

# Europäische Sprachmodelle: Herausforderungen, Initiativen, Lösungsansätze

Prof. Dr. Georg Rehm, DFKI GmbH  
georg.rehm@dfki.de



# Eine (sehr) kurze Geschichte der Erschließung – In drei Phasen

- Früher: Karteikarten erschließen Metadaten über Werke
- Dann: Einzug der digitalen Kataloge, die aber ebenfalls „nur“ Metadaten erschließen
- Kurze Zeit später: Suchmaschinen (Plural) erschließen eine heterogene Mischung aus Metadaten und Volltexten in diversen analogen und digitalen Bibliotheken sowie Katalogen



FACHVERANSTALTUNG

# KI IN BIBLIOTHEKEN: NEUE WEGE MIT GROSSEN SPRACHMODELLEN?

**Beginn** 07.12.2023, 13:00 Uhr  
**Ende** 08.12.2023, 13:00 Uhr  
**Wo** Frankfurt am Main



# KI in Bibliotheken – Eine Vision

- Große Sprachmodelle werden auch in Bibliotheken Einzug halten
- Bereits jetzt verbessern sie Basisfunktionen: OCR, OLR, Search, Discovery etc.
- In einigen Jahren dann die Revolution:

**Mittels ScienceGPT mit der wissenschaftlichen Literatur sprechen,  
mit allen wissenschaftlichen Publikationen in den Dialog treten,  
diese aktiv konsultieren, Forschung gemeinsam mit den Quellen betreiben!**

- Basis: Ein LLM wird auf Basis möglichst *aller* jemals erschienenen wissenschaftlichen Veröffentlichungen in allen Sprachen trainiert und mit Chat-Funktionalität ausgestattet.
- Voraussetzungen (Auswahl): Die gesamte wissenschaftliche Fachliteratur liegt digital vor (in allen Sprachen); die Rechteproblematik ist geklärt; LLMs halluzinieren nicht mehr; Faktualität liegt bei 100%; LLMs beherrschen Mathematik und alle Naturgesetze (durch das Training); die Ausgaben von LLMs liefern Quellenangaben; strukturiertes Wissen ist im LLM nutzbar.

# Überblick

- Einleitung
- Modelle und digitale Sprachgerechtigkeit
- Daten
- Wissenschaftliche Daten
- Schlussfolgerungen

# Modelle

# Kontext: Große Sprachmodelle – Large Language Models (LLMs)

- LLMs sorgen in allen Domänen für massive Disruptionen (GPT-3, ChatGPT, Gemini etc.)
- Sie stellen den bedeutendsten Durchbruch in der KI-Forschung der jüngsten Vergangenheit dar
- LLMs werden mit immens großen Datenmengen trainiert (Sprachdaten, d.h. Text)
- Sie nutzen Hunderte von Terabytes, oft Petabytes (Billionen Tokens), an Sprachdaten sowie auch Bild-, Video- und Audiodaten
- Besonderheit: Für fast alle Sprachen Europas existieren Herausforderungen

## BUSINESS

# ChatGPT Shows Just How Far Europe Lags in Tech

Analysis by Lionel Laurent | Bloomberg

February 21, 2023 at 2:12 a.m. EST



Comment 1



Gift Article



Share

Europe is where ChatGPT gets regulated, not invented. That's something to regret. As unhinged as the initial results of the artificial-intelligence arms race may be, they're also another reminder of how far the European Union lags behind the US and China when it comes to tech.



# Regelrechte Explosion des LT/NLP Markts: 439.85 Mrd. \$ bis 2030



## Natural Language Processing Market Size, Share & Trends Analysis Report By Component, By Deployment Model, By Enterprise Size, By Type, By Application, By End-use, By Region, And Segment Forecasts, 2023 - 2030

Report ID: GVR-4-68040-020-4 | Number of Pages: 100 | Format: Electronic (PDF)  
 Historical Range: 2017 - 2021 | Industry: [Technology](#)

<https://www.grandviewresearch.com/industry-analysis/natural-language-processing-market-report>

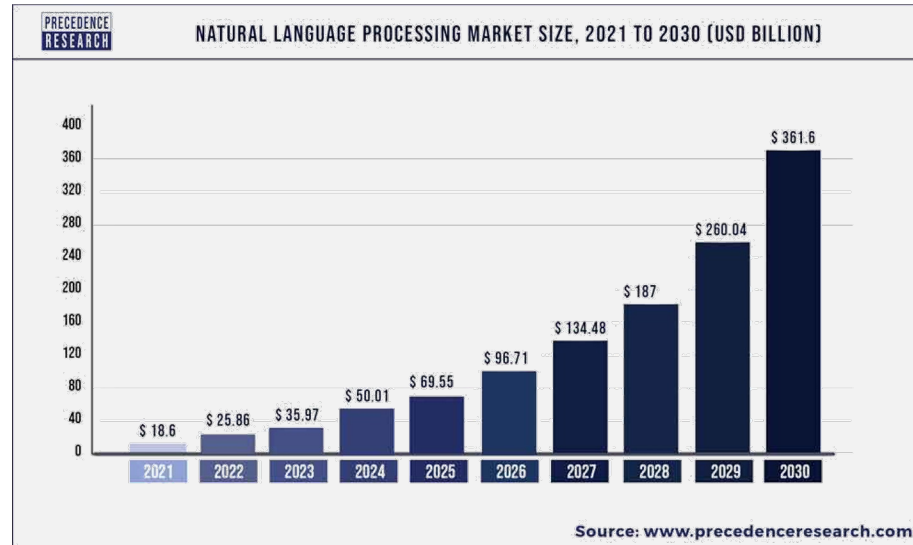
### Natural Language Processing Market Report Scope

Report Attribute	Details
Market size value in 2023	USD 40.98 billion
Revenue forecast in 2030	USD 439.85 billion
Growth rate	CAGR of 40.4% from 2023 to 2030
Base year for estimation	2022
Historical data	2017 - 2021
Forecast period	2023 - 2030
Quantitative units	Revenue in USD million and CAGR from 2022 to 2030

Players leading the NLP market include-

- 3M Co. (US)
- IBM Corporation (US)
- Hewlett-Packard Co. (US)
- Oracle Corporation (US)
- Apple Inc. (US)
- Microsoft Corporation (US)
- SAS Institute Inc. (US)
- Dolby Systems Inc. (US)
- Verint Systems Inc. (US)
- Net base Solutions Inc. (US)

US!



Ohne gezielte Interventionen der EU wird Europa in den kommenden Jahren weiter an den Rand gedrängt.

<https://www.precedenceresearch.com/natural-language-processing-market>

# Europäische Initiativen

- Es existieren diverse europäische Initiativen für die Entwicklung großer Sprachmodelle
  - Forschungsprojekte in nahezu allen Ländern, z.B. Spanien, Dänemark, Slovenien, Italien etc.
  - Firmen in vielen Ländern, z.B. Finnland (Silo.ai) oder Frankreich (Mistral)
  - EU-Projekte, z.B. HPLT
  - Übergreifende europäische Initiative: ALT-EDIC – dazu später mehr
- Wichtige Initiativen in Deutschland
  - Firmen: Aleph Alpha (Heidelberg) und andere
  - Forschungsprojekte: OpenGPT-X und andere
- Herausforderungen: Datenlage speziell für europäische Sprachen; High-Performance Compute Einrichtungen; Geschwindigkeit der zentralen Big-Tech-Player vs. Geschwindigkeit Europas

# Europaparlament und Europäische Kommission (2022/23)



Intergroup for Traditional Minorities, National Communities and Languages (März 2022)

## Intergroup for Traditional Minorities, National Communities and Languages

Dear Members of the Intergroup,  
Dear Colleagues,

Welcome to the next Minority Intergroup meeting, which will be held on Thursday, 5 May from 10.00 to 12.00. The format will be communicated early next week.

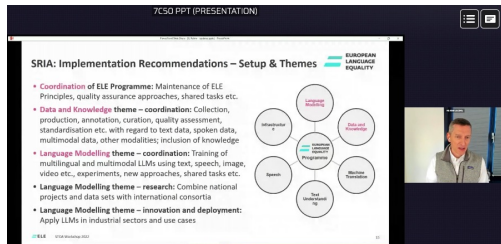
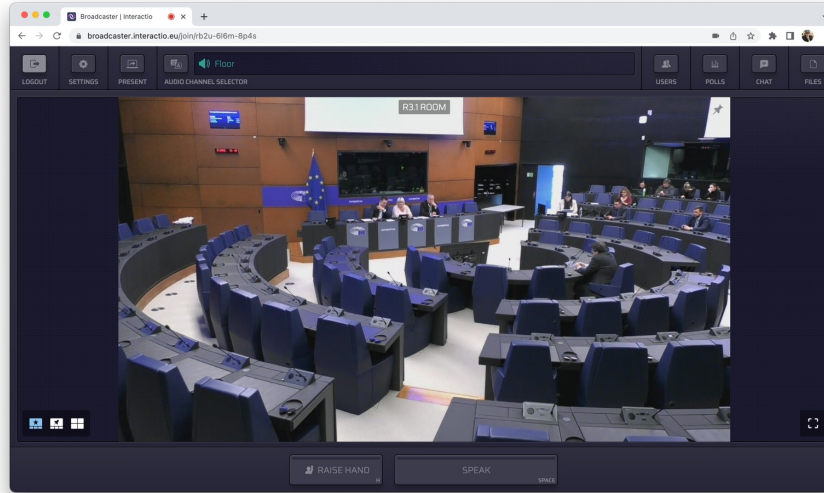
We have invited the UN Special Rapporteur on Minority Issues Dr. Frank La Rue to the occasion of celebrating this year the 30th anniversary of the adoption of the United Nations Declaration of the rights of persons belonging to national or ethnic, religious and linguistic minorities, which continues to be the most important UN instrument devoted to minority rights. He will speak about the role of the Declaration both at international and European level and look together at the future challenges and possibilities for improvement of minority rights and their enforcement.

