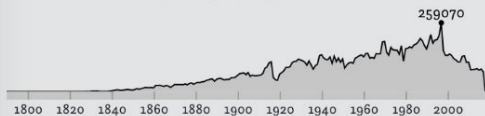SEARCH TEXT REUSE PASSAGES

start searching...

NUMBER OF PASSAGES PER YEAR ⓘ
Number of text reuse passages per year
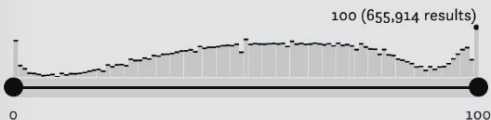
259070

1800  1820  1840  1860  1880  1900  1920  1940  1960  1980  2000

ADD NEW DATE FILTER ...

CLUSTER SIZE ⓘ
How to read histograms ⓘ

2 - 453 (15,756,994 results)

2                                              45,553

LEXICAL OVERLAP ⓘ

100 (655,914 results)

0                                              100

TIME SPAN IN DAYS ⓘ

0 - 727 (13,824,681 results)

72,700

TEXT REUSE

# impresso Text Reuse at Scale

SAVE ARTICLES TO COLLECTION ⌄

OVERVIEW    STATISTICS    6,121,859 CLUSTERS    VIEW 16,099,821 PASSAGES

**NEWSPAPER**
| | |
|---|---|
| L'Express | 3,152,955 |
| L'Impartial | 2,943,902 |
| Journal de Genève | 2,397,372 |
| Gazette de Lausanne | 2,197,976 |
| La Liberté | 1,991,592 |
| Le Peuple, La Sentinelle | 748,793 |
| Freiburger Nachrichten | 373,980 |
| Confédéré | 328,495 |
| Die Tat | 247,654 |
| Neue Zürcher Zeitung | 229,565 |

**COUNTRY**
| | |
|---|---|
| Switzerland | 14,186,834 |
| Luxembourg | 1,057,331 |

**TYPE**
| | |
|---|---|
| article | 13,758,049 |
| advertisement | 1,863,669 |
| page article | 229,565 |
| section | 48,166 |
| obituary | 13,106 |
| tables | 2,240 |
| buckets.type.NaN | 2,066 |
| weather news (other) | 552 |

Quantitative analysis of text reuse. Towards a methodology
Helsinki, 23-24.11.2023

**LANGUAGE**
| | |
|---|---|
| French | 13,945,214 |
| German | 1,295,522 |
| English | 2,663 |
| Luxembourgish | 766 |

**PERSON**
| | |
|---|---|
| United States | 186,999 |
| Les Anglais | 87,128 |
| Les Alliés | 53,866 |
| Paris | 50,322 |
| David Lloyd George | 46,276 |
| Les Brenets | 45,923 |
| Léon Blum | 44,533 |
| Les Chambres | 41,953 |

**LOCATION**
| | |
|---|---|
| France | 2,851,220 |
| Paris | 2,391,261 |
| Switzerland | 2,347,917 |
| Suisse, Moselle | 2,323,778 |
| Lausanne | 1,909,000 |
| Zürich | |
| Italy | |
| Berlin | |

Marten Düring & impresso team

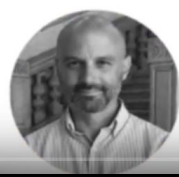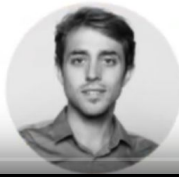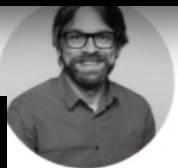# impresso - Media Monitoring of the Past II
*Beyond Borders: Connecting Historical Newspapers and Radio*

*How can semantic enrichments enhance the study of digitised historical newspapers?*

