# EUROPEAN LANGUAGE DATA SPACE

**Common European Language Data Space:**
**Developing a market for language data and services and benefitting from a joint European effort**

Firstname Lastname (Affiliation, Country)
name@domain.tld

# Large Language Models

# Context: Large Language Models

- Large language models are the most disruptive breakthrough in AI in recent history (BERT, GPT-3, ChatGPT, GPT-4 etc.)

- LLMs are based on vast amounts of training data

- LLMs use dozens, some even hundreds of terabytes of language and also image, video, audio etc. training data

- Europe's languages are vastly under-resourced, except English

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

- 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.

# 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**

Report ID: GVR-4-68040-020-4  |  Number of Pages: 100  |  Format: Electronic (PDF)

Historical Range: 2017 - 2021  |  Industry: Technology

Market Analysis Report

PDF

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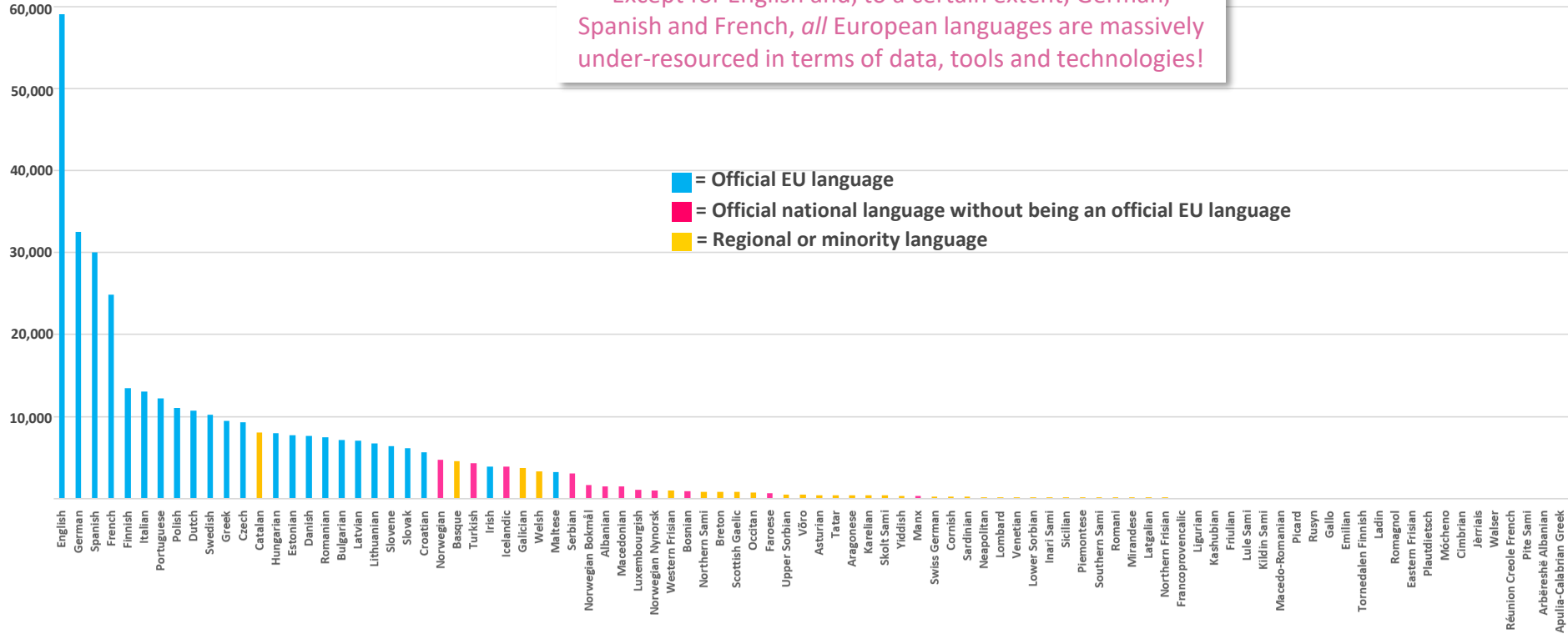- Verint Systems Inc. (US)
- Net base Solutions Inc. (US)

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

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

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

Common European Language Data Space

LDS

5

# Digital Language Equality?



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!

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

# EU Data Strategy & Data Spaces

- Data Spaces are an inherent part of the EU Data Strategy

- 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 approx. 20 official EU data space projects – focus on industry

# 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

# Establishment of the CELT and CELT+

| Establishment of the Centre of Excellence for Language Technologies (CELT) | Establishment of a Multi-Stakeholder Governance Body (CELT+) |
|---|---|
| Mission: strategic | Mission: tactical/operational |
| Members: governance representatives (from 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) |
| Nucleus of the future Language EDIC | Periphery of the future Language EDIC |
| 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 |
| Operational four months after inception meeting | Operational six months after inception meeting |
| Statutes ready in Q4 2023 | Statutes ready in Q4 2023/Q1 2024 |

# CELT and CELT+ Governance Bodies



**General Authorisation Framework (DSSC/EDIB)**

**Strategic Level**

**Tactical & Operational Level**

**Language Data Space – LDS**

## CELT
MS governments
**DEFINES**

## CELT+
Stakeholders from industry, academia, public administrations and NGOs
**IMPLEMENTS**

**GOVERNANCE building blocs**
- BUSINESS agreements
- OPERATIONAL agreements
- ORGANISATIONAL agreements

**TECHNICAL building blocs**
- data interoperability
- data sovereignty
- data value creation

# Previous Projects and Initiatives

- The four core partners – DFKI, ILSP, ELDA, TILDE – have been involved in many joint projects and initiatives, 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 in LDS** will be informed by ELG, ELRC-SHARE, META-SHARE.

# The LDS Decentralised Architecture

# Collaborations (selection)

- **DSSC** (Digital, EU; Community of Practice; Thematic Groups; Expert Groups)

- **OpenGPT-X** (Gaia-X; BMWK, Germany)

- **HPLT** (EU Horizon Europe)

- **DataBri-X** (IDSA; EU Horizon Europe)

- **European Language Grid** (ELG) – currently supported through OpenGPT-X, SciLake, DataBri-X – legal entity work in progress

- **European Language Equality** (ELE, EU PP/PA project)

- **INESData** (new language data space project in Spain; 65% of the 5M€ funding for industry for development of actual platform)

- **SciLake** (EOSC; EU Horizon Europe)

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

**Common European Language Data Space**

# Thank you!

Firstname Lastname (Affiliation, Country)
name@domain.tld

xx-xx-2023 Name of the event, City, Country
https://language-data-space.ec.europa.eu