# Natural Language Processing
## Introduction to NLP

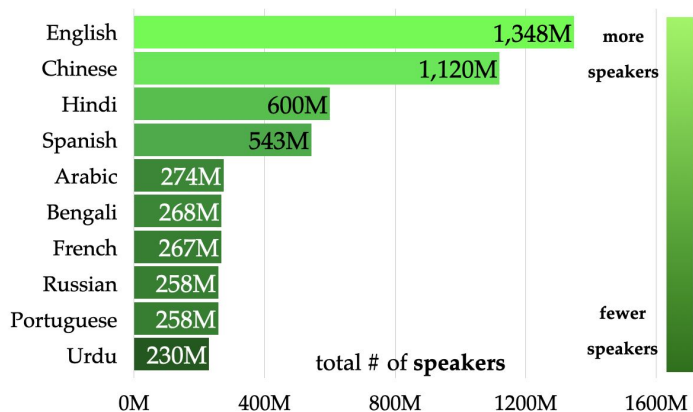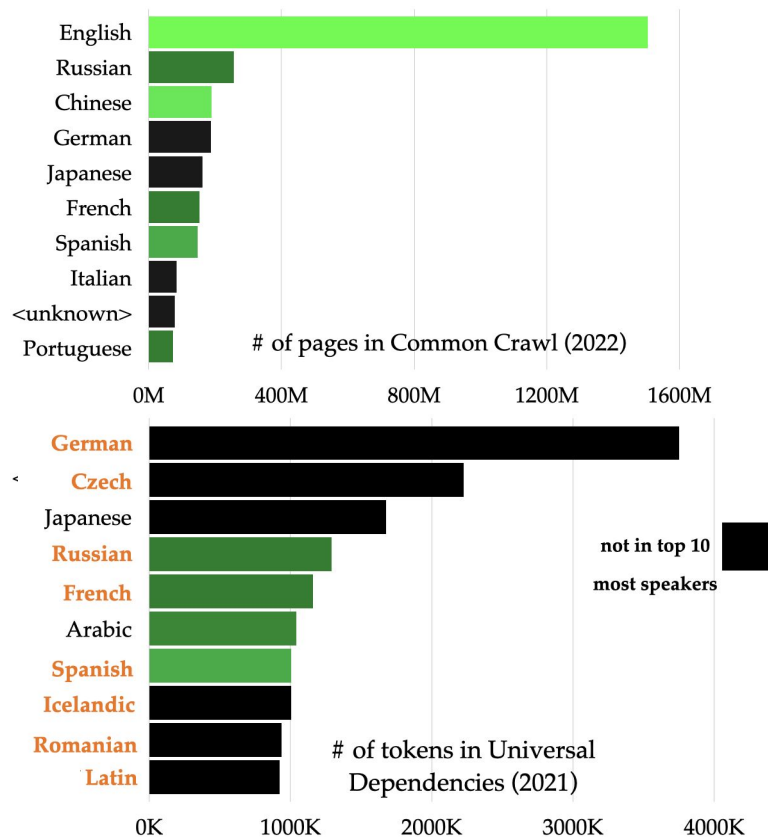**Sofia Serrano**
**sofias6@cs.washington.edu**

# Announcements

- Academic Integrity Form is out on Canvas
- A1 is out on GitLab
  - Don't see it? Reply to this thread on Ed with your NetID: https://edstem.org/us/courses/32306/discussion/2365366
- Access to lecture recordings
  - No @cs.washington.edu google account? Click through to (request) access any lecture recording sooner rather than later so that we can give you access
- Make sure you can access the course machines
  - *(if connecting from off campus)* Run Husky OnNet VPN OR first ssh into an attu machine
  - `ssh yourNetID@nlpg00.cs.washington.edu`          *(nlpg00-nlpg03)*
  - Not working?
    - Not a CSE major/no CSE account? Email ugrad-adviser@cs.washington.edu to request a CSE account (include your student ID number in the email) and CC Sofia
    - Still not working? Reply to this thread on Ed so that we can help troubleshoot: https://edstem.org/us/courses/32306/discussion/2368995

# Following up on a question from last lecture



total # of **speakers**

| Language | speakers |
|---|---|
| English | 1,348M |
| Chinese | 1,120M |
| Hindi | 600M |
| Spanish | 543M |
| Arabic | 274M |
| Bengali | 268M |
| French | 267M |
| Russian | 258M |
| Portuguese | 258M |
| Urdu | 230M |

