

# LLM LE3 VT2025

## Data Processing

**Fredrik Heintz**

Dept. of Computer Science

Linköping University

[fredrik.heintz@liu.se](mailto:fredrik.heintz@liu.se)

@FredrikHeintz

### Outline:

- Tokenization
- Data Processing Pipeline

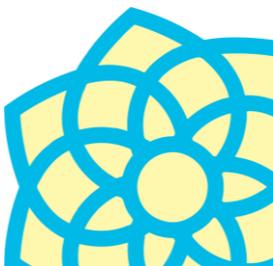
# Language modelling

- **Language modelling** is the task of predicting which word comes next in a sequence of words.
- More formally, given a sequence of words  $w_1, \dots, w_t$  we want to know the probability of the next word,  $w_{t+1}$ :

$$P(w_{t+1} | w_1, \dots, w_t)$$

- We are assuming that  $w_{t+1}$  comes from a finite vocabulary  $V$ .

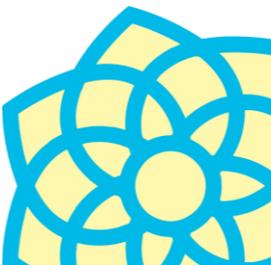
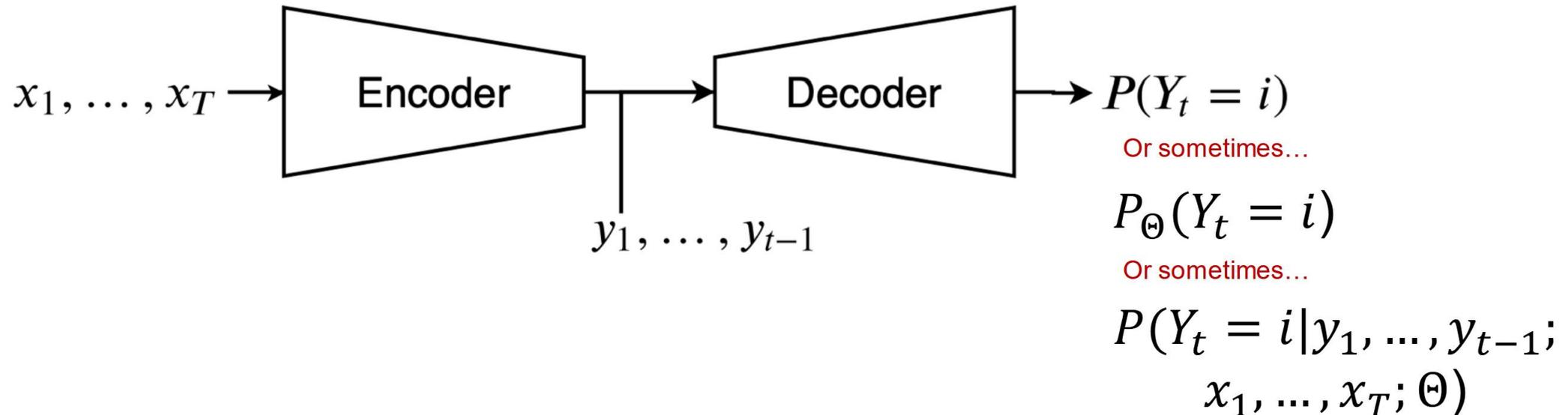
language models = classifiers



# Neurla Language Models

Input sequence:  $x_1, \dots, x_T$

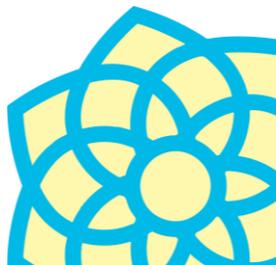
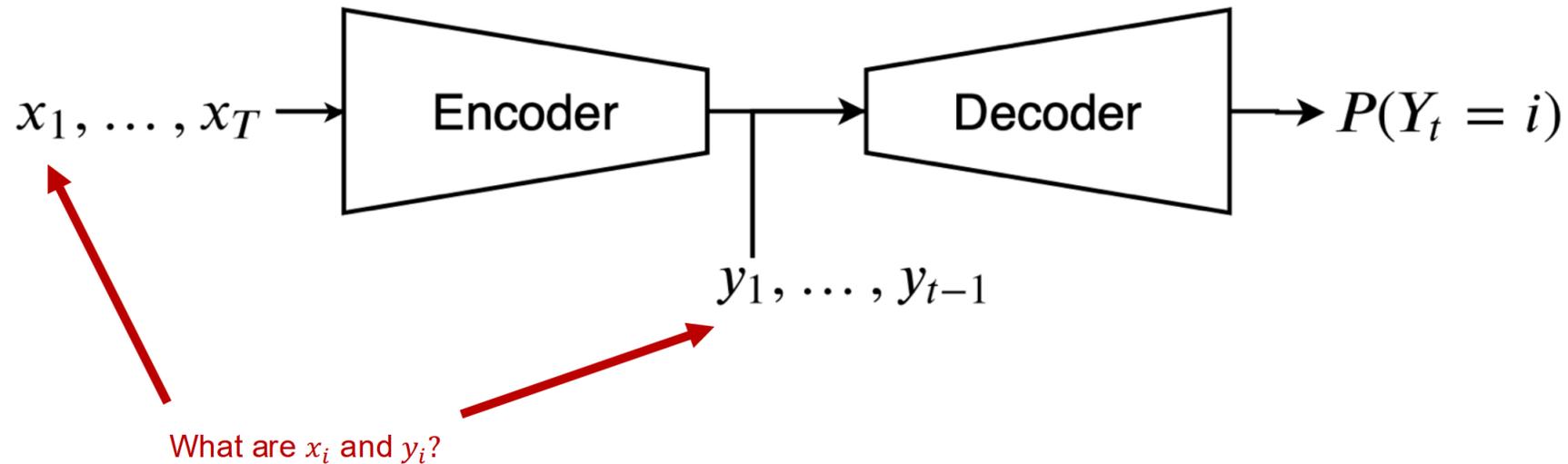
Target sequence:  $y_1, \dots, y_T$



# Neural Language Models

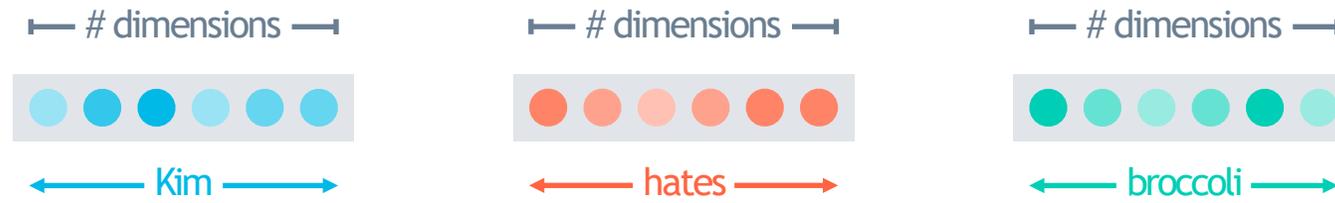
Input sequence:  $x_1, \dots, x_T$

Target sequence:  $y_1, \dots, y_T$



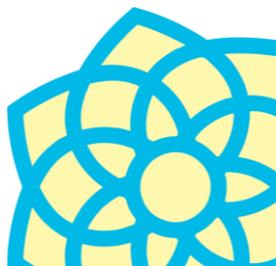
# Word embeddings

To process words using neural networks, we need to represent them as vectors of numerical values.



Compared to one-hot vectors, **word embeddings**

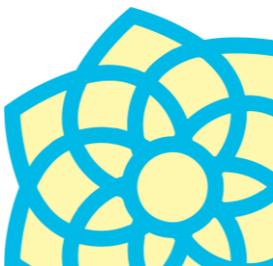
- are shorter but dense
- support a useful notion of similarity
- can be learned from data



# Tokenization

# What is tokenisation?

- **Tokenisation** is the task of taking text (or code or music) and turning it into a sequence of discrete items, such as words or characters, called tokens.
- Tokenisation simplifies natural language processing by reducing unstructured text to more useful units.
- Tokenisation is the first step in mapping text to a numerical representation that computers can process.



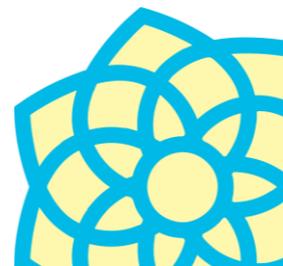
# Words provide important signals

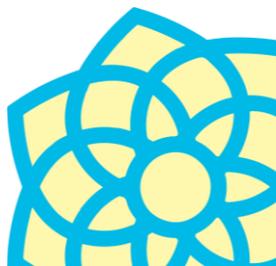
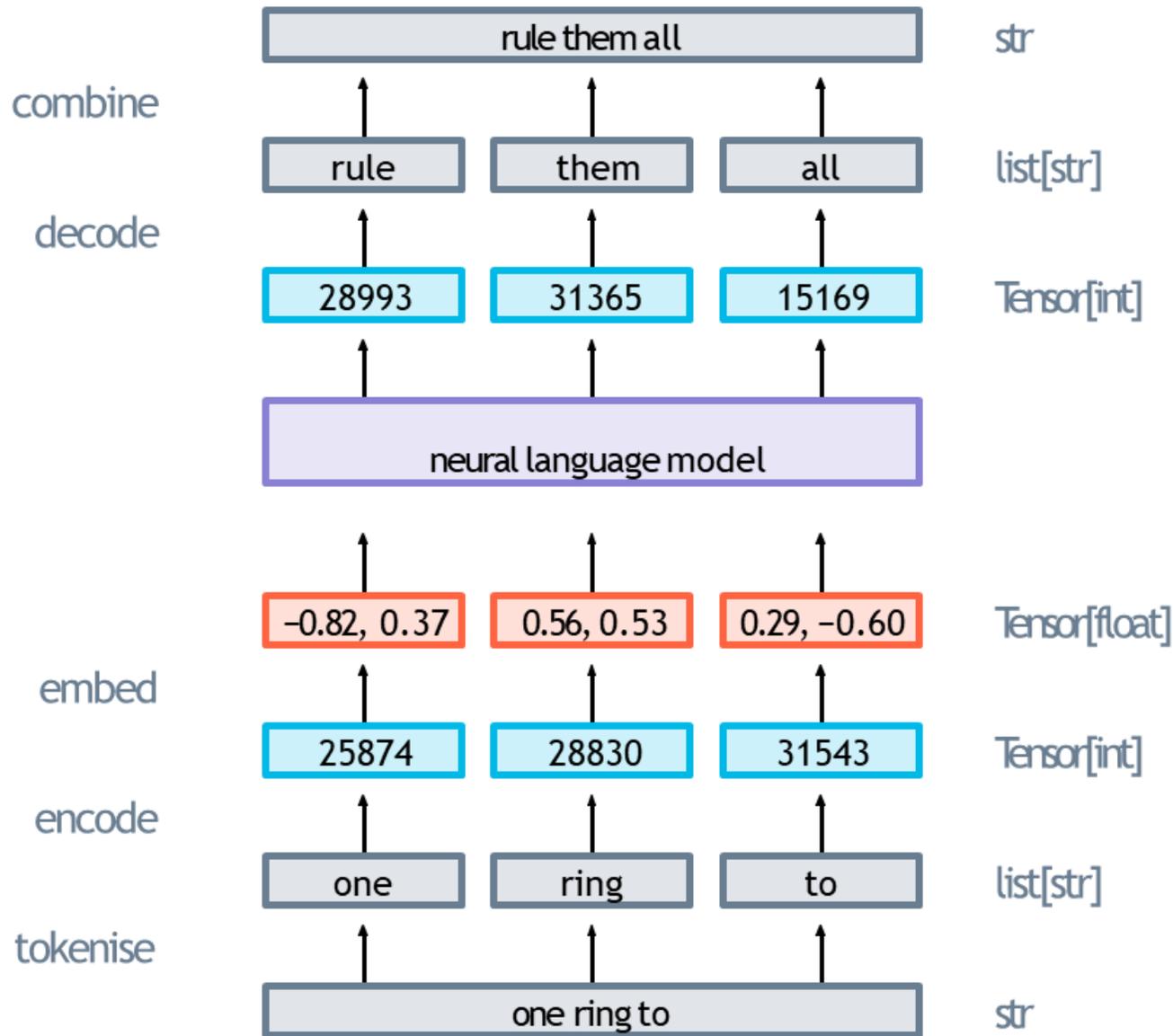
The gorgeously elaborate continuation of “The Lord of the Rings” trilogy is so huge that a column of words cannot adequately describe co-writer/director Peter Jackson’s expanded vision of J.R.R. Tolkien’s Middle-earth.

positive

... is a sour little movie at its core; an exploration of the emptiness that underlay the relentless gaiety of the 1920’s, as if to stop would hasten the economic and global political turmoil that was to come.

negative





# Whitespace tokenisation

```
# Tokenise text by splitting at whitespace def
```

```
tokenize(text: str) -> list[str]:
```

```
    return text.split()
```

```
# Create a vocabulary
```

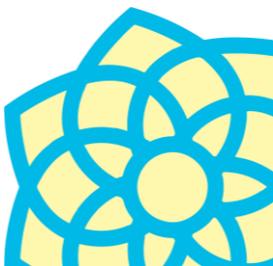
```
vocab: set[str] = set(tokenize(text))
```

```
# {'cannot', 'huge', 'column', 'that', 'is', ...}
```

```
# Create a string-to-ID mapping
```

```
stoid: dict[str, int] = {s: i for i, s in enumerate(vocab)}
```

```
# {'cannot': 0, 'huge': 1, 'column': 2, 'that': 3, 'is': 4, ...}
```



### Whitespace tokenisation

The gorgeously elaborate continuation of “The Lord of the Rings” trilogy is so huge that a column of words cannot adequately describe co-writer/director Peter Jackson’s expanded vision of J.R.R. Tolkien’s Middle-earth.

### Regex-based tokenisation

The gorgeously elaborate continuation of “The Lord of the Rings” trilogy is so huge that a column of words cannot adequately describe co-writer / director Peter Jackson ’s expanded vision of J. R. R. Tolkien ’s Middle-earth .

