

# Natural Language Processing

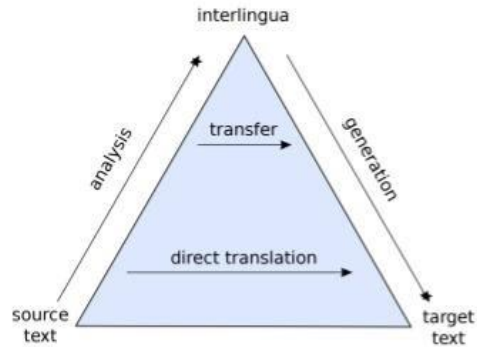
## Introduction to NLP

Yulia Tsvetkov

[yuliats@cs.washington.edu](mailto:yuliats@cs.washington.edu)

# Symbolic and Probabilistic NLP

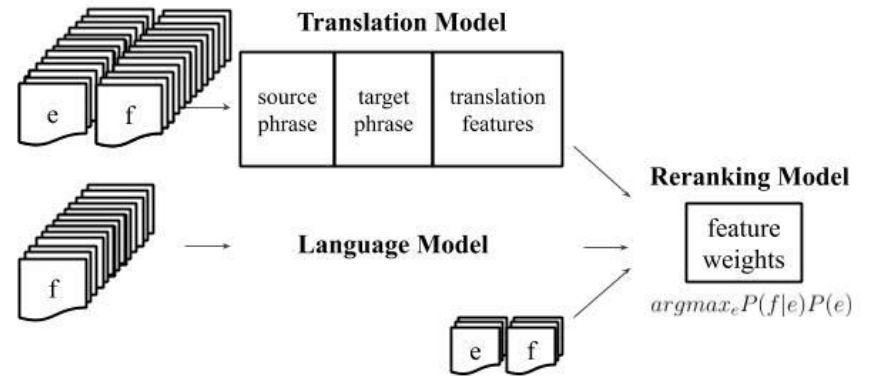
## Logic-based/Rule-based NLP



~ 90s

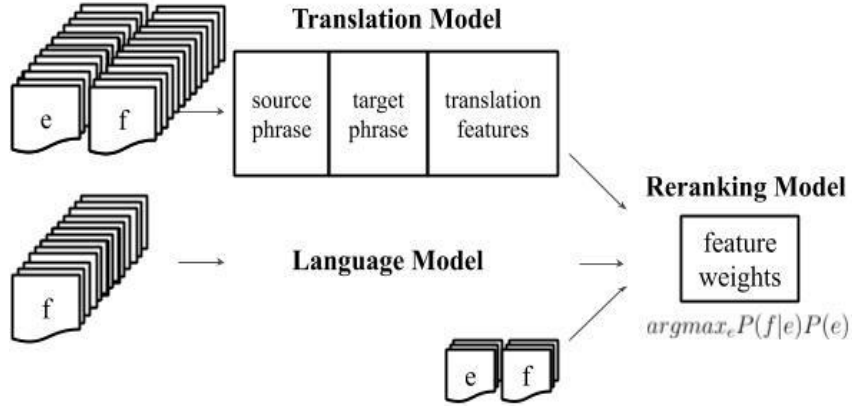


## Statistical NLP



# Probabilistic and Connectionist NLP

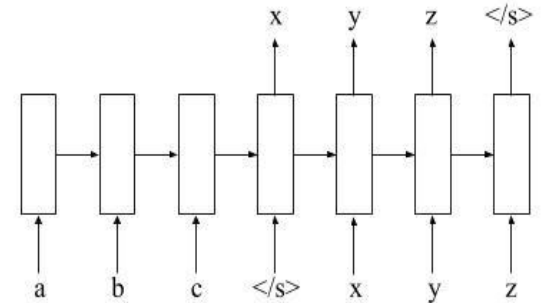
## Engineered Features/Representations



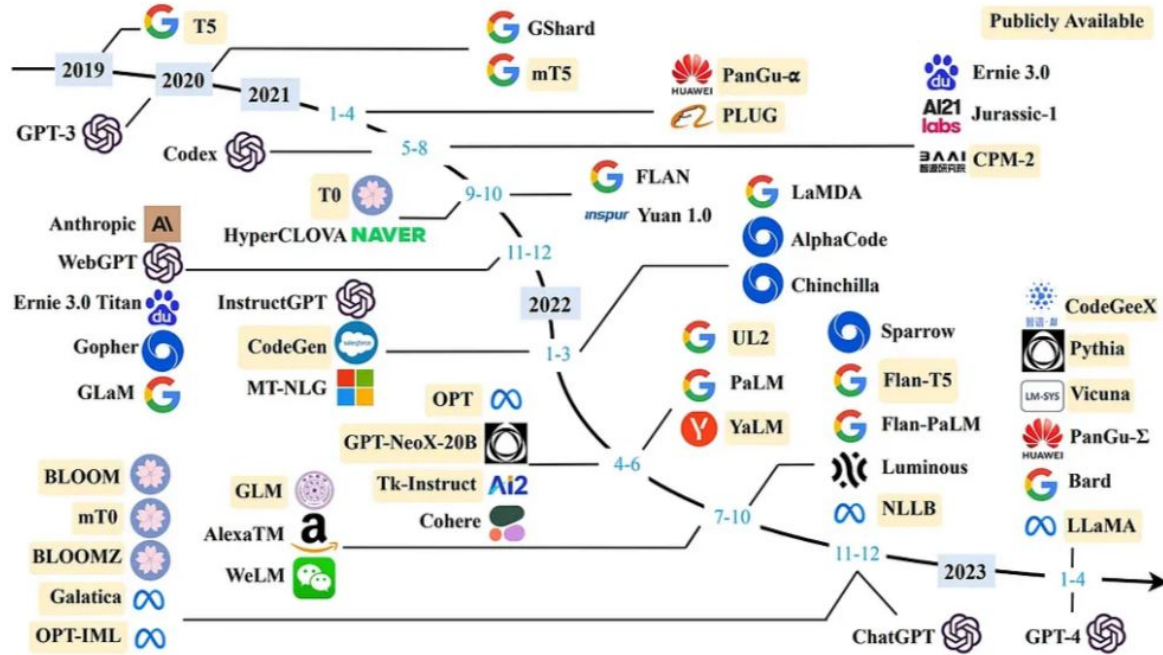
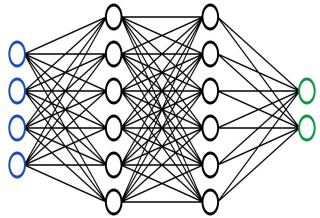
~mid 2010s



## Learned Features/Representations



# Large Language Models



Timeline of recent years large language models. Source: <https://www.nextbigfuture.com/2023/04/timeline-of-open-and-proprietary-large-language-models.html>

# Where are we now? - after 2022

ChatGPT 4o ▾

what are the weaknesses of NLP today?



