

DISCUSSION SESSION

Streamlining Data Publication: Automatic Metadata and Large Datasets in the Age of AI



Leibniz ScienceCampus
Digital Transformation
of Research

Anna Jacyszyn¹, Felix Bach¹, Tobias Kerzenmacher², Mahsa Vafaie¹

¹FIZ Karlsruhe - Leibniz Institute for Information Infrastructure

²KIT Institute of Meteorology and Climate Research Atmospheric Trace Gases and Remote Sensing (IMK-ASF)



Open Science Conference, 8-9 October 2025, Hamburg

Welcome

Session organisers



Anna Jacyszyn¹

DiTraRe coordinator,
dimension *Exploration
and Knowledge
Organisation*



Mahsa Vafaie^{1,3}

DiTraRe dimension
*Exploration and
Knowledge
Organisation*



**Tobias
Kerzenmacher²**

DiTraRe use case
*Publication of Large
Datasets*



Felix Bach¹

DiTraRe coordinator,
dimension *Tools and
Processes*



Kerstin Soltau¹

RADAR Product
Manager



**Stefan
Hofmann¹**

RADAR Full Stack
developer

¹FIZ Karlsruhe - Leibniz Institute for Information Infrastructure

²KIT Institute of Meteorology and Climate Research Atmospheric Trace Gases and Remote Sensing (IMKASF)

³KIT Institute of Applied Informatics and Formal Description Methods (AIFB)

Agenda

- **Ice-breaker** session
- Flash talks: **experts from different disciplines**
- **Fishbowl** discussion
- Workshop **summary**

Take notes with us and be part of the
session proceedings!



zbw.to/osc25-pad02

Ice-breaker session

Join at
slido.com
#2264 787



Flash Presentations

DiTraRe

Leibniz Science Campus *Digital Transformation of Research*

- Growth core to **establish new research branch.**



+



- Planned as a **4+4 years** project (start: September 2023).
- Funded by the Leibniz Association + FIZ KA + KIT.
- Analyse the process of **digitalisation of research.**
- **Multilevel interdisciplinary approach.** → We start with **4 specific use cases.**

DiTraRe Use Cases



1

Sensitive Data
in Sports
Science

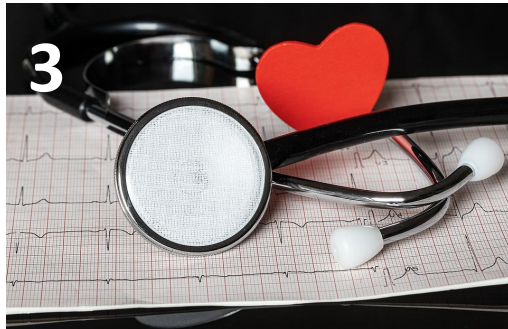
KIT Institute of Sports and Sports Science



2

Chemotion
Electronic Lab
Notebook

KIT Institute of Biological and Chemical Systems



3

AI in
Biomedical
Engineering

KIT Institute of Biomedical Engineering



4

Publication of
Large Datasets

KIT Institutes of Meteorology and Climate Research

DiTraRe dimensions

A. Reflection and Resonance

a dialogue between research and society, interactive process

B. Exploration and Knowledge Organisation

applied AI: represent, organise, and manage domain specific and procedural knowledge

C. Legal and Ethical Challenges

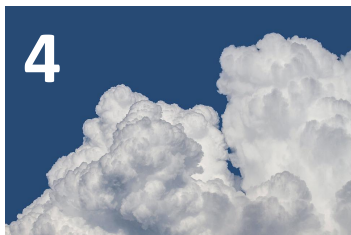
data ethics, data protection, copyright and data law

D. Tools and Processes

digital tools tailored precisely to the needs of researchers



This session



Use case
*Publication of Large
Datasets*

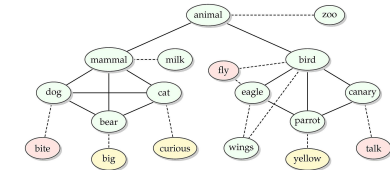


Dimension *Exploration
and Knowledge
Organisation*



Dimension *Tools and
Processes*

Discussion session
Streamlining Data
Publication: Automatic
Metadata and Large
Datasets in the Age of AI



4. Looking Forward – Role of AI

It's not just about making data open – it's about making data usable, FAIR, and meaningful for re-use.

DiTraRe dimension: Exploration & Knowledge Org.

FIZ Karlsruhe – Leibniz Institute for Information Infrastructure



Online collection
Wiedergutmachung for National Socialist Injustice

Wiedergutmachung is the term used to describe the German government's efforts to take responsibility and make amends for crimes committed by the National Socialist regime. In material terms, this encompasses the return of stolen property, compensation payments and assistance measures.

