# CMSC 473/673 Introduction to Natural Language Processing

Instructor: Lara J. Martin (she/they)

TA: Omkar Kulkarni (he)

https://laramartin.net/NLP-class/

# Learning Objectives

By the end of the course, you will be able to…

1. Recall common tasks in NLP and formulate problems for them. (HW1)

2. Diagnose and setup appropriate evaluation metrics for a given problem, including determining what an appropriate baseline might be. (HW2)

3. Compare and contrast language models and other NLP methods. (HW2, Exam)

Knowledge Checks

4. Implement AI systems that use popular NLP toolkits and libraries. (Grad Assignment, Project)

5. Construct a literature review from state-of-the-art research. (Grad Assignment, Project)

6. Plan and create an NLP system for a particular task. (HW3, Project)

7. Identify ethical issues in NLP systems and consider how they might be mitigated. (HW3)

# Grades

| Assignment | 473 (undergrad) | 673 (grad) |
|---|---|---|
| Class Knowledge Checks | 15% | 10% |
| Homework 1 | 10% | 5% |
| Homework 2 | 15% | 15% |
| Homework 3 | 15% | 15% |
| Exam | 15% | 15% |
| Project | 30% | 30% |
| Grad Assignment | - | 10% |

- In-class checks so that I can see how you're doing with the material
- Not graded for accuracy
- Can be made up by the end of the semester

- 3 homework assignments
- NLP tasks, evaluation & neural networks, prompt engineering & NLP ethics
- First homework is worth less than the other two
- Can be worked on alone or in pairs

# Grades

| Assignment | 473 (undergrad) | 673 (grad) |
|---|---|---|
| Class Knowledge Checks | 15% | 10% |
| Homework 1 | 10% | 5% |
| Homework 2 | 15% | 15% |
| Homework 3 | 15% | 15% |
| Exam | 15% | 15% |
| Project | 30% | 30% |
| Grad Assignment | - | 10% |

- New for this semester
- I want to test your knowledge of NLP concepts

- Group project (around 3-5 people)
- You will come up with your own topic with my help

- Implementation or literature review

The lecture schedule will be updated as the term progresses.

| Date | Lecture Topic | Readings for this Lesson | Homework Due |
|---|---|---|---|
| Tue, Jan 27, 2026 | No Class - Snow Day | | |
| Thu, Jan 29, 2026 | No Class - Snow Day | | |
| Fri, Jan 30, 2026 | Waitlist Deadline | | |
| Tue, Feb 3, 2026 | What is NLP?<br>[slides] | • Jacob Eisenstein, NLP Chapter 1 | |
| Thu, Feb 5, 2026 | Examples of NLP Tasks | • Dan Jurafsky and James H. Martin, SLP Chapter 2<br>• Jacob Eisenstein, NLP Chapter 2.2 & 4.5 | |
| Mon, Feb 9, 2026 | Waitlists Deactivated | | |
| Tue, Feb 10, 2026 | Examples of NLP Tasks | | |
| Thu, Feb 12, 2026 | Machine Learning Basics | • Ian Goodfellow, Yoshua Bengio, and Aaron Courville, Deep Learning Chapter 5.1-5.3 (ML Basics) | |
| Fri, Feb 13, 2026 | Last Day to Change Schedule | | |

# Academic Integrity

- If you feel the need to cheat on the assignment to do well on it, please talk to me or Omkar first. We can work it out ahead of time, but once you cheat it's hard to do anything.

If you cheat or plagiarize, you…
- aren't learning anything
- wasting money paying for tuition
- can get an F on the assignment or even for the entire class

More details on course website

# If you want to use ChatGPT

- Make sure you're saying that you used it

- Provide your prompt and the original generation (along with how you edited it)

- Make sure that <u>you're not avoiding the learning objectives by using it</u>

- If you do not say you're using it and I notice, that is an academic integrity violation

- It's okay to use grammar tools (e.g., spell check or Grammarly) or small-scale prediction (e.g., next word prediction, tab completion), provided that they don't change the **substance** of your work

# Learning Objectives

Develop a working vocabulary of terms in the field of NLP

Recognize NLP systems in your daily life
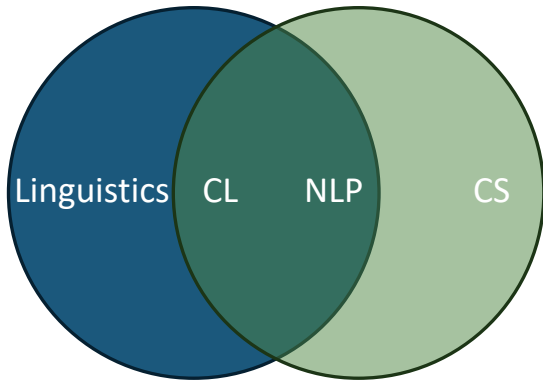
Define sub areas of linguistics

