# Bigger Smarter Data

# Extracting, Modeling and Linking Data for Literary History

Christof Schöch
(Trier University, Germany)

Korea University
Seoul, South Korea

23 May 2024

# Introduction

# Thanks

Korea University, as well as KADH (Korean Association for Digital Humanities).

The Ministry for Research and Education in Rhineland-Palatinate, Germany, for funding this research (Mining and Modeling Text, 2019-2023)

Thanks to all the project contributors: Maria Hinzmann, Matthias Bremm, Tinghui Duan, Anne Klee, Johanna Konstanciak, Julia Röttgermann and many others.
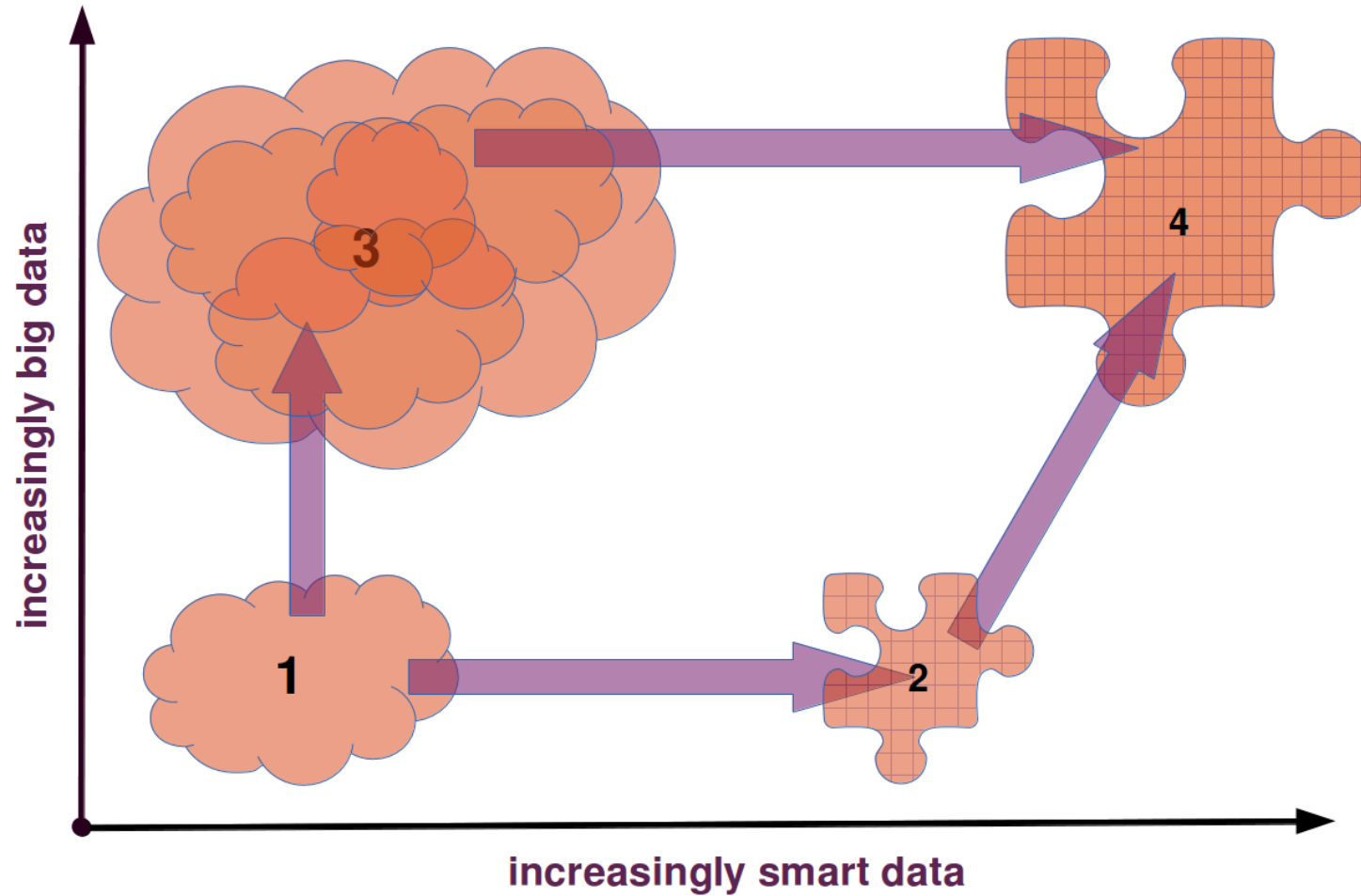
# Overview

