI+ +L + IWN+ + H +L I GKYL+F + LFL+++FL	
Sbjct	196	KEYRKYPKTLKQAAIDQLLFAEKETDYIWNLPDHPVLTISAGKYSFYSCLELLSMFL	
255			
Query	234	YKNEANHLPLKLYDFKNN----KSKQELFTFFSKKFTSQDIDAEEGHLIKFRDMIQYLN	
288		YK+EAN L+ K+ FK N KSK+E+FTFFSK+F S DID+EE L+KPRD+I YLN	
Sbjct	256	YKSEANQLISKIKGFKKNTTEEEKSKREIFTFFSKRFSMDIDSSEENQLVKPRDLILYLN	
315			
Query	289	HYPTAWNNDLKLESENKNKIMTTKLIDSIEFELNSNYP-----SFATDIQFKKEAKAF	
342		HYP AWN DL+L+S N MT KL IIE E+N ++P FAT +++ K	
Sbjct	316	HYPVAWNKDLELDS--NPAMTDKLSKIELEINRSFPLYEGNERFATFAKYQIWGKKH	
373			
Query	343	LFAS-NKKRNQTSFSNK--SYNEEIRHNPHIKQRDEIASALTPISFNVKEDKFKIFVK	
398		L S K+ SF+++ +Y E P +K ++A ++K K +F K	
Sbjct	374	LGKSIEKEYINASFTDEBITAYTYETDTCPELKAHKKLA-----DLKAAK-GLFGK	
424			
Query	399	KHVLEEYFPNSIGYEKFLEYNDPTEKEKEDFGLKLYSNPKTNKLIERIDNHKLVKSHGRN	
458		+ E +D +E L+ NP +KLI+RI+ +L S+GRN	
Sbjct	425	RKEKNE-----SDIKKTETSIRELQHEPNPIKDKLIQRIEKNLLTVSYGRN	
470			
Query	459	QDRFMDFSMRFLAENNYFGKDAFFKCYKFYDTQEQEFLQSNE-----NNDVVKFHKGK	
512		QDRFMDFS RFLAE NYFG+DA FK Y FY T EQ+ L+ E D +KPH+GK	
Sbjct	471	QDRFMDFSARFLAEINNYFGQDASFKMYHFYATDEQNSELEKYELPKDKKKYDSLKPHQGK	
530			
Query	513	VITYIKYEEHLKNYSYWDPPVEENNSMSVKISIGSEEKILKIQRNLMYFLENALY---	
569		+ +I Y+EHLK Y WD FV ENN++ +K+S E + IQR L+IY LE+AL	
Sbjct	531	LVHFISYKEHLKRYESWDDAFVIENNAIQLKLSFDGVENTVTIQRALLIYLLEDALRNIQ	
590			
Query	570	NENVENQGYKLVNYYRELKDV--EESIASLDLIXSNPDPFSKYKILPKRLLHNYAP	
626		N EN G +L+ YY K D+ ++ + D I+ P K+++KK+LP+RLL+NY+P	
Sbjct	591	NNTAENAGKQLLQYYSYSHNKADLSAFKQILTQQDSIE--PQQKTEFKKLLPRLLNLYSP	
648			
Query	627	A-KQDKAPENAFEILLKKADFREEQYKLLKKAHEKKNKEDFVKRNKGKQFKLHPIRKAC	
685		A + P ++ +L+KA E++Y L+ KA+ E N +DF+KRNKGKQFKL FIRKA	
Sbjct	649	AINHLQTPHSSLPLILEKALLAEKRYCSLVVKAAGAEGNYDDFIKRNKGKQFKLQFIRKAW	
708			
Query	686	QMMYFKEKYNTLKEGNAAFEKDPVIEKRKNKEHEFGHHKNLNI TREEFNDYCKWMPAFN	
745		+MYF+ Y L+ AA GHKK+ +I R+EFND+ ++MPAF	
Sbjct	709	NLMYFRNSY--LQNVQAA-----GHHKSFHIERDEFNDFSRYPFAFE	
748			

-continued

Query	746	GND5YKKYLRDLFSEKHFFDNQEQYKNLFESSVNLEAFYAKTKELFKKWIETNKP-TNNEN	
Sbjct	749	YK YL ++F +K FF+N E+K LF+S +LE Y KTK+ F+ W+ +N TN +	
804		EL5QYKYILNEMFEKKGFFENNEFKILFQSGTSLLENLYEKTQKFEIWLASNTAKTNKPD	
808			
Query	805	RYTLENYKNLILQKQVFINVYHFSKYLIDKNLNLSE--NNVIQYKSLNVEYLISDFYFQS	
Sbjct	809	Y L NY+ + FIN+ HF YL L ++ N I Y++L NV+YLI ++Y+	
863		NYHLNNYEQQFSNQLFFINLSHFINYLKSIGKQLTDANGQIIYEALNNVQYLIPEYYTD	
868			
Query	864	KLSIDQYKTCGKLFNKLKSNKLEDCLLYEIAYNYI--DKKNVHKI--DIQKILTSKIILT	
Sbjct	869	K + K+ KL+NKLK+ KLED LLYE+A Y+ DK+ K I K+LTS +	
919		KPERSESKSGNKLYNKLKATKLEDAALYEMAMCYLKADKQIADKAKHPITKLLTSDVEFN	
928			
Query	920	IN--DANTPYKISVFPNKLERYTEMIAIKNQNNLK--ARFLIDLPLYLSKNKIKKGD-S	
Sbjct	929	I + Y + VPF K++ + + K Q + K FL ++ YL +K KD	
974		ITNKEGIQLYHLLVFPKIDAFIGLKMHEQQDKKHPTSPFLANIVNYLE--LVKNDKDIR	
986			
Query	975	AGYEII----IKNDLEIEDINTINNKI INDSVKFTEVLMLEKEYFILKDKCILSK-NYID	
Sbjct	987	YE +K L +D+ I+ +I +S+KFT V +ELE+YFI K+ I+ K N ID	
1029		KTYEAFSTNPVKRTLTYDDLAKIDGHLISKS IKFTNVTLERLYFIFKESLIVKKGNNID	
1046			
Query	1030	NSEIPLSKQFSKQVWIKENENEIINRYNIACHFHLPLETPDNLLLNVQKFIKEELQ--N	
Sbjct	1047	I L+ + K+NE RN A HF +P +++D L+ + E FI E++ +	
1087		FKYIKGLRNYNNEKKNNE---GIRNKAPHFGIPDSKSYDQLIRDAEVMFIANEVKEPTH	
1102			
Query	1088	VSTINDLSKQEQYLILLFIKFKHNNFYLNLFNKNESKTIKNDKEVKKNRVLQKFINQVIL	
Sbjct	1103	+ DL+K + ++ HN+++ + + D + K+ QK+ +I	
1147		ATKYTDLNLKQLHTVCDKLMETVHNDYF-----SKEGDGKKKREAAQKYFENIIS	
1152			
Query	1148	KK	1149
Sbjct	1153	AK	1154

