# Secondments@Arquivo.pt and *new* research tools available

[daniel.gomes@fccn.pt](mailto:daniel.gomes@fccn.pt)

# Who are we?

**Free online** service to research the Past Web

Preserves **publicly accessible** information related with:

- Portugal
- **Research** and **Education** (international)

Governmental service provided by
the Ministry for Science and Technology (Portugal)

A digital research infrastructure

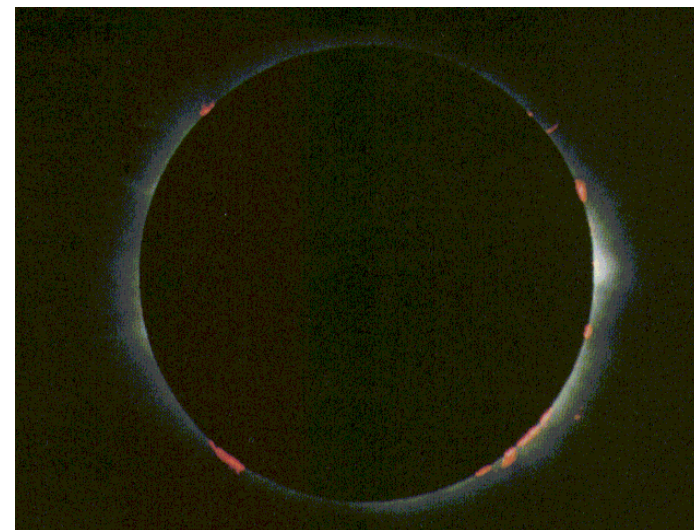# Arquivo.pt preserves national and international historical web content



**nautilus.fis.uc.pt**- **1993**
(oldest page)



**spacelink.nasa.gov – 1992**
(oldest image)

# CLEOPATRA Secondments

- Diego Alves (ESR14), "Robustness of Corpus-based Typological Strategies for Dependency Parsing", August 2022

- Swati (ESR12), "Cross-lingual news reporting bias", December 2022

➢ Get feedback from researchers with real questions about our research tools to improve them

# Robustness of Corpus-based Typological Strategies for Dependency Parsing

ESR14 Secondment

From 06/2022 to 08/2022

Diego Alves
Daniel Gomes

# Objectives

- Creation of datasets from texts extracted from the EAWP23 collection of the Arquivo.pt:
  - 2019 European Parliamentary Elections collection
  - 24 European Union official languages
- Conduct corpus-based typological analysis with syntactic information extracted from this dataset and compare to the state-of-the-art (i.e.: parallel corpora - PUD)
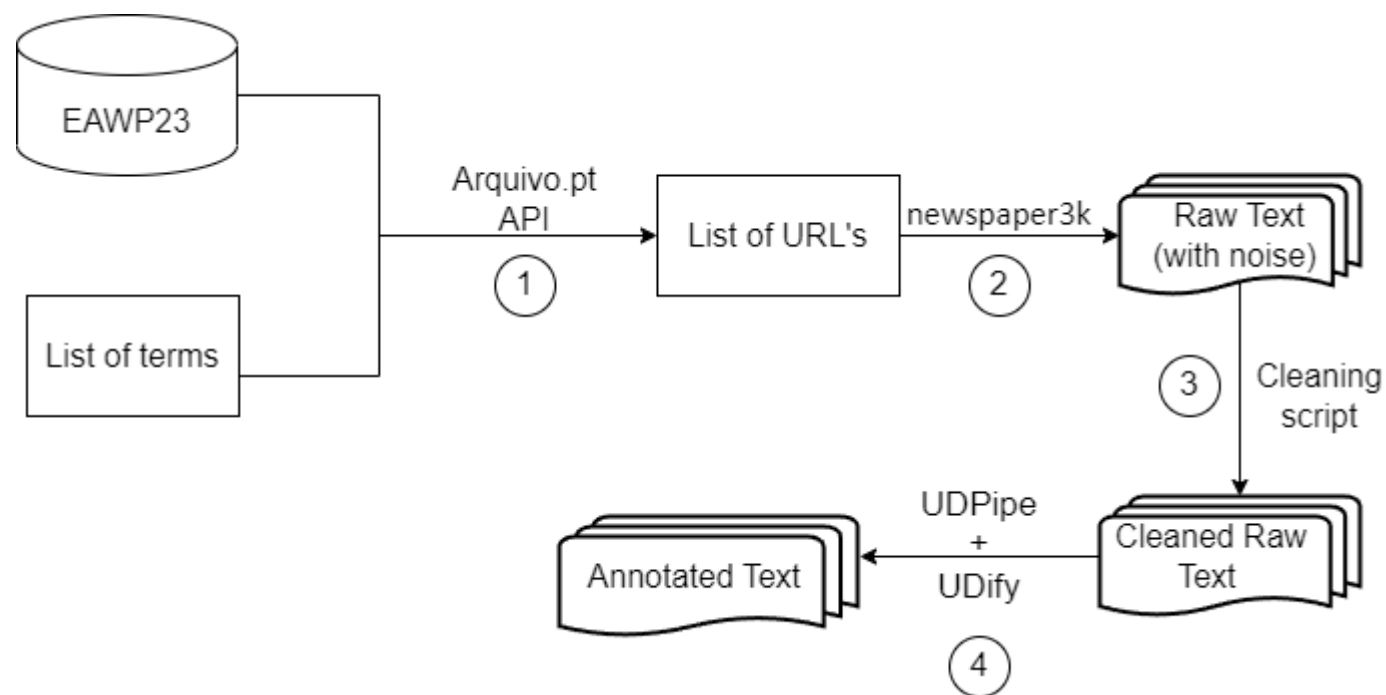
# Dataset creation

