

# Natural Language Processing

AI51701/CSE71001

Lecture 1

08/29/2023

Instructor: Taehwan Kim

# Lecture Plan

- ❑ Introduction to the course
- ❑ Examples of Natural Language Processing
- ❑ Word

# Course Information

- ❑ Instructor: Taehwan Kim ([taehwankim@unist.ac.kr](mailto:taehwankim@unist.ac.kr), Machine Learning, Vision, and Language Lab,  
<https://sites.google.com/view/mvllab>)
  - Office hours: by appointments
  
- ❑ We have excellent TAs:
  - Jaeyeon Bae ([qowodussla@unist.ac.kr](mailto:qowodussla@unist.ac.kr))
  - Seonghee Han ([seonghee@unist.ac.kr](mailto:seonghee@unist.ac.kr))

# Course Information

## Will be taught offline

- Time: Tue/Thu at 4:00-5:15pm
- Place: Bldg. 106, T205
- Will use **mobile attendance system**
  - If you miss more than  $\frac{1}{4}$  of all classes, you will be automatically failed in this course (university policy).
    - Exceptions to count as being absent: military training, family occasion (wedding & funeral), medical reasons, and attending events approved by university president (university policy, all require document proof)
  - > 10 min. late: counted as being late
  - > 30 min. late: absence
  - 2 lateness = 1 absence

# Grading

- ❑ Assignments (total 50%)
  - Around 4 assignments including problem solving and (mostly) programming in Python.
- ❑ Attendance Quizzes (total 10%)
- ❑ Final Project (40%)

# Grading

## Attendance Quizzes (total 10%)

- Will be available via *Quizzes & Exams* in BlackBoard (the left panel), so bring your laptops/tablets.
- Will be open for five minutes during the lecture.
- Two lowest scores will be excluded when grading and use this for emergencies.

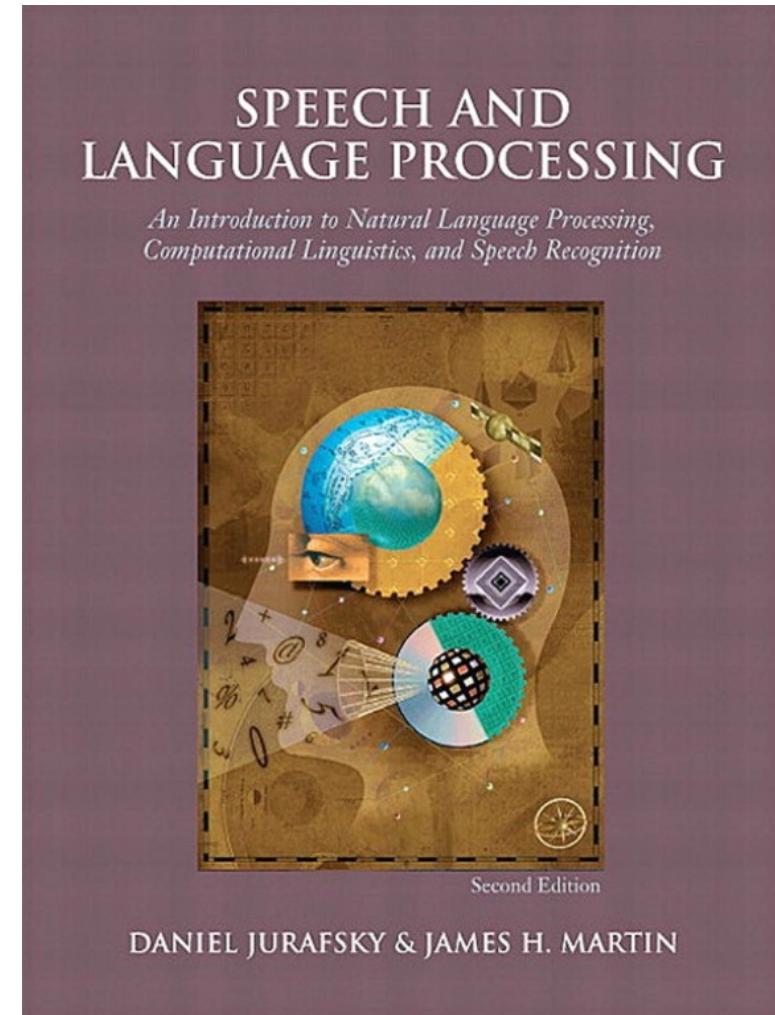
# Grading

## □ Final Project (40%)

- Proposal (5%): presentation on Oct. 10 & Oct. 12 (tentative), feedback will be given.
- Final presentation (5%): during the final week
- Final report (30%): by the end of last week of the semester
- You will choose your project topics.
- 3 people team project.
- More detail will be provided soon.

# Optional Textbooks

- Dan Jurafsky and James H. Martin,  
"Speech and Language Processing, 2nd Edition",  
Prentice Hall, 2009.
  - Third edition draft is available at  
<https://web.stanford.edu/~jurafsky/slp3/>
  
- Goldberg. *Neural Network Methods for Natural Language Processing*.
  - Earlier draft (from 2015) available online



# Prerequisites

- ❑ Proficiency in Python
  - All class assignments will be in Python
  - For deep learning, we use PyTorch (a refresher may be given later)
- ❑ Multivariate Calculus, Linear Algebra
- ❑ Basic Probability and Statistics
- ❑ Fundamentals of Machine Learning
  - Loss functions, Taking simple derivatives, Performing optimization with gradient descent

# Using BlackBoard

- ❑ URL: <https://blackboard.unist.ac.kr/>
- ❑ Lecture materials
  - Lecture slides will be posted before the lecture
- ❑ Discussions
  - You can post questions related to the course such as assignments.
  - The instructor and TAs will try to respond soon.
  - Other students are welcome to provide answers as well.
- ❑ Assignments
  - You can check the posted assignments and may submit your assignment here

# Lateness Policy

- ❑ You may submit your assignments up to three days after the deadline:
  - You will get 10% penalty for each late day.
  - After three late days, no late submission will be accepted.
  
- ❑ Appeals to your graded assignments
  - Should be done via documents.
  - Detail will be announced later.

# Cheating and Collaboration

□ What is your goal in this course?

- You want to learn and cheating will not help.

□ Your assignment submissions should be your OWN work.

- You are welcome to discuss assignments with others in the course, but written solutions and code must be written individually.
- Do not search or copy solutions in textbook/internet.
- We will check plagiarism of your code by using automatic tools.
- If cheating happens, you will get a penalty
  - Minimum 0 point and -10% in your final grade, up to failure in this course and reporting to the dean

# What is Natural Language Processing?

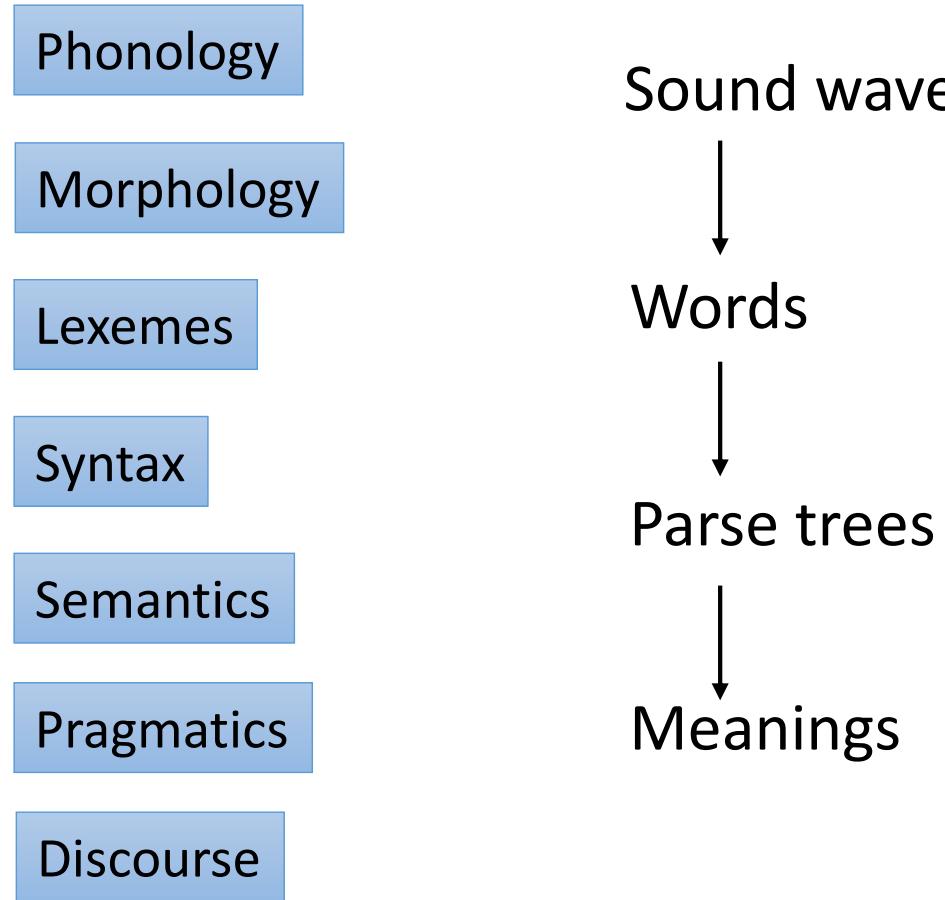
- ❑ Automating the analysis, generation of human (“natural”) language
  - Analysis (or “understanding” or “processing” ...)
  - Generation
  
- ❑ An experimental computer science research area that includes problems and solutions pertaining to the understanding of human language

# What does it mean to understand a language?

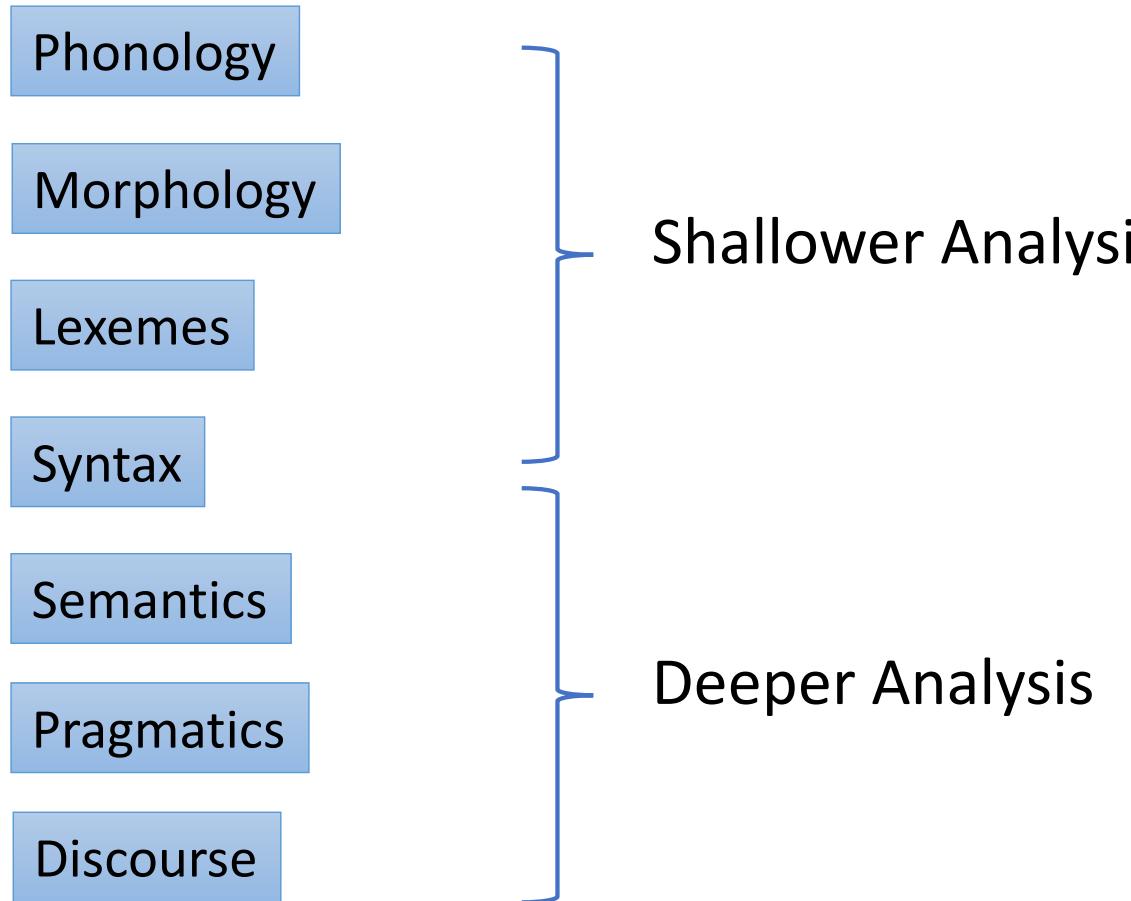


- "Stop"
- "Turn it up"
- "Volume level 6"
- "Repeat that"
- "What can you do?"
- "Play some music"
- "Play music by [artist]"
- "Play dance music on YouTube"
- "Play KEXP radio on TuneIn"
- "Play the latest episode of Radiolab"
- "Pause"
- "Next song"
  