>ref|WP_040310142.1| hypothetical protein [*Bacteroides coprosuis*]
 Length = 732
 Score = 297 bits (761), Expect = 7e-83, Method: Compositional matrix
 adjust.
 Identities = 213/668 (32%), Positives = 348/668 (52%), Gaps = 70/668 (10%)

Query	451	LVKSHGRNQRDFMDFSMRFLAENNYFGKDAFFKCYKFYDTQEQDEFLQSNENN-----D	
Sbjct	41	L +R QD+FM ++ R+LAE+NYFG++A FK Y+F +EQ++++ + N D	
504		LWSKYSRKQDKFMLYASRYLAESNYFGEEAMFKVYQFASNEEQEKYIVEAKQNLPKREYD	
100			
Query	505	DVKPHKGVTTYIKYEEHLKNYSYWDCCPFVEENNSMSVKISIGSEEKILKIQRNLMIIYFL	
Sbjct	101	+K+HKG++ Y Y HL +Y WD PFV ENN++ + + I E I+ IQR L+IYFL	
564		KLKYHKGRLVVYKSYHNLQEYPRWDYPFVVENNAIQIYVKILGEPWIVSIQRRLIIYFL	
160			
Query	565	ENALYNENVENQGYKLVNNYRELKKDVEESIASLDLISN---PDFKSKYKILPKRLL	
Sbjct	161	E+AL+++ E+ G L+ NY ++DV + ++N S +K+ P++L+	
621		EDALFSKKKESNGIALLQNYLPHHQDRVNRGLFVFKTGQTNLSTKEMSNLKLFPRKLI	
220			
Query	622	HNYAPAKQDKAPENAFETLLKKADFREEQYKLLKKAHEKKNKEDFVKNKGGKQKLFHFI	
Sbjct	221	+Y + + ++ +L + + +K + + GK KL +I	
681		QSYL-YEDNTGDMDSQVLSDTISINDTE-----KKGTKKILNLRVKGHLKLYI	
269			

-continued

Query	682	RKACQMMYFKEKYNTLKEGNAAFEKDPVIEKRKNKEHEFGHHKLNITREEFN DYCKWM	
741			
		RK ++YFK+ Y K+K GHK +IT++EF Y +WM	
Sbjct	270	RKVWNLIYFKDIY KDKAQRMGHKKPHITKDEFV FYTRWM	
309			
Query	742	FAFNGNSYKKYLRDLFSEKHFFDNQYK NLFESSVNLEAFYAKTKELFKKIETN-KPT	
800			
		++F SYK +L F +KHFF+N+E+K LF +S +++ Y +TK F KW N	
Sbjct	310	YSPESIPSYKDHLIQFFIKKHFFNNEEFKELFLN SSSIDELYLQTKRNF IKWSAHNVNSE	
369			
Query	801	NNENRYTLENYKNLILQKQVFINVYHFSKYL-IDKNLLNSENNVIQYK SLENVEYLI SDF	
859			
		E Y+LE+YK K ++INV HF +L +K + ++N +IQYK+L+N+ YLI F	
Sbjct	370	KKEKTYSLYEDYKLFPEKILYINVSHFISPLNQEKVIQKNDNGI IQYKALKNLSYLIKPF	
429			
Query	860	YFQSKLSIDQYKTCGKLFNKLKSNKLEDCLLYEIAYNYI-----DKKNVHKIDIQKILTS	
914			
		Y++ KL I+ YKT GK+FNKL+S KLEDCLLYEIAY Y+ ++ IQ	
Sbjct	430	YYKDKLEIEHYKTYGKVFNKLRSIKLEDCLLYEIAYRYLLNVTPSPFKYKQLIIQSFPKE	
489			
Query	915	KIILTINDANT PYKISVPPNKLERYTEMIAIKQN NLKAR-----FLIDL P	
960			
		K+ L +N + Y I VPF KL +I +N + A + +	
Sbjct	490	KVDLLVNAIYSFEIHNKKGAFIYSIQVPFLKLNELVCLIIY-RNSTKIAATNKEFLFLQIY	
548			
Query	961	LYLSKNKIKKKGDSAGYEIIK-NDLEI--EDINTINNKIINDSVKFTEVLM ELEYFI	
1016			
		YL K D Y + K N L++ +D+ I+ ++ T++L +LEK+I	
Sbjct	549	KYLVNYCKNKPVDYELYTVCYKFNQLKVLGYDDLLHFLKHIVKRGLQTLTQILTQLEKFLI	
608			
Query	1017	LKDKCILSKNYIDNSEIPSLKQFSKVIKENENEIINRYNIACHFHLPLLET PFDNLLLVN	
1076			
		+K+ + S+ + SK+ N ++ + R A +F +P + + ++L ++	
Sbjct	609	IKNNIQIDIQKQGTLSLSIYEC SKM---NNPQKLEDLRIKAINFDIPD TD-YPSILEHI	
664			
Query	1077	EQKFIKEE	1084
		E++FI +E	
Sbjct	665	EKQPIIKE	672
<p>>gb EGJ71546.1 hypothetical protein Bcop_1349 [<i>Bacteroides coprosuis</i> DSM 18011] Length = 680 Score = 288 bits (736), Expect = 5e-80, Method: Compositional matrix adjust. Identities = 208/656 (32%), Positives = 342/656 (52%), Gaps = 70/656 (11%)</p>			
Query	463	MDFSMRFLAENNYFGKDAFFKCYK FYDTQE QDEF LQSNENN DDVKPHKGVTTY	60
516			
		M ++ R+LAE+NYFG++A FK Y+F +EQ++++ + N D +K+HKG++ Y	
Sbjct	1	MLYASRYLAESNYPGEEAMFKVYQFASNEEQEKYIVEAKQNL PKREYDKLKYHKGR LVVY	
Query	517	IKYEEHLKNYSYWDPCPFVEENNSMSVKISIGSEEKILKIQRNLMIYPLEN ALYNENVENQ	
576			
		Y HL+ Y WD PFV ENN++ + + I E I+ IQR L+IYFLE+AL+++ E+	
Sbjct	61	KSYHNHLQEQYPRWDYPFVVENNAIQIYVKILGEPWIVSIQRRLIIYFLEDALF SKKESN	
120			
Query	577	GYKLVNMYRELKDDVEESIASLDLIKSN---PDFKSKYKKILPKRLLHNYAPAKQDKAP	
633			
		G L+ NY ++DV + ++N S +K+ P++L+ +Y + +	
Sbjct	121	GIALLONYLPHHQDRVNRGLFVFKTGQTNNLSTKEMS NLRKLPKRKLIQSYL-YEDNTGD	
179			
Query	634	ENAFETLLKADFREEQYKLLKKAHEKKNKEDFV KRNGKQFKLHPIRKACQMMYFPEK	
693			
		++ +L + + +K + + GK KL +IRK ++YFK+	
Sbjct	180	MDSPSQVLS DTS INDTE-----KKGTKKILNLFVGGKHLKLR YIRKVWNLIYFKDI	
229			

