

EUROPEAN LANGUAGE DATA SPACE



Country Workshop Malta

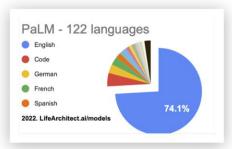
Khalid CHOUKRI (ELDA)

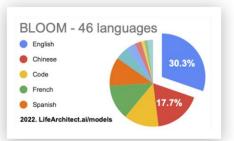
choukri@elda.org

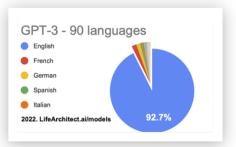
LDS Country Workshop in Malta https://language-data-space.ec.europa.eu

Context: Large Language Models (LLMs)

- Unprecedented capabilities: Large language models are the most disruptive breakthrough in AI in recent history (GPT-3, ChatGPT, GPT-4, Claude, Gemini etc.)
- LLMs are trained on vast amounts of data and optionally also image, video, audio etc. data, i.e., multimodal data
- 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 ...
- ... and to make a difference with European data and European stakeholders.
- Already now billions and billions are made but ...



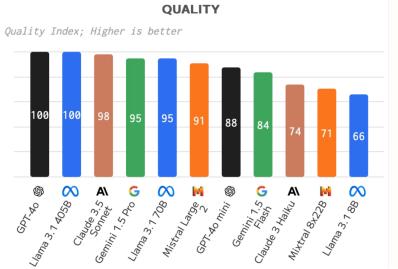






LLM News





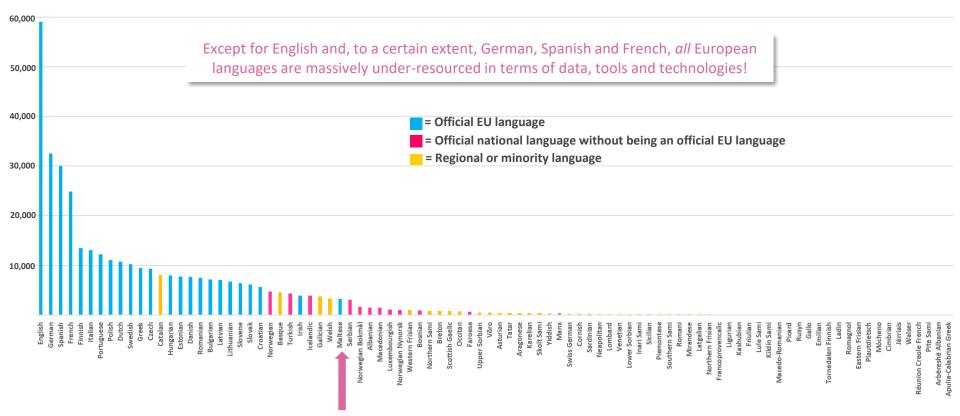
European Initiatives

- European initiatives for the development of LLMs
 - Large research projects in almost every country, e.g., Poland, Slovenia, Spain, Denmark, Italy, Germany etc.
 - Companies in many countries, e.g., Finland (Silo.ai), France (Mistral), Germany (Aleph Alpha)
 - EU-funded projects, e.g., HPLT, TrustLLM, OpenEuroLLM, LLMs4EU
 - New pan-European initiative: ALT-EDIC with an EC Support Action ALT-EDIC4EU
 - New EU initiative: AI Factories (tightly coupled with national HPC centres and EuroHPC JU)
- Challenges:
 - HPC facilities (amount, access and ease-of-use)
 - Speed of the big tech players in the US and Asia vs. speed of Europe
 - Availability of data for *all* European languages *right now the most crucial bottleneck by far*



