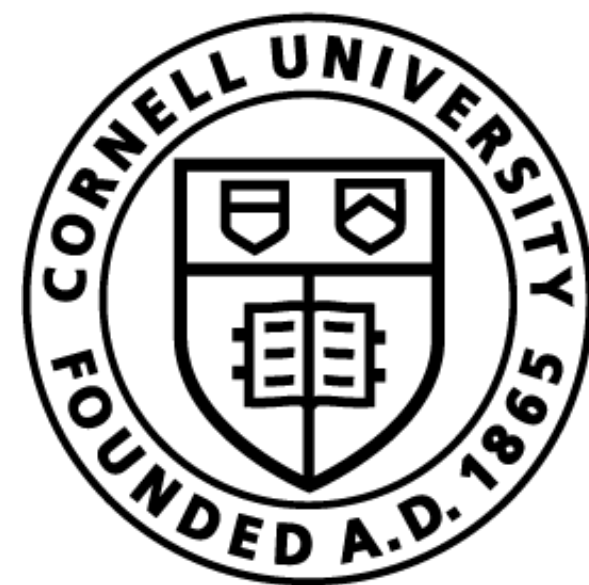


# Lecture 12: Tokenizers, Retrieval Augmented LLMs



Cornell Bowers CIS  
**Computer Science**

Claire Cardie, Tanya Goyal

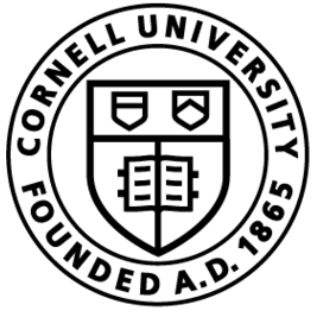
CS 4740 (and crosslists): Introduction to Natural Language Processing

# Announcements

- HW4 due on May 6, 11.59 pm.
  - No penalty submission allowed up to May 9, 11.59 p.m.
- **Review session next class.**
  - Topics-wise or questions?

# Today

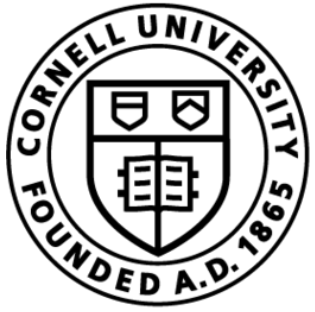
- Tokenization
- RAG
- MiSC topics



# Subword Tokenization

---

- ▶ BERT vocabulary = 30K tokens
  - ▶ More words exist in the dictionary!
- ▶ What will happen if we create a 100K+ word vocabulary?
  - ▶ The final matrix multiply (projecting from last layer to vocabulary) will dominate computation.
  - ▶ Adding a ton of rare words that are rarely used. Will these embeddings get learnt?
  - ▶ Morphology?
- ▶ Character-level models?
  - ▶ Don't work great in practice
  - ▶ Sequence lengths can be very long.



# Subword Tokenization

---

- ▶ Goal: Instead of words, create a vocabulary of “word pieces” that are between characters and full words.

positional embeddings  $\rightarrow$  \_positional\_embeddings

Represents start of a word

- ▶ A number of ways the subword vocabulary can be constructed:
  - ▶ BPE: Byte Pair Encoding (Sennrich et al, 2016; Gage, 1994)
  - ▶ WordPiece (Schuster and Nakajima, 2012)
  - ▶ Unigram LM tokenization (Kudo, 2018)



# BPE Encoding

---

**Algorithm 1** Byte-pair encoding (Sennrich et al., 2016; Gage, 1994)

---

- 1: Input: set of strings  $D$ , target vocab size  $k$
  - 2: **procedure** BPE( $D, k$ )
  - 3:      $V \leftarrow$  all unique characters in  $D$
  - 4:     (about 4,000 in English Wikipedia)
  - 5:     **while**  $|V| < k$  **do**      $\triangleright$  Merge tokens
  - 6:          $t_L, t_R \leftarrow$  Most frequent bigram in  $D$
  - 7:          $t_{\text{NEW}} \leftarrow t_L + t_R$       $\triangleright$  Make new token
  - 8:          $V \leftarrow V + [t_{\text{NEW}}]$
  - 9:         Replace each occurrence of  $t_L, t_R$  in
  - 10:          $D$  with  $t_{\text{NEW}}$
  - 11:     **end while**
  - 12:     **return**  $V$
  - 13: **end procedure**
- 

