# European Language Grid overview

Maria Giagkou
Institute for Language and Speech Processing, ATHENA Research and Innovation Centre
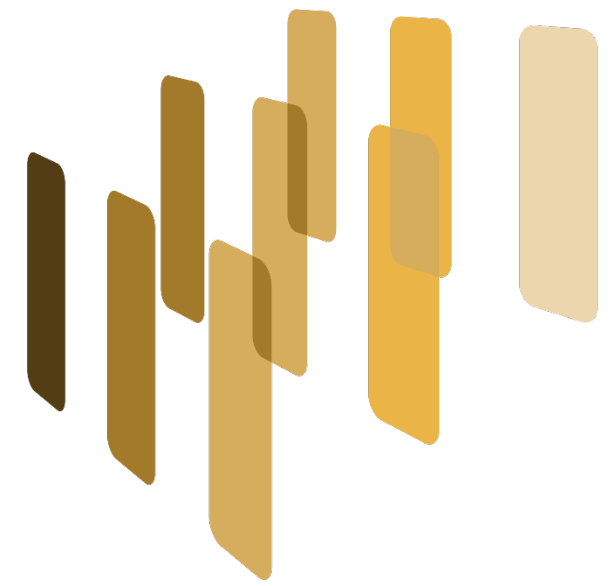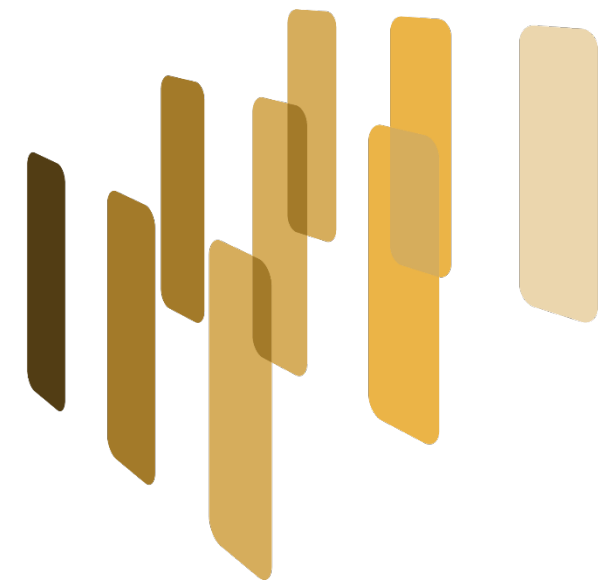
UniDive, WG3 meeting, September 8, 2023

# EUROPEAN LANGUAGE GRID

# European Language Grid **project**

- Project duration: 2019-2022

- Roots: META-NET Strategic Research Agenda for Multilingual Europe (2020)

  - 24 official EU languages, more than 60 official languages at the national and regional levels

  - Fragmented LT landscape: thousands for SMEs directly or indirectly active in the LT sector, hundreds of research and academic institutions

  - A "European Service Platform for Language Technologies" recommended

- Main aim: tackle the fragmentation of the European LT landscape

- Landmark outcome: the **European Language Grid platform**

# EUROPEAN LANGUAGE GRID



## European Language Grid **platform**

- Accessible at https://live.european-language-grid.eu/

- A platform for commercial and non-commercial LTs

- In support of developers and integrators of LTs as well as users/consumers of LT

- Enabling the European LT community to **upload** services and data sets, to **deploy** them and to connect with, and **use** resources made available by others

European Language Grid is a joint data, tools and services sharing platform

European Language Grid is a marketplace and **one-stop-shop** for the whole European Language Technology community

European Language Grid is the yellow pages of the European Language Technology community

