



# **Quasi-Parallel Corpora for Less-Resourced Languages: Parallelized Translations of Plato´s Faidon in Basque and Finnish.**

Koldo J. Garai

University of Eastern Finland (UEF) & University of the Basque Country (EHU)

Digital Research Data and Human Sciences (DRDHum).

2024/12/10

# OBJECTIVES

- To present a possible “easy” path to provide and enrich technologically less-resourced languages
  - Following the aims of the European Language Equality (ELE) program
  - It is a proposal to enlarge this investigation with more texts and non-Indo-European languages
- To present our outcomes to be tested regarding their efficiency and whether they meet the goals

# What has been done

- Two independent translations of Plato's Faidon
  - Finnish Faidoni by Calamnius 1887
  - Basque Faidon by Zaitegi 1975
- They have been parallelized and automatically annotated in Universal Dependencies
  - Now we look at the terms “less-resourced language” and “Quasi-parallel”

# Less-resourced language

## ■ European Language Equality

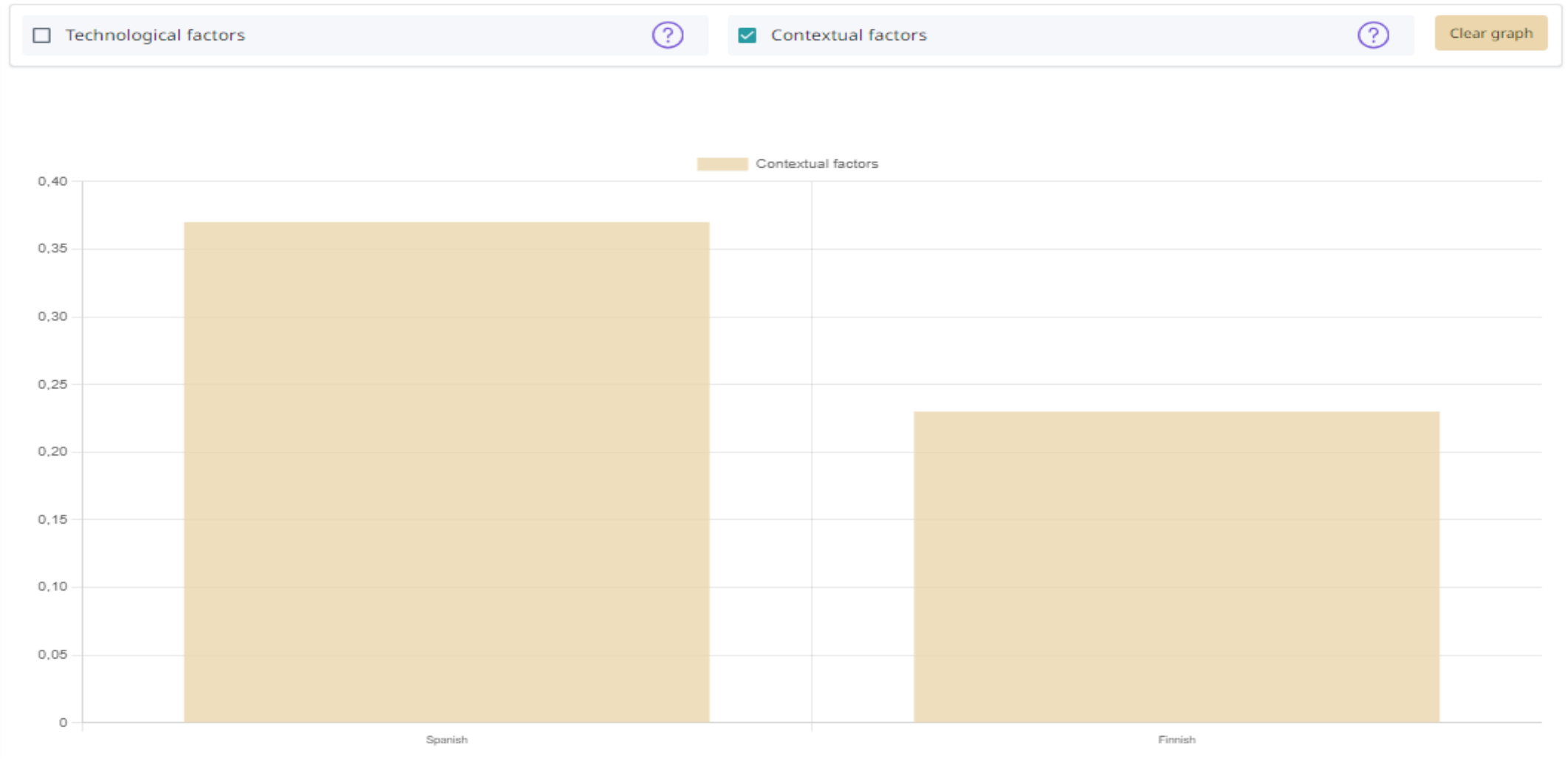
- the mission of the ELE Programme is to achieve language equality in Europe by 2030, reducing the technology gap with English and addressing the lack of available language data.
- “Digital Language Equality (DLE): all languages have the technological support and **situational context** necessary for them to continue to exist and to prosper as living languages in the digital age.” (Gaspari et al., 2023: 43).
- DLE Metric is a measure that reflects the digital readiness of a language and its contribution to the state of technology-enabled multilingualism”.
  - It considers two-factor branches: technological and contextual

# Technological & Contextual factors for metrics

- The technological factors (TFs) include language resources, tools, and services listed in the European Language Grid (ELG) catalog.
- The contextual factors (CFs) are based on societal, economic, and industrial conditions, and ecosystems of languages, to determine the potential for Language Resources and Technologies (LRTs), the extent to which a language has a **context that supports** the possibility of **evolving digitally** or not.

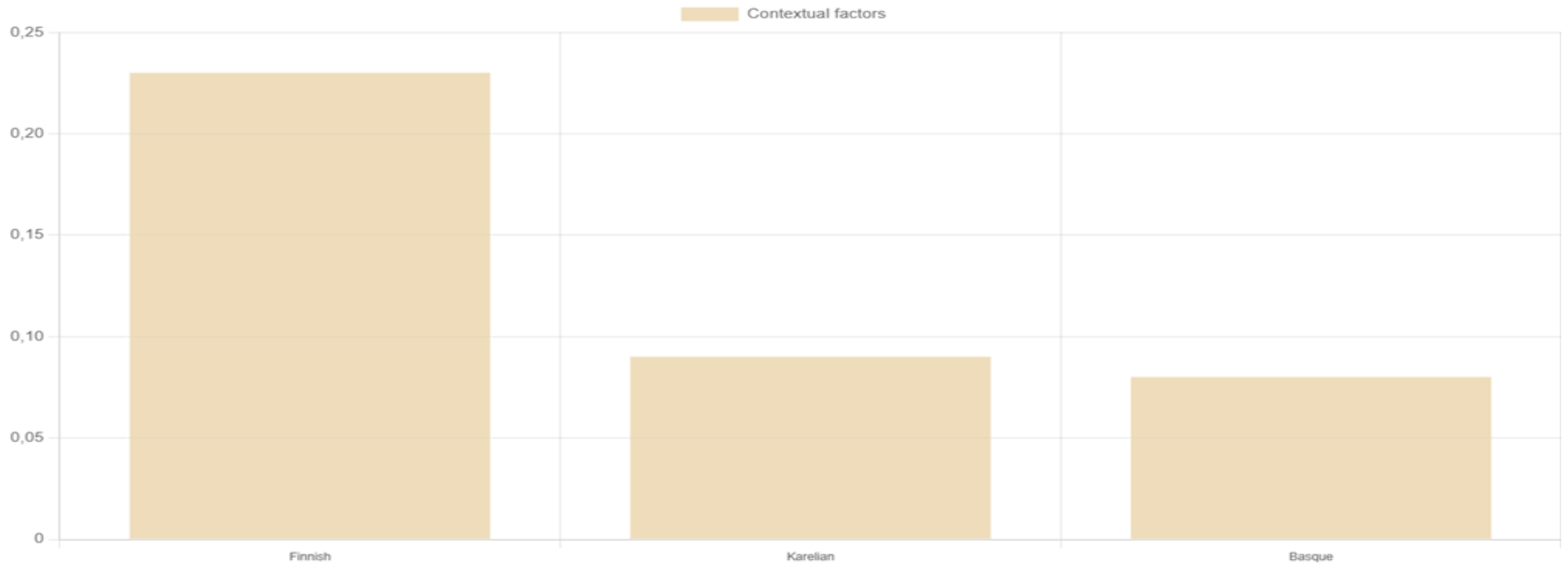
# Contextual factors: Spanish vs Finnish

