



Collections
Lexical
Resources
Editions
Infrastructure/
Operations

Posterslam

3. Text+ Plenary in Mannheim

Gefördert durch

DFG Deutsche
Forschungsgemeinschaft

Förderungsnummer 460033370

Teil der

nfdi Nationale
Forschungsdaten
Infrastruktur

<https://www.text-plus.org>



Agenda

- **Anwendungen und Methoden der Abgeleiteten Textformate (ATF) im Kontext von LLMs**

Florian Barth, José Calvo Tello, Keli Du, Philippe Genêt, Peter Leinen, Jörg Knappen, Thorsten Trippel, Andreas Witt

- ~~Decoding Discourse: Gender Dynamics in German Bundestag Debates (1949-2021)~~

~~Teresa Hailer~~

- Empowering AI Knowledge Management: A Community-Organizing Approach to Enhance Fidelity and Quality through Authority File Use in Metadata

Barbara Fischer

- Entwicklung von Transformer-basierten Modellen für historische Textnormalisierung

Yannic Bracke, Gregor Middell, Alexander Geyken

- Evaluation of LLMs to Support the Development of GermaNet

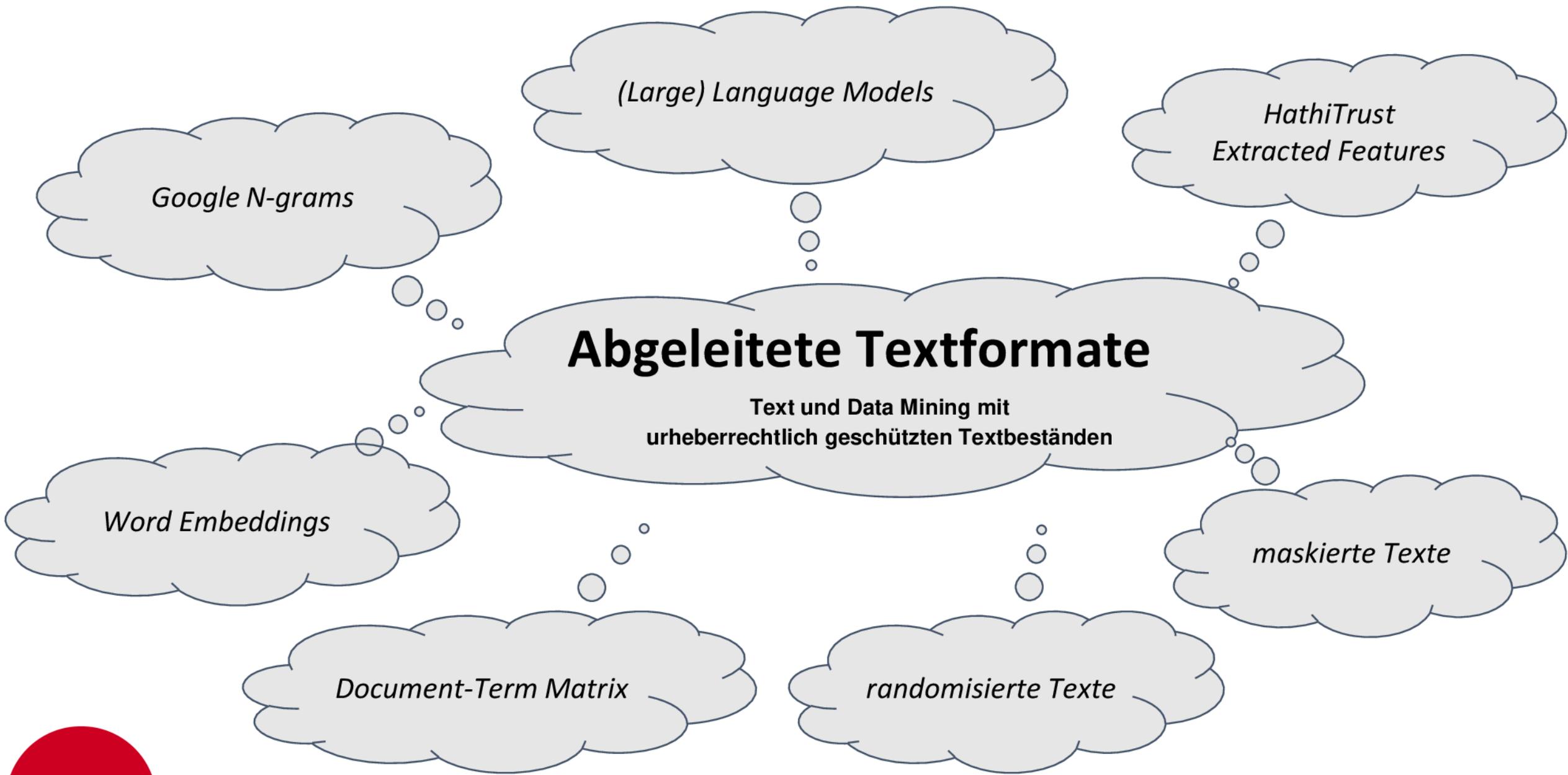
Reinhild Barkey, Erhard Hinrichs, Marie Hinrichs, Kimberly Sharp, Claus Zinn

- HERMES – Humanities Education in Research, Data, and Methods

Ruth Reiche, Andrea Rapp, Anna Schlander, Ksenia Stanicka-Brzezicka, Johanna Konstanciak

Collections
Lexical
Resources
Editions
Infrastructure/
Operations





Agenda

- Anwendungen und Methoden der Abgeleiteten Textformate (ATF) im Kontext von LLMs
Florian Barth, José Calvo Tello, Keli Du, Philippe Genêt, Peter Leinen, Jörg Knappen, Thorsten Trippel, Andreas Witt
- ~~Decoding Discourse: Gender Dynamics in German Bundestag Debates (1949-2021)~~
Teresa Hailer
- **Empowering AI Knowledge Management: A Community-Organizing Approach to Enhance Fidelity and Quality through Authority File Use in Metadata**
Barbara Fischer
- Entwicklung von Transformer-basierten Modellen für historische Textnormalisierung
Yannic Bracke, Gregor Middell, Alexander Geyken
- Evaluation of LLMs to Support the Development of GermaNet
Reinhild Barkey, Erhard Hinrichs, Marie Hinrichs, Kimberly Sharp, Claus Zinn
- HERMES – Humanities Education in Research, Data, and Methods
Ruth Reiche, Andrea Rapp, Anna Schlander, Ksenia Stanicka-Brzezicka, Johanna Konstanciak

Collections
Lexical
Resources
Editions
Infrastructure/
Operations



KI-basierte Wissensgraphen besser machen

EMPOWERING AI KNOWLEDGE MANAGEMENT 

A COMMUNITY-ORGANIZING APPROACH TO ENHANCE FIDELITY AND QUALITY THROUGH AUTHORITY FILE USE IN METADATA

AI helps creating metadata automatically. Authority data minimizes the risks of AI hallucination and increases precision and recall results in metadata based operations. Even more, authority data may embrace the diversity of GLAM and science communities. Thus, the German National Library and its partners further an participative infrastructure that empowers a broad range of diverse communities to use, contribute and edit the Integrated Authority File (GND) in German-speaking countries.