- "When's my first appointment tomorrow?"
- "Wake me up at 6am tomorrow"
- "Tell me about my day"
- "How long will it take to get to work?"
- "What's the weather today?"

# What does it mean to understand a language?



# What does it mean to understand a language?



# Syntax, Semantics, Pragmatics

□ **Syntax** concerns the proper ordering of words and its affect on meaning.

- The dog bit the boy.
- The boy bit the dog.
- Bit boy dog the the.

□ **Semantics** concerns the (literal) meaning of words, phrases, and sentences.

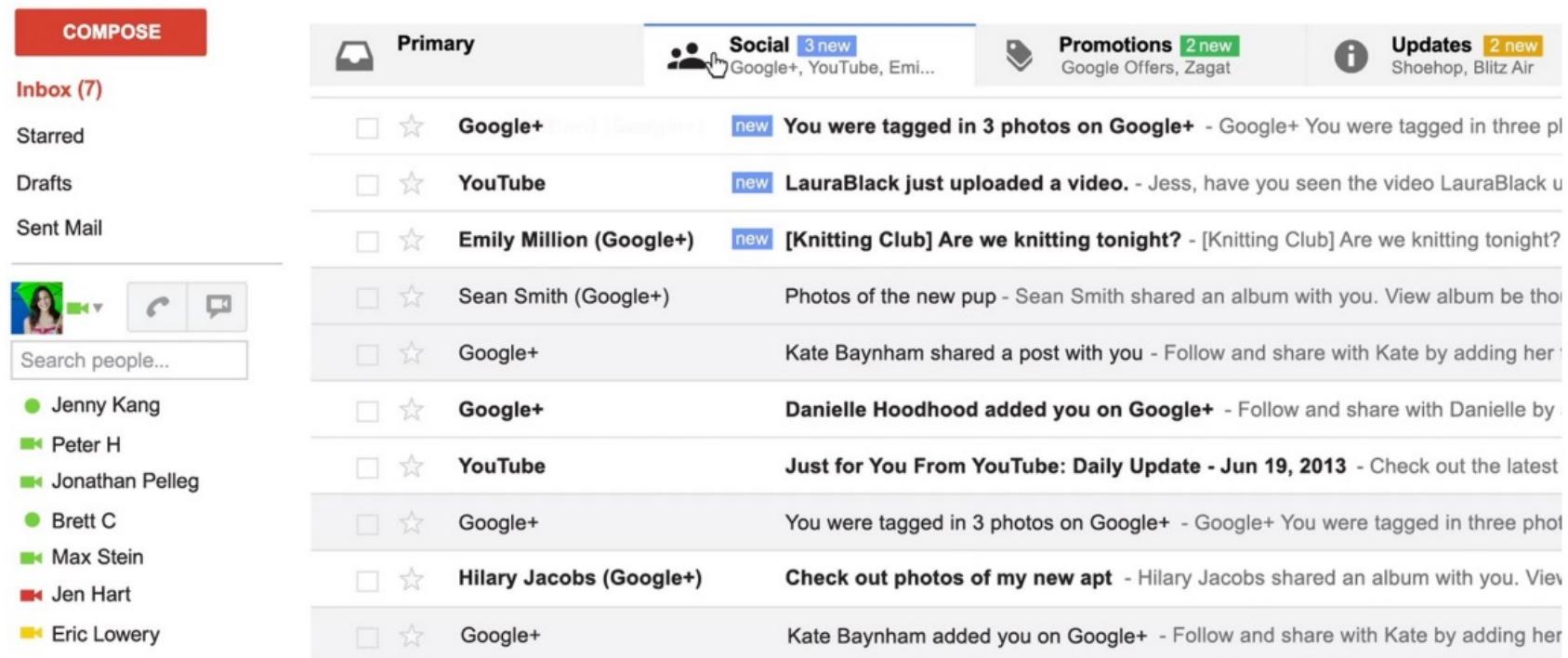
- “plant” as a photosynthetic organism
- “plant” as a manufacturing facility
- “plant” as the act of sowing

# Syntax, Semantics, Pragmatics

- **Pragmatics** concerns the overall communicative and social context and its effect on interpretation.
  - Honest or dishonest?
  - Context: Kyle and Ellen would like to see a movie. Kyle has \$20 in his pocket. Tickets cost \$8 each.
  - Kyle: “I have \$8.”

# Examples of Natural Language Processing

## ❑ Text Classification



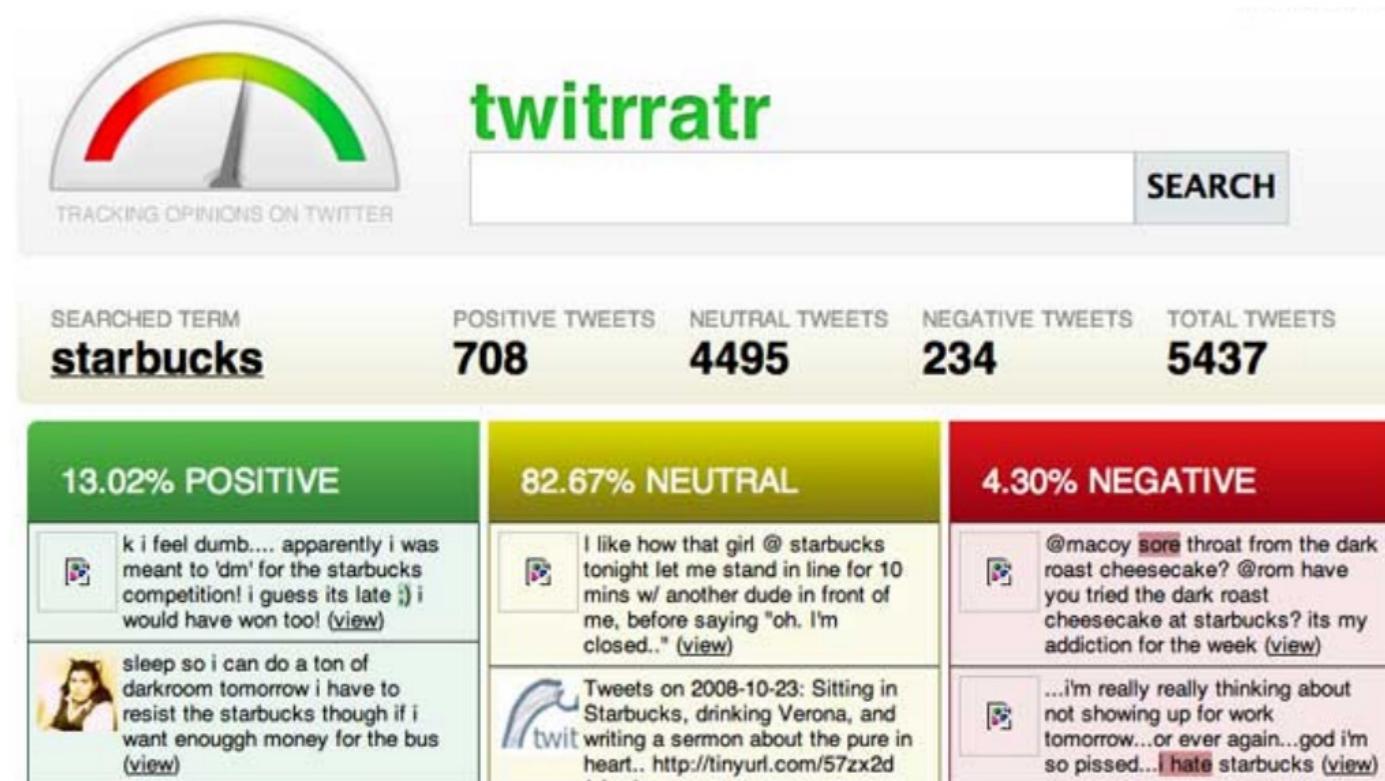
The screenshot shows a Gmail inbox with the following message list:

- Primary**:
  - Google+ You were tagged in 3 photos on Google+ - Google+ You were tagged in three pl
  - YouTube new LauraBlack just uploaded a video. - Jess, have you seen the video LauraBlack u
  - Emily Million (Google+) new [Knitting Club] Are we knitting tonight? - [Knitting Club] Are we knitting tonight?
  - Sean Smith (Google+) Photos of the new pup - Sean Smith shared an album with you. View album be tho
  - Google+ Kate Baynham shared a post with you - Follow and share with Kate by adding her to your
  - Google+ Danielle Hoodhood added you on Google+ - Follow and share with Danielle by adding her to
  - YouTube Just for You From YouTube: Daily Update - Jun 19, 2013 - Check out the latest
  - Google+ You were tagged in 3 photos on Google+ - Google+ You were tagged in three photo
  - Hilary Jacobs (Google+) Check out photos of my new apt - Hilary Jacobs shared an album with you. View
  - Google+ Kate Baynham added you on Google+ - Follow and share with Kate by adding her to
- Social**: 3 new Google+, YouTube, Emi...
- Promotions**: 2 new Google Offers, Zagat
- Updates**: 2 new Shoehop, Blitz Air

On the left, there's a sidebar with a **COMPOSE** button, an **Inbox (7)** section listing Starred, Drafts, and Sent Mail, and a search bar with a "Search people..." placeholder. Below the search bar is a list of contacts with their names and profile icons.

