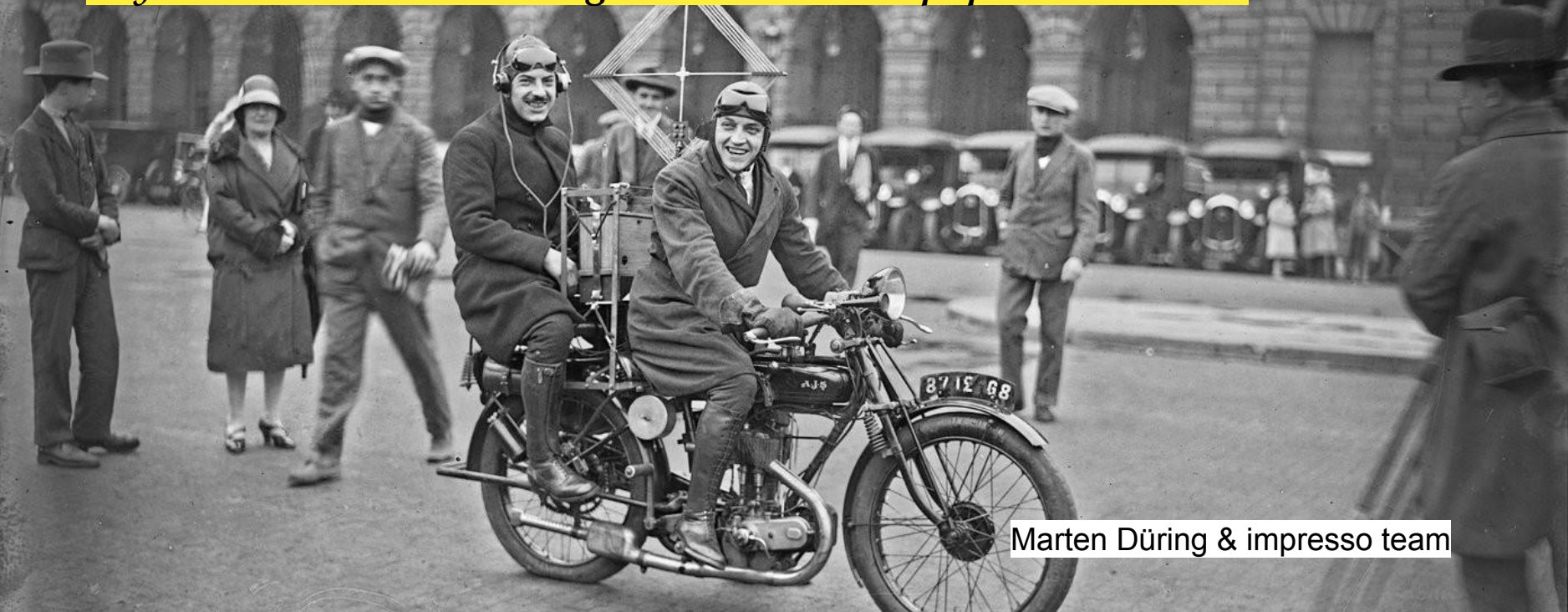# impresso - Media Monitoring of the Past I+II
*Beyond Borders: Connecting Historical Newspapers and Radio*

Marten Düring & impresso team

## Erstes Blatt.

### Serbische Wahlen.

### Erzbischof Dr. v. Abert †

### Türkisch-italienischer Krieg

### Zum Untergang der Titanic.

### Schloß Altenstein.

---

# Haute tension sur l'avenir électrique

*par Nicolas Honscher*



## L'ère des incertitudes

ÉNERGIE DE L'OUEST SUISSE (EOS)

## Christophe Babaiantz: le prix du renoncement à Kaiseraugst

---

# The landscape of historical newspapers



**European Newspapers - Titles Per Year by Country**

Legend: Wales, Switzerland, Sweden, Scotland, Norway, Netherlands, Italy, Isle of Man, Ireland, Germany, France, Finland, England, Denmark, Belgium, Austria

# The challenging landscape of historical newspapers

1. Institutional silos

2. Big and messy data

3. Noisy historical text

4. Visualisation and exploration

5. Digital Scholarship

# Newspapers as Data

'*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".

https://collectionsasdata.github.io/



Bibliothèque nationale du Luxembourg
Open Data

HOME    DATA ▾    TOOLS    API ▾

datasets contain XML (METS + ALTO), PDF, original TIFF and PNG files for every newspaper issue.