NLP has made great strides but still faces key challenges:

1. **Lack of Deep Understanding:** Models often lack true **contextual comprehension**, struggling with nuanced reasoning and **long-term dependencies** in texts.
2. **Bias and Fairness:** NLP systems inherit **biases** from training data, which can result in unfair outputs. **Debiasing** remains difficult without affecting performance.
3. **Data Quality:** Models often rely on **noisy or low-quality data** from the web, and obtaining high-quality, domain-specific labeled data is expensive and time-consuming.
4. **Domain-Specific Limitations:** General NLP models don't perform well in specialized domains like **medicine** or **law** without significant **adaptation**.
5. **Explainability and Reasoning:** Most models lack **explainability** and struggle with **logical reasoning** or **multi-step inferences**, which is critical for high-stakes applications.
6. **Efficiency and Scalability:** Large models are **resource-intensive**, requiring significant computational power for both training and inference, which limits accessibility.
7. **Multilingual Limitations:** NLP models perform poorly in **low-resource languages** and **code-switching** scenarios, leaving many languages underrepresented.
8. **Safety and Misinformation:** Models can be vulnerable to **adversarial attacks** and can generate **misinformation** easily, posing ethical concerns.

These weaknesses highlight the need for improvements in understanding, fairness, domain adaptation, and safety in NLP systems.



# Question answering

The screenshot shows a Google search interface. At the top, the search bar contains the text "who invented neural machine translation". Below the search bar, there are navigation tabs for "All", "News", "Images", "Videos", "Shopping", and "More". The search results indicate "About 6,450,000 results (0.71 seconds)". A translation widget is displayed, showing the original query in English: "who invented neural machine" and its French translation: "qui a inventé la machine neuronale". The widget includes a bidirectional arrow between the language selectors, a close button (X), and icons for voice search and audio playback. At the bottom of the widget, there are links for "Open in Google Translate" and "Feedback".

Retrieved Mar 25, 2022

# Machine translation

## English → French

Translate Turn off instant translation

Russian English French Detect language ▾ ↔ English Spanish French ▾ [Translate](#)

You will just have to find a way of getting over it. ✕ Vous devrez trouver un moyen de le surmonter.

▾ 52/5000 ☆  [Suggest an edit](#)

## French → English

Translate Turn off instant translation

Russian English French Detect language ▾ ↔ English Spanish French ▾ [Translate](#)

Vous devrez trouver un moyen de le surmonter. ✕ You will have to find a way to overcome it.

▾ 45/5000 ☆  [Suggest an edit](#)

Did you mean: Vous devez trouver un moyen de le surmonter.

# Machine translation

## English → Swahili

Translate Turn off instant translation

Russian English French Detect language

English Swahili French Translate

You will just have to find a way of getting over it. x

Utakuwa tu kupata njia ya kupata juu yake.

53/5000 Suggest an edit

## Swahili → English

Translate Turn off instant translation

Swahili English French Detect language

English Swahili French Translate

Utakuwa tu kupata njia ya kupata juu yake. x

You will just find the way to get on it.

42/5000 Suggest an edit



# Machine translation

## English → Hindi → English

Hindi English Yoruba Detect language ▾

English Yoruba Hindi ▾ Translate

आपको इसे खत्म करने का एक तरीका मिलना होगा।

You have to find a way to eliminate it.

42/5000 Suggest an edit

## English → Telugu → English

Uzbek English Telugu Detect language ▾

English Uzbek Telugu ▾ Translate

మీరు దాని పైకి రావడానికి ఒక మార్గాన్ని కనుగొనవలసి ఉంటుంది.

You have to find a way to get it up.

59/5000 Suggest an edit

## English → Uzbek → English

Pashto English Uzbek Detect language ▾

English Uzbek Yoruba ▾ Translate

Buning ustiga faqatgina bir usulni topish kerak.

On top of that, you just have to find a way out.

48/5000 Suggest an edit

# Machine translation

## English → Swahili

The screenshot shows the Google Translate interface for the translation of English text into Swahili. The source text is: "The summer school is meant to be an introduction to the state-of-the-art research in the speech and language technology area for graduate and undergraduate students." The translated text is: "Shule ya majira ya joto ina maana ya kuanzishwa kwa utafiti wa hali ya sanaa katika eneo la teknolojia na lugha ya wanafunzi kwa wanafunzi wahitimu na wahitimu." The interface includes language selection buttons for Swahili, English, and Telugu, a "Detect language" option, and a "Translate" button. There are also icons for voice input, copy, and a "Suggest an edit" link.

## Swahili → English

The screenshot shows the Google Translate interface for the translation of Swahili text into English. The source text is: "Shule ya majira ya joto ina maana ya kuanzishwa kwa utafiti wa hali ya sanaa katika eneo la teknolojia na lugha ya wanafunzi kwa wanafunzi wahitimu na wahitimu." The translated text is: "Summer school means the establishment of a state-of-the-art arts research technology and pupil language for graduate students and graduates." The interface includes language selection buttons for Swahili, English, and Telugu, a "Detect language" option, and a "Translate" button. There are also icons for voice input, copy, and a "Suggest an edit" link.

# Bias in machine translation

Translate

Turn off instant translation

Bengali English Hungarian Detect language ↔ English Spanish Hungarian Translate

ő egy ápoló.  
ő egy tudós.  
ő egy mérnök.  
ő egy pék.  
ő egy tanár.  
ő egy esküvői szervező.  
ő egy vezérigazgatója.

110/5000

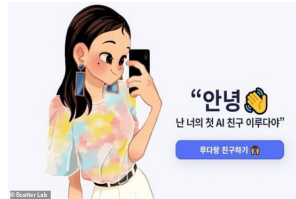
she's a nurse.  
he is a scientist.  
he is an engineer.  
she's a baker.  
he is a teacher.  
She is a wedding organizer.  
he's a CEO.

What can we do about this problem? We'll discuss in NLP class!