*impresso* team

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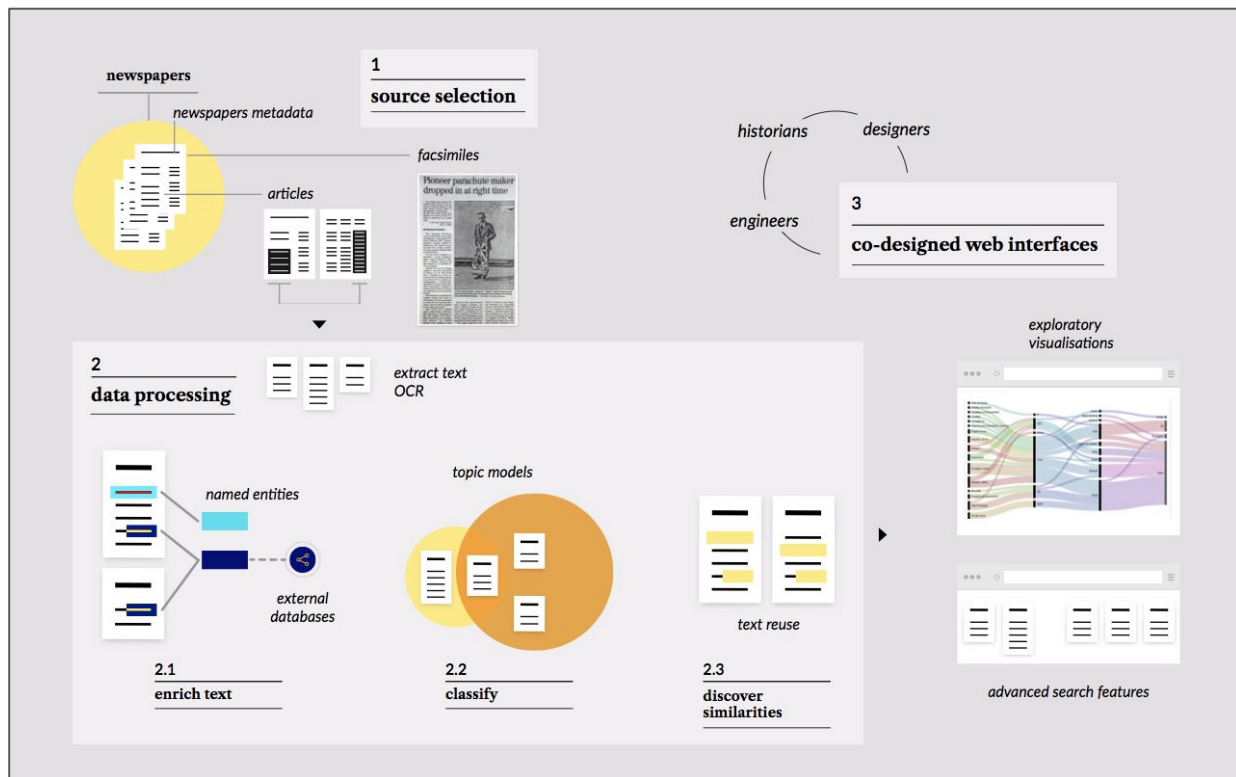Phillip Stroëbel
Martin Volk
Thijs van Beek
Lars Wieneke

+ a team of historical advisors and associated researchers in (media) history

# 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?

# Research interests

**Discipline**
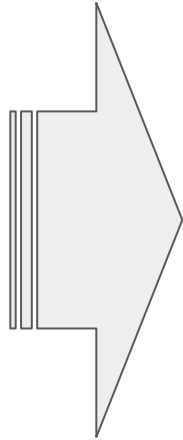
Media history

Gender and social norms

Political/cultural history

Social sciences

Social history

…
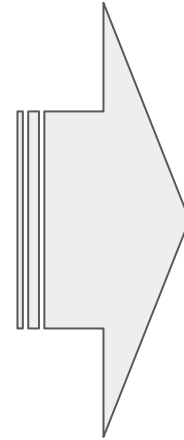
**Interests**

Layout

Social norms

Public opinion

Knowledge horizons

Biographical data

….