**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)

**DEV PACK**

## 3GB

of digitised newspapers

✓1 month of news
✓26 newspaper issues
✓112 pages
✓Luxemburger Wort (1877)
✓Public Domain, CC0 (See copyright notice)
✓Best for getting started with Big Data

⬇ DOWNLOAD (ZIP)

**SAMPLE PACK**

## 1GB

of digitised newspapers

✓11 different newspaper titles
✓1 issue per newspaper
✓News between 1845 and 1877
✓Public Domain, CC0 (See copyright notice)
✓Best for testing different newspapers and metadata

⬇ DOWNLOAD (ZIP)

https://data.bnl.lu

# TROVE

▾  2. Show the total number of articles per year

In another notebook, I look at different ways of visualising Trove newspaper searches over t
set the  q  parameter to a single space.

```
[6]:  # Set the q parameter to a single space to get ALL THE ARTICLES
      params["q"] = " "
```

Now we can find the total number of newspaper articles in Trove.

```
[7]:  # Get the JSON data from the Trove API using our parameters
      data = get_results(params)

      # Navigate down the JSON hierarchy to find the total results
      total = int(data["response"]["zone"][0]["records"]["total"])

      # Print the results
      print("There are currently {:,} articles in Trove!".format(total))
```

There are currently 233,666,567 articles in Trove!

Ok, that's not all that useful. What would be more interesting is to show the total number of a
details in this notebook but, in short, we have to loop through the decades from 1800 to 201

These two functions do just that.

https://glam-workbench.net/trove-newspapers/

OPEN A GLAM LAB

https://glamlabs.io/

**ANALYSING THE EDITIONS OF** *LES FLEURS DU MAL* **DE BAUDELAIRE FROM DATA.BNF.FR**

This notebook shows how to exploit the editions of *Les fleurs du mal* de Baudelaire using network graphs from data.bnf.fr.

⚙ Binder   👁 Preview

**COMPUTER VISION APPLIED TO SMITHSONIAN OPEN ACCESS**

This notebook introduces how to explore Smithsonian Open Access to apply computer vision methods in face detection.

⚙ Binder   👁 Preview

**ACCESSING EUROPEANA IIIF API**

This notebook extracts a dataset from the Europeana IIIF API. It performs an automatic search, retrieving the manifests from the IIIF server to create a dataset with the metadata as a CSV file.

⚙ Binder   👁 Preview

https://data.cervantesvirtual.com/glam-jupyter-notebooks

*In which areas of research is this method/approach currently being used and developed, and what are its merits?*


*What types of data and research interests fit well with this method/approach?*

# Interdisciplinary projects on newspapers

How can semantic indexing enhance the study and exploration of historical newspapers?

# The team

Estelle Bunout
Simon Clematide
Marten Duering
Maud Ehrmann
Andreas Fickers
Daniele Guido
Frédéric Kaplan
Peter Makarov
Matteo Romanello
Gerold Schneider
Paul Schroeder
Benoit Seguin
Phillip Stroëbel
Martin Volk
Thijs van Beek
Lars Wieneke

*+ historical advisors*
*+ associated historians*

# Objectives and research questions

1. How to adapt NLP tools to historical texts?

2. How to explore complex and vast amounts of data?

3. What is the impact of new tooling on digital scholarship?

# Target audience



Historians without specific digital literacy but curiosity to go beyond keyword search and the motivation to learn how to use a research tool.