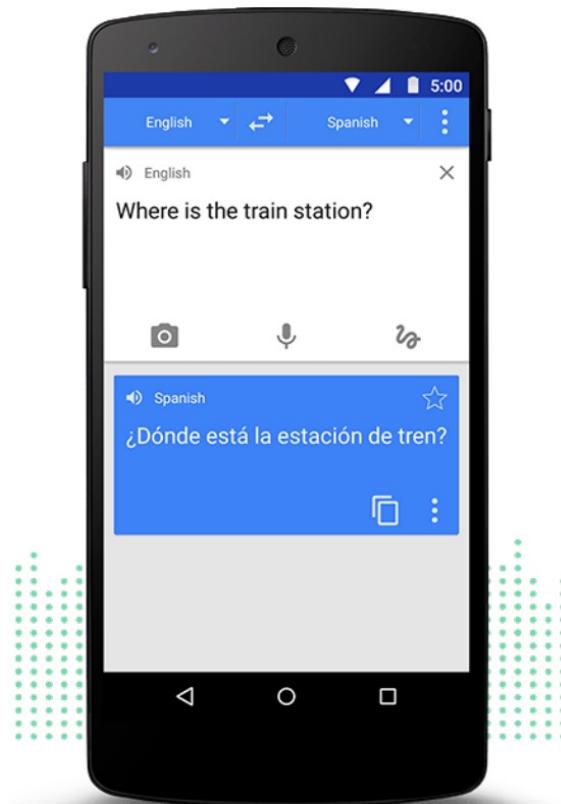
# Examples of Natural Language Processing

## ❑ Sentiment Analysis



# Examples of Natural Language Processing

## ❑ Machine Translation



# Examples of Natural Language Processing

## ❑ Question Answering

“Alexa, who was President when Barack Obama was nine?”

“Alexa, how’s my commute?”

“Alexa, what’s the weather?”

“Alexa, did the 49ers win?”



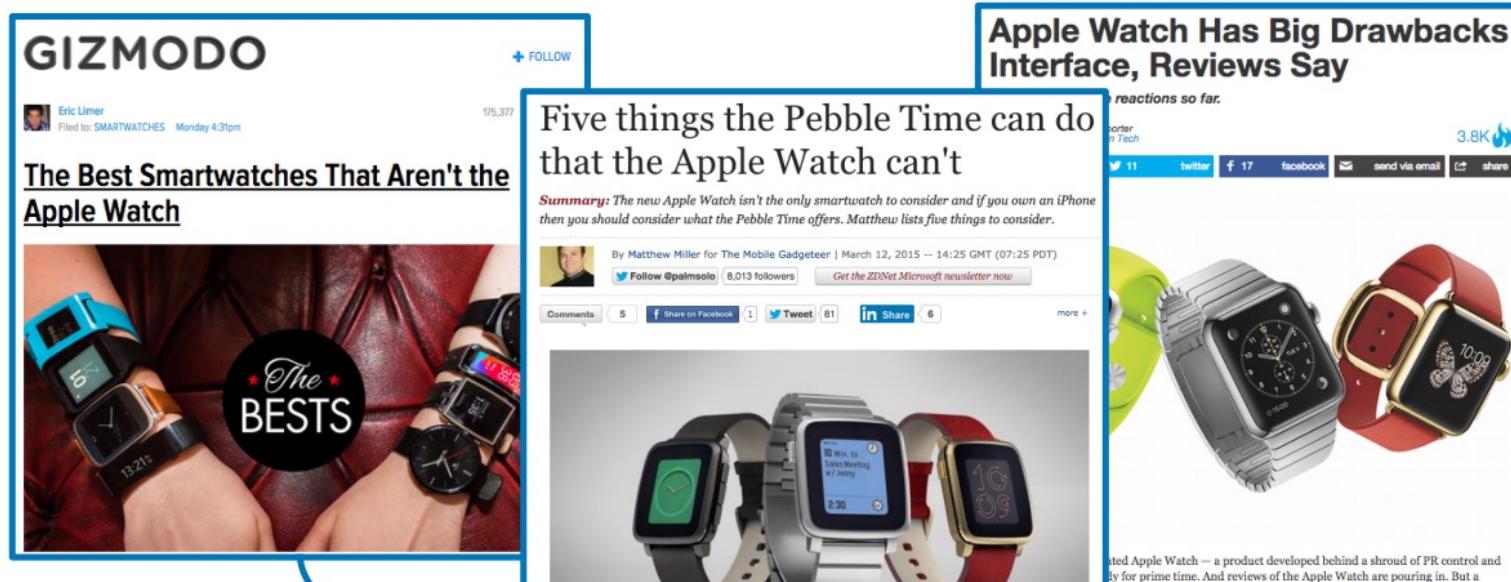
# Examples of Natural Language Processing

## ❑ Dialogue Systems



# Examples of Natural Language Processing

## ❑ Summarization



The Apple Watch has drawbacks. There are other smartwatches that offer more capabilities.

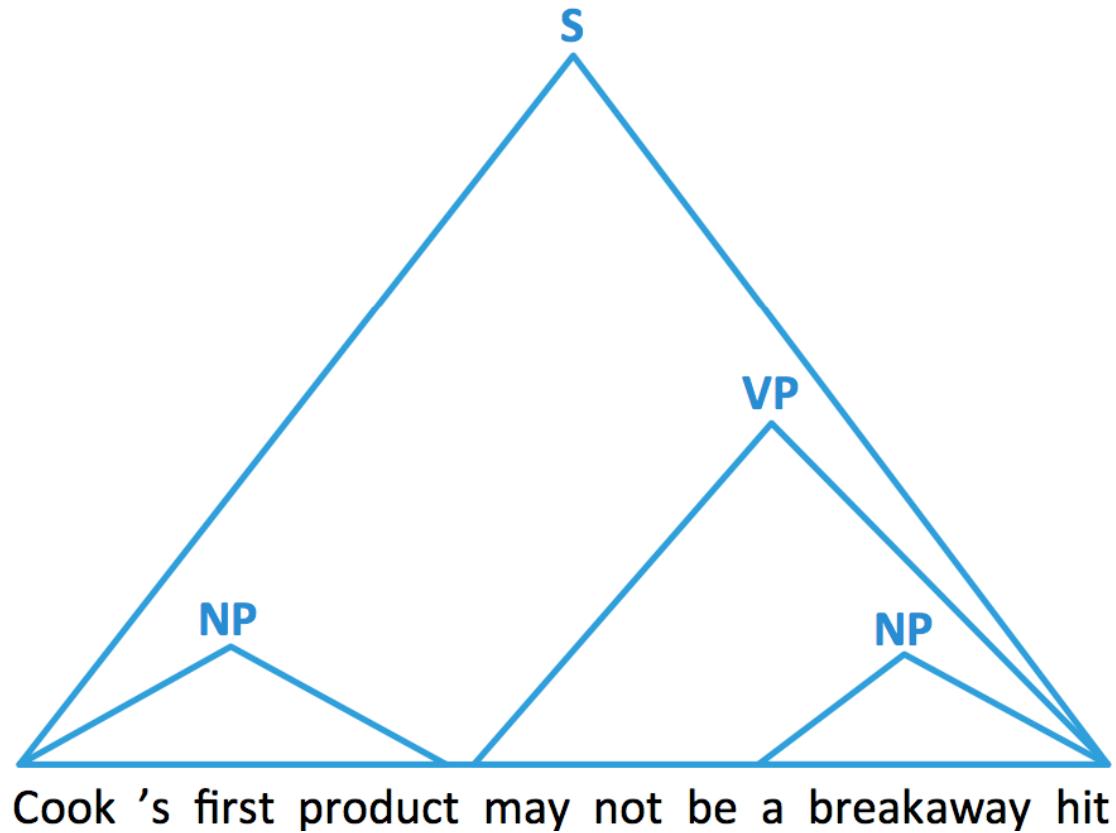
# Examples of Natural Language Processing

## ❑ Part-of-Speech Tagging

determiner	verb (past)	prep.	proper noun	proper noun	poss.	adj.	noun
Some	questioned	if	Tim	Cook	's	first	product
modal	verb	det.	adjective	noun	prep.	proper noun	punc.
would	be	a	breakaway	hit	for	Apple	.

# Examples of Natural Language Processing

- ❑ Syntactic Parsing



# Examples of Natural Language Processing

# ☐ Named Entity Recognition

Some questioned if Tim Cook's first product would be a breakaway hit for Apple.

# Examples of Natural Language Processing

## ❑ Coreference Resolution

Some questioned if Tim Cook's first product would be a breakaway hit for Apple.

It's the company's first new device since he became CEO.

# Examples of Natural Language Processing

## ❑ Reading Comprehension

Once there was a boy named **Fritz** who loved to draw. He drew everything. In the morning, **he drew a picture of his cereal with milk**. His papa said, “Don’t draw your cereal. Eat it!”

After school, Fritz drew a picture of his bicycle. His uncle said, “Don’t draw your bicycle. Ride it!”

...

What did Fritz draw first?

- A) the toothpaste
- B) his mama
- C) cereal and milk**
- D) his bicycle

# Examples of Natural Language Processing

## ❑ Sentence Similarity

Input	Output
Other ways are needed. We must find other ways.	4.4
Pakistan bomb victims' families end protest Pakistan bomb victims to be buried after protest ends	2.6
I absolutely do believe there was an iceberg in those waters. I don't believe there was any iceberg at all anywhere near the Titanic.	1.2

# Examples of Natural Language Processing

## ❑ Word Prediction

he bent down and searched the large container, trying to find anything else hidden in it other than the \_\_\_\_\_

# Other Language Technologies

- ❑ Speech processing
- ❑ Information retrieval / web search
- ❑ Knowledge representation / reasoning
- ❑ Multimodal: Image-to-text (image/video captioning, VQA)

# Related Fields

- ❑ Artificial Intelligence
- ❑ Machine Learning
- ❑ Linguistics
- ❑ Cognitive science
- ❑ Logic
- ❑ Data science
- ❑ Political science
- ❑ Education
- ❑ Economics
- ❑ ...many more

# Computational Linguistics vs. Natural Language Processing

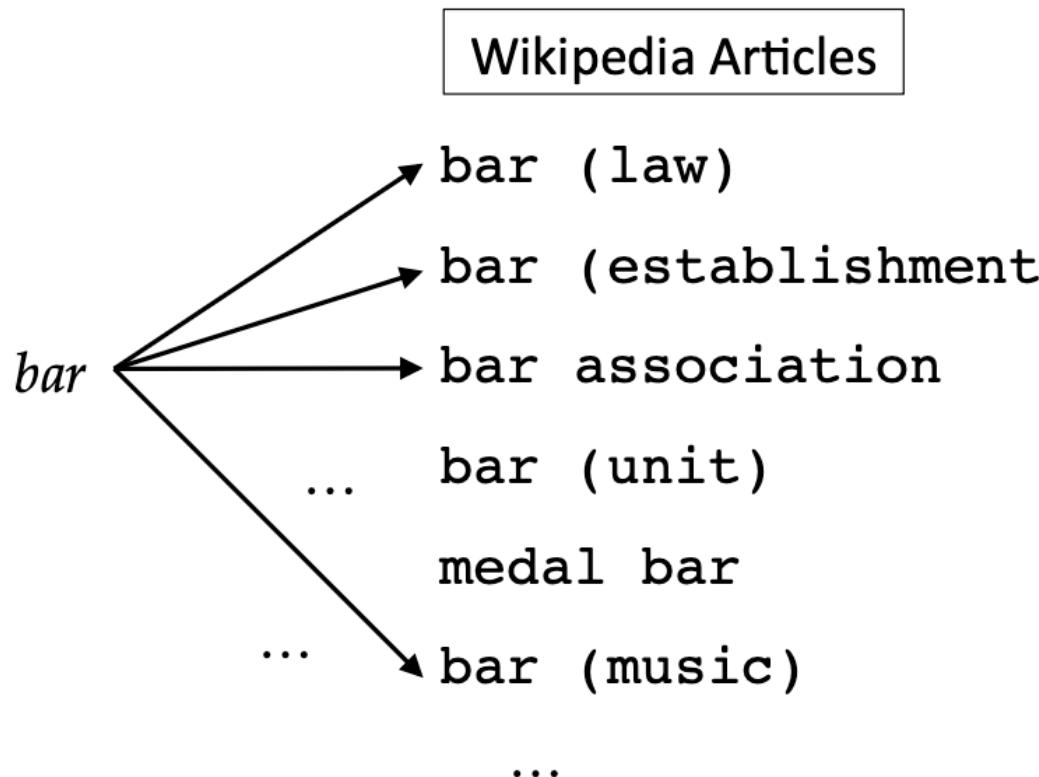
- ❑ Many people think of the two terms as **synonyms**
- ❑ Computational linguistics is about **linguistics**
  - Computational linguistics is more inclusive; more likely to include sociolinguistics, cognitive linguistics, and computational social science
- ❑ Natural language processing is about **processing**
  - NLP is more likely to use machine learning and involve engineering / system-building

