

## EUROPEAN LANGUAGE DATA SPACE



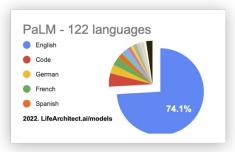
### **LDS Country Workshop Slovakia**

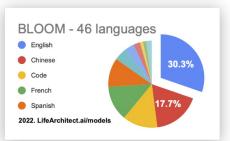
Dr. Simon Ostermann (DFKI GmbH, Germany) – LDS Project Manager and Site Manager Saarbrücken simon.ostermann@dfki.de

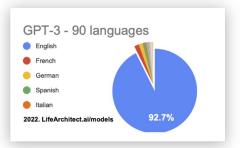
07-11-2024 LDS Country Workshop Slovakia https://language-data-space.ec.europa.eu

### **Context: Large Language Models (LLMs)**

- Unprecedented capabilities: LLMs are the most disruptive breakthrough in AI in recent history (GPT-3, ChatGPT, GPT-4, Claude, Gemini etc.)
- Trained on vast amounts of data (+ images, videos, audio, i.e., multimodal data)
- LLMs are getting larger and larger: more data, more parameters, more compute. Scaling laws: larger models outperform smaller models.
- Multilingualism makes everything much harder (data imbalance): Europe's languages are vastly under-resourced, except English
- Unprecedented opportunities:
  - The global LT/NLP market is expected to reach 439.85B\$ by 2030
  - The global Gen AI market is expected to reach 1.3T\$ by 2032
- A concerted effort for the collection of data for all European languages is very much needed to be able to develop LLMs according to our needs and cultures
- 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







Europe is where ChatGPT gets regulated, not invented. That's something to regret. As unhinged as the initial results of the artificialintelligence 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.



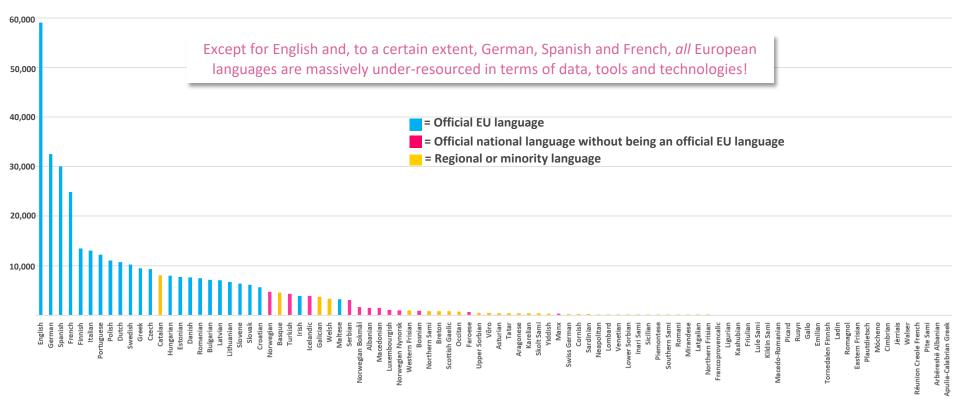
### **European Initiatives**

- 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:
  - Availability of data for European languages
  - HPC facilities
  - Speed of the big tech players in the US and Asia vs. speed of Europe



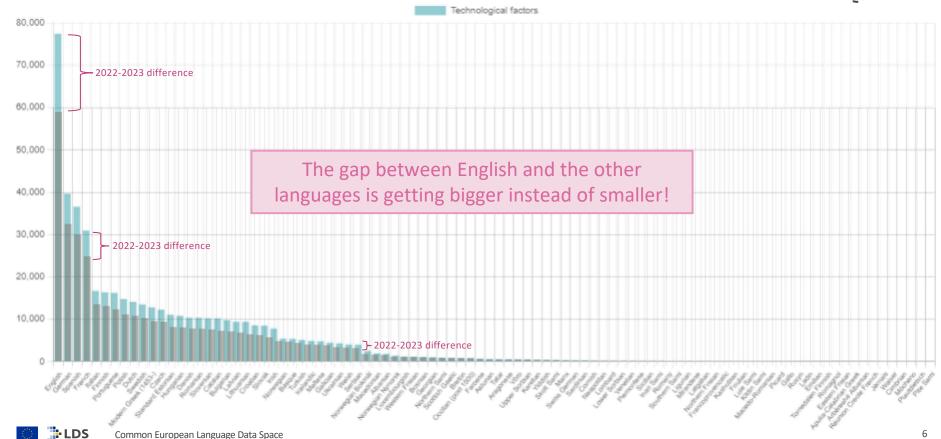
### **Digital Language Equality Metric: Technological Scores**