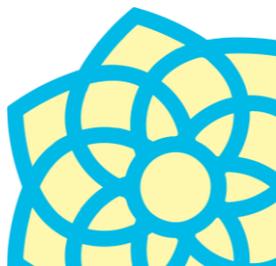
```
re.findall(r"[A-Za-z]\.| \w+(?:-\w+)*| '\w+| [^\w\s]+", text)
```

single letters  
followed by a period

whole words, incl.  
hyphenated words

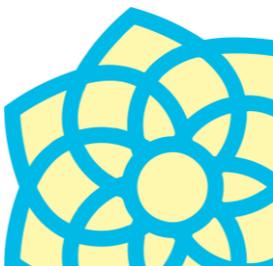
genitives (’s) and  
contractions (’ve)

punctuation, other  
non-word characters

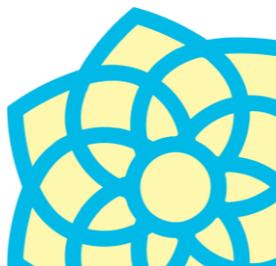
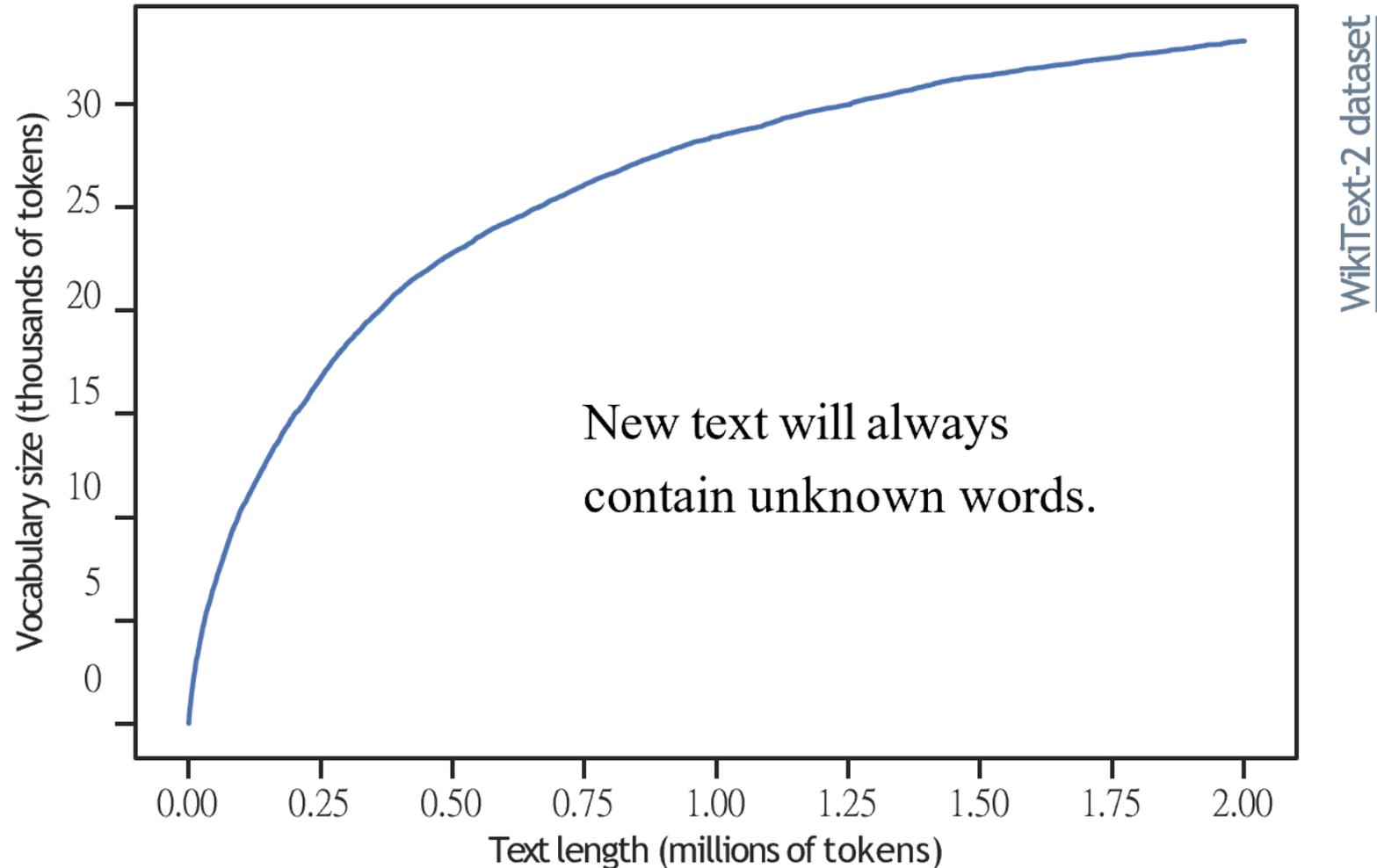


# Text normalisation

- **Text normalisation** refers to the process of converting text into a more useful, standard form.
- Standard techniques include case normalisation, harmonisation of spelling variants, lemmatisation, and removing punctuation.  
Harmonisation: *color* → *colour*. Lemmatization: *runs, ran, running* → *run*
- Text normalisation was once a critical step in NLP tasks but is no longer as widely used today.



# The challenge of unknown words – Heaps' law



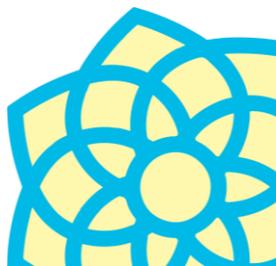
# Dealing with unknown words

- **Step 1:** Build the vocabulary as usual.  
often combined with a frequency threshold
- **Step 2:** Augment the vocabulary with a special token, such as [UNK] to represent unknown words.
- **Step 3:** When processing new text, replace any out-of-vocabulary (oov) word with the special [UNK] token.

*The quokka is adorable. →The [UNK] is adorable. (Assuming quokka is oov.)*



By Ena Music - Own work, CC BY-SA 4.0, [Link](#)



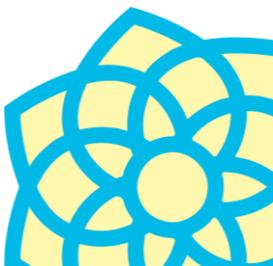
# But what is a word, anyway?

- There are many languages that do not adhere to the same concept of a “word” as English and Swedish.
- **Chinese** is written without spaces between characters. Identifying word boundaries is challenging.

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

- **Inuktitut** allows entire sentences to be expressed as single words by combining multiple morphemes.

*tusaatsiarunnanngittualuujunga* — “I cannot hear very well.”



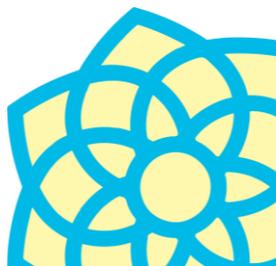
# Tokenization is language dependent

Table 1: Tokenizer comparisons between original LLaMA and Chinese LLaMA.

	Length	Content
<b>Original Sentence</b>	28	人工智能是计算机科学、心理学、哲学等学科融合的交叉学科。
<b>Original Tokenizer</b>	35	‘_’, ‘人’, ‘工’, ‘智’, ‘能’, ‘是’, ‘计’, ‘算’, ‘机’, ‘科’, ‘学’, ‘、’, ‘心’, ‘理’, ‘学’, ‘、’, ‘0xE5’, ‘0x93’, ‘0xB2’, ‘学’, ‘等’, ‘学’, ‘科’, ‘0xE8’, ‘0x9E’, ‘0x8D’, ‘合’, ‘的’, ‘交’, ‘0xE5’, ‘0x8F’, ‘0x89’, ‘学’, ‘科’, ‘。’
<b>Chinese Tokenizer</b>	16	‘_’, ‘人工智能’, ‘是’, ‘计算机’, ‘科学’, ‘、’, ‘心理学’, ‘、’, ‘哲学’, ‘等’, ‘学科’, ‘融合’, ‘的’, ‘交叉’, ‘学科’, ‘。’

LLaMA tokenizer is **unfriendly** to Chinese

Yiming Cui. et.al. EFFICIENT AND EFFECTIVE TEXT ENCODING FOR CHINESE LLAMA AND ALPACA. <https://arxiv.org/pdf/2304.08177.pdf>

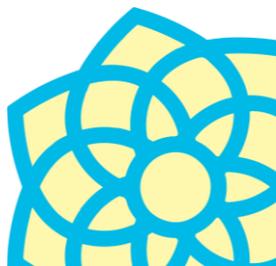


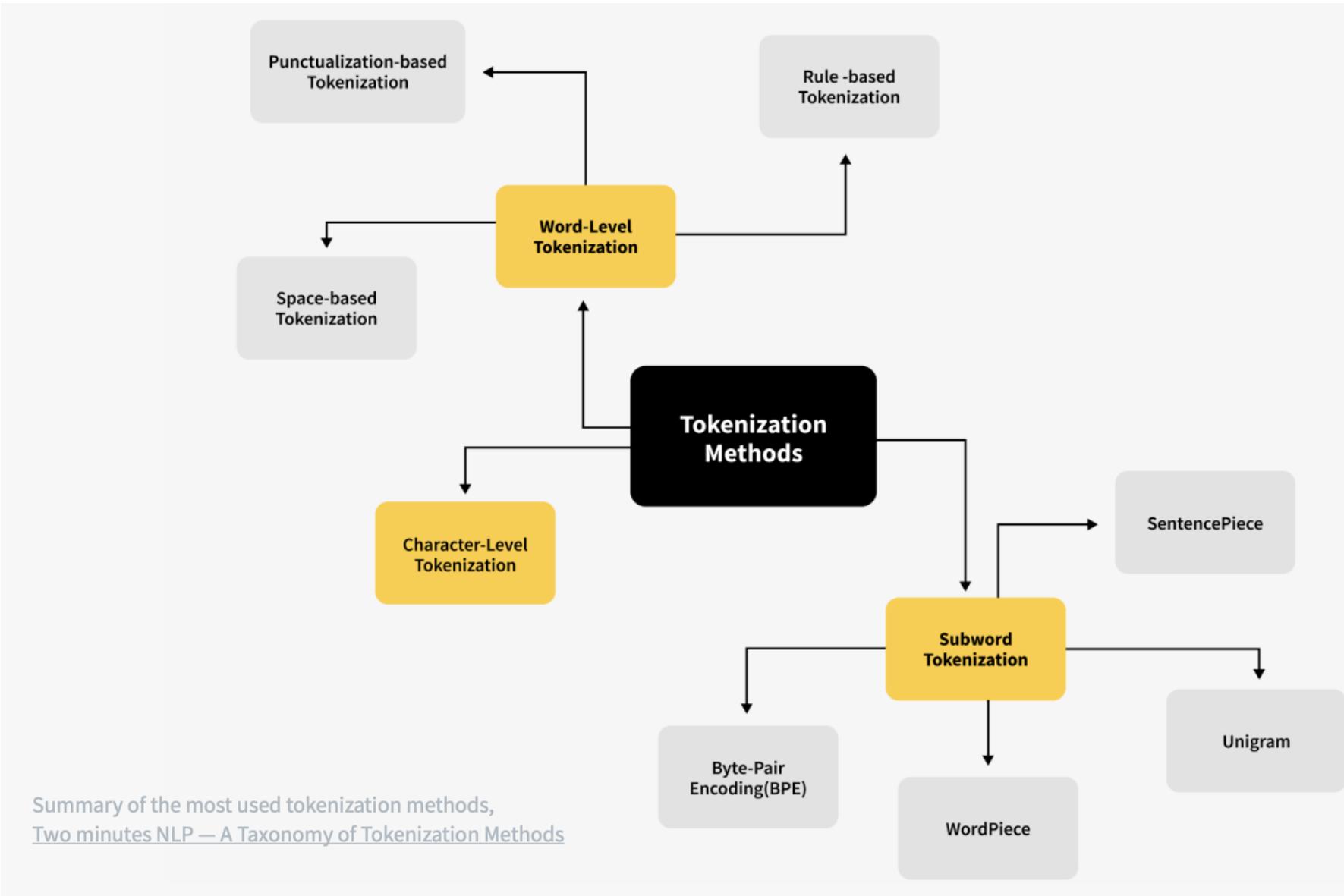
# Target representations for tokenisation

- **Option 1:** Tokenise into words  
But: concept of “word” not universal; unknown words
- **Option 2:** Tokenise into individual characters  
But: may be too small a unit for learning
- **Option 3:** Tokenise into subwords  
Intuition: words are composed of morphemes

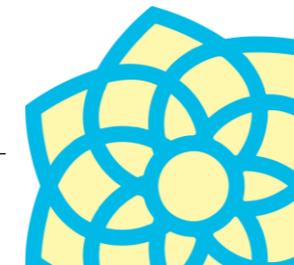
Let’s tokenize: “A hippopotamus ate my homework.”

Vocab Type	Example	Ex. length
character-level	['A', ' ', 'h', 'i', 'p', 'p', 'o', 'p', 'o', 't', 'a', 'm', 'u', 's', ' ', ' ', 'a', 't', 'e', ' ', ' ', 'm', 'y', ' ', ' ', 'h', 'o', 'm', 'e', 'w', 'o', 'r', 'k', '.']	31
subword-level	['A', 'hip', '##pop', '##ota', '##mus', 'ate', 'my', 'homework', '.']	9
word-level	['A', 'hippopotamus', 'ate', 'my', 'homework']	5



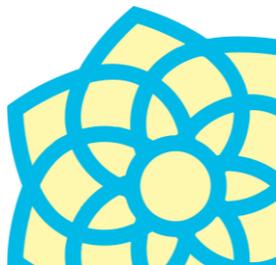


Summary of the most used tokenization methods,  
[Two minutes NLP – A Taxonomy of Tokenization Methods](#)



Tokenization Methods	Word-based tokenization	Character-based tokenization	Subword-based tokenization
<b>Example Tokenizers</b>	Space tokenization (split sentences by space); rule-based tokenization (e.g. Moses, spaCy)	Character tokenization (simply tokenize on every character)	Byte-Pair Encoding (BPE); WordPiece; SentencePiece; Unigram (tokenizing by parts of a word vs. the entirety of a word; see table above)
<b>Considerations</b>	<ul style="list-style-type: none"> <li>• <b>Downside:</b> Generates a very large vocabulary leading to a huge embedding matrix as the input and output layer; large number of out-of-vocabulary (OOV) tokens; and different meanings of very similar words</li> <li>• Transformer models normally have a vocabulary of less than 50,000 words, especially if they are trained only on a single language</li> </ul>	<ul style="list-style-type: none"> <li>• Lead to much smaller vocabulary; no OOV (out of vocabulary) tokens since every word can be assembled from individual characters</li> <li>• <b>Downside:</b> Generates very long sequences and less meaningful individual tokens, making it harder for the model to learn meaningful input representations. However, if character-based tokenization is used on non-English language, a single character could be quite information rich (like “mountain” in Mandarin).</li> </ul>	<ul style="list-style-type: none"> <li>• Subword-based tokenization methods follow the principle that frequently used words should not be split into smaller subwords, but rare words should be decomposed into meaningful subwords</li> <li>• <b>Benefit:</b> Solves the downsides faced by word-based tokenization and character-based tokenization and achieves both reasonable vocabulary size with meaningful learned context-independent representations.</li> </ul>

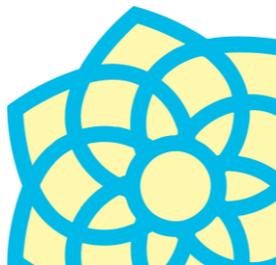
LooongLLaVA  
LGUer



# Subword-based Tokenization Methods

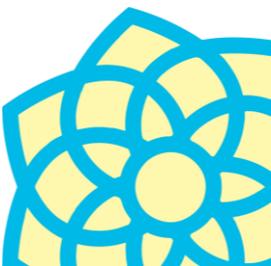
- **Byte-Pair Encoding** [[Gage 1994](#)]
  - Originally used in machine translation
- **WordPiece**
- **Unigram**
- **SentencePiece**

Subword-based Tokenization Methods	Byte-Pair Encoding (BPE)	WordPiece	Unigram	SentencePiece
<b>Description</b>	<p>One of the most popular subword tokenization algorithms. The Byte-Pair-Encoding works by starting with characters, while merging those that are the most frequently seen together, thus creating new tokens. It then works iteratively to build new tokens out of the most frequent pairs it sees in a corpus.</p> <p>BPE is able to build words it has never seen by using multiple subword tokens, and thus requires smaller vocabularies, with less chances of having “unk” (unknown) tokens.</p>	<p>Very similar to BPE. The difference is that WordPiece does not choose the highest frequency symbol pair, but the one that maximizes the likelihood of the training data once added to the vocabulary (evaluates what it loses by merging two symbols to ensure it's worth it)</p>	<p>In contrast to BPE / WordPiece, Unigram initializes its base vocabulary to a large number of symbols and progressively trims down each symbol to obtain a smaller vocabulary. It is often used together with SentencePiece.</p>	<p>The left 3 tokenizers assume input text uses spaces to separate words, and therefore are not usually applicable to languages that don't use spaces to separate words (e.g. Chinese). SentencePiece treats the input as a raw input stream, thus including the space in the set of characters to use. It then uses the BPE / Unigram algorithm to construct the appropriate vocabulary.</p>
<b>Considerations</b>	<p>BPE is particularly useful for handling rare and out-of-vocabulary words since it can generate subwords for new words based on the most common character sequences.</p> <p>Downside: BPE can result in subwords that do not correspond to linguistically meaningful units.</p>	<p>WordPiece can be particularly useful for languages where the meaning of a word can depend on the context in which it appears.</p>	<p>Unigram tokenization is particularly useful for languages with complex morphology and can generate subwords that correspond to linguistically meaningful units. However, unigram tokenization can struggle with rare and out-of-vocabulary words.</p>	<p>SentencePiece can be particularly useful for languages where the meaning of a word can depend on the context in which it appears.</p>



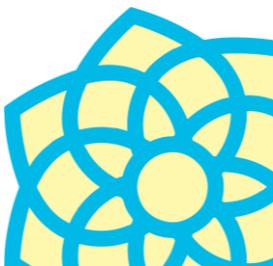
# Byte-Pair Encoding (BPE) [Gage 1994]

- **Byte Pair Encoding (BPE)** is an algorithm for learning subword tokens from text.
- **Step 1:** Encode the text into a sequence of bytes. Initialise the token vocabulary with all single bytes.
- **Step 2:** Create a new token by merging the most frequent pair of consecutive tokens. Add the new token to the vocabulary.
- Repeat the previous step as long as the token vocabulary does not exceed a predefined maximum size.