# ELG catalogue in numbers (2023/09/05)

**15.199 Language Resources & Technologies** (LRTs), including

- **As per type:**

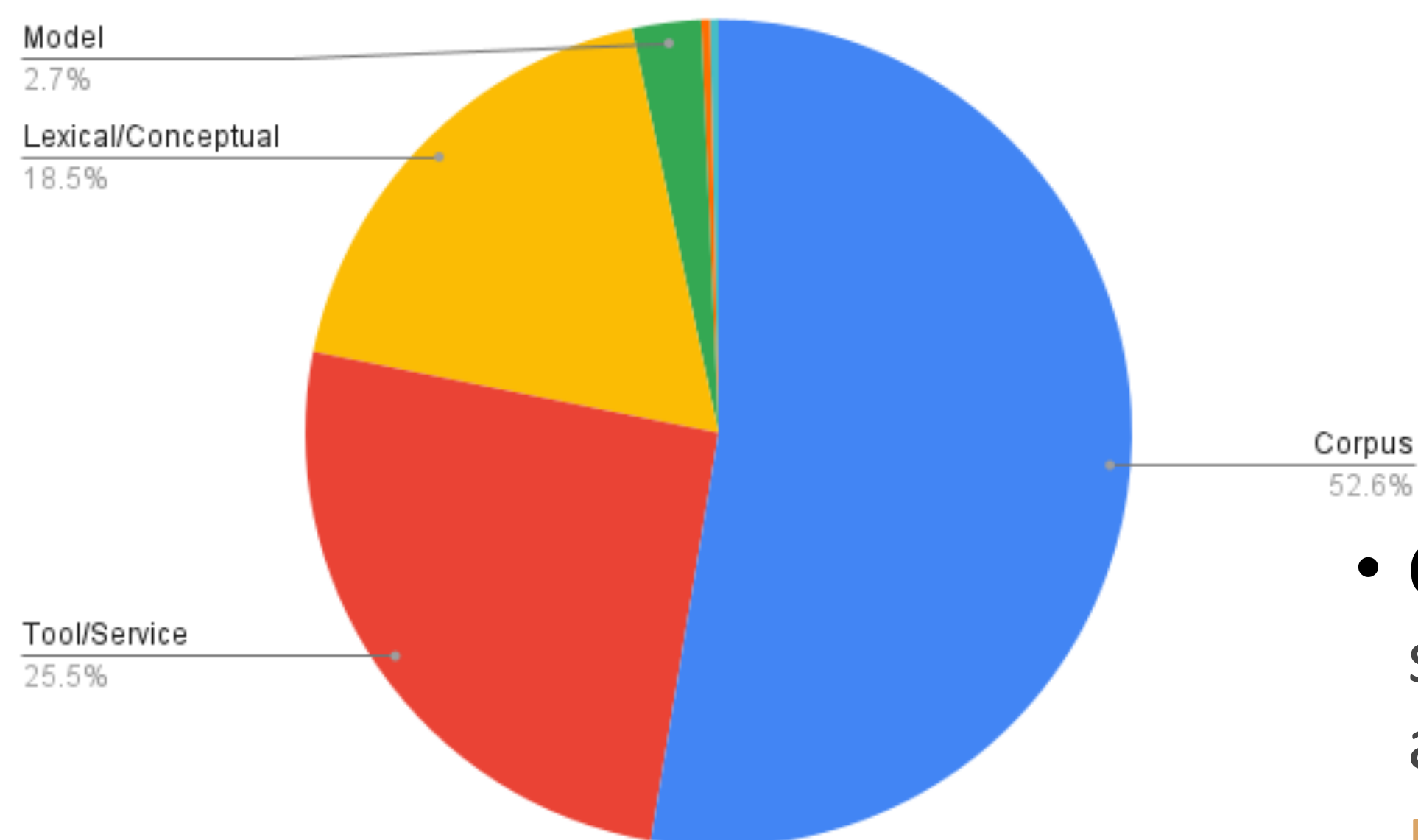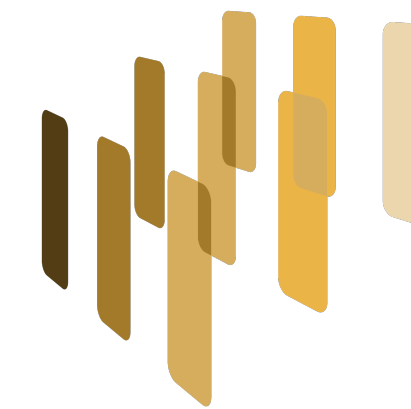  - 11.319 data resources, e.g. corpora (collections of data), models, lexical/conceptual resources (e.g., lexica, ontologies, terminologies, etc.), computational grammars

  - 3.880 language processing tools and services (e.g., MT services, speech synthesis/analysis tools, IE services, etc.),

- Of which, 1307 LRTs integrated in the ELG infrastructure: 1187 services running in the ELG cloud and 120 datasets uploaded and hosted in ELG; remaining records are "metadata-only records" providing information on the LRTs and re-directing consumers to the location where they can access them

- Related entities: 1775 Organizations involved in LT & 513 Projects

Model
2.7%

Lexical/Conceptual
18.5%

Corpus
52.6%

Tool/Service
25.5%

# ELG catalogue contents: Sources

- **Bridges** with existing platforms and infrastructures
  - Mainly in terms of metadata-based descriptions
  - Based on **open protocols** (OAI-PMH), or **APIs** offered by the platform or infrastructure providers
  - Respecting their own policies
- ELG also as infrastructural arm of ELE
  - using a mixture of automatic and collaborative population of the ELG catalogue

Browsing the ELG catalogue

Downloading a resource

Testing an MT service

Testing a dependency parser

# Provider's grid

- Providers can create new items
  - by validating and uploading a schema compliant metadata record (single or batch)
  - using an **interactive metadata editor**

