## CSCI 384 Computational Linguistics
Fall 2023

Introduction and warm-up unit:

- ▶ General introduction to the course (**today**)
- ▶ Lab: Python and relevant libraries (Friday)

Today:

- ▶ Computational linguistics in the context of related fields
- ▶ Syllabus and course details
- ▶ Basic vocabulary of computational linguistics

## Tokens, types, and lemmas

> *I rose to saw off the still rose that I saw still grew near the still.*

16 tokens.

Types occurring once: to, off, that, grew, near.

Types occurring twice: I, rose, saw, the.

Type occurring three times: still

Distinct lexemes/lemmas: I, rise, to, saw (verb-cut), off, the, still (adjective), rose (noun), that, see, still (adverb), grow, near, still (noun).

Coming up:

- ► Learn Python
- ► Take Python quiz (Thurs, Aug 24)
- ► Do Python warm-up assignment (Mon, Aug 28)
- ► Read J&M, Sections 2.(0–4) (Mon, Aug 28)

Next time: Introduction to Python and NLTK *in the lab*.