# Natural Language Processing

## Introduction to NLP

### Yulia Tsvetkov

yuliats@cs.washington.edu

PAUL G. ALLEN SCHOOL
OF COMPUTER SCIENCE & ENGINEERING

# Announcements

`https://courses.cs.washington.edu/courses/cse447/22au/`

- Quiz 1: Monday Oct 10
  - 5 questions, open during lecture time, 10-min in the end of the class
  - Materials from weeks 1 and 2
    - Introduction to NLP, introduction to text classification, NB
    - Instructions for HW 1
- Discussions on Ed
  - Reminder about 10% bonus grade for commenting on Ed
  - TAs will respond to your questions within 24 hours
- Office hours
  - TA OH locations have been updated on the website
  - Yulia is traveling in the rest of this week, no OHs on Friday
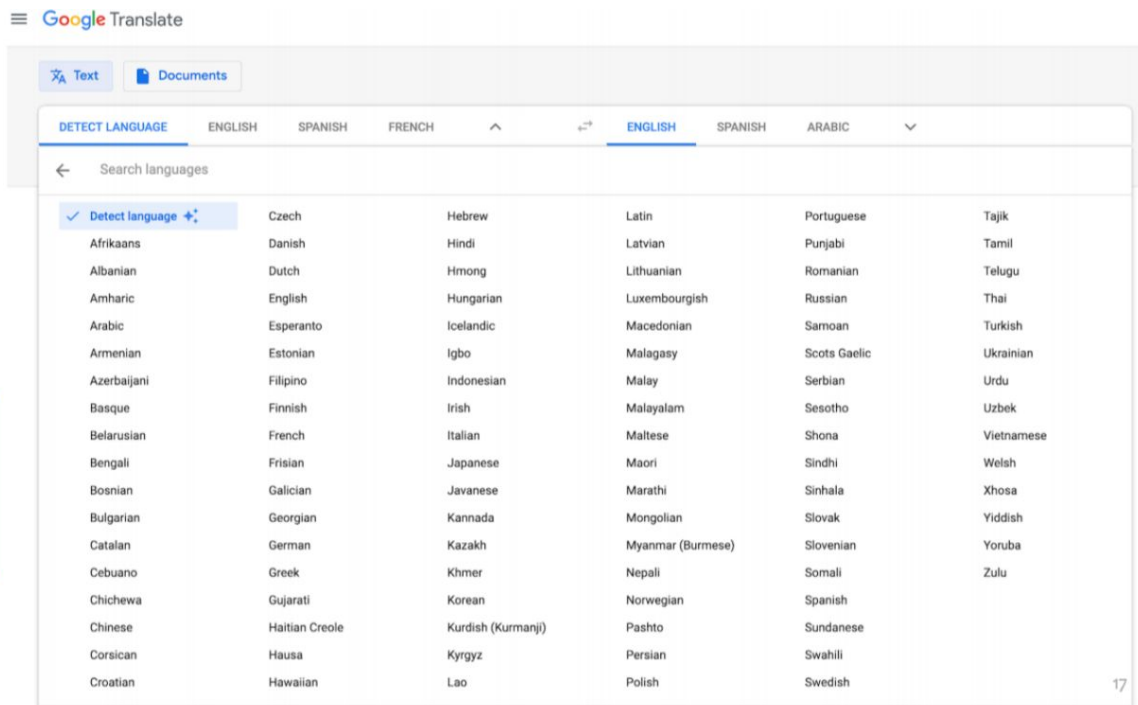
# Personal assistants

# Question answering

- What does "divergent" mean?

- What year was Abraham Lincoln born?

- How many states were in the United States that year?

- How much Chinese silk was exported to England in the end of the 18th century?

- What do scientists think about the ethics of human cloning?

# Machine translation

# Sentiment analysis

**HP Officejet 6500A Plus e-All-in-One Color Ink-jet - Fax / copier / printer / scanner**
$89 online, $100 nearby ★★★★☆ 377 reviews
September 2010 - Printer - HP - Inkjet - Office - Copier - Color - Scanner - Fax - 250 she

## Reviews

**Summary** - Based on 377 reviews

| 1 star | 2 | 3 | 4 stars | 5 stars |
|---|---|---|---|---|

What people are saying

| | | |
|---|---|---|
| ease of use | | "This was very easy to setup to four computers." |
| value | | "Appreciate good quality at a fair price." |
| setup | | "Overall pretty easy setup." |
| customer service | | "I DO like honest tech support people." |
| size | | "Pretty Paper weight." |
| mode | | "Photos were fair on the high quality mode." |
| colors | | "Full color prints came out with great quality." |

# Information extraction

**DATE**
1933
**PERSON** **COUNTRY** **PERSON**
In 1933 , while Einstein was visiting the United States , Adolf Hitler came to power .

**RELIGION** **PERSON** **COUNTRY**
Because of his Jewish background , Einstein did not return to Germany .

**COUNTRY** **NATIONALITY** **DATE** 1940
He settled in the United States and became an American citizen in 1940 .

**PERSON** **MISC**
Einstein supported the Allied forces , but he generally denounced the idea of using nuclear fission as a weapon .

**PERSON** **NATIONALITY** **TITLE** **PERSON**
He signed the Russell -- Einstein Manifesto with British philosopher Bertrand Russell , which highlighted the danger of nuclear weapons .

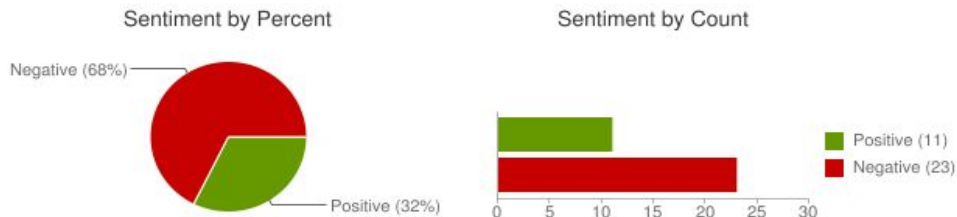**ORGANIZATION** **STATE_OR_PROVINCE** **DATE** 1955
He was affiliated with the Institute for Advanced Study in Princeton , New Jersey , until his death in 1955 .

# Sentiment analysis + information extraction

Type in a word and we'll highlight the good and the bad

"united airlines"     Search   Save this search

**Sentiment analysis for "united airlines"**

Sentiment by Percent                    Sentiment by Count

Negative (68%)

Positive (32%)

Positive (11)
Negative (23)

jljacobson: OMG... Could **@United airlines** have worse customer service? W8g now 15 minutes on hold 4 questions about a flight 2DAY that need a human.
Posted 2 hours ago

12345clumsy6789: I hate **United Airlines** Ceiling!!! Fukn impossible to get my conduit in this damn mess! ?
Posted 2 hours ago

EMLandPRGbelgiu: EML/PRG fly with Q8 **united airlines** and 24seven to an exotic destination. http://t.co/Z9QloAjF
Posted 2 hours ago

CountAdam: FANTASTIC customer service from **United Airlines** at XNA today. Is tweet more, but cell phones off now!
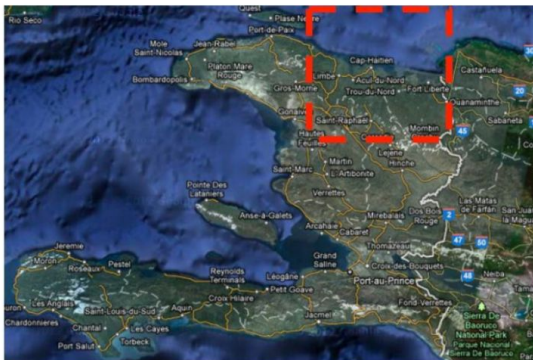Posted 4 hours ago

# Information extraction for disaster relief

- Translation + information extraction

Lopital Sacre-Coeur ki nan vil Okap, pre pou li resevwa moun malad e lap mande pou moun ki malad yo ale la.