SEARCH ARTICLES · SEARCH IMAGES · NGRAMS

GROUP BY ARTICLE ˅

ORDER BY RELEVANCE ˅   DISPLAY AS LIST TILES

"arnhem OR amhem(2 more)" ˅

front · armée · guerre · ennemi · n ˅   article ˅

French ✕

add keyword to search 🔍 ⚲

○ Frontpage   FIND SIMILAR WORDS ⓘ

**PUBLICATION DATE**
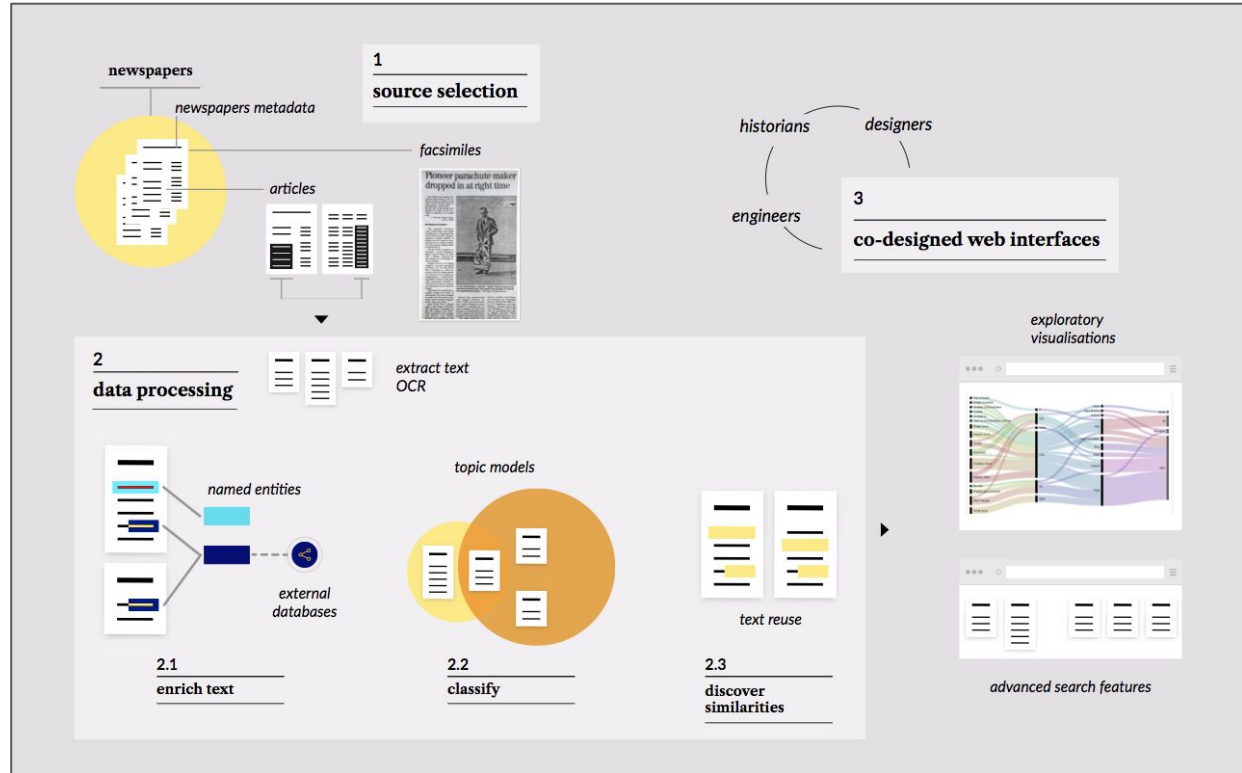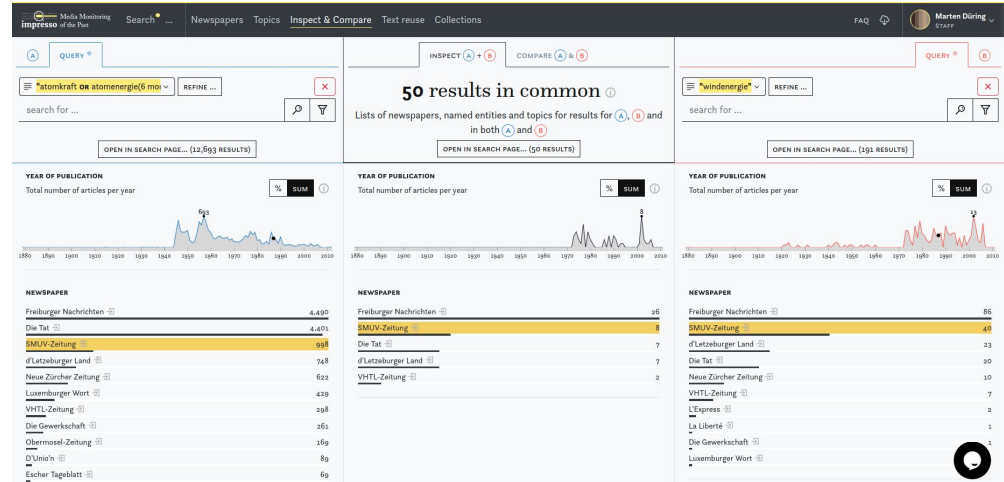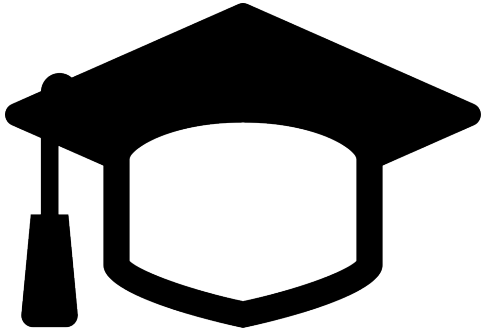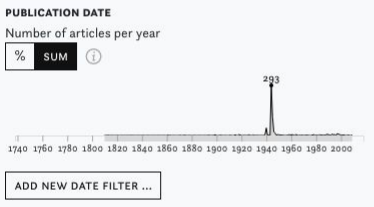Number of articles per year

% SUM ⓘ

293

1740 1760 1780 1800 1820 1840 1860 1880 1900 1920 1940 1960 1980 2000

ADD NEW DATE FILTER ...

**FILTER BY CONTENT LENGTH** ⓘ

65

0                    9,900

**FILTER BY LANGUAGE OF ARTICLES (1 OPTION)**   RESET
results are filtered when:

WRITTEN IN ˅
☑ French(689 results) ⧉

**FILTER BY NEWSPAPER TITLES (14 OPTIONS)** ⓘ   ⊖
check one or more newspaper to filter results

☐ Journal de Genève (141 results) ⧉
☐ L'Impartial (123 results) ⧉
☐ La Liberté (112 results) ⧉

---

**689** articles found containing **arnhem** or **amhem** or **arnheim** or **arnehm** - tagged as article; with topic FRONT · ARMÉE · GUERRE · ENNEMI · NORD; written in French