Kompetenzzentrum
Trier Center for Digital Humanities

# Bigger Smarter Data:
# Linked Open Data

Kompetenzzentrum
Trier Center for Digital Humanities

# Three Modes of Data (and Digital Humanities)

- Qualitative DH:

    - Datasets are typically small, curated, heavily annotated, flawless, specialized ('smart data')

    - Prototype: digital scholarly editions, e.g. *Faust Edition*

- Quantitative DH:

    - Datasets are typically large, scraped, unannotated, with errors and biases, generic ('big data')

    - Prototype: Large Language Models trained on lots of text, e.g. ChatGPT

Kompetenzzentrum
Trier Center for Digital Humanities

# Third Way: Bigger Smarter Data

# Background: What is Machine Learning?

- Fundamentally, ML involves detecting relations between features and labels

  - Features that we can observe in data

  - Labels, classes, or values that are relevant to our research

- We use this approach primarily for information retrieval

  - We start from a text collection

  - We may annotate part of the data, or use an existing model

  - And then derive labels, classes, values from the text

Kompetenzzentrum
Trier Center for Digital Humanities

# Background: What is Linked Open Data?

# Multilingualism

# Background: What is Literary History?

- Goals of literary history

  - Collecting and documenting knowledge of literary history

  - Providing explanations for the development of literature

- Organizational principles

  - Nations, periods, movements/currents, genres

  - Authors and works

  - Similarities and differences, continuities and change

- Functions

  - Explanations for literary developments

  - a. by cultural or socio-historical context

  - b. by internal dynamics of the literary system

Kompetenzzentrum
Trier Center for Digital Humanities

# The project 'Mining and Modeling Text'

# Literary History in Linked Open Data

- Building blocks

  - Subjects, including persons (author, etc.) and works (primary text, scholarly literature, etc.)

  - Objects, including works, but also themes, locations, protagonists, literary genre, etc.

  - Predicates, as required, including: author_of, about, sameAs etc.

  - Qualifications, e.g: Source (with type, date, URL)

- Some exemplary statement types

  - Bibliographic: `[person] author_of [work]`

  - Contentual: `[work] about [theme]`

  - Formal: `[work] narrative_form [type]`

  - and many more.

Kompetenzzentrum
Trier Center for Digital Humanities

# Wikidata for Literary History

- Idea: Create a "Wikidata for the history of literature"

  - Literary history information system

  - LOD-based, with explorative interface and SPARQL endpoint

  - Approach of an "atomization" of the historical knowledge

  - Linking with other knowledge systems (taxonomies, standard data, knowledge bases)

  - Key values: human and machine readable, open, collaborative, multilingual

- Compared to Wikidata

  - Focused on one domain (French novel, 1750-1800)

  - Better coverage / higher density of information for this domain

  - Development of a systematic ontology

  - much smaller: 300k vs. 1.5 billion statements

Kompetenzzentrum
Trier Center for Digital Humanities

# Mining: Information Extraction

# Pillar 1: *Bibliographie du genre romanesque français*

**59.25**    **VOLTAIRE, François-Marie Arouet de**

Candide ou l'Optimisme, traduit de l'allemand de Mr. le docteur Ralph
1759, in-12
BN
AL 1759 II 203-210; AT 1761 (1759); CorrL mars 1759
Bengesco Dufrenoy Gay Morize Q
Il paraît y avoir eu jusqu'à une vingtaine d'éditions datées de 1759. Sur la question de la véritable édition *princeps,* voir Bengesco; Morize; I.O. Wade, *Voltaire et Candide,* Princeton, 1959; B. Gagnebin, ds *Bulletin du bibliophile,* 1960, pp. 22-31; J.-D. Candaux, ds *Studies on Voltaire,* XVIII, 1961, pp. 173-178.

