

# CMSC 25700/35100: Natural Language Processing Lecture 2: Words and Word Representations

Chenhao Tan  
University of Chicago  
@ChenhaoTan, [chenhao@uchicago.edu](mailto:chenhao@uchicago.edu)

# Logistics

- Quizzes (I got 60 responses, you are all quick learners!!!)
- A1 is released
- Office hours (Mourad's office hour is 5-6 on Tuesday now)
- Coding tutorials will happen at 11am on Jan 21 and Jan 28

# Basic unit of language

- Words
- But what are words and how do computers know words?

**Tokenization:** convert sequence of characters into sequence of tokens

“Oh!” said Lydia stoutly, “I am not afraid; for though I am the youngest, I’m the tallest.”

(Austen, 1813)

“Oh!” said Lydia stoutly, “I am not afraid; for though I \_am\_ the youngest, I’m the tallest.”

(Austen, 1813)

“Oh!”  
said  
Lydia  
stoutly,  
“I  
am

not  
afraid;  
for  
though  
I  
\_am\_

the  
youngest,  
I’m  
the  
tallest.”

“Oh!” said Lydia stoutly, “I am not afraid; for though I \_am\_ the youngest, I’m the tallest.”

(Austen, 1813)

“ Oh ! ” said Lydia stoutly , “ I am not afraid ; for though I \_ am \_ the youngest , I ’ m the tallest . ”

- Most tokenizers are rule-based
- several conventions:

|        | Penn Treebank | Moses   |
|--------|---------------|---------|
| don't  | do n't        | don 't  |
| aren't | are n't       | aren 't |
| can't  | ca n't        | can 't  |
| won't  | wo n't        | won 't  |

- important to be consistent across NLP systems, match tokenization of external tools/resources
- see `nltk.tokenize` (also for sentence tokenization)

- Chinese, Japanese, Thai: no spaces between words
- Tokenization becomes highly non-trivial!

姚明进入总决赛  
“Yao Ming reaches the finals”

- Multiple conventions:

姚明 进入 总决赛  
“YaoMing reaches finals”

Chinese Treebank

姚 明 进入 总 决赛  
“Yao Ming reaches overall finals”

Peking University

Example from Chen et al. (2017), cited in Jurafsky & Martin (SLP3)



“ oh ! ” said lydia stoutly , “ i am not afraid ; for  
though i \_ am \_ the youngest , i ’m the tallest . ”

|   |     |   |        |   |          |
|---|-----|---|--------|---|----------|
| 3 | i   | 1 | !      | 1 | oh       |
| 2 | ,   | 1 | .      | 1 | said     |
| 2 | _   | 1 | ;      | 1 | stoutly  |
| 2 | am  | 1 | afraid | 1 | tallest  |
| 2 | the | 1 | for    | 1 | though   |
| 2 | “   | 1 | lydia  | 1 | youngest |
| 2 | ”   | 1 | not    | 1 | 'm       |

|   |     |   |        |   |          |
|---|-----|---|--------|---|----------|
| 3 | i   | 1 | !      | 1 | oh       |
| 2 | ,   | 1 | .      | 1 | said     |
| 2 | _   | 1 | ;      | 1 | stoutly  |
| 2 | am  | 1 | afraid | 1 | tallest  |
| 2 | the | 1 | for    | 1 | though   |
| 2 | "   | 1 | lydia  | 1 | youngest |
| 2 | "   | 1 | not    | 1 | 'm       |

**type:** a unique word (an entry in a vocabulary or dictionary)

**token:** an instance of a type in a corpus

|   |     |   |        |   |          |
|---|-----|---|--------|---|----------|
| 3 | i   | 1 | !      | 1 | oh       |
| 2 | ,   | 1 | .      | 1 | said     |
| 2 | _   | 1 | ;      | 1 | stoutly  |
| 2 | am  | 1 | afraid | 1 | tallest  |
| 2 | the | 1 | for    | 1 | though   |
| 2 | "   | 1 | lydia  | 1 | youngest |
| 2 | "   | 1 | not    | 1 | 'm       |

two types of counts:

type count =

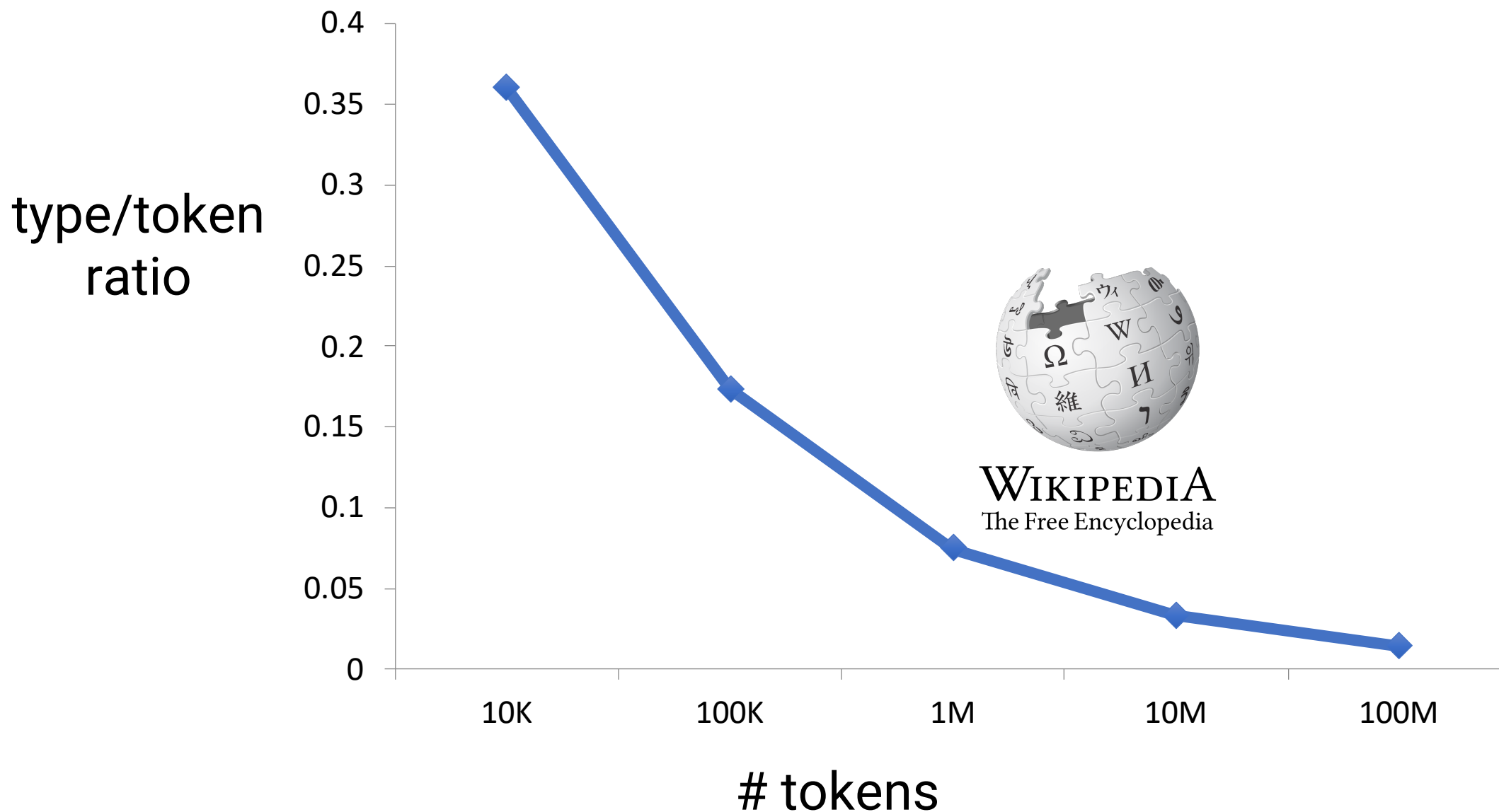
token count =

useful statistic: **type/token ratio**

(here,  $21/29 = 0.724$ )

How does the type/token ratio change when adding more data?