- EAWP23
  - Semi-automatic selection of relevant online content (40 relevant terms + collaborative list)
  - Crawling web content (6 crawls using different configurations)
  - 99 million URLs (4.8 TB)

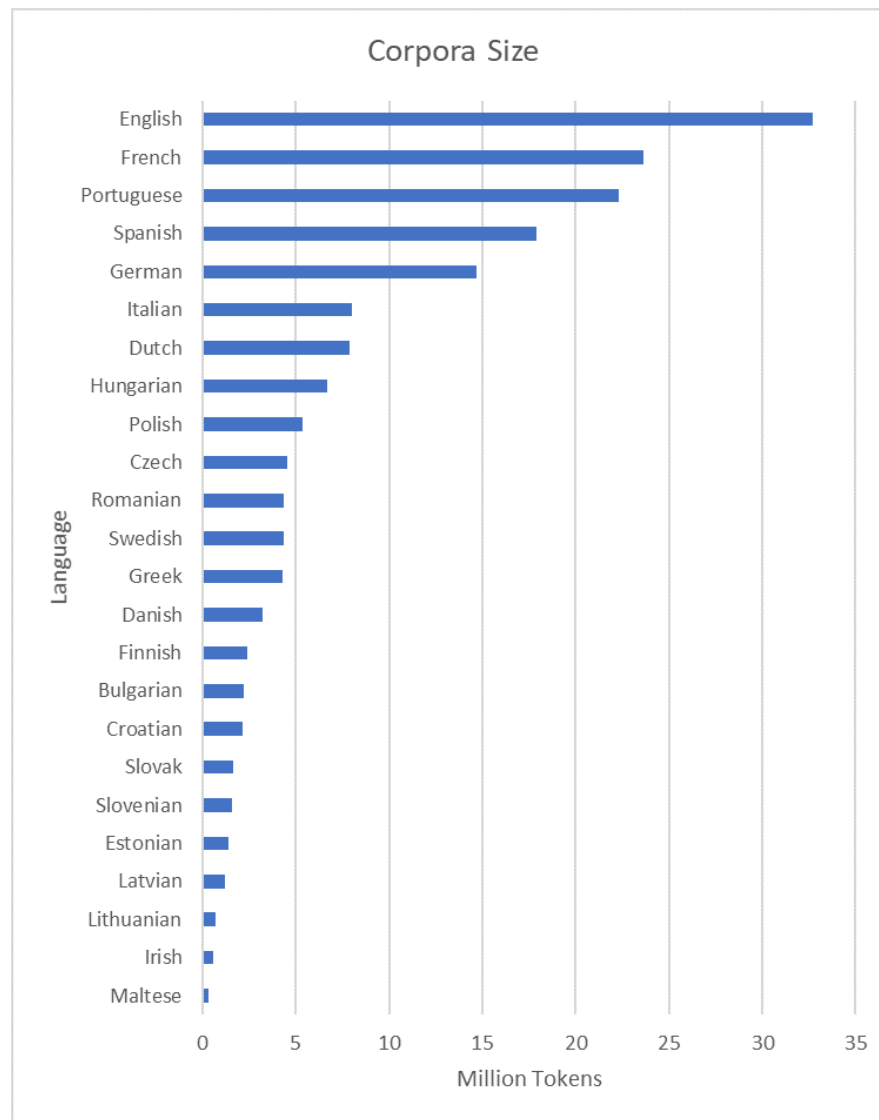# Dataset creation



Raw text corpus:

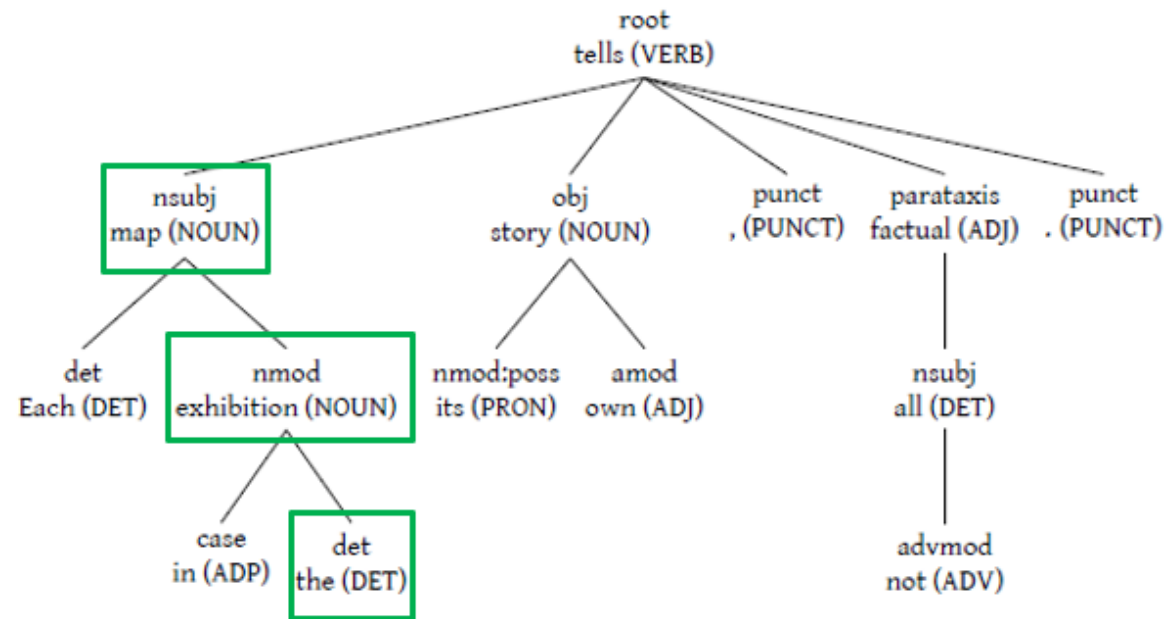https://dataverse.harvard.edu/dataset.xhtml?persistentId=doi:10.7910/DVN/WCGNHU

