# CLARIN and Libraries

9 May 2022 - 10 May 2022
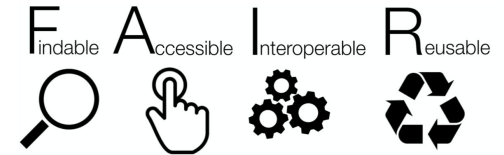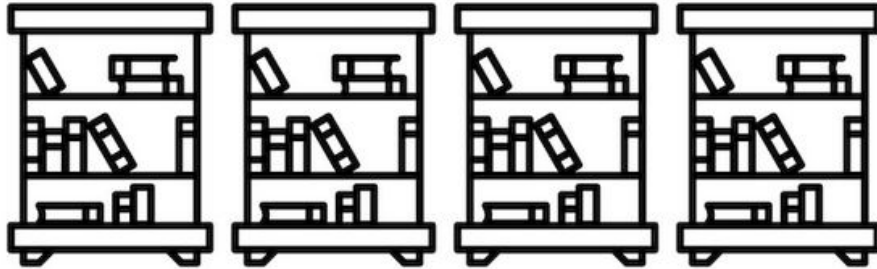
# DATA-KBR-BE: Data-level access to digitised collections for digital humanities research

**Sally Chambers, KBR, Royal Library of Belgium and Ghent Centre for Digital Humanities 🐦 @schambers3**
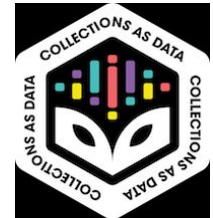
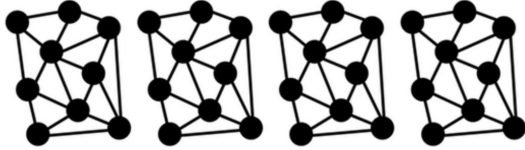# Collections as Data: "Always Already Computational" - Final Report



Padilla, Thomas, Allen, Laurie, Frost, Hannah, Potvin, Sarah, Russey Roke, Elizabeth, & Varner, Stewart. (2019, May 22). Final Report --- Always Already Computational: Collections as Data (Version 1). Zenodo. http://doi.org/10.5281/zenodo.3152935

**Always Already Computational: Collections as Data final report and project deliverables:** https://osf.io/mx6uk/wiki/home/

# Collections as Data: Part to Whole

https://collectionsasdata.github.io/part2whole/

**Collections as Data: Part to Whole** aims to foster the development of broadly viable models that support implementation *and* use of **collections as data**. This effort is made possible by The Andrew W. Mellon Foundation (**see grant narrative**).

Over a period of three years, *Part to Whole* **will fund and programmatically support** two cohorts. Cohorts will be comprised of project teams jointly led by librarians and disciplinary scholars. Project teams will develop models that support collections as data implementation and holistic reconceptualization of services and roles that support scholarly use. Collections as data produced by project activity will exhibit high research value, demonstrate the capacity to serve underrepresented communities, represent a diversity of content types, languages, and descriptive practices, and arise from a range of institutional contexts.

**Wittmann, R., Neatrour, A., Cummings, R., & Myntti, J. (2019). From Digital Library to Open Datasets.** *Information Technology and Libraries*, *38* (4), 49-61. https://doi.org/10.6017/ital.v38i4.11101

*Articles*

## From Digital Library to Open Datasets: Embracing a "Collections as Data" Framework

Rachel Wittmann, Anna Neatrour, Rebekah Cummings, and Jeremy Myntti

**ABSTRACT**

*This article discusses the burgeoning "collections as data" movement within the fields of digital libraries and digital humanities. Faculty at the University of Utah's Marriott Library are developing a collections as data strategy by leveraging existing Digital Library and Digital Matters programs. By selecting various digital collections, small- and large-scale approaches to developing open datasets are explored. Five case studies chronicling this strategy are reviewed, along with testing the datasets using various digital humanities methods, such as text mining, topic modeling, and GIS (geographic information system).*

# Newspapers as Data

**Bibliothèque nationale du Luxembourg**
**Open Data**

## Historical Newspapers

**https://data.bnl.lu/data/historical-newspapers/**

### STARTER PACK

## 250MB

of digitised newspapers

✓ 5 days of news
✓ 5 newspaper issues
✓ 22 pages
✓ D'Wäschfra (1868)
✓ Public Domain, CC0 (See copyright notice)
✓ Best for getting started & developing

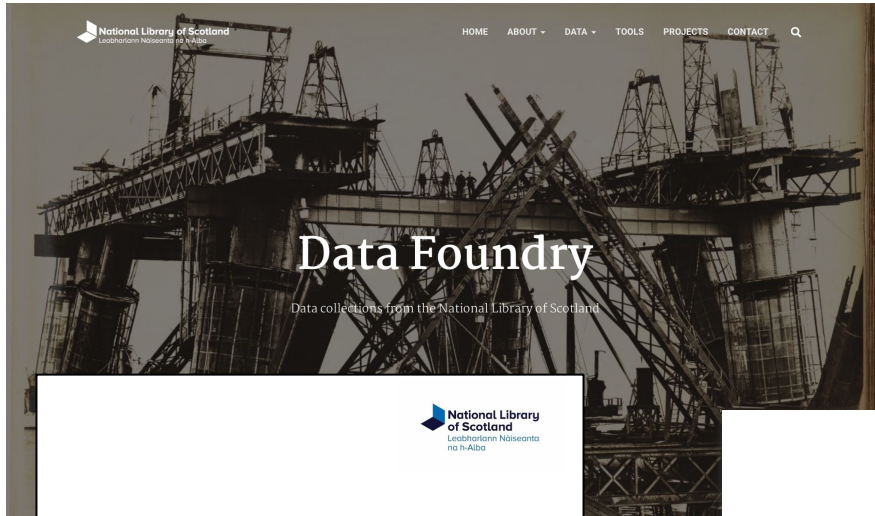⬇ DOWNLOAD (ZIP)

### TEXT ANALYSIS PACK

## 2GB

of processed newspapers data

✦ 38 years of news (1841-1878)
✦ 23663 processed newspaper issues
✦ 97285 processed pages
✦ 510505 extracted articles
✦ Public Domain, CC0 (See copyright notice)
✦ Best for getting started with text analysis

⬇ DOWNLOAD (ZIP)

# National Library of Scotland: Data Foundry

Ames, S. (2021) "Transparency, Provenance and Collections as Data: The National Library of Scotland's Data Foundry." LIBER Quarterly 31 (1): 1–13. https://doi.org/10.18352/lq.10371.

