

# Herramientas Digitales para las Humanidades Digitales en la e-infraestructura CLARIN

Mikel Iruskietea  
HiTZ zentroa-Ixa Taldea  
UPV/EHU

[www.clarin.eu](http://www.clarin.eu)  
[www.clarin-es.org](http://www.clarin-es.org)  
<http://ixa2.si.ehu.es/clarink>

Creación de un Proyecto en  
Humanidades Digitales basado  
en el análisis de textos:  
Modelado y Procesamiento



Creación de un proyecto en  
humanidades digitales basado  
en el análisis de textos:  
modelado y procesamiento

# Esquema

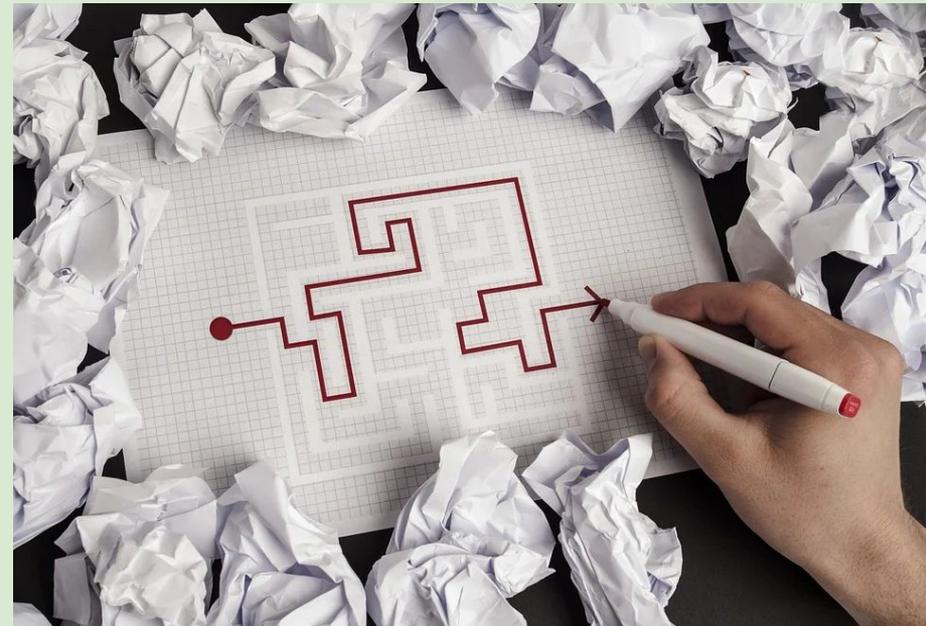
1. Resumen de la sesión
2. Introducción: e-Infraestructuras de investigación
  - a. Acciones conjuntas de infraestructuras
3. Justificación
4. Casos prácticos
  - a. Herramientas de análisis textual
  - b. Herramientas de análisis de voz
  - c. Casos de uso prácticos
  - d. Otros recursos
5. Conclusiones



# Mapa

## Resumen

- **Infraestructuras de investigación**
- Interoperabilidad (dentro y fuera)
- Justificación de infraestructuras
- Casos de uso



# Resumen de la sesión

## Objetivo

- *European Open Science Cloud*: EOSC
- Ciencia de principios FAIR: Encontrable, Accesible, Interoperable y Reutilizable
- Interoperabilidad en infraestructuras (CLARIN)

## Método

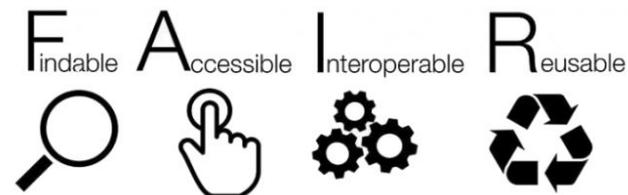
- Construir infraestructura que no se desarrollará en Europa para
  - Lenguas oficiales y cooficiales del estado
  - ALL-LT-in-ONE-URL: todos los recursos, todos los servicios

## Ejemplos

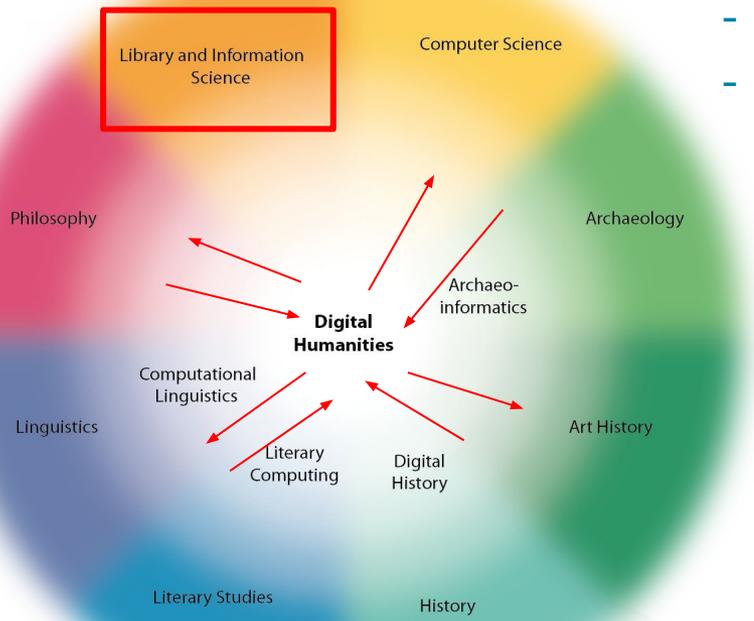
- Casos de uso y herramientas (sencillas) para investigar en infraestructuras europeas que se inter-comunican



**EUROPEAN OPEN  
SCIENCE CLOUD**



# Humanidades Digitales e infraestructuras



- Disciplinas: **Colaboración**
- Estudio de métodos/preguntas de investigación
  - **Método:** Cómo influyen las herramientas/recursos cuando se cuantifican algunos conceptos literarios (*Distant Reading*)
  - **Adaptar:** recursos de otros dominios
  - Modelado de datos, **metadatos**, bases de datos... (cómo se diseña el corpus literario y cuál es la pregunta de investigación)

- Colaboración e interdisciplinaridad
  - Humanistas digitales
  - Críticos literarios
  - Lingüistas
  - Informáticos

Teaching CLARIN  
in times of corona



CLARIN en el currículum universitario:  
<https://labur.eus/scuzN>

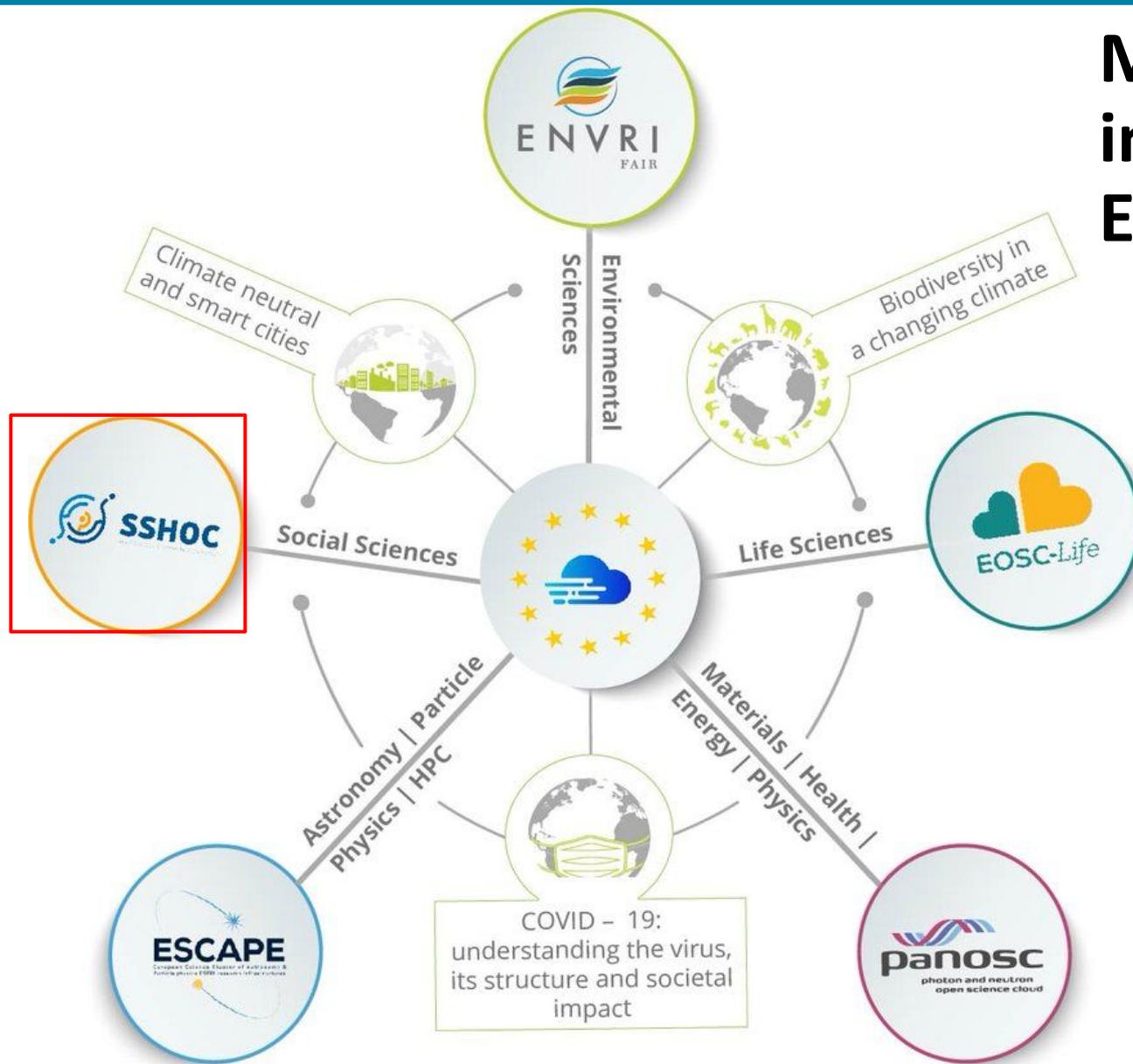
# Justificación de las infraestructuras

## Ciencia: abierta de impacto y reproducible:

- Fragmentación en la investigación de CCSS
- Recursos de CCSS dispersados en repositorios
- Poco re-uso en la investigación de CCSS
- Poca interdisciplinaridad
- Impacto social limitado

