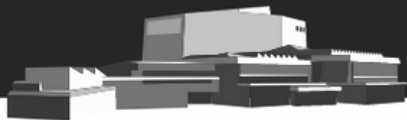




Staatsbibliothek  
zu Berlin  
Preußischer Kulturbesitz



# Beyond word clouds

## Practical applications in challenging cultural contexts

Jana Götze | Staatsbibliothek zu Berlin - Preußischer Kulturbesitz  
CLASP Seminar | 10 June 2026  
credit for many slides: Clemens Neudecker  
[jana.goetze@sbb.spk-berlin.de](mailto:jana.goetze@sbb.spk-berlin.de)

# Let's start with a quiz

Fill the blank!

Berlin is the capital of  
\_\_\_\_\_

FRG / Germany (1990 - now)

GDR / East Germany (1949-1989)

Soviet Occupation Zone (1945-1949)

Third Reich (1933-1945)

Weimar Republic (1918-1933)

Kingdom of Prussia (1701-1918)

Margraviate of Brandenburg (1618-1701)

# Agenda

- I. Background
- II. NLP data, tasks and challenges
- III. Related ongoing work
- IV. Opportunities

# Staatsbibliothek zu Berlin – Berlin State Library (SBB, “Stabi”)

- is part of the **Prussian Cultural Heritage Foundation** - Stiftung Preußischer Kulturbesitz (SPK)
- is **open freely to the public** 7 days/week in two sites
- is collecting **scientific literature** in all languages and from all times and countries since 1661
- has a collection of ca. 12M books with an annual **growth** of approx. 100k titles
- **Digitized Collections** provide access to >240,000 digitized documents under Public Domain license
- **Stabi Lab** for experiments, events, datasets, digital humanities



**Staatsbibliothek  
zu Berlin**  
Preußischer Kulturbesitz



# What cultural heritage organisations do

## **Provide data and services to the user**

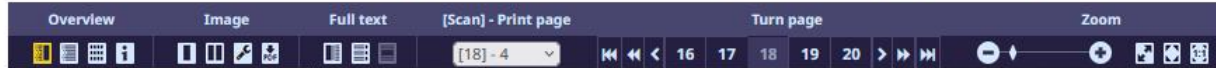
- Indexation, analysis and enhancement of digital data
- Create new services for users and scientists
- Growing (mostly open) data can serve as training data
- Expert knowledge about material and content creates quality data
- Sensitivity to data quality in creating, maintaining, and using data
- Transparency, data privacy and responsibility wrt AI are not just buzz words (public institution!)



GLAM

# What digital access looks like

Hermann, Paul: Frieda das Kind des Seiltänzers : eine Erzählung für die liebe Jugend ; mit 3 ... , 1877



## Compact Table of Contents

Page: [1]

### Rosen und Dornen

Page: [1]

#### Frieda das Kind des Seiltänzers

Page: [1]

#### Binding

Page: [6]

#### Engraved titlepage

Page: [9]

#### Title page

Page: [11]

#### Vorrede.

Page: [13]

#### Inhalt.

Page: [15] - 1

#### [Texte und Illustrationen]

— 4 —

auf eine bessere Versorgung zu machen. Er hatte den Nachtwächterposten auf dem Herrenhofe im Auge, den Diana, die zahnlose alte Wolfskünderin nicht mehr recht ordentlich versehen konnte. Dies Stückchen Gnadenbrod war ihm auch schon in nahe Aussicht gestellt, doch mußte er sich noch zuvor erst einer Prüfung unterwerfen und nachweisen, ob er wohl das Zeug für diesen wichtigen Posten hätte. Es hatten manche seiner Kameraden, verschiedentlich recht bunt und lebhaft in der Wollse gefärbt, und nicht so pechschwarz, zottelig wie er, um diese Anstellung sich längst beworben! Da mußte denn die Wahl

# What digital access looks like

Stabi

Digitalisierte Sammlungen

Hilfe DE | EN

## Simons, Gustav: Küchensünden und Volksgesundheit, 1905

Übersicht Bild Volltext [Scan] - Druckseite Blättern Zoom

— 76 —

wird oder aber man schiebt Geld, welches man im Inlande so gut gebrauchen könnte, dem Auslande in die Tasche.

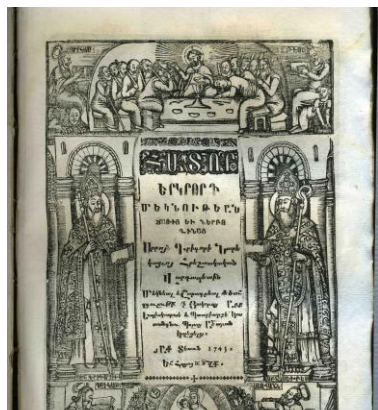
Gesetzt nun den Fall, ich kaufe nicht für diese Mark Zigarren, sondern z. B. Haselnüsse, so fördere ich erstens meine und der Meinigen Gesundheit und Schönheit, denn Nüsse sind gesund und erhalten das Gebiß infolge des aufzuwendenden Kaudruckes blank auch ohne Zahnbürsten, zumal wenn man mit den Kindern um die Wette noch als Erwachsener die Nüsse mit den Zähnen knackt. Man erspart zudem ein Industrieprodukt, die Butter, denn fettreiche Nüsse und Vollbrot sind allein schon eine vollwertige Mahlzeit. Man schädigt durch solchen Genuß auch keinen Volksgenossen, hat keine Verstopfung und braucht deshalb keinen Pillendreher und Pillenverschreiber. Während ich als Zigarrenkäufer einen Deutschen in die giftige Zigarrenfabrikluft hineindränge, locke ich durch Nußeinkauf einen Nußstaudenzüchter in Gottes freie Natur. Der Lebensreformer Ruf lautet doch „heraus aus der Bude, und nicht hinein in die Bude.“ Viele Nußkäufer machen durch die Befriedigung ihrer Bedürfnisse aus Deutschland einen Garten Eden, viele Nikotinschwelger locken einen Schornstein aus der Muttererde, einen sogenannten Industriespargel und machen höchstens die Amerikaner in Kuba oder sonst wo noch unverschämter. Gutes und Schlechtes liegt also in der Einkaufs-

76

wird oder aber man schiebt Geld, welches man im Inlande so gut gebrauchen könnte, dem Auslande in die Tasche.

Gesetzt nun den Fall, ich kaufe nicht für diese Mark Zigarren, sondern z. B. Haselnüsse, so fördere ich erstens meine und der Meinigen Gesundheit und Schönheit, denn Nüsse sind gesund und erhalten das Gebiß infolge des aufzuwendenden Kaudruckes blank auch ohne Zahnbürsten, zumal wenn man mit den Kindern um die Wette noch als Erwachsener die Nüsse mit den Zähnen knackt. Man erspart zudem ein Industrieprodukt, die Butter, denn fettreiche Nüsse und Vollbrot sind allein schon eine vollwertige Mahlzeit. Man schädigt durch solchen Genuß auch keinen Volksgenossen, hat keine Verstopfung und braucht deshalb keinen Pillendreher und Pillenverschreiber. Während ich als Zigarrenkäufer einen Deutschen in die giftige Zigarrenfabrikluft hineindränge, locke ich durch Nußeinkauf einen Nußstaudenzüchter in Gottes freie Natur. Der Lebensreformer Ruf lautet doch „heraus aus der Bude, und nicht hinein in die Bude.“ Viele Nußkäufer

# The variety of data is large



- occidental, oriental and East Asian manuscripts
- geographical maps, globes and atlases
- music scores
- autographs and bequests
- historical and modern printed matter
- 650,000+ newspaper issues
- 170,000+ photographs