AI and community collaboration on metadata in GLAM and science

The GND Knowledge Graph

-  informing and building community
-  tools for exploring
-  community empowerment
-  interoperable APIs

Organizing Communities

DFG DataCite DEUTSCHE BIBLIOTHEK HELMHOLTZ Open Science TIB UNIVERSITÄT BIELEFELD PID

1. GND-Normdaten integrieren
2. Precision & Recall verbessern
3. weniger Halluzinationen
4. mehr demokratische Kontrolle bei der Wissensorganisation durch partizipative Gestaltung der Normdaten
5. Integration von mehr Communities ermöglichen

Agenda

~~• Decoding Discourse: Gender Dynamics in German Bundestag Debates (1949-2021)~~

~~Teresa Hailer~~

• Empowering AI Knowledge Management: A Community-Organizing Approach to Enhance Fidelity and Quality through Authority File Use in Metadata

Barbara Fischer

• **Entwicklung von Transformer-basierten Modellen für historische Textnormalisierung**

Yannic Bracke, Gregor Middell, Alexander Geyken

• Evaluation of LLMs to Support the Development of GermaNet

Reinhild Barkey, Erhard Hinrichs, Marie Hinrichs, Kimberly Sharp, Claus Zinn

• HERMES – Humanities Education in Research, Data, and Methods

Ruth Reiche, Andrea Rapp, Anna Schlander, Ksenia Stanicka-Brzezicka, Johanna Konstanciak

• KI-gestützte Workflows im Umgang mit gesprochenen sprachlichen Daten

Alina Hemmer

Collections

Lexical

Resources

Editions

Infrastructure/

Operations



Transnormer – ein Transformer-basierter Ansatz für die Normalisierung historischer Texte

Yannic Bracke, Gregor Middell, Alexander Geyken (BBAW)

Nemlich man nimmt einen hölzernen Teller, theilet denselben in die Helffte, und zeichnet nach justen Winckeln beykommende fünf Linien vor sich, auff sel

Nemlich man nimmt einen **hölzernen** Teller, theilet denselben in die **Helffte**, und zeichnet nach justen Winckeln **beykommende** **fünf** Linien vor sich, auff ...



Nämlich man nimmt einen **hölzernen** Teller, teilt denselben in die **Hälfte**, und zeichnet nach justen Winckeln **beikommende** **fünf** Linien vor sich, auf ...

Quelle: Fleming, Hans Friedrich von: Der Vollkommene Teutsche Jäger. Bd. 1. Leipzig, 1719. In: Deutsches Textarchiv, https://www.deutschestextarchiv.de/fleming_jaeger01_1719/425

Agenda

- Empowering AI Knowledge Management: A Community-Organizing Approach to Enhance Fidelity and Quality through Authority File Use in Metadata
Barbara Fischer
- Entwicklung von Transformer-basierten Modellen für historische Textnormalisierung
Yannic Bracke, Gregor Middell, Alexander Geyken
- **Evaluation of LLMs to Support the Development of GermaNet**
Reinhild Barkey, Erhard Hinrichs, Marie Hinrichs, Kimberly Sharp, Claus Zinn
- HERMES – Humanities Education in Research, Data, and Methods
Ruth Reiche, Andrea Rapp, Anna Schlander, Ksenia Stanicka-Brzezicka, Johanna Konstanciak
- KI-gestützte Workflows im Umgang mit gesprochen sprachlichen Daten
Alina Hemmer
- Korpusproduktion in Zeiten großer Sprachmodelle
Thomas Eckart, Christopher Schröder, Erik Körner, Felix Helfer, Frank Binder

Collections
Lexical
Resources
Editions
Infrastructure/
Operations





- GermaNet mit 215,000 lexikalischen Einträgen, organisiert in 167,163 synsets
- Für Verben gibt es immer Beispielsätze
 - für Nomen und Adjektive aber selten
- Verwende LLMs zur Anreicherung
 - sind auf großen Datenmengen trainiert
 - generierte Sätze spiegeln die Wahrscheinlichkeit eines Satzes im Korpus wider
 - Beispielsätze häufig typisch
 - bessere Datengrundlage als Sprachexperten
- Evaluation der von ChatGPT erzeugten Beispielsätze
 - Stärken und Schwächen (Spoiler: Frequenzeffekte und Polysemie)

Respekt n. (Gefühl)

Paraphrase: *leichte Angst vor etwas*

Beispiel: *Er hatte vor Schlangen großen Respekt.*

Hyperonym: *Angst, Bammel, Muffe, Schiss...*

Hyponym: *Heidenrespekt*

Agenda

- Entwicklung von Transformer-basierten Modellen für historische Textnormalisierung
Yannic Bracke, Gregor Middell, Alexander Geyken
- Evaluation of LLMs to Support the Development of GermaNet
Reinhild Barkey, Erhard Hinrichs, Marie Hinrichs, Kimberly Sharp, Claus Zinn
- **HERMES – Humanities Education in Research, Data, and Methods**
Ruth Reiche, Andrea Rapp, Anna Schlander, Ksenia Stanicka-Brzezicka, Johanna Konstanciak
- KI-gestützte Workflows im Umgang mit gesprochen sprachlichen Daten
Alina Hemmer
- Korpusproduktion in Zeiten großer Sprachmodelle
Thomas Eckart, Christopher Schröder, Erik Körner, Felix Helfer, Frank Binder
- Legal Linguistic Memos mit Large Language Models: Automatisierte Erfassung und Klassifizierung von Sachverhaltsbeschreibungen im Familienrecht
Margret Mundorf