COMPARE ...   SAVE / EXPORT ˅   ☐

### La progression des Aîiiês en Allemagne

**Journal de Genève** ⧉ WEDNESDAY, APRIL 18, 1945 — p.2
Personal use

La progression des Aîiiês en Allemagne Q. o. allié, 17.-(AFP.) Le communiqué de mardi annonce : Des forces alliées débouchant de la région d'Arnhem on

LOCATIONS Claude Dallemagne, Zutphen, Gary Nord, Dessau, Leipzig, Halle, Saxony-Anhalt, Borna of Croatia, Altenburg, E. O. Plauen, Hof, Germany, Nuremberg, Heilbronn, Freudenstadt, Offenburg, Royan, Gironde, Berlin, Suisse, Moselle, Laufenburg, Germany, Lahr

| : Des forces alliées débouchant de la région d'Arnhem ont avancé au delà d'Otterloo. Elles ont fait également

Kriegsgriffe in CH ✕   Arnhem 1944, filtered by topics ✕

VIEW   ADD TO COLLECTION ... ˅

### Les forces alliées ont réalisé de nouveaux progrès à l'Ouest

**Journal de Genève** ⧉ SATURDAY, APRIL 7, 1945 — p.2
Personal use

Les forces alliées ont réalisé de nouveaux progrès à l'Ouest C. Q allié, 6. — (A. F. P.) Au nord de Niraègue. les patrouilles alliées ont traversé le

LOCATIONS Zutphen, Almelo, Diepholz, Petershagen Nord railway station, Warburg, Hamm, Dortmund, Fulda, Deventer, Hamelin, Berlin, Lingen, Germany, Rheine, Recklinghausen, Siegen, Gotha (town)

| de Niraègue. les patrouilles alliées ont traversé le Rhin Inférieur à l'ouest et à l'est d'Arnhem. Une vive

Kriegsgriffe in CH ✕

VIEW   ADD TO COLLECTION ... ˅

### DANS LE CORRIDOR ALLIÉ EN HOLLANDE Les Britanniques ont colmaté une ré[...]

**Journal de Genève** ⧉ TUESDAY, SEPTEMBER 26, 1944 — p.2
Personal use

DANS LE CORRIDOR ALLIÉ EN HOLLANDE Les Britanniques ont colmaté une récente brèche faite par les S. S. Auprès de la 2 ᵐᵉ armée britannique, 26. — (îer

LOCATIONS Eindhoven, Brussels, Turnhout, Berlin, Nancy Wake, Veghel, Deurne, Netherlands, Metz, Moselle
PEOPLE United Kingdom, Eisenhower Doctrine

| : Le corridor menant aux troupes portées d'Arnhem, qui avait été coupé, hier, par les Allemands, a de nouveau
dans la région d'Arnhem, où nous avons réussi à faire passer quelques renforts sur la berge septentrionale du Lek

« ‹ 1 2 3 4 ... › »

# Example codesign: Word embeddings



find similar words ×

Enlarge you search! Type **one word** and obtain a list of surrounding words

| atom | German ⇕ | 50 ⇕ |

*Click on one of the following words to update your search*

atom  atomkraft  nuklear  atomare  nukleare
atomaren  atomenergie  atomarer  atombombe
thermonuklearen  nuklearen  atombomben
wasserstoffbombe  plutonium  nuklearer  atomischen
wasserstoffbomben  atomund  atomares

Word embeddings?

Keyword suggestions?

OCR mistakes!

find similar words ×

Enlarge you search! Type **one word** and obtain a list of surrounding words

| nucleaire | French ⇕ | 50 ⇕ |

*Click on one of the following words to update your search*

nucleaire  nucieaire  nucleaires  nuclöaire  nueleaire
nucieaires  nuoleaire  thermonucleaire  hydrogene
nuclöaires  nucleai  atorniques  cleaire  reacteur
thermonucleaires  reacteurs  nucliaire  lhydrogene
cleaires  megatonnes  teleguides  lenergie