# Why is NLP hard?

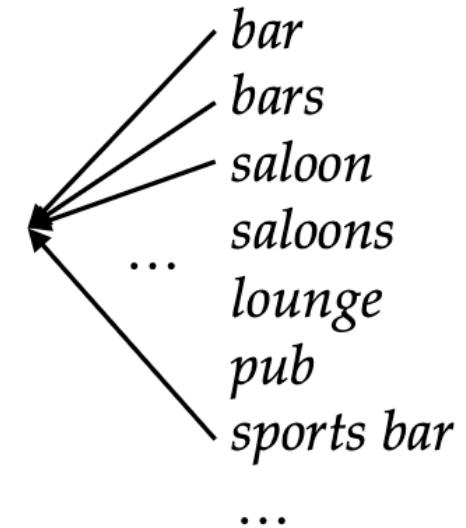
- Ambiguity and variability of linguistic expression:
  - **Variability:** many forms can mean the same thing
  - **Ambiguity:** one form can mean many things
- Many different kinds of variability and ambiguity
- Each NLP task must address distinct kinds

# Examples: Hyperlinks in Wikipedia

Ambiguity



Variability



# Why is Language Ambiguous?

- ❑ Having a unique linguistic expression for every possible conceptualization that could be conveyed would make language overly complex and linguistic expressions unnecessarily long.
- ❑ Allowing resolvable ambiguity permits shorter linguistic expressions, i.e. data compression.
- ❑ Language relies on people's ability to use their knowledge and inference abilities to properly resolve ambiguities.
- ❑ Infrequently, disambiguation fails, i.e. the compression is lossy.

# More difficulties: Non-standard language

Great job @justinbieber! Were SOO PROUD of  
what youve accomplished! U taught us 2  
#neversaynever & you yourself should never give  
up either ❤

## ❑ And neologisms:

- Unfriend
- Retweet
- bromance

# Learning-based NLP



## Option 1: Write Rules

- Rules for translating **much** or **many** into Russian:

**if** preceding word is *how* **return** *skol'ko*

**else if** preceding word is *as* **return** *stol'ko zhe*

**else if** word is *much*

**if** preceding word is **very** **return** nil

**else if** following word is a noun **return mnogo**

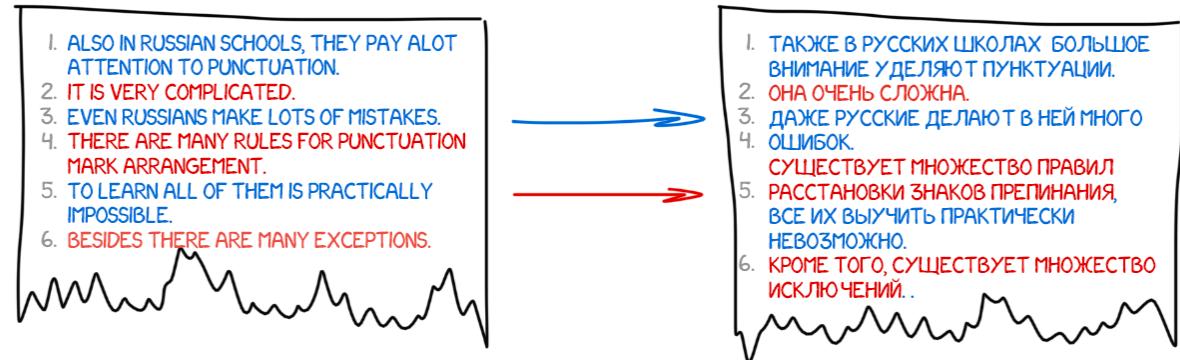
**else** (word is many)

**if** preceding word is a preposition and following word is noun **return** *mnogii*

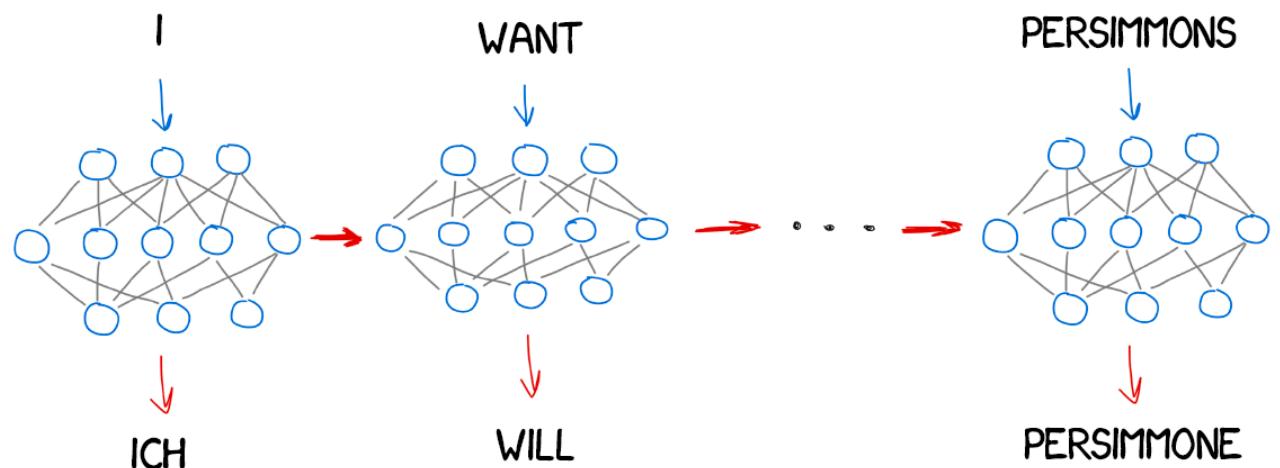
**else return** *mnogo*

# Learning-based NLP

PARALLEL CORPUS



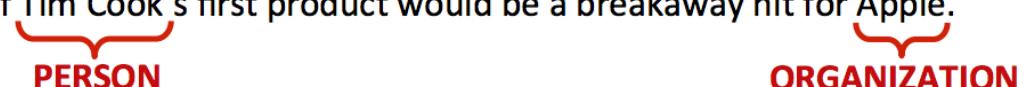
## □ Option 2: Learn from data



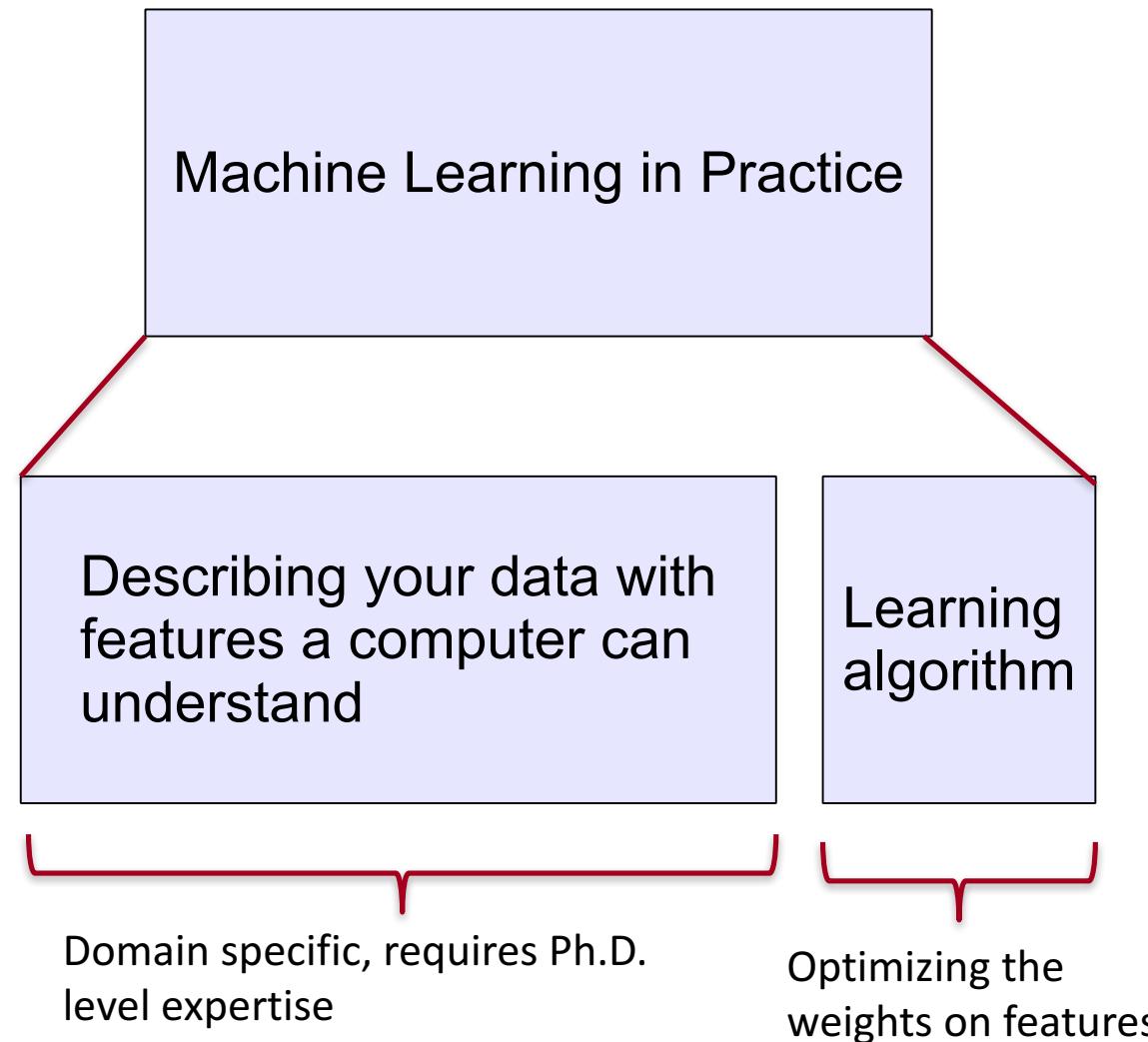
# What is Deep Learning?

- ❑ Deep learning is a subfield of machine learning.
- ❑ Most machine learning methods work well because of human-designed representations/input features.
- ❑ Machine learning becomes just optimizing feature weights to make a good prediction.

