

CS 489/698: Introduction to Natural Language Processing

Lecture 2: Words and Morphology

Instructor: Freda Shi

rhs@uwaterloo.ca

January 7th, 2026

What is a word?

Lexical Semantics

“A single distinct **meaningful** element of speech or writing, used with others (or sometimes alone) to form a sentence and typically shown with a space on either side when written or printed.”

[Source: Oxford Languages]

What is a word?

Things in dictionaries?

“One of the most prolific areas of change and variation in English is vocabulary; **new words** are constantly being coined to name or describe new inventions or innovations, or to better identify aspects of our rapidly changing world... Most general English dictionaries are designed to include only those words that **meet certain criteria** of usage across wide areas and over extended periods of time...”

[Source: Merriam-Webster; <https://www.merriam-webster.com>]

What is a word?

Things between spaces and punctuation?

This is English: The cat is cute.

This is French: Le chat est mignon.

This is Spanish: El gato es lindo.

This is Chinese: 猫很可爱。

This is Japanese: 猫はかわいい。

This is Thai: แมวน่ารักจัง.

What is a word?

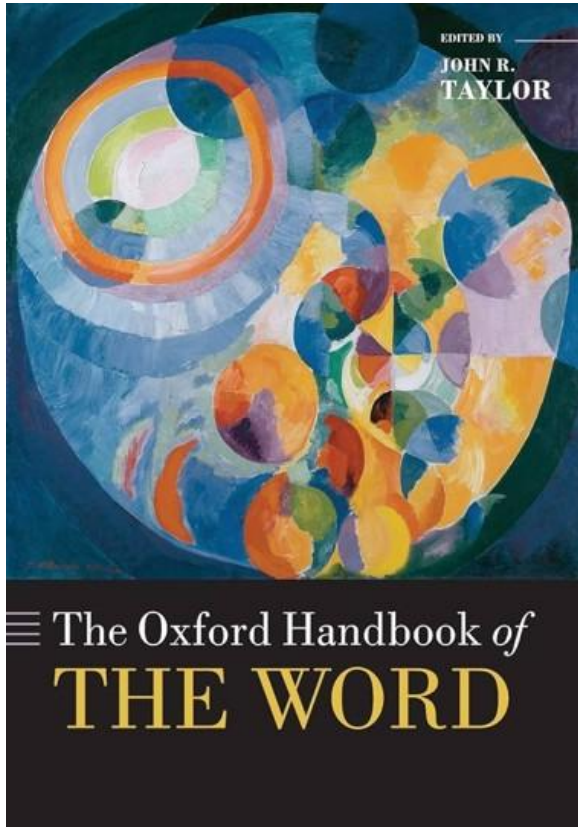
Smallest unit that can be uttered in isolation?

- You could utter this word in isolation: *unimpressively*.
- Also this one: *impress*.
- Probably also these when you talk about morphology: *un*, *ive*, *ly*.

Are they all words?

What is a word?

Each of the above points captures some, but very likely not all aspects of what a word is.



42 chapters.

Nearly 900 pages.

Covers a lot of aspects of what makes a word word, “to anyone who shares a fascination with words.”

Outline of Today's Lecture

Linguistic morphology

- The study of internal structures of words.

Lexical semantics

- The study of meanings of words.

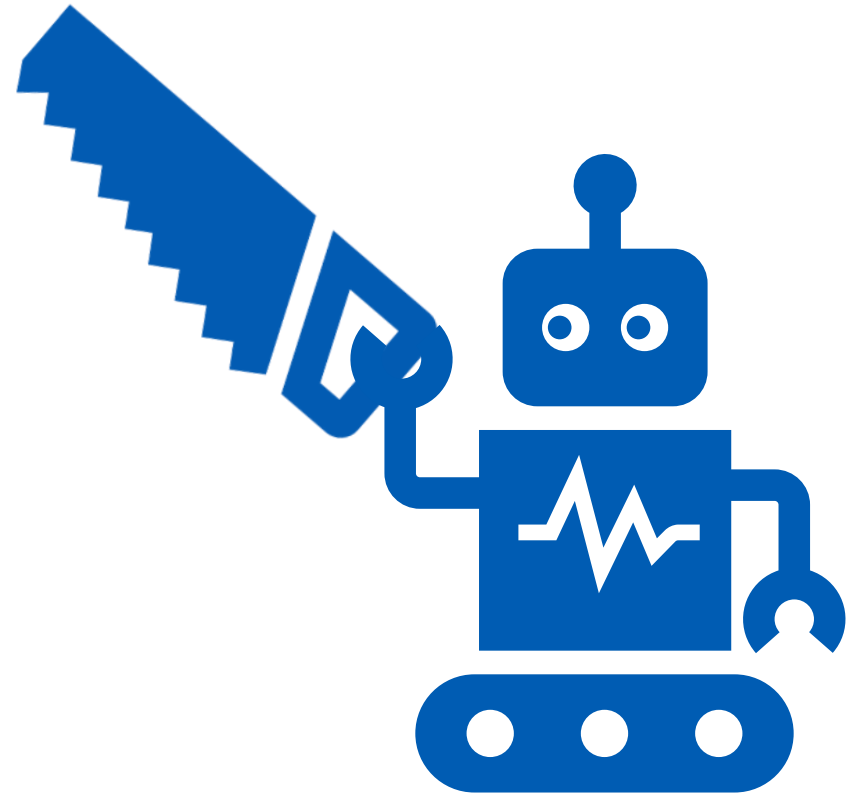
Word tokenization

- The process of splitting texts into “words” (tokens).