# more data $\rightarrow$ lower type/token ratio





WIKIPEDIA  
The Free Encyclopedia

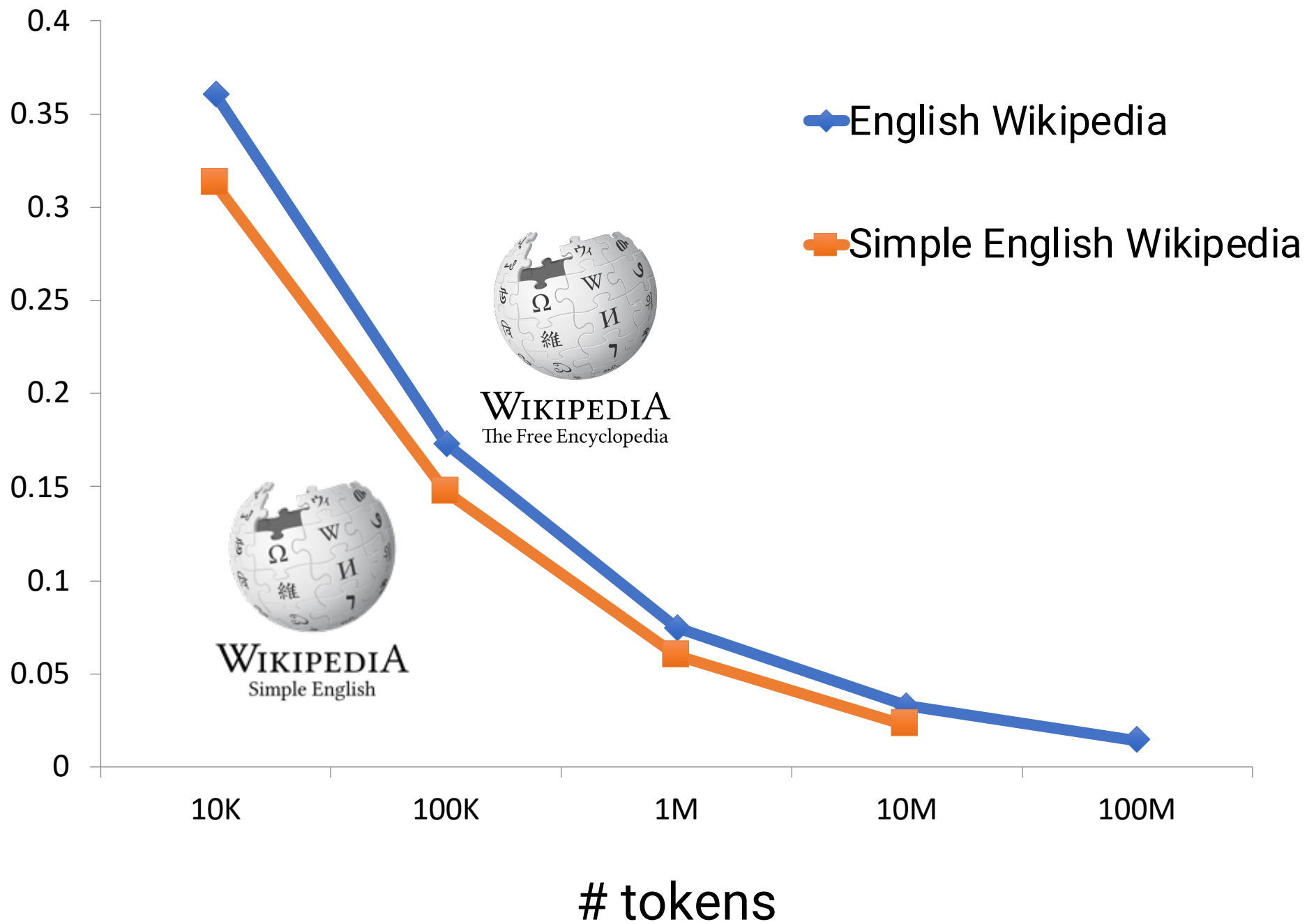
vs.



WIKIPEDIA  
Simple English

Which has a higher type/token ratio?