# ELG metadata schema: main features

- **rich structured model** aiming to describe LRTs properties throughout their lifecycle, from design to (re-)use

- three versions:
  - maximal (full description) with required, recommended and optional properties
  - minimal (intended to support discovery) with required properties
  - relaxed (intended only for harvesting from more generic sources) with strictly required properties

- implementation: XSD re-using elements and values from OWL ontologies and controlled vocabularies in SKOS
  - MS-OWL (derived from the META-SHARE model)
  - OMTD-SHARE (catering for annotation/extraction types, data formats, methods & "LT taxonomy" used for service functions, intended application and LT areas for organisations and projects)

- Mappings with DataCite, DCAT, DC, (partially) schema.org

# ELG minimal metadata schema



**LANGUAGE RESOURCE / TECHNOLOGY**

**IDENTITY**
- Resource name
- Description
- Version

**CATEGORIES**
- Keyword

**CONTACT**
- Additional information

**DOCUMENTATION**

**RELATED LRT'S**

**TOOL/SERVICE**

**CATEGORIES**
- Function

**TECHNICAL**
- Language dependent
- Input content resource
  - Resource type
  - Language *
- Output resource *
  - Resource type
  - Language *

**EVALUATION**

**DISTRIBUTION**

**TECHNICAL**
- Software distribution form
- Private
- Docker download location *
- Download location *
- Access location *
- Execution location *
- Web service type *
- Licence

**DATA**

**DATA**

- [Corpus](#)
- [Model](#)
- [Lexical/Conceptual resource](#)
- [Tool/Service (not ELG-compatible)](#)

# ELG Support



- **Documentation**: https://european-language-grid.readthedocs.io
  - structured with the user in mind
    - Using, Contributing, Validating
  - examples and detailed technical guidance
  - recommendations on split of metadata records, data packaging and integration of services
  - continuously updated
- **Schema documentation** and ready-to-use templates and examples: https://gitlab.com/european-language-grid/platform/ELG-SHARE-schema

ELG tutorial video that explains how to make resources and tools available (https://youtu.be/29-V2EyMn4E).

# European Language Equality (ELE)

**Objective:** *development of a strategic research, innovation and deployment agenda to achieve digital language equality in Europe by 2030*

**Runtime:** 18 + 12 months (ELE1 and ELE2)

Ended 30 June 2023

**http://www.european-language-equality.eu**

# ELE Some of the main outcomes and findings

- 35 language reports, >40 languages, ~100 authors (D1.4-D1.36, D1.38-D1.39) (https://european-language-equality.eu/deliverables/)
- Strategic Research, Innovation and Implementation Agenda and Roadmap (https://european-language-equality.eu/agenda/)
- ELE book (https://link.springer.com/content/pdf/10.1007/978-3-031-28819-7.pdf?pdf=button)

# ELE Some of the main outcomes and findings

- Evidence-based investigation of the level of technology support per language
- Aggregate (meta)data through multiple strategies, including **manual population by ELE partners**
- Import these metadata to the ELG Catalogue and use the contents of the Catalogue to:
  - Compute the Digital Language Equality metric
  - DLE computations and visualisations through the DLE dashboard (https://live.european-language-grid.eu/catalogue/dashboard)

# Dashboard walkaround

- DLE technological and contextual scores
- Cross-language comparisons
  - As per resource types and features
  - Through
    - Histograms
    - Heatmaps
    - Radial bars
- Within-language comparisons
- Evolution of resources over time

# Language resources and technologies

- Corpora and lexical resources are the most numerous resource types, across all languages

- All languages seem to be better supported with translation technologies and text processing tools among all tools and services functions