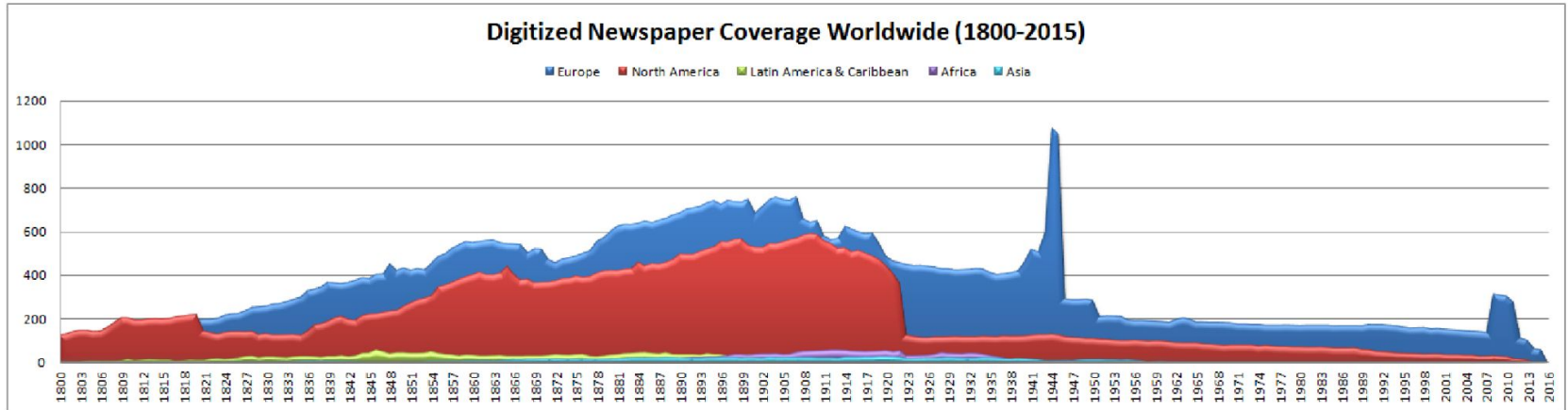**Newspaper Data**

Adverts

Opinion pieces

Press agency text

Images

Classifieds

News

Obituaries

….

# The big tip of a hidden iceberg



Digitized Newspaper Coverage Worldwide (1800-2015)

# The challenging landscape of historical newspapers

1.  Institutional silos

2.  Big and messy data

3.  Noisy historical text

4.  Visualisation and exploration
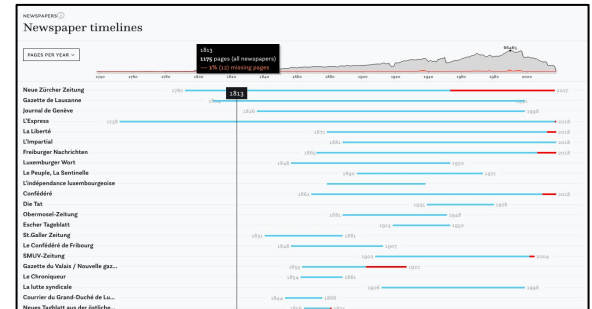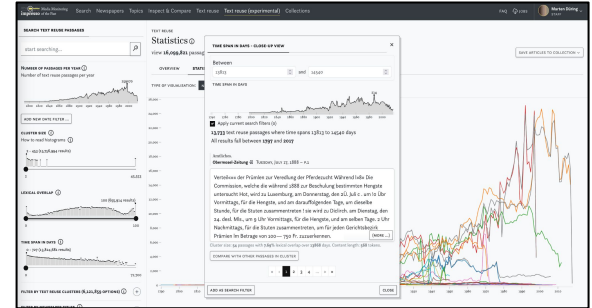
5.  Digital Scholarship

# The impresso app

**Based on** digitised newspaper collections provided by partners.

**Target audience** are historians and other humanities scholars.

**Iterative exploration** from multiple perspectives, not an analytical workbench.

**Interfaces** for overview, exploration, limits & opportunities, flaws in the data.

# About Text reuse at Scale

# impresso text reuse data

Using passim and applied to 76 newspapers from CH and LU; yielded 16 mio passages in 6 mio clusters



TYPE OF VISUALISATION: NUMBER OF TEXT REUSE CLUSTERS OVER TIME, BY NEWSPAPER ⌄

Romanello, Matteo. *Text re-use detection in a nutshell*, Blog post, impresso, 2018 <http://impresso-project.ch/news/2018/06/12/tradingzone-tr.html>.

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

*How can we enable scalable reading of text reuse data within the impresso web app?*

# Workshop

Funded by a mini grant by C2DH to conduct experimental research.

