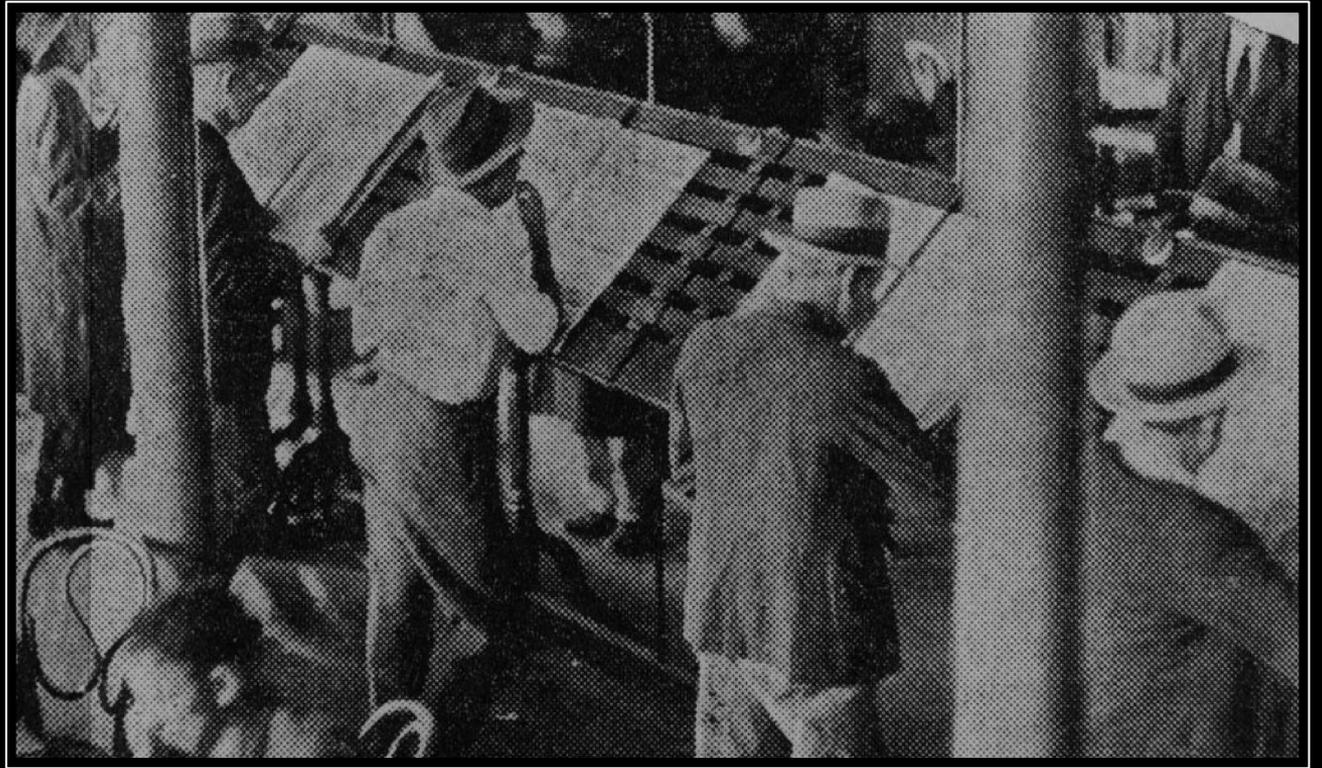# impresso Text Reuse at Scale
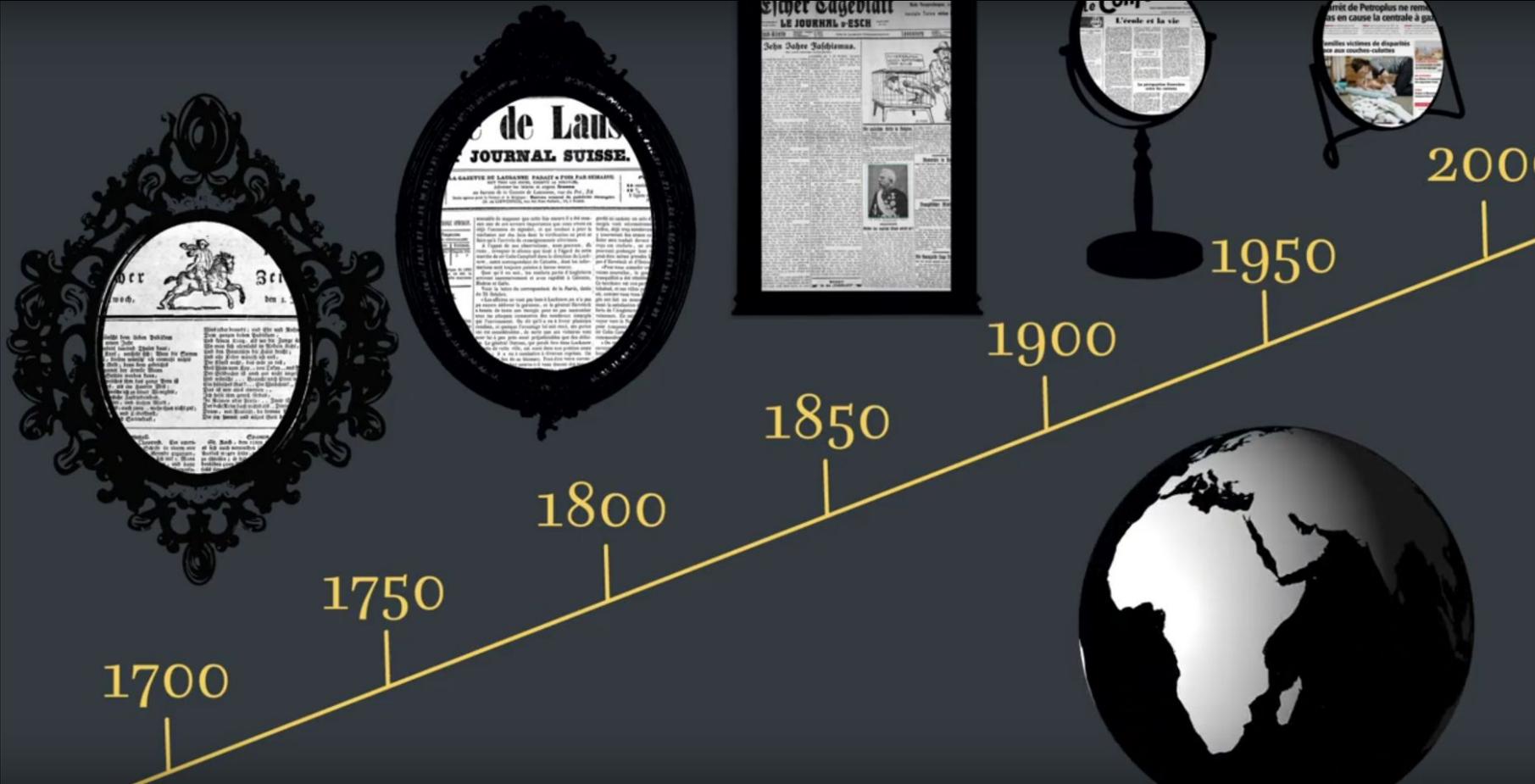
Kampagnen, Wiederverwertung und Plagiate in historischen Zeitungsammlungen

Marten Düring & impresso team

1. Ein paar Worte über historische Zeitungen

## Haute tension sur l'avenir électrique

par Nicolas Hirtzberer

### L'ère des incertitudes

## Christophe Babaiantz: le prix du renoncement à Kaiseraugst

# Forschungsinteressen: Geschichte der Medien und Medien als Quellen

**Fachdisziplinen**

Mediengeschichte

Gender studies

Kulturgeschichte

Sozialwissenschaften

Sozialgeschichte

…

**Ziele**

Layoutanalyse

Evolution von Genres

"Viralität"

Soziale Normen

Meinungsbildung

Wissenshorizonte

Biographien

Ernährung

….

**Zeitungselemente**

Werbung

Kolumnen

Agenturnachrichten

Bilder

Kleinanzeigen

Radioprogramme

Todesanzeigen

….

# The big tip of a hidden iceberg



Digitized Newspaper Coverage Worldwide (1800-2015)

Centre for Research Libraries. 2015. *The "State of the Art". A Comparative Analysis of Newspaper Digitization to Date*.

# Historische Zeitungen als Herausforderung

1. Institutionelle Silos

2. Big and messy data

3. Noisy historical text

4. Visualisierung und Exploration

5. Digitale Forschungskultur

Impresso

Critical content mining of 200 years of historical newspapers

Inwieweit dienen semantische Anreicherungen der Analyse und Exploration historischer Zeitungen?

# Das 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*

*+ beratende HistorikerInnen*
*+ assoziierte WissenschaftlerInnen*

engineer

web dev

designer / web dev

NLP/DH

(digital) historian

NLP/DH

NLP

NLP

(digital) historian

designer / web dev

(digital) historian

web dev

DH/robotics

# Ziele und Forschungsfragen

1. Wie passt man NLP-Werkzeuge an historischen Text an?

2. Wie erforscht man große und komplex Datenbestände?

3. Welche Rückwirkungen hat dies auf historische Forschung?

# Impressos Zielgruppe



HistorikerInnen ohne besondere digital Vorbildung aber mit der Neugier, Stichwortsuchen zu überwinden und der Motivation, ein neues Tool kennenzulernen.

Gianni Sarconi, The Master of Numbers (2006)



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

# Impresso II

09/2023 - 02/2027



impresso
Media
Monitoring
of the Past

Beyond Borders: Connecting Historical Newspapers and Radio

**Ziele**
- Anreicherung und Integration von **Zeitungs- und Radioquellen in einem semantischen Raum;**
- Korpuserweiterung auf Westeuropa zusammen mit 20 Projektpartnern;
- Interfaces für die Exploration und datengetriebene Analyse;
- Fallstudien in (Medien) Geschichte mit dem Leitthema "influences."