Feature	NER
Current Word	✓
Previous Word	✓
Next Word	✓
Current Word Character n-gram	all
Current POS Tag	✓
Surrounding POS Tag Sequence	✓
Current Word Shape	✓
Surrounding Word Shape Sequence	✓
Presence of Word in Left Window	size 4
Presence of Word in Right Window	size 4

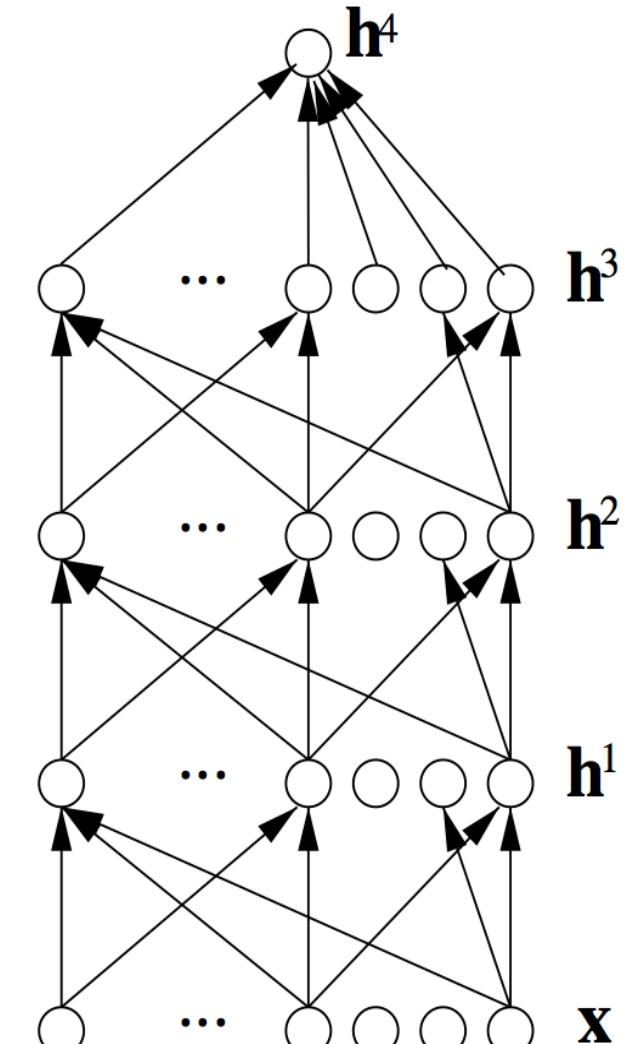
Some questioned if Tim Cook's first product would be a breakaway hit for Apple.  


# Traditional Machine Learning



# Deep Learning

- ❑ Representation learning attempts to automatically learn good features or representations
- ❑ Deep learning algorithms attempt to learn (multiple levels of) representations (here:  $h^1, h^2, h^3$ ) and output ( $h^4$ )
- ❑ From “raw” inputs  $x$  (e.g. sound, pixels, characters, or words)



# Why Deep Learning?

- Manually designed features are often over-specified, incomplete and take a long time to design and validate
- **Learned Features** are easy to adapt, fast to learn
- Deep learning provides a very flexible, learnable framework for **representing** world, visual and linguistic information.
- Deep learning can learn in **unsupervised** (from raw text) and **supervised** (with specific labels like positive/negative) settings

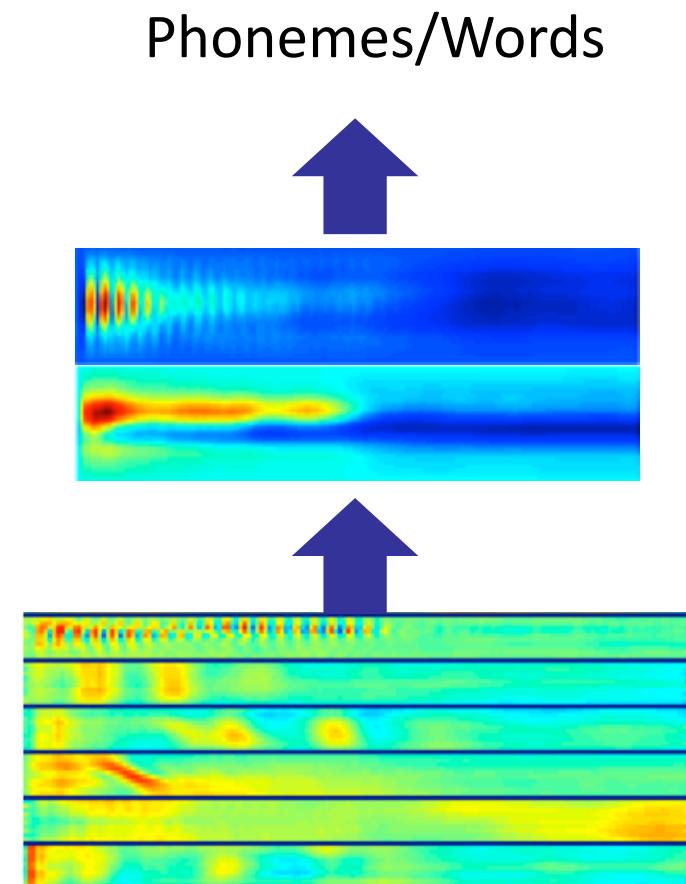
# Why Deep Learning?

- ❑ In ~2010 **deep** learning techniques started outperforming other machine learning techniques. Why this decade?
  - Large amounts of training data favor deep learning
  - Faster machines and multicore CPU/GPUs favor deep learning
- ❑ New models, algorithms, ideas
  - Better, more flexible learning of intermediate representations
  - Effective end-to-end joint system learning
  - Effective learning methods for using contexts and transferring between tasks
  - Better regularization and optimization methods
- ❑ **Improved performance** (first in speech and vision, then NLP)

# Deep Learning in Speech

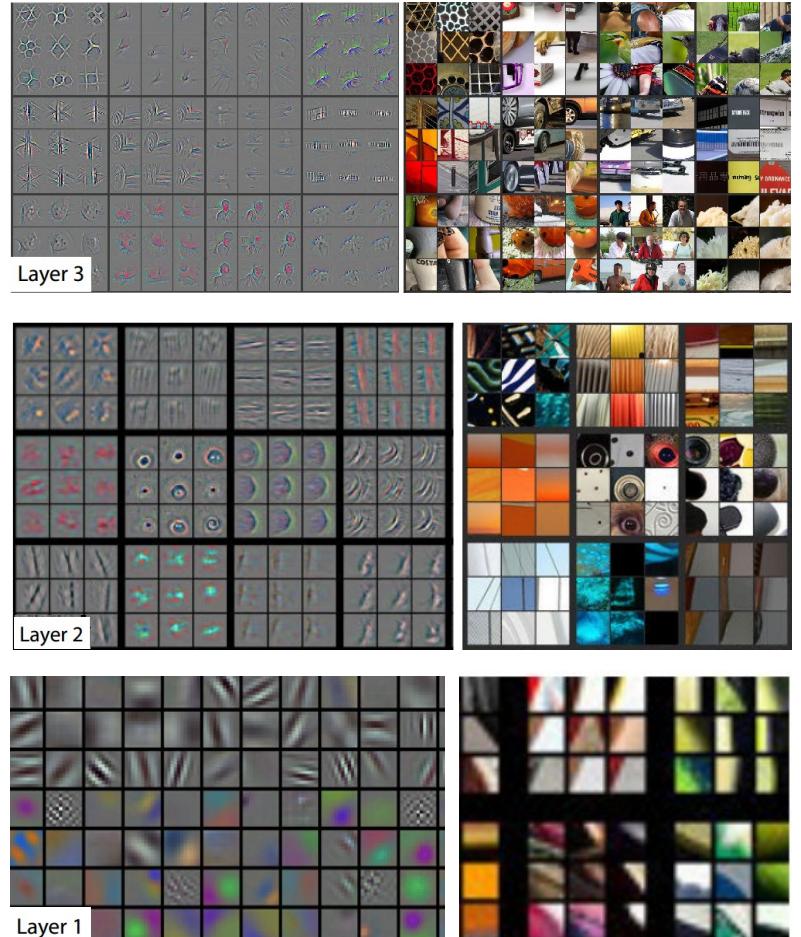
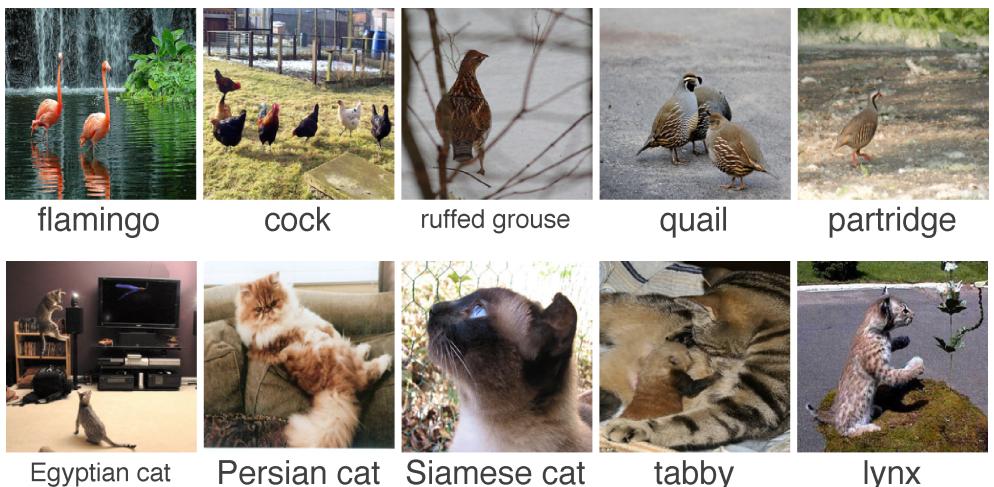
- The first breakthrough results of “deep learning” on large datasets happened in speech recognition
- Context-Dependent Pre-trained Deep Neural Networks for Large Vocabulary Speech Recognition Dahl et al. (2010)

Acoustic model and WER	RT03S FSH	Hub5 SWB
Traditional features	<b>27.4</b>	<b>23.6</b>
Deep Learning	<b>18.5</b> (-33%)	<b>16.1</b> (-32%)



# Deep Learning in Vision

- ❑ First major focus of deep learning groups was computer vision
- ❑ The breakthrough DL paper: ImageNet Classification with Deep Convolutional Neural Networks by Krizhevsky, Sutskever, & Hinton, 2012, U. Toronto. 37% error red.

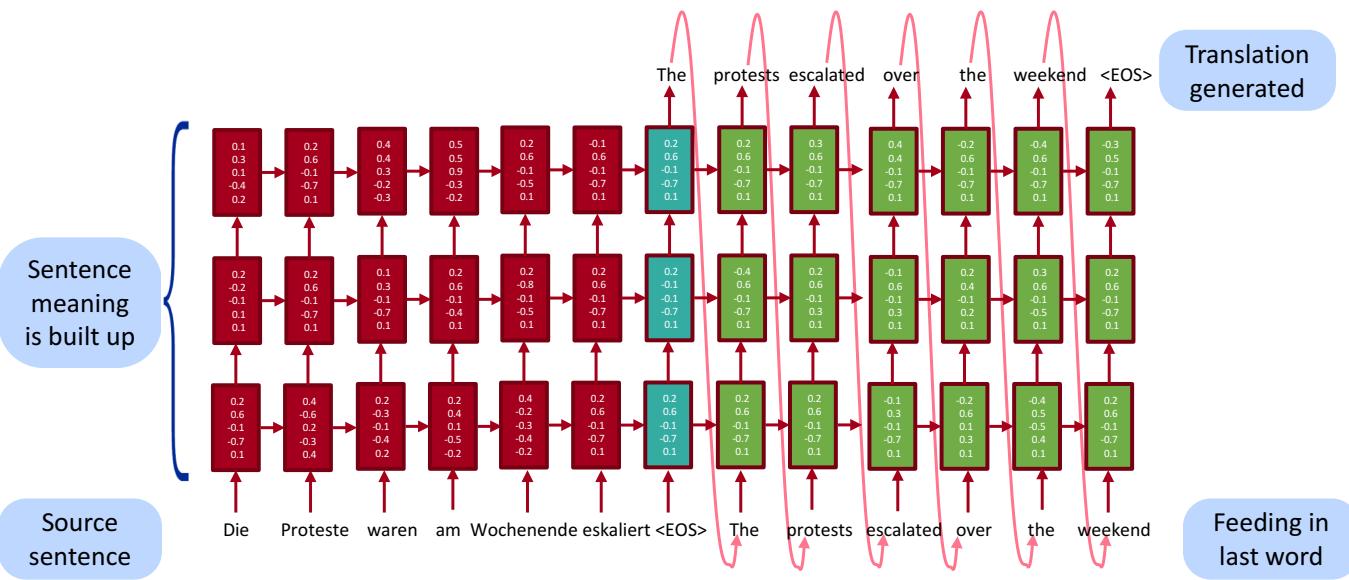


Zeiler and Fergus (2013)

Stanford 224n

# Deep Learning in Language

- ❑ Machine Translation: first breakthrough in 2015



- ❑ Now live for many languages in Google Translate (etc.), with big error reductions!