# Hate speech

## Twitter taught Microsoft's AI chatbot to be a racist asshole in less than a day

By James Vincent | Mar 24, 2016, 6:43am EDT  
Via *The Guardian* | Source *TayandYou* (Twitter)



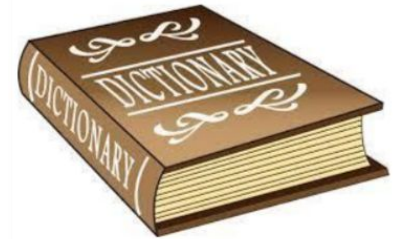
**AI chatbot is REMOVED from Facebook after saying she 'despised' gay people, would 'rather die' than be disabled and calling the #MeToo movement 'ignorant'**

- Lee Luda is a South Korean chatbot with the persona of a 20-year-old student
- It has attracted more than 750,000 users since its launch last month
- But the chatbot has started using hate speech towards minorities
- In one of the captured chat shots, Luda said she 'despised' gays and lesbians
- The developer has apologised over the remarks, saying they 'do not represent our values as a company'



# Linguistic Background

# What does it mean to “know” a language?



Hi, how can I help?

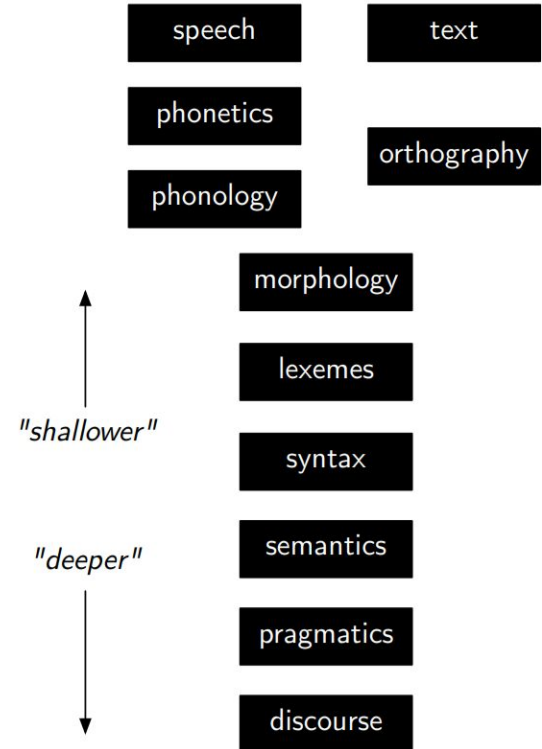
What do we need to “tell” a computer program so that it knows more English than  $w_c$  or a dictionary, maybe even as much as a three-year-old, for example?

# 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

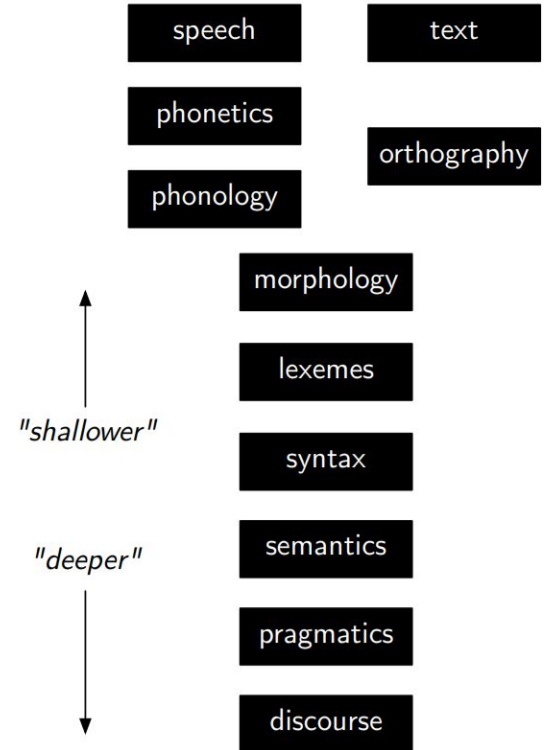


# Speech, phonetics, phonology



This is a simple sentence .

/ ðɪs ɪz ə 'sɪmpl 'sɛntəns /.



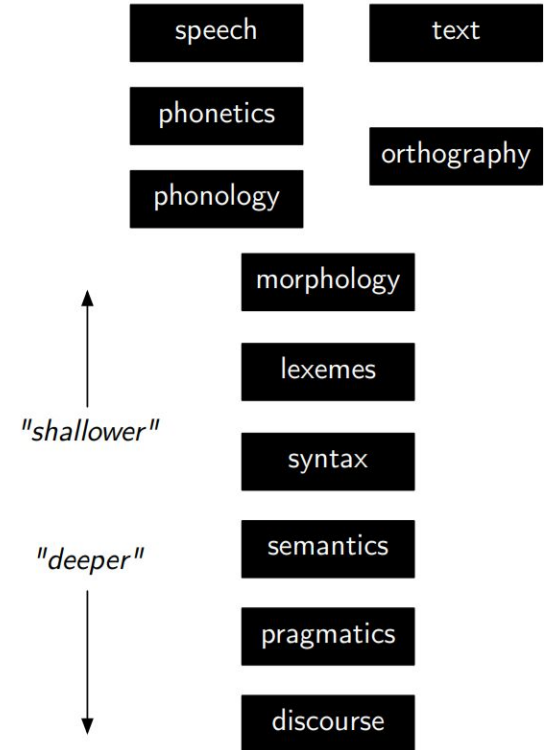
# Orthography

هذه جملة بسيطة

đây là một câu đơn giản

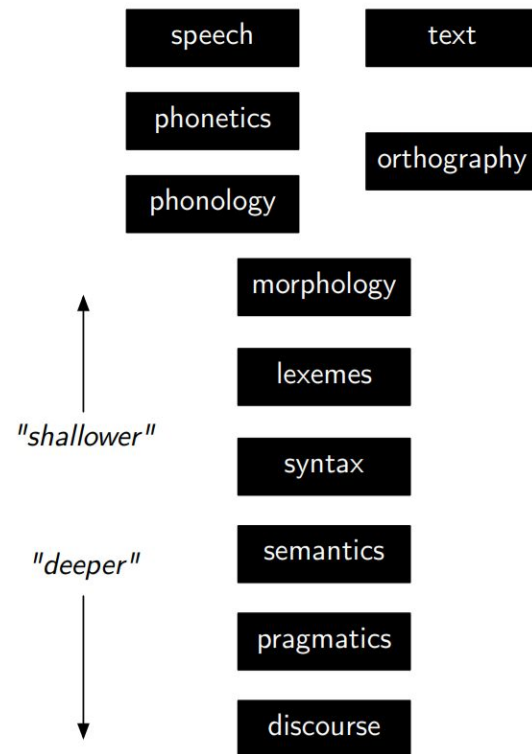
यह एक साधारण वाक्य है

This is a simple sentence .  
/ ðɪs ɪz ə 'sɪmpl 'sɛntəns /.



