



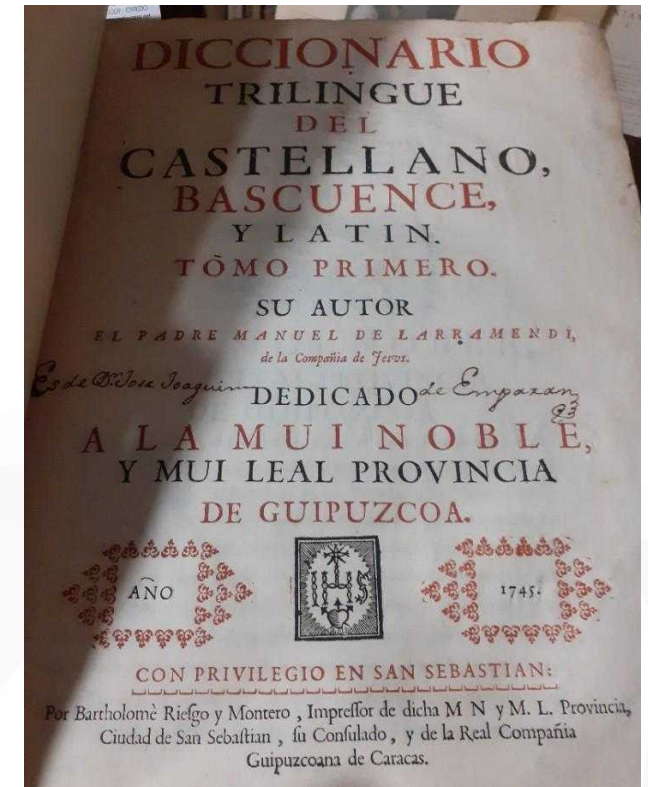
A workflow for historical dictionary digitization: Larramendi's Trilingual Dictionary



David Lindemann david.lindemann@ehu.eus
Mikel Alonso mikelalon@gmail.com

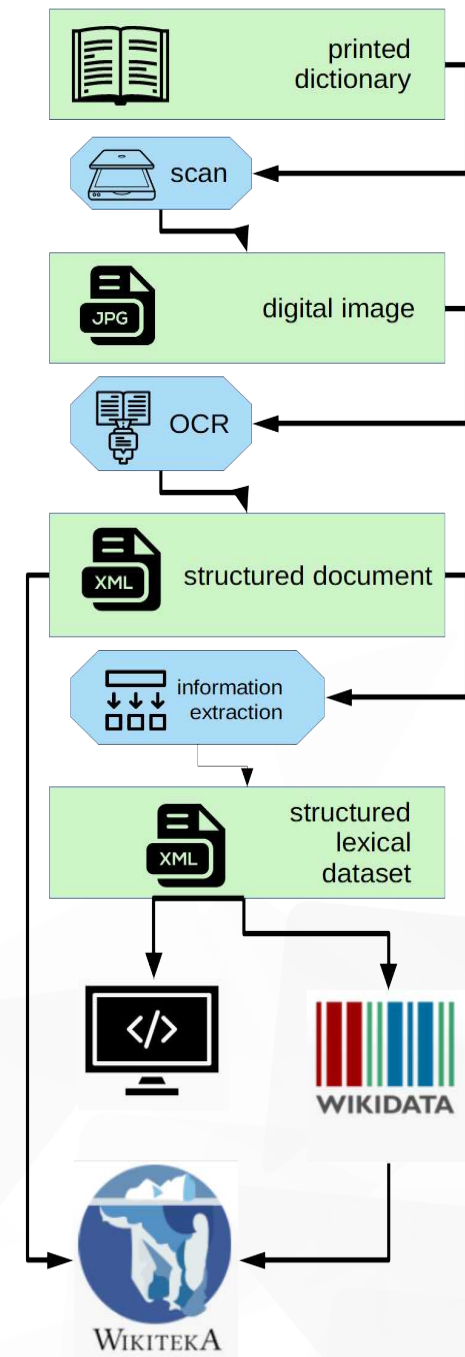
Introduction

- ▼ Larramendi's *Diccionario Trilingüe del Castellano, Bascuence y Latín* (LAR, 1745)
 - ▼ In Basque, *Hiztegi Hirukoitza* (the 'triple dictionary')
 - ▼ First published general Basque dictionary.
 - ▼ Brings significant shift in Basque Lexicography periodization (Urgell 2002).
 - ▼ During a century and a half main reference work for Basque
 - ▼ Subject to in-depth philological analysis (Urgell 1998a; 1998b, among others).
 - ▼ Until today only available on paper, and as collection of scanned images.
- ▼ UEU/Wikimedia Basque Country short-term research grant (2020)
 - ▼ Workflow design using Wikimedia platforms and freely available tools



