# Lecture 8: Linguistic Data and Annotation

LING 1340/2340: Data Science for Linguists

Na-Rae Han

# Objectives

▶ **Your term project**

- ◆ Plan submitted, repo created!

▶ **Linguistic annotation**

- ◆ Types of linguistic annotation
- ◆ Annotation formats

# Your term project

▶ Everyone's project repo is at our GitHub org.

▶ First progress report is due next!

  ◆ Focus on **data**: sourcing, curation and cleaning

▶ Managing your data

  ◆ You will be manipulating and processing your data.

  ◆ Should you include your data set in your GitHub repo?

    ◆ Depends on your license!

# Linguistic data, corpus, annotation, coding



▶ **How language exists in the world:**

▶ **How it is captured, becomes language DATA:**

- ◆ RAW data: spoken, written, visual
- ◆ Experimental data

▶ **Added (layers of) information:**

- ◆ Metadata
- ◆ Transcription, time alignment
- ◆ Acoustic measurements
- ◆ Annotation
- ◆ Coding

What are the differences?

# 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
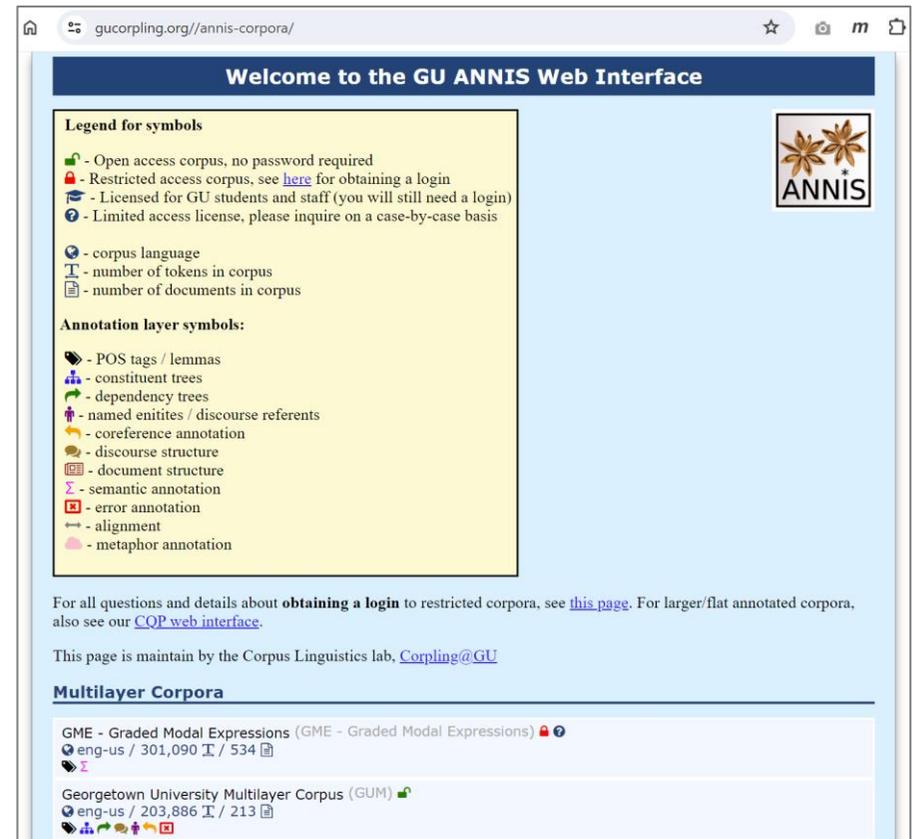      - ☐ https://corpus-tools.org/annis/
  - ◆ GU's ANNIS Web Interface: https://gucorpling.org/annis-corpora/
    - ◆ Example queries on this page (scroll down): https://gucorpling.org/gum/search.html

# Why annotate?

Why annotate text with linguistic information?

▸ Development and testing of linguistic theories, analysis

        ← 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
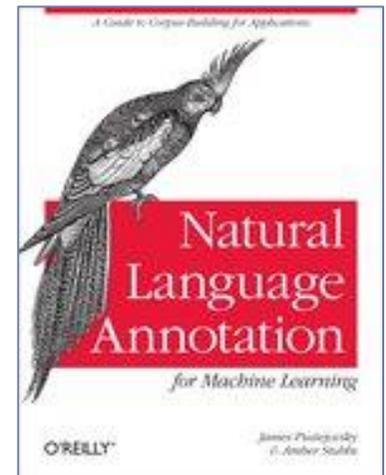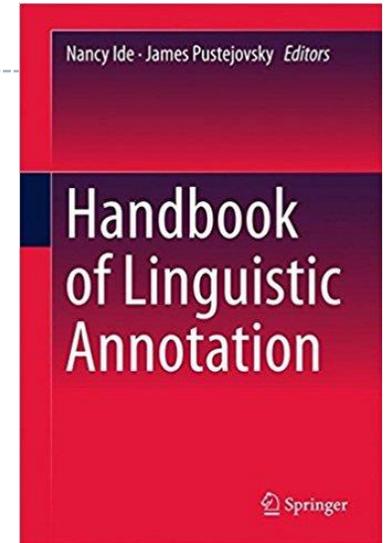
# What are linguists' roles in all this?

▶ **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.

# All about Linguistic Annotation

▶ *Handbook of Linguistic Annotation* (2017)

  ◆ Nancy Ide, James Pustejovsky (eds)

  ◆ https://link.springer.com/chapter/10.1007/978-94-024-0881-2_1

  ◆ Offers in-depth coverage on the topic of linguistic annotation

▶ *Natural Language Annotation for Machine Learning* (2012)

  ◆ James Pustejovsky, Amber Stubbs

  ◆ https://www.oreilly.com/library/view/natural-language-annotation/9781449332693/ch01.html

# Wrapping up

▶ To-do #9

 ◆ Explore APLS (Archive of Pittsburgh Language and Speech)

▶ Friday

 ◆ Guest presentation by Maya Asher, introduction to APLS project

▶ Your project

 ◆ Work on it! Focus on DATA.