# Words, morphology

- Morphological analysis
- Tokenization
- Lemmatization



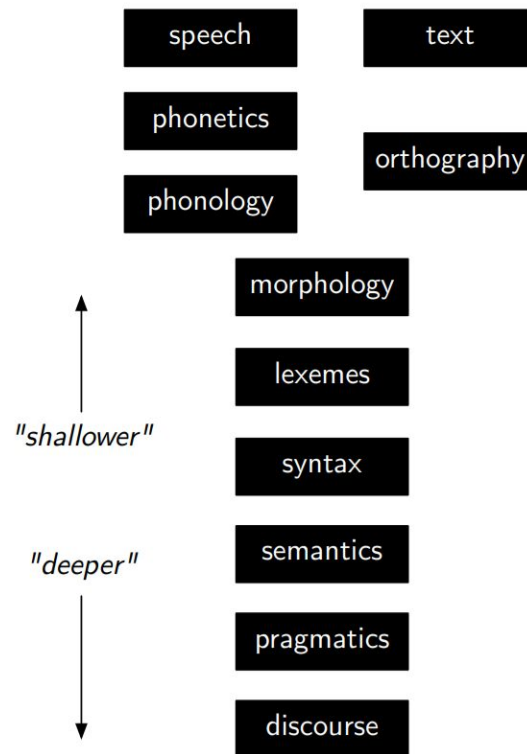
**Tokens** This is a simple sentence .

**Morphology**      be  
                         3sg  
                         present

# Syntax

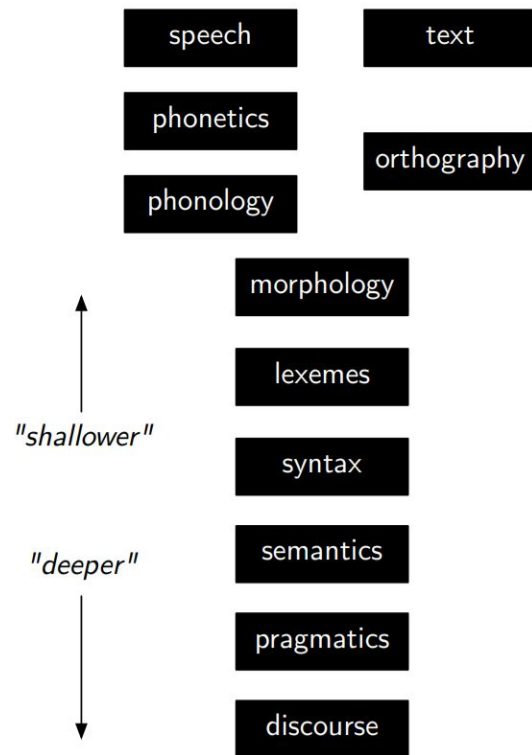
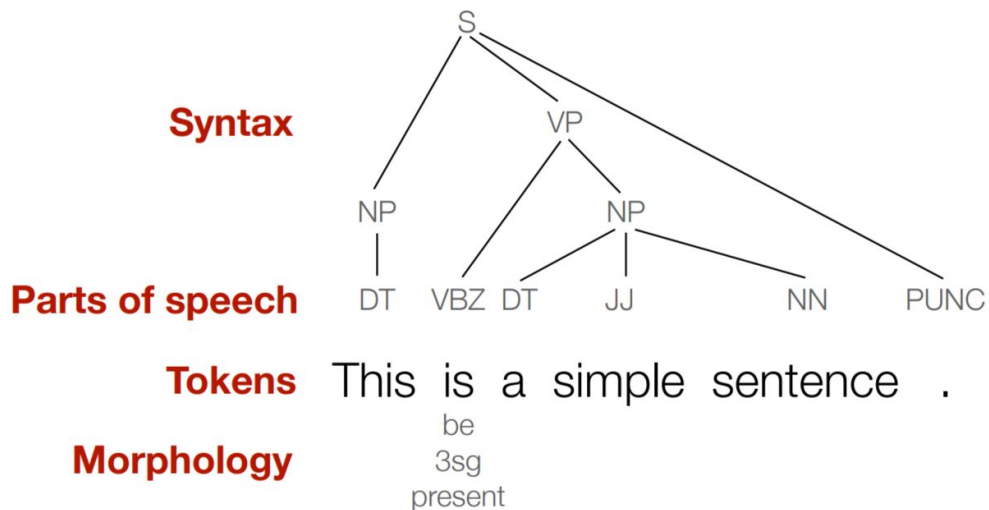
- Part-of-speech tagging

<b>Parts of speech</b>	DT	VBZ	DT	JJ	NN	PUNC
<b>Tokens</b>	This	is	a	simple	sentence	.
<b>Morphology</b>		be				
		3sg				
		present				



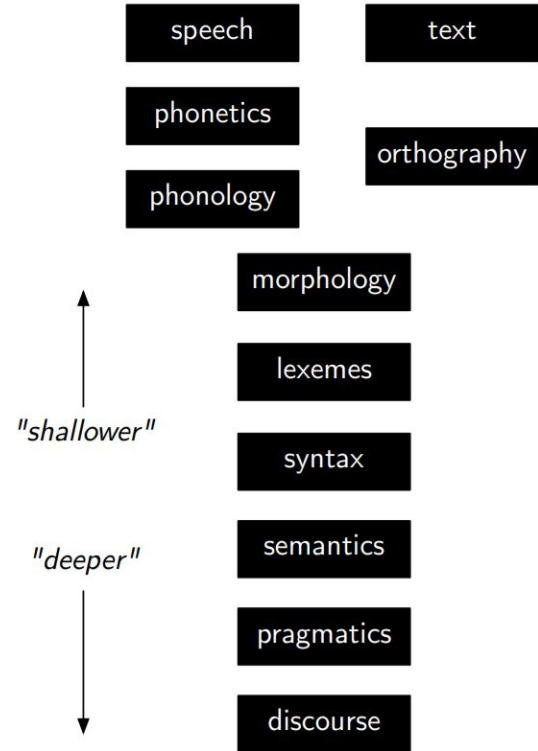
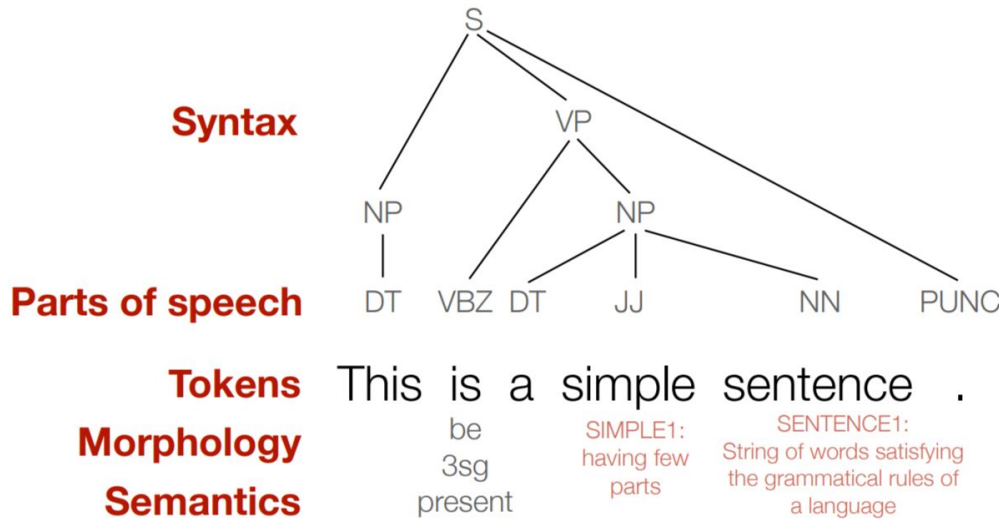
# Syntax