### **DLE Metric: 2022 vs. 2023**





## **EU Data Strategy**



### META SHARE LEARN + DISCOVER + PARTICIPATE + CONNECT + 1 LOGIN Search & exchange language resources Share your own resources! META-SHARE is an open and secure network of repositories for sharing and exchanging language data, tools and related web services earch the META-SHARE inventory number of text corpora downloads

## **Long History of Language Data Sharing**





### **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
- The Common European Language Data Space is one of the 15 official EU data space projects with a strong focus on industry



## Data initiatives - EU-level

Geo- information	Construction	Energy	Space	Public Adminis- tration	Research/ Education	Automotive	Manu- facturing	Mobility	Health	Agriculture	Climate	Finance	Culture	Media	Language	Smart cities & commu- nities	Tourism
	Construction	EU Energy Data Space intNET OMEGA-X EDDIE Enershare Synergies Data cellar	Space	Adminis-		Automotive		EU Mobility Data Space PrepDSpace 4Mobility Deploy EMDS*	EU Health Data Space  MyHealth@ EU  Support for HDABs  Healthdata @EU piliot  Central Services for Health Data@EU  PaTHED  Support for SNOMED CT  Capacity building for prim-sec. Use cases  Joint Action for primary uses  Joint action for secondary uses  Data Quality & Utility	EU Agriculture Data Space AgriData Space Divine Crack Sense ScaleAg Data AgData Value 4Growth* Dig4Live*	EU Green Deal Data Space GREAT ADAGD B-Cubed FAIRICUBE USAGE	EU Financial Data Space  Digital Europe	EU Cultural Data Space  Deployment Eureka3D SDCulture DE-BIAS  AIAEuropeana	Media EU Media Data Space TEMS	EU Language Data Space Digital Europe	& commu-	EU Tourism Data Space  DATES  DEST
					AqualNFRA  Blue-Cloud 2026  OSCARS  EVERSE  OSTRAIIs*  EOSC  Beyond  EOSC- ENTRUST  SIESTA*				EUCAIM  GDI		Project in the	n EU COM SWD(	non European Da 2024) 21 final. Fo	or timeline (2022	2-24) see p.56 ff.	in <u>EU COM SWD</u> I	(2024) 21 final.



### **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 infrastructure, openness, promotion
- Stakeholders: industry, research, public administration, cultural associations, NGOs and citizens
- The four core partners 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)







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









### **Consortium and Subcontractors**

Lead Partner and 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	3рс	DE					
CLARIN ERIC	CLARIN	NL					
Big Data Value Association (Data, Al and Robotics) AISBL	BDVA	BE					

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

12



### **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
- The first couple of EDICs are currently under development including the ALT-EDIC
- Coordinated by the French Ministry of Culture
- Close 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!