- The least populated technologies are HCI, NLG and image/video processing

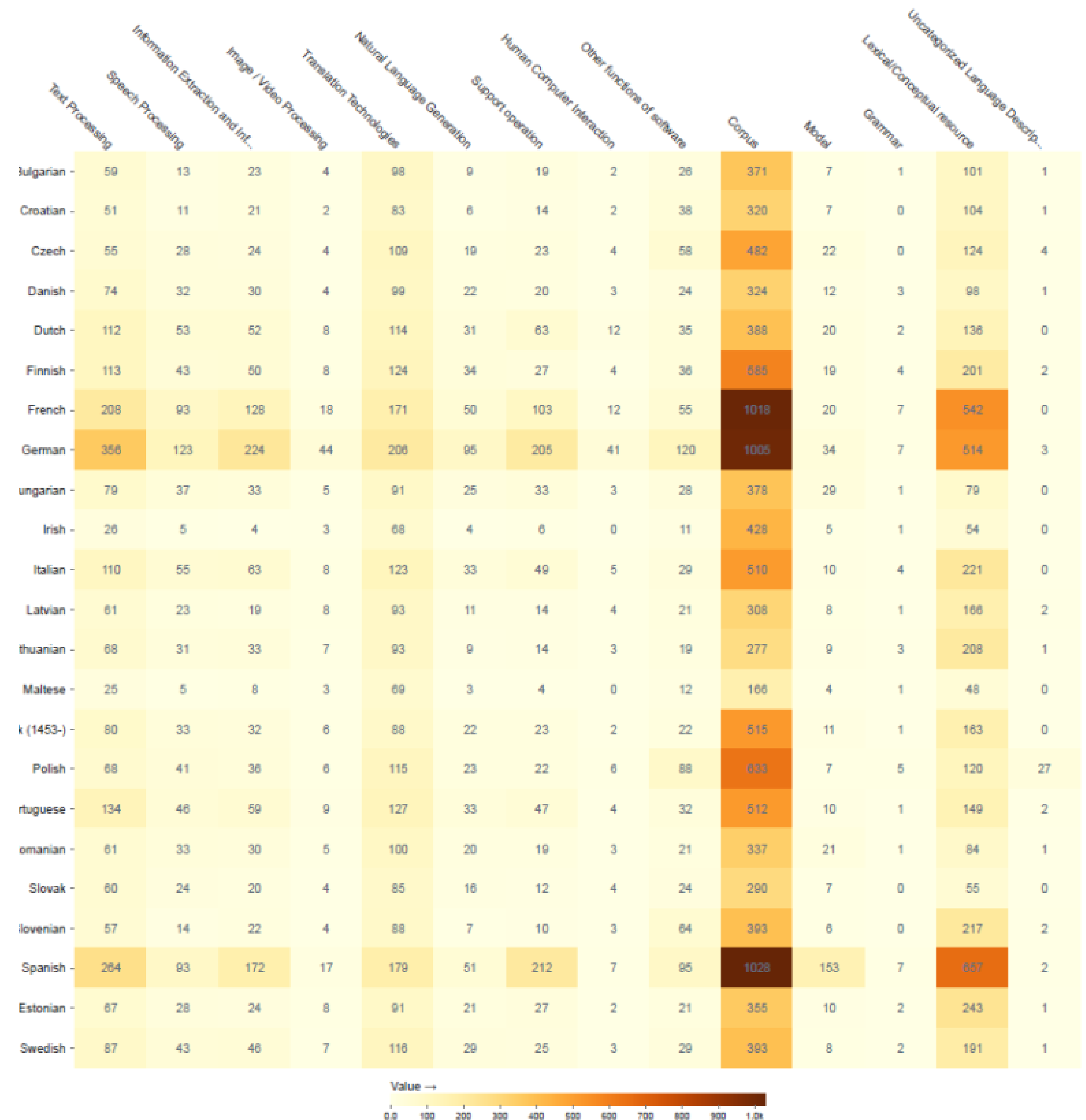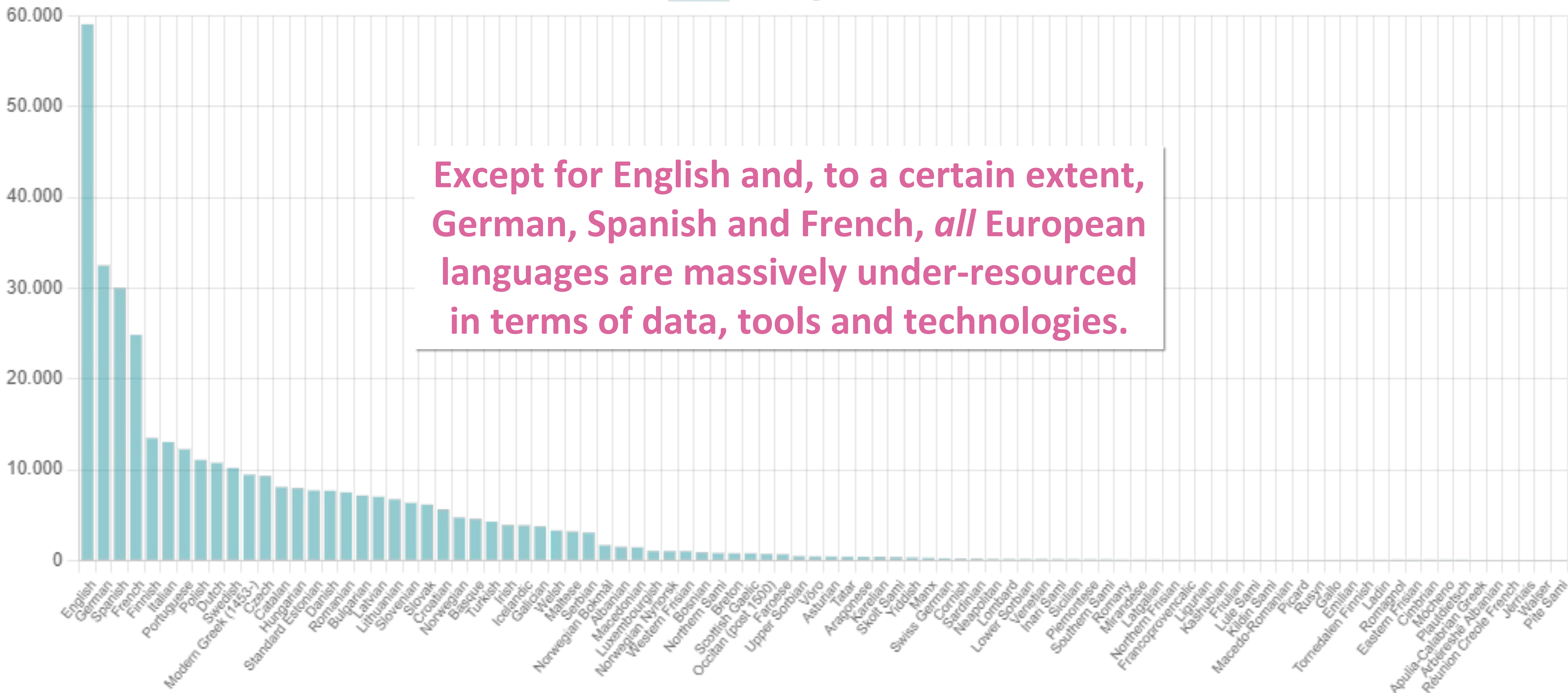| | Text Processing | Speech Processing | Information Extraction and Inf... | Image / Video Processing | Translation Technologies | Natural Language Generation | Support operation | Human Computer Interaction | Other functions of software | Corpus | Model | Grammar | Lexical/Conceptual resource | Uncategorized Language Descrip... |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Bulgarian | 59 | 13 | 23 | 4 | 98 | 9 | 19 | 2 | 26 | 371 | 7 | 1 | 101 | 1 |
| Croatian | 51 | 11 | 21 | 2 | 83 | 6 | 14 | 2 | 38 | 320 | 7 | 0 | 104 | 1 |
| Czech | 55 | 28 | 24 | 4 | 109 | 19 | 23 | 4 | 58 | 482 | 22 | 0 | 124 | 4 |
| Danish | 74 | 32 | 30 | 4 | 99 | 22 | 20 | 3 | 24 | 324 | 12 | 3 | 98 | 1 |
| Dutch | 112 | 53 | 52 | 8 | 114 | 31 | 63 | 12 | 35 | 388 | 20 | 2 | 136 | 0 |
| Finnish | 113 | 43 | 50 | 8 | 124 | 34 | 27 | 4 | 36 | 585 | 19 | 4 | 201 | 2 |
| French | 208 | 93 | 128 | 18 | 171 | 50 | 103 | 12 | 55 | 1018 | 20 | 7 | 542 | 0 |
| German | 356 | 123 | 224 | 44 | 206 | 95 | 205 | 41 | 120 | 1005 | 34 | 7 | 514 | 3 |