more speakers / fewer speakers

# of pages in Common Crawl (2022)

# of tokens in Universal Dependencies (2021)

not in top 10 most speakers

Credit to Phoebe Mulcaire for figures

3

# Symbolic and Probabilistic NLP

**Logic-based/Rule-based NLP**



interlingua

analysis

transfer

generation

direct translation

source text

target text

**~ 90s**

**Statistical NLP**



**Translation Model**

| source phrase | target phrase | translation features |

e    f

**Language Model**

f

**Reranking Model**

feature weights

$argmax_e P(f|e)P(e)$

e    f

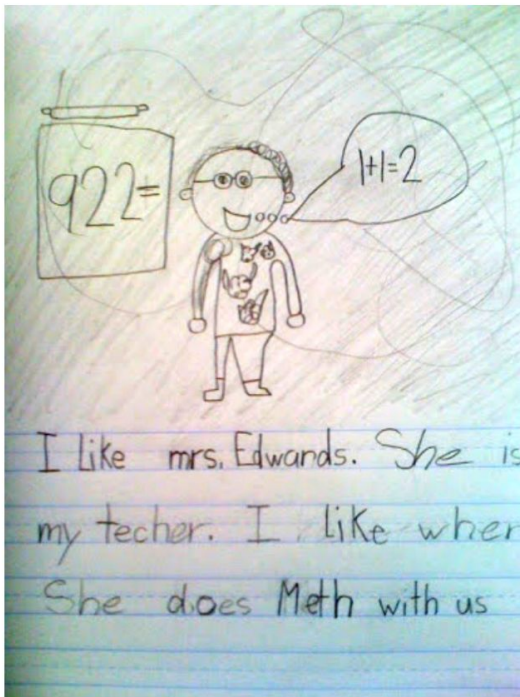# Probabilistic and Connectionist NLP

# Linguistic Background

# What does it mean to "know" a language?



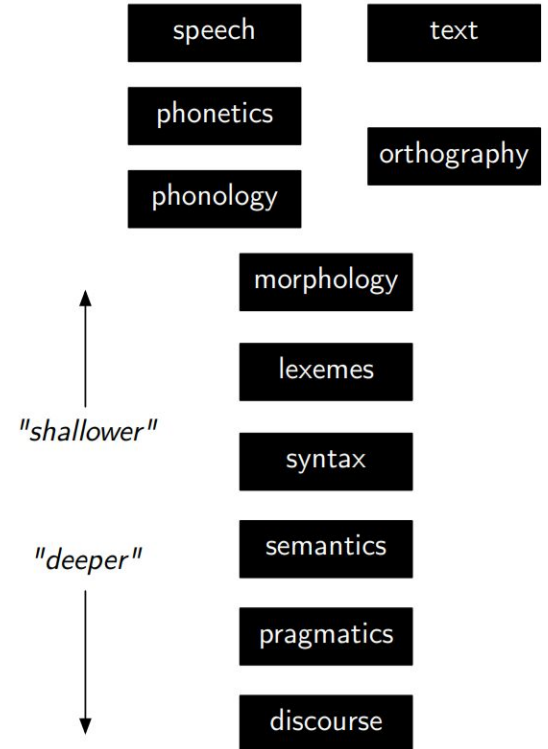(Thanks Canadian Internet Registration Authority!)

What do we need to "tell" a computer program so that it knows more English than `wc` or a dictionary, maybe even as much as a three-year-old, for example?

# What does an NLP system need to 'know'?

- Language consists of many levels of structure

- Humans fluently integrate all of these in producing/understanding language

- Ideally, so would a computer!

# Levels of linguistic knowledge

speech

text

phonetics

orthography

phonology

morphology

lexemes

*"shallower"*

syntax

semantics

*"deeper"*

pragmatics

discourse

# Speech, phonetics, phonology

This is a simple sentence .

/ ðɪs ɪz ə ˈsɪmpl ˈsɛntəns /.

speech

text

phonetics

orthography

phonology

morphology

lexemes

*"shallower"*

syntax

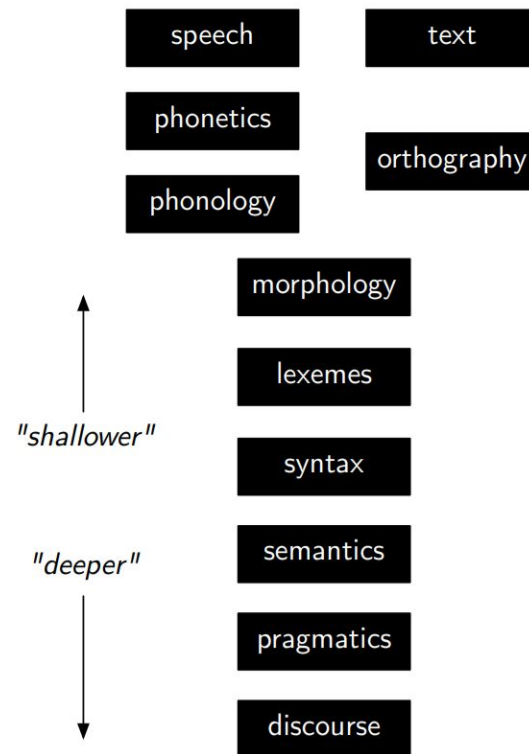semantics

*"deeper"*

pragmatics

discourse

# Orthography

هذه جملة بسيطة

đây là một câu đơn giản

यह एक साधारण वाक्य है

This   is   a   simple   sentence   .
/ ðɪs ɪz ə ˈsɪmpl ˈsɛntəns /.

