# TERMIKNOWLEDGE

EURA LEX XX 22

# Insights into data acquisition and data preparation for an online resource on COVID-19 terminology
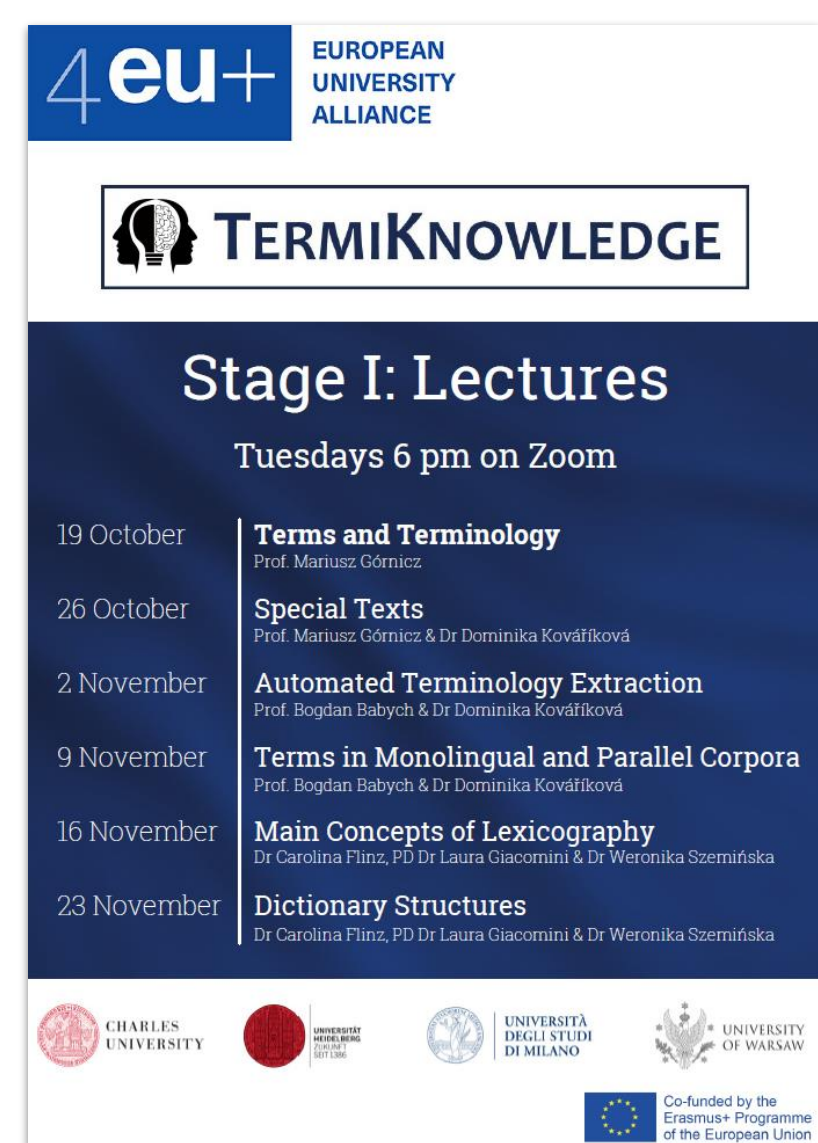
**Carolina Flinz** (carolina.flinz@unimi.it) – **Laura Giacomini** (laura.giacomini@iued.uni-heidelberg.de) – **Weronika Szemińska** (w.szeminska@uw.edu.pl)

## OBJECTIVES

- TermiKnowledge is an **international course in creating terminological resources** which took place between October 2021 and June 2022.

- Participants were **BA and MA students in translation, (corpus) linguistics**, **language and literature**

- The **goal** of TermiKnowledge was to teach students how to create a corpus, extract terms manually and using software, identify their equivalents in other languages, design a multilingual knowledge base and create terminological entries.

- The **4EU+ shared competencies** acquired during the course were:
  - data literacy – corpus work teaches students how to cope with specific terminological problems, use corpus tools and verify data;
  - multilingualism – the project offered an opportunity to compare concept systems and terminology across languages;
  - critical thinking – the students were involved in taking decisions during the entire lexicographic process (extent of the domain, headwords, access structure, entry structure etc.).

- The **multilingual knowledge base has been published on-line** for public use: https://terminology.mimuw.edu.pl/.