Distinguish between types and tokens

Define featurization & other ML terminology

Define some "classification" terminology

Distinguish between different text classification tasks

# Computational Linguistics
# =?
# Natural Language Processing

The computational **study** of language

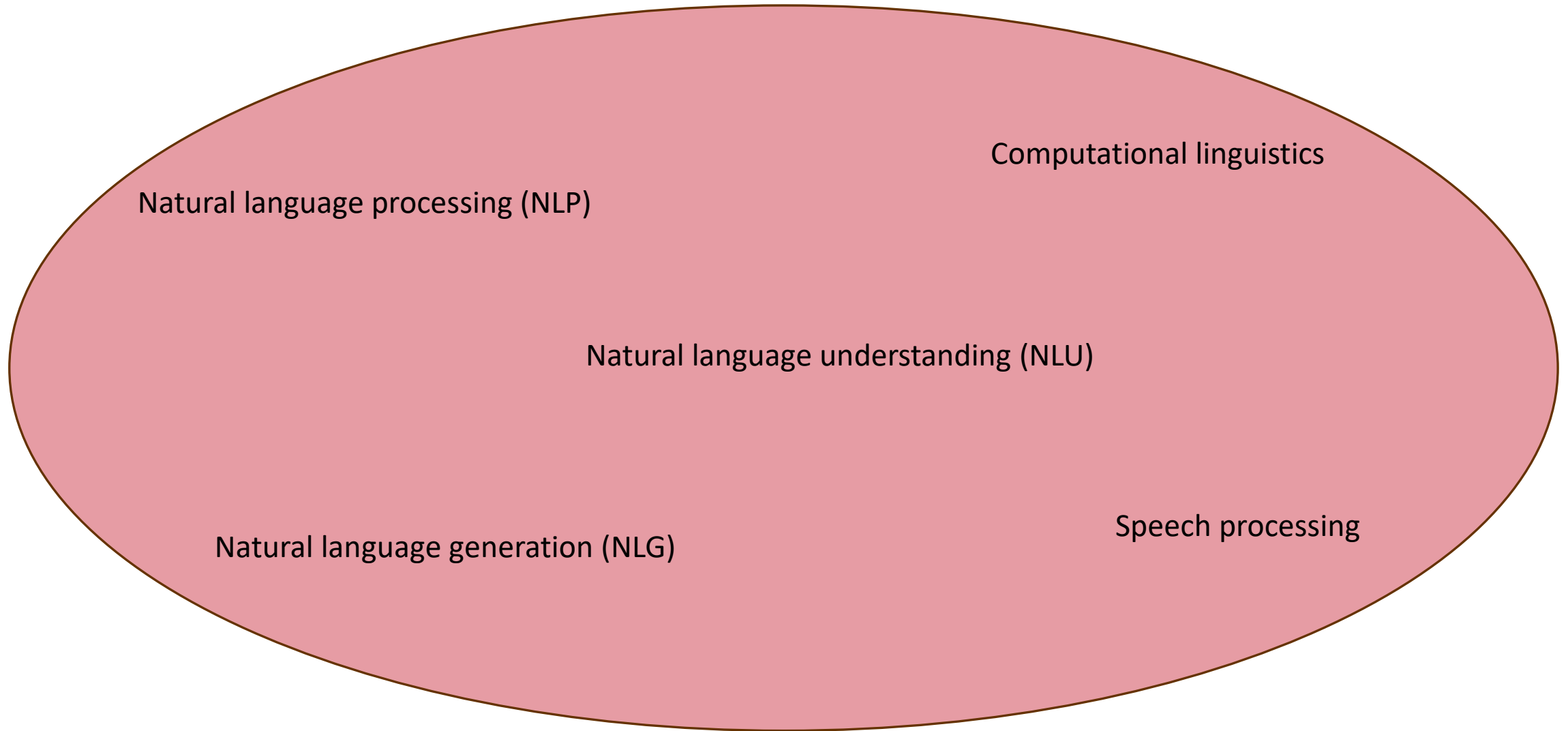# Computational Linguistics

# ≈

# Natural Language Processing

The computational **use** of language

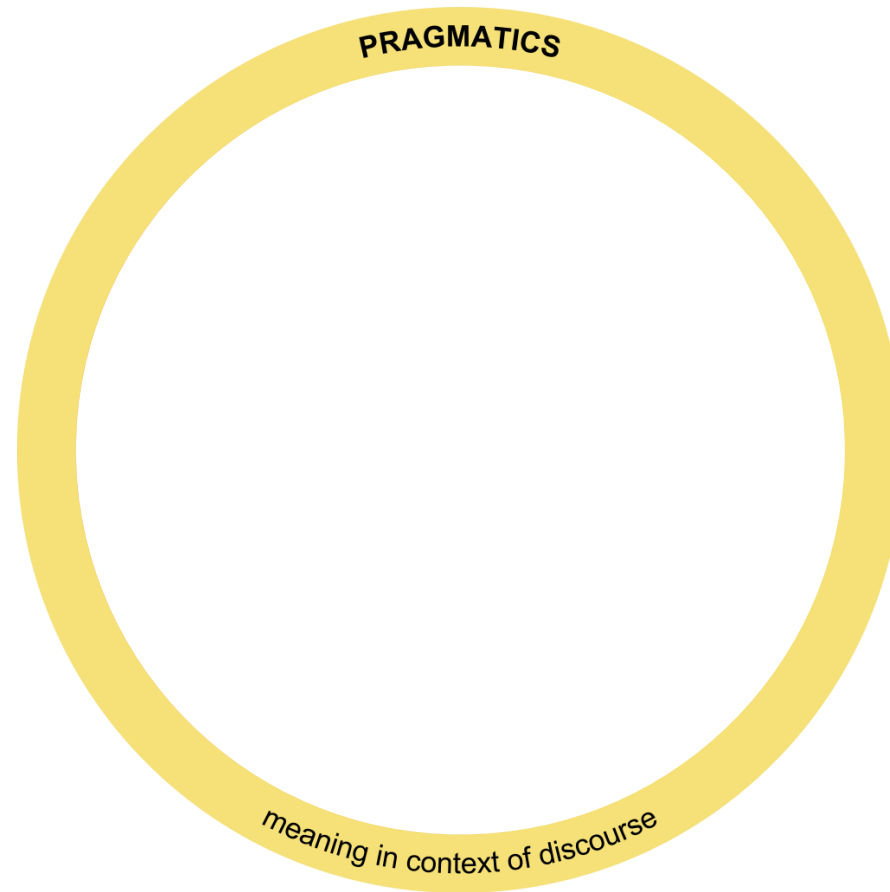Association for Computational Linguistics
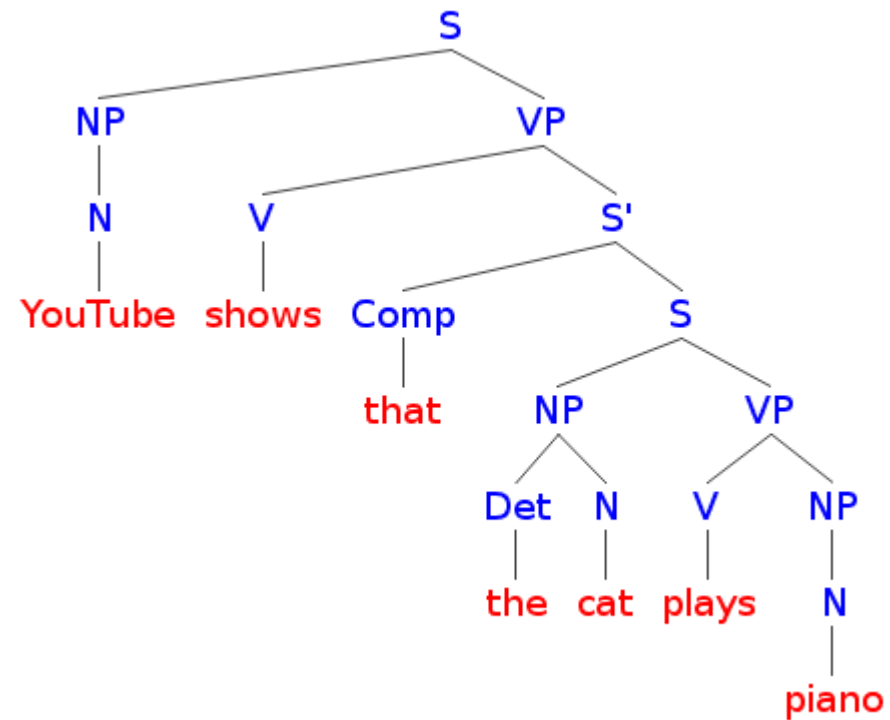
Language technologies

Computational linguistics

Natural language processing (NLP)

Natural language understanding (NLU)

Speech processing

Natural language generation (NLG)

# Linguistics

The study of language



PRAGMATICS

meaning in context of discourse

https://en.wikipedia.org/wiki/Morphology_(linguistics)#/media/File:Major_levels_of_linguistic_structure.svg

# Semantics

Meaning

$$S' = NP'(VP)$$
$$= \exists z[\text{mortal}(z) \wedge \text{loves}(z)(\text{Thetis})]$$

$$NP' = \text{Name}'$$
$$= \lambda P P(\text{Thetis})$$

$$VP' = \lambda x NP'(\lambda y V'(y)(x))$$
$$= \lambda x(\exists z[\text{mortal}(z) \wedge \text{loves}(z)(x)])$$

$$\text{Name}' = \lambda P P(\text{Thetis}) \quad V' = \text{loves}$$

$$NP' = \text{Det}'(N')$$
$$= \lambda Q(\exists z[\text{mortal}(z) \wedge Q(z)])$$

$$\text{Def}' = \lambda P \lambda Q(\exists z[P(z) \wedge Q(z)]) \quad N' = \text{mortal}$$