- **historical data**
- **heterogeneous data**
- **a lot of data**

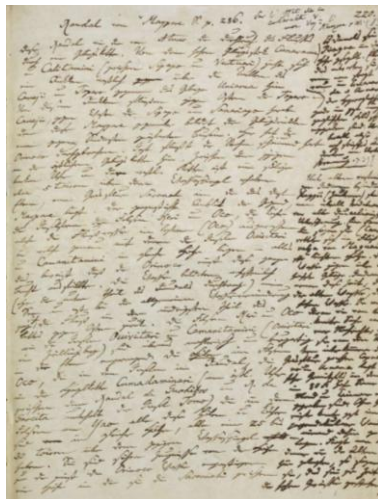


# Digitized collections



1477

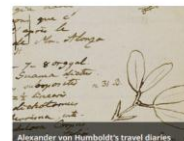
240,794 digitized objects



1800



The first European prints of the 19th century from the Stab's collections have already been widely digitized. More than 1,200 titles can be found in this database. Well-known and illustrated prints such as the Gutenberg Bible are currently on display in the Photo Gallery in the original - but also many other titles from the 19th century. Please check the Research Catalogue for more information.



Alexander von Humboldt's travel diaries



Orientalia from the Diet Library



E.T.A. Hoffmann



Type specimens from German-speaking foundries



Kriegsbilder



Type specimens from German-speaking foundries

Der Generaldirektor und leitende Mitarbeiter waren in zahlreichen bibliothekarischen und des Bibliothekswesens betreffenden Fachgremien und Organisationen der DDR teils leitend, teils aktiv mitarbeitend tätig und übten Funktionen in der International Federation of Library Associations (IFLA) und in der Association Internationale de Bibliothèques Musicales (AIBM) aus.

Beherrschendes Thema vieler Beratungen wie von Belegschaftsversammlungen waren die Rekonstruktion des Gebäudes sowie die Erweiterung der Magazine und der Lesesaalplätze. Im Mai 1976 begannen Instandsetzungsarbeiten an den Fassaden und zwar zunächst an der Front der Universitätsstraße. Sie wurden in der Folgezeit unkontinuierlich weitergeführt.

Die überaus angespannte Platzsituation in den Hauptmagazinen führte zu einer Reihe von Notmaßnahmen und zur Suche nach Räumlichkeiten für die Auslagerung von Buchbeständen. Nachdem im April 1975 die stählerne Dachkonstruktion über der Ruine des alten Kuppelensalaus aus Sicherheitsgründen gesprengt werden mußte und Abbruch- und Transportarbeiten einsetzten, die für viele der Mitarbeiter Beunruhigungen durch Lärm und Staub brachten, konnten schließlich grundlegende Entscheidungen getroffen werden. Die neue Variante, Ergebnis jahrelanger Diskussionen, Entwürfe und Kapazitätsberechnungen, sah die Errichtung von vier Büchertürmen mit einer Nutzfläche für etwa 2,4 Millionen Bände und eines Zwischenbaus vor sowie ein zusätzliches Funktionsgebäude außerhalb der Bibliothek. Der Baubeginn sollte sich aber noch weitere Jahre hinziehen!

## 1.2. Wissenschaftliches Leben

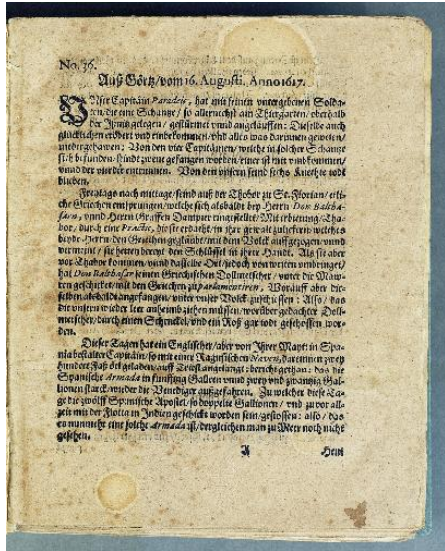
Mitarbeiter der DSB haben durch Forschungsarbeiten, fachwissenschaftliche und populärwissenschaftliche Veröffentlichungen sowie Vorträge bei Veranstaltungen oder auf Kongressen im In- und Ausland wesentlich zur Förderung des geistigen und kulturellen Lebens beigetragen und Arbeitsergebnisse und Bestände der Bibliothek für Fachkollegen und die Öffentlichkeit erschlossen.

Stellvertreter für viele dieser Beiträge seien die Mitarbeit an der Festschrift „Studien zur Buch- und Bibliotheksgeschichte, Hans Lüffing zum 70. Geburtstag“ erwähnt, das Handschrifteninventar „Der Nachlaß Hans Delbrück“, die Publikation „America in maps“ oder die Veröffentlichung von Band 7 der Konversationshefte Ludwig van Beethovens.

Weiten Wiederhall fand auch die Konferenz aus Anlaß des 75jährigen Bestehens der Arbeit am Gesamtkatalog der Wiegendrucke im November 1979. Fünf von den Referenten, die Forschungsergebnisse zur Buchdruck-, Wissenschafts- und Kulturgeschichte

2001

# Newspapers (ZEFYS portal)



1617



Berlinerische Nachrichten von Staats- und gelehrten Sachen  
1776



al-Mu'ayyad  
1897










Dziennik Berliński  
1939

# Stabi has 4 goals to serve the user (with ML methods)

- Provide **search and retrieval** across all collections, for text and image content
- Research and develop **open source technologies** for cultural heritage
- Provide collections as **open and machine-readable** digital data (“Collections as Data”)
- **Responsibly use AI** and curate digital data under consideration of problematic content based on ethical, legal and social criteria



# There is also contemporary material

-  **DETR-crowd is all you need** ▼  
Liu Weijia ; Zishen Zheng ; Ke Fan ; Kun He ; Taiqiu Huang ; Weijia Liu ; Xianlun Ke ; Yuming Xu  
In: Современные инновации, системы и технологии; 3(2023), 2 ☆  
[Online access](#)  🔖
- 
-  **Rotation is All You Need: Cross Dimensional Residual Interaction for Hyperspectral Image Classification** ▼  
Xin Qiao ; Swalpa Kumar Roy ; Weimin Huang  
In: IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing; 16(2023), Seite 5387-5404 ☆  
[Online access](#)  (via journal page) 🔖
- 
-  **Attention (to Virtuosity) Is All You Need: Religious Studies Pedagogy and Generative AI** ▼  
Jonathan Barlow ; Lynn Holt  
In: Religions; 15(2024), 9, p 1059 ☆  
[Online access options](#) 🔖
- 
-  **Cross attention is all you need: relational remote sensing change detection with transformer** ▼  
Kaixuan Lu ; Xiao Huang ; Ruiheng Xia ; Pan Zhang ; Junping Shen  
In: GIScience & Remote Sensing; 61(2024), 1 ☆  
[Online access](#)  (via issue) 🔖