Collections
Lexical
Resources
Editions
Infrastructure/
Operations



HERMES

Humanities Education in Research, Data, and Methods



FORSCHEN



Bring-your-own-data-Lab



Data Challenges



Forschungsstudienprogramm

LERNEN



Data Carpentries



Transferwerkstatt



Open Educational Resources

VERNETZEN



Kommunikationsformate



HERMES-Hub



WissKomm Academy – HERMES Edition

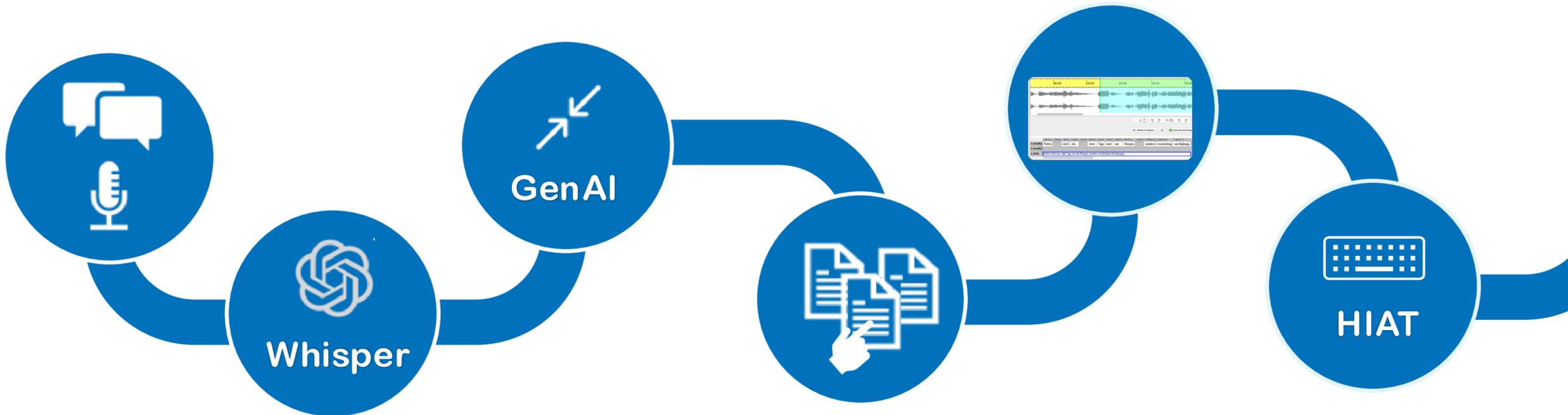


Agenda

- Evaluation of LLMs to Support the Development of GermaNet
Reinhild Barkey, Erhard Hinrichs, Marie Hinrichs, Kimberly Sharp, Claus Zinn
- HERMES – Humanities Education in Research, Data, and Methods
Ruth Reiche, Andrea Rapp, Anna Schlander, Ksenia Stanicka-Brzezicka, Johanna Konstanciak
- **KI-gestützte Workflows im Umgang mit gesprochen sprachlichen Daten**
Alina Hemmer
- Korpusproduktion in Zeiten großer Sprachmodelle
Thomas Eckart, Christopher Schröder, Erik Körner, Felix Helfer, Frank Binder
- Legal Linguistic Memos mit Large Language Models: Automatisierte Erfassung und Klassifizierung von Sachverhaltsbeschreibungen im Familienrecht
Margret Mundorf
- LLOD-isierung des Madras Tamil Lexicon: Modellierung eines Wörterbuchs einer in der IT-basierten Linguistik unterrepräsentierten außereuropäischen Sprache als Linguistic Linked Open Data
Liudmila Olalde, Thomas Malten, Frank Grieshaber

Collections
Lexical
Resources
Editions
Infrastructure/
Operations





KI-gestützte Workflows im Umgang mit gesprochenen sprachlichen Daten

Text+ Plenary, 10.10.2024

Kristin Bührig
Marcel Fladrich
Alina Hemmer

Agenda

- HERMES – Humanities Education in Research, Data, and Methods
Ruth Reiche, Andrea Rapp, Anna Schlander, Ksenia Stanicka-Brzezicka, Johanna Konstanciak
- KI-gestützte Workflows im Umgang mit gesprochen sprachlichen Daten
Alina Hemmer
- **Korpusproduktion in Zeiten großer Sprachmodelle**
Thomas Eckart, Christopher Schröder, Erik Körner, Felix Helfer, Frank Binder
- Legal Linguistic Memos mit Large Language Models: Automatisierte Erfassung und Klassifizierung von Sachverhaltsbeschreibungen im Familienrecht
Margret Mundorf
- LLOD-isierung des Madras Tamil Lexicon: Modellierung eines Wörterbuchs einer in der IT-basierten Linguistik unterrepräsentierten außereuropäischen Sprache als Linguistic Linked Open Data
Liudmila Olalde, Thomas Malten, Frank Grieshaber
- More uniformity and more diversity at the same time: LLMs and a 21st century standardisation paradox
Christian Mair

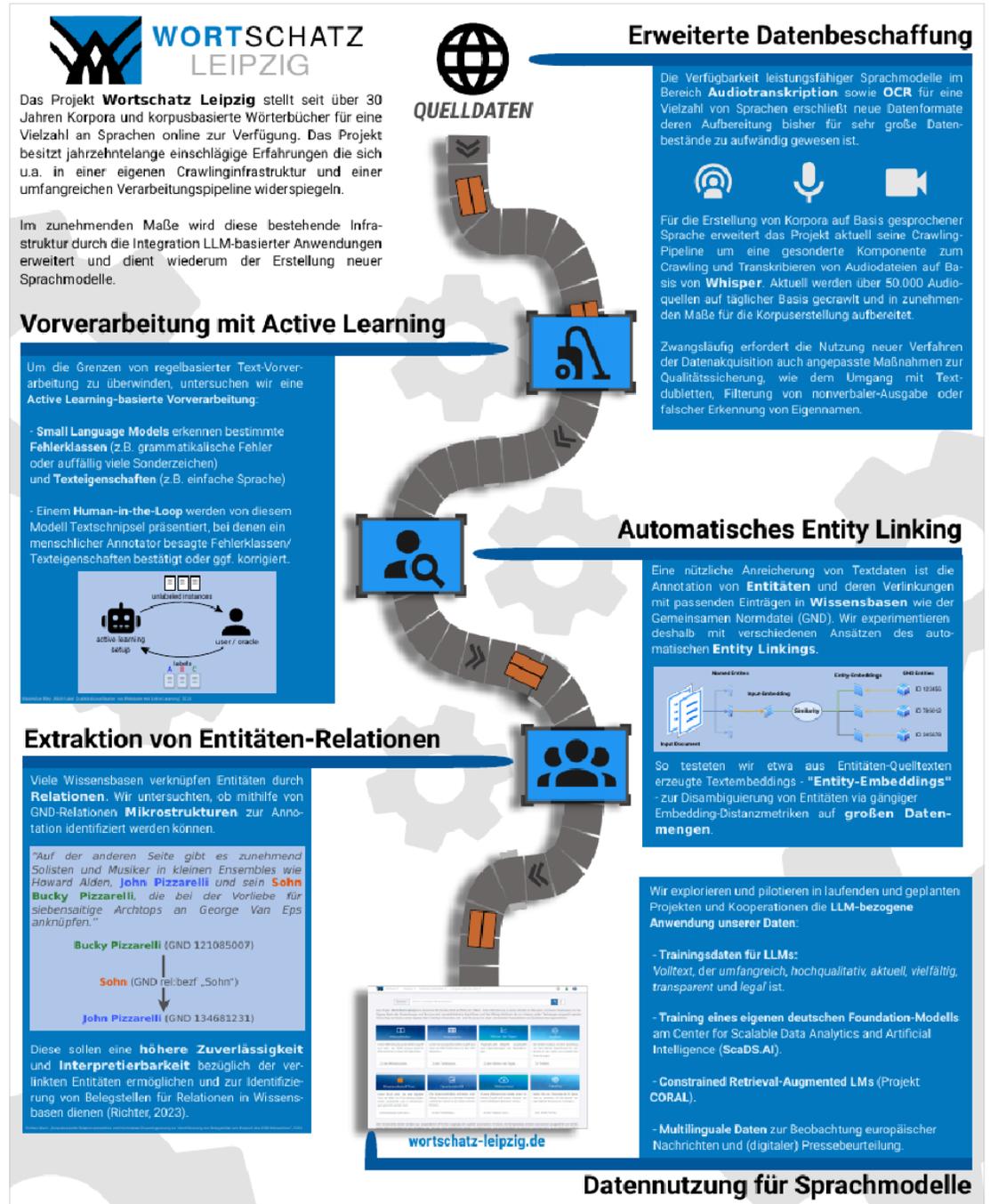
Collections
Lexical
Resources
Editions
Infrastructure/
Operations