type/token  
ratio





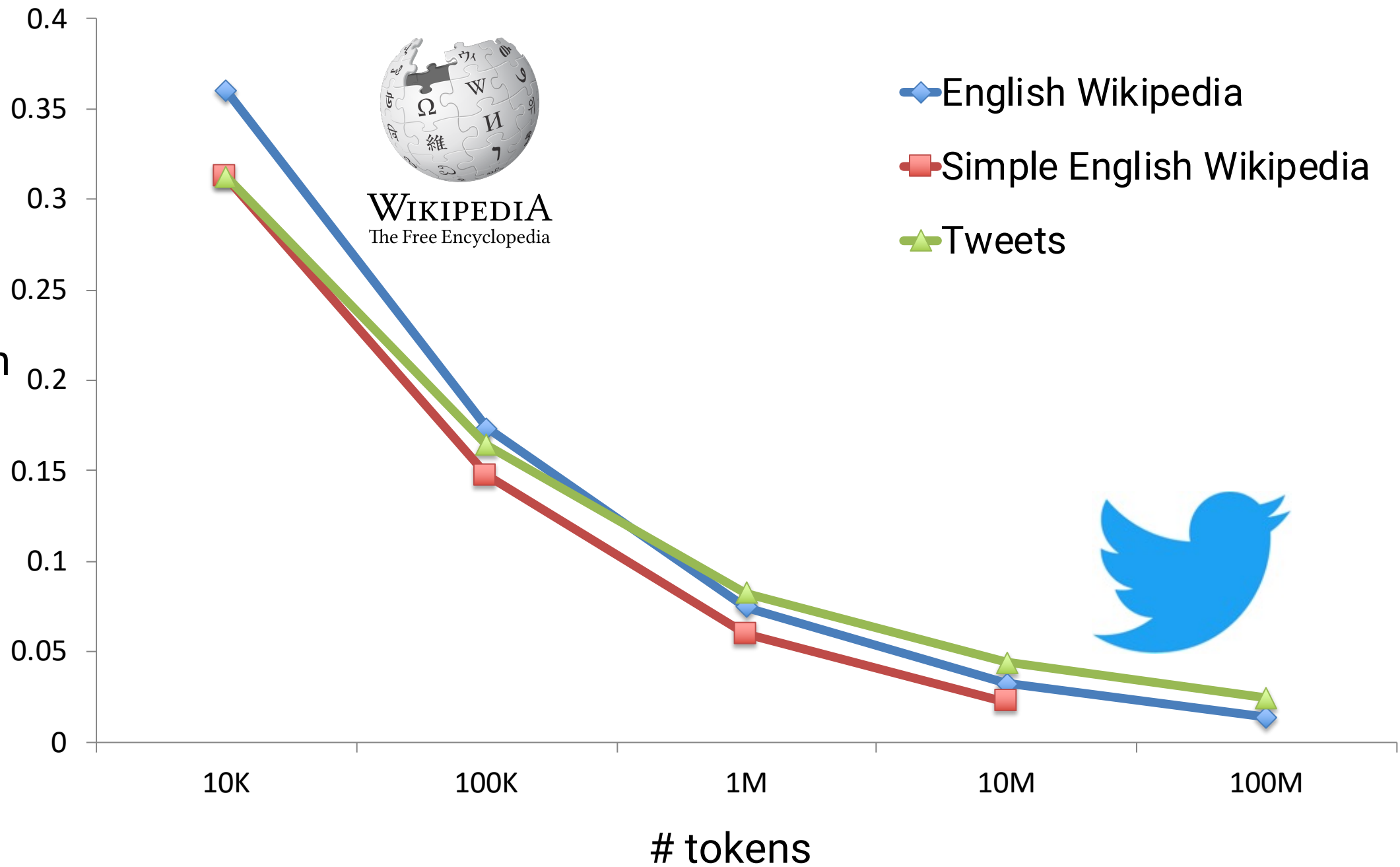
**WIKIPEDIA**  
The Free Encyclopedia

**VS.**





type/token  
ratio



|                  |                    |                       |
|------------------|--------------------|-----------------------|
| 224571 really    | 38 really2         | 12 reaaaaaally        |
| 1189 rly         | 37 reaaaaaally     | 12 rreally            |
| 1119 realy       | 35 reallyyyyy      | 11 reaallyy           |
| 731 rlly         | 31 reely           | 11 realllllyyy        |
| 590 realllly     | 30 reallllyyy      | 11 reeeallly          |
| 234 realllly     | 27 reaaly          | 11 reeeeeallly        |
| 216 reallyy      | 27 reallllyy       | 10 reaaaly            |
| 156 relly        | 26 reallllyyyy     | 10 reallyreallyreally |
| 146 reallllly    | 25 reallllllllly   | 9 r)eally             |
| 132 rily         | 22 reaaallly       | 9 really-really       |
| 104 reallyyy     | 21 really-         | 9 reallys             |
| 89 reallllllly   | 19 reeaally        | 9 reeeeeeeally        |
| 89 reeeally      | 18 reallllyyy      | 8 realky              |
| 84 reaaally      | 16 reaaaaallly     | 8 reallyyyyyyy        |
| 82 reaally       | 15 reaallly        | 8 reallyyyyyyyy       |
| 72 reeeeeally    | 15 reallllllllllly | 8 reeeaaally          |
| 65 reaaaaally    | 15 reallllyy       | 7 r3ally              |
| 57 reallyyyy     | 15 reallyreally    | 7 raelly              |
| 53 rilly         | 15 realyy          | 7 reaaaaaaally        |
| 50 reallllllllly | 14 reallllyyyy     | 7 reallllllllllllllly |
| 48 reeeeeeeally  | 14 reeeeeeeally    | 7 realllllllyyy       |
| 41 reeally       | 13 reeeaaally      | 7 reeeeeaaally        |

|                        |                    |                            |
|------------------------|--------------------|----------------------------|
| 7 reeeealy             | 5 rrly             | 3 reali y                  |
| 7 reeeeeeeeeally       | 5 rrrreally        | 3 realllllllllllllllllllly |
| 7 relaly               | 4 reaaaaly         | 3 realllllllyy             |
| 6 r-e-a-l-l-y          | 4 reaaalllly       | 3 realllllllyyyy           |
| 6 r-really             | 4 reaaallllyy      | 3 realllllllyyyyyyy        |
| 6 reaaaaaallly         | 4 reaalllly        | 3 reallllyyyyyy            |
| 6 realllllllllllly     | 4 reaallllyyy      | 3 realluy                  |
| 6 reallllyyyyyy        | 4 realllllllllyyyy | 3 really)                  |
| 6 realyl               | 4 reallllllyyyy    | 3 reallyl                  |
| 6 reeeaaaally          | 4 reeeaaaally      | 3 reallyyyyyyyyyy          |
| 6 reeeaaallly          | 4 reeealy          | 3 reeeaaallly              |
| 6 reeeaaalllyyy        | 4 reeeeeeeeeeeally | 3 reeaalllly               |
| 5 reaaaaaallly         | 4 rllly            | 3 reeaalllyyy              |
| 5 reaaaaalllly         | 3 r34lly           | 3 reeaaly                  |
| 5 reaallllyy           | 3 r]eally          | 3 reeaallly                |
| 5 realllllllllllllllly | 3 reaaaaaaaally    | 3 reealy                   |
| 5 realllllllllllllly   | 3 reaaaaaly        | 3 reeeaaallllyyy           |
| 5 reeallyyy            | 3 reaaaalllly      | 3 reeeaaallly              |
| 5 reeeaaaallly         | 3 reaaaalllyy      | 3 reeeaaaaaally            |
| 5 reeeeaally           | 3 reaaallly        | 3 reeeeaalllly             |
| 5 reeeeeeeally         | 3 reaalllly        | 3 reeeeealllly             |
| 5 relly                | 3 reaallyyy        | 2 reaaaaaaaaaally          |

|                            |                      |                              |
|----------------------------|----------------------|------------------------------|
| 2 reaaaaaaaaaally          | 2 really/            | 2 reeely                     |
| 2 reaaaaaaaaallly          | 2 reallyyyyyyyyy     | 2 rellys                     |
| 2 reaaaaaaalllllyyy        | 2 reallyyyyyyyyyyyyy | 2 rellyy                     |
| 2 reaaaaaallllly           | 2 realyyy            | 2 reqally                    |
| 2 reaaaaaallllly           | 2 reaqlly            | 2 rlyyy                      |
| 2 reaaalllllyyy            | 2 reeaaally          | 2 rlyyyy                     |
| 2 reaaalllllyyy            | 2 reeaallly          | 2 rreeaallyy                 |
| 2 reaaalllyyy              | 2 reeaalllyy         | 2 rrreally                   |
| 2 reaalllllyy              | 2 reeaallyy          | 1 r-r-r-really               |
| 2 reaalllllyyy             | 2 reeallyy           | 1 r3aly                      |
| 2 reaallyyy                | 2 reeeallyy          | 1 r3ly                       |
| 2 reaalyy                  | 2 reeeaaaalllyyy     | 1 raaahhhlllaaayyyy          |
| 2 realllllllllllllllllllly | 2 reeeaaaally        | 1 raeally                    |
| 2 realllllllllllllllllly   | 2 reeeaaaallly       | 1 re-e-e-eally               |
| 2 reallllllllyy            | 2 reeeaaaalllyyyy    | 1 re-eaaaaaaly               |
| 2 reallllllllyyyyy         | 2 reeeeaalllyyy      | 1 re-he-he-he-ealy           |
| 2 reallllllllyyyyy         | 2 reeeeaallyy        | 1 re-he-he-heeeeeally        |
| 2 reallllllyy              | 2 reeeeeeaaallllly   | 1 re3ally                    |
| 2 reallllllyyyyyyy         | 2 reeeeeeaaally      | 1 rea(1)ly                   |
| 2 reallllyyyyyyy           | 2 reeeeeeaaally      | 1 reaaaaaaaaaaaaaaaaaaaaally |
| 2 reallyyyyyyy             | 2 reeeeeeallly       | 1 reaaaaaaaaaaaaaaaaaaaaally |
| 2 really*                  | 2 reeeeeealy         | 1 reaaaaaaaaaaaaaaaaaallllly |