Annotated corpus:

https://dataverse.harvard.edu/dataset.xhtml?persistentId=doi:10.7910/DVN/ULWS1K

# Dataset Creation



Corpora Size — a horizontal bar chart showing Million Tokens by Language:
- English ~33
- French ~23.5
- Portuguese ~22
- Spanish ~18
- German ~14.5
- Italian ~8
- Dutch ~7.5
- Hungarian ~6.5
- Polish ~5.5
- Czech ~4.5
- Romanian ~4.5
- Swedish ~4.5
- Greek ~4.5
- Danish ~3
- Finnish ~2.5
- Bulgarian ~2.5
- Croatian ~2.5
- Slovak ~2
- Slovenian ~2
- Estonian ~1.5
- Latvian ~1.5
- Lithuanian ~1
- Irish ~0.5
- Maltese ~0.5

# Typological Analysis

- Best identified quantitative approach for dependency parsing improvement:
  - MarsaGram linear patterns

- MarsaGram is a tool for exploring treebanks, it extracts context-free grammars (CFG) from annotated data-sets

- Linear property:
  - C_A_precedes_B
  - element A (POS - deprel) precedes element B (POS - deprel) in a subtree ruled by element C (POS)

- Language vectors:
  - Features: MarsaGram patterns
  - Values: Frequency in the corpus

# Typological Analysis

- NOUN_precede_DET-det_NOUN-nmod



"Each map in the exhibition tells its own story, not all factual."

# Typological Analysis

PUD

Arquivo.pt

**Cluster Dendrogram**

**Cluster Dendrogram**

d
hclust (*, "ward.D2")

d
hclust (*, "complete")

# Conclusions and Perspectives

- Differences in the typological classification specially concerning Germanic languages

- To be evaluated:
  - Deeper analysis of the bias introduced by the automatic annotation
  - Impact on the choice of best language pairs for dependency parsing improvement

# *New* research tools available!

daniel.gomes@fccn.pt

# Search **texts** from the past in any language or about any event

# Arquivo.pt search query logs available for research *new!*



- Unique resource for innovative research
- Information retrieval over historical web data is unexplored ground
- Only web archive that supports full-text search
- Only web archive that supports image search

# APIs and Bulk download of web-archived content *new!*

# Four **A**pplication **P**rogramming **I**nterfaces to automatically search and access web-archived content

**Arquivo.pt APIs**

Your new app



### APIs
Daniel Gomes edited this page 11 days ago · 33 revisions

### APIs specific to Arquivo.pt that enable the full exploration of its functions

- Arquivo.pt API (Full-text & URL search)
- Image Search API v1.1 (beta version)

### APIs based on international standards to enable interoperability among web archives and code reuse

- CDX-server API (URL search)
- Memento API (URL search)

**arquivo.pt/api**

# 🔗 Bulk download of web-archived resources

If you need to download a large amount of web-archived resources, such as all the URLs archived from a large website along time, we suggest the following methodology:

1. Analyse the Arquivo.pt collections so that you may choose those which may contain the most interesting web-archived data for your use case. If you have any doubt, contact us.

2. Download the CDXJ index files, (what is CDXJ?) of the Arquivo.pt collections you selected to process. For this purpose, analyse the "column A: Collection ID" and the corresponding CDXJ index files on "column H: Collection CDXJ File");

3. Create a list of selected URLs to be downloaded, extracted from the CDXJ index files (e.g. using Linux grep command);

4. Download the web-archived resources for the list of selected URLs from Arquivo.pt by using the above APIs or, by building links to directly access the web-archived resources. These links are available on the Technical details of the Options top-right menu when accessing a web-archived page. For instance, for the URL http://publico.pt/ with timestamp 20120201160355 extracted from the CDXJ index file, build the following links to download the:

- original file of the web-archived page (loses replay quality because the original internal links are **not rewritten** to reference web-archived images or stylesheets): https://arquivo.pt/noFrame/replay/20120201160355id_/http://publico.pt/
- web-archived page without the Arquivo.pt UI frame (internal links are rewritten to reference web-archived resources): https://arquivo.pt/noFrame/replay/20120201160355/http://publico.pt/

**arquivo.pt/api#bulk**

# 1. Analyze the [Arquivo.pt collections](#) to choose the most relevant

List of the collections preserved by Arquivo.pt (publicly available)