<https://www.archivportal-d.de/themenportale/wiedergutmachung>



The Bundeszentralkartei (BZK) Collection

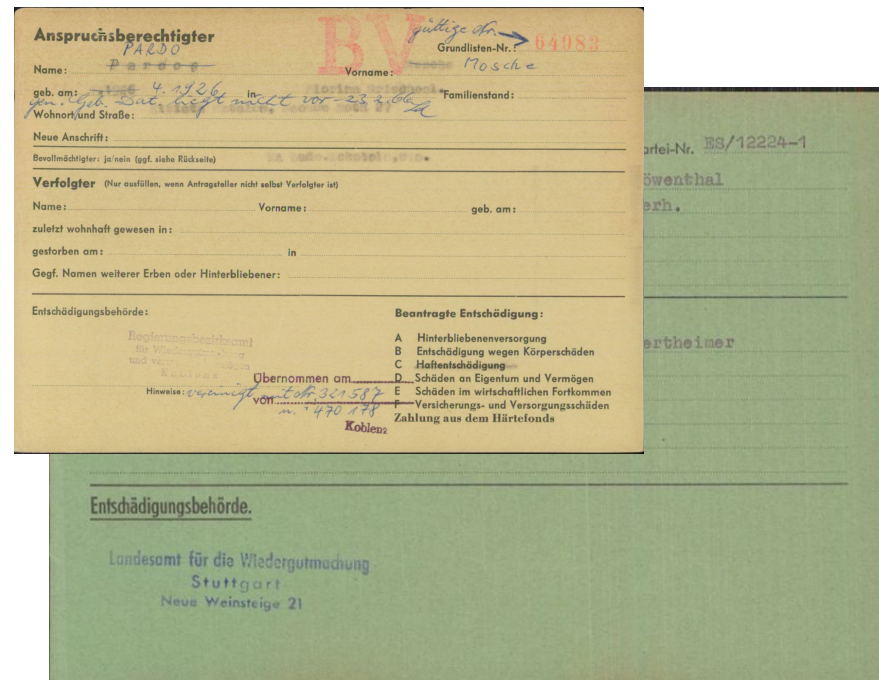
- Central Federal Index for Compensation of National Socialist injustices
- Central card file of most of applications for compensation in the Federal Republic of Germany
- about 1.9 million cards with basic information on applicants, persecuted persons and proceedings
- Kept from 1950s until today
- Basis for the **person search in Online Collection Wiedergutmachung**

<u>Anspruchsberechtigter.</u>		Kartei-Nr. <u>BS/12224-1</u>
Name: <u>Friedrich,</u>	geborene: <u>Löwenthal</u>	
Vorname: <u>Rosa</u>	Familienstand: <u>verh.</u>	
geb. am: <u>26.5.98</u>	in: <u>ESlingen</u>	
wohnhaft in (Ort, Straße, Land): <u>ESlingen (Neckar), Neckarstr. 85</u>		
<u>Tochter d. Erblasserin</u>		
<u>Verfolgter.</u> (Nur ausfüllen, wenn Anspruchsberechtigter nicht selbst Verfolgter ist.)		
Name: <u>Löwenthal,</u>	geborene: <u>Wertheimer</u>	
Vorname: <u>Jette</u>		
geb. am: <u>22.8.73</u>	in: <u>Bauerbach (Baden)</u>	
gest. am: <u>31.12.43</u>	in: <u>f.tot erklärt.</u>	
Namen weiterer Erben oder Hinterbliebener:		
<u>Entschädigungsbehörde.</u>		
Landesamt für die Wiedergutmachung Stuttgart Neue Weinsteige 21		

Source: Landesarchiv NRW – Abteilung Rheinland – BR 3015 ZK-Nr. 64083, 190667, 15932, scho8, 67800/II/6095

The Bundeszentralkartei (BZK) Collection

- Central Federal Index for Compensation of National Socialist injustices
- Central card file of most of applications for compensation in the Federal Republic of Germany
- about 1.9 million cards with basic information on applicants, persecuted persons and proceedings
- Kept from 1950s until today
- Basis for the **person search in Online Collection Wiedergutmachung**