[illegible]

|                                     |                                |                                   |
|-------------------------------------|--------------------------------|-----------------------------------|
| 1 reallyyyyyyyyyyy                  | 1 reeeaaallllyyy               | 1 reeeeaalllly                    |
| 1 really🙄                           | 1 reeeaaallly                  | 1 reeeeaallly                     |
| 1 realoly                           | 1 reeeaaalllyy                 | 1 reeeeaalllyy                    |
| 1 realys                            | 1 reeeaaaly                    | 1 reeeeaalllyyy                   |
| 1 realyyyyy                         | 1 reeeaalllly                  | 1 reeeeeeaaaaallllllly            |
| 1 real•ly                           | 1 reeeaallyy                   | 1 reeeeeeaaaaally                 |
| 1 reawly                            | 1 reeeaallyyy                  | 1 reeeeeeaaalllly                 |
| 1 ree-hee-heally                    | 1 reeeaalllly                  | 1 reeeeeeaaalllyyy                |
| 1 reeaaaaaaaaaaaaaaaaaalllllllllyyy | 1 reeeaalllly                  | 1 reeeeeeaaallly                  |
| 1 reeaaaaaaaaalllllllly             | 1 reeeaallllyy                 | 1 reeeeeeallllyyy                 |
| 1 reeaaaaalllllly                   | 1 reeeaallllyyy                | 1 reeeeeealy                      |
| 1 reeaaaallly                       | 1 reeeaallllyyyy               | 1 reeeeeeaaaaaally                |
| 1 reeaaalllly                       | 1 reeeaallyyy                  | 1 reeeeeeaaaaallllyyy             |
| 1 reeaaallllyyy                     | 1 reeeaallyyyy                 | 1 reeeeeeaaaaaally                |
| 1 reeaaallyy                        | 1 reeeallys                    | 1 reeeeeeaaalllly                 |
| 1 reeaaaly                          | 1 reeeaaaaaaaaaalllllllyyyyyyy | 1 reeeeeeaaally                   |
| 1 reeaalllllyyy                     | 1 reeeeeeaaaaaalllllllyyyy     | 1 reeeeeeaaaly                    |
| 1 reeaalllllyyy                     | 1 reeeeeeaaaaaallllllly        | 1 reeeeeealllly                   |
| 1 reeaallyyy                        | 1 reeeeeeaaaaalllly            | 1 reeeeeealllyyy                  |
| 1 reeaallly                         | 1 reeeeeeaaaaallly             | 1 reeeeeeaaaaallly                |
| 1 reeaalllyyy                       | 1 reeeeeeaaaaalllyyy           | 1 reeeeeeaaally                   |
| 1 reeaalllyyyy                      | 1 reeeeeeaaalllly              | 1 reeeeeeaaallyy                  |
| 1 reeeaaaaaaaly                     | 1 reeeeeeaaalllyyyyyy          | 1 reeeeeeaaaly                    |
| 1 reeeaaaaaally                     | 1 reeeeeeaaallllyy             | 1 reeeeeeaaaaaaaalllllyyyyyyy     |
| 1 reeeaaaalllly                     | 1 reeeeeeaaallly               | 1 reeeeeeaaaaaaaalllllllllyyyyyyy |
| 1 reeeaaaallllyyy                   | 1 reeeeeeaaallyy               | 1 reeeeeeaaaaaaaalllllllllyyyyyyy |
| 1 reeeaaaallllyyyy                  | 1 reeeeeeaaaly                 | 1 reeeeeeaaaaaaaalllyyy           |
| 1 reeeaaaallly                      | 1 reeeeeeallly                 | 1 reeeeeeaaaaaally                |
| 1 reeeaaaalllyyy                    | 1 reeeeeealllyy                | 1 reeeeeeaaaaaallyy               |
| 1 reeeaaaaly                        | 1 reeeeeealy                   | 1 reeeeeeaaallly                  |
| 1 reeeaaalllllyyy                   | 1 reeeeealllly                 | 1 reeeeeeaaally                   |

1 reeeeeeeeeeeaaally  
1 reeeeeeeeeeeaaally  
1 reeeeeeeeeeeaaally  
1 reeeeeeeeeeeaaally  
1 reeeeeeeeeeeeeeeaaaally  
1 reeeeeeeeeeeeeeeeeaaaally  
1 reeeeeeeeeeeeeeeeeaaally  
1 reeeeeeeeeeeeeeeeeesallllllllllllllllllllllllllly  
1 reeeeeeeelly  
1 reeeeeely  
1 reeeeeelly  
1 reeeeeely  
1 reeeellllllyyy  
1 reelllllyy  
1 reellly  
1 reelllyy  
1 reheheally  
1 relally  
1 rellllllly  
1 relllly  
1 rellyrell  
1 rellyz  
1 rieeely  
1 rlllllly  
1 rllllllyy  
1 rlllly  
1 rlllyrllly  
1 rllly  
1 rlllyy  
1 rlyrlyrly

1 rlyy  
1 rraarreellyy  
1 rreaalllyyy  
1 rreaally  
1 rreeaaalllllyyyy  
1 rreeaalllllyy  
1 rreeaallyyy  
1 rreealy  
1 rreeeaaaaalllllyyyy  
1 rreeeaaalllllyyy  
1 rreeeeeeaaaaallllllllyyyyyy  
1 rreeeeeeallly  
1 rreeeeely  
1 rrrealllyyy  
1 rrreeeaaalllyyy  
1 rrreeealllyyy  
1 rrreeeeaaalllllyy  
1 rrreeeeallly  
1 rrrlyyy  
1 rrrreeeally  
1 rrrreeeeeeaaaaallllllllyyyyyy  
1 rrrrreeeeaaalllly  
1 rrrrreeeeaaalllyy  
1 rrrrrreally  
1 rrrrrreal  
1 rrrrrrrreeeeeeaaaaallllllyyyyyy  
1 rrrrrrrrrreally  
1 rrrrrrrrrrrrrrrrrreeeeeeeeeeeeeeaaaaaallllllllyyyyyy