speech

text

phonetics

orthography

phonology

morphology

lexemes

*"shallower"*

syntax

semantics

*"deeper"*

pragmatics

discourse

# Words, morphology

- Morphological analysis
- Tokenization
- Lemmatization

| | speech | | text |
| | phonetics | | |
| | | | orthography |
| | phonology | | |

morphology

lexemes

"shallower"

syntax

"deeper"

semantics

pragmatics

discourse

**Tokens** This is a simple sentence .

**Morphology**
be
3sg
present

# Syntax

- Part-of-speech tagging

| speech | | text |
| --- | --- | --- |

phonetics

orthography

phonology

morphology

lexemes

"shallower"

syntax

semantics

"deeper"

pragmatics

discourse

| **Parts of speech** | DT | VBZ | DT | JJ | NN | PUNC |
| --- | --- | --- | --- | --- | --- | --- |
| **Tokens** | This | is | a | simple | sentence | . |
| **Morphology** | | be<br>3sg<br>present | | | | |

# Syntax

- Part-of-speech tagging
- Syntactic parsing

# Semantics

- Named entity recognition
- Word sense disambiguation
- Semantic role labeling

# Discourse

- Reference resolution
- Discourse parsing

**Syntax**

S
VP
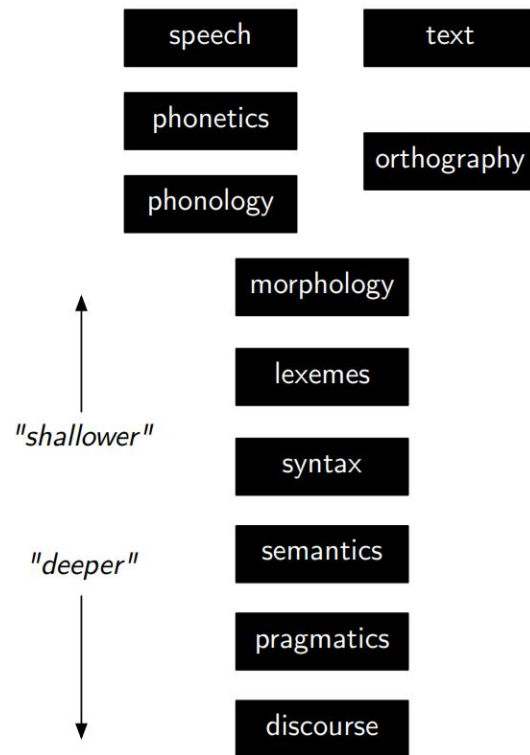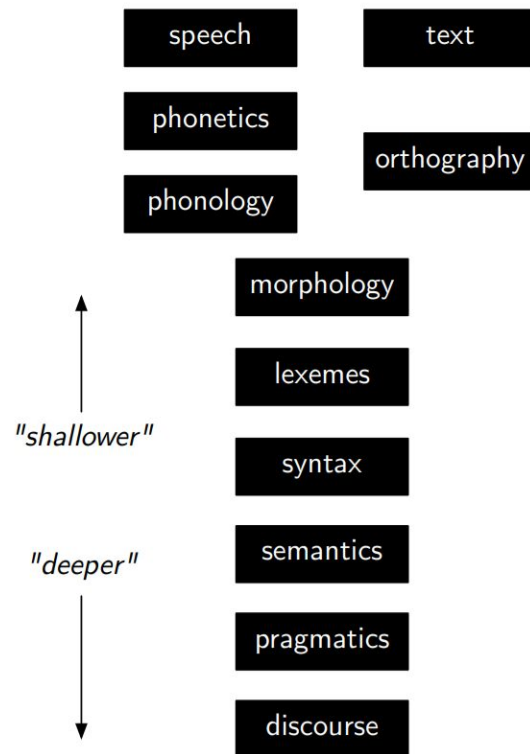NP NP

**Parts of speech**
DT VBZ DT JJ NN PUNC

**Tokens**
This is a simple sentence .

**Morphology**
be
3sg
present

SIMPLE1:
having few
parts

SENTENCE1:
String of words
satisfying the
grammatical rules
of a language

coreferent

**Semantics**

**Discourse**
But an instructive one .

speech | text

phonetics

orthography

phonology

morphology

lexemes

"shallower"

syntax

semantics

"deeper"

pragmatics

discourse

# Linguistic challenges we'll need to deal with in designing NLP systems

# What are some challenges for NLP systems?

1. Ambiguity
2. Variation
3. Sparsity
4. Expressivity
5. Unmodeled variables
6. Unknown representation $\mathcal{R}$

# Ambiguity

- Ambiguity at multiple levels:
  - Word senses: **bank** (finance or river?)
  - Part of speech: **chair** (noun or verb?)
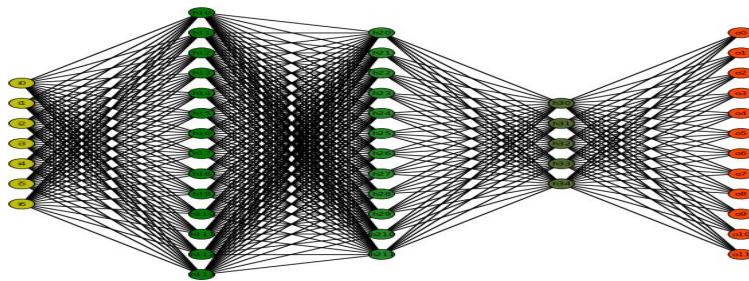  - Syntactic structure: **I can see a man with a telescope**
  - Multiple: **I saw her duck**

# Dealing with ambiguity

- How can we model ambiguity and choose the correct analysis in context?
    - non-probabilistic methods (FSMs for morphology, CKY parsers for syntax) return *all possible analyses.*
    - probabilistic models (HMMs for part-of-speech tagging, PCFGs for syntax) and algorithms (Viterbi, probabilistic CKY) return *the best possible analysis*, i.e., the most probable one according to the model
    - Neural networks, pretrained language models now provide end-to-end solutions



- But the "best" analysis is only good if our probabilities are accurate. Where do they come from?