colder

replayed

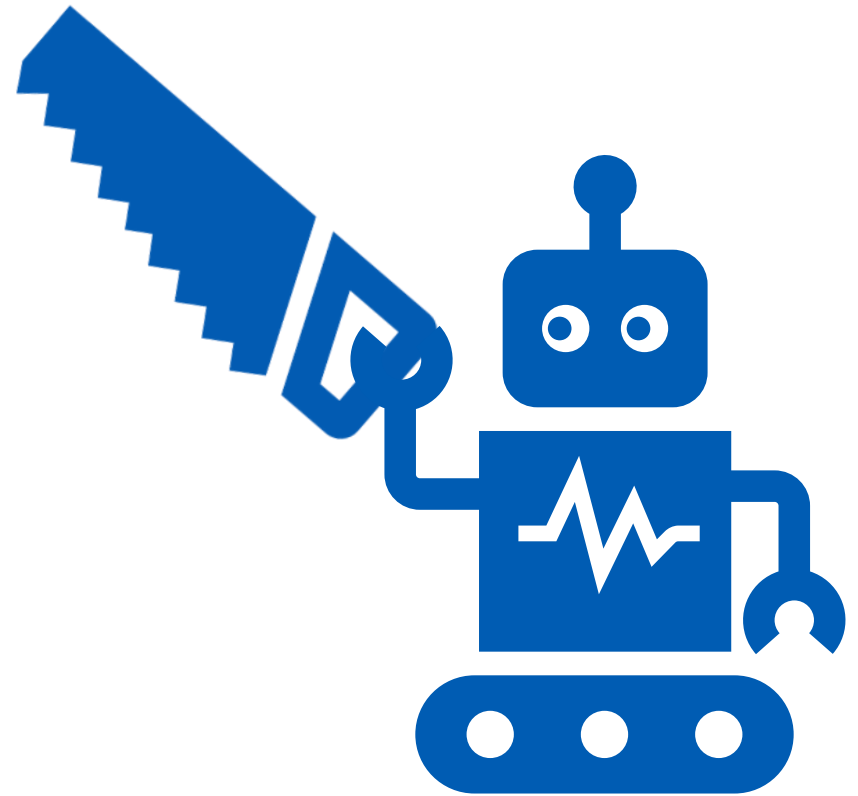
gameplay



cold|er

re|play|ed

game|play



Morphology

The study of how words are built from smaller meaning-bearing units.

Types of morphemes:

- **Stem**: a core meaning-bearing unit.
- **Affix**: a piece that attaches to a stem, adding some function or meaning.

Prefix: **un**happy, **pre**define

Suffix: cat**s**, walk**ed**

Infix: (Malay) Gigi (*teeth*) → **Ger**igi (*tooth blade*)

Circumfix: (German) mach (root of *machen*; *to make*) → **ge**mach**t** (*made*; *past participle*)

Interfix: speed**o**meter

[See more in Chap. 6.2 of Doner. The Linguistic analysis of word and sentences structures]

Types of Word Formation

Inflection: adding morphemes to a word to indicate grammatical information.

- *walk* → *walk***ed**
- *cat* → *cat***s**

Derivation: adding morphemes to a word to create a new word with a different meaning.

- *happy* → *happi***ness**
- *define* → *pre***define**

Compounding: combining two or more words to create a new word.

- *key* + *board* → *keyboard*
- *law* + *suit* → *lawsuit*
- *book* + *case* → *bookcase*

Isolating Language

In languages like Classical Chinese, Vietnamese, and Thai

- Each word form typically consists of one single morpheme.
- There is little morphology other than compounding.

➤ Inflection

们: 我们, 你们, 他们

mén: wǒmén, nǐmén, tāmén

plural: we, you (pl.), they

➤ Derivation

家: 艺术家

jiā: yì shù jiā, artist

➤ Compound

高
gāo

+

地

ground, land
земля

档

grade, quality
сорт, качество

速

speed
скорость

=

高地

gāodì
highland
возвышенность

高档

gāodàng
high quality
высококачественный

高速

gāosù
high speed
скоростной

Morphological Decomposition

Usually, morphological decomposition is simply splitting a word into its morphemes:

walked = walk + ed

greatness = great + ness

But it can actually be a hierarchical structure:

unbreakable = un + (break + able)

internationalization = (((inter + nation) + al) + iz[e]) + tion

There is ambiguity in hierarchical decomposition!

The door is unlockable.

Morphology in NLP

Individual tasks that address morphology:

- **Lemmatization**: putting words/tokens in a standard format.
 - **Lemma**: canonical/dictionary form of a word.
 - **Wordform**: fully inflected or derived form of a word as it appears in text.

wordform	lemma
run	run
ran	run
running	run

Morphology in NLP

Individual tasks that address morphology:

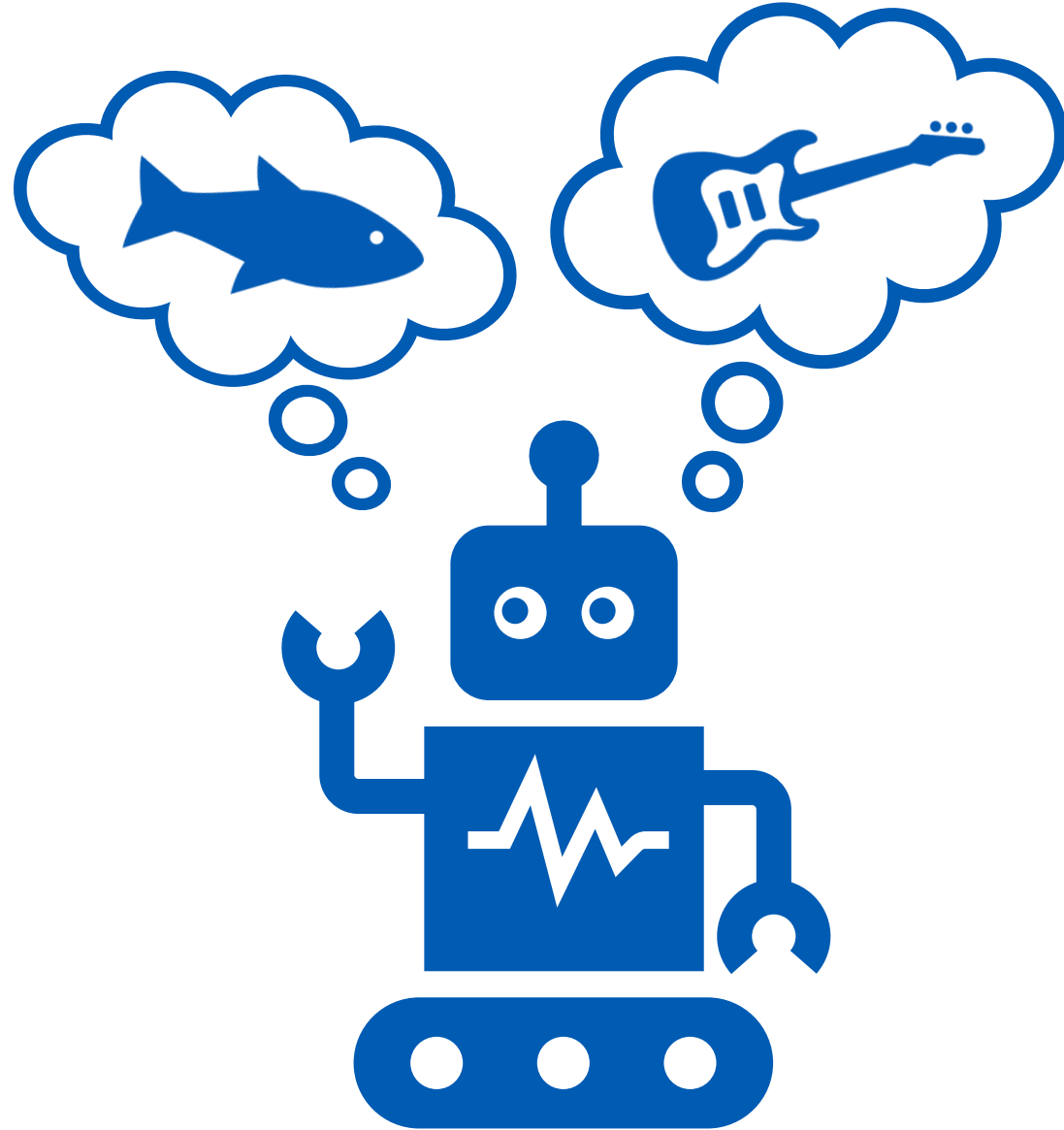
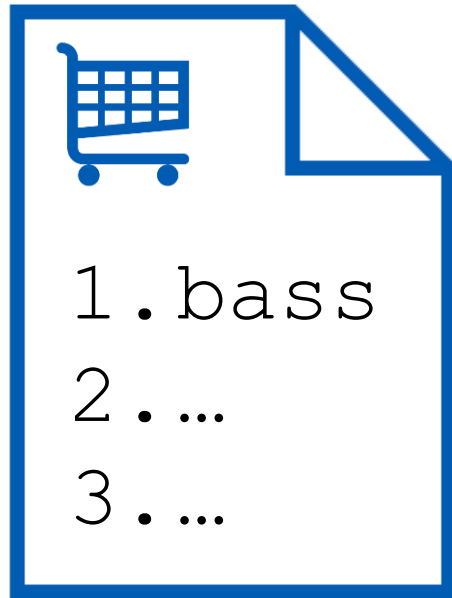
- Stemming: reducing words to their stems (approximately) by removing affixes. More conventional engineering-oriented approach used in applications such as retrieval.