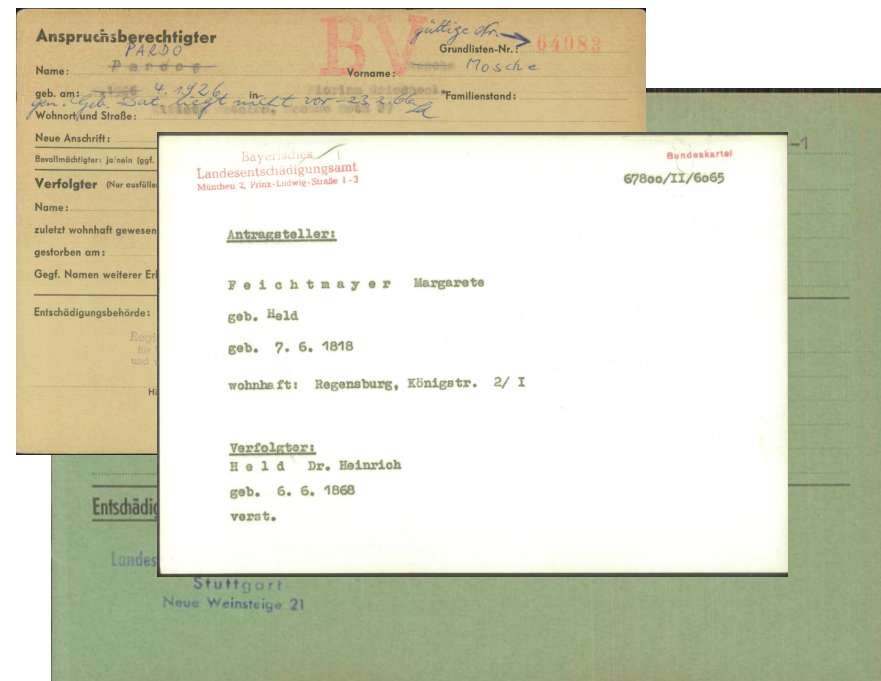


Anspruchsberechtigter
PA 600
Name: P... Vorname: Mosche
geb. am: 4. 12. 1926 in: 1912 Grundlisten-Nr.: 64083
Wohnort und Straße: ... Familienstand:
Neue Anschrift:
Bevollmächtigter: ja/nein (ggf. siehe Rückseite) ...
Verfolger (Nur ausfüllen, wenn Antragsteller nicht selbst Verfolger ist)
Name: ... Vorname: ... geb. am: ...
zuletzt wohnhaft gewesen in: ...
gestorben am: ... in: ...
Gegf. Namen weiterer Erben oder Hinterbliebener: ...
Entschädigungsbehörde: ...
Beantragte Entschädigung:
A Hinterbliebenenversorgung
B Entschädigung wegen Körperschäden
C Haftentschädigung
D Schäden an Eigentum und Vermögen
E Schäden im wirtschaftlichen Fortkommen
F Versicherungs- und Versorgungsschäden
Zahlung aus dem Härtefonds
Übernommen am: 30. 3. 1982
von: 470 178
Koblenz
Entschädigungsbehörde:
Landesamt für die Wiedergutmachung
Stuttgart
Neue Weinsteige 21

Source: Landesarchiv NRW – Abteilung Rheinland – BR 3015 ZK-Nr. 64083, 190667, 15932, scho8, 67800/II/6095

The Bundeszentralkartei (BZK) Collection

- Central Federal Index for Compensation of National Socialist injustices
- Central card file of most of applications for compensation in the Federal Republic of Germany
- about 1.9 million cards with basic information on applicants, persecuted persons and proceedings
- Kept from 1950s until today
- Basis for the **person search in Online Collection Wiedergutmachung**



Anspruchsberechtigter
PA 600
Name: P...
geb. am: 4. 11. 1926
Wohnort und Straße: ...
Neue Anschrift: ...
Bevollmächtigter: ja/nein (ggf.)
Verfolger (Bitte ausfüllen)
Name: ...
zuletzt wohnhaft gewesen
gestorben am: ...
Gegf. Namen weiterer Erf...
Entschädigungsbehörde: ...

Verfolger
Name: ...
zuletzt wohnhaft gewesen
gestorben am: ...
Gegf. Namen weiterer Erf...
Entschädigungsbehörde: ...

Antragsteller:
Feichtmayer Margarete
geb. Held
geb. 7. 6. 1918
wohnhaft: Regensburg, Königstr. 2/ I

Verfolgter:
Held Dr. Heinrich
geb. 6. 6. 1908
verst.

Landesentschädigungsamt
München 2, Prinz-Luitwig-Strasse 1-3
Stuttgart
Neue Weinsteige 21

Bundeskartei
67800/II/6065

Source: Landesarchiv NRW – Abteilung Rheinland – BR 3015 ZK-Nr. 64083, 190667, 15932, scho8, 67800/II/6095

The Bundeszentralkartei (BZK) Collection

- Central Federal Index for Compensation of National Socialist injustices
- Central card file of most of applications for compensation in the Federal Republic of Germany
- about 1.9 million cards with basic information on applicants, persecuted persons and proceedings
- Kept from 1950s until today
- Basis for the **person search in Online Collection Wiedergutmachung**

The image shows three overlapping forms from the Bundeszentralkartei (BZK) collection. The top form is a 'Verfolgter' (Persecuted) card for Margarete Feichtmayer, with handwritten details including birth date (26.7.86), place of birth (Braunschweig), and address (Hamburg). The middle form is a 'Verfolgter' card for Helga Schreiber, with handwritten details including birth date (26.7.86), place of birth (Braunschweig), and address (Hamburg). The bottom form is a 'Verfolgter' card for Martha, Luise, and Christine Bräunschweig, with handwritten details including birth date (26.7.86), place of birth (Braunschweig), and address (Hamburg). The forms are filled out with handwritten information in black ink.