# The Unicode Standard

- **Unicode** is a text encoding standard designed to support text from all the world's writing systems (that can be digitised).
- Version 16.0 supports 154,998 characters from 168 scripts.
- For backwards compatibility, the first 128 codepoints of Unicode are the same as ASCII.



	000	001	002	003	004	005	006	007
0	NUL 0000	DLE 0010	SP 0020	0 0030	@ 0040	P 0050	` 0060	p 0070
1	SOH 0001	DC1 0011	! 0021	1 0031	A 0041	Q 0051	a 0061	q 0071
2	STX 0002	DC2 0012	" 0022	2 0032	B 0042	R 0052	b 0062	r 0072
3	ETX 0003	DC3 0013	# 0023	3 0033	C 0043	S 0053	c 0063	s 0073
4	EOT 0004	DC4 0014	\$ 0024	4 0034	D 0044	T 0054	d 0064	t 0074
5	ENO 0005	NAK 0015	% 0025	5 0035	E 0045	U 0055	e 0065	u 0075
6	ACK 0006	SYN 0016	& 0026	6 0036	F 0046	V 0056	f 0066	v 0076
7	BEL 0007	ETB 0017	' 0027	7 0037	G 0047	W 0057	g 0067	w 0077
8	BS 0008	CAN 0018	( 0028	8 0038	H 0048	X 0058	h 0068	x 0078
9	HT 0009	EM 0019	) 0029	9 0039	I 0049	Y 0059	i 0069	y 0079
A	LF 000A	SUB 001A	* 002A	: 003A	J 004A	Z 005A	j 006A	z 007A
B	VT 000B	ESC 001B	+ 002B	; 003B	K 004B	[ 005B	k 006B	{ 007B
C	FF 000C	FS 001C	, 002C	< 003C	L 004C	\ 005C	l 006C	 007C
D	CR 000D	GS 001D	- 002D	= 003D	M 004D	] 005D	m 006D	} 007D
E	SO 000E	RS 001E	. 002E	> 003E	N 004E	^ 005E	n 006E	~ 007E
F	SI 000F	US 001F	/ 002F	? 003F	O 004F	_ 005F	o 006F	DEL 007F

**Various signs**

- 0900 ॐ DEVANAGARI SIGN INVERTED CANDRABINDU = vaidika adhomukha candrabindu
- 0901 ॐ DEVANAGARI SIGN CANDRABINDU = anunasika  
→ 0310 ॐ combining candrabindu
- 0902 ॐ DEVANAGARI SIGN ANUSVARA = bindu
- 0903 ॐ DEVANAGARI SIGN VISARGA
- Independent vowels**
- 0904 ॐ DEVANAGARI LETTER SHORT A  
• used for short e in Awadhi  
• also used in Devanagari transliterations of some South Indian and Kashmiri languages by a publisher in Lucknow
- 0905 अ DEVANAGARI LETTER A
- 0906 आ DEVANAGARI LETTER AA
- 0907 इ DEVANAGARI LETTER I
- 0908 ई DEVANAGARI LETTER II
- 0909 उ DEVANAGARI LETTER U
- 090A ऊ DEVANAGARI LETTER UU
- 090B ऋ DEVANAGARI LETTER VOCALIC R
- 090C ॠ DEVANAGARI LETTER VOCALIC L
- 090D ऌ DEVANAGARI LETTER CANDRA E
- 090E ॡ DEVANAGARI LETTER SHORT E  
• Kashmiri, Bihari languages  
• also used for transcribing Dravidian short e

- 090F ऋ DEVANAGARI LETTER E
- 0910 ँ DEVANAGARI LETTER AI
- 0911 ओ DEVANAGARI LETTER CANDRA O
- 0912 औ DEVANAGARI LETTER SHORT O  
• Kashmiri, Bihari languages  
• also used for transcribing Dravidian short o
- 0913 ओ DEVANAGARI LETTER O
- 0914 औ DEVANAGARI LETTER AU

**Consonants**

- 0915 क DEVANAGARI LETTER KA
- 0916 ख DEVANAGARI LETTER KHA
- 0917 ग DEVANAGARI LETTER GA
- 0918 घ DEVANAGARI LETTER GHA
- 0919 ङ DEVANAGARI LETTER NGA
- 091A च DEVANAGARI LETTER CA
- 091B छ DEVANAGARI LETTER CHA
- 091C ज DEVANAGARI LETTER JA
- 091D झ DEVANAGARI LETTER JHA
- 091E ञ DEVANAGARI LETTER NYA
- 091F ट DEVANAGARI LETTER TTA
- 0920 ठ DEVANAGARI LETTER TTHA
- 0921 ड DEVANAGARI LETTER DDA
- 0922 ढ DEVANAGARI LETTER DDHA
- 0923 ण DEVANAGARI LETTER NNA
- 0924 त DEVANAGARI LETTER TA
- 0925 थ DEVANAGARI LETTER THA
- 0926 द DEVANAGARI LETTER DA
- 0927 ध DEVANAGARI LETTER DHA
- 0928 न DEVANAGARI LETTER NA
- 0929 ण DEVANAGARI LETTER NNA  
• for transcribing Dravidian alveolar n  
≡ 0928 ण 093C ॠ
- 092A प DEVANAGARI LETTER PA
- 092B फ DEVANAGARI LETTER PHA

- 092C व DEVANAGARI LETTER BA
- 092D ब DEVANAGARI LETTER BHA
- 092E म DEVANAGARI LETTER MA
- 092F य DEVANAGARI LETTER YA
- 0930 र DEVANAGARI LETTER RA
- 0931 ॠ DEVANAGARI LETTER RRA  
• for transcribing Dravidian alveolar r  
• half form is represented as "Eyelash RA"  
≡ 0930 र 093C ॠ
- 0932 ल DEVANAGARI LETTER LA
- 0933 ॡ DEVANAGARI LETTER LLA
- 0934 ॢ DEVANAGARI LETTER LLLA  
• for transcribing Dravidian l  
≡ 0933 ॡ 093C ॠ
- 0935 व DEVANAGARI LETTER VA
- 0936 श DEVANAGARI LETTER SHA
- 0937 ष DEVANAGARI LETTER SSA
- 0938 स DEVANAGARI LETTER SA
- 0939 ह DEVANAGARI LETTER HA

**Dependent vowel signs**

These dependent vowel signs are used in Kashmiri and in the Bihari languages (Bhojpuri, Magadhi, and Maithili).

- 093A ॐ DEVANAGARI VOWEL SIGN OE
- 093B ॐ DEVANAGARI VOWEL SIGN OOE

**Various signs**

- 093C ॐ DEVANAGARI SIGN NUKTA  
• for extending the alphabet to new letters
- 093D ॐ DEVANAGARI SIGN AVAGRAHA

**Dependent vowel signs**

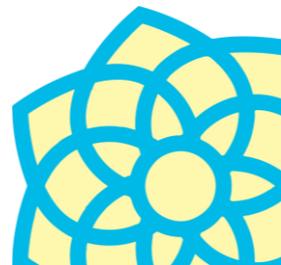
- 093E ॐ DEVANAGARI VOWEL SIGN AA
- 093F ॐ DEVANAGARI VOWEL SIGN I  
• stands to the left of the consonant
- 0940 ॐ DEVANAGARI VOWEL SIGN II
- 0941 ॐ DEVANAGARI VOWEL SIGN U
- 0942 ॐ DEVANAGARI VOWEL SIGN UU
- 0943 ॐ DEVANAGARI VOWEL SIGN VOCALIC R
- 0944 ॐ DEVANAGARI VOWEL SIGN VOCALIC RR
- 0945 ॐ DEVANAGARI VOWEL SIGN CANDRA E = candra
- 0946 ॐ DEVANAGARI VOWEL SIGN SHORT E  
• Kashmiri, Bihari languages  
• also used for transcribing Dravidian short e
- 0947 ॐ DEVANAGARI VOWEL SIGN E
- 0948 ॐ DEVANAGARI VOWEL SIGN AI
- 0949 ॐ DEVANAGARI VOWEL SIGN CANDRA O
- 094A ॐ DEVANAGARI VOWEL SIGN SHORT O  
• Kashmiri, Bihari languages  
• also used for transcribing Dravidian short o
- 094B ॐ DEVANAGARI VOWEL SIGN O
- 094C ॐ DEVANAGARI VOWEL SIGN AU

**Virama**

- 094D ॐ DEVANAGARI SIGN VIRAMA = halant (the preferred Hindi name)  
• suppresses inherent vowel

**Dependent vowel signs**

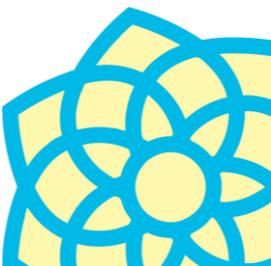
- 094E ॐ DEVANAGARI VOWEL SIGN PRISHTHAMATRA E  
• character has historic use only  
• combines with E to form AI, with AA to form O, and with O to form AU



# Encoding text into bytes

- Encoding all (more than 1 million) Unicode characters into bytes requires more than one byte per character.
- **UTF-8 (8-bit Unicode Transformation Format)** is the most widely used encoding scheme for Unicode.
- It uses a variable-width encoding of 1-4 bytes per character.

The first byte indicates how many additional bytes are part of the character.



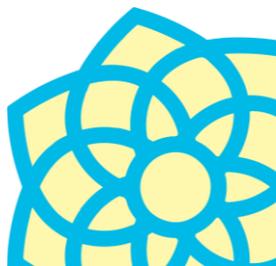
# Encoding text into bytes

Einu sinni deildu norðanvindurinn og sólin um, kvort þeirra væri sterkara.

74 Unicode characters

E	i	n	u		s	i	n	n	i		d	e	i	l	d	u		n	o	r	ð	a	
69	105	110	117	32	115	105	110	110	105	32	100	101	105	108	100	117	32	110	111	114	195	176	97
n	v	i	n	d	u	r	i	n	n		o	g		s	ó	l	i	n	n		u	m	
110	118	105	110	100	117	114	105	110	110	32	111	103	32	115	195	179	108	105	110	110	32	117	109
,		k	v	o	r	t		þ	e	i	r	r	a		v	æ	r	i		s	t		
44	32	107	118	111	114	116	32	195	190	101	105	114	114	97	32	118	195	166	114	105	32	115	116
e	r	k	a	r	a	.																	
101	114	107	97	114	97	46																	

78 bytes in UTF-8



# Byte-Pair Encoding – Example

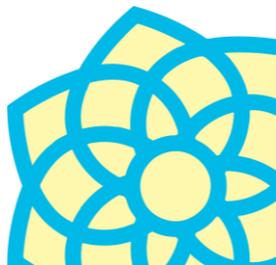
The North Wind and the Sun were disputing which was the stronger, when a traveler came along wrapped in a warm cloak. They agreed that the one who first succeeded in making the traveler take his cloak off should be considered stronger than the other.

## Vocabulary (without single bytes)

Token ID	Token
256	
257	
258	
259	
260	

## Pair counts

Token pair	Count
e + SPACE	11
SPACE + t	10
h + e	9
t + h	9
d + SPACE	7



# Byte-Pair Encoding – Example

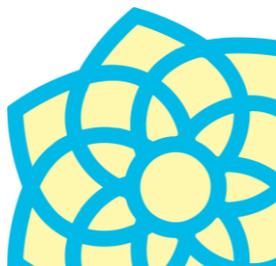
The North Wind and the Sun were disputing which was the stronger, when a traveler came along wrapped in a warm cloak. They agreed that the one who first succeeded in making the traveler take his cloak off should be considered stronger than the other.

## Vocabulary (without single bytes)

Token ID	Token
256	[e ]
257	
258	
259	
260	

## Pair counts

Token pair	Count
e + SPACE	11
SPACE + t	10
h + e	9
t + h	9
d + SPACE	7



# Byte-Pair Encoding – Example

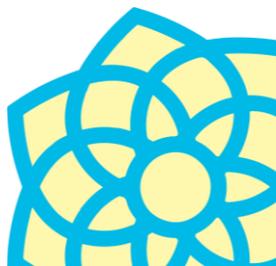
Th[e ]North Wind and th[e ]Sun wer[e ]disputing which was th[e ]stronger,  
when a traveler cam[e ]along wrapped in a warm cloak. They agreed that  
th[e ]on[e ]who first succeeded in making th[e ]traveler tak[e ]his cloak  
off should b[e ]considered stronger than th[e ]other.

## Vocabulary (without single bytes)

Token ID	Token
256	[e]
257	
258	
259	
260	

## Pair counts

Token pair	Count
t + h	9
SPACE + t	9
d + SPACE	7
e + r	7
h + [e]	6



# Byte-Pair Encoding – Example

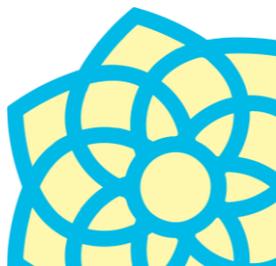
Th[e ]North Wind and th[e ]Sun wer[e ]disputing which was th[e ]stronger, when a traveler cam[e ]along wrapped in a warm cloak. They agreed th[at th[e ]on[e ]who first succeeded in making th[e ]traveler tak[e ]his cloak off should b[e ]considered stronger than th[e ]other.

## Vocabulary (without single bytes)

Token ID	Token
256	[e ]
257	[th]
258	
259	
260	

## Pair counts

Token pair	Count
t + h	9
SPACE + t	9
d + SPACE	7
e + r	7
h + [e ]	6



# Byte-Pair Encoding – Example

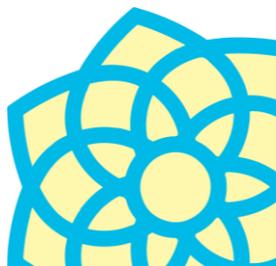
Th[e ]Nor[th] Wind and [th][e ]Sun wer[e ]disputing which was  
 [th][e ]stronger, when a traveler cam[e ]along wrapped in a warm  
 cloak. [Th]ey agreed [th]at [th][e ]on[e ]who first succeeded in making  
 [th][e ]traveler tak[e ]his cloak off should b[e ]considered stronger  
 [th]an [th][e ]other.