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**

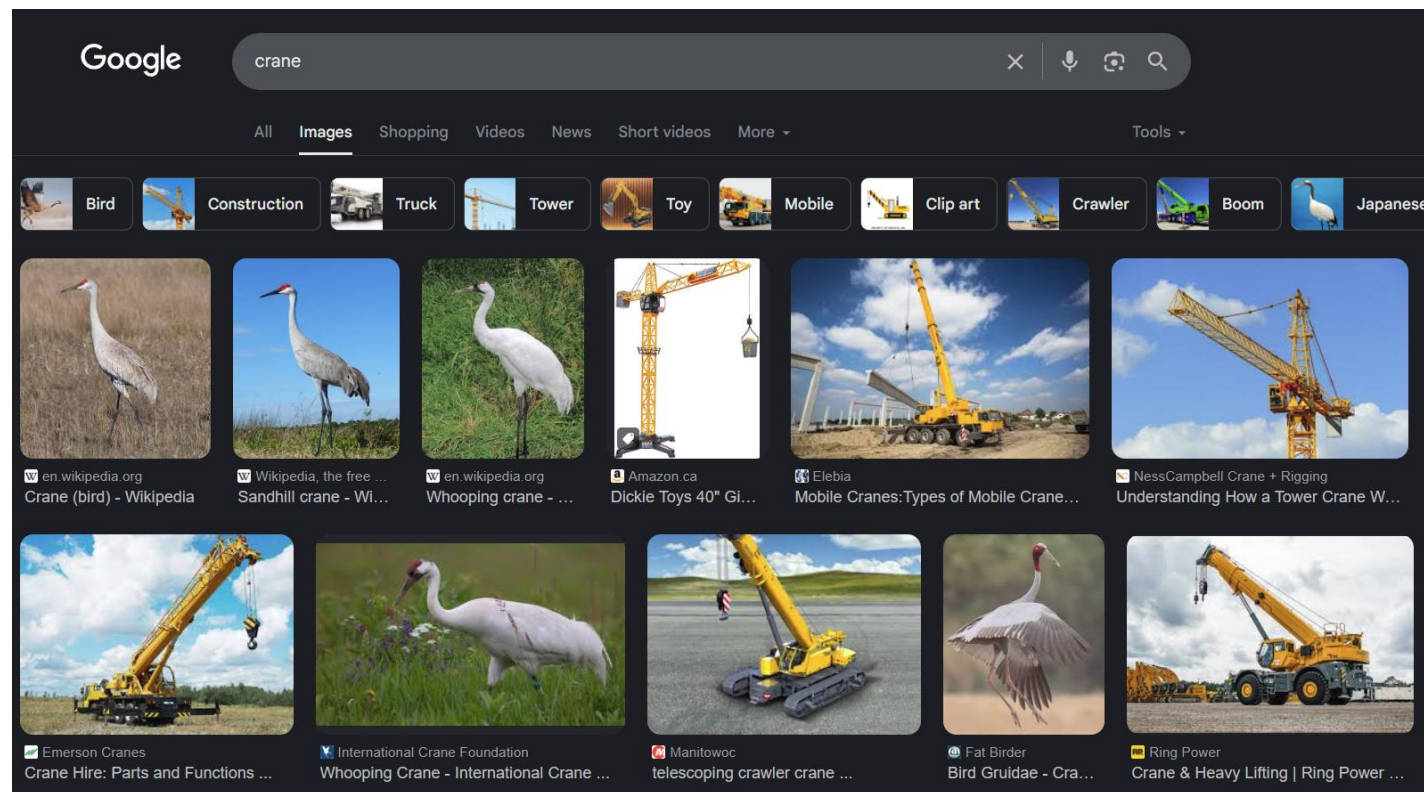
Lexical Semantics



Variability and Ambiguity in Words

Lemmatization and stemming tackles the problem of variability---multiple forms could share the same or similar meanings.

On the other hand, one wordform could refer to multiple meanings.



Polysemy vs. Homonymy

- **Polysemy**: a word has multiple **related** meanings.
 - She is a **star**.
 - (pointing to the sky) The **star** is shining.
- **Homonymy**: a word has multiple meanings originated from different sources.
 - I need to go to the **bank** as I don't have enough cash.
 - I am sitting on the **bank** of the river.

Question: which one do you think is the case for crane?

Synonyms

Synonyms (informal definition): words that have the same meanings according to some criteria.

- *couch vs. sofa*
- *big vs. large*
- *water vs. H_2O*
- There are very few (or no) examples of perfect synonymy.
- Synonymy is a relation between senses rather than words.
- *How big is the plane?*
- *How large is the plane?*
- *Miss Nelson became a kind of big sister to Benjamin.*
- *Miss Nelson became a kind of large sister to Benjamin. (*)*

Antonyms

Antonyms: senses that are opposite with respect to (at least) one dimensionality of meaning.

- *dark and light*
- *dark and bright*
- *hot and cold*
- *in and out*

Hyponymy/Hypernymy, and Meronym/Holonym

- Sense A is a **hyponym** of sense B if A is more specific, denoting a subclass of B.
- Conversely, B is a hypernym of A.
 - *dog* is a hyponym of *animal*
 - *corgi* is a hyponym of *dog*
- Sense A is a meronym of sense B if A is a part of B.
- Conversely, B is a holonym of A.
 - *hand* is a meronym of *body*
 - *Finger* is a meronym of *hand*

The WordNet database: <https://wordnet.princeton.edu>

Word Sense Disambiguation

- Word-Sense Disambiguation (WSD): the task of determining which sense of a word is used in a particular context, given a set of predefined possible senses.
- Relatedly, word sense induction (WSI) requires clustering word usages into senses without predefined ground truths.

Default solution (as of 2026): encode the context of words with a pretrained model, and train a neural network to predict the sense.

Or... prompting a pretrained language model.

The Role of Word Senses

A practical question: We now have powerful neural language models, which do not distinguish word senses. Is WSD still a meaningful task?

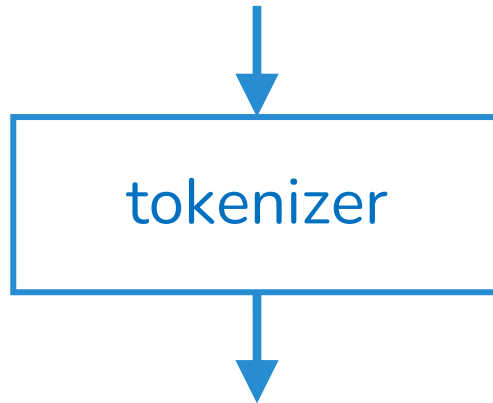
A philosophical question in lexical semantics: Do discrete word senses even “exist”?

[Li. 2024. *Semantic minimalism and the continuous nature of polysemy*. Mind and Language]

Tokenization

- Tokenization: the process that converts running text (i.e., a sequence of characters) into a sequence of tokens.

"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 . "

Conventions in Rule-Based Tokenizers

	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

It is important to check and ensure consistency when comparing results across different tokenizers.

See `nltk.tokenize`, which also works for sentence tokenization.

[<https://www.nltk.org/api/nltk.tokenize.html>]

Tokenization across Languages

There is no explicit whitespace between words in some languages, and tokenization becomes highly nontrivial in these cases.

