

CHR 2025

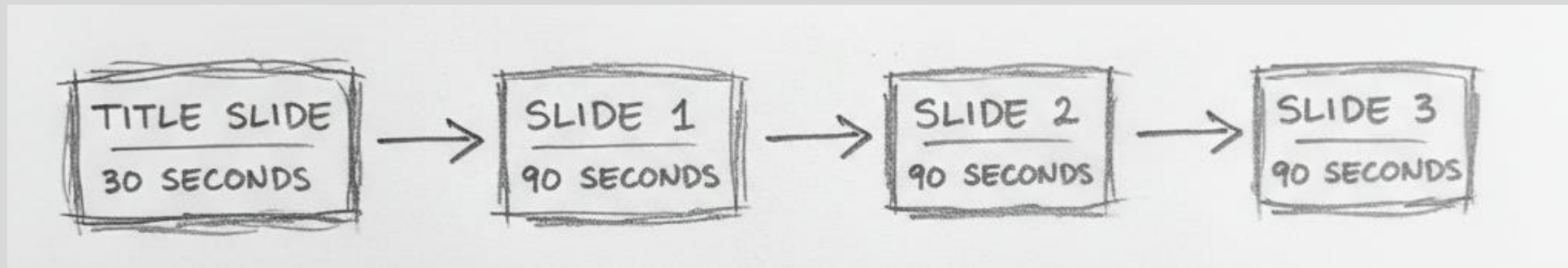
Lightning Talks

Chair: Taylor Arnold

Welcome!

While all forms of papers had increased submissions this year, the lightning talks had by far the largest growth and ultimately lowest acceptance rate.

We have a fantastic set of 17 presentations. To make sure everyone has time to present, we will use the following format with automatically advancing slides:

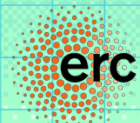


A static PDF version of these slides is available on the CHR 2025 website.

Bringing together close reading questions and distant reading methods in the analysis of archived web

Victor Harbo Johnston, Helle Strandgaard Jensen and Sasch Berg Bogeberg

WEB CHILD



European Research Council
Established by the European Commission

Contact Us

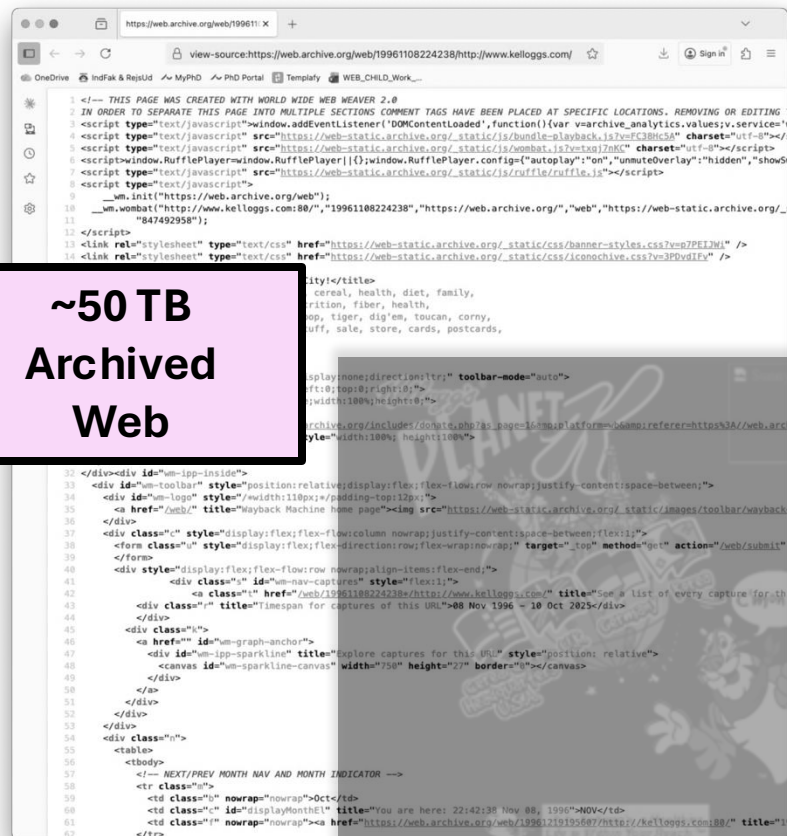


Sascha Berg **Bøgebjerg**
Helle Strandgaard **Jensen**
Victor Harbo **Johnston**
Aarhus University

Print



~50 TB
Archived
Web



Oral History

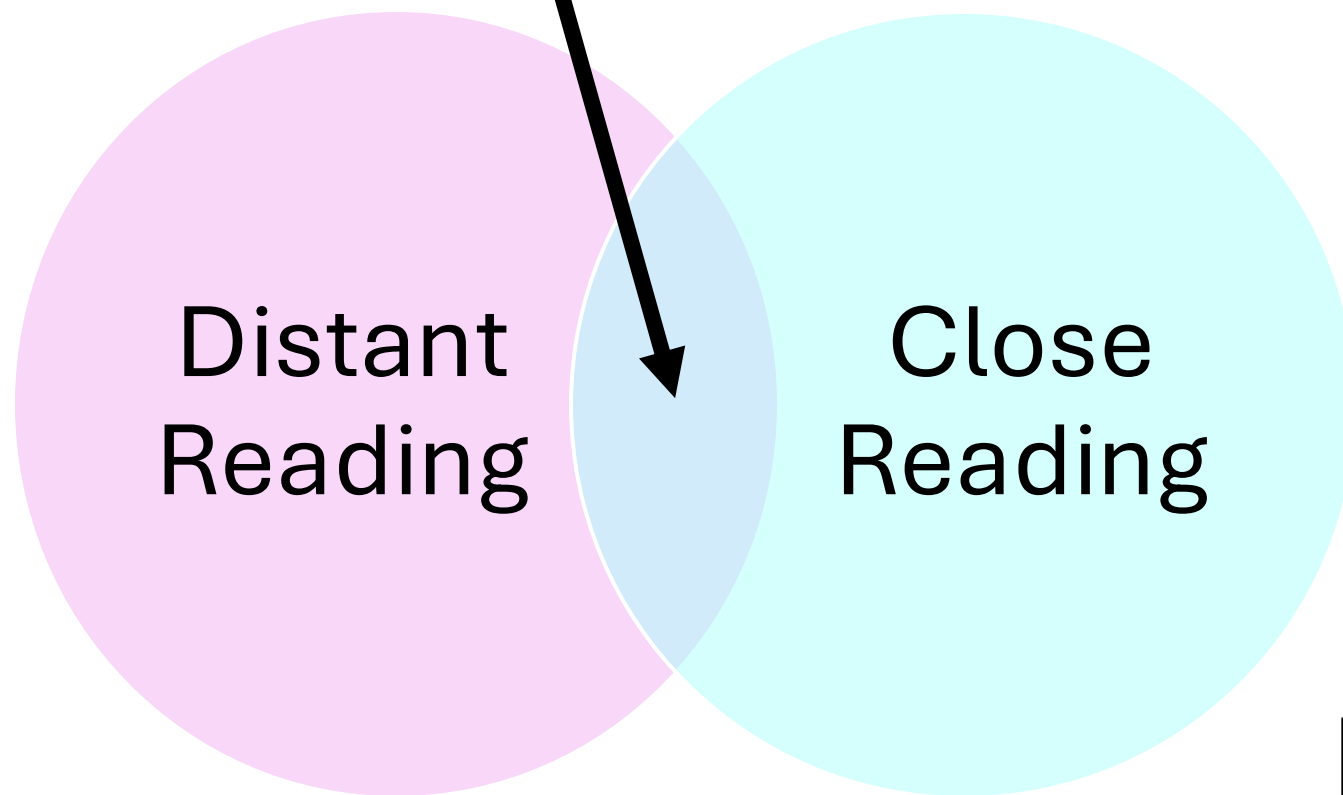


Survey



WEB CHILD

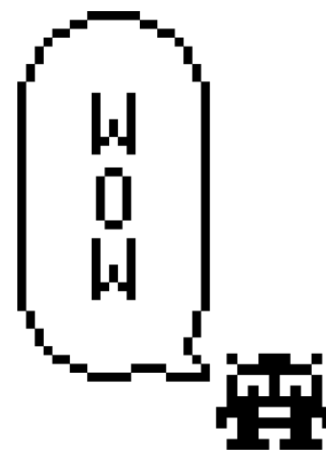
Scalable Reading



Contact Us



Sascha Berg **Bøgebjerg**
Helle Strandgaard **Jensen**
Victor Harbo **Johnston**
Aarhus University



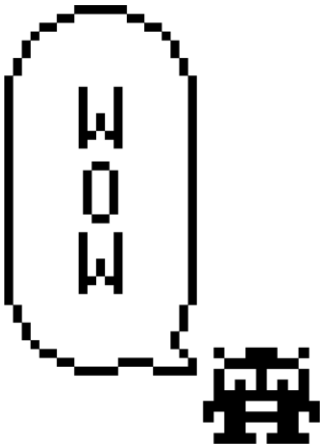
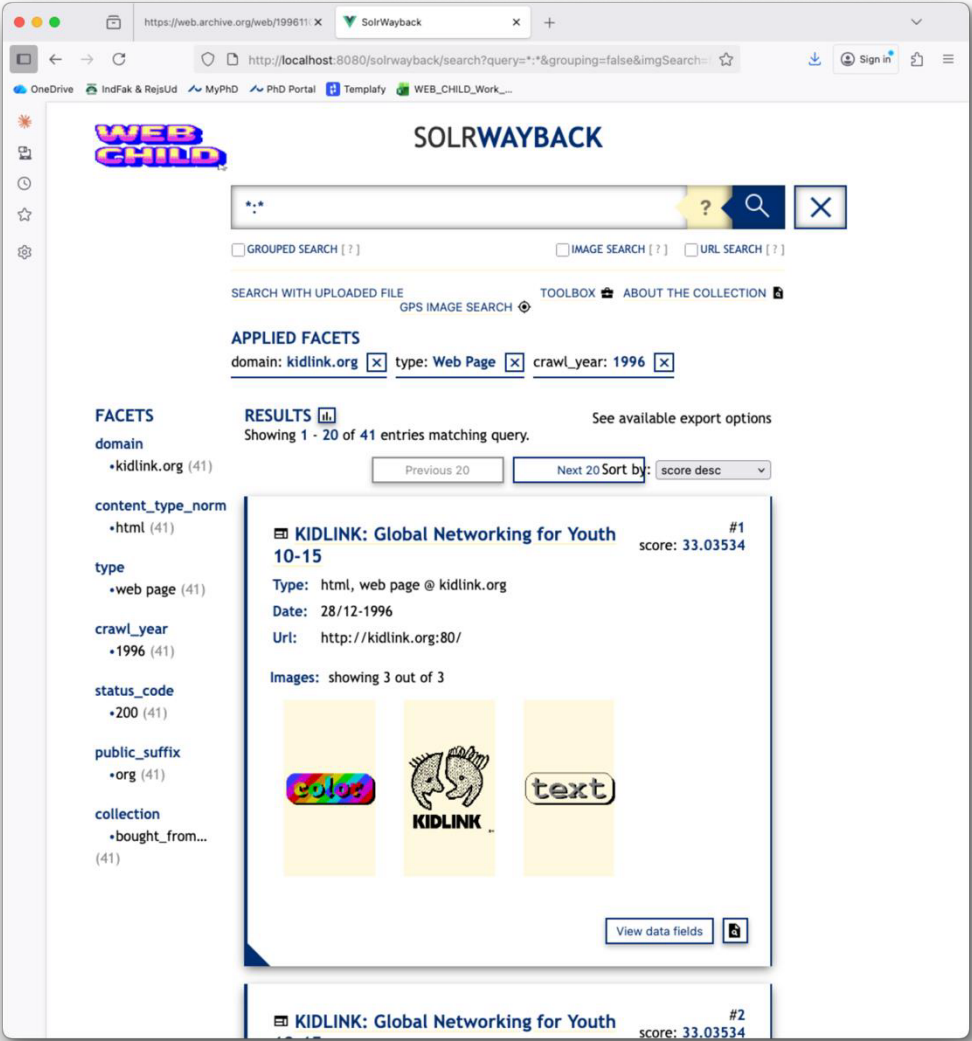
WEB CHILD

Contact Us



Sascha Berg **Bøgebjerg**
Helle Strandgaard **Jensen**
Victor Harbo **Johnston**
Aarhus University

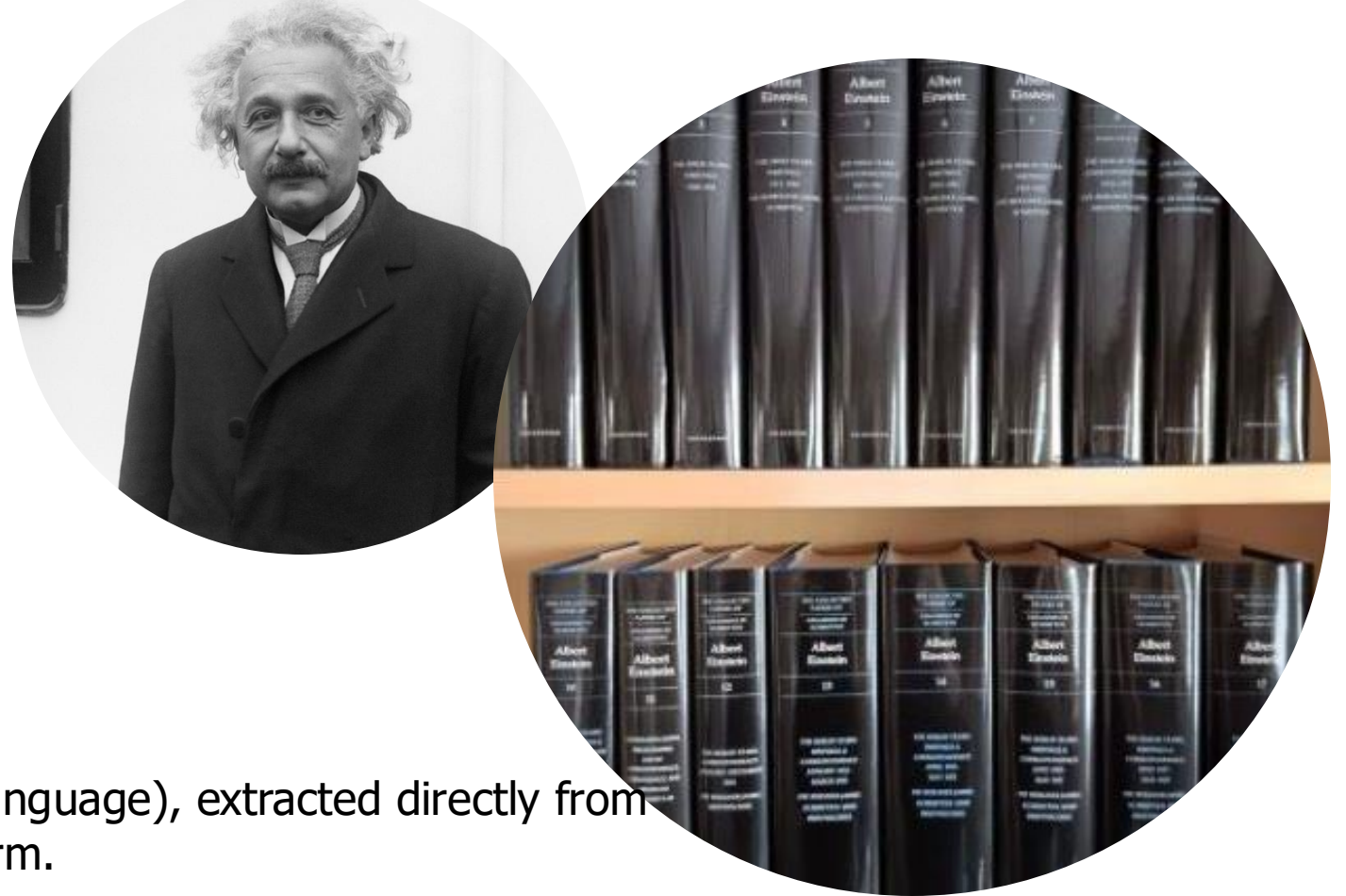
Preliminary Solutions



Einstein AI: Contextual Retrieval from the Collected Papers of Albert Einstein Using RAG and GraphRAG Architectures

Florin-Stefan Morar

- *The Collected Papers of Albert Einstein* (CPAE), published by Princeton University Press. The definitive scholarly edition, established in 1977 and publishing since 1987.
- A massive ongoing project with over 30,000 unique documents total. Currently spans 17 volumes (up to 1930).
- Our dataset focuses on the “Early Berlin Years”, Includes *Annus Mirabilis* (1905), the completion of General Relativity (1915), and Einstein’s rise to fame (1919)
- Primary Text: The German Edition (original language), extracted directly from the open-access *Digital Einstein Papers* platform.
- Translations: English versions were integrated via custom OCR of the supplementary English volumes to ensure cross-lingual retrieval.



Research question comes from an effort to bridge history of science and LLMs

How did Einstein's personal life and social context lead to the development of his scientific ideas?