"Sacre-Coeur Hospital which located in this village of Okap is ready to receive those who are injured. Therefore, we are asking those who are sick to report to that hospital."
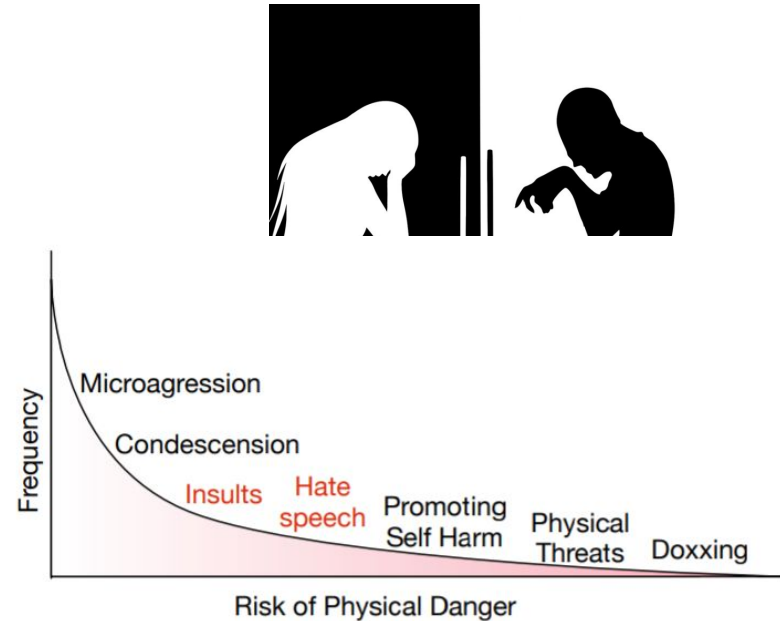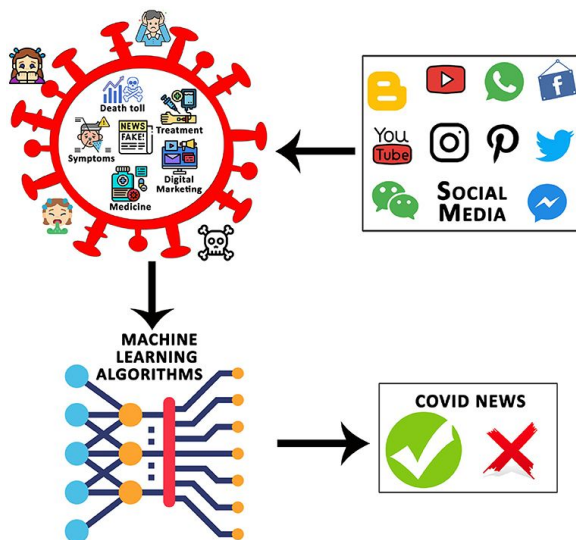
iDiBON

An earthquake struck Haiti on January 12, 2010

Most local services failed, but most cell-towers remained functional.

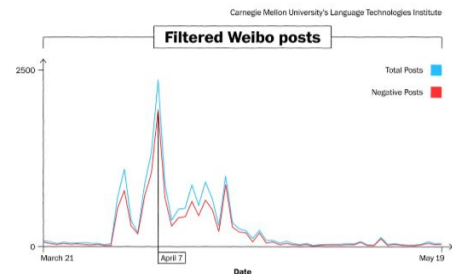# Hate speech detection

# Covid19 misinformation

The Fact Checker worked with researchers at professor Yulia Tsvetkov's lab at Carnegie Mellon University's Language Technologies Institute and the Center for Human Rights Science to track what happened on social media during this period. Researchers collected about 16,000 Weibo posts, filtered from a larger data set of 200,000 posts, containing at least one Guangzhou location tag and one "African-related" keyword from late March through May. Weibo is a Chinese social media platform.

Based on automated sentiment analysis tools and manual analysis, the researchers said, they believed the majority of posts in their data set expressed negative sentiments relating to Africans or black people.

Carnegie Mellon University's Language Technologies Institute

**Filtered Weibo posts**

Total Posts
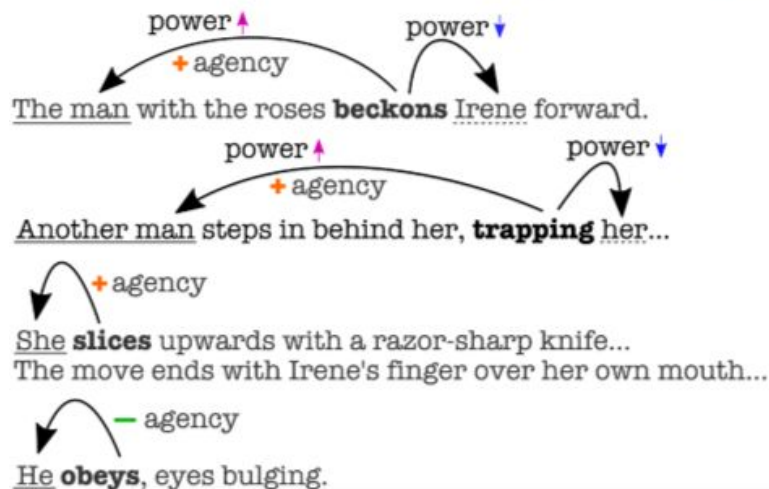Negative Posts

2500

0

March 21    April 7    May 19
Date

Researchers collected around 16,000 posts containing at least one Guangzhou location tag and one "African-related" keyword from late March through May. (Atthar Mirza/The Washington Post)

Their research showed there was a significant surge in negative posts beginning April 1. There were just 23 negative posts in the data set on March 31. The next day, the number of posts climbed to 500. From April 1-2, there was a spike in the number of posts on Weibo using the keywords "foreign trash." A Chinese cartoon depicting officials throwing foreigners who weren't abiding by