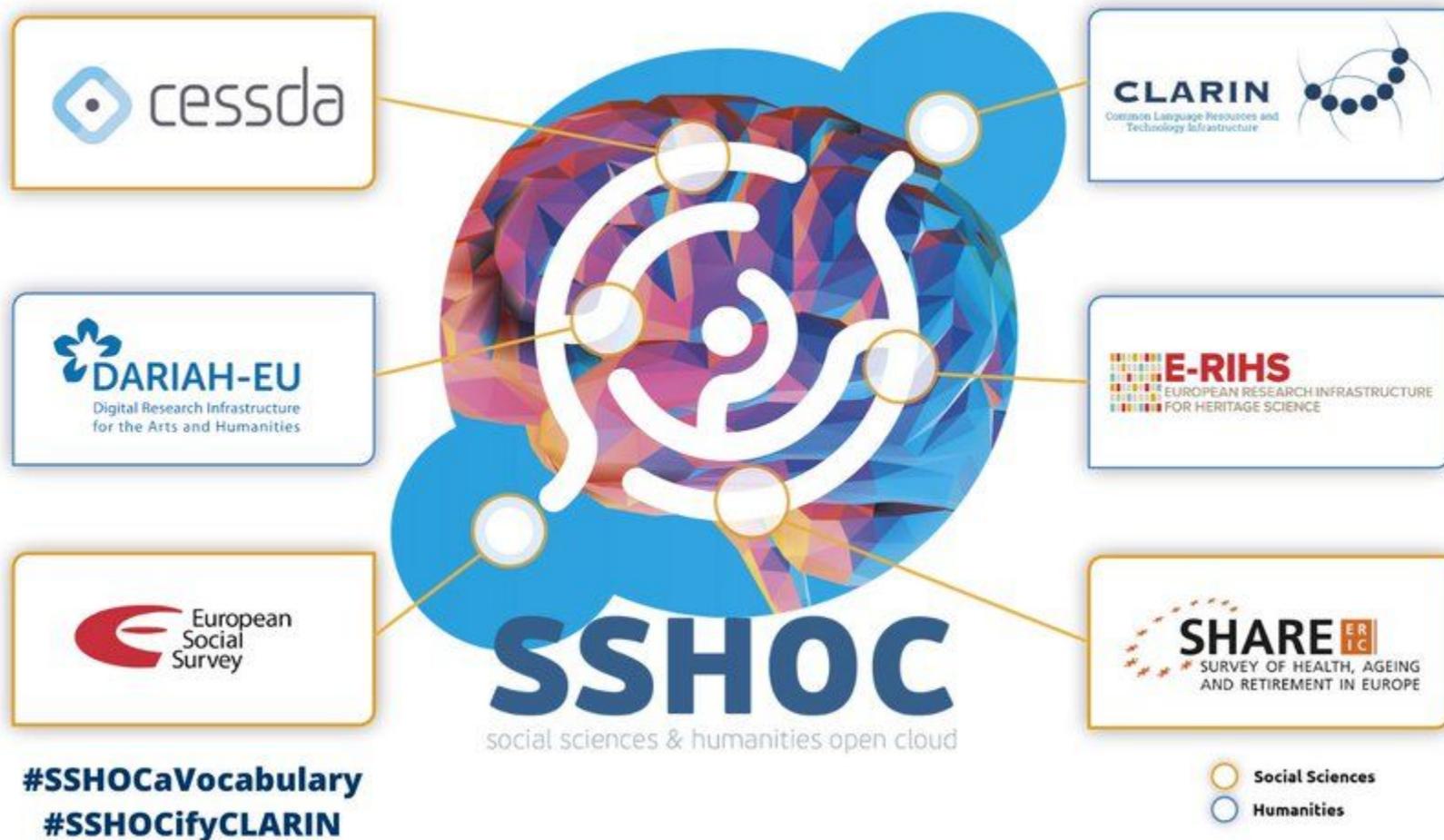


# Mapa de las infraestructuras EOSC

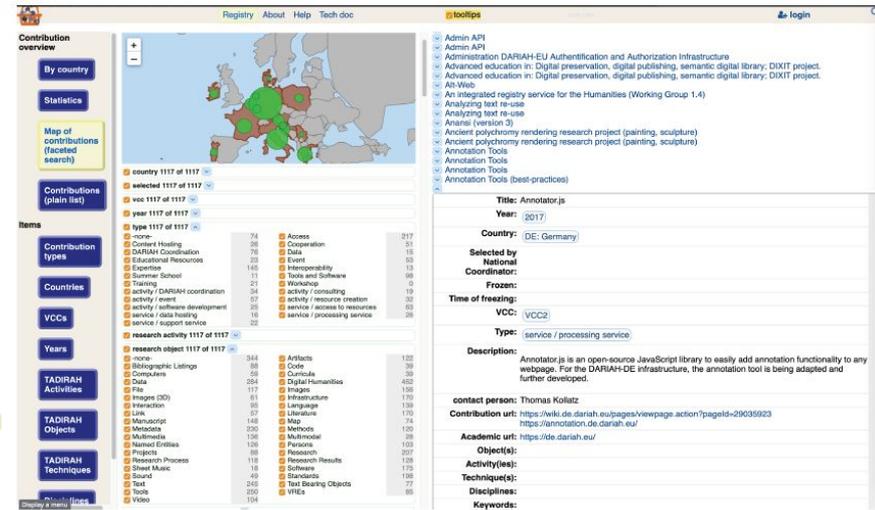
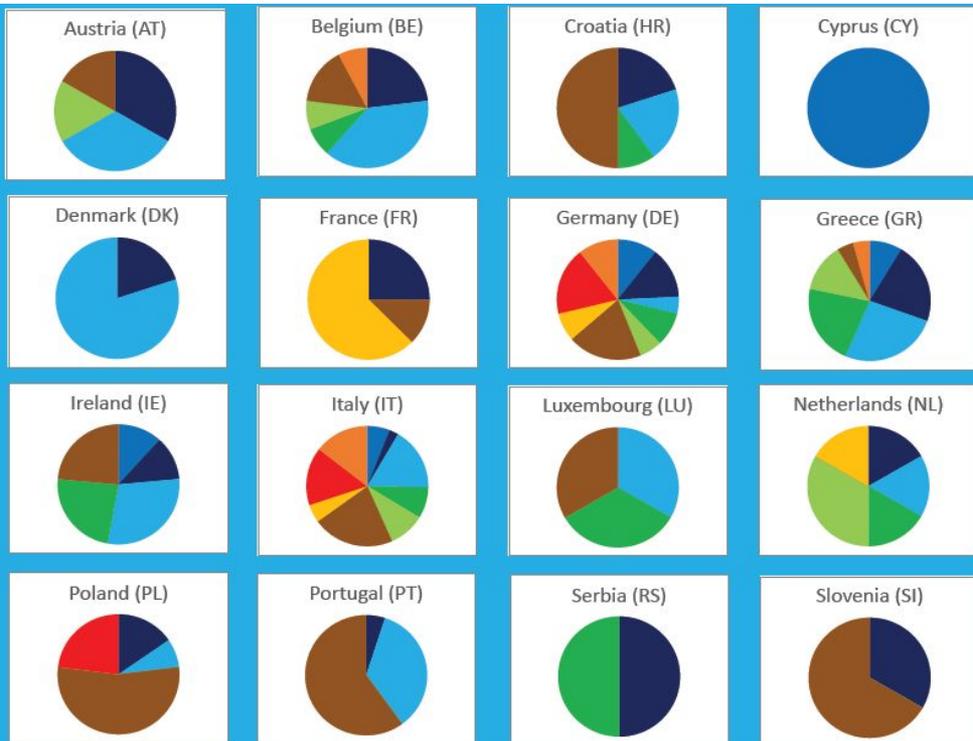


[#SSHOC](#) es una de las acciones de INFRAEOSC 04-2018, que consolida y conecta las e-infraestructuras europeas en **European Open Science Cloud**.

# SSHOC: Conexión de las e-infraestructuras europeas en HD y CCSS



# ¿Cómo contribuyen las e-infraestructuras?



# Creación de infraestructuras y red estratégica

- **INTELE**: red estratégica para la participación oficial en infraestructuras europeas CLARIN y DARIAH
  - **Impulsar** la investigación en humanidades y ciencias sociales
  - **Impulsar** los proyectos y programas **internacionales**
- **CLARIN**: **C**ommon **L**anguage **R**esources and Technology **I**nfrastructure
  - ESFRI ERIC (2012) y ESFRI Landmark (2016)
- **DARIAH**: **D**igital **R**esearch **I**nfrastructure for the **A**rts and **H**umanities
  - ESFRI ERIC (2014) y ESFRI Landmark (2016)



CLARIN

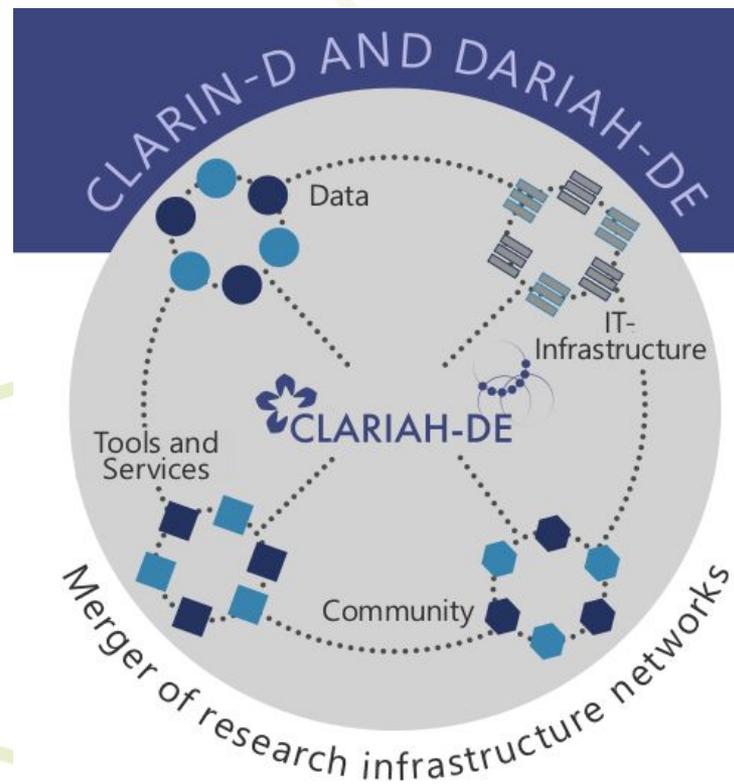


DARIAH-EU

# INTELE: Red estratégica para la promoción de las infraestructuras de tecnologías del lenguaje en eHumanidades y ciencias sociales



- (1) **Impulsar actividades** de promoción de las infraestructuras CLARIN y DARIAH
- (2) **Conectar grupos investigadores** que tengan interés para participar en dichas infraestructuras europeas
- (3) Elaborar un **catálogo de herramientas y casos de uso** para castellano y lenguas cooficiales (euskera, catalán, gallego)
- (4) Elaborar un informe para la **reevaluación positiva** de dichas infraestructuras



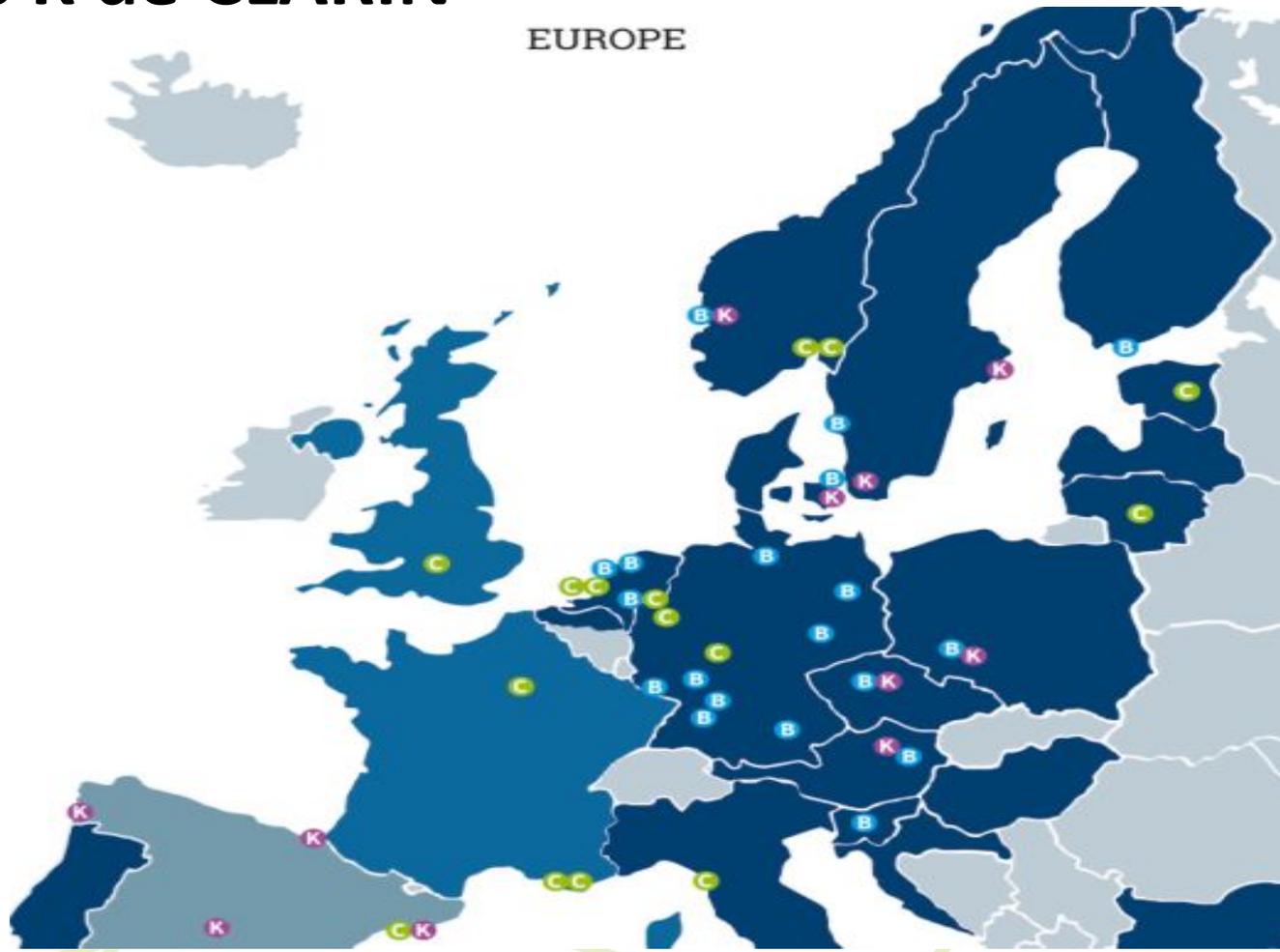
The content-related and technological foundations created by CLARIN-D and DARIAH-DE will be aligned, integrated, further developed and jointly maintained in CLARIAH-DE.



# Nodos y Centros K de CLARIN



- ERIC members
- Observers
- Countries with participating centres
- Centre Providing Data
- Centre Providing Metadata
- Knowledge Centre



# ¿Por qué una infraestructura?

## TL y situación de las lenguas

- 5 clases de situaciones en **2012**
  - Situación excelente:
  - Situación adecuada: inglés
  - Situación media: alemán **castellano**, francés holandés ...
  - Situación en desarrollo: **euskara**, gallego, **catalán**, esloveno, ...
  - Situación pobre: irlandés (gaélico), islandés, rumano, ...

- Investigar más rápido y con mayor calidad
  - Tener más tiempo para investigar
    - - tiempo programando
    - - tiempo creando recursos
    - + impacto social
    - + reutilización...

- Proyecto Europeo EUROPEAN LANGUAGE EQUALITY (ELE) 2020

- Desarrollar una agenda estratégica de investigación e innovación, y una hoja de ruta para lograr la igualdad total de las lenguas europeas en el ámbito digital para 2030.

META-NET

Offin Noyce (coord.) · Clara Sureda Mollat

THE BASQUE EUSKARA  
LANGUAGE ARO  
IN THE DIGITALEAN  
DIGITAL AGE

Arrokazuak Herriak  
Eus Herriak  
Igor Oñativita  
Kajal Sorribila  
Herriak: Oñativita, Sorribila  
Igor Oñativita  
Herriak: Oñativita, Sorribila  
Igor Oñativita  
Herriak: Oñativita, Sorribila