# Partner

***National or state libraries (holding digitised newspaper collections)***

Bibliothèque Nationale Suisse, BN
Bibliothèque Nationale du Luxembourg, BNL
Österreichische Nationalbibliothek, ONB
Staatsbibliothek zu Berlin, SBB
The British Library (BL)
Bibliothèque nationale de France, BnF
Staats- und Universitätsbibliothek Hamburg, HUB
Bibliothèque royale de Belgique/Koninklijke Bibliotheek van België, KBR
Koninklijke Bibliotheek, KB

***Newspapers***

Le Temps
Neue Zürcher Zeitung

***Audiovisual heritage institutions and archives (holding digitised radio collections)***

Radio Television Suisse (RTS)
Österreichischer Rundfunk, ORF
British Broadcasting Corporation (BBC)
DeutschlandRadio
Institut National de l'Audiovisuel, INA
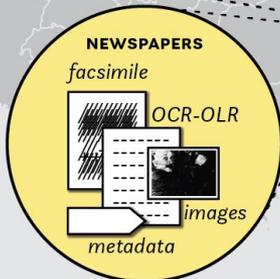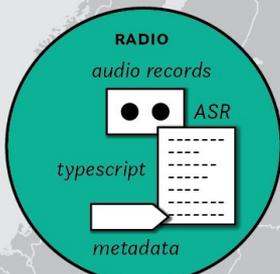Nederlands Instituut voor Beeld en Geluid, NISV

***Research Networks***

Entangled Media Histories Research Network for European media historians (EMHIS)
Memoriav, the Swiss network for audiovisual cultural heritage preservation
infoclio.ch

# Source collection

**1**

*European media archives*

AUSTRIA

BELGIUM

FRANCE

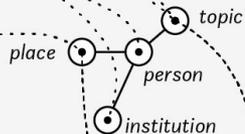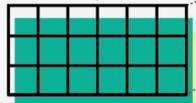GERMANY

LUXEMBOURG

THE NETHERLANDS

SWITZERLAND

UK

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

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

# Media processing

**2**

*Enriching & connecting*

SEMANTIC ENRICHMENT
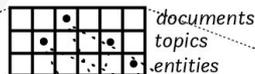ACROSS LANGUAGES
ACROSS MEDIA

*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*

**USER-ORIENTED API**

**API**

COMPATIBLE DATA

*community of historians*

**IMPRESSO WEB APP**

*historians*
*designers*
*engineers*

**IMPRESSO DATA LAB**
*executable notebooks*

CO-DESIGNED INTERFACES

CO-DEVELOPED METHODS

# 3. Text reuse

# Über Text Reuse

Text reuse ist die "the meaningful reiteration of text, usually beyond the simple repetition of common language" (Romanello et al. 2014).

Warum text reuse?

- Um festzustellen, ob eine digitale Bibliothek mehrere Ausgaben desselben Werks bzw. derselben Werke enthält.
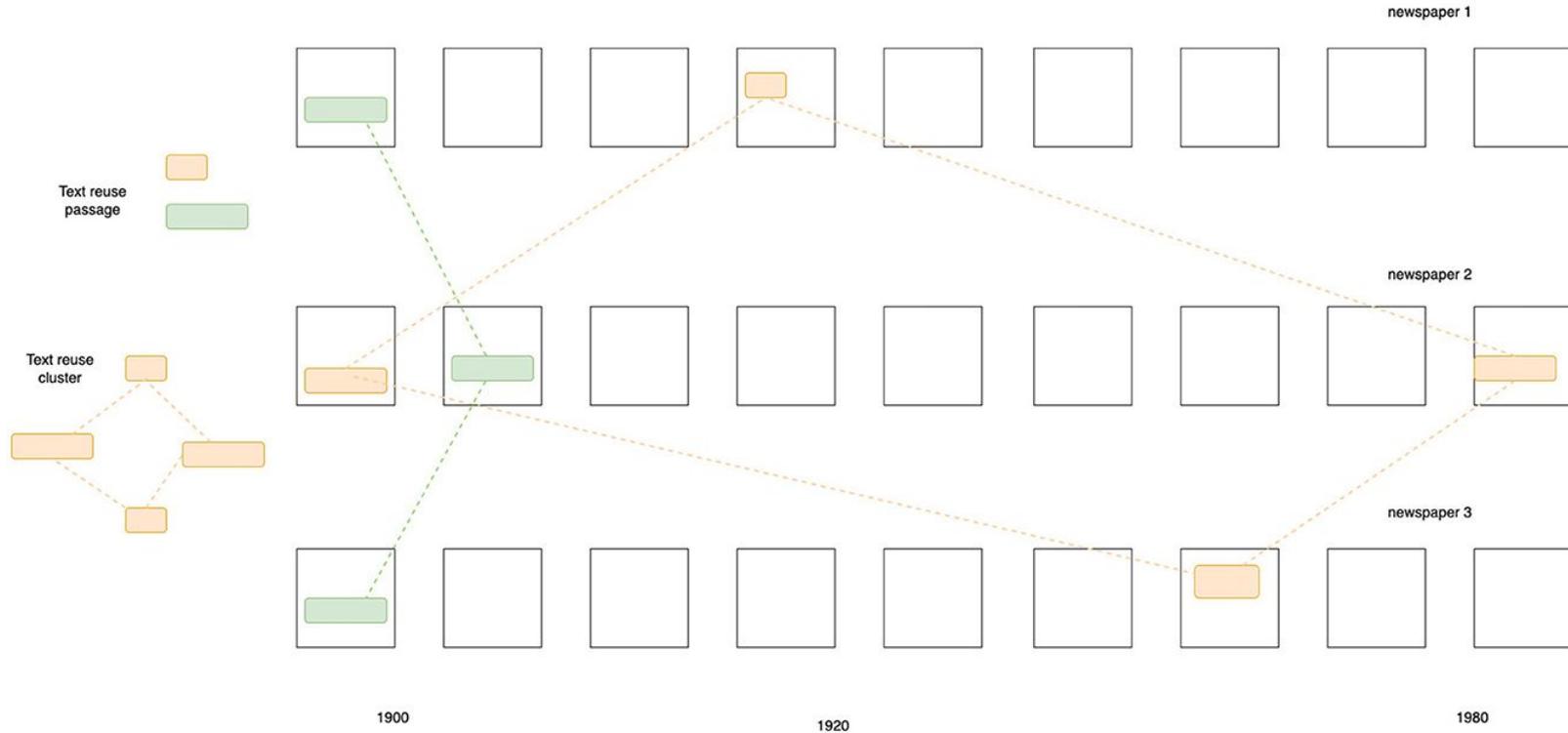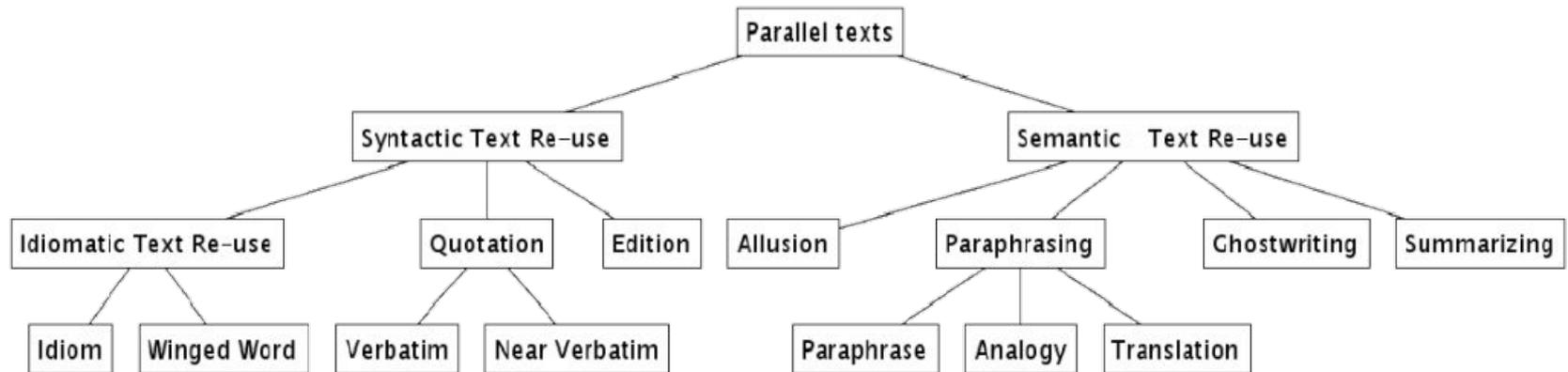
- Um Zitate in einem Text zu finden, vorausgesetzt, die Zielwerke sind bekannt (z. B. Zitate aus der Bibel in der englischen Literatur des 17. Jahrhunderts finden) Untersuchung der Viralität und Verbreitung von Texten (z. B. Viral Texts von Cordell und Smith für historische Zeitungen)

- Um Duplikate innerhalb einer Textsammlung zu identifizieren (und möglicherweise herauszufiltern), bevor weitere Verarbeitungsschritte durchgeführt werden (z. B. topic modeling, wie von Schofield et al. (2017) dargestellt)

Matteo Romanello and Simon Hengchen, "Detecting Text Reuse with Passim," *Programming Historian* 10 (2021), https://doi.org/10.46430/phen0092.

