# Lecture 11: Linguistic Annotation

LING 1340/2340: Data Science for Linguists

Na-Rae Han

# Objectives

- AMR review

- Linguistic annotation

  - Types of linguistic annotation

  - Annotation formats

  - Annotation tools

  - How to plan and run an annotation project

    - An anatomy of annotation project

# AMR example

▸ Guidelines:

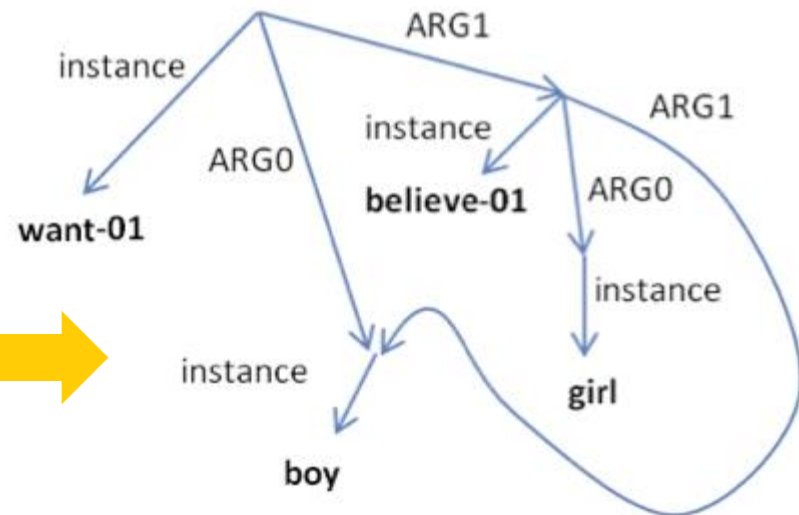◆ https://github.com/amrisi/amr-guidelines/blob/master/amr.md

The boy desires the girl to believe him.

The boy desires to be believed by the girl.

The boy has a desire to be believed by the girl.

The boy's desire is for the girl to believe him.

The boy is desirous of the girl believing him.



```
(w / want-01
    :ARG0 (b / boy)
    :ARG1 (b2 / believe-01
                :ARG0 (g / girl)
                :ARG1 b))
```

# AMR annotated corpora

▸ [https://amr.isi.edu/download.html](https://amr.isi.edu/download.html)

▸ *The Little Prince* by Antoine de Saint-Exupéry is annotated in AMR in full.

- ◆ English
- ◆ Chinese

▸ Why build such corpora?

# Linguistic annotation: what types?

▸ **What types of linguistic annotation have we seen so far?**

▸ GUM: **The Georgetown University Multilayer Corpus**

  ◆ https://gucorpling.org/gum/index.html

  ◆ A corpus with *all* levels of linguistic knowledge annotated!!

# *Using* GUM

- ▶ How to explore and use the GUM corpus?
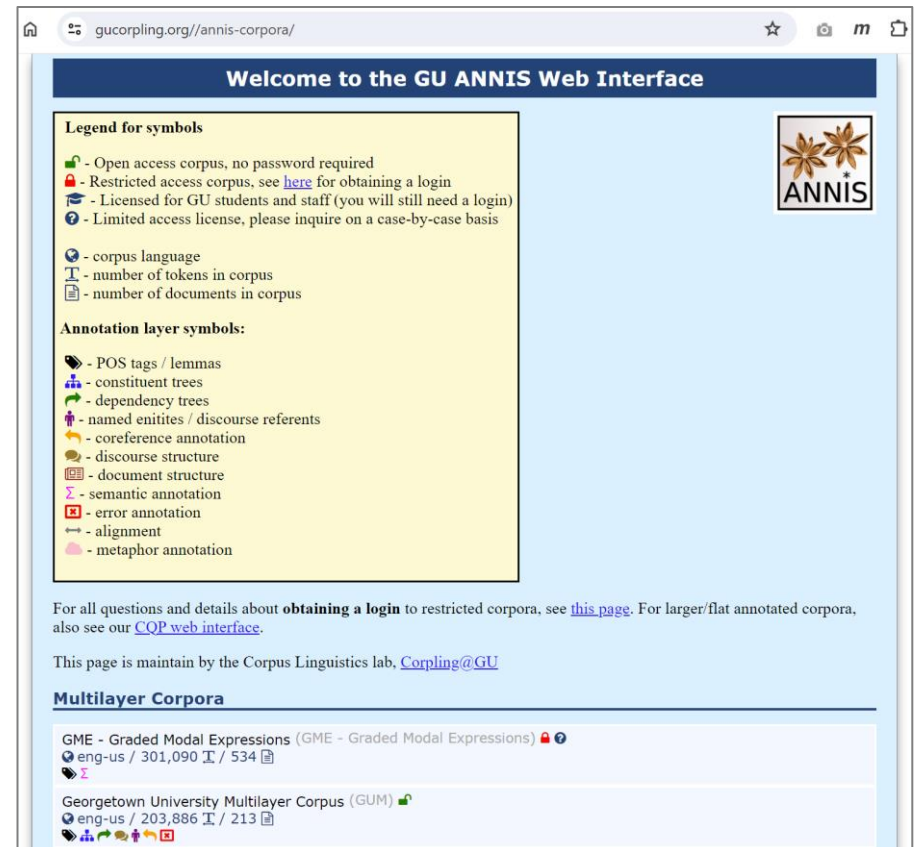    - ◆ Download the source on GitHub: https://github.com/amir-zeldes/gum then process it yourself

    - ◆ Or: use the **ANNIS** interface
        - ◆ A web browser-based search and visualization architecture for complex multilayer linguistic corpora with diverse types of annotation
        - ◆ https://corpus-tools.org/annis/
        - ◆ GU ANNIS Web Interface: https://gucorpling.org/annis-corpora/

# Why annotate?

Why annotate text with linguistic information?

▶ **Development and testing of linguistic theories**

⬅ Assists empirical linguistic inquiries

▶ **Develop and evaluate (statistically based) NLP technologies**

⬅ Becomes the basis of "language models" in NLP applications

⬅ Linguistic annotation represents linguistic knowledge of humans that AI agents learn through machine learning, which they then mimic

# What are linguists' roles in all this?

- Doing the annotation
  - Linguistics undergrads and grads make excellent annotators.

- Leading annotation projects
  - Design annotation schemes
  - Develop annotation guidelines
  - Train and supervise annotators
  - An example https://www.ldc.upenn.edu/sites/www.ldc.upenn.edu/files/penn-etb-2-style-guidelines.pdf

- As part of the NLP community, help keep linguistic knowledge representation in balance with engineering-side considerations

- Be a USER of linguistically annotated data by conducting empirical research
  - An example: https://web.stanford.edu/~bresnan/qs-submit.pdf

- Increasingly: Be a community-minded steward of language data. Address concerns of ethics and representation.

# Wrapping up

▶ **Next class**

- ◆ An anatomy of annotation project

- ◆ Annotation wrap

- ◆ Machine learning: regression

▶ **Your project**

- ◆ Progress Report #1 due this Friday!