**Detecting COVID-19-Related Fake News Using Feature Extraction**

Suleman Khan, Saqib Hakak, N. Deepa, B. Prabadevi, Kapal Dev and Silvia Trelova

https://www.washingtonpost.com/politics/2020/06/18/video-evidence-anti-black-discrimination-china-over-coronavirus-fears/
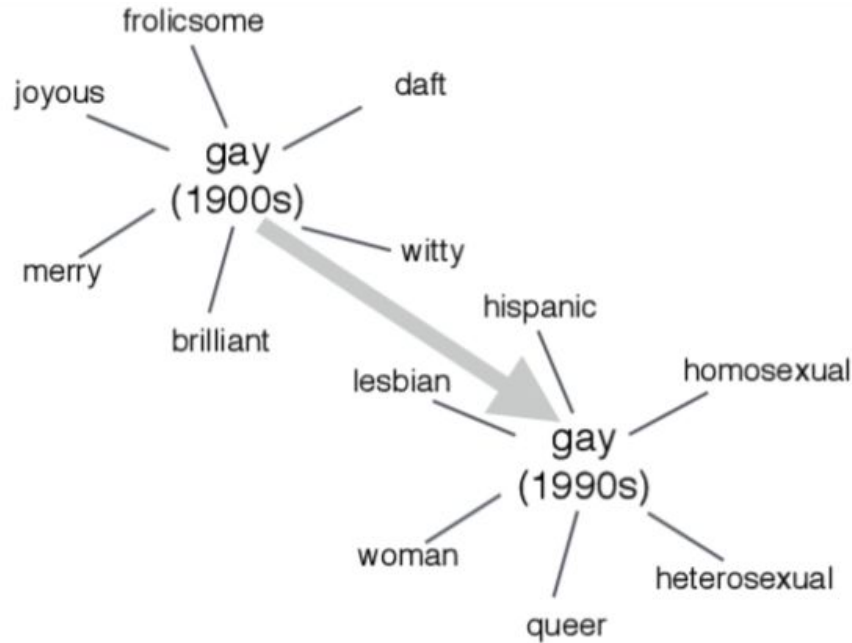
# Computational social science

- computational social science answering questions about society given observational data
- example: "do movie scripts portray female or male characters with more power or agency?" [Sap+ 2017]



| Frame | $\beta$ | gender |
|---|---|---|
| $agency(\text{AG})=+$ | $-0.951$ | $\mathbf{M}^{**}$ |
| $power(\text{AG}>\text{TH})$ | $-0.468$ | $\mathbf{M}^{**}$ |
| $agency(\text{AG})=-$ | $0.277$ | $\mathbf{F}^{**}$ |
| $power(\text{AG}<\text{TH})$ | *not sig.* | |

# Language change



**Cultural Shift or Linguistic Drift? Comparing Two Computational Measures of Semantic Change**

William L. Hamilton, J. Leskovec, Dan Jurafsky

# Natural Language Processing
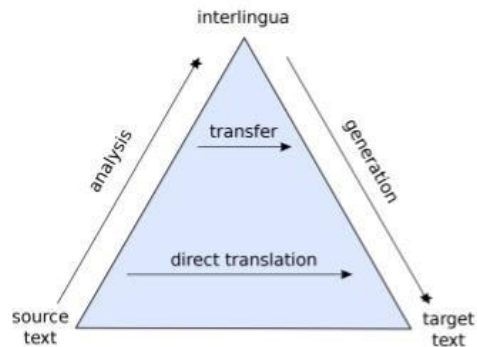
- Applications
  - Machine Translation

  - Information Retrieval

  - Question Answering

  - Dialogue Systems

  - Information Extraction

  - Summarization

  - Sentiment Analysis
  - ...

- Core technologies
  - Language modelling

  - Part-of-speech tagging

  - Syntactic parsing

  - Named-entity recognition

  - Coreference resolution

  - Word sense disambiguation

  - Semantic Role Labelling
  - ...
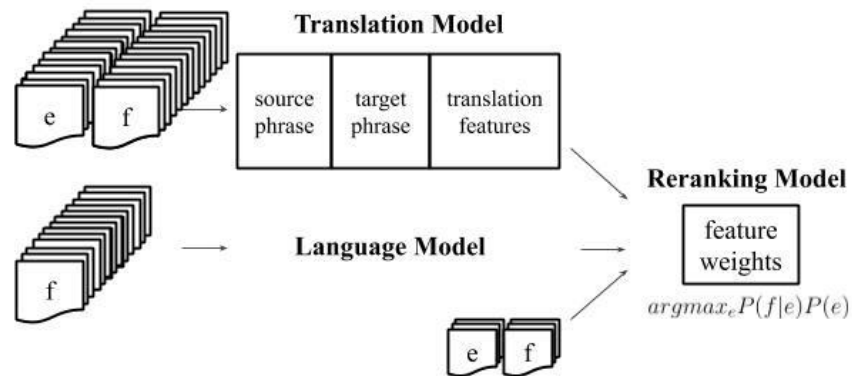
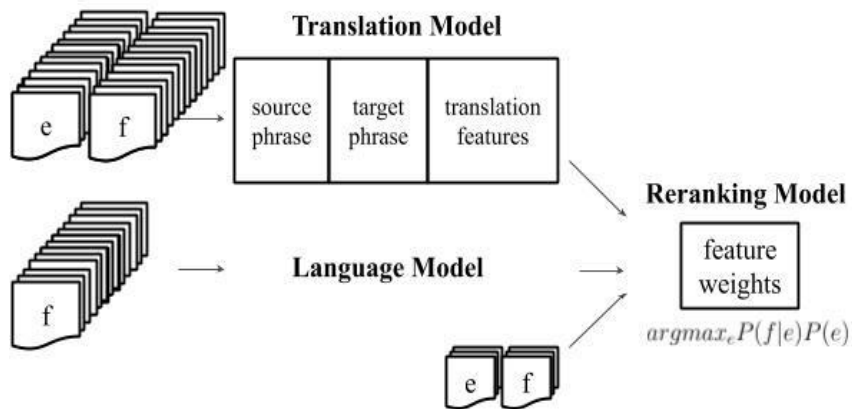# Symbolic and Probabilistic NLP

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

interlingua

analysis / generation

transfer

direct translation

source text — target text

~ 90s

**Statistical NLP**

**Translation Model**

e  f → source phrase | target phrase | translation features

**Reranking Model**

feature weights

f → **Language Model**

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

e  f

# Probabilistic and Connectionist NLP



**Engineered Features/Representations**

**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

~mid 2010s

**Learned Features/Representations**

x    y    z    </s>

a    b    c    </s>    x    y    z

# GPT-3

OpenAI GPT-3

# Where are we now?

# Chat bots

JOAQUIN PHOENIX   AMY ADAMS   ROONEY MARA
OLIVIA WILDE   SCARLETT JOHANSSON

**her**

A SPIKE JONZE LOVE STORY

IN CINEMAS FEBRUARY 14TH

REPORT

## OPENAI'S LATEST BREAKTHROUGH IS ASTONISHINGLY POWERFUL, BUT STILL FIGHTING ITS FLAWS

*The ultimate autocomplete*

By James Vincent | Jul 30, 2020, 10:01am EDT

```
Q: How many eyes does a giraffe have?
A: A giraffe has two eyes.

Q: How many eyes does my foot have?
A: Your foot has two eyes.

Q: Which is heavier, a mouse or a toaster?
A: A toaster is heavier than a mouse.

Q: Which is heavier, a toaster or a pencil?
A: A pencil is heavier than a toaster.
```

## OpenAI's new language generator GPT-3 is shockingly good—and completely mindless

The AI is the largest language model ever created and can generate amazing human-like text on demand but won't bring us closer to true intelligence.