# Construcción de la infraestructura: CLARIAH

National Roadmap for Large-Scale Research Infrastructure

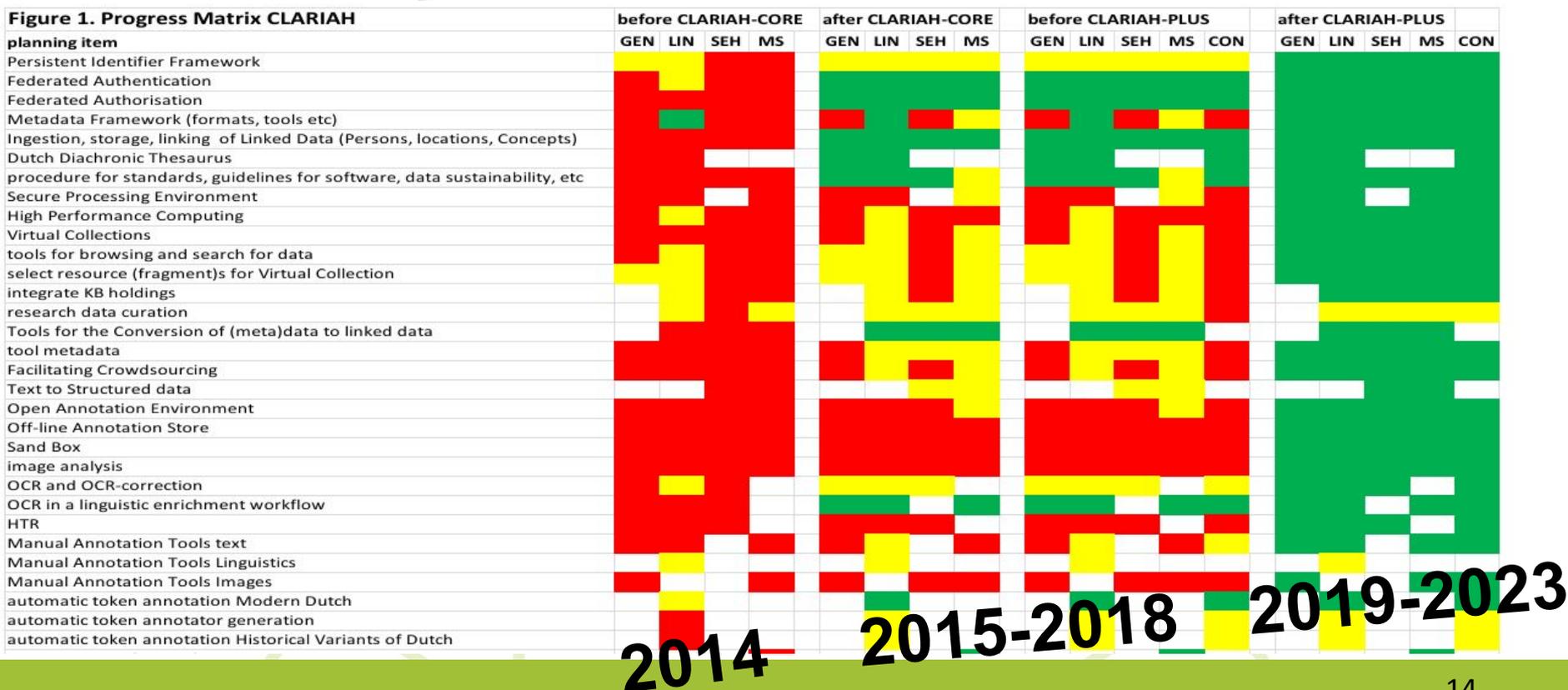
1 General information

GENeric functionality, LINGuistics, Socio-Economic  
History and Media Studies, CONtent of texts: history, literary

Fuente:

<https://www.clariah.nl/over/bestanden/downloads/send/10-folders/166-clariah-plus>

Figure 1. Progress Matrix CLARIAH



# Personas trabajando en la infraestructura CLARIN

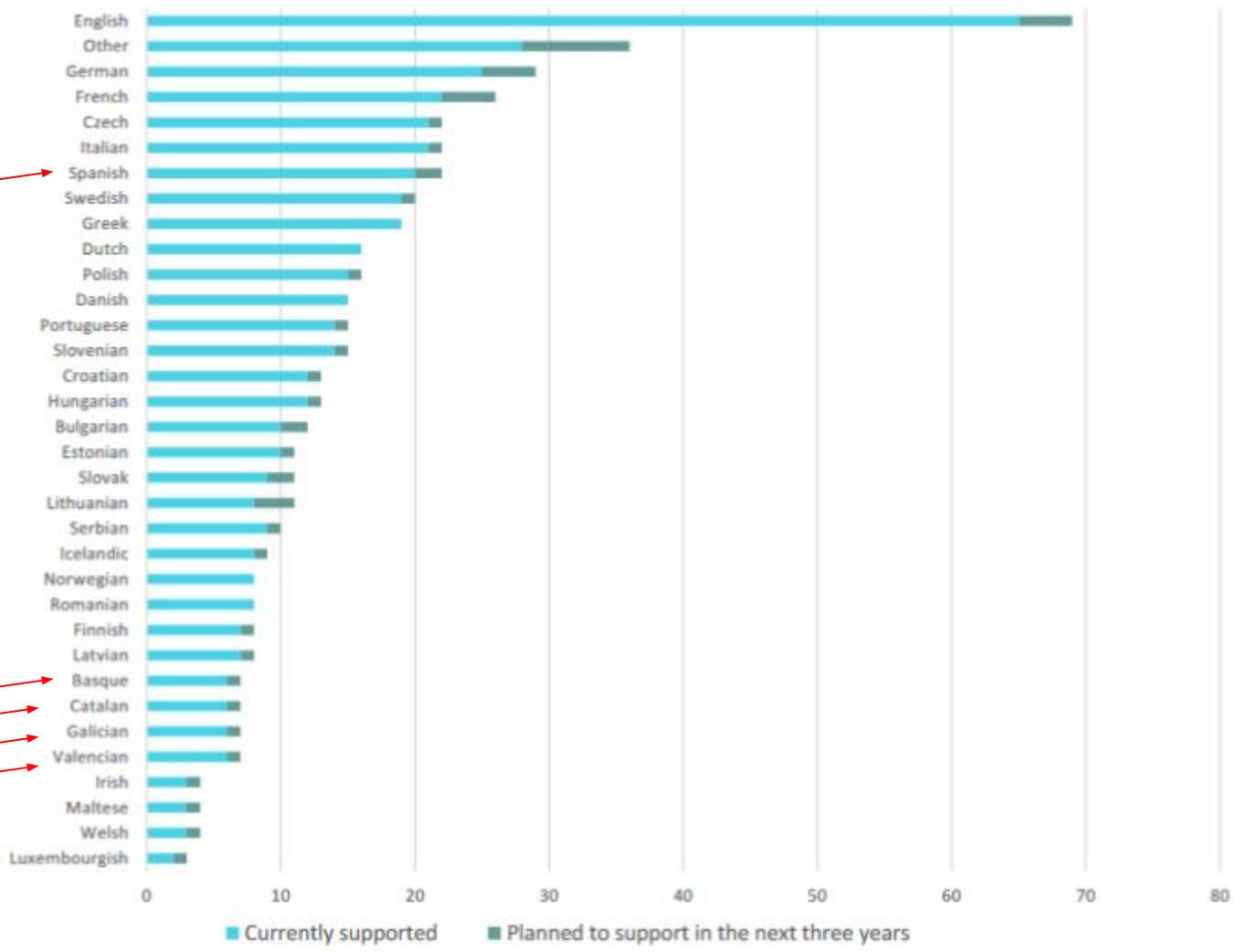


Figure 4: Number of respondents that already work with the language and/or plan to process it in the upcoming three years

# Áreas en las que se está trabajando en CLARIN

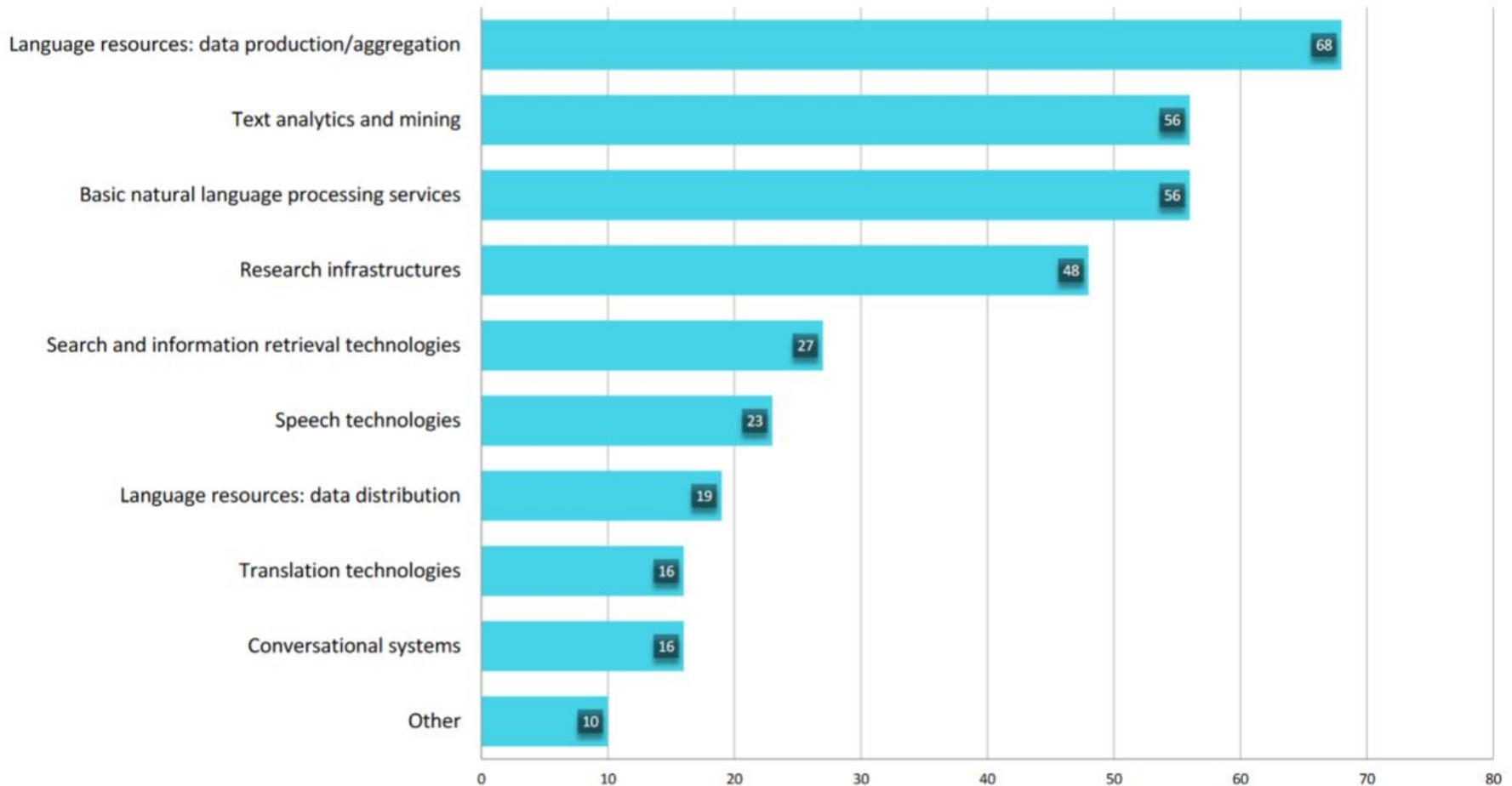


Figure 2: LT areas in which the respondents conduct research or develop tools and services

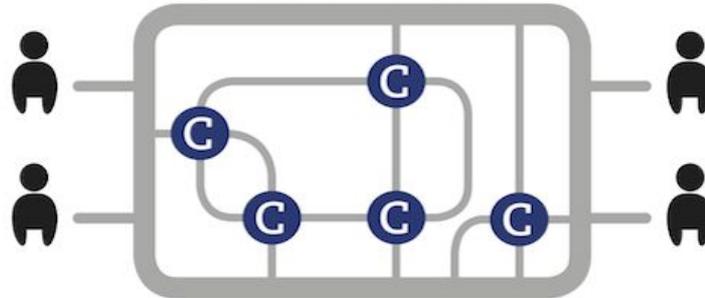


## Recursos para Humanidades y Ciencias Sociales

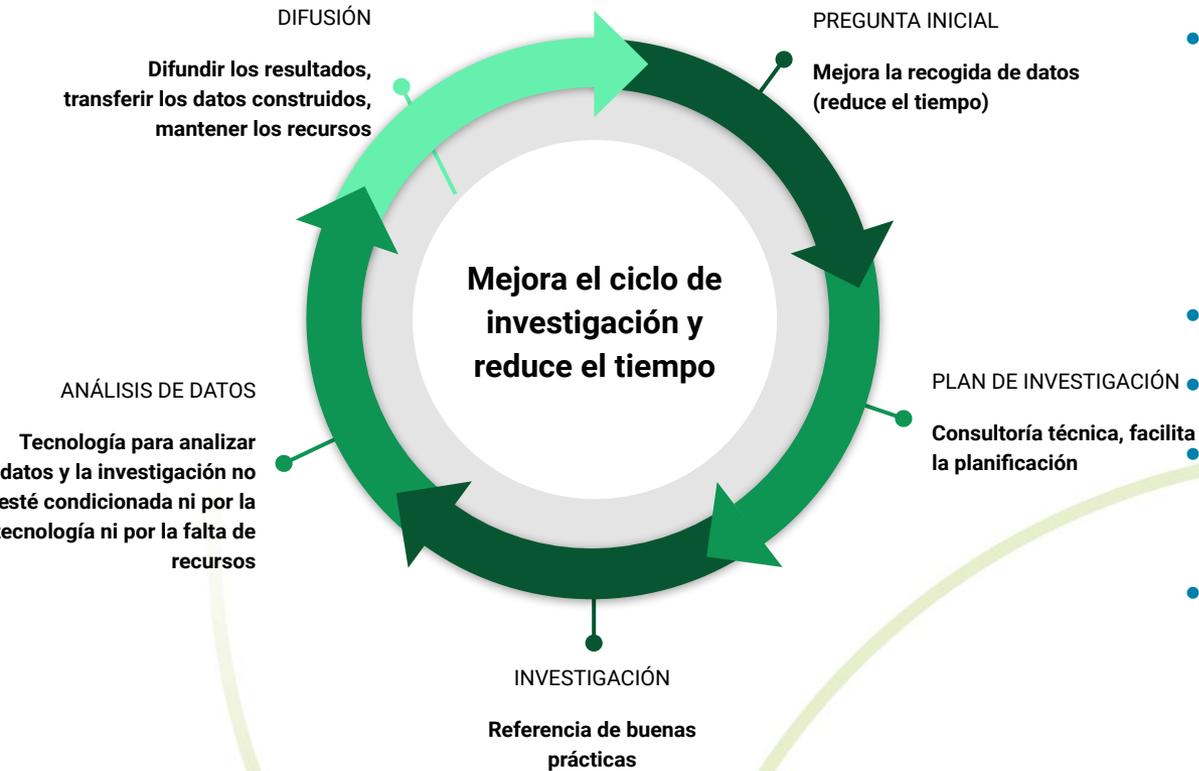
- **Facilitar el uso** de las TL

- Desde una URL
  - Los datos pueden estar en sitios diferentes
- Datos sobre la lengua
  - Texto y vídeo
- Herramientas avanzadas e interoperables
  - buscar, analizar, combinar y crear...

services to researcher



# La infraestructura más que un proyecto



- Acceso confederado a todos los recursos y datos desde una única web
- **Estándares**  
**Protocolos** comunes  
Ayuda para el cambio de paradigma
- Diseño de los recursos **estratégicos**