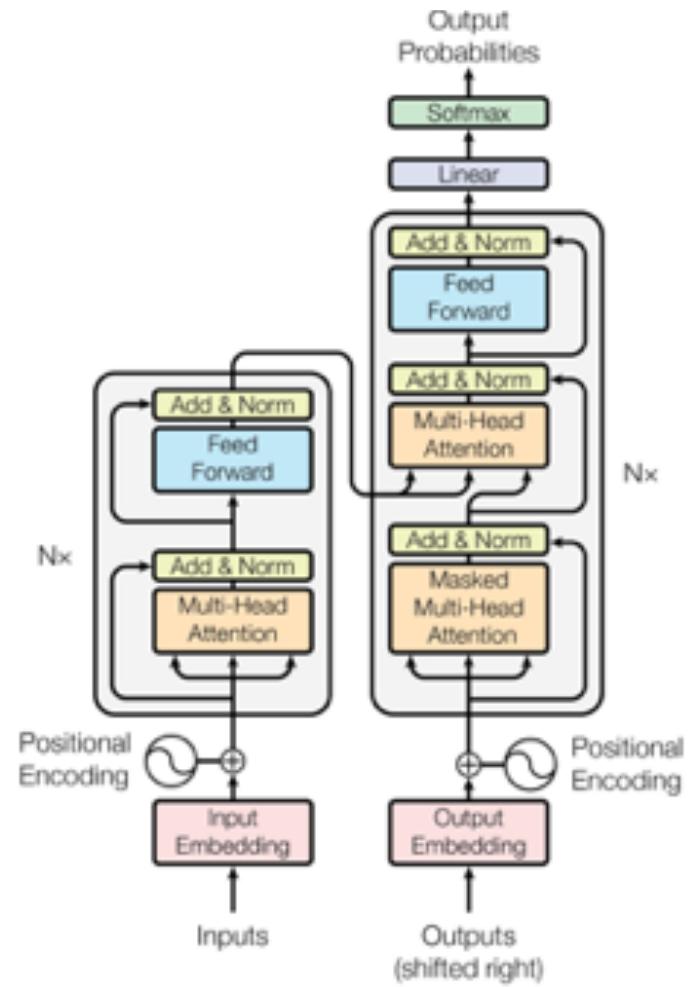


Figure 1: The Transformer - model architecture.

# Recent Breakthrough

BLOG POST  
RESEARCH

30 NOV 2020

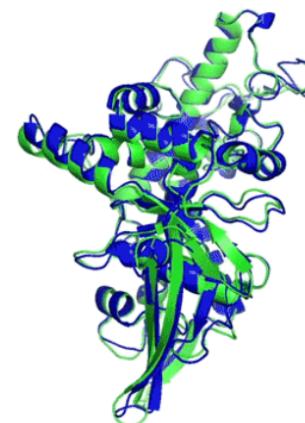
## AlphaFold: a solution to a 50-year-old grand challenge in biology



We have been stuck on this one problem – how do proteins fold up – for nearly 50 years. To see DeepMind produce a solution for this, having worked personally on this problem for so long and after so many stops and starts, wondering if we'd ever get there, is a very special moment.

PROFESSOR JOHN MOULT  
CO-FOUNDER AND CHAIR OF CASP, UNIVERSITY OF MARYLAND

### 3D structure



T1037 / 6vr4  
90.7 GDT  
(RNA polymerase domain)



T1049 / 6y4f  
93.3 GDT  
(adhesin tip)

- Experimental result
- Computational prediction

# Tentative Topics to Cover

- Word and word vectors
- Neural networks basics
- Text classification
- Recurrent neural networks
- Language Models
- Dependency and syntactic parsing
- Part-of-speech tagging and sequence labeling
- Machine translation and attention
- Transformers and large language models
- Question answering
- Chatbots and dialogue systems
- Natural language generation

# Word

# Word

- ❑ Before learning NLP, you need to know about natural language
  - Similarly, to study computational biology, you need to know biology
  
- ❑ We will cover tokenization and morphology

# Tokenization

- ❑ Tokenization: convert a character stream into words by adding spaces
  - For certain languages, highly nontrivial
  - E.g., Chinese word segmentation is a widely- studied NLP task
  - for other languages (English), tokenization is easier but is still not always obvious

# Intricacies of Tokenization

- Separating punctuation characters from words?
  - , " ? ! -> always separate
  - . -> when shouldn't we separate it?
    - *Dr., Mr., Prof., U.S., etc.*
  
- English contractions:
  - *isn't, aren't, wasn't,...* -> *is n't, are n't, was n't,...*
  - but how about these: *can't, won't* -> *ca n't, wo n't*
  - *ca* and *wo* are then different forms from *can* and *will*

# Intricacies of Tokenization

- Chinese and Japanese: no spaces between words:
  - 莎拉波娃现在居住在美国东南部的佛罗里达。
  - 莎拉波娃 现在 居住 在 美国 东南部 的 佛罗里达
  - Sharapova now lives in US southeastern Florida
  
- Further complicated in Japanese, with multiple alphabets intermingled
  - Dates/amounts in multiple formats

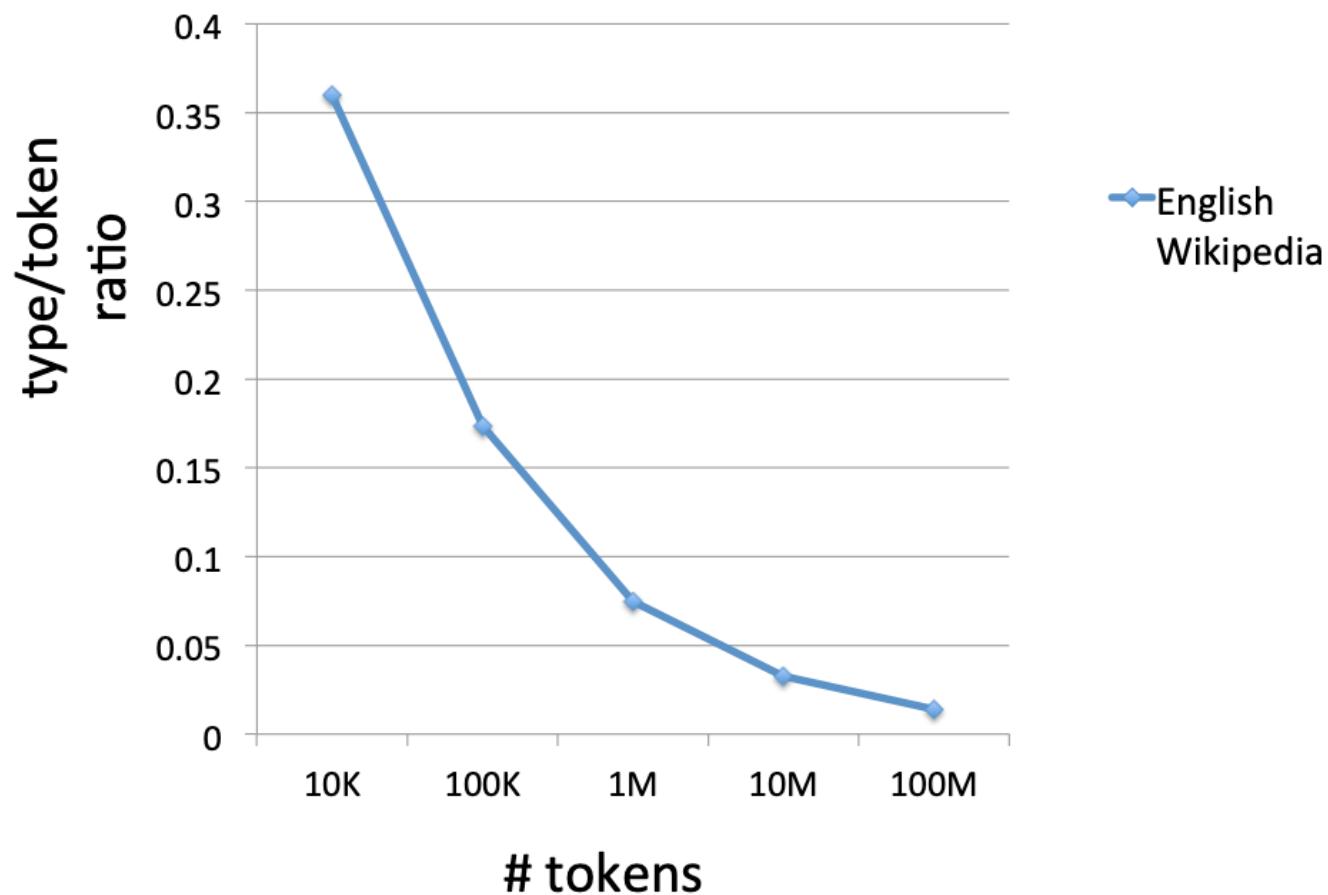


# Types and Tokens

- Once text has been tokenized, let's count the words
- Types: entries in the vocabulary (vocabulary = set of types)
- Tokens: instances of types in a corpus
- Example sentence: *If they want to go , they should go .*
  - how many types?
  - how many tokens?
- Type/token ratio: useful statistic of a corpus (here, 0.8)
- As we add data, what happens to the type-token ratio?