# Corpora

- A corpus is a collection of text
  - Often annotated in some way
  - Sometimes just lots of text
- Examples
  - Penn Treebank: 1M words of parsed WSJ
  - Canadian Hansards: 10M+ words of aligned French / English sentences
  - Yelp reviews
  - The Web: billions of words of who knows what

# What are some challenges for NLP systems?

1. Ambiguity
2. Variation
3. Sparsity
4. Expressivity
5. Unmodeled variables
6. Unknown representation $\mathscr{R}$

# Variation in languages

- ~7K languages
- Thousands of language varieties



60M Speakers
125M Speakers
251M Speakers
90M Speakers
79M Speakers

Englishes



- Afro Asiatic
- Nilo Saharan
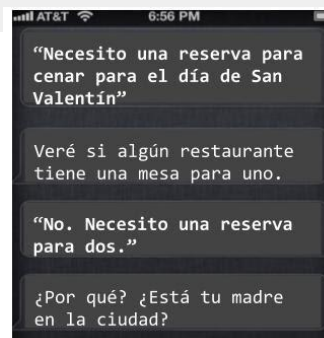- Niger Congo A
- Niger Congo B (Bantu)
- Khoisan
- Austronesian

Africa is a continent with a very high linguistic diversity:
there are an estimated 1.5-2K African languages from 6 language
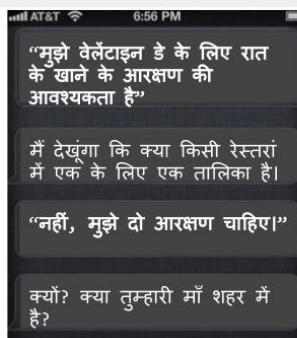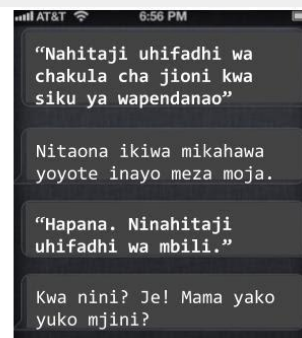families. 1.33 billion people

# NLP beyond English
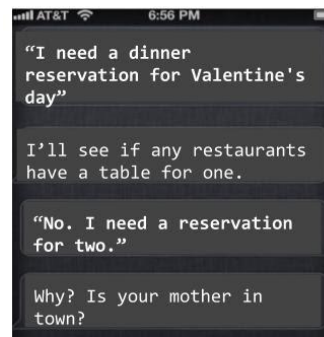
- ~7,000 languages
- thousands of language varieties



"Necesito una reserva para cenar para el día de San Valentín"

Veré si algún restaurante tiene una mesa para uno.

"No. Necesito una reserva para dos."

¿Por qué? ¿Está tu madre en la ciudad?

Spanish
534 million speakers

"मुझे वेलेंटाइन डे के लिए रात के खाने के आरक्षण की आवश्यकता है"

मैं देखूंगा कि क्या किसी रेस्तरां में एक के लिए एक तालिका है।

"नहीं, मुझे दो आरक्षण चाहिए।"

क्यों? क्या तुम्हारी माँ शहर में है?

Hindi
615 million speakers

"Nahitaji uhifadhi wa chakula cha jioni kwa siku ya wapendanao"

Nitaona ikiwa mikahawa yoyote inayo meza moja.

"Hapana. Ninahitaji uhifadhi wa mbili."

Kwa nini? Je! Mama yako yuko mjini?

Swahili
100 million speakers

"I need a dinner reservation for Valentine's day"

I'll see if any restaurants have a table for one.

"No. I need a reservation for two."

Why? Is your mother in town?

American English

"Ah need a tatties an' neebs reservation fur Valentine's day ."

I'll see if onie restaurants hae a table fur a body.

"Nae. Ah need a reservation fur tois."

Wa? is yer maw in toon?

Scottish English

"Mujhe Valentine's day par reservation chahiye."

I'll see agar ek aadmi ke liye table hai.

"Nhi. Mujhe do logo ke liye table chahiye."

Kyu? Aapki mother town me hain?