<https://live.e>

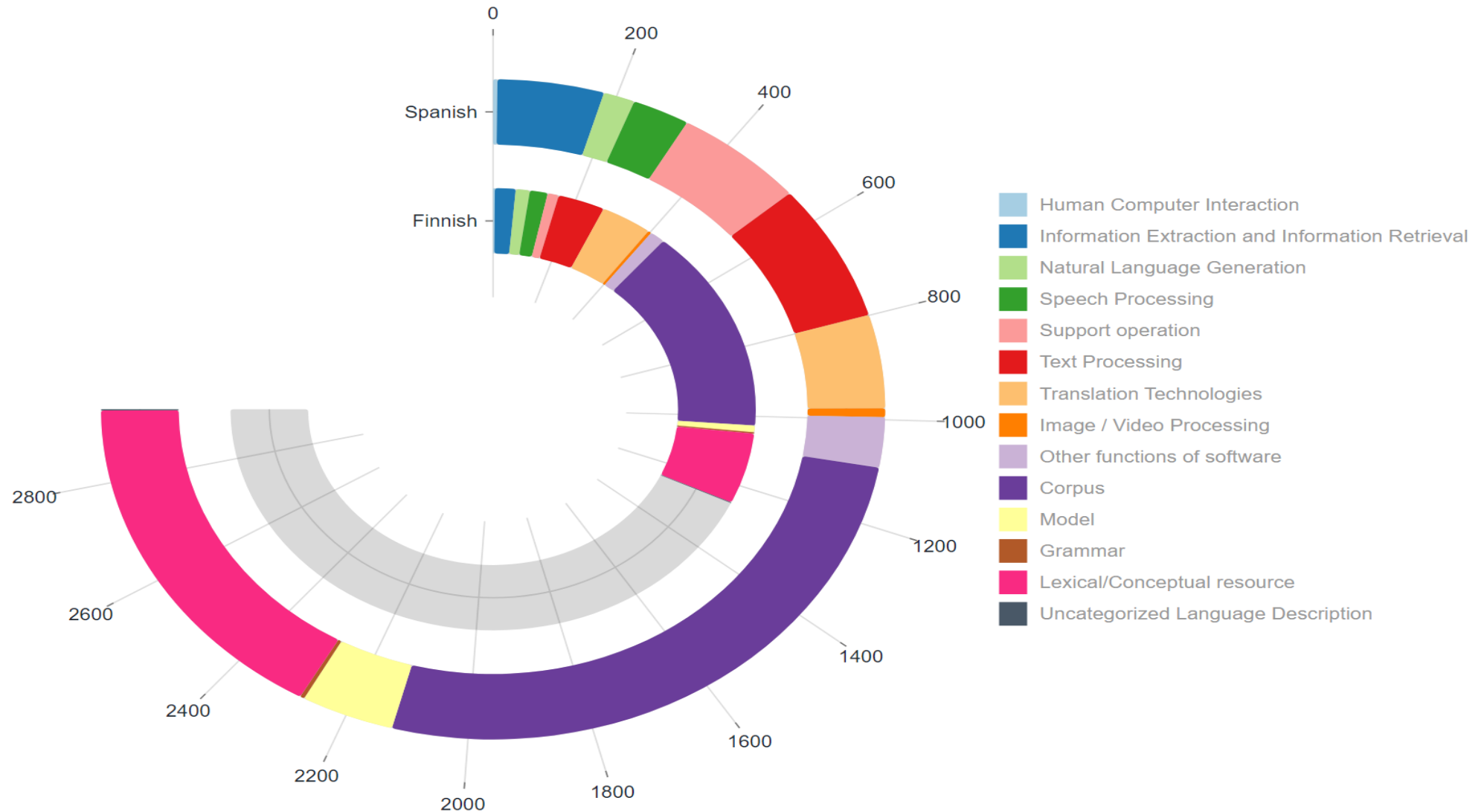


<https://live.european-language-grid.eu/catalogue/dashboard>

# Contextual Factors: Finnish, Karelian, Basque

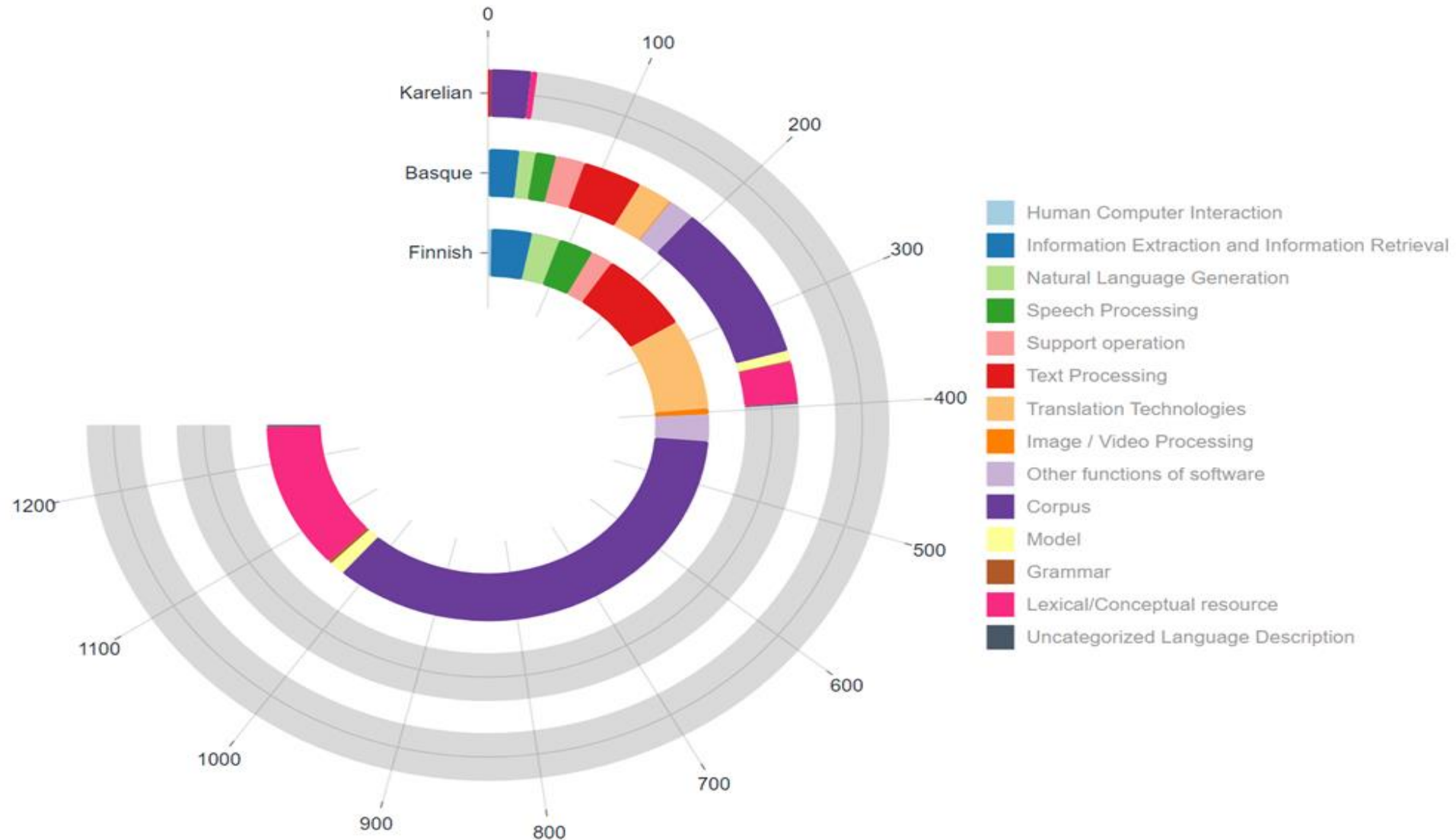


# TFs: Spanish / Finnish





# TFs: Finnish / Karelian / Basque



# Parallel vs Comparable Corpora: Quasi-Parallel?

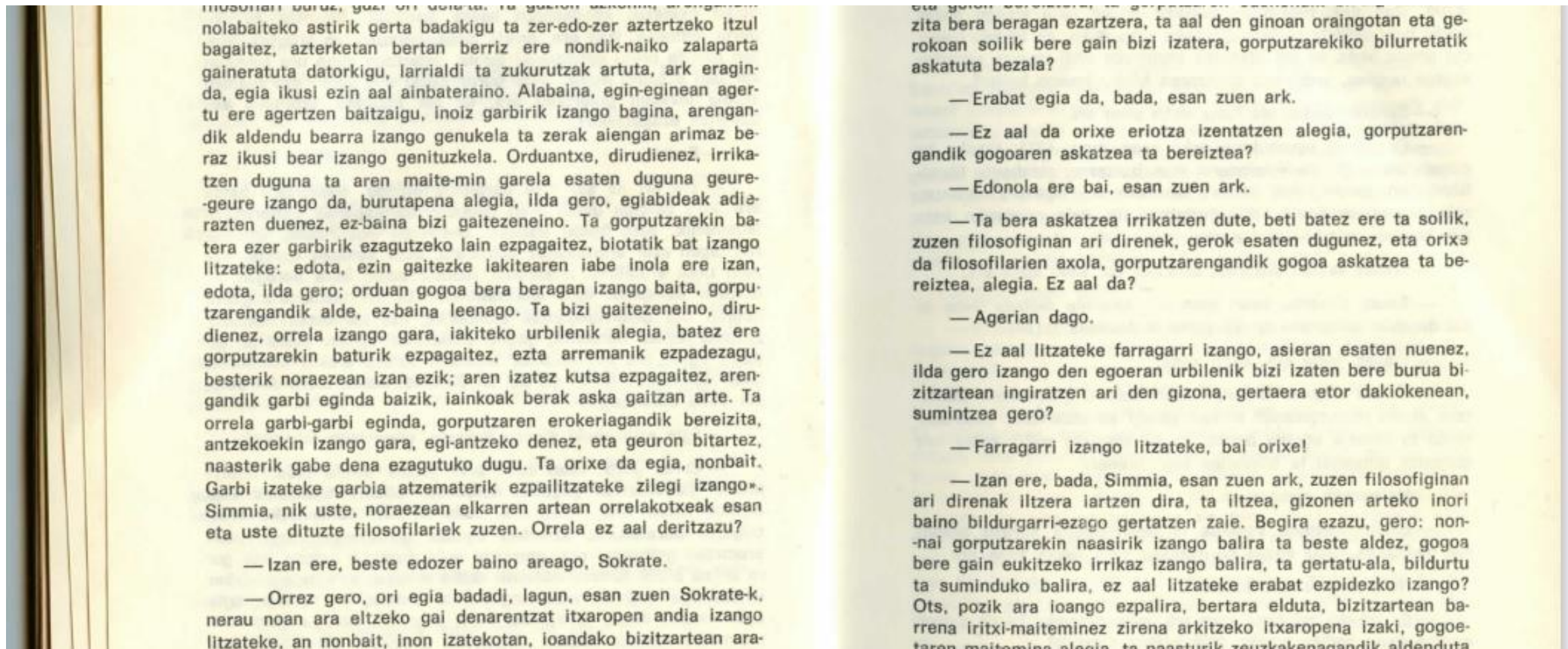
- Type A: Source texts plus translations ( $L1 > L2$ ). This is the parallel corpora, strictly speaking.
- Type B: The comparability is given by the design of the same sampling frame. ( $L1 + L1'$ ), similar balance and representativeness; the same proportions of the texts of the same genres, same sampling period, in the same domains, and not translations of each other.
- A combination of A and B . This also is included in the comparable category (McEnery and Xiao, 2018)

# Process

- 1) Scanning the printed version and editing the OCR errors to have it in a machine-readable format of the Basque text,
- 2) Aligning both texts,
- 3) Standardizing the Basque version,
- 4) Annotating both texts, original Finnish and standardized Basque.



# Digitizing: The text





# Digitizing: the process

- 1) Transform the **hard copy** text into txt: scanning and OCR reading
- 2) Process the txt: find **clusters** of the misread tokens (Word Smith Tools 8. Scott, Michael, 2021); **hyphens** are part of the token.
- 3) Extract the list of tokens, including the frequencies, to single out OCR mistakes from recurrent forms.