Goals

- ▼ Propose a workflow for historical dictionary digitization using Wikimedia platforms and free tools
- ▼ Predict workload...
 - ▼ for training set compilation (OCR, information extraction)
 - ▼ for result validation (covering the whole dictionary)
- ▼ Publish LAR contents on Wikisource and Wikidata platforms



Optical Character Recognition (OCR)

▼ Kraken

- ▼ machine-learning based
- ▼ freely available
- ▼ deals with font style recognition
- ▼ result export in ALTO XML format

▼ Steps

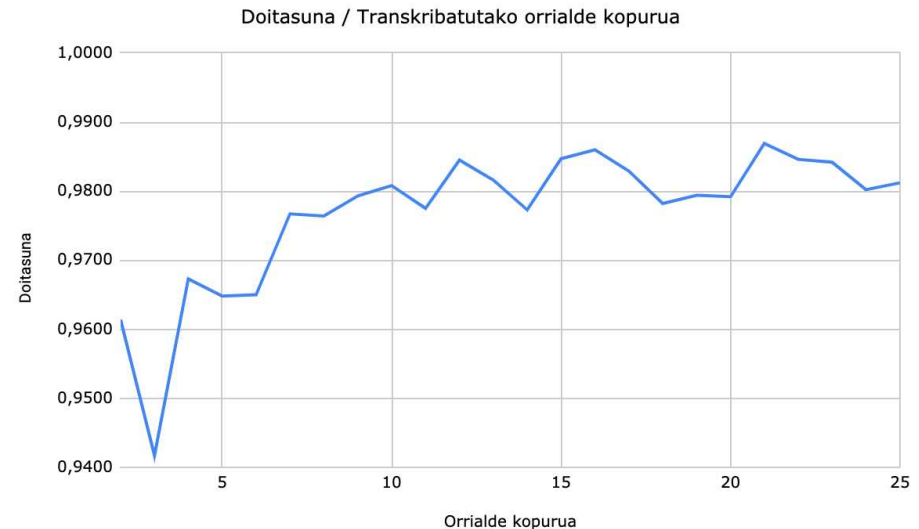
- ▼ image pre-processing
- ▼ transcription of a training set
- ▼ application of the trained model to the whole content

V E. 371
Lat. Vertibilitas.
Vertible , *aldacoya* ; *giracoya*. Lat. Vertibilis.
Vertical , *bugaindarra*. Lat. Verticalis.
Vertice , *bugaina*. Lat. Vertex , icis.
Verticidad , veafe *vertibilidad*.
Vertigo , vertiginoso , veafe *vaguido*.
Vespero , *illunabarreco izarra*. Lat. Vesperus.
Vespertino , *arratsaldeco* , *arrastezguicoa*. Lat. Vespertinus.
Vesquir , antiquado , lo mismo que *vivir*.
Veste , lo mismo que *vestido*.
Vestido , *soñecoa* , *jazcaya* , *jaunzcaya* , *aldagarria* , *filda* , *abillamendua*. Lat. Vestis , vestitus.
Vestidura , lo mismo. Lat. Indumentum.
Vestigio , *aztarnà* , *sená* , *hatzá*. Lat. Vestigium.
Vestiglo , monstruo formidable , *mamuza*. Lat. Spectrum horridum.
Vestimenta , vestimento , veafe *vestido*.
Vestir , *janci* , *jaunci* . *bestitu*. Lat. Vestire , inducere.

VE.
Lat. Vertibilitas.
371
Vertible , @aldacoya , @giracoya Lat. Vertibilis.
Vertical , @bugaindarra Lat. Verticalis.
Vertice , @bugaina Lat. Vertex , icis.
Verticidad , veafe @vertibilidad.
Vertigo , vertiginoso , veafe @vaguido.
Vespero , @illunabarreco @izarra Lat. Vesperus.
Vespertino , @arratsaldeco , @arrastezguicoa Lat. Vespertinus.
Vesquir , antiquado , lo mismo que vivir.
Veste , lo mismo que @vestido.
Vestido , @soñecoa , @jazcaya , @jaunzcaya , @aldagarria , @filda , @abillamendua Lat. Vestis , vestitus.
Vestidura , lo mismo Lat. Indumentum.
Vestigio , @aztarnà , @sená , @hatzá Lat. Vestigium.

