

Selected Infrastructure Activities around Language Data at DFKI

Prof. Dr. Georg Rehm

DFKI GmbH and Humboldt-Universität zu Berlin

georg.rehm@dfki.de



06 June 2024 – Inria and ECDF Collaboration Event – Berlin, Germany

Outline

- European Language Data Space (LDS)
- European Language Grid (ELG)
- NFDI for Data Science and Artificial Intelligence (NFDI4DS)

European Language Data Space

Relevant Observations

- Large language models are the most disruptive breakthrough in AI in recent history (BERT, GPT-3, ChatGPT, GPT-4 etc.)
- LLMs are trained on vast amounts of training data (language data)
- LLMs use dozens, some even hundreds of terabytes (trillions of tokens) of language and also image, video, audio etc. training data
- Europe's languages are vastly under-resourced, except English
- In 2024, language data is more relevant than ever
- A concerted effort for the collection of (very large amounts of) language data for all European languages is very much needed
- The global NLP/LT/Gen-AI market is in the hundreds of billions of US-\$ already and expected to grow to 439.85B US-\$ by 2030 – no significant players from Europe
- Already now billions and billions are made but ...

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.

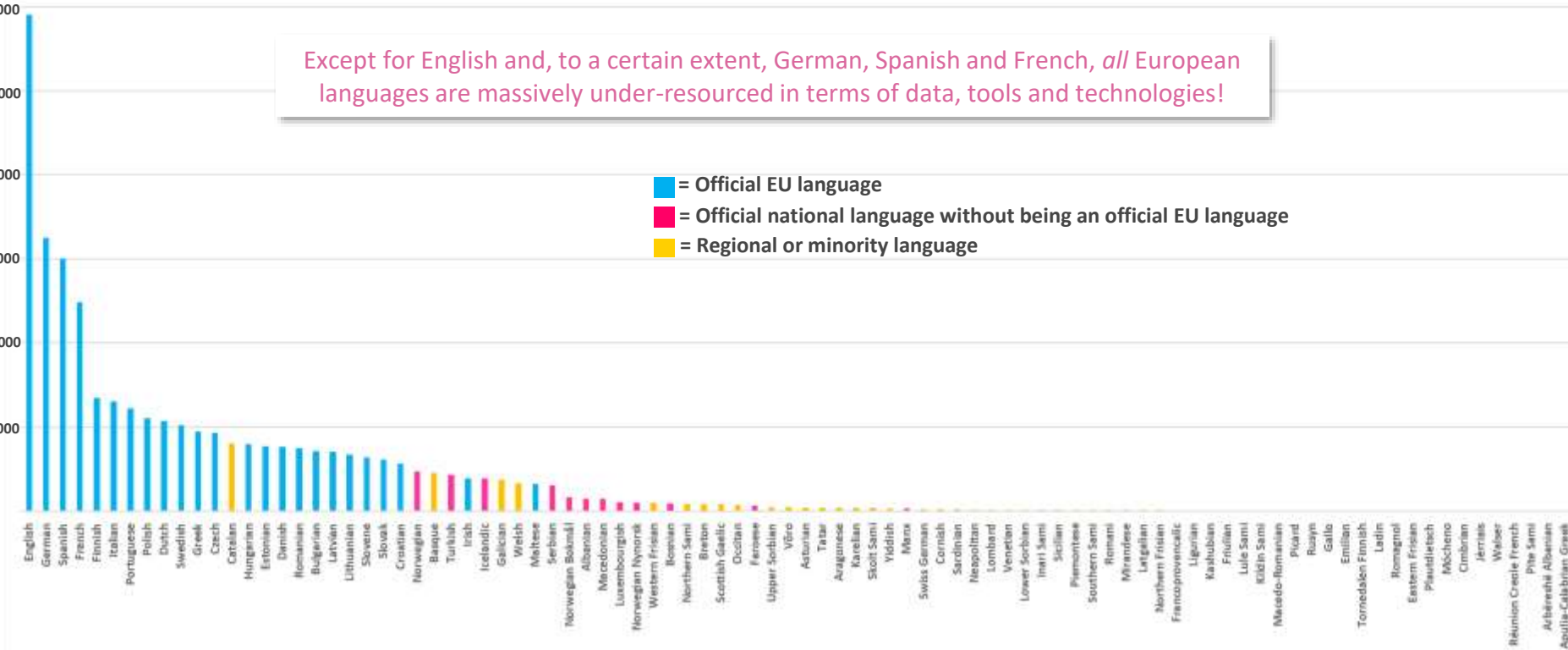
European Initiatives

- European initiatives for the development of LLMs
 - Large research projects in almost every country, e.g., Spain, Denmark, Italy, Germany etc.
 - Companies in many countries, e.g., Finland (Silo.ai), France (Mistral), Germany (Aleph Alpha)
 - EU and nationally funded projects, e.g., HPLT, TrustLLM
 - New pan-European initiative: ALT-EDIC
- Challenges:
 - HPC facilities
 - Speed of the big tech players in the US and Asia vs. speed of Europe
 - Availability of data for European languages

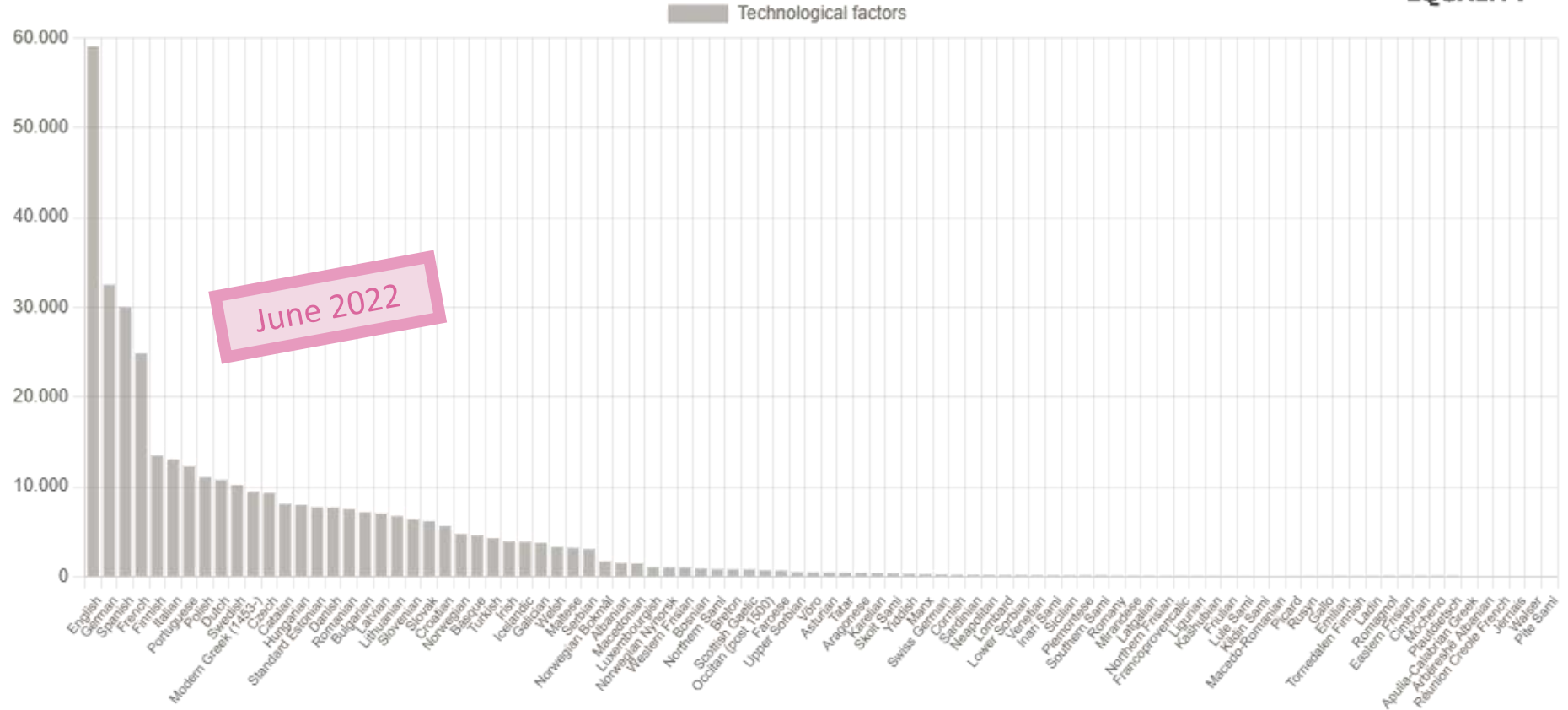
Digital Language Equality Metric: Technological Scores