# Add similar

# Ngrams

# Filter by topic

# Inspect & Compare
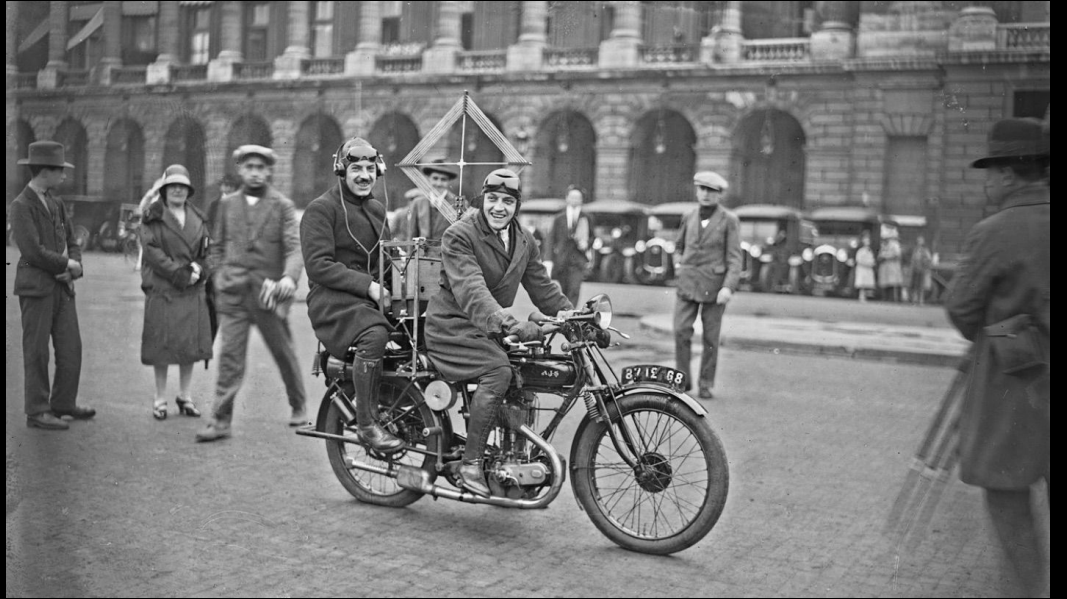
Gianni Sarconi, The Master of Numbers (2006)


Jennifer Kuhns, Glass mosaic of Alice Paul, suffragist. 1885-1977 (2019)

*How could this method/approach contribute to the study of historical processes of knowledge evolution?*

Outlook: What is next?

# Text Reuse at Scale

Opportunity: integration of TR data with other semantic enrichments.

1. Expert workshop in Nov 2022 produced research objectives, tasks & mockups.

2. New prototype component.

3. User evaluation.

4. Research paper [preprint].

5. Revised prototype to be integrated in main app.



b) Overview tab: Distribution of metadata and enrichments in a set of clusters.

c) Statistics tab: Distribution of Text Reuse measures over time and newspapers.

a) Search and filters.

d) Passages tab: Passages within cluster c62714.

e) Passages tab: Comparative view.

a) Search and filters.

b) Overview tab: Distribution of metadata and enrichments in a set of clusters.

c) Statistics tab: Distribution of Text Reuse measures over time and newspapers.

d) Passages tab: Passages within cluster c62714.

e) Passages tab: Comparative view.

# Example: When did nuclear power ads circulate?

# Editing press agency content



**COMPARE TEXT REUSE PASSAGES**     ✕

MONDAY, NOVEMBER 28, 1910 "M. Winston Churchill, ministre du commerce, pronon..." TEXTREUSEPASSAGE

Compare the passage below     with # [ 2 ] of 3 passages   BY DATE (DESC) ⌄

lia crise britaiiiilquc.
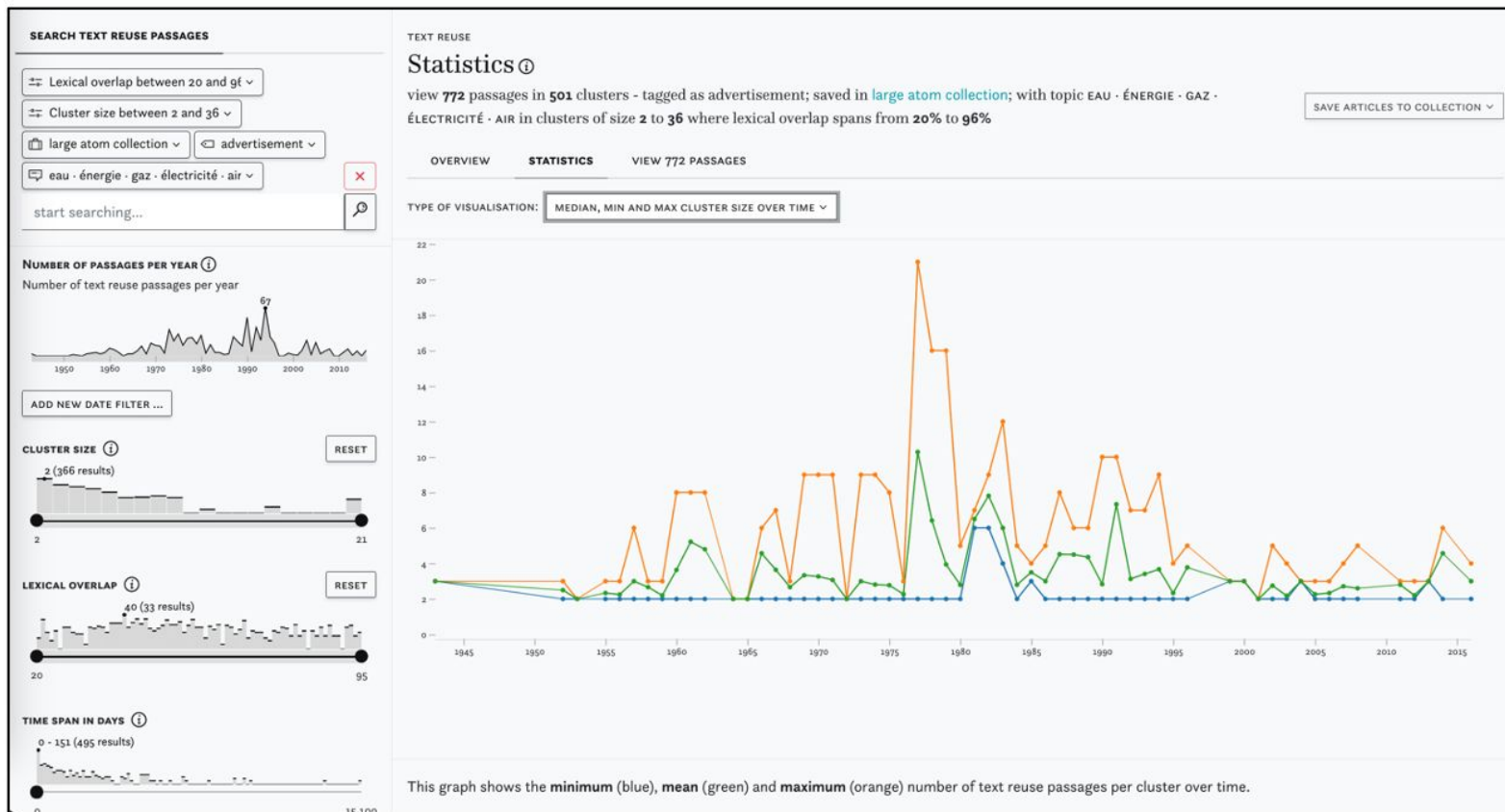**Gazette de Lausanne** ⇥ MONDAY, NOVEMBER 28, 1910 – P.2