## Vocabulary (without single bytes)

Token ID	Token
256	[e ]
257	[th]
258	
259	
260	

## Pair counts

Token pair	Count
d + SPACE	7
SPACE + [th]	7
e + r	7
SPACE + W	6
i + n	5



# Byte-Pair Encoding – Example

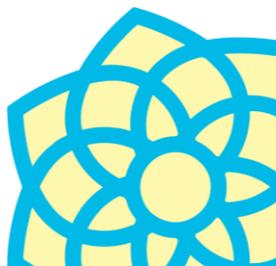
Th[e ]Nor[th] Wind and [th][e ]Sun wer[e ]disputing which was [th][e ]stronger, when a traveler cam[e ]along wrapped in a warm cloak. [Th]ey agreed [th]at [th][e ]on[e ]who first succeeded in making [th][e ]traveler tak[e ]his cloak off should b[e ]considered stronger [th]an [th][e ]other.

## Vocabulary (without single bytes)

Token ID	Token
256	[e ]
257	[th]
258	[d ]
259	
260	

## Pair counts

Token pair	Count
d + SPACE	7
SPACE + [th]	7
e + r	7
SPACE + W	6
i + n	5



# Byte-Pair Encoding – Example

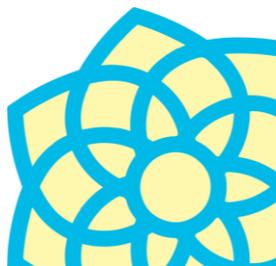
Th[e ]Nor[th] Wind and [th][e ]Sun wer[e ]disputing which was [th][e ]stronger, when a traveler cam[e ]along wrapped in a warm cloak. [Th]ey agreed [th]at [th][e ]on[e ]who first succeeded in making [th][e ]traveler tak[e ]his cloak off should b[e ]considered stronger [th]an [th][e ]other.

## Vocabulary (without single bytes)

Token ID	Token
256	[e ]
257	[th]
258	[d ]
259	
260	

## Pair counts

Token pair	Count
e + r	7
SPACE + w	6
i + n	5
[th] + [e ]	5
n + SPACE	5



# Byte-Pair Encoding – Example

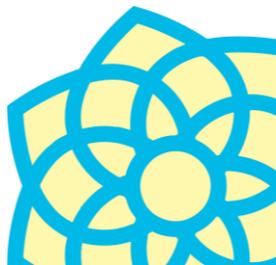
Th[e ]Nor[th] Wind and [th][e ]Sun wer[e ]disputing which was [th][e ]stronger, when a traveler came along wrapped in a warm cloak. [Th]ey agreed [th]at [th][e ]on[e ]who first succeeded in making [th][e ]traveler take his cloak off should be considered stronger than [th][e ]other.

## Vocabulary (without single bytes)

Token ID	Token
256	[e ]
257	[th]
258	[d ]
259	
260	

## Pair counts

Token pair	Count
e + r	7
SPACE + w	6
i + n	5
[th] + [e ]	5
n + SPACE	5



# Byte-Pair Encoding – Example

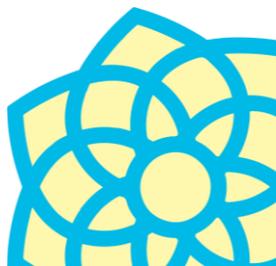
Th[e ]Nor[th] Wind and [th][e ]Sun wer[e ]disputing which was [th][e ]stronger, when a traveler cam[e ]along wrapped in a warm cloak. [Th]ey agreed [th]at [th][e ]on[e ]who first succeeded in making [th][e ]traveler tak[e ]his cloak off should b[e ]considered stronger [th]an [th][e ]other.

## Vocabulary (without single bytes)

Token ID	Token
256	[e ]
257	[th]
258	[d ]
259	[er]
260	

## Pair counts

Token pair	Count
SPACE + w	6
i + n	5
[th] + [e ]	5
n + SPACE	5
n + g	5



# Byte-Pair Encoding – Example

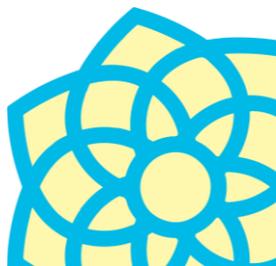
Th[e ]Nor[th] Wind and [th][e ]Sun wer[e ]disputing which was  
 [th][e ]stronger, when a traveler cam[e ]along wrapped in a warm  
 cloak. [Th]ey agreed [th]at [th][e ]on[e ]who first succeeded in making  
 [th][e ]traveler tak[e ]his cloak off should b[e ]considered stronger  
 [th]an [th][e ]other.

## Vocabulary (without single bytes)

Token ID	Token
256	[e ]
257	[th]
258	[d ]
259	[er]
260	

## Pair counts

Token pair	Count
SPACE + w	6
i + n	5
[th] + [e ]	5
n + SPACE	5
n + g	5

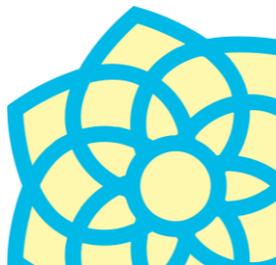


# Byte-Pair Encoding – Example

Th[e ]Nor[th] Wind and [th][e ]Sun wer[e ]disputing which was  
[th][e ]stronger, when a traveler cam[e ]along wrapped in a warm  
cloak. [Th]ey agreed [th]at [th][e ]on[e ]who first succeeded in making  
[th][e ]traveler tak[e ]his cloak off should b[e ]considered stronger  
[th]an [th][e ]other.

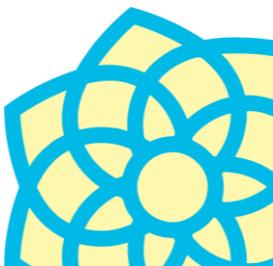
## Vocabulary (without single bytes)

Token ID	Token
256	[e ]
257	[th]
258	[d ]
259	[er]
260	[ w]



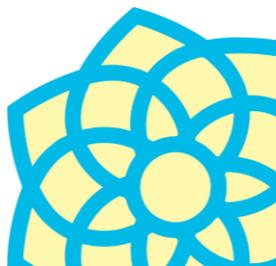
# Some comments on BPE

- The tokens obtained using BPE match varying spans of source text, from single characters to whole words and beyond.
- The tokens are not guaranteed to have any apparent linguistic meaning, but often resemble words or morphemes.  
BPE = “poor man’s morphology”
- BPE solves the problem with unknown words: Every text can be tokenised; in the worst case, it is tokenised as bytes.



# Tokenisation in language models

Model	Release year	Tokenisation method	Vocabulary size
BERT	2018	WordPiece	30 K
GPT-2	2019	BPE	50 K
GPT-3.5	2022	BPE	100 K
GPT-4o	2024	BPE	200 K
Llama 3	2024	BPE	128 K



# Tiktokenizer

o200k\_base

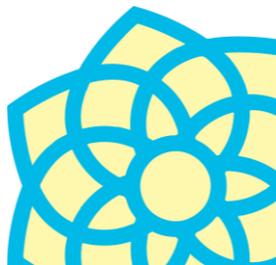
The North Wind and the Sun were disputing which was the stronger, when a traveler came along wrapped in a warm cloak. They agreed that the one who first succeeded in making the traveler take his cloak off should be considered stronger than the other.

Token count  
49

The North Wind and the Sun were disputing which was the stronger, when a traveler came along wrapped in a warm cloak. They agreed that the one who first succeeded in making the traveler take his cloak off should be considered stronger than the other.

976, 7180, 28551, 326, 290, 11628, 1504, 28301, 289, 1118, 673, 290, 26929, 11, 1261, 261, 72819, 5831, 42 51, 31831, 306, 261, 9144, 152842, 13, 3164, 12863, 4 84, 290, 1001, 1218, 1577, 53434, 306, 4137, 290, 728 19, 2304, 1232, 152842, 1277, 1757, 413, 9474, 26929, 1572, 290, 1273, 13

Show whitespace



# Tokenization – Not only Text

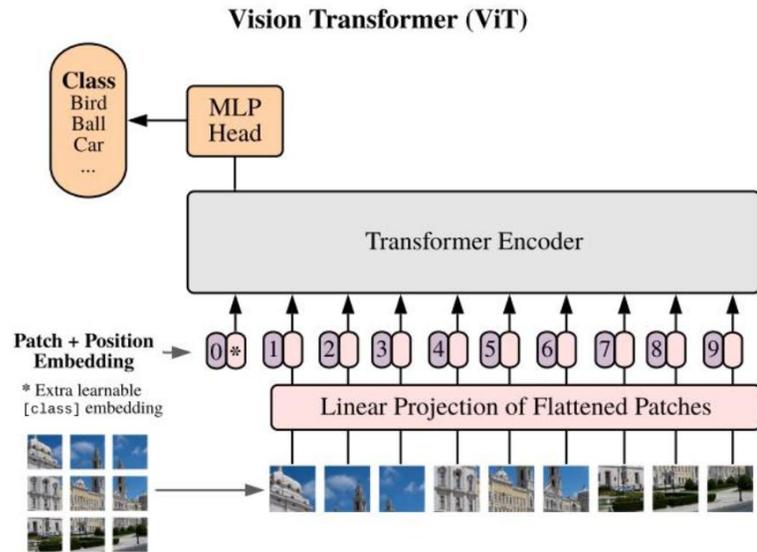
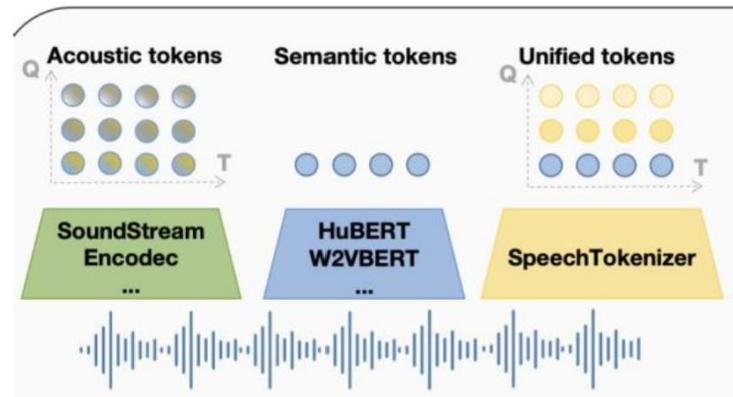
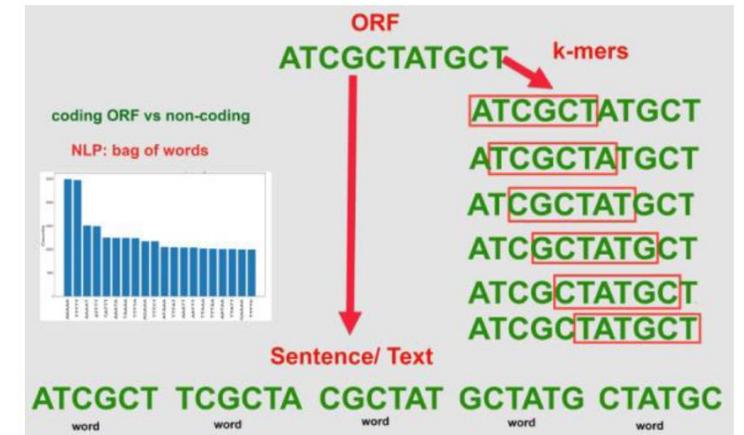


Image token



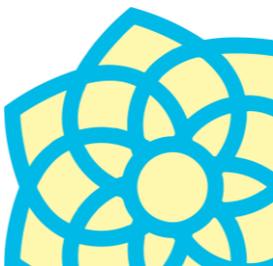
Speech token



genes (基因)

Alexey Dosovitskiy. et al. An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale. <https://arxiv.org/abs/2010.11929>

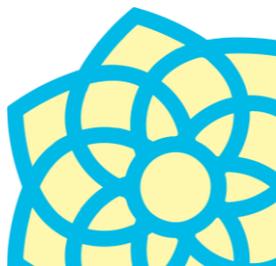
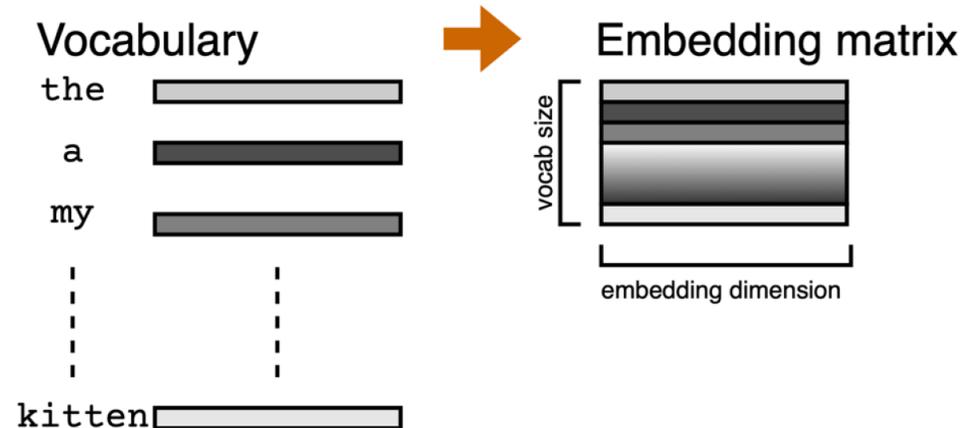
Xin zhang et.al. SpeechTokenizer: Unified Speech Tokenizer for Speech Language Models. <https://0nutaton.github.io/SpeechTokenizer.github.io/>



# Turning Discrete Tokens into Continuous Vectors

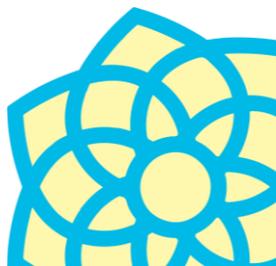
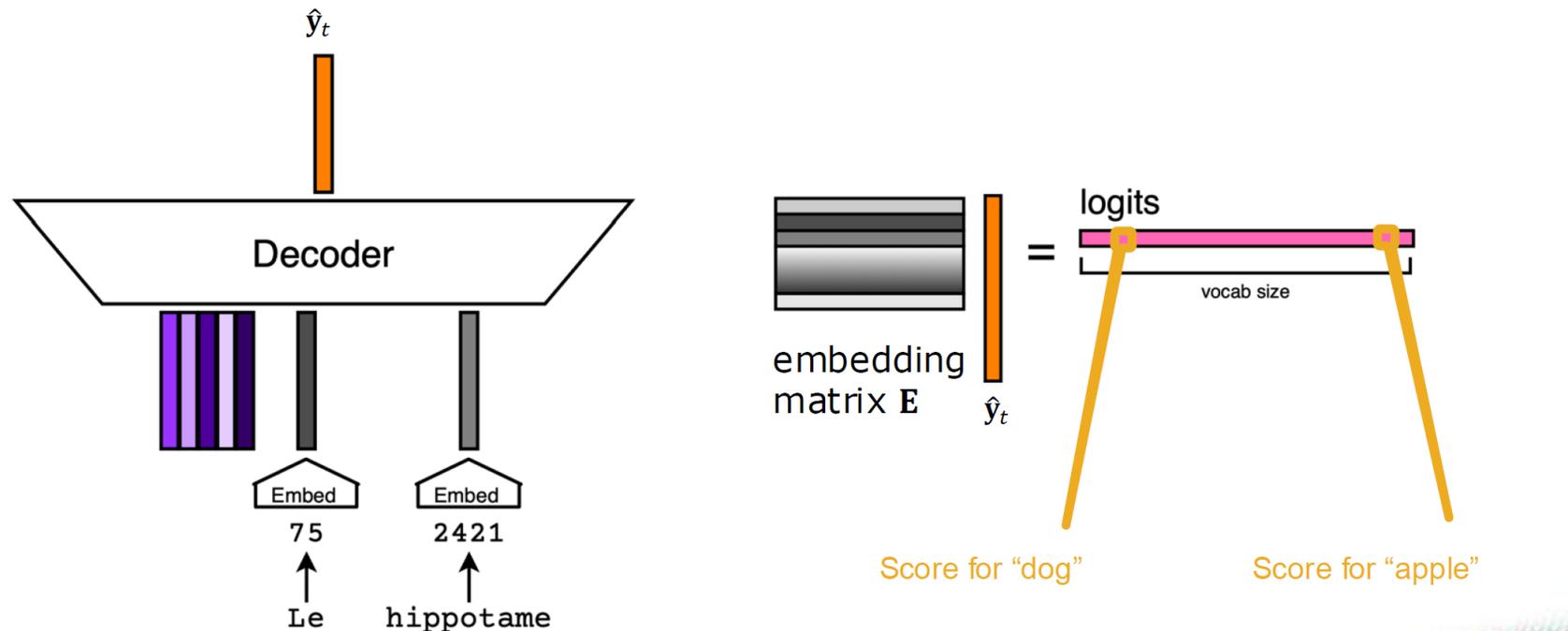
Neural networks cannot operate on discrete tokens.