Except for English and, to a certain extent, German, Spanish and French, *all* European languages are massively under-resourced in terms of data, tools and technologies!

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



DLE Metric: 2022 vs. 2023 (1/3)



EU Data Strategy & Data Spaces

- Data Spaces are an inherent part of the EU Data Strategy
- Data Spaces will help to establish a data economy in Europe
- Various data economy and data infrastructure initiatives in Europe with slightly different goals and individual positioning but conceptual, technical, legal and operational overlap:
 - Data Spaces Business Alliance (DSBA): Gaia-X, IDSA, FIWARE, BDVA
 - EU: DSSC (incl. DSBA), Simpl, approx. 20 data spaces eventually
- The Common European Language Data Space is one of the 14 official EU data space projects with a strong focus on industry

Common European Language Data Space



- Type of action: procurement (CNECT/LUX/2022/OP/0026)
- Budget: 6M€ (+ 2M€ if renewed)
- Runtime: 36 months (+ 12 months if renewed)
- Objective: Develop and deploy a European platform and marketplace for the collection, creation, sharing and re-use of multilingual and multimodal language data
- Salient features: governance framework, technical architecture and infrastructure, openness, promotion
- Stakeholders: industry, research, public administration, cultural associations, NGOs and citizens

Consortium and Subcontractors

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

Plus legal experts (Delcade, France) and approx. 30 organisations for the logistics of multiple country workshops

Previous Projects and Initiatives

- The four core partners – DFKI, ILSP, ELDA, TILDE – have been involved in many projects, including:
- **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)

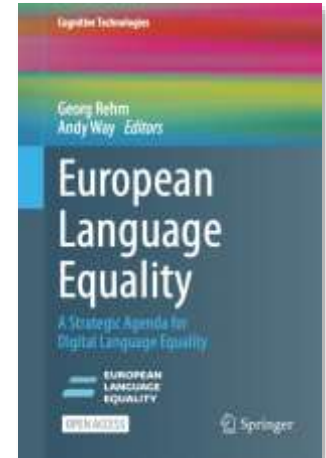
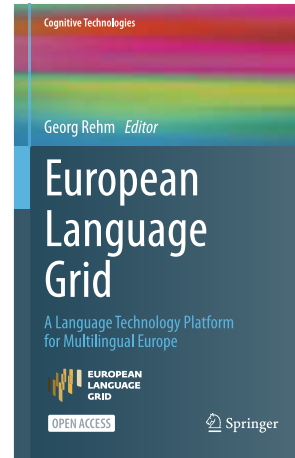
META NET



**EUROPEAN
LANGUAGE
GRID**

**EUROPEAN
LANGUAGE
EQUALITY**

The **technical development work in LDS** will be informed by ELG, ELRC-SHARE, META-SHARE.



Classes of Data

Class of Data	Typical Size	Providers	Integration into LDS	Relevance for LLMs
Regular Corpora and Language Resources	Small (MB, GB)	Primarily NLP/LT research: ELG, META-SHARE, CLARIN, ELRA, ELDA etc.	Can be easily integrated by connecting the repositories to LDS	Usually very high quality data and thus relevant for LLMs but not as base data
Web Crawls	Very big (TB, PB)	Common Crawl (and OSCAR-processed CC dumps), Internet Archive dumps etc.	Challenge due to their size (hard to transfer, hard to preprocess, hard to store; must be close to the HPC)	Indispensable due to their size and coverage – but: high level of noise, massive need for pre-processing
New, fresh data from industry and other organisations	Arbitrary size, ideally as large as possible	Publishing houses, media companies, libraries, call centres, broadcasters etc.; also: Media Data Space	Can be easily integrated by connecting these organisations to LDS	Especially high quality data or domain-specific data or data covering specific languages and thus highly relevant for LLMs

Alliance for Language Technologies EDIC (ALT-EDIC)

- European Digital Infrastructure Consortium (EDIC): a new legal entity type in the EU
- Two EDICs have recently been established including the ALT-EDIC (not-for-profit association in France)
- Coordinated by the French Ministry of Culture
- Collaboration between: ALT-EDIC Working Group, EC, LDS
- ALT-EDIC action plan will concentrate on:
 - 1. Data
 - 2. Existing language models
 - 3. New language models
 - 4. Evaluation, certification, normalization
 - 5. Ecosystem
 - 6. EDIC implementation
- We expect many synergies between LDS, ALT-EDIC, DSSC, Simpl, other data spaces and other projects!