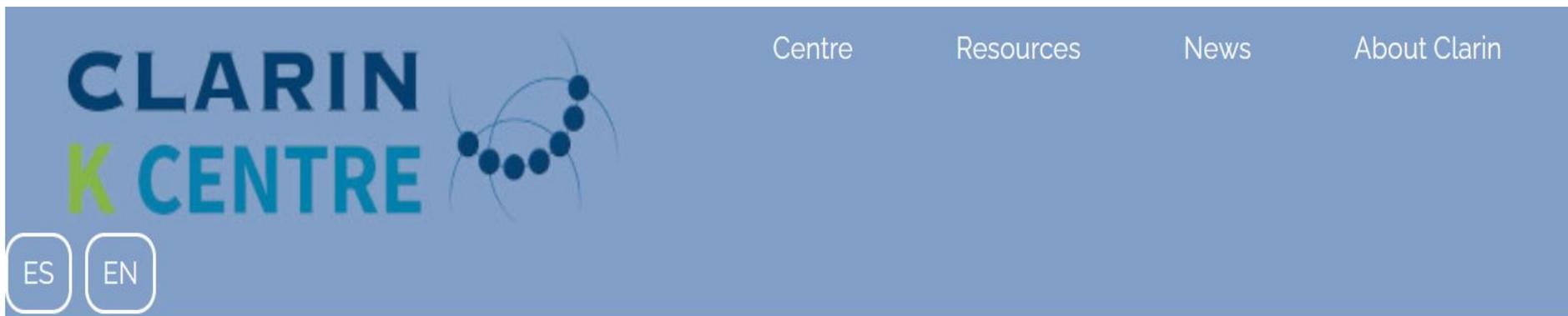
- Corpora: abiertos y **públicos**
- Abiertos solo para la **academia**
- Únicamente para **autorizados**

PUB

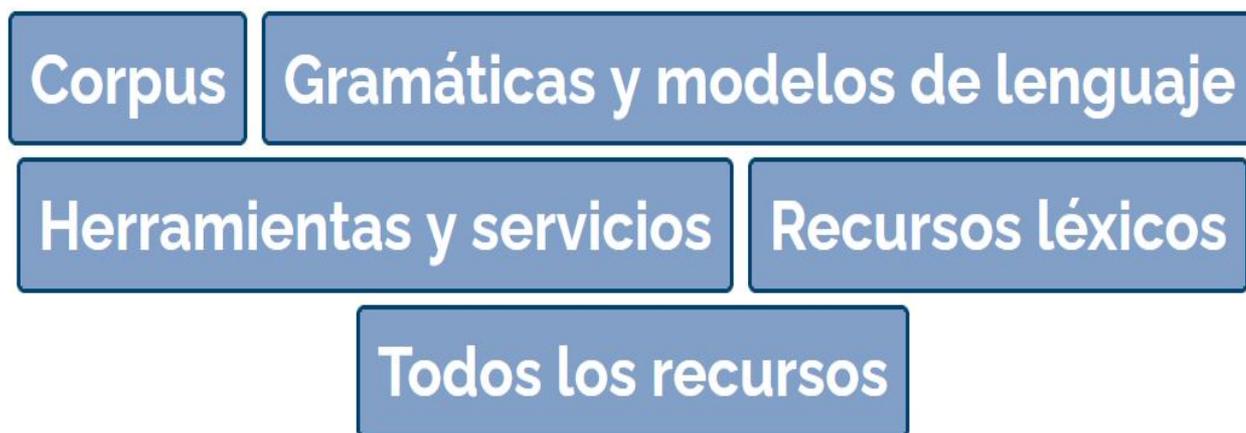
AKA

AUT

# Catálogos de recursos y herramientas



## Recursos y servicios



<https://ixa2.si.ehu.eus/clarin-es/recursos>

# Posibilidades del usuario en infraestructuras

1. Descubrir datos en
  - Repositorios de datos
2. Archivos avanzados de datos
  - Archivos del grupo de investigación o infraestructura
  - Archivos personales
3. Análisis de datos y procesamiento del lenguaje

## Relación usuario/comunidad con infraestructura

4. Se ofrecen los servicios a la comunidad
5. Se evalúa y se ajusta la tecnología
6. Se piden servicios y casos de uso a la comunidad

# Resource families: facilitando el acceso a los datos resumiendo por tipo de datos

## Corpora

- Computer-mediated communication corpora
- Corpora of academic texts
- Historical corpora
- L2 learner corpora
- Legal corpora
- Literary corpora
- Manually annotated corpora
- Multimodal corpora
- Newspaper corpora
- Parallel corpora
- Parliamentary corpora
- Reference corpora
- Spoken corpora

## Lexical Resources

- Lexica
- Dictionaries
- Conceptual Resources
- Glossaries
- Wordlists

## Tools

- Normalization
- Named entity recognition
- Part-of-speech tagging and lemmatization
- Tools for sentiment analysis

# CLARIN: experiencias de usuarios

<https://zenodo.org/record/4288980#.X9YDS7N7mDI>

## CLARIN through the eyes of the researchers

**Tour de CLARIN**   
Volume III



# Servicio ad hoc del centro CLARIN K

Publicado en el Tour de CLARIN

- Tesis en *Basque Center on Cognition, Brain and Language* (BCBL)
- Tema: “My PhD work focuses on the amount of exposure to each language within bilingual contexts, and how it shapes language acquisition at a cognitive and neural level”
  - **Herramientas:**
  - [ANALHITZA](#)
  - <https://switchboard.clarin.eu>



<https://www.clarin.eu/blog/tour-de-clarin-interview-jose-perez-navarro>

# Uso del centro CLARIN-K IMPACT-CKC

## Interview | **Mikel Iruskieta**



Mikel Iruskieta is a computational linguist who is part of the Ixa Research Group and the Didactics of Language and Literature Department at the University of the Basque country. He has collaborated with the CLARIN IMPACT-CKC Knowledge Centre, which helped him and his colleagues digitize Basque texts.

### Could you briefly describe your academic and research background?

<

My current research focuses on the didactics and analysis of Basque, mostly regarding discourse parsing and evaluation of discourse structure. For the last five years, I have mainly worked on adapting language technologies for teaching and learning purposes. With that goal, I have created and now co-lead a postgraduate programme in Basque (University Specialist in ICT and Digital Competences in Education, Continuing Education and Language Teaching), a research group working in Digital Humanities and Education. Our aim is to build a research community that will conduct research and teach in Basque by adopting a critical approach and using language technologies in a pedagogical context. In this postgraduate programme, my colleagues and I are developing a new framework of the socio-tech pedagogy for Basque that will cover the following topics:

<http://clarin-es.org/tour-de-clarin-vol-iii/>

- The Basics of Technology and Pedagogy;
- Formal Education and Technology;
- Continuing Education and Technology;
- Language Teaching and Technology Development;
- Society and Education, Opportunities and Risk of Technology;
- E-learning: Approaches and resources; and
- Digital Research: Methods and resources.

>

### Does the fact that Basque is a language isolate have any bearing on the development of language tools tailored to it?

<

The history and current situation of the Basque language are both complex and interesting. Basque has a relatively small community of speakers (751,700 active and 1,185,500 passive speakers) which lives in contact with three powerful language communities, namely Spanish and French (as official languages in the Basque Country) and English (as a foreign language). It is also not supported enough by official language policies. As a result, Basque is still considered an under-resourced language. In this context, the work of the Ixa Group for NLP is highly valuable. They have developed basic resources for Basque (as well as for other languages) which are used by the research community, for example IXApipes (a modular set of NLP tools which provide easy access to NLP technology for several languages that can be used or exploit its modularity to pick and change different components) and ANALHITZA (a web service to analyse Basque, Spanish and English texts without needing any technical experience). Many more basic and advanced tools and resources for Basque can be found on the website of the HITZ: Basque Center for Language Technology.

>

### How did you get involved with the IMPACT K-Centre and how did they help you with your research?

<

I learned about the IMPACT K-Centre when they joined CLARIN. Because I was working on several different digitization projects for Basque and for Spanish, I immediately got in touch with them and asked for their help. Isabel Martínez Sempere, the manager of IMPACT, helped me solve a digitization issue that I encountered when I was analysing the most frequently occurring words in *Pulgarcito*, which is a Cuban children's magazine written in Spanish from 1919 to 1920. This magazine consists of very diverse materials, such as drawings and handwritten texts, which are

## Analizadores: Language Resource Switchboard

Para encontrar la herramienta adecuada para tu tipo de datos lingüísticos

- <https://switchboard.clarin.eu/>

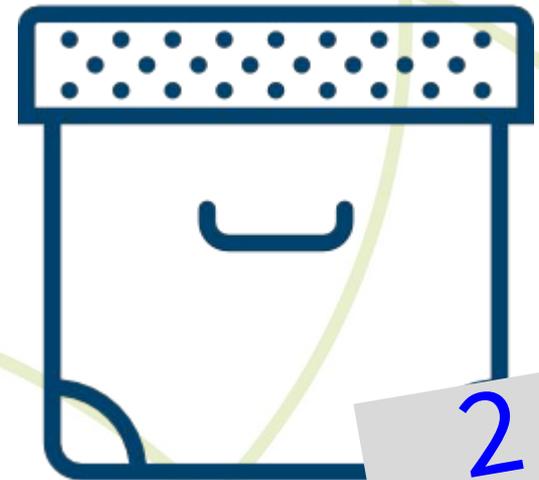


1

## Servicios de depósito

Para depositar y mantener corpora y recursos:

- [www.clarin.eu/content/depositing-services](http://www.clarin.eu/content/depositing-services)



2

## Recursos de lenguaje

Corpus y metadatos: en grandes cantidades y para búsquedas rápidas:

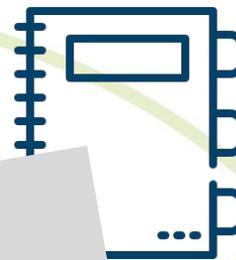
- [vlo.clarin.eu/#tour](http://vlo.clarin.eu/#tour)
- [contentsearch.clarin.eu](http://contentsearch.clarin.eu)
- <https://labur.eus/gZ1ld>



3



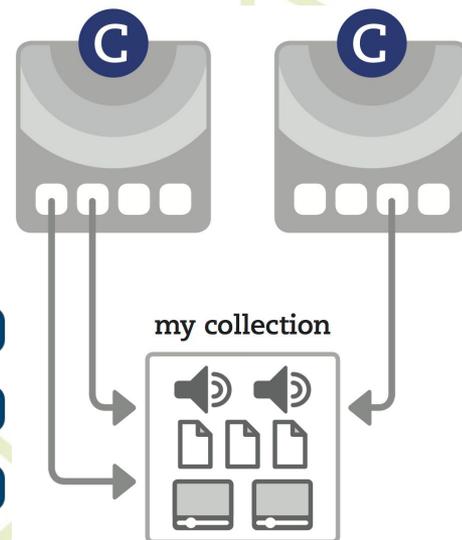
4



## Colecciones virtuales

Para crear corpus virtuales y poder mencionarlos (replicabilidad):

- <https://www.clarin.eu/content/virtual-collections>



# Servicios de depósito, recursos, analizadores y corpus virtual enlazados: CLARIN Virtual collections



## Euskarazko haurren corpusak

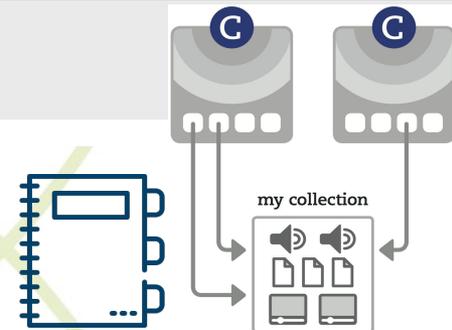
### General

Name: Euskarazko haurren corpusak

### Resources

Reference	Actions
Frogs French Iduguine Corpus	...
Basque SotoValle Corpus - 040505	...
Basque Luque Corpus - 33cas3	...
Haur Hezkuntzako ipuin-bilduma HDL 11304/f27f5e92-af01-4a37-a6d9-82cf14afa160	...

State	Type	Created	
private	extensional	2021-06-29	

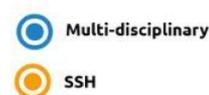


# SSH Open Marketplace y CLARIN

@SSHOpenCloud Objetivo: crear un escenario sostenible y perdurable para compartir y optimizar los datos y servicios en CCSS



#SSHOCaVocabulary  
#SSHOCifyCLARIN



# Obtener y transformar: casos de uso digitalización: herramientas y servicios en CLARIN

¿Qué hacer con textos digitalizados?

- Fondos digitales de bibliotecas
- Tu libro a estudiar digitalizado
- Centro de competencia IMPACT



# IMPACT CLARIN K-centre y BNE

## Vida de Lazarillo de Tormes

 **IMPACT DATASET BROWSER**

This resource is property of

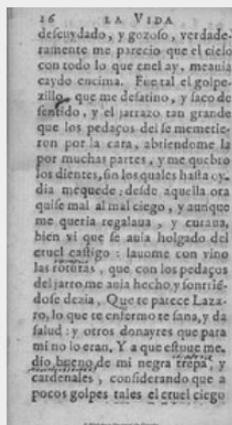


and distributed by the Impact Centre of Competence



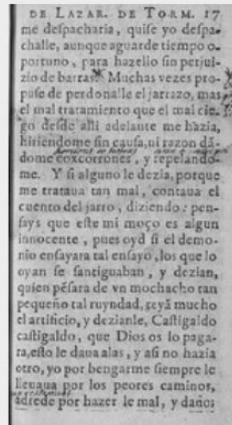
440435

[TIFF](#) [XML](#)



440436

[TIFF](#) [XML](#)



440437

[TIFF](#) [XML](#)



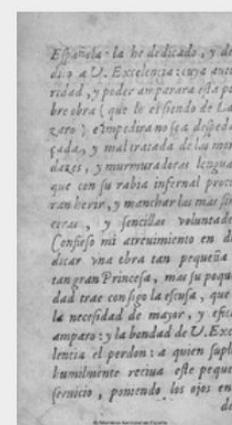
440438

[TIFF](#) [XML](#)



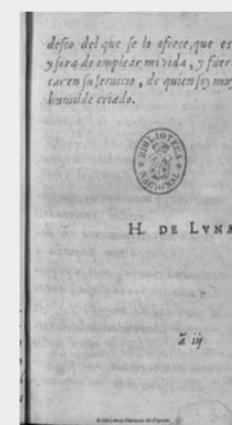
440439

[TIFF](#) [XML](#)



440440

[TIFF](#) [XML](#)



440441

[TIFF](#) [XML](#)

# Libro en castellano

Objetivo: analizar un libro digitalizado con texto escrito a mano y a máquina

- Publicado en [CLARIN](#)
- Descarga: wget [imagenes.sld.cu/download/pulgarcito/volumen-2.pdf](http://imagenes.sld.cu/download/pulgarcito/volumen-2.pdf)

PULGARCITO

VOL. II - NUM. I - ENERO 1920 - 20 CTs.