-continued

Query	694	YNTLKEGNAAPEKKDPVIEKRKNKEHEFGHHKLNLI TREEFNDYCKWMFAFNGNDSYKKY	
753		Y	
		K+K GHHK +IT++EF Y +WM++F SYK +	
Sbjct	230	Y-----KDKAQRMGHHKKFHI TKDEFVFTYTRWMYSFESIPSYKDH	
269			
Query	754	LRDLFSEKHFFDNQEQYKNLFESSVNLEAFYAKTKELFKKIETNKPTNN-ENRYTLENYK	
812		L F +KHF+N+E+K LF +S +++ Y +TK F KW N + E Y+LE+YK	
Sbjct	270	LIQFFIKKHFNNNEEFKELFLNSSSIDELYLQTKRNFIKWSAHNVNSEKKEKTYSLLEDYK	
329			
Query	813	NLILQKQVFINVYHFSKYL-IDKNLLNSENNVIQYKSLNVEYLI SDFYFQSKLSIDQYK	
871		K ++INV HF +L +K + ++N +IQYK+L+N+ YLI FY++ KL I+ YK	
Sbjct	330	LFPEFKILYINVSHFISFLNQEKVIQKNDNGIIQYKALKNLSYLIKPFYKDKLEIEHYK	
389			
Query	872	TCGKLFNKLKSNKLEDCLLYEIAYNYI DKKNVHKIDIQKILTSKIIILTINDANT-	
925		T GK+FNKL+S KLEDCLLYEIAI Y+ ++ IQ K+ L +N +	
Sbjct	390	TYGKVFNKLRSIKLEDCLLYEIAIRYLLNVTSPFPKYQLIIQSFPKEKVDLVNAIYSF	
449			
Query	926	-----PYKISVPFNKLERYTEMIAIKQNQLKAR----FLIDLPLYLSKNKIKKGGK	
972		Y I VPF KL +I +N + A + + YL K	
Sbjct	450	EIHNNKGAFIYSIQVPLKLNELVCLIIY-RNSTKIAATNKEFLFIQIYKYLNVNYCKNKPV	
508			
Query	973	DSAGYEIIK-NDLEI---EDINTINNKI INDSVKFTEVLMELEKYFILKDKCILSKNYI	
1028		D Y + K N L++ +D+ I+ ++ T++L +LEK+ I+K+ +	
Sbjct	509	DYELYTVCYKFNQLKVLGYDDLHFLKHIVKRGQLTQILTQLEKFLIKNNIQIDIQKQ	
568			
Query	1029	DNSEIPSLKQFSKVIKENENEI INYRNIACHFHLPLETFDNLNLLNVEQKFIKKE	1084
		S+ + SK+ N ++ + R A +F +P + + ++L ++E++FI +E	
Sbjct	569	GTLSLSLIYEC SKM--NPNQKLEDLRKAINFDIPD TD-YPSILEHIEKQFIKKE	620
>ref WP_006744704.1 hypothetical protein [Bacteroides coprosuis]			
gb EGJ71547.1 hypothetical protein Bcop_1350 [Bacteroides coprosuis DSM 18011]			
Length = 370			
Score = 122 bits (305), Expect = 7e-26, Method: Compositional matrix adjust.			
Identities = 101/328 (31%), Positives = 170/328 (52%), Gaps = 30/328 (9%)			
Query	23	KGIAAPITEKALDNIK-----SKQKNDLNKEA--RLHYFSIGHSFKQIDTKKVPDYVLIE	75
		KG+ AP++E A N + +++ DL+ +A ++ FSIG +FK +K++F+Y L +	
Sbjct	18	KGLLAPLSEIAYRNFELVCDNNEKSGDLSLQAISSIIYQFSIGQTFKSHNIKQLFNQYQLND	77
Query	76	ELKDEKPLKFITLQK-DFFTKEFSIKLQKLINSIRNINNHVHNFNDINL--NKIDSNVF	132
		E P K+++LQK F E + L KL+++IR +N +Y HN + + N I +	
Sbjct	78	EKDRFVPTKYLSLQKKQFIDNEIASTLLKLVSAIRELNQNYTHNLEPLRIGNNII TPQII	137
Query	133	HFLKESFELAIIEKYYKVNK-----KYPLDNEIVL---FLKELFIK DEN--TALLNYF	180
		FL + FE++II K K +Y D+ +L F+ LF + + T +	
Sbjct	138	EFLHDLFEVSIIMLLNKS VKDRNFTSQYNEDSLNLLLKEFILTLPFPETS YTTTEEIEVL	197
Query	181	TNLSKDEAIEYILTFTITENKIWINNEHNI LNIEKGKYLTFEAMLPLITIFLYKNEANH	240
		SKDE I +L F + W +N ++ I+ GKYL+F + LF+ ++FL+K +	
Sbjct	198	RKKS KDELINL E V L F F D V Q S V Y Q W K V V N P L M M P I Q C G K Y L S F T S C L F I N S L F L F K E D T K I	257
Query	241	LLPKLYDFKNKNSK-----QELFTFFSKKFTSQDIDAEEGHLIKFRDMIQYLNHYPTAWN	295
		+ P + F NK+ Q F+ F+ +T I ++ + K +D+I+YLN YP+ W	
Sbjct	258	IFPN-FPFLKNKTDEBTQVLQMFPSLFANPYTLHYIPSOHLRIAKHKDII EYLNLYPSFPWQ	316