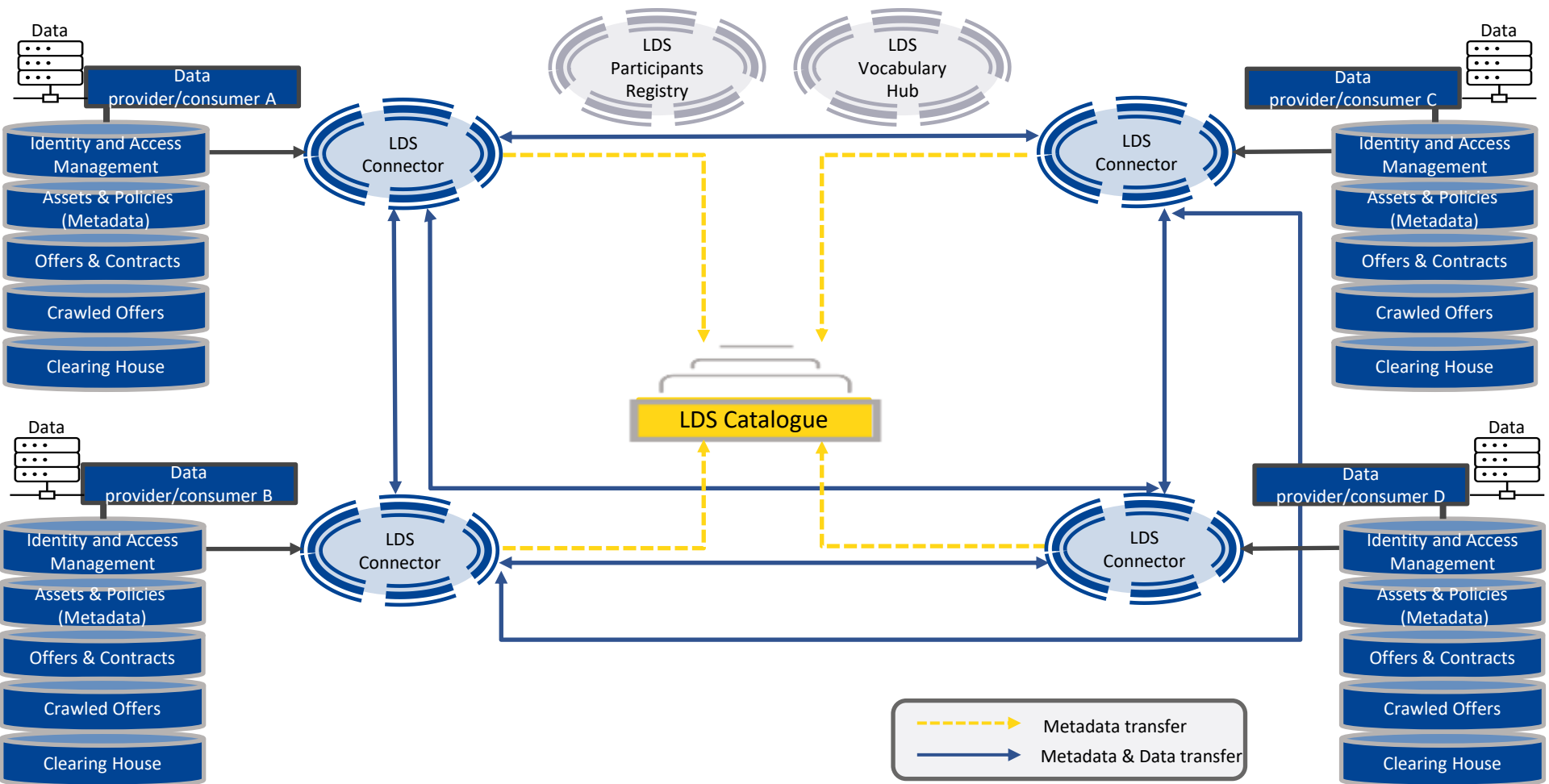
Long History of Language Data Sharing in the European NLP/CL Community

The screenshot shows the META-SHARE website. At the top left is the logo "META-SHARE". To the right are navigation links: "LEARN", "DISCOVER", "REGISTER", "CONTACT", and "LOGIN". The main heading is "Search & exchange language resources". Below this is a sub-heading: "META-SHARE is an open and secure network of repositories for sharing and exchanging language data, tools and related web services." To the right of this text is a button that says "Share your own resources!" and another button below it that says "JOIN OUR NETWORK NOW!". There is a search bar with the placeholder text "Search for NLP (to 4444) resources". Below the search bar are four statistics: "4,481 users", "2,887 languages resources", "92% text corpora", and "27,630 number of downloads".

The screenshot shows the CLARIN Virtual Language Observatory (VLO) homepage. At the top left is the logo "Virtual Language Observatory" and navigation links: "Search", "Contributors", and "Help". At the top right is the CLARIN logo. The main heading is "CLARIN Virtual Language Observatory". Below this is a sub-heading: "Welcome to the VLO!". The main text says: "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." There are two buttons: "See all records" and "Take a quick tour". At the bottom is a search bar with the placeholder text "Search through 1,234,567 records" and a search icon.

The screenshot shows the ELRC-SHARE Repository homepage. At the top right is the logo "European Language Resource Centre for the Americas". Below this is a banner image of many flags. The main heading is "ELRC-SHARE Repository". Below this is a search bar with the placeholder text "Type in your keywords, please...". Below the search bar is a sub-heading: "Welcome to the ELRC-SHARE repository!". Below this is a paragraph of text: "The ELRC-SHARE repository is used for discovering, storing, searching and accessing Language Resources that are collected through the European Language Resource Centre and deposited under the Creative Commons Attribution 3.0 License." Below this is another paragraph of text: "If you want to contribute resources, all you need to do is register your user in http://www.elrc.org and get your ID card and upload your data with a simple form."

The screenshot shows the ELRA website search results page. At the top left is the ELRA logo. Below this is a search bar and a "Search" button. Below the search bar is a list of search results. The first result is "2000 CMC, Shared Task - Arabic & Czech". Below this is a list of search results for "2000 CMC, Shared Task - Arabic & Czech". The second result is "2000 CMC, Shared Task - Ten Languages". Below this is a list of search results for "2000 CMC, Shared Task - Ten Languages".



Current LDS
Prototype



MANAGEMENT

[Home](#)

[History](#) ▼

[Storage Solutions](#) ▼

OPERATIONS

[Assets](#) ▼

[Policies](#) ▼

[Offers](#) ▼

LDS Connector Management Panel

Here you can create and manage your assets, your policies and your offers and review your contract agreements.

Assets
15

Create A New
Asset

View My
Assets

Policies
16

Define A New
Policy

View My
Policies

Offers
6

Create Offers

View Offers

Create new data asset

Create a new asset

Select language (optional)
Language
ENGLISH

Basic properties
Title, short description, version, ...

Privacy properties
anonymization or sensitive data details

Language
A language of the resource

Type properties
media type, linguality type, annotation type, corpus subclass, ...

IprHolder properties
ipr holder or creator details...

Documentation
related documentation

Temporal properties
time constraints

Identifiers
identifier details

Distribution
media type, format, ...

Data address
base url, type, ...

Details

Privacy

Language

Type

IPR holder

is documented by

Temporal Coverage

Identifiers

distribution

Data address

CREATE ASSET

Current LDS
Prototype

Adjust and attach policy

POLICY CLASS

Current LDS
Prototype

✓ Interval-restricted Data Usage Policy Class
allows data usage for a specified time period

Purpose-restricted Data Usage Policy Class
allows data usage for a specific declared purpose

pending

Connector-restricted Data Usage Policy Class
allows data usage for a specific connector

Perpetual Data Sale (Payment once) Policy Class
allows data usage after payment is completed (for datasets requiring once-off payment)

pending

pending

Location Restriction of the participant for Data Usage Policy Class
restricts data usage to participants in a specific location

Attribution Data Policy Class
allows distribution of data with mandatory attribution

pending

pending

Share Alike derivatives Policy Class
allows distribution of derivatives with a compatible

Attach Policy When Distribute to a third-party Policy Class
allows distribution of data to a third-party with the specified attached policy

pending

pending

Derivatives not allowed Policy Class
distribution of derivatives is not allowed

pending

Create and publish offer

Current LDS
Prototype

Progress: 1 Select Asset (checked) — 2 Assign Policy (checked) — 3 Review & Publish (active)

Buttons: Previous, Publish

Name	Description	LR type
Arab-Andalusian music corpus	This repository contains Arab-Andalusian corpus collected in the CompMusic project. The following files are available for 164 concert recordings (overall playabl...	Corpus
AcCompl-it Dataset	The AcCompl-It dataset comprehends the Complexity and the Acceptability Datasets. The first data set is composed of 2,530 Italian sentences annotated with human...	Corpus