Source: Landesarchiv NRW – Abteilung Rheinland – BR 3015 ZK-Nr. 64083, 190667, 15932, scho8, 67800/II/6095

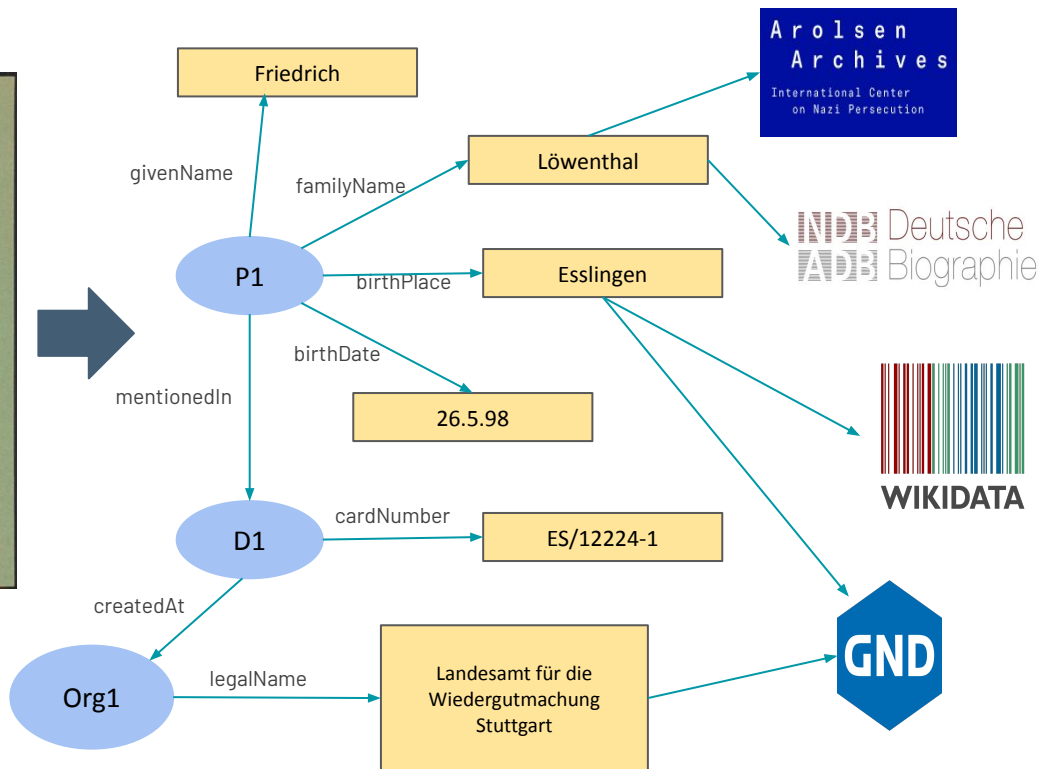
The BZK Knowledge Graph

Anspruchsberechtigter. Kartei-Nr. ES/12224-1

Name: Friedrich, geborene: Löwenthal
 Vorname: Rosa, Familienstand: verh.
 geb. am: 26.5.98 in: Esslingen
 wohnhaft in (Ort, Straße, Land): Esslingen (Neckar), Neckarstr. 85
Tochter d. Erblasserin

Verfolgter. (Nur ausfüllen, wenn Anspruchsberechtigter nicht selbst Verfolgter ist.)
 Name: Löwenthal, geborene: Wertheimer
 Vorname: Jette
 geb. am: 22.8.73 in: Bauerbach (Baden)
 gest. am: 31.12.43 in: f. tot erklärt.
 Namen weiterer Erben oder Hinterbliebenen:

Entschädigungsbehörde.
 Landesamt für die Wiedergutmachung
 Stuttgart
 Neue Weinsteige 21



“Traditional” Information Extraction

Rule-based IE



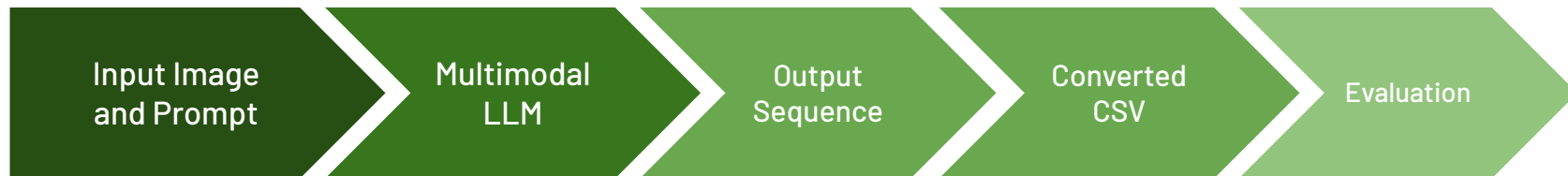
NER-based IE



Limitations

- ❑ **Error Propagation**
 - ❑ Text Recognition
 - ❑ Layout Analysis
- ❑ **Template-based**
 - ❑ Generalisability
 - ❑ Laborious with >40 layout types

End-to-end MLLM-based Information Extraction



Anspruchsberechtigter
Name Wissel Geburtsdatum 19.6.67 Geburtsort Karlsruhe
Vatername P. Fischer Familienname Wissel
Wohnort u. Straße Alt. Am Oben Schloß 14a Postfach 100 18
Verfänger Wie soll sein, sein Anspruchsberechtigter soll selber Verfänger sein
Name Wissel Vatersname Fischer / Fischer geb. Fischer
gebürtlich in Wormsheim St. 16
Cogn. Nomes weiterer Eltern: ?
Entscheidungsbeförder:
Freie und Hausmutter Hamburg
Staatshilfe
bei der Wohnungsgemeinschaft
Hamburg 1, Alsterdorf Straße 9
(Postfach)
DOPPELMELDUNG
Stufen
gemeldet u. am: 6. April 1958

<image>\nPlease provide the following information as you can see on the image as a Python dictionary [schema]...

