# Word Sense Disambiguation and Crowdsourcing
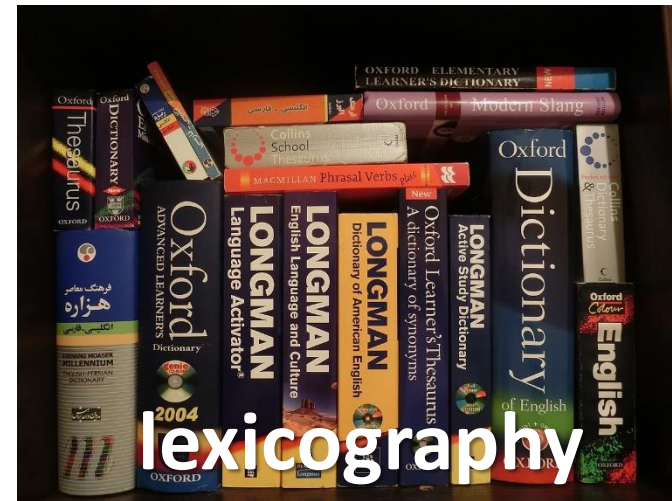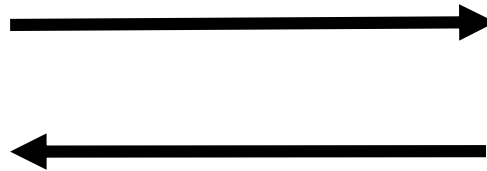
## Roberto Navigli

Dept. of Computer Science

Sapienza University of Rome

February 18th, 2019 - ELEXIS observer event, Vienna

# Lexicographic data for Natural Language Processing and vice versa

- We aim to show the **impact of lexicographic data for NLP**

- A paradigmatic task is Word Sense Disambiguation thanks to its need to leverage lexical-semantic knowledge resources

- However, we also aim to show that **NLP can help lexicography**

**NLP**

**lexicography**

# Multilingual Word Sense Disambiguation and Entity Linking

- **Word Sense Disambiguation:** automatic assignment of senses to words in context

- **Entity Linking:** automatic assignment of named entities to mentions in context

# Multilingual Word Sense Disambiguation and Entity Linking

- **Objective 1:** Develop algorithms that will use ELEXIS lexicographic resources to bootstrap disambiguation in a dozen languages

- **Objective 2:** Show high performance in many languages
  - Quantitative evaluation based on standard multilingual datasets (SemEval 2013; 2015 on multilingual WSD; Entity Linking datasets)
  - Perform validation in multiple languages and with different sense inventories: demonstrate **high-quality sense annotations**

# Challenges in WSD and Entity Linking

- **Issues:**
  - **The knowledge acquisition bottleneck:**
    - **Supervised approaches** suffer from **lack of annotated data** (only English and little else)
    - **Knowledge-based approaches** need computational lexicons, semantic collocations, graph-like dictionary structure, etc.
  - **Reference inventories**
    - WordNet is too fine grained
    - Wikipedia is too rich
- The ELEXIS dictionary matrix will prove **important benefits for both issues**

elexis

# Workplan (1/2)

- **Textual data** from:
  - the Universal Dependencies project (POS tagged)
  - the *TenTen corpora from Lexical Computing
  - Semantically-annotated corpora from partners

- **Phase 1a (October 2018/February 2019):**
  - **Algorithms:** Babelfy (Uniroma1) + Wikifier (JSI)
  - **Inventory:** use existing inventories (BabelNet, Wikipedia)
  - **Validation:** show the data to lexicographers in ELEXIS + observers
  - **Goal:** prepare the framework

- **Phase 1b (February 2019/June 2019):**
  - **Disambiguation** of the corpora + analysis

# BabelNet: a shared multilingual inventory of meanings

- **Multilingual**: the same concept in tens of languages

- It **integrates different kinds of open resources**, such as WordNet, Wikipedia, Wikidata, Wiktionary, etc.

- **Wide coverage**: 284 languages and 16 million entries!

- Used by **more than 800 universities and research centers**!