## ORGANISATION

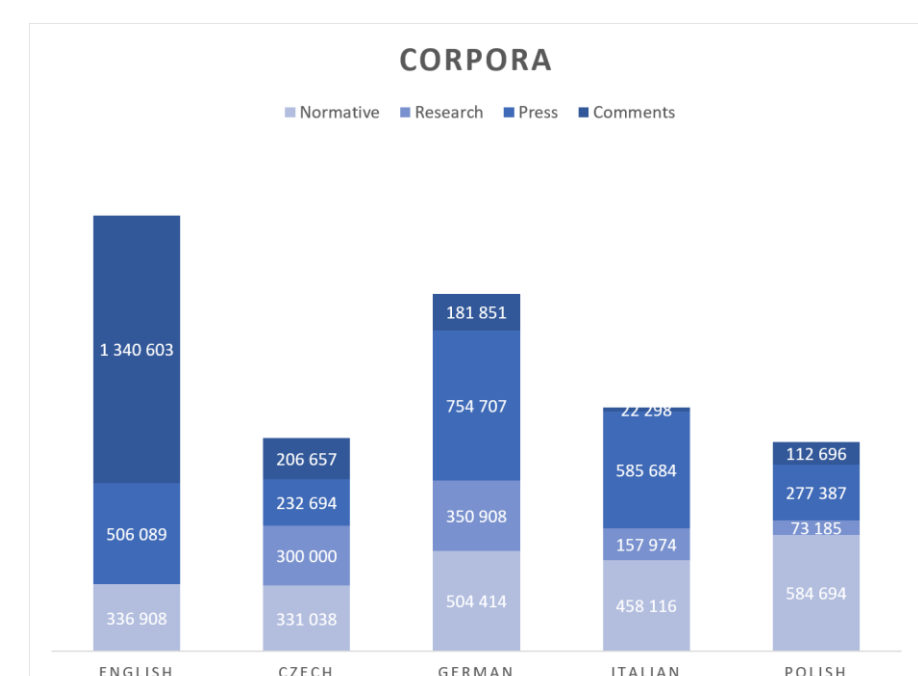- TermiKnowledge was run under the **4EU+ Alliance** by an international team of **experts in terminology, corpus linguistics, translation and knowledge representation** from **four European universities** (University of Warsaw – UW, Heidelberg University – HU, Charles University in Prague – CU, and Milan University – MU):

  - Mariusz Górnicz (PI, UW), Weronika Szemińska (UW), Laura Giacomini (HU), Bogdan Babych (HU), Dominika Kováříková (CU), Carolina Flinz (MU), Jerzy Tyszkiewicz (UW)
  - student assistants: Weronika Stefańska (UW), Pia Kruse (HU), Petr Louda (CU), Rita Luppi (MU)

- The project included **two courses**:
  - Semester I from mid-October 2021 to March 2022
  - Semester II from March 2022 to June 2022

- Each course consisted in
  - on-line lectures and research-based classes,
  - collaborative practical tasks on the creation of terminological resources as well as
  - short mobilities involving final works towards publication.

- Assessment was based on portfolios.



## THE LEXICOGRAPHIC PROCESS (SEMESTER I)

| a) CORPUS COMPILATION | b) TERMINOLOGY WORK | c) ABSTRACT MICROSTRUCTURE | d) ENTRY COMPILATION |

### a) CORPUS COMPILATION

- five languages: English, Czech, Italian, German, and Polish
- four corpora per language:
  - **Normative** – legal regulations such as EU directives, national regulations, and medical guidelines
  - **Research** – published research papers
  - **Press** – general press articles
  - **Comments** – texts in online comments sections under press articles



\* The already existing *Covid-19* corpus available on SketchEngine served as the English Research corpus. With its 224,061,570 words, it was by far the largest corpus in the project – and had to be omitted in the chart.

### b) TERMINOLOGY WORK

- Keyword extraction
- Keyword analysis and validation
- Contrastive analysis of keywords
- Headword selection (similar across corpora)
- Data preparation:
  - decision on the form of the headword
  - search for synonyms and variants
  - collocation extraction
  - search for definitions and encyclopaedic information
  - search for examples
  - preparation of notes

### c) ABSTRACT MICROSTRUCTURE

**TERM**
Contents
Related terms
NORMATIVE:
- Synonyms and variants
- Definition
- Encyclopaedic information
- Examples
- Collocations
- Keyness
- Note

RESEARCH:
- Synonyms and variants
- Definition/description
- Examples
- Collocations
- Keyness

PRESS:
- Synonyms and variants
- Description
- Examples
- Collocations
- Keyness

COMMENTS:
- Synonyms and variants
- Description
- Examples
- Collocations
- Keyness

### d) ENTRY COMPILATION

- The entries were compiled and stored in a dedicated application designed by a team of IT students from UW.
- The app was based on the familiar MediaWiki engine.
- All entries are subject to the CC-BY licence.
- The students worked on specific entry sections corresponding to the working group they belonged to / the type of corpus they worked on, i.e. Normative, Research, Press, or Comments.
- All students could propose related terms, which were then verified and added by the instructors or student assistants.

### ACCESS STRUCTURE



### ENTRY STRUCTURE



## CHALLENGES AND DISCOVERIES



What competences have you acquired during the TermiKnowledge project?

**REFLECTIONS AFTER SEMESTER I**
- student portfolios
- MentiMeter survey during mid-term meeting in Heidelberg

**CHANGES IN SEMESTER II**
- organisation
  - more workshops
  - no working groups
- primary data
  - one corpus per language
  - only normative and press texts
- access structure
  - separate entry lists for each language
  - list of all English entries with equivalent other language entries
  - links to all equivalent entries at the top of each article
- microstructure
  - no entry sections
- implementation
  - more refined term relations