-continued

Query	296	NDLKLESENKKNKIMTTKLIDSIIIEFELN	323
		L S++ M+ L D +IE E N	
Sbjct	317	E--ALNSQSPCFPMShLLKDFLIERECN	342

Database: All non-redundant GenBank CDS translations + PDB + SwissProt + PIR + PRF
 excluding environmental samples from WGS projects
 Posted date: Oct. 11, 2015 2:03 PM
 Number of letters in database: 26,613,539,447
 Number of sequences in database: 73,055,898

Lambda	K	H
0.317	0.135	0.383
Gapped		
Lambda	K	H
0.267	0.0410	0.140

Matrix: BLOSUM62

Gap Penalties: Existence: 11, Extension: 1
 Number of Sequences: 73055898
 Number of Hits to DB: 2969008310
 Number of extensions: 138028101
 Number of successful extensions: 489804
 Number of sequences better than 100: 2601
 Number of HSP's better than 100 without gapping: 0
 Number of HSP's gapped: 487533
 Number of HSP's successfully gapped: 2931
 Length of query: 1150
 Length of database: 26613539447
 Length adjustment: 162
 Effective length of query: 988
 Effective length of database: 14778483971
 Effective search space: 14601142163348
 Effective search space used: 14601142163348

T: 11
 A: 40
 X1: 16 (7.3 bits)
 X2: 38 (14.6 bits)
 X3: 64 (24.7 bits)
 S1: 41 (20.4 bits)
 S2: 85 (37.4 bits)
 ka-blk-alpha gapped: 1.9
 ka-blk-alpha ungapped: 0.7916
 ka-blk-alpha_v gapped: 42.6028
 ka-blk-alpha_v ungapped: 4.96466
 ka-blk-sigma gapped: 43.6362

Example 4: Generation of Group29 Mutants with Enhanced Specificity

[1206] Recently a method was described for the generation of Cas9 orthologs with enhanced specificity (Slaymaker et al. 2015). This strategy can be used to enhance the specificity of group29 orthologs. Primary residues for mutagenesis are all positive charges residues within the HEPN domain, since this is the only known structure in the absence of a crystal and we know that specificity mutants in RuvC worked in Cas9. An example of a conserved Arginine residue within HEPN domain is 8204.

[1207] Additional candidates are positive charged residues that are conserved between different orthologs. These can be used to generate group29 mutants with enhanced specificity.

Example 5: Design of RNA-Guide *E. coli* Screen and Validation

[1208] RNA-guide *E. coli* screen was designed to determine the likelihood of a guide working against RNA. Targeting guides were identified by tiling along a number of essential genes in *E. coli*. Specifically, a library of 55,700 targeting guides were created based on 45 essential genes

that are unique to their operon. From this, it could be assumed that any knockdown observed was likely due to cis-interference and not trans-interference, which is possible with the collateral effect. Taking the entire coding region, plus nucleotides into the 5' and 3' UTRs, a total of 54,600 unique sites in the *E. coli* genome were targeted. To establish baseline depletion, eleven hundred (1100) non-targeting guides ("dummy" guides) were also selected by generating randomized 30 nucleotide sequences that did not map to anywhere in the *E. coli* genome. The relative depletion of these dummy guides was used to determine significant levels of depletion among targeting guides.

[1209] As shown in FIG. 50, a protein plasmid and a guide plasmid were prepared, wherein the protein plasmid comprises an RNA-targeting protein, e.g., the large group 29 protein (Large Grp29 Protein), e.g., *Bergeyella zoohelcum* (BZ) or *Prevotella buccae* (PB) protein, or no Cas protein and a first antibiotic resistance gene, e.g., AmpR gene, and the guide plasmid comprises a targeting guide or a non-targeting guide and a second antibiotic resistance gene, e.g., CmR gene. From the library of targeting guides, a guide plasmid library was generated. A protein plasmid and a guide plasmid were co-electroporated at 50/50 ratio and

co-introduced into electrocompetent cell, e.g., MegaX DH10B T1R cells. The transformed cells were then allowed to grow for 1 hour at 37° C. and subsequently were incubated overnight on Chloramphenicol and Carbenicillin agar plates. The colonies grown on Chloramphenicol and Carbenicillin agar plates were collected 11 hours post transformation, and DNA was prepared for sequencing (maxiprep, NGS library prep, and loaded onto NextSeq). The experiment consisted of three bioreps of four plates per condition (control condition, i.e., no Cas protein, and experimental condition, i.e., with BZ or PB). Cross correlations of guide abundances are shown in FIGS. 51 (with BZ) and 52 (with PB).

[1210] Targeting guide depletions in alternating colors by gene for the 45 essential genes are shown in FIG. 53 for BZ screen and FIG. 54 for PB screen. Red horizontal line represents the safely depleted line, i.e. depletion above the line represents depletion above non-targeting guide. Non-targeting (dummy) guide depletions are shown in black. Depletion levels for targeting guide that are common to both PB and BZ (intersecting guides) are shown in FIG. 55. Normalizing depleted guide efficacy across each gene with both BZ and PB indicated that targeted positions are not spatially biased, as shown in FIGS. 56 (BZ) and 57 (PB).