Apache License, Version 2.0

European Language Grid



EUROPEAN LANGUAGE GRID

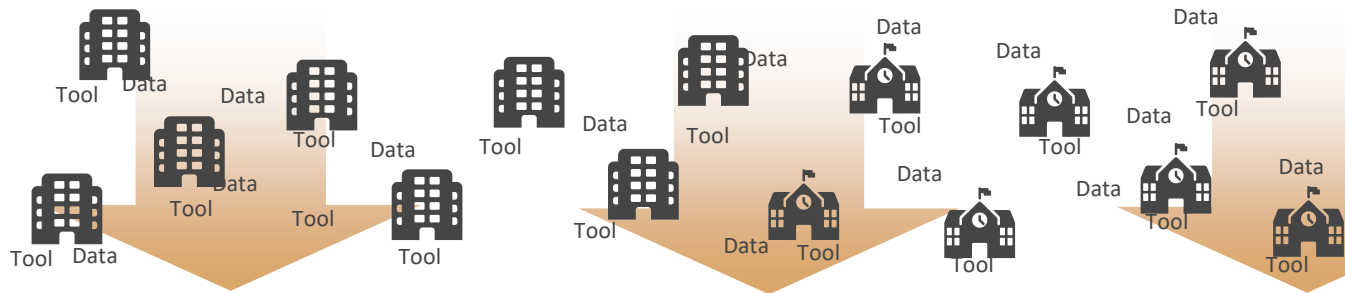
Objectives (Selection)

1. Establish the ELG as the primary LT platform and marketplace in Europe to tackle the fragmentation of the European LT landscape.
2. ELG as a platform for commercial and non-commercial, industry-related LTs (functional and non-functional).
3. Enable the European LT community to upload services and data sets, to deploy them and to connect with, and make use of those resources made available by others.
4. Enable businesses to grow and benefit from scaling up.
5. Unleash enormous potential for innovation.



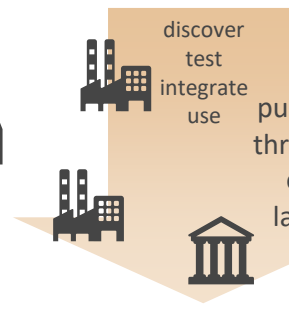
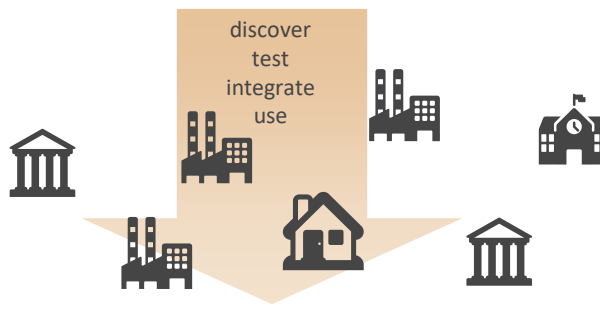
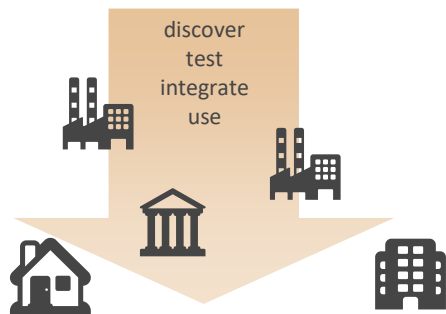
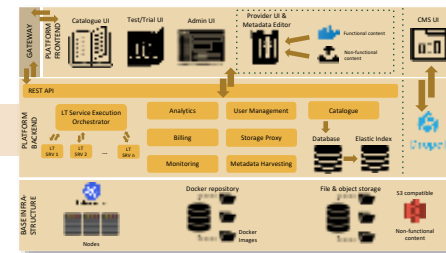
Kick-off meeting, 22/23 January 2019

Developers of Language Technologies: Companies, Universities, Research Centres (approx. 1750-2000 organisations in total in Europe)

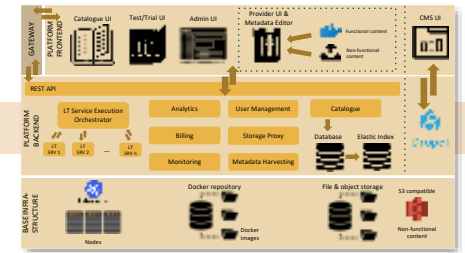


Make available their LT, NLP and speech tools, services and components as well as their data sets, corpora etc.

EUROPEAN LANGUAGE GRID



Discover, download, purchase, use (for example, through the ELG Python SDK or through the web UI) language technologies or language resources



ELG is a **joint tool and resource sharing technology platform** as well as **marketplace** for the whole European LT community (approx. 1750-2000 organisations)

ELG makes all European Language Technologies and Language Resources available in a one-stop-shop.

ELG Release 3 (May 2024)

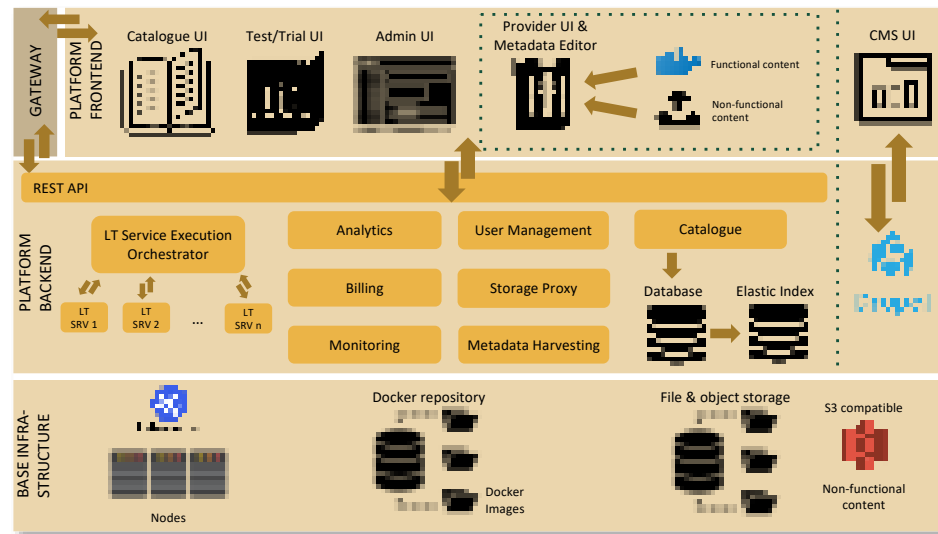
17,800+ Resources without organisations

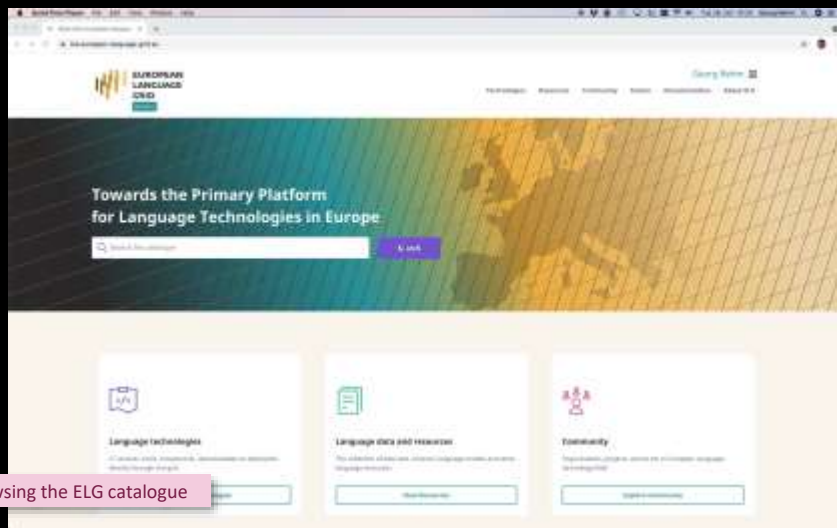
- 8090 corpora and data sets
- 3860 services and tools (1142 services integrated)
- 2830 lexical/conceptual resources
- 513 models, grammars, lang. descriptions
- 2000 organisations (research orgs., companies)

Users can connect to the ELG cloud platform via ELG APIs, remote APIs, ELG GUI, Python SDK, download of containers or source code.