## https://data.nls.uk

https://data.nls.uk/download/national-library-of-scotland-open-data-publication-plan.pdf

# Digital Library of the Caribbean as Data



**https://dataverse.fiu.edu/dataverse/dloc-as-data**

# CLARIN Resource Families

**Resource families**

CLARIN

About ⌄    Language Resources ⌄    Learn & Exchange ⌄    Events    News    Contact
🇺🇦 CLARIN and Ukraine

Home  /  Language Resources  /  Resource Families

**Resource families**

## Resource Families

### Introduction

The aim of the CLARIN Resource Families initiative is to provide a
available language resources in the CLARIN infrastructure for res
sciences and human language technologies. The overviews are r
the listings  are sorted by language.

The listings for each family include the most important metadata
size, text sources, time periods, annotations and licences as well
concordancers. In addition to the resources found in the CLARIN
other existing valuable language resources which have not yet b

Currently, overviews are available for 12 corpora families, 5 families of lexical resources, and 4 tool families. See below. For information about applying for funding for small projects that can help to extend the scope of the initiative, see https://www.clarin.eu/content/clarin-resource-families-project-funding.

**Corpora**
- Computer-mediated communication corpora
- Corpora of academic texts
- Historical corpora
- L2 learner 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

## https://www.clarin.eu/resource-families

# DATA-KBR-BE

**Facilitating data-level access to KBR's digitised and born-digital collections for digital humanities research**
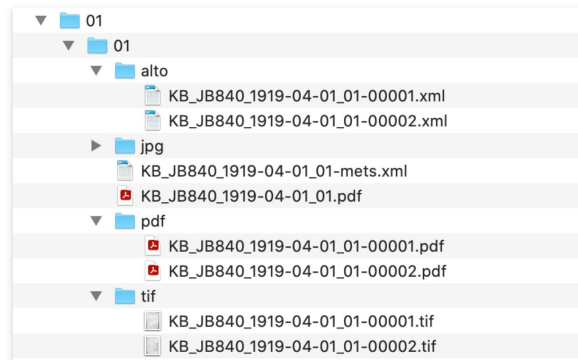
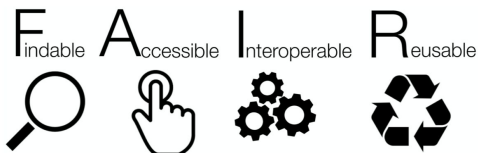**https://www.kbr.be/en/projects/data-kbr-be/**

# Digitised Historical Newspapers as Data

## Collections as Data

Providing data-level access to digital collections is a primary challenge for undertaking digital humanities research. In the United States, the flagship initiatives, 'Always Already Computational: Collections as Data' and 'Collections as Data: Part to Whole', define 'Collections as Data' as a "conceptual orientation to collections that renders them as ordered information, stored digitally, so that they are inherently amenable to computation". The initiative was established to document, exchange experience and share knowledge to encourage cultural heritage institutions to implement 'collections as data' in their own institutions. **DATA-KBR-BE will kick-start the implementation of 'Collections As Data' in Belgium.**

**https://collectionsasdata.github.io**



**F**indable **A**ccessible **I**nteroperable **R**eusable

**Providing data-level access to the underlying files of digitised and born-digital cultural heritage resources to facilitate data analysis by means of tools and methods developed in the field of digital humanities**
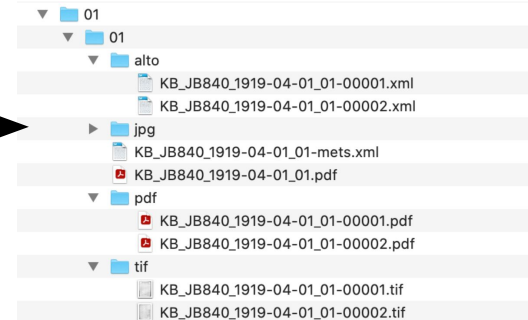
**https://www.kbr.be/en/projects/data-kbr-be/**