Les suffragistes font de l'action directe
**L'indépendance luxembourgeoise** ⇥ TUESDAY, NOVEMBER 29, 1910 – P.2

M. Winston Churchill, ministre du commerce, prononçait ven- dredi soir à Bradford a été coupé de plu- sieurs interruptions de suffragistes et de suffragettes qui ont été expulsés à tour de rôle. Comme l'orateur rentrait de Bradford à Londres, il fut attaqué dans le train par un individu, qui a essayé de le frapper avec une cravache en disant : « Voilà pour toi, chien ! » Deux agents de police parèrent le coup et s'emparè- rent de l'homme après une lutte violente. A la gare de Londres, trois femmes ont également essayé de frapper M. Churchill ; elles en ont été empêchées par les agents.

M. Winston Churchill, revenant à Londres, après avoir prononcé un dis- cours à Bradford, où les suffragettes et leurs partisans avaient déjà violemment manifesté contre lui, a été attaqué dans le train par un individu, qui a essayé de le frapper avec une cravache en disant: «Voilà pour toi, chien!» Deux agents de police, qui accompagnaient M. Churchill, parèrent le coup, et s'em- parèrent de l'homme, après une lutte violente. On croit que l'assaillant est un des suffragistes expulsés de la réunion où M. Winston Churchill venait de parler. A la gare de Londres, trois femmes ont également essayé de frapper M. Chur- chill ; elles en ont été empêchées par les agents.
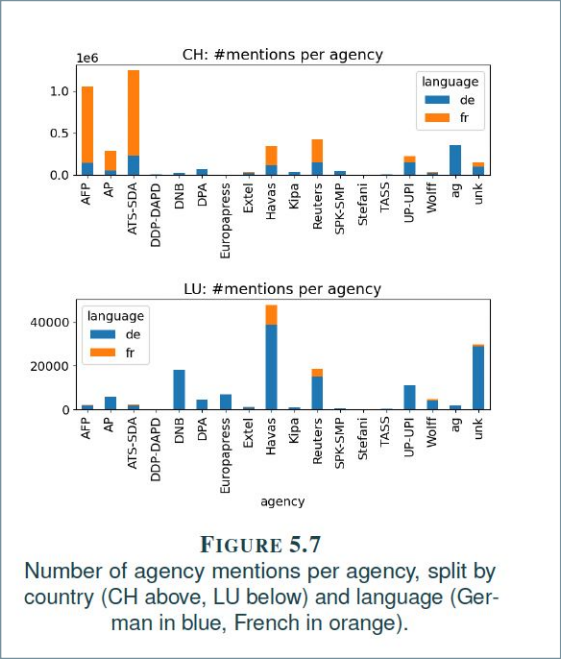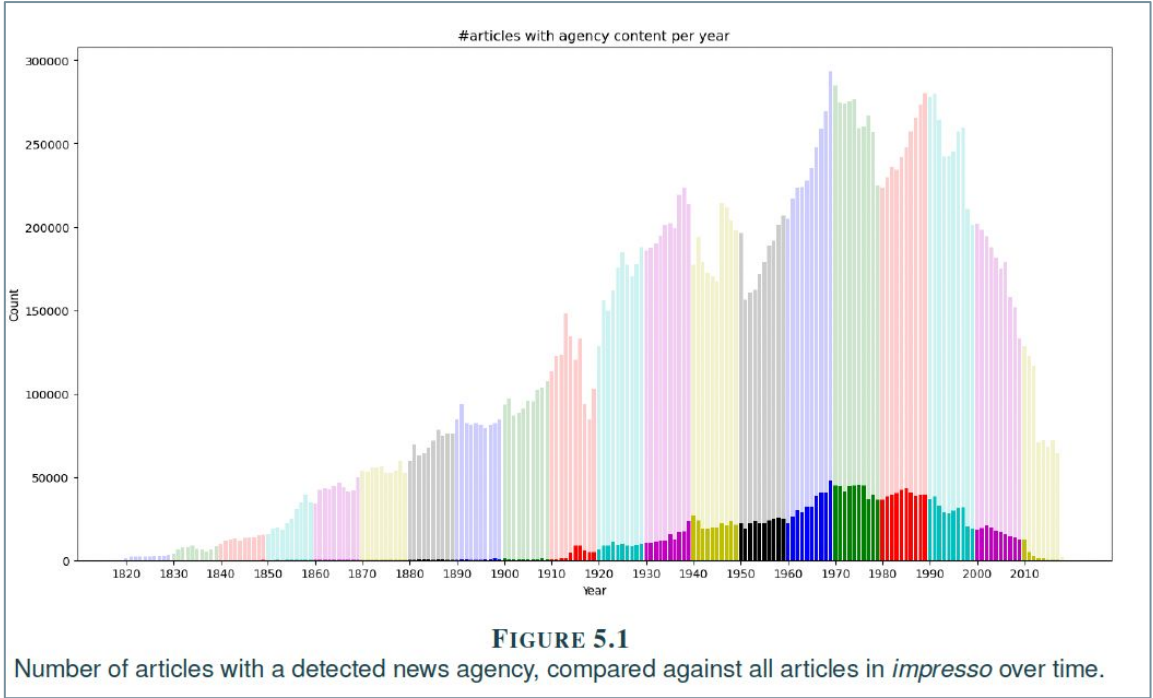
# Where did the news come from?