# Text reuse Passagen und Cluster

# Typen von Text Reuse



Parallel texts
- Syntactic Text Re-use
  - Idiomatic Text Re-use
    - Idiom
    - Winged Word
  - Quotation
    - Verbatim
    - Near Verbatim
  - Edition
- Semantic Text Re-use
  - Allusion
  - Paraphrasing
    - Paraphrase
    - Analogy
    - Translation
  - Ghostwriting
  - Summarizing

Marco Büchler, Historical text reuse: what is it?,
https://www.etrap.eu/historical-text-re-use/

# Typen von Text Reuse

Rosson et al. Reception Reader: Exploring Text Reuse in Early Modern British Publications, 2023,
https://openhumanitiesdata.metajnl.com/articles/10.5334/johd.101

| TYPE OF REUSE | EXAMPLES | POSSIBLE RESEARCH QUESTIONS |
|---|---|---|
| Quotes | Latin, biblical, famous quotes. | What was the process of quotes from Lucretius becoming epigraphs over time? |
| Reprints of longer passages | Reused sections or fragments from essays or treatises appearing in works by different authors. | What was the distribution of Hume's essays outside of his published works? |
| Modified reuse | Modified reuse of a specific work in another work. | How did Clarendon's *History of the Rebellion* feature in other historical works? |
| Verse reprints | Reprinting of poetry in unexpected or uncommon locations. | How did Dryden's poetry spread outside of known collections? |
| Unattributed reuse | Hidden or obscured reuse of texts. | Can we gain a broader understanding of the reception of Hume's essays by exploring their use in other works without proper attribution? |
| Artefacts | Imprint of publisher, advertisement. | What was the distribution pattern of advertising for Hume's *Treatise* in printed books in the eighteenth century? |

# Beispiele: Job-Anzeigen - intendierte Zirkulation

# Beispiele: Publikation und Korrektur (modified reuse)



COMPARE TEXT REUSE PASSAGES ✕

WEDNESDAY, OCTOBER 13, 1982 "Le solde sera acquis en Italie, en France et en Su..." TYPES_TEXTREUSEPASSAGE

Compare the passage below with # 2 of 5 passages BY DATE (DESC)

Assez d'électricité cet hiver
**L'Express** ⊡ WEDNESDAY, OCTOBER 13, 1982 – P.2

Rectificatif de l'ATS
**Journal de Genève** ⊡ THURSDAY, OCTOBER 14, 1982 – P.8

Le solde sera acquis en Italie, en France et en Suisse alémanique, ce qui représente environ dix millions de francs sur une année.

le solde sera acquis en Italie, France et Suisse alémanique, ce qui représente environ 100 millions de francs sur une année,

## Rectificatif de l'ATS

Dans la nouvelle intitulée « Electricité : assez de courant pour cet hiver », publiée dans nos éditions d'hier, l'Agence télégraphique suisse s'est trompée dans les zéros. A la fin du deuxième paragraphe, il faut lire que « le solde sera acquis en Italie, France et Suisse alémanique, ce qui représente environ 100 millions de francs sur une année, et non pas 10 millions comme noté par erreur). (Réd.)

# Beispiele: Nachdrucke zu Anlässen (reprints)

# Beispiele: Zitate (quotes)

**Left panel:**

COMPARE TEXT REUSE PASSAGES ✕

THURSDAY, MAY 26, 1966 "Bundesverfassung heisst es : « Im Namen Gottes des..." TYPES_TEXTREUSEPASSAGE

Compare the passage below — with # [1] of 6 passages — BY DATE (DESC)

Was meint Exuperantius?
**Die Tat** ⊡ THURSDAY, MAY 26, 1966 – P.15

für «junge
**Freiburger Nachrichten** ⊡ SATURDAY, JANUARY 13, 1996 – P.20

Bundesverfassung heisst es : « Im Namen Gottes des Allmächtigen ! Die Schweizerische Eidgenossenschaft , In der Ab- sicht , den Bund der Eidgenossen zu befestigen , die Einheit , Kraft und Ehre der schweizerischen Nation zu erhalten und zu fördern , hat nachste- hende Bundesverfassung angenommen — . » Und hier haben Sie den offiziellen Titel unseres schweizerischen Staatswesens : Schweizerische Eidgenossenschaft .

der Bun- desverfassung : « Im Namen Gottes des Allmächtigen ! Die Schweizerische Eid- genossenschaft , in der Absicht , den Bund der Eidgenossen zu festigen , die Einheit , Kraft und Ehre der schweizeri- schen Nation zu erhalten und zu för- dern , hat nachstehende Bundesverfas- sung angenommen . »

**Right panel:**

COMPARE TEXT REUSE PASSAGES ✕

SATURDAY, OCTOBER 5, 1996 "de la Bible (Luc 10 : 27 : « Tu aimeras le Seigneu..." TYPES_TEXTREUSEPASSAGE

Compare the passage below — with # [1] of 43 passages — BY DATE (DESC)

Sachons être confiants en Dieu
**L'Express** ⊡ SATURDAY, OCTOBER 5, 1996 – P.23

Choisir son camp
**L'Express** ⊡ TUESDAY, JANUARY 13, 2015 – P.27

de la Bible (Luc 10 : 27 : « Tu aimeras le Seigneur, ton Dieu, de tout ton cœur, de toute ton âme, de toute ta force, et de toute ta pen- sée ; et ton prochain comme toi- même ») et que l'on

« Tu aimeras le Sei- gneur, ton Dieu, de tout ton cœur, de toute ton âme, de toute ta force, et de toute ta pensée ; et ton prochain comme toi-même. » (Luc 10.27)

# Text reuse and historical research objectives

1. (Trans-) national media ecosystems

2. Newspaper content as *bricolage*

3. Historicising virality

4. Tracing historical events

5. Capturing historical Zeitgeist

# (Trans-) nationale Medien-Ökosysteme

*Goal: Highlight the connectivity of historical media across borders and languages. Which ideological, commercial, and financial structures made this work and how did they shape media?*

*TR: (How) does content flow through the international network of newspapers?*

Focus e.g. on cut-paste practices, the relevance of information infrastructure, geography, cross-border cultural affinities.



Beelen, K. Digitising newspapers press directories to understand the landscape of historical newspapers

- Smith, D. A., Cordell, R., and Mullen, A. (2015). Computational methods for uncovering reprinted texts in antebellum newspapers. *Am. Liter. Hist.* 27, E1–E15. doi: 10.1093/alh/ajv029
- Keck, J., Oiva, M., and Fyfe, P. (2022). Lajos kossuth and the transnational news: a computational and multilingual approach to digitized newspaper collections. *Media History* 29, 287–304. doi: 10.1080/13688804.2022.2146905
- Salmi, H., Rantala, H., Vesanto, A., and Ginter, F. (2019). "The long-term reuse of text in the finnish press, 1771–1920," in *Proceedings of the Digital Humanities in the Nordic Countries 4th Conference*, eds. C. Navarretta, M. Agirrezabal, and B. Maegaard (Copenhagen, Denmark: CEUR Workshop Proceedings), 253–273.
- Paju, P., Salmi, H., Rantala, H., Lundell, P., and Marjanen, Vesanto, A. (2022). "Textual migration across the baltic sea: Creating a database of text reuse between Finland and Sweden," in *Proceedings of the 6th Digital Humanities in the Nordic and Baltic Countries Conference (DHNB 2022), CEUR Workshop Proceedings*, eds. K. Berglund, M. La Mela, and I. Zwart (Aachen: CEUR-WS.org), 361–369.
- Beelen, K. Digitising newspapers press directories to understand the landscape of historical newspapers, 2023, https://livingwithmachines.ac.uk/digitising-newspapers-press-directories-to-understand-the-landscape-of-historical-newspapers/.
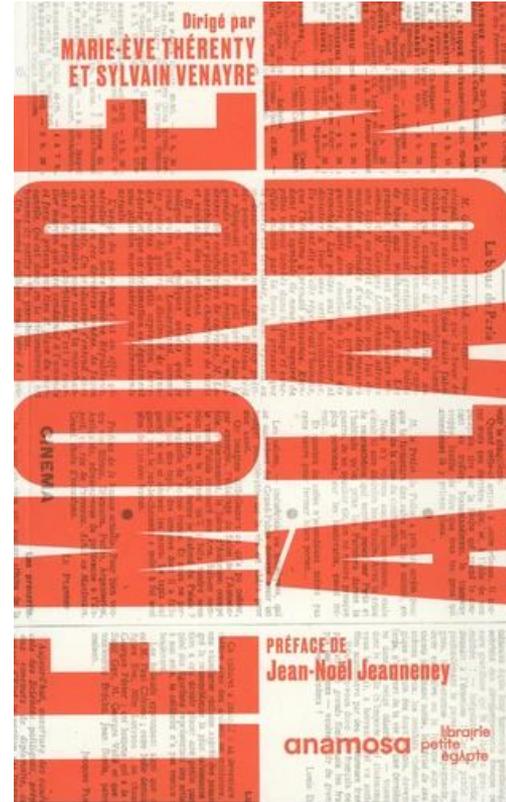
# Newspaper content as bricolage

*From which textual resources were newspapers compiled?*

*TR: How did formulaic ways of representing information emerge and spread (e.g. weather reports)? Can we determine the origins of content? Which content was influential?*

Focuses on parallel developments across titles and the emergence of common practices.

- Walma, L. W. B. (2015). Filtering the "news:" Uncovering morphine's multiple meanings on delpher's dutch newspapers and the need to distinguish more article types. *Tijdschrift voor Tijdschriftstudies*. 38, 61–78. doi: 10.18352/ts.345
- Thèrenty, M.-E., and Venayre, S. (2021). *Le monde à la une. Une histoire de la presse par ses rubriques*. Anamosa, illustrated èdition edition. doi: 10.3917/anamo.there.2021.02

# Historicising virality

*Which conditions (geography, type of information, infrastructure) enable rapid dissemination?*

Focus is on the breadth and speed with which content spreads. Pioneering work of Paju et al. who define a virality score based on the number of titles within a cluster, the number of unique printing locations, and the distance in days between the first and last passage publication date.