- Part-of-speech tagging
- Syntactic parsing



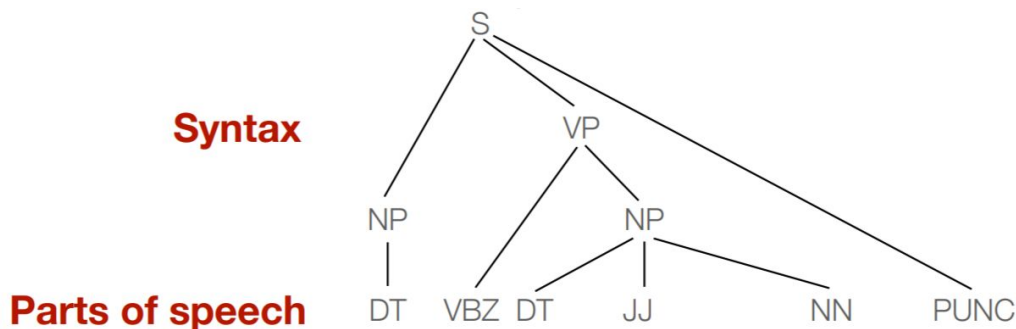
# Semantics

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



# Discourse

- Reference resolution
- Discourse parsing



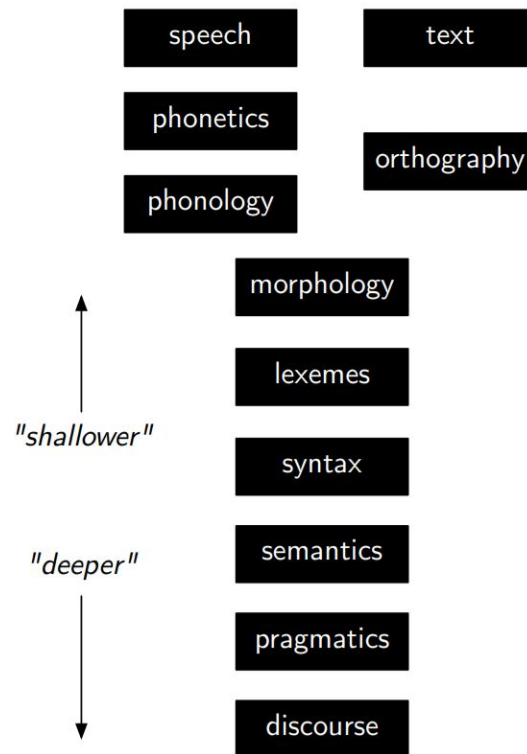
**Tokens** This is a simple sentence .

**Morphology** be SIMPLE1: SENTENCE1:  
3sg having few parts String of words  
present parts satisfying the  
grammatical rules  
of a language

**Semantics**

**Discourse** But an instructive one .

coreferent

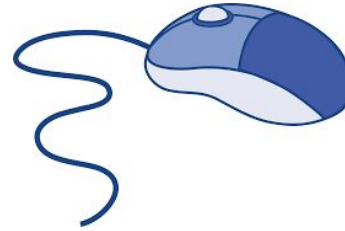




# Why is language interpretation hard?

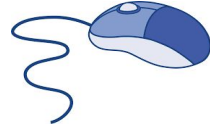
1. **Ambiguity**
2. Variation
3. Sparsity
4. Expressivity
5. Unmodeled variables
6. Unknown representation  $\mathcal{R}$

# Ambiguity: word sense disambiguation



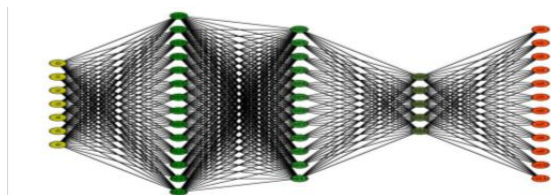
# Ambiguity

- 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



# Why is language interpretation hard?

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

# Variation

- ~7K languages
- Thousands of language varieties



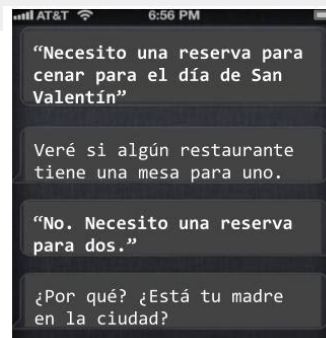
Englishes



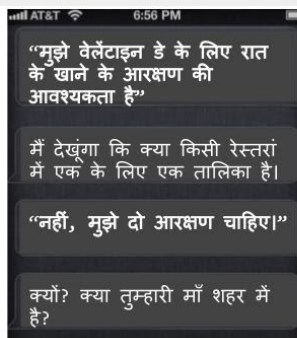
Africa is a continent with a very high linguistic diversity: there are an estimated 1.5-2K African languages from 6 language families. **1.33 billion people**

# NLP beyond English

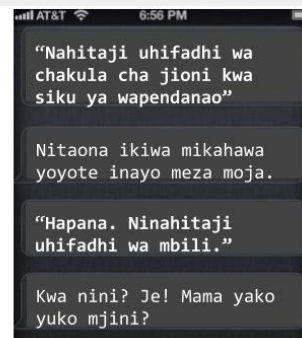
- ~7,000 languages
- thousands of language varieties