*3e personne; Europe, Amérique; Candide, Cunégonde, Pangloss, Martin; voyages, aventures romanesques, désastres; thèmes philosophiques, ton satirique.*

Autres éditions:
— s.l., 1759. Bengesco donne 10 éditions s.l. 1759; Morize en cite 12; selon Besterman il y aurait une vingtaine d'éditions portant la date de 1759.
— Londres, 1759 (Bengesco, Morize)
— s.l., 1760 (Morize donne une édition; Bengesco en donne deux)
— s.l., 1761 (Bengesco)
— Genève, 1761 (Morize)
— ds *Seconde suite des Mélanges,* 1761 (Bengesco, Morize)
— Aux Délices, 1763 (Bengesco, Morize)

# Pillar 2: primary literature (novels)



- Corpus of 200 French novels (1750-1800)

- Coding: in XML-TEI, with metadata, according to ELTeC schema

- Analysis methods: Topic modeling, NER, stylometry, etc.

*Collection of Eighteenth-Century French novels (1750-1800)*, ed. Julia Röttgermann. See

Kompetenzzentrum
Trier Center for Digital Humanities

# Pillar 3: Scholarly Literature



- Annotation Guidelines (based on the data model)

- Manual annotations (using INCEpTION)

- Linking of INCEpTION with MiMoTextBase and Wikidata => disambiguation

- Creation of statements about authors and works (genres, themes, etc.)

Kompetenzzentrum
Trier Center for Digital Humanities

# Modeling: Data Modeling

# Modular Data Model

- Module 1: Theme

- Module 2: Space

- Module 3: Narrative form

- Module 4: Literary work

- Module 5: Author

- Module 6: Mapping

- Module 7: Referencing

- Module 8: Versioning & publication

- Module 9: Terminology

- Module 10: Bibliography

- Module 11: Scholarly literature

Kompetenzzentrum
Trier Center for Digital Humanities

# Example: The module on themes

MODULE 1: theme

[Bibliographie du genre romanesque français]

[BGRF_Matching-Table]

[topic labels and concepts 11-2020]

topic model

https://github.com/MiMoText/[...]

https://zenodo.org/record/[... ]

[scholarly work]

rdfs:range

stated in

rdfs:range rdfs:range rdfs:range

described at URL

instance of

described at URL

topic

literary work

https://schema.org/about

[wdt:P921=main subject]

[author]

rdfs:domain

topic interest

[author]

rdfs:range

[author]

rdfs:range

[literary work]

rdfs:range

rdfs:range

[spatial concept]

[thematic concept]

exact match

http://xmlns.com/foaf/spec/#term_topic_interest

exact match
close match

rdfs:domain

owl:disjointWith

about

rdfs:domain=
statement / claim

rdfs:range

[topic model]

part of

[topic]

instance of

image

[image: wordl]

represented by

related to

[skos:altLabel
@en | @fr | @de]

[string]

rdfs:range

[author]

rdfs:range

[work]

rdfs:range

thematic concept

instance of

[thematic concept]

exact match

[wikidata identifier]

[spatial concept]