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

Word Types vs. Word Tokens

" 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

Word Types vs. Word Tokens

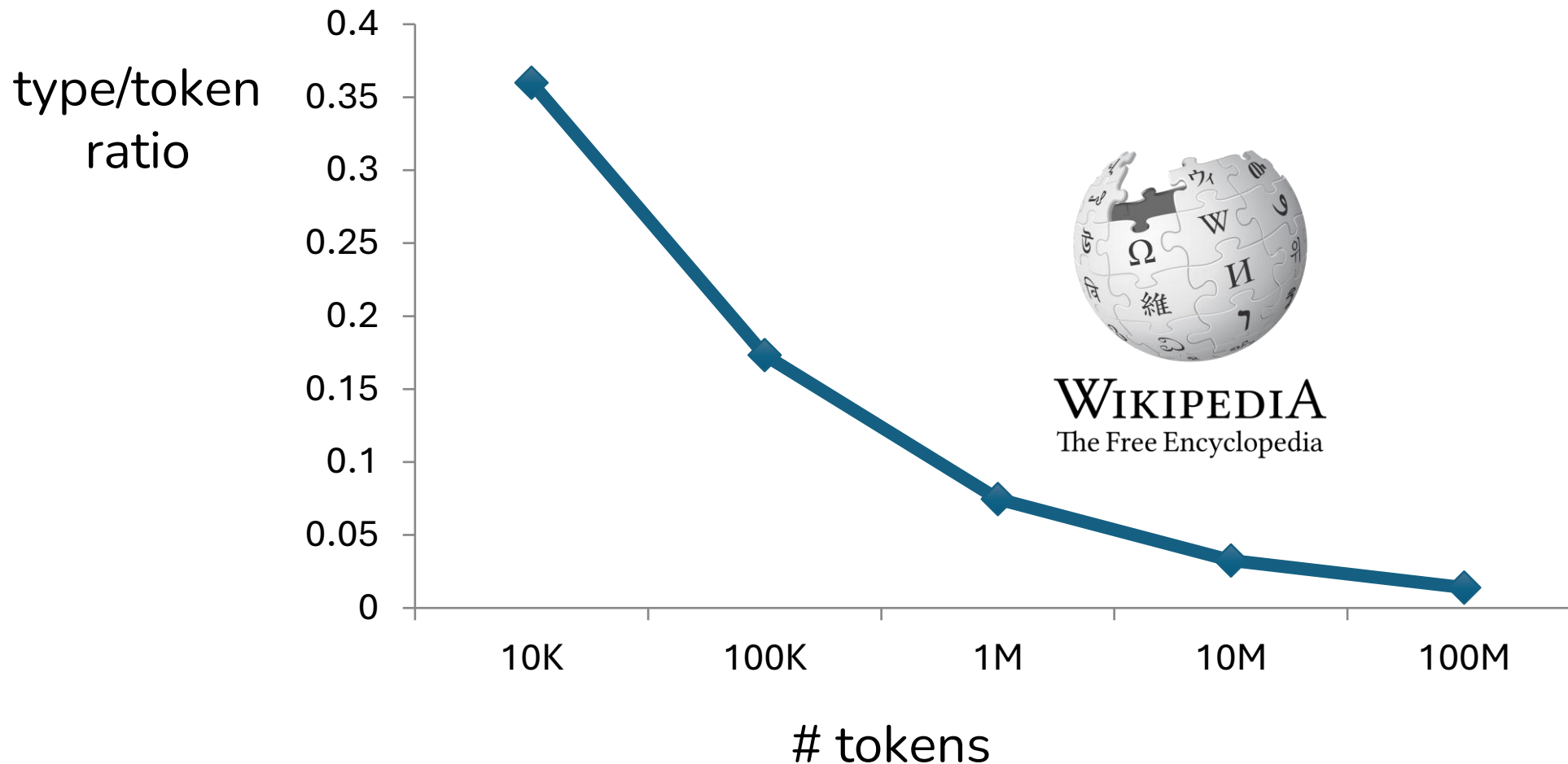
3	i	1	!	1	oh
2	,	1	.	1	said
2	_	1	;	1	stoutly
2	am	1	afraid	1	tallest
2	the	1	for	1	though
2	"	1	lydia	1	youngest
2	"	1	not	1	'm

Type: a unique word (an entry in a vocabulary or dictionary) – 21 types.

Token: an instance of a type in the text – 29 tokens.

Type/Token Ratio

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



Type/Token Ratio: Wikipedia vs. Twitter

How do the type/token ratio curves compare between Wikipedia and Twitter?