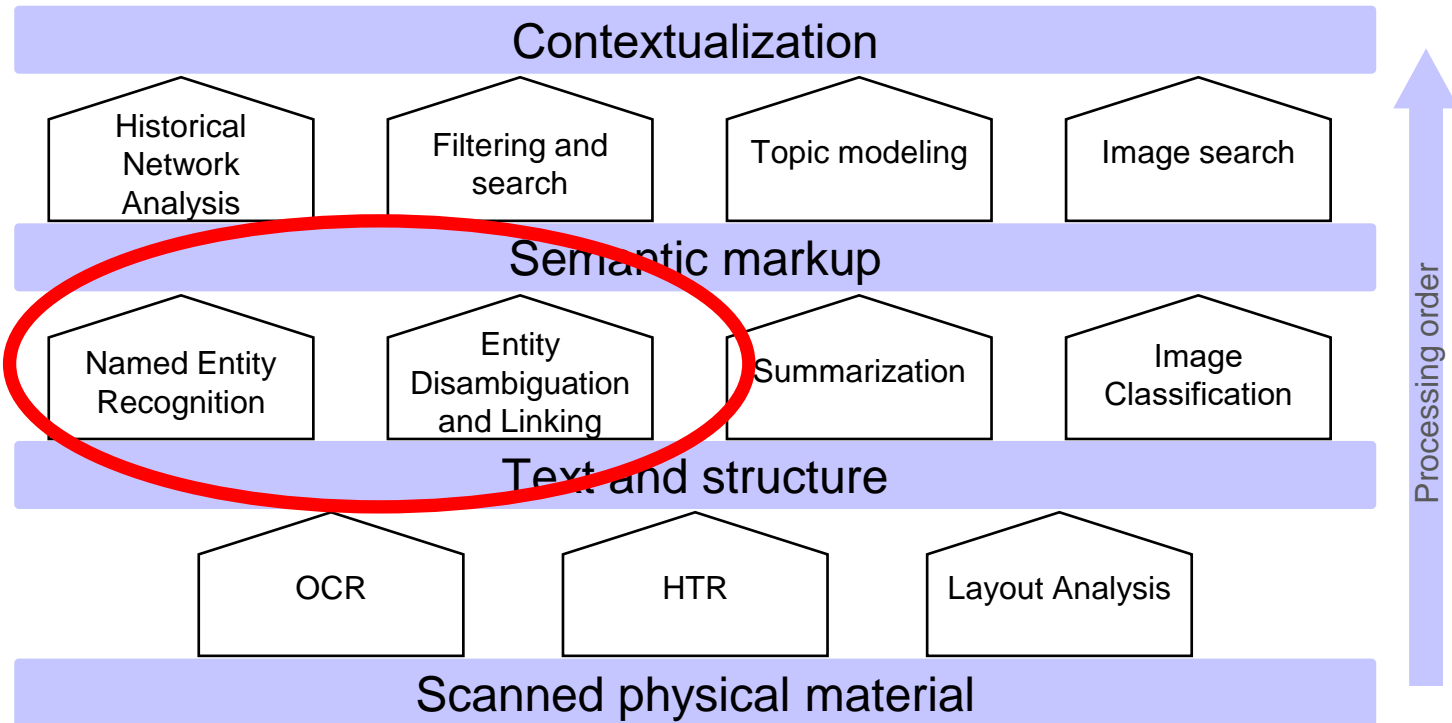
# The main target is the digital humanities researcher

- Modern methods vs. historic material
  - models and methods must work on a variety of material (type, age, quality)
- Results as a means to an end
  - search and filtering
  - topic modeling
  - network analysis
- Deployment
  - usable tools, incl. frontend
  - must process a huge volume of data
  - public institution has limited infrastructure (finances, maintaining, hardware)



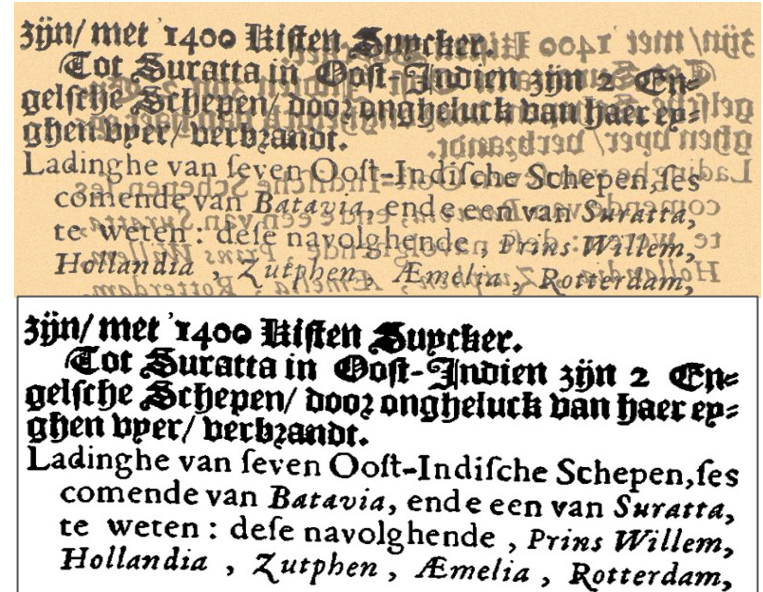
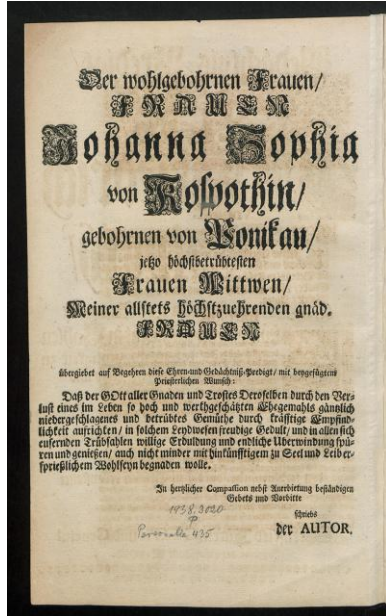
Discovery vs. modeling

# From scanned image to contextualized information



# Prep 1: Image Enhancement and Binarization

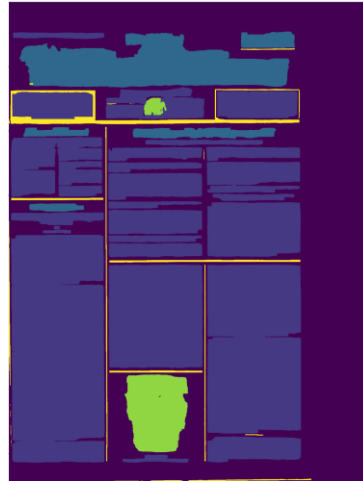
A hybrid CNN-Transformer model for Historical Document Image Binarization, 2023. <https://doi.org/10.1145/3604951.3605508> | [https://github.com/qurator-spk/sbb\\_binarization](https://github.com/qurator-spk/sbb_binarization)



# Prep 2: Document Layout Analysis (Segmentation)

Document Layout Analysis with Deep Learning and Heuristics, 2023.

<https://doi.org/10.1145/3604951.3605513> | <https://github.com/qurator-spk/eynollah>



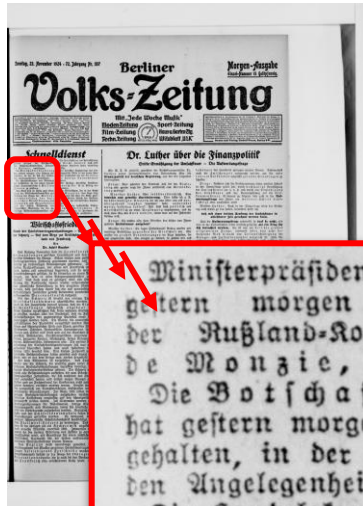
- Background
- Main text
- Header
- Image
- Separator



# Prep 3: Text Recognition (OCR, HTR)

OCR-D: An end-to-end open source OCR framework for historical printed documents, 2019.