[1211] Sequence logos were created showing the nucleotide frequency for the left protospacer flanking sequence (LPFS) and the right protospacer flanking sequence (RPFS). As shown in FIG. 58, the nucleotide G, A, or U is more preferred at position 10 (L1) with BZ protein. In other words, little or no cleavage was observed when C is at L1. Similarly, PB prefers A, G, or U nucleotide at L1. For the RPFS, cleavage was observed when either A, G, U or C was present at R2 (position 42) or R3 (position 43) with both BZ and PB. FIG. 59 shows the protospacer flanking sequence (PFS) wheel illustrating nucleotide frequencies for L1, R2, and R3 for BZ. FIG. 60 shows a protospacer flanking sequence (PFS) wheel illustrating nucleotide frequencies for L1, R2, and R3 for PB. Mean depletion for RPFS (R2 and R3) at fixed LPFS (L1) is shown in FIG. 61.

[1212] For validation screen, 4 types of sites in Ampicillin (162 sites)/Kanamycin (162 sites) were selected and a library of 322 sites were generated. Abundance of guides under experimental condition and control condition were determined and plotted as shown in FIGS. 62 (PB) and 63 (BZ). Experimental condition represents screening with Cas protein (either BZ or PB) with the library and target antibiotic resistance gene, and control condition indicates screening with only library and target antibiotic resistance gene (no Cas protein). Results of validation screens for BZ and PB are shown in FIGS. 62 and 63.

Example 6: Computational Sequence Analysis

[1213] From complete compiled Ensembl Release 27 genomes (Yates et al., 2016), CRISPR repeats were identified using PILER-CR (Edgar, 2007). Proteins within 10 kb of identified CRISPR arrays were clustered into loci, with loci rejected if more than one protein of size 700 amino acids or larger or if either Cas1 or Cas2 present. For candidate Class 2 effectors, only proteins in these remaining loci of size 900 aa to 1800 aa were selected. These candidate effectors were subjected to the BLASTP program (Camacho et al., 2008) searched against the NCBI non-redundant (NR) protein sequence database with an E-value cutoff of 1e-7. All discovered proteins were then grouped into putative families

via a nearest-neighbor grouping with the same E-value cutoff. Only putative families with at least ten candidate effectors and more than 50% of candidate effectors within 10 kb of CRISPR arrays were considered. HHpred (Remmert et al., 2012) and existing CRISPR locus classification rules were used to classify each family, leaving Cas13b as the only unclassified family. Within this family, truncated or suspected partially sequenced effectors were discarded, leaving 105 loci, and 81 with a non-redundant protein. Multiple sequence alignments on these 81 proteins (as well as the accessory Csx27 and Csx28 proteins) were performed using BLOSUM62 (Henikoff and Henikoff, 1992) to identify the HEPN domains and to sort the loci into phylogenetic trees. Vienna RNAfold (Lorenz et al., 2011) was used to predict secondary structure of each direct repeat, whose transcriptional orientation was chosen as identical to that of Cas13b in its locus. CRISPRTarget (Biswas et al., 2013) was used to search the spacers in each locus against NCBI phage and plasmid genomes. Weblogos were generated for all unique direct repeats and protospacer flanking sequences (Crooks et al., 2004). TMHMM Server v. 2.0 (Moller et al., 2001) was used to predict the transmembrane helices in Csx27 and Csx28.

Example 7: RFP-Tagged Protein Fluorescent Imaging

[1214] One Shot Stb13 Chemically Competent *E. coli* were transformed with RFP (negative control) or RFP fused to the N- or C-termini of Csx27 of *Bergeyella zoohelcum* or Csx28 of *Prevotella buccae*. Clones were cultured up in 5 mL of antibiotic LB overnight, then spun down at 5000 g and resuspended in PBS with 1% methanol-free formaldehyde. After 30 minutes fixation, cells were washed once with PBS and then diluted 1:2 in PBS. 5 uL of sample was pipetted onto a silane-coated slide, which was covered with a coverslip. Fluorescent imaging was performed in a 63x objective microscope with oil immersion.

Example 8: Bacterial RNA-Sequencing

[1215] RNA was isolated and prepared for sequencing using a previously described protocol (Shmakov et al., 2015; Heidrich et al., 2015). For native *Bergeyella zoohelcum* ATCC 43767 RNA sequencing, we repeated the experiment with a modified protocol, omitting TAP prior to 5' adapter ligation, to promote enrichment of processed transcripts originating from the CRISPR array.

[1216] The prepared cDNA libraries were sequenced on a MiSeq (Illumina). Reads from each sample were identified on the basis of their associated barcode and aligned to the appropriate RefSeq reference genome using BWA (Li and Durbin, 2009). Paired-end alignments were used to extract entire transcript sequences using Galaxy (<https://usegalaxy.org>), and these sequences were analyzed using Geneious 8.1.8.

Example 9: *E. coli* Essential Gene Experiment

[1217] The intersection of two *E. coli* DH10B strain essential gene studies (Gerdes et al., 2003; Baba et al., 2006) was taken, and further pared down to 45 genes by only selecting genes exclusive to their respective operons. Over all genes 54,600 spacers were designed to tile at single resolution across the coding region, as well as to extend 60 nt into the 5' UTR and 3'UTR. In addition, 1100 non-

targeting, randomly generated spacers with no precise match to the *E. coli* DH10B strain genome were added to the library as a non-targeting negative control. The library of spacers was cloned into a *Bergeyella zoohelcum* or *Prevotella buccae* direct repeat-spacer-direct repeat backbone containing a chloramphenicol resistance gene using Golden Gate Assembly (NEB) with 100 cycles, and then transformed over five 22.7 cm×22.7 cm chloramphenicol LB Agar plates. Libraries of transformants were scraped from plates and DNA extracted using the Macherey-Nagel Nucleobond Xtra Midiprep Kit (Macherey-Nagel). 50 ng of library plasmid and equimolar gene plasmid containing an ampicillin resistance gene (bzcas13b, bzcas13b & bzcsx27, pbcas13b, pbcas13b & pbcxs28, empty vector pBR322) were transformed into MegaX DH10B™ T1R Electrocomp™ Cells (ThermoFisher) according to manufacturer's protocol, with four separate 22.7 cm×22.7 cm carbenicillin-chloramphenicol LB Agar plates per bioreplicate, and three bioreplicates per condition (twelve transformations total per condition). Eleven hours post-transformation, libraries of transformants were scraped from plates and DNA extracted using the Macherey-Nagel Nucleobond Xtra Maxiprep Kit (Macherey-Nagel).