Paju, P., Salmi, H., Rantala, H., Lundell, P., and Marjanen, Vesanto, A. (2022). "Textual migration across the baltic sea: Creating a database of text reuse between Finland and Sweden," in *Proceedings of the 6th Digital Humanities in the Nordic and Baltic Countries Conference (DHNB 2022), CEUR Workshop Proceedings*, eds. K. Berglund, M. La Mela, and I. Zwart (Aachen: CEUR-WS.org), 361–369.

## Textual Migration Across the Baltic Sea: Creating a Database of Text Reuse Between Finland and Sweden

Petri Paju[1], Hannu Salmi[1], Heli Rantala[1], Patrik Lundell[2], Jani Marjanen[3] and Aleksi Vesanto[1]

[1] University of Turku, Department of Cultural History, Turku, FI-20014, Finland
[2] Örebro University, School of Humanities, Örebro, SE-70182, Sweden
[3] University of Helsinki, Department of Digital Humanities, Helsinki, FI-00014, Finland

**Abstract**

In this paper, we present a database and an interface on text reuse between newspapers and journals published in the Swedish language in Sweden and Finland during the 1645–1918 time frame. Using two national, digital newspaper collections, we detected their textual similarities with a computational method to study the textual migration, i.e., information flows, between the two countries. For purposes of this project, we developed a database of detected clusters of text reuse and an online interface to search, examine and analyse the transnational movement of information. The database, *Text Reuse in the Swedish-language Press, 1645–1918*, is accessible online and includes texts from over 1,100 newspapers and journals published at approximately 150 locations at various times during the 274-year time frame.

## 1. Introduction

This short paper presents a database and an interface on text reuse among Swedish-language newspapers and journals during the 1645–1918 time frame. The database and interface were built as part of the project, *Information flows across the Baltic Sea: Swedish-language press as a cultural mediator, 1771–1918*. The database and the accompanying project aim to study information flows, particularly between Sweden and Finland from the period when present-day Finland was part of the Swedish kingdom to the establishment of Finland as a Grand Duchy in the Russian Empire after 1809 and until the Independence of Finland in 1917 and Civil War in 1918. Even after the 1809 separation, news and other texts circulated because of the common cultural heritage and shared language, i.e., Swedish. The border was relatively easy to cross, and newspapers circulated between Sweden and Finland regularly. However, because their national histories eventually diverged, these press materials have been preserved, processed, and siloed in two national libraries. Still, print media digitisation makes it possible to study overlaps in large collections of texts and see how information was spread across the Baltic Sea.

In our project, textual migration was traced using a method based on the software BLAST, which can be applied to text-reuse detection. With the method, we detected every text passage with 300 or more characters of similarity and combined these passages into reuse clusters. We included Swedish-language papers published in Sweden and Finland, but excluded the Finnish-language press in Finland, as textual migration within Finland has been studied in previous publications [1, 2]. To strengthen the

# Tracing historical events

*How do individual media titles situate events in the political, economic, social, and cultural context relevant to them? How does this affect the perception of these events by their audiences?*

Focus is on the appropriation of often global news on the local scale.



**FIGURE 2.**
The bar plot of the number of reprinted articles mentioning Kossuth in OcEx corpus shows that the American tour established Kossuth as an international celebrity.

Keck, J., Oiva, M., and Fyfe, P. (2022). Lajos kossuth and the transnational news

- Oiva, M., Nivala, A., Salmi, H., Latva, O., Jalava, M., Keck, J., et al. (2020). Spreading news in 1904. *Media History* 26, 391–407. doi: 10.1080/13688804.2019.1652090
- Keck, J., Oiva, M., and Fyfe, P. (2022). Lajos kossuth and the transnational news: a computational and multilingual approach to digitized newspaper collections. *Media History* 29, 287–304. doi: 10.1080/13688804.2022.2146905

# Tracing historical events

# Capturing historical Zeitgeist

*Historical media partially capture the attitudes, norms, beliefs, moods and feelings of past generations, or Zeitgeist.*

Focus is on texts which were produced independently but still share certain characteristics.

- Verheul, J., Salmi, H., Riedl, M., Nivala, A., Viola, L., Keck, J., et al. (2022). Using word vector models to trace conceptual change over time and space in historical newspapers, 1840–1914. *Dig. Human. Quart*. 16, 7445. Available online at: https://www.digitalhumanities.org/dhq/vol/16/2/000550/000550.html
- Paasikivi, S., Salmi, H., Vesanto, A., and Ginter, F. (2022). Infectious media: Cholera and the circulation of texts in the finnish press, 1860–1920. *Media Hist*. 29, 17–38. doi: 10.1080/13688804.2022.2054408
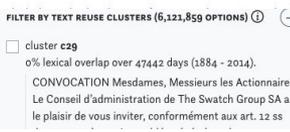
# Tasks to support research objectives

| Task | Title | Level | Support |
|------|-------|-------|---------|
| 1 | Obtain an overview of text reuse in a corpus, collection or query | Corpus | Yes |
| 2 | Obtain an overview of a single cluster | Cluster | Yes |
| 3 | Compare passages | Passage | Yes |
| 4 | Compare clusters | Cluster | Yes |
| 5 | Identify different types of text reuse | Corpus | Yes |
| 6 | Generate research corpora based on text reuse clusters | Corpus | Yes |
| 7 | Identify connections | Corpus | Partial |
| 8 | Detect and trace virality | Corpus | No |
| 9 | Search for passages | Passage | No |
| 10 | De-duplicate content | Corpus | No |
| 11 | Export of text reuse data | All | Planned |

# Temporalities in text reuse data

| Type | Description | Measures | Examples |
|------|-------------|----------|----------|
| Duration | The time period which is covered by a cluster ranging from the earliest to the latest publication date of individual passages. | Publication date | Paju et al.'s notions of fast and slow text reuse fall into this category. |
| Virality | The speed (measured in days) and breadth of text reuse passages spreading within a corpus. Speed corresponds to time passed (e.g., days) whereas breadth corresponds to the number of publications which contain a passage at a given point in time. | Publication date, number of publications | News of the sinking of the Titanic or the destruction of the Hindenburg Zeppelin traveled around the world within days or weeks. |
| Rhythm | Pattern with which text reuse passages appear over time. | Distance between publication dates | Reprints of articles on the occasion of their anniversary, e.g., on the occasion of the bombing of Hiroshima. |

# Capture the characteristics of text reuse through filters

| Measure | Description | Implementation in interface prototype | |
|---|---|---|---|
| Passages per year | Number of passages counted in a given year. | Line chart which displays the count of passages per year for a given query or filter operation. This gives a first indication, during which years text reuse occurred more commonly. Time sliders and precise date entry allow users to filter for exact date ranges to inspect. | **NUMBER OF PASSAGES PER YEAR** ⓘ<br>Number of text reuse passages per year<br>*1815*<br>1800 1820 1840 1860 1880 1900 1920 1940 1960 1980 2000 |
| Cluster size | The number of passages contained in a cluster. | Histogram which shows the distribution of text reuse cluster sizes and indicates the highest score. The histogram groups clusters of size n and displays their sum. This gives a first indication of averages as well as outliers. Sliders can be used to specify a cluster size range of interest. Filtering by cluster size allows to exclude or explicitly focus on outliers but different cluster sizes may also correspond to different types of content. | **CLUSTER SIZE** ⓘ<br>How to read histograms ⓘ<br>2 - 453 (15,756,994 results)<br>2　　　　45,553 |
| Lexical overlap | The percentage of unique tokens that all passages in a cluster have in common. All text was lowercased and punctuation was stripped. | Histogram which shows the distribution of lexical overlap in percent and indicates the largest number of clusters for a given score. Extremely low lexical overlap decreases the chance to discover meaningful text reuse whilst extremely high overlap will only reveal near-copies of content and may be too restrictive for some purposes. | **LEXICAL OVERLAP** ⓘ<br>100 (655,914 results)<br>0　　　　100 |
| Time span | The time window covered by documents in the cluster, measured in number of days. | Histogram which shows the gap between the earliest publication date of an article in a text reuse cluster and the latest measured in days and indicates the largest number of passages for a given score. This is an efficient approach to discover or filter for instances of slow, mid-range and rapid text reuse. The histogram groups clusters by the number of days in between publication dates and displays their sum. | **TIME SPAN IN DAYS** ⓘ<br>0 - 727 (13,824,681 results)<br>0　　　　72,700 |
| Text reuse clusters | Clusters store text segments (or passages) that are reused in different units of a corpus. | List of text reuse clusters which match a given query, sorted by number of passages. Each cluster is characterized with basic information (passages count, lexical overlap, time periods and years covered) as well as a snippet preview of the passage. Clusters are sorted by the number of matching passages. Clusters can be selected manually for further inspection in the Text Reuse app or in other *impresso* components such as Search. | **FILTER BY TEXT REUSE CLUSTERS (6,121,859 OPTIONS)** ⓘ ⊖<br>☐ cluster **c29**<br>0% lexical overlap over 47442 days (1884 - 2014).<br>CONVOCATION Mesdames, Messieurs les Actionnaires, Le Conseil d'administration de The Swatch Group SA a le plaisir de vous inviter, conformément aux art. 12 ss |

# 3. Demo *impresso* Text reuse at Scale

# Types of text reuse in newspapers seen through filters

|  | **Passages per year** | **Cluster size** | **Lexical overlap** | **Time span (=duration)** |
|---|---|---|---|---|
| **Reprints of historical materials** | Few | Small | Very high | Very large |
| **Co-publication** | Many | Rather small | Very high | Very short |
| **Advertising campaign** | Many | Large | High | Varies |

# Hands-on part and Challenge: Search and lexical overlap

1. **Go to** https://bit.ly/impresso-tr
2. **Go to** https://dev.impresso-project.ch/app/ and login
3. **Search** for "titanic"
4. **Exercise**: Familiarise yourself with all the different filters and views. Which questions emerge?
5. **Lexical overlap**: Go to Passages and experiment with the different sorting choices. What do you observe?