# Disambiguation: Babelfy

- We used Babelfy for disambiguating the Wikipedia corpus

- Why?
  - The first (and only) system that performs Word Sense Disambiguation (common nouns, verbs, adjectives, adverbs) and Entity Linking (names) **jointly**

# Disambiguation: Babelfy

- We used Babelfy for disambiguating the Wikipedia corpus
- Why?
  - The first (and only) system that performs Word Sense Disambiguation (common nouns, verbs, adjectives, adverbs) and Entity Linking (names) **jointly**
  - **Knowledge-based:** does not need millions of sentences annotated in each language
  - Works in **arbitrary languages** (284 languages)
  - Can disambiguate **texts written in mixed languages** (language-agnostic setting)

fy **B** abelfy

"Word sense disambiguation and entity linking together!"

elexis

Allianz Arena - Munich

hitter#n#2 - strikers

Thomas Needham

Mario Balotelli - Mario

Mario Kempes - Mario

Thomas Hitzlsperger - Thomas

Thomas Strakosha - Thomas

Mario Vasili - Mario

Thomas Flath - Thomas

Mario Gómez - Mario

Central defender - strikers

striker#n#1 - strikers

Bayern Munich - Munich

Bayern Munich Junior Team - Munich

Munich#n#1 - Munich

Thomas Pfannkuch - Thomas

forward#n#1 - strikers

Aquinas#n#1 - Thomas

Thomas Ravelli - Thomas

Thomas Fernandez - Thomas

José Mourinho - Mario

Thomas Müller - Thomas

Mario Götze - Mario

Mario Mandžukić - Mario

History of FC Bayern Mun

Mario Erb - Mario

TSV 1860 München - Munich

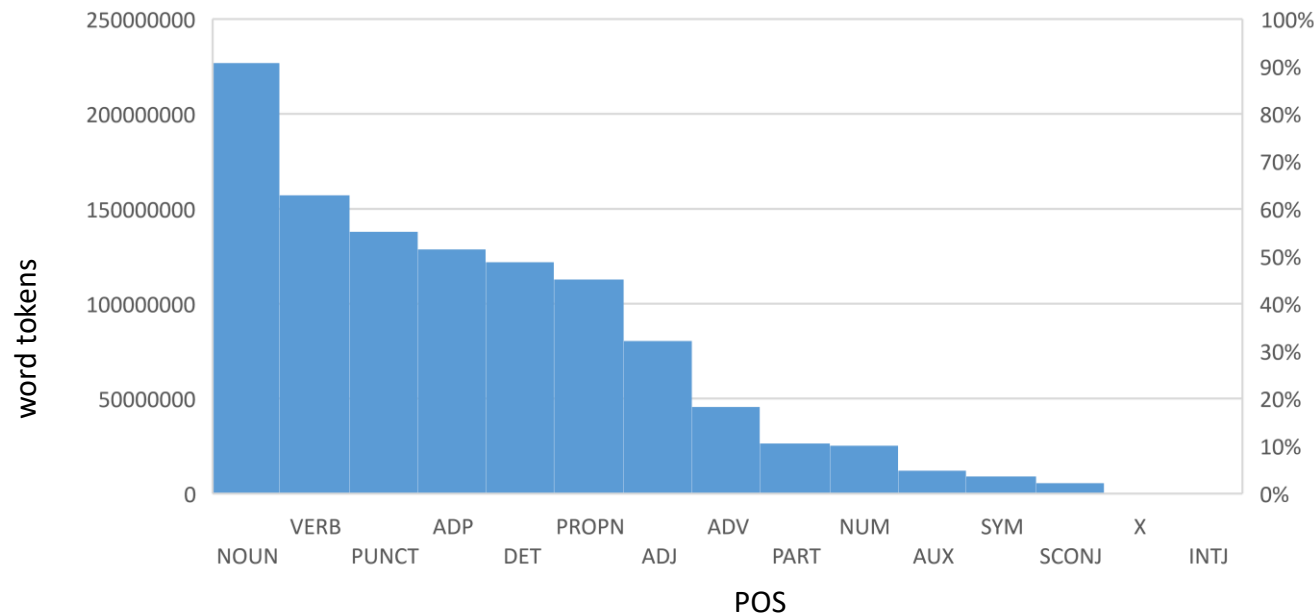# Statistics on the Babelfied English Datasets from UD