- 4) Organize the tokens **alphabetically**
  - The errors should be in the least frequent words,
  - the character sequences should cluster around the more frequent correctly read words.
- 5) Create the **concordances** per token and break the concordances alphabetically for each key token
  - We needed a wider context than single words to **disambiguate homo-form** cases.
  - Locate the reference in the .txt text,
  - Keep the scanned text also available to visually double-check the token.
- 7) Evaluate the speed of correction in a sample time, i.e., one hour of work.



# Digitizing: the dictionaries

- Rationale:
  - Create a dictionary of **systematic OCR errors**
    - It could be used for same typeface with more texts.
- Outcome:
  - The **corrected original** text
  - A golden rule

AALDE	AAIDE
AALDEAK	AAIDEAK
AALDE-KLDE	AAIDE-KIDE
AALEGLN	AALEGIN
AALEGLNAK	AALEGINAK
AALEGLNKA	AALEGINKA
AAILK	AALIK
AALLK	AALIK
AAZURLK	AAZURIK
ADLERAZL-BEARREAN	ADIERAZI-BEARREAN
ALERBIDE	AGERBIDE
AGRTUKO	AGERTUKO
AGERTREN	AGERTZEN
ALN	AIN
ALNAKOA	AINAKOA
ALNBATEAN	AINBATEAN
AINRECTE	AINRECTE

# Alignment

- The leading **intuition**: to create "fake" Translation Memories using CAT aligners (TM).
- The **problem**: the segmenting algorithms look at punctuation and string length
- The solution: manual aligning combining CAT tools, Aligners, and direct and inverse automatic translation
  - Content based aligning
- The question: is the product reusable in further texts?