| Hungarian | 79 | 37 | 33 | 5 | 91 | 25 | 33 | 3 | 28 | 378 | 29 | 1 | 79 | 0 |
| Irish | 26 | 5 | 4 | 3 | 68 | 4 | 6 | 0 | 11 | 428 | 5 | | 54 | |
| Italian | 110 | 55 | 63 | 8 | 123 | 33 | 49 | 5 | 29 | 510 | 10 | 4 | 221 | 0 |
| Latvian | 61 | 23 | 19 | 8 | 93 | 11 | 14 | | 21 | 308 | 8 | 1 | 166 | 2 |
| Lithuanian | 68 | 31 | 33 | 7 | 93 | 9 | 14 | 3 | 19 | 277 | 9 | 3 | 208 | 1 |
| Maltese | 25 | 5 | 8 | 3 | 69 | 3 | 4 | 0 | 12 | 166 | 4 | | 48 | 0 |
| k (1453-) | 80 | 33 | 32 | 6 | 88 | 22 | 23 | 2 | 22 | 515 | 11 | 1 | 163 | 0 |
| Polish | 68 | 41 | 36 | 6 | 115 | 23 | 22 | 6 | 88 | 633 | 7 | 5 | 120 | 27 |
| Portuguese | 134 | 46 | 59 | 9 | 127 | 33 | 47 | 4 | 32 | 512 | 10 | 1 | 149 | 2 |
| Romanian | 61 | 33 | 30 | 5 | 100 | 20 | 19 | 3 | 21 | 337 | 21 | 1 | 84 | 1 |
| Slovak | 60 | 24 | 20 | 4 | 85 | 16 | 12 | 4 | 24 | 290 | 7 | 0 | 55 | 0 |
| Slovenian | 57 | 14 | 22 | 4 | 88 | 7 | 10 | 3 | 64 | 393 | 6 | 0 | 217 | 2 |
| Spanish | 264 | 93 | 172 | 17 | 179 | 51 | 212 | 7 | 95 | 1028 | 153 | 7 | 657 | 2 |
| Estonian | 67 | 28 | 24 | 8 | 91 | 21 | 27 | 2 | 21 | 355 | 10 | 2 | 243 | 1 |
| Swedish | 87 | 43 | 46 | 7 | 116 | 29 | 25 | 3 | 29 | 393 | 8 | 2 | 191 | 1 |