- Total number of word tokens: 488515

- Total number of word types: 26892

- Total number of disambiguated word tokens: 85094

- Total number of disambiguated word types: 24144

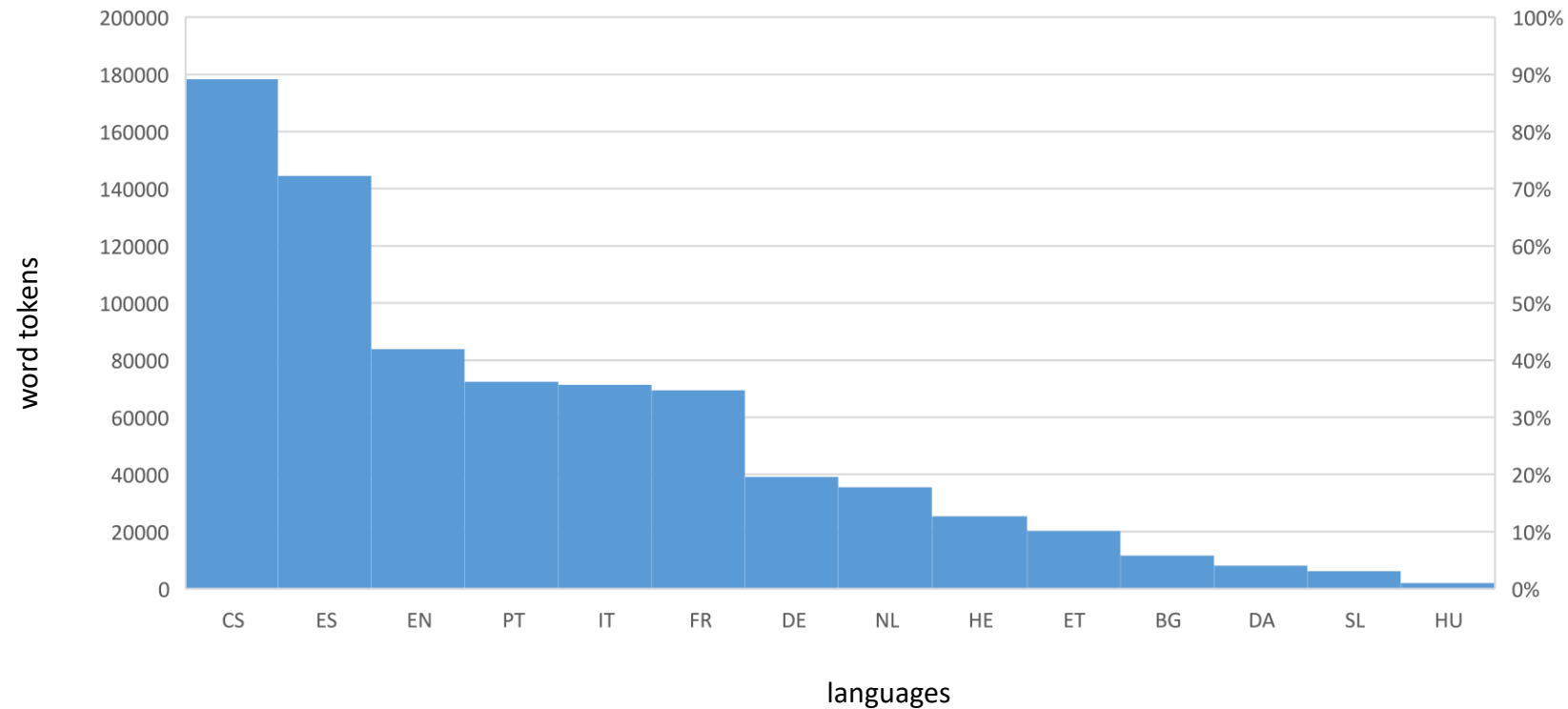# Statistics on the Babelfied English Datasets from the *TenTen Corpora
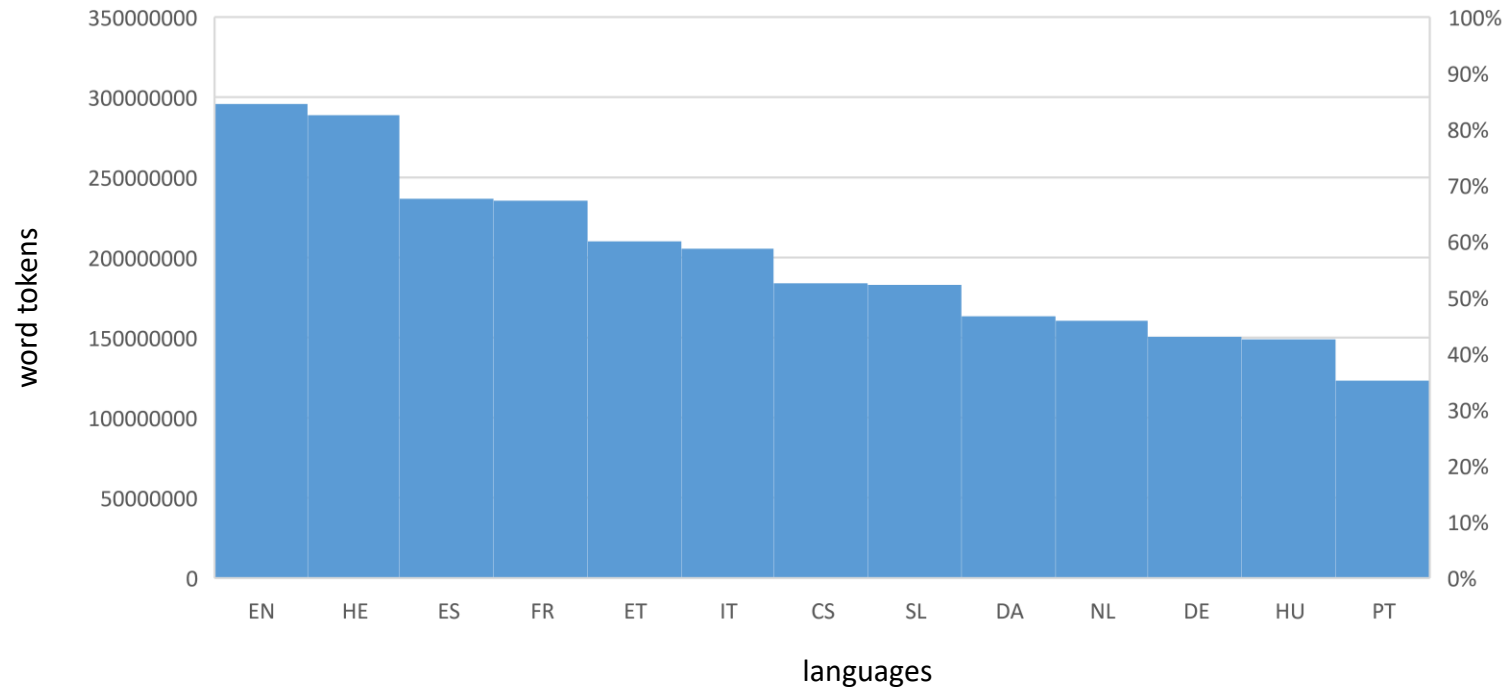
- Total number of word tokens: 1170202338