Optical Character Recognition (OCR)

- ▼ Model training: Inclusion of more transcribed dictionary columns into the training set, until desired precision rate reached
- ▼ After 25 pages (i.e. 50 columns), we estimated that precision would not significantly grow
- ▼ Infrequent characters not properly recognized: Inclusion of additional 14 columns (where these appear)
- ▼ With a semi-automatic transcription of less than 4% of the content, OCRed version of the whole dictionary with 98,5% precision
- ▼ Significantly better than TXT versions available before



ALTO to Wikitext

ALTO XML format produced by Kraken:

- Represents characters, font style, and text box position on page.
- The latter is used for a rule-based approach of segmentation of dictionary text into entries (indented lines). We allow only indented lines starting with capital letter of the corresponding alphabet section.

ALTO to Wikitext:

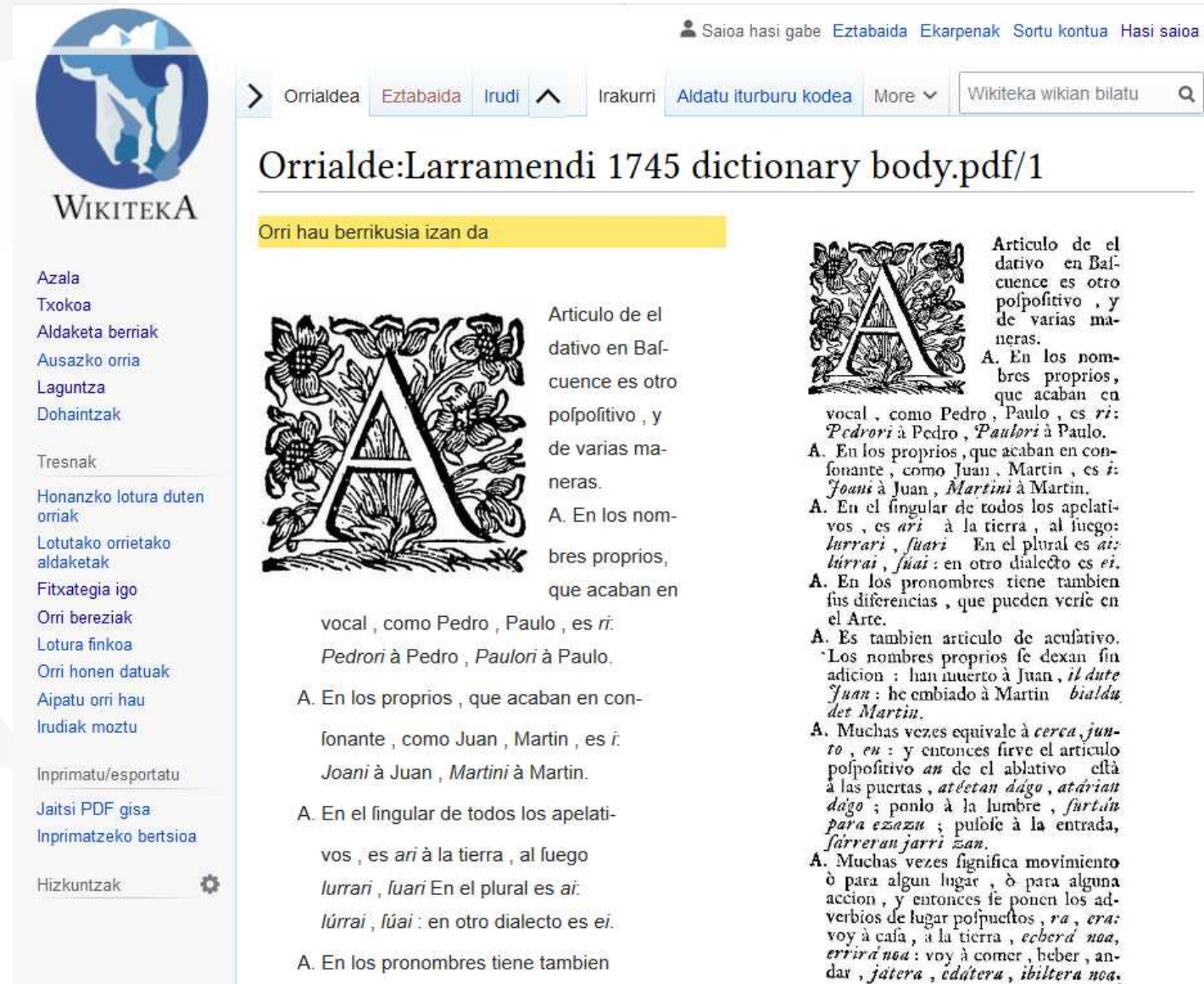
- Preserves font style
- Line indent represented by colon(s)
- Anchor template for headwords (in draft: first word of negative indented line starting with correct letter)


```
:::::A B.
::bandono, ''utziera'', ''lajaera'', ''lagaera''.
:::Lat. Derelictio. Envn total abandono,
::''utziera'' ''gucizcoán''.
:{{sarrera|abanicarse}}Abanicarse, ''aizatu'', ''aizatzea''.
Lat. Fla-
::bello ventum facere , movere.
:{{sarrera|abanico}}Abanico, ''aitzequiña'', ''aizeguillea'',
''aize-''
::''emallea''. Lat. Flabellum.
:{{sarrera|abanillo}}Abanillo, diminutivo de abano.
:{{sarrera|abanino}}Abanino , moda de que vñaban las Damas
::de Palacio , y era vn pedazo de Gaffa
::blanca , atraveffada en el escote de el
::Jubon , ''abaninoa''. Lat. Lineus colli
::amictus.
:{{sarrera|abno}}Abno, es vn Abanico grande , que col-
::gaba de el techo, y meneado con cuer-
::da, haze ayre, yahuyenta las moças.
::''aizequin'' ''aundia''. Lat. Flabrum , i.
:{{sarrera|abaratar}}Abaratar , ''merquetú'' , ''merqué''
''ifini'',
:::: eee re ve
:{{sarrera|abaratado}}Abaratado demafiadamente,
''merquetue-''
::''quia'' , ''merquequi'' ''ifinia''.
```