Example 10: *E. coli* Essential Gene Analysis

[1218] Prepared DNA libraries were sequenced on a Next-Seq (Illumina), with reads mapped to the input library of spacers. Spacer depletions were calculated as the read abundance of a spacer in the empty vector condition divided by read abundance in each gene plasmid condition. Mean depletions over three bioreplicates were calculated, with a maximum coefficient of variation of 0.2 and minimum spacer read abundance of $\frac{1}{3}N$, where $N=55,700$, as a quality-control filter. Weblogos of the strongly depleted (top 1% depleted) spacers were generated (Crooks et al., 2004), and from each identified PFS heatmaps of the ratio of safely depleted ($>5\sigma$ above mean depletion of non-targeting spacers) spacers to all spacers in the screen were generated. For spatial analysis via empirical cumulative distribution functions, safely depleted spacers were aggregated across the first or last 250 nt of genes.

[1219] For secondary structure analysis, the RNA accessibility model from Vienna RNAplfold (Bernhart et al., 2006) was utilized. RNAplfold calculates through a moving average of RNA folds the probability that a region u of RNA is unpaired given its cis sequence context in a four-parameter model, where W is the moving average window length in nucleotides, L is the maximum permissible pairing distance between nucleotides in the window, and u_{start} and u_{end} are the start and end of the region u , respectively. To apply this model to our data, spacers from the *E. coli* essential gene screen were separated into training/testing cohorts of five or more, each represented by a unique permissible PFS and gene and containing at least one spacer in the top 2% of depleted spacers from the screen (to enhance predictive signal). These cohorts were randomly divided into a training set (~80%) and a testing set (~20%). For optimizing a secondary structure-mediated model of efficient spacer design, the percent of cohorts for which the top spacer is accurately predicted or falls in the top 3 depleted spacers in a cohort were selected as objective functions top 1 and top 3 accuracy. The two objective functions on the training data set were optimized, first by fixing W and L while varying u_{start} and u_{end} , then by fixing u_{start} and u_{end} and varying W

and L (FIG. 66B). In the case of bzcas13b with bzcsx27, as well as that of bzcas13b alone, the optimized parameters were found to be approximately $W=240$, $L=180$, $u_{start}=16$, and $u_{end}=30$. The performance of this RNAplfold model relative to 10^6 Monte Carlo simulations performed on the testing data set were gauged, and empirical P-values of less than $1e-2$ for top 1 accuracy, and less than $1e-5$ for top 3 accuracy were found. Similar predictive power applied to pbcas13b with pbcxs28, as well as to pbcas13b alone.

Example 11: Kanamycin Validation Screen Experiment

[1220] A total of 160 kanamycin-targeting spacers was selected, 42 of which contain both PFS rules, 47 of which contain one rule, and 71 of which contain no rules, to which 162 non-targeting control spacers were added. The library of spacers was cloned into either a bzcas13b and *Bergeyella zoohelcum* direct repeat-spacer-direct repeat backbone or simply a *Bergeyella zoohelcum* direct repeat-spacer-direct repeat backbone containing a chloramphenicol resistance gene using Golden Gate Assembly (NEB) with 100 cycles, and then transformed over one 22.7 cm×22.7 cm carbenicillin LB Agar plate. The two cloned library plasmids were then re-transformed with over a 22.7 cm×22.7 cm chloramphenicol LB Agar plate or a 22.7 cm×22.7 cm kanamycin-chloramphenicol LB Agar plate. Libraries of transformants were scraped from plates and DNA extracted using the Qiagen Plasmid Plus Maxi Kit (Qiagen). 100 ng of library DNA and 100 ng of pMAX-GFP (Lonza), containing a kanamycin resistance gene were added to 50 μ L of chemically competent 10-beta cells (NEB) and transformed according to the manufacturer's protocol.

Example 12: Kanamycin Validation Screen Analysis

[1221] Prepared DNA libraries were sequenced on a Next-Seq (Illumina), with reads mapped to the input library of spacers. For normalizing the abundance of spacers of two separate clonings, the corrected experimental read abundance of a given spacer was calculated as the read abundance of that spacer in the bzcas13b plasmid (kanamycin-chloramphenicol transformation) multiplied by the ratio of the read abundance ratio of that spacer in the non-bzcas13b plasmid (chloramphenicol-only transformation) to the read abundance ratio of that spacer in the bzcas13b plasmid (chloramphenicol-only transformation).

Example 12: MS2 Phage Drop Plaque Assay

[1222] Individual spacers for MS2 interference were ordered as complementary oligonucleotides containing overhangs allowing for directional cloning in between two direct repeat sequences in vectors containing cas13b. 10 μ M of each complementary oligo were annealed in 10 \times PNK Buffer (NEB), supplemented with 10 mM ATP and 5 units of T4PNK (NEB). Oligos were incubated at 37 C for 30 min., followed by heating to 95 C for 5 min. and then annealed by cooling to 4 C. Annealed oligos were then diluted 1:100 and incubated with 25 ng of Eco31I digested cas13b vector in the presence of Rapid Ligation Buffer and T7 DNA ligase (Enzymatics). Individual plasmids were prepared using the QIAprep Spin Miniprep Kit (Qiagen), sequence confirmed and then transformed into C3000 (ATCC 15597) cells made competent using the Mix & Go E.

coli Transformation Kit (Zymo). In the case of experiments using csx27 or csx28, C3000 cells harboring these plasmids were made competent and then transformed with cas13b DR-spacer-DR plasmids. Following transformation, individual clones were picked and grown overnight in LB containing the appropriate antibiotics. The following morning, cultures were diluted 1:100 and grown to an OD₆₀₀ of 2.0, then mixed with 4 mL of antibiotic containing Top Agar (10 g/L tryptone, 5 g/L yeast extract, 10 g/L sodium chloride, 5 g/L agar) and poured on to LB-antibiotic base plates. 10 fold serial-dilutions of MS2 phage were made in LB and then spotted onto hardened top agar with a multi-channel pipette. For assessing interference levels in FIG. 4, samples were blinded using a key and the lowest dilution of phage at which plaque formation occurred was compared to pACYC by eye, where the lowest dilution of MS2 that formed plaques on pACYC was set to 1.