Approach

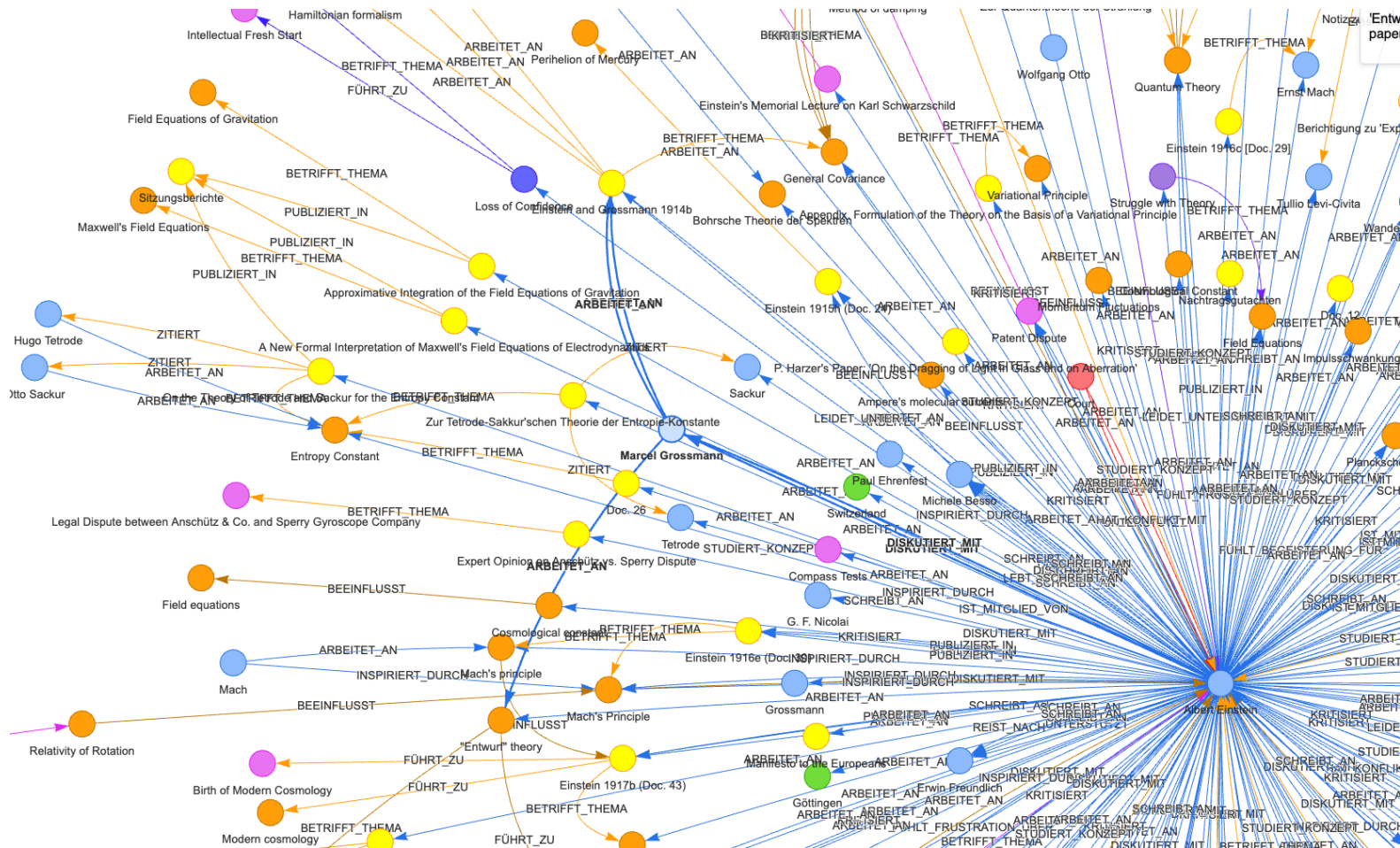
- A Custom Hybrid Pipeline for Historical Reasoning
- **Foundation:** Adapted the Microsoft GraphRAG architecture (specifically "Local Search" patterns) but re-engineered for the Google ecosystem. Replaced default extractors with Google Gemini 2.5 Pro to handle the complex, multilingual schema. Custom pipeline using NetworkX for graph management, FAISS for vector indexing, and SentenceTransformers (all-MiniLM-L6-v2) for embeddings.
- **Graph building schema:** A dual ontology, to simultaneously extract rigorous intellectual data (the "Science") and nuanced interpersonal dynamics. This allows us to map not just what Einstein discovered, but who he struggled with, what he felt about his work, and how his ideas evolved through social interaction. The Graph is built out of ENTITIES (person, organization, institution, location concept, method, theory) and RELATIONSHIPS (Interpersonal: colab with, wrote to etc.; Institutional: studied at etc.)
- Technical process: 1200 char chunking, Gemini 2.5 pro to extract entities using schema, NetworkX graph and stored as .pickle file for retrieval
- **How the Query Works (The "Graph-Aware" Retrieval):**
 - Vector Seeding: The user's query is embedded and matched against the vector index to find the top-k "Seed Nodes" (specific relevant document chunks).
 - Graph Expansion: The system uses the graph structure to "walk" from these seed nodes to their neighbors (1-2 hops).
 - *Result:* This retrieves semantically linked content (e.g., a letter *about* a concept) that lacks shared keywords, solving the "lexical gap."
 - Community Clustering: Retrieved nodes are clustered (using DBSCAN) to group thematically related evidence before the LLM synthesizes the final answer.

```
##ENTITY TYPES (Use these categorically)
- ##PERSON: Key individuals (e.g., Albert Einstein, Mileva Marić, Marcel Grossmann)
- ##ORGANIZATION: Schools, workplaces, organizations (e.g., ETH Zürich, Swiss Patent Office, University of Zurich)
- ##INSTITUTION: Similar to organization but for formal institutions (e.g., Kaiser Wilhelm Institute)
- ##LOCATION: Cities, countries, regions (e.g., Ulm, Bern, Zurich, Princeton)
- ##THEORY: Specific theories (e.g., Special Relativity, General Relativity, Quantum Theory)
- ##CONCEPT: Scientific concepts and ideas (e.g., spacetime, photoelectric effect, quantum mechanics, ether)
- ##METHOD: Mathematical or experimental methods (e.g., tensor calculus, thought experiment)
- ##PUBLICATION: Books, papers, journals (e.g., "Zur Elektrodynamik bewegter Körper", Annalen der Physik)
- ##PAPER: Specific research papers (e.g., "On the Electrodynamics of Moving Bodies")
- ##EXPERIMENT: Scientific experiments (e.g., Michelson-Morley experiment)
- ##MILESTONE: Significant life or professional milestones (e.g., graduation, marriage, Nobel Prize, World War II)
- ##DATE: Specific dates or time periods (e.g., 1885, March 14, 1953)
- ##FAMILY_MEMBER: Family relationships (e.g., Hans Albert Einstein, Elsa Einstein)
- ##HONOR: Honors and awards (e.g., Nobel Prize in Physics)
- ##BOOK: Published books (e.g., "The Meaning of Relativity")
- ##LETTER: Correspondence (e.g., "Letter to Mileva Marić, 1893")
- ##PATENT: Patents (e.g., "Electromagnetic Pump Patent")
- ##FIELD: Academic fields (e.g., Physics, Mathematics)
- ##DISCOVERY: Discoveries (e.g., Brownian motion, Bose-Einstein condensate)
- ##STRUGGLE: Abstract themes (e.g., financial struggle, jobsearch, antisemitism, scientific collaboration)
- ##SENTIMENT: Emotional states expressed (e.g., Frustriert, Zuversichtlich, Einsam) - use sparingly, only when explicitly stated
```

```
##RELATIONSHIP TYPES (Use English or German as appropriate)
- ##interpersonal:
  - collaborated_with, discussed_with, wrote_to, had_conflict_with, supports, loves, misses, meets
  - assisted_by, supervised_by, mentored_by, mentee_of, advisor_of, student_of, teacher_of, tutor
- ##institutional:
  - studied_at, works_for, is_member_of, taught_at, lectured_at
  - student_of, assistant_of, tutor_of, lecturer_of
- ##intellectual:
  - influenced_by, authored_by, proposed_by, criticized, inspired_by, studies_concept, works_on, published_in, cites
  - co-discovered, co-authored, co-proposed, co-criticized, co-inspired, co-studies, co-works, co-publishes, co-cites
- ##experiential:
  - traveled_to, born_in, lives_in, travels_to
  - left_from, moved_to
- ##social/romantic:
  - knew_about, influenced, concerns, themes, contributes_to, based_on, supported_by, opposed_by
  - friend_of, mentor_of, mentee_of, colleague_of
- ##emotional (use sparingly):
  - feels_frustration_about, feels_enthusiasm_for, suffers_from
  - finds_frustrating_about, finds_inspiring_about, finds_challenging_about, finds_interesting_about
- ##scientific:
  - discovered, invented, developed_by, experimented_on, referenced_in, reviewed_by, translated_by, patented_by
- ##family:
  - parent_of, child_of, spouse_of, sibling_of
- ##other:
  - related_to, married, received, supervised_by, mentored, corresponded_with, mentor_of, founded, presented_at, attended_by
```

```
##EXAMPLES
- ##PERSON:
  - name: "Albert Einstein"
  - type: "PERSON"
  - description: "A physicist who developed the theory of general relativity."
  - roles: ["Physicist"]
  - relationships: ["collaborated_with", "wrote_to", "had_conflict_with", "assisted_by", "supervised_by", "mentored_by", "mentee_of", "advisor_of", "student_of", "teacher_of", "tutor"]
- ##ORGANIZATION:
  - name: "ETH Zürich"
  - type: "ORGANIZATION"
  - description: "A Swiss federal technical university."
  - roles: ["University"]
  - relationships: ["studied_at", "works_for", "is_member_of", "taught_at", "lectured_at"]
- ##LOCATION:
  - name: "Ulm"
  - type: "LOCATION"
  - description: "A city in Germany."
  - roles: ["City"]
  - relationships: ["born_in", "lives_in", "travels_to"]
- ##THEORY:
  - name: "Special Relativity"
  - type: "THEORY"
  - description: "A theory of physics developed by Albert Einstein."
  - roles: ["Theory"]
  - relationships: ["influenced_by", "authored_by", "proposed_by", "criticized", "inspired_by", "studies_concept", "works_on", "published_in", "cites"]
- ##CONCEPT:
  - name: "Spacetime"
  - type: "CONCEPT"
  - description: "A concept in physics representing the four-dimensional continuum of space and time."
  - roles: ["Concept"]
  - relationships: ["influenced_by", "authored_by", "proposed_by", "criticized", "inspired_by", "studies_concept", "works_on", "published_in", "cites"]
- ##METHOD:
  - name: "Tensor Calculus"
  - type: "METHOD"
  - description: "A mathematical framework used in general relativity."
  - roles: ["Method"]
  - relationships: ["influenced_by", "authored_by", "proposed_by", "criticized", "inspired_by", "studies_concept", "works_on", "published_in", "cites"]
- ##PUBLICATION:
  - name: "Zur Elektrodynamik bewegter Körper"
  - type: "PUBLICATION"
  - description: "A scientific paper published by Einstein."
  - roles: ["Paper"]
  - relationships: ["influenced_by", "authored_by", "proposed_by", "criticized", "inspired_by", "studies_concept", "works_on", "published_in", "cites"]
- ##PAPER:
  - name: "On the Electrodynamics of Moving Bodies"
  - type: "PAPER"
  - description: "A scientific paper published by Einstein."
  - roles: ["Paper"]
  - relationships: ["influenced_by", "authored_by", "proposed_by", "criticized", "inspired_by", "studies_concept", "works_on", "published_in", "cites"]
- ##EXPERIMENT:
  - name: "Michelson-Morley experiment"
  - type: "EXPERIMENT"
  - description: "A scientific experiment conducted by Michelson and Morley."
  - roles: ["Experiment"]
  - relationships: ["influenced_by", "authored_by", "proposed_by", "criticized", "inspired_by", "studies_concept", "works_on", "published_in", "cites"]
- ##MILESTONE:
  - name: "Nobel Prize"
  - type: "MILESTONE"
  - description: "A significant life or professional milestone."
  - roles: ["Milestone"]
  - relationships: ["influenced_by", "authored_by", "proposed_by", "criticized", "inspired_by", "studies_concept", "works_on", "published_in", "cites"]
- ##DATE:
  - name: "1885"
  - type: "DATE"
  - description: "A specific date or time period."
  - roles: ["Date"]
  - relationships: ["influenced_by", "authored_by", "proposed_by", "criticized", "inspired_by", "studies_concept", "works_on", "published_in", "cites"]
- ##FAMILY_MEMBER:
  - name: "Hans Albert Einstein"
  - type: "FAMILY_MEMBER"
  - description: "A family relationship."
  - roles: ["Family Member"]
  - relationships: ["influenced_by", "authored_by", "proposed_by", "criticized", "inspired_by", "studies_concept", "works_on", "published_in", "cites"]
- ##HONOR:
  - name: "Nobel Prize in Physics"
  - type: "HONOR"
  - description: "A honor and award."
  - roles: ["Honor"]
  - relationships: ["influenced_by", "authored_by", "proposed_by", "criticized", "inspired_by", "studies_concept", "works_on", "published_in", "cites"]
- ##BOOK:
  - name: "The Meaning of Relativity"
  - type: "BOOK"
  - description: "A published book."
  - roles: ["Book"]
  - relationships: ["influenced_by", "authored_by", "proposed_by", "criticized", "inspired_by", "studies_concept", "works_on", "published_in", "cites"]
- ##LETTER:
  - name: "Letter to Mileva Marić, 1893"
  - type: "LETTER"
  - description: "A letter."
  - roles: ["Letter"]
  - relationships: ["influenced_by", "authored_by", "proposed_by", "criticized", "inspired_by", "studies_concept", "works_on", "published_in", "cites"]
- ##PATENT:
  - name: "Electromagnetic Pump Patent"
  - type: "PATENT"
  - description: "A patent."
  - roles: ["Patent"]
  - relationships: ["influenced_by", "authored_by", "proposed_by", "criticized", "inspired_by", "studies_concept", "works_on", "published_in", "cites"]
- ##FIELD:
  - name: "Physics"
  - type: "FIELD"
  - description: "An academic field."
  - roles: ["Field"]
  - relationships: ["influenced_by", "authored_by", "proposed_by", "criticized", "inspired_by", "studies_concept", "works_on", "published_in", "cites"]
- ##DISCOVERY:
  - name: "Brownian motion"
  - type: "DISCOVERY"
  - description: "A discovery."
  - roles: ["Discovery"]
  - relationships: ["influenced_by", "authored_by", "proposed_by", "criticized", "inspired_by", "studies_concept", "works_on", "published_in", "cites"]
- ##STRUGGLE:
  - name: "Financial struggle"
  - type: "STRUGGLE"
  - description: "An abstract theme."
  - roles: ["Struggle"]
  - relationships: ["influenced_by", "authored_by", "proposed_by", "criticized", "inspired_by", "studies_concept", "works_on", "published_in", "cites"]
- ##SENTIMENT:
  - name: "Frustriert"
  - type: "SENTIMENT"
  - description: "An emotional state expressed."
  - roles: ["Sentiment"]
  - relationships: ["influenced_by", "authored_by", "proposed_by", "criticized", "inspired_by", "studies_concept", "works_on", "published_in", "cites"]
```


Some results



<https://einstein-service-465121314425.us-central1.run.app/>

^^^ simple RAG

Vs. GraphRAG [example](#)

Sample response:

Note that for this experimental version, to reduce cost, graph construction was limited to 1000 nodes, 34536 edges, 6.6% of the data. A full-fledged implementation would have 15,237 nodes and 500000+ edges. Even so the results are very promising.

Ask Einstein AI

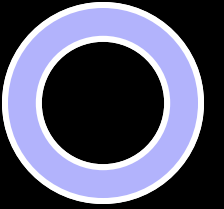
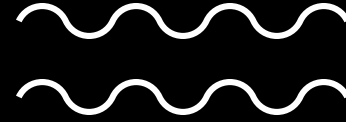
Your Question: How did Einstein develop the theory of relativity?

Sources and References

- Source 1 (merged_ainstein_content, Volume Einstein papers_v01, Page, 011): Referenced 1 time. Discusses the publication of Einstein's letter on "Relativistic Theory of the Non-symmetric Field," and Einstein's interest in publishing it. [View Paper](#)
- Source 2 (merged_ainstein_content, Volume Einstein papers_v01, Page, 281): Referenced 1 time. Discusses Einstein's familiarity of Maxwell's theory to the special theory of relativity and details his disagreement with Hertz's work on electrodynamics. [View Paper](#)
- Source 3 (merged_ainstein_content, Volume Einstein papers_v01, Page, 307): Referenced 1 time. Discusses Einstein's early work in electrodynamics phenomena, his family's involvement in electrodynamics, and his "particular inclination towards physics" from a young age. [View Paper](#)
- Source 4 (merged_ainstein_content, Volume Einstein papers_v010, Page, 048): Referenced 2 times. Discusses Einstein's disagreements with critics of his theory of relativity, and Einstein, who struggled with the concepts of universal time and the constancy of the speed of light. [View Paper](#)
- Source 5 (merged_ainstein_content, Volume Einstein papers_v01, Page, 036): Referenced 1 time. Highlights Einstein's dedication to theoretical physics, his interest in experimental work, and his work at ETH. [View Paper](#)

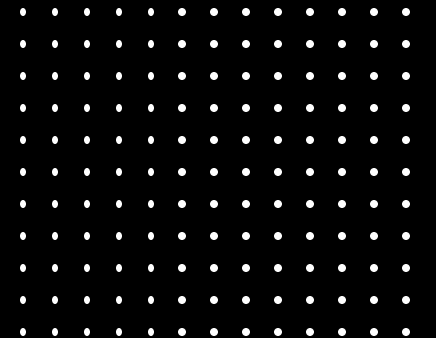
Modeling Intertextuality: An Ontological Framework for Literary Studies

Laura Untner



Intertextuality is messy, but we need structure.

Some ontologies already provide this structure:
... mainly for scholarly publishing (BIBO, CiTO, etc.)
... or for specific corpora or projects (SAWS,
Hypermedia Dante, OntoPoetry, GOLEM,
MiMoText, etc.)