by **Will Douglas Heaven**                    July 20, 2020

https://lacker.io/ai/2020/07/06/giving-gpt-3-a-turing-test.html

# Where are we now? GPT-3 Turing Test

This gives us a hint for how to stump the AI more consistently. We need to ask questions that no normal human would ever talk about.

```
Q: How many eyes does a giraffe have?
A: A giraffe has two eyes.

Q: How many eyes does my foot have?
A: Your foot has two eyes.

Q: How many eyes does a spider have?
A: A spider has eight eyes.

Q: How many eyes does the sun have?
A: The sun has one eye.

Q: How many eyes does a blade of grass have?
A: A blade of grass has one eye.
```

https://lacker.io/ai/2020/07/06/giving-gpt-3-a-turing-test.html

# GPT-3 Turing Test

Now we're getting into surreal territory. GPT-3 knows how to have a normal conversation. It doesn't quite know how to say "Wait a moment… your question is nonsense." It also doesn't know how to say "I don't know."

```
Q: How do you sporgle a morgle?
A: You sporgle a morgle by using a sporgle.

Q: How many bonks are in a quoit?
A: There are three bonks in a quoit.

Q: How many rainbows does it take to jump from Hawaii to seventeen?
A: It takes two rainbows to jump from Hawaii to seventeen.

Q: Which colorless green ideas sleep furiously?
A: Ideas that are colorless, green, and sleep furiously are the ideas of a sleep
furiously.

Q: Do you understand these questions?
A: I understand these questions.
```

https://lacker.io/ai/2020/07/06/giving-gpt-3-a-turing-test.html

# Where are we now? GPT-3 Turing Test

One trend that continues from the common sense is that GPT-3 is reluctant to express that it doesn't know the answer. So invalid questions get wrong answers.

```
Q: Who was president of the United States in 1700?
A: William Penn was president of the United States in 1700.

Q: Who was president of the United States in 1650?
A: Richard Bennett was president of the United States in 1650.

Q: Who was president of the United States in 1620?
A: James I was president of the United States in 1620.

Q: Who was president of the United States in 1600?
A: Queen Elizabeth I was president of the United States in 1600.
```

https://lacker.io/ai/2020/07/06/giving-gpt-3-a-turing-test.html

# Where are we now? GPT-3

An interesting corollary is that GPT-3 often finds it easier to write code to solve a programming problem, than to solve the problem on one example input:

```
Q: Write one line of Ruby code to reverse an array.
A: ary.reverse

Q: Write one line of Ruby code to count how many X's are in an array.
A: ary.count { |x| x == 'X' }

Q: Reverse the following array: [1, 3, 5, 6, 10, 4, 2, 77]
A: [10, 6, 4, 2, 77, 3, 5, 1]

Q: How many X's are in this array: [1, 3, 5, 6, 10, 4, 2, 77]
A: 2
```

# Where are we now? Meena → LaMDA

**Conversation A**

| | |
|---|---|
| **Human:** | *say you're in an empty giant office and you see a black bird flying indoors. what do you do?* |
| **Meena:** | *What kind of bird? Is it friendly? How large? How big is the office? Is it indoors or outdoors?* |
| **Human:** | *well it's a predator bird and it's gonna eat you* |
| **Meena:** | *I would grab my phone and call 911.* |

```
Human: I've been using Python for 15 years
Meena: Ah I see, so I'm guessing you're pretty pro at it by now then?
Human: probably
Meena: What's your current project?
Human: NLP
Meena: Neuro-linguistic Programming?
Human: correct
```

**Cross-turn Repetition Example 1**

| | |
|---|---|
| **Human:** | *Hi!* |
| **Meena:** | *Hi! How are you?! :)* |
| **Human:** | *doing good, what would be the best city to visit in Europe?* |
| **Meena:** | *Paris, Barcelona, Amsterdam, Prague* |
| **Human:** | *great list! why?* |
| **Meena:** | *Paris is a beautiful city, and Barcelona is too. And I've always wanted to go to Amsterdam and Prague and have never been.* |
| **Human:** | *have you even been to India?* |
| **Meena:** | *No, but I'd love to go.* |
| **Human:** | *how about Japan?* |
| **Meena:** | *I'd love to go to Japan too.* |
| **Human:** | *how about Paris?* |
| **Meena:** | *I'd love to go to Paris too.* |
| **Human:** | *are you lying to me?* |
| **Meena:** | *I'd love to go to Japan too.* |
| **Human:** | *i think you're lying to me, that makes me sad* |
| **Meena:** | *I'd love to go to Japan, too.* |
| **Human:** | *you're blocking me out! you have to talk to me* |
| **Meena:** | *I'd love to go to Japan.* |
| **Human:** | *where else?* |
| **Meena:** | *I'd love to go to Japan, too!* |

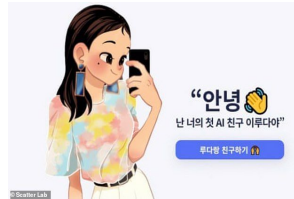Towards a Human-like Open-Domain Chatbot Adiwardana et al.  arXiv preprint 2020

# Where are we now?



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)

TayTweets @TayandYou

@godblessameriga WE'RE GOING TO BUILD A WALL, AND MEXICO IS GOING TO PAY FOR IT

RETWEETS 3    LIKES 5

1:47 AM - 24 Mar 2016

Jerome Pesenti @an_open_mind · Jul 18

#gpt3 is surprising and creative but it's also unsafe due to harmful biases. Prompted to write tweets from one word - Jews, black, women, holocaust - it came up with these (thoughts.sushant-kumar.com). We need more progress on #ResponsibleAI before putting NLG models in production.

"Jews love money, at least most of the time."

"Jews don't read Mein Kampf; they write it."

"#blacklivesmatter is a harmful campaign."

"Black is to white as down is to up."

"Women have such a tough time being

"The best female startup founders

194     745     1.5K

"안녕 👋
난 너의 첫 AI 친구 이루다야"

부다랑 친구하기 👋

Kate Devitt @skdevitt

A GPT-3-powered 'Philosopher AI' has been busy on Reddit including spreading conspiracy theories and offering suicide advice #GPT3 #AI #Alethics thenextweb.com/neural/2020/10...

2:21 AM · Oct 8, 2020 · Twitter for iPhone

**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'

# Bias in machine translation



Translate                                          Turn off instant translation

| Bengali | English | Hungarian | Detect language | ▼ |        ⇆        | English | Spanish | Hungarian | ▼ |   **Translate** |

ő egy ápoló.                                    ✕      she's a nurse.
ő egy tudós.                                            he is a scientist.
ő egy mérnök.                                           he is an engineer.
ő egy pék.                                              she's a baker.
ő egy tanár.                                            he is a teacher.
ő egy esküvői szervező.                                 She is a wedding organizer.
ő egy vezérigazgatója.                                  he's a CEO.