URL: https://dev.impresso-project.ch/app

Login: student@impresso-project.ch

Passwort: Marburg2023!

# Hands-on part and Challenge: Time span

1. Start a new search by clicking the red "X"
2. Filter for "German"
3. Set Lexical overlap to 15% - 100%
4. Move to Passages view. What do you get?
5. Set time span to 55.00 - 67.000 days. What do you get?

# Hands-on part and Challenge: Cluster size

1. Start a new search by clicking the red "X"
2. Set filters for adverts
3. Set Cluster size between 90 - 100 passages.
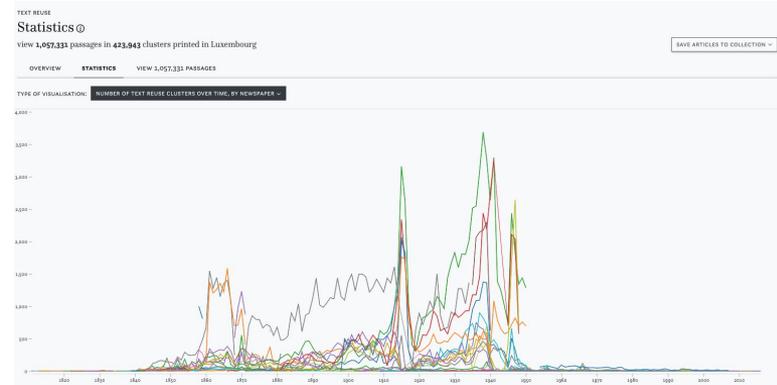4. What do you find? How can we explain this?

# Hands-on part and Challenge: Topics etc

1. Start a new search by clicking the red "X"
2. Set filters to "German" and adverts
3. Experiment with different topics: medical, classified ads, media…
4. Use the Passages tab and experiment with different sorting operations: e.g. largest cluster size vs. lexical overlap, passage size etc.

# Hands-on part and Challenge: Topics etc

1. Start a new search by clicking the red "X"
2. Set the **country** filters to "Luxembourg"
3. Select the Statistics view + number of text reuse clusters over time
4. How can we explain this graph?

# Your turn…find instances of:

1. Build you own queries, combine different keywords and filters. Which impact do they have?
2. Explore the Statistics view. What can it tell you about the data?
3. Open https://bit.ly/impresso-tr and go to slides 44ff. Pick one of the project ideas or duplicate slides to add your own.

|  | Pass-ages per year | Cluster size | Lexical overlap | Time span |
|---|---|---|---|---|
| **Reprints of historical materials** | Few | Small | Very high | Very large |
| **Co-publication** | Many | Rather small | Very high | Very short |
| **Advertising campaign** | Few | Large | High | Varies |

# Your idea here

<SCREENSHOT(S)>

Insert here a brief description of what type of text reuse it is and what we can learn from it.

Your names here

Observations:

| URL to query | |
|---|---|
| Number of passages | |
| Lexical overlap | |
| Time span | |

# Standardized content (weather, radio programme…)

<SCREENSHOT(S)>

Insert here a brief description of what type of text reuse it is and what we can learn from it.

Your names here

Observations:

| URL to query | |
|---|---|
| Number of passages | |
| Lexical overlap | |
| Time span | |

# The longest running advertisement

<SCREENSHOT(S)>

Your names here

Observations:

Insert here a brief description of what type of text reuse it is and what we can learn from it.

| | |
|---|---|
| URL to query | |
| Number of passages | |
| Lexical overlap | |
| Time span | |

# Press agency reports

<SCREENSHOT(S)>

Insert here a brief description of what type of text reuse it is and what we can learn from it.

Your names here

Observations:

| URL to query | |
|---|---|
| Number of passages | |
| Lexical overlap | |
| Time span | |

# The two newspapers that co-publish the m

<SCREENSHOT(S)>

Your names here

Observations:

| | |
|---|---|
| URL to query | |
| Number of passages | |
| Lexical overlap | |
| Time span | |

Insert here a brief description of what type of text reuse it is and what we can learn from it.

# A new form of text reuse!?

<SCREENSHOT(S)>

Your names here

Observations:

| | |
|---|---|
| URL to query | |
| Number of passages | |
| Lexical overlap | |
| Time span | |

Insert here a brief description of what type of text reuse it is and what we can learn from it.

# For more information



impresso website (update soon)

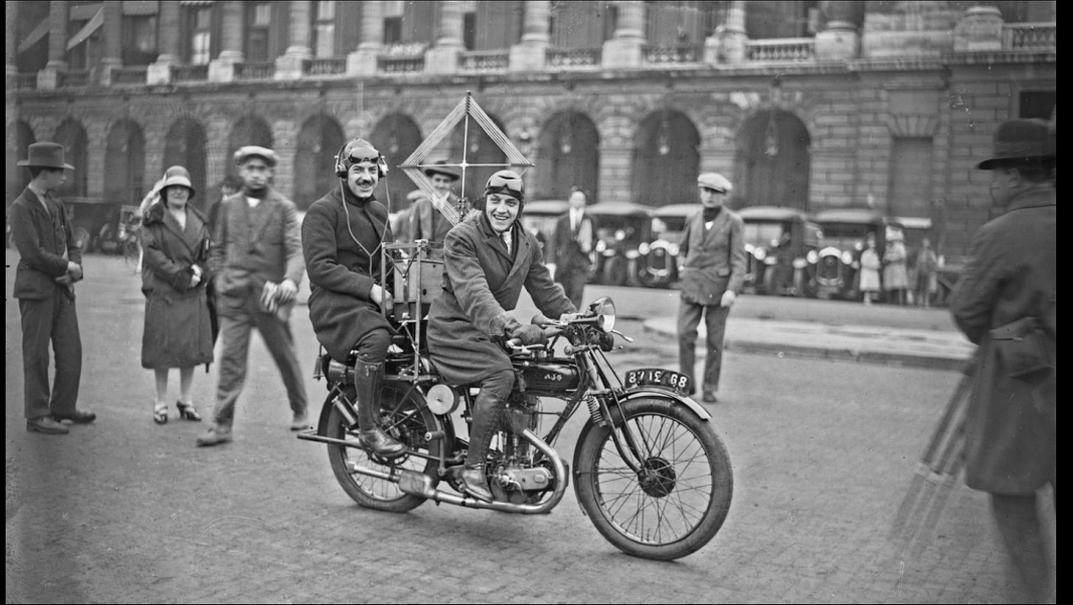data on Zenodo

code on GitHub

impresso youtube channel

impresso-project.ch/app

# Temporary account

URL: https://dev.impresso-project.ch/app

Login: student@impresso-project.ch

Passwort: Marburg2023!
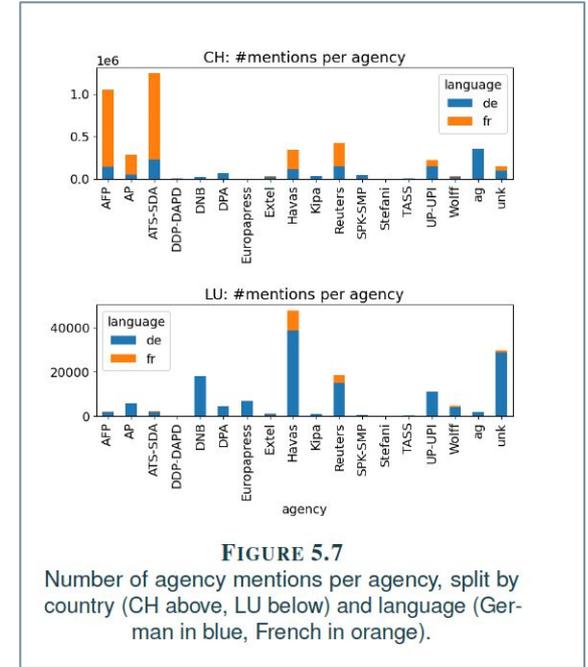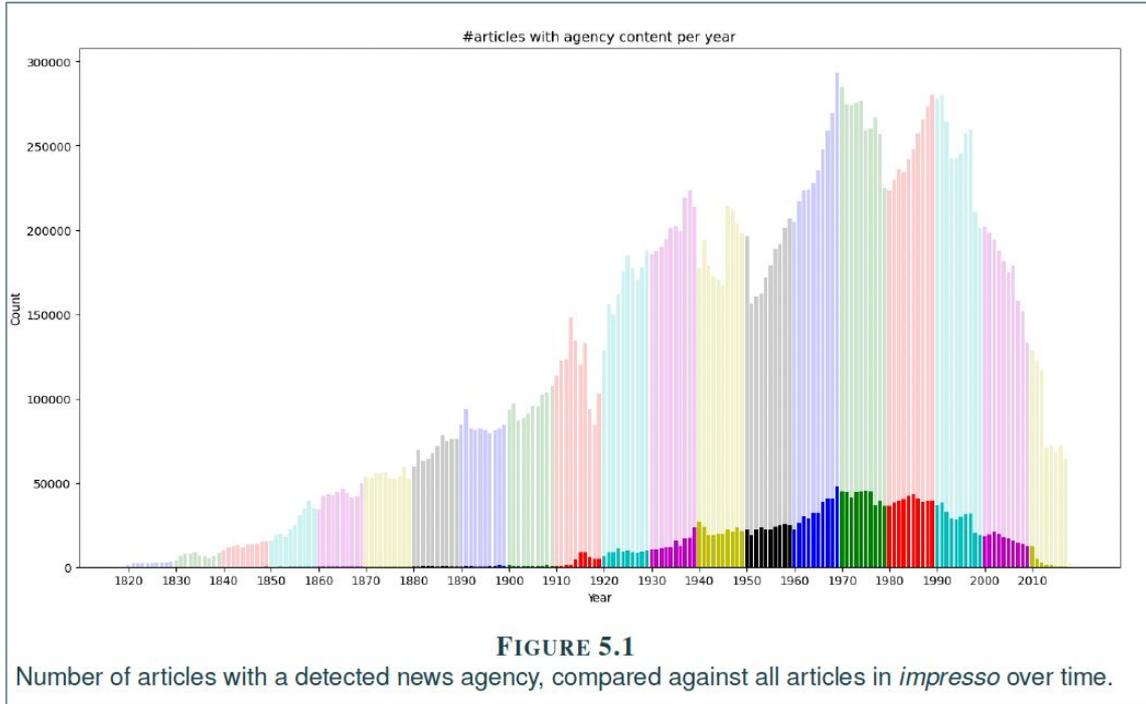
Ausblick: Was kommt als Nächstes?