GHENT UNIVERSITY

KBR

# Digitised Historical Newspapers as Data



https://www.kbr.be/en/projects/data-kbr-be/

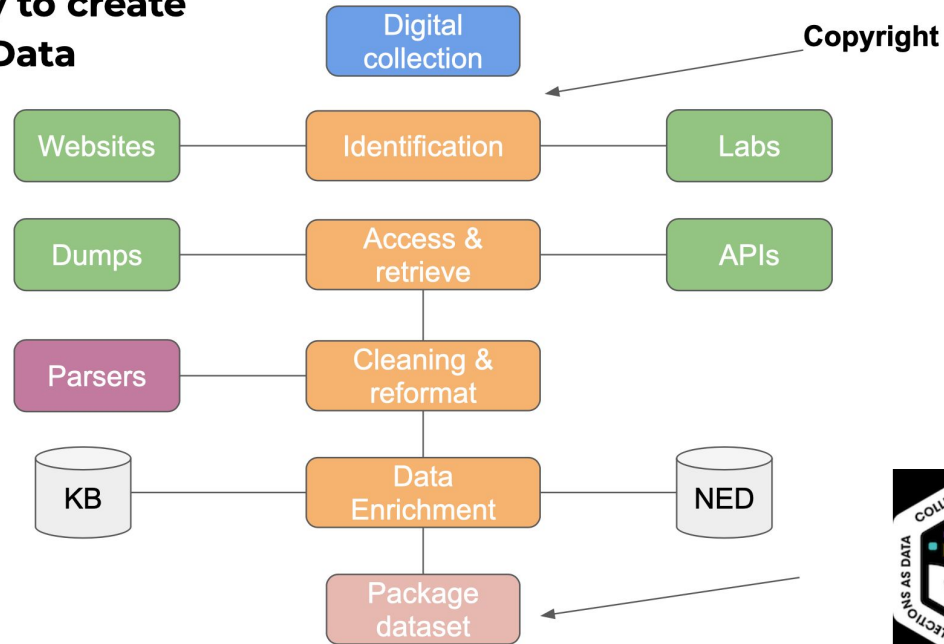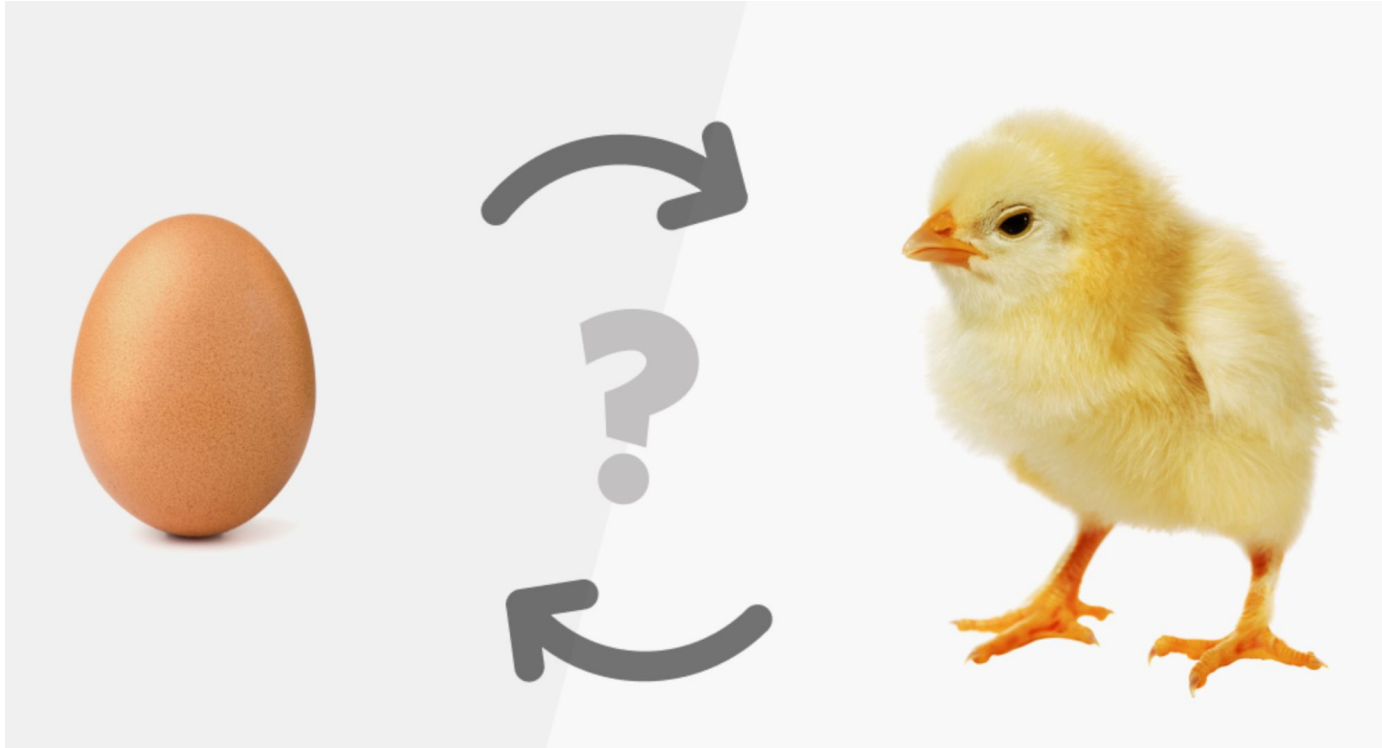https://www.kbr.be/en/projects/digital-research-lab/

# Collections as Data Methodology



A methodology to create Collections as Data

Candela, G., Sáez, M. D., Escobar Esteban, Mp., & Marco-Such, M. (2020). *Reusing digital collections from GLAM institutions.* Journal of Information Science: https://doi.org/10.1177/0165551520950246 & http://rua.ua.es/dspace/handle/10045/109460

# Collections or Corpora?

# What is a corpus … ?

**McGillivray, B., Poibeau, T., & Ruiz Fabo, P. (2020). Digital Humanities and Natural Language Processing:"Je t'aime… Moi non plus". DHQ, Digital Humanities Quarterly. 2020, 14.2**
http://www.digitalhumanities.org/dhq/vol/14/2/000454/000454.html

# Collections as Data @ KBR

Candela, G., Sáez, M. D., Escobar Esteban, Mp., & Marco-Such, M. (2020). *Reusing digital collections from GLAM institutions*. Journal of Information Science: https://doi.org/10.1177/0165551520950246 & http://rua.ua.es/dspace/handle/10045/109460

# Interdisciplinary Research Scenarios

- **Collective Action Belgium** led by GhentCDH ↗, focuses on social history in the Interbellum and World War Two period and aims to trace the dynamics of contention, strikes, demonstrations and other forms of collective action in Belgium as reported in Belgian newspapers;
- **The feuilleton in Belgium**, led by ACDC ↗, focuses on literary studies in the period 1830–1930 and aims to map the publication of literature in Belgian newspapers across the first century of the Belgian nation state;
- **History of Belgian Journalism**, led by ULB ↗ and KBR, focuses on media history from 1886 until now and aims to trace the history of Belgian journalism through the lens of critical discourses about journalism as in Belgian newspapers.