During the second topic on the agenda Prof. Dr. Georg Rehm, Principal Researcher and Research Fellow at the German Research Centre for Artificial Intelligence (DFKI) in Berlin and Co-ordinator of the ELE project will present the European Language Equality (ELE) Project. With a large consortium consisting of 52 partners covering research and industry and all major pan-European initiatives, the European Language Equality (ELE) EU project develops a strategic research, innovation and implementation agenda with recommendations, as well as a roadmap with concrete steps, for achieving digital language equality in Europe by 2030. The ELE project is a response to the European Parliament resolution "Language equality in the digital age", which was passed by the EP in a landslide vote – 592 votes in favour and only 45 against – in September 2018. The Professor will present the methodology followed in the project, the current state of preliminary results with regard to language barriers in Europe and the stark imbalance in terms of technology support of Europe's languages. Based on these preliminary findings, the ELE project team is currently developing a strategic research, innovation and implementation agenda, to be presented to the European Institutions in June/July 2022.

Interpretation in English, French, Spanish and Hungarian will be available during the meeting.

We look forward to meeting you in Strasbourg.

Best regards,  
Kings Gál  
Francisco Altamir  
Loreda Vinciguerra  
Co-Chairs of the Minority Intergroup

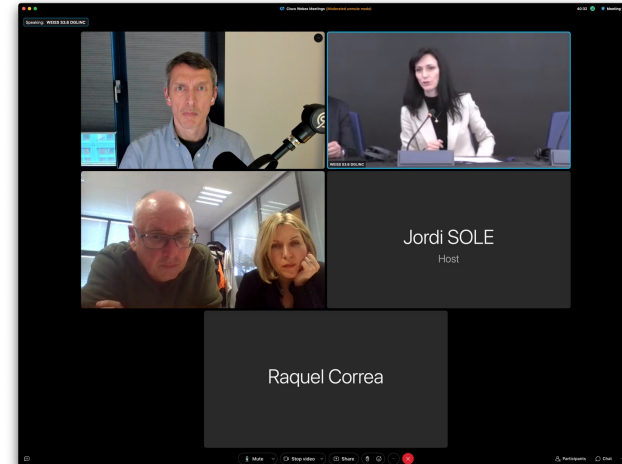


**STOA Workshop**  
Towards full digital language equality in a multilingual European Union  
8 Nov 2022, 09:30-11:30 CET  
Room SPAAK 7C50 & Online via Interactio

**Background**

At least 21 European languages are in danger of digital extinction due to a severe lack of technology support, concluded the [STOA 2012](#) reports prepared by a group of more than 230 experts from all over Europe. For the past decade the introduction of novel technologies (automated translation) has precipitated a revolution in digital language services, allowing for ever faster and more accurate automatic speech recognition (ASR) and machine translation (MT) tools for a wide range of general technology support between the five most spoken EU languages (English, French, German, Spanish and Italian), and the emerging technologies. This digital inequality further increases when regional and minority languages are considered, leading to a dearth of online technological support, both in spoken audio (voice) and written text form. As digital services become an ever-increasing part of our lives, such digital language inequalities could eventually threaten the digital survival of EU languages.

The STOA examines the requirements for research and development investment of language technologies in the context of EU multilingualism. It will present the results of the EU project [Language Equality in the Digital Age \(LEDA\)](#), which proposes a roadmap towards achieving full digital language equality by 2030. A panel discussion will ensue, where policymakers will join experts from academia and the industry to discuss challenges and opportunities for digital language equality in the EU. This event is the third in a series of STOA events on language technologies in the EU: the first, in [2019](#) and the second in [2021](#), a both on STOA 2012 on "Language equality in the digital age", which led to the European Parliament's [Resolution](#) of the same name.



Besprechung mit Commissioner Mariya Gabriel (März 2023)



STOA Workshop "Towards full digital language equality in a multilingual European Union" (November 2022)



# European Language Equality (ELE) 1 und 2

<http://www.european-language-equality.eu>

Koordinator: DCU (ADAPT Centre)

Co-Koordinator: DFKI



META-FORUM 2022 – 8./9. Juni, Brüssel, Belgien



**Konsortium:** 52 Partner (Kerngruppe: DCU – ADAPT Centre, DFKI, Charles University, ILSP, University of the Basque Country)

**Ziel:** *Entwicklung einer strategischen Forschungsagenda, um in Europa bis 2030 für digitale Sprachgerechtigkeit zu sorgen*

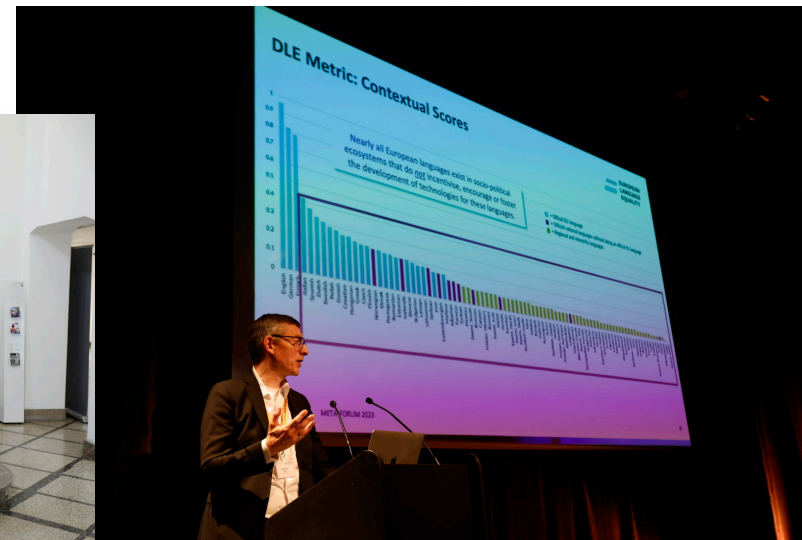
**Laufzeit:** 18 Monate • Januar 2021 – Juni 2022