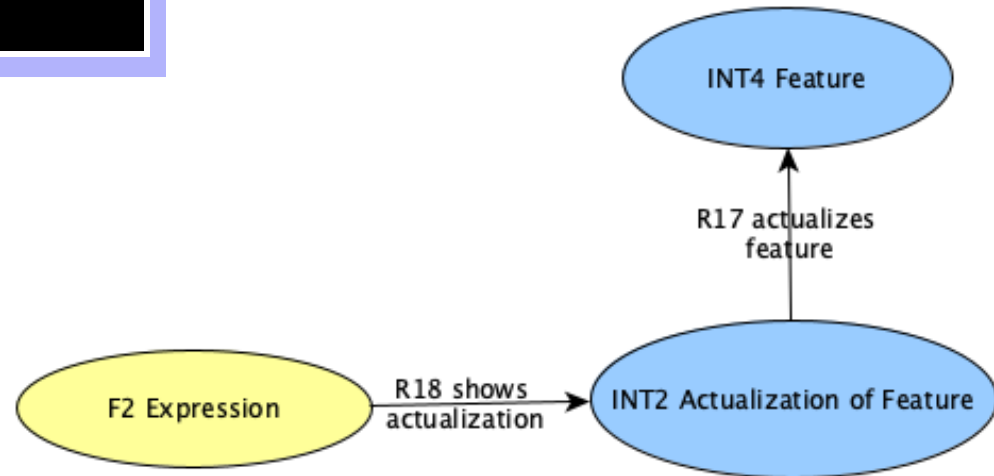
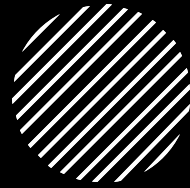


INTRO: The Intertextual, Interpictorial, and Intermedial Relations Ontology

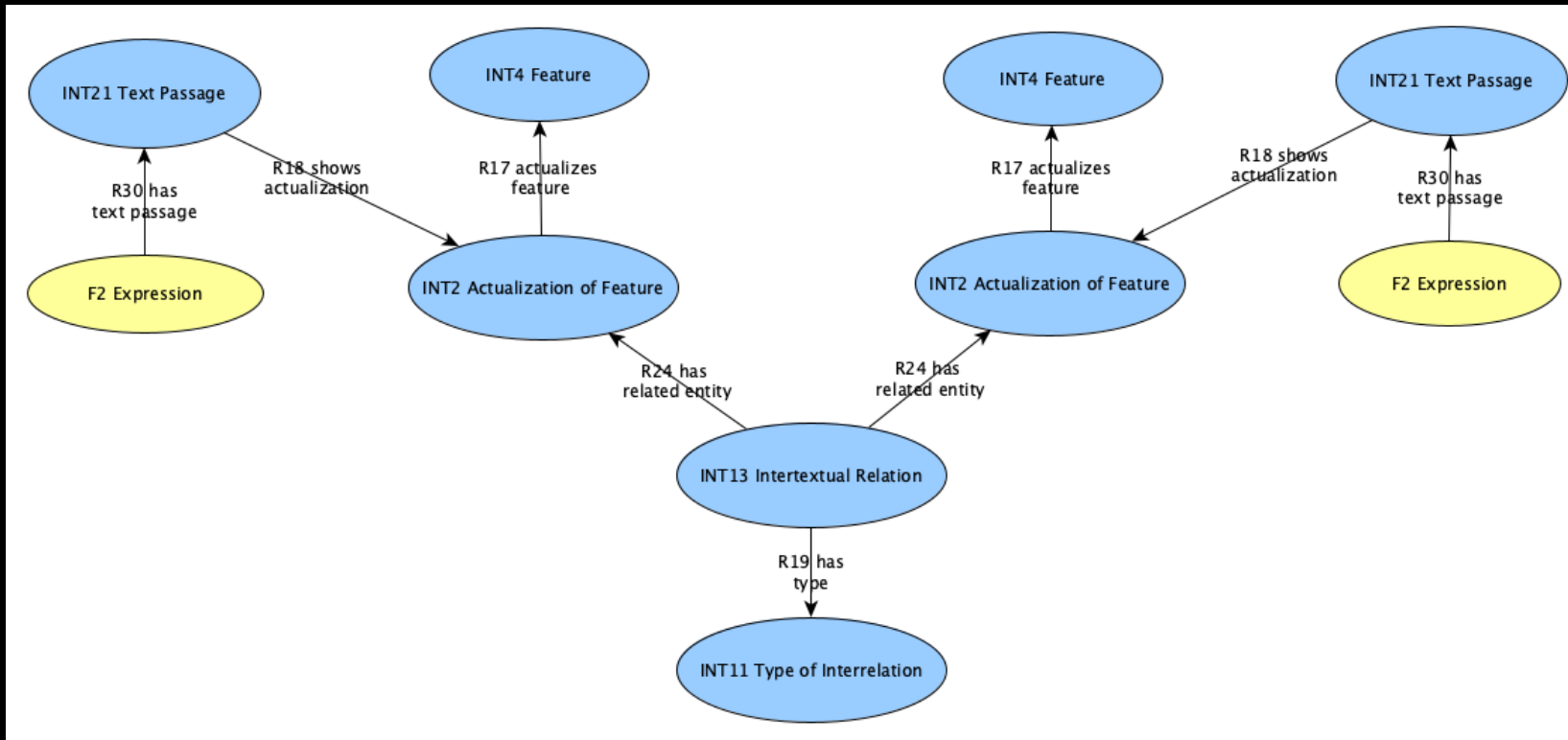
... was originally developed by Bernhard Oberreither.

... is built on CIDOC CRM and LRMoo.

... is text-centric.



>> <https://boberreither.github.io/INTRO/> <<



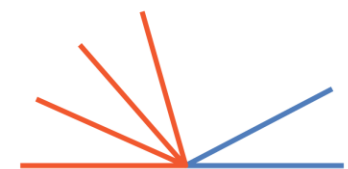
Current projects

Laura Untner (2025): From Wikidata to CIDOC CRM: A Use Scenario for Digital Comparative Literary Studies. In: Journal of Open Humanities Data [accepted].

Sappho Digital: <https://sappho-digital.com>

Automating the Study of Digital Literary Memory: A Multilingual LLM Pipeline for Wikipedia-Based Cultural Analysis

Botond Szemes
University of Tartu
DigiTS Research Group



DigiTS
Center for Digital
Text Scholarship



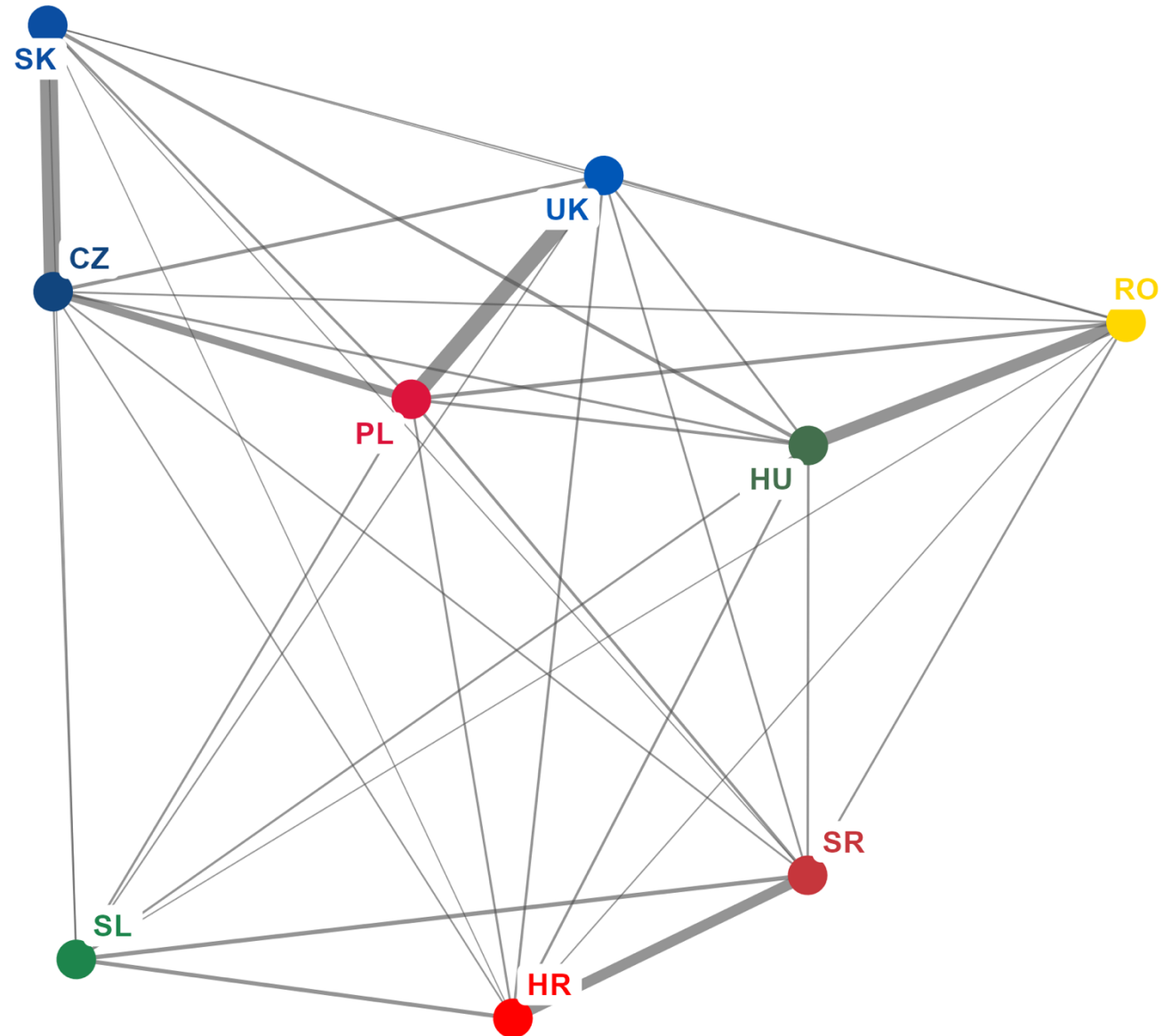
**Funded by
the European Union**

- Clean Wikidata query results
- Who is Polish, Hungarian etc.? Based on the language of publications
- What is literature? Before 1800: everything. After 1800: just fiction in the modern sense + drama, poetry.
- Manual cleaned dataset for the Visegrad region (Czech, Hungarian, Polish, Slovak) as gold standard
- Test an LLM pipeline:
 - send the Wikipedia articles in a given language to a local version of a large model with the same prompt via API
- **llama3.3:70B** overall **F1 = 0.89**
- In the extended dataset cross validation of results from different languages

LLM	Author	Wiki	F1 Score	Precision	Recall
llama3.1:8b	Cz	Hu	0,84	0,74	0,96
	Cz	Pl	0,80	0,72	0,91
	Cz	Sk	0,78	0,70	0,89
	Hu	Cz	0,81	0,73	0,92
llama3.3:70b	Cz	Hu	0,92	0,91	0,93
	Cz	Pl	0,87	0,87	0,87
	Cz	Sk	0,84	0,82	0,85
	Hu	Cz	0,90	0,92	0,89
gemma3:27b	Cz	Hu	0,78	0,71	0,86
	Cz	Pl	0,85	0,81	0,90
	Cz	Sk	0,76	0,72	0,79
	Hu	Cz	0,82	0,91	0,74
gpt-oss:120b	Cz	Hu	0,87	0,88	0,86
	Cz	Pl	0,89	0,91	0,88
	Cz	Sk	0,85	0,86	0,84
	Hu	Cz	0,92	0,98	0,87

Network of Literary Memory

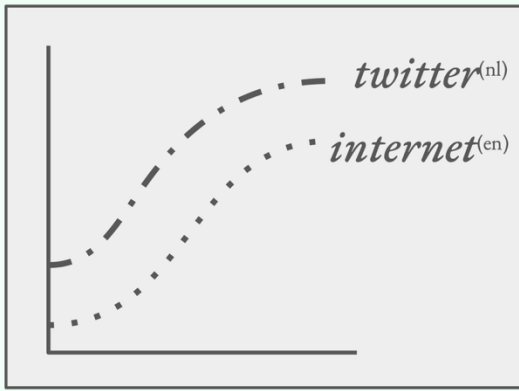
Weighted by summed memory score



Measuring the Synchronicity of Historical European Parliamentary Discourse, 1949-2018

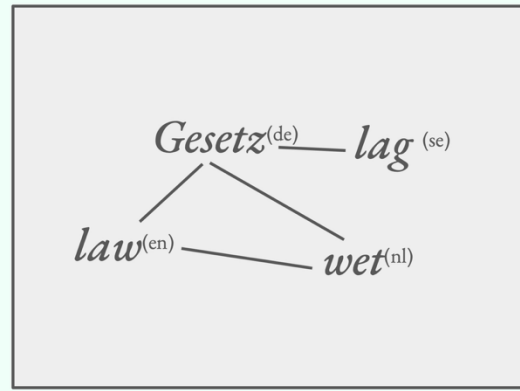
Ruben Ros and Risto Turunen

How did *British*, *Dutch*, *Swedish*, and *German* parliamentary discourse develop between 1945 and 2018?



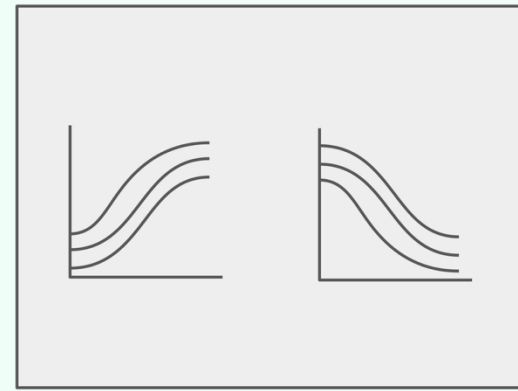
Find bilingual **word pairs** with high correlation in frequency time series.

Spearman correlation coefficient + Dynamic Time Warping



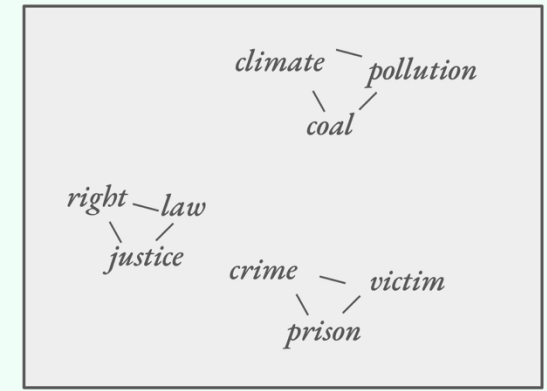
Filter **concepts**: multilingual networks of words with high internal frequency correlation.

Contextual translation + word alignment = distributions of translations



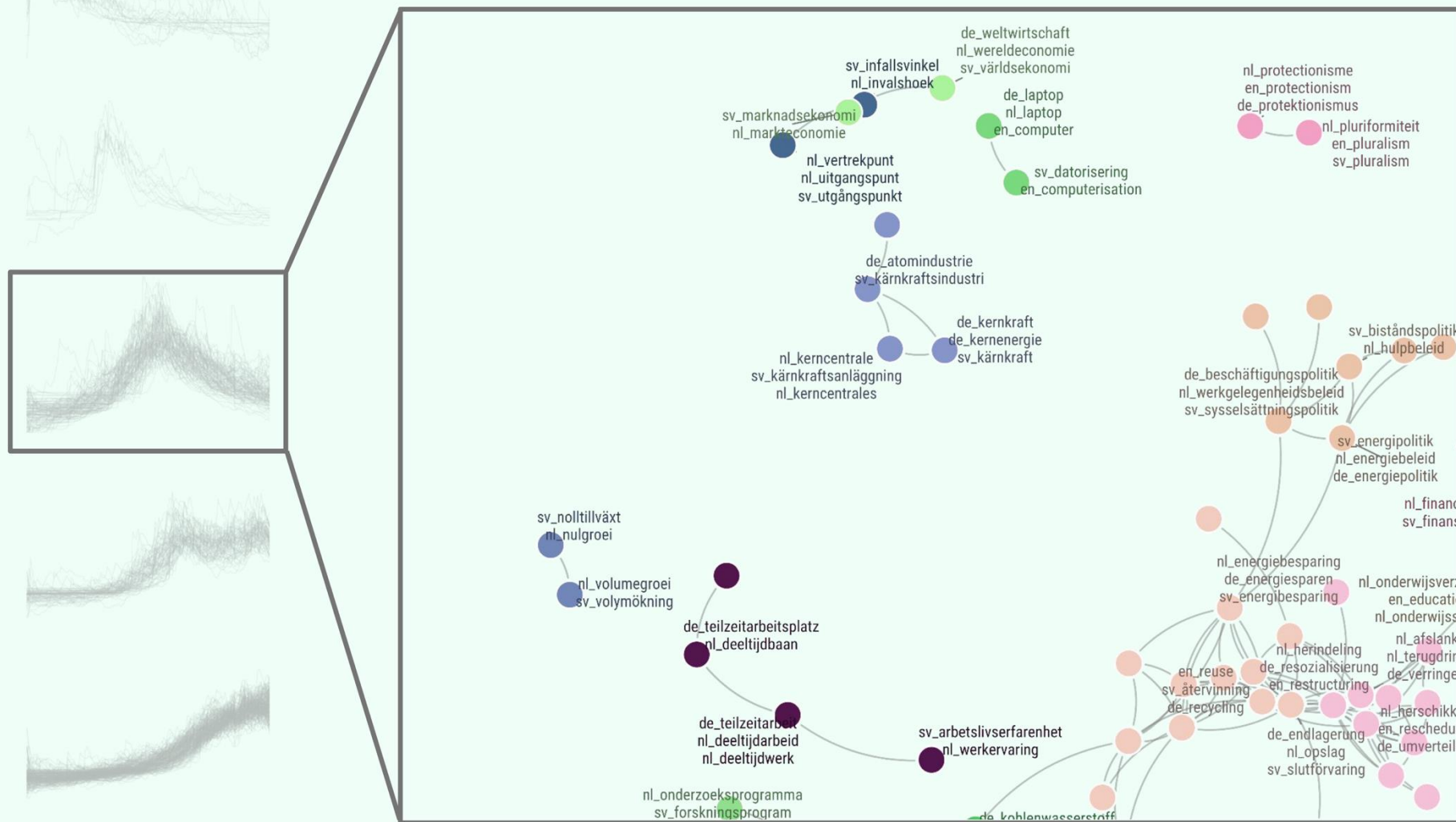
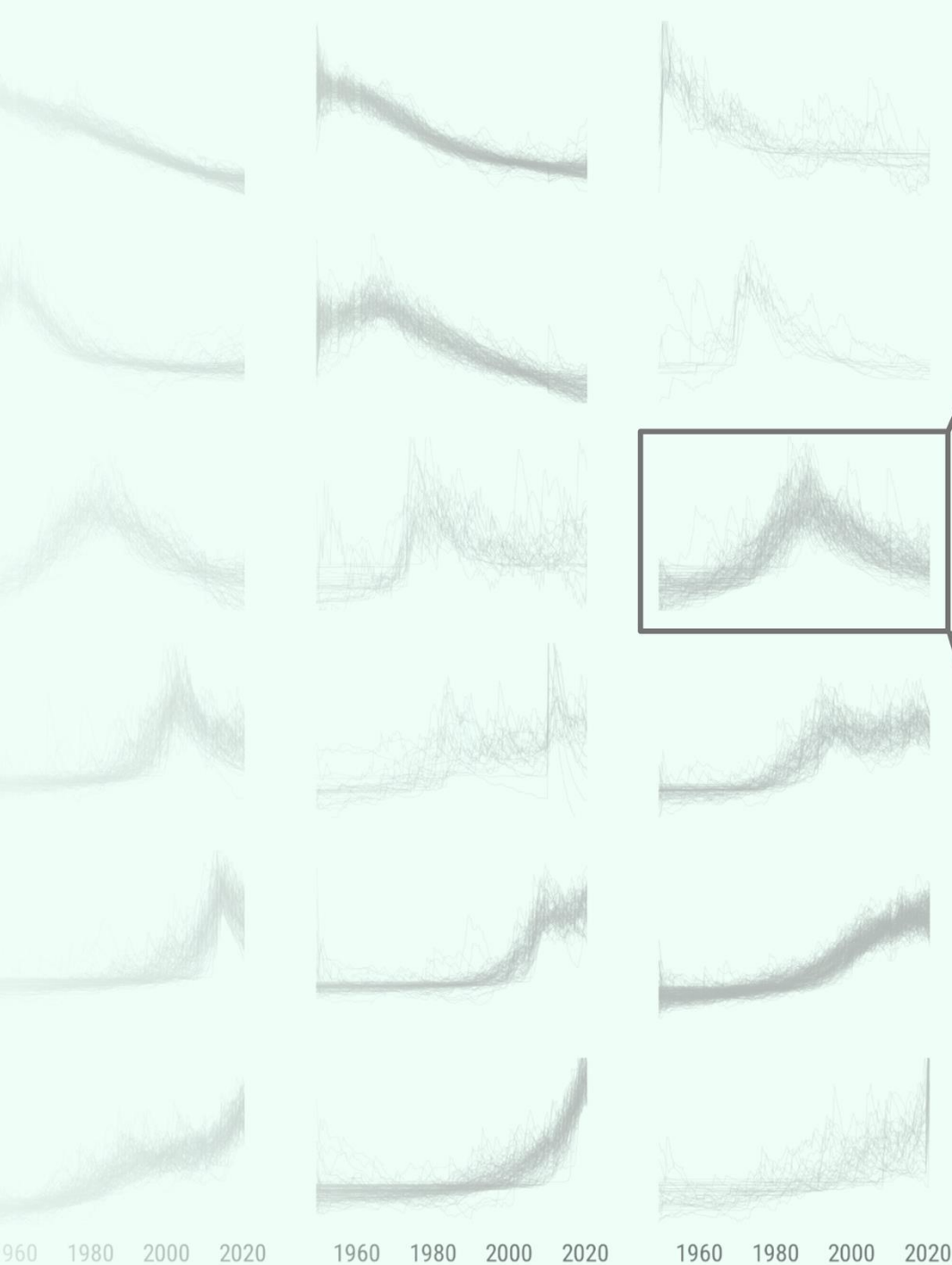
Cluster averaged concept **frequency time-series**

K-means clustering with normalized time series.