Workshop with 10 NLP, history, design experts in Nov 2022 which produced:

- A set of historical research objectives based on case studies
- Focus on the integration of TR with other semantic enrichments
- Three mockups to support envisioned tasks

# Outputs

Mockups, objectives and tasks inspired:

- Text reuse at Scale interface for impresso Web app
- Press agency detection project (impresso + Lea Marxen)
- User evaluation (Zenodo)
- Paper: *impresso* Text Reuse at Scale (Frontiers in Big Data, 6/2023)

# Prios and desiderata

**Main Priorities**

1. Scalable reading of text reuse data
2. Integration in impresso Web app and its component to enable general exploration
3. Integration with other semantic enrichments for more precision
4. ~~Keep Frontiers deadline~~

**Main Desiderata**

1. Passage-based search
2. Passage-based comparison
3. ~~Cluster-level view~~
4. TR data export

# Side project: 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.

# Side project: News agency mentions



#articles with agency content per year

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



CH: #mentions per agency

LU: #mentions per agency

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

# Side project: In Swiss newspapers



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

*Why do we want to work TR in newspapers?*

# Types of Text Reuse



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

# Types of 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? |

| | **Research questions** | **Purpose of TR** | **Examples** | **Match with Rosson et al.** |
|---|---|---|---|---|
| **Media ecosystems** | How did historical media function? How did transnational information flow shape historical media? | TR as indicator of content flow across titles and borders. | Train and telegraph lines predict the travel of news items. | |
| **Bricolage** | Where does content in historical newspapers come from? | TR to identify the patterns and fragments which constitute content and their evolution over time. | Death notices, weather reports, event anniversaries | Reprints, verse reprints |
| **Virality** | Which content spreads? | TR as an indicator of the distribution of content within media across time, space, titles. | Paju et al. 2022 | |
| **Events** | How did journalists shape the content they published? | TR to reconstruct the changes in texts spreading across media along editorial lines. | Adjustments of press agency content: title, adjectives, omissions, additions. | Modified reuse |
| **Zeitgeist** | How did ideas co-evolve? | TR to trace co-evolving ideas. | "A shared way" to write text, design adverts | quotes |
| **Unattributed reuse\*** | Which content circulated without attribution? | TR as evidence for undeclared reuse. | Plagiarism | Unattributed reuse |
| **Data cleaning\*** | When does TR reduce data quality? | TR as indicators of unwanted duplicates | Mastheads, co-publication | Artefacts |

|  | Research questions | Purpose of TR | Examples |
|---|---|---|---|
| **Media ecosystems** | How did historical media function? How did transnational information flow shape historical media? | TR as indicator of content flow across titles and borders. | Train and telegraph lines predict the travel of news items. |
| **Bricolage** | Where does content in historical newspapers come from? | TR to identify the patterns and fragments which constitute content and their evolution over time. | Death notices, weather reports, event anniversaries |
| **Virality** | Which content spreads? | TR as an indicator of the distribution of content within media across time, space, titles. | Paju et al. 2022 |
| **Events** | How did journalists shape the content they published? | TR to reconstruct the changes in texts spreading across media along editorial lines. | Adjustments of press agency content: title, adjectives, omissions, additions. |
| **Zeitgeist** | How did ideas co-evolve? | TR to trace co-evolving ideas. | "A shared way" to write text, design adverts |
| **Unattributed reuse*** | Which content circulated without attribution? | TR as evidence for undeclared reuse. | Plagiarism |
| **Data cleaning*** | When does TR reduce data quality? | TR as indicators of unwanted duplicates | Mastheads, co-publication |

|  | Research questions | Purpose of TR | Examples |
|---|---|---|---|
| **Media ecosystems** | How did historical media function? How did transnational information flow shape historical media? | TR as indicator of content flow across titles and borders. | Train and telegraph lines predict the travel of news items. |
| **Bricolage** | Where does content in historical newspapers come from? | TR to identify the patterns and fragments which constitute content and their evolution over time. | Death notices, weather reports, event anniversaries |
| **Virality** | Which content spreads? | TR as an indicator of the distribution of content within media across time, space, titles. | Paju et al. 2022 |
| **Events** | How did journalists shape the content they published? | TR to reconstruct the changes in texts spreading across media along editorial lines. | Adjustments of press agency content: title, adjectives, omissions, additions. |
| **Zeitgeist** | How did ideas co-evolve? | TR to trace co-evolving ideas. | "A shared way" to write text, design adverts |
| **Unattributed reuse*** | Which content circulated without attribution? | TR as evidence for undeclared reuse. | Plagiarism |
| **Data cleaning*** | When does TR reduce data quality? | TR as indicators of unwanted duplicates | Mastheads, co-publication |