# How many words are there?

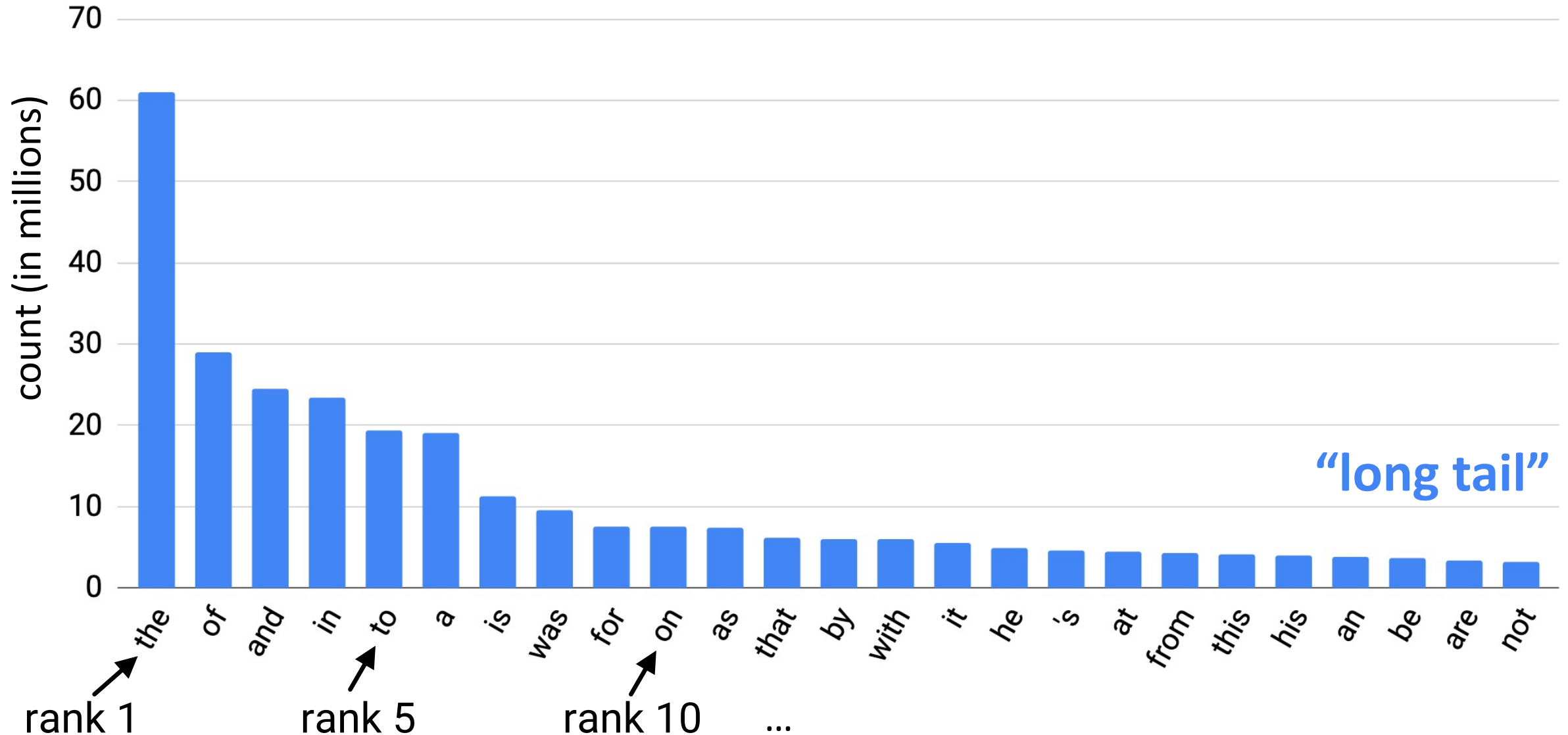
- size of vocabulary continues to grow as you collect more data
- you'll never find all the words



# How are words distributed?

|        |            |   |  |
|--------|------------|---|--|
| 224571 | really     | 1 | rreeeeeeaaaaaalllllllllyyyyyy                            |
| 1189   | rly        | 1 | rreeeeeeallly  |
| 1119   | realy      | 1 | rreeeeeeely  |
| 731    | rllly      | 1 | rrreallyyy   |
| 590    | reallly    | 1 | rrreeeaaalllyyy  |
| 234    | realllly   | 1 | rrreeealllyyy  |
| 216    | reallyy    | 1 | rrreeeeaaalllllyy  |
| 156    | relly      | 1 | rrreeeeallly   |
| 146    | reallllly  | 1 | rrrlyyy  |
| 132    | rily       | 1 | rrrreeeally  |
| 104    | reallyyy   | 1 | rrrreeeeeeaaaaaalllllllllyyyyyy                          |
| 89     | realllllly | 1 | rrrrreeeeaaallly   |
| 89     | reeeally   | 1 | rrrrreeeeaaalllyy  |
| 84     | reaaally   | 1 | rrrrrreally  |
| 82     | reaally    | 1 | rrrrrrealy   |
| 72     | reeeeally  | 1 | rrrrrrreeeeeeaaaaalllllyyyyyyy                           |
| 65     | reaaaally  | 1 | rrrrrrrrrrreally   |
| ...    |            | 1 | rrrrrrrrrrrrrrrrrrrrreeeeeeeeeeeeeeaaaaaalllllllllyyyyyy |

**Zipf's law:** frequency of a word is inversely proportional to its rank in the word frequency list



# The Long Tail

- there are so many word types!
- but words have **internal structure**

# Morphology in NLP

- NLP problems that address morphology:

lemmatization

stemming

# Terminology

- **lemma**

- canonical/dictionary form of a word
- words with same lemma have same stem, part of speech, rough semantics

- **word form**

- full inflected or derived form of a word as it appears in text

| word form | lemma |
|-----------|-------|
| run       | run   |
| ran       | run   |
| running   | run   |

# Lemmatization

- **lemmatization**: convert wordform to lemma

am, is, are → be

car, cars, car's, cars' → car

the boy's cars are different colors



the boy car be different color

- mostly about finding the correct dictionary entry,  
but this may depend on the context

# Stemming

- **stemming**: reduce words to their stems by removing affixes
  - usually implemented with language-specific, manually-designed rules
  - commonly used in information retrieval
  - example:

Caillou is an average, imaginative four-year-old boy with a love for forms of transportive machinery such as rocket ships and airplanes.



Caillou is an **averag imagin** four year old **boi** with a love for **form** of **transport machineri** such as rocket **ship** and **airplan**

# Porter's Algorithm

(the most common English stemmer)

## Step 1a

|      |      |          |          |
|------|------|----------|----------|
| sses | → ss | caresses | → caress |
| ies  | → i  | ponies   | → poni   |
| ss   | → ss | caress   | → caress |
| s    | → ∅  | cats     | → cat    |

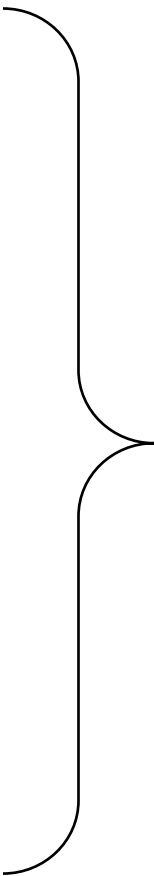
## Step 1b

|          |     |           |           |
|----------|-----|-----------|-----------|
| (*v*)ing | → ∅ | walking   | → walk    |
|          |     | sing      | → sing    |
| (*v*)ed  | → ∅ | plastered | → plaster |
| ...      |     |           |           |




# Idiosyncrasies of the Porter Stemmer

sever  
severed  
severing  
several  
severe  
severely  
severity



→ sever

wit  
wits  
witness  
witnesses  
witnessing



→ wit

# Lemmatization vs. Stemming

- Lemmatization
  - viewed as an NLP task
  - solved with dictionary look-up, possibly with machine learning
- Stemming uses manually-defined rules
  - simple and fast, but limited due to reliance on rules
  - may conflate words erroneously
- Both may remove information
  - To mitigate, combine lemma/stem form with original form, use both!