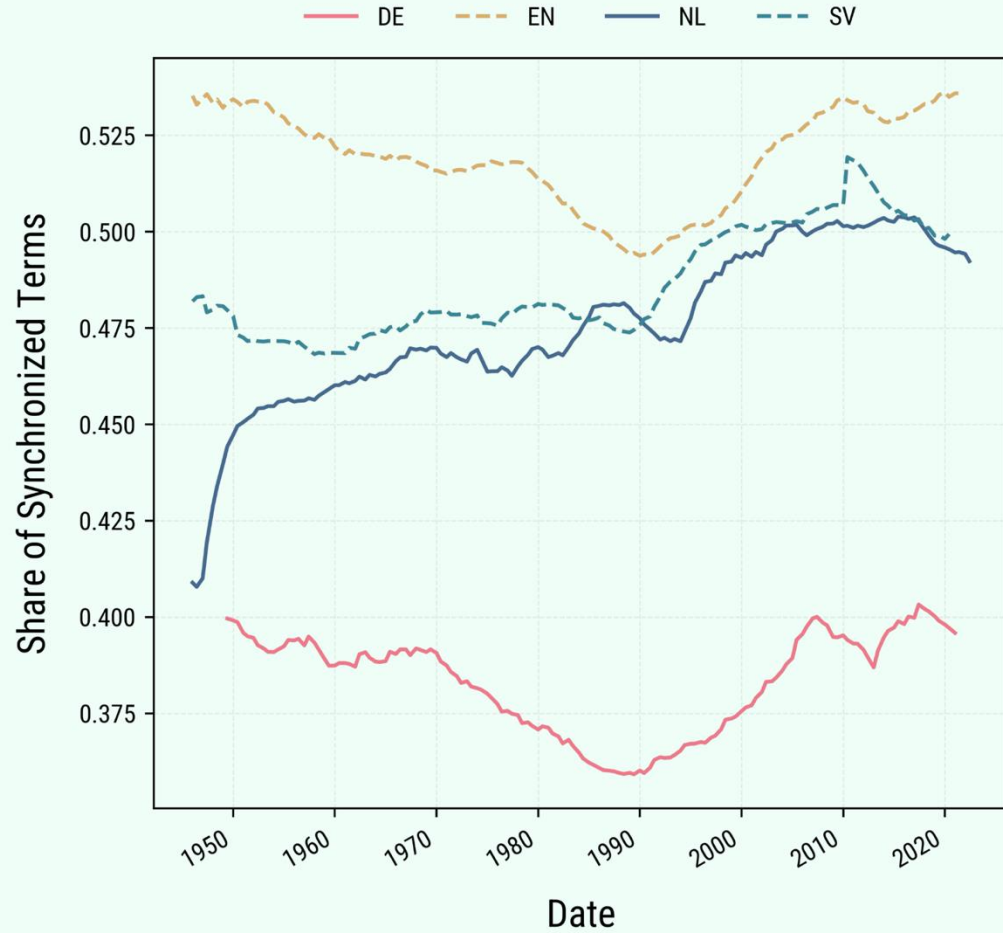


Find **semantic clusters** within time-series clusters

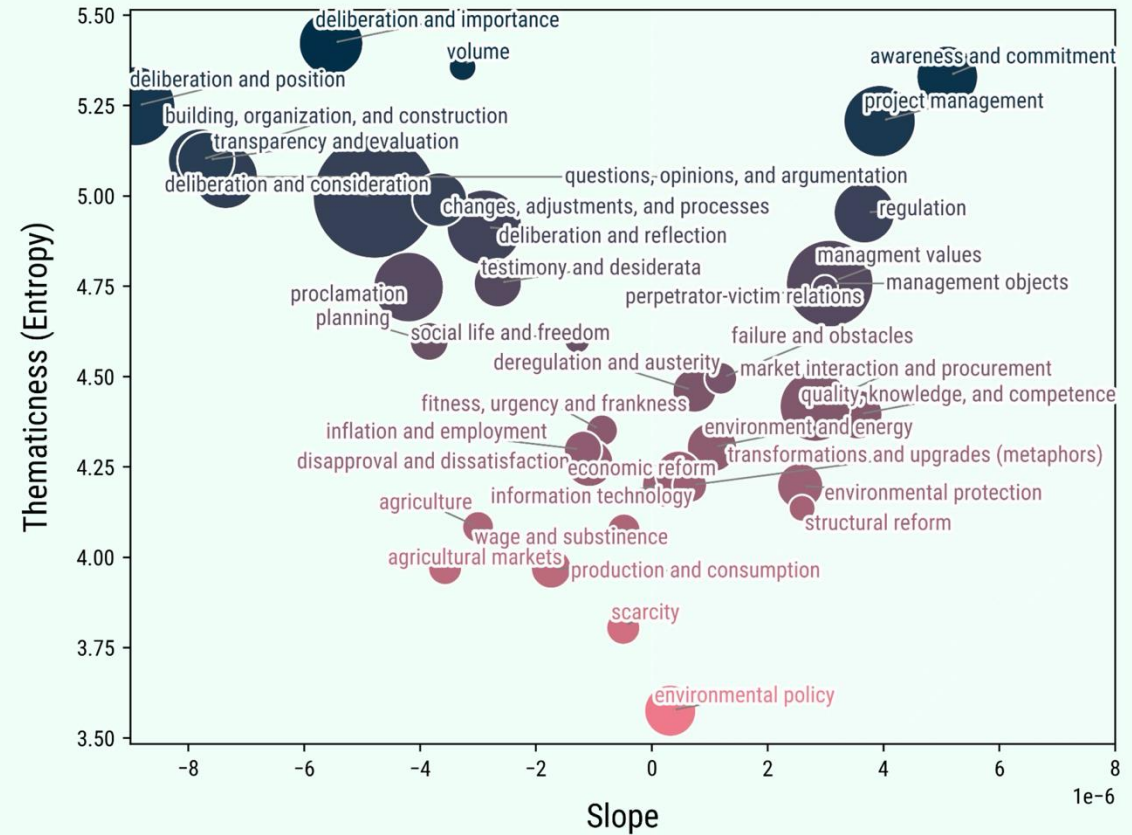
Type embeddings from adapted multilingual bert-base



Identifying semantic clusters within temporal clusters.



The **share of synchronized terms** increases in all countries, especially from the 1990s onwards.



The slope and entropy of clusters reveals shared trends in both **procedural style** and **thematic content**

Disorder or (self-)murder? Making sense of suicide in 19th-century British newspapers

Nilo Pedrazzini and Daniel C. S. Wilson

Disorder OR (self-)murder?

Making sense of suicide in 19th-century British newspapers

Nilo Pedrazzini (The Alan Turing Institute, UK)
Daniel CS Wilson (University College London, UK)

UNPRECEDENTED CASE.—*Mary Murgetts*, an interesting young woman, was indicted for throwing herself into the river, with intent to commit *felo de se*.
The prisoner pleaded guilty.
The COMMON SERGEANT told the prisoner that her offence was one forbidden in the decalogue, and by the common law of England declared to be a felony. He (the Common Sergeant) had been informed that down to the present time she entertained a determination to destroy herself; hence, the best thing the Court could do would be to respite her judgment until the next sessions, and he hoped ere that time the chaplain of the gaol would be able to impress upon her mind the sin and folly of self-destruction.
The session terminated on Thursday night in recognition of

1843

sin and folly

23 YEARS OF SUFFERING
—
CLAYTON WOMAN HANGS HERSELF.
—
A MERCIFUL VIEW.
—
After suffering from rheumatism for the long period of 23 years, Margaret Heys (53), a single woman, who resided with her niece at 19, Blackburn-road, Clayton-le-Moors, put an end to her sufferings on Monday afternoon by hanging herself. Owing to the rheumatism deceased was infirm, and she attended to the house whilst her niece went to the mill. She was all right at noon on Monday, but when a next door neighbour tried the door of the house during the afternoon she found it locked. About five o'clock she suspected that something untoward had happened, and she called deceased's nephew, who, on entering the house, found his aunt hanging by a rope in the front room. The body was immediately cut down and medical aid summoned, but

1907

suffering

Links to newspaper collections:



Living with
Machines (LwM)



Heritage Made
Digital (HMD14)

Disorder OR (self-)murder?

Making sense of suicide in 19th-century British newspapers

Nilo Pedrazzini (The Alan Turing Institute, UK)
Daniel CS Wilson (University College London, UK)

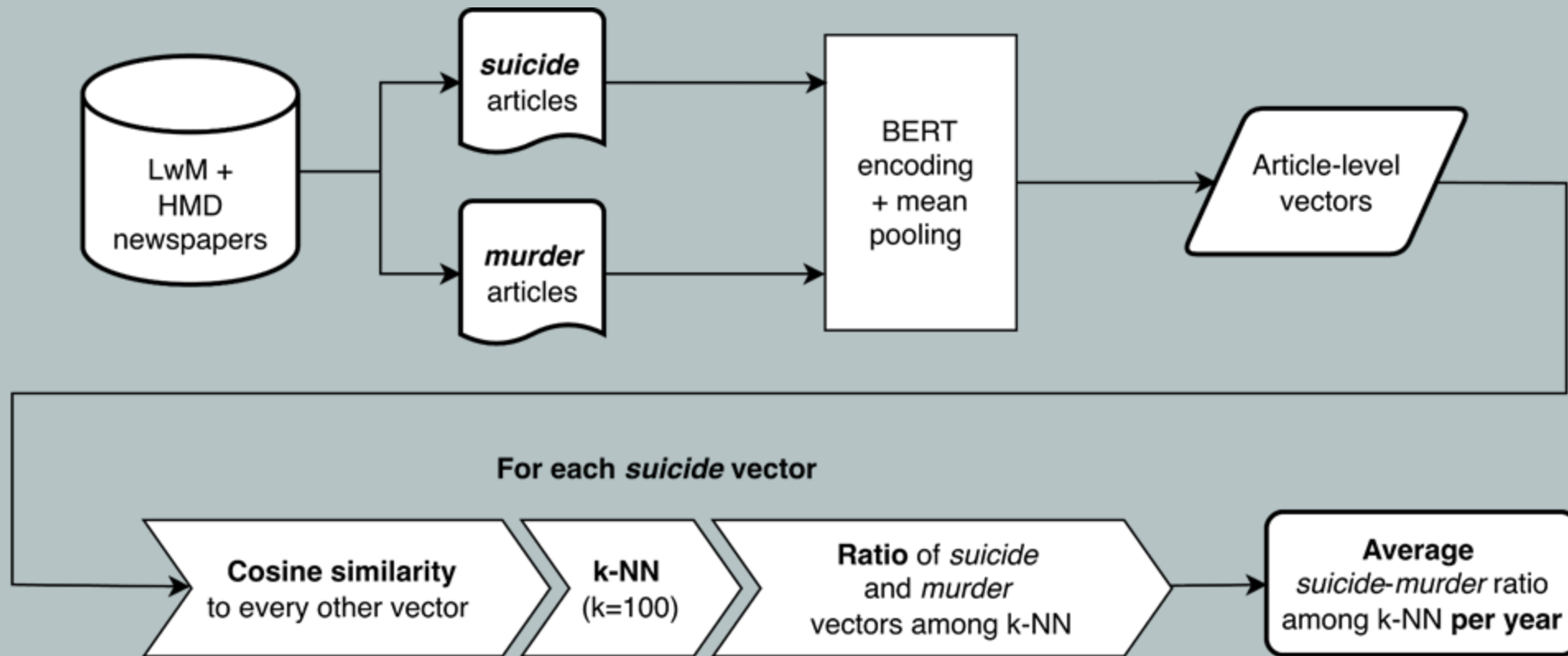


Figure 2. Flowchart of the method. Newspaper articles on *suicide* and *murder* are encoded with BERT into article-level vectors. For each suicide vector we compute cosine k-NN over all articles, derive the suicide-murder neighbour ratio, and average this ratio by year.

Disorder OR (self-)murder?

Making sense of suicide in 19th-century British newspapers

Nilo Pedrazzini (The Alan Turing Institute, UK)
Daniel CS Wilson (University College London, UK)

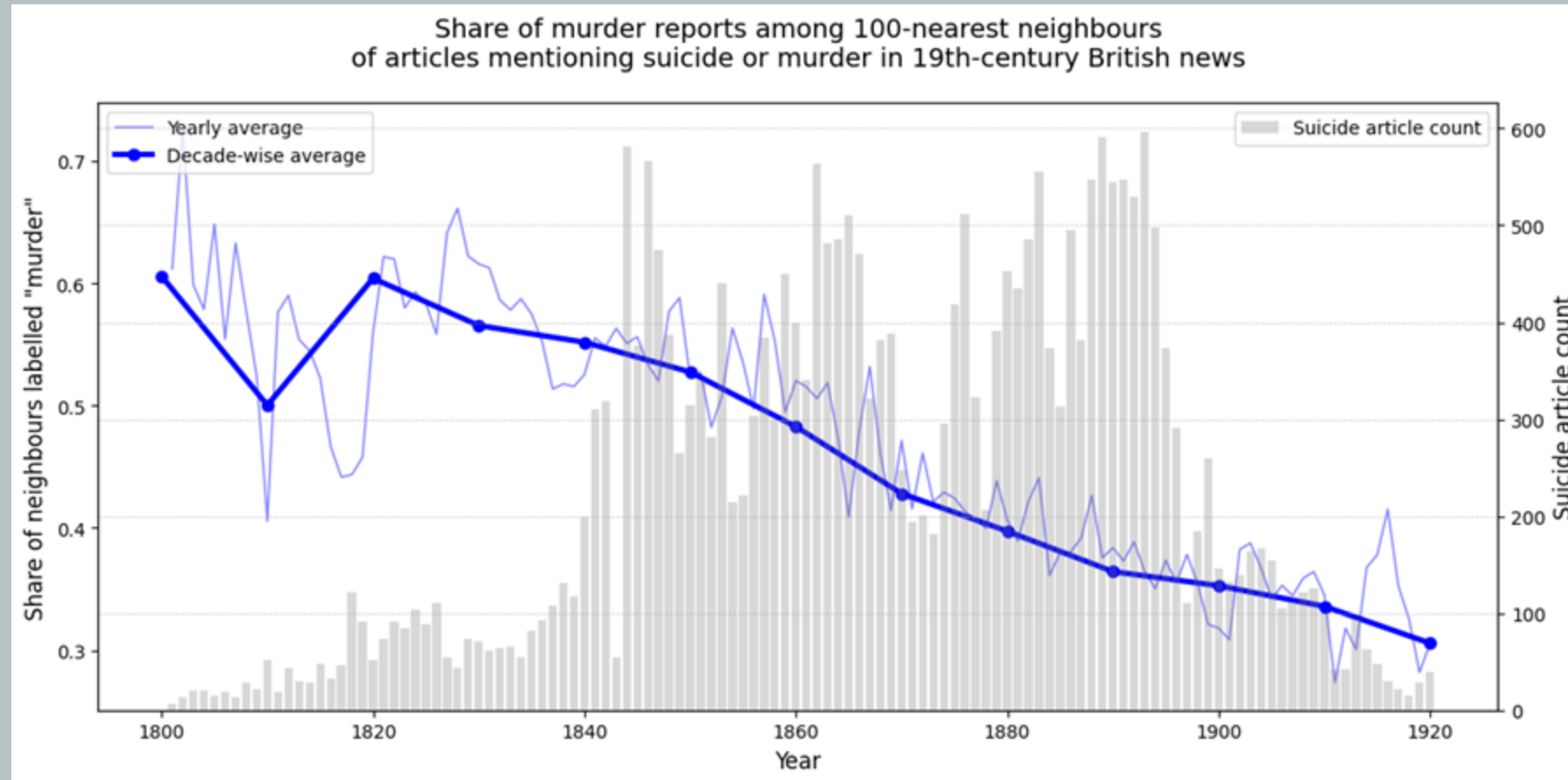


Figure 3. Average ratio of neighbours from *murder* articles among the 100 nearest neighbours of suicide articles (yearly and decade-wise averages, 1800–1920). Grey bars show the yearly count of suicide articles.