|  | **Research questions** | **Purpose of TR** | **Examples** |
|---|---|---|---|
| **Media ecosystems** | How did historical media function? How did transnational information flow shape historical media? | TR as indicator of content flow across titles and borders. | Train and telegraph lines predict the travel of news items. |
| **Bricolage** | Where does content in historical newspapers come from? | TR to identify the patterns and fragments which constitute content and their evolution over time. | Death notices, weather reports, event anniversaries |
| **Virality** | Which content spreads? | TR as an indicator of the distribution of content within media across time, space, titles. | Paju et al. 2022 |
| **Events** | How did journalists shape the content they published? | TR to reconstruct the changes in texts spreading across media along editorial lines. | Adjustments of press agency content: title, adjectives, omissions, additions. |
| **Zeitgeist** | How did ideas co-evolve? | TR to trace co-evolving ideas. | "A shared way" to write text, design adverts |
| **Unattributed reuse*** | Which content circulated without attribution? | TR as evidence for undeclared reuse. | Plagiarism |
| **Data cleaning*** | When does TR reduce data quality? | TR as indicators of unwanted duplicates | Mastheads, co-publication |

|  | **Research questions** | **Purpose of TR** | **Examples** |
|---|---|---|---|
| **Media ecosystems** | How did historical media function? How did transnational information flow shape historical media? | TR as indicator of content flow across titles and borders. | Train and telegraph lines predict the travel of news items. |
| **Bricolage** | Where does content in historical newspapers come from? | TR to identify the patterns and fragments which constitute content and their evolution over time. | Death notices, weather reports, event anniversaries |
| **Virality** | Which content spreads? | TR as an indicator of the distribution of content within media across time, space, titles. | Paju et al. 2022 |
| **Events** | How did journalists shape the content they published? | TR to reconstruct the changes in texts spreading across media along editorial lines. | Adjustments of press agency content: title, adjectives, omissions, additions. |
| **Zeitgeist** | How did ideas co-evolve? | TR to trace co-evolving ideas. | "A shared way" to write text, design adverts |
| **Unattributed reuse*** | Which content circulated without attribution? | TR as evidence for undeclared reuse. | Plagiarism |
| **Data cleaning*** | When does TR reduce data quality? | TR as indicators of unwanted duplicates | Mastheads, co-publication |

|  | Research questions | Purpose of TR | Examples |
|---|---|---|---|
| **Media ecosystems** | How did historical media function? How did transnational information flow shape historical media? | TR as indicator of content flow across titles and borders. | Train and telegraph lines predict the travel of news items. |
| **Bricolage** | Where does content in historical newspapers come from? | TR to identify the patterns and fragments which constitute content and their evolution over time. | Death notices, weather reports, event anniversaries |
| **Virality** | Which content spreads? | TR as an indicator of the distribution of content within media across time, space, titles. | Paju et al. 2022 |
| **Events** | How did journalists shape the content they published? | TR to reconstruct the changes in texts spreading across media along editorial lines. | Adjustments of press agency content: title, adjectives, omissions, additions. |
| **Zeitgeist** | How did ideas co-evolve? | TR to trace co-evolving ideas. | "A shared way" to write text, design adverts |
| **Unattributed reuse*** | Which content circulated without attribution? | TR as evidence for undeclared reuse. | Plagiarism |
| **Data cleaning*** | When does TR reduce data quality? | TR as indicators of unwanted duplicates | Mastheads, co-publication |