**Konsortium:** 7 Partner (Kerngruppe + EFNIL und ELEN)

**Ziel:** *Revision und Erweiterung der strategischen Forschungsagenda*

**Laufzeit:** 12 Monate • Juli 2022 – Juni 2023

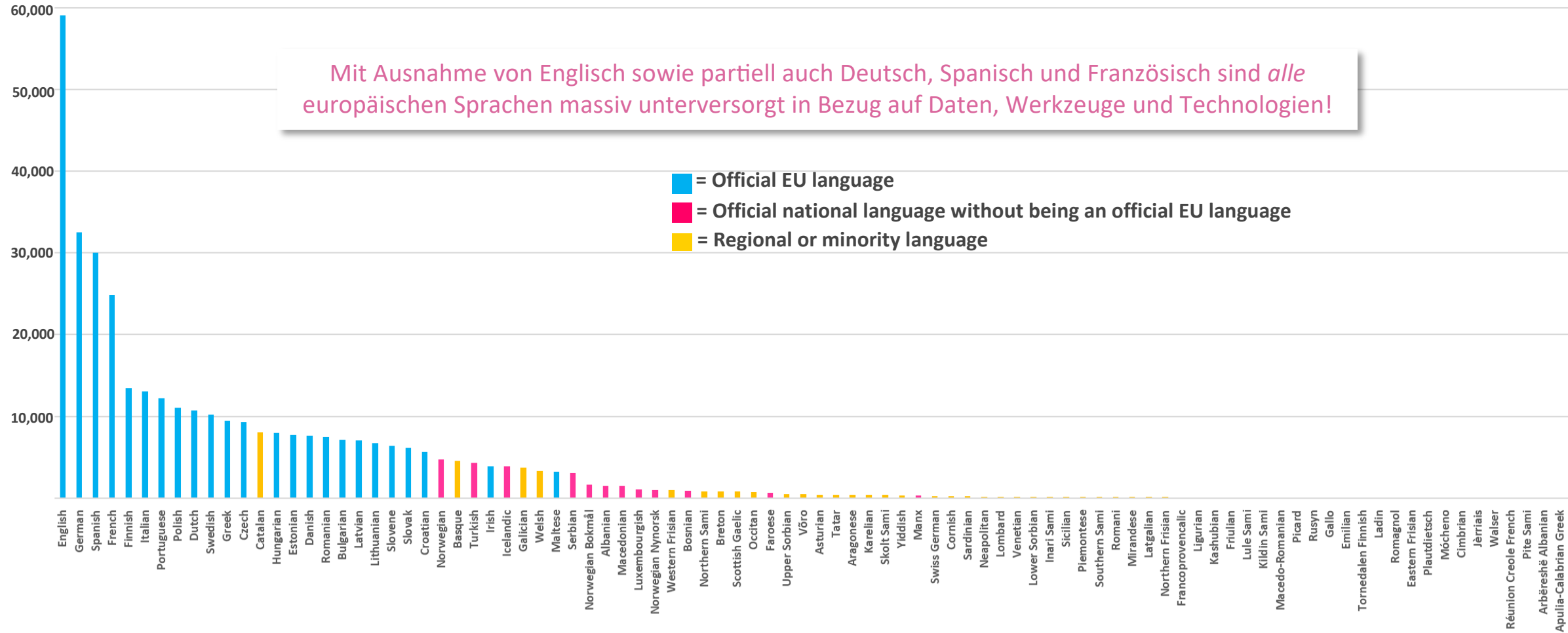


META  FORUM 2023  
27. Juni 2023

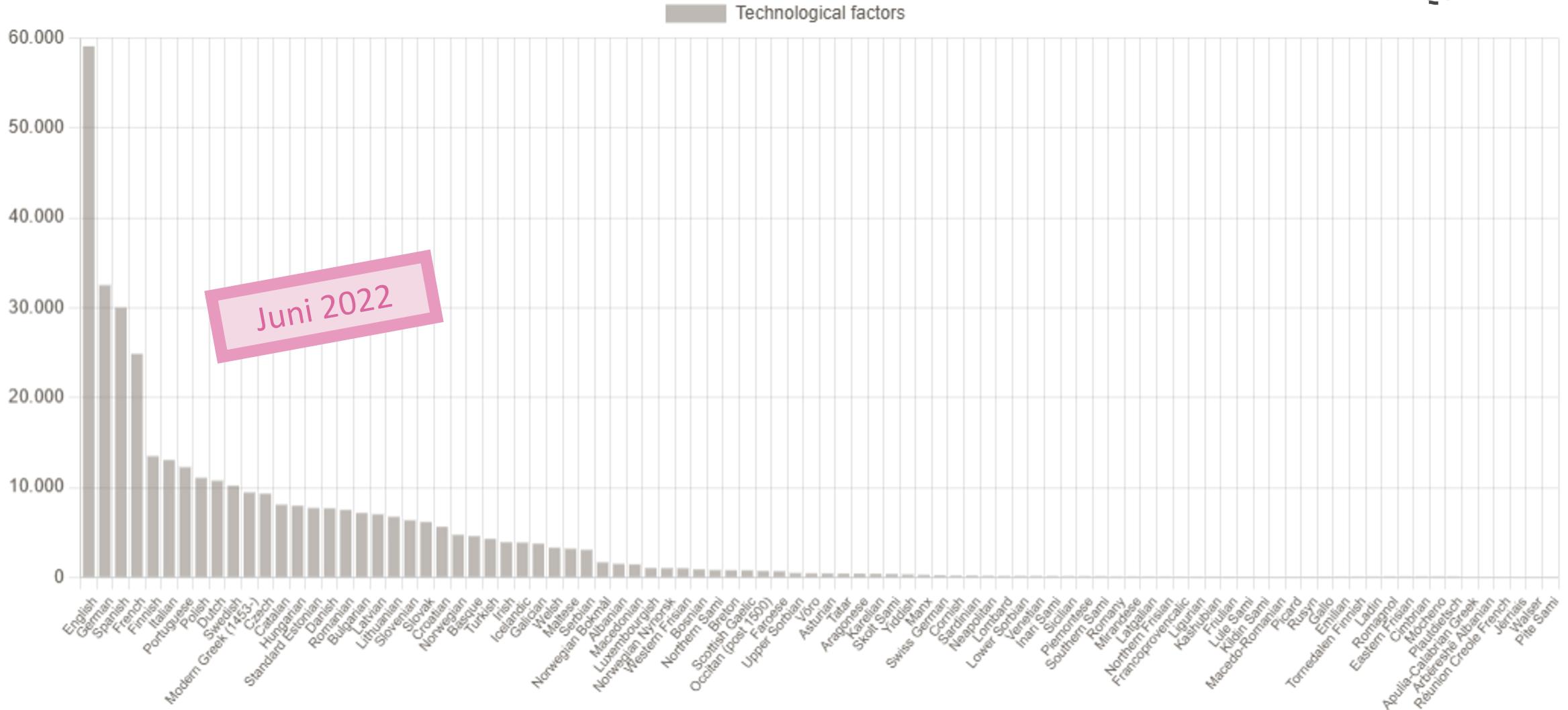
# Digital Language Equality Metrik

Mit Ausnahme von Englisch sowie partiell auch Deutsch, Spanisch und Französisch sind *alle* europäischen Sprachen massiv unterversorgt in Bezug auf Daten, Werkzeuge und Technologien!

- = Official EU language
- = Official national language without being an official EU language
- = Regional or minority language



# Digital Language Equality Metrik: 2022 vs. 2023 (1/3)



Juni 2022



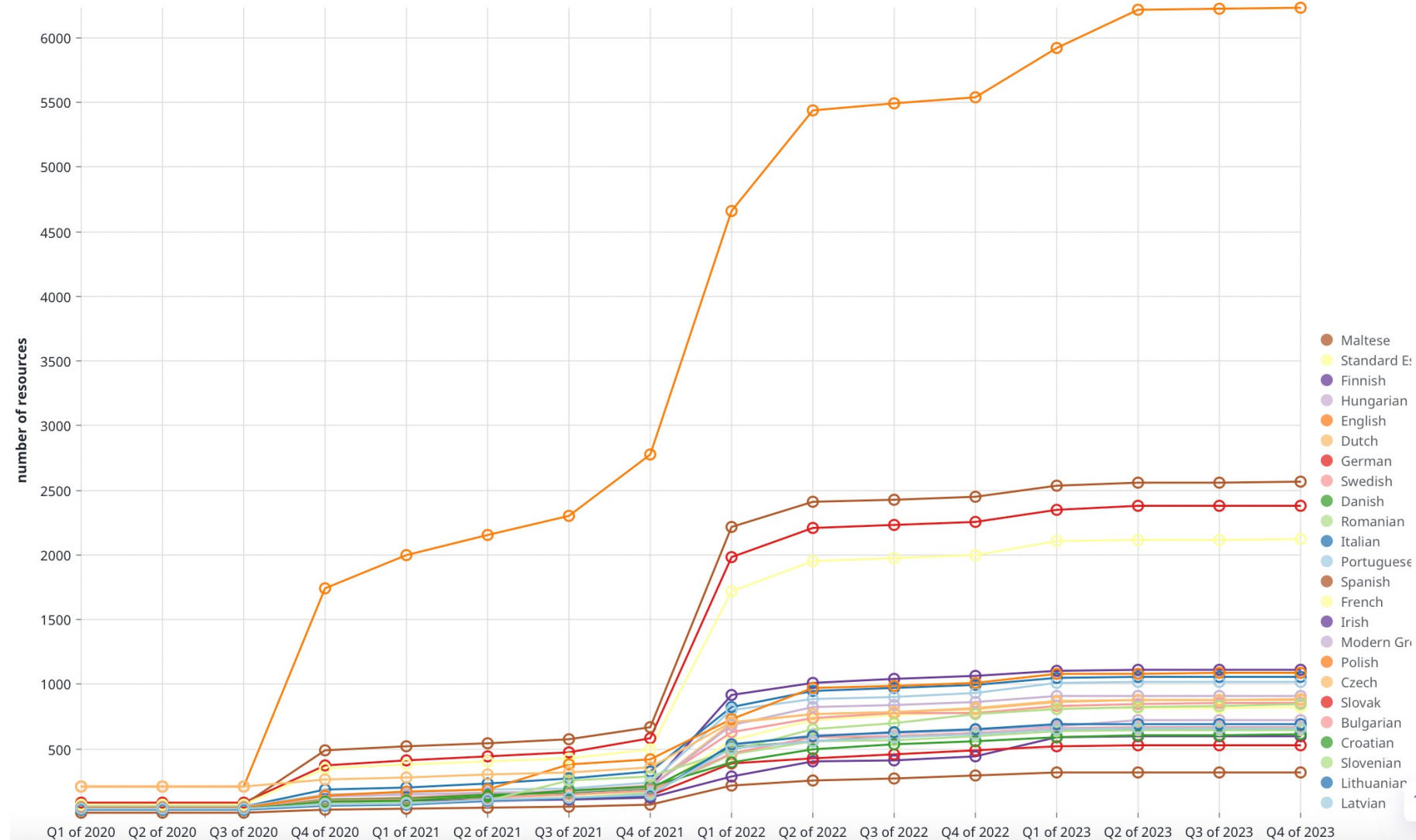


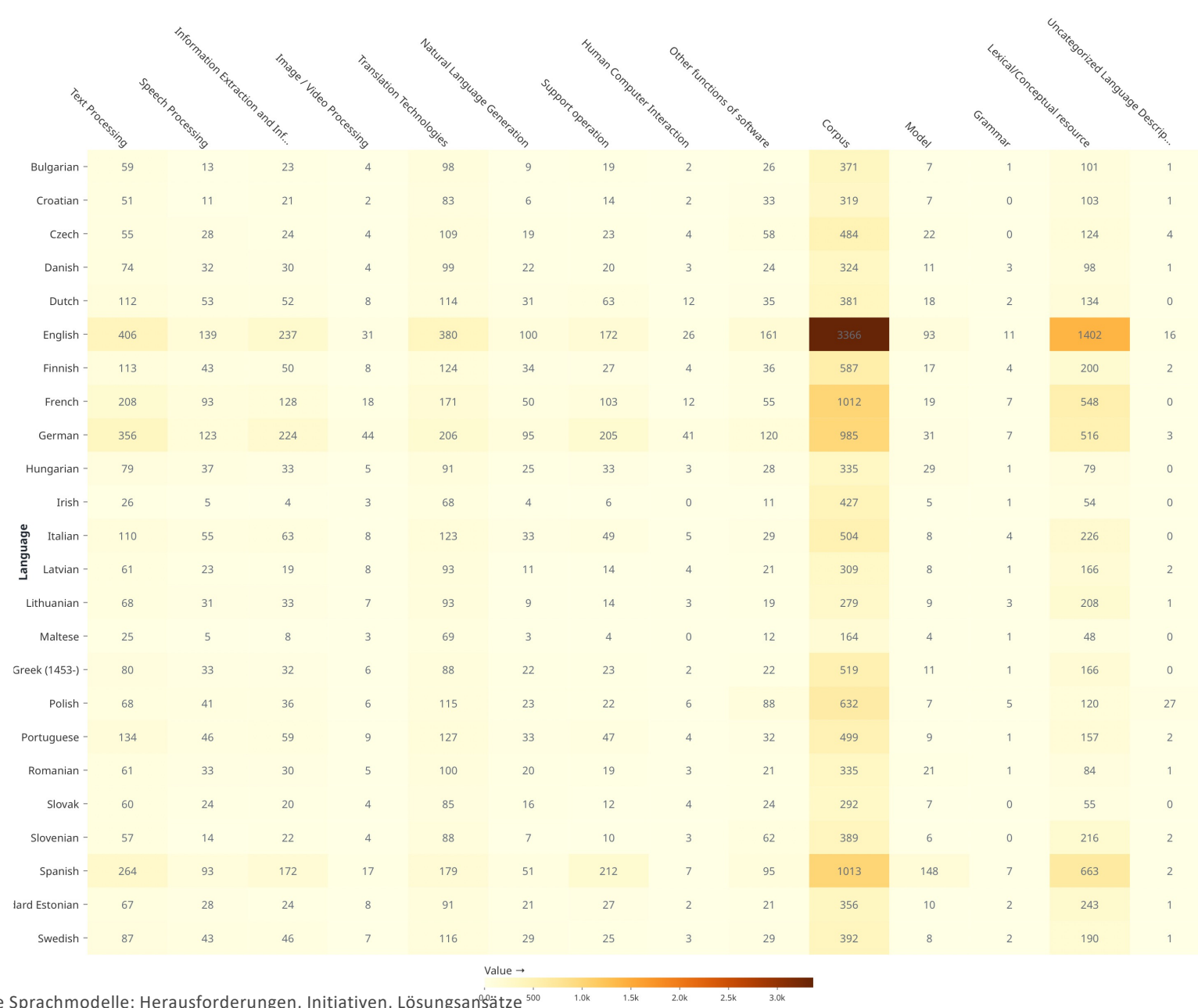
# Digital Language Equality Metrik: 2022 vs. 2023 (3/3)