Mapping Literary Networks through Epigraphs

Tomás Espino Barrera

L E T T R E S

DE DEUX AMANS,

Habitans d'une petite Ville
au pied des Alpes.

RECUEILLIES ET PUBLIÉES

PAR J. J. ROUSSEAU.

PREMIERE PARTIE.



A AMSTERDAM,

Chez MARC MICHEL REY.

MDCCLXI.

EPIMAPS - Mapping Epigraphical Networks

1. ELTeC corpus, aprox. 2000 novels
2. Epigraph extraction and enrichment (incl. sources)
 - 1st open-access multilingual epigraph database
3. Interactive mapping (Nodegoat)
4. Distant Reading
5. Case studies

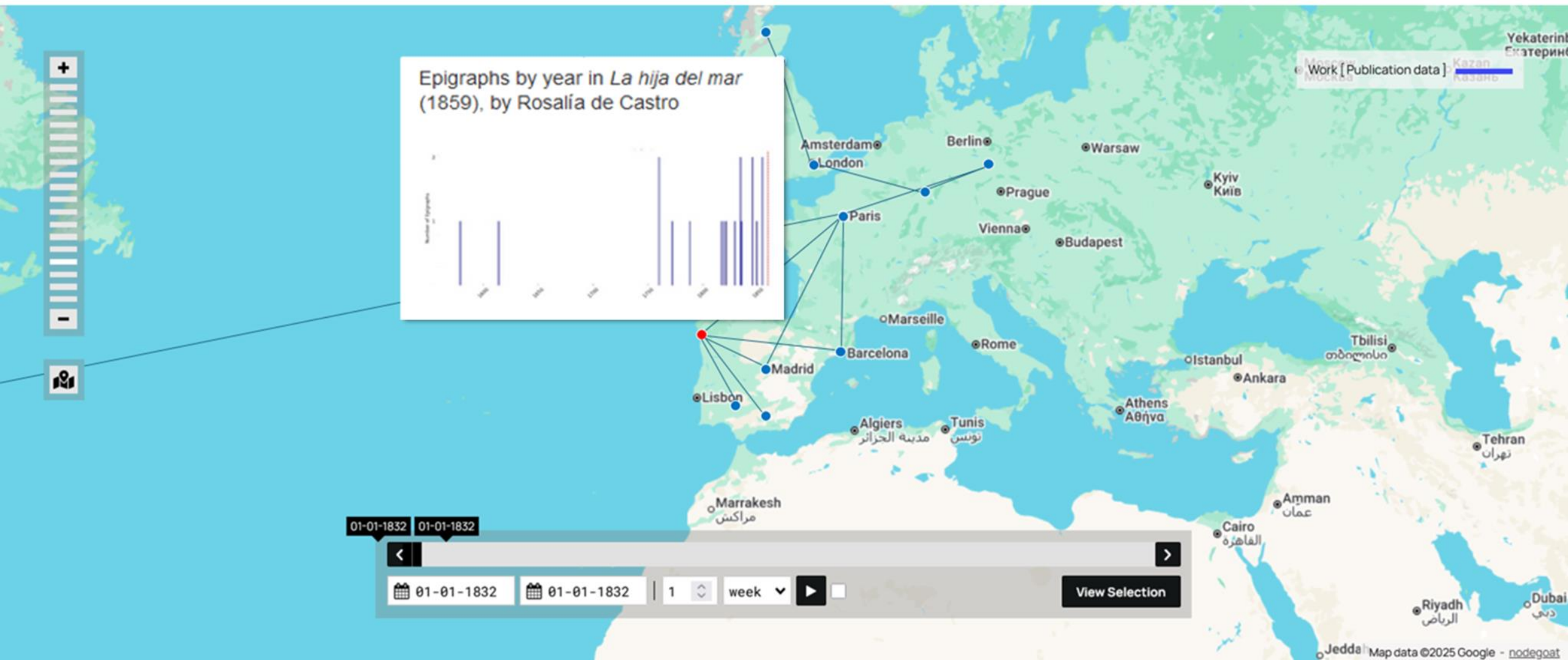
EPIMAPS

subgenre	narrative-j	size-categ	reprint-co	time-slot	Epigraphs:	Feuilleton	title	chapters	epigraph	e-author-r	e-author-r	e-author-id	e-title	e-title-id	e-publicati	e-original	e-translati	e-publicat
NA	NA	long	unspecified	1750	Y	N	Y	N	"O quantur Persius	"Pers."	wikidata:Q332785	<i>Saturae</i> 1. wikidata:Q87143921	Rome?					1st centur
NA	NA	short	unspecified	1750	Y	N	Y	N	"Haec legit Voltaire	"Ovid."	wikidata:Q9068	<i>Méropé</i> wikidata:Q3333324	Paris	"Hoc legite austeri, cri				1744
NA	NA	short	unspecified	1750	Y	N	Y	N	"Non semꝑ Phaedrus	"Phaedr."	wikidata:Q52166	<i>Fabulae Aë</i> wikidata:	Rome?					1st centur
NA	NA	long	unspecified	1760	Y	N	Y	N	"successer Tacitus	"Tacit."	wikidata:Q2161	<i>Histoires</i> 2 wikidata:Q1247073	Rome?					2nd centur
NA	NA	short	unspecified	1760	Y	N	Y	N	"Nulla Viro Catullus	"Cat."	wikidata:Q163079	LXIV "Argo wikidata:						1st centur
epistolary	NA	long	unspecified	1760	Y	N	Y	N	"Non la co Petrarch	"Petrarc"	wikidata:Q1401	<i>Canzonere</i> wikidata:Q777574	Avignon / Italy					14th C.
epistolary	NA	long	unspecified	1761	Y	N	Y	Y (2.12)	"O qual fia Metastasic	N/A	wikidata:Q29473	<i>Attilio Regi</i> wikidata:	Several locations (opera)					1750
NA	NA	short	unspecified	1760	Y	N	Y	N	"Des jeune Gresset, Je	"Gresset"	wikidata:Q942163	<i>Le Méchan</i> wikidata:Q3224937	Paris	selected fragments of				1747
didactic no	NA	long	unspecified	1760	Y	N	Y	N	"Sanabilibꝑ Seneca	"Sen"	wikidata:Q2054	"de ira. L. I wikidata:Q3704115	Rome?					1st centur
NA	NA	short	unspecified	1760	Y	N	Y	N	"J'admire l. N/A	N/A	N/A	N/A wikidata:	N/A					N/A
libertine no	NA	long	unspecified	1760	Y	N	Y	N	"Priape, so Piron, Alex	"Piron"	wikidata:Q983437	<i>Ode à Priap</i> wikidata:	N/A					ca. 1710
NA	NA	medium	unspecified	1760	Y	N	Y	N	"..... Ut nec Horace	"Hor."	wikidata:Q6197	<i>Ars poetica</i> wikidata:Q677997	Rome?					1st centur
NA	NA	long	unspecified	1760	Y	N	Y	N	"Tout ce qꝑ Self-epigra	N/A	N/A	"Tome II. F N/A	London?	"tout ce qui est hors d				1766
NA	NA	short	unspecified	1760	Y	N	Y	N	"Ecce spec Seneca	Senec.	wikidata:Q2054	<i>De provide</i> wikidata:Q2264770	Rome?					1st centur
NA	NA	short	unspecified	1760	Y	N	Y	N	"Non miroꝑ Seneca	Senec.	wikidata:Q2054	<i>De provide</i> wikidata:Q2264770	Rome?	"Ego uero non miror, s				1st centur
NA	NA	short	unspecified	1760	Y	N	Y	N	"La Vertu c N/A	N/A	N/A	N/A N/A	N/A					N/A
NA	NA	short	unspecified	1760	Y	N	Y	N	"Virtue car Richardsor	"Row."	wikidata:Q295941	<i>Clarissa</i> wikidata:Q980534	London					1748
NA	NA	short	unspecified	1760	Y	N	Y	N	"Si vuole à Guarini, Gi	"Guarini"	wikidata:Q542039	<i>Il pastor fic</i> wikidata:Q477990	Venice	"Si vuole à punto / Far				1590
NA	NA	short	unspecified	1760	Y	N	Y	N	"On peut tꝑ Boileau, Ni	"Boil."	wikidata:Q188857	Satire X [Sæ wikidata:	Paris					1694
NA	NA	medium	unspecified	1770	Y	N	Y	N	"Le Temps Leibniz, Go	"Leibnitz"	wikidata:Q9047	Different p N/A	Different p	"Le présent est gros d				ca. 1710-1
NA	NA	long	unspecified	1770	Y	N	Y	N	"One Almiꝑ Milton, Jo	"Milton"	wikidata:Q79759	<i>Paradise L</i> wikidata:Q28754	London		"Il est un s			1667
NA	NA	short	unspecified	1770	Y	N	Y	N	"Nunc scio Virgil	"Virg."	wikidata:Q1398	<i>Eclogues</i> wikidata:Q546203	Rome?					1st centur
libertine no	NA	medium	unspecified	1770	Y	N	Y	N	"La faute e N/A	N/A	N/A	N/A N/A	N/A					N/A
epistolary	NA	long	unspecified	1770	Y	N	N	Y (2.34)	"O serpent N/A	N/A	N/A	N/A N/A	N/A					N/A

EPIMAPS

You know my method. It is founded upon the observation of trifles.

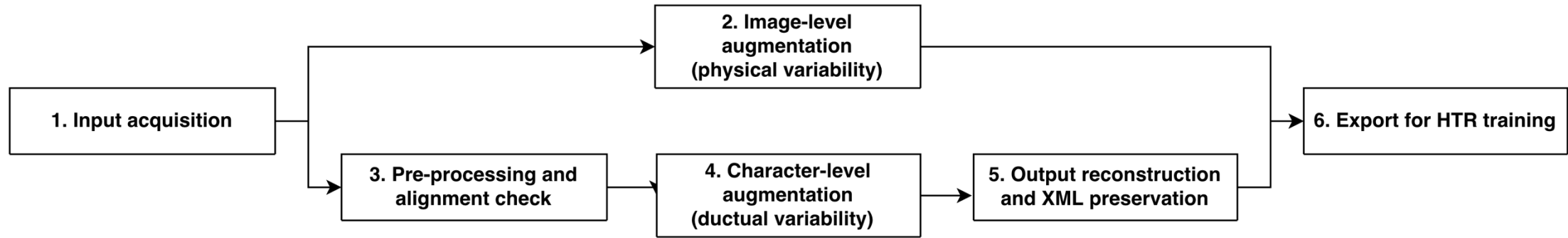
—Arthur Conan Doyle



Low-Cost Synthetic Data Generation for HTR Training: Evaluating a Multimodal Strategy for Historical Manuscript Processing

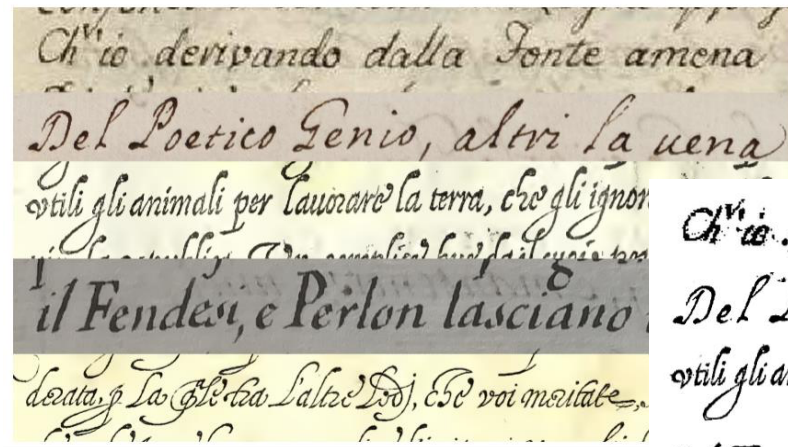
Serena Carlamaria Crespi and Carlos-Emiliano González-Gallardo

A Modular Approach for Data Augmentation



Training corpus

- 446 manuscript pages (images & ALTO transcriptions)
- Mostly poetic Italian codices (17th and 18th centuries)

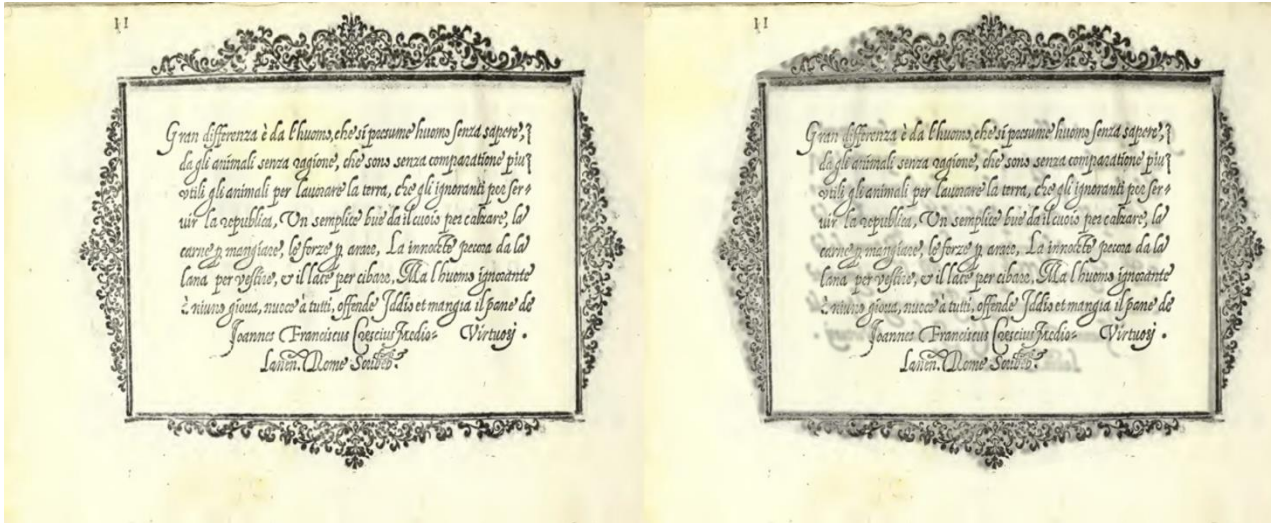


Ch'io derivando dalla Fonte amena
Del Poetico Genio, altri la uena
utili gli animali per l'auorare la terra, che gli ignoranti per ser
il Fender, e Perlon lasciano i posti
deata, p. La Gloria. L'altra Doj, Ede voi meritate, si può dire

Corpus & Augmentations

I. Image-level augmented corpus

- 1,338 manuscript pages (images & ALTO transcriptions)



Ink-bleeding simulation. Cresci, Essemplare, Presso Altobello Salicato, alla Libreria della Fortezza, Venice, fol. 2r. **Left:** original line before augmentation ; **right:** result after applying the ink-bleeding simulation.

II. Character-level augmented corpus

- 1,074 manuscript pages (images & ALTO transcriptions)



Example of Bézier-based ductual augmentation on the same Cresci Essemplare

Beyond the Statistics: Migration to a Kyiv Suburb through the Lens of the 1897 Census

Konstantin Mogarichev, Tetiana Shyshkina and Maria Volkova

ПЕРВАЯ ВСЕОБЩАЯ ПЕРЕПИСЬ НАСЕЛЕНИЯ
РОССИЙСКОЙ ИМПЕРИИ, 1897 Г.

ИЗДАНИЕ ЦЕНТРАЛЬНОГО СТАТИСТИЧЕСКОГО КОМИТЕТА МИНИСТЕРСТВА ВНУТРЕННИХЪ ДѢЛЪ

ПОДЪ РЕДАКЦІЕЮ Н. А. ТРОЙНИЦКАГО.

XVI. КІЕВСКАЯ ГУБЕРНІЯ.



1904.

- Вы чѣмъ занимаетесь, сударыня?
- Какой вы не деликатный, однако.