**Thetis**     **loves**     **a**     **mortal**

# Syntax

Grammar



https://allthingslinguistic.com/post/100617668093/how-to-draw-syntax-trees-part-3-type-1-a

# Phonology

Processing of sounds

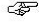

tsunami

↓

sunami

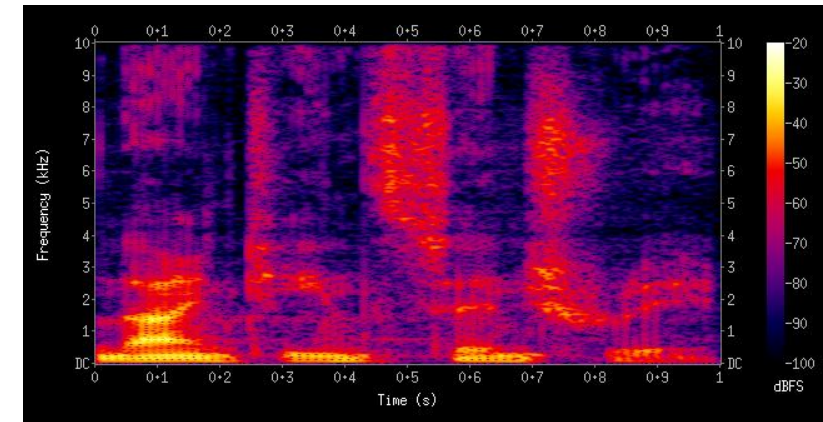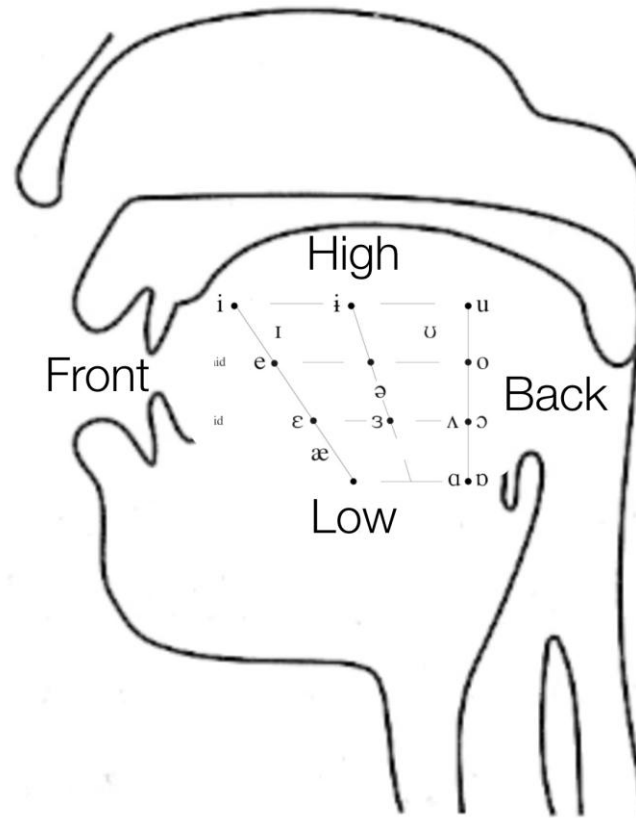https://upload.wikimedia.org/wikipedia/commons/a/a5/Tsunami_by_hokusai_19th_century.jpg

| /ðɪs/ *this* | | Dep | *Coda | Max |
|---|---|---|---|---|
| a. ☞ [dɪs] | | | * | |
| b. ☞ [dɪ] | | | | * |
| c. [dɪ.sə] | | *! | | |

https://pubs.asha.org/doi/10.1044/0161-1461%282001/022%29

# Phonetics

Physical production/understanding of sounds



https://wstyler.ucsd.edu/talks/l111_3_phonetics_review_handout.html



https://en.wikipedia.org/wiki/Spectrogram#/media/File:Spectrogram-19thC.png

# Back to CL vs NLP

Computational linguistics: Using computers to solve linguistic questions
- E.g., How does language X order their sentences? SVO, SOV, VOS...?
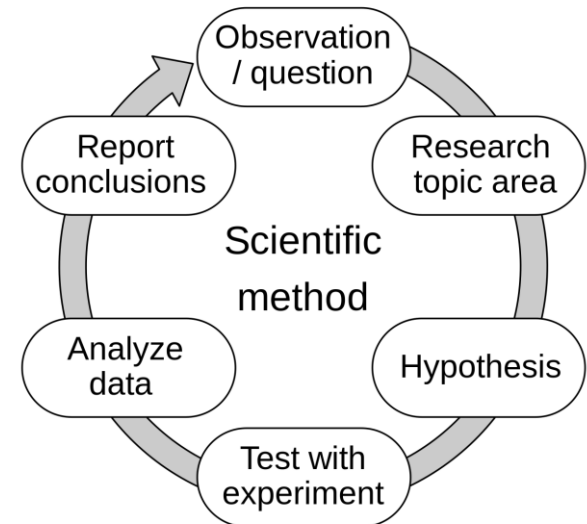
And this can inform NLP work
- E.g., How can we create a system that generates text in language X?