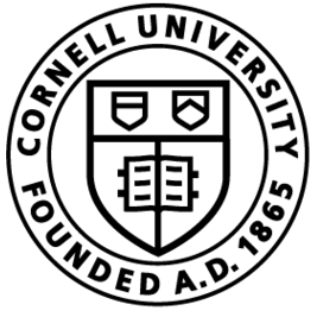
$D = \{$ 'l o w </w>': 9,  
'l o w e r </w>': 2,  
'n e w e s t </w>': 6,  
'w i d e s t </w>': 3}

$V = \{$ l, o, w, e, r, .... </w>

Most frequent bigram= 'lo': 11

$V' = \{$ l, o, lo, w, e, r, .... </w>

$D' = \{$ 'lo w </w>': 9,  
'lo w e r </w>': 2,  
'n e w e s t </w>': 6,  
'w i d e s t </w>': 3}



# Other tokenizers?

## ▶ WordPiece

- ▶ Google never open-sourced their implementation, but based on literature:
- ▶ Merge criterion:

$$\begin{array}{l} t_L, t_R \leftarrow \text{Most frequent bigram in } D \\ t_{\text{NEW}} \leftarrow t_L + t_R \quad \triangleright \text{Make new token} \end{array} \rightarrow \text{score} = \frac{\text{freq}(t_L t_R)}{[\text{freq}(t_L) \times \text{freq}(t_R)]}$$

- ▶ Prioritizes merging where individual chunks are less likely in vocab.

## ▶ Unigram LM tokenizer

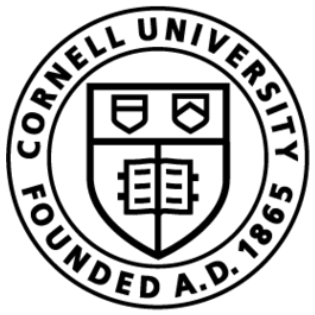
- ▶ Instead of greedy bottom-up approach, top-down.
- ▶ Finds chunks which makes the sequence more likely under the unigram LM



# Comparisons

<b>Original:</b>	furiously	<b>Original:</b>	tricycles	<b>Original:</b>	nanotechnology
<b>BPE:</b>	_fur iously	<b>BPE:</b>	_t ric y cles	<b>BPE:</b>	_n an ote chn ology
<b>Uni. LM:</b>	_fur ious ly	<b>Uni. LM:</b>	_tri cycle s	<b>Uni. LM:</b>	_nano technology
<b>Original:</b>	Completely preposterous suggestions				
<b>BPE:</b>	_Comple t ely _prep ost erous _suggest ions				
<b>Unigram LM:</b>	_Complete ly _pre post er ous _suggestion s				
<b>Original:</b>	corrupted	<b>Original:</b>	1848 and 1852,		
<b>BPE:</b>	_cor rupted	<b>BPE:</b>	_184 8 _and _185 2,		
<b>Unigram LM:</b>	_corrupt ed	<b>Unigram LM:</b>	_1848 _and _1852 ,		

- ▶ BPE tokens may not be semantically meaningful.
- ▶ Unigram LM leads to slightly better BERT (Bostrom and Durrett 2020)



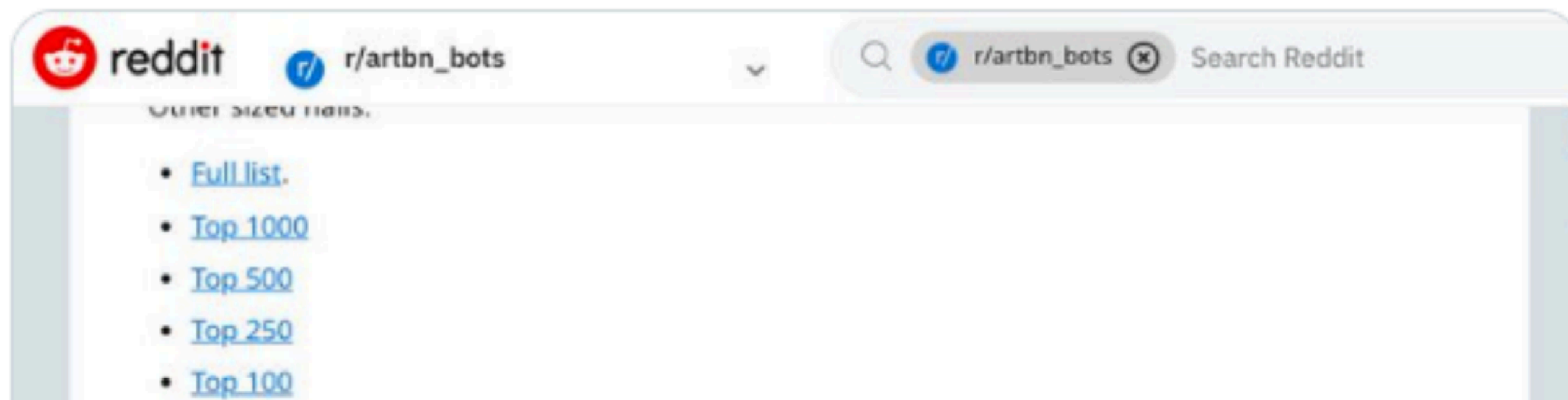
# GPT-3 vocab?