Value →

0.0   100   200   300   400   500   600   700   800   900   1.0k
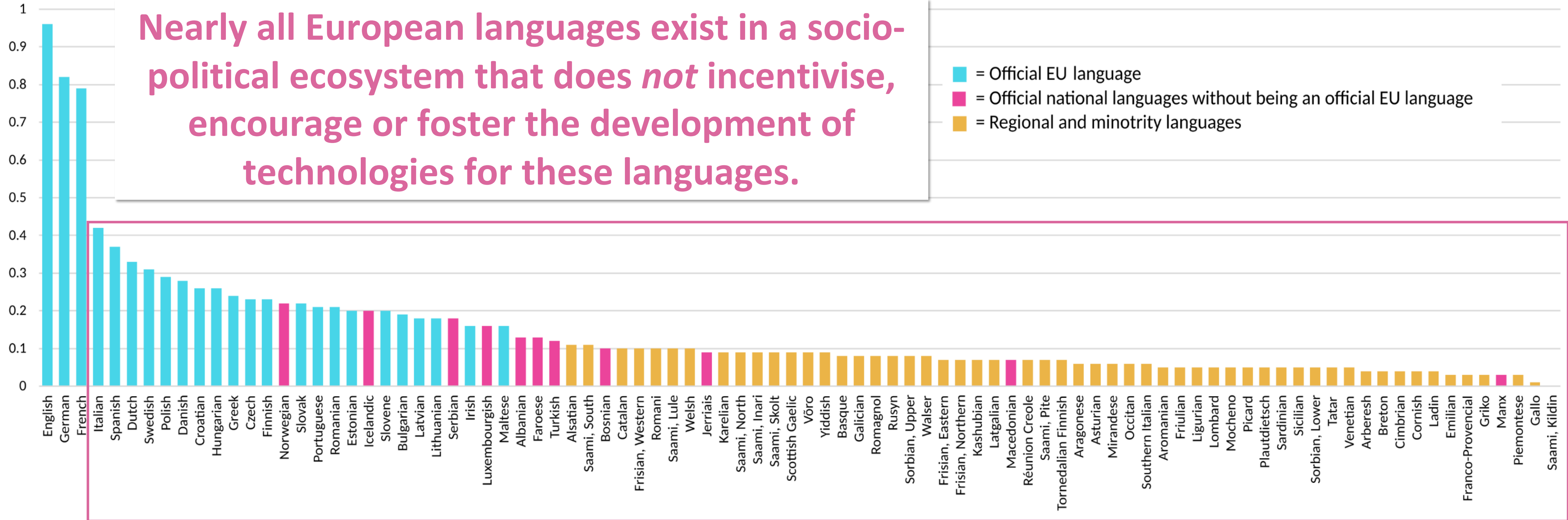
# DLE metric – technological scores