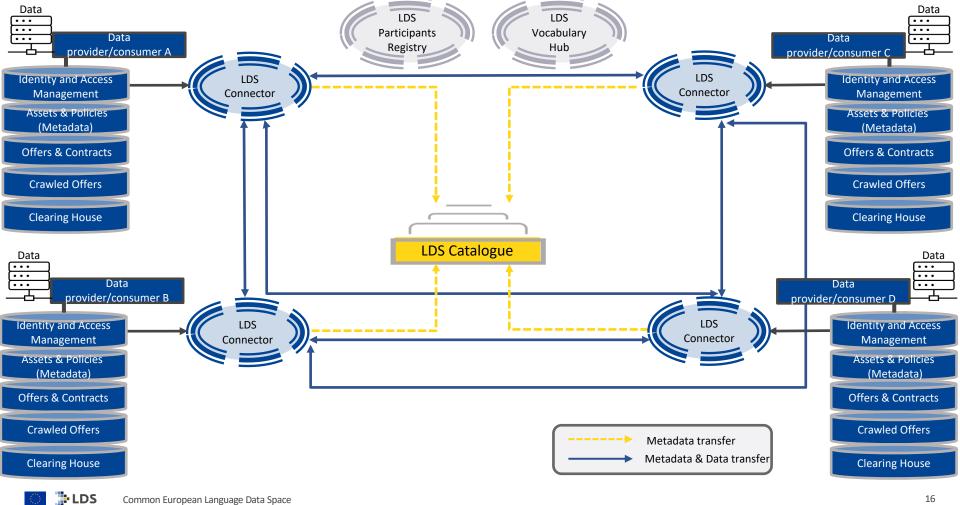
#### **ALT-EDIC Members**

- **17 Members States:** Bulgaria, Croatia, Czechia, Denmark, Finland, France, Greece, Hungary, Ireland, Italy, Latvia, Lithuania, Luxembourg, Netherlands, Poland, Slovenia and Spain;
- **8 observing Member States:** Austria, Belgium, Cyprus, Estonia, Malta, Portugal, Romania and Slovakia.