# TMX Editor: the segmentation issue

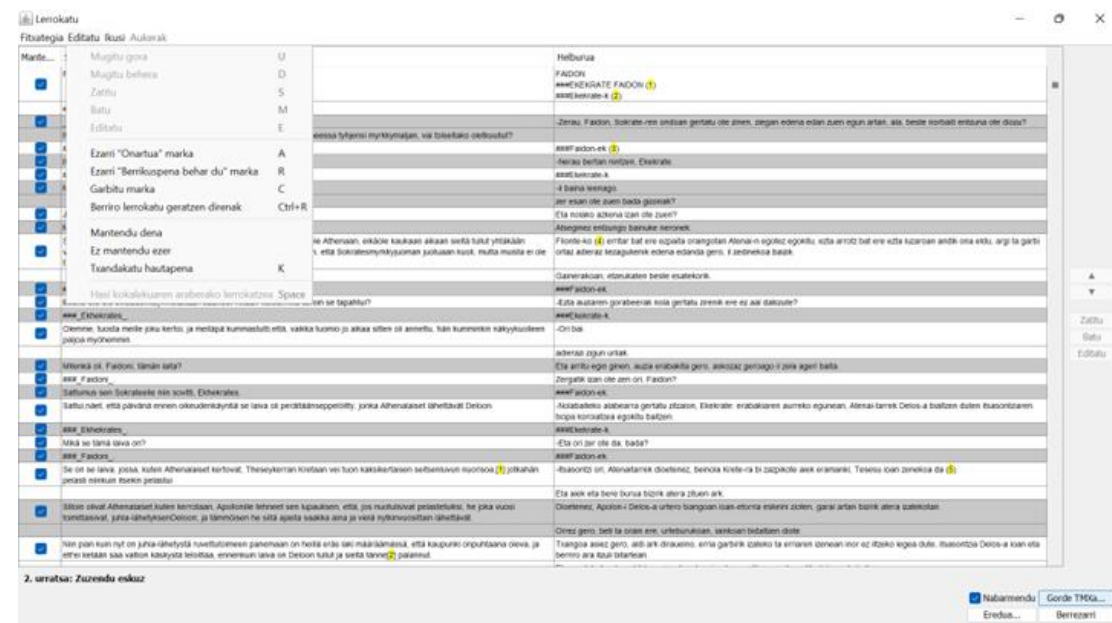
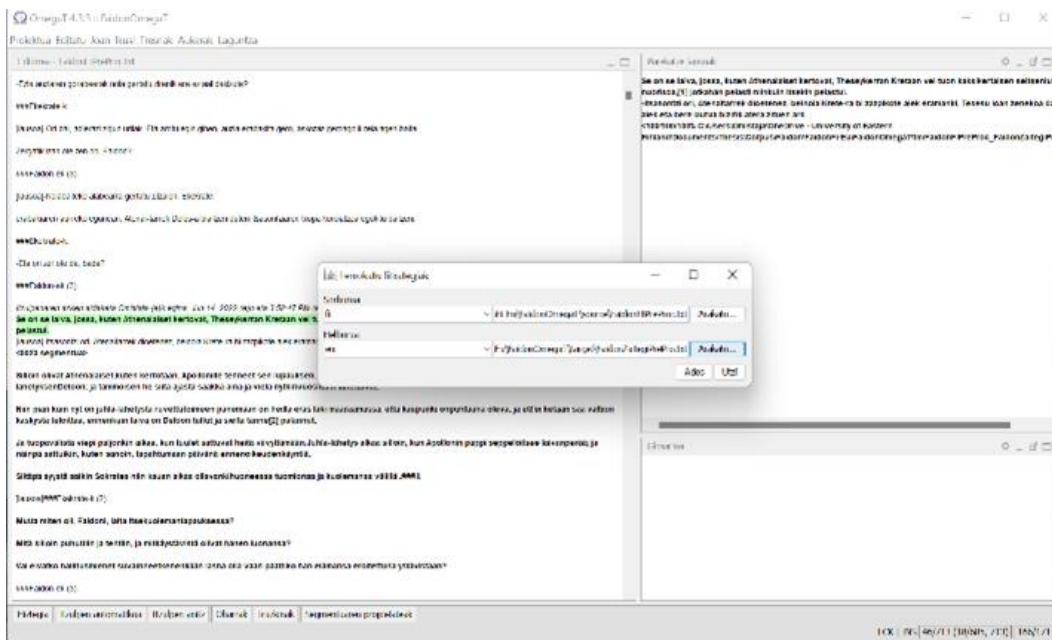
1	Faidoni	FAIDON
2	[L1> _Ekhekrates_	Ekhekrate-k [L2>
3	Itsekö, Faidoni, oitit Sokrateen luona sinä päivänä, jona hän vankhuoneessa tyhjensi myrkkymaljan, vai loisetako otel kuullut?	-Zerai, Faidon, Sokrate-ren ondoan gertalu ote zinen, ziegan edena edan zuen egun artan, ala, beste norbaiti entzuna ote diozu?
4	_Faidoni_	Faidon-ek [L3>
5	Itse siellä olin, Ekhekrates.	-Neraiu bertan nintzen, Ekhekrate.
6	_Ekhekrates_	Ekhekrate-k.
7	Mitä se siis olikaan, jota se mies ennen kuolemataan puhui?	-Ii baina leenago, zer esan ote zuen bada gizonak?
8	<b>Ja miten kuoli hän?</b>	<b>Eta nolako azkena izan ote zuen?</b>
9	Mieteltäni tuota kuulisin.	Atseginex entzungo bainuke neronek.
10	Sillä tähän aikaan ei ainoakaan Filialaisista kansalaisistani matkustele Athenaan, eikä ole kaukaan aikaan siellä hulfut yhtäkään vierasta, joka olisi tietänyt meille jotain varmaa noista kertoa, paitsi sen, että Sokrates myrkyttyoman juotuaan kuoli.	Filonite-ko [L4> erriat bat ere ezpaita oraingotan Atenai-n egotez egokitu, ezta arrotz bat ere ezta luzaroan andik ona eldu, argi ta garbi ortaz adieraz liezagukunik edena edanda gero, il zedinekoa balzik.
11	mutta muista ei ole tiedetty mitään jutella.	Gainerakoan, etzeukaten beste esatekorik.
12	_Faidoni_	Faidon-ek.
13	Ettekö siis ole oikeudenkäynnistäkään saaneet mitään tietoa, millä tavoin se tapahtui?	-Ezta auziaren gorabeerak nola gertalu zirenik ere ez aai dakizute?
14	_Ekhekrates_	Ekhekrate-k.
15	Olemme, luosta meille joku kertoi, ja meitäpä kummashutti, että, vaikka tuomio jo aikaa sitten oli annettu, hän kumminkin näkyy kuolleen paljoo myöhemmin.	-Ori bat, adierazi zigun urtiak. Eta arritu egin ginen, auzia erabakita gero, askozaz geroago il zela ageri baita.
16	Mitenkä oli, Faidoni, tämän laita?	Zergatik izan ote zen ori, Faidon?
17	_Faidoni_	Faidon-ek.
18	Sattumus sen Sokrateelle niin sovitii, Ekhekrates.	-Nolabaiteko atabearra gertalu zitzaion, Ekhekrate:
19	Sattui, näet, että päivänä ennen oikeudenkäyntiä se laiva oli perättään seppelöity, jonka Athenalaiset lähettivät Deloon.	erabakiaren zurreko egunean, Atenai-tarrek Delos-a bialtzen duten itsasontziaren txopa korotatzea egokitu baitzen.
20	_Ekhekrates_	Ekhekrate-k.
21	Mikä se tämä laiva on?	-Eta ori zer ote da, bada?
22	_Faidoni_	Faidon-ek.
23	Se on se laiva, jossa, kuten Athenalaiset kertovat, Thesey kerran Kretaan vei luon kaksikertaisen seitsemäluvan nuorisoo.[L1> Jotki hän pelasti niinkuin itsekin pelastui.	-Itsasontzi ori, Atenaitarrek dioetenez, beinola Krete-ra bi zazpikote aiek eramanki, Tesesu ioan zenekoa da [L5>. Eta aiek eta bere burua bizirik atera zitzen ark.
24	Silloin olivat Athenalaiset, kuten kerrotaan, Apollonille tehneet sen lupauksen, että, jos nuo tulisivat pelastetuiksi, he joka vuosi toimittaisivat juhla-lähetksen Deloon.	Dioetenez, Apolon-i Delos-a urtero txangoan ioan-eterria eskeini zioien, garai artan bizirik atera izatekotan.

CREATIONDATE="20220730T11:13:01Z" C:\Users\Ornitaja\OneDrive - University of Eastern Finland\Documents\Thesis\Corpus\Faidon\FaidonOrnitaja\FaidonOrnitaja	Erabat
jotka hallitsiväisyyttä noudattavat, eikä ole	-Eta zer? Orien arleko neurmidunek ere ez aai dute orae berorise izaten, nolabaiteko neurmigabekeriak zuer egin ditzaiean?
hallitsijasta.	Erabaitako dela esanaren, erabat nolako erabaitako gertatzen zaien izanda orae neurmitasunak buruz.
1 melkein mahdottomaksi, niin saabus lämpöön hillitsiväisyyteen tulee, heidän n, sillä pelään menettävänsä muita itävät he itsiään muuttamista nauhinnoista, hallitsijana.	Aien irrikan dauden beste atseginak uzteko bidur izaten direnez gero, baitzuegandik bagelu egiten dira, besteak gaina arfuta. Atseginak izatela egotza eragabekeria izentatzen dute, noski: atabaina, orretakeet, erabat atseginak garaituak izanda, beste atseginak garaitza gertatzen zaien.
in nauhinnojen alla orjaaminen, kuitenkin utamia hahuja hallitsivat oysttä että ovat	Ori, orraitzu esaten nuenaren antzeko da; nolabaik neurmigabekeriaren bidur zentzuraztea, alegia. -Egi-antzeko baita.
1, mitä istken sanoin, että he tavattaisa evät."	-Oimma zeroneko ori, ezta ori onbidetari buruz aidaaketa zuzena, atseginak atseginen aldera, atsekeabeak atsekeabeen aldera, bidurak biduraren aldera, andiagoa baitzuegaren aldera aidaaketa, dirua bezala.
1711 tämä viikkokausi oli siksi tuhtimman nauhinnoja saattaa, heidän istaan sekä suurempain pienempai istaan, vaan että ainoa ja oikea raha, jota isten, on viikkokausi, ja että kaikella sillä istaan ja myydään, on tiedonkoti oroa. Itseväisyys ja oikeameleisyys sekä yri on annustaan viikkokausi kansa ta ntot, pellet ja kaikki muut senkattavat. 1711 Mutta jos nämä kaikki eroitetaan pois isten vastaan, niin on demokraatin hyve kansain orjallinen, että ote siitä mitään ote on demokraatin erik puolesta kaitoa kyyri, joka annustaan ja annustasun kansa puolesta kaitoa.	1711 Egin-eginean, oritak gertak ezat ta kaitu egiten dira, senentzuean, neurmitasuna, biduraztean, ta erabat arfuta, egotza esanaren gogotaren lagun egotza, atseginak, bidurak eta arfuta gaiterituen beste gertak erabaita ta kaituta. Gogotaren aiek beronita ta aietzen bitaneri aidaaketa, aiatu esanena ezta nolabaiteko bidur-aitaketa buruz, eta egin-eginean bidur esanena, erabaitako dela egiten ere arfutatuta.
ään ne, jotka meille ovat mysteerit vaan itse toessa he jo ammooilla ajoilla a sillä opista, että joka puhdistamatonas ekskeen, saa saantassa maata, mutta joka lnyä saapuu, hän saa jumalien kanssa	Izan ere, aiatu gertak nolabaiteko gertaketa da, neurmitasuna, zindotasuna ta kenera, ta gogoteta bera garbitide besterik ezatela egituko da. Itzulu zedonak ezarri idikugun aiek ezpaitzen befauner beebiko gizen, dirudienez, egin-eginean aspaldidantik asmakartit itz egin baitzufen, alegia, Adea-a argitugabe ta ikangabe elizen dena aitan egongo dela; garbitunik eta eskainirik ara bertara elizen dena, bertuz, kaituen ondoen bizi izango dela.
181, Hyvän-kantajia on monta; mutta mie 1711: Tämmöisiä ovat miehen	Izan ere, aiatu baita, atsearrek aipatzen dituztenek dioetenez, "aitak eta noski atse-erabaitako baita erabaitako erabaita."



