# EUROPEAN LANGUAGE DATA SPACE

## European Language Data Space

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

29-01-2024 2nd Technology Workshop
https://language-data-space.ec.europa.eu

# Context: Large Language Models (LLMs)

- 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 (the only exception is English)

- A concerted effort for the collection of enormous amounts of language data for all European languages is very much needed

**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.

# Global LT/NLP Market is exploding: 439.85B$ by 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**

Market Analysis Report

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

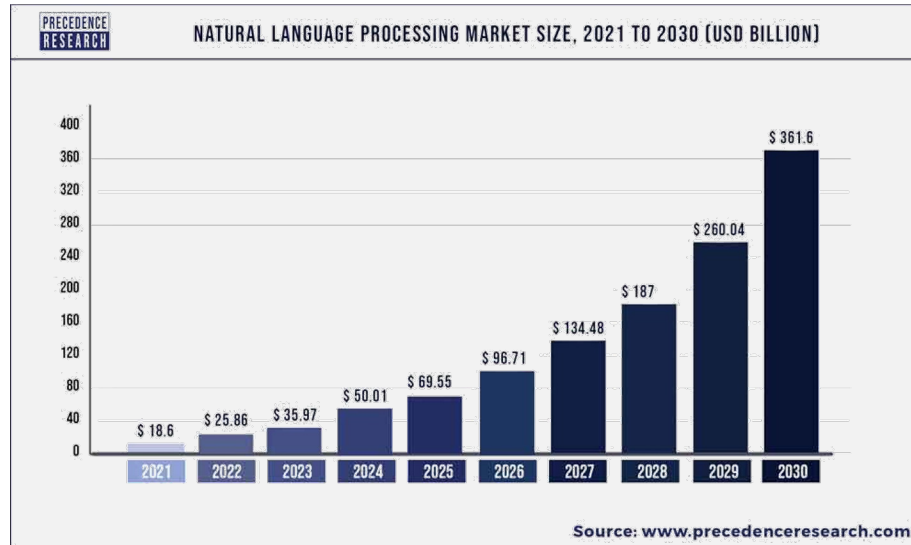### Natural Language Processing Market Report Scope

| 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)
- Dolbey Systems Inc. (US)
- Verint Systems Inc. (US)
- Net base Solutions Inc. (US)

All US!



PRECEDENCE RESEARCH — NATURAL LANGUAGE PROCESSING MARKET SIZE, 2021 TO 2030 (USD BILLION)

$18.6 (2021), $25.86 (2022), $35.97 (2023), $50.01 (2024), $69.55 (2025), $96.71 (2026), $134.48 (2027), $187 (2028), $260.04 (2029), $361.6 (2030)

Source: www.precedenceresearch.com

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

**Without a decisive intervention by the EU, Europe will be pushed further to the side lines in the global NLP market.**

# EU Data Strategy & Data Spaces

- Data Spaces are an inherent part of the EU Data Strategy and key instrument for the data economy
- The EU project DSSC is currently helping to develop the technical and operational specifications
  - IDSA, Gaia-X, FI-WARE, BDVA are members of DSSC – convergence of data space approaches!
- Important goals and properties of data spaces:
  - Technical sovereignty and data sovereignty
  - Establishing trust between participants
  - Data transactions can be free or for a certain price
  - Data providers/owners will be fully in control of their own data
  - Policies, standards, rules, contracts will be technically enforced
- The LDS is one of the currently 14 official EU data space projects – 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, selling 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

| 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 | 3pc | DE |
| Capgemini Deutschland GmbH | CapG | 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

| Centre of Excellence for Language Technologies (CELT) | Multi-Stakeholder Governance Body: LDS User Group |
|---|---|
| Mission: strategic | Mission: tactical/operational |
| Members: government representatives (from up to two different ministries from each Member State) | Members: companies, research centres and other organisations – also: identify other stakeholders (60%/40% private/public distribution of members) |
| Address collection, creation, sharing and re-use of data and models; aggregate initiatives and coordinate LDS governance scheme | Coordinate development of the blueprint, focus upon the technical building blocks and operational aspects of the LDS |

# Typical Use Cases of the LDS as a Data Space and Data Marketplace

**Organisations that develop NLP/LLM technologies**

• Check if data sets are available in LDS that cover a certain language and purchase as well as download those that address their need

• Check if data sets are available in LDS that cover a certain language and genre/register and purchase as well as download those that address their need

**Organisations that would like to provide language data**

• Make language data sets available, either for free or for a certain price; these organisations will *always* remain in control of their own data; terms of use can be based on templates or bespoke

**Data preprocessing through organisations that offer certain services**

• Certain data sets may need preprocessing before they are made available through LDS; such preprocessing services, e.g., anonymisation, can also be offered through LDS

# Large Language Models *for* Europe *made in* Europe – The Big Picture

**LLM Platform**
- User-friendly platform for pre-training, fine-tuning and deploying LLMs
- Can be realised on top of ELG (requires additional development work & HPC access)
- Huge gap that needs to be addressed soon



**Data Curation**
- High-performance tools for language data curation, filtering and annotation
- Filtering, quality and bias assessment, identification and analysis of various text and document properties
- Huge gap that needs to be addressed soon

**Federated Learning**
- Federated learning of LLMs by making use of the data sets provided in a decentralised way by LDS
- Huge gap that needs to be addressed soon

**NLP/LT Platform**
- General NLP and Language Technology Platform
- Provide general, i.e., non-LLM-based, NLP and LT services (e.g., anonymisation)
- Addressed by European Language Grid


EUROPEAN LANGUAGE GRID

**Compute**
- High-Performance Compute infrastructure for training LLMs
- Addressed by EuroHPC Joint Undertaking
- Simple yet crucial issue: storage capacity, esp. for multimodal data


EuroHPC Joint Undertaking

**Language Data**
- Secure and trusted sharing of language data covering all European languages
- Decentralised, i.e., federated
- Addressed by Language Data Space


EUROPEAN LANGUAGE DATA SPACE

# Next Steps

- LDS is in full swing: technical development, promotion, dissemination, governance etc.

- Collaboration with DSSC and collaboration with ALT-EDIC Working Group

- Collaboration with LLM projects such as HPLT and OpenGPT-X

- Collaboration with EuroHPC JU

- Collaboration with other data spaces, especially Media and Cultural Heritage

- Very important next step:

  - Raising awareness and adoption of LDS by industry and other organisations

  - Identify and make available new and fresh language data, especially from industry and ideally covering all European languages and modalities

# Subscribe to our newsletter!







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

**Common European Language Data Space**

# Thank you!

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

29-01-2024 2nd Technology Workshop
https://language-data-space.ec.europa.eu