KBR **Where time is treasured**

Belgian Science Policy Office
belspo

# From collections to corpora

| Year | Date | Title | Cause | Number of strikers | Description |
|------|------|-------|-------|--------------------|-------------|
| 1893 | 12–18 April | Belgian general strike of 1893 | Franchise reform | 200,000 | Successfully led to the establishment of universal male suffrage with plural votes.[3] Thirteen strikers were killed and socialist leaders were briefly arrested.[5] |
| 1902 | 10–20 April | Belgian general strike of 1902 | Franchise reform and an end to plural vote | 350,000 | Failed to achieve the abolition of the plural vote as Catholics and Liberals united to oppose constitutional reform. The Belgian Workers' Party had been reluctant to support the strike and it soon descended into violence in Brussels and parts of Wallonia. 12 workers and one policeman were killed. Union membership dropped sharply in its aftermath.[5] |
| 1913 | 14–24 April | Belgian general strike of 1913 | Franchise reform | 400,000 | Carefully planned to avoid the same problems as 1902, the strike gained the promise of electoral reform but its proposals were postponed by the outbreak of World War I and the subsequent German occupation. The policy was finally adopted in 1919.[3][6] |
| 1932 | 7 July–9 September | Belgian general strike of 1932 | Pay, working hours and unemployment insurance | | Began after a spontaneous strike by coal miners in the Borinage and involved Communist agitation following a severe decrease in living standards and real wages during the Great Depression. Two people were killed during the strike.[7] |
| 1936 | 2 June-2 July | Belgian general strike of 1936 | Working hours, paid holiday, union reforms | 500,000 | Broke out at the port of Antwerp and led to the creation of a National Labour Conference.[7] Although influenced by the French Popular Front and held against the backdrop of the Spanish Civil War, it was also supported by Catholic trade unions.[8] |
| 1950 | 24 July-3 August | Belgian general strike of 1950 | "Royal Question" | 700,000 | Chiefly active in Wallonia, the strike contributed to the abdication of King Leopold III on 1 August 1950. At least four strikers were killed. |

## https://en.wikipedia.org/wiki/General_strikes_in_Belgium

GHENT UNIVERSITY   KBR

# Which newspapers are digitised?

| | Title | Journal No. | Digitally availability | 1886 | 1893 | 1902 | 1913 | 1932 | 1936 | 1950 | Extraction |
|---|-------|-------------|------------------------|------|------|------|------|------|------|------|------------|
| 1 | Vooruit: socialistisch dagblad | JB 809 | 1884-1889; 1901-1902; 1911-1950 | Yes | No | Yes | Yes | Yes | Yes | Yes | 1913, 1950 |
| 2 | Het Volk : antisocialistisch dagblad | JB 785 | 1911-1916; 18 apr. 1918-31 maart 1927; 1931-1950. | No | No | No | Yes | Yes | Yes | Yes | 1913, 1950 |
| 3 | Vaderland | JB 310 | 16 maart 1910-30 dec. 1913. | No | No | No | Yes | No | No | No | 1913, 1950 |
| 4 | Le Peuple : organe quotidien de la démocratie socialiste | JB837 | 1885-1907; 1911-1914; 1918-1940; 1944-1950. | Yes | Yes | Yes | Yes | Yes | Yes | Yes | 1913, 1950 |
| 5 | Le Vingtième Siècle | JB729 | 6 juin 1895-13 mai 1940 | Yes | Yes | Yes | Yes | Yes | Yes | No | 1913, 1950 |
| 6 | La Meuse : journal de Liège et de la Province | JB638 | 1856; 1858-1882; 1884-5 août 1914; 28 nov. 1918-11 mai 1940; 9 sept. 1944-1950. | No | Yes | Yes | Yes | Yes | Yes | Yes | 1913, 1950 |

# Importance of historical context

"The newspaper was founded in Ghent in 1884 with the support of the socialist cooperative Vooruit. When the Parti ouvrier belge - Belgische Werkliedenpartij ("Belgian Workers Party", POB-BWP, 1885) was founded, it was recognised as its official organ for the Flemish part of the country. The Ghent socialist leader Edward Anseele was editor-in-chief, but he also worked as a typographer. *Vooruit* was published under German censorship during the two world wars. Having reached its peak in the 1950s, it began a slow decline. It was succeeded by the daily *De Morgen* in 1978."

**Vooruit: socialistisch dagblad**

# DATA-KBR-BE Data Extraction 1



- 3 Research Scenarios

- 12 Newspaper Titles

- 3 years: 1885, 1913, 1950

- 18 'drop offs'

(list of titles for extraction)

# Methods for Data Sharing: KBR Send Data



**https://datasend.kbr.be**

# Methods for Data Sharing



https://pro.europeana.eu/page/harvesting-and-downloads#downloads

# Interdisciplinary collaboration with data scientists



**Dilawar Ali, Kenzo Milleville, Alec Van de Broeck & Steven Verstockt, IDLab, UGent**

# NewspAIper Demonstrator

**Dilawar Ali, Kenzo Milleville, Alec Van de Broeck & Steven Verstockt, IDLab, UGent**
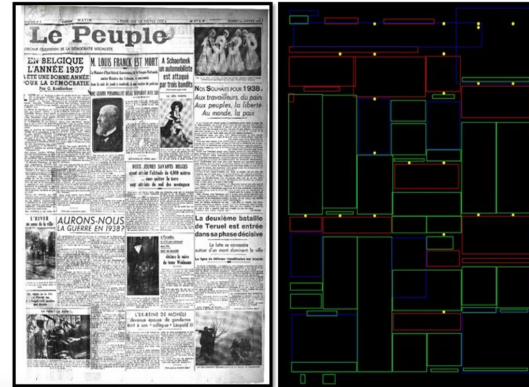
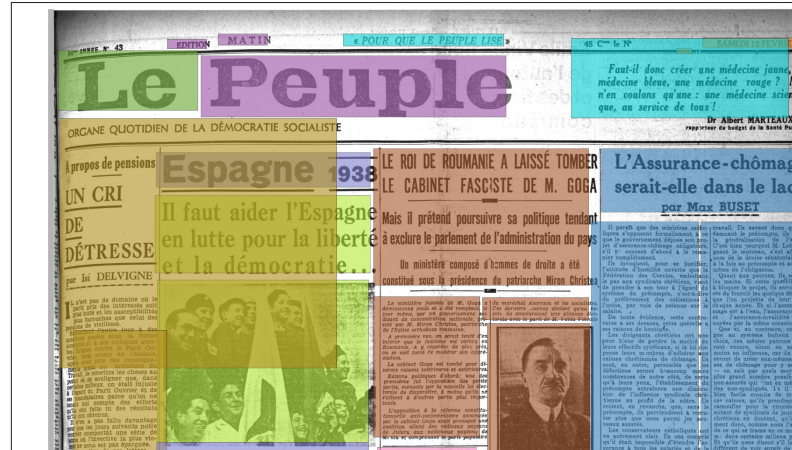