— *What is your occupation, madam?*

— *How indiscreet you are, really.*

Illustration published in Budilnik, No. 5, 1897

Rethinking Census 1897: Demiivka Individual Census Sheets

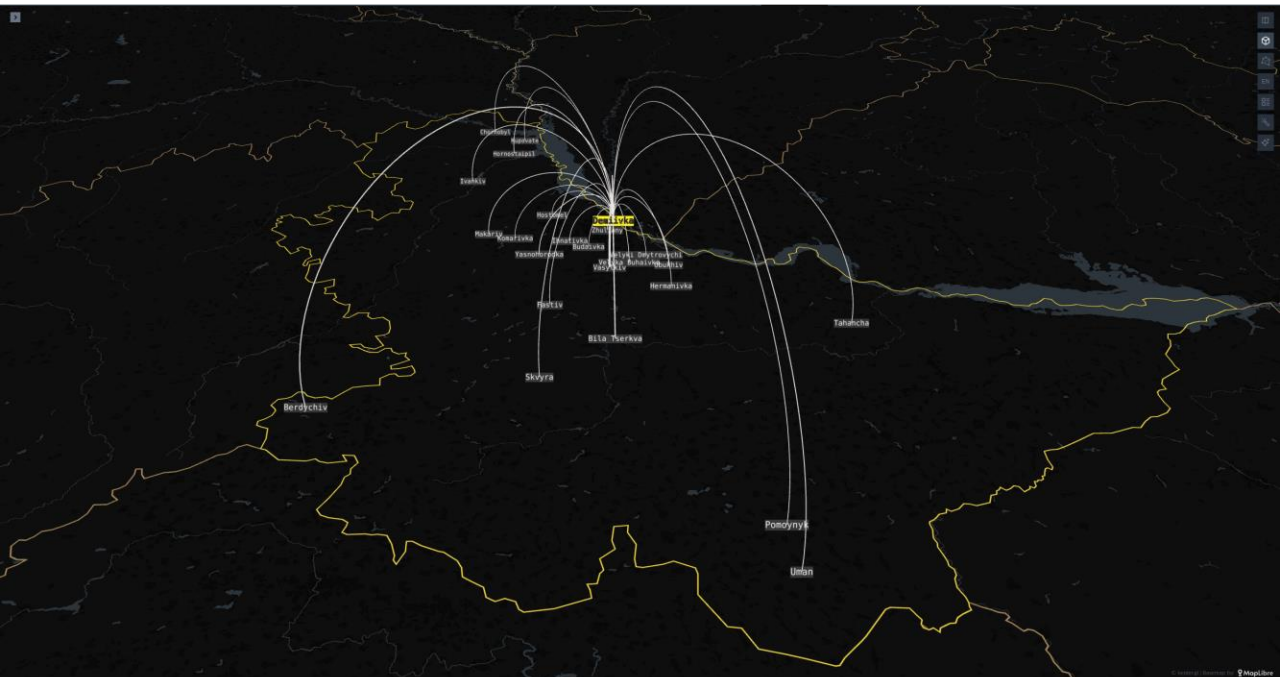
[illegible]

[illegible]

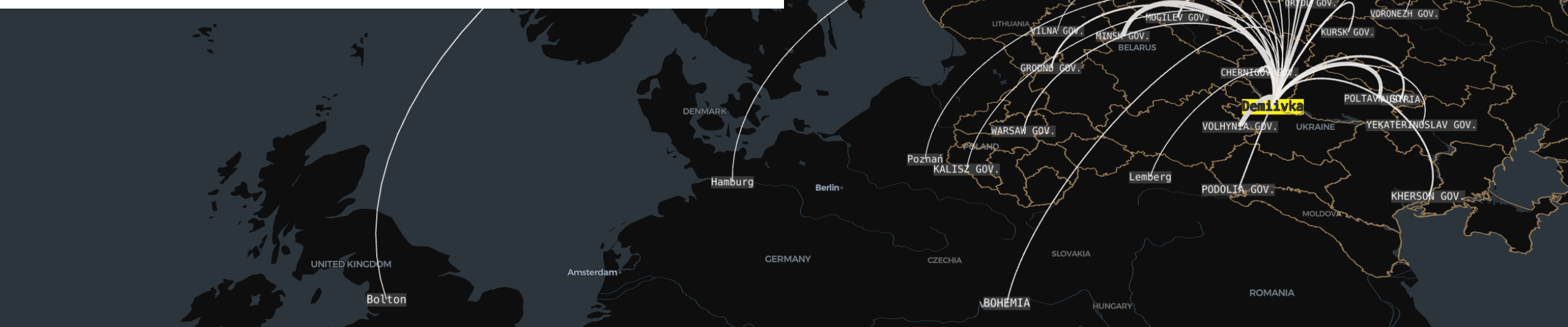
Artemii Plekhanov

[illegible]

Migration to Demiivka from Beyond Kyiv Gubernia (male + female)



Migration to Demiivka within Kyiv Gubernia (male+female)



Speculative Reconstruction and the Ethics of the Fragment: Early Experiments with Generative AI in Art History

Katarina Mohar

Fragments, AI, and Speculative Reconstruction

- Generative AI used to explore **visual hypotheses**, not restorations
- Two micro-datasets: **Selo fresco fragments** and **four paintings by Almenak**
- Early Selo tests exposed the **base model's generic "medievalness"**
- Reveals a core issue: AI fills **absence** with the *probable*, not the *historically specific*

original, Selo



DreamBooth inpainting tests



What AI Learns — and What It Doesn't

- Fine-tuning captured **surface style** from minimal data
- But models struggled with **composition** and **narrative coherence**
- Outputs work as **multiple speculative possibilities**, not reconstructions
- Key insight: AI exposes the **range of interpretations** a fragment can support

Almenak: *The Peddler*



(late 17th century, oil/canvas,
National Gallery of Slovenia)

LoRa tests



Prompt: An elderly peddler showing trinkets to a peasant woman outside a cottage, depicted in the style of a 17th-century oil painting, with expressive gestures, ochre tones, and soft painterly edge in style of almnk

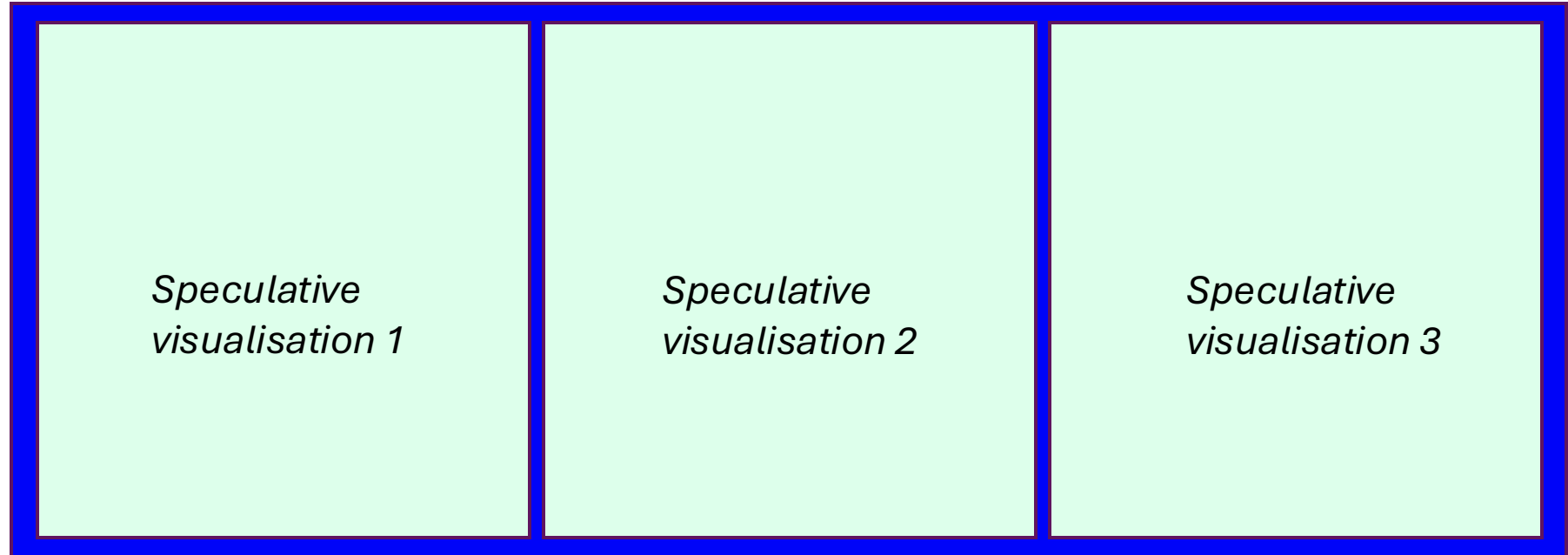
The Ethics of Completing Fragments

- Fragments carry meaning through **incompleteness**
- AI “completions” risk producing **false wholeness** or stylistic flattening
- Our guidelines for responsible use:
 - Label AI outputs as **speculative visualizations**
 - Always show them **alongside the fragment**
 - Present **multiple options**, never a single authoritative fill
- Aim: use AI to **clarify absence**, not erase it

original, Selo



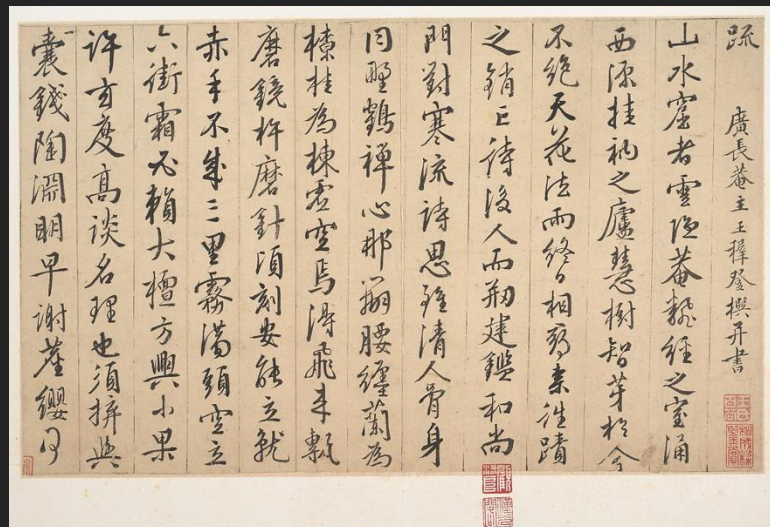
Hypothetical AI-generated fill



When Larger LLMs Aren't Enough: Word Segmentation in Historical Chinese Texts

Hao Tan

- “Word Segmentation” - What is the Problem Here?



“removespacesbetween
sentences this is a sentence
this is another sentence”

“東京大 学 院 生”

“東京大学 + 院生”

A graduate student at
the University of Tokyo

“東京 + 大学院生”

A graduate student
in Tokyo

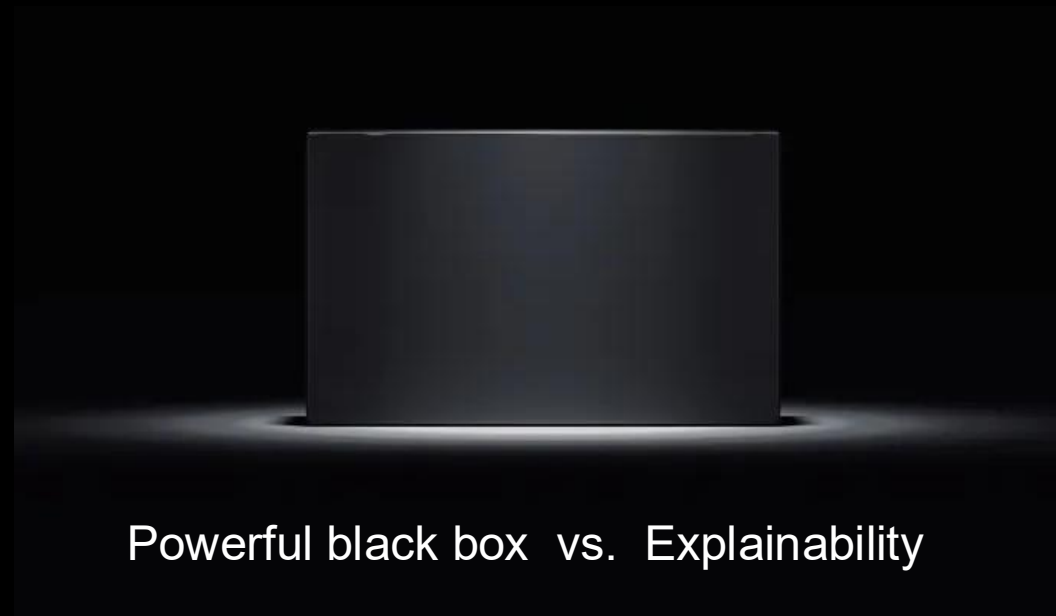
• Why LLMs Aren't Enough for This Problem?

-- Language evolves: "One Size Doesn't Fit All Eras"

	'yellow'	'river'	'enter'	'sea'	'mouth'
	黄	河	入	海	口
Modern Chinese	the Yellow River		estuary		
	黄	河	入	海	口
Medieval Chinese	the Yellow River		enter	estuary	
	黄	河	入	海	口
Ancient Chinese	yellow	river	enter	sea	mouth
	黄	河	入	海	口

- **LLMs in Humanities:**
“What are the Hidden Assumptions Here?”

Digital humanities \longleftrightarrow Interpretive scholarship



Towards animal-centric affective analysis in poetry

Thomas Haider

Framework for Affective Analysis of Animals

Animals in NH German Poetry:
over 20k mentions (in 65k poems)

Taxon	Frequency	Taxon	Frequency	Taxon	Frequency
Nachtigall	1036	Hund	822	Adler	776
Pferd	653	Esel	521	Vieh	493
Wurm	462	Lamm	406	Hahn	378
Fisch	339	Schlangen	338	Bienen	333
Raben	315	Schlange	314	Tauben	308
Fuchs	299	Löwen	293	Hirsch	291
Brut	286	Wolf	265	Löwe	264
Taube	263	Reh	251	Schwan	250
Drachen	249	Schmetterling	236	Schimmel	233
Grillen	224	Aar	211	Schafe	208
Ungeheuer	198	Schwalbe	195	Kukuk	186
Kuh	184	Schaf	175	Bär	173

Table 1: Most frequent animal taxa in German poetry corpus.

Dimensions of Interest:

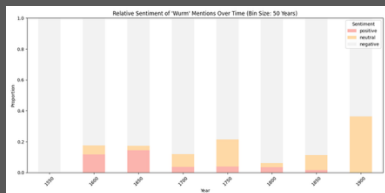
- Diegetic Function
- Agency
- Power
- Target Sentiment
- Representation/Connotation

Target Sentiment Classification (of Animals)

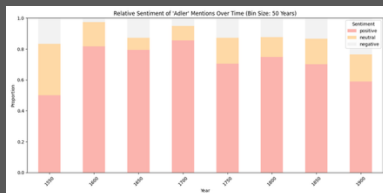
- LLMs in zero shot -> .7 F1 macro
- Few shot hurt
- Animal name as generic 'Tier' -> minus 7 points (sentiment encoded in symbol only)

Metric / Class	Qwen3	Mistral7B	DeepSeekV3	LLaMA3.3
Macro F1	0.659	0.337	0.692	0.696
Negative F1	0.797	0.226	0.762	0.778
Neutral F1	0.415	0.444	0.573	0.573
Positive F1	0.767	0.340	0.742	0.736
Negative Precision	0.724	0.875	0.752	0.690
Neutral Precision	0.595	0.292	0.526	0.623
Positive Precision	0.737	0.735	0.814	0.812
Negative Recall	0.885	0.130	0.773	0.891
Neutral Recall	0.319	0.924	0.630	0.531
Positive Recall	0.798	0.221	0.681	0.672

Table 2: Comparison of Model Performance on Target Sentiment Task



Worm



Eagle

- Why is an animal evaluated a certain way?
- Prompt DeepSeekV3 to interpret 'meaning' of animal in context -> connotation labels (representation)
- Calc. association measures sentiment & connotation
- Counterintuitive connotation -> non-canonical?

Top Keywords for Sentiment: POSITIVE (Adler)

	keyword	observed	expected	pmi	z_score	fisher_p
2	strength	351	302.76	0.2133	2.8888	0.0
3	transcendence	160	130.4	0.2952	2.6376	0.0
5	vision	40	29.98	0.4162	1.838	9e-06
15	victory	37	28.48	0.3777	1.603	0.000229
23	divine	60	49.46	0.2787	1.5083	0.000891

Top Keywords for Sentiment: NEGATIVE (Adler)

	keyword	observed	expected	pmi	z_score	fisher_p
0	vulnerability	21	3.18	2.7236	10.0219	0.0
1	destruction	10	1.48	2.7601	7.0247	0.0
4	retreat	7	1.02	2.7761	5.9188	7e-06
6	suffering	5	0.57	3.1386	5.8852	1.9e-05
7	predator	7	1.14	2.6241	5.509	2.1e-05

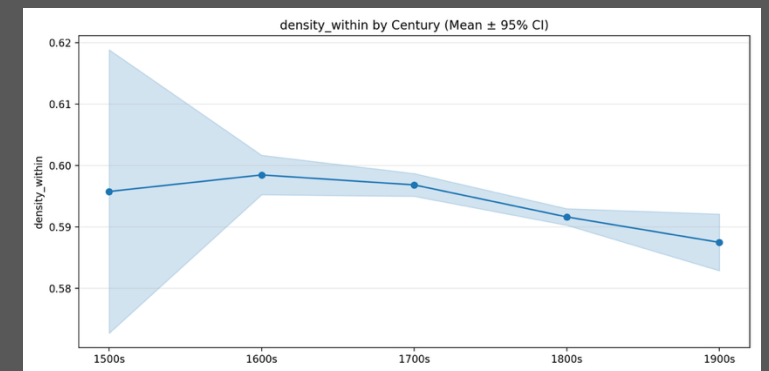
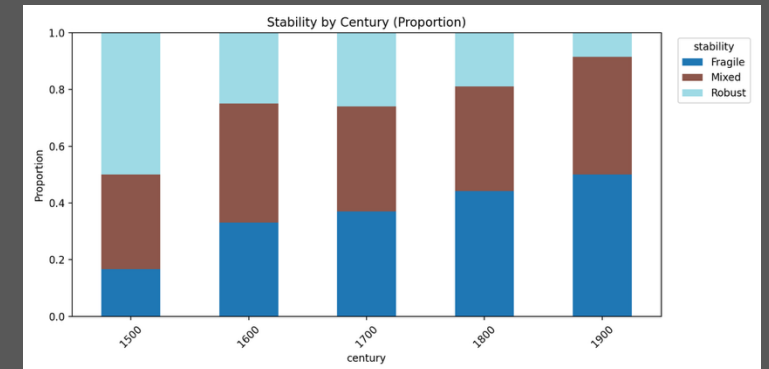
How do we evaluate connotation/representation?