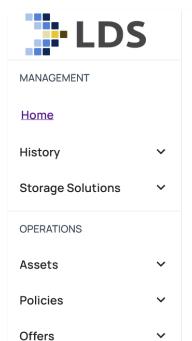


## The Common European Language Data Space - Implementation





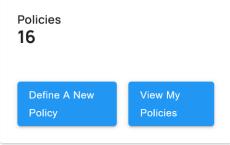


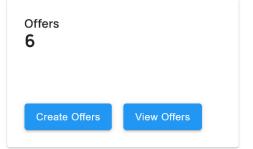


### LDS Connector Management Panel

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

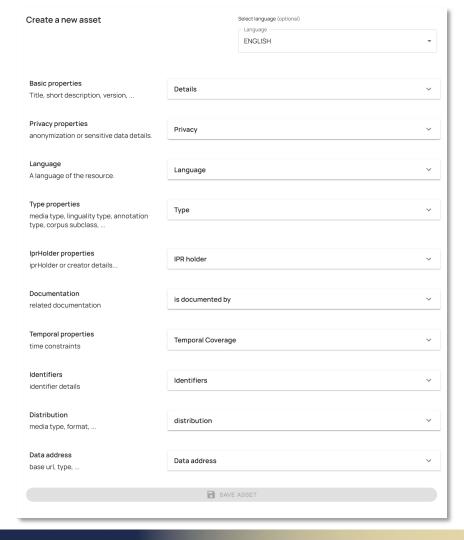








# Step 1: Create new data asset

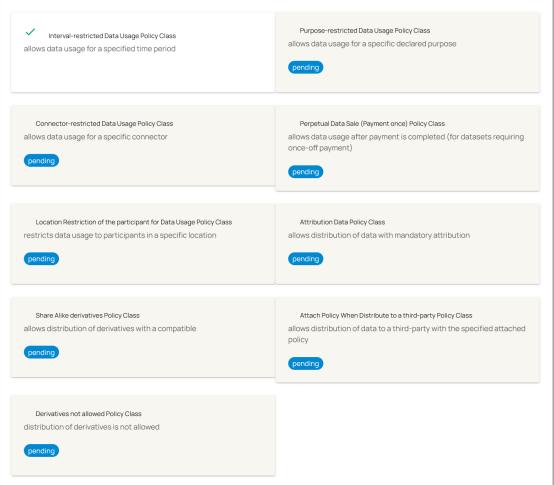






#### POLICY CLASS

# Step 2: Create and adjust policy

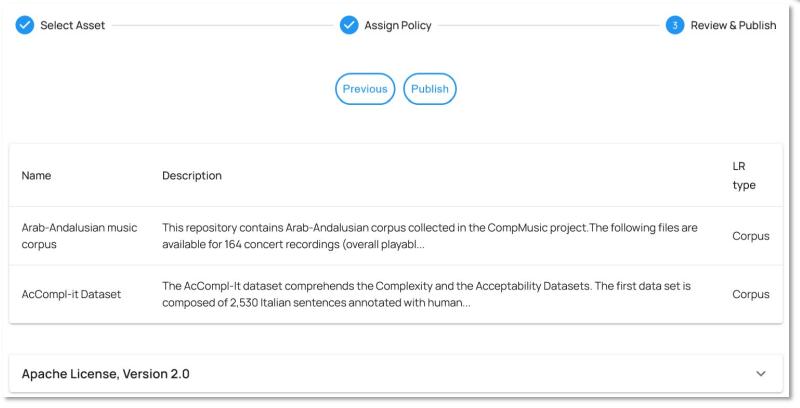




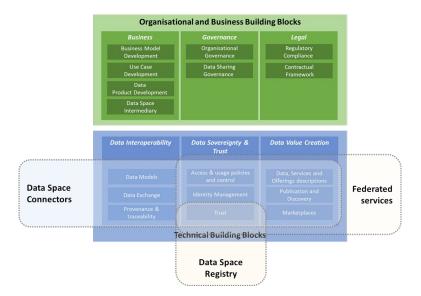


### Step 3: Create and publish offer



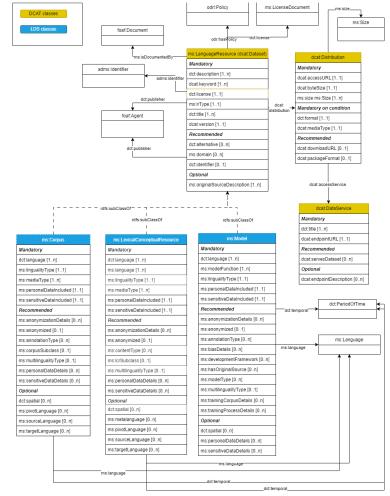






### **Build on Existing Solutions**

- Following DSSC (see above)
- Eclipse Data Space Components (EDC)
- DCAT-AP, Language DCAT-AP (see right), ODRL
- Mappers from existing platforms

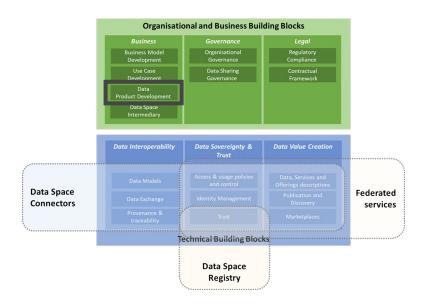


## The Common European Language Data Space – a Value Proposition



### DSSC Blueprint 1.0 – Data Product Development







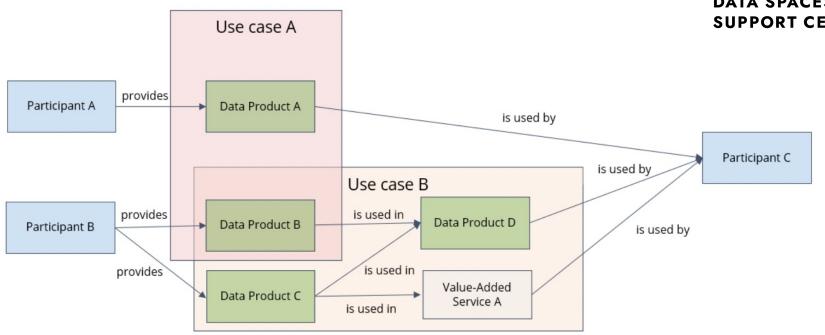
### **Language Data – Language Resources – Data Products**