**https://tw06v072.ugent.be/kbr/**

# NewspAIper Demonstrator: Document Layout Analysis & Article Separation



Page Viewer

Click on an article/image to show the recognized text and related articles/images

Le Peuple: Saturday, 12 February 1938, page 1



*Document layout analysis using BelgicaPress*

## https://tw06v072.ugent.be/kbr/

GHENT UNIVERSITY

KBR

# NewspAIper Demonstrator: Named Entity Recognition & Linking

## Page Viewer

Click on an article/image to show the recognized text and related articles/images

Le Peuple: Thursday, 24 February 1938, page 1



**Most similar articles**

Samedi matin, devant de nombreux délégués réunis à la Maison du Peuple de 'Bruxelles le camarade D.

TROISIÈME JOURNÉE DE DISCUSSION sur la Politique étrangère au Conseil Général du P.O.B. CORNEILLE ME

DEVANT LE PARLEMENT Dans un discours vivement acclamé Paul-Henri Spaak souligne les mérites de la po

**Article text:**

Au Conseil Général du P.O.B. Buset Paul-Henri Spaak ont introduit un large débat sur LA POLITIQUE EXTÉRIEURE DE LA BEL GIQUE Les délégués ont salué la mémoire d'Edouard Anseele LA DISCUSSION SERA POURSUIVIE MERCREDI PROCHAIN

## https://tw06v072.ugent.be/kbr/

# NewspAIper Demonstrator - Image similarity

**Page Viewer**

Click on an article/image to show the recognized text and related articles/images

Le Peuple: Tuesday, 1 March 1938, page 1

Most similar images:

Most similar images:

GHENT UNIVERSITY

KBR

https://tw06v072.ugent.be/kbr/

# Computer Vision and Machine Learning Approaches to Improve Searchability and Building Corpora using Historical Newspapers

DILAWAR ALI, KENZO MILLEVILLE, STEVEN VERSTOCKT, IDLab, Ghent University-imec, Ghent, Belgium

SALLY CHAMBERS, JULIE M. BIRKHOLZ, Ghent Centre for Digital Humanities, Ghent University, Ghent, Belgium and KBR- the Royal Library of Belgium, Belgium

Historical newspaper collections provide a wealth of information abou
opportunities for not only broadening the access of these collections in
engineering techniques to enhance the extraction of information such
These research domains are, however, dependent on the quality of t
layout analysis (e.g. for automatic article segmentation and image
depending on the techniques used by cultural heritage institutions
constantly re-digitize collections in line with the improvements of th
only in computer engineering but also in the (digital) humanities, to
vision and machine learning approaches to these challenges.

## 3.3 General Workflow

In this research project, we focused on metadata enrichment approaches using traditional computer vision and machine learning approaches. The general workflow is shown in figure 1.
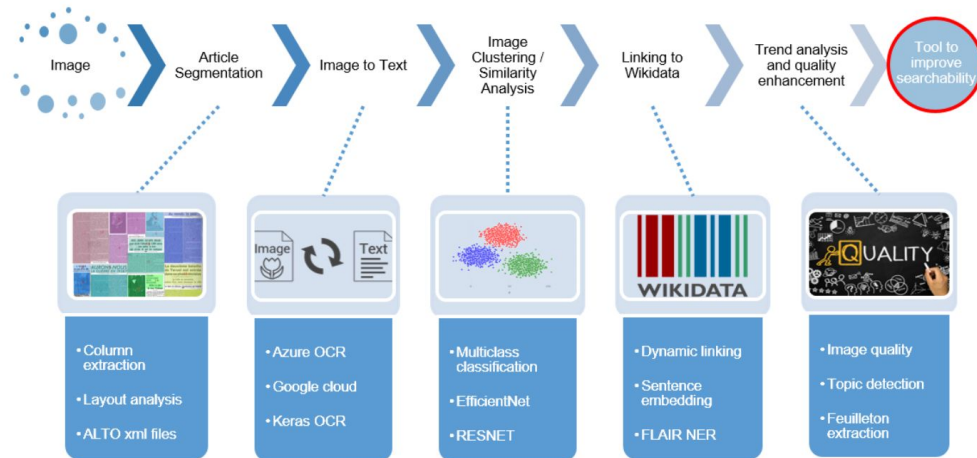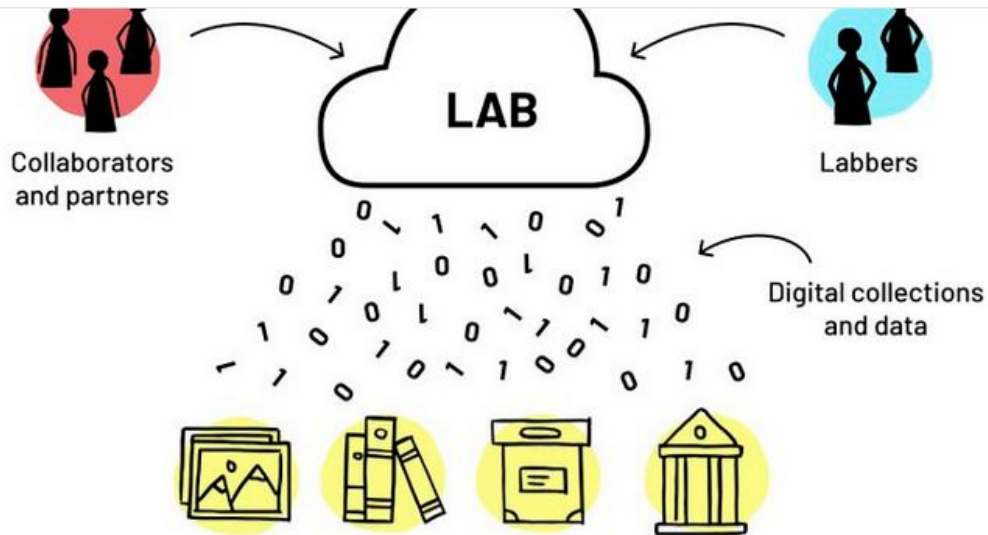