[rdfs:label
@en | @fr | @de]

part of

stated in

[string]

thematic vocabulary MiMoText

described at URL

https://github.com/MiMoText/

Kompetenzzentrum
Trier Center for Digital Humanities

# Example: The module on narrative location



'spatial vocabulary'
(=mmt:Q25)

campagne |
rural area |
ländlicher Raum

wdt:Q175185
WIKIDATA

described at URL

https://github.com/MiMoText/vocabularies/blob/main/Raumvokabular.tsv

stated in

rdfs:label @fr | de | en

exact match

part of

wdt:Q90
WIKIDATA

'BGRF'
(=mmt:Q1)

stated in

mmt:Q3223

'spatial concept'
(=mmt:Q26)

exact match
(P31)

reference URL
(=P18)

stated in

stated in

narrative_location

instance of

'La religieuse'
(=novel by Denis Diderot;
mmt:Q3730)

narrative_location

'Paris'
(=mmt:Q3508)

coordinate location

48.856944444444,2.3513888888889

"Paris, province"

narrative_location
_string

narrative_location
_string

"Longchamp"

narrative_location
_string

"Paris"

stated in

https://github.com/MiMoText/roman18/tree/master/NER-with-SpaCy

occurence in text

occurence in text

stated in

stated in

"22"

"30"

'NER_novels locations'
(= mmt:Q3730)

described at URL

Kompetenzzentrum
Trier Center for Digital Humanities

# Meta-Statements

| about | ⬍ libertinism | |
|---|---|---|
| | ▾ 2 references | |
| | stated in | Bibliographie du genre romanesque français |
| | stated in | BGRF_matching-table (03-2022) |
| | ⬍ correspondence | |
| | ▾ 2 references | |
| | stated in | Topic Model MMT 11-2020 |
| | stated in | topic labels and concepts (11-2020) |

Kompetenzzentrum
Trier Center for Digital Humanities

# Linking with Wikidata for 'federated queries'



MiMoTextBase

rdfs:label → "DIDEROT, Denis"

http://data.mimotext.uni-trier.de/wiki/Item:Q306

exact match

https://www.wikidata.org/wiki/Q448

Wikidata

"Denis Diderot"

rdfs:label

https://www.wikidata.org/wiki/Q448

date of birth → 5 October 1713

Federated Queries

(idea of visualization, see Abel 2019, 5)

Kompetenzzentrum
Trier Center for Digital Humanities

# Result: Queryable Database

Kompetenzzentrum
Trier Center for Digital Humanities

# The MiMoTextBase

Log in

| Main Page | Discussion | | Read | View source | View history | Search MiMoText 🔍 |

**Main page**
**Recent changes**
**Random page**
**Help about MediaWiki**

Tools

**What links here**
**Related changes**
**Special pages**
**Printable version**
**Permanent link**
**Page information**

In other languages

✏ Add links

## Main Page

**Mining and Modeling Text: Interdisciplinary applications, informational development, legal perspectives (MiMo Text)**

The acquisition of knowledge from large amounts of text and data which can no longer be handled by individuals is becoming increasingly important due to the possibilities of digitisation. For the humanities, this means in particular that digital full texts and rich metadata must not only be available, but must also be available in a form that promotes knowledge in the humanities.

The aim of the MiMoText project is therefore to establish an information network for the humanities fed from various sources, which, by making it available as Linked Open Data, is not only freely available and can be linked to other knowledge resources of the Semantic Web, but also offers innovative and efficient access possibilities to scientific information.