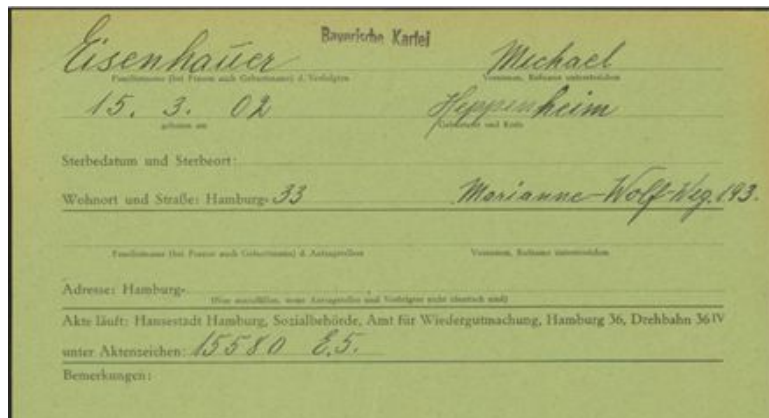
Open-source MLLM

```
{ "full_response": "python\n\nCompensationOffice1" :  
"Nordrhein-Westfalen",\n...
```

```
BZKNr: 1906 67
ApplicantFirstName: Pinkas
ApplicantLastName: Riesel
ApplicantBirthDate: 19.6.67
ApplicantBirthPlace: Kossov
VictimtFirstName:
Jakob/Jerzb
VictimLastName: Riesel
VictimDeathPlace: Warschau
```

Source: Landesarchiv NRW – Abteilung Rheinland – BR 3015 ZK-Nr. 190667

BZKOpen dataset and data schema



Keys for data extraction:

CompensationOffice1, BZKNr, ApplicantFirstName,
ApplicantLastName, ApplicantAltFirstName, ApplicantBirthName,
ApplicantAltLastName, ApplicantBirthDate, ApplicantBirthPlace,
ApplicantCurrentAddress, VictimFirstName, VictimLastName,
VictimAltFirstName, VictimBirthName, VictimAltLastName,
VictimBirthDate, VictimBirthPlace, VictimDeathDate,
VictimDeathPlace

Dataset Viewer

Auto-converted to Parquet

</> API

Subset (2)
normalized · 516 rows

Split (3)
train · 361 rows

Q

Search this dataset

image image · width (px)	CompensationOffice1 string · lengths	BZKNr string · lengths	Layout class string · classes	ApplicantFirstName string · lengths
<div><div><div></div><div></div></div><div>2.46k2.54k</div></div>	<div><div><div></div><div></div></div><div>0116</div></div>	<div><div><div></div><div></div></div><div>021</div></div>	<div><div><div></div><div></div></div><div>16 values</div></div>	<div><div><div></div><div></div></div><div>019</div></div>
<div><div><div></div></div></div>	Bayerisches Landesentschädigungsamt	BEG 069 657	BY-BE-Hauptphase	Etz
<div><div><div></div></div></div>	Entschädigungsamt Berlin	363801	BY-BE-Hauptphase	Cirl
<div><div><div></div></div></div>	Hess. Staatsministerium Der Minister des Inner...	12435	HE-Frühe-Phase	
<div><div><div></div></div></div>	Karlsruhe	EK 9048 DP 246	BY-HB-Frühe Phase	

<https://huggingface.co/datasets/MahsaVafaie/BZKopen>

Source: Landesarchiv NRW – Abteilung Rheinland – BR 3015 ZK-Nr. 15580 E.5

Information Extraction with MLLMs

Models

Donut
GPT-4o-mini
InternVL2
InternVL2_5

Metrics

Normalised Edit
Distance (NED)
Exact Matches (EM)
Partial Matches (PM)

Model	Size	NED	EM (t=0)	PM (t=1)	PM (t=3)
Donut-base	-	0.415	56%	57%	58%
Donut-base-finetuned	-	0.358	59%	61%	64%
GPT-4o-mini	-	0.184	72%	76%	79%
InternVL2.0	8B	0.382	53%	55%	60%
InternVL2.0	26B	0.431	48%	52%	57%
InternVL2.0	40B	0.158	76%	79%	83%
InternVL2.0-Llama3	76B	0.286	64%	67%	70%
InternVL2.0-finetuned	40B	0.173	74%	77%	81%
InternVL2.5	8B	0.340	57%	60%	64%
InternVL2.5	26B	0.311	60%	64%	69%
InternVL2.5	38B	0.080	83%	88%	91%
InternVL2.5	78B	0.139	77%	82%	84%
InternVL2.5-finetuned	38B	0.117	79%	84%	86%

Table 1. Performance of different transformer-based models on the BZKOpen dataset, including both pre-trained and fine-tuned variants. For fine-tuned models, the BZKOpen train set was used.