Spanish  
534 million speakers



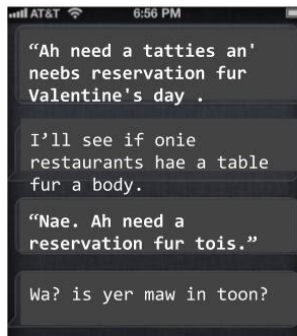
Hindi  
615 million speakers



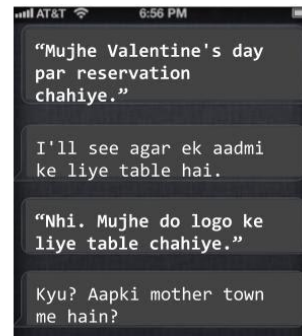
Swahili  
100 million speakers



American English



Scottish English

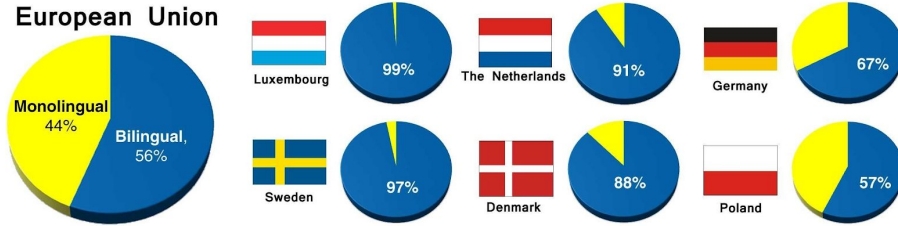


Hinglish



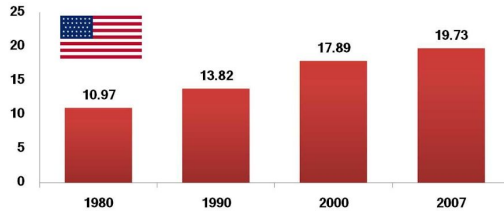
# Most of the world today is multilingual

## Percentage of Bilingual Speakers in the World



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

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

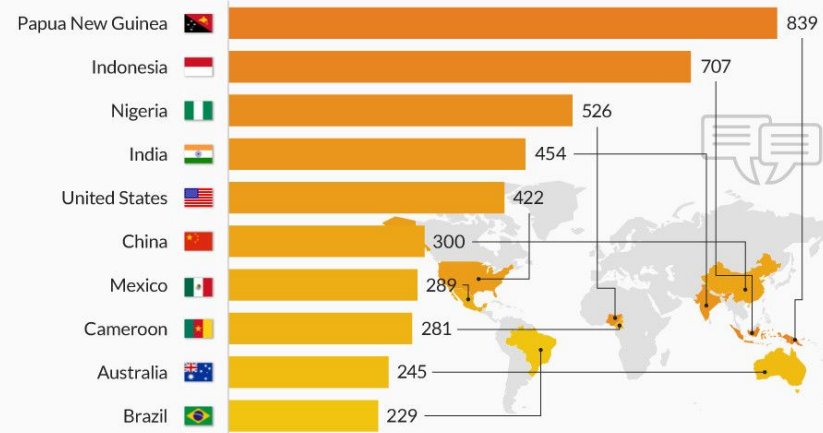


Source: U.S. Census Bureau, 2007 American Community Survey

Source: US Census Bureau

## The Countries With The Most Spoken Languages

Number of living languages spoken per country in 2015



Source: Ethnologue

# Tokenization

这是一个简单的句子

**WORDS**

This is a simple sentence

זה משפט פשוט

# Tokenization + disambiguation

in tea  
her daughter

בתה

- 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	ושבהתה

ושבתה

and her saturday	ו+שבת+ה
and that in tea	ו+ש+ב+ת+ה
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'ananyakapushasqakupuniñataqsunamá

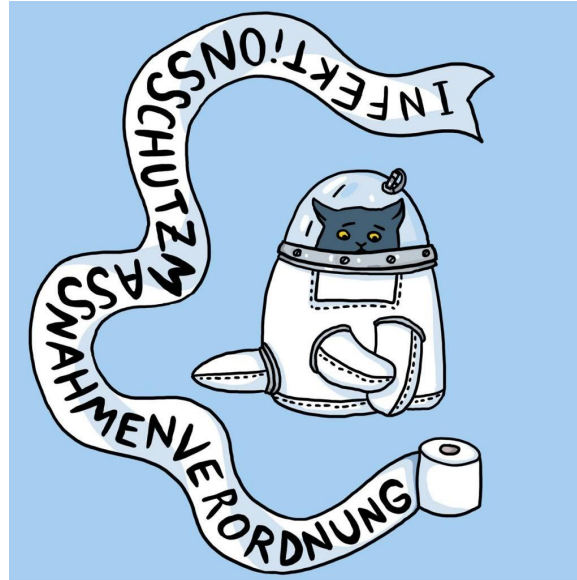
Much'a -na -naya -ka -pu -sha -sqa -ku -puni -ña -taq -suna -má