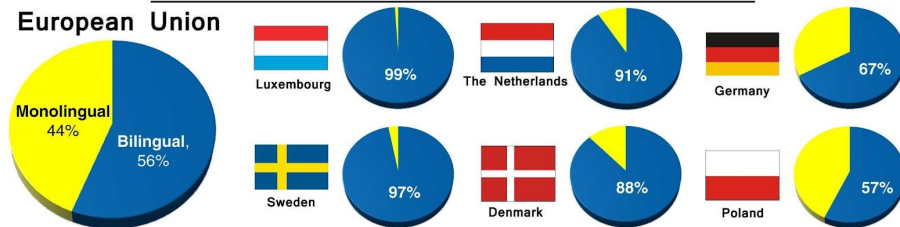Hinglish

# Most of the world today is multilingual



## Percentage of Bilingual Speakers in the World

European Union
- Monolingual 44%
- Bilingual, 56%

Luxembourg 99%
The Netherlands 91%
Germany 67%
Sweden 97%
Denmark 88%
Poland 57%

Source: European Commission, "Europeans and their Languages," 2006

Percentage of US Population who spoke a language other than English at home by year

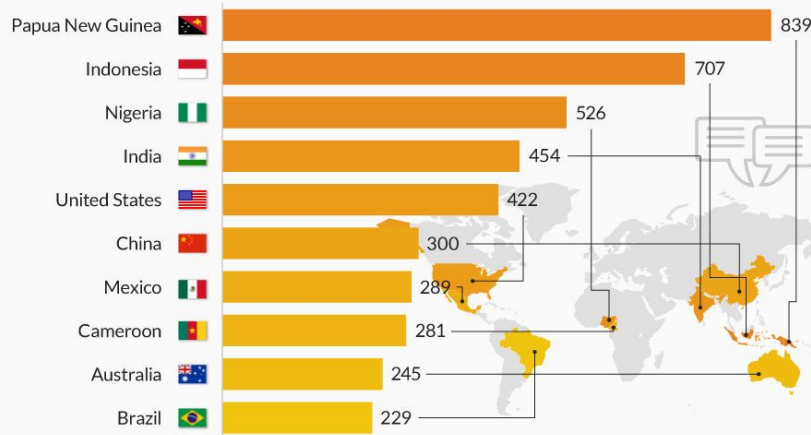- 1980: 10.97
- 1990: 13.82
- 2000: 17.89
- 2007: 19.73

Source: U.S. Census Bureau, 2007 American Community Survey

Source: US Census Bureau

## The Countries With The Most Spoken Languages
Number of living languages spoken per country in 2015

- Papua New Guinea: 839
- Indonesia: 707
- Nigeria: 526
- India: 454
- United States: 422
- China: 300
- Mexico: 289
- Cameroon: 281
- Australia: 245
- Brazil: 229

Source: Ethnologue

# Semantic analysis

- **Every language represents the world in a different way**
    - For example, it could depend on cultural or historical conditions



- Russian has very few words for colors, Japanese has hundreds
- Multiword expressions, e.g. happy as a clam, it's raining cats and dogs or wake up and metaphors, e.g. love is a journey are very different across languages

# Tokenization

这是一个简单的句子

**WORDS**   This   is   a   simple   sentence

זה   משפט   פשוט

# Tokenization + disambiguation

in tea        בתה
her daughter

- most of the vowels unspecified

| | |
|---|---|
| in tea | בתה |
| in the tea | בהתה |
| that in tea | שבתה |
| that in the tea | שבהתה |
| and that in the tea | ושבההתה |

ושבתה

| | |
|---|---|
| and her saturday | ו+שבת+ה |
| and that in tea | ו+ש+ב+תה |
| and that her daughter | ו+ש+בת+ה |

- most of the vowels unspecified
- particles, prepositions, the definite article, conjunctions attach to the words which follow them
- tokenization is highly ambiguous