# Data-Driven Segmenters

- Segment words into pieces (**subword units** or **wordpieces**) based on common character sequences in a dataset
- Most popular methods:
  - Byte pair encoding (BPE)
  - SentencePiece's unigram language model (LM)
- These are efficient and effective, but they don't necessarily correspond to morphology; splits may be arbitrary
- Very popular when using neural networks (BERT, GPT, etc)

# Byte Pair Encoding (BPE)

(Gage, 1994)

- Simple data compression technique
- Iteratively replaces most frequent pair of bytes in a sequence with a single, unused byte
- Sennrich et al. (2016) adapted BPE for segmenting words

# Byte Pair Encoding for Words

- “merge”: operation that combines two consecutive units into a single unit
  - initially, units are characters (e.g., `s` or `t`)
  - after merges, units become character sequences (e.g., `st` or `books`)
- greedy algorithm:
  - merge 2 units with the largest 2-unit sequence count, produce merged unit
  - replace all instances of that 2-unit sequence with the merged unit, recompute counts

# Example from movie review dataset (Stanford Sentiment Treebank):

word that was not in training set:

writer/director/producer

BPE segmenter based on training set



writ@@ er@@ /@@ direct@@ or@@ /@@  
producer

(recover original text by removing “@@ ”)

It wouldn't be my preferred way of spending 100  
minutes or \$ 7.00 .

likely good: “prefer” is the lemma of “preferred”



It wouldn't be my prefer@@ red way of sp@@ ending 100  
minutes or \$ 7@@ .@@ 00 .

maybe bad: “spending” is not related to “ending”



# Summary (1/3)

- to do NLP on some text, we need to preprocess it:
  - tokenize documents into sentences
  - tokenize sentences into tokens
- rule-based tokenizers exist for many languages
- for writing systems without whitespace, tokenization becomes complex



## Summary (2/3)

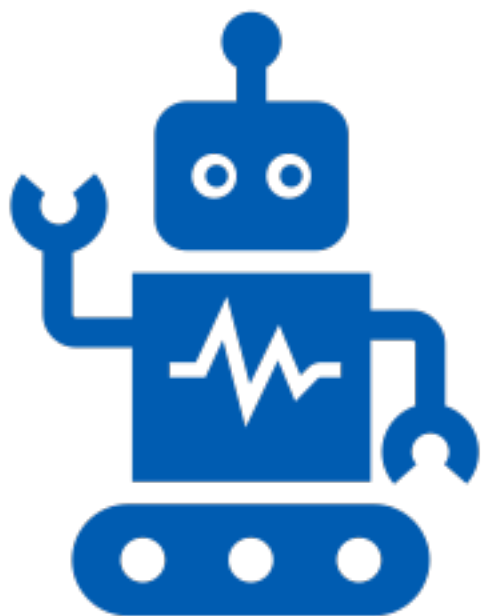
- useful terms: **type**, **token**, **type/token ratio**
- when adding data, number of types keeps increasing
- most types are extremely rare (Zipf's law)
- people can often understand a novel word type based on its internal structure and the context of its use

# Summary (3/3)

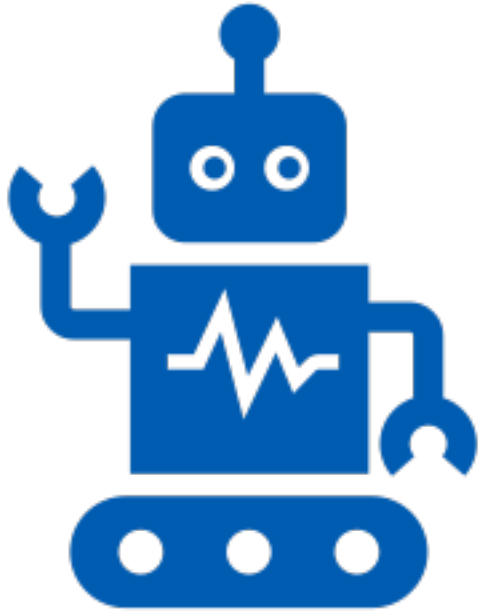
- **morphology**: study of how words are built from morphemes
- **lemmatization**: convert wordform to lemma (may depend on context)
- **stemming**: removing affixes from words to get stems (simple, rule-based)
- data-driven segmentation (BPE, unigram LM) splits words based on data, very common in deep learning

# How do represent words computationally?

荃者所以在鱼，得鱼而忘荃    Nets are for fish;  
Once you get the fish, you can forget the net.  
言者所以在意，得意而忘言    Words are for meaning;  
Once you get the meaning, you can forget the words  
庄子(Zhuangzi), Chapter 26



|     |      |         |       |        |             |
|-----|------|---------|-------|--------|-------------|
| cat | chef | chicken | civic | cooked | council ... |
| ↓   | ↓    | ↓       | ↓     | ↓      | ↓           |
| 17  | 91   | 253     | 104   | 5      | 6001 ...    |



cat    chef    chicken    civic    cooked    council ...

|      |      |      |      |      |      |
|------|------|------|------|------|------|
| ↓    | ↓    | ↓    | ↓    | ↓    | ↓    |
| 0.1  | -0.1 | -0.4 | 0.1  | -0.5 | 0.6  |
| 7.9  | 2.1  | 2.4  | 0    | -1.1 | -1.3 |
| 2.4  | 3.8  | 9.7  | -1.5 | 7.6  | 0    |
| -1.3 | -0.1 | -1.0 | 2.4  | -3.1 | 3.4  |
| 0.5  | 5.3  | 3.2  | 0.2  | 4.2  | -0.6 |

**“embeddings”**

# Motivation for word embeddings?

**Variability**

multiple forms,  
similar meaning

really really reallly



$\begin{bmatrix} 2.1 \\ -7.9 \\ 8.4 \\ -1.3 \end{bmatrix}$



$\begin{bmatrix} 2.3 \\ -6.1 \\ 7.8 \\ -0.8 \end{bmatrix}$



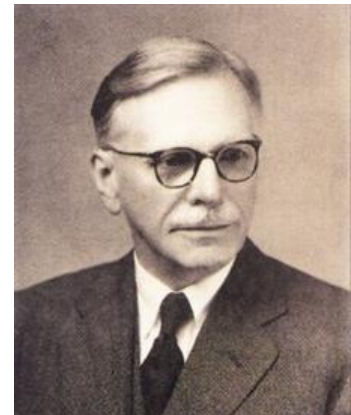
$\begin{bmatrix} 1.9 \\ -6.8 \\ 7.7 \\ -1.0 \end{bmatrix}$

How should we embed words?

**distributional hypothesis**: words that appear in similar contexts have similar meanings (Joos, 1950; Harris, 1954; Firth, 1957)

“You shall know a word by the company it keeps.”

–Firth (1957)



words we are computing vectors for:

|                   |         | cat | chef | chicken | civic | cooked | council |
|-------------------|---------|-----|------|---------|-------|--------|---------|
| context<br>words: | the     |     |      |         |       |        |         |
|                   | cat     |     |      |         |       |        |         |
|                   | chicken |     |      |         |       |        |         |
|                   | city    |     |      |         |       |        |         |
|                   | cook    |     |      |         |       |        |         |



... , the club may also employ a **chef** to prepare and cook food items .