Fig. 1. General workflow: Pipeline for developing a tool to improve searchability

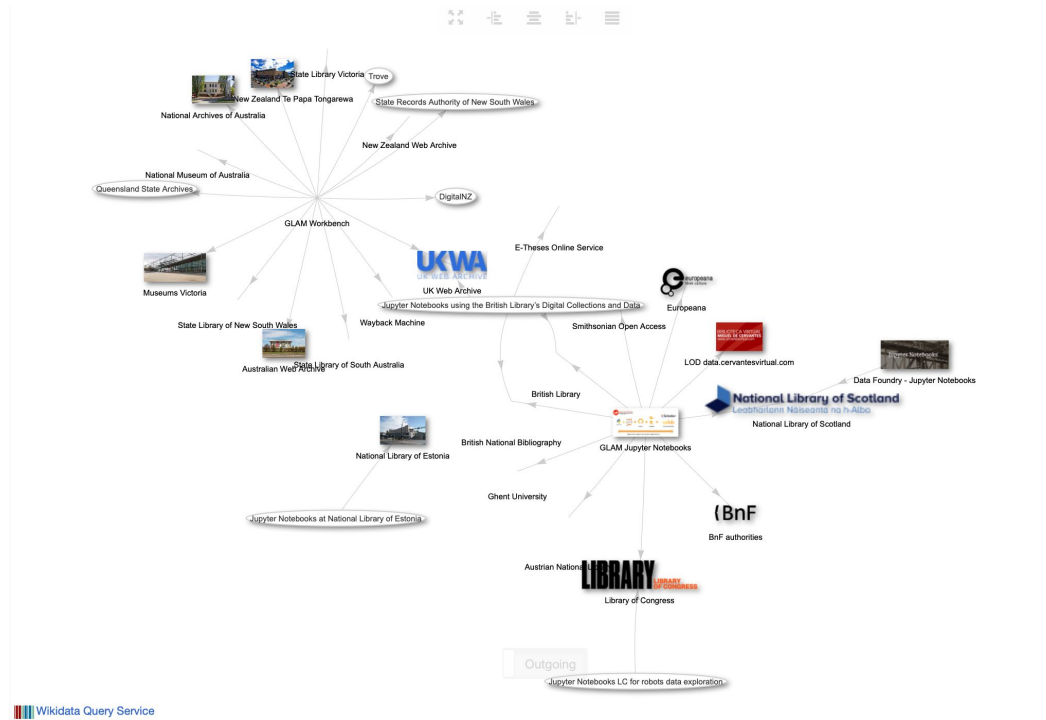International GLAM Labs Community

Collaborators and partners

LAB

Labbers

0 1 1 0 0 1
0 1 0 1 0 1 0
1 1 0 1 0 1 0
0 1 0 1 1 0 1 0
0 1 0

Digital collections and data

OPEN A GLAM LAB

**Chapter: Rethinking Collections as Data:**
**https://glamlabs.pubpub.org/pub/urhm68yr/release/1**

**https://glamlabs.io/books/**

# Computational access to digital collections



**International GLAM Labs Community: Computational Access to Digital Collections**
https://glamlabs.io/computational-access-to-digital-collections/

# Corpus Building: an interdisciplinary digital hermeneutics workflow



Oberbichler, S. et al. (2021) *Integrated interdisciplinary workflows for research on historical newspapers: Perspectives from humanities scholars, computer scientists, and librarians*. JASIST, August 2021.
https://doi.org/10.1002/asi.24565

# Historical Data Analysis using Jupyter Notebooks

## NLP Notebooks for Newspaper Collections

*A collection of notebooks for Natural Language Processing*

The following notebooks are aimed particularly at digital humanities scholars who use newspapers as a source. The focus lies on (topic-specific) collection building, a field that is becoming increasingly interesting with better article separation. Very specific research problems are addressed, such as building up collections with ambiguous keywords or working with certain genres. In order to best meet the needs of digital humanities scholars, NLP methods are adapted in new ways, the output is human-readable and the processed newspaper articles can be exported in the form of the original file. In addition, the notebooks allow the users to control the single steps and to choose what is best for their collection. While the NewsEye demonstrator offers the possibility to create datasets quickly and effectively, these notebooks offer possibilities to work on these collections according to specific questions.
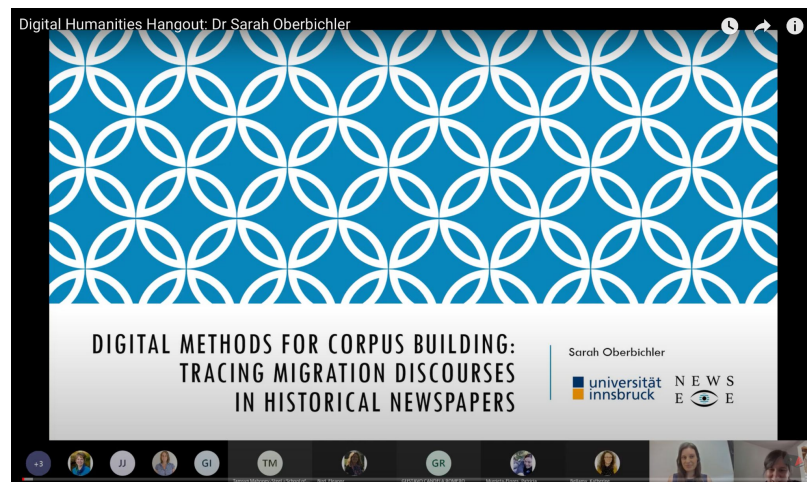
1. Text classification for topic-specific newspaper collections
2. Group similar newspaper articles
3. Discover a newspaper collection with diachronic Ngram clouds
4. Discourse in Spanish flu coverage
5. Topic Modeling and Uses of the Term Telegraph in the Context of Journalism