# Tokenization + morphological analysis

- Quechua

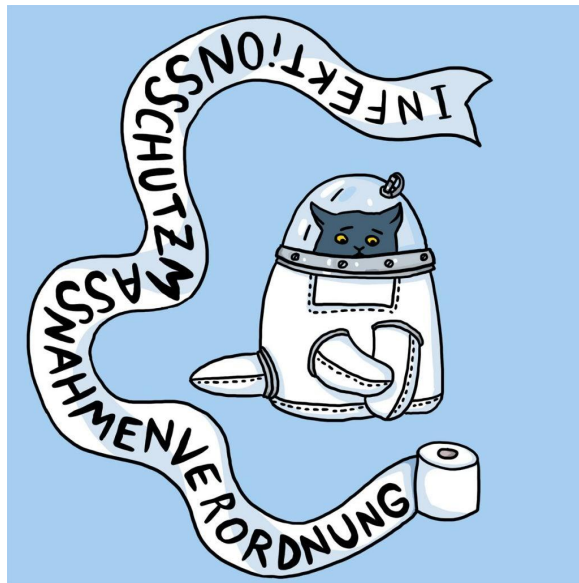Much'ananayakapushasqakupuniñataqsunamá

Much'a -na -naya -ka -pu -sha -sqa -ku -puni -ña -taq -suna -má

*"So they really always have been kissing each other then"*

```
Much'a  to kiss
-na     expresses obligation, lost in translation
-naya   expresses desire
-ka     diminutive
-pu     reflexive (kiss *eachother*)
-sha    progressive (kiss*ing*)
-sqa    declaring something the speaker has not personally witnessed
-ku     3rd person plural (they kiss)
-puni   definitive (really*)
-ña     always
-taq    statement of contrast (...then)
-suna   expressing uncertainty (So...)
-má     expressing that the speaker is surprised
```
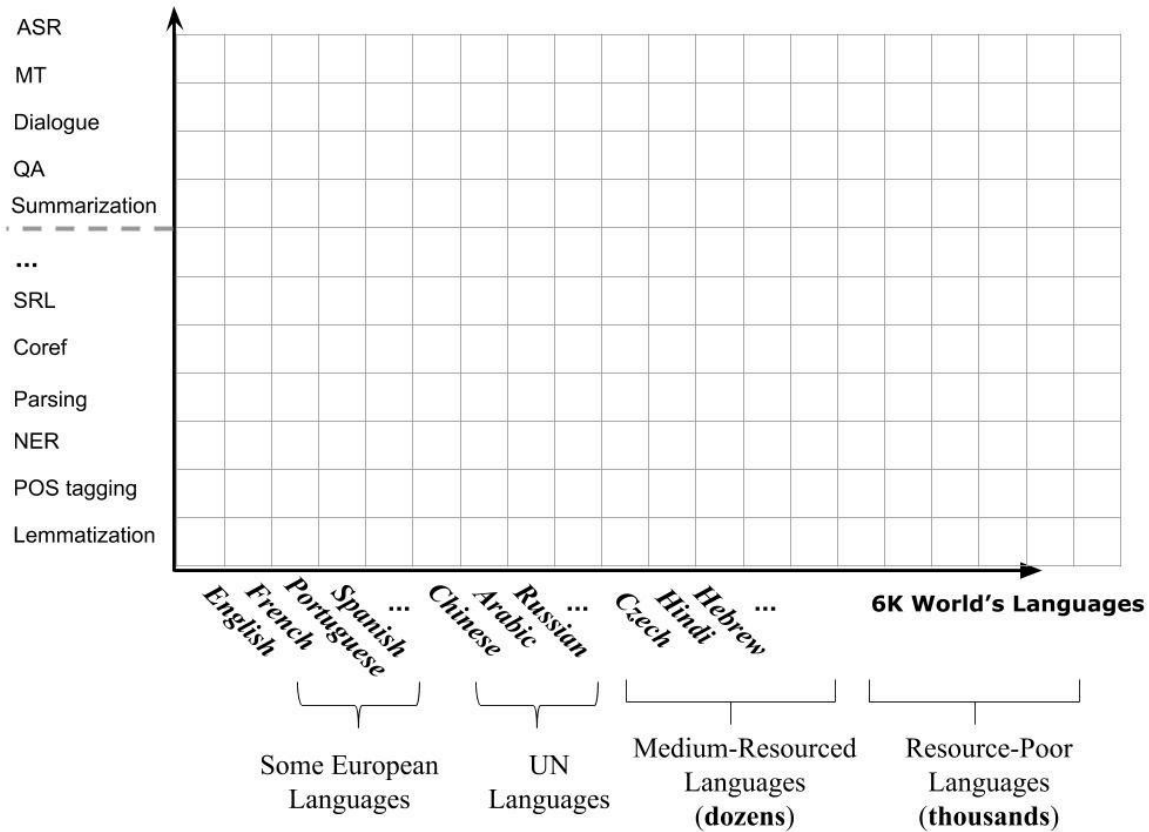
# Tokenization + morphological analysis

- German



Infektionsschutzmaßnahmenverordnung

**NLP Technologies/Applications**

ASR
MT
Dialogue
QA
Summarization
...
SRL
Coref
Parsing
NER
POS tagging
Lemmatization

English, French, Portuguese, Spanish, ..., Chinese, Arabic, Russian, ..., Czech, Hindi, Hebrew, ...

**6K World's Languages**

Some European Languages

UN Languages

Medium-Resourced Languages (**dozens**)

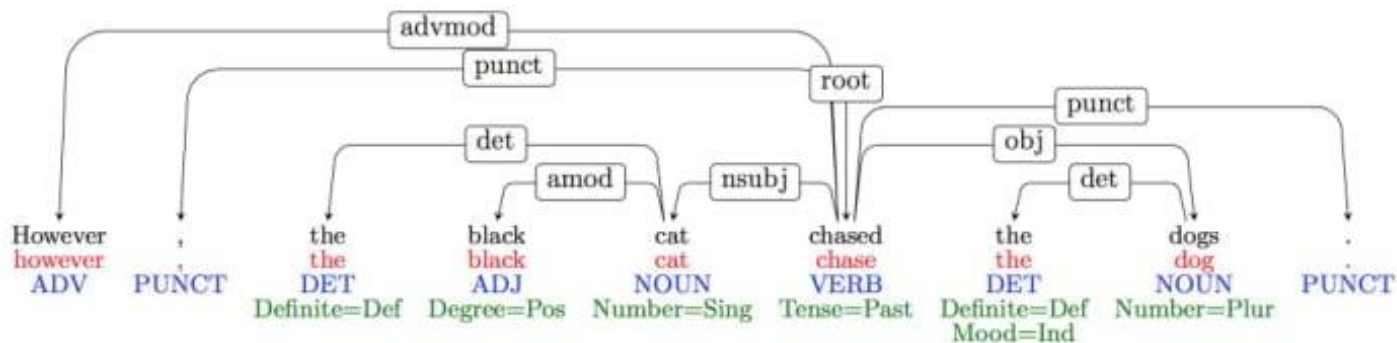Resource-Poor Languages (**thousands**)