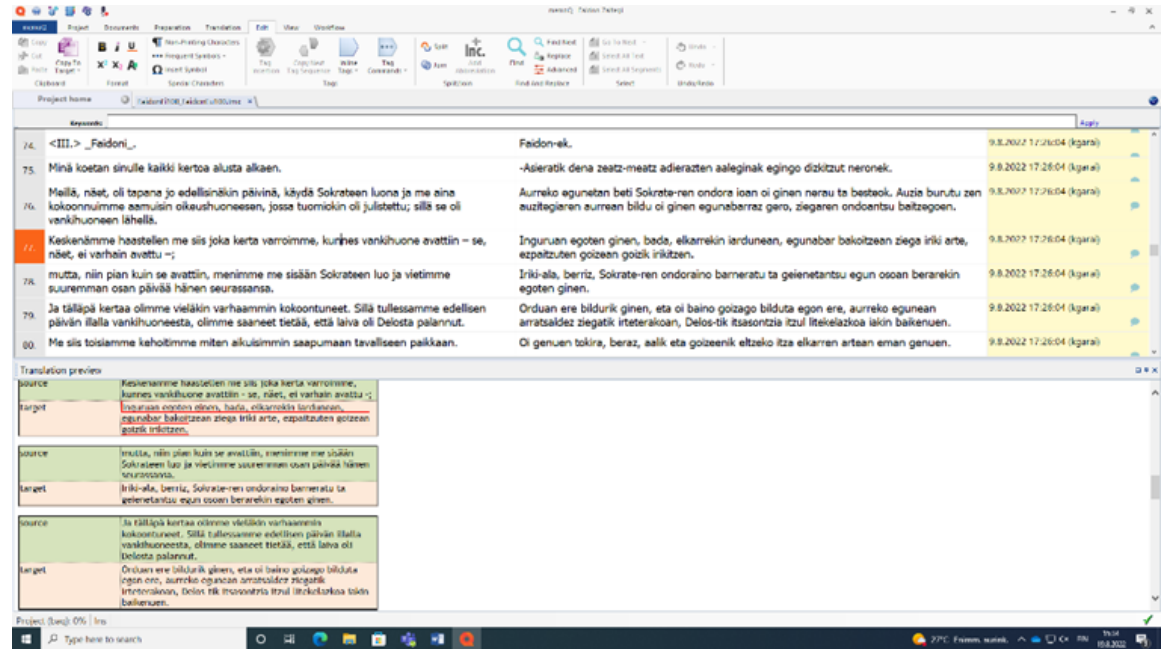
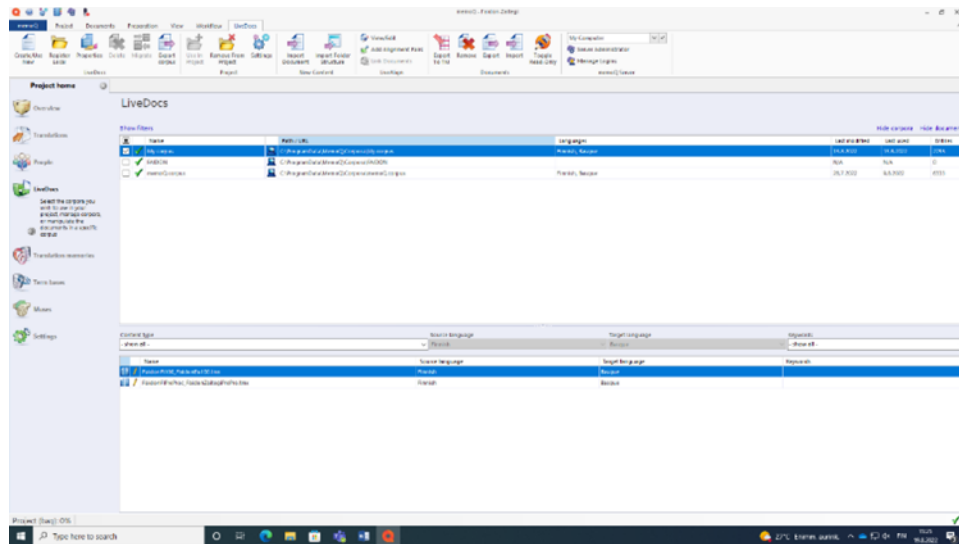


# Fake Translation Memories with CAT systems: OmegaT





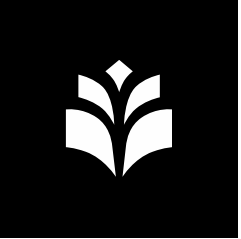
# Recovery by Livedocs (memoQ)





# Content based manual alignment

The screenshot displays the Google Translate web interface with two side-by-side windows. The left window shows the source text in Finnish: "Ja täällä kertaa olimme vieläkin varhaammin kokoontuneet. Sillä tullessamme edellisen päivän illalla vankihuoneesta, olimme saaneet tietää, että laiva oli Delosta palannut." The right window shows the target text in Euskara: "Eta oraingoan lehenago ere bildu ginen. Zeren bezperako arratsaldean kartzelatik etorri ginenean, ontzia Delosetik itzuli zela jakin genuen." The interface includes language selection menus (FINLANDIERA, EUSKARA, INGELESA), a character count (173 / 5.000), and navigation icons (Historia, Gordetakoak, Egin ekarpena). The bottom of the image shows a Windows taskbar with the date 19.8.2022 and temperature 27°C.



# Beyond punctuation and string length

The screenshot shows the memoQ software interface. The main window displays a list of source and target text segments. A 'Quick find' dialog box is open, showing search results for the word 'mutta'. Below the main list, a 'Translation preview' window shows a detailed comparison of source and target text for a specific segment.