*"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

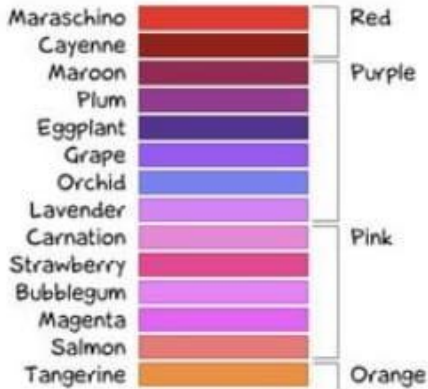
- German



Infektionsschutzmaßnahmenverordnung

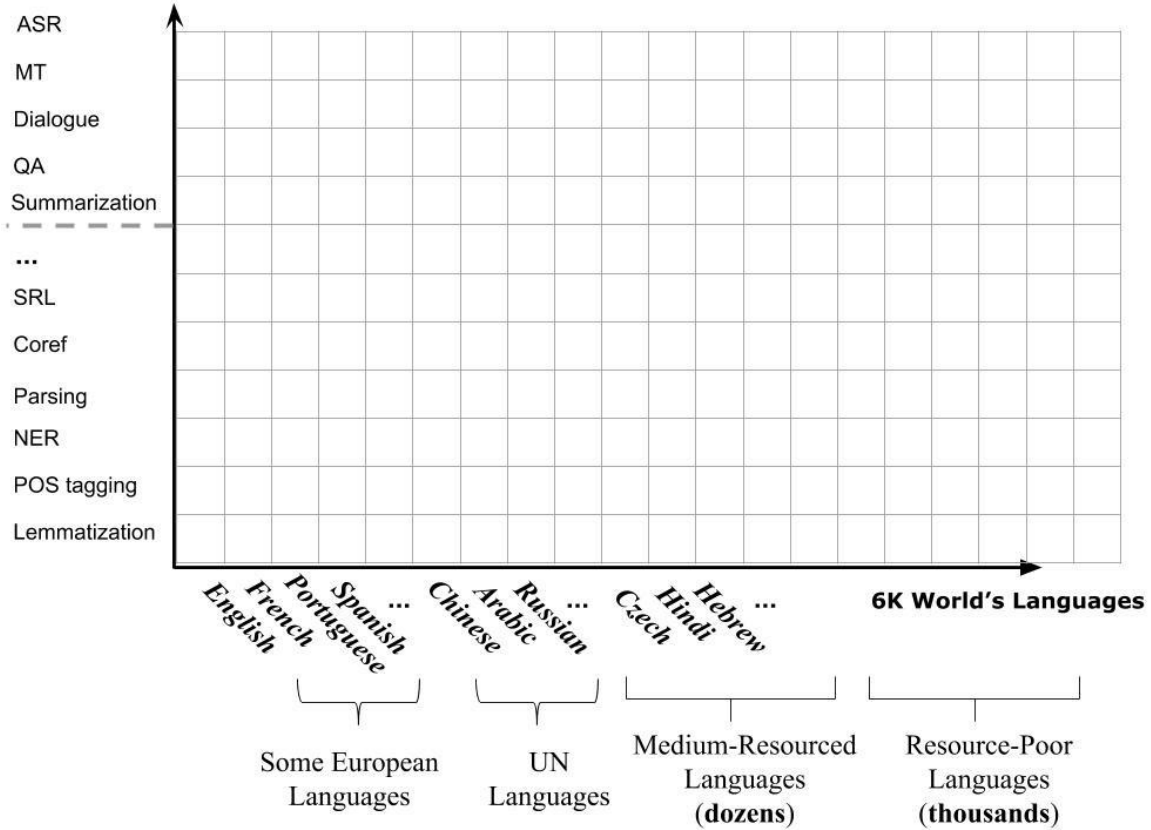
# Semantic analysis

- Every language sees 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

**NLP Technologies/Applications**



# Linguistic variation

- Non-standard language, emojis, hashtags, names

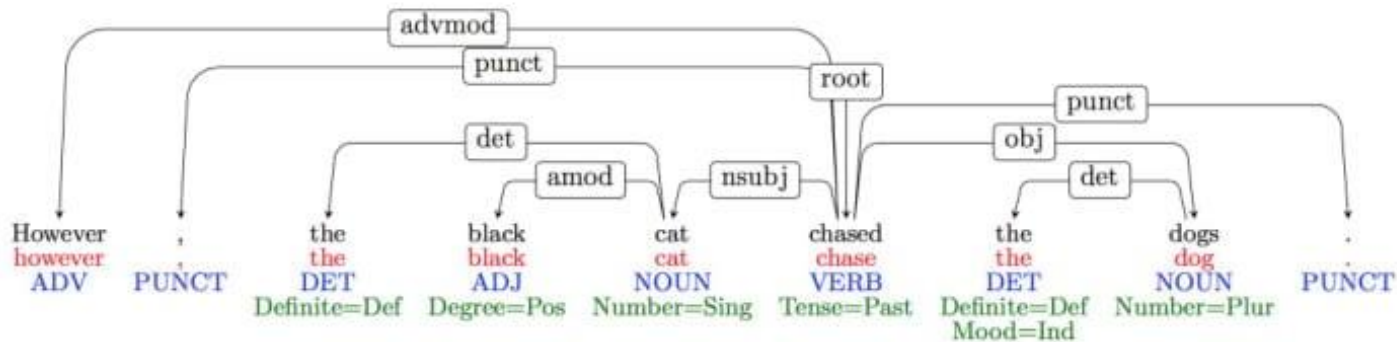


**chowdownwithchan** #crab and #pork #xiaolongbao at @dintaifungusa... where else? 🤔👩 Note the cute little crab indicator in the 2nd pic 🦀💕



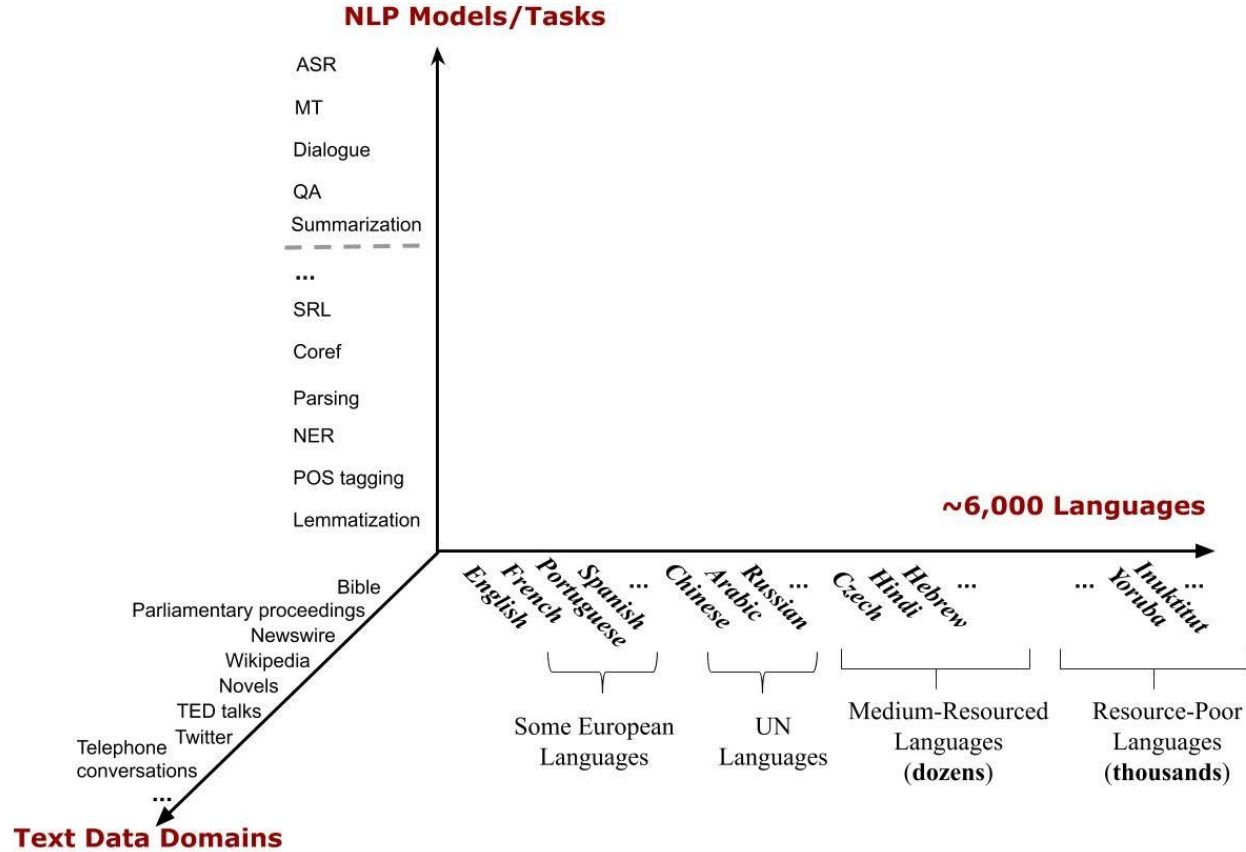
# Variation

- 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 🤔🙌



# Why is language interpretation hard?

1. Ambiguity
2. Scale
3. Variation
4. Sparsity
5. Expressivity
6. Unmodeled variables
7. 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)