Or not...
- E.g., Let's feed a model a bunch of text so that it can generate text in language X.

How do we solve any of these problems?

Data!

# Where does the data come from?

Corpus (plural: corpora)
◦ Literally a "body" of text

Languages with few corpora are called "low-resource languages"
◦ This might not mean the language is endangered!

We can collect corpora in a few different ways:
◦ Curation: data tagged & organized by experts
◦ Internet: data "scraped" from open-access sources (Wikipedia, Reddit)
   ◦ Or data collected with permission from closed sources (Facebook, texts) – more rare
◦ Elicitation: carefully getting participants to produce language (lab studies, crowdsourcing, field studies)
◦ Pre-existing corpora

**!** Facebook has gotten into trouble several times for using data or manipulating people's feeds without their permission

# Benchmarking

Collecting & publishing corpora is helpful for…
- Replication
- Improving performance

# Benchmarking

Your task

If you want people to work on your problem, make it easy for them to get started and to measure their progress. Provide:

- Test data, for evaluating the final systems
- Development data, for measuring whether a change to the system helps, and for tuning parameters
- An evaluation metric (formula for measuring how well a system does on the dev or test data)
- A program for computing the evaluation metric
- Labeled training data and other data resources
- A prize? – with clear rules on what data can be used

# What does the data look like?

Curated data (and some collected data) are usually labeled, especially when made for a particular **task**

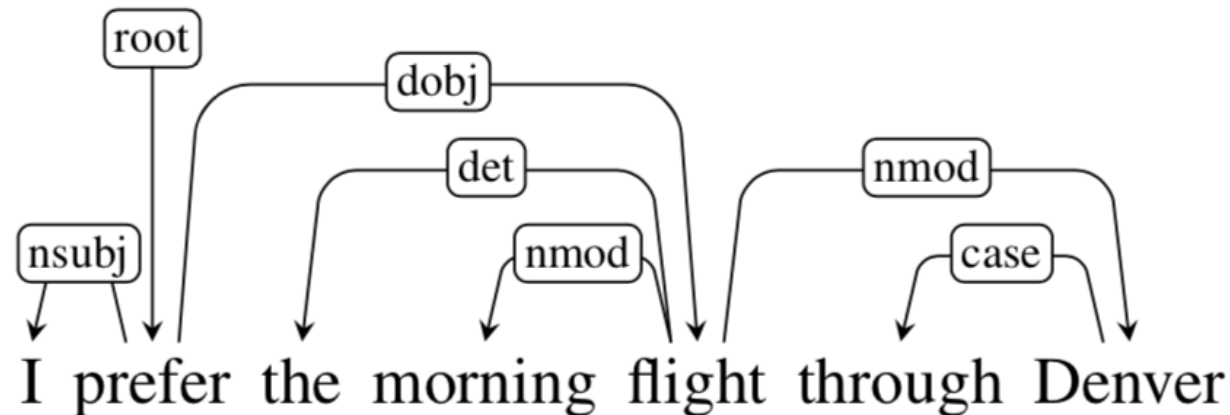◦ E.g., Universal dependencies (https://universaldependencies.org/)

# Current UD Languages

Information about language families (and genera for families with multiple branches) is mostly taken from WALS Online (IE = Indo-European).