A B.
Abandono, *utziera*, *lajaera*, *lagaera*.
Lat. Derelictio. En vn total abandono,
utziera gucizcoán.
Abanicarse, *aizatu*, *aizatzea*. Lat. Fla-
bello ventum facere, movere.
Abanico, *aitzequiña*, *aizeguillea*, *aize-
emallea*. Lat. Flabellum.
Abanillo, diminutivo de abano.
Abanino, moda de que vñaban las Damas
de Palacio, y era vn pedazo de Gaffa
blanca, atraveffada en el escote de el
Jubon, *abaninoa*. Lat. Lineus colli
amictus.
Abano, es vn Abanico grande, que col-
gaba de el techo, y meneado con cuer-
da, haze ayre, yahuyenta las moças,
aizequin aundia. Lat. Flabrum, i.
Abaratar, *merquetú*, *merqué ifini*.
Lat. Minoris vendere, venire. Vide
barato.
Abaratado demafiadamente, *merquetue-
guia*, *merquequi ifinia*.
Abarca, calzado de cuero. Es del Baf-
cunce: dixofe así, por la semejanza,
que tiene con la barca, *au barca*, y

Wikiteka (Basque Wikisource platform)

- Collaborative validation of OCR results
- Active community
- Two-step validation



Wikiteka logo:  WIKITEKA

Navigation: > Orrialdea **Eztabaida** Irudi  Irakurri Aldatu iturburu kodea More

Orrialde:Larramendi 1745 dictionary body.pdf/1

Orrri hau berrikusia izan da

A Artículo de el dativo en Baf-cuence es otro pospositivo , y de varias maneras.

A. En los nombres propios, que acaban en vocal , como Pedro , Paulo , es *ri*: *Pedrori* à Pedro , *Paulori* à Paulo.

A. En los propios , que acaban en consonante , como Juan , Martin , es *i*: *Joani* à Juan , *Martini* à Martin.

A. En el singular de todos los apelativos , es *ari* à la tierra , al fuego *lurrari* , *luari* En el plural es *ai*: *lúrrai* , *lúai* : en otro dialecto es *ei*.

A. En los pronombres tiene tambien sus diferencias , que pueden verçie en el Arte.

A. Es tambien articulo de acusativo. Los nombres propios se dexan sin adición : han muerto à Juan , *il dute Juan* : he embiado à Martin *bialdu det Martin*.

A. Muchas vezes equivale à *cerca* , *junto* , *en* : y entonces sirve el articulo pospositivo *an* de el ablativo está à las puertas , *atetan dago* , *atavian dago* ; pono à la lumbre , *furtu para ezazu* ; pulso à la entrada , *sárreran jarri zan*.