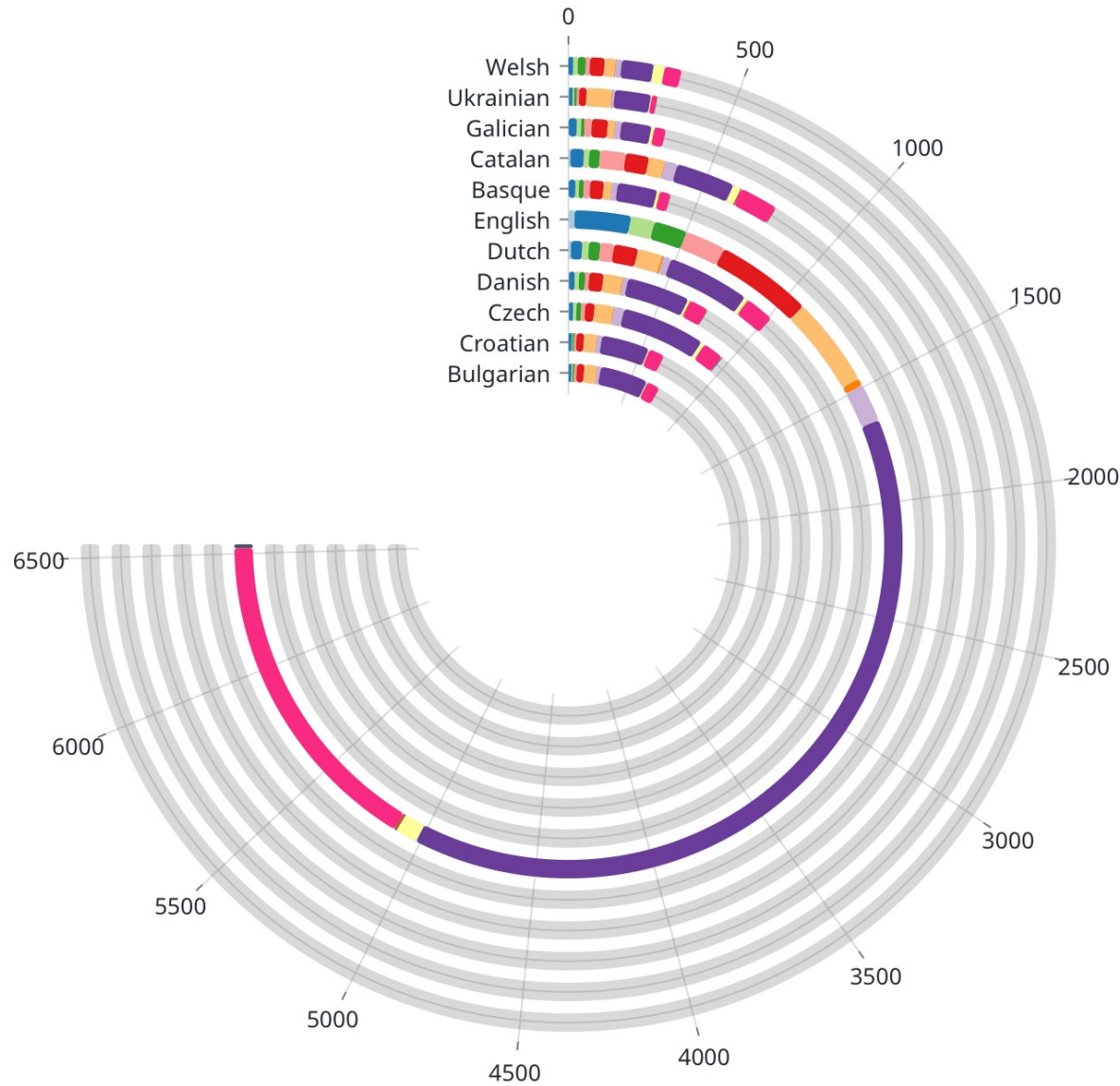
# Zeitliche Entwicklung (alternative Visualisierung)

- Alle Sprachen machen Fortschritte, aber ...
- Englisch macht deutlich mehr Fortschritte im ersten Halbjahr 2023
- Die Lücke wird größer und nicht etwa kleiner!





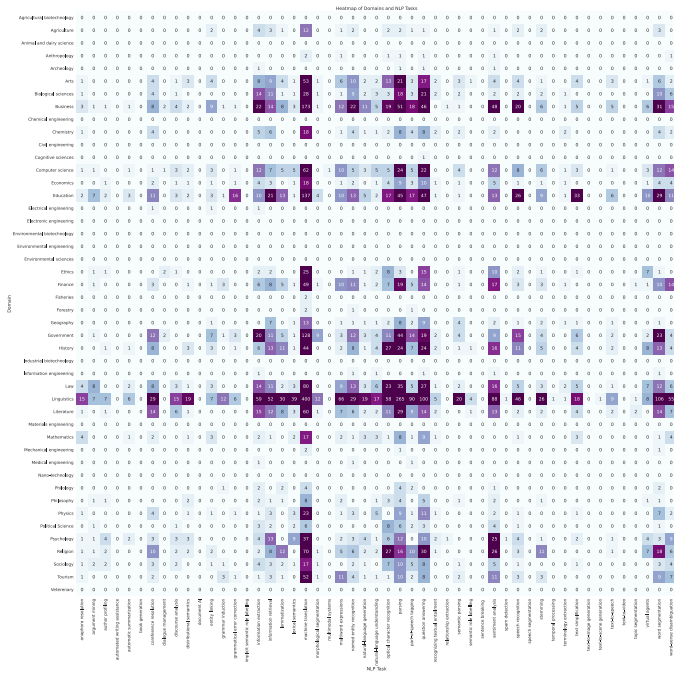
Erstellt mit dem ELE-Dashboard  
<https://live.european-language-grid.eu/catalogue/dashboard>



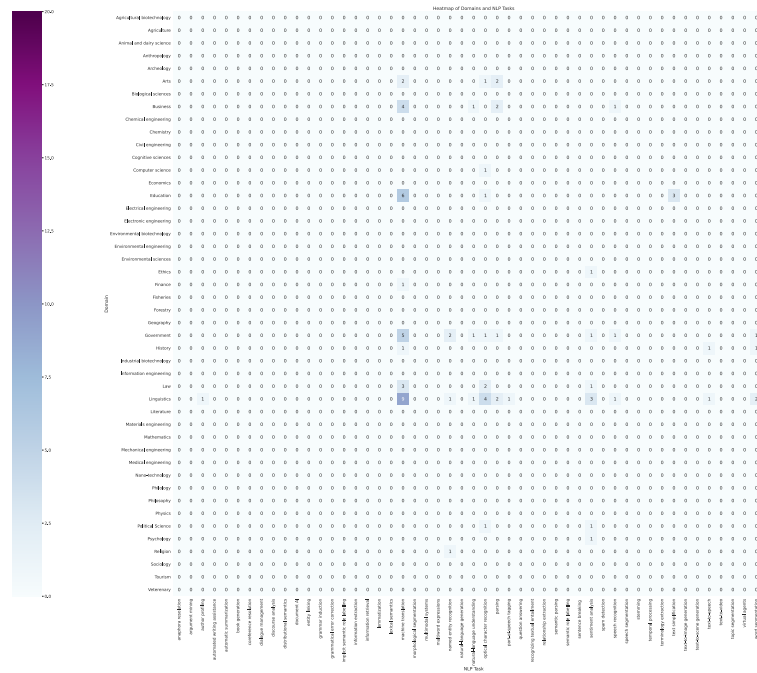
- Human Computer Interaction
- Information Extraction and Information Retrieval
- Natural Language Generation
- Speech Processing
- Support operation
- Text Processing
- Translation Technologies
- Image / Video Processing
- Other functions of software
- Corpus
- Model
- Grammar
- Lexical/Conceptual resource
- Uncategorized Language Description

Erstellt mit dem ELE-Dashboard  
<https://live.european-language-grid.eu/catalogue/dashboard>

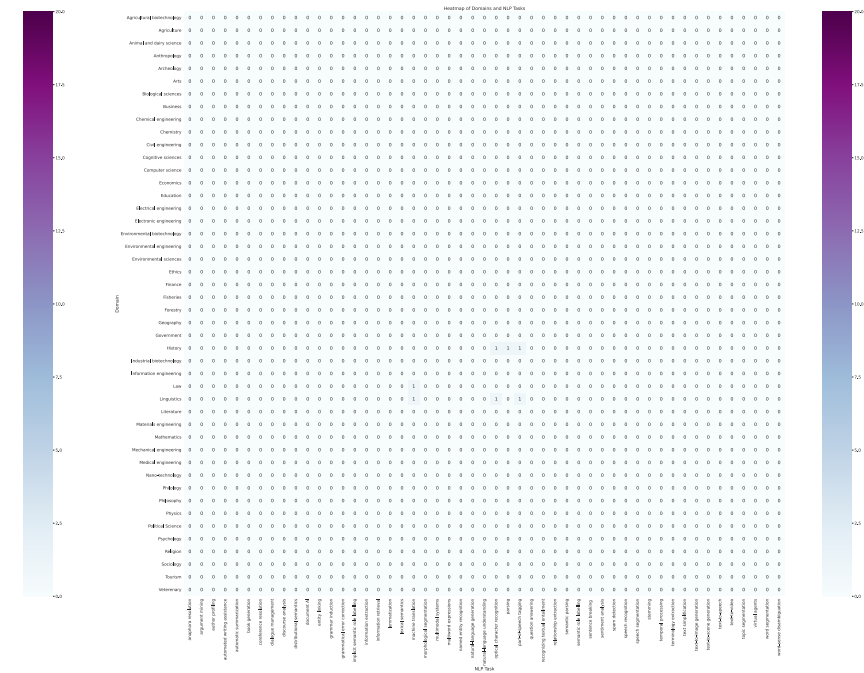
# Digital Language *Inequality* (Analyse der ACL Anthology)



Englisch



Dänisch



Färöisch

Diego Alves, Marko Tadić: “EuLTDom2023 – European LT Domains 2023” (2023), ELE FSTP Projekt

# OpenGPT-X: Große Sprachmodelle für das Deutsche

OpenGPT-X entwickelt große KI-Sprachmodelle, die neue datenbasierte Business-Lösungen ermöglichen und sich speziell an europäische Bedürfnisse richten.