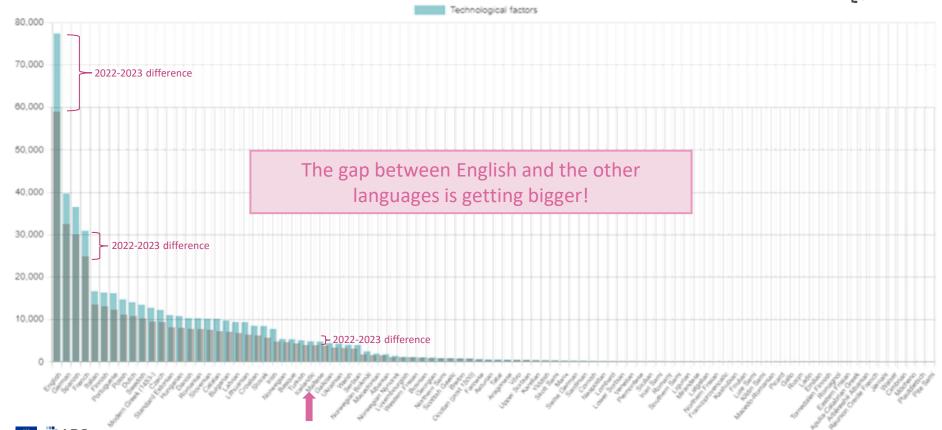
Digital Language Equality (DLE) Metric: Technological Scores





DLE Metric: 2022 vs. 2023





EU Data Strategy & Data Spaces

- Data Spaces are an inherent part of the EU Data Strategy
- Data Spaces will help establish and grow the data economy in Europe
 - Organisations can sell and monetise their data, i.e., they can benefit from its value
- Convergence of the data economy initiatives in Europe (Gaia-X, IDSA, FIWARE, BDVA).
- Important initiatives in Europe:
 - DSSC blueprint and specification development, community coordination and events
 - Simpl development of Open Source data space components
 - Approx. 20++ data space initiatives with the EU's official mandate
 - The Common European Language Data Space is one of the official EU data space projects with a strong focus on industry
- European investment (both national and EU) already more than 2 billion Euros



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
		EU Energy Data Space intNET OMEGA-X EDDIE Enershare Synergies Data cellar			EU Skills Data Space DS4Skills EDGE-Skills EU Open Science Cloud Skills4EOSC EOSC Focus FAIR-IMPACT RDA TIGER FAIRCORE 4EOSC AMEOSC EuroScience Gateway FAIR-EASE		EU Manu- facturing Data Space Jace 4.0 SMARTENA NCE UNDERPIN*	EU Mobility Data Space PrepDSpace 4Mobility Deploy EMDS*	EU Health Data Space MyHealth@ EU Support for HDABs Healthdata @EU pilot Central Services for Health Data@EU PaTHED Supprt for SNOMED CT Capacity building for prim+sec. Use cases Joint Action	EU Agriculture Data Space AgriData Space Divine Crack Sense ScaleAg Data AgData Value 4Growth* Dig4Live*	EU Green Deal Data Space GREAT AD4GD B-Cubed FAIRICUBE USAGE	EU Financial Data Space Digital Europe	EU Cultural Data Space Deploy- ment Eureka3D SDCulture DE-BIAS AIAEuropeana	EU Media Data Space TEMS	EU Language Data Space Digital Europe		EU Tourism Data Space DATES DFST
					RAISE SciLake SCILAKE EOSC4 Cancer GraspOS CRAFT-OA AquaINFRA Blue-Cloud 2026 OSCARS EVERSE OSTRAIIS* EOSC Beyond EOSC- ENTRUST SIESTA*				for primary uses Joint action for secondary uses Data Quality & Utility Label Dev. EUCAIM GDI		Project in the	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 trustworthy and secure European infrastructure and marketplace for the collection, creation, sharing and re-use of multilingual and multimodal language data
- LDS is a marketplace for all organisations commercial and public.
- LDS provides helpdesks for legal and for technical questions.
- 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.
- Technical development informed by ELG, ELRC-SHARE, META-SHARE.