Example 13: DNA Targeting

[1223] A 34 nt target sequence consisting of a 30 nt protospacer and a permissive PFS a permissive PFS (5'-G, 3'-AAA) was cloned into pUC19 into two locations. For the transcribed target, the target sequence was cloned into the coding strand of the bla gene, in frame immediately after the start codon, with the G of the start codon serving as the 5' PFS. For the non-transcribed target the identical target sequence (protospacer and PFS) were cloned into the AatII site of pUC19, so that the protospacer appears on the non-transcribed strand with respect to the pBla and pLac promoters. To determine interference, 25 ng of the ampicillin resistant target plasmid and 25 ng of the chloramphenicol resistant bzcas13b or empty vector (pACYC) were added to 5 uL of NovaBlue GigaSingle cells (Novagen). The cells were incubated for 30 minutes on ice, heatshocked for 30 seconds at 42 C and incubated on ice for 2 minutes. Then, 95 uL of SOC was added to cells and they were incubated with shaking at 37 C for 90 minutes, before plating the entire outgrowth (100 uL) on plates containing both chloramphenicol and ampicillin.

METHODS REFERENCES

- [1224] 1. Baba T, Ara T, Hasegawa M, Takai Y, Okumura Y, Baba M, et al. Construction of *Escherichia coli* K-12 in-frame, single-gene knockout mutants: the Keio collection. *Molecular systems biology*. 2006; 2:2006 0008. doi: 10.1038/msb4100050. PubMed PMID: 16738554; PubMed Central PMCID: PMC1681482.
- [1225] 2. Bernhart S H, Hofacker I L, Stadler P F. Local RNA base pairing probabilities in large sequences. *Bioinformatics*. 2006; 22(5):614-5. doi: 10.1093/bioinformatics/btk014. PubMed PMID: 16368769.
- [1226] 3. Biswas A, Gagnon J N, Brouns S J, Fineran P C, Brown C M. CRISPRTarget: bioinformatic prediction and analysis of crRNA targets. *RNA biology*. 2013; 10(5): 817-27. doi: 10.4161/rna.24046. PubMed PMID: 23492433; PubMed Central PMCID: PMC3737339.
- [1227] 4. Camacho C, Coulouris G, Avagyan V, Ma N, Papadopoulos J, Bealer K, et al. BLAST+: architecture and applications. *BMC bioinformatics*. 2009; 10:421. doi: 10.1186/1471-2105-10-421. PubMed PMID: 20003500; PubMed Central PMCID: PMC2803857.
- [1228] 5. Crooks G E, Hon G, Chandonia J M, Brenner S E. WebLogo: a sequence logo generator. *Genome research*. 2004; 14(6):1188-90. doi: 10.1101/gr.849004. PubMed PMID: 15173120; PubMed Central PMCID: PMC419797.
- [1229] 6. Edgar R C. PILER-CR: fast and accurate identification of CRISPR repeats. *BMC bioinformatics*. 2007; 8:18. doi: 10.1186/1471-2105-8-18. PubMed PMID: 17239253; PubMed Central PMCID: PMC1790904.
- [1230] 7. Gerdes S Y, Scholle M D, Campbell J W, Balazsi G, Ravasz E, Daugherty M D, et al. Experimental determination and system level analysis of essential genes in *Escherichia coli* MG1655. *Journal of bacteriology*. 2003; 185(19):5673-84. PubMed PMID: 13129938; PubMed Central PMCID: PMC193955.
- [1231] 8. Heidrich N, Dugar G, Vogel J, Sharma C M. Investigating CRISPR RNA Biogenesis and Function Using RNA-seq. *Methods in molecular biology*. 2015; 1311:1-21. doi: 10.1007/978-1-4939-2687-9_1. PubMed PMID: 25981463.
- [1232] 9. Henikoff S, Henikoff J G. Amino acid substitution matrices from protein blocks. *Proceedings of the National Academy of Sciences of the United States of America*. 1992; 89(22):10915-9. PubMed PMID: 1438297; PubMed Central PMCID: PMC50453.
- [1233] 10. Li H, Durbin R. Fast and accurate short read alignment with Burrows-Wheeler transform. *Bioinformatics*. 2009; 25(14):1754-60. doi: 10.1093/bioinformatics/btp324. PubMed PMID: 19451168; PubMed Central PMCID: PMC2705234.
- [1234] 11. Lorenz R, Bernhart S H, Honer Zu Siederdisen C, Tafer H, Flamm C, Stadler P F, et al. Vienna RNA Package 2.0. *Algorithms for molecular biology: AMB*. 2011; 6:26. doi: 10.1186/1748-7188-6-26. PubMed PMID: 22115189; PubMed Central PMCID: PMC3319429.
- [1235] 12. Moller S, Croning M D, Apweiler R. Evaluation of methods for the prediction of membrane spanning regions. *Bioinformatics*. 2001; 17(7):646-53. PubMed PMID: 11448883.
- [1236] 13. Remmert M, Biegert A, Hauser A, Soding J. HHblits: lightning-fast iterative protein sequence searching by HMM-HMM alignment. *Nature methods*. 2012; 9(2):173-5. doi: 10.1038/nmeth.1818. PubMed PMID: 22198341.
- [1237] 14. Shmakov S, Abudayyeh O O, Makarova K S, Wolf Y I, Gootenberg J S, Semenova E, et al. Discovery and Functional Characterization of Diverse Class 2 CRISPR-Cas Systems. *Molecular cell*. 2015; 60(3):385-97. doi: 10.1016/j.molcel.2015.10.008. PubMed PMID: 26593719; PubMed Central PMCID: PMC4660269.
- [1238] 15. Yates A, Akanni W, Amode M R, Barrell D, Billis K, Carvalho-Silva D, et al. Ensembl 2016. *Nucleic acids research*. 2016; 44(D1):D710-6. doi: 10.1093/nar/gkv1157. PubMed PMID: 26687719; PubMed Central PMCID: PMC4702834.
- [1239] While preferred embodiments of the present invention have been shown and described herein, it will be obvious to those skilled in the art that such embodiments are provided by way of example only. Numerous variations, changes, and substitutions will now occur to those skilled in the art without departing from the invention. It should be understood that various alternatives to the embodiments of the invention described herein may be employed in practicing the invention. It is intended that the following claims define the scope of the invention and that methods and structures within the scope of these claims and their equivalents be covered thereby.