JUQUEMOS HOY A...



LOS PATINES

Interview at CLARIN:

<https://www.clarin.eu/blog/what-impact-k-centre-can-do-you>

# Comparación entre IMPACT vs Transcribus

<p>CUANDO UN NIÑO          &lt;&lt;SI POEAA?          M\$&amp;ECE,          UMBETRAFO          conminas y ca</p>	<p>c          d          Di          d d ded le          des</p>
<p><b>IMPACT</b></p>	<p><b>Transcribus</b></p>



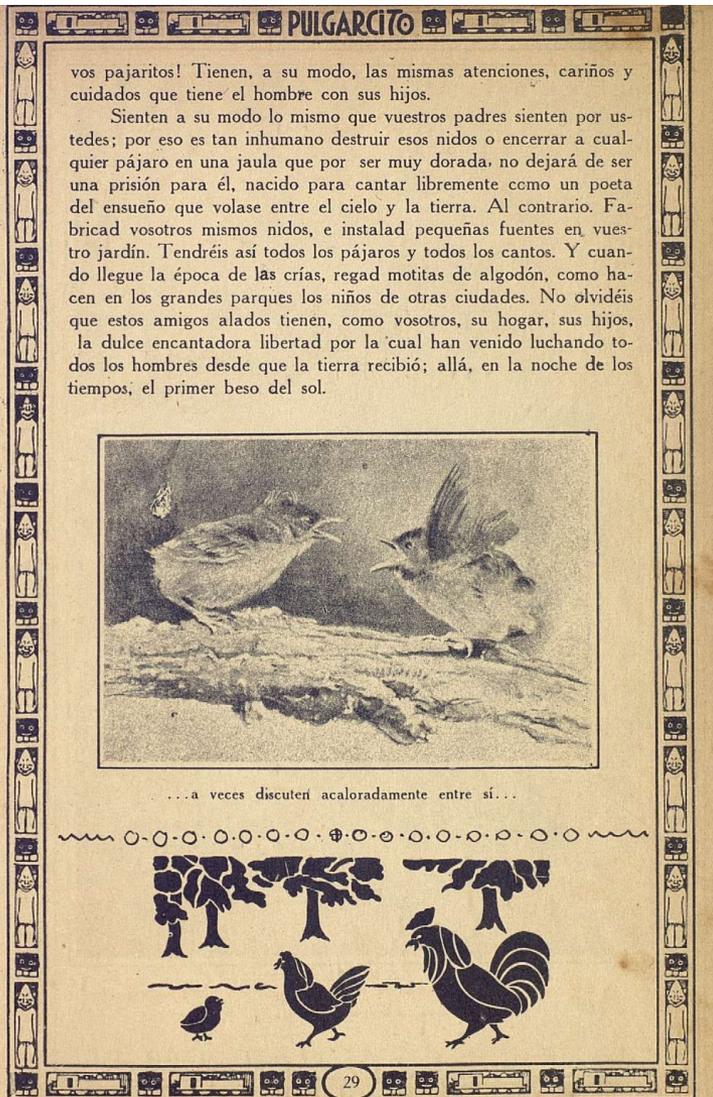
U ab ← → ?! Unclear

---

c  
 d  
 Di  
 d d ded le  
 des

<https://transcribus.eu/lite/collection/89934/doc/617293/detail/5>

# PDF2TXT (IMPACT)



vos pajaritos! Tienen, a su modo, las mismas atenciones, cariños y cuidados que tiene el hombre con sus hijos. Sienten a su modo lo mismo que vuestros padres sienten por ustedes; por eso es tan inhumano destruir esos nidos o encerrar a cualquier pájaro en una jaula que por ser muy dorada, no dejará de ser una prisión para él, nacido para cantar libremente como un poeta del ensueño que volase entre el cielo y la tierra. Al contrario. Fabricad vosotros mismos nidos, e instalad pequeñas fuentes en vuestro jardín. Tendréis así todos los pájaros y todos los cantos. Y cuando llegue la época de las crías, regad motitas de algodón, como hacen en los grandes parques los niños de otras ciudades. No olvidéis que estos amigos alados tienen, como vosotros, su hogar, sus hijos, la dulce encantadora libertad por la cual han venido luchando todos los hombres desde que la tierra recibió; allá, en la noche de los tiempos, el primer beso del sol.

vos pajaritos! Tienen, a su modo, las mismas atenciones, cariños y 4 cuidados que tiene el hombre con sus hijos. Sienten a su modo lo mismo que vuestros padres sienten por ustedes; por eso es tan inhumano destruir esos nidos o encerrar a cualquier pájaro en una jaula que por ser muy dorada, no dejará de ser una prisión para él, nacido para cantar libremente como un poeta del ensueño que volase entre el cielo y la tierra. Al contrario. Fabricad vosotros mismos nidos, e instalad pequeñas fuentes en vuestro jardín. Tendréis así todos los pájaros y todos los cantos. Y cuando llegue la época de las crías, regad motitas de algodón, como hacen en los grandes parques los niños de otras ciudades. No olvidéis que estos amigos alados tienen, como vosotros, su hogar, sus hijos, la dulce encantadora libertad por la cual han venido luchando todos los hombres desde que la tierra recibió; allá, en la noche de los tiempos, el primer beso del sol. ...a veces discuten acaloradamente entre sí...

O-O-O-O'O'O-O - \$-0.0-0.0-0-0 -

# Analizar y presentar: Casos de uso análisis textual en CLARIN

- Interoperabilidad y el análisis textual
- Interoperabilidad y la transcripción de videos
- Otros casos de uso



# Uso de las infraestructuras

## 1. Colección virtual y ejemplo OCR:

- VTL: <https://wlt.pcss.pl/>

## 2. Texto escrito:

- Eudat: FAIR data across borders and disciplines
- Interoperable con el [Switchboard](#) de CLARIN
- Detección de lengua y propone herramientas/métodos
- Analiza el texto y ofrece técnicas de visualización

## 3. De voz a texto:

- Servicio [BAS](#) de CLARIN
- Decenas de lenguas y variaciones
- Múltiples formatos de salida para seguir con la investigación: TXT, SCV, PRAAT, Video...

## 4. Traducción automática:

- Lindat: <https://lindat.mff.cuni.cz/services/translation/>

## 5. Otros ejemplos

# Cobertura por lenguas de la tecnología CLARIN

<https://vlo.clarin.eu> hay más de 1 millón de componentes (herramientas, corpus ...)

- English (154.928)
- German (143.271)
- Dutch (117.312)
- Danish (109.962)
- Slovenian (73.495)
- Polish (40.466)
- French (24.432)
- Afrikaans (7.870)
- ...

- Objetivo y compromiso con
  - FAIR
  - Sustainable Development



ALL LTs in ONE url

Recursos en las lenguas de CLARIN-es

- **Spanish; castilian (14.444)**
- Catalan; valencian (1.364)
- Basque (498)
- Galician (216)

DATA TYPE	
spoken	99558
speech	4455
writing	1543
gestures	36 1307
pointing-gestures	454
facial-expressions	452
emotional-state	451

9 INDUSTRY, INNOVATION AND INFRASTRUCTURE



4 QUALITY EDUCATION



8 DECENT WORK AND ECONOMIC GROWTH



16 PEACE, JUSTICE AND STRONG INSTITUTIONS



11 SUSTAINABLE CITIES AND COMMUNITIES



17 PARTNERSHIPS FOR THE GOALS



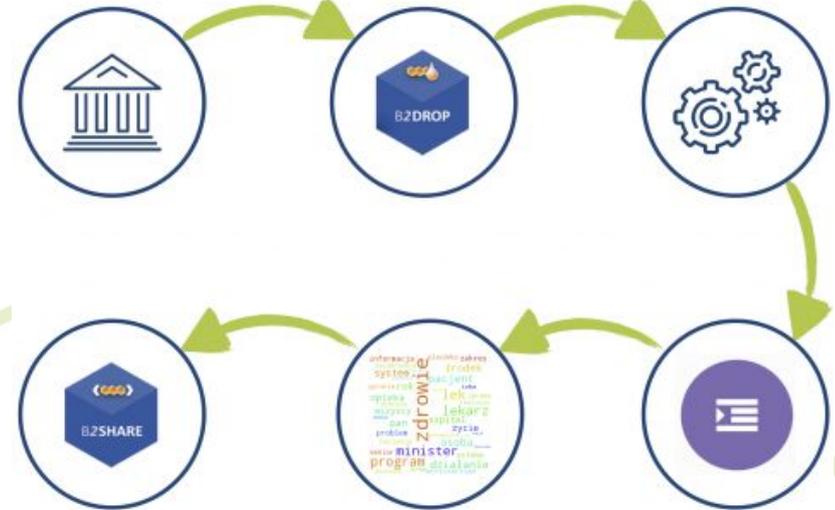
5 GENDER EQUALITY



# Investigar en la nube en CLARIN: principios FAIR

Ejemplo en euskera en la EOSC y EuDAT

- Datos:
  - Corpus de cuentos infantiles
- Análisis sintáctico
- Análisis en la nube
- Publicación persistente



Interoperable



Original en inglés:

<https://www.clarin.eu/showcase/eosc-portal-demonstration>

[https://www.youtube.com/watch?v=YvZ9Y\\_uyr7M](https://www.youtube.com/watch?v=YvZ9Y_uyr7M)

# Colección virtual ad hoc

<https://b2drop.eudat.eu/apps/files/?dir=/UCM%20presentaci%C3%B3n&fileid=27166334>

## Colección virtual Edad de Plata

### General

Name: Colección virtual Edad de Plata  
Type: EXTENSIONAL  
Creation date: 2021-10-08  
Description: Presentación UCM  
Purpose: REFERENCE  
Reproducibility: INTENDED  
Keywords: 

- UCM
- Edad
- de
- Plata

### Resources

Reference

ELTeC

CONSSA

Audio sobre novelas

T1-unamuno

T2-Fernan-Caballero

Taxi-1

Taxi-2

ELTeC (<https://github.com/COST-ELTeC/ELTeC-spa>)

CONSSA (<https://github.com/cligs/conssa>)

Audio sobre novelas (<http://contenidosdigitales.uned.es/fez/view/intecca:VideoCMAV.5a6f2b5bb1111f57648b4cb3>)

Unamuno (<https://b2drop.eudat.eu/s/dLBDMkDBrs3HFT7>)

Caballero (<https://b2drop.eudat.eu/s/Hz5cP4aip5wobDP>)

Save Collection Cancel

# Laboratorio virtual: <https://wlt.pcss.pl/>



**Title:** Seis horas dentro de un taxi

**Author:** Andrés Carranque de Ríos



1

Transcription advancement  
 0%



2

Transcription advancement  
 0%

### Transcribe the page automatically

This function creates an automatic transcription of the selected page

Recognition range

Recognition profile  

- Latin
- Polish
- English
- German
- French
- Latin
- Russian

# Herramientas automáticas y revisión manual

No.	Page name	Number of rows	Verified	Correct	Incorrect	Last edition	Verification
1	 <a href="#">Seis horas dentro de un taxi_Ca rranque_p.1.pdf</a>	87	0	0	0	2021-11-15 10:29	<a href="#">Verify</a>
2	 <a href="#">Seis horas dentro de un taxi_Ca rranque_p.2.pdf</a>	87	0	0	0	mikel.iruskietia@ehu.eus 2021-11-15 10:56	<a href="#">Verify</a>

WLT Transcription object: Seis horas dentro de un taxi

Transcribe automatically Swap panels Correct transcription Download Previous page

Fold Transcribed object Edit selection Fold Transcription



1 .

2 tomóvil se apea de su coche y empieza a gritar. Esto mismo

3 realiza | hombre joven, con aire de estudiante. Le pego con el codo

4 al conduc-

5 mi «jefe», y hasta el grupo de curiosos que comienza a rodearnos.

6 Por tor, y el auto que.la frenado. La dama. nos dice desde el

7 interior:

8 fin se arregla todo buenamente y regresamos al taxi. Entoaces, el

9 – "Vaya hacia la plaza de Santa Ana... No corra.

10 asombro s pinta en. nuestras caras. Resulta que cel señor vencia-

11 Pasada la plaza de Antón Martín noto que el chófer mo hace sc-

12 ble y la jovencita no aparecen por ningún lado. nas de que mire el

13 retrovisor. Observo por el espejito, y «o que la

14 –¡Atiza!—exclama el chófer—. Aquí había algún lío.. El viejo –

15 dama entrega un billete al adolescente, Después se efe t; i1 agia-

16 de los bigotes debe ser un Landrü. BOB - yep ni xi cERdiacct eras

17 SESS \*

# Herramientas de monitorización/verificación

Title: 6 Seis horas dentro de un taxi

Author: Andrés Carranque de Ríos

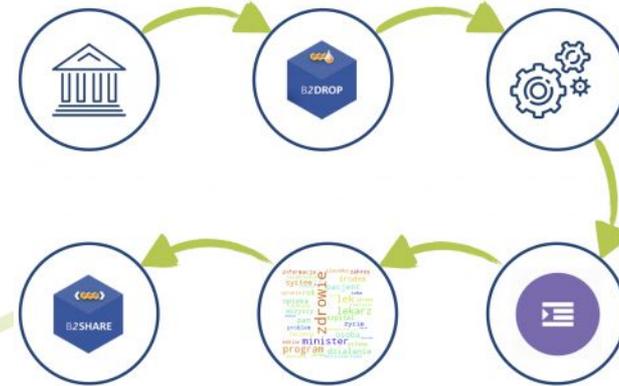
Advanced verification: Pages to verify at this property: 1 Verified rows: 87 / 174 Correct: 83 ✓ Incorrect: 4 ✗

No.	Page name	Number of rows	Verified	Correct	Incorrect	Last edition	Verification
1	 <a href="#">Seis_horas_dentro_de_un_taxi_Carranque_p.1.pdf</a>	87	0	0	0	2021-11-15 10:29	<a href="#">Verify</a>
2	 <a href="#">Seis_horas_dentro_de_un_taxi_Carranque_p.2.pdf</a>	87	87	83	4	<a href="mailto:mikel.iruskiet@ehu.eus">mikel.iruskiet@ehu.eus</a> 2021-11-15 10:59	<a href="#">Verify</a>

# Investigar en la nube en CLARIN: principios FAIR

Ejemplo de compilación de novelas de varias fuentes en la infraestructura de EuDAT

- Datos:
  - “La edad de plata” por Jose Calvo Tello
  - ELTeC por Borja Navarro
- Análisis sintáctico
- Análisis en la nube
- Publicación persistente



Interoperable



# Recursos en la nube para texto



1. Descargar texto: [Gitzip](#)
2. Corpus en [Eudat](#)
3. Analizar con un clic en [Switchboard](#)
4. Elegir un recurso para el análisis de texto
5. ...