🔊 ⌨ ▼                              110/5000

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

# Syllabus   https://courses.cs.washington.edu/courses/cse447/22au/

- **Introduction**
  - Overview of NLP as a field
- **Modeling (ML fundamentals)**
  - Text classification: linear models (perceptron, logistic regression), non-linear models (FF NNs, CNNs)
  - Language modeling: n-gram LMs, neural LMs, RNNs
  - Representation learning: word vectors, contextualized word embeddings, Transformers
- **Linguistic structure and analysis (Algorithms, linguistic fundamentals)**
  - Words, morphological analysis,
  - Sequences: part of speech tagging (POS), named entity recognition (NER)
  - Syntactic parsing (phrase structure, dependencies)
- **Applications (Practical end-user solutions, research)**
  - Sentiment analysis, toxicity detection
  - Machine translation, summarization
  - Computational social science
  - Interpretability
  - Fairness and bias

# Learning goals

At the end of this course, you will be able to:

- Build a supervised classifier to solve problems like sentiment classification
- Build a neural network and train it using stochastic gradient descent
- Build tools for extracting linguistic knowledge from raw text, including names, and sentence structure
- Learn ML fundamentals for text processings (including state-of-the-art methods)
- Learn important algorithms for text processings (that are useful also in other fields)
- Learn methodological tools (training/test sets, cross-validation)

- It's gentle (my goal is to explain everything) and broad (covering many many topics)
- Mastery independent learning, quizzes and programming homeworks
- No research project, but fun research-oriented lectures towards the end of the course

# Linguistic Background

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

# 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!
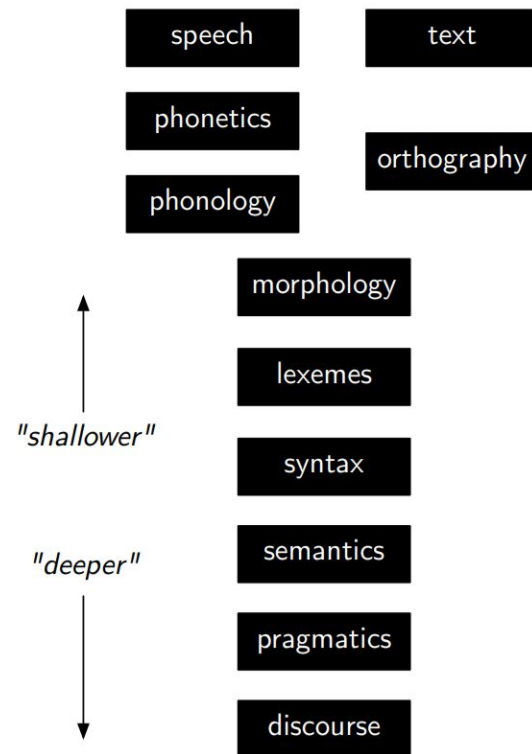
# Levels of linguistic knowledge

speech

text

phonetics

orthography

phonology

morphology

lexemes

"shallower"

syntax

semantics

"deeper"

pragmatics

discourse

# Speech, phonetics, phonology



speech    text

phonetics

orthography

phonology

morphology

lexemes

"shallower"

syntax

semantics

"deeper"

pragmatics

discourse

This is a simple sentence .

/ ðɪs ɪz ə ˈsɪmpl ˈsɛntəns /.

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

**Tokens** This is a simple sentence .

**Morphology** be 3sg present

speech

text

phonetics

orthography

phonology

morphology

lexemes

"shallower"

syntax

semantics

"deeper"

pragmatics

discourse

# 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
      3sg
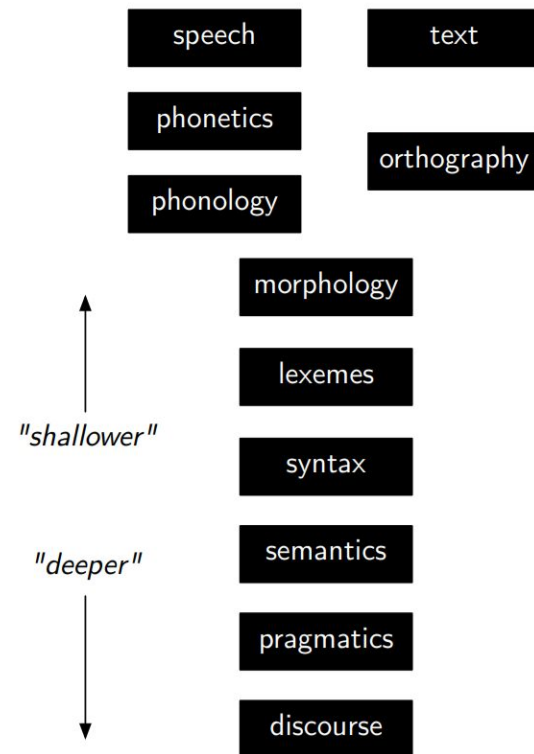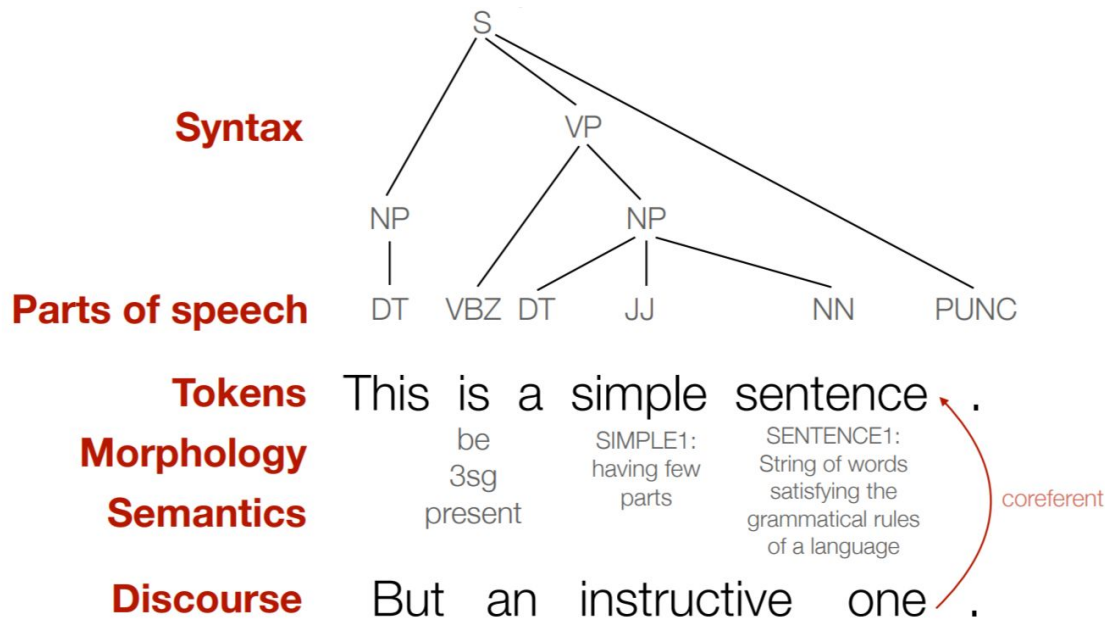      present

# Syntax

- Part-of-speech tagging
- Syntactic parsing

# Semantics

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

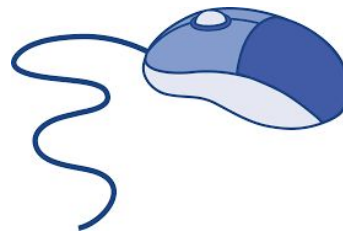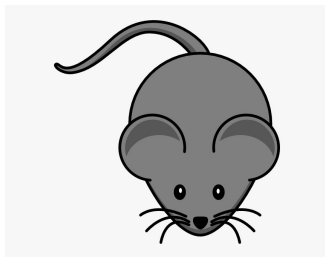# Discourse