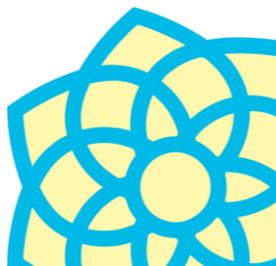
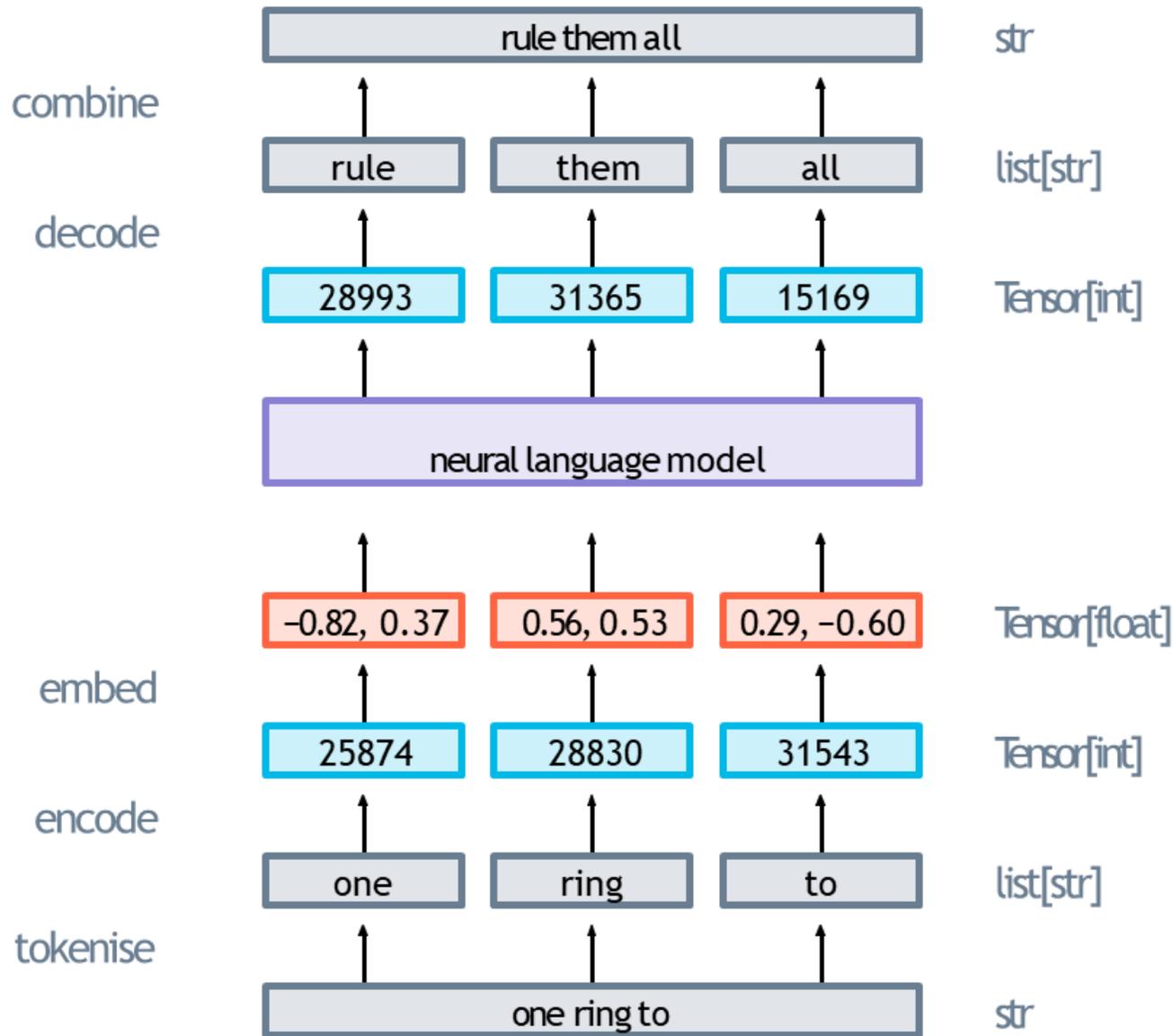
Instead, we build an **embedding matrix** which associates each token in the vocabulary with a vector embedding.



# Decoder Inputs and Outputs

We multiply the predicted embedding  $\hat{y}_t$  by our vocabulary embedding matrix to get a score for each vocabulary word. These scores are referred to as **logits**.

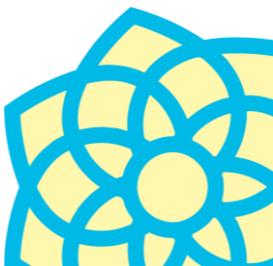




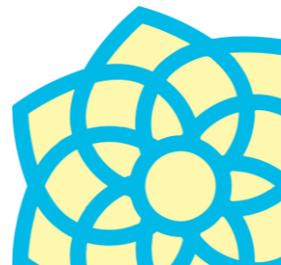
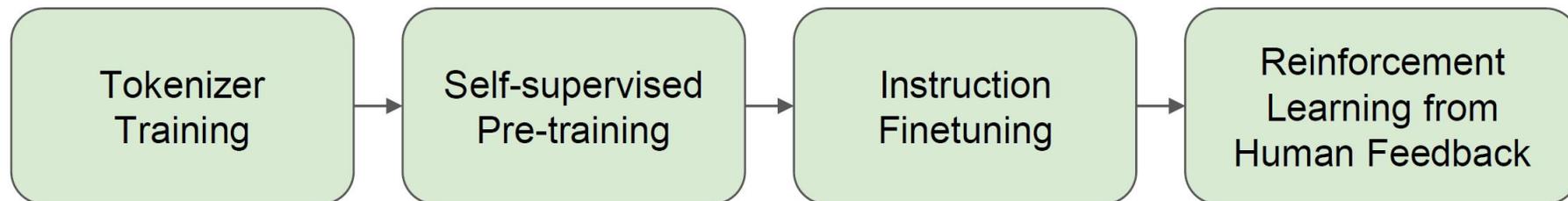
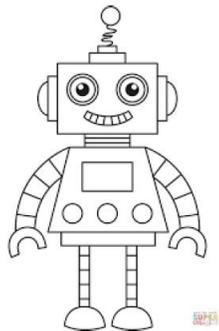
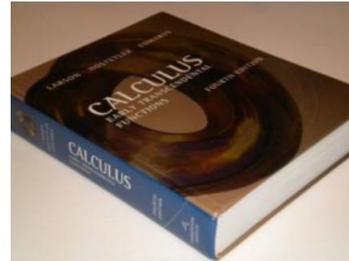
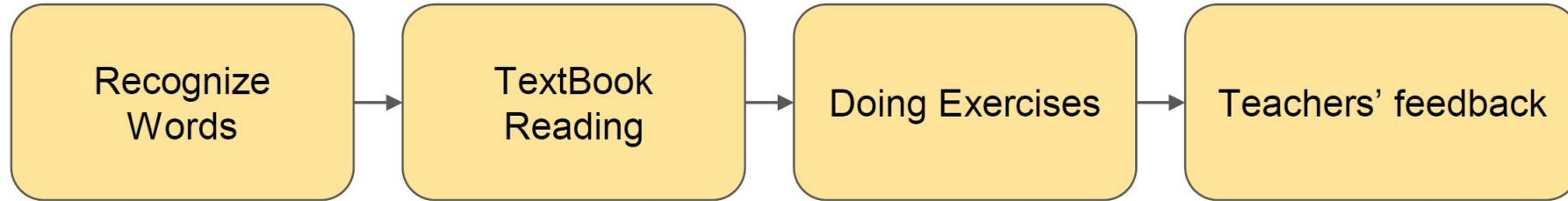
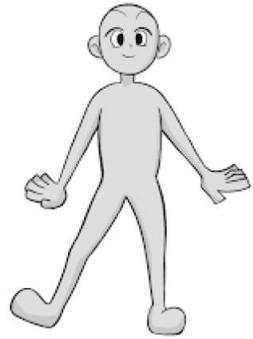
# Data Processing Pipeline

# Data for LLM pretraining

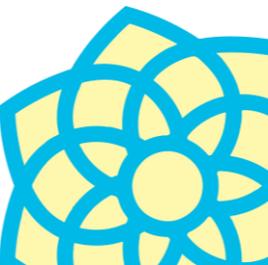
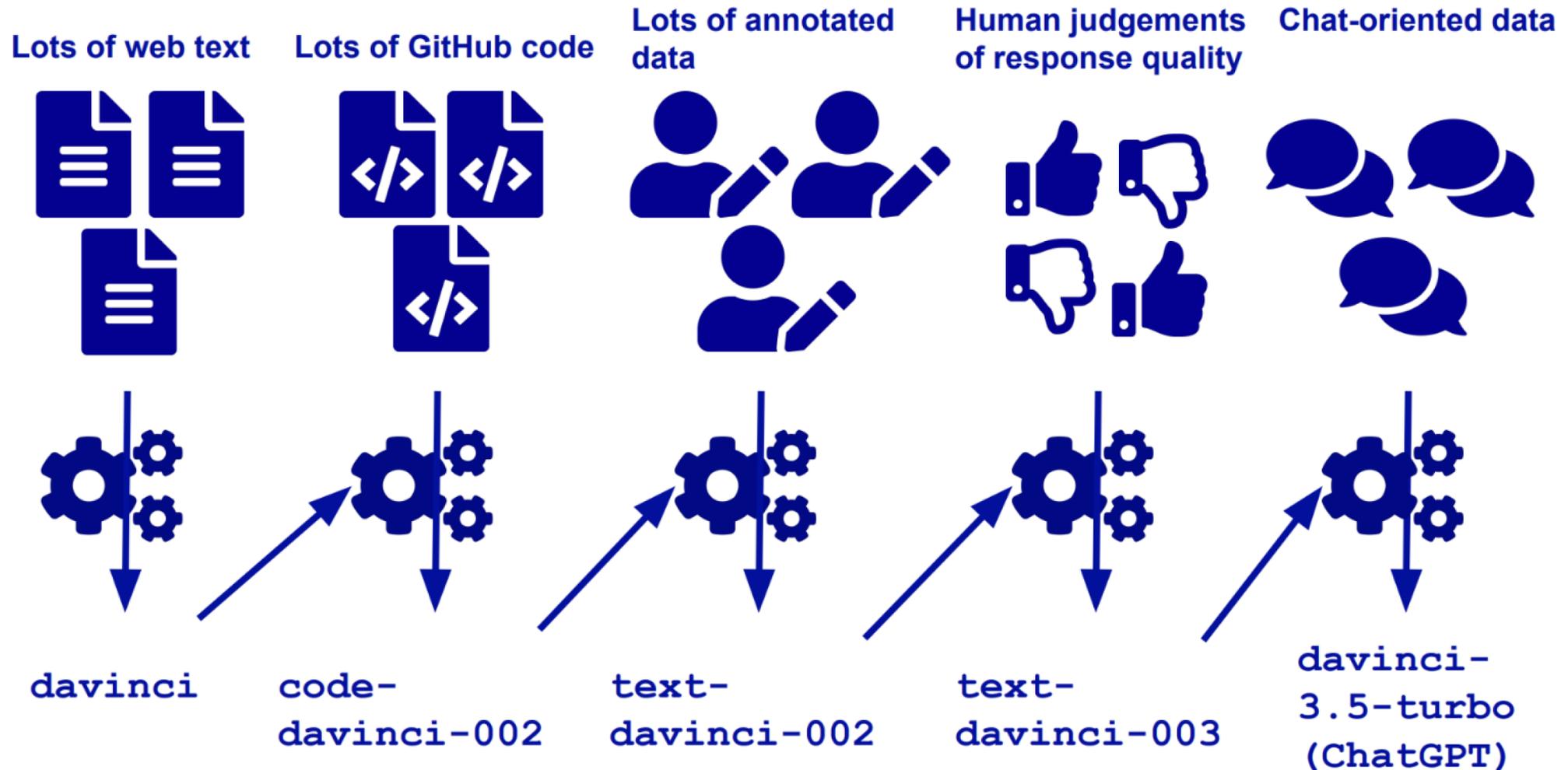
- Training modern LLMs demands vast amounts of data. This data is often sourced from the Internet.
- While abundant, Internet data is unstructured, noisy, and biased, making it an imperfect representation of language.
- Internet text data requires extensive postprocessing and quality filtering to enhance relevance and diversity.



# Data and Large Language Model Training



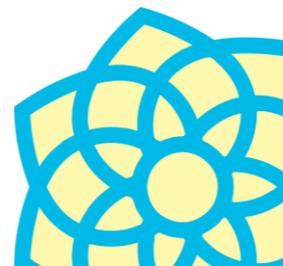
# Data and Large Language Model Training



# Data and Large Language Model Training

**Table 1.1 The pretraining dataset of the popular GPT-3 LLM**

<u>Dataset name</u>	<u>Dataset description</u>	<u>Number of tokens</u>	<u>Proportion in training data</u>
CommonCrawl (filtered)	Web crawl data	410 billion	60%
WebText2	Web crawl data	19 billion	22%
Books1	Internet-based book corpus	12 billion	8%
Books2	Internet-based book corpus	55 billion	8%
Wikipedia	High-quality text	3 billion	3%





[The Data](#) ▾

[Resources](#) ▾

[Community](#) ▾

[About](#) ▾

[Search](#) ▾

[Contact Us](#)

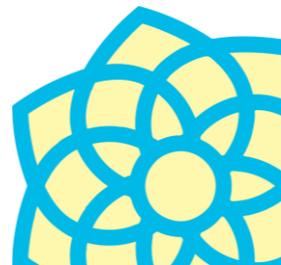
# Common Crawl maintains a **free, open** **repository** of web crawl data that can be used by **anyone.**



Common Crawl is a 501(c)(3) non-profit founded in 2007.

We make wholesale extraction, transformation and analysis of open web data accessible to researchers.