- Recursos
  - Para el castellano: 4
  - Para el inglés: 16
  - Para el alemán: 13
  - Para el polaco: 26



1. Constituency Parsing
2. Coreference Resolution
3. Dependency Parsing
4. Distant Reading
5. Extraction of Polish terminology
6. Inclusion detection
7. Keyword Extractor
8. Lemmatization
9. Machine Translation
10. Metadata Processing
11. Morpho-syntactic tagger
12. Morphological Analysis
13. Named Entity Recognition
14. Named Entity Relation Detection
15. Part-Of-Speech Tagging
16. Sentiment Analysis
17. Shallow Parsing
18. Spatial expression detection
19. Speech Recognition
20. Stylometry
21. TF, IDF, TF-IDF calculation
22. Text Analytics
23. Text Enhancement
24. Text Summarization
25. Tokenisation
26. Topic Modelling
27. Visualisation of Geographic Data
28. Word sense disambiguation

# Ejemplo de análisis de TEI+XML del corpus ELTeC



## Resources

SPA2013\_OrtegaYFrias\_ElDuende.xml 1.12 MiB

Mediatype

application/tei+xml

## Matching Tools

▼ Distant Reading



> Open Voyant Tools

Voyant Tools

Cirrus Terms Links Reader TermsBerry Trends Document Terms

**El duende de la corte : edición ELTeC(Ortega y Frias, Ramón (1825-1883))**

El duende de la corte  
o  
Memorias de un fraile  
Novela histórica original  
de

Relative Frequencies

Document Segments (El duende de la corte...)

Summary Documents Phrases Contexts Bubblelines Correlations

This corpus has 1 document with 178,052 total words and 14,073 unique word forms. Created now.  
Vocabulary Density: 0.079  
Average Words Per Sentence: 15.2  
Most frequent words in the corpus: a (4049); no (3542); más (1316); me (1076); mi (909)

Document	Left	Term	Right
1) El du...	fácilmente sin necesidad de pedirlos	a	la imaginación del poeta, cuyas
1) El du...	vez se arranca una lágrima	a	los ojos, un suspiro al
1) El du...	una historia lo que voy	a	referir. No he tenido que
1) El du...	lectura de algunos párrafos convenció	a	mi amigo de que había
1) El du...	de cedérmelo. [1] Así vino	a	mis manos esta historia, y

4,049 context expand

Voyant Tools Stéfan Sinclair & Geoffrey Rockwell (© 2021) Privacy v 2.4 (M55)

LINDAT Repository Corpus Search TreeQuery Treex More Apps About

LINDAT/CLARIN Services / UDPipe

## UDPipe

About Run REST API Documentation

UDPipe is a trainable pipeline for tokenization, tagging, lemmatization and dependency parsing of CoNLL-U files. UDPipe is language-agnostic and can be trained given annotated data in CoNLL-U format. Trained models are provided for nearly all UD treebanks. UDPipe is available as a binary for Linux/Windows/OS X, as a library for C++, Python, Perl, Java, C#, and as a web service. Third-party CRAN packages also exist.

UDPipe is a free software distributed under the Mozilla Public License 2.0 and the linguistic models are free for non-commercial use and distributed under the CC BY-NC-SA license, although for some models the original data used to create the model may impose additional licensing conditions. UDPipe is versioned using Semantic Versioning.

Copyright 2017 by Institute of Formal and Applied Linguistics, Faculty of Mathematics and Physics, Charles University, Czech Republic.

Description of the available methods is available in the API Documentation and the models are described in the UDPipe User's Manual.

Service

The service is freely available for testing. Respect the CC BY-NC-SA licence of the models – **explicit written permission of the authors is required for any commercial exploitation of the system.** If you use the service, you agree that data obtained by us during such use can be used for further improvements of the systems at UFAL. A comments and reactions are welcome.

Model:  UD 2.5 (description)  UD 2.4 (description)  UD 2.0 (description)  UD 1.2 (description)

basque-bdt-ud-2.5-191206

Actions:  Tag and Lemmatize  Parse

1

Title (Plain Text)	Char: UDPipe tokenizer	Char: UDPipe tagger	Char: UDPipe parser
BABARRUN ALE MAGIKOAK	Document Type: CONLL-U conllu.forms conllu.misc Language: Basque	conllu.lemmas conllu.upostags conllu.xpostags conllu.feats	conllu.heads conllu.deprels

Andoni izeneko mutiko bat baserrian bizi zen bere amarekin bakar-bakarrik

Calling UDPipe tagger ...

2

Process Input

Output Text Show Table Show Trees

Save Output File

```
# newdoc
# newpar
# sent_id = 1
# text = BABARRUN ALE MAGIKOAK
1 BABARRUN Babarrun PROP_N 3 nmod SpacesBefore=\\n
2 ALE ale PROP_N 1 flat
3 MAGIKOAK MAGIKOAK PROP_N Case=Erg|Definite=Def|Number=Sing 0 root SpacesAfter=\\n\\n

# newpar
# sent_id = 2
# text = Andoni izeneko mutiko bat baserrian bizi zen bere amarekin bakar-bakarrik.
1 Andoni Andoni PROP_N Case=Dat|Definite=Def|Number=Sing 6 iobj
2 izeneko izen NOUN 3 nmod
3 mutiko mutiko NOUN 6 nsubj
4 bat bat NUM NumType=Card 3 nummod
5 baserrian baseri NOUN Animacy=Inan|Case=Ine|Definite=Def|Number=Sing 6 obl
6 bizi bizi ADJ Case=Abs|Definite=Ind 0 root
7 zen izan AUX Aspect=Prog|Mood=Ind|Number|abs|=Sing|Person|abs|=3 6 aux
8 bere bera DET Case=Gen|Number=Sing 9 nmod
9 amarekin ama NOUN Case=Com|Definite=Def|Number=Sing 6 obl
10 bakar-bakarrik bakar-bakarrik ADV 6 advmod SpaceAfter=No
11 . PUNCT 6 punct
```

3

Process Input

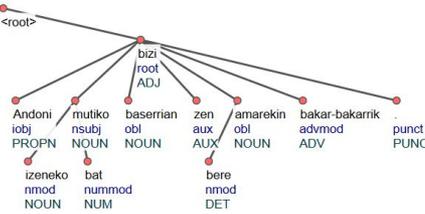
Output Text Show Table

Save Tree as SVG

Previous 1 2 3 4 5 6 7 8 9 10 11 12 ... Next

Andoni izeneko mutiko bat baserrian bizi zen bere amarekin bakar-bakarrik .

4



GO TO EUDAT WEBSITE

SEARCH RECORDS FOR...

RECORDS = C1f2BA03D9584E04AF81E9026B53471E

### Haur Hezkuntzako ipuin-bilduma

by Iruskietia, Mikel;  
Mar 14, 2020

Description: Euskal Herriko ikastolen elkartearen lantzen diren ipuinen bilduma

Disciplines: 1.2.1 → Linguistics → Languages;

Keywords: haur hezkuntza; ipuinak

DOI: 10.23728/bzshare.c1f2ba03d9584e04af81e9026b53471e Copy

PID: 11304/f27f5e92-af01-4a37-a6d9-82cfa4af160 Copy

Name	Size
01-3U-ST-Arrowall-LUZEZA.txt	1,36KB

Basic meta...  
Open Access...  
License...  
Attribution-ShareAlike

6

TUNDRA FileBank\_ee3a1b14-d284-4579-994a-b965c61aabe7

Treebanks Tutorial About Old TUNDRA CLARIN-D

Query

Enter either a TIGERSearch query, or simply a word in quotation marks.

Visualization

Andoni izeneko mutiko bat baserrian bizi zen bere amarekin bakar-bakarrik .

5

# Consecuencias de...

## THE NETWORK IN EUROPE



The CLARIN, DARIAH and CLARIAH research infrastructures are active on a national and the European level. CLARIN and DARIAH both have the status of a European Research Infrastructure Consortium (ERIC).

1. Constituency Parsing
2. Coreference Resolution
3. Dependency Parsing
4. Distant Reading
5. Extraction of Polish terminology
6. Inclusion detection
7. Keyword Extractor
8. Lemmatization
9. Machine Translation
10. Metadata Processing
11. Morpho-syntactic tagger
12. Morphological Analysis
13. Named Entity Recognition
14. Named Entity Relation Detection
15. Part-Of-Speech Tagging
16. Sentiment Analysis
17. Shallow Parsing
18. Spatial expression detection
19. Speech Recognition
20. Stylometry
21. TF, IDF, TF-IDF calculation
22. Text Analytics
23. Text Enhancement
24. Text Summarization
25. Tokenisation
26. Topic Modelling
27. Visualisation of Geographic Data

	Polaco	Alemán	Inglés	Castel
1. Constituency Parsing		x	x	
2. Coreference Resolution	x			
3. Dependency Parsing	x	x	x	x
4. Distant Reading	x	x	x	x
5. Extraction of Polish terminology	x			
6. Inclusion detection	x			
7. Keyword Extractor	x			
8. Lemmatization		x	x	
9. Machine Translation		x	x	
10. Metadata Processing				
11. Morpho-syntactic tagger	x		x	
12. Morphological Analysis	x	x	x	
13. Named Entity Recognition	x	x	x	x
14. Named Entity Relation Detection				
15. Part-Of-Speech Tagging	x	x	x	
16. Sentiment Analysis	x			
17. Shallow Parsing	x			
18. Spatial expression detection	x			
19. Speech Recognition				
20. Stylometry				
21. TF, IDF, TF-IDF calculation	x			
22. Text Analytics	x	x	x	x
23. Text Enhancement			x	
24. Text Summarization	x			
25. Tokenisation				
26. Topic Modelling				
27. Visualisation of Geographic Data				

# Recursos en la nube para la voz

1. Descargar [este vídeo](#) del congreso
2. Analizar con un clic en [BAS](#)
3. Observar los resultados

1. Mary TTS
2. ASR
3. TextAlign
4. Pipeline without ASR
5. Pho2Syl
6. Chunker
7. AnnotConv
8. G2P
9. OCTRA - online text transcription system.
10. AudioEnhance
11. WebMAUS General
12. Chunk Preparation
13. Coala
14. WebMINNI
15. WebMAUS Basic
16. Anonymizer
17. TextEnhance
18. Formant Analysis
19. Subtitle
20. EMU Magic
21. Voice Activity Detection
22. EMU webApp - online labeling of speech data and more.
23. Pipeline with ASR
24. SpeakDiar

Congreso de los Diputados

Sesión Plenaria  
Sesión nº 9  
19/02/2020



✕ Cerrar

## Asuntos

### PREGUNTAS.

► PREGUNTA del Diputado D. PABLO CASADO BLANCO que formula al Excmo. Sr. Presidente del Gobierno: ¿Ha sufrido más de 3 millones de españoles? (Núm.Exp. 180/000026)

Casado Blanco, Pablo (GP)

Sánchez Pérez-Castejón, Pedro (GS) (Presidente del Gobierno)

Casado Blanco, Pablo (GP)

Sánchez Pérez-Castejón, Pedro (GS) (Presidente del Gobierno)

PREGUNTA de la Diputada D<sup>a</sup> INÉS ARRIMADAS GARCÍA que formula al Excmo. Sr. Presidente del Gobierno: ¿Va usted a velar por todos los españoles? (Núm.Exp. 180/000031)

Arrimadas García, Inés (GCs)

Sánchez Pérez-Castejón, Pedro (GS) (Presidente del Gobierno)

Arrimadas García, Inés (GCs)

Sánchez Pérez-Castejón, Pedro (GS) (Presidente del Gobierno)

PREGUNTA de la Diputada D<sup>a</sup> CAYETANA ÁLVAREZ DE TOLEDO que formula al Excmo. Sr. Presidente del Gobierno: ¿Va usted a velar por todos los españoles? (Núm.Exp. 180/000026)

Álvarez de Toledo Peralta-Ramos, Cayetana (GP)

Calvo Poyato, Carmen (GS) (Vicepresidenta Primera del Gobierno)

PREGUNTA del Diputado D. PABLO CASADO BLANCO, del Grupo Parlamentario Popular en el Congreso, que formula al Excmo. Sr. Presidente del Gobierno: ¿Ha renunciado el Gobierno a reducir el desempleo que sufren más de 3 millones de españoles?

► Sesión Completa

► Ver orden del día

# Transcripciones automáticas

Original

(El señor MINISTRO DE INCLUSIÓN, SEGURIDAD SOCIAL Y MIGRACIONES (Escrivá Belmonte): Señora López Álvarez,) en primer lugar, permítame que me dirija hacia la pantalla, aunque le dé la espalda, para que se me pueda oír y responda por el señor Marlaska en términos solidarios, ya que puede ser contestada por el Gobierno. La contestación es clara: no solamente el Gobierno español, yo diría que ningún Gobierno, al menos europeo, en ningún caso fomenta la inmigración irregular, en ningún caso. Lo que hace es intentar evitar que ocurra. (Aplausos).