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.
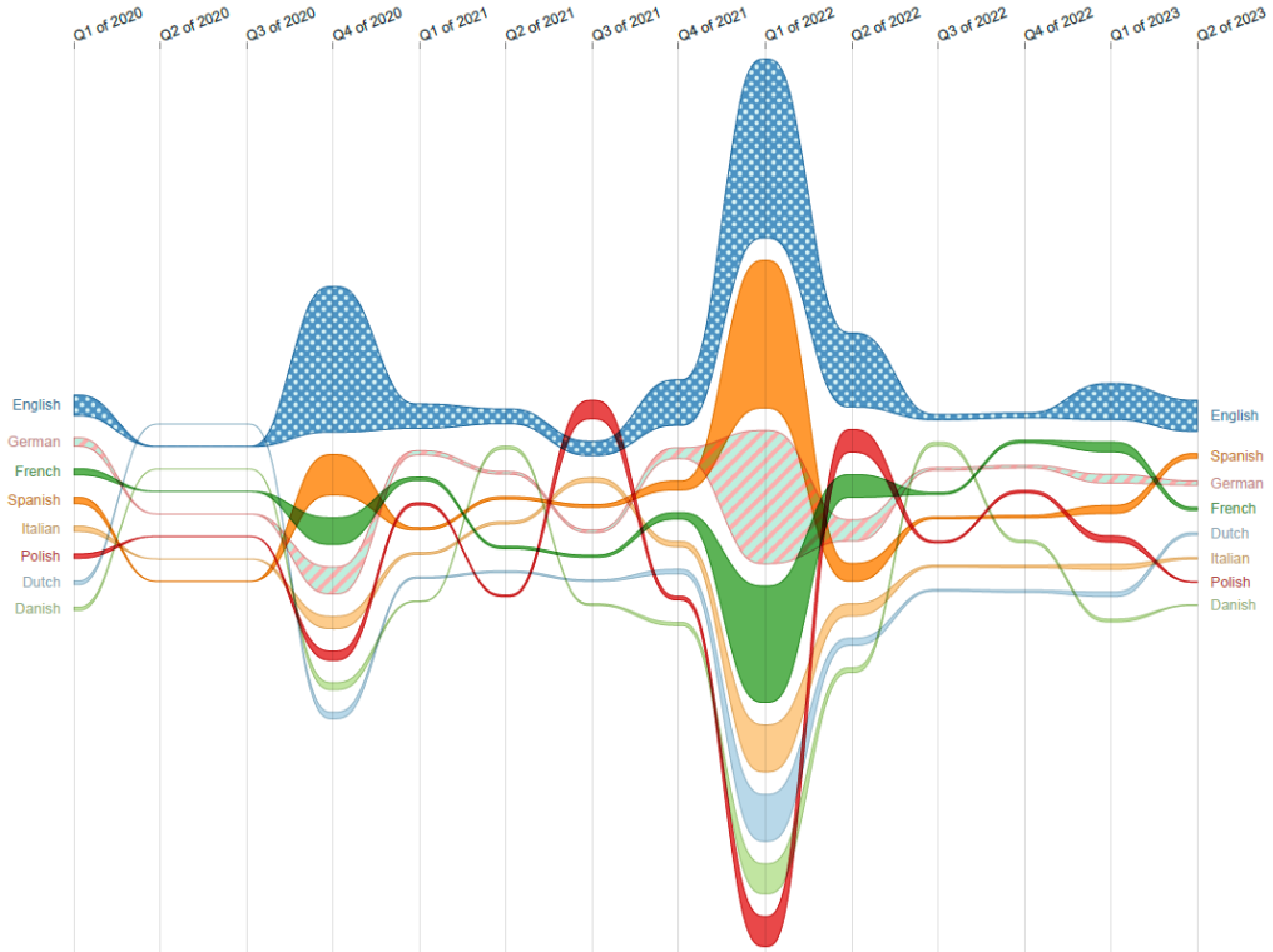
# DLE metric – contextual scores



Nearly all European languages exist in a socio-political ecosystem that does *not* incentivise, encourage or foster the development of technologies for these languages.