SEQUENCE LISTING

The patent application contains a lengthy "Sequence Listing" section. A copy of the "Sequence Listing" is available in electronic form from the USPTO web site (<http://seqdata.uspto.gov/?pageRequest=docDetail&DocID=US20170211142A1>). An electronic copy of the "Sequence Listing" will also be available from the USPTO upon request and payment of the fee set forth in 37 CFR 1.19(b)(3).

What is claimed:

1. A method for identifying novel nucleic acid modifying effectors, comprising:

identifying a putative nucleic acid modifying loci from a set of nucleic acid sequences encoding the putative nucleic acid modifying enzyme loci that are within a defined distance from a conserved genomic element of the loci, that comprise at least one protein above a defined size limit, or both;

grouping the identified putative nucleic acid modifying loci into subsets comprising homologous proteins;

identifying a final set of candidate nucleic acid modifying loci by selecting nucleic acid modifying loci from one or more subsets based on one or more of the following; subsets comprising loci with putative effector proteins with low domain homology matches to known protein domains relative to loci in other subsets, subsets comprising putative proteins with minimal distances to the conserved genomic element relative to loci in other subsets, subsets with loci comprising large effector proteins having a same orientations as putative adjacent accessory proteins relative to large effector proteins in other subsets, subset comprising putative effector proteins with lower existing nucleic acid modifying classifications relative to other loci, subsets comprising loci with a lower proximity to known nucleic acid modifying loci relative to other subsets, and total number of candidate loci in each subset.

2. The method of claim **1**, wherein the set of nucleic acid sequences is obtained from a genomic or metagenomic database.

3. The method of claim **4**, wherein the genomic or metagenomic database comprises prokaryotic genomic or metagenomic sequences.

4. The method of claim **1**, wherein the defined distance from the conserved genomic element is between 1 kb and 25 kb.

5. The method of claim **1**, wherein said conserved genetic element comprises a repetitive element.

6. The method of claim **5**, wherein the repetitive element is a CRISPR array.

7. The method of claim **6**, wherein the defined distance is within 10 kb of the CRISPR array.

8. The method of any one or claims **1** to **7**, wherein the defined size limit is greater than 200 amino acids.

9. The method of claim **8**, wherein the defined size limit is greater than 700 amino acids.

10. The method of claim **1**, wherein said conserved genomic elements are identifying using a repeat or pattern finding analysis of the set of nucleic acids.

11. The method of claim **10**, wherein said repeat or pattern finding method comprises PILER-CR.

12. The method of claim **1**, wherein said grouping step is based, at least in part, on results of a domain homology search or an HHpred protein domain homology search.

13. The method of claim **1**, wherein the defined threshold is a BLAST nearest-neighbor cut-off value of 0 to 1e-7.

14. The method of claim **1**, further comprising a filtering step that includes only loci with putative proteins between 900 and 1800 amino acids.

15. The method of any one of claims **1** to **14**, further comprising experimentally validating the nucleic acid modifying function of the candidate nucleic acid modifying effectors by generating a set of nucleic acid constructs encoding the nucleic acid modifying effectors and performing one or more biochemical validation assays.

16. The method of claim **20**, wherein said validation comprises PAM validation in bacterial colonies, in vitro cleavage assays, the Surveyor method, experiments in mammalian cells, PFS validation, or a combination thereof.

17. The method of any one of claims **1** to **16**, further comprising preparing a non-naturally occurring or engineered composition comprising one or more proteins from the identified nucleic acid modifying loci.

18. The method of claim **17**, wherein the identified loci comprise a Class 2 CRISPR effector.

19. The method of claim **17**, wherein the identified loci lack Cas1 or Cas2.

20. The method of claim **17**, wherein the identified loci comprise a single effector.

21. The method of claim **17**, wherein loci comprise a single large effector protein that is greater than 900, or greater than 1100 amino acids in length.

22. The method of claim **21**, wherein the single large effector comprises at least one HEPN domain.

23. The method of claim **22**, wherein the at least one HEPN domain is near a N- or C-terminus of the effector protein.

24. The method of claim **23**, wherein the at least one HEPN domain is located in an interior position of the effector protein.

25. The method of claim **22**, wherein the single large effector protein comprises a HEPN domain at the N- and C-terminus and two HEPN domains internal to the protein.

26. The method of claim **20**, wherein the loci further comprises one or two small putative accessory proteins within 2 kb to 10 kb of the CRISPR array.

27. The method of claim **26**, wherein the small accessory protein comprise multiple predicted transmembrane domains.

28. The method of claim 27, wherein the small accessory protein comprises four predicted transmembrane domains.

29. The method of claim 26, wherein small accessory protein comprises at least one HEPN domain.

30. The method of claim 26, wherein the small accessory protein comprises at least one HEPN domain and a at least one transmembrane domain.

31. The method of claim 17, wherein the loci comprise no additional proteins out to 25 kb from the CRISPR array.

32. The method of claim 17, wherein the CRISPR array comprises direct repeat sequences comprising about 36 nucleotides in length.

33. The method of claim 32, wherein the direct repeat comprises a GTTG/GUUG at the 5' end that is reverse complementary to a CAAC at the 3' end.

34. The method of claim 17, wherein the CRISPR array comprises spacer sequences comprising about 30 nucleotides in length.

35. The method of claim 17, wherein said loci lack a small accessory protein.

* * * * *