Line	Source	Target	Date
116.	"Emme ainakaan mitään tarkkaa, Sokrates."	-Argi ta garbi ezer ere ez, Sokrate.	
117.	"No niin, kuulonhan mukaan minäkin näistä tulen puhumaan; ja mitä olen sattunut näistä kuulemaan, sen mielelläni teille puhun.	-Nik neuk ere omenka bakarrik ortaz esaten dut; eta ent gabe esango dizuet.	
118.	Niinhan lieneekin sopivinta, että, kun ollaan aikeissa lähteä tuonne pois, mietitään ja keskustellaan sinne-lähdöstä, minkälaisen luulemme lähdön olevan;	Ain zuzen, batez ere ara ioateko denari, arako bidea nolakoa izango den uste dugun aztertea ta elezarretan adieraztea baitagokio.	9.8.2022 17:26:04 (kgarai)
119.	mitäpä muuta meillä olisikaan tekemistä koko ajan päivän laskuun saakka?"	Zer besterik egin daiteke eguzkia etzaterainoko bitarte ontan? <19>.	9.8.2022 17:26:04 (kgarai)
120.	<VI.> "Mutta millähän syyllä itsensä surmaaminen väitetään vääräksi, Sokrates?"	-Zergatik esaten dute, bada, eztela egundo norberak bere burua iltzea zilegi, Sokrate?"	9.8.2022 17:26:04 (kgarai)
121.	Sillä, niinkuin äsken kuulustelitkin, olen jo ennen sekä Filolaolta, kun hän meidän luonamme oleskeli, että muutamilta muiltakin kuullut, ett'ei sovi sitä tehdä;	Dagoneko, nik neuk, oraintsu zerorrek galdetzen zenidanez, egotez gure artean lerrokatzen zen Filolao-gandik, baita beste zenbaitzuegandik ere entzuna daukat, eztela orrelakorik egin bear.	9.8.2022 17:26:04 (kgarai)
122.	mutta mitään tarkkaa näistä en ole keneltäkään koskaan kuullut."	Baina, argi ta garbi ortaz eztiot egundaino inori entzun, arean.	9.8.2022 17:26:04 (kgarai)

Translation preview

source	target
mutta mitään tarkkaa näistä en ole keneltäkään koskaan kuullut."	Baina, argi ta garbi ortaz eztiot egundaino inori entzun, arean.
"Ole vaan hyvillä mielin," sanoi Sokrates; "saanet sen plankin kuulla. Mutta sinusta arvattavasti kuuluu kummalta, että tämä väite yksin kaikista muista on ehdoton ja että se ei, niinkuin kaikki muu, riipu yksityisestä ihmisestä, niin että muutamasti ja muutamille olisi parempi kuolla kuin elää.	-Ots, gogoz ekin bear diogu, esan zuen. Bear bada, laster arritzekoa deritzakezunik entzun aal izango baituzu, beste guztien artean bakun auxe bakarrik badadi, beinik-bein, eta besteetan bezala, inoiz gizonari egokitu eypadakio, iltzea bizi izatea baino obe ote denentz iakitea.
Vaan mitä nyt niihin tulee, joille olisi parempi kuolla, niin arvattavasti sinusta kummalta kuuluisi, jos näillä ihmisillä	

# Standardization

- The Basque text was produced before the adoption of a written standard.
- The NLPs work usually only with modern spelling
  - Using original (pre-standard) spelling brings annotation errors.
- The “standardizer” has two components:
  - 1) a dictionary of **non-lemmatized occurrences** built by hand
  - 2) the **language-independent program**: it substitutes the tokens in the input text by reading the dictionaries and creating a new output text.

# Dictionary (1)

- Problem: **how to discriminate and extract** those tokens that need to be standardized reliably and systematically
  - Are the spell-checkers enough?
- Procedure:
  - When the actualization is obvious, we add directly the standardized token in the chart, helped by the spelling corrector Xuxen (IXA group)
  - When there is a doubt check it in the Basque Language Academy dictionary (Euskaltzaindiaren Hiztegia, n.d.);
  - When the word does not appear in the Basque language academy, just “trust the text”.

# Dictionaries

- The tokens are not lemmatized intentionally
  - For possible losses in grammatical information (multi word items, etc.)
  - What is **possible** vs. **probable** in language
    - Eventually the number of tokens will stabilize.
    - The substitution by rules can be prone to errors.

Word	WORD3
AAIDE	AHAIDE
AAIDEAK	AHAIDEAK
AAIDE-KIDE	AHAIDE-KIDE
AAL	AHAL
AALA	AHALA
AALBAITEKO	ALBAITEKO
AALEAN	AHALEAN
AALEGIN	AHALEGIN
AALEGINAK	AHALEGINAK
AALEGINEZ	AHALEGINEZ
AALEGINKA	AHALEGINKA
AALIK	AHALIK
AALTSUA	AHALTSUA
AALTSUAGORLK	AHALTSUAGORIK
AANTZ	AHANTZ
AANZTERAINO	AHANZTERAINO
AAZTEA	AHAZTEA
AAZTU	AHAZTU
AAZTURIK	AHAZTURIK
ABAR-AGADUNIK	ABAR-HAGADUNIK
ABERE-SORRETARA	ABERE SORRETARA
ADES	HADES
ADESA	HADESA
ADES-A	HADESA
ADES-AREN	HADESAREN
ADES-EKO	HADESEKO
ADES-EN	HADESEN
ADES-ERUNTZ	HADESERANTZ
ADES-KO	HADESKO
ADIERAZI-BEARREAN	ADIERAZI BEHARREAN
ADIMEN-GABEKERIA	ADIMEN GABEKERIA
ADIMEN-GABEKOA	ADIMEN GABEKOA
AGITZ	HAGITZ
AGORTU-EZINEKO	AGORTU EZINEKO
AIDEAK	AIREAK
AIEI	HAIEI
AIEK	HAIEK
AIEKIN	HAIEKIN
AIEN	HAIEN
AIENGAN	HAIENGAN
AIENGANDIK	HAIENGANDIK
AIGITO-N	EGIPTON

# Sorting tokens by length

- To avoid errors:
  - Substitution could happen inside the token, if using rules

WORD2	WORD3	# occurs	Length
armoni-gabekoarekin	harmonia gabekoarekin	1	19
erakutsi-aaleginean	erakutsi ahaleginean	1	19
gorputz-inarkunetan	gorputz inarkunetan	1	19
biotz-pozgarritzat	bihotz pozgarritzat	1	18
eztirautenetaruntz	ez dirautenetarantz	1	18
iauretxe-lapurreta	jauretxe lapurreta	1	18
litzatekenearenera	litzatekenearenera	1	18
adierazi-bearrean	adierazi beharrean	1	17
armoniagabearekin	harmonia gabearekin	1	17
bildurgarri-ezago	beldurgarri ezago	1	17
ta	eta	568	2
Ta	Eta	62	2
oi	ohi	23	2
il	<b>hil</b>	<b>22</b>	2
au	hau	21	2
an	<b>han</b>	<b>16</b>	2
io	jo	10	2
ar	har	9	2
An	<b>Han</b>	<b>3</b>	2
oe	ohe	2	2
as	has	1	2
eo	eho	1	2
ll	<b>Hil</b>	<b>1</b>	2
lo	Jo	1	2
Oi	Ohi	1	2



# The "inclusive standardizer"

- The program should perform these tasks:
  - 1- To sort the extracted list of tokens of the dictionary by decreasing number of characters, from the largest word (e.g. "zuzengabeagoarengandik") to the shortest (e.g. "jo").
  - 2- To find the token of the left column of the dictionary in the text and substitute it with the correspondent token of the column in the right.
  - 3- To repeat this process until no more instances of the left token are found in the plain text.

# The program

- The steps the program should take are these:
  - Open the CSV file and read it into a dictionary where the keys are the original words, and the values are the standard words.
  - Open the text file and read it into a string.
  - Loop through the dictionary and use the replace method to replace all instances of the original words with the standard words.
  - Write the corrected text back to the file.

# Universal Dependencies: practical analysis

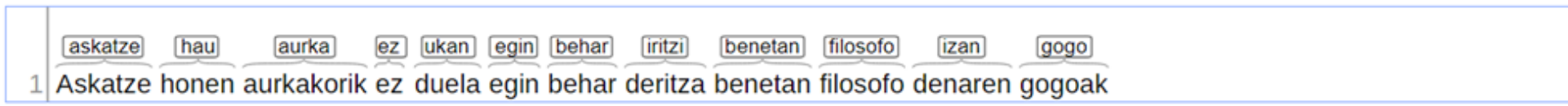
- Askatze honen aurkakorik ez duela egin behar deritza benetan filosofo denaren gogoak.
- Aska-tze[verb nominalization] / hau-EN[gen] / aurka-ko[gen-locative]-rik[partitive] / ez / du[Aux]-eLa[completive-object] / egin / behar / d-eritz-a [3 argument defective synthetic verb] / benetan / filosofo / da[Aux]-EN[relative]-a[det]-rEN[gen] / gogo-a[det]-K[ergative].
- Freeing / this-of[gen] / against-of-(ellipsis)-“any”[part] / not / has-that / do / must / be-lieves / truly / philosopher / is-that[relative]-the-of[gen] / mind-the[erg]
- The mind of that who is a true philosopher believes that it must not do (anything) against this freeing.

# Stanza

## Part-of-Speech (XPOS):

XPOS is not available for this language at this time.