Gefördert durch:



Partner



# OpenGPT-X

- Große Sprachmodelle sind mittlerweile eine Schlüsseltechnologie
  - Die erfolgreichsten LLMs stammen aus den USA sowie China.
  - Rechtslage z.B. bzgl. Verarbeitung von Daten, Server-Standorten etc. oftmals unklar.
  - Die rapide zunehmende Relevanz von LLMs bedeutet für Europa:
    - Digital Souveränität: Unabhängigkeit von Technologien und Daten sicherstellen.
    - Mehr Innovation und Wettbewerbsfähigkeit in Europa.
- OpenGPT-X entwickelt offene und europäische Sprachmodelle.
  - Europäischer Datenschutz wird vollständig berücksichtigt.
  - Steigert Innovation und stärkt Europas Wettbewerbsfähigkeit („LLMs made in Europe“).
  - Förderung durch das BMWK (Januar 2022 bis Dezember 2024).

# OpenGPT-X: Colossal OSCAR

- Mittlerweile ist klar: Die Trainingsdaten sind entscheidend.
- Eine Hauptquelle sind Webdaten (z.B. Webcrawls von CommonCrawl), da sie – vergleichsweise – unkompliziert riesige Datenmengen verfügbar machen (im Petabyte-Bereich).
- Herausforderung: Was sind gute bzw. schlechte Webdaten?
- OSCAR Projekt (Dr. Pedro Ortiz): Annotierte und gefilterte Webdaten für das Training von Sprachmodellen (z.B. für BLOOM von BigScience)
- Neue Veröffentlichung: Colossal OSCAR (10x mehr Daten – bis zu eine Billionen Tokens inkl. Qualitätsannotationen)
- Mehr als 120TB Daten (150+ Sprachen, unkomprimiert)



OSCAR Projekt  
<https://oscar-project.org/>



# OpenGPT-X: Im-datasets

- Neben Webdaten sind kuratierte Datensätze entscheidend für die Fähigkeiten von Sprachmodellen: Qualität statt Masse.
- Kuratierte Datensätze für Sprachmodelle werden bisher nicht einheitlich veröffentlicht (unterschiedliche Formate, Plattformen, Repositorien etc.)
- Das Software-Framework Im-datasets (Dr. Malte Ostendorff) macht viele Datensätze über eine einheitliche Schnittstelle verfügbar.
- Inkl. diverser Werkzeuge für Filterung, Sampling etc.
- Aktuell mehr als 400 Datensätze aus +60 Quellen in 32 Sprachen
- Publikation in Arbeit



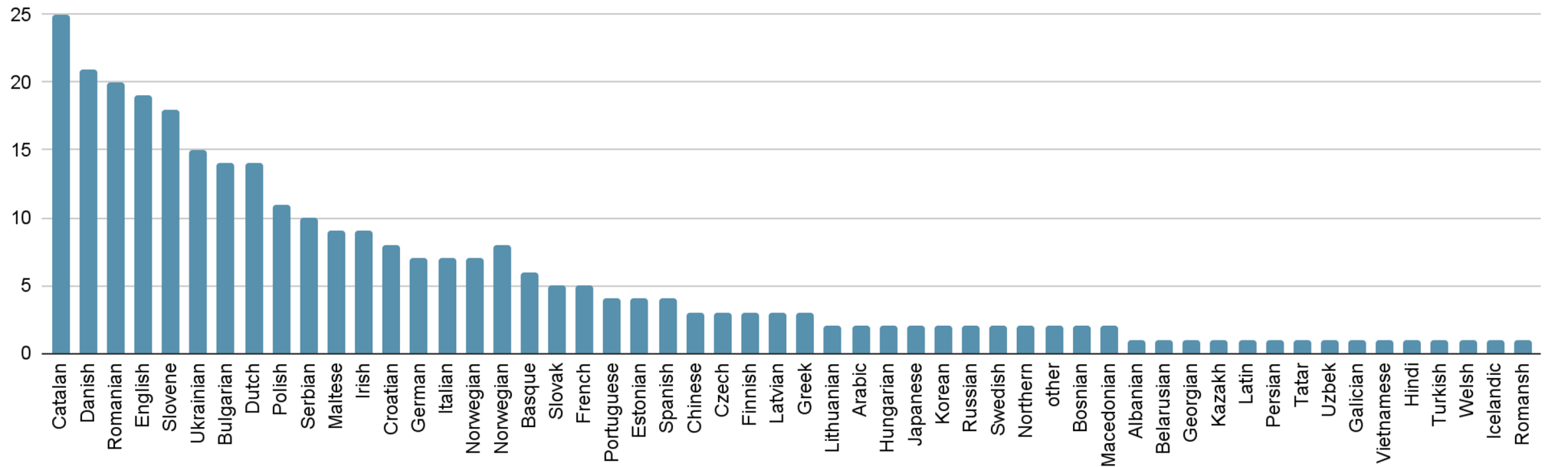
Im-datasets  
[github.com/malteos/Im-datasets](https://github.com/malteos/Im-datasets)



Zusammenarbeit mit unseren europäischen Partnern

# OpenGPT-X: Im-datasets

## Number of Entries by Language

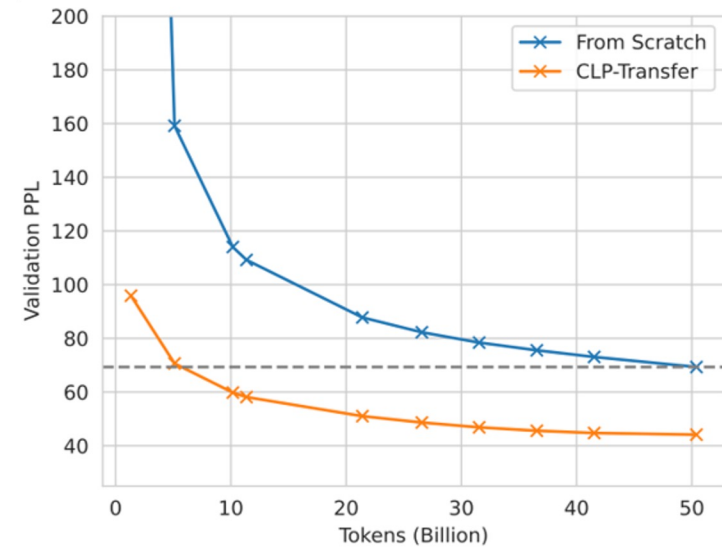


# OpenGPT-X: Transfer Learning

- Das “Open” in OpenGPT-X steht auch für die Zugänglichkeit von LLM-Technologie (geringere Barrieren in Bezug auf Ressourcenbedarf)
- Transfer Learning: Bestehende Ressourcen „recylen“ um die Effizienz zu steigern (z.B. Open Source LLMs wie BLOOM anpassen).
- Cross-lingual & Progressive Transfer Learning [1]:
  - Bestehende Sprachmodelle auf eine neue Sprache anpassen (z.B. Englisch → Deutsch)
- Ergebnis: Bis zu 80% Reduktion des Trainingsaufwands.
- Größtes deutsches Open-Source Sprachmodell (bei Publikation):

[BLOOM-CLP-GERMAN \(6.4B parameters\)](#)

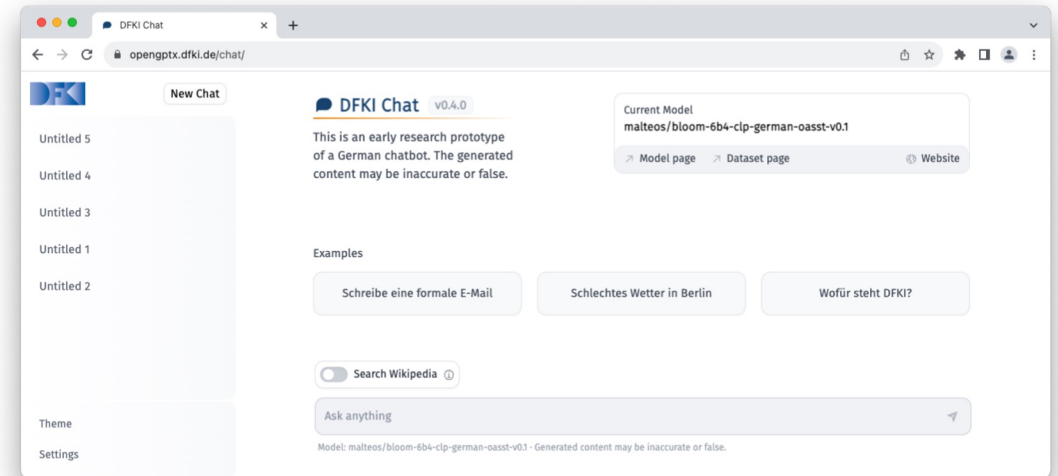
[1] Malte Ostendorff und Georg Rehm. “Efficient Language Model Training through Cross-Lingual and Progressive Transfer Learning”. In: *PML4DC@ICLR 2023*. 05 Mai 2023.