<https://doi.org/10.1145/3322905.3322917> | <https://github.com/OCR-D/core>



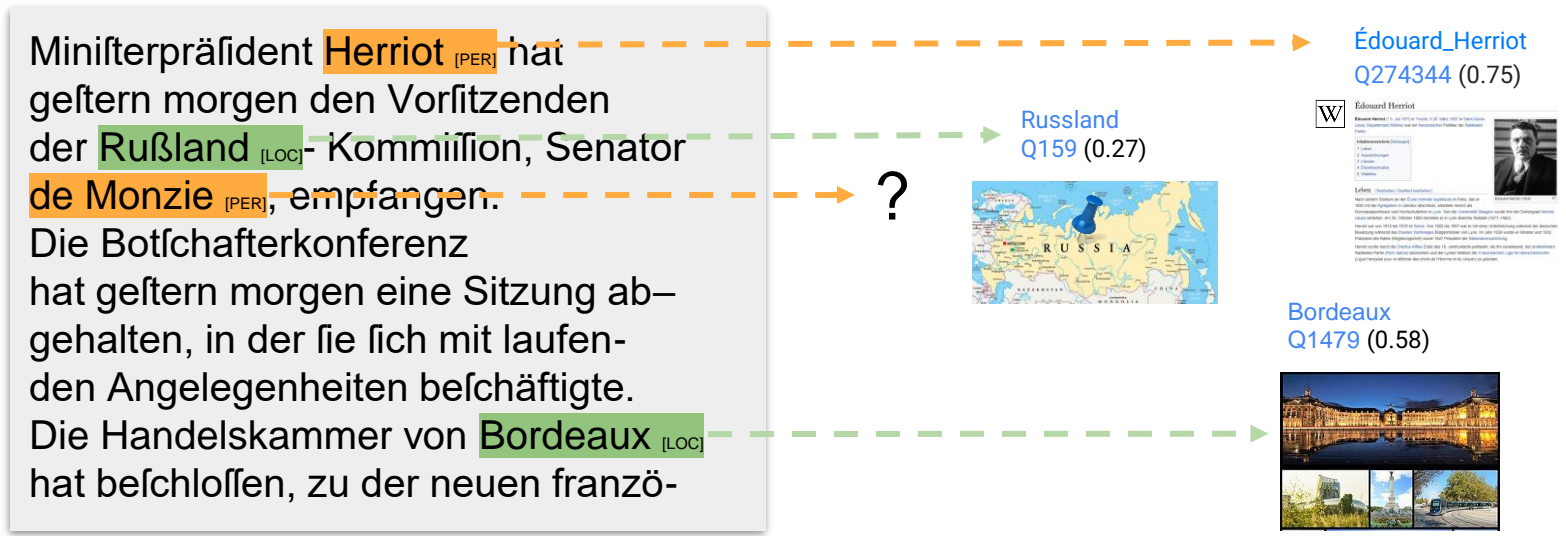
Ministerpräsident Herriot hat gestern morgen den Vorsitzenden der Rußland-Kommission, Senator de Monzie, empfangen. Die Botschafterkonferenz hat gestern morgen eine Sitzung abgehalten, in der sie sich mit laufenden Angelegenheiten beschäftigte. Die Handelskammer von Bordeaux hat beschlossen, zu der neuen franzö-

Ministerpräsident Herriot hat die Feierlichkeiten zur Ueberführung geiern morgen den Vorsitzenden der Autlanü-Kommission, Senator d'Almeida, empfangen. Die Botschafterkonferenz hat gestern morgen eine Sitzung abgehalten, in der sie sich mit laufenden Angelegenheiten beschäftigte. Die Handelskammer von Bordeaux hat beschlossen, zu der neuen französischen Inlandsanleihe 1 Million Francs zu zeichnen. Aus Genf sind in Sofia zwei Delegierte der Balkan-Kommission zur Prüfung der Frage der Massenwanderung der bulgarischen Bevölkerung aus Thrazien und Mazedonien und der letzten Beschwerde der bulgarischen Regierung an die zuständige Balkan-Kommission eingetroffen.

Ministerpräsident Herriot hat gestern morgen den Vorsitzenden der Rußland-Kommission, Senator de Monzie, empfangen. Die Botschafterkonferenz hat gestern morgen eine Sitzung abgehalten, in der sie sich mit laufenden Angelegenheiten beschäftigte. Die Handelskammer von Bordeaux hat beschlossen, zu der neuen französischen Inlandsanleihe 1 Million Francs zu zeichnen. Aus Genf sind in Sofia zwei Delegierte der Völkerbundskommission zur Prüfung der Frage der Massenwanderung der bulgarischen Bevölkerung aus Thrazien und Mazedonien und der letzten Beschwerde der bulgarischen Regierung an die zuständige Völkerbundskommission eingetroffen.

# Named Entity Recognition and Entity Linking

BERT for Named Entity Recognition in Contemporary and Historic German, 2019. [https://konvens.org/proceedings/2019/papers/KONVENS2019\\_paper\\_4.pdf](https://konvens.org/proceedings/2019/papers/KONVENS2019_paper_4.pdf) | [https://github.com/qurator-spk/sbb\\_ner](https://github.com/qurator-spk/sbb_ner) | [https://github.com/qurator-spk/sbb\\_ned](https://github.com/qurator-spk/sbb_ned)

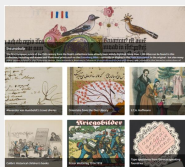


# Named Entity Recognition and Linking Data

- mostly newspaper data, German or multilingual (de/fr/en)
- historical vs. contemporary

## Unlabeled data

- DC-SBB: Digitized Collections
- 1470 – 1945
- 2,333,647 pages



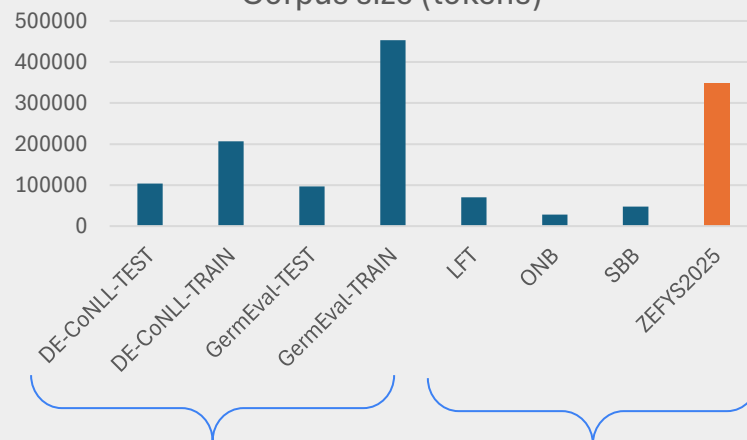
## Knowledge bases

- Wikidata
- Wikipedia
- GND – Integrated Authority File



## Annotated data

Corpus size (tokens)

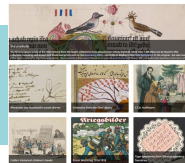


- CoNLL 2003, Frankfurter Rundschau 1992
- German Wikipedia and various online newspapers
- LFT, 1926
- ONB, 1710 – 1873
- SBB, 1872 – 1930
- ZEFYS, 1830 – 1940

# 1) Named Entity Recognition

Multilingual BERT-Base model

+unsupervised pre-training  
on historic text



- How well do models perform on historic text when they are trained on contemporary data?
  - Do we need historic ground truth?
- Does unsupervised training on historic text help?

A lot worse than on the contemporary data.

Yes, for the historical datasets, results improve irrespective of the ground truth.