- Which meaning comes from from the symbol (animal name), the context, what's encoded in the model?
- What does a reader need to know about the symbol (to understand the poem)?
- What's the signal in the context?
- Create synthetic ground truth?
- Context stays constant -> Inject various Animals
- -> DeepSeek 'interprets'

- Here stood once a **cow** and thought
- Here stood once a **dragon** and thought
 - **contemplation**; simplicity; existence; reflection; ordinary
 - mythical; transience; memory; decay; **contemplation**

Jaccard coef.
Keyword overlap
(thresholded) ->
context matters
less and less

Topic diversity
increases (more
contexts)
(density ↓)

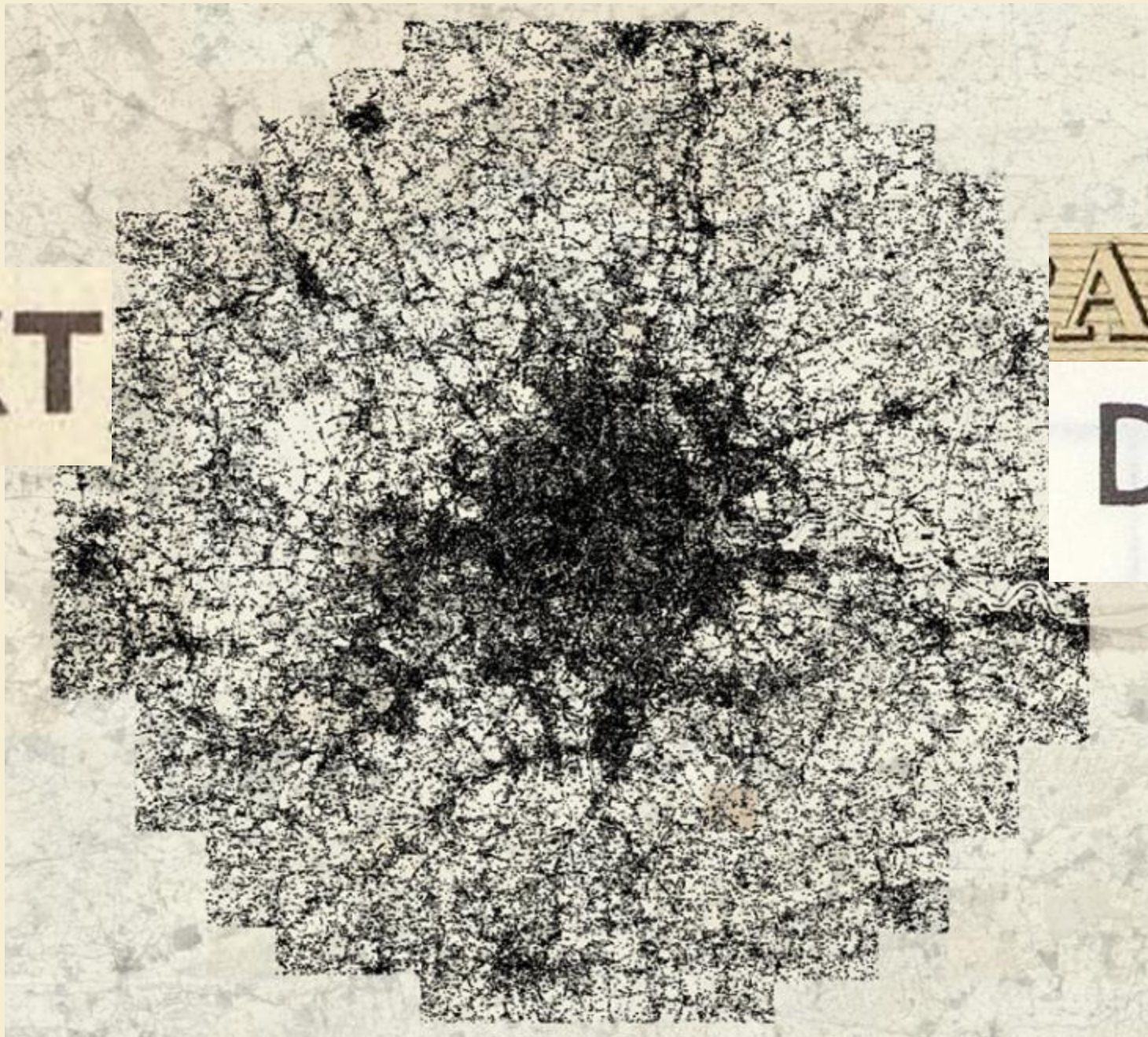


Neighbourhood Walks: A New Semantic Topology for Historical Map Text

Katherine McDonough, Kaspar Beelen and Daniel C. S. Wilson

MAP

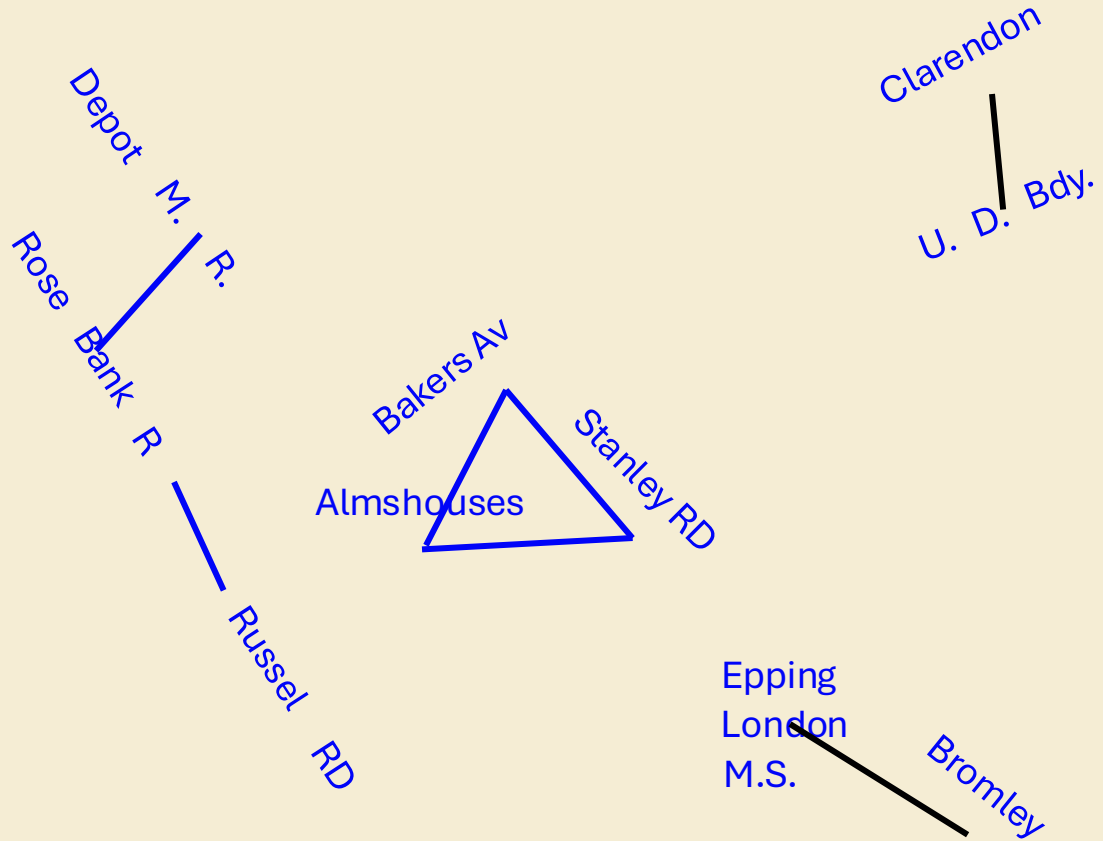
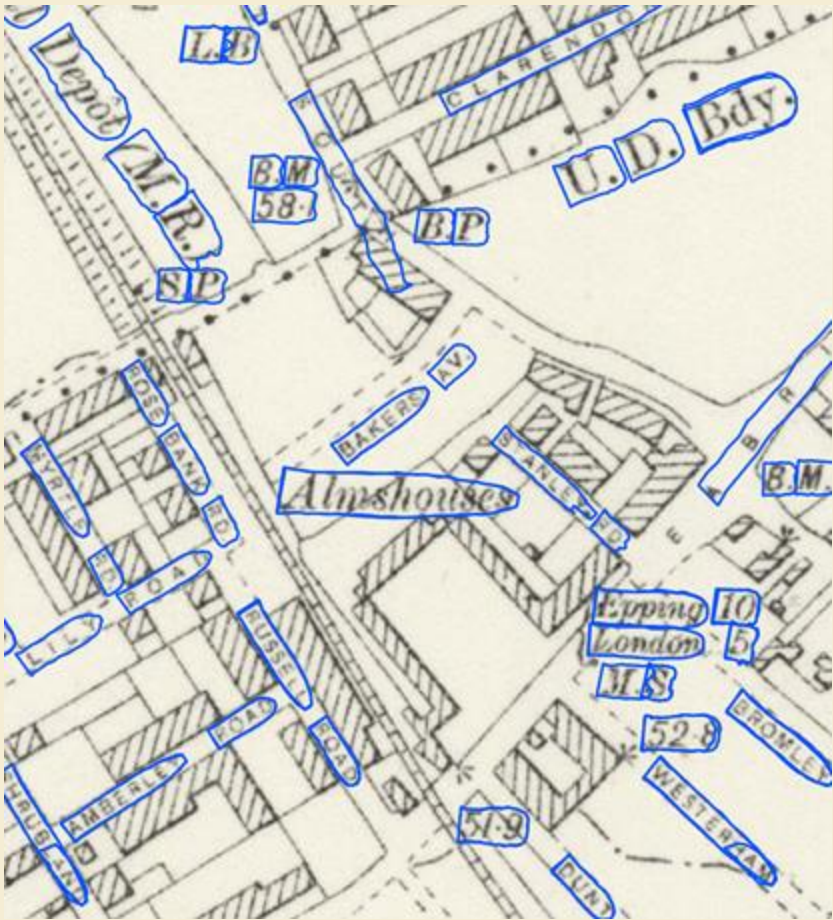
TEXT



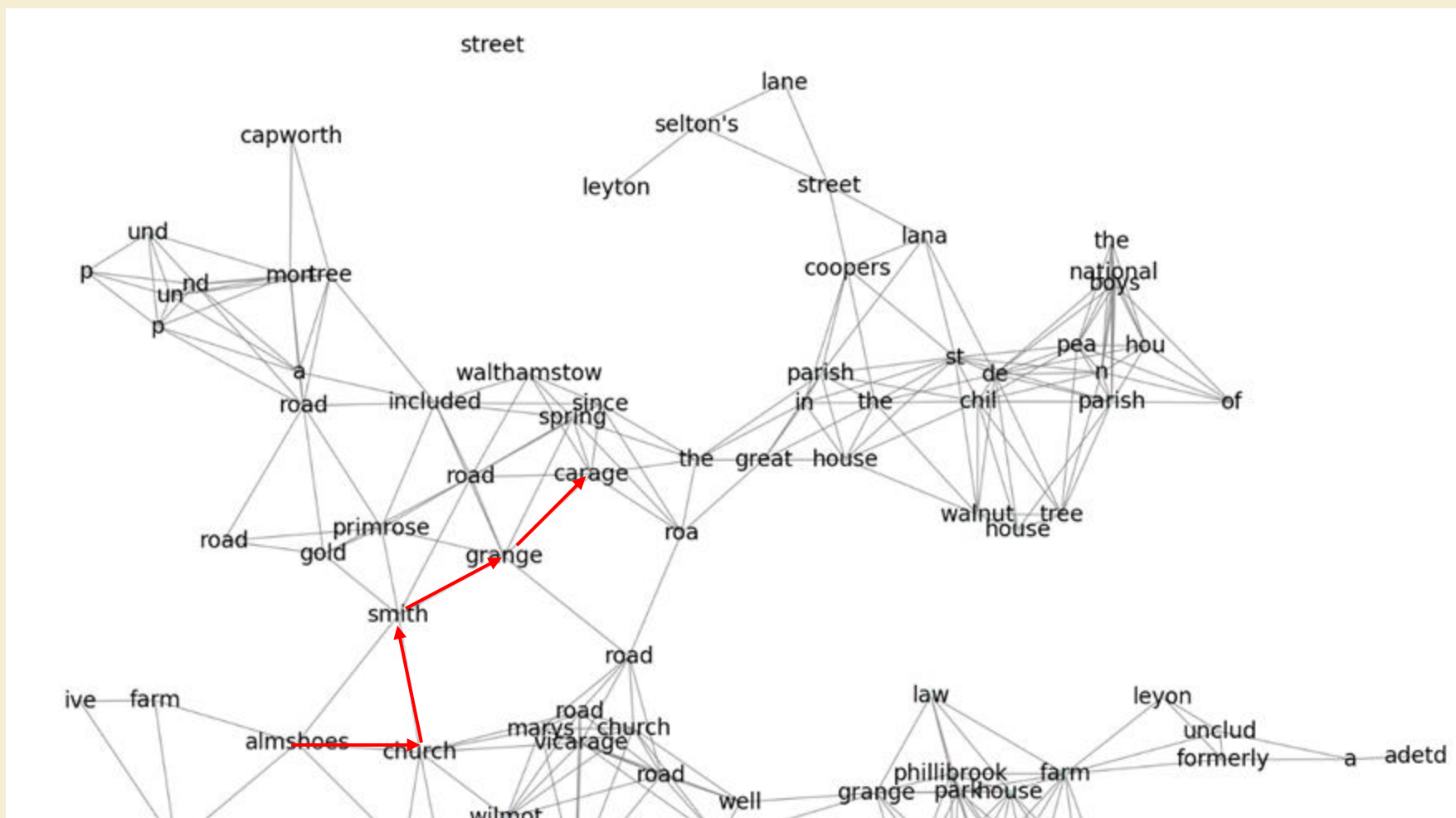
AS

DATA

Maps 2 Labels 2 Networks



Networks 2 Walks



walk: almsho[us]es \Rightarrow smith \Rightarrow orange \Rightarrow [g]arage \Rightarrow ...

Where Empires End: Tracing the Geography of a “Soaring Spirit” in Poetry

Antonina Martynenko, Artjoms Šeļa and Petr Plecháč



From A to B: Soaring view in poetry in six languages

Poet's gaze:

'From Malta's temples to the gates of Rome'

Six European poetic corpora (17–20th c.)
Czech, German, English, French, Russian, Slovenian

1086 poetic formulas



Method

1

Extract
from A – to B
formula from a poem

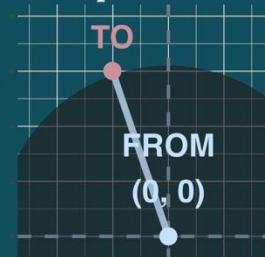
2

Plot & calculate
Haversine distance

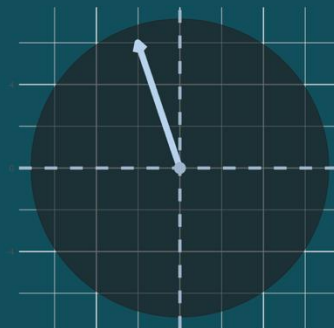


3

Transpose to FROM(0,0)



4

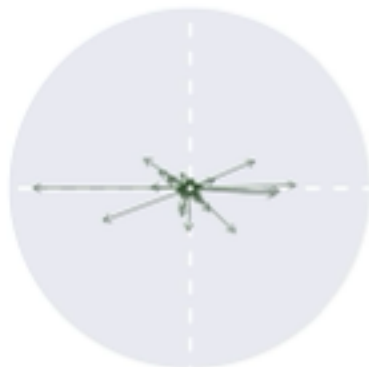


East <—> West gaze of empire

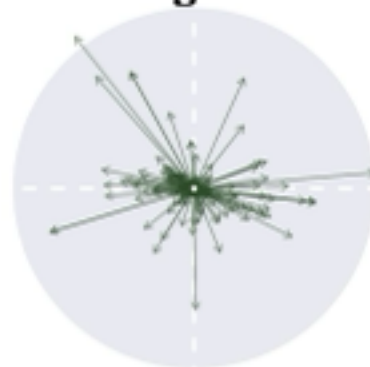
Czech



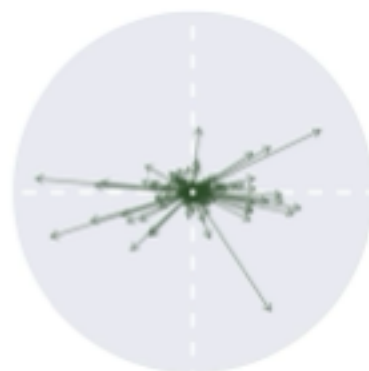
German



English



French



Russian



Slovenian

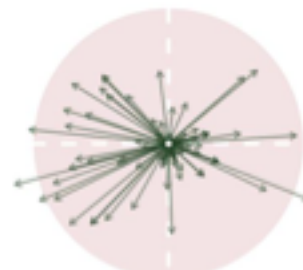


Longer distances (>1000 km)

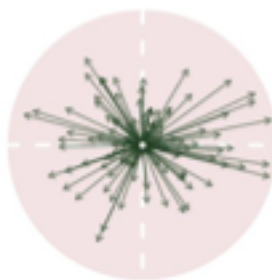
East <—> West direction (except English!)

"Younger" national literatures explore more **local** spaces

fr



en



Shorter distances
do not have
particular directions



Shared symbolic
borders and centers
of European poetry

Locations appeared in
4, 5, or all 6 corpora and
their *from* directions
(normalised distances)

Rapid Cultural Analytics Using LLMs: A Case of Dreams

Andres Karjus

Tallinn University
Estonian Business School
University of Tartu

@andreskarjus on LinkedIn/Bluesky/X

- The Estonian “It’s coming together!” NGO (Hakkab looma) collected dreams and aspirations at various events across the country
- 873 groups of people (~5503 people total; btw Estonia is 1.3M)
- Not very standardized data collection though... lots of text files

****Ages:**** 16 16 17

****Gender:****

Boys: 1 man, 1 man, and another man

Girls: no girls!

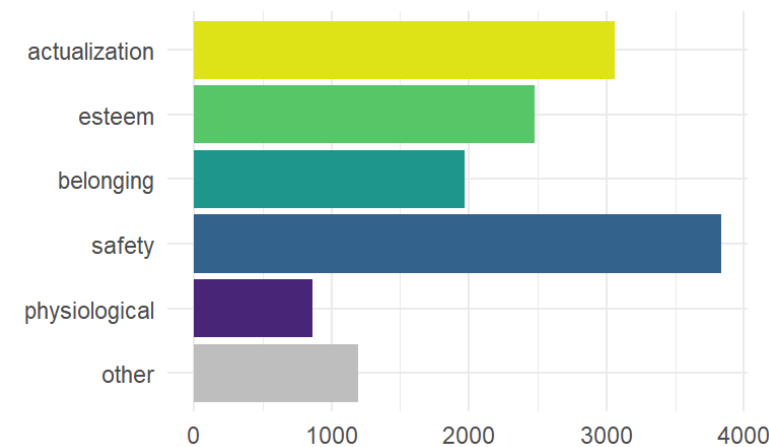
****Place of residence:**** First man – Milky Way, solar system, Earth. Other men: Narva

****Question 1:**** Answers: Right now I dream of living abroad. As a kid I used to dream of playing video games all day. Dominus hat in roblox. To live in America.