Korpusproduktion in Zeiten großer Sprachmodelle

Integration von LLMs in die Verarbeitungspipeline des Wortschatz Leipzig

Thomas Eckart, Felix Helfer, Erik Körner, Frank Binder, Christopher Schröder



Agenda

- KI-gestützte Workflows im Umgang mit gesprochenen sprachlichen Daten
Alina Hemmer
- Korpusproduktion in Zeiten großer Sprachmodelle
Thomas Eckart, Christopher Schröder, Erik Körner, Felix Helfer, Frank Binder
- **Legal Linguistic Memos mit Large Language Models: Automatisierte Erfassung und Klassifizierung von Sachverhaltsbeschreibungen im Familienrecht**
Margret Mundorf
- LLOD-isierung des Madras Tamil Lexicon: Modellierung eines Wörterbuchs einer in der IT-basierten Linguistik unterrepräsentierten außereuropäischen Sprache als Linguistic Linked Open Data
Liudmila Olalde, Thomas Malten, Frank Grieshaber
- More uniformity and more diversity at the same time: LLMs and a 21st century standardisation paradox
Christian Mair
- Perspektiven des Einsatzes von LLM in Text+
Florian Barth, Philippe Genêt, Erik Körner, Peter Leinen, Pia Schwarz, Claus Zinn

Collections
Lexical
Resources
Editions
Infrastructure/
Operations





Pitch

Agenda

Collections
Lexical
Resources
Editions
Infrastructure/
Operations

- Korpusproduktion in Zeiten großer Sprachmodelle
Thomas Eckart, Christopher Schröder, Erik Körner, Felix Helfer, Frank Binder
- Legal Linguistic Memos mit Large Language Models: Automatisierte Erfassung und Klassifizierung von Sachverhaltsbeschreibungen im Familienrecht
Margret Mundorf
- **LLOD-isierung des Madras Tamil Lexicon: Modellierung eines Wörterbuchs einer in der IT-basierten Linguistik unterrepräsentierten außereuropäischen Sprache als Linguistic Linked Open Data**
Liudmila Olalde, Thomas Malten, Frank Grieshaber
- More uniformity and more diversity at the same time: LLMs and a 21st century standardisation paradox
Christian Mair
- Perspektiven des Einsatzes von LLM in Text+
Florian Barth, Philippe Genêt, Erik Körner, Peter Leinen, Pia Schwarz, Claus Zinn
- ~~Reaching for the stars: Integration von LLMs in komplexe automatisierte Workflows~~
Daniela Schneider



LLOD-isierung des Madras Tamil Lexicon

Modellierung eines Wörterbuchs einer in der IT-basierten Linguistik unterrepräsentierten außereuropäischen Sprache als Linguistic Linked Open Data

Worum geht es? Um ein Tamil-Englisch Wörterbuch.

Was ist Tamil? Eine südindische Sprache nicht indoeuropäischen Ursprungs.

Wo wird Tamil gesprochen? Hauptsächlich in Südindien und Sri Lanka.

Wie viele Leute sprechen Tamil? Erstsprache von ca. 70-80 Millionen Menschen.

Was ist das Tamil Lexicon? Ein von 1924-1939 in sechs Bänden und einem Nachtragsband erschienenenes historisch-literarisches Wörterbuch (Madras University).

Ist es nicht bereits online? Ja, aber in eingeschränkter Funktionalität in proprietärem Format.

Was machen wir? Eine Neumodellierung der digital vorliegenden Daten als LLOD.

Ist das alles? Nein, wir machen noch eine API.

Wer hat was davon? Alle.

Wo erfahre ich mehr davon? Am Posterstand.

தமிழ்

Agenda

Collections

Lexical

Resources

Editions

Infrastructure/

Operations

- Legal Linguistic Memos mit Large Language Models: Automatisierte Erfassung und Klassifizierung von Sachverhaltsbeschreibungen im Familienrecht

Margret Mundorf

- LLOD-isierung des Madras Tamil Lexicon: Modellierung eines Wörterbuchs einer in der IT-basierten Linguistik unterrepräsentierten außereuropäischen Sprache als Linguistic Linked Open Data

Liudmila Olalde, Thomas Malten, Frank Grieshaber

- **More uniformity and more diversity at the same time: LLMs and a 21st century standardisation paradox**

Christian Mair

- Perspektiven des Einsatzes von LLM in Text+

Florian Barth, Philippe Genêt, Erik Körner, Peter Leinen, Pia Schwarz, Claus Zinn

- ~~Reaching for the stars: Integration von LLMs in komplexe automatisierte Workflows~~

~~Daniela Schneider~~

- SwineBad: Tabellenextraktion und Informationsstrukturierung aus dem Swinemünder Badeanzeiger

Steffen Steiner, Frank Krüger



Christian Mair

Uniformity and Diversity at the Same Time: LLMs and the 21st Century Standardisation Paradox

- Die Geisteswissenschaften sollen aufhören, große Sprachmodelle nur mit Sorgenfalten zu betrachten.
- Sie schenken uns fantastische neue Forschungsthemen:
 - LLMs und Sprachplanung für intelligente Mehrsprachigkeit
 - LLMs und Diskriminierung
 - LLMs und die kritische Ethnographie der *clickworker sweatshops* von den Philippinen bis Westafrika
 - ...
- **Und, Thema des Posters:**
 - Wird die britische Schriftnorm die LLMs überleben?

<https://www.etsy.com/listing/1230349109/hanging-union-jack-united-kingdom-draped> -->