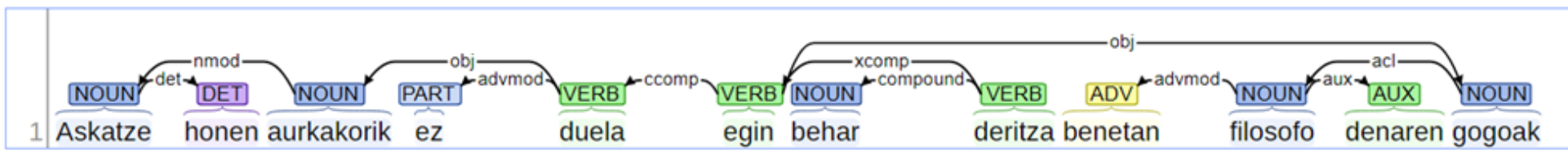
## Lemmas:



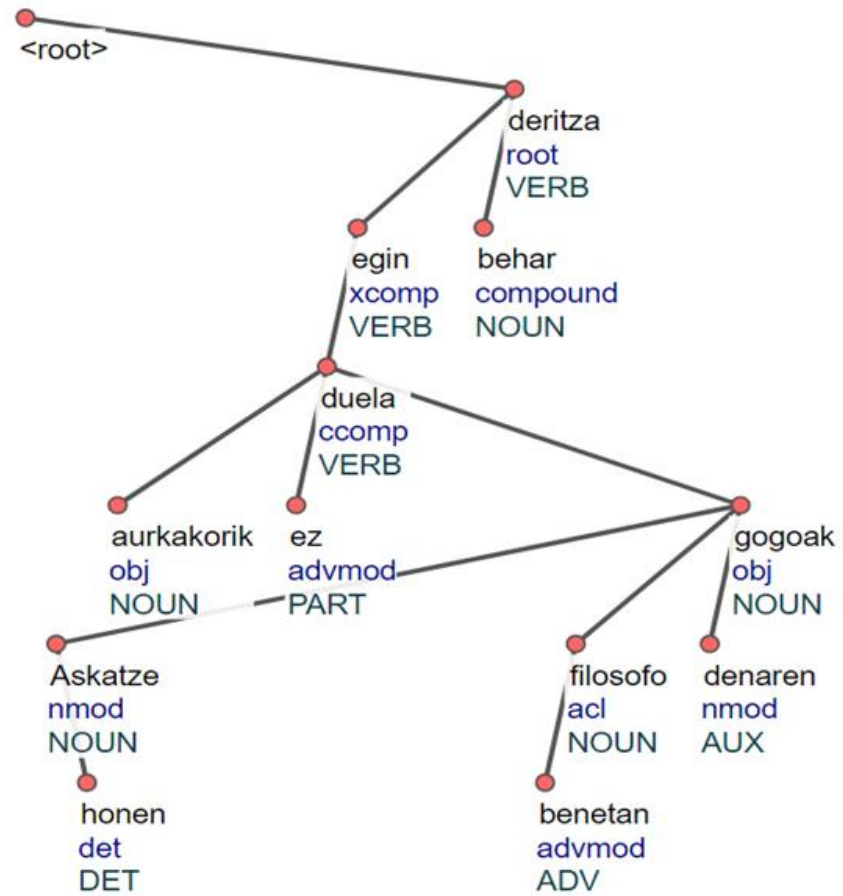
## Named Entity Recognition:

NER is not available for this language at this time.

## Universal Dependencies:



# UDPipe Lindat





# UDPipe/MaltIxa comparison. Lemmas and UPoS

Id	Form	Lemma1	Lemma2	UPoSTag	Category	Subcategory
1	Askatze	<u>askatze</u>	<u>askatu</u>	<u>NOUN</u>	<u>ADI</u>	ADI_SIN
2	honen	hau	hau	DET	DET	DET_ERKARR
3	aurka- korik	aurkako	aurkako	NOUN	IZE	IZE_ARR
4	ez	ez	ez	PART	PRT	PRT
5	duela	ukan	ukan	VERB	ADT	ADT
6	egin	egin	egin	VERB	ADI	ADI_SIN
7	behar	behar	behar	NOUN	IZE	IZE_ARR
8	deritza	iritzi	iritzi	VERB	ADT	ADT
9	benetan	benetan	benetan	ADV	ADB	ADB_ARR
10	filosofo	filosofo	filosofo	NOUN	IZE	IZE_ARR
11	denaren	<u>izan</u>	<u>dena</u>	<u>AUX</u>	<u>DET</u>	DET_ORO
12	gogoak	gogo	gogo	NOUN	IZE	IZE_ARR
13	.	.	.	PUNCT	PUNT	PUNT_PUNT



# UDPipes/MaltIxa Features comparison

Id	Form	Head1	Head2	DepRel1	Relation 2	Feats1	Feats2
1	Askatze	3	5	nmod	xmod	-	ADM:ADIZE
2	honen	1	3	det	ncmod	Case=Gen Definite=Def Number=Sing	KAS:GEN NUM:S
3	aurkakorik	5	5	obj	ncobj	Case=Par Definite=Ind	KAS:PAR
4	ez	5	5	advmod	ncmod	Polarity=Neg	-
5	duela	6	0	ccomp	ROOT	Aspect=Prog Mood=Ind Number[abs]=Sing Number[erg]=Sing Person[abs]=3 Person[erg]=3 VerbForm=Fin	ERL:KONPL ASP:PN T
6	egin	8	7	xcomp	xmod	VerbForm=Part	ADM:PART
7	behar	8	5	compound	ncobj	Case=Abs Definite=Ind	KAS:ABS
8	deritza	0	7	root	mw	Aspect=Prog Mood=Ind Number[abs]=Plur Number[erg]=Sing Person[abs]=3 Person[erg]=3 VerbForm=Fin	ASP:PNT
9	benetan	10	12	advmod	ncmod	-	-
10	filosofo	12	12	acl	ncmod	Animacy=Inan Case=Abs Definite=Ind	-
11	denaren	10	10	cop	detmod	Aspect=Prog Case=Gen Definite=Def Mood=Ind Number=Sing Number[abs]=Sing Person[abs]=3 VerbForm=Fin	KAS:GEN NUM:S
12	gogoak	6	5	nsubj	ncobj	Animacy=Inan Case=Abs Definite=Def Number=Plur	KAS:ABS NUM:P
13	.	8	8	advmod	ncmod	-	-

# Conclusions

## ■ Outcomes

- Dictionaries:
  - Systematized OCR errors with their corresponding corrections
  - Standardization dictionaries for that period
- Original text in machine readable format.
- Aligned texts Finnish-(original) Basque
- Aligned texts Original vs Standard Basque
- Automatically annotated texts, Finnish and Basque



# .. And Questions:

- If hand-curation is expensive,
  - to which extent is this process “an easy path” to enrich less-resourced languages?
  - To which extent are the content based aligned segments reusable in other similar future texts of the same domain?

# References

- Aranzabe, M. J., Atutxa, A., Bengoetxea, K., Diaz, A., de Ilarraza, Goenaga, I., Gojenola, K., & Uria, L. (2015). Automatic Conversion of the Basque Dependency Treebank to Universal Dependencies. *14th International Workshop on Treebanks and Linguistic Theories*, 9.
- Bengoetxea, K., & Gojenola, K. (2010). Application of different techniques to dependency parsing of Basque. *Proceedings of the NAACL HLT 2010 First Workshop on Statistical Parsing of Morphologically-Rich Languages*, 31–39.
- Estarrona, A., Etxeberria, I., Etxepare, R., Padilla-Moyano, M., & Soraluze, A. (2020). Dealing with dialectal variation in the construction of the Basque historical corpus. *Proceedings of the 7th Workshop on NLP for Similar Languages, Varieties and Dialects*, 79–89. <https://aclanthology.org/2020.vardial-1.8>
- Estarrona, A., Etxeberria, I., Soraluze, A., Etxepare, R., & Padilla-Moyano, M. (2021). The first annotated corpus of historical Basque. *Digital Scholarship in the Humanities*, fqab066. <https://doi.org/10.1093/lc/fqab066>
- *Euskaltzaindiaren Hiztegia—Bilaketa*. (n.d.). Retrieved May 21, 2023, from [https://www.euskaltzaindia.eus/index.php?option=com\\_hiztegianbilatu&task=hasiera&Itemid=1693&lang=eu-ES](https://www.euskaltzaindia.eus/index.php?option=com_hiztegianbilatu&task=hasiera&Itemid=1693&lang=eu-ES)
- *EUSTAGGER: Lemmatizer/tagger for Basque*. (2022, October 25). <https://ixa2.si.ehu.eus/eustagger/>
- Gaspari, F., Grützner-Zahn, A., Rehm, G., Gallagher, O., Giagkou, M., Piperidis, S., & Way, A. (2023). Digital Language Equality: Definition, Metric, Dashboard. In *European Language Equality* (pp. 39–73). Springer International Publishing. [https://doi.org/10.1007/978-3-031-28819-7\\_3](https://doi.org/10.1007/978-3-031-28819-7_3)
- Hovy, E. (2022). Chapter 21. Corpus Annotation. In R. Mitkov (Ed.), *The Oxford Handbook of Computational Linguistics* (pp. 508–517). Oxford University Press. <https://doi.org/10.1093/oxfordhb/9780199573691.013.011>
- IXA Taldea. (2004). *Ixa taldeko etiketen eskuliburua*. [https://garaterm.ehu.es/sites/default/files/dokumentuak/4083/Etiketen\\_eskuliburua\\_IXA.pdf](https://garaterm.ehu.es/sites/default/files/dokumentuak/4083/Etiketen_eskuliburua_IXA.pdf)
- Lapshinova-Koltunski, E., Popović, M., & Koponen, M. (2022). DiHuTra: A Parallel Corpus to Analyse Differences between Human Translations. *Proceedings of the 23rd Annual Conference of the European Association for Machine Translation*, 337–338. <https://aclanthology.org/2022.eamt-1.58>