MiMoTextBase was built as part of the project "Mining and Modeling Text" (2019-2023). It is implemented using a Wikibase infrastructure and integrates various data from heterogeneous sources. Note that the project is ongoing and the contents and structure of the MiMoTextBase will continually be further developed.
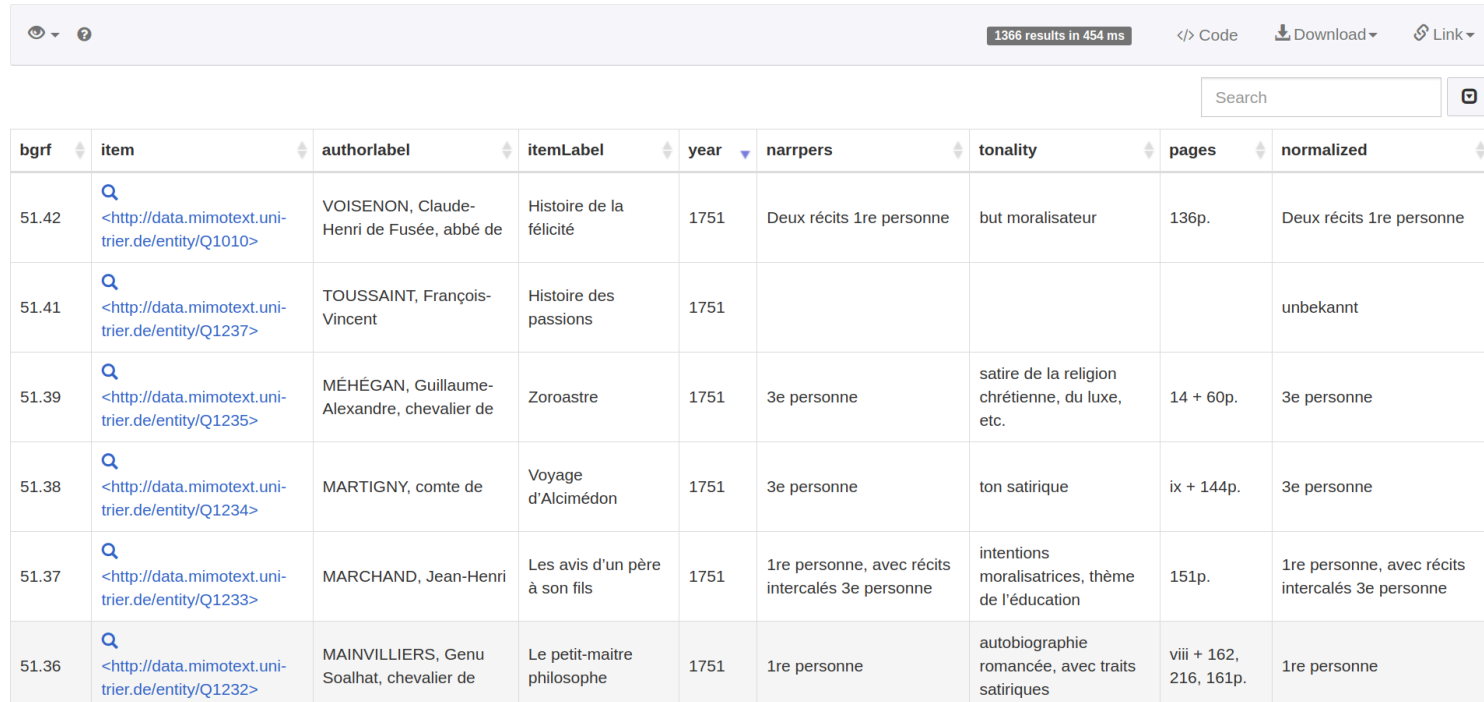
**Key starting points**

- Tutorial and further information about the MiMoTextBase: https://docs.mimotext.uni-trier.de⧉
- SPARQL endpoint: https://query.mimotext.uni-trier.de⧉
- MiMoTextBase (current page): https://data.mimotext.uni-trier.de⧉
- Project homepage: https://mimotext.uni-trier.de/en⧉
- Ontology repository: https://github.com/MiMoText/ontology⧉

**Example pages**

- An author entry (Révéroni Saint-Cyr): http://data.mimotext.uni-trier.de/wiki/Item:Q851⧉
- A title entry (Liaisons dangereuses): http://data.mimotext.uni-trier.de/wiki/Item:Q1053⧉
- A thematic concept (travel): http://data.mimotext.uni-trier.de/wiki/Item:Q3126⧉
- A spatial concept (Geneva): http://data.mimotext.uni-trier.de/wiki/Item:Q3478⧉