## News Agency Recognition

- Construction of an annotated dataset (27 agencies, ca 2000 articles, fr and de);

- Training and evaluation of models to recognise news agencies;

- Application on the whole *impresso* corpus;

- First analyses.

# News agency mentions



**FIGURE 5.1**
Number of articles with a detected news agency, compared against all articles in *impresso* over time.



**FIGURE 5.7**
Number of agency mentions per agency, split by country (CH above, LU below) and language (German in blue, French in orange).

# In Swiss newspapers



**FIGURE 5.8**
News agency mentions in Swiss newspapers over time, split by language (French above, German below).
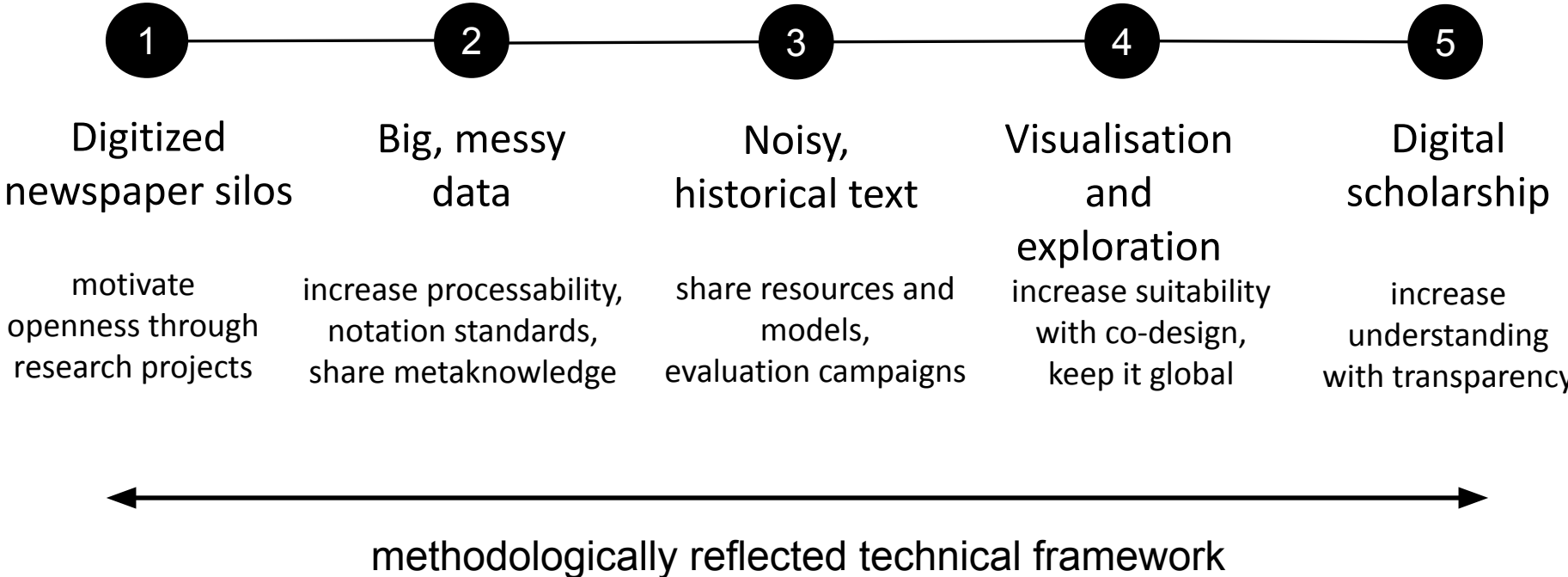
# In Luxembourgish newspapers (1940-1944)



**FIGURE 5.15**

The development of agency mentions in Luxembourg in the years 1930-1949, split by language (French above, German below). The bars show the distribution of the different agency mentions (left y-axis), while the blue line indicates the general share of articles with agency mentions in the Luxembourgish corpus (right y-axis).