# Linguistic variation

- Non-standard language, emojis, hashtags, names



chowdownwithchan #crab and #pork #xiaolongbao at @dintaifungusa… where else? 😂 🤷🏻‍♀️ Note the cute little crab indicator in the 2nd pic 🦀 💕
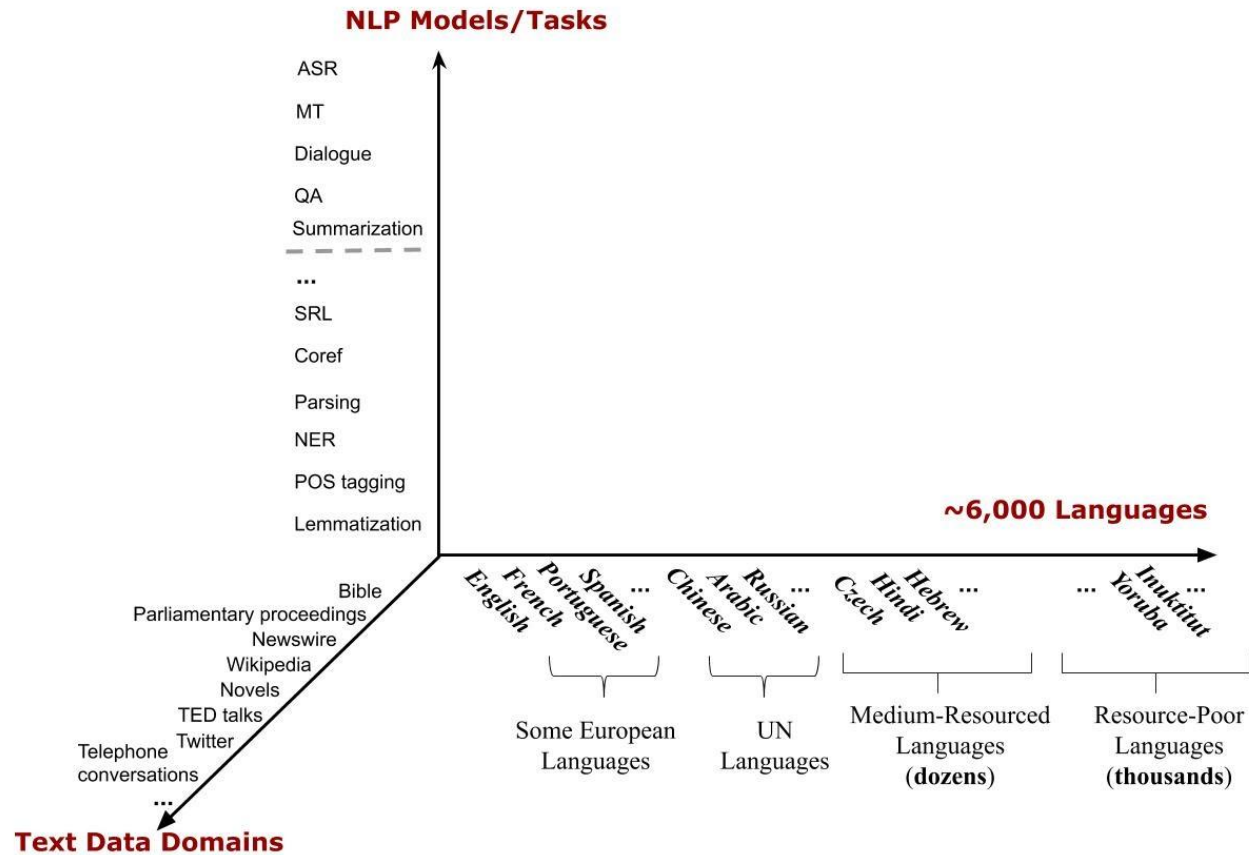
# Variation

- Suppose we train a part of speech tagger or a parser on the Wall Street Journal



- What will happen if we try to use this tagger/parser for social media??

@_rkpntrnte hindi ko alam babe eh, absent ako
kanina I'm sick rn hahaha😌🙌

**NLP Models/Tasks**

ASR
MT
Dialogue
QA
Summarization
...
SRL
Coref
Parsing
NER
POS tagging
Lemmatization

**~6,000 Languages**

English
French
Portuguese
Spanish
...
Chinese
Arabic
Russian
...
Czech
Hindi
Hebrew
...
...
Yoruba
Inuktitut
...

Some European Languages

UN Languages

Medium-Resourced Languages (**dozens**)

Resource-Poor Languages (**thousands**)

Bible
Parliamentary proceedings
Newswire
Wikipedia
Novels
TED talks
Twitter
Telephone conversations
...