Lead Partner and Coordinator						
Deutsches Forschungszentrum für Künstliche Intelligenz GmbH						
Partners and Operation Leads						
R.C. "Athena", Institute for Language and Speech Processing						
Evaluations and Language Possurees Distribution Agency						

Main Subcontractors

3pc GmbH Neue Kommunikation

ig Data Value Association (Data, Al and Robotics) AISBI

LDS Launch Conference – 19 March 2025 – Villers-Cotterêts, France







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;
- **9 observing Member States:** Austria, Belgium, Cyprus, Estonia, **Malta**, Portugal, Romania and Slovakia, and the Flanders (region).



Long History of Language Data Sharing META SHARE LEARN + DISCOVER + PARTICIPATE + CONNECT + 1 LOGIN ❷ Virtual Language Observatory Search Contributors Help Search & exchange language resources Share your own resources! META-SHARE is an open and secure network of repositories for sharing and exchanging **CLARIN Virtual Language Observatory** language data, tools and related web services Welcome to the VLO! earch the META-SHARE Inventor 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 number of text corpora downloads Search through 1,030,321 records 1096 Language Resources (Page 1 of 55) Media Type 2006 CoNLL Shared Task - Arabic & Czech LANGUAGE Catalogue S Documentation & Media S About S 8 Czech consists of dependency treebanks used as part of the CoNLL 2006 shared task on multi-lingual dependency parsing. The Conference on earning (CoNLL) is accompanied every year by a shared task intended to promote natural lan Type in your keywords, please...

Language Technologies

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

Browse ELG and find the LT services,

are looking for.

resources, developers and providers you

3884

2812

510

513



Common European Language Data Space

The ELRC-SHARE repository is used for documenting storing browsing and accessing Language Resources that are of

If you want to contribute resources, all you have to do is register (new user) or login (returing user) and go on to describe

Coordination and considered useful for feeding the CEF Automated Translation (CEFAT) platform

Welcome to the ELRC-SHARE repository!

CLARIN

Cart total Vew cart Register Logis

Order by Resource Name A-Z

- Ten Languages 🚿

Classes of Data

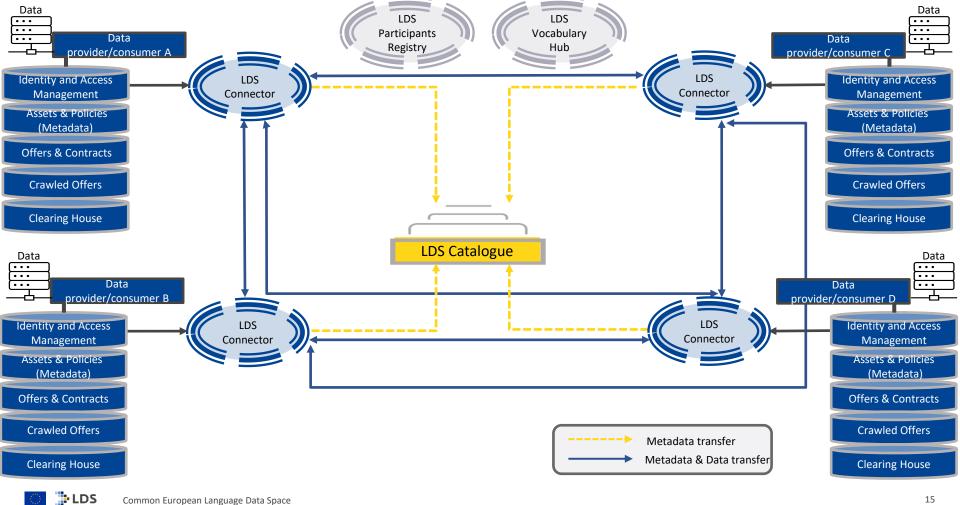
Class of Data	Typical Size	Providers	Integration into LDS	Relevance, especially for LLM Development
Regular Corpora and Language Resources	Small (MB, GB)	Primarily NLP/LT research: ELG, META-SHARE, CLARIN, ELRA, ELDA etc.	Can be easily integrated by connecting the repositories to LDS	Usually very high 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	High quality data; domain-specific data; data covering specific languages; raw data; processed data; evaluation data; data relevant for LLM development etc.