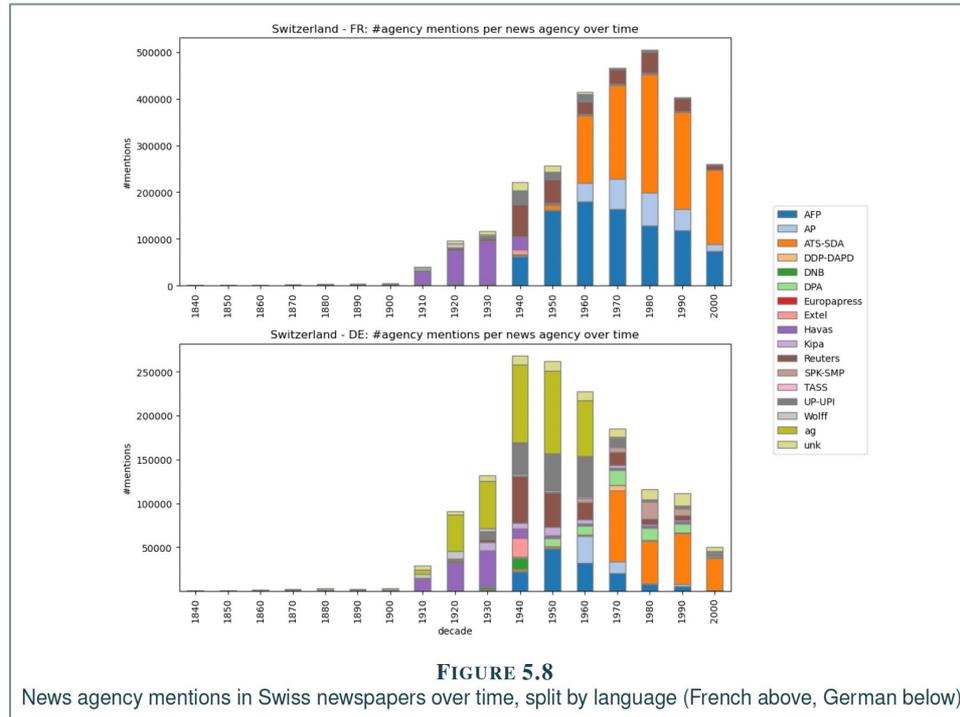# Neue Perspektiven auf Agenturnachrichten

**Erkennung von Presseagenturen**

- Auf Basis eines Trainingskorpus (27 Agenturen, ca 2000 Artikel, fr und de);
- Training und Evaluierung von untersch. Sprachmodellen;
- Eweiterung auf den ganzen *impresso* Korpus;
- Erste Analysen;
- Integration in die App.

Media Monitoring of the Past impresso — Search — Newspapers Topics Inspect & Compare Text reuse — FAQ

SEARCH ARTICLES    SEARCH IMAGES    NGRAMS

GROUP BY ARTICLE ∨    ORDER BY RELEVANCE ∨    DISPLAY

"conference" ∨    Reuters AND AFP ∨    ✕

add keyword to search

16,641 articles found containing conference gathered by Reuters and AFP

FILTER BY LANGUAGE OF ARTICLES (2 OPTIONS)    ⊕

FILTER BY NEWSPAPER TITLES (10 OPTIONS) ⓘ    ⊕

FILTER BY ARTICLE TYPE (1 OPTION)    ⊕

FILTER BY COUNTRY OF PUBLICATION (2 OPTIONS)    ⊕

FILTER BY ACCESS RIGHT (4 OPTIONS)    ⊕

FILTER BY ARCHIVE (4 OPTIONS)    ⊕

FILTER BY PERSON (10,266 OPTIONS) ⓘ    ⊕

FILTER BY LOCATION (5,778 OPTIONS) ⓘ    ⊕

FILTER BY TOPIC (164 OPTIONS) ⓘ    ⊕

FILTER BY NEWS AGENCY (12 OPTIONS)    RESET

WRITTEN BY ∨

☑ Reuters(16,641 results)

☑ AFP(16,641 results)

☐ ATS SDA (3,361 results)

☐ UP UPI (1,430 results)

☐ AP (785 results)

**s-fffft ^ffl^^ J|I •«fjifjr v : _ _T mWk ...**
Die Tat  WEDNESDAY, MAY 1, 1968 – p.13
Personal use

Zürich . ( UP ) 624 , 8 Millionen Franken zahlten die ln der Schweiz tätigen LebenherungsgeselL schaffen lm Jahre 1967 aus , jeden Tag durch

LOCATIONS Zürich District, Switzerland, NATO, ATP World Tour Finals, Switzerla team, Lagos, Nigerian records in athletics, Pound sterling, Montevideo
PEOPLE ATP World Tour Finals

. Mai an gewähren die « Atlantic Passenger Steamship Conference Lines » 1968

VIEW    ADD TO COLLECTION ... ∨

**Erdbeben registriert Pasadena (Kaliforni...**
Freiburger Nachrichten  WEDNESDAY, JUNE 16, 1954 – p.8
Personal use

Erdbeben registriert Pasadena ( Kalifornien ) , 16 . Juni . ag . ( Reuter . ) Das kali technologische Institut registrierte am Dienstag zwei

LOCATIONS Pasadena, California, University of Freiburg, Salaberry-de-Valleyfield Switzerland
PEOPLE Ernest Hill

, und Abbe J . D . Cadieux vom Pressedienst CCC ( Conference Catholique befindet , reiste

VIEW    ADD TO COLLECTION ... ∨

**ten auf Ministerebene befassen. Diese Ko...**
Die Tat  MONDAY, JUNE 2, 1958 – p.2
Personal use

ten auf Ministerebene befassen . Diese Konferenz soll am 15 . September in Mon

# Erwähnungen von Presseagenturen



**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 der schweizerischen Presse…



**FIGURE 5.8**
News agency mentions in Swiss newspapers over time, split by language (French above, German below).

# …in der luxemburgischen Presse (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

## 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.

# For more information



impresso website (update soon)

data on Zenodo

code on GitHub

impresso youtube channel

impresso-project.ch/app

# Text reuse and historical research objectives

1.  (Trans-) national media ecosystems

2.  Newspaper content as *bricolage*

3.  Historicising virality

4.  Tracing historical events

5.  Capturing historical Zeitgeist

# (Trans-) nationale Medien-Ökosysteme

*Goal: Highlight the connectivity of historical media across borders and languages. Which ideological, commercial, and financial structures made this work and how did they shape media?*

*TR: (How) does content flow through the international network of newspapers?*

Focus e.g. on cut-paste practices, the relevance of information infrastructure, geography, cross-border cultural affinities.



Beelen, K. Digitising newspapers press directories to understand the landscape of historical newspapers

- Smith, D. A., Cordell, R., and Mullen, A. (2015). Computational methods for uncovering reprinted texts in antebellum newspapers. *Am. Liter. Hist*. 27, E1–E15. doi: 10.1093/alh/ajv029
- Keck, J., Oiva, M., and Fyfe, P. (2022). Lajos kossuth and the transnational news: a computational and multilingual approach to digitized newspaper collections. *Media History* 29, 287–304. doi: 10.1080/13688804.2022.2146905
- Salmi, H., Rantala, H., Vesanto, A., and Ginter, F. (2019). "The long-term reuse of text in the finnish press, 1771–1920," in *Proceedings of the Digital Humanities in the Nordic Countries 4th Conference*, eds. C. Navarretta, M. Agirrezabal, and B. Maegaard (Copenhagen, Denmark: CEUR Workshop Proceedings), 253–273.
- Paju, P., Salmi, H., Rantala, H., Lundell, P., and Marjanen, Vesanto, A. (2022). "Textual migration across the baltic sea: Creating a database of text reuse between Finland and Sweden," in *Proceedings of the 6th Digital Humanities in the Nordic and Baltic Countries Conference (DHNB 2022), CEUR Workshop Proceedings*, eds. K. Berglund, M. La Mela, and I. Zwart (Aachen: CEUR-WS.org), 361–369.
- Beelen, K. Digitising newspapers press directories to understand the landscape of historical newspapers, 2023, https://livingwithmachines.ac.uk/digitising-newspapers-press-directories-to-understand-the-landscape-of-historical-newspapers/.
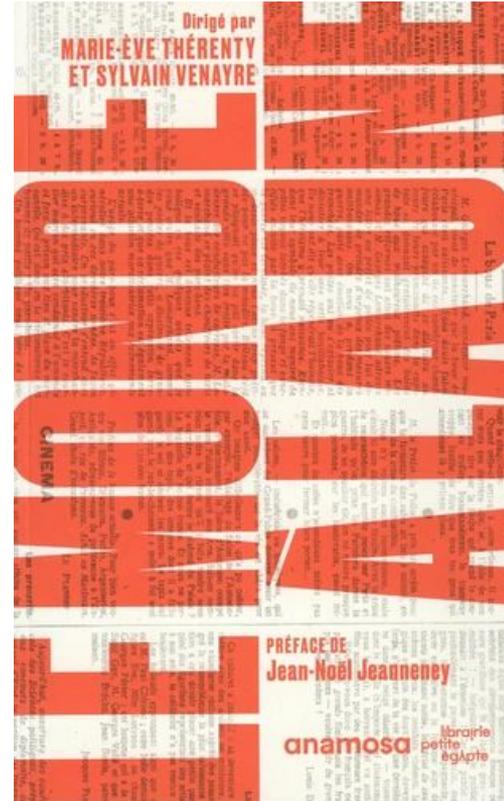
# Newspaper content as bricolage

*From which textual resources were newspapers compiled?*

*TR: How did formulaic ways of representing information emerge and spread (e.g. weather reports)? Can we determine the origins of content? Which content was influential?*

Focuses on parallel developments across titles and the emergence of common practices.