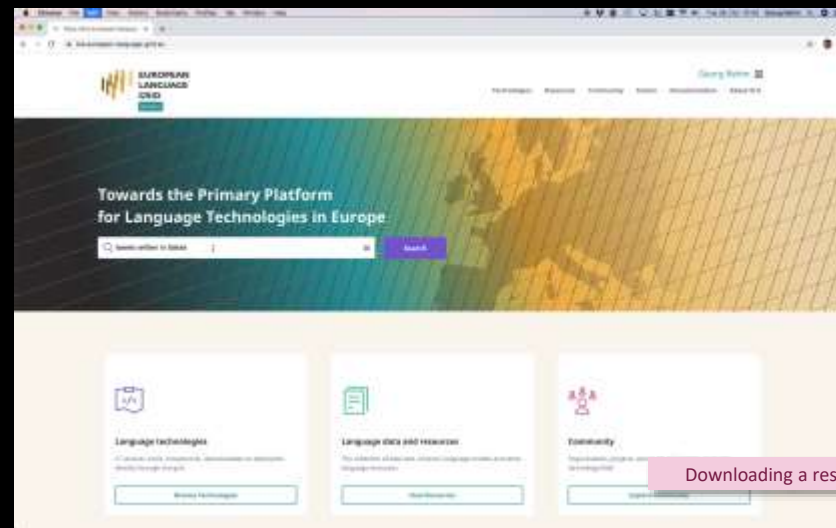
+ ELG / ELE	(10081)	+ Hugging Face	(385)
+ ELRC-SHARE	(3099)	+ CLARIN-PL	(351)
+ ELRA Catalogue of Language Resources	(1178)	+ Quantum Stat Datasets	(262)
+ Zenodo	(1107)	+ LREC Shared LR (ELRA)	(144)
+ LINDAT/CLARIAH-CZ	(596)	+ META-SHARE/ILSP	(69)
+ CLARIN.SI	(527)	+ META-SHARE/DFKI	(2)