The screenshot shows the 'Text Completion Demo' interface for OpenGPT-X. It features the logos for 'openGPT-X' and 'Deutsches Forschungszentrum für Künstliche Intelligenz GmbH'. Under 'Text completion', the available model is 'BLOOM-CLP German 6.4B'. The 'Your prompt' field contains '(How to use the model?)'. Below this, there is a section to 'Definiere folgende Wörter' with definitions for 'Universität', 'Regierungschef', and 'Wirtschaft'. The 'Examples' section has a text input field with the placeholder 'Enter a prompt or try the examples ...'. The 'Settings' section shows 'Sample mode' selected, with a note that it provides imaginative completions. A 'Generate' button is located at the bottom of the settings. At the bottom of the interface, the 'Model prediction' section is visible with the placeholder text 'Please provide a prompt to get an output...'

# OpenGPT-X: Transfer Learning

- Zentrale Funktion eines Sprachmodells: Vorhersage des wahrscheinlichsten nächsten Wortes (next word prediction).
- Chatbots oder virtuelle Assistenten (ChatGPT etc.) sind Sprachmodelle, die für Chat optimiert wurden.
- DFKI Chat:
  - Optimierte Variante von BLOOM-CLP-German (siehe vorherige Folie).
  - Integration von Retrieval-Augmented Generation
  - Bei Bedarf können aktuelle Zusatzinformationen aus Wikipedia genutzt werden.



DFKI Chat  
<https://opengptx.dfki.de/chat>

# Daten

# Existierende LT/NLP-Plattformen, -Infrastrukturen, -Repositorien

**META-SHARE** LEARN DISCOVER PARTICIPATE CONNECT LOGIN

Search & exchange language resources

META-SHARE is an open and secure network of repositories for sharing and exchanging language data, tools and related web services

Share your own resources!

[JOIN OUR NETWORK NOW](#)

Already a member? [Log in](#)

Search the META-SHARE inventory

OR [LEARN MORE](#)

4,481 users | 2,887 language resources | 32% text corpora | 27,630 number of downloads

Virtual Language Observatory Search Contributors Help CLARIN

## CLARIN Virtual Language Observatory

Welcome to the VLO!

Use the **search bar** below to start searching through hundreds of thousands of language resources, or **continue** to browse everything and use **facets** to narrow down to your area of interest or discover new resources.

[See all records](#) [Take a quick tour](#)

Search through 1,030,321 records

European Language Resource Coordination  
Coordinating Europe Facility  
ELRC-SHARE

## ELRC-SHARE Repository

Type in your keywords, please...

Welcome to the ELRC-SHARE repository!

The ELRC-SHARE repository is used for documenting, storing, browsing and accessing [Language Resources](#) that are collected through the [European Language Resource Coordination](#) and considered useful for feeding the [CEF Automated Translation \(CEF AT\)](#) platform.

If you want to contribute resources, all you have to do is [register](#) (new user) or [login](#) (returning user) and go on to describe and upload your data with a simple form.

EUROPEAN LANGUAGE RESOURCE ASSOCIATION

1096 Language Resources (Page 1 of 55)

2006 CoNLL Shared Task - Arabic & Czech

Arabic, Czech  
ID: ELRA-W0087  
ISLRN: 798.485.294.792.1

2006 CoNLL Shared Task - Arabic & Czech consists of dependency treebanks used as part of the CoNLL 2006 shared task on multi-lingual dependency parsing. The Conference on Computational Natural Language Learning (CoNLL) is accompanied every year by a shared task intended to promote natural lan...

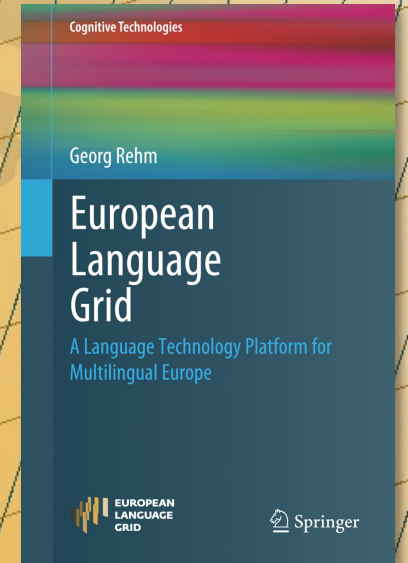
2006 CoNLL Shared Task - Ten Languages

Bulgarian, Danish, Dutch, Finnish, German, Japanese, Portuguese, Slovenian, Spanish, Croatian, Swedish, Turkish  
ID: ELRA-W0088  
ISLRN: 978.227.532.844.8

Discover, try out, use and download LT services and resources for all European languages.

Browse ELG and find the LT services, resources, developers and providers you are looking for.

Search



ELG RELEASE 3.0 (DECEMBER 2023) — GRID WORKING — NO MAINTENANCE SCHEDULED

**8015**  
Corpora

**3894**  
Tools & Services

**2822**  
Conceptual Resources

**511**  
Models & Grammars

**1779**  
Organizations

**514**  
Projects

# EU Datenstrategie und Datenräume (Data Spaces)

- Data Spaces sind inhärenter Teil der EU-Datenstrategie
- Data Spaces sind auch ein Instrument, um eine Datenwirtschaft in Europa zu etablieren
- Diese Daten können für diverse Zwecke verwendet werden, u.a. auch für das Training von LLMs
- In Europa existieren verschiedene Initiativen bzgl. Datenwirtschaft und Dateninfrastrukturen, die ähnliche Ziele und Positionen haben, sich aber teilweise unterscheiden:
  - Data Spaces Business Alliance (DSBA): Gaia-X, IDSA, FIWARE, BDVA – *Konvergenz!*
  - EU: DSSC (incl. DSBA), SIMPL, 14 offizielle EU Data Space-Projekte
- Der Common European Language Data Space ist einer dieser 14 offiziellen EU Data Spaces.



# Common European Language Data Space

- Procurement Projekt (CNECT/LUX/2022/OP/0026)
- Budget: 6M€ (+ 2M€, sofern die EU die Option der Verlängerung wählt)
- Laufzeit: 36 Monate (+ 12 Monate)
- Ziel: Entwicklung und Bereitstellung einer europäischen Plattform mit Marktplatz für die Sammlung, Erstellung, Teilen und Nutzung mehrsprachiger und multimodaler Sprachdaten
- Wichtige Komponenten: Governance, Architektur, Infrastruktur, Offenheit, Dissemination
- Stakeholder: Industrie, Forschung, öffentliche Verwaltung, Kultureinrichtungen, NGOs etc.



**EUROPEAN  
LANGUAGE  
DATA SPACE**

# Konsortium und Unterauftragnehmer

Koordinator		
Deutsches Forschungszentrum für Künstliche Intelligenz GmbH	DFKI	DE
Partner		
R.C. "Athena", Institute for Language and Speech Processing	ILSP	GR
Evaluations and Language Resources Distribution Agency	ELDA	FR
TILDE	TILDE	LV
Unterauftragnehmer		
3pc GmbH Neue Kommunikation	3pc	DE
Capgemini Deutschland GmbH	CapG	DE
CLARIN ERIC	CLARIN	NL
Big Data Value Association (Data, AI and Robotics) AISBL	BDVA	BE

Außerdem: Rechtsexperten (Delcade, Frankreich) und ca. 30 Organisationen, die die vorgesehenen Workshops in den Ländern organisieren.

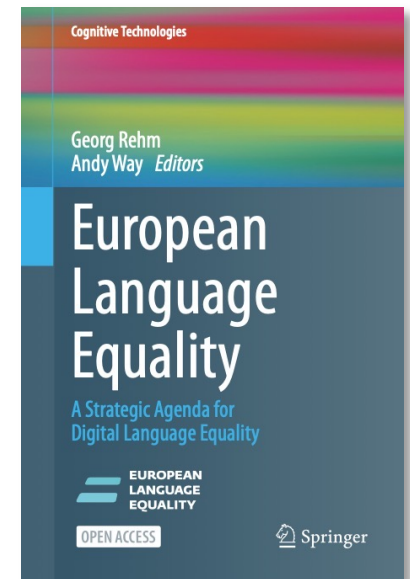
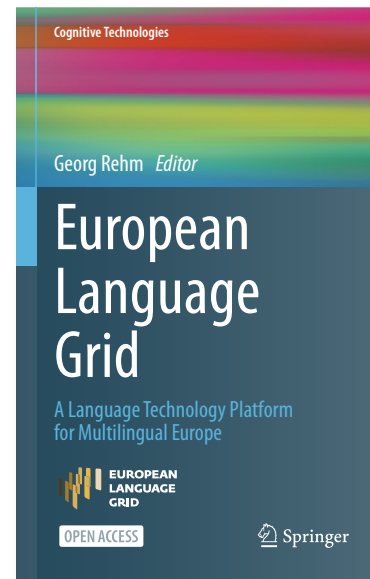
# Frühere Projekte und Initiativen

- Die Kerngruppe – DFKI, ILSP, ELDA, TILDE – arbeitet seit Jahren in diversen Projekten bzgl. Infrastruktur und Plattformen zusammen:
- **META-NET** (FP7, 2010-2013)
  - META-SHARE
- **ELRC** (CEF, 2014-2023)
  - ELRC-SHARE
- **ELG** (H2020, 2019-2022)
  - ELG Cloud Platform
- **ELE** (PP/PA, 2021-2023)