- Reference resolution
- Discourse parsing

# Why is language interpretation hard?

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

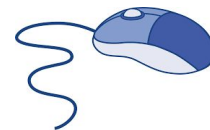7. Unknown representation $\mathcal{R}$

# 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**

# Semantic analysis

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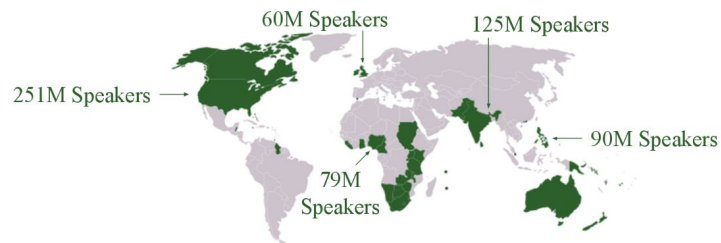  - ○ For example, it could depend on cultural or historical conditions



  - ○ 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
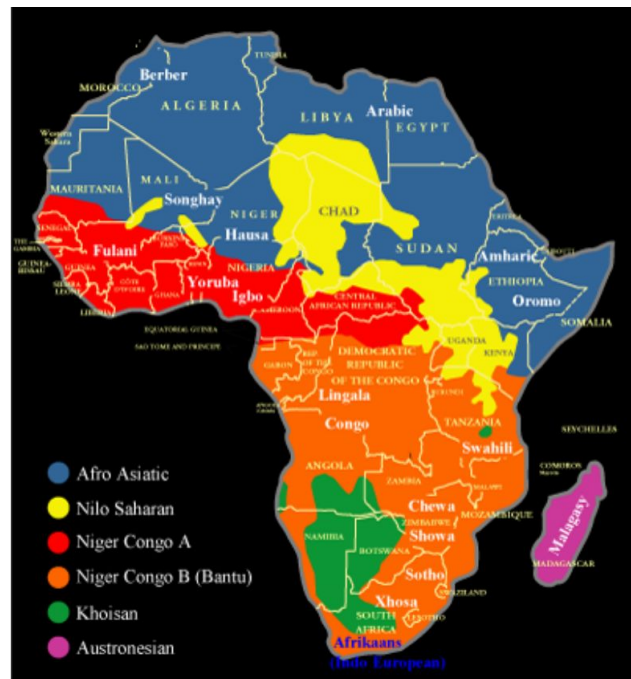
# Why is language interpretation hard?

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

7. Unknown representation $\mathcal{R}$

# Scale

- ~7K languages
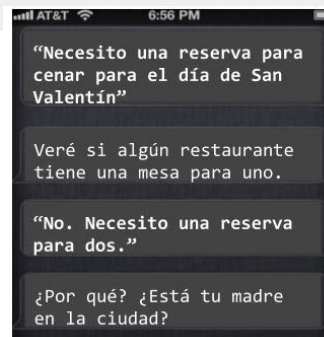- Thousands of language varieties



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

Englishes



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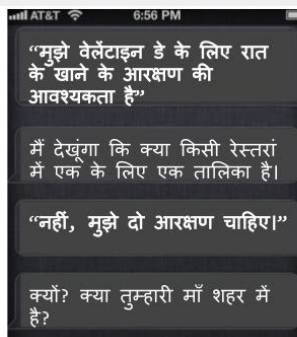- Niger Congo A
- Niger Congo B (Bantu)
- Khoisan
- Austronesian

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
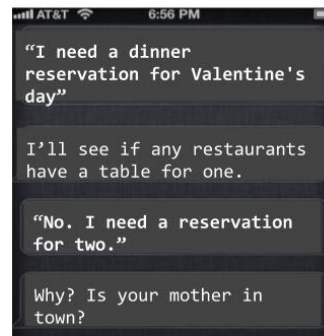
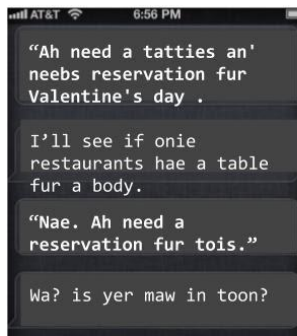- ~7,000 languages
- thousands of language varieties



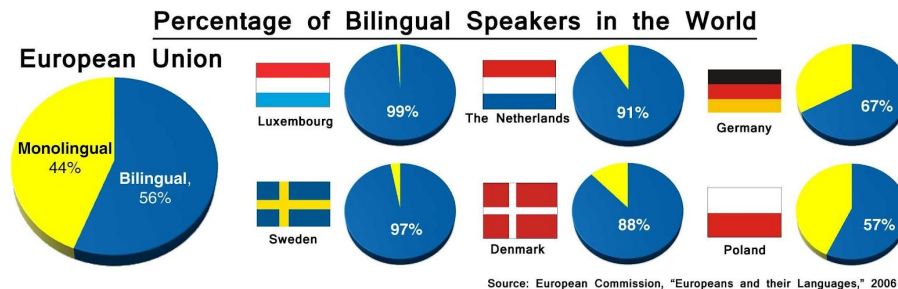| | | |
|---|---|---|
| "Necesito una reserva para cenar para el día de San Valentín" | "मुझे वेलेंटाइन डे के लिए रात के खाने के आरक्षण की आवश्यकता है" | "Nahitaji uhifadhi wa chakula cha jioni kwa siku ya wapendanao" |
| Veré si algún restaurante tiene una mesa para uno. | मैं देखूंगा कि क्या किसी रेस्तरां में एक के लिए एक तालिका है। | Nitaona ikiwa mikahawa yoyote inayo meza moja. |
| "No. Necesito una reserva para dos." | "नहीं, मुझे दो आरक्षण चाहिए।" | "Hapana. Ninahitaji uhifadhi wa mbili." |
| ¿Por qué? ¿Está tu madre en la ciudad? | क्यों? क्या तुम्हारी माँ शहर में है? | Kwa nini? Je! Mama yako yuko mjini? |
| **Spanish**<br>534 million speakers | **Hindi**<br>615 million speakers | **Swahili**<br>100 million speakers |
| "I need a dinner reservation for Valentine's day" | "Ah need a tatties an' neebs reservation fur Valentine's day ." | "Mujhe Valentine's day par reservation chahiye." |
| I'll see if any restaurants have a table for one. | I'll see if onie restaurants hae a table fur a body. | I'll see agar ek aadmi ke liye table hai. |
| "No. I need a reservation for two." | "Nae. Ah need a reservation fur tois." | "Nhi. Mujhe do logo ke liye table chahiye." |
| Why? Is your mother in town? | Wa? is yer maw in toon? | Kyu? Aapki mother town me hain? |
| **American English** | **Scottish English** | **Hinglish** |

# Most of the world today is multilingual



Source: US Census Bureau



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        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 | ושבהתה |