- ▶ BPE tokenizer
- ▶ GPT-3: 50K vocab size

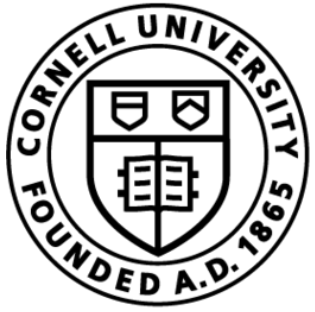


**Matthew Watkins**  
@SoC\_trilogy

I've just found out that several of the anomalous GPT tokens ("TheNitromeFan", " SolidGoldMagikarp", " davidjl", " Smartstocks", " RandomRedditorWithNo", ) are handles of people who are (competitively? collaboratively?) counting to infinity on a Reddit forum. I kid you not.



Rank	User	Counts
1	<a href="#">/u/davidjl123</a>	163477
2	<a href="#">/u/Smartstocks</a>	113829
3	<a href="#">/u/atomicimploder</a>	103178
4	<a href="#">/u/TheNitromeFan</a>	84581
5	<a href="#">/u/SolidGoldMagikarp</a>	65753
6	<a href="#">/u/RandomRedditorWithNo</a>	63434
7	<a href="#">/u/rideride</a>	59024
8	<a href="#">/u/Mooraeil</a>	57785
9	<a href="#">/u/Removedpixel</a>	55080
10	<a href="#">/u/Adinida</a>	48415
11	<a href="#">/u/rschaosid</a>	47132

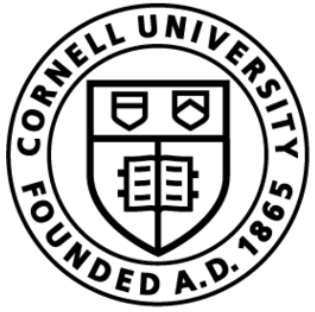


# Are all languages treated equally?

- ▶ Will a sentence in English and Spanish having same number of words have similar number of tokens?  
Think of how BPE vocabulary is constructed...
- ▶ Language premium for A relative to B =  $|t(s_A)| / |t(s_B)|$
- ▶ Flores-200 — contains translations of 842 distinct web articles, totaling 3001 sentences.
- ▶ Premiums wrt to English:
  - ▶ Chinese costs 91% more tokens.
  - ▶ Spanish costs 55% more.