Agenda

Collections

Lexical

Resources

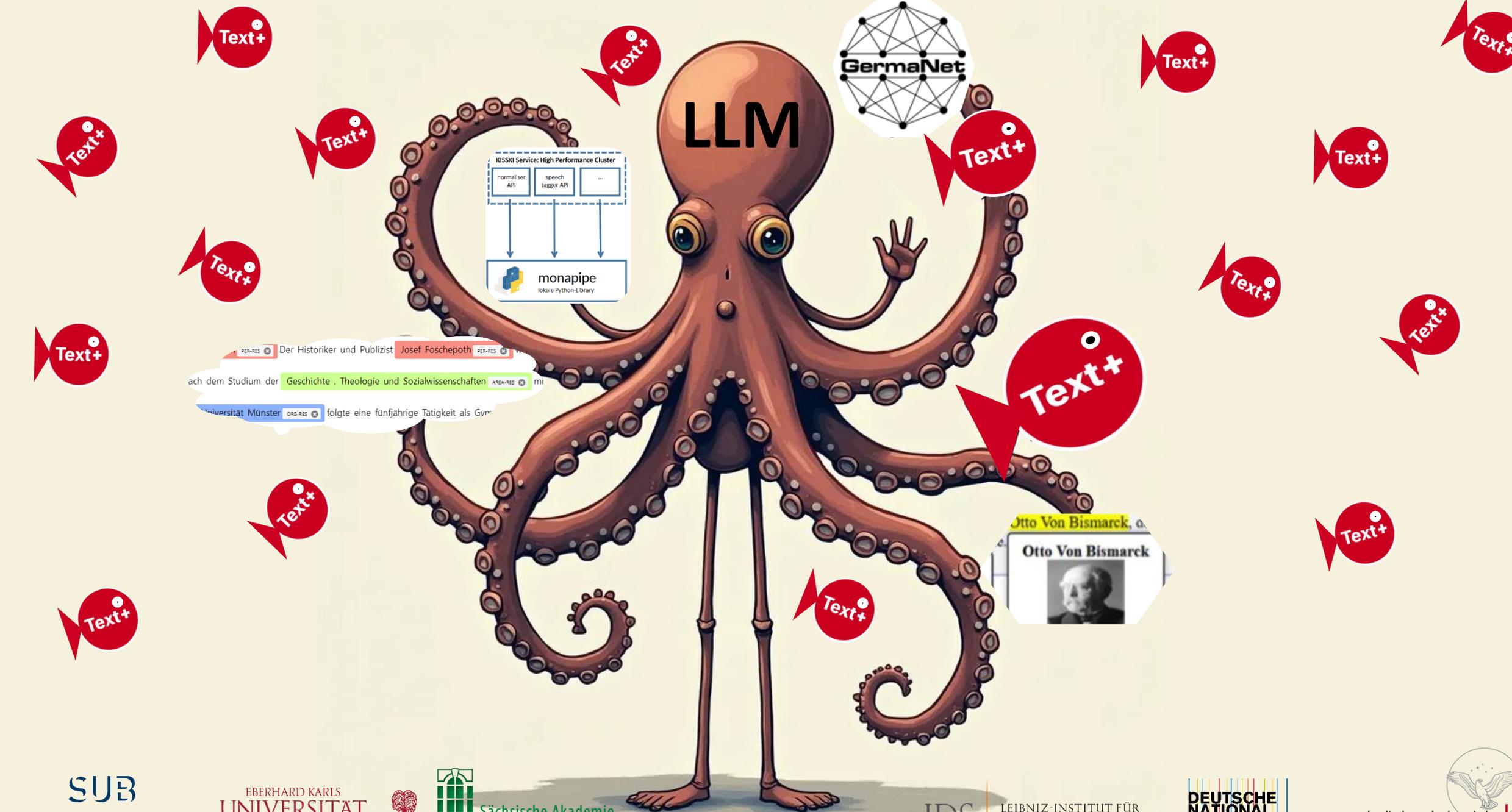
Editions

Infrastructure/

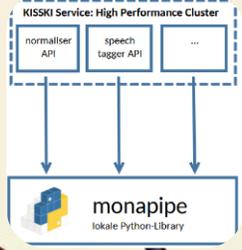
Operations

- LLOD-isierung des Madras Tamil Lexicon: Modellierung eines Wörterbuchs einer in der IT-basierten Linguistik unterrepräsentierten außereuropäischen Sprache als Linguistic Linked Open Data
Liudmila Olalde, Thomas Malten, Frank Grieshaber
- More uniformity and more diversity at the same time: LLMs and a 21st century standardisation paradox
Christian Mair
- **Perspektiven des Einsatzes von LLM in Text+**
Florian Barth, Philippe Genêt, Erik Körner, Peter Leinen, Pia Schwarz, Claus Zinn
- ~~Reaching for the stars: Integration von LLMs in komplexe automatisierte Workflows~~
Daniela Schneider
- SwineBad: Tabellenextraktion und Informationsstrukturierung aus dem Swinemünder Badeanzeiger
Steffen Steiner, Frank Krüger
- SwissGB4Science - ein Volltext Korpus für die Forschung
Eric Dubey, Matteo Lorenzini, Martin Reisacher, Tim Rüdiger

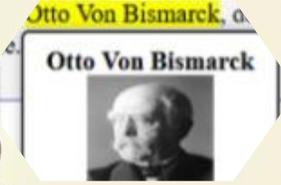




LLM



Der Historiker und Publizist Josef Foschepoth
ach dem Studium der Geschichte , Theologie und Sozialwissenschaften
Universität Münster folgte eine fünfjährige Tätigkeit als Gym



Agenda

- Perspektiven des Einsatzes von LLM in Text+
Florian Barth, Philippe Genêt, Erik Körner, Peter Leinen, Pia Schwarz, Claus Zinn
- ~~Reaching for the stars: Integration von LLMs in komplexe automatisierte Workflows~~
~~Daniela Schneider~~
- **SwineBad: Tabellenextraktion und Informationsstrukturierung aus dem Swinemünder Badeanzeiger**
Steffen Steiner, Frank Krüger
- SwissGB4Science - ein Volltext Korpus für die Forschung
Eric Dubey, Matteo Lorenzini, Martin Reisacher, Tim Rüdiger
- Synthetische Datensätze in der CLS
Daniel Kababgi, Emilie Sitter, Robin Martin Aust, Marie-Christine Boucher, Berenike Herrmann
- Text+ LLM Service
Alexander Steckel, Umut Basaran, Stefan Buddenbohm, Maik Wegener, Philipp Wieder

Collections
Lexical
Resources
Editions
Infrastructure/
Operations



SwineBad

Tabellenextraktion und Informationsstrukturierung aus dem Swinemünder Badeanzeiger



Nr.	Name und Wohnort	Beruf	Wohnung	Personenanzahl
970	Herrn Hermann Appelbaum	Elektrotechniker	Maaß, Lotsenstr. 81	1
971	Herrn Hermann Appelbaum	Elektrotechniker	Maaß, Lotsenstr. 81	1
972	Herrn Hermann Appelbaum	Elektrotechniker	Maaß, Lotsenstr. 81	1
973	Herrn Hermann Appelbaum	Elektrotechniker	Maaß, Lotsenstr. 81	1
974	Herrn Hermann Appelbaum	Elektrotechniker	Maaß, Lotsenstr. 81	1
975	Herrn Hermann Appelbaum	Elektrotechniker	Maaß, Lotsenstr. 81	1
976	Herrn Hermann Appelbaum	Elektrotechniker	Maaß, Lotsenstr. 81	1
977	Herrn Hermann Appelbaum	Elektrotechniker	Maaß, Lotsenstr. 81	1
978	Herrn Hermann Appelbaum	Elektrotechniker	Maaß, Lotsenstr. 81	1
979	Herrn Hermann Appelbaum	Elektrotechniker	Maaß, Lotsenstr. 81	1
980	Herrn Hermann Appelbaum	Elektrotechniker	Maaß, Lotsenstr. 81	1
981	Herrn Hermann Appelbaum	Elektrotechniker	Maaß, Lotsenstr. 81	1
982	Herrn Hermann Appelbaum	Elektrotechniker	Maaß, Lotsenstr. 81	1
983	Herrn Hermann Appelbaum	Elektrotechniker	Maaß, Lotsenstr. 81	1
984	Herrn Hermann Appelbaum	Elektrotechniker	Maaß, Lotsenstr. 81	1
985	Herrn Hermann Appelbaum	Elektrotechniker	Maaß, Lotsenstr. 81	1
986	Herrn Hermann Appelbaum	Elektrotechniker	Maaß, Lotsenstr. 81	1
987	Herrn Hermann Appelbaum	Elektrotechniker	Maaß, Lotsenstr. 81	1
988	Herrn Hermann Appelbaum	Elektrotechniker	Maaß, Lotsenstr. 81	1
989	Herrn Hermann Appelbaum	Elektrotechniker	Maaß, Lotsenstr. 81	1
990	Herrn Hermann Appelbaum	Elektrotechniker	Maaß, Lotsenstr. 81	1
991	Herrn Hermann Appelbaum	Elektrotechniker	Maaß, Lotsenstr. 81	1
992	Herrn Hermann Appelbaum	Elektrotechniker	Maaß, Lotsenstr. 81	1
993	Herrn Hermann Appelbaum	Elektrotechniker	Maaß, Lotsenstr. 81	1
994	Herrn Hermann Appelbaum	Elektrotechniker	Maaß, Lotsenstr. 81	1
995	Herrn Hermann Appelbaum	Elektrotechniker	Maaß, Lotsenstr. 81	1
996	Herrn Hermann Appelbaum	Elektrotechniker	Maaß, Lotsenstr. 81	1
997	Herrn Hermann Appelbaum	Elektrotechniker	Maaß, Lotsenstr. 81	1
998	Herrn Hermann Appelbaum	Elektrotechniker	Maaß, Lotsenstr. 81	1
999	Herrn Hermann Appelbaum	Elektrotechniker	Maaß, Lotsenstr. 81	1
1000	Herrn Hermann Appelbaum	Elektrotechniker	Maaß, Lotsenstr. 81	1



```
{  
  "Nummer": "970",  
  "Vorname": "Hermann",  
  "Nachname": "Appelbaum",  
  "Titel": null,  
  "Beruf": "Elektrotechniker",  
  "Sozialer Stand": null,  
  "Begleitung": null,  
  "Wohnort": "Berlin",  
  "Wohnung": "Maaß, Lotsenstr. 81",  
  "Personenanzahl": "1"  
},
```