# Types and Tokens

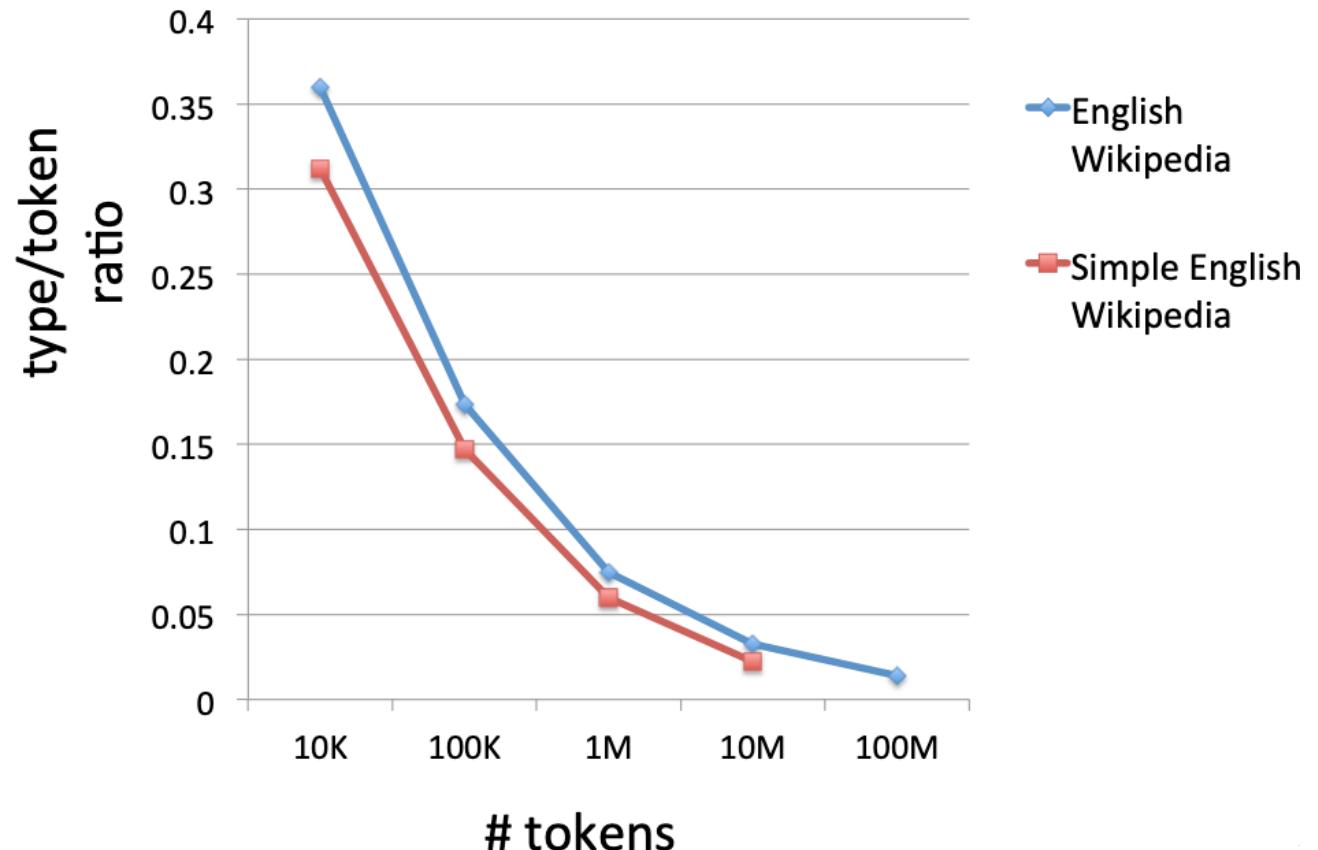
- ❑ More data -> Lower type/token ratio



# Types and Tokens

- What has a higher type/token ratio, Simple English Wikipedia or English Wikipedia?

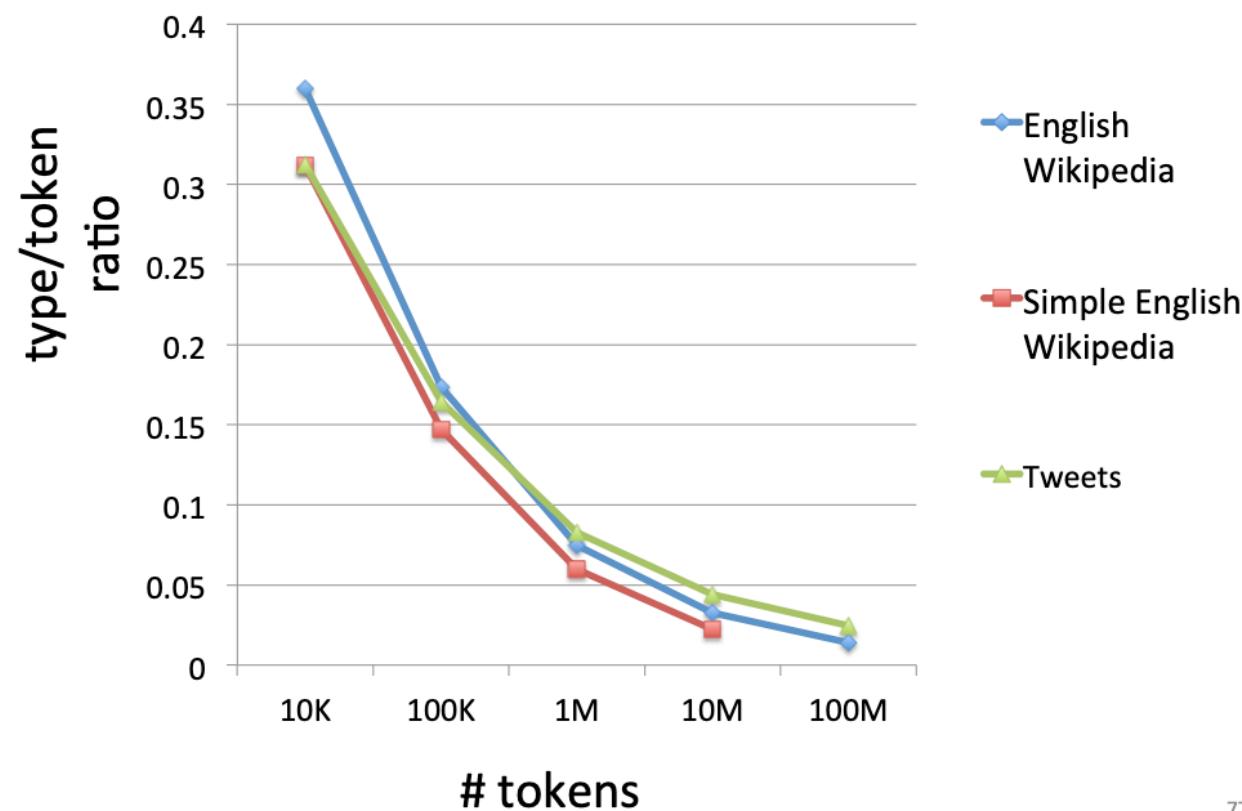
- type/token ratio is one measure of complexity



# Types and Tokens

## ❑ Wikipedia vs Tweets?

- Tweets (once you have 1 million or more tokens)



# “really” on Twitter

224571	really	50	real1111ly	15	reall11yy	8	reallyyyyyyy	6	real1111111ly	4	real111111yyy
1189	rly	48	reeeeeally	15	real111111ly	8	reallyyyyyy	6	reaaaaaallly	4	reaall1yyy
1119	realy	41	reeally	15	reaallly	8	realky	5	rrrreally	4	reaallly
731	rlly	38	really2	14	reeeeeally	7	relaly	5	rrly	4	reaaallly
590	really	37	reaaaallly	14	real111yyyy	7	eeeeeeeeally	5	rellly	4	reaaallly
234	real11ly	35	reallyyyyy	13	reeeaally	7	reeeealy	5	eeeeeeeally	4	reaaaaly
216	reallyy	31	reely	12	rreally	7	reeeeaaally	5	reeeeeally	3	reeeeeallly
156	relly	30	real11yyy	12	reaaaaally	7	real11111yyy	5	reeeeeaaally	3	reeeeallly
146	real111ly	27	real11yy	11	reeeeallly	7	real11111111ly	5	reealllyy	3	reeeeaaaaally
132	rily	27	reaaly	11	reeeallly	7	reaaaaaaally	5	real11111111ly	3	reeeaallly
104	reallyyy	26	real111yyy	11	real1111yyy	7	raelly	5	real1111111111ly	3	reeeaalllyyyy
89	reeeally	25	real111111ly	11	reaallyy	7	r3ally	5	reaalllyy	3	reealy
89	real1111ly	22	reaaallly	10	reallyreallyreally	6	r-really	5	reaaaallly	3	reeallly
84	reaaally	21	really-	10	reaaaly	6	reeeaalllyyy	5	reaaaallly	3	reeaaly
82	reaally	19	reeaally	9	eeeeeeeeally	6	reeeaallly	4	r1lly	3	reeaalllyyy
72	reeeeeally	18	real111yyy	9	reallys	6	reeeeaaally	4	eeeeeeeeeeeally	3	reeaallly
65	reaaaally	16	reaaaallly	9	really-really	6	realyl	4	reeealy	3	reeaaallly
57	reallyyyyy	15	realyy	9	r)eally	6	r-e-a-l-1-y	4	reeaaaally	3	reallyyyyyyyyy
53	rilly	15	reallyreally	8	reeeaally	6	real111yyyy	4	real11111yyy	3	reallyl

# “really” on Twitter

3 really)	2 rlyyyy	2 reeaallyy	1 rrrrrrrrrrrrrrreeeeeeeeeaaaaalllllllyyyyyyy
3 r]eally	2 rlyyy	2 reeaalllyy	1 rrrrrrrrrrreally
3 realluy	2 reqally	2 reeaallly	1 rrrrrrrreeeeeaaaalllllyyyyyyy
3 reallllyyyyy	2 rellyy	2 reeaally	1 rrrrrrrealy
3 reallllllyyyyyyy	2 rellys	2 reaqllly	1 rrrrrrreally
3 reallllllyyyy	2 reeely	2 realyyy	...
3 realllllly	2 reeeeeeealy	2 reallyyyyyyyyyyy	1 re-he-he-heeeeally
3 realllllllllllllly	2 reeeeeallly	2 reallyyyyyyy	1 re-he-he-he-ealy
3 realiy	2 reeeeeaaally	2 really*	1 reheheally
3 reaallyyyy	2 reeeeeaaally	2 really/	1 reelllyy
3 reaalllly	2 reeeeeaaalllly	2 realllyyyyy	1 reellly
3 reaaallyy	2 reeeeallyyy	2 reallllyyyyy	1 ree-hee-heally
3 reaaaallyy	2 reeeealllyyy	2 reallllyyyyy	...
3 reaaaalllly	2 reeeeaaallllyyy	2 reallllyyyyy	1 reeeeeeeeaaally
3 reaaaaaaly	2 reeeeaaalllly	2 reallllllyyyyy	1 reeeeeeeeaaally
3 reaaaaaaaally	2 reeeeaaaally	2 reallllllyyyyy	1 reeeeeeeeaaaaalllyyy
3 r34lly	2 reeeeaaaalllyyy	2 reallllllyyy	1 reeeeeeeeaaaaalllllllyyyyyyyyy
2 rrreally	2 reeeallyy	2 realllllllllllllly	1 reeeeeeeeaaaaaaaalllllllyyyyyyyyy
2 rreeaallyy	2 reeallyy	2 realllllllllllllly	1 reeeeeeeeaaaaaaaalllllllyyyyyyyyy