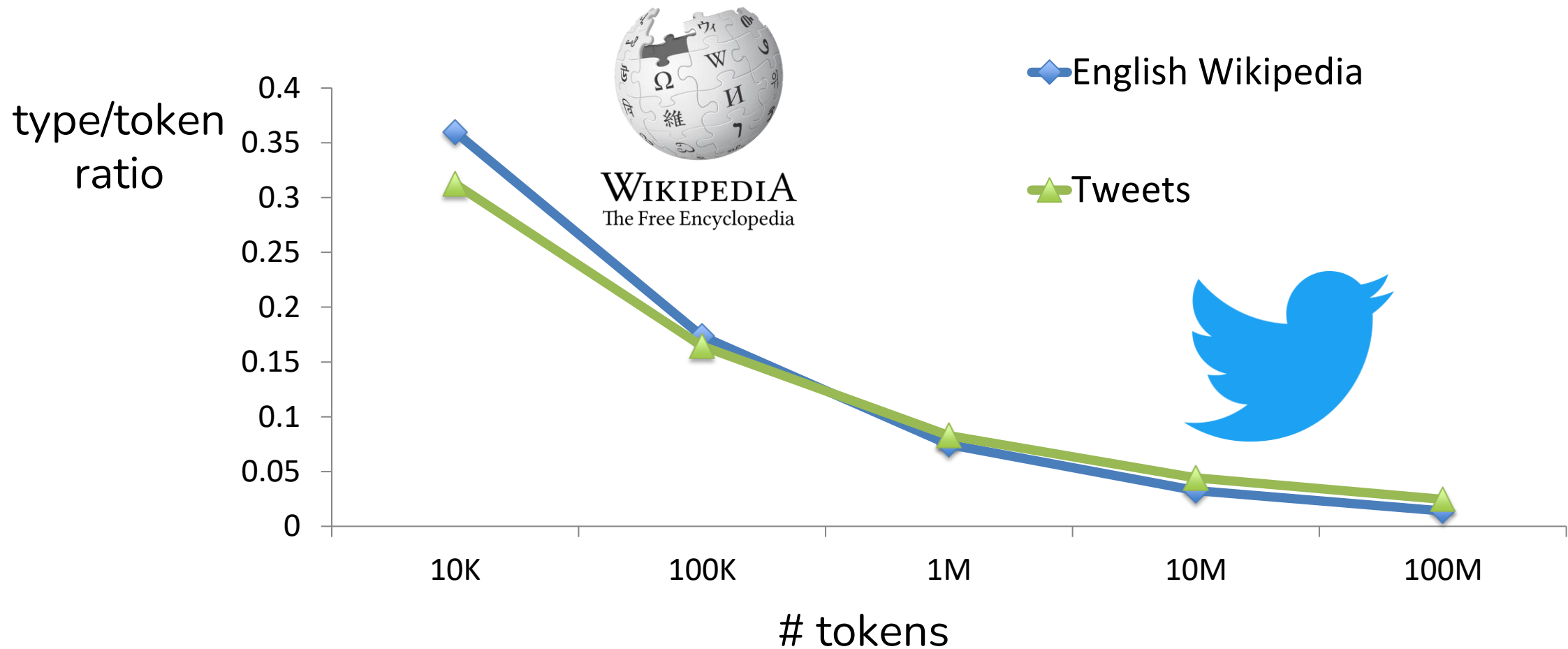
WIKIPEDIA
The Free Encyclopedia

VS.



Type/Token Ratio: Wikipedia vs. Twitter

How do the type/token ratio curves compare between Wikipedia and Twitter?



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 reallllly	27 reaaly	11 reeeeallly
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 realllllllyyyy	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 reeaallllyyy
5 reaaaaalllly	3 r34lly	3 reeaaly
5 reaallllyy	3 r]eally	3 reeaallly
5 realllllllllllllllllly	3 reaaaaaaaally	3 reealy
5 realllllllllllllly	3 reaaaaaly	3 reeeaaallllyyy
5 reeallyyy	3 reaaaalllly	3 reeeaaallly
5 reeeaaaallly	3 reaaaalllyy	3 reeeeaaaaally
5 reeeeaally	3 reaaalllyy	3 reeeeaalllly
5 reeeeeeeeeally	3 reaalllly	3 reeeeealllly
5 relly	3 reaallyyyy	2 reaaaaaaaaaally

2 reaaaaaaaaally	2 really/	2 reeely
2 reaaaaaaaaallly	2 reallyyyyyyyy	2 rellys
2 reaaaaaalllllyyy	2 reallyyyyyyyyyyyy	2 rellyy
2 reaaaaallllly	2 realyyy	2 reqally
2 reaaaaalllly	2 reaqlly	2 rlyyy
2 reaaalllllyyy	2 reeaaally	2 rlyyyy
2 reaaalllllyyy	2 reeaalllly	2 rreeaallyy
2 reaaallllyyy	2 reeaallllyy	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 reeeeeaaaallllyyy	1 raaahhhlllaaayyyy
2 realllllllllllllllllllly	2 reeeeeaaaally	1 raeally
2 realllllllllllllllllly	2 reeeeeaaalllly	1 re-e-e-eally
2 reallllllllyy	2 reeeeeaaallllyyyy	1 re-eaaaaaaly
2 reallllllllyyyyyy	2 reeeeeallllyyy	1 re-he-he-he-ealy
2 reallllllllyyyyyy	2 reeeeeallyyy	1 re-he-he-heeeeeally
2 reallllllyy	2 reeeeeaaaallllly	1 re3ally
2 reallllllyyyyyyy	2 reeeeeaaaally	1 rea(1)ly
2 reallllyyyyyyy	2 reeeeeeaally	1 reaaaaaaaaaaaaaaaaaally
2 reallyyyyyyy	2 reeeeeaally	1 reaaaaaaaaaaaaaaaaaally
2 really*	2 reeeeeeealy	1 reaaaaaaaaaaaaaaaaaalllly

[illegible]

1 reallyyyyyyyyyyy	1 reeeaaallllyyy	1 reeeeaalllly
1 really🙄	1 reeeaaallyy	1 reeeeaallly
1 realoly	1 reeeaaalllyyy	1 reeeeaalllyyy
1 realys	1 reeeaaaly	1 reeeeaalllyyyy
1 realyyyyy	1 reeeaaalllly	1 reeeeeeaaaaaalllllllly
1 real•ly	1 reeeaaallyy	1 reeeeeeaaaaaally
1 reawly	1 reeeaaalllyyy	1 reeeeeeaaalllly
1 ree-hee-heally	1 reeeaallllly	1 reeeeeeaaallllyyy
1 reeaaaaaaaaaaaaaaaaaallllllllllyyy	1 reeeaalllly	1 reeeeeeaaallly
1 reeaaaaaaaaalllllllly	1 reeeaallllyy	1 reeeeeeaallllyyy
1 reeaaaaalllllly	1 reeeaallllyyy	1 reeeeeealy
1 reeaaaallly	1 reeeaallllyyyy	1 reeeeeeaaaaaally
1 reeaaalllly	1 reeeaallllyyy	1 reeeeeeaaaaaalllllyyyy
1 reeaaallllyyy	1 reeeaallllyyyy	1 reeeeeeaaaaaally
1 reeaaallyyy	1 reeeallys	1 reeeeeeaaaalllly
1 reeaaaly	1 reeeeaaaaaaalllllllyyyyyyy	1 reeeeeeaaally
1 reeaalllllyyy	1 reeeeaaaaaaalllllllyyyy	1 reeeeeeaaaly
1 reeaallllyyy	1 reeeeaaaaaallllllly	1 reeeeeeaallllly
1 reeaallyyy	1 reeeeaaaaaalllly	1 reeeeeeaallllyyy
1 reeaallly	1 reeeeaaaaaallly	1 reeeeeeaaaaallly
1 reeaalllyyy	1 reeeeaaaaaalllyyy	1 reeeeeeaaally
1 reeaalllyyyy	1 reeeeaaaaaallly	1 reeeeeeaaallyy
1 reeeaaaaaaaly	1 reeeeaaaaalllllyyyyyy	1 reeeeeeaaaly
1 reeeaaaaaally	1 reeeeaaaaallllyy	1 reeeeeeaaalllly
1 reeeaaaallllly	1 reeeeaaaaalllyy	1 reeeeeeaaalllyyy
1 reeeaaaallllyyy	1 reeeeaaaaallyy	1 reeeeeeaaaaallly
1 reeeaaaallllyyyy	1 reeeeaaly	1 reeeeeeaaallyy
1 reeeaaaallly	1 reeeeaallly	1 reeeeeeaaallyy
1 reeeaaaalllyyy	1 reeeeaalllyy	1 reeeeeeaaally
1 reeeaaaaly	1 reeeeaaly	1 reeeeeeaaalllly
1 reeeaaaalllllyyyy	1 reeeeaalllly	1 reeeeeeaaallyy