- Tabellen ankommender Badegäste 1910–1932
- Software Pipeline zur Extraktion und Strukturierung

1. Segmentierung der Tabellen
2. OCR der Tabellendaten
3. LLM-basierte Korrektur und Zusammenfassung
4. LLM-basierte Strukturierung

Agenda

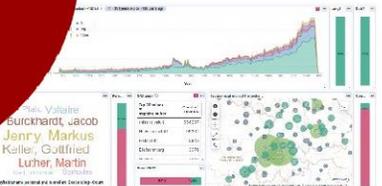
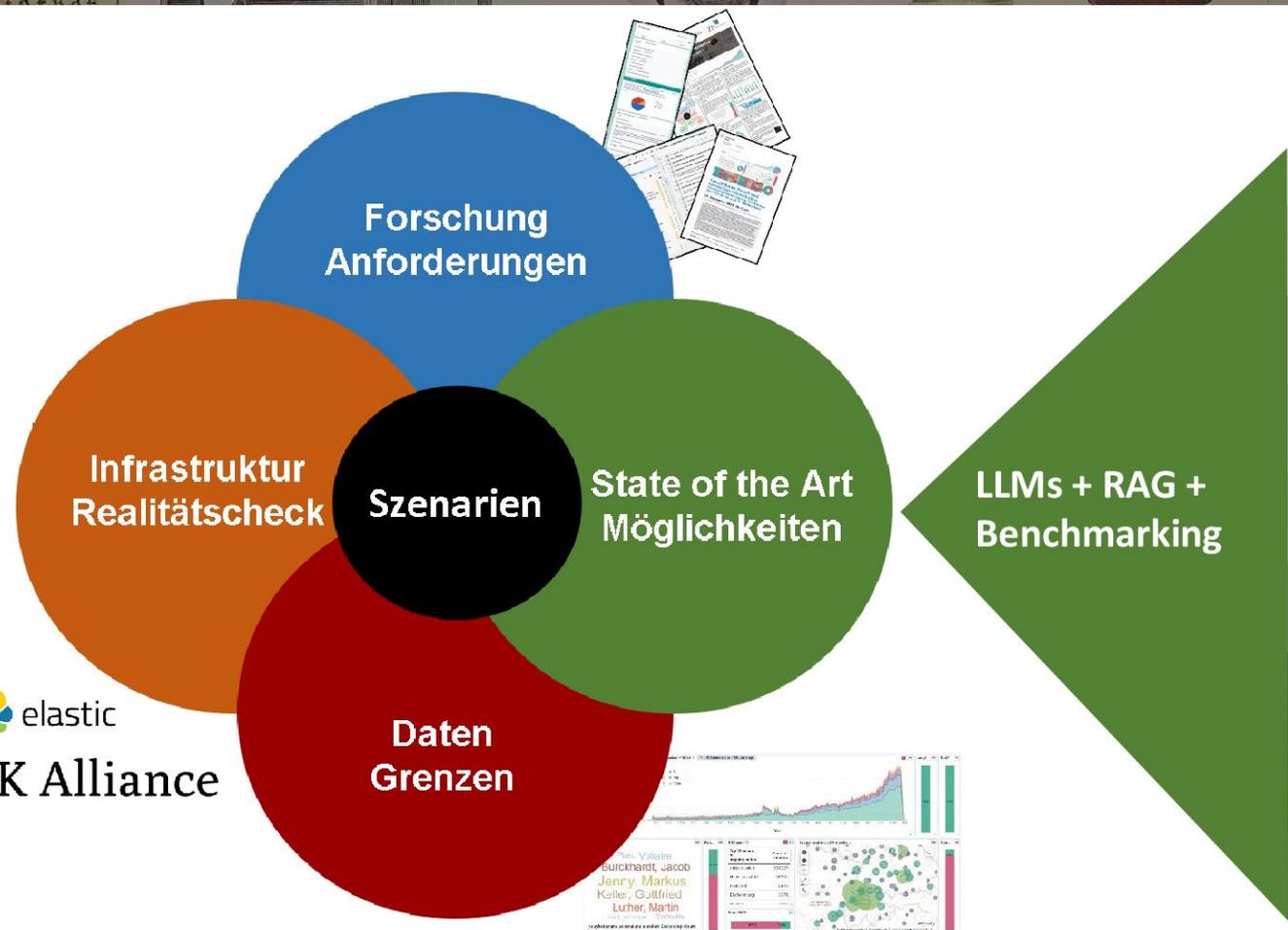
- ~~Reaching for the stars: Integration von LLMs in komplexe automatisierte Workflows~~
Daniela Schneider
- ~~SwineBad: Tabellenextraktion und Informationsstrukturierung aus dem Swinemünder Badeanzeiger~~
Steffen Steiner, Frank Krüger
- **SwissGB4Science - ein Volltext Korpus für die Forschung**
Eric Dubey, Matteo Lorenzini, Martin Reisacher, Tim Rüdiger
- Synthetische Datensätze in der CLS
Daniel Kababgi, Emilie Sitter, Robin Martin Aust, Marie-Christine Boucher, Berenike Herrmann
- Text+ LLM Service
Alexander Steckel, Umut Basaran, Stefan Buddenbohm, Maik Wegener, Philipp Wieder
- Wissen, wen man fragt – Agentic RAG für Automatisches Question Answering in der Domäne deutscher Grammatik
Christian Lang, Ngoc Duyen Tanja Tu, Roman Schneider

Collections
Lexical
Resources
Editions
Infrastructure/
Operations



Swiss Google Books for Research

90 Millionen historische Seiten für die Forschung?



Anreicherungen

Gemini

```

[{"id": "1", "text": "Einige Beispiele für die Verwendung von Gemini."},
{"id": "2", "text": "Weitere Informationen zu den verschiedenen Modellen."},
{"id": "3", "text": "Zusätzliche Ressourcen und Links."}
    
```

Zusammenfassungen

Zusammenfassung und Themenordnung des Textes:

Der Text handelt von der politischen Entwicklung in der Schweiz im 16. Jahrhundert, insbesondere in dem Aufbruch und der Konsolidierung der Helvetik und der Schweiz. Im Fokus steht dabei der Vergleich der Ansätze von **Martin Luthers** und **Ulrich Zwingli**, dem Konflikt zwischen den beiden Reformatoren in der Zürcher Bibel und der daraus resultierenden Spaltung der Protestanten.