■ = Official EU language
■ = Official national languages without being an official EU language
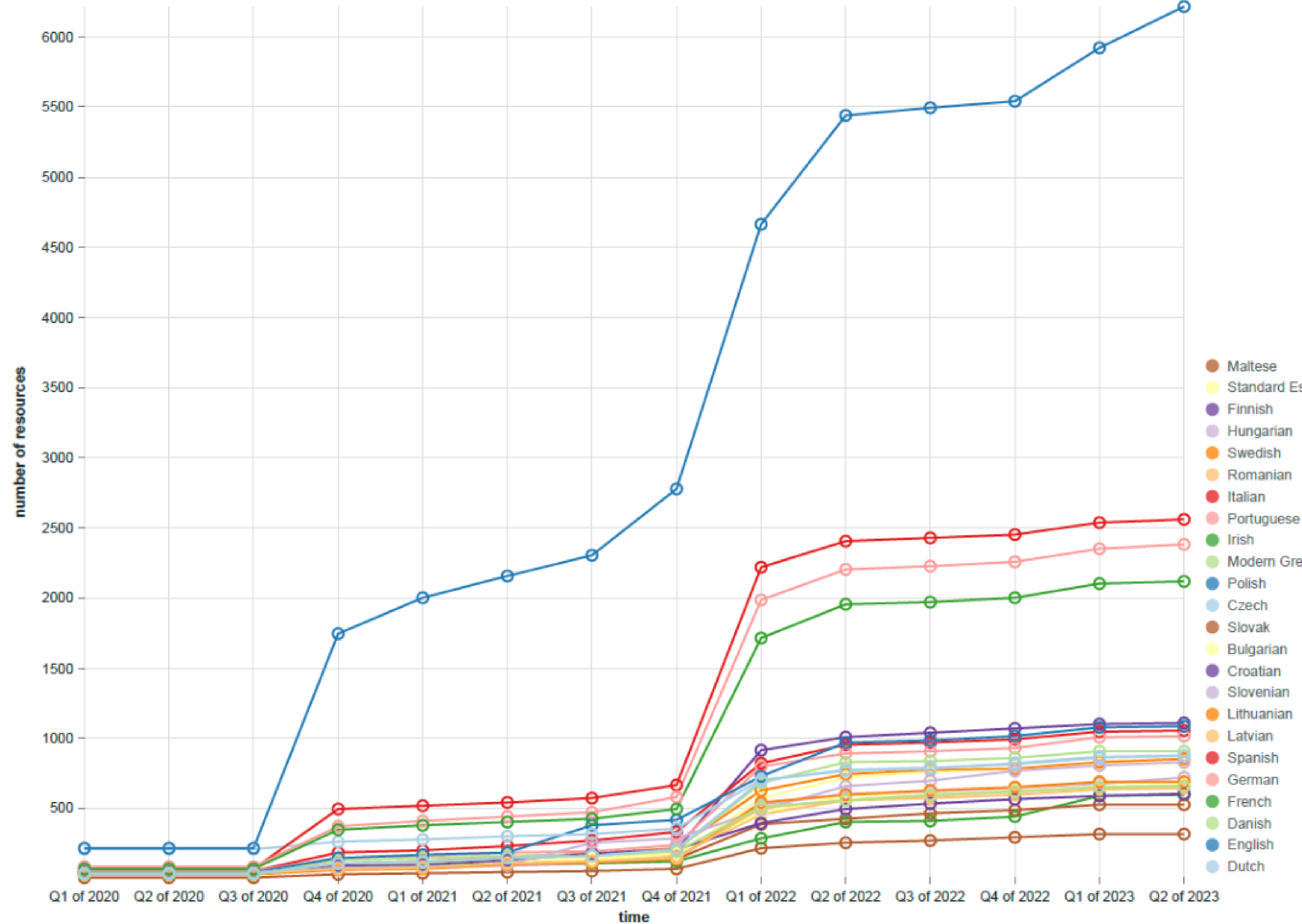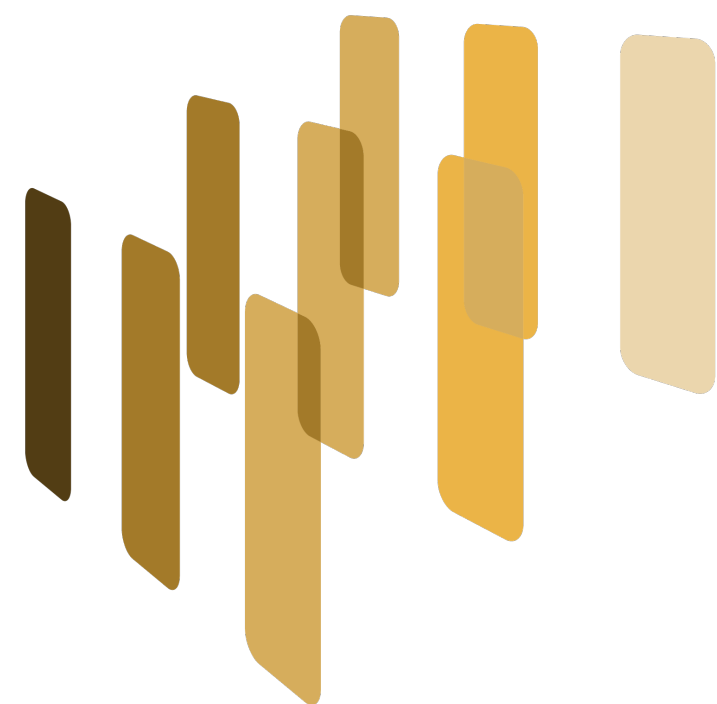■ = Regional and minotrity languages

# Evolution over time (1/2)

# Evolution over time (2/2)

- All languages are progressing

- English shows steeper progress in Q1-Q2 2023

- Distance has increased



ELG

# EUROPEAN LANGUAGE GRID

# EUROPEAN LANGUAGE EQUALITY

## Thank you!

**Maria Giagkou**

Institute for Language and Speech Processing / "Athena" Research and Innovation Center

www.ilsp.gr / www.athenarc.gr

mgiagkou@athenarc.gr