IBM

abajo en primer lugar el hermitage hija hacia la pantalla y no puedo a la espalda para que se me pueda oír y respondo por él señor smart laska entrenó solidarios cuyo poder con el gobierno la contestación es clara y no solamente el gobierno español ya que cualquier gobierno al menos europeo en ningún caso fomenta la inmigración irregular ningún caso lo que se intenta evitar que ocurra qué

European Media Lab

señaló que cada año en 1º lugar que permita dirija hacia la pantalla y no pudo dar la espalda para que se me pueda unir y respaldado el señor más hasta que en los aviaros como puede ser colobiano la contestación es clara el no solamente el gobierno español y al igual que el Gobierno al menos europeo y en ningún caso fomenta la inmigración irregular en ningún caso lo que hace es intentar evitar que ocurran

# Transcripciones automáticas

Original

(El señor MINISTRO DE INCLUSIÓN, SEGURIDAD SOCIAL Y MIGRACIONES (Escrivá Belmonte): Señora López Álvarez,) en primer lugar, permítame que me dirija hacia la pantalla, aunque le dé la espalda, para que se me pueda oír y responda por el señor Marlaska en términos solidarios, ya que puede ser contestada por el Gobierno. La contestación es clara: no solamente el Gobierno español, yo diría que ningún Gobierno, al menos europeo, en ningún caso fomenta la inmigración irregular, en ningún caso. Lo que hace es intentar evitar que ocurra. (Aplausos).

PRAAT textgrid



File type = "ooTextFile"

Object class = "TextGrid"

xmin = 0

xmax = 26.302000

tiers? <exists>

(...)

intervals [12]:

xmin = 3.472000

xmax = 3.982000

text = "permita"

intervals [13]:

xmin = 3.982000

xmax = 4.542000

text = "dirija"

intervals [14]:

xmin = 4.542000

xmax = 4.822000

text = "hacia"

<https://youtu.be/7II-gOShtFA>

# Transcripciones bilingües euskera-castellano



EU ES



Pre

Albisteak eta ekitaldiak · Ekitaldiak eta gertaerak

2020 ira 14

## EUSKO LEGEBILTZARRAREN 40. URTEURRE

*Eusko Legebiltzarrak Euskal Herriko Unibertsitatearen (UPV-EHU) udako ik*

Lekua MIRAMAR JAUREGIA

Datak: leh, 26/10/2020 - art, 27/10/2020

Ordua: 10:00 -18:00



urteurrena  
aniversario  
1980 - 2020

**EUSKO LEGEBILTZARRAREN 40. URTEURRENA:  
ATZERANZKO BEGIRADA**

**40 ANIVERSARIO DEL PARLAMENTO VASCO:  
UNA MIRADA RETROSPECTIVA**



FUENTE:

<https://www.legebiltzarra.eus/portal/eu/web/eusko-legebiltzarra/noticias-y-eventos/actos-y-eventos/-/buscador/content/40-aniversario-del-parlamento-vasco-una-mirada-retrospectiva>

callGoogleASR: egun on guztioi eta ongi etorri abestia eusko legebiltzarrak euskal herriko unibertsitatearen udako ikastaroen baitan antolatu duen 2000 goiko ikastaro honetara eskerrak eman nahi dizkizuet jardunaldi hauetan parte hartu duzuen guztioi hizlari partehartzaile antolatzaileei ere gehiago covid-19 da gure bizitzak etengabe baldintzatzen dituen une honetan ikastaro hau horren lekuko eusko legebiltzarraren 40. urteurrena atzerako begirada da ikasturte honetarako aukeratutako gaia ezin ziteken besterik izan izan ere aurten 40 (...) izan gara eta legebiltzarrak horretan paper garrantzitsua izan du **en estos dos legislaturas el parlamento vasco se ha ido construyendo y consolidando dia a dia del mismo modo que este pueblo nuestro pueblo se ha ido reconstruyendo la trayectoria de la camara ha sido y es fiel reflejo de la evolucion social la presencia de la mujer (...)** eta konpromisoz aurre egiteko zuen ekarpenak helburu horretan lagunduko dugula sinetsita berriz ere eskerrak eman nahi dizkizuet guztioi

[https://clarin.phonetik.uni-muenchen.de/BASWebServices/data/2021.05.28\\_00.47.01\\_234C2EA21FEC82BA69904D1D91E42A41/Eusko-legebiltzarra.txt](https://clarin.phonetik.uni-muenchen.de/BASWebServices/data/2021.05.28_00.47.01_234C2EA21FEC82BA69904D1D91E42A41/Eusko-legebiltzarra.txt)

# OH Portal (CLARIN)

<https://clarin.phonetik.uni-muenche.n.de/apps/TranscriptionPortal/>

0 1 0 0 0 0

Help Statistics Feedback

OCTRA: plain.par , Language: eng-GB , Audio duration: 00:58



OCTRA v1.4.3 (url) Dictaphone Editor Linear Editor 2D-Editor

TRN Werkzeuge Exportieren DE

TASTENKOMBINATIONEN [ALT + 8]

ÜBERSICHT [ALT + 0]

HILFE

Audio player interface with waveform and transcription text. The transcription text is: "however the list number one the birch ...", "...canoes slid on the smooth planks do the ...", and "...sheet to the dark blue background it's e...".

0 1 0 0 0 0 Help Statistics Feedback

Harvard2.TextGrid , Language: eng-GB , Audio duration: 00:58

TextGrid interface showing waveform and segmentation data. The interface includes a search bar and a list of segments with their start and end times.

Webinar: <https://youtu.be/X6bFGJpMjVQ>

# Curso: Oral Archives for Sociolinguistic Research

Objetivo: curso de sociolingüística para mostrar las posibilidades y los desafíos de los archivos de historia oral

Herramientas:

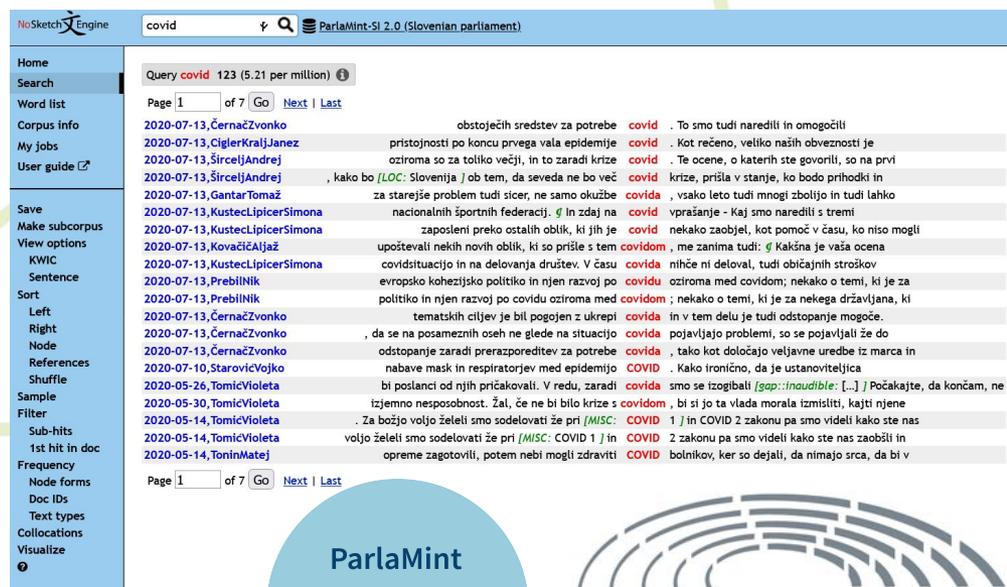
1. Orthographic transcripts with Octra
2. ASR
3. Chunker
4. G2P
5. MAUS segmentation
6. Pho2Syl

**Servicio CLARIN:  
BASWebServices**

# ParlaMint corpus multilingüe (texto escrito)

1. Obtener datos de parlamentos y sus metadatos
2. Convertirlos al esquema de ParlaMint
3. Anotación lingüística (UDpipe y NERC)
4. Hacer corpus disponibles a través de concordantes (noSketch Engine / KonText) y Parlameter

<https://www.clarin.si/noske/parlamint.cgi/>



The screenshot displays the ParlaMint search interface. The search bar contains the query 'covid' and the corpus is identified as 'ParlaMint-SI 2.0 (Slovenian parliament)'. The search results show a frequency of 123 (5.21 per million) for the query. The interface includes a sidebar with navigation options like Home, Search, Word list, Corpus info, My jobs, and User guide. The main content area lists search results with dates, names, and snippets of text. The results are paginated, showing page 1 of 7.

Date	Name	Snippet	Frequency
2020-07-13	ČrnačZvonko	obstoječih sredstev za potrebe	covid
2020-07-13	ČiglerKrajJanez	pristojnosti po koncu prvega vala epidemije	covid
2020-07-13	ŠirčeljAndrej	oziroma so za toliko večji, in to zaradi krize	covid
2020-07-13	ŠirčeljAndrej	, kako bo [LOC: Slovenija ] ob tem, da seveda ne bo več	covid
2020-07-13	GantarTomaž	za starejše problem tudi sicer, ne samo okužbe	covida
2020-07-13	KustecLipicerSimona	nacionalnih športnih federacij. ¶ In zdaj na	covid
2020-07-13	KustecLipicerSimona	zaposleni preko ostalih oblik, ki jih je	covida
2020-07-13	KovačičAljaž	upoštevali nekih novih oblik, ki so prišle s tem	covidom
2020-07-13	KustecLipicerSimona	covidstuaucioj in na delovanja društev. V času	covida
2020-07-13	PrebilNIK	evropsko kohezijsko politiko in njen razvoj po	covidu
2020-07-13	PrebilNIK	politiko in njen razvoj po covidu oziroma med	covidom
2020-07-13	ČrnačZvonko	tematskih ciljev je bil pogojen z ukrepi	covida
2020-07-13	ČrnačZvonko	, da se na posameznih oseh ne glede na situacijo	covida
2020-07-13	ČrnačZvonko	odstopanje zaradi prerazporeditev za potrebe	covida
2020-07-10	StarovičVojko	nabave mask in respiratorjev med epidemijo	covida
2020-05-26	TomičVioleta	bi poslanci od njih pričakovali. V redu, zaradi	covida
2020-05-30	TomičVioleta	izjemno nesposobnost. Žal, če ne bi bilo krize s	covidom
2020-05-14	TomičVioleta	Za božjo voljo želeli smo sodelovati že pri [MISC: COVID	COVID
2020-05-14	TomičVioleta	voljo želeli smo sodelovati že pri [MISC: COVID 1 ] in	COVID
2020-05-14	TorinMatej	opreme zagotovili, potem nebi mogli zdraviti	COVID



# Preguntas de investigación

1. ¿Hay alguna diferencia en los discursos en uno u otro género?
2. ¿Es determinante la edad para diferenciar la forma en el que hablan?
3. ¿Influye la experiencia en la manera de hablar?
4. ¿Hay diferencias entre la izquierda/centro/derecha en los discursos?
5. ¿Se puede diferenciar el gobierno o la oposición por la manera de hablar?
6. ¿Hay cambios en la temática del parlamento, cuando la derecha influye en la agenda mediática?
7. ¿Habla de la misma manera el partido de los verdes en los diferentes parlamentos o depende de la correlación de fuerzas entre la izquierda y la derecha?
8. ¿Cuál es el texto más complejo la traducción o el original?
9. ¿Utilizan los políticos una u otra lengua dependiendo del tema?

# Materiales para el uso práctico del corpus

- 1 Introduction
- 2 Instructions for use
- 3 Corpora and concordancers
  - 3.1 Corpora
  - 3.2 Concordancers
- 4 Parliamentary records
  - 4.1 Parliamentary discourse
  - 4.2 Faithfulness of the records
  - 4.3 Know your research dataset
- 5 Language and gender
- 6 Corpus analysis
  - 6.1 The siParl 2.0 corpus
  - 6.2 TASK 1: Representation of women in the Slovenian Parliament
    - 6.2.1 Creating subcorpora
    - 6.2.2 Using frequency lists
    - 6.2.3 Comparative analysis
  - 6.3 TASK 2: Issues addressed by women
    - 6.3.1 Extracting keywords
    - 6.3.2 Analysing concordances
    - 6.3.3 Comparative analysis
  - 6.4 TASK 3: Topics related to women
    - 6.4.1 Working with frequencies
    - 6.4.2 Extracting collocations
    - 6.4.3 Comparative analysis

## Voices of the Parliament A Corpus Approach to Parliamentary Discourse Research

»Prvič, sem političarka in  
ne politik, drugič pa ...«

Korpusni pristop  
k raziskovanju  
parlamentarnega  
diskurza



<https://sidih.github.io/voices/toc.html>

# Casos de uso de impacto

- Corpus: Norwegian Newspaper Corpus (en CLARIN Resource Family 'Newspaper Corpora')
- Herramienta: **Corpuscle** corpus management and search system, desarrollado en CLARINO Bergen Centre

De Smedt says: 'The useful thing about Corpuscle is that it has powerful search possibilities and that it allows the download of the data in a very usable format. Those are the advances of Corpuscle.'

***'This is the first study to demonstrate the effect of such a spelling change in various Norwegian media sources.'***

Koenraad De Smedt

# Traducción automática

<https://lindat.mff.cuni.cz/services/translation/>

Preguntas a la historia que fue la Edad de Plata de la cultura española, una época de gran auge y esplendor, de la cultura española es la desde los tiempos de la revolución llamada gloriosa desde 1868 hasta el intento de golpe de Estado y la Guerra Civil de 1936-39 da de plata con las décadas posteriores a 1902 y hasta el mismo 1936, pero esto parece constituir un desenfoque histórico pues a nombres fundamentales como los de Giner de los Ríos Benito Pérez Galdós Clarín Santiago Ramón y Cajal Marcelino Menéndez Pelayo Manuel Bartolomé Cossío etcétera. públicas el sufragio universal constituyeron logros de la época aunque a la vez no hay que olvidar que fueron asimismo años de hambre y miseria de violencia en la vida diaria de luchas de clases muy sombrías Francisco Abad profesor de historia de la lengua española UNED Radio 5

**Czech-Ukrainian Translation**  
The Czech-Ukrainian and Ukrainian-Czech translation is available at <https://lindat.cz/translation/>

The translation service is available for *personal and non-commercial use* (see terms of use for more details).