- Walma, L. W. B. (2015). Filtering the "news:" Uncovering morphine's multiple meanings on delpher's dutch newspapers and the need to distinguish more article types. *Tijdschrift voor Tijdschriftstudies*. 38, 61–78. doi: 10.18352/ts.345
- Thèrenty, M.-E., and Venayre, S. (2021). *Le monde à la une. Une histoire de la presse par ses rubriques*. Anamosa, illustrated èdition edition. doi: 10.3917/anamo.there.2021.02

# Historicising virality

*Which conditions (geography, type of information, infrastructure) enable rapid dissemination?*

Focus is on the breadth and speed with which content spreads. Pioneering work of Paju et al. who define a virality score based on the number of titles within a cluster, the number of unique printing locations, and the distance in days between the first and last passage publication date.

Paju, P., Salmi, H., Rantala, H., Lundell, P., and Marjanen, Vesanto, A. (2022). "Textual migration across the baltic sea: Creating a database of text reuse between Finland and Sweden," in *Proceedings of the 6th Digital Humanities in the Nordic and Baltic Countries Conference (DHNB 2022), CEUR Workshop Proceedings*, eds. K. Berglund, M. La Mela, and I. Zwart (Aachen: CEUR-WS.org), 361–369.

## Textual Migration Across the Baltic Sea: Creating a Database of Text Reuse Between Finland and Sweden

Petri Paju[1], Hannu Salmi[1], Heli Rantala[1], Patrik Lundell[2], Jani Marjanen[3] and Aleksi Vesanto[1]

[1] *University of Turku, Department of Cultural History, Turku, FI-20014, Finland*
[2] *Örebro University, School of Humanities, Örebro, SE-70182, Sweden*
[3] *University of Helsinki, Department of Digital Humanities, Helsinki, FI-00014, Finland*

**Abstract**
In this paper, we present a database and an interface on text reuse between newspapers and journals published in the Swedish language in Sweden and Finland during the 1645–1918 time frame. Using two national, digital newspaper collections, we detected their textual similarities with a computational method to study the textual migration, i.e., information flows, between the two countries. For purposes of this project, we developed a database of detected clusters of text reuse and an online interface to search, examine and analyse the transnational movement of information. The database, *Text Reuse in the Swedish-language Press, 1645–1918*, is accessible online and includes texts from over 1,100 newspapers and journals published at approximately 150 locations at various times during the 274-year time frame.

**Keywords**
computational history, text reuse, historical newspaper, digital collections, database construction, Finland, Sweden, transnational history, information flow

### 1. Introduction

This short paper presents a database and an interface on text reuse among Swedish-language newspapers and journals during the 1645–1918 time frame. The database and interface were built as part of the project, *Information flows across the Baltic Sea: Swedish-language press as a cultural mediator, 1771–1914*. The database and the accompanying project aim to study information flows, particularly between Sweden and Finland from the period when present-day Finland was part of the Swedish kingdom to the establishment of Finland as a Grand Duchy in the Russian Empire after 1809 and until the Independence of Finland in 1917 and Civil War in 1918. Even after the 1809 separation, news and other texts circulated because of the common cultural heritage and shared language, i.e., Swedish. The border was relatively easy to cross, and newspapers circulated between Sweden and Finland regularly. However, because their national histories eventually diverged, these press materials have been preserved, processed, and siloed in two national libraries. Still, print media digitisation makes it possible to study overlaps in large collections of texts and see how information was spread across the Baltic Sea.

In our project, textual migration was traced using a method based on the software BLAST, which can be applied to text-reuse detection. With the method, we detected every text passage with 300 or more characters of similarity and combined these passages into reuse clusters. We included Swedish-language papers published in Sweden and Finland, but excluded the Finnish-language press in Finland, as textual migration within Finland has been studied in previous publications [1, 2]. To strengthen the

# Tracing historical events

*How do individual media titles situate events in the political, economic, social, and cultural context relevant to them? How does this affect the perception of these events by their audiences?*

Focus is on the appropriation of often global news on the local scale.



**FIGURE 2.**
The bar plot of the number of reprinted articles mentioning Kossuth in OcEx corpus shows that the American tour established Kossuth as an international celebrity.

Keck, J., Oiva, M., and Fyfe, P. (2022). Lajos kossuth and the transnational news

- Oiva, M., Nivala, A., Salmi, H., Latva, O., Jalava, M., Keck, J., et al. (2020). Spreading news in 1904. *Media History* 26, 391–407. doi: 10.1080/13688804.2019.1652090
- Keck, J., Oiva, M., and Fyfe, P. (2022). Lajos kossuth and the transnational news: a computational and multilingual approach to digitized newspaper collections. *Media History* 29, 287–304. doi: 10.1080/13688804.2022.2146905

# Tracing historical events

# Capturing historical Zeitgeist

*Historical media partially capture the attitudes, norms, beliefs, moods and feelings of past generations, or Zeitgeist.*

Focus is on texts which were produced independently but still share certain characteristics.

- Verheul, J., Salmi, H., Riedl, M., Nivala, A., Viola, L., Keck, J., et al. (2022). Using word vector models to trace conceptual change over time and space in historical newspapers, 1840–1914. *Dig. Human. Quart*. 16, 7445. Available online at: https://www.digitalhumanities.org/dhq/vol/16/2/000550/000550.html
- Paasikivi, S., Salmi, H., Vesanto, A., and Ginter, F. (2022). Infectious media: Cholera and the circulation of texts in the finnish press, 1860–1920. *Media Hist*. 29, 17–38. doi: 10.1080/13688804.2022.2054408

# Tasks to support research objectives

| Task | Title | Level | Support |
|------|-------|-------|---------|
| 1 | Obtain an overview of text reuse in a corpus, collection or query | Corpus | Yes |
| 2 | Obtain an overview of a single cluster | Cluster | Yes |
| 3 | Compare passages | Passage | Yes |
| 4 | Compare clusters | Cluster | Yes |
| 5 | Identify different types of text reuse | Corpus | Yes |
| 6 | Generate research corpora based on text reuse clusters | Corpus | Yes |
| 7 | Identify connections | Corpus | Partial |
| 8 | Detect and trace virality | Corpus | No |
| 9 | Search for passages | Passage | No |
| 10 | De-duplicate content | Corpus | No |
| 11 | Export of text reuse data | All | Planned |

# Temporalities in text reuse data

| Type | Description | Measures | Examples |
|------|-------------|----------|----------|
| Duration | The time period which is covered by a cluster ranging from the earliest to the latest publication date of individual passages. | Publication date | Paju et al.'s notions of fast and slow text reuse fall into this category. |
| Virality | The speed (measured in days) and breadth of text reuse passages spreading within a corpus. Speed corresponds to time passed (e.g., days) whereas breadth corresponds to the number of publications which contain a passage at a given point in time. | Publication date, number of publications | News of the sinking of the Titanic or the destruction of the Hindenburg Zeppelin traveled around the world within days or weeks. |
| Rhythm | Pattern with which text reuse passages appear over time. | Distance between publication dates | Reprints of articles on the occasion of their anniversary, e.g., on the occasion of the bombing of Hiroshima. |

# Capture the characteristics of text reuse through filters

| Measure | Description | Implementation in interface prototype | |
|---|---|---|---|
| Passages per year | Number of passages counted in a given year. | Line chart which displays the count of passages per year for a given query or filter operation. This gives a first indication, during which years text reuse occurred more commonly. Time sliders and precise date entry allow users to filter for exact date ranges to inspect. |  |
| Cluster size | The number of passages contained in a cluster. | Histogram which shows the distribution of text reuse cluster sizes and indicates the highest score. The histogram groups clusters of size n and displays their sum. This gives a first indication of averages as well as outliers. Sliders can be used to specify a cluster size range of interest. Filtering by cluster size allows to exclude or explicitly focus on outliers but different cluster sizes may also correspond to different types of content. |  |
| Lexical overlap | The percentage of unique tokens that all passages in a cluster have in common. All text was lowercased and punctuation was stripped. | Histogram which shows the distribution of lexical overlap in percent and indicates the largest number of clusters for a given score. Extremely low lexical overlap decreases the chance to discover meaningful text reuse whilst extremely high overlap will only reveal near-copies of content and may be too restrictive for some purposes. |  |
| Time span | The time window covered by documents in the cluster, measured in number of days. | Histogram which shows the gap between the earliest publication date of an article in a text reuse cluster and the latest measured in days and indicates the largest number of passages for a given score. This is an efficient approach to discover or filter for instances of slow, mid-range and rapid text reuse. The histogram groups clusters by the number of days in between publication dates and displays their sum. |  |
| Text reuse clusters | Clusters store text segments (or passages) that are reused in different units of a corpus. | List of text reuse clusters which match a given query, sorted by number of passages. Each cluster is characterized with basic information (passages count, lexical overlap, time periods and years covered) as well as a snippet preview of the passage. Clusters are sorted by the number of matching passages. Clusters can be selected manually for further inspection in the Text Reuse app or in other *impresso* components such as Search. |  |

# 3. Demo *impresso* Text reuse at Scale

# Types of text reuse in newspapers seen through filters

|  | **Passages per year** | **Cluster size** | **Lexical overlap** | **Time span (=duration)** |
|---|---|---|---|---|
| **Reprints of historical materials** | Few | Small | Very high | Very large |
| **Co-publication** | Many | Rather small | Very high | Very short |
| **Advertising campaign** | Many | Large | High | Varies |

# Hands-on part and Challenge: Search and lexical overlap

1. **Go to** https://dev.impresso-project.ch/app/ and login
2. **Go to** https://bit.ly/impresso-tr
3. **Search** for "titanic"
4. **Exercise**: Familiarise yourself with all the different filters and views. Which questions emerge?
5. **Lexical overlap**: Go to Passages and experiment with the different sorting choices. What do you observe?