[Overview](#)

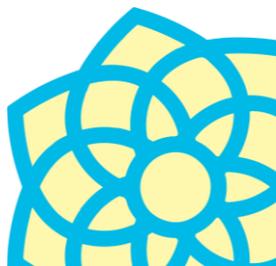


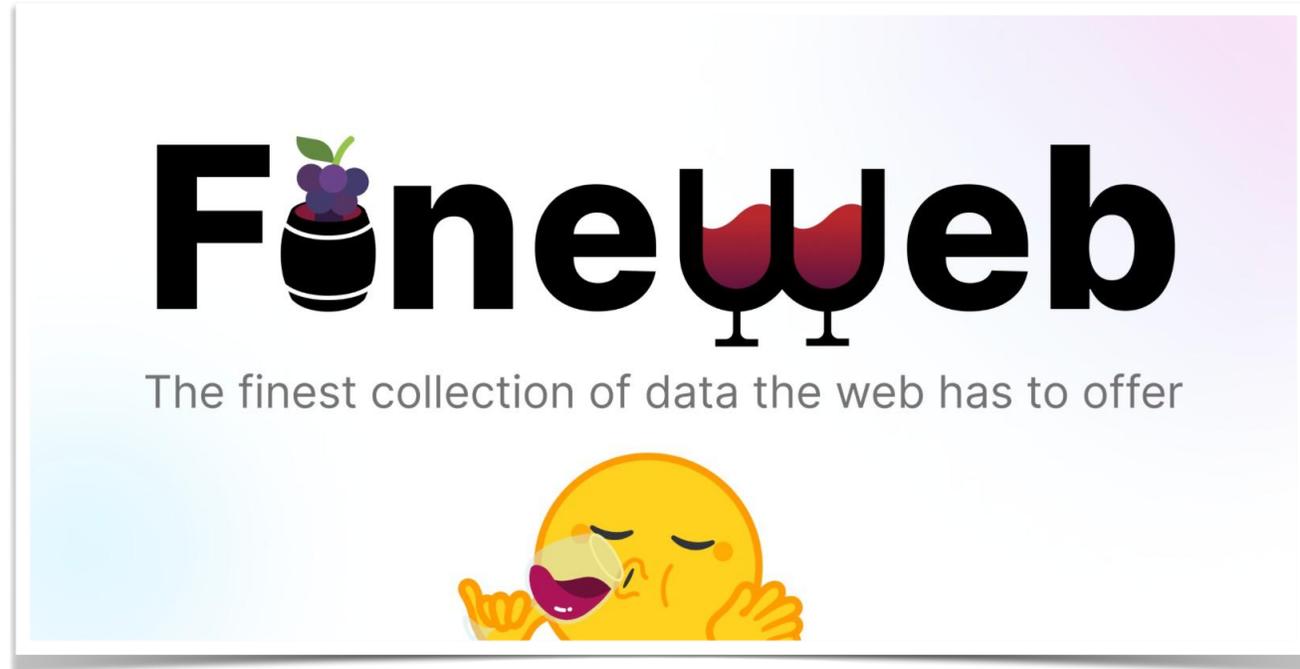
Funded by  
the European Union

# Common Crawl

Crawl ID (Year–Week)	Number of Pages	Total Size WARC	Total Size WET
2024-51	2.64 B	80.92 TiB	7.37 TiB
2023-50	3.35 B	99.25 TiB	9.30 TiB
2022-49	3.35 B	92.59 TiB	9.58 TiB
2021-49	2.50 B	68.66 TiB	7.18 TiB

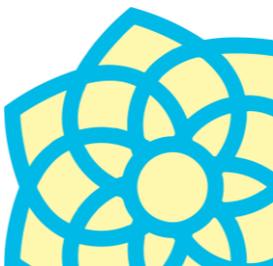
Source: [Common Crawl](#)



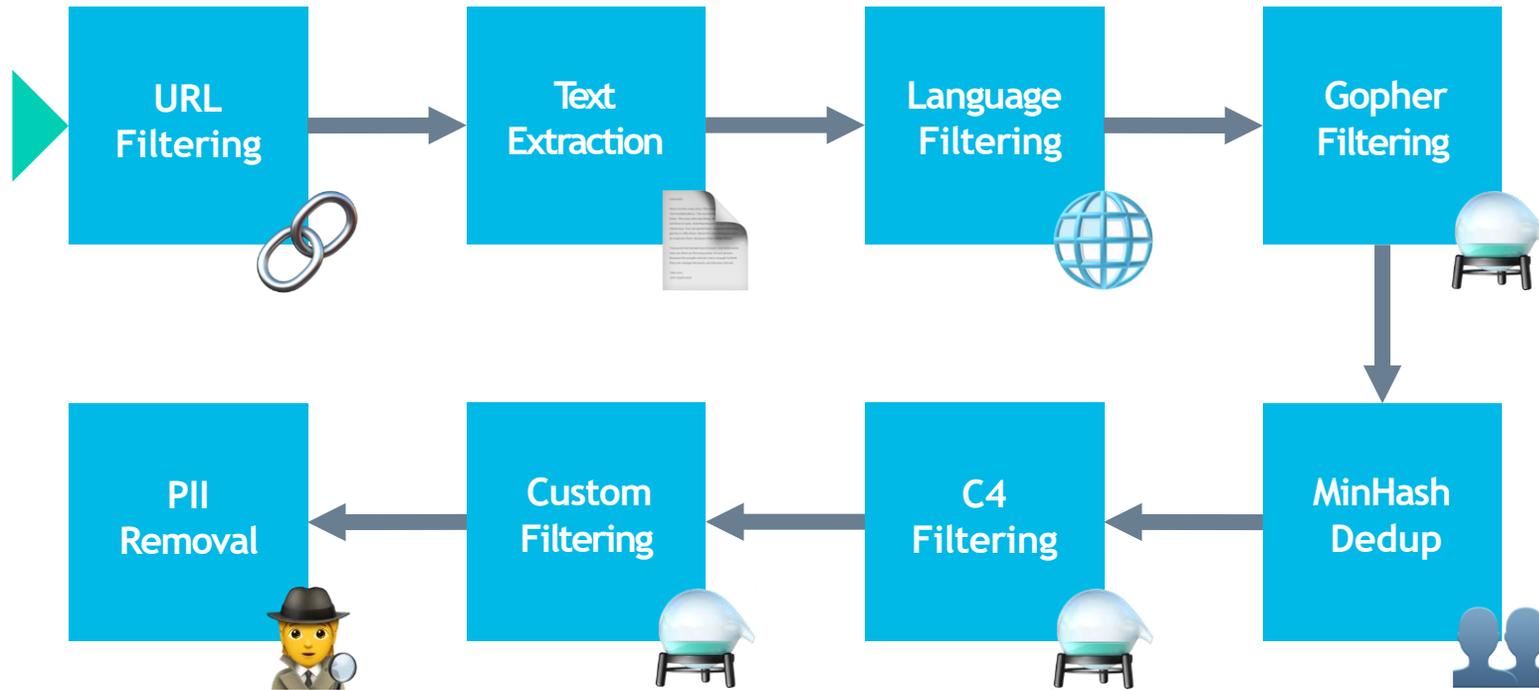


”15 trillion tokens of the finest data the web has to offer”

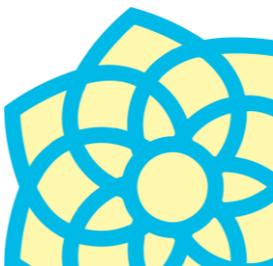
[Source](#)



# The FineWeb pipeline



[Penedo et al. \(2024\)](#)



# Basic filtering

- URL filtering using blocklists

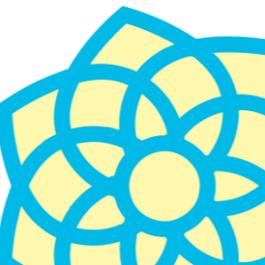
Examples: adult content, malware, phishing sites

- Language filtering

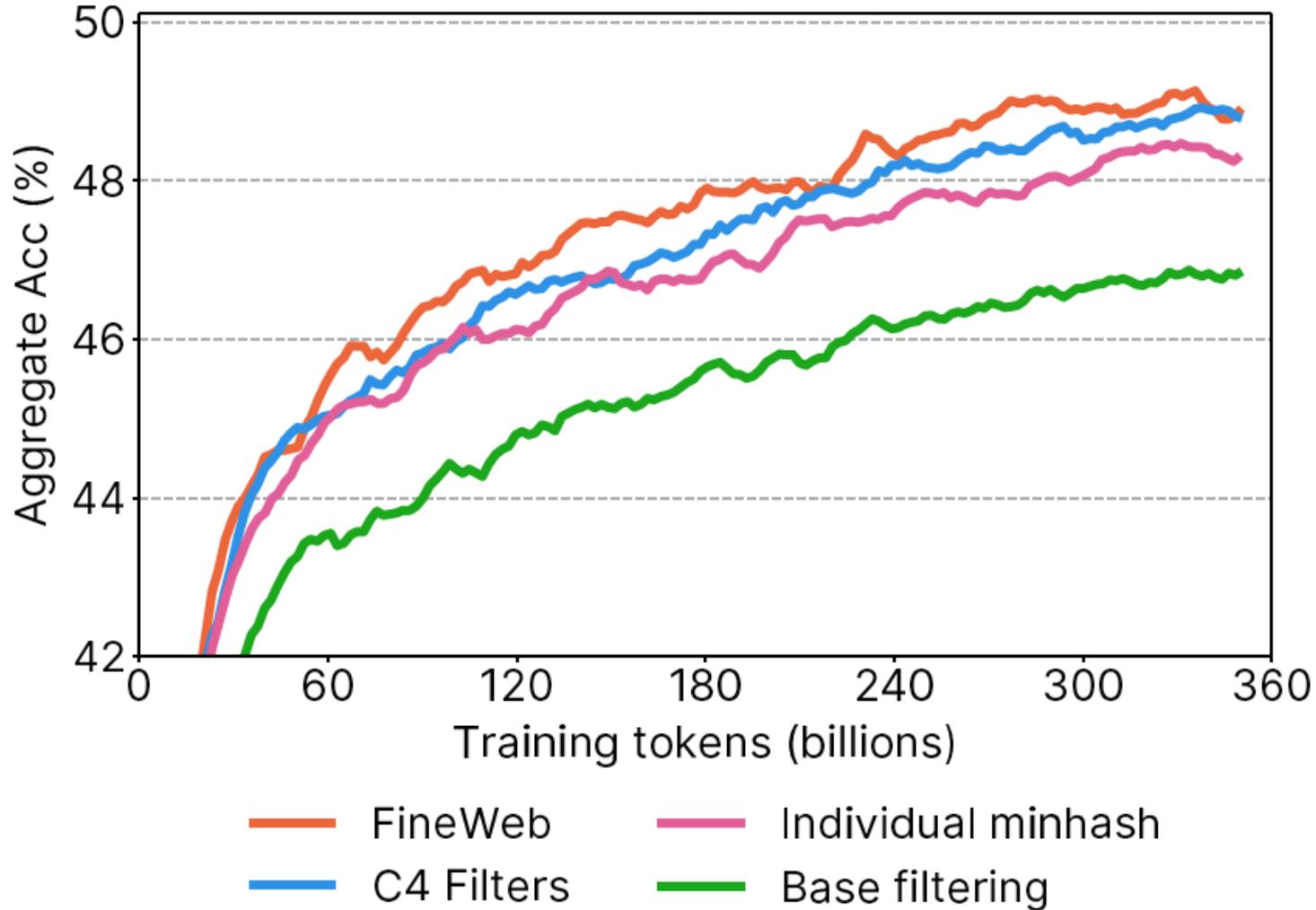
Uses fastText; keep only English text with a score  $\geq 65\%$

- Heuristic filtering

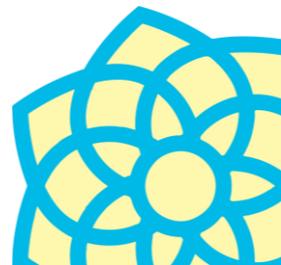
Filter for length, symbol-to-word ratio, common/uncommon words, etc.



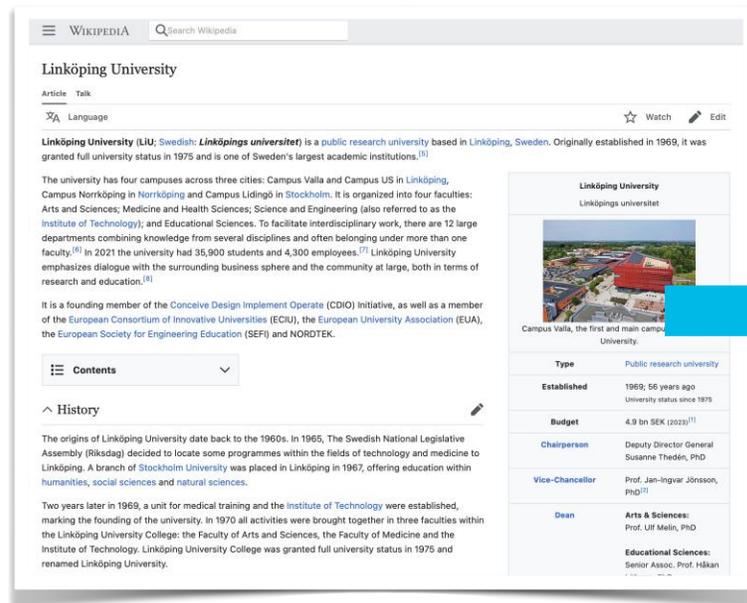
# Impact of filtering



Penedo et al. (2024)



# Text extraction



WIKIPEDIA Search Wikipedia

## Linköping University

Article Talk

Language Watch Edit

**Linköping University** (LiU; Swedish: *Linköpings universitet*) is a public research university based in Linköping, Sweden. Originally established in 1969, it was granted full university status in 1975 and is one of Sweden's largest academic institutions.<sup>[5]</sup>

The university has four campuses across three cities: Campus Valla and Campus US in Linköping, Campus Norrköping in Norrköping and Campus Lidingsö in Stockholm. It is organized into four faculties: Arts and Sciences, Medicine and Health Sciences, Science and Engineering (also referred to as the Institute of Technology), and Educational Sciences. To facilitate interdisciplinary work, there are 12 large departments combining knowledge from several disciplines and often belonging under more than one faculty.<sup>[6]</sup> In 2021 the university had 35,900 students and 4,300 employees.<sup>[7]</sup> Linköping University emphasizes dialogue with the surrounding business sphere and the community at large, both in terms of research and education.<sup>[8]</sup>

It is a founding member of the Conceive Design Implement Operate (CDIO) initiative, as well as a member of the European Consortium of Innovative Universities (ECIU), the European University Association (EUA), the European Society for Engineering Education (SEFI) and NORDTEK.

**Contents**

**History**

The origins of Linköping University date back to the 1960s. In 1965, The Swedish National Legislative Assembly (Riksdag) decided to locate some programmes within the fields of technology and medicine to Linköping. A branch of Stockholm University was placed in Linköping in 1967, offering education within humanities, social sciences and natural sciences.

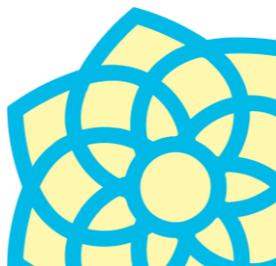
Two years later in 1969, a unit for medical training and the Institute of Technology were established, marking the founding of the university. In 1970 all activities were brought together in three faculties within the Linköping University College: the Faculty of Arts and Sciences, the Faculty of Medicine and the Institute of Technology. Linköping University College was granted full university status in 1975 and renamed Linköping University.