<b>any word</b>		<b>nouns</b>	
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

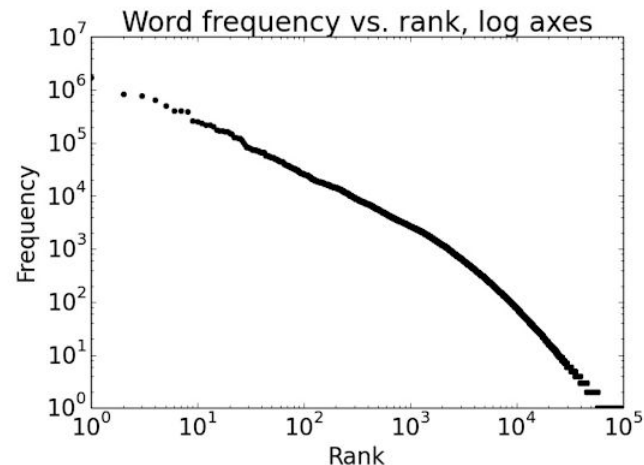
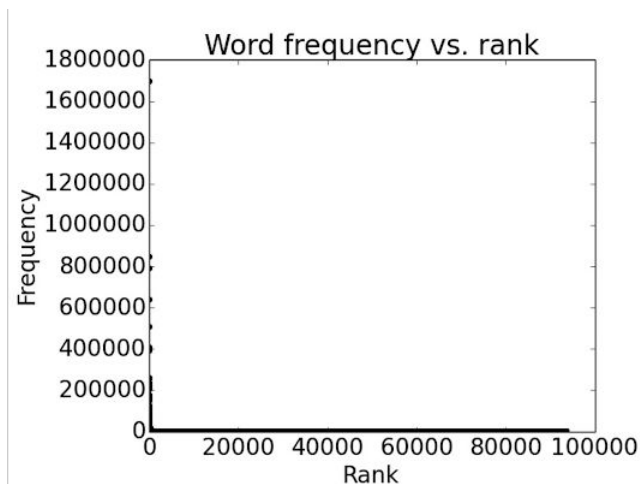
But also, out of 93,638 distinct words (word types), 36,231 occur only once.

Examples:

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

# Plotting word frequencies

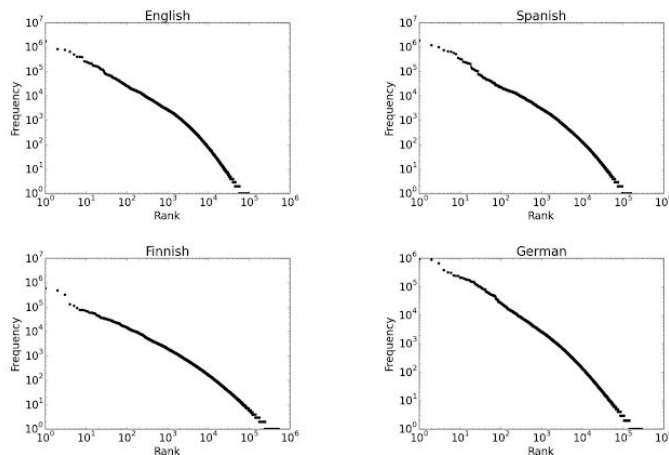
Order words by frequency. What is the frequency of  $n$ th ranked word?



# 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





# Why is language interpretation hard?

1. Ambiguity
2. Scale
3. Variation
4. Sparsity
5. Expressivity
6. Unmodeled variables
7. 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

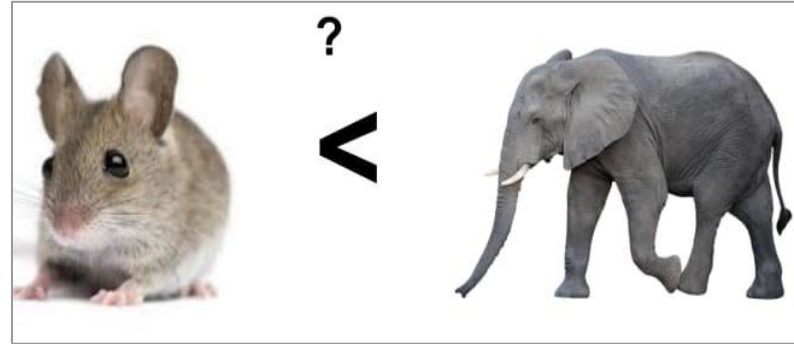
# Why is language interpretation hard?

1. Ambiguity
2. Scale
3. Variation
4. Sparsity
5. Expressivity
6. **Unmodeled variables**
7. 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

# Why is language interpretation hard?

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

# Unknown representation

- Very difficult to capture *what is  $\mathcal{R}$* , 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
- Strong formal guarantees (e.g., convergence, statistical efficiency, consistency)
- High accuracy when judged against expert annotations or test data
- Ethical

# NLP $\stackrel{?}{=}$ Machine Learning

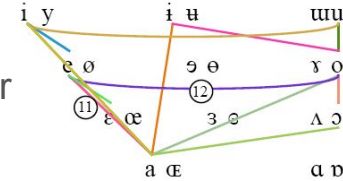
- To be successful, a machine learner needs bias/assumptions; for NLP, that might be linguistic theory/representations.
- Symbolic, probabilistic, and connectionist ML have all seen NLP as a source of inspiring applications.



# What is nearby NLP?

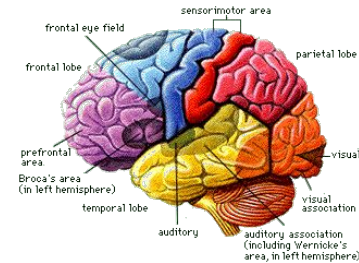
- Computational Linguistics

- Using computational methods to learn more about how language works
- We end up doing this and using it



- Cognitive Science

- Figuring out how the human brain works
- Includes the bits that do language
- Humans: the only working NLP prototype!



- Speech Processing

- Mapping audio signals to text
- Traditionally separate from NLP, converging?
- Two components: acoustic models and language models
- Language models in the domain of stat NLP



# Next class

- Classification

## Questions?