A. Muchas vezes significa movimiento ò para algun lugar , ò para alguna accion , y entonces se ponen los adverbios de lugar pospuestos , *ra* , *era*: voy à casa , à la tierra , *esberd noa* , *errid noa* : voy à comer , heber , andar , *jatera* , *edatera* , *ibiltera noa*.

Information extraction

<lemma_Spanish>

<definition_Spanish>

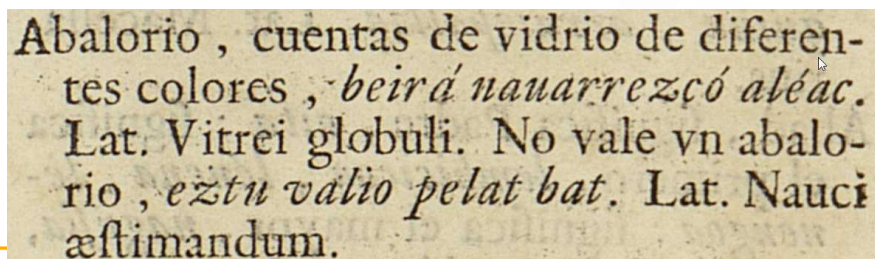
Monopolio, vn convenio prohibido de mercaderes, para vender à tanto los generos, *ambàsala*. Lat. Monopolium.

<TE_Basque>

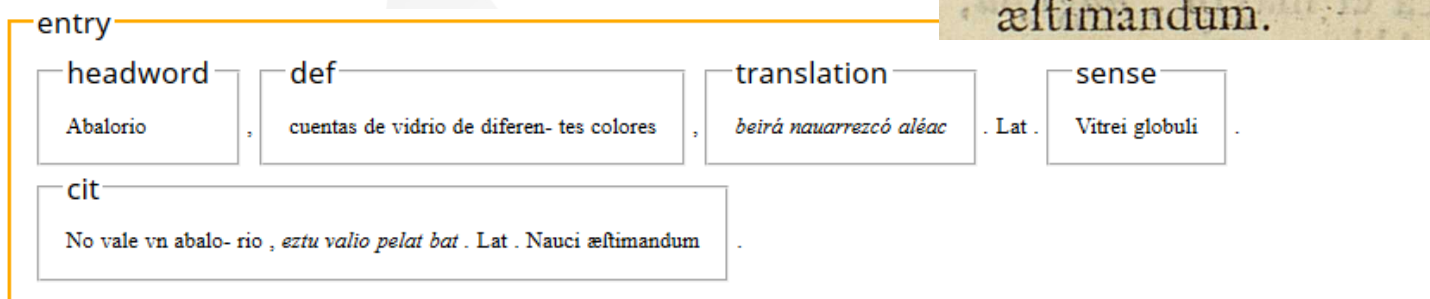
<TE_Latin>

Information extraction

- ▼ ELEXIFIER is a ML-based tool for information extraction from dictionary plain text. As Elexis observer institution, we take part in Elexifier beta evaluation, as early adopters.
 - ▼ TEI-Lex0, with (at the time of preparing the dataset) a still limited tagset
 - ▼ Training set: Manual annotation of entry structure on 10 pages.
 - ▼ Annotated microstructure elements: headword, Basque equivalent, latin equivalent, cross-reference, definition (Spanish), "cit" element for other content



Abalorio , cuentas de vidrio de diferentes colores , *beirá nauarrezcó aléac.* Lat. Vitrei globuli. No vale vn abalorio , *eztu valio pelat bat.* Lat. Nauci æstimandum.



Information extraction

▼ ELEXIFIER results:

- ▼ Spanish headwords identified with very high precision.
- ▼ Basque equivalents also identified with high precision, but many false positives (other items in italics font style also taken as Basque equivalents, such as cross-references)
- ▼ Latin equivalents identified with low precision.

▼ Strategy for improving results:

- ▼ Use a fully validated transcription (once wikisource community has finished).
- ▼ Use a full TEI-Lex0 tagset
- ▼ Manually annotate a larger training set.

