# Text+

# Utilise and Preserve Text- and Language-Based Research Data

Discussion of the proposal for the funding second phase 2026–31

27.11.2025 Bonn

1. What is **Text+?**

2. Where does **Text+** stand now?

3. **Text+**: the road until 2031

# 1. What is **Text+?**

# 2. Where does Text+ stand now?

# 3. Text+: the road until 2031

# Covering the **entire humanities research landscape**
# Services for the **entire spectrum of text-based research**
# Strong combination of **university and non-university institutions**

**Andrea Rapp**

**Hanna Fischer**

**Andreas Witt**

**Philipp Wieder**

**Philippe Genêt**

**21 universities, 8 academies of sciences and humanities, 12 computing centres, libraries, and non-university research institutions**

6 task areas with co-spokespersons from universities and non-university research institutions:

- **Library and Archive Holdings as Research Data**: Elke Teich (Saarland University) and Philippe Genêt (DNB)
- **Language Corpora and Elicited Data**: Angelika Zirker (Tübingen University) and Andreas Witt (IDS Mannheim)
- **Lexical Resources**: Hanna Fischer (Marburg University) and Alexander Geyken (BBAW)
- **Editions**: Andrea Rapp (Technical University Darmstadt) and Andreas Speer (Academy of Sciences North-Rhine Westphalia)
- **Infrastructure Operations**: Philipp Wieder (GWDG), Vivien Petras (Humboldt University Berlin) and Regine Stein (VZG)
- **Administration**: Philipp Wieder (GWDG) and Andreas Witt (IDS Mannheim)

Text+

# Covering the **entire humanities research landscape**
# Services for the **entire spectrum of text-based research**
# Strong combination of **university and non-university institutions**



**Andrea Rapp**

**Hanna Fischer**

**Andreas Witt**

**Philipp Wieder**

**Philippe Genêt**

**21 universities, 8 academies of sciences and humanities, 12 computing centres, libraries, and non-university research institutions**

6 task areas with co-spokespersons from universities and non-university research institutions:

- **Library and Archive Holdings as Research Data**: Elke Teich (Saarland University) and Philippe Genêt (DNB)
- **Language Corpora and Elicited Data**: Angelika Zirker (Tübingen University) and Andreas Witt (IDS Mannheim)
- **Lexical Resources**: Hanna Fischer (Marburg University) and Alexander Geyken (BBAW)
- **Editions**: Andrea Rapp (Technical University Darmstadt) and Andreas Speer (Academy of Sciences North-Rhine Westphalia)
- **Infrastructure Operations**: Philipp Wieder (GWDG), Vivien Petras (Humboldt University Berlin) and Regine Stein (VZG)
- **Administration**: Philipp Wieder (GWDG) and Andreas Witt (IDS Mannheim)

Text+

5

1. What is Text+?

2. **Where does Text+ stand now?**

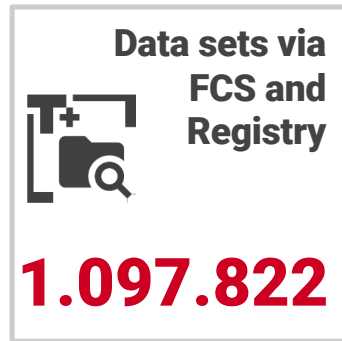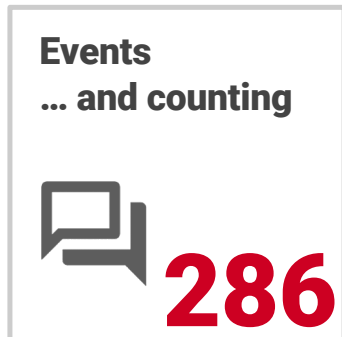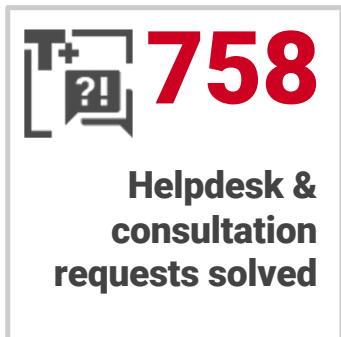3. Text+: the road until 2031

# Achievements for the community (2021–26)

- Empowering community's access to text and language data

- Making research data FAIR and globally reusable

- Providing trusted services and expert consulting

- Strengthening the community through training and outreach

- Driving innovation with AI-ready data

- Leading in standardisation and authority data

- Ensuring sustainability of services

- Actively engaging in OneNFDI and Europe

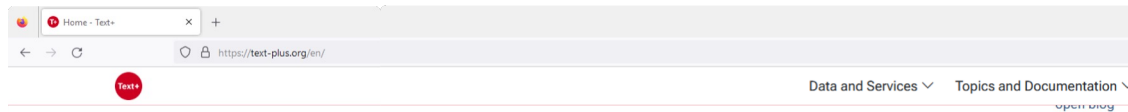- Achieving full interoperability with Base4NFDI

# Achievements for the community (2021–26)

- Empowering community's access to text and language data

- **Making research data FAIR and globally reusable**

- Providing trusted services and expert consulting

- **Strengthening the community through training and outreach**

- **Driving innovation with AI-ready data**

- Leading in standardisation and authority data

- Ensuring sustainability of services

- Actively engaging in OneNFDI and Europe

- Achieving full interoperability with Base4NFDI

# Achievements in numbers (2021–26)

**Text+ Centres … and counting**
**33**

**Data sets via FCS and Registry**
**1.097.822**

**115**
**Tools & services for the community**

**Webportal visits**
**34.400**

**Text+ Data Space**

**758**
**Helpdesk & consultation requests solved**

**Events … and counting**
**286**

**528**
**Publications via the project bibliography**

**Community activities**

Text+

# Text+ portal as entry point to the portfolio



## For

- Diversity of text-based academic disciplines
- Heterogenous research data encompassed and represented in the data domains