Information Extraction with MLLMs

Model	Size	Prompting Strategy	NED	EM (t=0)	PM (t=1)	PM (t=3)
GPT-4o-mini	-	ZS	0.184	72%	76%	79%
GPT-4o-mini	-	1FS	0.093	83%	87%	90%
GPT-4o-mini	-	2FS	0.080	84%	88%	91%
GPT-4o-mini	-	5FS	0.074	85%	89%	92%
InternVL2.0	40B	ZS	0.158	76%	79%	83%
InternVL2.0	40B	1FS	0.160	76%	79%	82%
InternVL2.0	40B	2FS	0.134	78%	82%	84%
InternVL2.0	40B	5FS	0.117	79%	82%	85%
InternVL2.5	38B	ZS	0.080	83%	88%	91%
InternVL2.5	38B	1FS	0.070	84%	89%	92%
InternVL2.5	38B	2FS	0.060	86%	90%	93%
InternVL2.5	38B	5FS	0.062	86%	90%	92%

Table 2. Performance of zero-shot (ZS) and few-shot (FS) prompting with different numbers of shots and different prompts for the InternVL2.5-38B and GPT-4o-mini models.

Takeaways

	Classical ML/DL	Generative AI (LLMs)
Efficiency		✓
Adaptability		✓
Reusability		✓
Reliability	✓	
Reproducibility	✓	
Data Privacy	✓	✓*

DiTraRe dimension: Tools and Processes

FIZ Karlsruhe – Leibniz Institute for Information Infrastructure

- About me: Felix Bach, Head of Research Data Department at FIZ Karlsruhe
 - Co-Spokesperson NFDI4Chem,
 - Scientific Coordinator of DiTraRe and
 - PI for dimension: Tools and Processes
- Focus of Dimension T&P:
 - Develop, test, and integrate digital tools and infrastructures for trusted workflows
 - Bridge between technical infrastructures (repositories, ELNs, AI services) and disciplinary practices
- Tools in place:
 - RADAR – repository service for research data, handling large and complex datasets
- Achievements so far:
 - Integration of ELNs (e.g., Chemotion) and repositories (e.g., RADAR4Chem) into FAIR workflows
 - First implementations of AI-supported metadata curation and interoperability solutions
- Future directions / in progress:
 - Efficient big data access (UL, DL, single file access within large data packages, viewer integration)
 - Distributed dataset strategies and interfaces
 - Exploitation of disciplinary metadata/terminologies
 - AI-based automatic FAIRness assessment and automatic metadata generation/enhancement

DiTraRe dimension: Tools and Processes

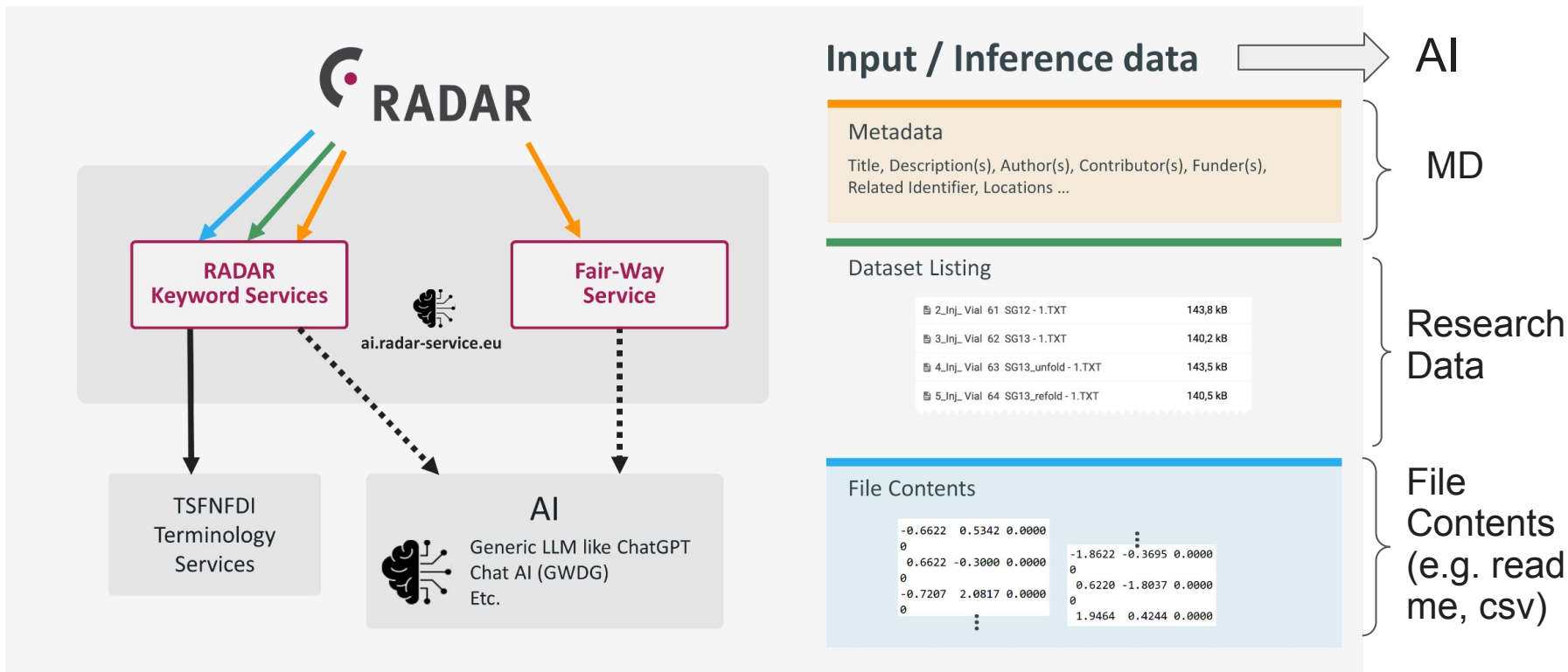
FIZ Karlsruhe – Leibniz Institute for Information Infrastructure