BERT models outperform previous biLSTM+CRF results but cannot reach best results

# 1) Named Entity Recognition numbers

2019

2025

pre-train:		none	DC-SBB
train	test	$F_1$	$F_1$
GermEval + CoNLL	CoNLL	<b>80.2</b>	79.4
	GermEval	<b>88.0</b>	85.7
	LFT	55.1	<b>55.2</b>
	ONB	58.6	<b>60.1</b>
	SBB	64.1	<b>65.1</b>

gual-cased (Riedl and Padó, 2018)

5-fold cross validation on	pre-train	precision	recall	$F_1$	$F_1$
ONB	Newspaper (1703-1875)	-	-	-	-
	DC-SBB+GermEval + CoNLL	81.5 ±1.8	87.8 ±1.4	<b>84.6 ±1.5</b>	-
	DC-SBB + GermEval	81.6 ±2.5	87.5 ±1.6	84.5 ±1.8	-
	DC-SBB + CoNLL	81.7 ±2.8	87.5 ±1.9	84.5 ±2.3	-
	DC-SBB	81.8 ±2.3	87.1 ±2.1	84.3 ±2.0	-
	GermEval	80.8 ±2.1	85.4 ±1.2	83.0 ±1.4	78.56
	GermEval + CoNLL	80.0 ±1.5	84.7 ±1.6	82.3 ±1.5	-
	CoNLL	79.1 ±2.5	84.5 ±2.1	81.7 ±2.2	76.17
	none	78.0 ±2.4	84.1 ±1.9	80.9 ±2.0	73.31

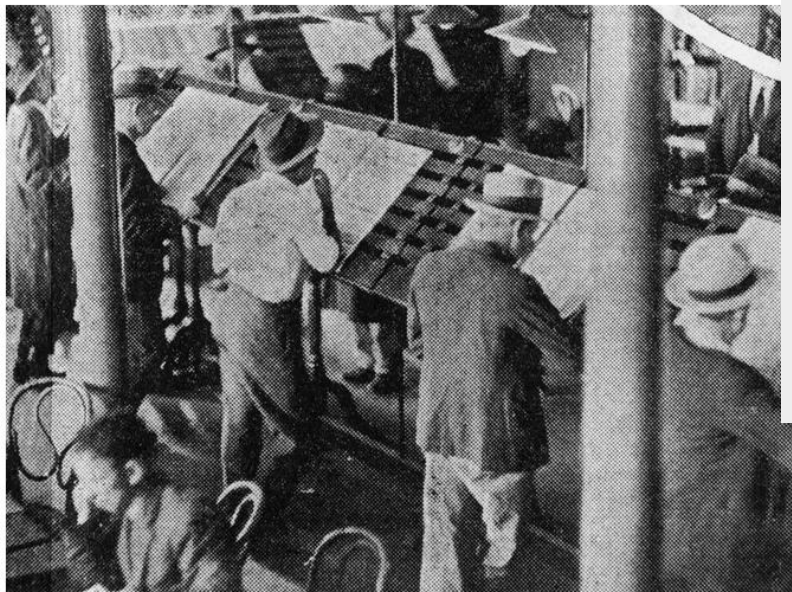
model	$F_1^{TE}$	$F_1^{PZ}$	$F_1^{PH}$	$F_1^{PER}$	$F_1^{LOC}$	$F_1^{ORG}$
<b>Europeana-ÖNB</b>						
Europ-ELECTRA	0.84	<b>0.86</b>	<b>0.88</b>	0.82	0.94	0.44
XLm-RoBERTa	0.84	<b>0.88</b>	0.86	0.85	0.92	0.13
hmBERT	0.84	<b>0.86</b>	<b>0.86</b>	0.78	0.92	0.50
RoBERTa	0.81	0.84	<b>0.87</b>	0.81	0.92	0.20
GBERT	0.83	<b>0.84</b>	0.83	0.75	0.92	0.40
GermanBERT	0.82	0.83	<b>0.84</b>	0.73	0.92	0.53
hmBERT-mini	0.70	0.76	<b>0.78</b>	0.68	0.86	0.28
DistilBERT	0.70	0.74	<b>0.75</b>	0.61	0.86	0.15
hmBERT-tiny	0.50	<b>0.63</b>	<b>0.65</b>	0.48	0.78	0.38

Training on the new larger ZEFYS dataset

## 2) Linking Entities

CLEF-HIPE-2020

### Named Entity Processing on Historical Newspapers



HIPE (Identifying Historical People, Places and other Entities) is a **evaluation campaign on named entity processing on historical newspapers** in **French, German and English**, which was organized in the context of the *impresso* project and run as a [CLEF 2020](#) Evaluation Lab.

Ministerpräsident **Herriot** [PER] hat  
gestern morgen den Vorsitzenden  
der **Rußland** [LOC]-Kommission, Senator  
**de Monzie** [PER], empfangen.  
Die Botschafterkonferenz  
hat gestern morgen eine Sitzung ab-  
gehalten, in der sie sich mit laufen-  
den Angelegenheiten beschäftigte.  
Die Handelskammer von **Bordeaux** [LOC]  
hat beschlossen, zu der neuen franzö-

Édouard Herriot  
Q274344 (0.75)



Russland  
Q159 (0.27)

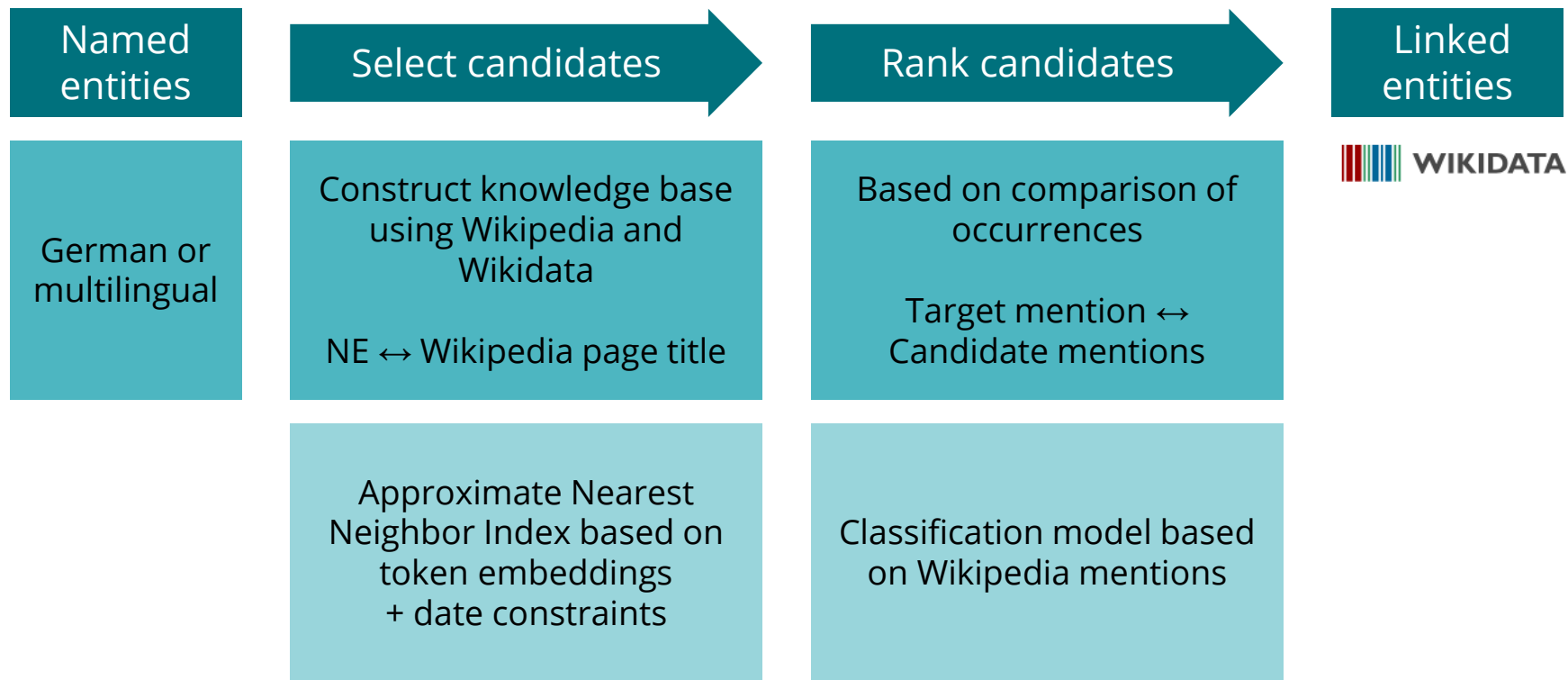


Bordeaux  
Q1479  
(0.58)