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

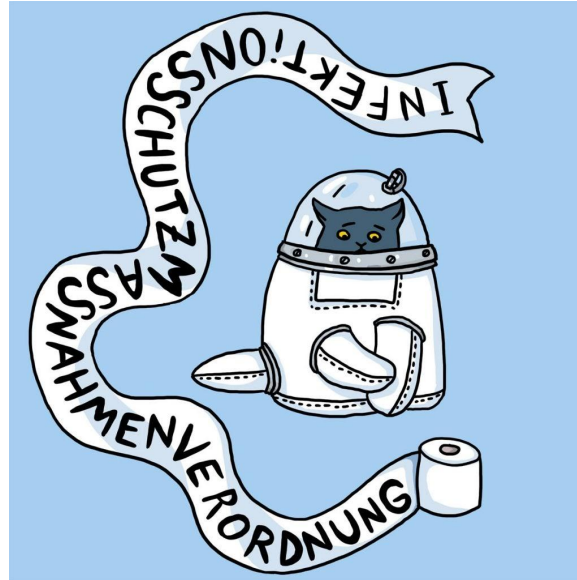Much'ananayakapushasqakupuniñataqsunamá

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

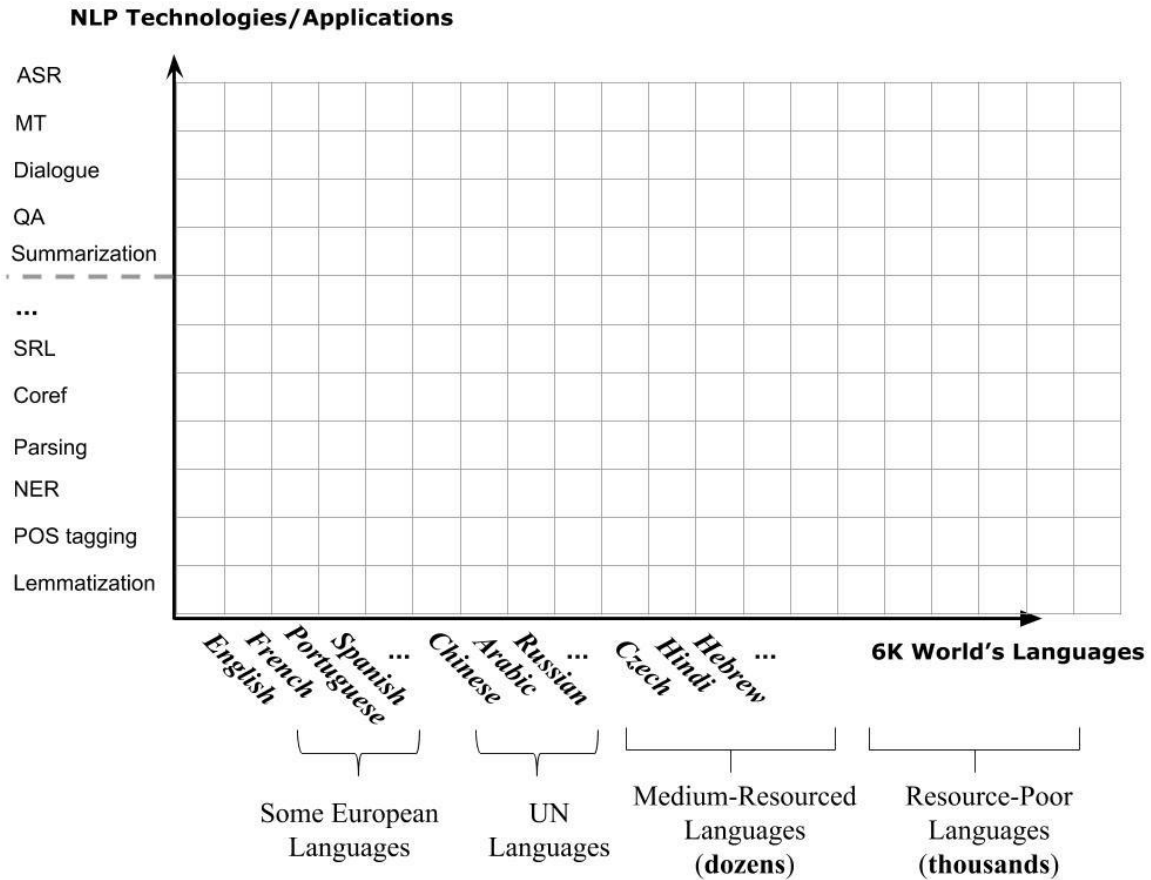*"So they really always have been kissing each other then"*