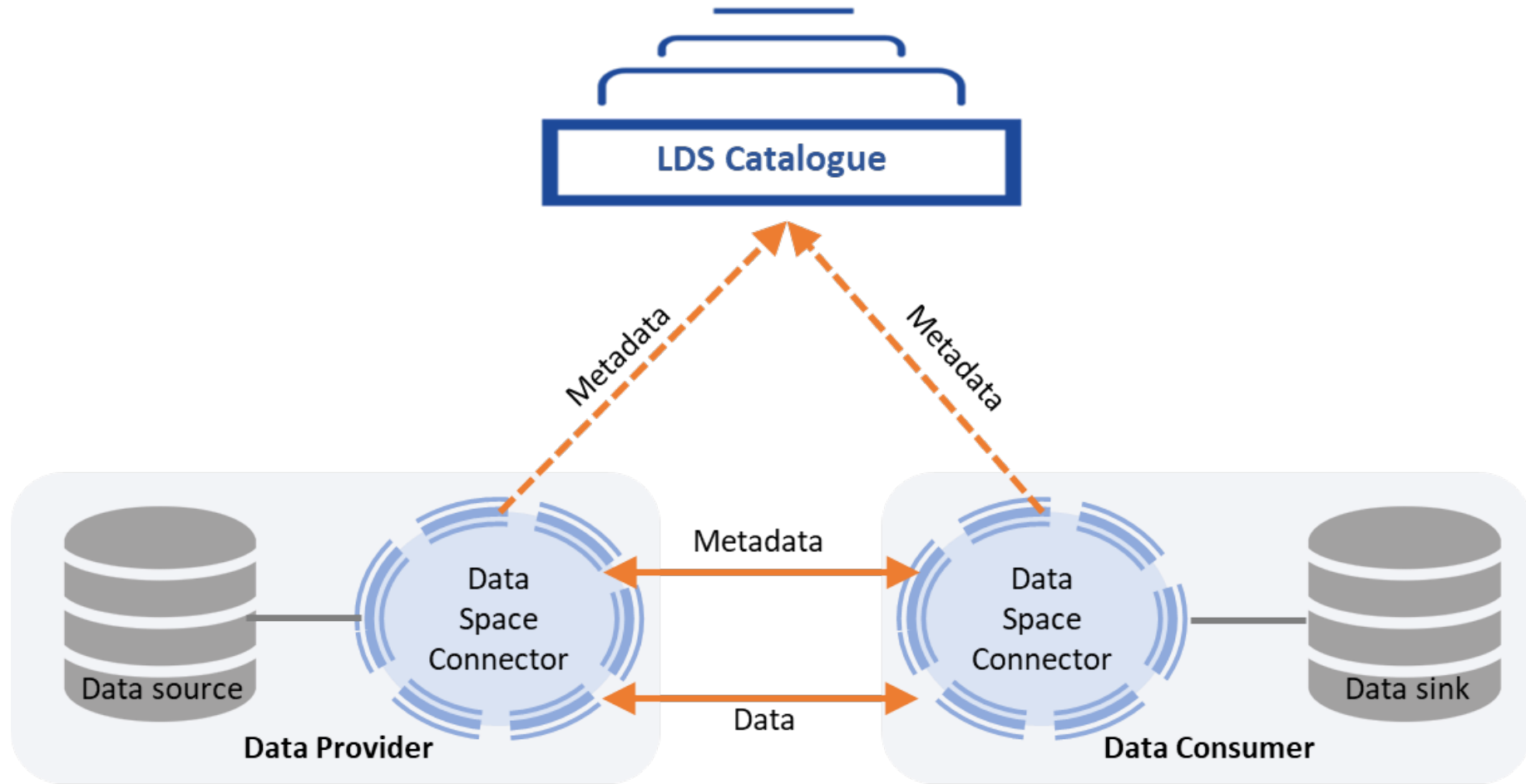
**META**  **NET**



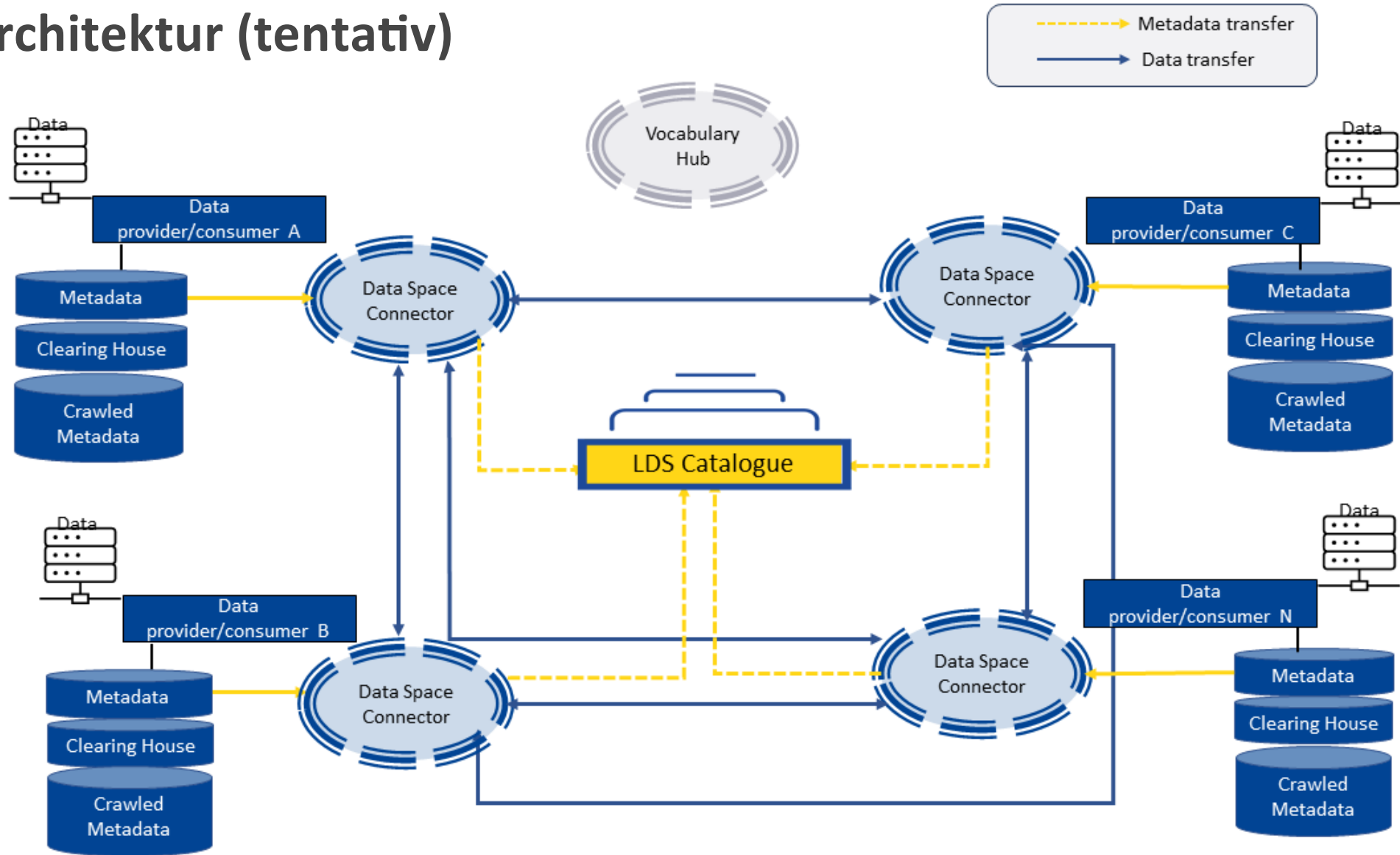
Die **dezentrale LDS-Infrastruktur** wird auf Erfahrungen aus ELG, ELRC-SHARE und META-SHARE basieren.



# LDS Datenflüsse



# LDS Architektur (tentativ)



# Unterschiedliche Klassen von Daten

Klasse	Typische Größe	Anbieter	Integration in LDS	Relevanz für LLMs
Reguläre Korpora & Sprachressourcen	Klein (MB, GB)	Primär NLP/LT-Forschung: ELG, META-SHARE, CLARIN, ELRA, ELDA etc.	Können einfach integriert werden, indem die Repositorien an LDS angebunden werden	Typischerweise qualitativ sehr hochwertige Daten und daher relevant für LLMs, aber nicht als "Basisdaten" für das Training
Web-Crawls	Sehr groß (TB, PB)	Common Crawl (und OSCAR-Dumps), Internet Archive-Dumps etc.	Herausforderung ist die Größe (aufwendig zu kopieren, zu prozessieren, zu speichern; müssen sehr nah an HPC sein)	Unabdingbar auf Grund ihrer Größe und Abdeckung – aber: viel Rauschen, großer Bedarf für Vorverarbeitung und Filterung
Neue, frische Daten aus der Industrie sowie von anderen Organisationen	Beliebige Größe, idealerweise so groß wie möglich	Verlage, Medienanbieter, Bibliotheken, Call-Centres, Radio- und Fernsehanstalten, andere Data Spaces	Können einfach integriert werden, indem die Organisationen an LDS angebunden werden	Insbesondere hochqualitative Daten oder domänenspezifische Daten oder Daten, die bestimmte Sprachen abdecken, daher sehr relevant für LLMs

# Wissenschaftliche Daten

# NFDI for Data Science and Artificial Intelligence

- NFDI: Nationale Forschungsdateninfrastruktur
- NFDI4DS: NFDI for Data Science and Artificial Intelligence (Laufzeit 5 + 5 Jahre)
- NFDI4DS entwickelt eine Forschungsdateninfrastruktur für DS und KI
- Basis: Existierende Infrastrukturen und Werkzeuge, die interoperabel gemacht und anschließend verbunden werden (inklusive u.a. European Language Grid)
- Langfristig wird die NFDI4DS-Infrastruktur mit der restlichen NFDI-Infrastruktur verbunden
- DFKI co-koordiniert zwei der fünf Task Areas:
  - Infrastructure and Services
  - Data Science and AI – Transfer and Application
- Forschungsthemen: Scholarly Information Processing, Informationsextraktion, Erstellung von Research Knowledge Graphs (RKGs)