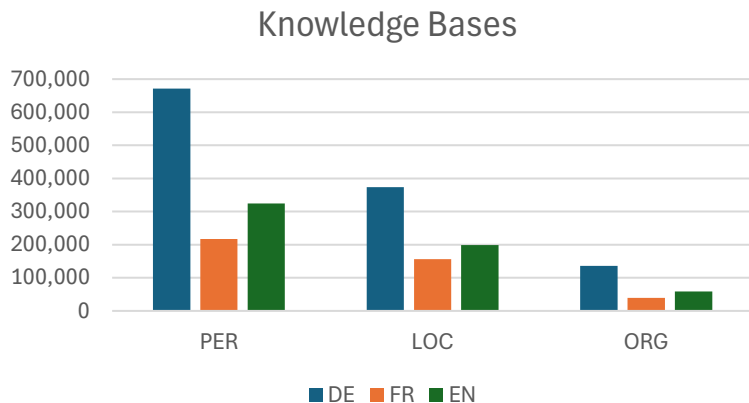
?

## 2) Linking Entities, continue with BERT models



## 2) Named Entity Linking results

- How well do models perform on historic text and for the different languages?



Results are similar for FR and DE even though the knowledge bases have different sizes (coverage of the test data is similar)

Lower scores for EN → OCR

	2020		2022	
	SBB	Best system	SBB	Best system
Language	F1	F1	F1	F1
German	0.389	0.534	0.506	0.506
French	0.407	0.598	0.596	0.602
English	0.141	0.531	0.393	0.546

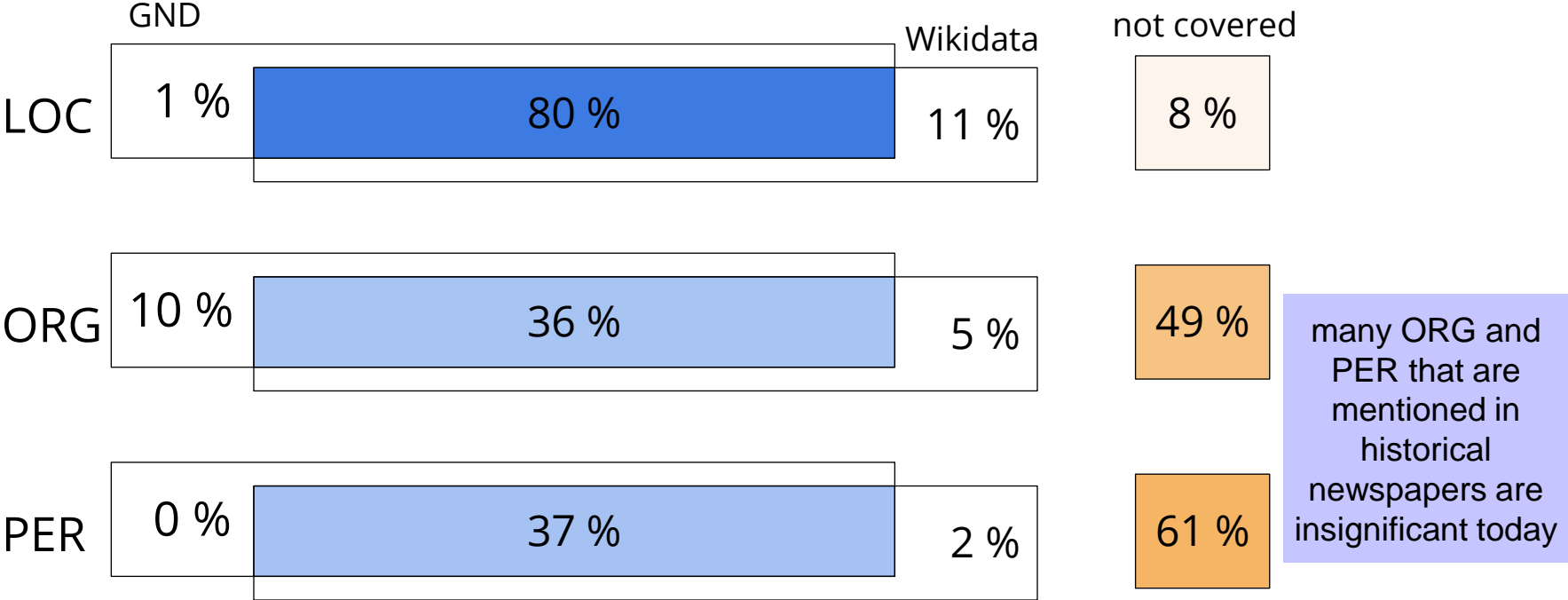
Effect of knowledge base size and quality is not clear!

coverage  
.86 – .99

# Contemporary knowledge bases cover the same entities

Overall, GND covers 68.1% / 51.8%, Wikidata covers 69.5% / 54.1%

N = 791 entities



# Knowledge base quality matters, too

- **duplicate entries**
- **missing entities (mismatching focus)**
  - GND focusses on authors, newspapers on public persons
- **incomplete entries**
  - “empty” entries contain attributes that are matched by the method
- **missing spelling variants**
  - newspaper used a spelling that is missing from the correct entry
- **entities with changing extension**

## **Reichstag** (Q160208)

parliament of Germany from 1871 to 1918  
16 statements, 25 sitelinks – 02:53, 25 May 2026

## **Reichstag** (Q878525)

parliament of the Third Reich from 1933 to 1945  
13 statements, 21 sitelinks – 12:14, 18 May 2026

## **Reichstag** (Q321246)

parliament of the North German Confederation 1867-1870  
14 statements, 11 sitelinks – 16:29, 27 January 2026

## **Lüneburger Heide** (Q61040297)

Item Discussion

special protection area in the EU defined by the bird

## **Lüneburger Heide** (Q60685764)

Item Discussion

protected area in the European Union defined by the habi

Do we choose the most extensive one? the one closest in time? only the one that matches in every detail?