RADAR adaption relevant for the use case “Publication of large Datasets”

- Integration of terminology services of Base4NFDI/TIB Hannover (TS4NFDI)
 - terminologies and ontologies for earth sciences
 - usable in keywords (type ahead, suggestion list pops up)
 - integrated widgets in metadata annotation
 - future work: exploit data file contents (e.g. NetCDF files contain reusable MD)
- Initial experimental implementation of automatic generation of metadata by LLMs/SLMs based on available text documents (project proposals, papers), title, abstract...
- gitHub and gitLab integration - data and SW import
- Optimisation of Big Data handling
 - WebDav ingest additionally to web-based ingest (failure-intolerant for long transfers)
 - Additional disc copy for random access (DL of selected files/folders in a data package)

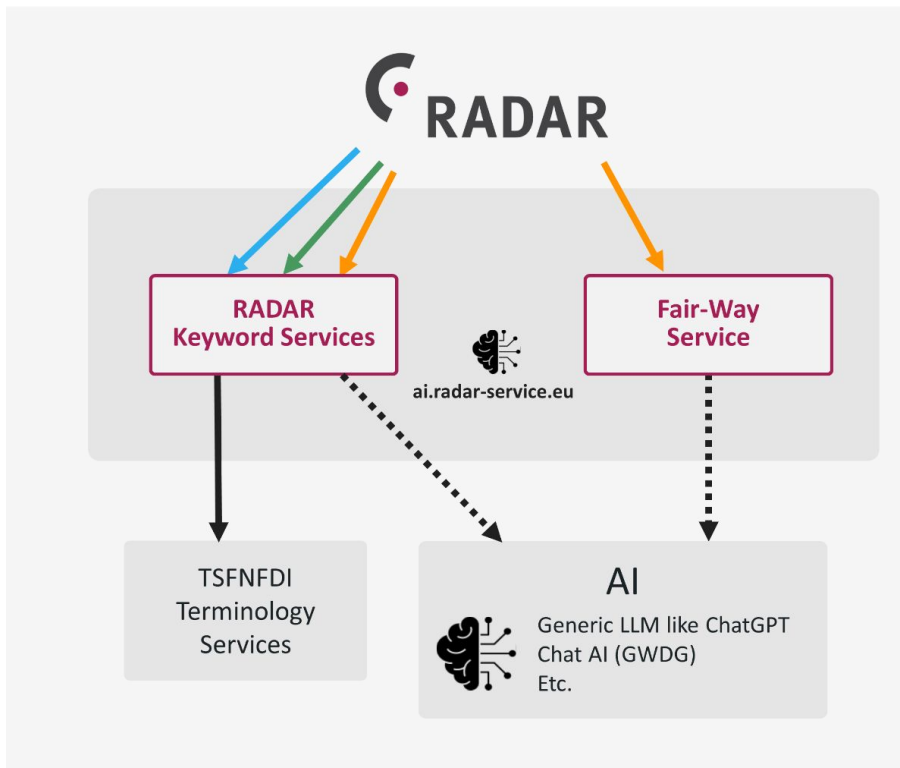
Enhancing FAIR Research Data Management with AI Support

a) Metadata enhancement



Enhancing FAIR Research Data Management with AI Support

b) FAIRness enhancement



FAIRness in RADAR

data must be published with rich MD!