Sources of data sets

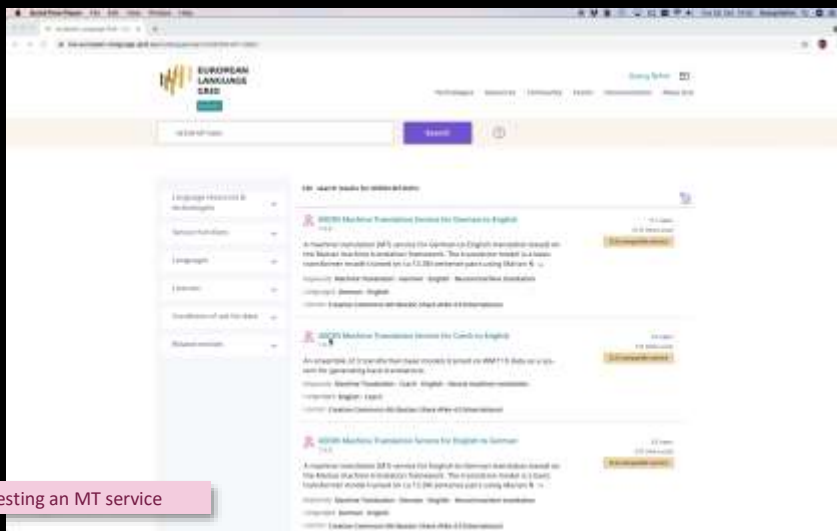




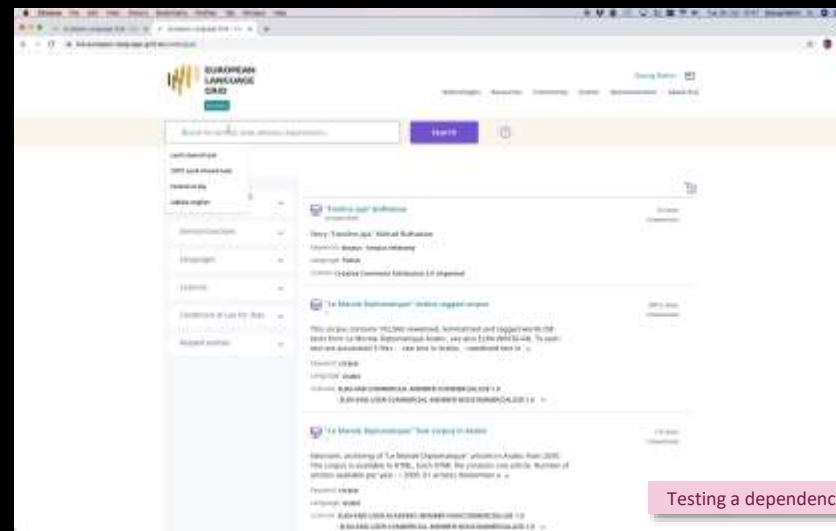
Browsing the ELG catalogue



Downloading a resource

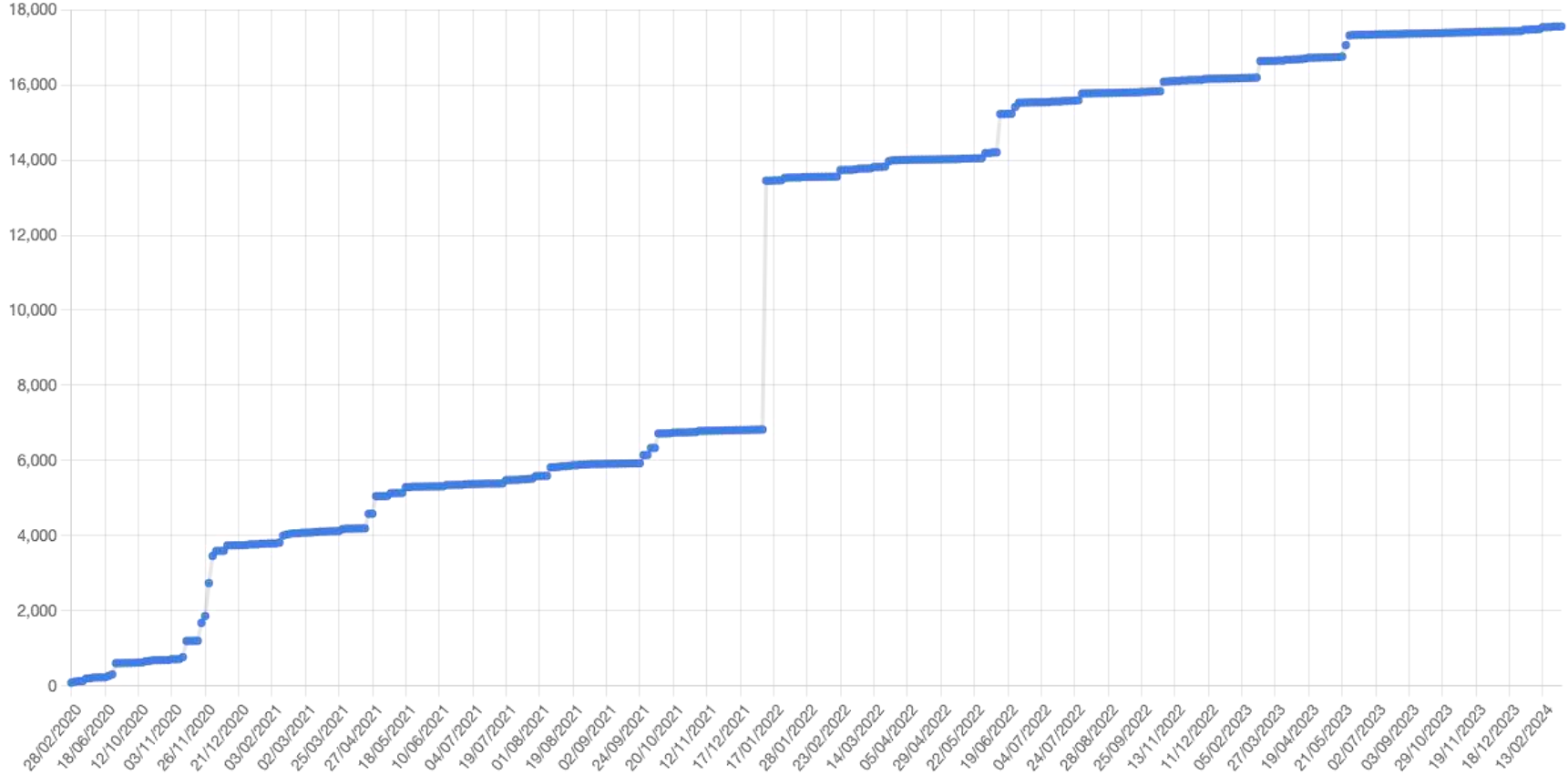


Testing an MT service

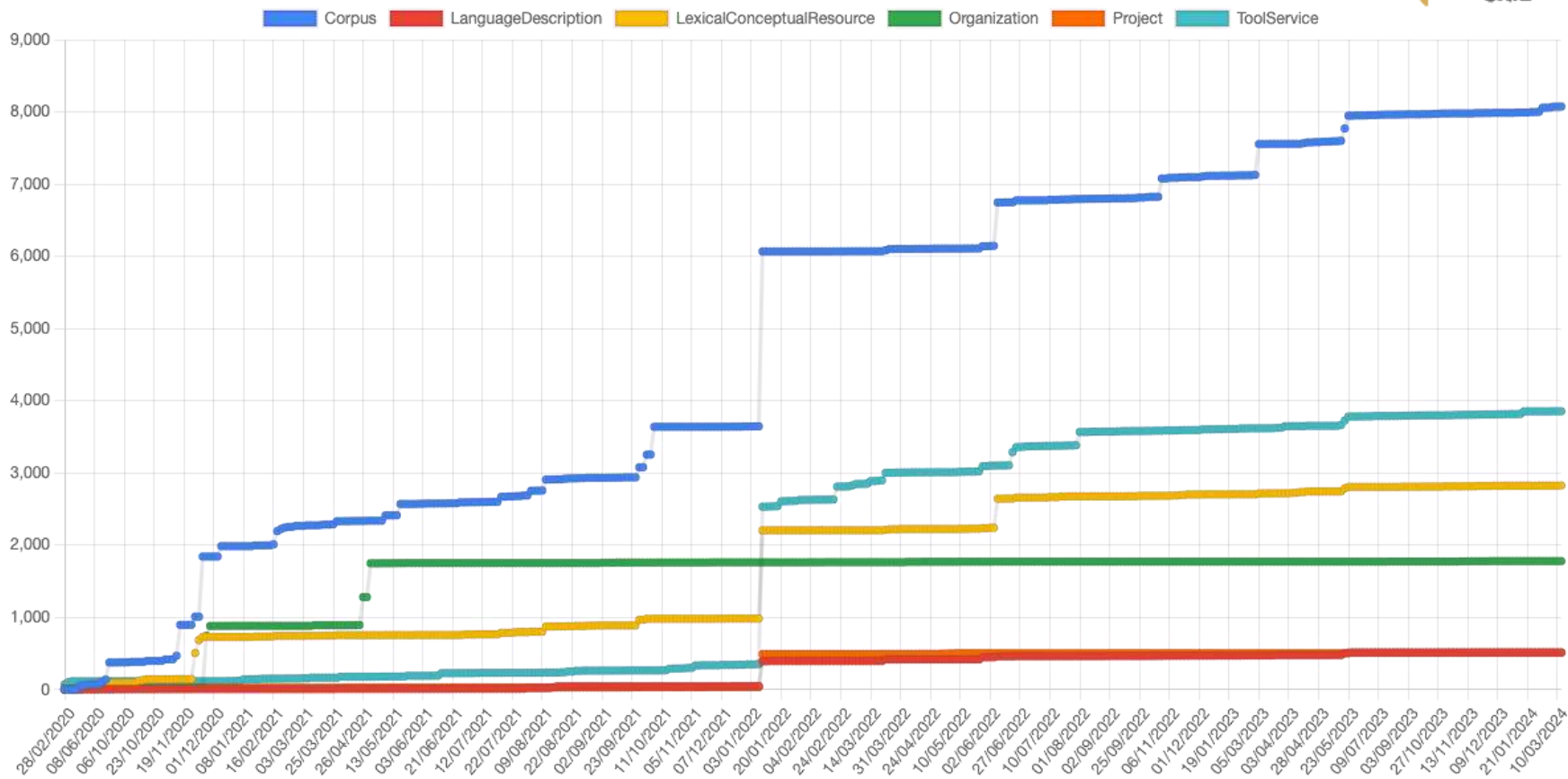