1 reeeeeeeeeeeaaally
1 reeeeeeeeeeeaaally
1 reeeeeeeeeeeaaally
1 reeeeeeeeeeeeeaaaally
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 rellyzy
1 rieeely
1 rlllllly
1 rllllllyy
1 rlllly
1 rlllyrllly
1 rllly
1 rllyyy
1 rlyrlyrly

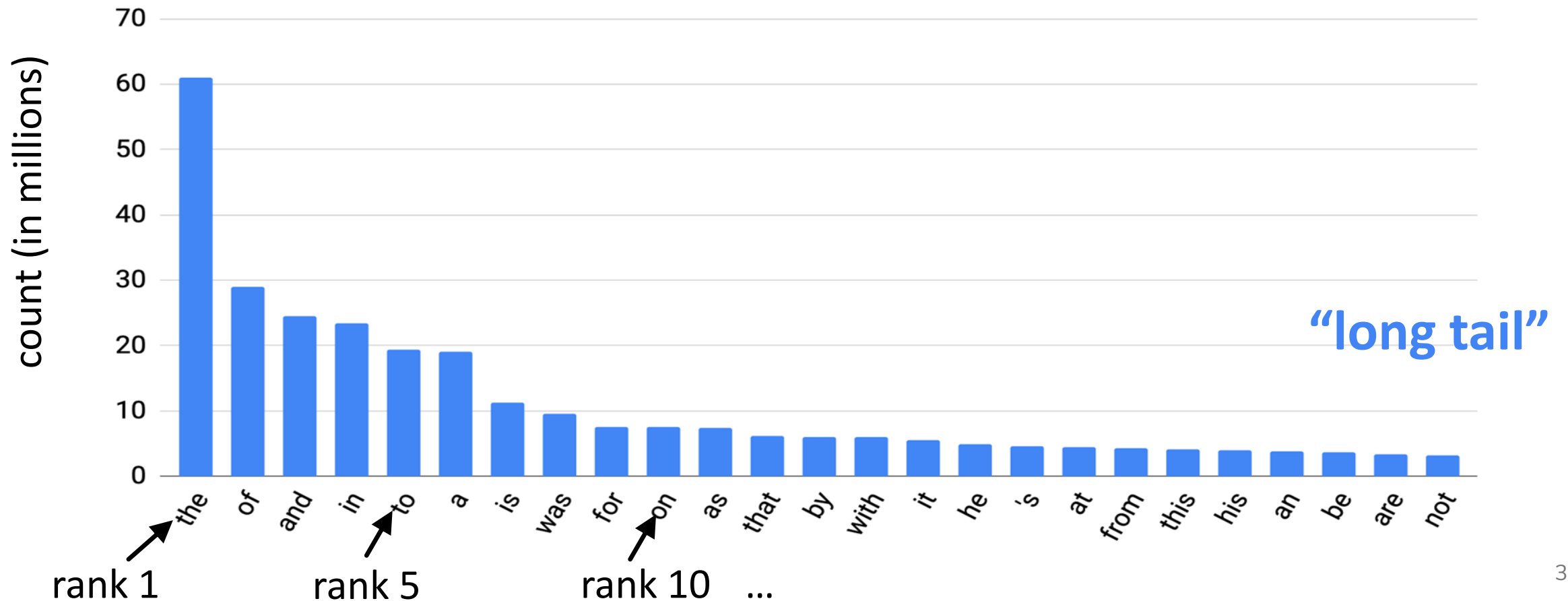
1 rlyy
1 rraarreellyy
1 rreaalllyyy
1 rreaally
1 rreeaaalllllyyyy
1 rreeaalllllyy
1 rreeaalllyyy
1 rreealy
1 rreeeaaaaalllllyyyy
1 rreeeaaalllllyyy
1 rreeeeeeaaaaallllllllyyyyyy
1 rreeeeeeallly
1 rreeeeeeely
1 rrrealllyyy
1 rrreeeaaalllyyy
1 rrreeealllyyy
1 rrreeeeaaalllllyy
1 rrreeeeallly
1 rrrlyyy
1 rrrreeeally
1 rrrreeeeeeaaaaallllllllyyyyyy
1 rrrrreeeeaaalllly
1 rrrrreeeeaaalllyy
1 rrrrrreally
1 rrrrrrealy
1 rrrrrrrreeeeeeaaaaallllllyyyyyyy
1 rrrrrrrrrreally
1 rrrrrrrrrrrrrrrrrreeeeeeeeeeeeeeaaaaaallllllllyyyyyy

How are words distributed?