| | | Name | | Size | | Family |
|---|---|---|---|---|---|---|
| ▶ | | Abaza | 1 | <1K | | Northwest Caucasian |
| ▶ | | Abkhaz | 1 | 6K | | Northwest Caucasian |
| ▶ | | Afrikaans | 1 | 49K | | IE, Germanic |
| ▶ | | Akkadian | 2 | 25K | | Afro-Asiatic, Semitic |
| ▶ | | Akuntsu | 1 | 1K | | Tupian, Tupari |
| ▶ | | Albanian | 2 | 4K | | IE, Albanian |
| ▶ | | Amharic | 1 | 10K | | Afro-Asiatic, Semitic |
| ▶ | | Ancient Greek | 3 | 456K | | IE, Greek |
| ▶ | | Ancient Hebrew | 1 | 39K | | Afro-Asiatic, Semitic |
| ▶ | | Apurina | 1 | <1K | | Arawakan |
| ▶ | | Arabic | 3 | 1,042K | | Afro-Asiatic, Semitic |
| ▶ | | Armenian | 2 | 94K | | IE, Armenian |
| ▶ | | Assyrian | 1 | <1K | | Afro-Asiatic, Semitic |
| ▶ | | Azerbaijani | 1 | <1K | | Turkic, Southwestern |
| ▶ | | Bambara | 1 | 13K | | Mande |
| ▶ | | Basque | 1 | 121K | | Basque |
| ▶ | | Bavarian | 1 | 15K | | IE, Germanic |
| ▶ | | Beja | 1 | 11K | | Afro-Asiatic, Cushitic |
| ▶ | | Belarusian | 1 | 305K | | IE, Slavic |
| ▶ | | Bengali | 1 | <1K | | IE, Indic |
| ▶ | | Bhojpuri | 1 | 6K | | IE, Indic |
| ▶ | | Bororo | 1 | 6K | | Bororoan |
| ▶ | | Breton | 1 | 10K | | IE, Celtic |
| ▶ | | Bulgarian | 1 | 156K | | IE, Slavic |
| ▶ | | Buryat | 1 | 10K | | Mongolic |
| ▶ | | Cantonese | 1 | 13K | | Sino-Tibetan, Chinese |
| ▶ | | Cappadocian | 2 | 4K | | IE, Greek |
| ▶ | | Catalan | 1 | 553K | | IE, Romance |
| ▶ | | Cebuano | 1 | 1K | | Austronesian, Central Philippine |
| ▶ | | Chinese | 7 | 309K | | Sino-Tibetan, Chinese |
| ▶ | | Chukchi | 1 | 6K | | Chukotko-Kamchatkan |
| ▶ | | Classical Armenian | 1 | 88K | | IE, Armenian |
| ▶ | | Classical Chinese | 2 | 433K | | Sino-Tibetan, Chinese |
| ▶ | | Coptic | 1 | 57K | | Afro-Asiatic, Egyptian |
| ▶ | | Croatian | 1 | 199K | | IE, Slavic |
| ▶ | | Czech | 6 | 2,252K | | IE, Slavic |
| ▶ | | Danish | 1 | 100K | | IE, Germanic |
| ▶ | | Dutch | 2 | 506K | | IE, Germanic |
| ▶ | | Egyptian | 1 | 14K | | Afro-Asiatic, Egyptian |
| ▶ | | English | 11 | 760K | | IE, Germanic |
| ▶ | | Erzya | 1 | 20K | | Uralic, Mordvin |

# What does the data look like?

Curated data (and some collected data) are usually labeled, especially when made for a particular **task**

◦ E.g., Universal dependencies (https://universaldependencies.org/)



https://medium.com/data-science-in-your-pocket/dependency-parsing-associated-algorithms-in-nlp-96d65dd95d3e
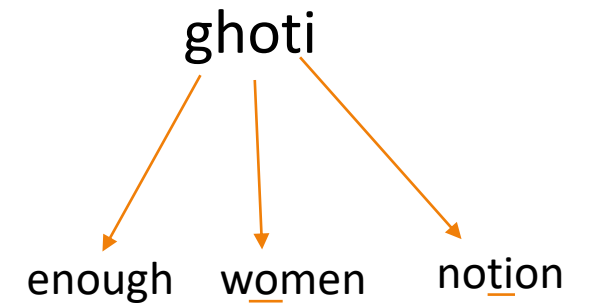
# Modalities

Text

Audio (speech)

Video (closed captioning, sign languages)

Pictures (handwriting recognition, image captioning)

Any of these can be labeled

TTS isn't straight forward. Unless you have information on how text is pronounced, an orthography (a writing system) by itself can be misleading.

ghoti

enough    women    notion

# What's in a word?

bat

https://www.freepngimg.com/download/bat/9-2-bat-png-hd.png
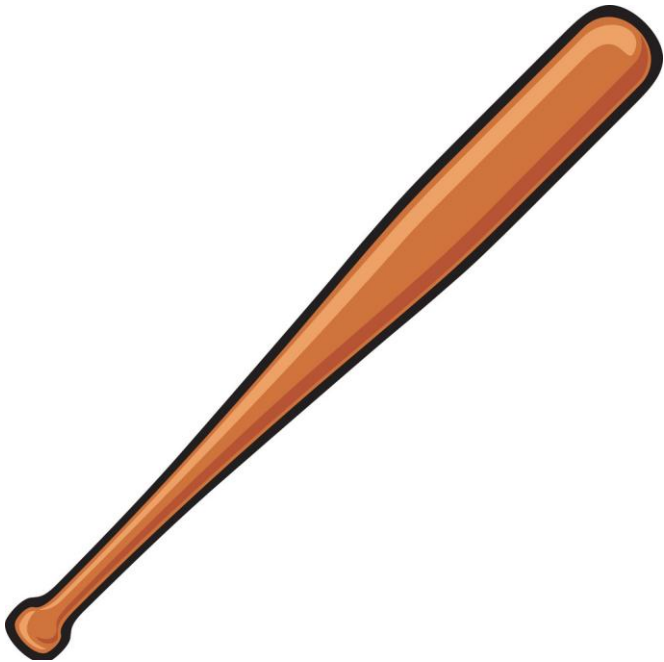
# What's in a word?