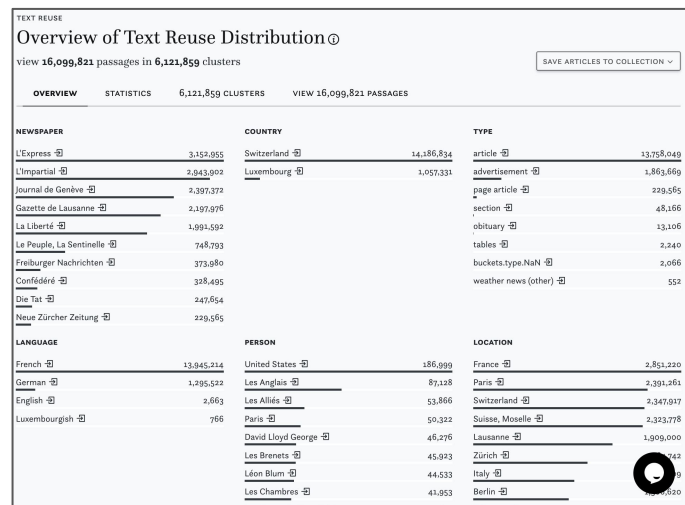|  | Research questions | Purpose of TR | Examples |
|---|---|---|---|
| **Media ecosystems** | How did historical media function? How did transnational information flow shape historical media? | TR as indicator of content flow across titles and borders. | Train and telegraph lines predict the travel of news items. |
| **Bricolage** | Where does content in historical newspapers come from? | TR to identify the patterns and fragments which constitute content and their evolution over time. | Death notices, weather reports, event anniversaries |
| **Virality** | Which content spreads? | TR as an indicator of the distribution of content within media across time, space, titles. | Paju et al. 2022 |
| **Events** | How did journalists shape the content they published? | TR to reconstruct the changes in texts spreading across media along editorial lines. | Adjustments of press agency content: title, adjectives, omissions, additions. |
| **Zeitgeist** | How did ideas co-evolve? | TR to trace co-evolving ideas. | "A shared way" to write text, design adverts |
| **Unattributed reuse*** | Which content circulated without attribution? | TR as evidence for undeclared reuse. | Plagiarism |
| **Data cleaning*** | When does TR reduce data quality? | TR as indicators of unwanted duplicates | Mastheads, co-publication |

|  | Research questions | Purpose of TR | Examples |
|---|---|---|---|
| **Media ecosystems** | How did historical media function? How did transnational information flow shape historical media? | TR as indicator of content flow across titles and borders. | Train and telegraph lines predict the travel of news items. |
| **Bricolage** | Where does content in historical newspapers come from? | TR to identify the patterns and fragments which constitute content and their evolution over time. | Death notices, weather reports, event anniversaries |
| **Virality** | Which content spreads? | TR as an indicator of the distribution of content within media across time, space, titles. | Paju et al. 2022 |
| **Events** | How did journalists shape the content they published? | TR to reconstruct the changes in texts spreading across media along editorial lines. | Adjustments of press agency content: title, adjectives, omissions, additions. |
| **Zeitgeist** | How did ideas co-evolve? | TR to trace co-evolving ideas. | "A shared way" to write text, design adverts |
| **Unattributed reuse*** | Which content circulated without attribution? | TR as evidence for undeclared reuse. | Plagiarism |
| **Data cleaning*** | When does TR reduce data quality? | TR as indicators of unwanted duplicates | Mastheads, co-publication |

*What do we want to do with text reuse data?*

| 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 | Cluster | Yes |
| 3 | Compare passages | Passages | 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 | Passages | No |
| 10 | De-duplicate content | Corpus | No |
| 11 | Export of text reuse data | All | Planned |

| 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 | Cluster | Yes |
| 3 | Compare passages | Passages | 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 | Passages | No |
| 10 | De-duplicate content | Corpus | No |
| 11 | Export of text reuse data | All | Planned |

**SEARCH TEXT REUSE PASSAGES**

- Lexical overlap between 76 and 10 ⌄
- "hiroshima" ⌄
- tr-nobp-all-v01-c180388965941 ⌄  ✕

hiroshima

**FILTER BY TEXT REUSE CLUSTERS (1 OPTION)** ⓘ    RESET

results are filtered when:

☑ cluster **c180388965941**

81.48% lexical overlap the same day (Thu, Aug 5, 2010).

Aucun président américain en exercice ne sest rendu dans les deux villes martyres . Mais d aucuns dans l archipel espèrent une visite de Barack Obama à Hiroshima en novembre .

(3 results)

| 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 | Cluster | Yes |
| 3 | Compare passages | Passages | 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 | Passages | No |
| 10 | De-duplicate content | Corpus | No |
| 11 | Export of text reuse data | All | Planned |

| 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 | Cluster | Yes |
| 3 | Compare passages | Passages | 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 | Passages | No |
| 10 | De-duplicate content | Corpus | No |
| 11 | Export of text reuse data | All | Planned |

| 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 | Cluster | Yes |
| 3 | Compare passages | Passages | 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 | Passages | No |
| 10 | De-duplicate content | Corpus | No |
| 11 | Export of text reuse data | All | Planned |