# NER and EL takeaways

- **Longterm goal and requirement:**
  - deliver robust and solid baseline performance even without optimization on the diverse material (age, heterogeneity, size)
- OCR is crucial
- NER results are stable, EL results are a starting point
- improve knowledge base construction
- exploit more constraints like time/date, e.g. publication dates

# There are more NLP tasks

## Multilingual search and summarization

Multilingual search

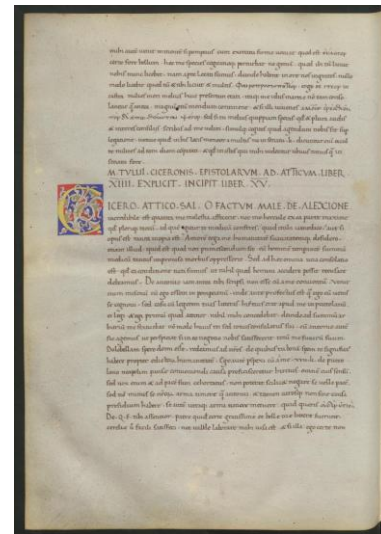
Summarization

Translation

The screenshot displays the SPUNK Assistant interface. On the left, there are search filters and a search history section. The main area shows search results for 'Socialism and the Rise of Social Thought in Modern Japan' and 'The Language of Revolution: Defining Socialism and 社会 in Public Media'. A 'CONCEPT SYNTHESIS' panel on the right provides a summary of the findings, mentioning 'transnational revolutionary ideas were budding against a backdrop of increasing state surveillance and the codification of "dangerous thought" through the Peace Preservation Law (article 6)'. At the bottom, there are buttons for 'Articles selected', 'Generate Synthesis', and 'Download PDF'.

## Hand-written text processing

- Text initials
- Finding duplicates
- Identify texts

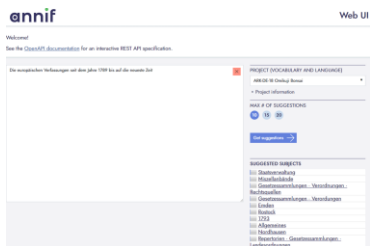


# Semi-automatic Subject Indexing

- Analyse requirements in different subject domains at Stabi
- Train and evaluate [Annif](#) models for ARK and BK

## ARK results (summary)

project id	eval params	P	R	F1	NDCC	F1@5 (ref / LB)	NDCC@5 (ref / LB)
ark-mit0-da-18		0.0020	0.0775	0.0206	0.0523	0.0257	0.0473
ark-rt01-da-18	1:1, 4:5	0.2384	0.1933	0.2081	0.2036	0.1535	0.3152
ark-ontokj-da-18	1:2, 0:5	0.3404	0.1991	0.3452	0.3938	0.2216	0.4886
ark-ontokj-de	1:2, 0:0	0.2425	0.4338	0.3120	0.4415	0.1846	0.4755
ark-fa01en-da-18	1:1, 0:0	0.2958	0.2285	0.2466	0.2434	0.1947	0.3275
ark-ontokj-de-01a-only	1:1, 0:0	0.4853	0.4582	0.4669	0.4678	0.3771	0.5243
ark-ontokj-de-01a-content	1:1, 0:0	0.4861	0.4587	0.4675	0.4683	0.2200	0.5103
ark-ontokj-lar-01a-content	1:1, 0:0	0.5755	0.4998	0.5170	0.5153	0.2876	0.6305
ark-ontokj-multiling-01a-content	1:1, 0:0	0.4779	0.4321	0.4604	0.4639	0.2197	0.5410
ark-ontokj-free-01a-content	1:1, 0:0	0.4228	0.6044	0.4103	0.4369	0.1885	0.4939
ark-ontokj-eng-01a-content	1:1, 0:0	0.3838	0.3694	0.3741	0.3795	0.1839	0.4702



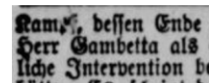
# Transcription and annotation

neat: neat annotation tool

[User Guide](#) | [Annotation Guidelines](#) | [Issues](#)



<< LOCATION	POSITION	TOKEN	NE-TAG	NE-EMB	ID >>	
	9	10	wäre	O	O	-
	10	11	,	O	O	-
	11	12	wenn	O	O	-
	12	13	nicht	O	O	-
	13	14	Herr	O	O	-
	14	15	Gambetta	B-PER	O	Q295090
	15	16	als	O	O	-
	16	17	deus	O	O	-
	17	18	ex	O	O	-
	18	19	machina	O	O	-
	19	20	erchiene	O	O	-
	20	21	wäre	O	O	-

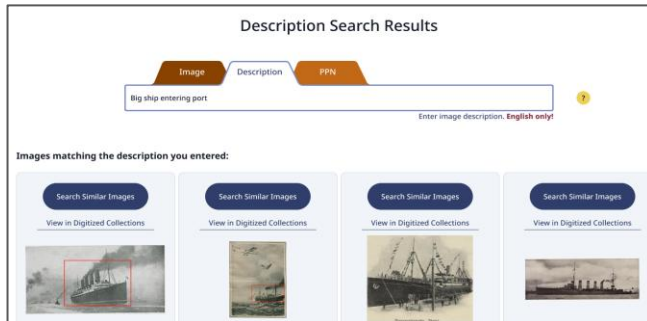


[entarge](#) | [full](#)

<https://github.com/qurator-spk/neat>

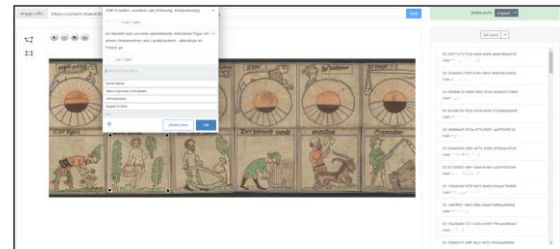
## Image (Similarity) Search

- Develop a multi-modal text-image-search
  - Find images by upload of an example image
  - Find images by a textual description (48 languages currently supported)
  - Find all images within a specific document



## Image Analysis for Special Subject Domains

- Develop web-based image annotation tool
- Adapt and train AI models specifically for
  - Iconographic images ([ICONCLASS](#))
  - Provenance features (e.g. printers marks)
  - Watermarks
  - Illustrations in children's books





But there are also great opportunities in using ML methods

develop methods for more efficient processing

create more ways of accessing data for the user

model data that changes with time and context

use diverse and trustworthy data



Culture for AI

AI for Culture

# Where to go from here

- multilinguality: SBB has documents in many languages
  - *search/retrieval, summarization*
- entity linking: enable historic research through entity search
  - *search for entities by ID*
- image processing
  - *search, comparison*
- derived text formats
  - *for material that cannot be accessed directly*
- pre-training with more rich, varied and contextualized data

# Outputs

- Open Source Software
  - [github.com/qurator-spk](https://github.com/qurator-spk)
  - Document Layout Analysis (Segmentation)
  - Optical Character Recognition
  - Image Similarity Search
  - Named Entity Recognition
  - Entity Disambiguation and Linking
- Models [huggingface.co/SBB](https://huggingface.co/SBB)
  - Document Layout Analysis
  - Named Entity Recognition
  - Entity Linking
- Datasets
  - Hugging Face [huggingface.co/SBB](https://huggingface.co/SBB)
  - Zenodo [zenodo.org/communities/stabi/](https://zenodo.org/communities/stabi/)

# OCR data



Datasets of Staatsbibliothek zu Berlin - Berlin State Library

Published June 26, 2019 | Version 1.0

Dataset  Open

## OCR fulltexts of the Digital Collections of the Berlin State Library (DC-SBB)

Labusch, Kai<sup>1</sup>  ; Zellhöfer, David<sup>1</sup> 

Show affiliations

The digital collections of the SBB contain 153,942 digitized works from the time period of 1470 to 1945.