Der Text bezieht sich auf die politische Entwicklung in der Schweiz, darunter den **Medienkampf** zwischen Kaiser Karl V. und König Franz I. von Frankreich, den **deutschen Überfall** auf Luzern und die **innerschweizerischen Konflikte** innerhalb des Heiligen Römischen Reiches.

Unterschiedliche Modelle

Agenda

- SwineBad: Tabellenextraktion und Informationsstrukturierung aus dem Swinemünder Badeanzeiger
Steffen Steiner, Frank Krüger
- SwissGB4Science - ein Volltext Korpus für die Forschung
Eric Dubey, Matteo Lorenzini, Martin Reisacher, Tim Rüdiger
- **Synthetische Datensätze in der CLS**
Daniel Kababgi, Emilie Sitter, Robin Martin Aust, Marie-Christine Boucher, Berenike Herrmann
- Text+ LLM Service
Alexander Steckel, Umut Basaran, Stefan Buddenbohm, Maik Wegener, Philipp Wieder
- Wissen, wen man fragt – Agentic RAG für Automatisches Question Answering in der Domäne deutscher Grammatik
Christian Lang, Ngoc Duyen Tanja Tu, Roman Schneider
- “Computer, was bedeutet ‘Tiki-Taka’?” Eine Studie zur Generierung von Definitionsparaphrasen für Bedeutungswörterbücher am Beispiel des DWDS
Alexander Geyken, Gregor Middell

Collections
Lexical
Resources
Editions
Infrastructure/
Operations



Synthetische Datensätze in den Computational Literary Studies

Eines der typischen Probleme in den CLS:

- Annotieren ist schwer und dauert
- Textstellen zum Annotieren finden ist schwer und dauert



Was tun!?

Lösungen müssen her:

- Sind LLMs eine Lösung?
- Falls nein: können wir zumindest mit LLMs weitere Trainingsdaten für ein eigenes Modell erzeugen?

**Falls ihr auf die letzte Frage
Antworten (und ein angeregtes
Gespräch) haben wollt, kommt
bei unserem Poster vorbei!**

Agenda

- SwissGB4Science - ein Volltext Korpus für die Forschung
Eric Dubey, Matteo Lorenzini, Martin Reisacher, Tim Rüdiger
- Synthetische Datensätze in der CLS
Daniel Kababgi, Emilie Sitter, Robin Martin Aust, Marie-Christine Boucher, Berenike Herrmann
- **Text+ LLM Service**
Alexander Steckel, Umut Basaran, Stefan Buddenbohm, Maik Wegener, Philipp Wieder
- Wissen, wen man fragt – Agentic RAG für Automatisches Question Answering in der Domäne deutscher Grammatik
Christian Lang, Ngoc Duyen Tanja Tu, Roman Schneider
- “Computer, was bedeutet ‘Tiki-Taka’?” Eine Studie zur Generierung von Definitionsparaphrasen für Bedeutungswörterbücher am Beispiel des DWDS
Alexander Geyken, Gregor Middell
- “fRAG Deine Daten doch selbst” – Potenziale des Einsatzes von Retrieval Augmented Generation für Forschungsdaten und Forschungsdateninfrastrukturen
Timm Lehmborg

Collections
Lexical
Resources
Editions
Infrastructure/
Operations



Text+

LLM Service: Text+ AI Assistant



AI as a Service

Text- and language-based humanities offer extensive use-cases for Large Language Models (LLMs). Through GWDG, an additional web service will be made available on the Text+ website providing access to open-source and custom fine-tuneable LLMs, including:

- (Meta) LLaMA
- Mixtral
- Qwen
- Codestral
- Chat-GPT

Interface



Base LLM selection, add, remove



Edit base LLM



Add custom LLM



Sources section in RAG

Features

- Free use of various open-source models
- Users can create, edit and delete custom LLMs
- Collaborations section allows users to invite collaborators to chat with the custom LLMs
- Retrieval-augmented generation (RAG) on personal documents (currently .pdf, .txt, .csv)
- Sources section on generated answers for users to check and enable citations
- Compliance with legislative requirements and user privacy interests
- No data leakage
- With Open AI's Chat-GPT, no single user related data is externally transmitted, as the current implementation makes all users appear as one
- Currently, the assistant is available for project participants who log in via Academic Cloud

The AI Assistant is based on GWDG's LLM Service whose architecture and features are described in detail in these articles:

- Decker, Hossain, Meisel: "The New GWDG LLM Service", [DFWG Nachrichten 03/2024](#)
- M. Azampour: "LLM Services and the AI Act of the EU", [DFWG Nachrichten 04-05/2024](#)

A plethora of useful information on AI adoption can be sourced in this white paper:

- <https://docs.kispi.de/white-paper/Barriers-to-AI-adoption-B-report.pdf>

Use Cases

- Retrieval Augmented Generation (RAG)
- Entity Linking
- Enrichment of GermaNet
- Query generation for search strings
- Data preprocessing for NER model
- Historical normalizations
- Runtime environment

In this poster session, the respective task force should exhibit a poster specifically on this topic:



Feedback



We believe in agile development. As such, the Text+ AI Assistant is a first installment of a service to our user base. Currently, it does not meet all of our requirements, but shows the way to the (future) use cases.

We depend on you to deliver feedback.



Try it out, and let us know about your experiences and suggestions directly or via the **Text+ helpdesk!**



text-plus.org/ai-assistant

Agenda

- Synthetische Datensätze in der CLS
Daniel Kababgi, Emilie Sitter, Robin Martin Aust, Marie-Christine Boucher, Berenike Herrmann
- Text+ LLM Service
Alexander Steckel, Umut Basaran, Stefan Buddenbohm, Maik Wegener, Philipp Wieder
- **Wissen, wen man fragt – Agentic RAG für Automatisches Question Answering in der Domäne deutscher Grammatik**
Christian Lang, Ngoc Duyen Tanja Tu, Roman Schneider
- “Computer, was bedeutet ‘Tiki-Taka’?” Eine Studie zur Generierung von Definitionsparaphrasen für Bedeutungswörterbücher am Beispiel des DWDS
Alexander Geyken, Gregor Middell
- “fRAG Deine Daten doch selbst” – Potenziale des Einsatzes von Retrieval Augmented Generation für Forschungsdaten und Forschungsdateninfrastrukturen
Timm Lehmberg
- “Nun sag', wie hast du's mit den LLMs?” – Antworten der Text+ Community auf die Gretchenfrage
Stine Ziegler, Philippe Genêt

Collections
Lexical
Resources
Editions
Infrastructure/
Operations



WISSEN, WEN MAN FRAGT...