- The NLP and Computational Linguistics community has been sharing language data since the 1990s
- Back then: annotated corpora, treebanks, grammars, lexicons, smaller language models
- The term "language resource" (LR) was established (data, documentation, evaluation, metrics etc.)
- language resource ≈ data product
- Of utmost importance now: In the LDS, identify and make available large amounts of language data to enable industry and research to pre-train large language models for Europe
- Typical availability of LRs *since the late 90s and early 2000s*:
  - For research purposes: free of charge
  - For commercial use: often for a certain fee. LRs developed by European research organisations were licensed by European but also large US tech companies, e.g., for online NLP services (Machine Translation). LDS: Streamline exchange of such data



### **LDS and Data Products**





DSSC Blueprint 1.0 – Data Product Development



### **Language Data Space – Value Proposition**

### **Data Providers**

- Additional revenue LDS as a marketplace
  - Sell data products
  - Find new customers
  - Extend or enrich datasets using AI/NLP services offered in the wider LDS ecosystem
- Legal compliance by design
  - Stay in control over use and access of data
  - Compliance with EU regulation and standards
- Limited effort
  - Keep existing infrastructure and workflows
  - Interoperability with other data spaces
  - Legal and technical helpdesks available
- Contribute to European LLMs: from and for Europe

### **Data Consumers**

- Buy or access data products to develop better services (including LLMs)
  - Multilingual data
  - Multimodal data
  - Domain-specific data
  - All European languages
  - Easy discoverability and access
- Limited effort: keep existing infrastructure
- Legal compliance by design
  - Compliance with EU regulation and standards
  - Transparency: emphasis on data provenance















○ D

### Data Products in LDS – Training Data for Generative AI and LLMs

- A few examples of recent agreements:
  - Reddit: \$60 million per year (Google)
  - Shutterstock: \$25-50 million (Apple)
  - Springer: Tens of millions (Open AI)
  - Offer for news publishers: \$1-5 million per year (Open AI)
  - Offer for owners of large datasets: \$50 million (Apple)
- Global market is **enormous** owners/providers of large amounts of content **are paid large sums by the US technology enterprises** that currently dominate the AI product landscape
- It's up to the data providers to establish offers and prices that make sense for them
- Our ambition: to establish LDS as a marketplace for European language data



### **Next Steps**

- LDS 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 relevant data spaces, especially Media and Cultural Heritage
  - EuroHPC
- Adoption of LDS by industry and other organisations → grow the LDS User Group
- Identify and make available new and fresh language data, especially from industry and covering all European languages and modalities



### **LDS User Group**

### Meetings and Conferences

- Inaugural meeting in March 2024
- Second meeting in June 2024
- Third Meeting yesterday
- Launch Conference: early 2025

### https://language-data-space.ec.europa.eu



Join the LDS user group



The European Language Data Space (LDS) user group members shall actively contribute to and take advantage of the LDS, bringing in their own requirements and validating the emerging LDS infrastructure.

If you are a stakeholder who is in need of language data or if you want to give the language data of your organisation a second life, potentially monetising it, you are welcome to join.

Click to join



### Communication

- Established a mailing list for the LDS user group
- The LDS user group will grow new members will be added to the mailing list

### • If you're interested, please get actively involved and join the LDS User Group!

- Validation of concepts, ideas, software; first test installations of the LDS connector (foreseen for Q4 2024); first trial exchanges of data; surveys etc.
- You can also help on a more substantial, in-depth level please approach us if you're interested.





**Common European Language Data Space** 

### Thank you!





A Common European Language Data Space – funded under contract LC-01936389 with the European Union. Dr. Simon Ostermann (DFKI GmbH, Germany) – LDS Project Manager and Site Manager Saarbrücken simon.ostermann@dfki.de

07-11-2024 LDS Country Workshop Slovakia https://language-data-space.ec.europa.eu