Collections | About

| Collection ID | Description | Collection Start date | Collection End date | Total files (Number of URLs collected) | Total seeds (Preserved sites) | Total volume of WARCS files (TB) | Collection CDXJ File publicly available |
|---|---|---|---|---|---|---|---|
| AWP1 | 1st complete crawl of the Portuguese web, mainly from the .PT domain, in 2008. | 2008-02-12 | 2008-03-06 | 56,046,288 | 154,787 | 1.60 | AWP1_15.cdxj |
| AWP2 | 2nd complete crawl of the Portuguese web, mainly from the .PT domain, in 2008. | 2008-03-11 | 2008-05-30 | 48,718,404 | - | 1.60 | AWP1_15.cdxj |

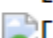# Cross-lingual collection on 2019 EU celections (EAWP23)



Created as a use case for cross lingual research

➢ Useful for under-resourced languages

➢ Parallel corpora because EU official websites are carefully translated

# 2. Download the correspondent CDXJ index files



**Index of /cdxj**

| [ICO] | Name | Last modified | Size | Description |
|---|---|---|---|---|
| [PARENTDIR] | Parent Directory | | - | |
| [ ] | AWP1_15.cdxj | 2022-07-11 14:22 | 482G | |
| [ ] | AWP8.cdxj | 2022-11-10 11:45 | 36G | |
| [ ] | AWP16_25.cdxj | 2022-07-11 14:44 | 1.1T | |
| [ ] | AWP26.cdxj | 2022-07-11 13:50 | 111G | |
| [ ] | AWP27.cdxj | 2022-07-11 13:47 | 90G | |
| [ ] | AWP28.cdxj | 2022-07-11 14:00 | 188G | |
| [ ] | EAWP23.cdxj | 2023-02-03 15:17 | 43G | |
| [ ] | EAWP24.cdxj | 2023-02-03 14:03 | 2.6G | |

https://arquivo.pt/cdxj

# 3. Select URLs to be downloaded from the CDXJ index files

pt,bvl)/avisos.html 19961013163232 {"status": "200", "url": "http://www.bvl.pt/avisos.html", "filename": "AWP-Roteiro-20090510220155-00000.arc.gz", "length": "0", "mime": "text/html", "offset": "5376161", "digest": "XLFEUTKLZGP2H4VK3RUKX7PI223I5MRH", "collection": "Roteiro"}

19961013163232 {"status": "200", "url": "http://www.bvl.pt/avisos.html"
"XLFEUTKLZGP2H4VK3RUKX7PI223I5MRH", "collection": "Roteiro"}

90510220155-00000.arc.gz", "length": "0", "mime": "text/html"

# 4. Download original file of the web-archived page

# 4. Download web-archived page without the Arquivo.pt UI frame



https://arquivo.pt/noFrame/replay/**20120201160355/http://publico.pt/**

# Tools to support the preservation of research data available online

# CitationSaver preserves citations to web resources *new!*



[arquivo.pt/citationsaver](arquivo.pt/citationsaver)

# **SavePageNow** to reference later! *new!*



arquivo.pt/savepagenow

# Learn more

# Tutorial: how to explore Arquivo.pt using Python *new!*



[arquivo.pt/tutorial](arquivo.pt/tutorial)

# Suggested publications

- Web archives as research infrastructure for digital societies: the case study of Arquivo.pt (2022), Daniel Gomes ***new!***

- The Anatomy of a Web Archive Image Search Engine (2022), André Mourão ***new!***

- Information Search in Web Archives (2014), Miguel Costa

- Learning Temporal-Dependent Ranking Models (2014), Costa et al.

➢ [arquivo.pt/publications](arquivo.pt/publications)

# Arquivo.pt **Annual** Awards





Support of the President of the Republic of Portugal

➢ **arquivo.pt/awards**

# Keep in touch!

**Subscribe our channels**:

➢ [arquivo.pt/subscribe](arquivo.pt/subscribe): mailing list in English

➢ [arquivo.pt/news](arquivo.pt/news): social networks and videos