**Kompetenzzentrum**
Trier Center for Digital Humanities

# SPARQL endpoint



- SPARQL = SPARQL Protocol and RDF Query Language

- Used to formulate queries

# MiMoText Base: Query for themes in novels

# Some sample queries: simple queries

- List of novels with information from BGRF

- Number of works per author (first 25))

- Themes of novels, in French and in English

Kompetenzzentrum
Trier Center for Digital Humanities

# Example queries: visualizations

- Number of novels per year

- Narrative forms over time (decades)

- Book history: print formats over time (5 years)

Kompetenzzentrum
Trier Center for Digital Humanities

# Sample queries: networked and federated

- Link with catalogue data from French National Library (using BNF id)

- Narrative locations of novels (map)

- Authors by birth year, with portrait)

- Alternative author names from Wikidata infobox

- Network of influences between authors (using 'influenced by')

- Querying MiMoText from Wikidata (it works both ways)

- Novels and basic information, from Wikidata

Kompetenzzentrum
Trier Center for Digital Humanities

# Sample queries: comparative queries

- Themes from topic modeling compared to themes in BGRF

- Themes from BGRF vs. Topic Modeling (in one query)

Kompetenzzentrum
Trier Center for Digital Humanities

# Conclusion

# Opportunities & challenges

- Opportunities

    - Linking heterogeneous data from different types of sources

    - Modeling, collecting and comparing contradictory statements

    - Transparency in knowledge production (sources)

- Challenges

    - Lack of consensus on relevant statement types in the discipline

    - Complexity reduction (triple structure)

    - Interoperability (tension 'Wikiverse' vs. OWL standard)

Kompetenzzentrum
Trier Center for Digital Humanities

# Lessons Learned

- Federated queries

  - Central element of the LOD vision

  - => Making it happen is not trivial (data model, infrastructure)

- Modeling meta statements

- Very important: perspectives / statements, not facts

- => Very different approaches in different technical contexts

- Exchange across communities

  - Literary Studies vs. Digital Humanities vs. Wikiverse

  - => is essential but needs more development

- There is still so much to do!

  - => We are continuing this effort in a new project called 'Linked Open Data in the Humanities' (LODinG)

Kompetenzzentrum
Trier Center for Digital Humanities

# Many thanks for your kind attention

# Further resources

- Tutorial: https://docs.mimotext.uni-trier.de

- SPARQL endpoint: https://query.mimotext.uni-trier.de

- MiMoTextBase: https://data.mimotext.uni-trier.de

- MiMoText Ontology: https://github.com/MiMoText/ontology

- Reference publication: 'Smart Modeling for Digital Literary History'

- Overview of visuals: mimotext.github.io/MiMoTextBase_Tutorial/visualizations.html

Kompetenzzentrum
Trier Center for Digital Humanities

# References

Martin, Mylne, and Frautschi. 1977. *Bibliographie Du Genre Romanesque Français, 1751-1800*. Mandell.

Röttgermann, Julia. 2024. "The Collection of Eighteenth-Century French Novels 1751-1800." *Journal of Open Humanities Data* 10 (1): 31. https://doi.org/10.5334/johd.201.

Schöch, Christof. 2013. "Big? Smart? Clean? Messy? Data in the Digital Humanities." *Journal of Digital Humanities* 2 (3): 1–19. https://journalofdigitalhumanities.org/2-3/big-smart-clean-messy-data-in-the-humanities/.

Schöch, Christof, Maria Hinzmann, Julia Röttgermann, Katharina Dietz, and Anne Klee. 2022. "Smart Modelling for Literary History." *International Journal of Humanities and Arts Computing* 16 (1): 78–93. https://doi.org/10.3366/ijhac.2022.0278.

Kompetenzzentrum
Trier Center for Digital Humanities