Oberbichler, S. (2020) NLP Notebooks for Newspaper Collections.
https://github.com/NewsEye/NLP-Notebooks-Newspaper-Collections

Oberbichler, S. (2021) Digital Methods for Corpus Building: Tracing Migration Discourses in Historical Newspapers, DH Hangout (Lancaster, Ghent, Lisbon), June 2021: https://www.youtube.com/watch?v=5w8X0qXP67M

# Europeana Research Community Café - Collections as Data

**Thomas Padilla's presentation**

# A European data space for cultural heritage

# Cultural Heritage Data is Humanities Research Data


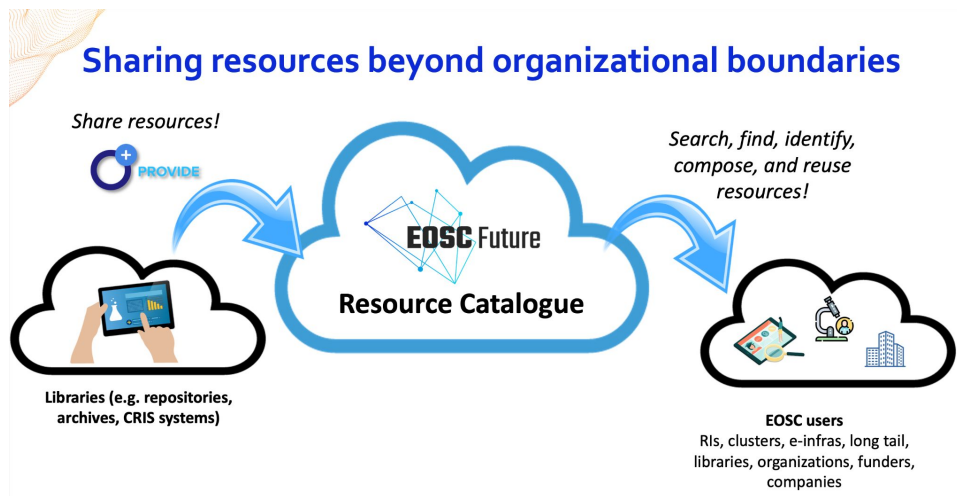
Tasovac, T., Chambers, S. and Tóth-Czifra, E. (2020) Cultural Heritage Data from a Humanities Research Perspective: A DARIAH Position Paper. https://hal.archives-ouvertes.fr/hal-02961317

# European Open Science Cloud (EOSC) and Libraries



Sharing resources beyond organizational boundaries

Share resources!

Search, find, identify, compose, and reuse resources!

PROVIDE

EOSC Future

**Resource Catalogue**

Libraries (e.g. repositories, archives, CRIS systems)

EOSC users
RIs, clusters, e-infras, long tail, libraries, organizations, funders, companies

**Libraries as EOSC providers**

*How and why libraries should become EOSC providers – Inge Van Nieuwerburgh, University of Gent*

*Current and upcoming functionalities for libraries – Paolo Manghi, OpenAIRE*

*A community presentation from Susanne Blumesberger, University of Vienna Library & Lisa Hönegger, AUSSDA*
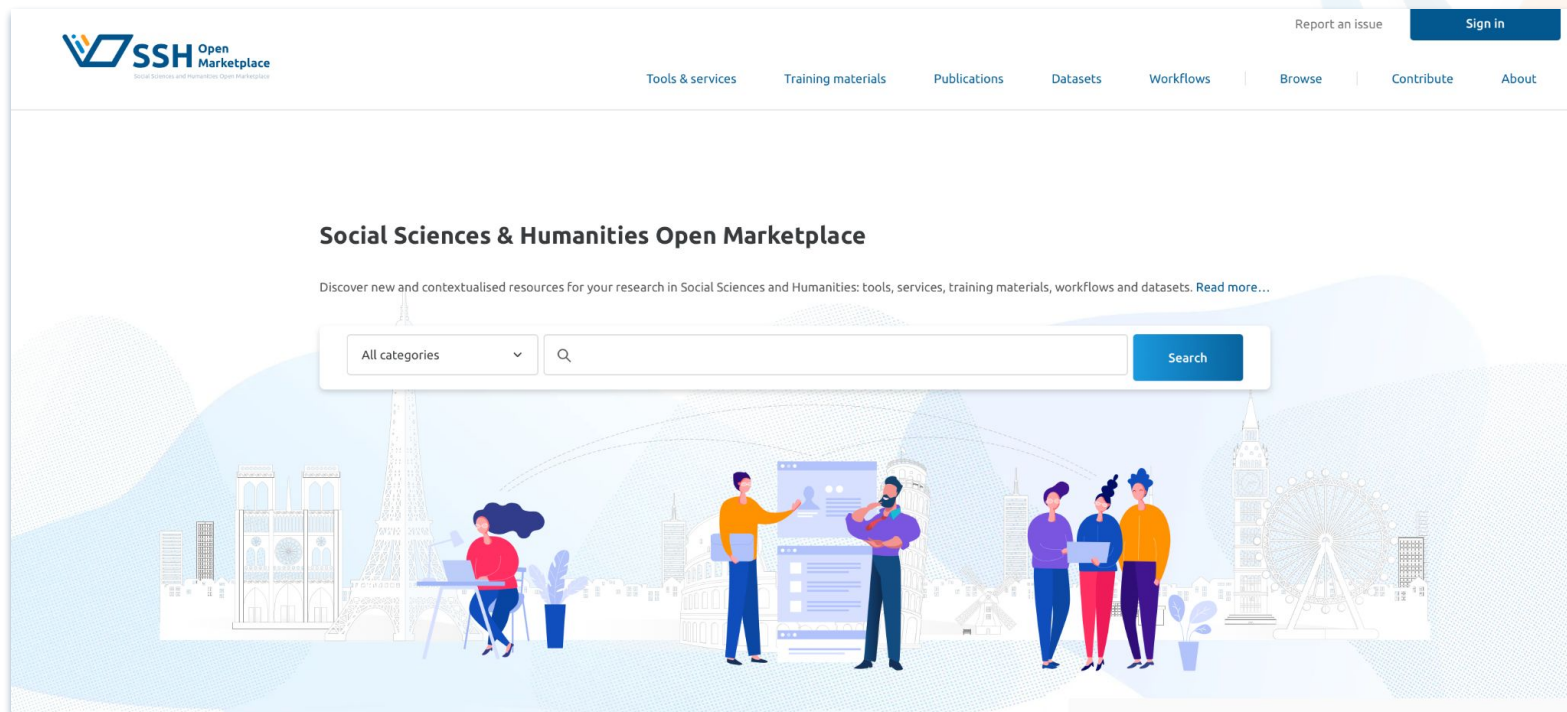
*Libraries as intermediaries – Sara Garavelli*