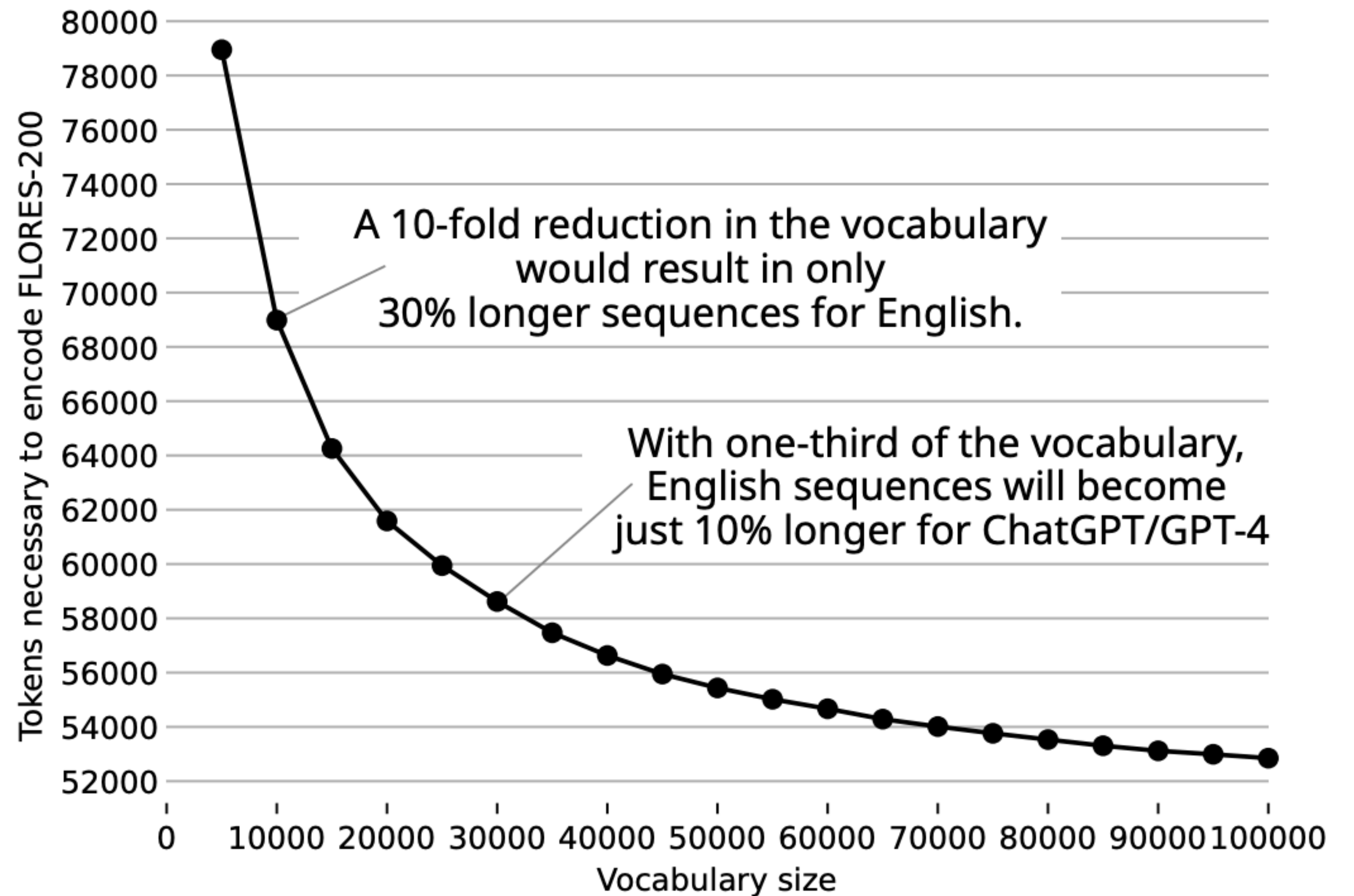
## Premiums with respect to English

	GPT-2 RoBERTa	ChatGPT GPT-4	FlanT5
Bulgarian	5.51	2.64	—
Burmese	16.89	11.70	—
Chinese (Simplified)	3.21	1.91	—
Dzongkha	16.36	12.33	—
English	1.00	1.00	1.00
French	2.00	1.60	1.60
German	2.14	1.58	1.37
Italian	2.01	1.64	2.18
Japanese	3.00	2.30	—
Jingpho	2.65	2.35	3.41
Maori	2.45	2.35	3.28
Norwegian Bokmål	1.86	1.56	2.24
Odia	13.38	12.48	—
Pangasinan	1.66	1.57	2.18
Portuguese	1.94	1.48	2.21
Romanian	2.48	1.88	1.50
Santali	12.86	12.80	—
Shan	18.76	15.05	—
Spanish	1.99	1.55	2.23
Standard Arabic	4.40	3.04	—
Tumbuka	2.78	2.57	3.29
Vietnamese	4.54	2.45	—



# Are all languages treated equally?

- ▶ Fairness implications?
  - ▶ Cost
  - ▶ Latency
  - ▶ Long context processing
- ▶ How much longer will English language tokenization be if we dedicate a fraction of English tokens to other languages?



cl100k\_base tokenizer used by ChatGPT, GPT-4