LDS Infrastructure – Basic Functionalities

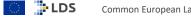
- Apply, i.e., become a participant in the LDS marketplace
- Discover data offers (= data products) through the joint Language Data Space (LDS) catalogue
- Negotiate access to a data product
- Transfer data product from provider to consumer
- Provide data products





Typical Workflow – Purchase Process

- Data provider prepares a new data asset and adds a license (obligatory) and policy (optional)
- Data provider asks LDS helpdesk for advice regarding certain technical or legal questions
- Data provider publishes an offer through their LDS Connector
- This offer is now visible throughout the whole LDS marketplace, i.e., in all LDS Connector instances
- Data consumers can now discover this new offer through their own LDS Connector
- **Specific data consumer** evaluates the new offer ...
- ... and decides to purchase the data offer from the data provider through the LDS Connector
- The corresponding data set is transferred from the provider's connector to the buyer's connector



Common European Language Data Space

16

Language Data Space Marketplace - Value Proposition

Data Providers

- Sell data products (= data sets, data offers)
 - 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 Albased 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
- Find new customers for services and products



Language Data – Language Resources – Data Products – Data Offers

- 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: identify and make available, through the LDS, 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
 - Many unique LRs developed by European research organisations were licensed by European but also large US tech companies, e.g., for online NLP services (Machine Translation)



Futurism

TAP RUNNETH DRY 11.13.23, 4:05 PM EST by MAGGIE HARRISON DUPRÉ

AI Companies Are Running Out of Training Data

The well is running dry.

/ Artificial Intelligence / Ai / Ai Industry / Ai Training





Data Products expected in LDS – Training Data for Gen. Al and LLMs

- A few examples of recent data 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 for data licenses
- It's up to the data providers to establish offers and prices that make sense for them
- Our ambition is to establish LDS as the marketplace for European language data



Technical Helpdesk – Legal Helpdesk – Business Development Helpdesk



- Legal helpdesk: provides legal advice and assistance to all stakeholders involved in the LDS. It provides guidance for legal questions related to the collection and sharing of language data, licensing schemes, IPR clearance, data protection requirements as well as confidentiality aspects. These may include:
 - questions related to any kind of utilisation of language data in the LDS;
 - legal queries, e.g., how to exploit the LDS by clearing specific legal aspects, etc.



- **Technical helpdesk:** provides technical advice and assistance to all stakeholders involved in the LDS, providing first-line support to technical questions. These may include:
 - general troubleshooting and assistance on the usage of the LDS;
 - technical support, e.g., validation of technical compliance of data assets
- Business Development helpdesk helps monetise language data with experts in business development, pricing and revenues strategies, licensing issues for valuable assets, etc.



Next Steps

- LDS is in full swing: technical development, promotion, dissemination, governance etc.
- Testing: LDS software, metadata model, license and policy templates
- Adoption of LDS by industry and other organisations → grow the LDS User Group
- Collaborations with
 - DSSC, Simpl and ALT-EDIC
 - European projects, e.g., HPLT, OpenGPT-X, OpenWebSearch
 - Other relevant data spaces, especially Media (TEMS) and Cultural Heritage (Europeana)
 - EuroHPC Joint Undertaking
- Identify and make available new, fresh, interesting, novel, relevant language data, especially from industry and covering all European languages and modalities



LDS User Group

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

Meetings and Conferences

- Inaugural meeting in March 2024
- Second meeting in June 2024
- Third meeting in November 2024
- The LDS User Group is constantly growing
- Testing the LDS connector (currently 10+ instances)



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



LDS User Group (as of 05 March 2025)					
Individual Members	Unique Organisations				
198	158				

- Please join the LDS User Group or encourage members of your networks to join!
 - Validation of concepts, ideas, software; first test installations of the LDS connector (currently ongoing!); first trial exchanges of data; surveys etc.





Common European Language Data Space

Thank you!



A Common European Language Data Space – funded under contract LC-01936389 with the European Union. Khalid CHOUKRI (ELDA)

choukri@elda.org

LDS Country Workshop in Malta https://language-data-space.ec.europa.eu