Testing a dependency parser

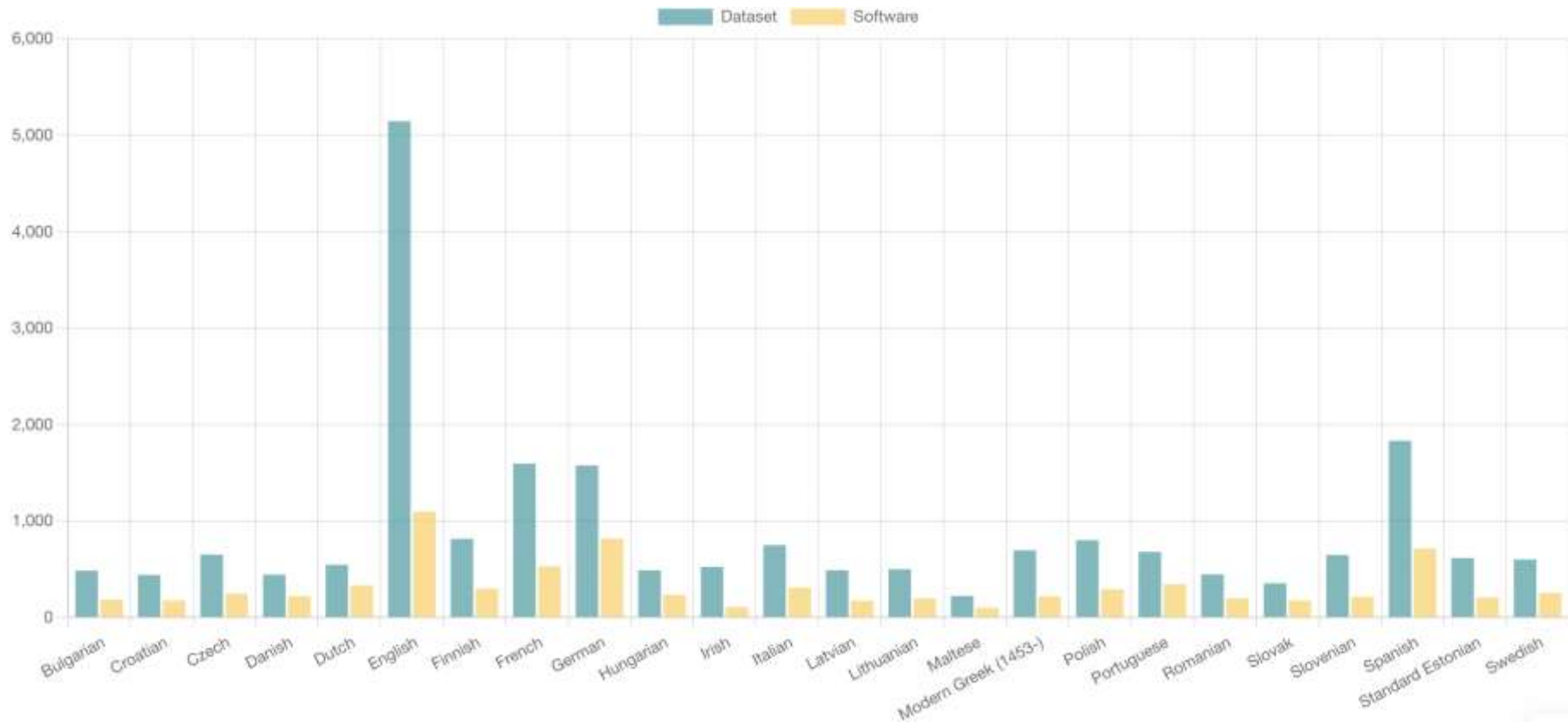
Number of Resources over Time



Number of Resources over Time

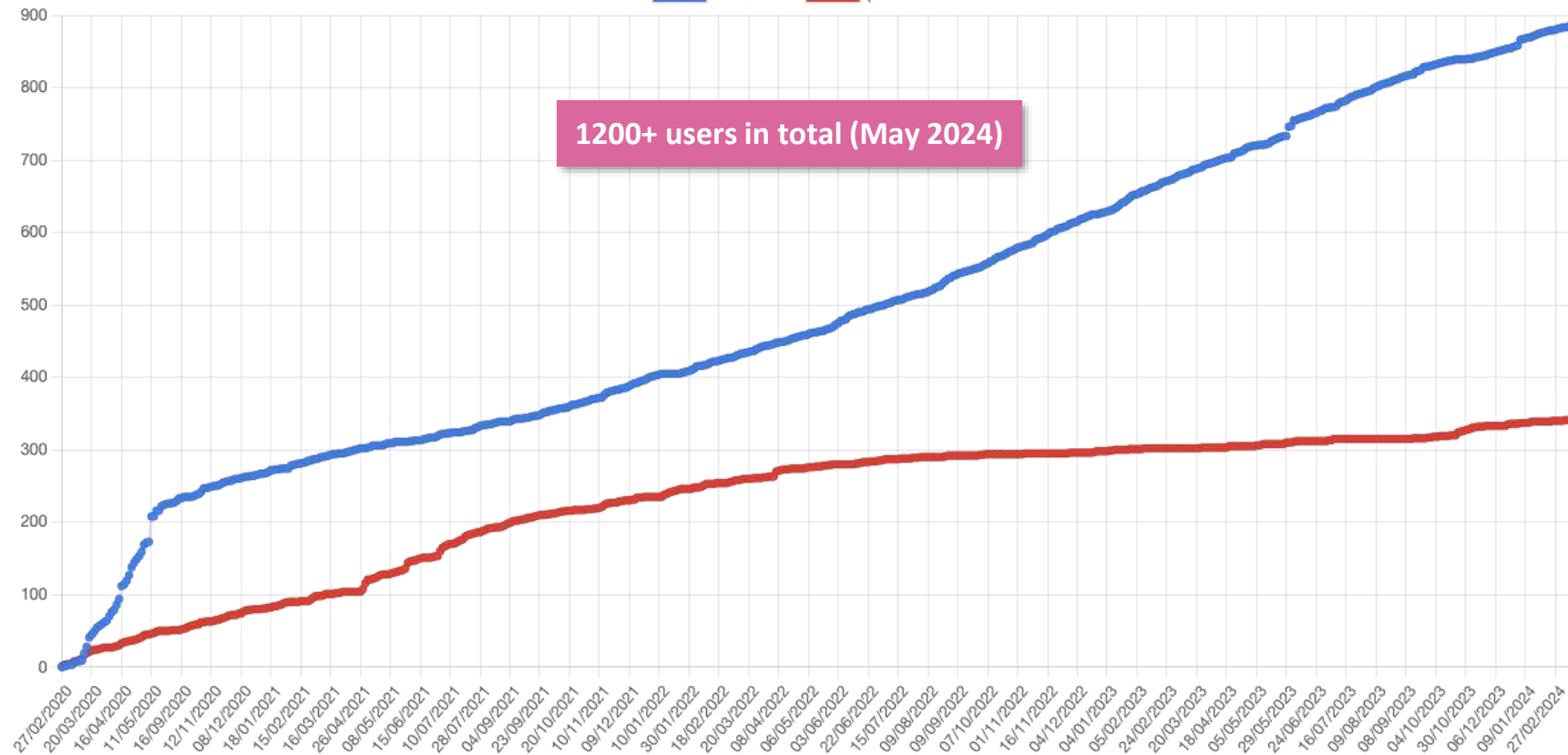


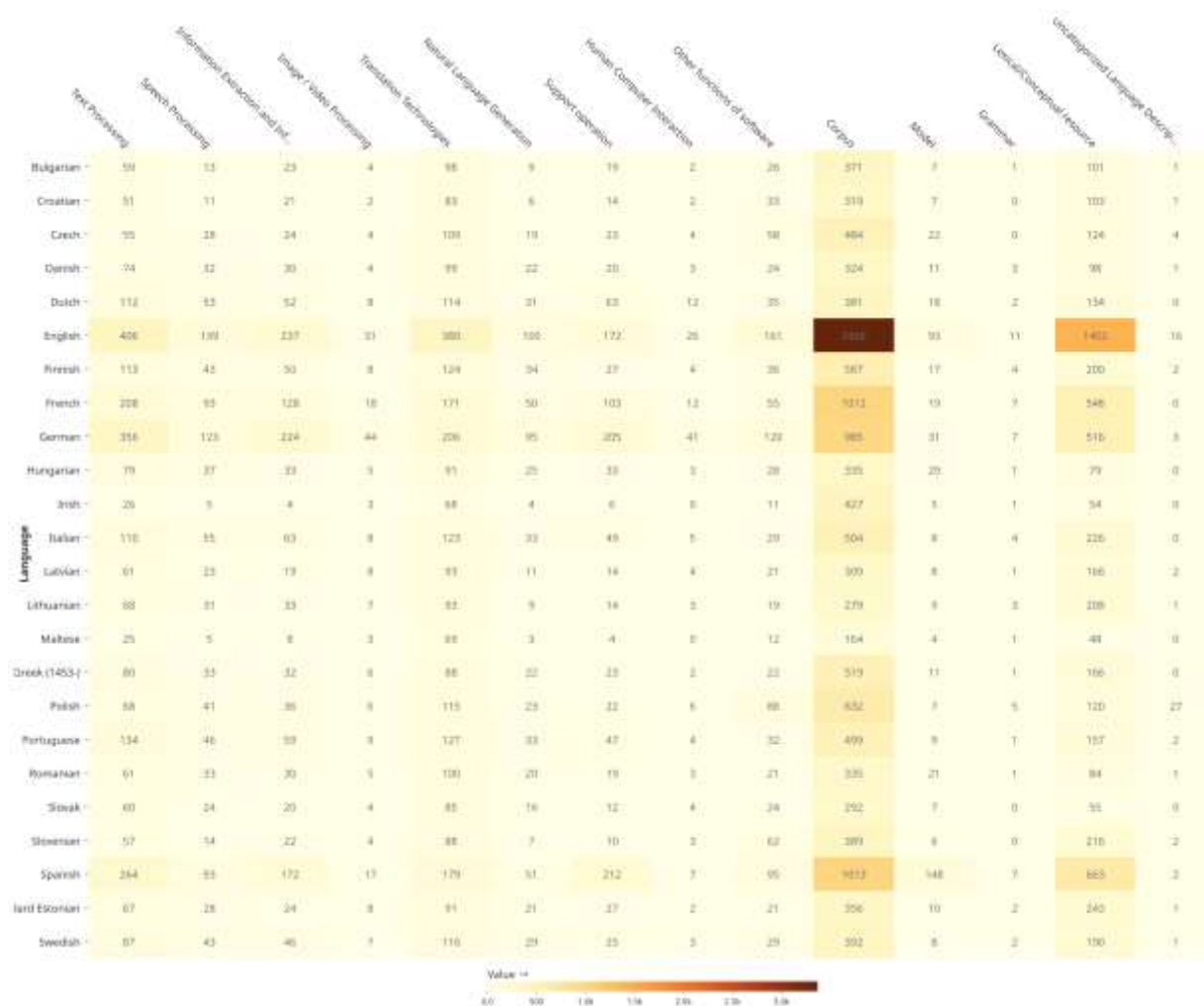
Number of LRTs available in ELG (24 EU languages only)