224571	really	1	rreeeeeeaaaaaalllllllllyyyyyy
1189	rly	1	rreeeeeeallly
1119	realy	1	rreeeeeeely
731	rlly	1	rrreallyyy
590	reaallly	1	rrreeeaaalllyyy
234	reaallly	1	rrreeealllyyy
216	reallyy	1	rrreeeeaaalllllyy
156	relly	1	rrreeeeallly
146	reaalllly	1	rrrllyyy
132	rily	1	rrrreeeally
104	reallyyy	1	rrrreeeeeeaaaaaalllllllllyyyyyy
89	reaallllly	1	rrrrreeeeaaallly
89	reeeally	1	rrrrreeeeaaalllyy
84	reaaally	1	rrrrrreally
82	reaally	1	rrrrrrealy
72	reeeeally	1	rrrrrrreeeeeeaaaaalllllyyyyyyy
65	reaaaally	1	rrrrrrrrreally
...		1	rrrrrrrrrrrrrrrrrrrrreeeeeeeeeeeeeeaaaaaalllllllllyyyyyy

Zipf's Law

Frequency of a word is (roughly) inversely proportional to its rank in the word frequency list.

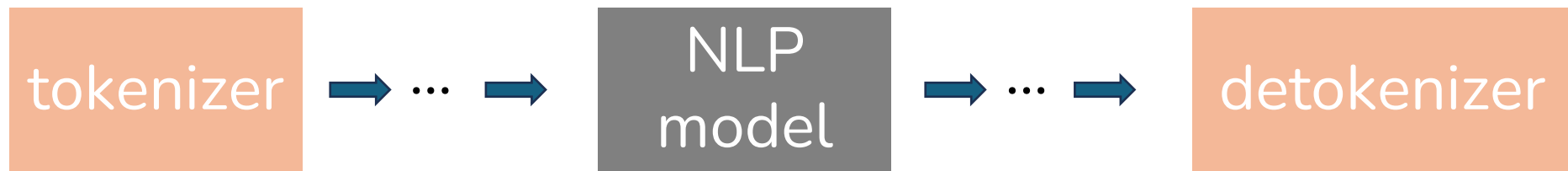


Tokenization in Modern NLP Systems

There are so many word types, but the words have shared internal structures and meanings (recall what we've talked about in morphology).

Modern NLP systems always convert tokens into numerical indices for further processing. Can we do better than assigning each word a unique index?

Raw text input



Raw text output

Data-Driven Subword-Based Tokenizers

Data-driven tokenizers offer an option that **learns** the tokenization rules from data, tokenizing texts into **subword units** (a.k.a, wordpieces) using statistics of character sequences in the dataset.

Two most popular methods:

- Byte Pair Encoding (BPE): Gage (1994), Sennrich et al. (2016).
- SentencePiece: Kudo (2018).

[Gage, P. (1994). A new algorithm for data compression. *The C Users Journal*, 12(2), 23-38.]

[Sennrich, R., Haddow, B., & Birch, A. (2016). Neural machine translation of rare words with subword units. In *Proceedings of ACL* (pp. 1715-1725).]

[Kudo, T. (2018). Subword Regularization: Improving Neural Network Translation Models with Multiple Subword Candidates. In *Proceedings of ACL* (pp. 66-75).]

Byte Pair Encoding

Originally introduced by Gage (1994) for data compression, and later adapted (and revived) by Sennrich et al. (2016) for NLP.

Key idea: **merge** symbols with a greedy algorithm.

Initialize the vocabulary with the set of characters, and iteratively merge the most frequent pair of symbols to extend the vocabulary.

Training Corpus

c a t

c a t s

c o n c a t e n a t i o n

c a t e g o r i z a t i o n

Byte Pair Encoding

Training Corpus

c a t
c a t s
c o n c a t e n a t i o n
c a t e g o r i z a t i o n

Initial Vocabulary

a c e g i n o r s t z

Count Symbol-Pair Frequencies

<a t>: 6, <c a>: 4, <o n>: 3, ...

Update Vocabulary

a c e g i n o r s t z a t

Byte Pair Encoding

Training Corpus

c a t
c a t s
c o n c a t e n a t i o n
c a t e g o r i z a t i o n

Current Vocabulary

a c e g i n o r s t z a t

Count Symbol-Pair Frequencies

<c at>: 4, <o n>: 3, ...

Update Vocabulary

a c e g i n o r s t z a t cat

Repeat this process until the vocabulary reaches the desired size.

Byte Pair Encoding

The BPE proposal is not optimal in terms of compression rate under the same vocabulary size.

Training Corpus

c a t
c a t s
c o n c a t e n a t i o n
c a t e g o r i z a t i o n

Terminate Vocabulary

a c e g i n o r s t z at cat

A Better Vocabulary

a c e g i n o r s t z cat on

Apply Trained BPE to New Corpus

In addition to having the tokens, we will also need to know the “merge rules.” Starting from individual characters, and merge following the rules.

Terminate Vocabulary

a c e g i n o r s t z at cat

Merge Rules

a t → at

c at → cat

Word to be tokenized

c a t e g o r y

Result

cat e g o r y

If there is unknown character, add a new term (unlikely in real practice).

SentencePiece Tokenization

Kudo (2018): find the vocabulary for a unigram language model that maximizes the likelihood of the training corpus.

$$P(\langle x_1, x_2, \dots, x_n \rangle) = \prod_{i=1}^n P(x_i)$$

$$P(x_i) := \text{optional-smoothing} \left(\frac{\text{count}(x_i)}{\sum_{x \in V} \text{count}(x)} \right)$$

The Google command-line toolkit that implements this algorithm and some others (including BPE): <https://github.com/google/sentencepiece>

Byte-Level BPE

A practical question: How do large language models tokenize texts from different languages, with a unified tokenizer and fixed vocabulary size?

That's great 👍

54 68 61 **74** 2019 73 20 67 72 65 61 **74** 20 **1F44D**

All in hexadecimal!

Prepend zeros to fix the length of tokens (to ensure the unique decoding), and do BPE on the bit/multi-bit level.

Next

Word Embeddings