**Source**  
English

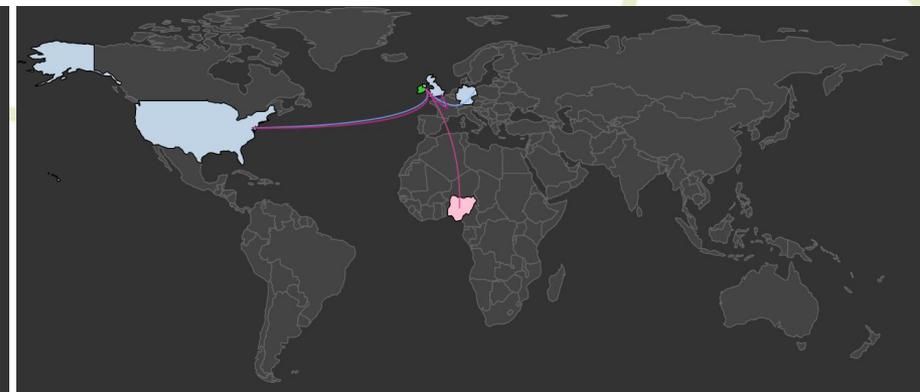
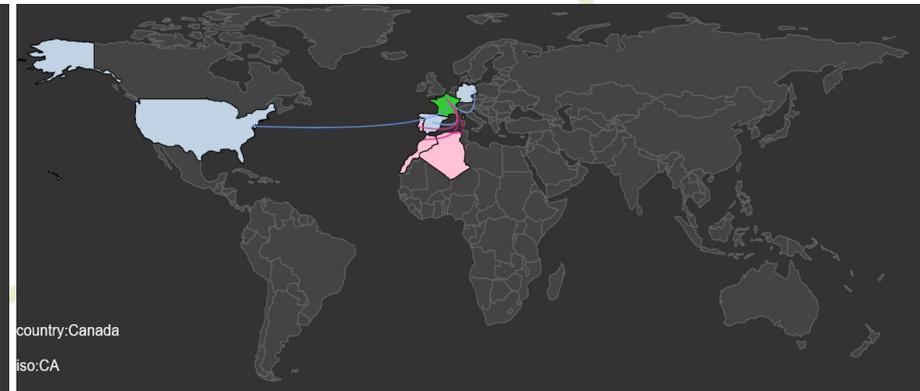
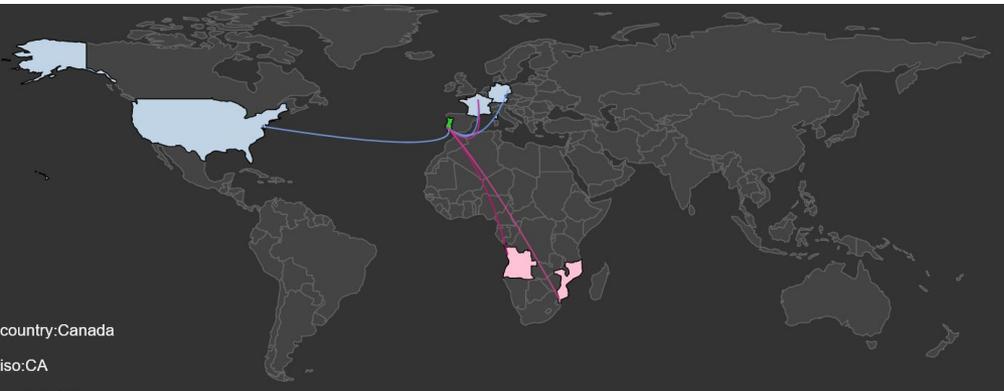
advanced

**Input sentences**

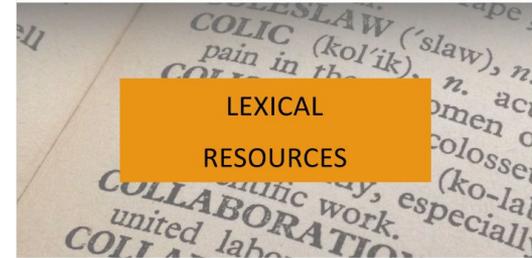
**Target**  
Czech  
Czech  
French  
German  
Hindi  
Polish  
Russian

# Textual Emigration Analysis

1. Historia, literatura y lingüística computacional.
2. Corpus wikipedia (texto)



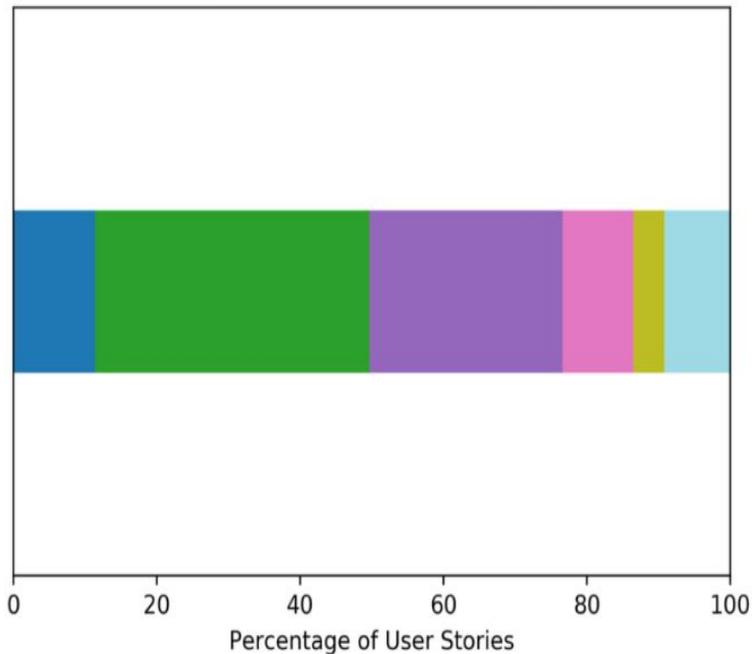
# Text+ User Stories (definir necesidades) DARIAH



## Text+ User Stories sorted by DFG Subject Areas

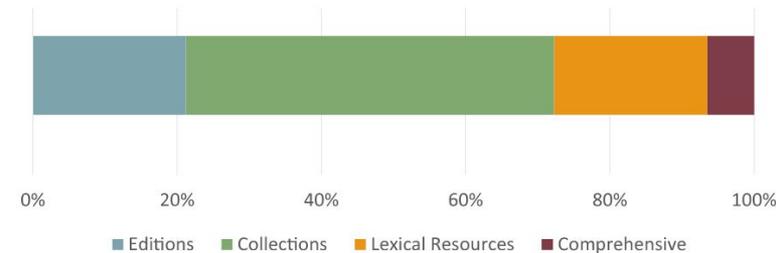
Information about the DFG subject areas can be found on the [webpages of the DFG](#).

Percentage of User Stories per Subject Area



Subject Areas

- Ancient Cultures
- Linguistics
- Literary Studies
- Social and Cultural Anthropology, Religious Studies, etc.
- Philosophy
- Other areas



Distribution bases on 120 user stories on 24 August 2020.

# Ejemplos prácticos de INTELE



**"Facilitando el acceso computacional a colecciones digitales"**

7 de junio de 2021  
16:30h CET  
(Online)

María Dolores Sáez  
Gustavo Candela  
María Pilar Escobar

Formulario de inscripción: [http://ixa2.si.ehu.es/intele/form\\_enlace\\_webinar](http://ixa2.si.ehu.es/intele/form_enlace_webinar)

**"Text Analysis for Spanish"**

17 de mayo de 2021  
14:30h CET  
(Online)

Quinn Dombrowski

**"Distant Reading for European Literary History"**

11 de junio de 2021  
14:00h CET  
(Online)

Rosario Arias  
Borja Navarro  
Christof Schöch

**"Programming Historian: Un proyecto colaborativo para poner la programación al alcance de los humanistas"**

25 de marzo de 2021  
14:30h CET  
(Online)

Jennifer Isasi  
Riva Quiroga

Formulario de inscripción: [http://ixa2.si.ehu.es/intele/form\\_enlace\\_webinar](http://ixa2.si.ehu.es/intele/form_enlace_webinar)

"Facilitando el acceso computacional a colecciones digitales". (Biblioteca Virtual Miguel de Cervantes)

Sesión práctica: [github.com/hibernator11/notebook-ph](https://github.com/hibernator11/notebook-ph)

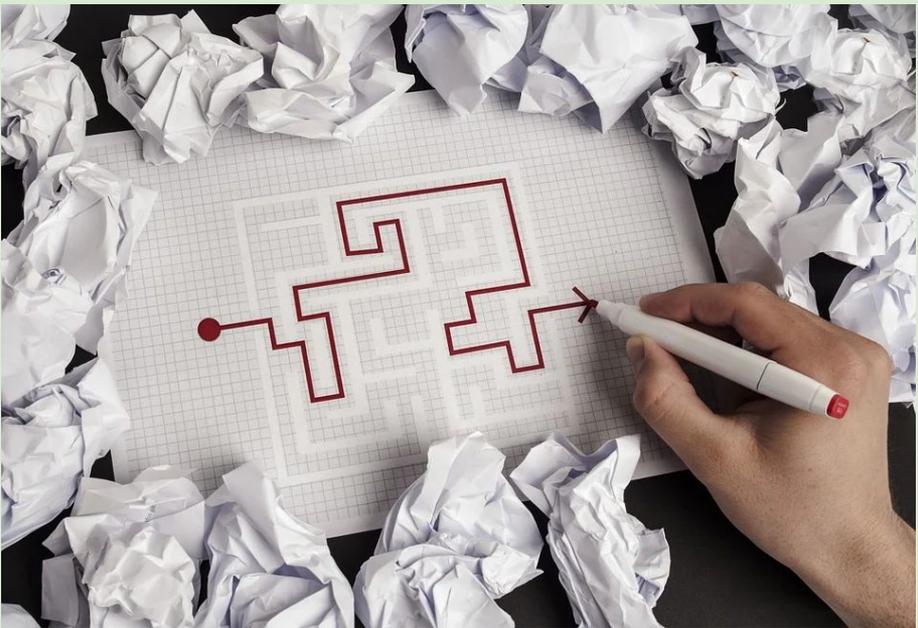
"Análisis de textos para el español"  
<https://github.com/quinnanya/intro-to-nlp-es>

"Distant Reading for European Literary History"  
Sesión práctica: [github.com/bncolorado/Processing-ELTeC-corpus](https://github.com/bncolorado/Processing-ELTeC-corpus)

Programming Historian:  
<https://programmingshistorian.org/es/lecciones/analisis-de-sentimientos-r>  
Parte práctica: <https://rstudio.cloud/project/2342606>

# Mapa

Conclusiones

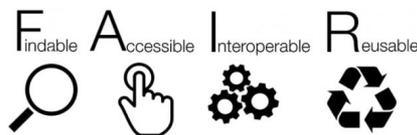


# Conclusiones de la sesión

- **Ciencia abierta: de impacto y reproducible (CLARIN)**
- Interoperabilidad en (y entre) infraestructuras (EOSC)
  - Construir infraestructura que no se desarrollará en Europa para
    - Lenguas oficiales y cooficiales del estado
  - Evitar la fragmentación: ALL-LT-in-ONE-URL
  - + re-uso de los datos + prosumidores + impacto social
- + Casos de uso y herramientas (sencillas) para investigar



**EUROPEAN OPEN  
SCIENCE CLOUD**



# Referencias y enlaces de interés

- Bel, N. Gonzalez-Blanco, E. Irukieta, M. (2016). [CLARIN Centro-K-español](#). *Procesamiento del Lenguaje Natural* 57: 151-154. ISSN: 1135-5948.
- Irukieta, M. Bel, N. (2017). [CLARIN-K Centre Spain: una infraestructura orientada usuario](#). LINHD-UNED. Escuela de Verano HD.
- Krauwer, S., & Hinrichs, E. (2014). The CLARIN research infrastructure: resources and tools for e-humanities scholars. In *Proceedings of the Ninth International Conference on Language Resources and Evaluation (LREC-2014)* (pp. 1525-1531). European Language Resources Association (ELRA).
- Váradi, T., Wittenburg, P., Krauwer, S., Wynne, M., & Koskenniemi, K. (2008). CLARIN: Common language resources and technology infrastructure. In *6th International Conference on Language Resources and Evaluation (LREC 2008)*.
- Otegi, A. Imaz, O. Díaz de Ilarraza, A. Irukieta, M. Uria, L. (2017). [ANALHITZA: a tool to extract linguistic information from large corpora in Humanities research](#). *Procesamiento del Lenguaje Natural* 58, pp. 77-84.
- CLARIN: <https://www.clarin.eu/>
- DARIAH: <https://www.dariah.eu/>
- INTELE: <http://ixa2.si.ehu.eus/intele/>

# Eskerrik asko, gracias

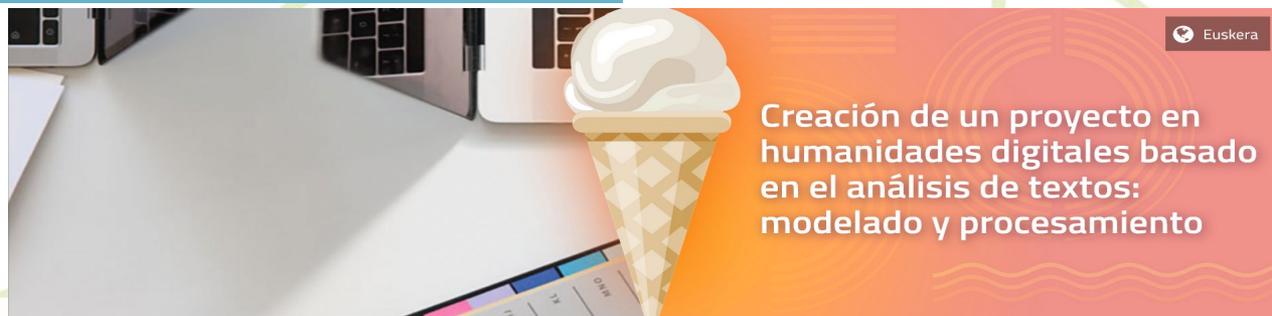
- Preguntas??

# Herramientas Digitales para las Humanidades Digitales en la e-infraestructura CLARIN

Mikel Iruskietea  
HiTZ zentroa-Ixa Taldea  
UPV/EHU

[www.clarin.eu](http://www.clarin.eu)  
[www.clarin-es.org](http://www.clarin-es.org)  
<http://ixa2.si.ehu.es/clarink>

Creación de un Proyecto en  
Humanidades Digitales basado  
en el análisis de textos:  
Modelado y Procesamiento



Creación de un proyecto en  
humanidades digitales basado  
en el análisis de textos:  
modelado y procesamiento