At the time of publication, 28,909 works have been OCR-processed resulting in 4,988,099 full-text pages. For each page with OCR text, the language has been determined by *langid* (Lui/Baldwin 2012).

corpus-entropy.pkl	entropy rate per document page
corpus-language.pkl	language per document page
corpus.zip	fulltext corpus (extracts to .txt format)
de_corpus.zip	German sub-corpus (extracts to .txt format)
selection_de.pkl	Selection list of German documents
xml2csv_alto.csv	fulltext corpus per document page (incl.OCR word confidences)

### Sources

Marco Lui and Timothy Baldwin. 2012. *Langid.py*:

An off-the-shelf language identification tool. In Proceedings of the ACL 2012 System Demonstrations,

ACL '12, pages 25–30, Stroudsburg, PA, USA. Association for Computational Linguistics

<https://zenodo.org/records/3257041>

# NER models

README Apache-2.0 license

## NER - Demo

Task: Named Entity Recognition

Model: DC-SBB + CONLL + GERMEVAL

Input text:

Die Staatsbibliothek zu Berlin (ab 1661: Churfürstliche Bibliothek; ab 1701: Königliche Bibliothek; ab 1918: Preußische Staatsbibliothek; ab 1954: Deutsche Staatsbibliothek)[2] ist eine Einrichtung der Stiftung Preußischer Kulturbesitz, einer durch Bundesgesetz errichteten rechtsfähigen Stiftung des öffentlichen Rechts mit Sitz in Berlin. Die Bibliothek sammelt für den Spitzenbedarf der Forschung wissenschaftlich relevante Literatur aus allen Zeiten, allen Ländern und in allen Sprachen. Sie ist eine der größten Bibliotheken Deutschlands und darüber hinaus eine der größten der Erde.

Go

Ergebnis:

Die Staatsbibliothek zu Berlin ( ab 1661 : Churfürstliche Bibliothek ; ab 1701 : Königliche Bibliothek ; ab 1918 : Preußische Staatsbibliothek ; ab 1954 : Deutsche Staatsbibliothek ) [ 2 ] ist eine Einrichtung der Stiftung Preußischer Kulturbesitz , einer durch Bundesgesetz errichteten rechtsfähigen Stiftung des öffentlichen Rechts mit Sitz in Berlin . Die Bibliothek sammelt für den Spitzenbedarf der Forschung wissenschaftlich relevante Literatur aus allen Zeiten , allen Ländern und in allen Sprachen . Sie ist eine der größten Bibliotheken Deutschlands und darüber hinaus eine der größten der Erde . Die Staatsbibliothek zu Berlin ist die größte wissenschaftliche Universalbibliothek im deutschen Sprachraum . Zu den bedeutendsten Unterstützern der Bibliothek gehört die Deutsche Forschungsgemeinschaft ( DFG ) mit Sitz in Bonn .

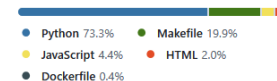
Legende:

- [Person]
- [Ort]
- [Organisation]
- [keine Named Entity]

### Contributors 3

- labusch Kai Labusch
- cneud Clemens Neudecker
- joergleh joergleh

### Languages



[https://github.com/qurator-spk/sbb\\_ner](https://github.com/qurator-spk/sbb_ner)

# API access

## digital.staatsbibliothek-berlin.de

### OAI 2.0 Request Results

More information to this OAI interface and the sets are available in the Lab of the SBB: <https://lab.sbb.berlin/dc>

---

[Identify](#) | [ListRecords \(oai\\_dc\)](#) | [ListRecords \(mets\)](#) | [ListSets](#) | [ListMetadataFormats](#) | [ListIdentifiers](#)

---

You are viewing an HTML version of the XML OAI response. To see the underlying XML use your web browsers view source option. More information about this XSLT is at the [bottom of the page](#).

<b>Datestamp of response</b>	2026-05-23T15:20:13Z
------------------------------	----------------------

<b>Request URL</b>	https://oai.sbb.berlin/
--------------------	-------------------------

### OAI Error(s)

The request could not be completed due to the following error or errors.

<b>Error Code</b>	badVerb
-------------------	---------

Required verb argument not found.

---

[Identify](#) | [ListRecords \(oai\\_dc\)](#) | [ListRecords \(mets\)](#) | [ListSets](#) | [ListMetadataFormats](#) | [ListIdentifiers](#)

---

### About the XSLT

An XSLT file has converted the [OAI-PMH 2.0](#) responses into XHTML which looks nice in a browser which supports XSLT such as Mozilla, Firebird and Internet Explorer. The XSLT file was created by [Christopher Gutteridge](#) at the University of Southampton as part of the [GNU EPrints system](#), and is freely redistributable under the [GPL](#).

If you want to use the XSL file on your own OAI interface you may but due to the way XSLT works you must install the XSL file on the same server as the OAI script, you can't just link to this copy.

For more information or to download the XSL file please see the [OAI to XHTML XSLT homepage](#).

<https://oai.sbb.berlin/oai>

# Summary

Libraries have real data and real users

Curated data that is not available elsewhere  
+ expert knowledge

The material is a challenge at every level of processing (age, heterogeneity, size)

**Libraries can be a testbed for NLP methods**

Ideas and suggestions welcome!



# References

Illustrations: Katerina Limpitsouni / [undraw.co](https://undraw.co)

- Alkemade, H. et al. (2023). *Datasheets for Digital Cultural Heritage Datasets*.
- Labusch, K., & Neudecker, C. (2020). *Named Entity Disambiguation and Linking on Historic Newspaper OCR with BERT*. CLEF 2020.
- Labusch, K., & Neudecker, C. (2022). *Entity Linking in Multilingual Newspapers and Classical Commentaries with BERT*. CLEF 2022.
- Labusch, K., & Neudecker, C. (2023). *Gauging the Limitations of Natural Language Supervised Text-Image Metrics Learning by Iconclass Visual Concepts*. HIP '23.
- Labusch, K. et al. (2024). *Automatisierte semantische Anreicherung von historischen Texten*. b.i.t.online, 27(3), 232–241.
- Labusch, K. et al. (2019). *BERT for Named Entity Recognition in Contemporary and Historical German*. KONVENS 2019.
- Neudecker, C. et al. (2021). *A survey of OCR evaluation tools and metrics*. The 6th International Workshop on Historical Document Imaging and Processing, 13–18.
- Neudecker, C. et al. (2021). *Methoden und Metriken zur Messung von OCR-Qualität für die Kuratierung von Daten und Metadaten*. Qualität in der Inhaltserschließung.
- Neudecker, C. (2022). *Cultural Heritage as Data: Digital Curation and Artificial Intelligence in Libraries*. Proc. of the Third Conference on Digital Curation Technologies.
- Menzel, S. et al. (2021). *Named Entity Linking mit Wikidata und GND – Das Potenzial handkuratierter und strukturierter Datenquellen für die semantische Anreicherung von Volltexten*. Qualität in der Inhaltserschließung.
- Rezanezhad, V. et al. (2023). *A hybrid CNN-Transformer model for Historical Document Image Binarization*. HIP '23.
- Rezanezhad, V. et al. (2023). *Document Layout Analysis with Deep Learning and Heuristics*. HIP '23.
- Schaefer, R., & Neudecker, C. (2020). *A Two-Step Approach for Automatic OCR Post-Correction*. In Proc. of the 4th Joint SIGHUM Workshop on Computational Linguistics for Cultural Heritage, Social Sciences, Humanities and Literature.
- Schneider, S. et al. (2025). *ZEFYS2025: A German Historical Newspaper Dataset for Named Entity Recognition and Entity Linking*. KONVENS 2025.
- Zellhöfer, D. (2019). *Multimodal Datasets of the Berlin State Library*. In Proc. of the Poster Session of the 2nd Conference on Language, Data and Knowledge (LDK).