DA / LAR (+ OEH)

21693	gafar	Gafas de ballesta, <i>tiracacoac</i> . Lat. Unci, orum.
		Gafas, <i>antojos</i> , vease.
		Gafas, en los trucos, <i>cacola</i> . Lat. Uncinata tabella.
21694	gafas	Gafar, <i>arripata</i> . Lat. Uncis arripere.
21695	gafedad	Gafedad, <i>lepra</i> , vease.
		Gafedad de pies, y manos, <i>oñetaco</i> , <i>ta esquetaco elbarria</i> , <i>macurdea</i> . Lat. Nervorum in manibus, pedibulque contratio.
21696	gafete	Gafez, lo mismo que <i>gafedad</i> .
21697	gafez	Gafo, <i>leproso</i> , vease.
21698	gafo	Gafo de pies, y manos, <i>elbarria</i> , <i>macurra</i> . Lat. Manibus, ac pedibus captus.
21699	gagates	Gagates, lo mismo que <i>azabache</i> , vease.
21700	gage	Gage, prenda, señal, viene de el Balcuence <i>gaisca</i> , que significa lo mismo, y se compone de <i>gai</i> , <i>gaia</i> , materia, señal, instrumento para algo, y <i>chea</i> , pequeño. Lat. Oppigneratio.
21701	gages	Gage, gages, <i>gaisca</i> , <i>farjac</i> . Lat. Merces, stipendium.
21702	gago	Gago, lo mismo que <i>gangofo</i> , vease.
21703	gaita	Gaita, instrumento musico de cuero hinchado, y flauta, <i>chirolarria</i> . Lat. Tibia utricularia.
21704	gaiteria	Gaita de rueda, teclas, y cuerdas, que comunmente tocan los ciegos, <i>boltrula</i> . Lat. Pfalterium fidicinum.
21705	gaitero	Gaita, de tamboritero, vease <i>flauta</i> .
21706	gajo	Gaita, melecina, <i>ciriscaria</i> , <i>ayuda</i> , <i>ajuda muscuria</i> . Lat. Clysterium, ij.
21707	gajoso	Gaiteria, traje alegre, de diversos colores, <i>chorajancia</i> . Lat. Fatilis ornatus.
21708	gala	Gaitero, <i>chirolarria</i> . Lat. Tibicen.
		Gaitero, vestido de gaiteria, <i>chorajanciduna</i> . Lat. Hilari ornatu gestiens.
		Gajo, rama cortada, es de el Balcuence, <i>gajoa</i> , o <i>gaisoa</i> , que significa pobre en sentido de digno de compasion, y tambien enfermo, y la rama cortada, como enferma, y se seca, y muere, y pobre rama, que la cortan sin que, ni para que, y por esto <i>gajoa</i> , <i>gaisoa</i> , <i>Adar ebajua</i> , <i>tantas epaquia</i> . Lat. Ramus dissectus.
		Gajo de vbas, &c. <i>cordoca</i> , <i>chordoquea</i> , <i>lacaña</i> . Lat. Scapus.
		Gajolo, <i>cordocaduna</i> , <i>chordocat sua</i> , <i>lacainduna</i> . Lat. Scapis plenus.
		Gala, es voz Balcogada, como tambien los

Gafas de ballesta, *tiracacoac*. Lat. Unci, orum.
 Gafas, *antojos*, vease.
 Gafas, en los trucos, *cacola*. Lat. Uncinata tabella.

Antojos, anteojos, *beguordeac*, *beacayac*, *miserac*. Lat. Conspicilia, orum.

BEGI-ORDE (L ap. **A**; **SP**, **Lar**, **Hb**, **H**).
 "Lunettes, biserak (**Ax**)" **SP**. "Antojos, anteojos, **begiordeak**, **beakaiak**, **miserak**" **Lar**. "Begietakoak ou begiordeak (**SP**), lunettes" **Dv**. "Begi ordeak, lunettes; id. begi-orde thailuak (**Ax**)" **H**.

Mapping historical to modern standard spelling (ES, EU)

<i>LAR</i>	<i>DA</i>	<i>matching normalized lemma-sign</i>
Obsession	obsesión	obsesion
Hueffo	huesso	hueso
Occiffion	occisión	occision
Atràs	atras	atras

▼ Spanish

- ▼ Remove diacritics, historical 'long s' and double 's' > 's'

▼ Basque

- ▼ Sequence of regex operations, discussed in Alonso (2021)

▼ examples:

- ▼ *begiordeac* > *begiordeak* > *begiorde*
- ▼ *astigarra* > *astigar*
- ▼ *astotua* > *astotu*