LLM & Prompt (e.g. ChatGPT)

- (costs up to 0,20€ / request)
- results of mixed quality
- Privacy issues

FAIR-Way Service

- Open Source
- FAIRsFAIR metric (standard)
- good results

*F-UJI best tool (very good metrics) but not the quickest



Fishbowl Discussion

Fishbowl discussion



Take notes with us:
zbw.to/osc25-pad02

Thank you for attending!

- **Symposium** on Digitalisation of Research: 2-3 December 2025, Karlsruhe
 - Interdisciplinary panels, discussions with experts!
- **Interdisciplinary Colloquium** on Digitalisation of Research
 - Monthly presentations of experts, join us online or in person!
- Stay **connected** with DiTraRe
 - Website: www.ditrare.de/en
 - Email: ditrare@fiz-karlsruhe.de
 - LinkedIn: www.linkedin.com/company/ditrare
 - Mastodon: sigmoid.social/@DiTraRe
 - YouTube: www.youtube.com/@DiTraRe
 - Zenodo: zenodo.org/communities/ditrare



www.ditrare.de/en





New Publication: CIKM'25

End-to-end Information Extraction from Archival Records with Multimodal Large Language Models

Mahsa Vafaie
mahsa.vafaie@fiz-karlsruhe.de
FIZ Karlsruhe – Leibniz Institute for
Information Infrastructure,
Hermann-von-Helmholtz-Platz 1
76344, Eggenstein-Leopoldshafen
Germany

Sven Hertling
sven.hertling@fiz-karlsruhe.de
Data and Web Science Group
University of Mannheim, Germany
FIZ Karlsruhe – Leibniz Institute for
Information Infrastructure,
Hermann-von-Helmholtz-Platz 1
76344, Eggenstein-Leopoldshafen
Germany

Inger Banse-Strobel
inger.banse-strobel@bundesarchiv.de
Bundesarchiv, Potsdamer Str. 1
56075 Koblenz, Germany

Kevin Dubout
k.dubout@bundesarchiv.de
Bundesarchiv, Potsdamer Str. 1
56075 Koblenz, Germany

Harald Sack
harald.sack@fiz-karlsruhe.de
FIZ Karlsruhe – Leibniz Institute for
Information Infrastructure,
Hermann-von-Helmholtz-Platz 1
76344, Eggenstein-Leopoldshafen
Germany

ABSTRACT

Semi-structured Document Understanding presents a challenging research task due to the significant variations in layout, style, font, and content of documents. This complexity is further amplified when dealing with *born-analogue* historical documents, such as digitised archival records, which contain degraded print, handwritten annotations, stamps, marginalia and inconsistent formatting resulting from historical production and digitisation processes. Traditional approaches for extracting information from semi-structured documents rely on manual labour, making them costly and inefficient. This is partly due to the fact that within document collections, there are various layout types, each requiring customised optimisation to account for structural differences, which substantially increases the effort needed to achieve consistent quality. The emergence of Multimodal Large Language Models (MLLMs) has significantly advanced Document Understanding by enabling flexible, prompt-based understanding of document images, needless of OCR outputs or layout encodings. Moreover, the encoder-decoder architectures have overcome the limitations of encoder-only models, such as reliance on annotated datasets and fixed input lengths. However, there still remains a gap in effectively applying these models in real-world scenarios. To address this gap, we first introduce BZKOpen, a new annotated dataset designed for key information extraction from historical German index cards. Furthermore,

we systematically assess the capabilities of several state-of-the-art MLLMs—including the open-source InternVL2.0 and InternVL2.5 series, and the commercial GPT-4o-mini—on the task of extracting key information from these archival documents. Both zero-shot and few-shot prompting strategies are evaluated across different model configurations to identify the optimal conditions for performance. Interestingly, our results reveal that increasing model size does not necessarily lead to better performance on this dataset. Among all models tested, the open-source InternVL2.5-38B consistently achieves the most robust results, outperforming both larger InternVL models and the proprietary alternative. We further provide practical insights into prompt engineering and inference settings, offering guidance for applying MLLMs to real-world key information extraction tasks. Additionally, we highlight the need for more ground truth datasets that include a wider range of historical documents with varying quality and in multiple languages, in order to fully explore the potentials and limitations of MLLMs for key information extraction from historical records.

CCS CONCEPTS

• Information systems → Information extraction.

KEYWORDS

Multimodal Large Language Models, Document Understanding, Key Information Extraction, Digital Cultural Heritage

ACM Reference Format:

Mahsa Vafaie, Sven Hertling, Inger Banse-Strobel, Kevin Dubout, and Harald Sack. 2025. End-to-end Information Extraction from Archival Records with Multimodal Large Language Models. In *Proceedings of The 34th ACM International Conference on Information and Knowledge Management (CIKM '25)*. ACM, New York, NY, USA, 9 pages. <https://doi.org/XXXXXX.XXXXXXXX>

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.
CIKM '25, November 10–16, 2025, Seoul, South Korea.

© 2025 Copyright held by the owner/authors. Publication rights licensed to ACM.
ACM ISBN 978-1-xxxx-xxxx-x/YY/MM
<https://doi.org/XXXXXX.XXXXXXXX>