URL: https://dev.impresso-project.ch/app

Login: student@impresso-project.ch

Passwort: Marburg2023!

# Hands-on part and Challenge: Time span

1. Start a new search by clicking the red "X"
2. Filter for "German"
3. Set Lexical overlap to 15% - 100%
4. Move to Passages view. What do you get?
5. Set time span to 55.00 - 67.000 days. What do you get?

# Hands-on part and Challenge: Cluster size

1. Start a new search by clicking the red "X"
2. Set filters for adverts
3. Set Cluster size between 90 - 100 passages.
4. What do you find? How can we explain this?

# Hands-on part and Challenge: Topics etc

1. Start a new search by clicking the red "X"
2. Set filters to "German" and adverts
3. Experiment with different topics: medical, classified ads, media…
4. Use the Passages tab and experiment with different sorting operations: e.g. largest cluster size vs. lexical overlap, passage size etc.

# Hands-on part and Challenge: Topics etc

1. Start a new search by clicking the red "X"
2. Set the **country** filters to "Luxembourg"
3. Select the Statistics view + number of text reuse clusters over time
4. How can we explain this graph?

# Your turn…find instances of:

1. Build you own queries, combine different keywords and filters. Which impact do they have?
2. Explore the Statistics view. What can it tell you about the data?
3. Open https://bit.ly/impresso-tr and go to slides 44ff. Pick one of the project ideas or duplicate slides to add your own.

| | Pass-ages per year | Cluster size | Lexical overlap | Time span |
|---|---|---|---|---|
| **Reprints of historical materials** | Few | Small | Very high | Very large |
| **Co-publication** | Many | Rather small | Very high | Very short |
| **Advertising campaign** | Few | Large | High | Varies |

# Your idea here

<SCREENSHOT(S)>

Insert here a brief description of what type of text reuse it is and what we can learn from it.

Your names here

Observations:

| | |
|---|---|
| URL to query | |
| Number of passages | |
| Lexical overlap | |
| Time span | |

# Standardized content (weather, radio programme…)

<SCREENSHOT(S)>

Insert here a brief description of what type of text reuse it is and what we can learn from it.

Your names here

Observations:

| | |
|---|---|
| URL to query | |
| Number of passages | |
| Lexical overlap | |
| Time span | |

# The longest running advertisement

<SCREENSHOT(S)>

Your names here

Observations:

| | |
|---|---|
| URL to query | |
| Number of passages | |
| Lexical overlap | |
| Time span | |

Insert here a brief description of what type of text reuse it is and what we can learn from it.

# Press agency reports

<SCREENSHOT(S)>

Insert here a brief description of what type of text reuse it is and what we can learn from it.

Your names here

Observations:

| | |
|---|---|
| URL to query | |
| Number of passages | |
| Lexical overlap | |
| Time span | |

# The two newspapers that co-publish the m

<SCREENSHOT(S)>

Your names here

Observations:

| URL to query | |
|---|---|
| Number of passages | |
| Lexical overlap | |
| Time span | |

Insert here a brief description of what type of text reuse it is and what we can learn from it.

# A new form of text reuse!?

<SCREENSHOT(S)>

Insert here a brief description of what type of text reuse it is and what we can learn from it.

Your names here

Observations:

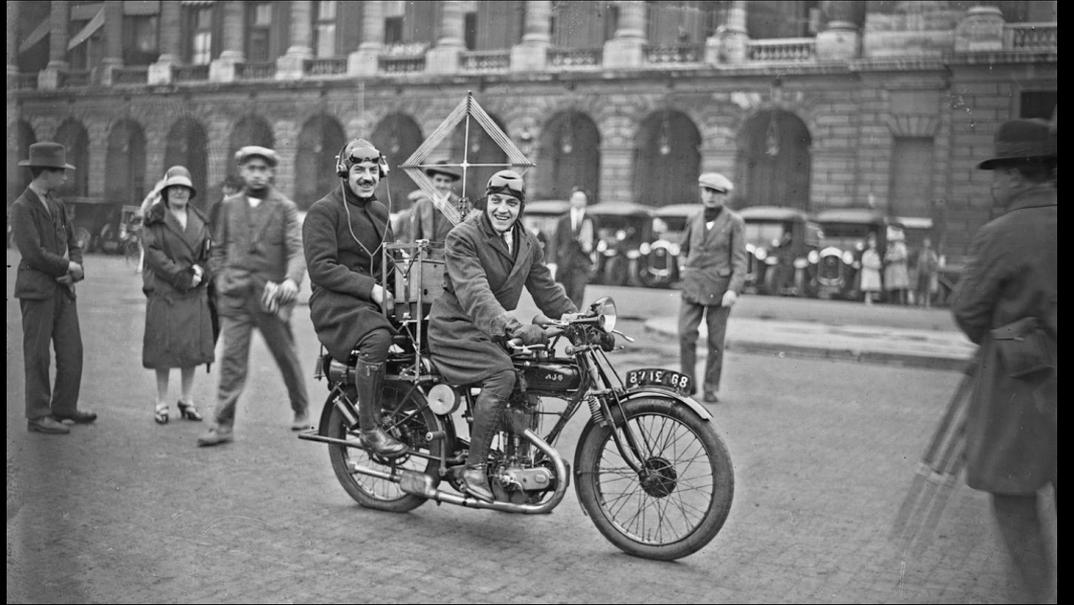| | |
|---|---|
| URL to query | |
| Number of passages | |
| Lexical overlap | |
| Time span | |

# Temporary account

URL: https://dev.impresso-project.ch/app

Login: student@impresso-project.ch

Passwort: Marburg2023!

Ausblick: Was kommt als Nächstes?

# Neue Perspektiven auf Agenturnachrichten

**Erkennung von Presseagenturen**

- Auf Basis eines Trainingskorpus (27 Agenturen, ca 2000 Artikel, fr und de);

- Training und Evaluierung von untersch. Sprachmodellen;

- Eweiterung auf den ganzen *impresso* Korpus;

- Erste Analysen;

- Integration in die App.

# Erwähnungen von Presseagenturen



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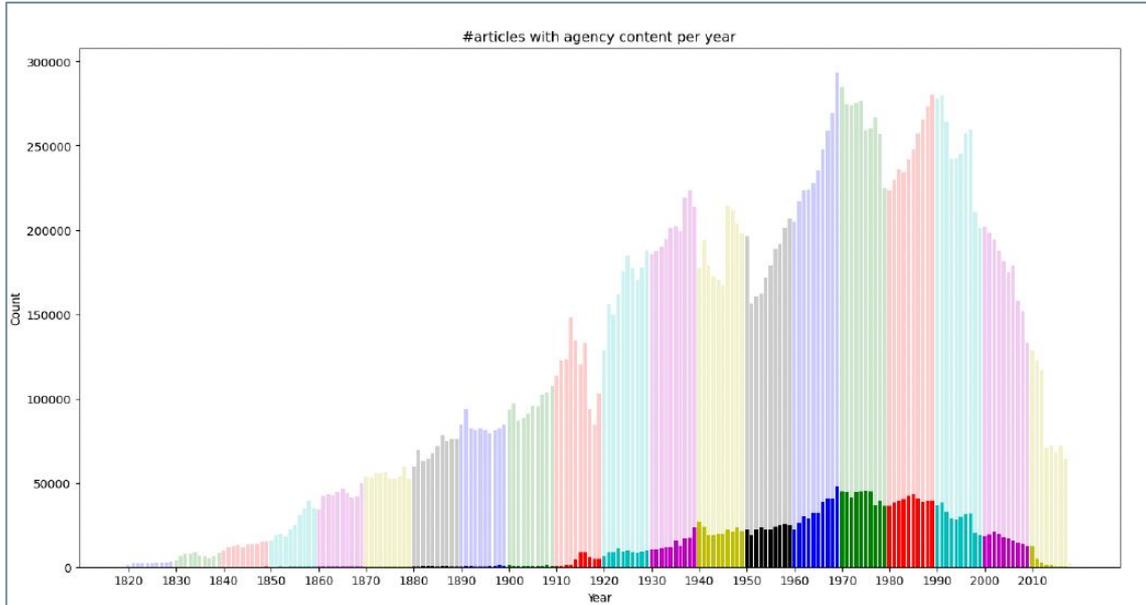Number of articles with a detected news agency, compared against all articles in *impresso* over time.



**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 der schweizerischen Presse…



**FIGURE 5.8**
News agency mentions in Swiss newspapers over time, split by language (French above, German below).

# …in der luxemburgischen Presse (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?

# 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.