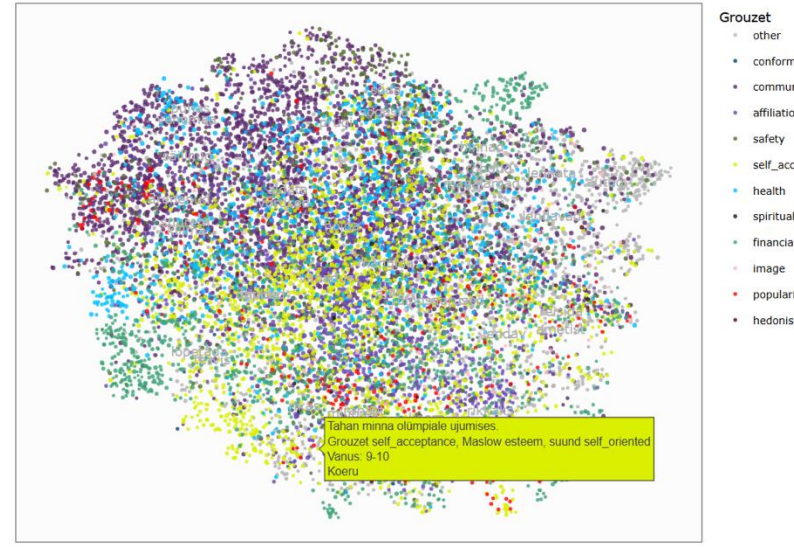
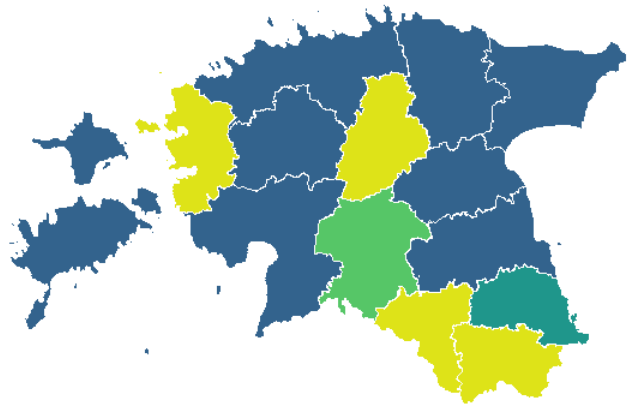
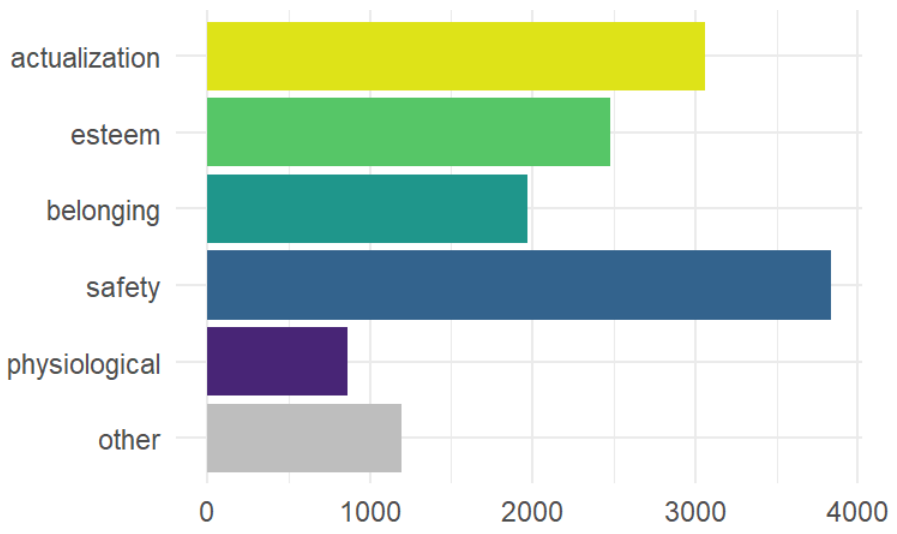
****Question 2:**** In 2050 there’s new tech like super fast cars, and Minecraft is round.

- The Estonian “It’s coming together!” (Hakkab looma) NGO collected dreams and aspirations at various events across the country
- 873 groups of people (~5503 people total; btw Estonia is 1.3M)
- No standardized data collection procedures though, lots of text files

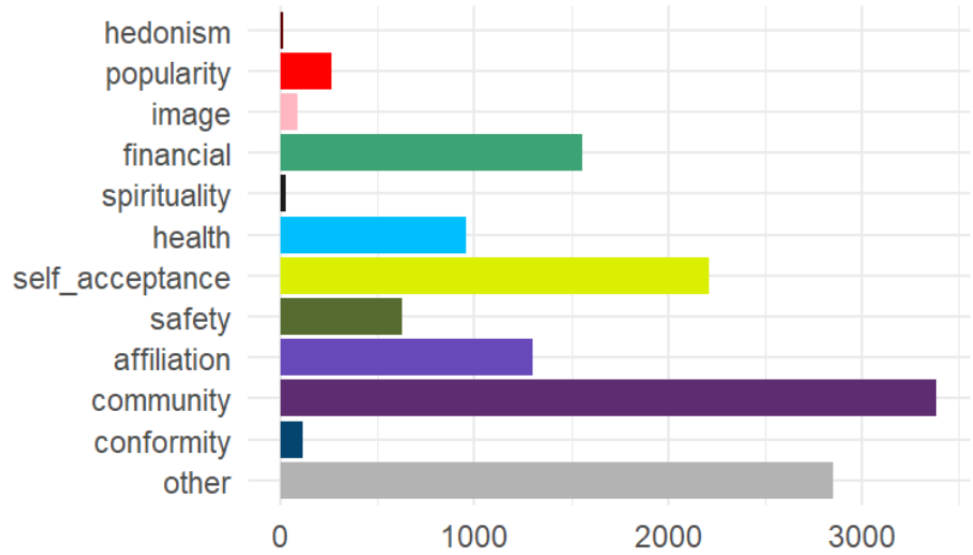
- A little inter-sector collaboration (with Innar Liiv at Taltech)
- Challenge: how to structure and analyze?
- Solution: LLM, structured outputs (GPT-4.1 + Pydantic)
- Parse, unitize and classify - all in one pass
- Result: 13407 dream units (1 to 145 per group)
- Classified into several standard psychological taxonomies + topics



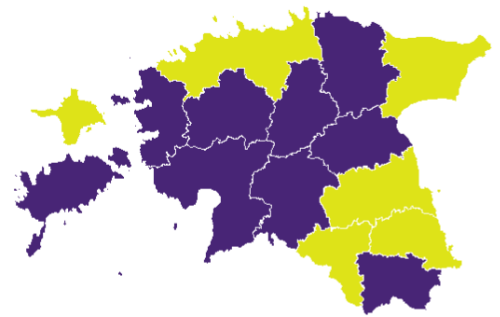
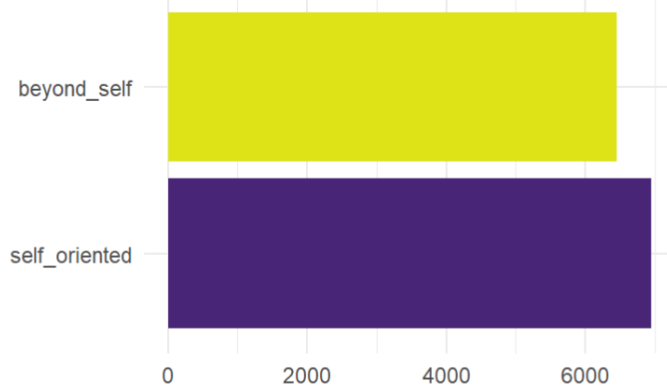
Maslow's needs:



Grouzet's goals; top: community feeling, to improve the world



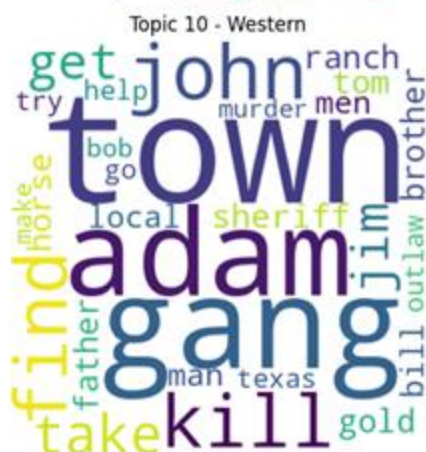
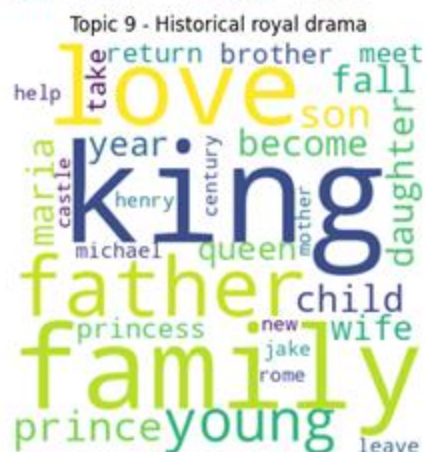
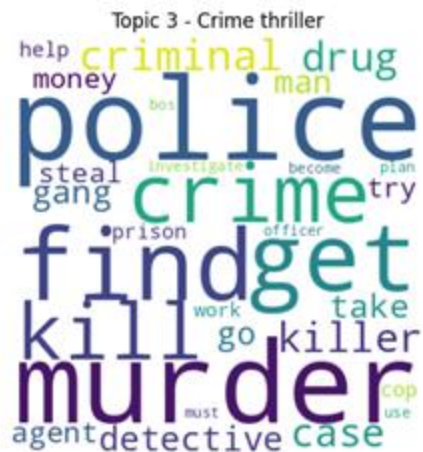
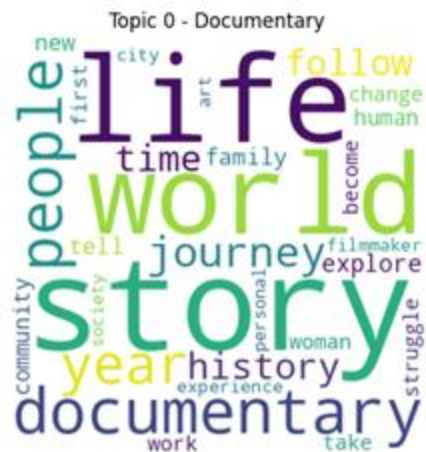
Damon's goals: self-oriented or beyond self

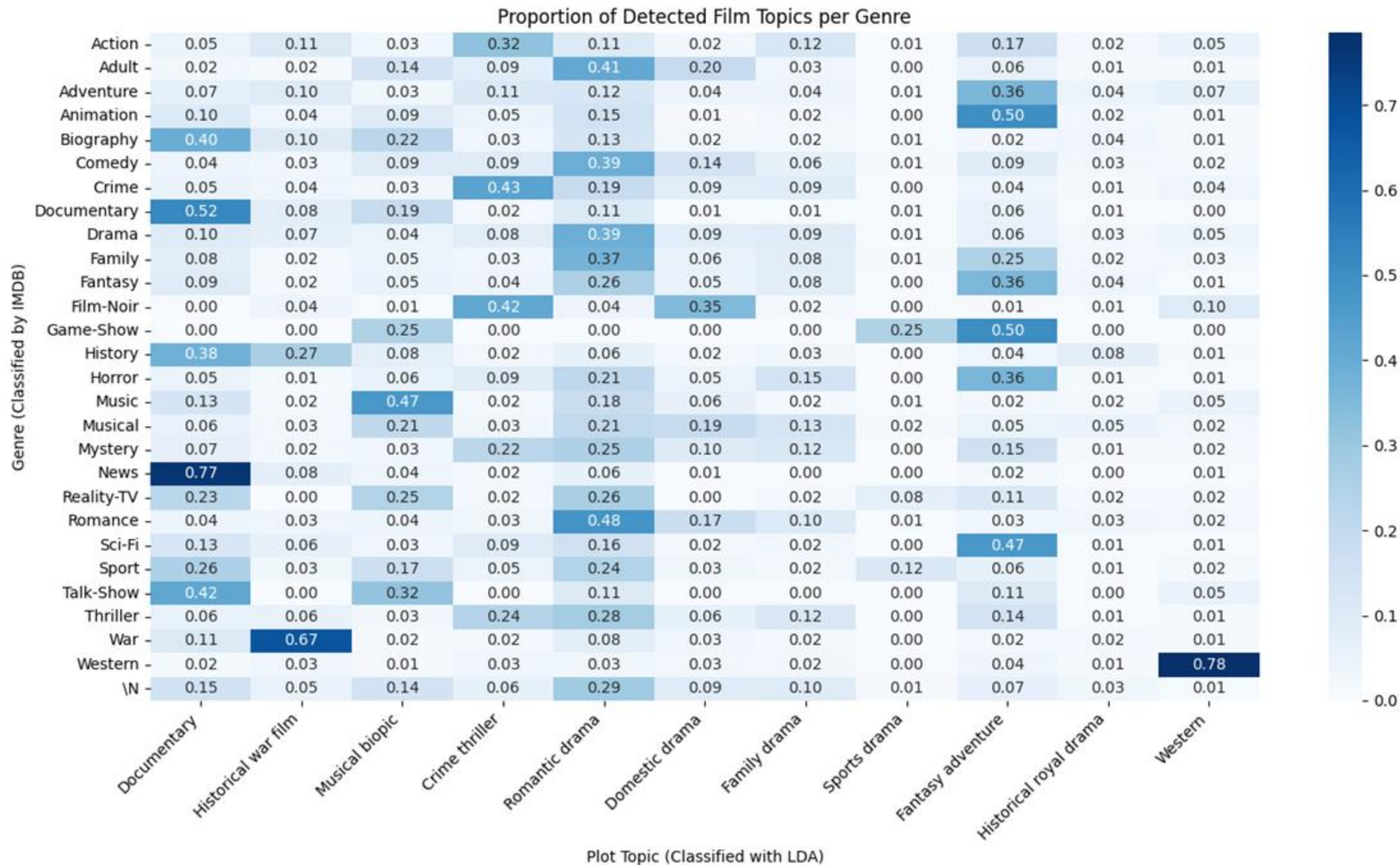


Also, we're hiring postdocs, ERC level salary, remote possible

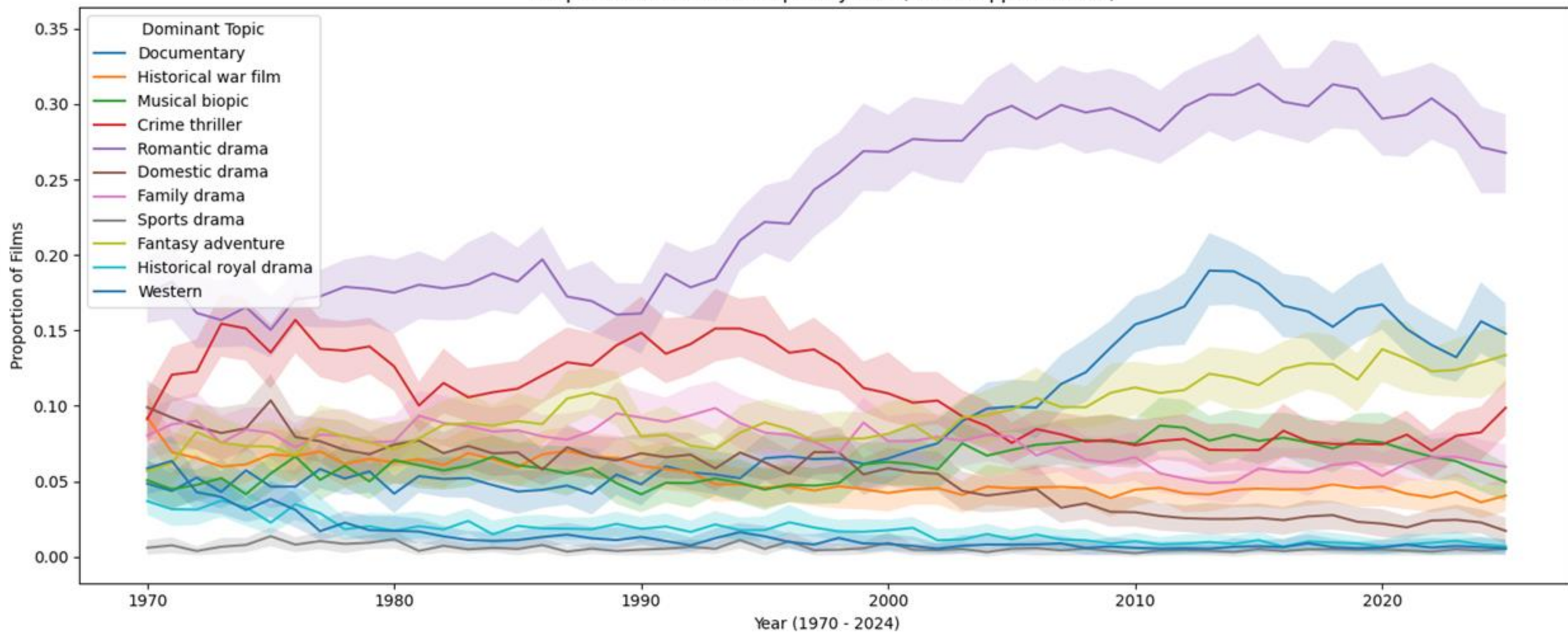
A Diachronic Analysis of Cinematic Trends and Their Reception

Jan Tvrz





Proportion of Dominant Topics by Year (Bootstrapped 95% CI)

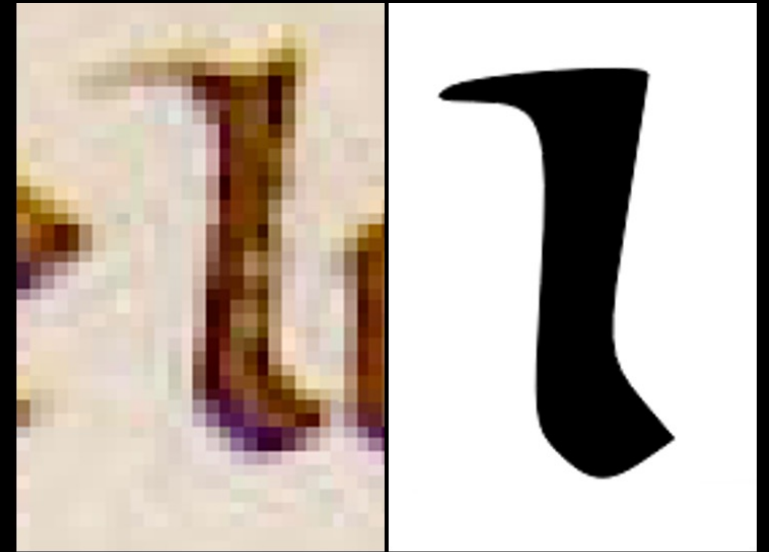


Quill2Vec: A Tool for Vector Manipulation of Medieval Latin Script

Mart Herman Gerrit Makkink

- New and innovative pipeline
 - Based on SVGs
 - For palaeography
- Convert individual
 - Characters
 - Pen strokes
 - Part of strokes
- Scalable Vector Graphics (SVG)
 - Mathematically expressed shapes
 - XML-like
- GUI-based
 - No programming required

WHEN TO USE QUIL2VEC?



TECHNICAL

- Front end is built in Qt
- All processes are done through Python
- CV2 for image preprocessing
- Potrace for image tracing