... is up to remy , linguini , and the **chef** colette to cook for many people ...

... cooking program the cook and the **chef** with simon bryant , who is ...

|         | chef |
|---------|------|
| the     | 0    |
| cat     | 0    |
| chicken | 0    |
| city    | 0    |
| cook    | 0    |

... , the club may also employ **a chef to** prepare and cook food items .

... is up to remy , linguini , and **the chef colette** to cook for many people ...

... cooking program the cook and **the chef with** simon bryant , who is ...

window size:  $w = 1$

|         | chef |
|---------|------|
| the     | 2    |
| cat     | 0    |
| chicken | 0    |
| city    | 0    |
| cook    | 0    |

... , the club may also employ a **chef** to prepare and cook food items .

... is up to remy , linguini , and the **chef** colette to cook for many people ...

... cooking program the cook and the **chef** with simon bryant , who is ...

window size:  $w = 4$

|         | chef |
|---------|------|
| the     | 3    |
| cat     | 0    |
| chicken | 0    |
| city    | 0    |
| cook    | 3    |

words we are computing vectors for:

|                   |         | cat   | chef | chicken | civic | cooked | council |
|-------------------|---------|-------|------|---------|-------|--------|---------|
| context<br>words: | the     | 24708 | 7410 | 7853    | 16486 | 3463   | 316380  |
|                   | cat     | 2336  | 14   | 23      | 0     | 1      | 36      |
|                   | chicken | 23    | 21   | 1640    | 1     | 181    | 7       |
|                   | city    | 116   | 89   | 62      | 943   | 7      | 27033   |
|                   | cook    | 12    | 113  | 34      | 6     | 34     | 51      |

- once we have word vectors, we can compute similarities!
- many ways to define similarity of two vectors
- a simple way: **dot product** (also called inner product):

$$\mathbf{u} \cdot \mathbf{v} = \mathbf{u}^\top \mathbf{v} = \sum_i u_i v_i$$

$\mathbf{u}$  = a vector

$u_i$  = entry i in the vector

- dot product is large when the vectors have very large values in the same dimensions

with dot product as similarity function, let's find the most similar words ("nearest neighbors") to each word:

| nearest<br>neighbors | cat     | chef    | chicken | civic   | cooked  | council |
|----------------------|---------|---------|---------|---------|---------|---------|
|                      | council | council | council | council | council | council |
|                      | cat     | cat     | cat     | cat     | cat     | cat     |
|                      | civic   | civic   | civic   | civic   | civic   | civic   |
|                      | chicken | chicken | chicken | chicken | chicken | chicken |
|                      | chef    | chef    | chef    | chef    | chef    | chef    |
|                      | cooked  | cooked  | cooked  | cooked  | cooked  | cooked  |

with dot product as similarity function, let's find the words ("nearest neighbors") to each word:

| council |        |
|---------|--------|
| the     | 316380 |
| cat     | 36     |
| chicken | 7      |
| city    | 27033  |
| cook    | 51     |

|   | chef    | chicken | civic   | cooked  | council |
|---|---------|---------|---------|---------|---------|
| l | council | council | council | council | council |
|   | cat     | cat     | cat     | cat     | cat     |
|   | civic   | civic   | civic   | civic   | civic   |
| n | chicken | chicken | chicken | chicken | chicken |
|   | chef    | chef    | chef    | chef    | chef    |
|   | cooked  | cooked  | cooked  | cooked  | cooked  |

- dot product is large when vectors have large values in same dimensions, doesn't control for vector length

- vector length:  $\|\mathbf{u}\| = \sqrt{\sum_i u_i^2}$

$\mathbf{u}$  = a vector

$u_i$  = entry i in the vector

- **cosine similarity:**  $\frac{\mathbf{u}^\top \mathbf{v}}{\|\mathbf{u}\| \|\mathbf{v}\|}$

this is the cosine of the angle between the two vectors!



now using cosine similarity:

nearest  
neighbors

| cat     | chef    | chicken | civic   | cooked  | council |
|---------|---------|---------|---------|---------|---------|
| cat     | chef    | chicken | civic   | cooked  | council |
| chef    | civic   | cooked  | council | chef    | civic   |
| cooked  | cooked  | chef    | chef    | civic   | chef    |
| civic   | council | civic   | cooked  | council | cooked  |
| council | cat     | council | cat     | cat     | cat     |
| chicken | chicken | cat     | chicken | chicken | chicken |

# Improving Counts

- counts of common words (“the”) are very large, but not very useful
- many ways proposed for improving raw counts

# Improving Counts

- counts of common words (“the”) are very large, but not very useful
- many ways proposed for improving raw counts
- **tf-idf** (“term frequency-inverse document frequency”)
  - **tf** is just the count of the two words
  - **idf** is computed for the context word
    - document frequency = number of documents in which a word appears
    - inverse of document frequency will be small for common words
  - vector entry is then the product of tf and idf

# Pointwise Mutual Information (PMI)

- consider two random variables,  $X$  and  $Y$
- do two events  $X = x$  and  $Y = y$  occur together more often than if they were independent?

$$\text{pmi}(x, y) = \log_2 \frac{p_{X,Y}(x, y)}{p_X(x) p_Y(y)}$$

- if they are independent,  $\text{PMI} = 0$

# PMI for Word Vectors

- for word vectors,  
     $X$  is the **center word**  
     $Y$  is the **context word**

$$\text{pmi}(x, y) = \log_2 \frac{p_{X,Y}(x, y)}{p_X(x) p_Y(y)}$$

- each probability can be estimated using counts we already computed!

$\#(x, y)$  = co-occurrence count of  $x$  and  $y$

$N$  = total count

$$p_{X,Y}(x, y) = \frac{\#(x, y)}{N}$$

$$p_X(x) = \frac{\sum_y \#(x, y)}{N}$$

$$p_Y(y) = \frac{\sum_x \#(x, y)}{N}$$

# Top co-occurrence counts with “chicken”

|       |         |      |       |     |       |
|-------|---------|------|-------|-----|-------|
| 14464 | ,       | 1525 | or    | 508 | pork  |
| 7853  | the     | 1225 | for   | 500 | meat  |
| 6276  | and     | 1061 | 's    | 481 | be    |
| 5931  | .       | 940  | fried | 479 | he    |
| 5213  | a       | 906  | on    | 452 | such  |
| 3963  | of      | 889  | was   | 445 | his   |
| 3282  | in      | 869  | that  | 417 | at    |
| 2520  | to      | 828  | are   | 405 | soup  |
| 2438  | "       | 777  | by    | 389 | made  |
| 2339  | is      | 746  | from  | 384 | rice  |
| 2127  | with    | 710  | it    | 375 | but   |
| 1818  | (       | 600  | beef  | 350 | has   |
| 1745  | )       | 590  | which | 330 | fish  |
| 1640  | chicken | 557  | also  | 325 | other |
| 1594  | as      | 531  | an    | 318 | this  |