bats

# What's in a word?

Fledermaus
*flutter mouse*

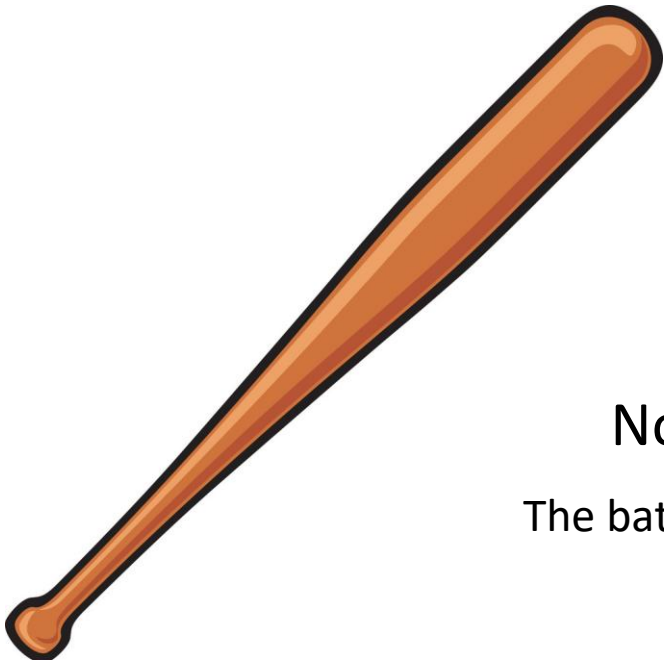# What's in a word?



bat

# What's in a word?

bat

Noun?

The bat was heavy.

Verb?

They bat 1000.

# What's in a word?

):

INTRO TO INTRO TO NLP

# What's in a word?

my leg is hurting nasty ):

# What's in a word?

add two cups (a pint): bring to a boil

# Tokens vs Types

The film got a great opening and the film went on to become a hit .

**Vocabulary:** the words (items) you know

**Type:** an element of the vocabulary.

**Token:** an instance of that type in running text.

How many of types & tokens appear in the above sentence?

# Tokens vs Types

**Types**
- The
- film
- got
- a
- great
- opening
- and
- the
- went
- on
- to
- become
- hit
- .

**Tokens**
- The
- film
- got
- a
- great
- opening
- and
- the
- ~~film~~
- went
- on
- to
- become
- ~~a~~
- hit
- .

# For your {task} how do you define tokens?

Sometimes:

1. They're defined for you by the *dataset creator*

What usually happens when you input a word that your writing/texting program doesn't recognize?
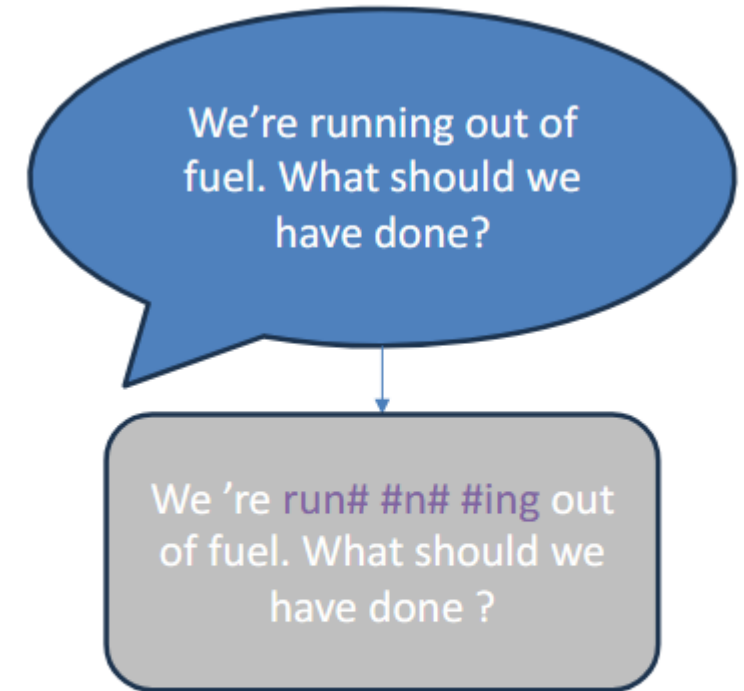
We're running out of fuel. What should we have done?