Index of LAR: ES ⇨ EU

abarcar

- Autoridades: abarcar
- Larramendi (Arau): abarcar https://eu.wikisource.org/wiki/Hiztegi_Hirukoitza/A#abarcar
⇒ itzulpen hautagaiak: abarcá, quien mucho abarca, poco aprieta
- Larramendi (Arau): abarcar2 https://eu.wikisource.org/wiki/Hiztegi_Hirukoitza/A#abarcar2
⇒ itzulpen hautagaiak: abarcatu, gauz ascó eguitea, artzea, asco abarcatzea, da guichi estutzea, sco eguin nai, guichi eguin bay
- Larramendi (IA): abarcar
- Larramendi (IA): abarcar2
⇒ itzulpen hautagaiak: abarcatu, gauz ascó eguitea, artzea

- ▼ DA (Diccionario de Autoridades) lemmata:
 - ▼ Out of 41,968, 27,825 found in LAR; 4,875 LAR lemmata not found in DA
- ▼ LAR lemmata and Basque equivalent candidates:
 - ▼ Rule-based extraction:
 - ▼ 30,045 unique equivalent candidates for 32,044 Spanish headwords
 - ▼ Machine-learning-based extraction:
 - ▼ 27,125 unique equivalent candidates for 29,932 Spanish headwords

Index of LAR: ES ⇨ EU

abarcar

- Autoridades: abarcar
- Larramendi (Arau): abarcar https://eu.wikisource.org/wiki/Hiztegi_Hirukoitza/A#abarcar
⇒ itzulpen hautagaiak: abarcá, quien mucho abarca, poco aprieta
- Larramendi (Arau): abarcar2 https://eu.wikisource.org/wiki/Hiztegi_Hirukoitza/A#abarcar2
⇒ itzulpen hautagaiak: abarcatu, gauz ascó eguitea, artzea, asco abarcatzea, da guich
- Larramendi (IA): abarcar
- Larramendi (IA): abarcar2
⇒ itzulpen hautagaiak: abarcatu, gauz ascó eguitea, artzea

abarcas

- Larramendi (Arau): abarcas https://eu.wikisource.org/wiki/Hiztegi_Hirukoitza/A#abarcas
⇒ itzulpen hautagaiak: abarcac
- Larramendi (IA): abarcas
⇒ itzulpen hautagaiak: abarcac

abarcon

- Autoridades: abarcon

abarquillado

- Autoridades: abarquillado

abarquillar

- Autoridades: abarquillar

→ Abarcas , abarcac. Lat. Pero, onis.

Abarcar , viene de el Balfuence *abarcá*,
que le calza , dando muchas bueltas
à la pierna , como abrazandola con
vna cuerda muy larga: y de aqui, *quien
mucho abarca , poco aprieta*.

Abarcar, *abarcatu , gauz ascó eguitea*,
artzea. Lat. Complector, eris. Quien
mucho abarca , poco aprieta : *alco
abarcatzea , da guichi estutzea*.
lco eguin nai , guichi eguin bay.
Lat. Dum complecti vis multa pauca.
tenes.

- ▼ Rule-based headword extraction:
Headwords connected to wikisource
navigation anchors

Index of LAR: ES ⇨ EU

→ Abarcas , abarcac. Lat. Pero, onis.

Abarcar , viene de el Balcuence abarcá,
que le calza , dando muchas bueltas
à la pierna , como abrazandola con
vna cuerda muy larga: y de aqui, *quien
mucho abarca , poco aprieta.*

Abarcar, abarcatu , gauz alcó eguitea,
artzea. Lat. Complector, eris. Quien
mucho abarca , poco aprieta : *alco
abarcatzea , da guichi estutzea.*
lco eguin nai , quichi eguin bay.
Lat. Dum complecti vis multa pauca.
tenes.

mucho abarca , poco aprieta.
Abarcar, abarcatu , gauz alcó eguitea,
artzea. Lat. Complector, eris. Quien
mucho abarca , poco aprieta : *asco
abarcatzea , da guichi estutzea.*
Asco eguin nai , guichi eguin bay.
Lat. Dum complecti vis multa pauca.
tenes.
Abarracarse , meterse dentro de las bar-
racas , *echoletan , chaoletan sartu:*

- ▼ If a transcription error is found, it can be directly corrected (logged into wikisource)

Index of LAR: EU ⇄ ES

abarizketa Lar: abarizqueta (arauz) Lar: [abarizqueta] (IAz)

⇒ carrascal https://eu.wikisource.org/wiki/Hiztegi_Hirukoitza/C#carrascal

abarka ezpartuzkoa Lar: abarca ezpartuzcoa (arauz) Lar: [abarca ezpartuzcoa] (IAz)