# Words with largest PMI with “chicken”

|      |          |     |            |     |           |
|------|----------|-----|------------|-----|-----------|
| 10.2 | fried    | 7.0 | robot      | 6.1 | pig       |
| 9.7  | chicken  | 6.9 | burger     | 6.0 | breeds    |
| 9.3  | pork     | 6.8 | recipe     | 6.0 | vegetable |
| 9.0  | beef     | 6.6 | vegetables | 6.0 | potato    |
| 8.7  | soup     | 6.6 | potatoes   | 5.9 | goose     |
| 7.8  | sauce    | 6.6 | goat       | 5.9 | dixie     |
| 7.7  | curry    | 6.5 | eggs       | 5.9 | kung      |
| 7.6  | cooked   | 6.4 | cow        | 5.9 | pie       |
| 7.5  | lamb     | 6.4 | pizza      | 5.8 | menu      |
| 7.4  | dish     | 6.4 | rice       | 5.8 | steamed   |
| 7.3  | shrimp   | 6.3 | ribs       | 5.8 | tastes    |
| 7.3  | egg      | 6.3 | tomatoes   | 5.7 | beans     |
| 7.2  | sandwich | 6.2 | cheese     | 5.7 | butter    |
| 7.2  | dishes   | 6.2 | duck       | 5.7 | barn      |
| 7.2  | meat     | 6.1 | chili      | 5.7 | breed     |

words we are computing vectors for:

|                   |         | cat   | chef | chicken | civic | cooked | council |
|-------------------|---------|-------|------|---------|-------|--------|---------|
| context<br>words: | the     | 24708 | 7410 | 7853    | 16486 | 3463   | 316380  |
|                   | cat     | 2336  | 14   | 23      | 0     | 1      | 36      |
|                   | chicken | 23    | 21   | 1640    | 1     | 181    | 7       |
|                   | city    | 116   | 89   | 62      | 943   | 7      | 27033   |
|                   | cook    | 12    | 113  | 34      | 6     | 34     | 51      |

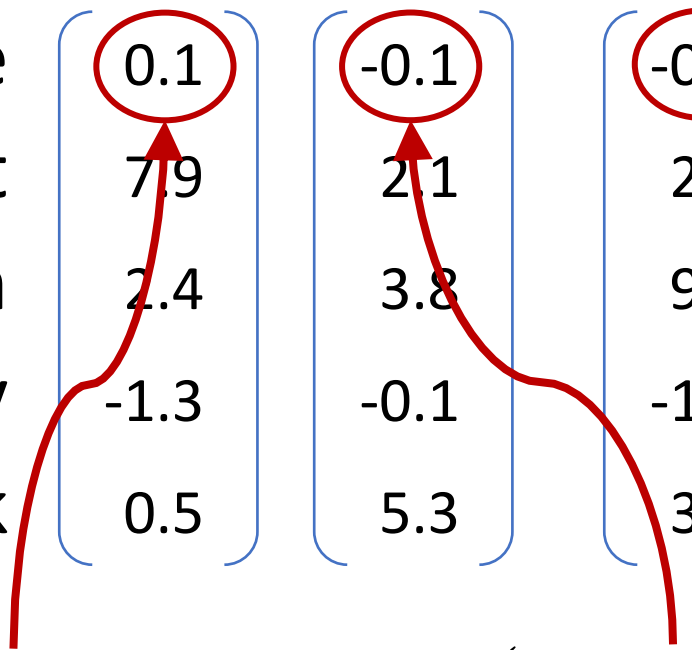


words we are computing vectors for:

context words:

|         | cat  | chef | chicken | civic | cooked | council |
|---------|------|------|---------|-------|--------|---------|
| the     | 0.1  | -0.1 | -0.4    | 0.1   | -0.5   | 0.6     |
| cat     | 7.9  | 2.1  | 2.4     | 0     | -1.1   | -1.3    |
| chicken | 2.4  | 3.8  | 9.7     | -1.5  | 7.6    | 0       |
| city    | -1.3 | -0.1 | -1.0    | 2.4   | -3.1   | 3.4     |
| cook    | 0.5  | 5.3  | 3.2     | 0.2   | 4.2    | -0.6    |

pmi(cat, the) pmi(chef, the) ...



using counts:

nearest  
neighbors

| cat     | chef    | chicken | civic   | cooked  | council |
|---------|---------|---------|---------|---------|---------|
| cat     | chef    | chicken | civic   | cooked  | council |
| chef    | civic   | cooked  | council | chef    | civic   |
| cooked  | cooked  | chef    | chef    | civic   | chef    |
| civic   | council | civic   | cooked  | council | cooked  |
| council | cat     | council | cat     | cat     | cat     |
| chicken | chicken | cat     | chicken | chicken | chicken |

using PMIs:

nearest  
neighbors

| cat     | chef    | chicken | civic   | cooked  | council |
|---------|---------|---------|---------|---------|---------|
| cat     | chef    | chicken | civic   | cooked  | council |
| chicken | chicken | cooked  | council | chicken | civic   |
| chef    | cooked  | chef    | chef    | chef    | chicken |
| cooked  | cat     | cat     | cat     | cat     | chef    |
| civic   | council | council | chicken | council | cooked  |
| council | civic   | civic   | cooked  | civic   | cat     |

# Positive PMI?

- some have found benefit by truncating PMI at 0 (“positive PMI”)
- negative PMI: words occur together less than we would expect, i.e., they are anticorrelated
- these anticorrelations may need more data to reliably estimate
- however, negative PMIs do seem reasonable!

## Largest PMIs:

|      |          |
|------|----------|
| 10.2 | fried    |
| 9.7  | chicken  |
| 9.3  | pork     |
| 9.0  | beef     |
| 8.7  | soup     |
| 7.8  | sauce    |
| 7.7  | curry    |
| 7.6  | cooked   |
| 7.5  | lamb     |
| 7.4  | dish     |
| 7.3  | shrimp   |
| 7.3  | egg      |
| 7.2  | sandwich |

## PMIs close to zero:

|        |              |
|--------|--------------|
| 0.003  | climbed      |
| 0.003  | detailing    |
| 0.002  | turkish      |
| 0.002  | oaks         |
| 0.001  | productivity |
| 0.000  | swing        |
| -0.001 | structures   |
| -0.001 | thirteenth   |
| -0.001 | commentators |
| -0.001 | palmer       |
| -0.002 | obstacles    |
| -0.003 | horns        |
| -0.003 | burning      |

## Smallest PMIs:

|      |              |
|------|--------------|
| -4.6 | users        |
| -4.6 | data         |
| -4.7 | discussion   |
| -4.7 | museum       |
| -4.7 | below        |
| -4.8 | editors      |
| -4.8 | railway      |
| -4.8 | committee    |
| -4.8 | elected      |
| -4.9 | championship |
| -5.0 | archive      |
| -5.3 | edits        |
| -6.1 | deletion     |

# Summary (1/2)

- **word vectors**: the use of a vector of numbers to represent a word
- **distributional word vectors**: the use of distributional statistics (e.g., word co-occurrence counts) in defining word vectors
- different window sizes lead to different kinds of similarity in the vectors
- many options for computing similarity of two vectors:
  - **dot product** is a simple starting point
  - **cosine similarity** accounts for vector length and works better for word vectors
- simple distributional vectors are ok, but modifying counts can help

## Summary (2/2)

- one method: scale down counts for common context words, usually with a variant of **inverse document frequency (idf)**
- another method: **pointwise mutual information (PMI)**
- PMI provides information about whether two events are independent
- view center and context word slots as random variables; words are events/outcomes of those random variables
- use  $\text{PMI}(\text{center word}, \text{context word})$  values in place of counts for better word vectors