We 're run# #n# #ing out of fuel. What should we have done ?

*why?*
- *scaleably handling novel words*
  - *linguistic reasons*
- *historical reasons / technical debt*

# For your {task} how do you define tokens?

Sometimes:

1. They're defined for you by the *dataset creator*

2. They're defined by the *model*

We're running out of fuel. What should we have done?

We 're run# #n# #ing out of fuel. What should we have done ?
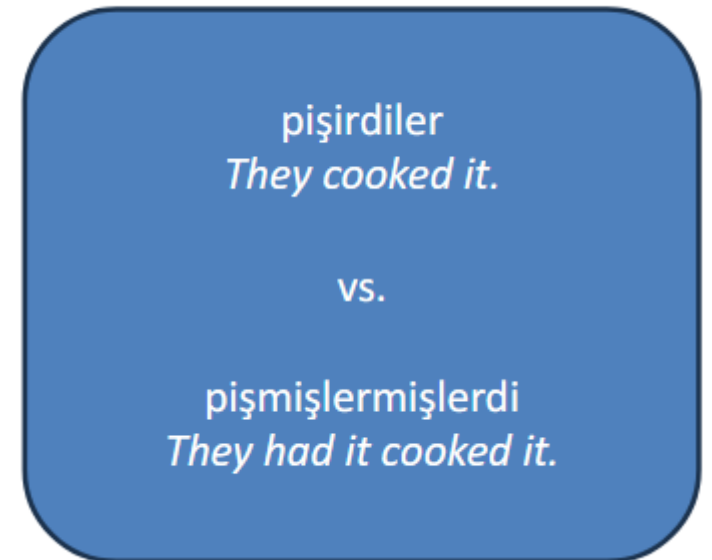
*(why? scaleably handling novel words)*
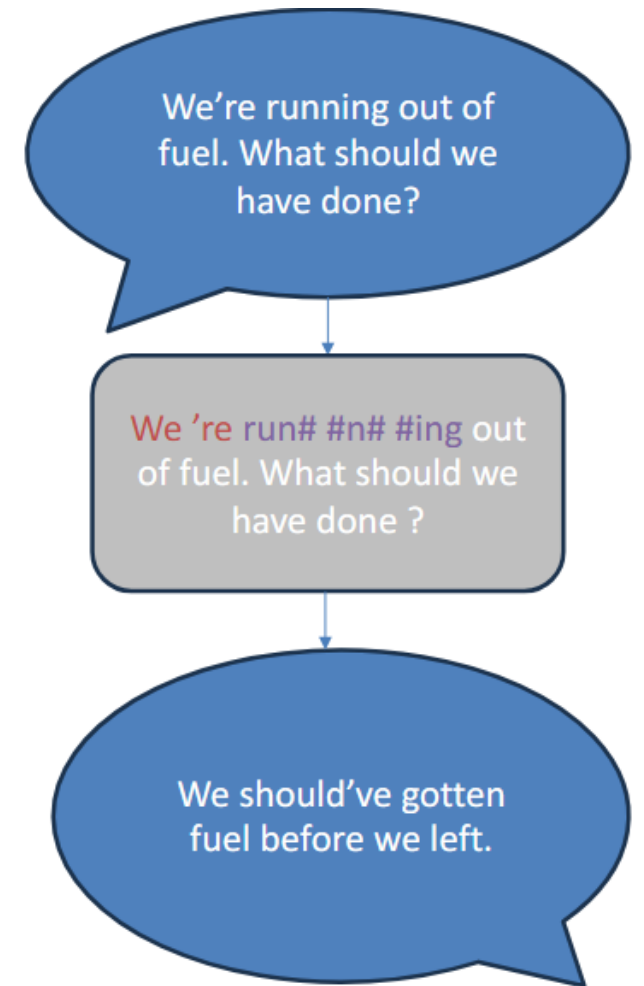
# For your {task} how do you define tokens?

Sometimes:

1. They're defined for you by the *dataset creator*

2. They're defined by the *model*

3. It might be part of the *research problem itself*

pişirdiler
*They cooked it.*

vs.

pişmişlermişlerdi
*They had it cooked it.*

# For your {task} how do you define tokens?

Sometimes:

1. They're defined for you by the *dataset creator*

2. They're defined by the *model*

3. It might be part of the *research problem itself*

4. They're defined by the *end user*
   1. You'll need to handle points 1 and/or 2 on-the-backend…
   2. and then reversing the process to present output to the user

# What are some NLP applications that you see in your daily life?

- ChatGPT/chatbots

- Speech recognition

- Dialog agents (Siri/Alexa)

- Translation

- Auto-correct/ Grammar correction

- Auto-complete (search engines, email)

- Search engine agent

- Code assistants

- Email summarizer

# True or False

The following sentence has the same number of types as tokens (i.e., # types = # tokens)
The dog caught the frisbee .