# Reliable Semantic Indexing of Historical Newspapers at Scale: Are We There Yet?

**1** — Digitized newspaper silos

motivate openness through research projects

**2** — Big, messy data

increase processability, notation standards, share metaknowledge

**3** — Noisy, historical text

share resources and models, evaluation campaigns

**4** — Visualisation and exploration

increase suitability with co-design, keep it global

**5** — Digital scholarship

increase understanding with transparency

methodologically reflected technical framework

# What's next ? (from the community)

What are the main challenges we need to address in relation with historical newspapers?

- Document processing
- Text and image processing
- **Evaluation** (digitisation and content mining)
- Exploration of enriched, **global** collections
- **Working with data**
- **Workflows**
- **Criticism, inclusivity**
- Legal matters



Report from Dagstuhl Seminar 22292

Computational Approaches to Digitised Historical Newspapers

Edited by
Maud Ehrmann[1], Marten Düring[2], Clemens Neudecker[3], and Antoine Doucet[4]

1 EPFL - Lausanne, CH, maud.ehrmann@epfl.ch
2 University of Luxembourg, LU, marten.during@uni.lu
3 Staatsbibliothek zu Berlin, DE, clemens.neudecker@sbb.spk-berlin.de
4 University of La Rochelle, FR, antoine.doucet@univ-lr.fr

— Abstract —
Historical newspapers are mirrors of past societies, keeping track of the small and great history and reflecting the political, moral, and economic environments in which they were produced. Highly valued as primary sources by historians and humanities scholars, newspaper archives have been massively digitised in libraries, resulting in large collections of machine-readable documents and, over the past half-decade, in numerous academic research initiatives on their automatic processing. The Dagstuhl Seminar 22292 "Computational Approaches to Digitised Historical Newspaper" gathered researchers and practitioners with backgrounds in natural language processing, computer vision, digital history and digital library involved in computational approaches to historical newspapers with the objectives to share experiences, analyse successes and shortcomings, deepen our understanding of the interplay between computational aspects and digital scholarship, and discuss future challenges. This report documents the program and the outcomes of the seminar.

# What's next ? (from impresso)



Beyond Borders: Connecting Historical Newspapers and Radio

**Main objectives**
- Enrichment and integration of **newspaper and radio sources in a single semantic space;**
- Expand the corpus to Western Europe;
- Develop interfaces for exploratory and computational research;
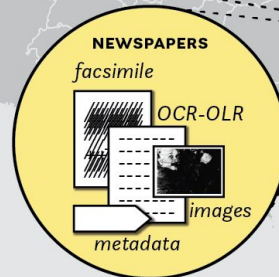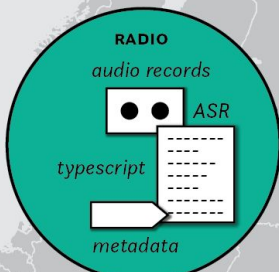- Conduct case studies in (media) history, theme "influences."

09/2023 - 02/2027

# Source collection
**1**

*European media archives*

AUSTRIA

BELGIUM

FRANCE

GERMANY

LUXEMBOURG

THE NETHERLANDS

SWITZERLAND

UK

**RADIO**
*audio records*
*ASR*
*typescript*
*metadata*

**NEWSPAPERS**
*facsimile*
*OCR-OLR*
*images*
*metadata*

# Media processing
**2**

*Enriching & connecting*

SEMANTIC ENRICHMENT
ACROSS LANGUAGES
ACROSS MEDIA

*dense vector representations*

*topic*
*place*
*person*
*institution*

*interview*
*advertisement*
*radio schedule*

**EXTERNAL KNOWLEDGE**

# Media exploration
**3**

*Connected and comparable enriched media sources*

**PROJECT DATA**
*documents*
*topics*
*entities*

*process data using impresso API*

**EXTERNAL DOCUMENTS**

*compatible enrichments*

**API**

**USER-ORIENTED API**

COMPATIBLE DATA

*community of historians*

**IMPRESSO WEB APP**

*historians*
*designers*
*engineers*

**IMPRESSO DATA LAB**
*executable notebooks*

CO-DESIGNED INTERFACES

CO-DEVELOPED METHODS

# For more information



impresso website (update soon)

data on Zenodo

code on GitHub

impresso youtube channel

impresso-project.ch/app