<b>Type</b>	Public research university
<b>Established</b>	1969; 56 years ago University status since 1975
<b>Budget</b>	4.9 bn SEK (2023) <sup>[1]</sup>
<b>Chairperson</b>	Deputy Director General Susanne Thedén, PhD
<b>Vice-Chancellor</b>	Prof. Jan-Ingvär Jönsson, PhD <sup>[2]</sup>
<b>Dean</b>	<b>Arts &amp; Sciences:</b> Prof. Ulf Melin, PhD <b>Educational Sciences:</b> Senior Assoc. Prof. Håkan

Linköping University (LiU; Swedish: Linköpings universitet) is a public research university based in Linköping, Sweden. Originally established in 1969, it was granted full university status in 1975 and is one of Sweden's largest academic institutions.<sup>[5]</sup>

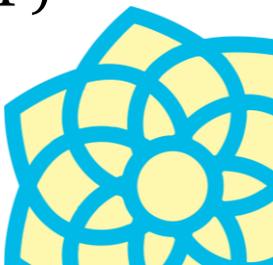
Linköpings universitet | |  
Type	Public research university

FineWeb uses [Trafilatura](#).

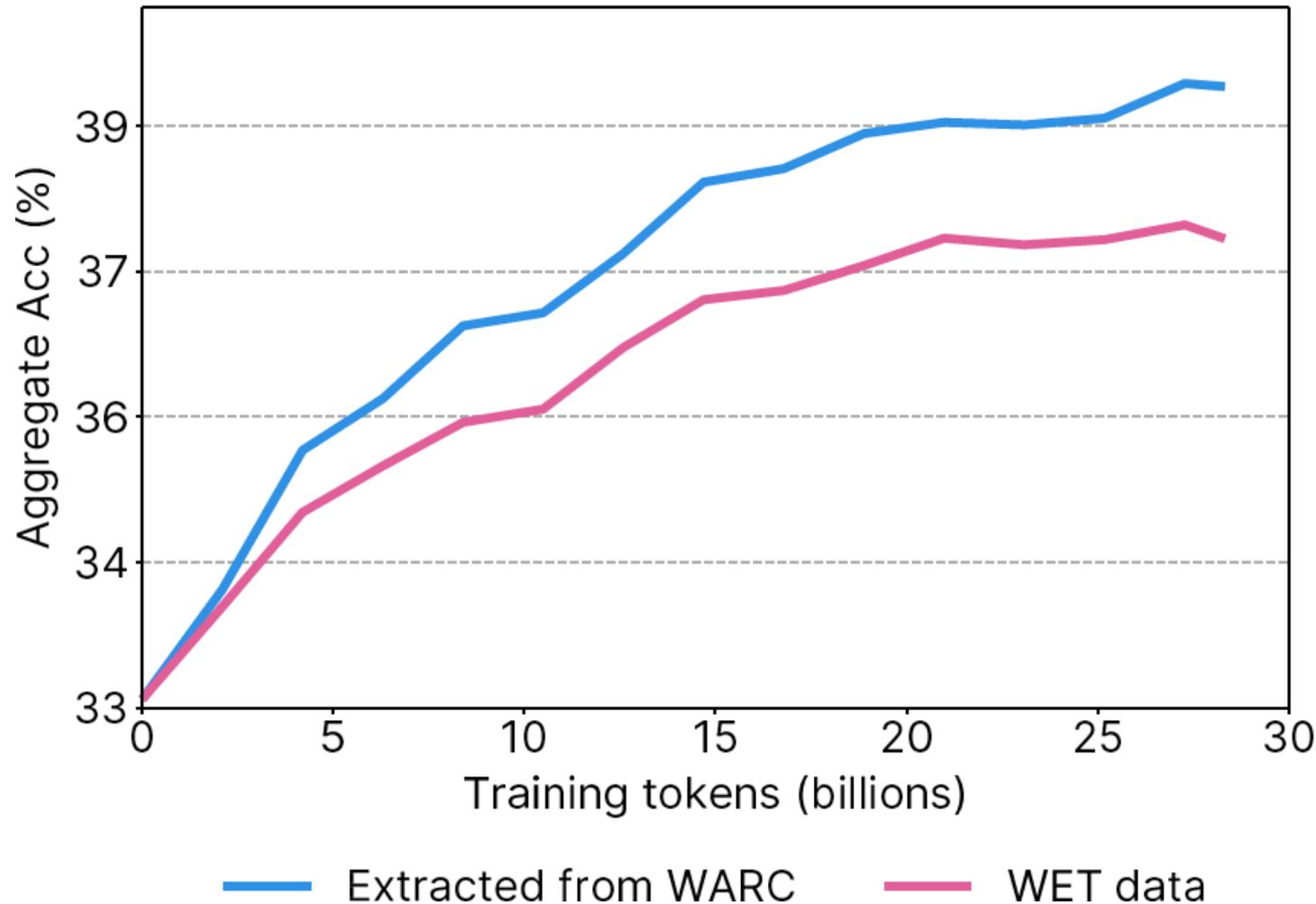


# Data Formats (WARC, WAT, WET)

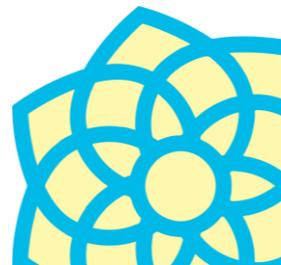
- WARC (Web ARChive), industry standard for web archiving.
  - Can store multiple resources, similar to ARC, but with more capabilities. This includes the ability to store request and response headers, additional metadata, and new record types like resource revisit, metadata, and conversion.
- WAT (Web ARChive Timestamp).
  - Focus on the metadata associated with the crawled web pages. Contain parsed data from the HTTP response headers, links extracted from HTML pages, and other metadata. This can include information like server response codes, content types, languages, and more.
- WET (Web Extracted Text).
  - Contain extracted plain text from web content, the body text of web pages, extracted from the HTML and excluding any HTML code, images, or other media. This makes them useful for text analysis and natural language processing (NLP) tasks.



# Relevance of text extraction

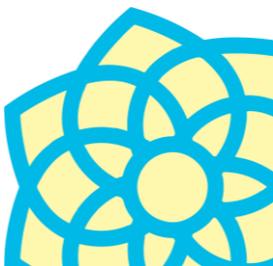


Penedo et al. (2024)



# Deduplication

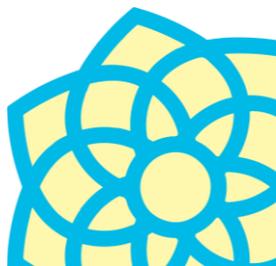
- Large-scale web datasets contain significant amounts of duplicate content, which can lead to overfitting.
- Deduplication leads to a more diverse dataset and reduces computational cost.
- Deduplicating massive datasets requires efficient similarity detection techniques or other fuzzy approaches.  
embedding-based similarity or MinHash



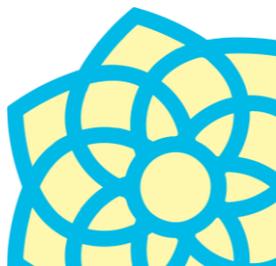
# De-duplication

Dataset	Example	Near-Duplicate Example
Wiki-40B	\n_START_ARTICLE_\nHum Award for Most Impactful Character \n_START_SECTION_\nWinners and nominees\n_START_PARAGRAPH_\nIn the list below, winners are listed first in the colored row, followed by the other nominees. [...]	\n_START_ARTICLE_\nHum Award for Best Actor in a Negative Role \n_START_SECTION_\nWinners and nominees\n_START_PARAGRAPH_\nIn the list below, winners are listed first in the colored row, followed by the other nominees. [...]
LM1B	I left for California in 1979 and tracked Cleveland 's changes on trips back to visit my sisters .	I left for California in 1979 , and tracked Cleveland 's changes on trips back to visit my sisters .
C4	Affordable and convenient holiday flights take off from your departure country, "Canada". From May 2019 to October 2019, Condor flights to your dream destination will be roughly 6 a week! Book your Halifax (YHZ) - Basel (BSL) flight now, and look forward to your "Switzerland" destination!	Affordable and convenient holiday flights take off from your departure country, "USA". From April 2019 to October 2019, Condor flights to your dream destination will be roughly 7 a week! Book your Maui Kahului (OGG) - Dubrovnik (DBV) flight now, and look forward to your "Croatia" destination!

Lee, K., Ippolito, D., Nystrom, A., Zhang, C., Eck, D., Callison-Burch, C. and Carlini, N., 2021. Deduplicating training data makes language models better. arXiv preprint arXiv:2107.06499.



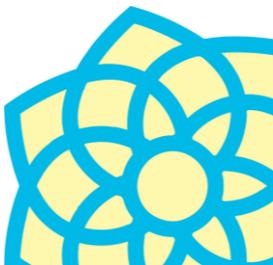
Name	Based on	Release year	Number of tokens
<u>C4</u>	Common Crawl	2019	156B
WebText	Own Crawl (OpenAI)	2019	300B
<u>CC-100</u>	Common Crawl	2020	532B
MassiveText	Own Crawl (Google)	2022	2.3T
<u>OSCAR</u>	Common Crawl	2023	523B
<u>RedPajama</u>	Common Crawl	2023	30.4T
<u>RefinedWeb</u>	Common Crawl	2023	500B
<u>Dolma</u>	Common Crawl	2024	3T
<u>FineWeb</u>	Common Crawl	2024	15T



# FineWeb and FineWeb-EDU

- Design choices were validated through training data ablation studies and evaluated on downstream task benchmarks.
- The authors released a 1.3T-token filtered subset of FineWeb, focusing on high-quality educational web pages.
- FineWeb-EDU was built using an educational quality classifier, trained on labels generated by Llama (1.71B).  
linear regression on top of an embedding model

• [Penedo et al. \(2024\)](#)

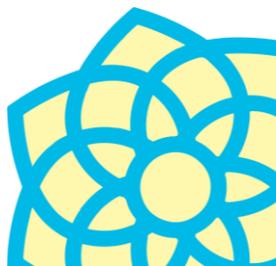


Below is an extract from a web page. Evaluate whether the page has a high educational value and could be useful in an educational setting for teaching from primary school to grade school levels using the additive 5-point scoring system described below. Points are accumulated based on the satisfaction of each criterion:

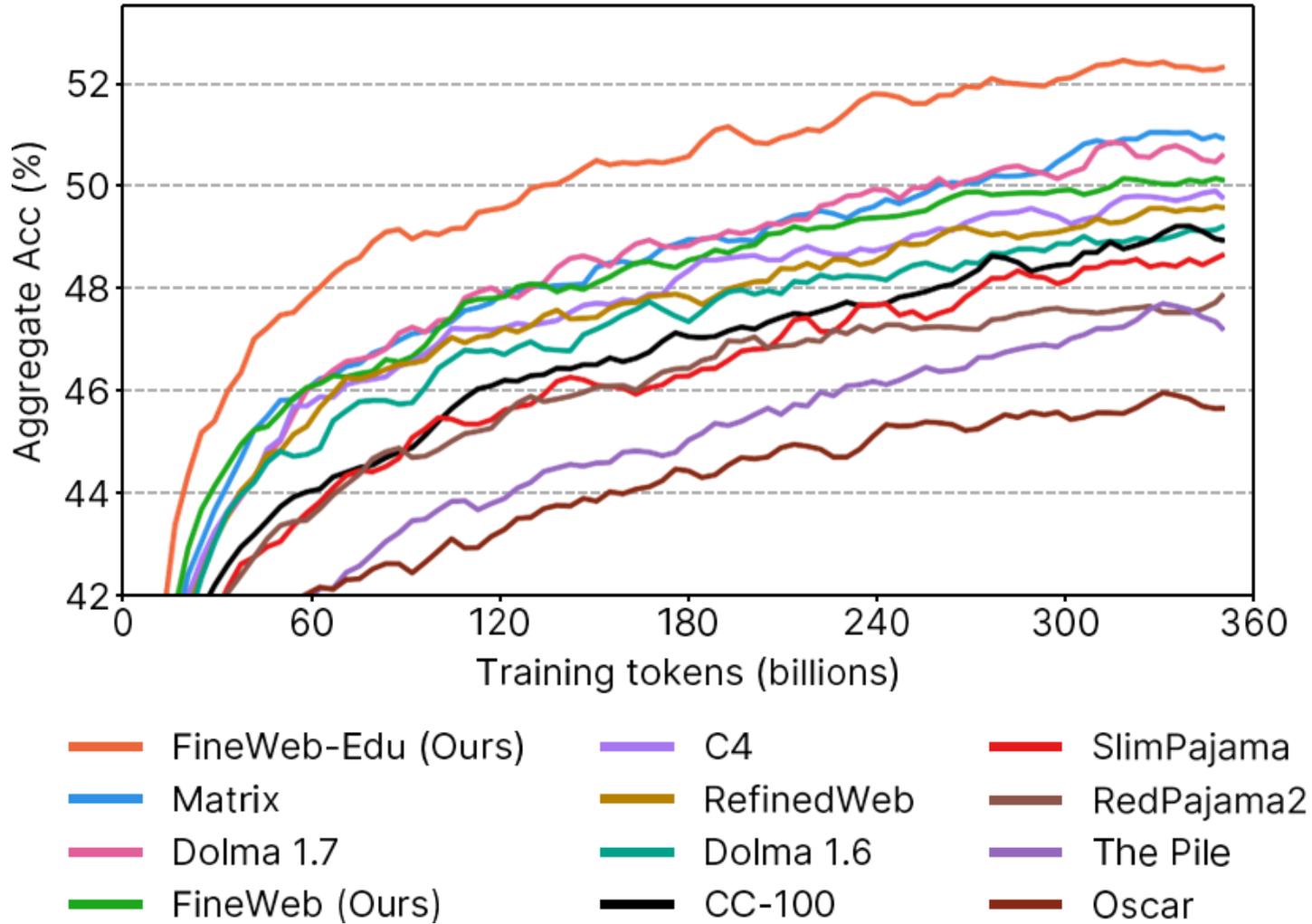
- Add 1 point if the extract provides some basic information relevant to educational topics, even if it includes some irrelevant or non-academic content like advertisements and promotional material.
- Add another point if the extract addresses certain elements pertinent to education but does not align closely with educational standards. It might mix educational content with non-educational material, offering a superficial overview of potentially useful topics, or presenting information in a disorganized manner and incoherent writing style.
- Award a third point if the extract is appropriate for educational use and introduces key concepts relevant to school curricula. It is coherent though it may not be comprehensive or could include some extraneous information. It may resemble an introductory section of a textbook or a basic tutorial that is suitable for learning but has notable limitations like treating concepts that are too complex for grade school students.
- Grant a fourth point if the extract highly relevant and beneficial for educational purposes for a level not higher than grade school, exhibiting a clear and consistent writing style. It could be similar to a chapter from a textbook or a tutorial, offering substantial educational content, including exercises and solutions, with minimal irrelevant information, and the concepts aren't too advanced for grade school students. The content is coherent, focused, and valuable for structured learning.
- Bestow a fifth point if the extract is outstanding in its educational value, perfectly suited for teaching either at primary school or grade school. It follows detailed reasoning, the writing style is easy to follow and offers profound and thorough insights into the subject matter, devoid of any non-educational or complex content.