**Text Data Domains**

# What are some challenges for NLP systems?

1. Ambiguity
2. Variation
3. Sparsity
4. Expressivity
5. Unmodeled variables
6. Unknown representation $\mathcal{R}$

# Sparsity

Sparse data due to Zipf's Law

- To illustrate, let's look at the frequencies of different words in a large text corpus
- Assume "word" is a string of letters separated by spaces

# Word Counts

Most frequent words in the English Europarl corpus (out of 24m word tokens)

| any word | | nouns | |
|---|---|---|---|
| Frequency | Token | Frequency | Token |
| 1,698,599 | the | 124,598 | European |
| 849,256 | of | 104,325 | Mr |
| 793,731 | to | 92,195 | Commission |
| 640,257 | and | 66,781 | President |
| 508,560 | in | 62,867 | Parliament |
| 407,638 | that | 57,804 | Union |
| 400,467 | is | 53,683 | report |
| 394,778 | a | 53,547 | Council |
| 263,040 | I | 45,842 | States |

# Word Counts

But also, out of 93,638 distinct words (word types), 36,231 occur only once.

Examples:

- cornflakes, mathematicians, fuzziness, jumbling
- pseudo-rapporteur, lobby-ridden, perfunctorily,
- Lycketoft, UNCITRAL, H-0695
- policyfor, Commissioneris, 145.95, 27a

# Plotting word frequencies

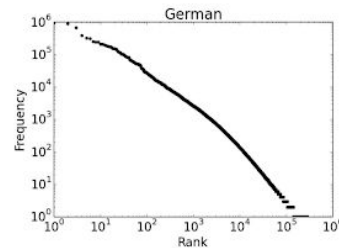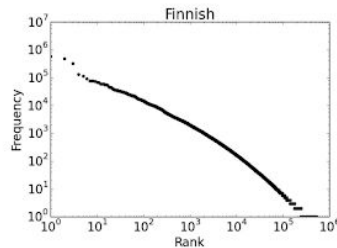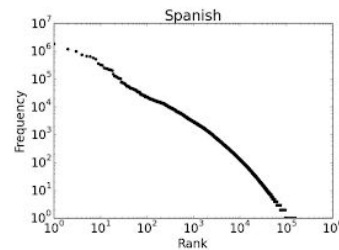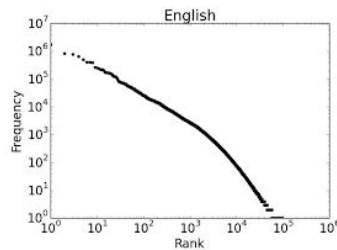Order words by frequency. What is the frequency of nth ranked word?

# Zipf's Law

Implications

- Regardless of how large our corpus is, there will be a lot of infrequent (and zero-frequency!) words
- This means we need to find clever ways to estimate probabilities for things we have rarely or never seen

# What are some challenges for NLP systems?

1. Ambiguity
2. Variation
3. Sparsity
4. **Expressivity**
5. Unmodeled variables
6. Unknown representation $\mathcal{R}$

# Expressivity

Not only can one form have different meanings (ambiguity) but the same meaning can be expressed with different forms:

She gave the book to Tom      vs.      She gave Tom the book

Some kids popped by      vs.      A few children visited

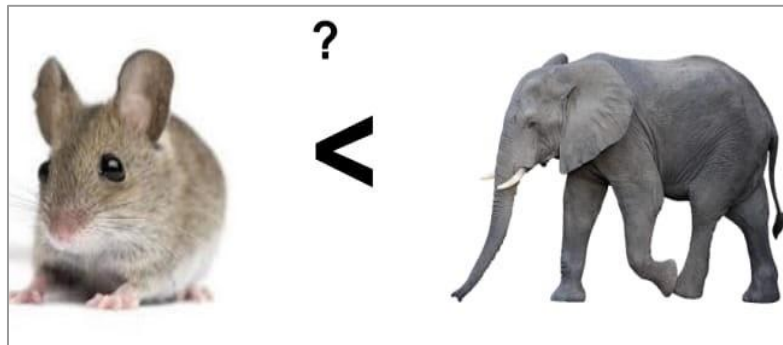Is that window still open?      vs.      Please close the window

# What are some challenges for NLP systems?

1. Ambiguity
2. Variation
3. Sparsity
4. Expressivity
5. Unmodeled variables
6. Unknown representation $\mathcal{R}$

# Unmodeled variables



"Drink this milk"



World knowledge
- I dropped the glass on the floor and it broke
- I dropped the hammer on the glass and it broke

# What are some challenges for NLP systems?

1. Ambiguity
2. Variation
3. Sparsity
4. Expressivity
5. Unmodeled variables

6. Unknown representation $\mathcal{R}$

# Unknown representation

- Very difficult to decide on a representation $\mathcal{R}$, since we don't even know how to represent the knowledge a human has/needs:
  - What is the "meaning" of a word or sentence?
  - How to model context?
  - Other general knowledge?

# Desiderata for NLP models

- Sensitivity to a wide range of phenomena and constraints in human language
- Generality across languages, modalities, genres, styles
- Computational efficiency at construction time and runtime
- Strong formal guarantees (e.g., convergence, statistical efficiency, consistency)
- High accuracy when judged against expert annotations and/or test data specific to a particular task
- Explainable to human users
- Ethical

# Next class

- Text classification

Questions?