- **Ausgangslage:** Auch in der Domäne Grammatik mildert RAG Probleme wie Halluzinationen ab und verbessert die Antwortqualität. Dies gilt jedoch **nicht für alle Fragetypen** (vgl. Lang et al. 2024).
 - **Strategie:** Dynamische **Steuerung** der Antwortprozesse in Abhängigkeit des Fragetyps
- ↓
- **Umsetzung:** Prototyp einer Agentic RAG-Chain: Generierendes LLM fungiert u. a. als **Router** zur passgenauen Anbindung externer Ressourcen für unterschiedliche Fragetypen.

Agenda

- Synthetische Datensätze in der CLS
Daniel Kababgi, Emilie Sitter, Robin Martin Aust, Marie-Christine Boucher, Berenike Herrmann
- Text+ LLM Service
Alexander Steckel, Umut Basaran, Stefan Buddenbohm, Maik Wegener, Philipp Wieder
- Wissen, wen man fragt – Agentic RAG für Automatisches Question Answering in der Domäne deutscher Grammatik
Christian Lang, Ngoc Duyen Tanja Tu, Roman Schneider
- **“Computer, was bedeutet ‘Tiki-Taka’?” Eine Studie zur Generierung von Definitionsparaphrasen für Bedeutungswörterbücher am Beispiel des DWDS**
Alexander Geyken, Gregor Middell
- “fRAG Deine Daten doch selbst” – Potenziale des Einsatzes von Retrieval Augmented Generation für Forschungsdaten und Forschungsdateninfrastrukturen
Timm Lehmberg
- “Nun sag', wie hast du's mit den LLMs?” – Antworten der Text+ Community auf die Gretchenfrage
Stine Ziegler, Philippe Genêt

Collections
Lexical
Resources
Editions
Infrastructure/
Operations



LLMs für Wörterbuchdefinitionen

Alexander Geyken und Gregor Middell (BBAW)



Aufgabe: definiere *Mausohr*

1. eine Pflanze mit herzförmigen Blättern, die an die Form von Mäuseohren erinnern. (ChatGPT)



Damit zusammenhängende Fragen:

- Halluzination oder nicht?
- Erkennen das die Nutzerinnen und Nutzer?
- Gibt es leichtere oder schwerere Wörter für LLMs (Kriterien)?
- Prompting Strategien und weitere Optimierungen für LLMs?

↪ Poster: "Computer, was bedeutet Tiki-Taka?"

Agenda

- Synthetische Datensätze in der CLS
Daniel Kababgi, Emilie Sitter, Robin Martin Aust, Marie-Christine Boucher, Berenike Herrmann
- Text+ LLM Service
Alexander Steckel, Umut Basaran, Stefan Buddenbohm, Maik Wegener, Philipp Wieder
- Wissen, wen man fragt – Agentic RAG für Automatisches Question Answering in der Domäne deutscher Grammatik
Christian Lang, Ngoc Duyen Tanja Tu, Roman Schneider
- “Computer, was bedeutet ‘Tiki-Taka’?” Eine Studie zur Generierung von Definitionsparaphrasen für Bedeutungswörterbücher am Beispiel des DWDS
Alexander Geyken, Gregor Middell
- **“fRAG Deine Daten doch selbst” – Potenziale des Einsatzes von Retrieval Augmented Generation für Forschungsdaten und Forschungsdateninfrastrukturen**
Timm Lehmberg
- “Nun sag', wie hast du's mit den LLMs?” – Antworten der Text+ Community auf die Gretchenfrage
Stine Ziegler, Philippe Genêt

Collections
Lexical
Resources
Editions
Infrastructure/
Operations



A surreal illustration of a chaotic office. In the center, a man with a white, featureless face and a blue suit with a red tie stands with his hands raised in a gesture of surprise or exasperation. He is surrounded by a vast sea of papers, some floating in the air and others covering the floor. Numerous small, white, alien-like creatures with large black eyes and thin limbs are scattered throughout the scene, some sitting on desks, some on shelves, and some on the floor. The office is filled with wooden desks, filing cabinets, and shelves, all cluttered with papers and folders. The lighting is warm, with yellow pendant lamps hanging from the ceiling. The overall atmosphere is one of overwhelming information and chaos.

fRAG Deine Daten
doch selbst!

Agenda

- Synthetische Datensätze in der CLS
Daniel Kababgi, Emilie Sitter, Robin Martin Aust, Marie-Christine Boucher, Berenike Herrmann
- Text+ LLM Service
Alexander Steckel, Umut Basaran, Stefan Buddenbohm, Maik Wegener, Philipp Wieder
- Wissen, wen man fragt – Agentic RAG für Automatisches Question Answering in der Domäne deutscher Grammatik
Christian Lang, Ngoc Duyen Tanja Tu, Roman Schneider
- “Computer, was bedeutet ‘Tiki-Taka’?” Eine Studie zur Generierung von Definitionsparaphrasen für Bedeutungswörterbücher am Beispiel des DWDS
Alexander Geyken, Gregor Middell
- “fRAG Deine Daten doch selbst” – Potenziale des Einsatzes von Retrieval Augmented Generation für Forschungsdaten und Forschungsdateninfrastrukturen
Timm Lehmborg
- “Nun sag', wie hast du's mit den LLMs?” – Antworten der Text+ Community auf die Gretchenfrage
Stine Ziegler, Philippe Genêt

Collections
Lexical
Resources
Editions
Infrastructure/
Operations

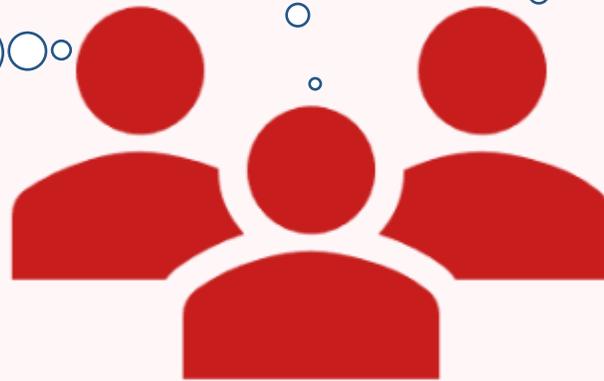


„Nun sag', wie hast du's mit denn LLMs?“

Ich wünsche mir ...

Mich stört ...

Ich nutze große Sprachmodelle für ...



Antworten der Text+ Community auf die Gretchenfrage



Vielen Dank an alle
Vortragenden und Zuhörenden!

Die Postersession findet ab
sofort bis 18:00 Uhr in O 138
(Fuchs-Petrolub-Saal) statt.

Collections
Lexical
Resources
Editions
Infrastructure/
Operations

Gefördert durch

DFG Deutsche
Forschungsgemeinschaft

Förderungsnummer 460033370

Teil der

nfdi Nationale
Forschungsdaten
Infrastruktur

Die vorliegende Präsentation wurde im Rahmen des Konsortiums Text+ im Kontext der Arbeit des Vereins Nationale Forschungsdateninfrastruktur (NFDI) e.V. verfasst. NFDI wird von der Bundesrepublik Deutschland und den 16 Bundesländern finanziert, und das Konsortium Text+ wird gefördert durch die Deutsche Forschungsgemeinschaft (DFG) – Projektnummer 460033370. Die Autor:innen bedanken sich für die Förderung sowie Unterstützung. Ein Dank geht außerdem an alle Einrichtungen und Akteur:innen, die sich für den Verein und dessen Ziele engagieren.

<https://www.text-plus.org>