Presentation – recording

**Presentation:**
https://eoscfuture.eu/wp-content/uploads/2022/04/ProviderDays_Libraries_Presentation2Inge_clean.pdf
**Recording:** https://www.youtube.com/watch?v=_vxHJNU78Pw

# Social Sciences and Humanities Open Marketplace



**https://marketplace.sshopencloud.eu**

SSHOC
social sciences & humanities open cloud

# SSHOC Open Marketplace: from tools to workflows

## Create a dictionary in TEI

This scenario sets out the best practices for creating a born-digital dictionary, especially with the TEI (Text Encoding Initiative). However, buiding a standardized lexicographical dataset is not only a data format problem, it is also an intellectual and technical process where one has to choose how to model their data, and with which tools operate in order to create an easy-to-use and sustainable resource.
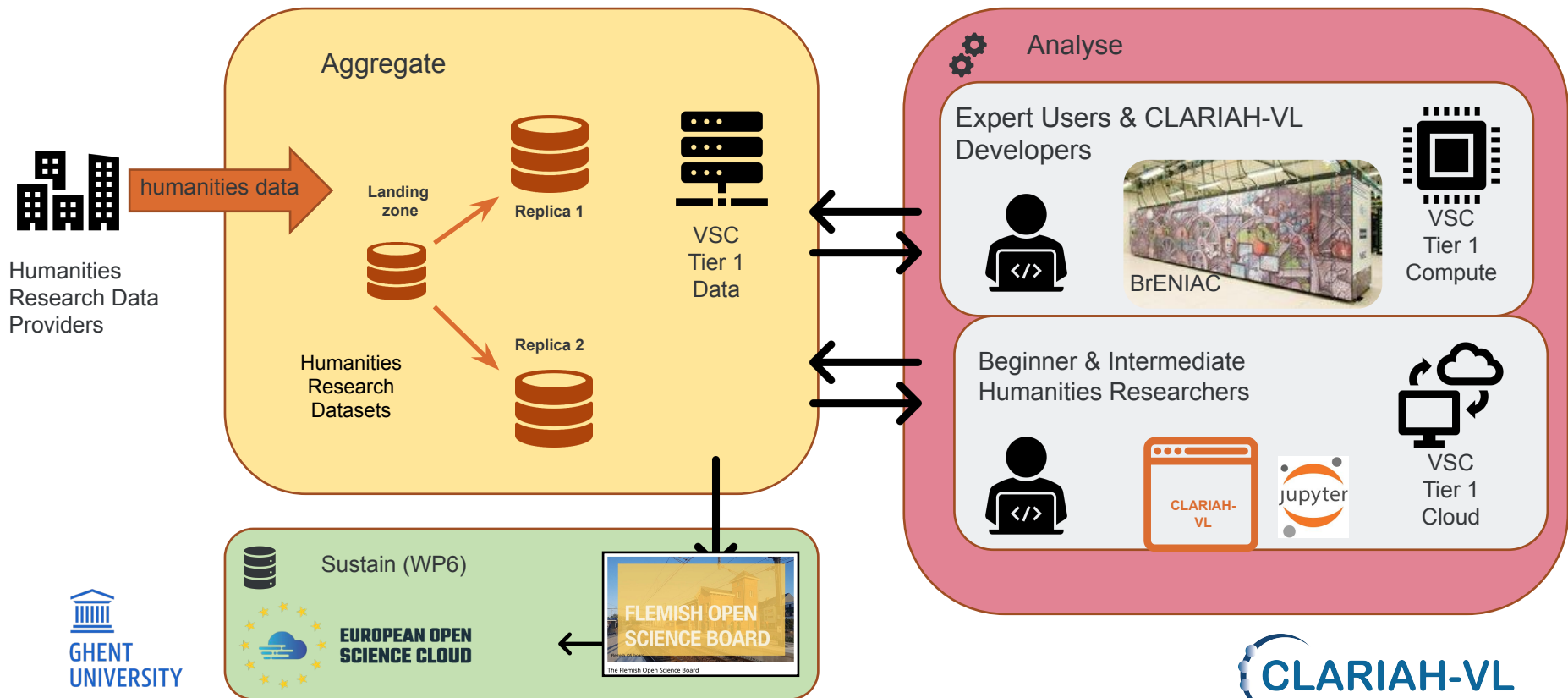


Photo by Joshua Hoehne on Unsplash

## Workflow steps (5)

| 1 | Build the model of the dictionary | Expand ▼ |
| 2 | Create a corpus of useful resources for the dictionary. | Expand ▼ |
| 3 | Setting up the editing environment | Expand ▼ |
| 4 | Setting up the publishing environment | Expand ▼ |
| 5 | Provide for long-term availability | Expand ▼ |

**https://marketplace.sshopencloud.eu/workflow/4qFarh**

**SSHOC**
social sciences & humanities open cloud

# Data-driven research in the arts and humanities

# CLARIN and Libraries

9 May 2022 - 10 May 2022

# Thanks for listening!

**Sally Chambers**
Sally.Chambers@kbr.be|Sally.Chambers@UGent.be |
**KBR, Royal Library of Belgium and**
**Ghent Centre for Digital Humanities, Belgium**