# References

- *LF Aligner*. (2023, May 7). SourceForge. <https://sourceforge.net/projects/aligner/>
- Li, B., Gaussier, E., & Yang, D. (2018). Measuring bilingual corpus comparability. *Natural Language Engineering*, 24(4), 523–549. <https://doi.org/10.1017/s1351324917000481>
- Lindén, K., & Dyster, W. (2023). Language Report Finnish. In G. Rehm & A. Way (Eds.), *European Language Equality: A Strategic Agenda for Digital Language Equality* (pp. 135–138). Springer International Publishing. [https://doi.org/10.1007/978-3-031-28819-7\\_15](https://doi.org/10.1007/978-3-031-28819-7_15)
- *LiveDocs—memoQ*. (2024, November 15). Translation Software - memoQ. <https://www.memoq.com/tools/livedocs>
- *Maltixa | Ixa taldea*. (n.d.). Retrieved June 23, 2022, from <https://ixa.ehu.eus/node/4457?language=en>
- McEnery, T. (2022). Corpora. In R. Mitkov (Ed.), *The Oxford Handbook of Computational Linguistics* (pp. 494–507). Oxford University Press. [https://doi.org/10.1093/oxfordhb/9780199573691.013.47\\_update\\_001](https://doi.org/10.1093/oxfordhb/9780199573691.013.47_update_001)
- McEnery, T., & Xiao, R. (2018). *Parallel and Comparable Corpora: What is Happening?* (pp. 18–31). Multilingual Matters. <https://doi.org/10.21832/9781853599873-005>
- memoQ (Director). (2019, September 18). *Hidden Treasures in memoQ* [Video recording]. <https://www.youtube.com/watch?v=OjwGEppw92M>
- *META-NET White Paper Series—META Multilingual Europe Technology Alliance*. (n.d.). Retrieved October 24, 2023, from <http://www.meta-net.eu/whitepapers/overview>
- Moseley, C., Nicolas, A., & Unesco. (2010). *Atlas of the world's languages in danger* (3rd ed., entirely rev., enl.updated, p. 1 atlas (154, 62 pages) : 29 color maps; 20 x 29 cm + 1 map (86 x 119 cm folded to 29 x 18 cm).) [Map]. Unesco Pub. Paris.

# References

- Nivre, J., de Marneffe, M.-C., Ginter, F., Hajič, J., Manning, C. D., Pyysalo, S., Schuster, S., Tyers, F., & Zeman, D. (2020). *Universal Dependencies v2: An Evergrowing Multilingual Treebank Collection* (No. arXiv:2004.10643). arXiv. <https://doi.org/10.48550/arXiv.2004.10643>
- *OmegaT - The Free Translation Memory Tool—OmegaT*. (n.d.). OmegaT - The Free Translation Memory Tool. Retrieved August 19, 2022, from <http://omegat.org/>
- Pace-Sigge, M. T. L. (2013). The concept of Lexical Priming in the context of language use. *ICAME*, 37, 149–173.
- Plato. (1978). *Platon. IV., Kriton eta Faidon* (I. Zaitegi Plazaola, Trans.). Euskaltzaindia.
- Plato, 427? BCE-347? BCE. (2006). *Faidoni Platonin keskustelma Sokrateen viimeisistä hetkistä jasielun kuolemattomuudesta* (J. W. (Johan W. Calamnius, Trans.)). <https://www.gutenberg.org/ebooks/19210>
- *Project Gutenberg*. (n.d.). Project Gutenberg. Retrieved March 14, 2022, from <https://www.gutenberg.org/>
- Rehm, G., & Way, A. (Eds.). (2023a). *European Language Equality: A Strategic Agenda for Digital Language Equality*. Springer International Publishing. <https://doi.org/10.1007/978-3-031-28819-7>
- Rehm, G., & Way, A. (2023b). European Language Equality: Introduction. In G. Rehm & A. Way (Eds.), *European Language Equality: A Strategic Agenda for Digital Language Equality* (pp. 1–10). Springer International Publishing. [https://doi.org/10.1007/978-3-031-28819-7\\_1](https://doi.org/10.1007/978-3-031-28819-7_1)
- Salminen, T. (2010). Europe and the Caucasus. In *Atlas of the world's languages in danger*. UNESCO. <https://researchportal.helsinki.fi/en/publications/europe-and-the-caucasus>
- Sarasola, K., Aldabe, I., Ilarraza, A. D. de, Estarrona, A., Farwell, A., Hernáez, I., & Navas, E. (2023). Language Report Basque. In G. Rehm & A. Way (Eds.), *European Language Equality: A Strategic Agenda for Digital Language Equality* (pp. 95–98). Springer International Publishing. [https://doi.org/10.1007/978-3-031-28819-7\\_5](https://doi.org/10.1007/978-3-031-28819-7_5)

# References

- Scott, M., & Tribble, C. (2006). *Textual Patterns: Key Words and Corpus Analysis in Language Education / Mike Scott and Christopher Tribble* (Issue v. 22). John Benjamins Publishing Co.
- Scott, Michael. (2021). *WordSmith Tools* (Version 8.0.0.61) [Windows].
- *TMX editor*. (2015, September 16). SourceForge. <https://sourceforge.net/projects/tmxeditor/>
- *UDPipe*. (2022). [C++]. ÚFAL. <https://github.com/ufal/udpipe> (Original work published 2016)
- *UNESCO WAL*. (n.d.). UNESCO WAL. Retrieved October 21, 2023, from <https://en.wal.unesco.org/>
- *Xuxen | Ixa taldea*. (n.d.). Retrieved May 30, 2022, from <http://ixa.ehu.eus/node/4463?language=en>
- Zaitegi, Iokin. (1974). Platon Euskeratzen. *Euskera: Euskaltzaindiaren lan eta agiriak = Trabajos y actas de la Real Academia de la Lengua Vasca = Travaux et actes de l'Academie de la Langue basque*, XIX(1), 358–361.
- Zaitegi Plazaola, I. 1906-1979., Lafitte, P. 1901-1983., & Zaitegi Plazaola, Iokin. (1962). *Platon'eneko atarian*. [s.n.].
- Zeman, D., Dušek, O., Mareček, D., Popel, M., Ramasamy, L., Štěpánek, J., Žabokrtský, Z., & Hajič, J. (2014). HamleDT: Harmonized multi-language dependency treebank. *Language Resources and Evaluation*, 48(4), 601–637. <https://doi.org/10.1007/s10579-014-9275-2>
- Zeman, D., Nivre, J., Abrams, M., Ackermann, E., Aepli, N., Aghaei, H., Agić, Ž., Ahmadi, A., Ahrenberg, L., Ajede, C. K., Akkurt, S. F., Aleksandravičiūtė, G., Alfina, I., Algom, A., Alnajjar, K., Alzetta, C., Andersen, E., Antonsen, L., Aoyama, T., ... Ziane, R. (2023). Universal Dependencies 2.12. <http://universaldependencies.org/>. <https://lindat.mff.cuni.cz/repository/xmlui/handle/11234/1-5150>