Media Monitoring impresso of the Past    Search ...    Newspapers    T

SEARCH TEXT REUSE PASSAGES

Text reuse time span between 55,: ∨
Lexical overlap between 13 and 87 ∨
vie · monde · mort · foi · peuple OF ∨    ×

hiroshima

NUMBER OF PASSAGES PER YEAR (i)
Number of text reuse passages per year

14

1820  1840  1860  1880  1900  1920  1940  1960  1980  2000

ADD NEW DATE FILTER …

CLUSTER SIZE (i)
How to read histograms (i)

4 (24 results)

2                                    44

LEXICAL OVERLAP (i)    RESET

18 (10 results)

16                                    85

# Job adverts - intentional distribution

# Rectification - 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          Rectificatif de l'ATS
L'Express  WEDNESDAY, OCTOBER 13,      Journal de Genève  THURSDAY,
1982 – P.2                             OCTOBER 14, 1982 – P.8

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

# Event anniversaries - reprints

# Quotes



**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 . »

**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)

# Edits of 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.                            Les suffragistes font de l'action directe
**Gazette de Lausanne** ⤢  MONDAY, NOVEMBER 28, 1910 –   **L'indépendance luxembourgeoise** ⤢  TUESDAY,
P.2                                                NOVEMBER 29, 1910 – P.2

M. Winston Churchill, ministre du commerce,        M. Winston Churchill, **revenant à Londres, après avoir**
prononçait ven- dredi soir à Bradford a été coupé de   prononcé un dis- cours à Bradford, où les suffragettes
plu- sieurs interruptions de suffragistes et de       et leurs partisans avaient déjà violemment manifesté
suffragettes qui ont été expulsés à tour de rôle.      contre lui, a été attaqué dans le train par un individu,
Comme l'orateur rentrait de Bradford à Londres, il fut   qui a essayé de le frapper avec une cravache en
attaqué dans le train par un individu, qui a essayé de   disant: «Voilà pour toi, chien!» Deux agents de police,
le frapper avec une cravache en disant : « Voilà pour   **qui accompagnaient M. Churchill, parèrent** le coup, et
toi, chien ! » Deux agents de police parèrent le coup et   s'em- parèrent de l'homme, après une lutte violente.
s'emparè- rent de l'homme après une lutte violente. A   **On croit que l'assaillant est un des suffragistes**
la gare de Londres, trois femmes ont également          **expulsés de la réunion où M. Winston Churchill**
essayé de frapper M. Churchill ; elles en ont été       **venait de parler.** A la gare de Londres, trois femmes
empêchées par les agents.                               ont également essayé de frapper M. Chur- chill ; elles
                                                        en ont été empêchées par les agents.

| 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 | Cluster | Yes |
| 3 | Compare passages | Passages | 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 | Passages | No |
| 10 | De-duplicate content | Corpus | No |
| 11 | Export of text reuse data | All | Planned |

**Create new collection** ✕

Please note: Collections are currently limited to 10.000 items. ⓘ
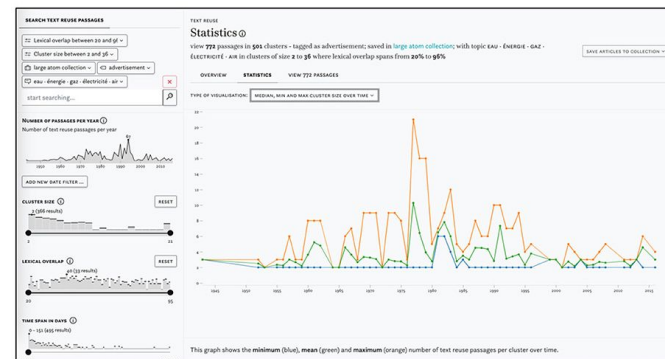
COLLECTION NAME

collection of bible quotes

DESCRIPTION

with topic vie · monde · mort · foi · peuple or problème · fait · question · exemple · monde where lexical overlap spans from

CANCEL   CREATE

| 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 | Cluster | Yes |
| 3 | Compare passages | Passages | 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 | Passages | No |
| 10 | De-duplicate content | Corpus | No |
| 11 | Export of text reuse data | All | Planned |

| 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 | Cluster | Yes |
| 3 | Compare passages | Passages | 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 | Passages | No |
| 10 | De-duplicate content | Corpus | No |
| 11 | Export of text reuse data | All | Planned |

| 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 | Cluster | Yes |
| 3 | Compare passages | Passages | 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 | Passages | No |
| 10 | De-duplicate content | Corpus | No |
| 11 | Export of text reuse data | All | Planned |