## Through

- RDM services across entire research data lifecycle
- Training and consulting: fostering digital literacy
- Transfer to other disciplines and domains

## Responsive to

- Tectonic shift caused by the availability of AI

**Services for the community**

5 basic services integrations

Data depositing form

many other services

> 20 base models in the LLM service

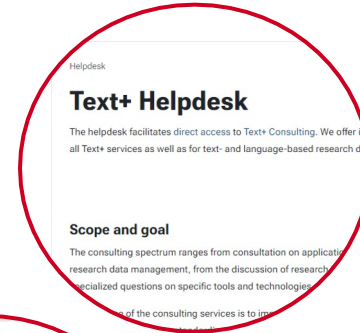>115 services in the catalogue

>230 events

FCS

Helpdesk

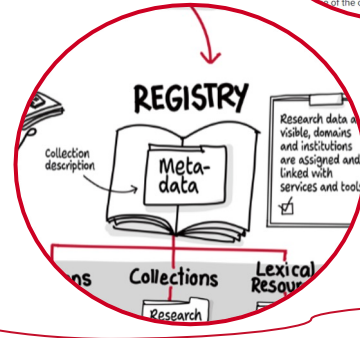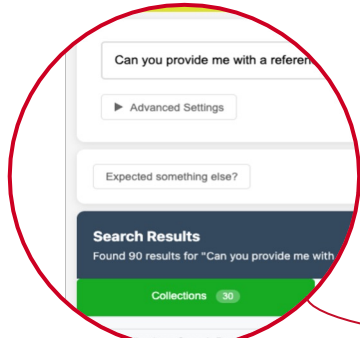GND-Agency Text+

3 data domains in the registry

AI based search (experimental)

# 1. What is Text+?

# 2. Where does Text+ stand now?

# 3. **Text+**: the road until 2031

# Strategic themes of **Text+**

- Language/text as dominant interface to data

- Dynamic interplay between structured and unstructured data

- Human-centred research culture

- Longevity of the research data infrastructure

# Work programme of the data domains

**M1 Federated Data Infrastructure**

- Text+ Data Space
- Federated resource catalogue and Federated Content Search
- Quality Standards for Research Data (FAIR)

**M2 Interfaces between Data from Data Domains**

- Entity linking and authority files (GND)
- Interoperability and knowledge graph integration
- Tracing and citing data and tools

**M3 Enabling Text and Language Data for AI**

- Automatic annotation
- High quality and curated data for LLMs
- Derived text formats

**M4 Implicit and Explicit Structures of Language Data**

- Capturing structures in unstructured data
- Evaluation and proving of generative AI
- Knowledge representation of implicit and explicit structures

**M5 Community Activities**

- Dissemination and community outreach
- RDM training and education
- Consulting and helpdesk

# Example activity: Leveraging **Library Holdings for AI**

- Train models and scripts for automatic classification tasks

- Provide large scale, high quality, curated data for LLM training (copyright-free)

- Publish derived text formats (DTF) of in-copyright material

  [Task Area Library and Archive Holdings as Research Data, Measure 3 Enabling text and language data for AI]

Text+

# Work programme of infrastructure operations

**T+ DATA**
**IO-M1**
**Federated Research Data Infrastructure**

**T+ PROCESSES**
**IO-M2**
**Service Integration w/ Data Domains**

**T+ PLATFORM**
**IO-M3**
**Platform and Core Services**

**T+ ECOSYSTEM**
**IO-M4**
**Coherence**

- Four measures supporting two core project goals:
  - Technical gateway to the Text+ Data Space
  - Empowering text-based research

- DATA: FCS & Registry, data integration with data centres, data enrichment

- PROCESSES: UX/UI of services, technology stack coherence

- PLATFORM: enabling and supporting role for the project (e.g. LLM), Portal, Base4NFDI, IT service management

- ECOSYSTEM: alignment with NFDI, EOSC, standards, indicators and impact alignment with data domains and administration

Text+

# Example activity:
# Basic services integration

- IO integrates services of Base4NFDI in the architecture of Text+ such as Terminology Service (TS4NFDI), Knowledge Graph Infrastructure (KGI4NFDI), Persistent Identifiers (PID4NFDI), Jupyter4NFDI and Identity and Access Management (IAM4NFDI).

- IO aligns genuine Text+ services such as LLM or search & retrieval (FCS, registry) with the overall (upcoming) architecture of OneNFDI.

- By these activities, Text+ leverages the uptake of its services in the larger NFDI community.

[Task Area Infrastructure Operations , Measure 3 Platform and Measure 4 Ecosystem]

# What **Text+ will achieve** until 2031

- Offering stable, reliable, sustainable services that are actively used by our communities
- Making fine-tuned LLMs and high quality, curated data for LLM training available
- Enabling automatic transcriptions and annotations
- Collaborating with specialised information services, learned societies, and other experts

[Text+ Vision] **The text- and language-oriented humanities and social sciences extensively <span style="color:red">leverage the possibilities of digitisation in their research, teaching, and transfer</span>, establishing a common data culture.**

1. What is Text+?

2. Where does Text+ stand now?

3. Text+: the road until 2031

[Text+ Vision] **The text- and language-oriented humanities and social sciences extensively leverage the possibilities of digitisation in their research, teaching, and transfer, establishing a common data culture**.

1. What is Text+?

2. Where does Text+ stand now?

3. Text+: the road until 2031