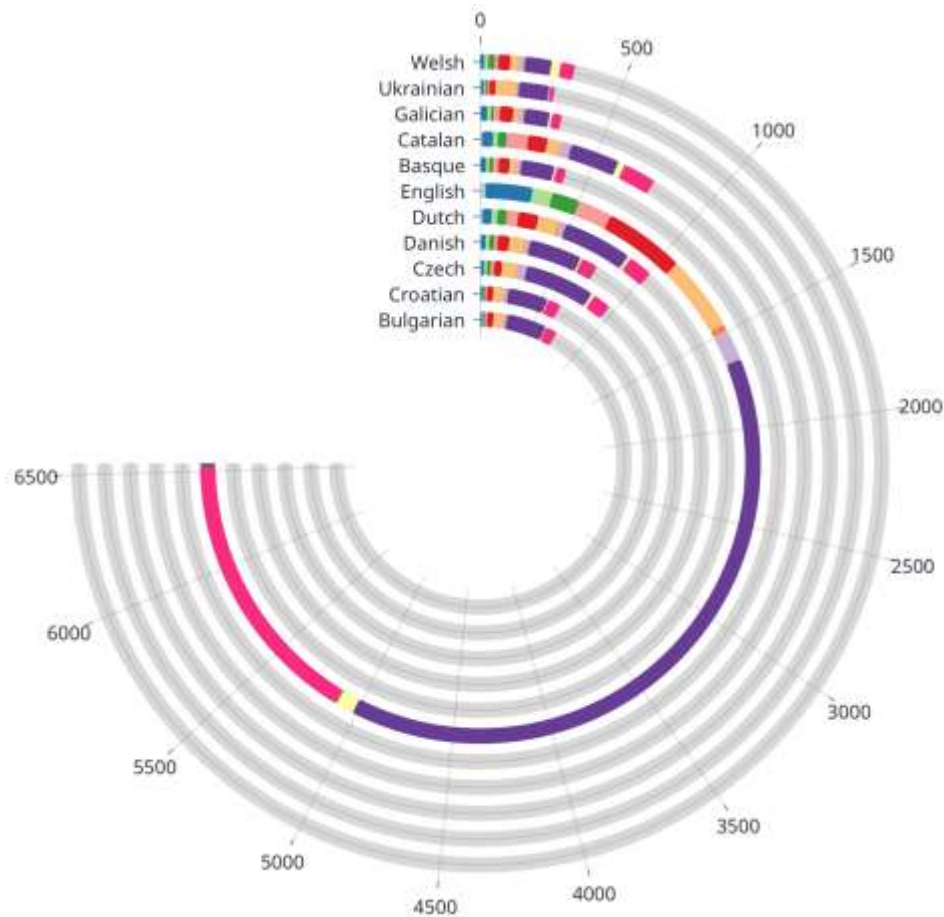


Number of Users (Consumers of LT, Providers of LT)

consumers providers







- 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

Produced using the ELE Dashboard
<https://live.european-language-grid.eu/catalogue/dashboard>

European Language Grid and Data Spaces

- Interoperability between ELG and **LDS**: installing the LDS connector within ELG → discoverability of ELG resources through LDS.
- Interoperability between the emerging **DataBri-X** toolbox and ELG using a similar methodology, which will be compliant with IDSA.
- In **OpenGPT-X**, we adapt ELG to Gaia-X, e.g., by serving SDs to federated catalogues.



NFDI for Data Science and AI

NFDI for Data Science and Artificial Intelligence

(max. runtime of project: 10 years)

- NFDI: German National Research Data Infrastructure
- NFDI4DS develops a research data infrastructure for DS and AI
- Foundation: Multiple existing tools and infrastructures – including ELG and LDS
- The NFDI4DS infrastructure will eventually also be connected to and made interoperable with the emerging national NFDI infrastructure
- Research topics: scholarly information processing, knowledge extraction and the construction of research knowledge graphs such as ORKG
- DFKI co-coordinates two Task Areas: Infrastructure & Services, Transfer & Application
- Most recent event: Natural Scientific Language Processing Workshop (NSLP 2024), 27 May 2024, proceedings to be published with Springer (LNAI) soon.
- Upcoming event: 4th Int. Workshop on Scientific Knowledge – Representation, Discovery, and Assessment (Sci-K 2024), 11/12 Nov. 2024, Baltimore.





Natural Scientific Language Processing Workshop
(NSLP 2024), 27 May 2024



<https://nfdi4ds.github.io/nslp2024/>



Summary and Conclusions

Summary and Conclusions

• European Language Grid



- ELG is the most comprehensive NLP/LT platform in Europe.
- There's still an extremely strong imbalance in terms of technology support of Europe's languages – all languages except English under-resourced!
- ELG can be used to measure the LT support for Europe's languages.
- The EU concentrates on data spaces and establishing a data economy in Europe.
- We're integrating ELG into various data space initiatives.
- A follow-up project is very much needed to adapt ELG to modern state of the art technologies, especially LLMs.

• European Language Data Space



- Project is in full swing: technical development, promotion, dissemination, governance etc.
- Collaborations with DSSC, Simpl and ALT-EDIC; European projects, e.g., HPLT, OpenGPT-X, OpenWebSearch; other data spaces, especially Media and Cultural Heritage
- Adoption of LDS by industry and other organisations → grow the LDS User Group
- Make available new and fresh language data, esp. from industry, covering all European languages and modalities
- ***Still the key open challenge: interoperability between platforms and infrastructures!***





Thank you!



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

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

06-06-2024 Inria and ECDF Event
<https://language-data-space.ec.europa.eu>