ANGLETERRE

**Journal de Genève** ⊡ WEDNESDAY, SEPTEMBER 12, 1883 – P.2

> Le dernier courrier du Japon nous ap- prend qu'un inceudie a détruit eu grande partie la prison de Hiroshima. Soixante uu détenus ont été bru es vifs, et 156 out été plus ou moius grièvement blessés. Profitant de. la. confusion, 120 prisonniers se sont enfuis, et 15 d'entre eux. sf. ulemtuit ont pu être recapturés.
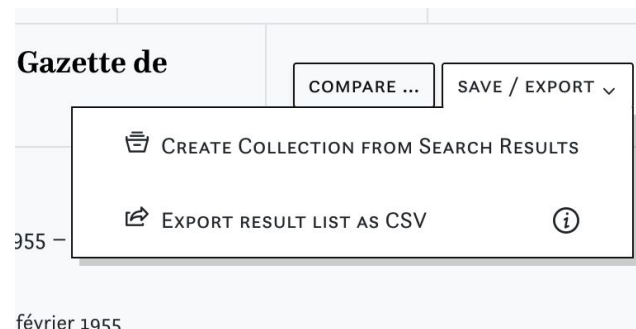
Cluster size: **3** passages with **58.33%** lexical overlap over **3** days.
Content length: **317** tokens.

COMPARE WITH OTHER PASSAGES IN CLUSTER

| 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 | Cluster | Yes |
| 3 | Compare passages | Passages | 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 | Passages | No |
| 10 | De-duplicate content | Corpus | No |
| 11 | Export of text reuse data | All | Planned |

| 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 | Cluster | Yes |
| 3 | Compare passages | Passages | 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 | Passages | No |
| 10 | De-duplicate content | Corpus | No |
| 11 | Export of text reuse data | All | Planned |

**Gazette de**

COMPARE ...    SAVE / EXPORT ⌄

⊟ CREATE COLLECTION FROM SEARCH RESULTS

↱ EXPORT RESULT LIST AS CSV    ⓘ

955 —

février 1955

# 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. |

Future development of the application should focus on these tasks / features / overall improvements: *

I would repeat what I wrote earlier--more historical/humanistic contextualization is needed to help users understand how the mechanics of text reuse, as represented in the variables users can use to filter data, reflect distinct historical genres, texts, technologies of reproduction, etc.

# 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><br>1815 |
| 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><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><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><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 |

# Data-driven text reuse classification?

|  | 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 |

# Finally..

Do we need a(nother) universal typology of historical TR types?

Do TR types have distinct properties (lex overlap, temporality, cluster/passage size, links to other sem enrichments) across source types?

Can we classify them analogue to the way we can predict content types?
E.g. quotes, unattributed reuse, modified reuse, rhythmic

Impresso use case: "Which types of TR can I find in my collection of articles?"

Can we link texts between between spoken and written text and languages?

How do we do TRD in the LLM age?

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

09/2023 - 02/2027

Main objectives

1.  **Enrichment and integration** of newspaper and radio sources
2.  **Expand** the corpus to Western Europe
3.  **Develop** interfaces for exploratory and computational research
4.  **Conduct** case studies in (media) history, theme "influences"

# Partners

**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

# 1 Source collection

## European media archives

AUSTRIA

BELGIUM

FRANCE

GERMANY

LUXEMBOURG

THE NETHERLANDS

SWITZERLAND

UK

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

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

# 2 Media processing

## Enriching & connecting

SEMANTIC ENRICHMENT
ACROSS LANGUAGES
ACROSS MEDIA

*dense vector representations*

topic
place
person
institution

interview
advertisement
radio schedule

**EXTERNAL KNOWLEDGE**

# 3 Media exploration

## 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

**IMPRESSO WEB APP**

historians
designers
engineers

**IMPRESSO DATA LAB**
*executable notebooks*

*community of historians*

CO-DESIGNED INTERFACES

CO-DEVELOPED METHODS

# Thank you



impresso website (update soon)



data on Zenodo



code on GitHub



impresso youtube channel



impresso-project.ch/app