⇒ alborga https://eu.wikisource.org/wiki/Hiztegi_Hirukoitza/A#alborga

abarkak Lar: abarcac (arauz) Lar: [abarcac] (IAz)

🔍 Sarasola: abarka 🔍 Wikidata: abarka 🔍 OEH: abarka <https://www.euskaltzaindia.eus/index.php?option=>

⇒ abarcas https://eu.wikisource.org/wiki/Hiztegi_Hirukoitza/A#abarcas

abarkatu Lar: abarcatu (arauz) Lar: [abarcatu] (IAz)

⇒ abarcar2 https://eu.wikisource.org/wiki/Hiztegi_Hirukoitza/A#abarcar2

→ Carralcal , abarizqueta , tartacadia ,
zarbaztaga. Lat. Illicetum.

Carralco , carralca , es voz Balcogada

carralcó , *garralcó* significa mucha
llama , y el carralco es oportuno par
ello. *Abarrá tartacá* , *zarbazta*.

Lat. Ilex , cis , prinus.

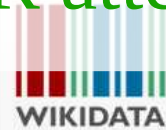
Carralpada , bebida compuesta de vino,

agua , miel , *eztiedaria*. Lat. Vinum
melle

- ▼ Reversed dictionary Basque to Spanish:
 - ▼ Spanish equivalents linked to wikisource anchors
- ▼ Mapping of Basque items to modern Basque dictionaries:
 - ▼ ML extracted Basque items:
 - ▼ 38,300 candidates, 11,551 found in modern dictionaries
 - ▼ Rule-based approach:
 - ▼ 60,193 candidates, 15,152 found in modern dictionaries

LAR attestations in Wikidata

- ▶ "attested in" property, "attested as" and "described at URL" qualifiers
- ▶ Wikidata Basque lexemes: linked to senses, forms, and other lexical resources
- ▶ To start with:
Both methods agree, and Sarasola cites as LAR item
(1,179 lemmata)
- ▶ POS mapping
- ▶ Problems with letter H, and hyphenization
- ▶ To be continued...



Azala
Komunitatearen ataria
Project chat
Elementu berria sortu
Aldaketa berriak
Ausazko orria
Query Service
Gertukoak
Laguntza
Dohaintzak

Lexicographical data
Lexema berria sortu
Aldaketa berriak
Random Lexeme

Tresnak
Honanzko lotura duten orriak
Lotutako orrietako aldaketak
Orri bereziak
Lotura finkoa
Orri honen datuak
Aipatu orri hau
Concept URI

Inprimatu/esportatu
Jaitsi PDF gisa
Inprimatzeko bertsioa

Lexeme Eztabaida

(L51983) **abarka** aldatu

eu

Language euskara
Lexical category substantibo

Adierazpenak

ahoskera audio aldatu

LL-Q8752 (eus)-Xabier Cañas-abarka.wav
1,1s; 97 KB
⚠
→ 0 erreferentzia
+ Erreferentzia gehitu
+ balioa gehitu

Ahotsak Lexeme *ingeles* aldatu

abarka
→ 0 erreferentzia
+ Erreferentzia gehitu
+ balioa gehitu

Elhuyar Hiztegien identifikatzailea aldatu

Heu000675_0
→ 0 erreferentzia
+ Erreferentzia gehitu
+ balioa gehitu

non jasoa aldatu

Larramendiren Hiztegi Hirukoitza
honela jasoa abarcac
dagokion URLa https://eu.wikisource.org/wiki/Hiztegi_Hirukoitza/A#abarcas
→ 0 erreferentzia
+ Erreferentzia gehitu
+ balioa gehitu

Links

- ▼ LAR at Wikidata: <http://www.wikidata.org/entity/Q65216433>
- ▼ LAR at Wikisource: https://eu.wikisource.org/wiki/Hiztegi_Hirukoitza
- ▼ Transcription guidelines: https://eu.wikisource.org/wiki/Laguntza:Orrialdeen_egoera
https://en.wikisource.org/wiki/Help:Page_status
- ▼ Produced dictionaries: <http://lexbib.org/larramendi>
- ▼ Code: <https://github.com/dlindem/LBLR/tree/master/Larramendi>

Eskerrik asko / Thank you!

*Funding acknowledgements:
UEU/EWKE Digital Humanities
Grant (2020),
Basque Government (IT 1344-19)*

*David Lindemann david.lindemann@ehu.es
Mikel Alonso mikelalon@gmail.com*