# NFDI for Data Science and Artificial Intelligence



- Es existieren diverse Initiativen zur Repräsentation wissenschaftlichen Wissens in Form von Knowledge Graphs (Research Knowledge Graphs, RKGs), z.B. ORKG (TIB Hannover)
- In NFDI4DS konzentrieren wir uns auf die Verarbeitung wissenschaftlicher Publikationen
- Domäne Computerlinguistik, d.h. Publikationen der Association for Computational Linguistics (ACL)
- Aktuell zwei Ziele:
  - Klassifikation beliebiger wissenschaftlicher Artikel in ihre Fachrichtungen
  - Extraktion von Informationen aus Publikationen und Integration dieser Informationen in RKGs
- Vision: Robuste Extraktion derartiger Informationen aus beliebigen wissenschaftlichen Artikeln und ihre Integration in wissenschaftliche Wissensgraphen
- *RKGs repräsentieren wissenschaftliche Ergebnisse, Daten und Artefakte explizit (symbolisch)*
- *LLMs repräsentieren wissenschaftliche Ergebnisse, Daten und Artefakte implizit (subsymbolisch)*
- Ein hybrider Ansatz erscheint sinnvoll, d.h. *LLMs für RKGs* und *RKGs für LLMs*

# Schlussfolgerungen

# Zusammenfassung

- LLMs sind nicht für Europas Sprachen und ihre linguistischen und kulturelle Bedürfnisse optimiert.
- Das DFKI arbeitet in diversen Projekten an europäischen Sprachdaten und Sprachmodellen.
- Erste Modelle wurden bereits veröffentlicht – weitere folgen Anfang 2024.
- Der Language Data Space läuft auf Hochtouren: Technologie, Dissemination, Governance etc.
- Zusammenarbeit mit allen relevanten Initiativen: DSSC, HPLT, EuroHPC, OpenWebSearch etc.
- Von zentraler Bedeutung: Nutzung des LDS durch Industrie und andere Organisationen, die in sinnvoller Weise Daten bereitstellen und teilen können sowie auch wollen
- Ziel: Identifizierung und Bereitstellung neuer und „frischer“ Sprachdaten, speziell aus der Industrie sowie von anderen Organisationen, die interessante Daten haben, und idealerweise auch viele bzw. alle europäischen Sprachen sowie auch alle Modalitäten abdecken

# Zusammenfassung

- Modelle:
  - Diverse LLM-Forschungsprojekte in Europa – jeweils mit zwei Jahren administrativem Vorlauf
  - Diverse LLM-Firmen in Europa versuchen, mit OpenAI zu konkurrieren – Herausforderung!
- Daten:
  - In Europa ist digitale Sprachgerechtigkeit noch in weiter Ferne
  - Common European Language Data Space
  - Erste nationale Language Data Spaces entstehen, z.B. in Spanien
- Wissenschaftliche Daten:
  - Diverse RKGs/SKGs werden aktiv mit Inhalten gefüllt
- *Erstmals institutionelle europäische Kollaboration!* (ab Anfang 2024):
  - ALT-EDIC

# Alliance for Language Technologies EDIC (ALT-EDIC)

- European Digital Infrastructure Consortium (EDIC): Ein neuer Typ juristische Person in der EU
- 15+ EU-Mitgliedsstaaten sind aktuell an diesem Multi-Country Project (MCP) beteiligt
- Koordiniert vom französischen Ministerium für Kultur
- Enge Zusammenarbeit zwischen: ALT-EDIC Arbeitsgruppe, EC, LDS
- ALT-EDIC Aktionsplan konzentriert sich auf:
  1. Daten;
  2. Existierende und neue Sprachmodelle;
  3. Evaluatierung, Zertifizierung, Normalisierung;
  4. Ökosystem;
  5. EDIC-Implementierung

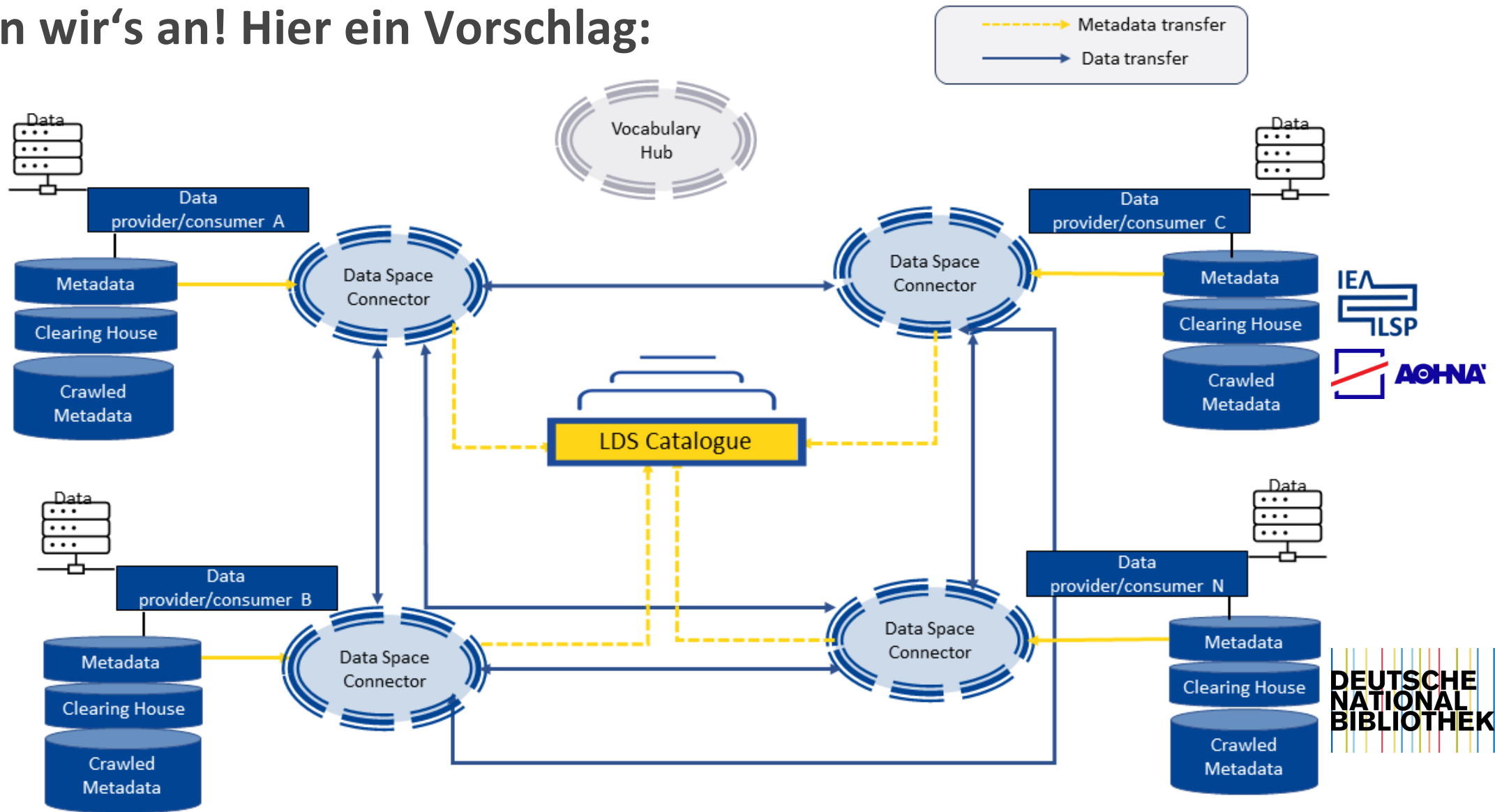
# KI in Bibliotheken – Eine Vision – Schlussfolgerungen

**Mittels ScienceGPT mit der wissenschaftlichen Literatur sprechen,  
mit allen wissenschaftlichen Publikationen in den Dialog treten,  
diese aktiv konsultieren, Forschung gemeinsam mit den Quellen betreiben!**

- In dieser Vision liefern Bibliotheken nicht mehr „nur“ Hinweise auf relevante Literatur (und die Literatur selbst), sondern Nutzerinnen können sich aktiv mit der wissenschaftlichen Literatur „unterhalten“, extrem effizient recherchieren und kreativ inspirieren lassen.
- Ist das realistisch? Ja!
- Werden sich Bibliotheken grundsätzlich transformieren ... müssen? Vermutlich!
- Durch systematische Datenbereitstellung und Datennutzung können Bibliotheken und wissenschaftliche Literatur eine Schlüsselrolle für die nächsten Durchbrüche bei LLMs spielen.

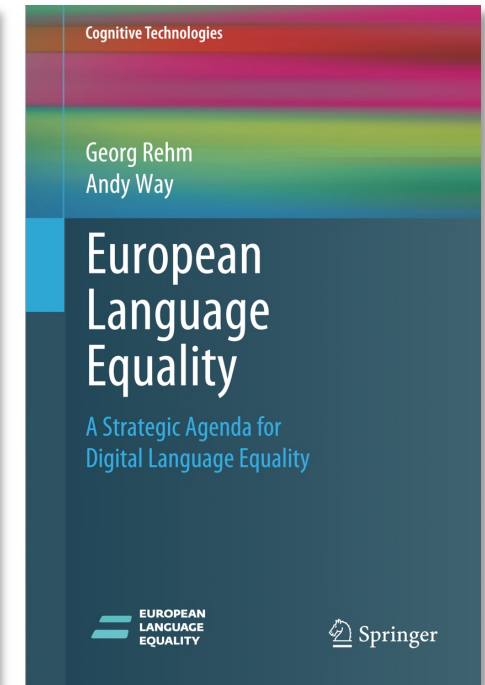
# Packen wir's an! Hier ein Vorschlag:

dfki  
ai





**Herzlichen Dank!**



A Common European Language Data Space – funded under contract LC-01936389 with the European Union.

Prof. Dr. Georg Rehm (DFKI GmbH)  
georg.rehm@dfki.de

07.12.2023 – KI in Bibliotheken: Neue Wege mit großen Sprachmodellen? – DNB, Frankfurt  
<https://language-data-space.ec.europa.eu>