- Total number of word types: 5454205

- Total number of disambiguated word tokens: 295783220

- Total number of disambiguated word types: 844705

# Overall statistics on words Babelfied across languages of UD Corpora

# Overall statistics on words Babelfied across languages of *TenTen Corpora

# Qualitative analysis + decisions about lexicographic evaluation

- We are in the process of defining the evaluation process

- Shall we check if the best option was chosen among those in the inventory?

- Annotations (potential options: correct, so and so, wrong?)

- Evaluation measures

# Workplan (2/2)

- **Phase 2 (June 2019-2021):**
  - **Algorithms:** New algorithm with multilingual sense embeddings + semantic graphs (Uniroma1)
  - **Inventory:** use the ELEXIS dictionary matrix from WP2
  - **Validation:** show the data to lexicographers in ELEXIS + observers
  - **Goals:**
    1. show we can now disambiguate in arbitrary languages with reputable dictionaries
    2. show improvements coming from the dictionary matrix resulting from WP2

# Multilingual semantic parsing

- Semantic parsing is the task of mapping sentences to a formal representation
  - Abstract Meaning Representation (AMR)
  - Universal Conceptual Cognitive Annotation (UCCA)
  - CCG-based like Discourse Representation Structures (DRS)

**The boy wants to visit New York City**

```
(w / want-01
   :ARG0 (b / boy)
   :ARG1 (g / visit-01
            :ARG0 b
            :ARG1 (c / city
                     :name (n / name
                               :op1 "New"
                               :op2 "York"
                               :op3 "City")))))
```

# Multilingual semantic parsing

- Semantic parsing is the task of mapping sentences to a formal representation

- **Objective 1:** develop algorithms for semantic parsing in multiple languages which take advantage of ELEXIS lexicographic data

- **Objective 2:** exploit bilingual and multilingual data to innovate semantic parsing algorithms

- **Expectations from other partners:** creation of a multilingual test set benchmark (à la SemEval) for the task; curation/validation of verb frames for parsing in different languages based on ELEXIS data

elexis

# Lexical-semantic analytics for NLP

- Based on analytics computed on the ELEXIS resources for words, phrases, collocations, senses, domains, etc. we will explore three directions:

  - **T3.3.1 Sense clustering:** semi-automatic algorithms to group fine-grained sense distinctions, also across languages

| race#n (WordNet) | |
| --- | --- |
| #1 | Any competition (→ contest). |
| #2 | People who are believed to belong to the same genetic stock (→ group). |
| #3 | A contest of speed (→ contest). |
| #4 | The flow of air that is driven backwards by an aircraft propeller (→ flow). |
| #5 | A taxonomic group that is a division of a species; usually arises as a consequence of geographical isolation within a species (→ taxonomic group). |
| #6 | A canal for a current of water (→ canal). |

| race#n (ODE) | |
| --- | --- |
| #1.1 | **Core:** SPORT A competition between runners, horses, vehicles, etc. ● RACING A series of such competitions for horses or dogs ● A situation in which individuals or groups compete (→ contest) ● ASTRONOMY The course of the sun or moon through the heavens (→ trajectory). |
| #1.2 | **Core:** NAUTICAL A strong or rapid current (→ flow). |
| #1.3 | **Core:** A groove, channel, or passage. ● MECHANICS A water channel ● Smooth groove or guide for balls (→ indentation, conduit) ● FARMING Fenced passageway in a stockyard (→ route) ● TEXTILES The channel along which the shuttle moves. |
| #2.1 | **Core:** ANTHROPOLOGY Division of humankind (→ ethnic group). ● The condition of belonging to a racial division or group ● A group of people sharing the same culture, history, language ● BIOLOGY A group of people descended from a common ancestor. |
| #3.1 | **Core:** BOTANY, FOOD A ginger root (→ plant part). |

# Lexical-semantic analytics for NLP

- Based on analytics computed on the ELEXIS resources for words, phrases, collocations, senses, domains, etc. we will explore three directions:

  - **T3.3.2 Domain labeling of text:** ELEXIS resources shown to improve domain labeling across languages

# Lexical-semantic analytics for NLP

- Based on analytics computed on the ELEXIS resources for words, phrases, collocations, senses, domains, etc. we will explore three directions:
  - **T.3.3.3 Diachronic distribution of senses:** sense frequency ranking over time across resources (Most Frequent Sense is a strong baseline in WSD)

# Challenges in lexical-semantic analytics

- **Sense clustering:**
  - Fine granularity
  - Not obvious what a good cluster of senses is

- **Domain labeling of text:**
  - Elicit information from ELEXIS resources (what is a good set of domain labels? which resources provide domain-specific content?)
  - Work in dozens of languages

- **Diachronic distribution of senses:**
  - Create reliable distributions of senses in many languages
  - Leverage such distributions in WSD and Entity Linking

- **See interaction with WP4**

elexis

# Crowdsourcing and gamification

- Objectives:
  - Validating the output of WP2 (links between resources, the dictionary matrix)
  - Validating and improving the data produced by WP3 and WP4

- Proposal for a crossword game developed jointly with Babelscape
- Other crowdsourcing efforts are on-going

- Goal: collect experiences from the consortium and observers

# Questions?