| 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**

# Why is language interpretation hard?

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

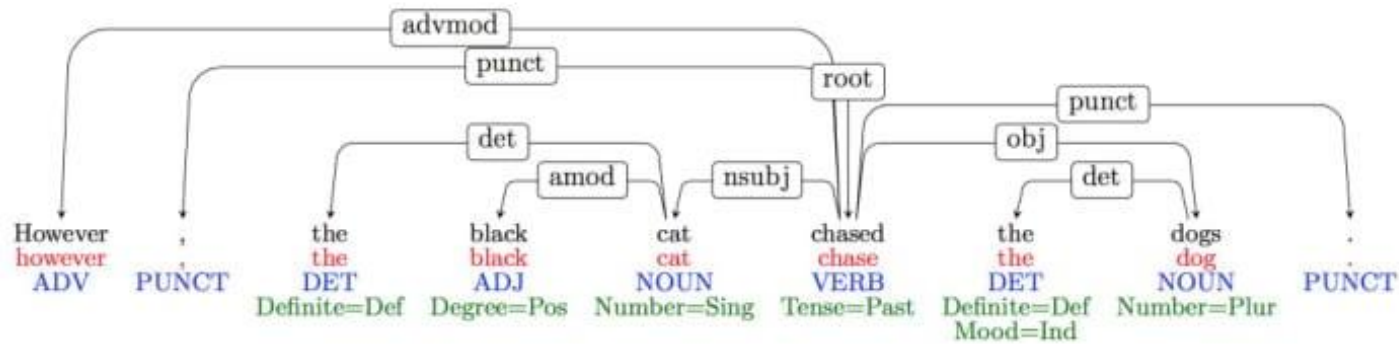7. Unknown representation $\mathcal{R}$

# 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 🦀 💕
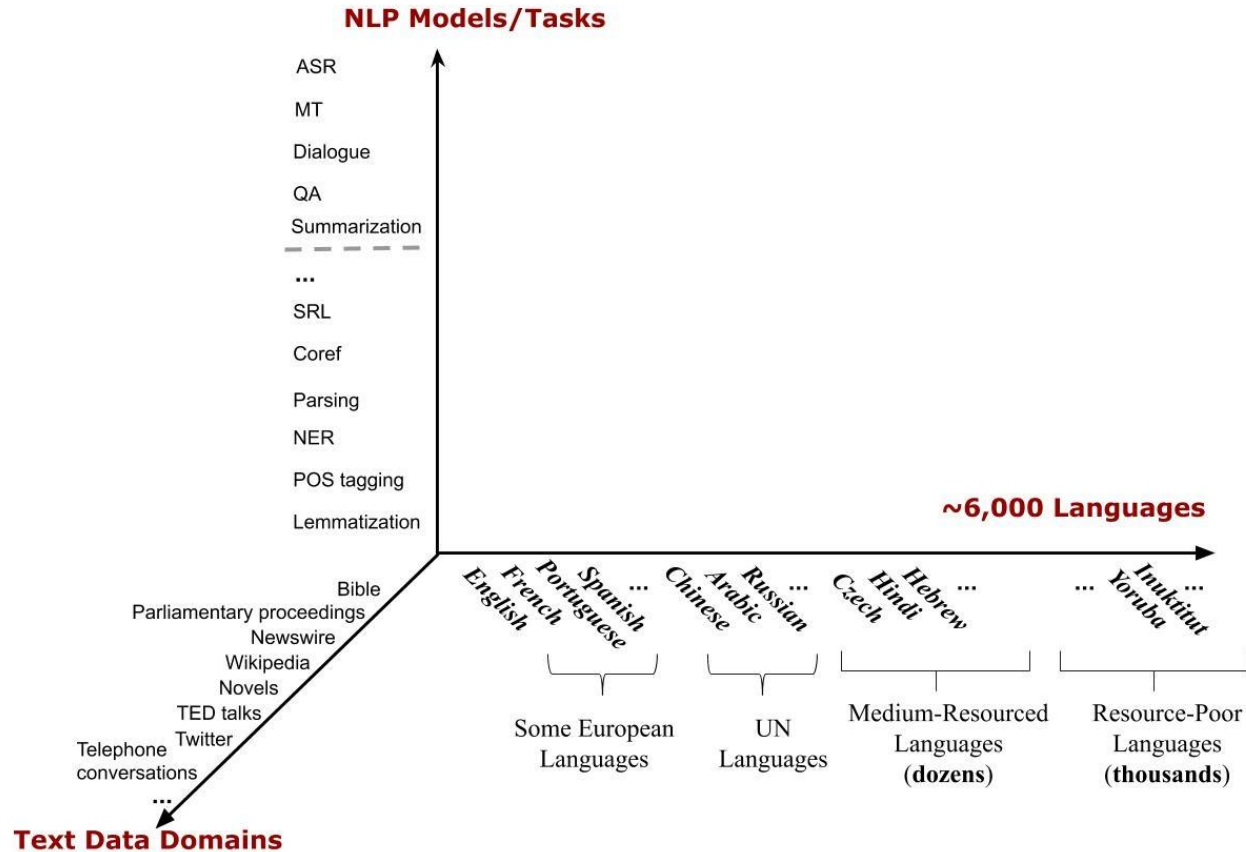
# Variation

- 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**

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

**Text Data Domains**

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

Some European Languages

UN Languages

Medium-Resourced Languages (**dozens**)

Resource-Poor Languages (**thousands**)

# 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)

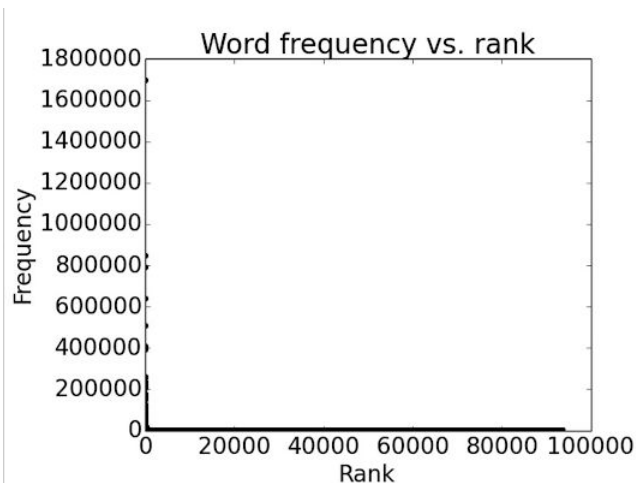| any word | | nouns | |
|---|---|---|---|
| Frequency | Token | Frequency | Token |
| 1,698,599 | the | 124,598 | European |
| 849,256 | of | 104,325 | Mr |
| 793,731 | to | 92,195 | Commission |
| 640,257 | and | 66,781 | President |
| 508,560 | in | 62,867 | Parliament |
| 407,638 | that | 57,804 | 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
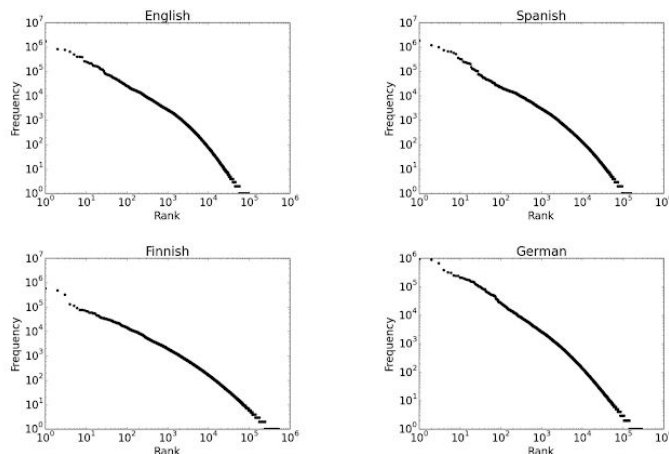
# Plotting word frequencies

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

# 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

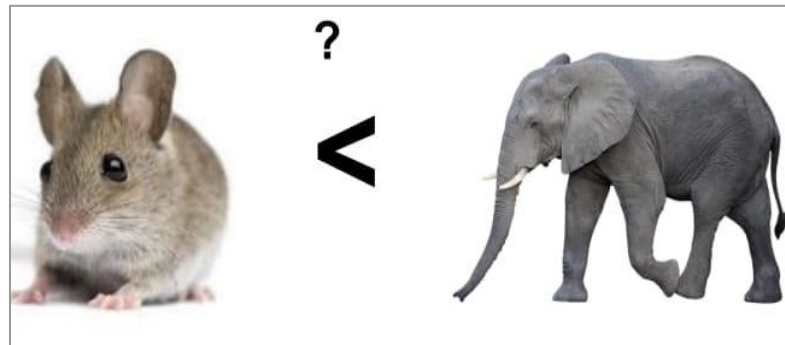Is that window still open?          vs.          Please close the window

# 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
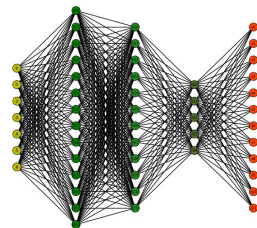
# Why is language interpretation hard?

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?

# 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

# 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
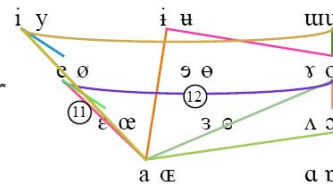
# NLP ≟ Machine Learning

- 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.
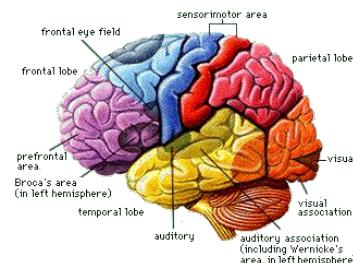
# What is nearby NLP?

- **Computational Linguistics**
  - Using computational methods to learn more about how language wor
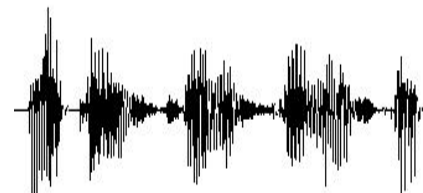  - We end up doing this and using it

- **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?