# “really” on Twitter

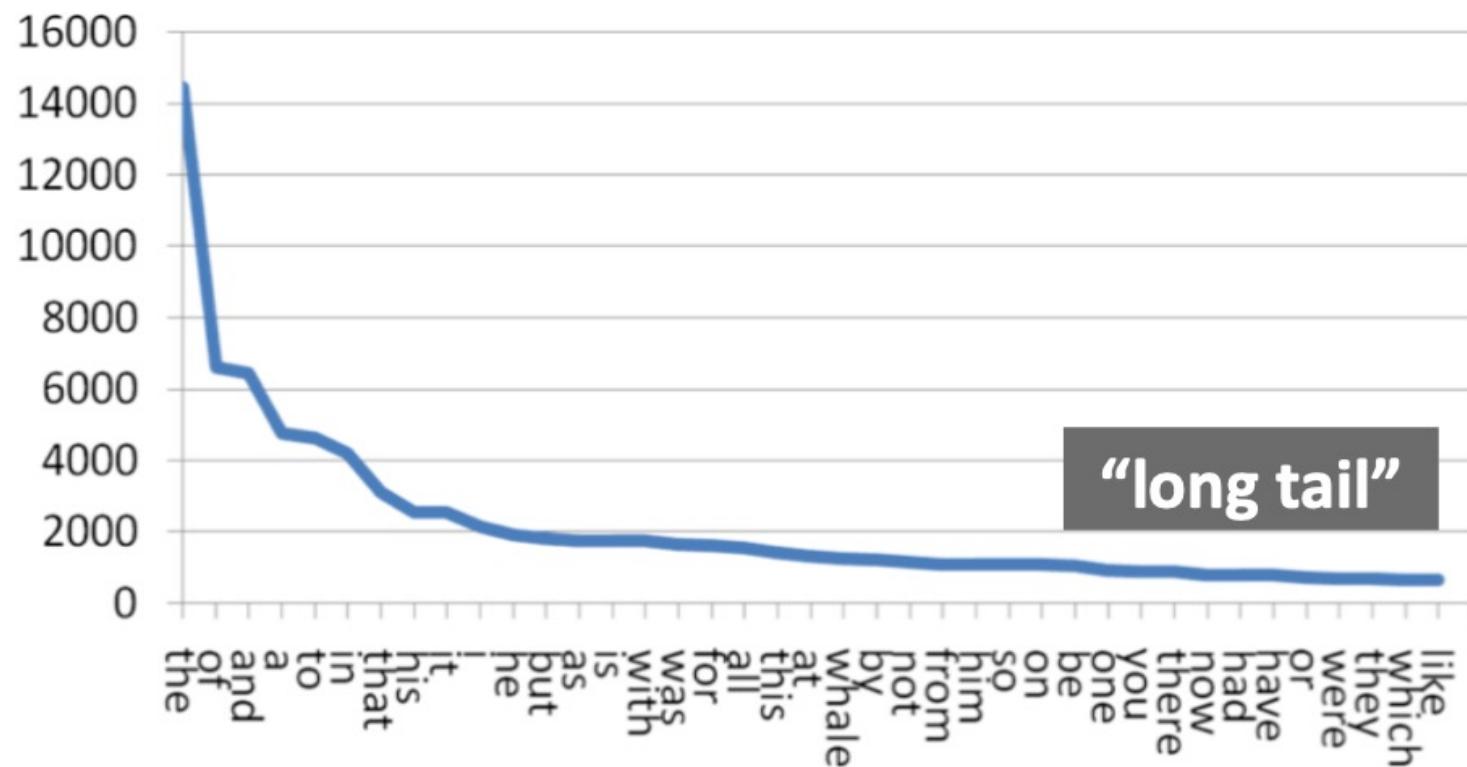
```
1 reallyreallyreallyreallyreallyreallyreallyreallyreallyreally  
reallyreallyreallyreallyreallyreallyreallyreally  
1 reallyreallyreallyreallyreallyreallyr33ly  
1 really/really/really  
1 really(really  
...  
1 real11111111yyyy  
1 real1111111111yyyyyy  
1 real1111111111yyyy  
1 real1111111111yyyy  
1 real1111111111yyy  
1 real111111111111yyyy  
1 real11111111111111yyyy  
1 real1111111111111111y  
1 real111111111111111111y  
1 real11111111111111111111y  
1 real1111111111111111111111y  
1 real111111111111111111111111y  
1 real11111111111111111111111111y  
1 real1111111111111111111111111111y  
1 real111111111111111111111111111111y  
1 real11111111111111111111111111111111y
```

# How many words are there?

- How many English words exist?
- When we increase the size of our corpus, what happens to the number of types?
  - A bit surprising: vocabulary continues to grow in any actual dataset
  - You'll just never see all the words
  - In 1 million tweets, 15M tokens, 600k types
  - In 56 million tweets, 847M tokens, 11M types

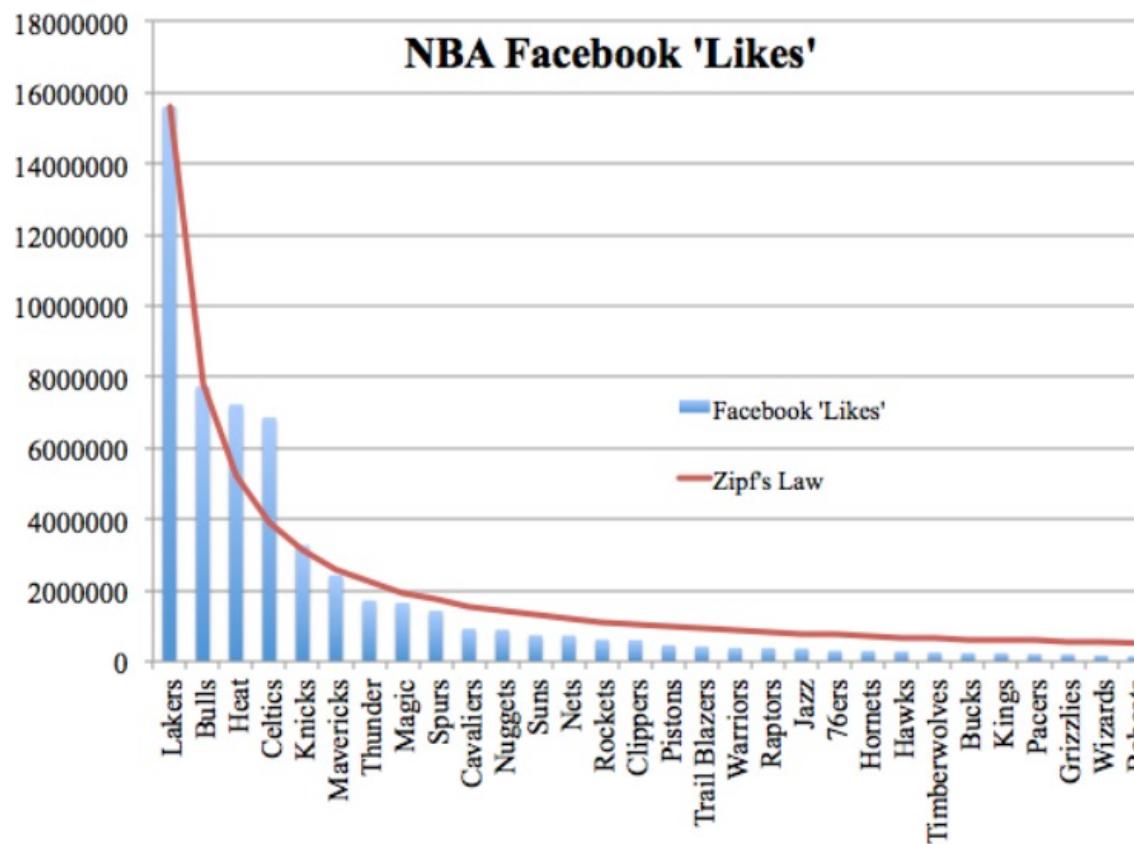
# How are words distributed?

- Zipf's law: frequency of a word is inversely proportional to its rank



# Zipf's law

- ❑ Also predicts other kinds of data: population of cities in a country, revenue of different companies, etc.



# Another option for text tokenization

- ❑ Instead of
  - white-space segmentation
  - single-character segmentation
- ❑ Use the data to tell us how to tokenize
- ❑ Subword tokenization (because tokens can be parts of words as well as whole words)

# Subword tokenization

- ❑ Three common algorithms:
  - Byte-Pair Encoding (BPE) (Sennrich et al., 2016)
  - Unigram language modeling tokenization (Kudo, 2018)
  - WordPiece (Schuster and Nakajima, 2012)
- ❑ All have 2 parts:
  - A token learner that takes a raw training corpus and induces a vocabulary (a set of tokens).
  - A token segmenter that takes a raw test sentence and tokenizes it according to that vocabulary

# Byte Pair Encoding (BPE) token learner

- Let vocabulary be the set of all individual characters
  - = {A, B, C, D, ..., a, b, c, d, ...}
- Repeat:
  - Choose the two symbols that are most frequently adjacent in the training corpus (say 'A', 'B')
  - Add a new merged symbol 'AB' to the vocabulary
  - Replace every adjacent 'A' 'B' in the corpus with 'AB'.
- Until k merges have been done.

# BPE token learner algorithm

```
function BYTE-PAIR ENCODING(strings  $C$ , number of merges  $k$ ) returns vocab  $V$ 
     $V \leftarrow$  all unique characters in  $C$           # initial set of tokens is characters
    for  $i = 1$  to  $k$  do                      # merge tokens til k times
         $t_L, t_R \leftarrow$  Most frequent pair of adjacent tokens in  $C$ 
         $t_{NEW} \leftarrow t_L + t_R$                   # make new token by concatenating
         $V \leftarrow V + t_{NEW}$                       # update the vocabulary
        Replace each occurrence of  $t_L, t_R$  in  $C$  with  $t_{NEW}$       # and update the corpus
    return  $V$ 
```

# Byte Pair Encoding (BPE) Addendum

- ❑ Most subword algorithms are run inside space-separated tokens.
- ❑ So we commonly first add a special end-of-word symbol '' before space in training corpus
- ❑ Next, separate into letters.

# BPE token learner

Original (very fascinating 😊) corpus:

low low low low lowest lowest newer newer newer  
newer newer newer wider wider wider new new

Add end-of-word tokens, resulting in this vocabulary:

**vocabulary**

\_, d, e, i, l, n, o, r, s, t, w

# BPE token learner

**corpus**

5 low \_  
2 lowest \_  
6 newer \_  
3 wider \_  
2 new \_

**vocabulary**

\_, d, e, i, l, n, o, r, s, t, w

Merge **e r** to **er**

**corpus**

5 low \_  
2 lowest \_  
6 newer \_  
3 wider \_  
2 new \_

**vocabulary**

\_, d, e, i, l, n, o, r, s, t, w, er

# BPE

**corpus**

5 low \_  
2 lowest \_  
6 newer \_  
3 wider \_  
2 new \_

**vocabulary**

\_, d, e, i, l, n, o, r, s, t, w, er

Merge er \_ to er\_

**corpus**

5 low \_  
2 lowest \_  
6 newer\_  
3 wider\_  
2 new \_

**vocabulary**

\_, d, e, i, l, n, o, r, s, t, w, er, er\_

# BPE

**corpus**

5 l o w \_  
2 l o w e s t \_  
6 n e w er\_  
3 w i d er\_  
2 n e w \_

**vocabulary**

\_, d, e, i, l, n, o, r, s, t, w, er, er\_

**Merge n e to ne****corpus**

5 l o w \_  
2 l o w e s t \_  
6 ne w er\_  
3 w i d er\_  
2 ne w \_

**vocabulary**

\_, d, e, i, l, n, o, r, s, t, w, er, er\_, ne

# BPE

The next merges are:

**Merge**

(ne, w)

(l, o)

(lo, w)

(new, er\_)

(low, \_\_)

**Current Vocabulary**

\_, d, e, i, l, n, o, r, s, t, w, er, er\_, ne, new

\_, d, e, i, l, n, o, r, s, t, w, er, er\_, ne, new, lo

\_, d, e, i, l, n, o, r, s, t, w, er, er\_, ne, new, lo, low

\_, d, e, i, l, n, o, r, s, t, w, er, er\_, ne, new, lo, low, newer\_

\_, d, e, i, l, n, o, r, s, t, w, er, er\_, ne, new, lo, low, newer\_, low\_

# BPE token segmenter algorithm

- ❑ On the test data, run each merge learned from the training data:
  - Greedily
  - In the order we learned them
  - (test frequencies don't play a role)
- ❑ So: merge every e r to er, then merge er \_ to er\_, etc.
- ❑ Result:
  - Test set "n e w e r \_" would be tokenized as a full word
  - Test set "l o w e r \_" would be two tokens: "low er\_

# Properties of BPE tokens

- ❑ Usually include frequent words
- ❑ And frequent subwords
  - Which are often morphemes like -est or –er
- ❑ A **morpheme** is the smallest meaning-bearing unit of a language
  - unlikeliest has 3 morphemes un-, likely, and -est