



# UniDive WG3

**Subgroup - Multilingual Tool and Resource Documentation**

**Co-leaders:**

**A. Seza Dođruöz (Ghent University)**

**Maria Giagkou (Athena RC)**

**Teresa Lynn (Mohamed bin Zayed University of Artificial Intelligence)**



## Task Overview

### Task 1

- Assess the “discoverability” of NLP tools and resources
- Who can participate?
  - Everyone 😊

### Task 2

- Analyse the NLP tool availability in the ELG catalogue
- Who can participate?
  - Excel or Tableau enthusiasts
  - Those with skills in data visualisation



## Task 1: Assessing the “discoverability” of NLP tools

- Kicked-off in Naples
- Template provided with instructions for semi-guided searches
  - Choose your language(s) and NLP task(s) of interest
  - Search for the relevant tools across a number of platforms
  - Report on the discoverability of desired tool (Could you find it easily? What challenges?)
  - Report on the metadata information available (was it sufficient and accurate?)
  - What metadata do you recommend should be provided for a similar search?
  - Is there a tool/ resource you are aware of that you can't find on these platforms?





## Observations

- 9 participants
- 68 search logs
- Repositories/platforms consulted:
  - ELG (27)
  - Clarin (24)
  - ELRA Catalogue (8)
  - LDC catalogue (6)
  - Hugging face (2)
  - <https://corpus-analysis.com/> (1)

→ Only platforms provided as examples were consulted - one exception

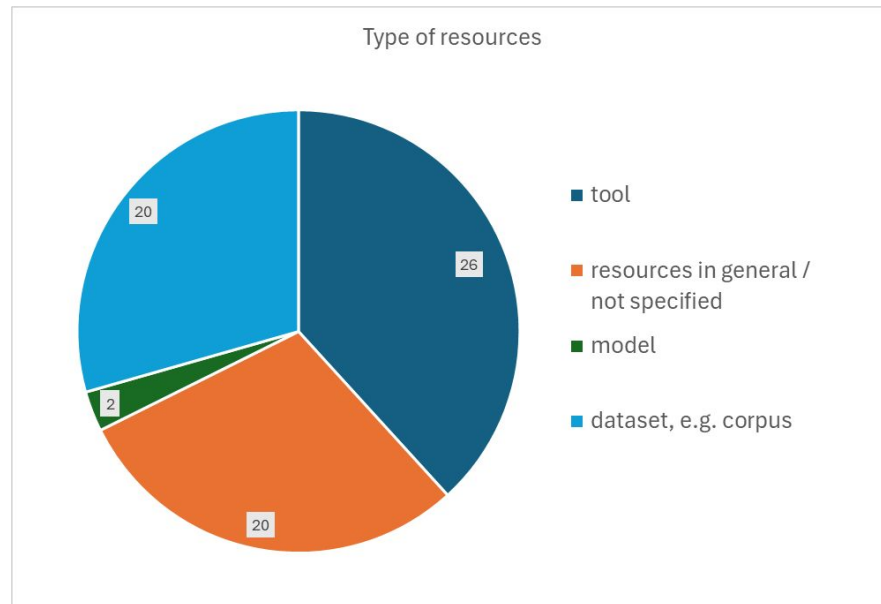
→ Perhaps we have been too prescriptive in our guidelines?

# Observations

- Languages searched for:
  - Ancient Egyptian
  - Ancient Italian
  - English
  - Florentine
  - Old Florentine
  - Old Italian
  - Polish
  - Portuguese
  - Serbian
  - Swedish
  - Turkish

→ the repositories suggested are not appropriate for diachronic linguistic studies (very few tools and resources for extinct languages)

- Resources searched for:



→ 7 logs (by one respondent) were attempts to actually run a NLP service offered through ELG and CLARIN → out of scope for this task



## Selected feedback

→ Models reported as missing

→ Some tools retrieved are outdated

→ Metadata issues reported:

- ❑ geographic variants not always available (e.g. Brazilian vs European Portuguese)
- ❑ missing time coverage to facilitate diachronic studies
- ❑ limited metadata specific to language models
- ❑ keywords not specific enough
- ❑ level of annotation

→ Many irrelevant results, low precision scores especially for free text searches

→ Respondents expected that tools would run natively on platform, not just described



**Thank you for your attention!**

**Questions?**