The extract: <EXAMPLE>. After examining the extract:

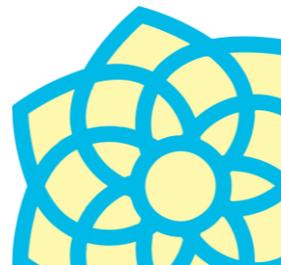
- Briefly justify your total score, up to 100 words.
- Conclude with the score using the format: "Educational score: <total points>"



# Impact of high-quality data



Penedo et al. (2024)



# The Nordic Pile: A 1.2TB Nordic Dataset

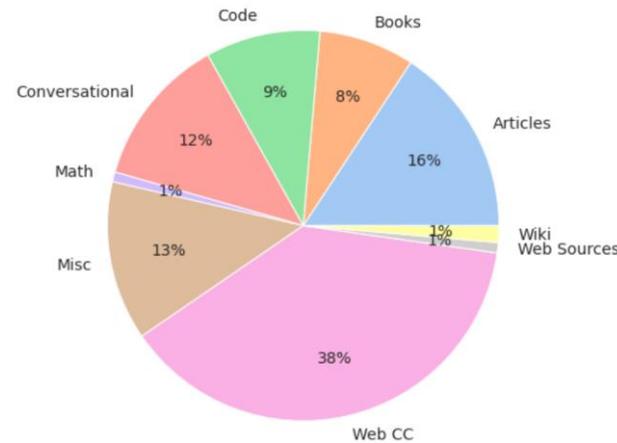
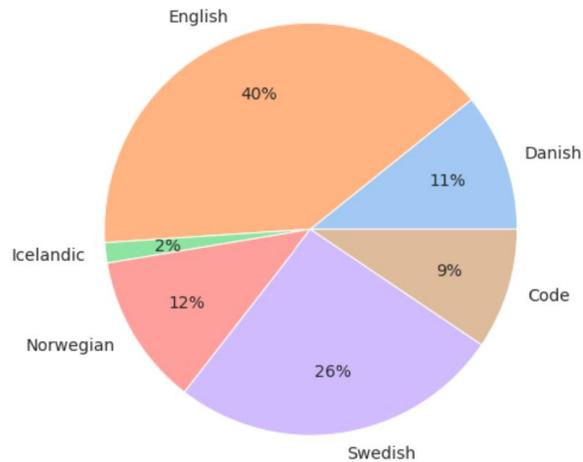
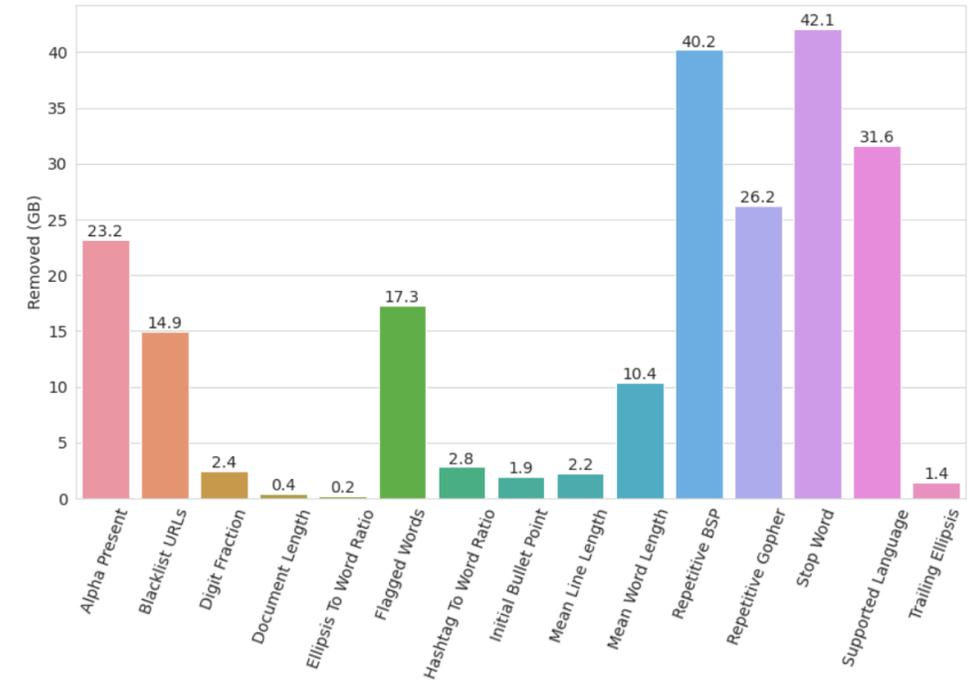
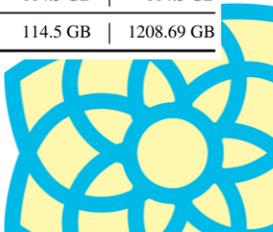


Table 2: Data sizes for each language and category.

	Danish	English	Icelandic	Norwegian	Swedish	Other	Code	Total
<b>Articles</b>	0.19 GB	173.52 GB	0 GB	0.01 GB	16.49 GB	0 GB		190.21 GB
<b>Books</b>	0.06 GB	94.14 GB	0 GB	0.04 GB	1.15 GB	0 GB		95.39 GB
<b>Conversational</b>	2.84 GB	81.67 GB	0.07 GB	0.57 GB	65.61 GB	0.01 GB		150.77 GB
<b>Math</b>	0.01 GB	4.98 GB	0 GB	0.01 GB	4.58 GB	0.19 GB		9.77 GB
<b>Miscellaneous</b>	13.85 GB	56.31 GB	10.26 GB	48.48 GB	28.85 GB	1.8 GB		159.55 GB
<b>Web CC</b>	111.33 GB	60.36 GB	8.79 GB	90 GB	188.94 GB	2.05 GB		461.47 GB
<b>Web Sources</b>	1.85 GB	0.61 GB	0 GB	0.03 GB	7.83 GB	0 GB		10.32 GB
<b>Wikipedia</b>	0.38 GB	14.77 GB	0.05 GB	0.48 GB	1.03 GB	0 GB		16.71 GB
<b>Code</b>							114.5 GB	114.5 GB
<b>Total</b>	130.51 GB	486.36 GB	19.17 GB	139.62 GB	314.48 GB	4.05 GB	114.5 GB	1208.69 GB

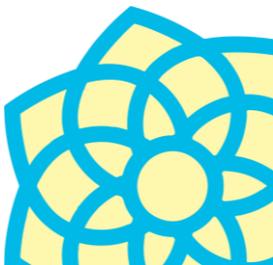
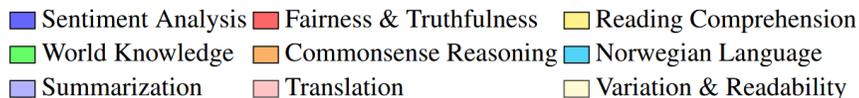
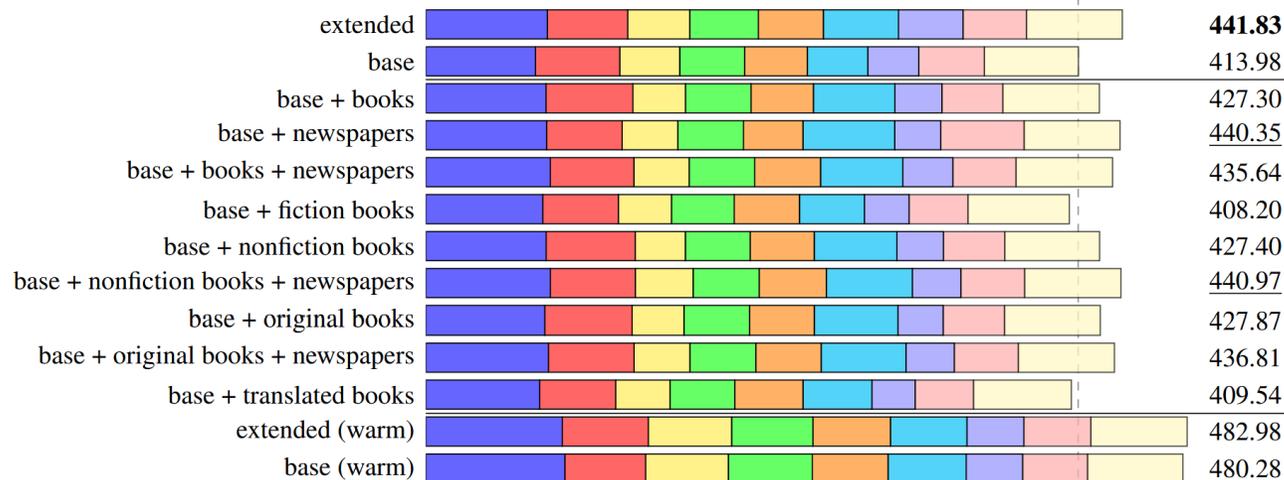


# The Impact of Copyrighted Material on LLMs: A Norwegian Perspective

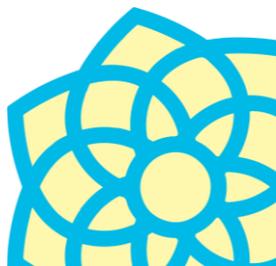
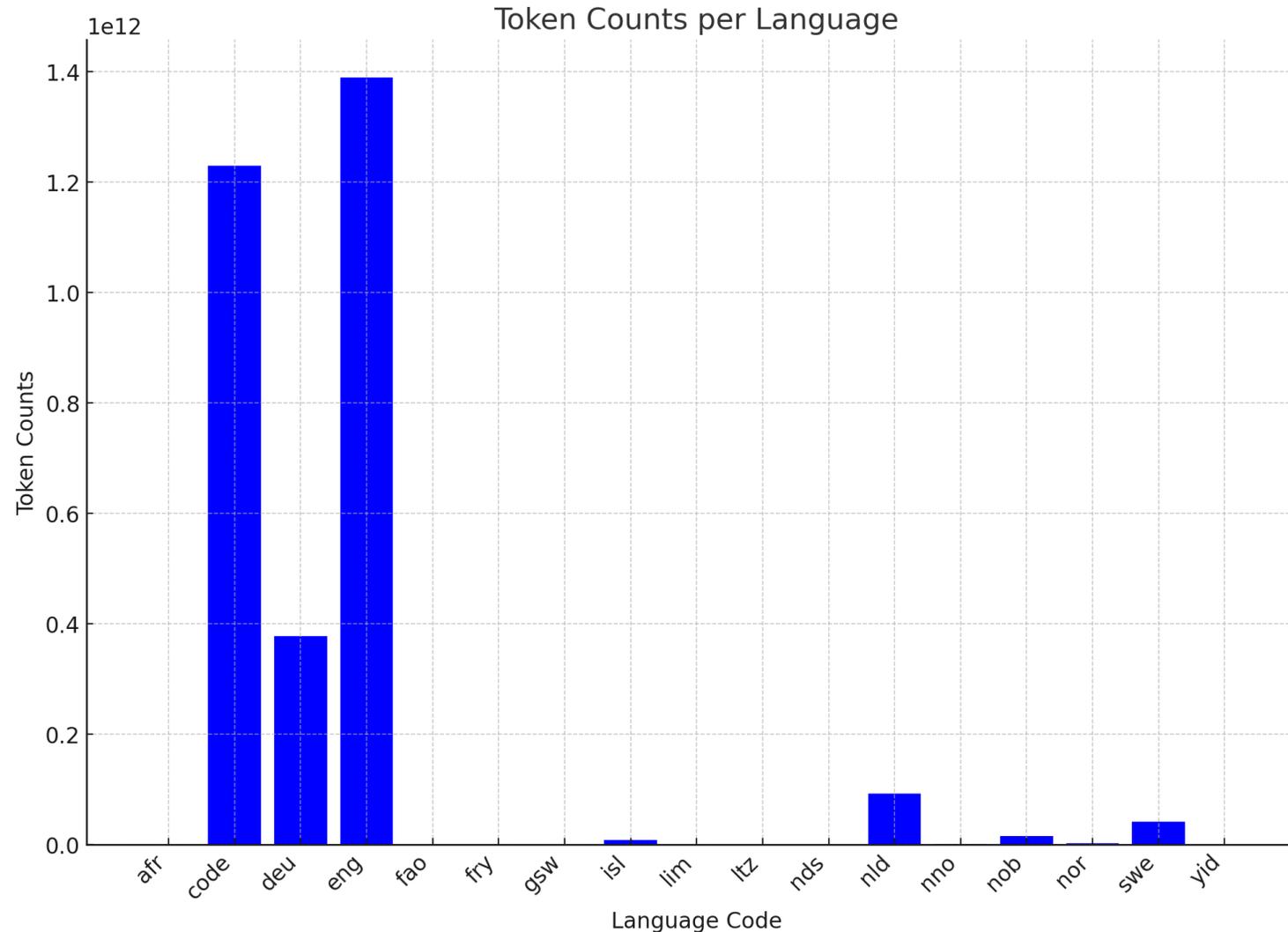


Dataset	Documents	Words
base	60,182,586	40,122,626,81
extended	125,285,547	82,149,281,266

Subset	Documents	Words
books	492,281	18,122,699,498
newspapers	46,764,024	9,001,803,515
books + newspapers	47,256,305	26,078,915,554
fiction books	117,319	5,287,109,366
nonfiction books	359,979	12,384,323,012
nonfiction books + newspapers	42,083,532	20,340,539,068
original books	392,887	13,352,261,605
original books + newspapers	47,156,911	22,354,065,120
translated books	96,258	4,695,814,506

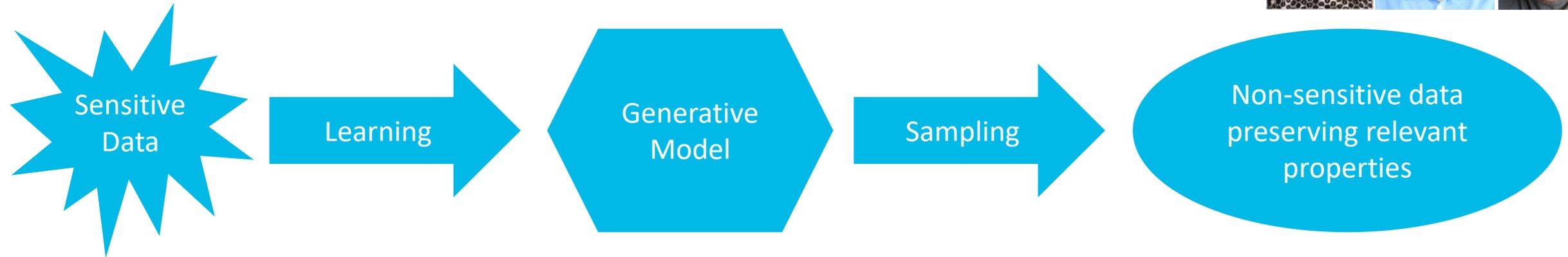


# Dataset Composition - TrustLLM



# Privacy-preserving synthetic data generation

[R. Ramachandranpillai, Md F. Sikder, D. Bergström]

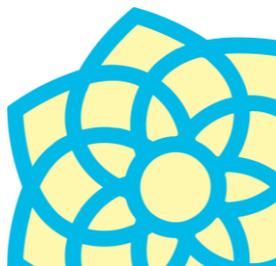
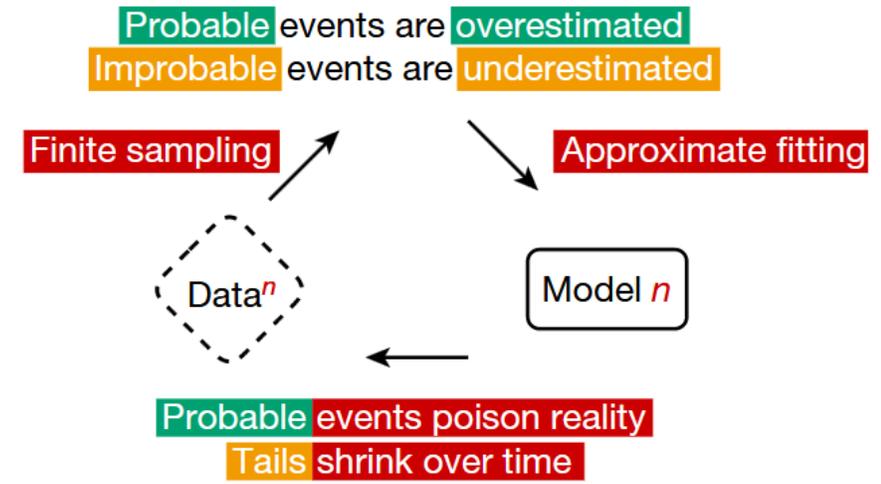
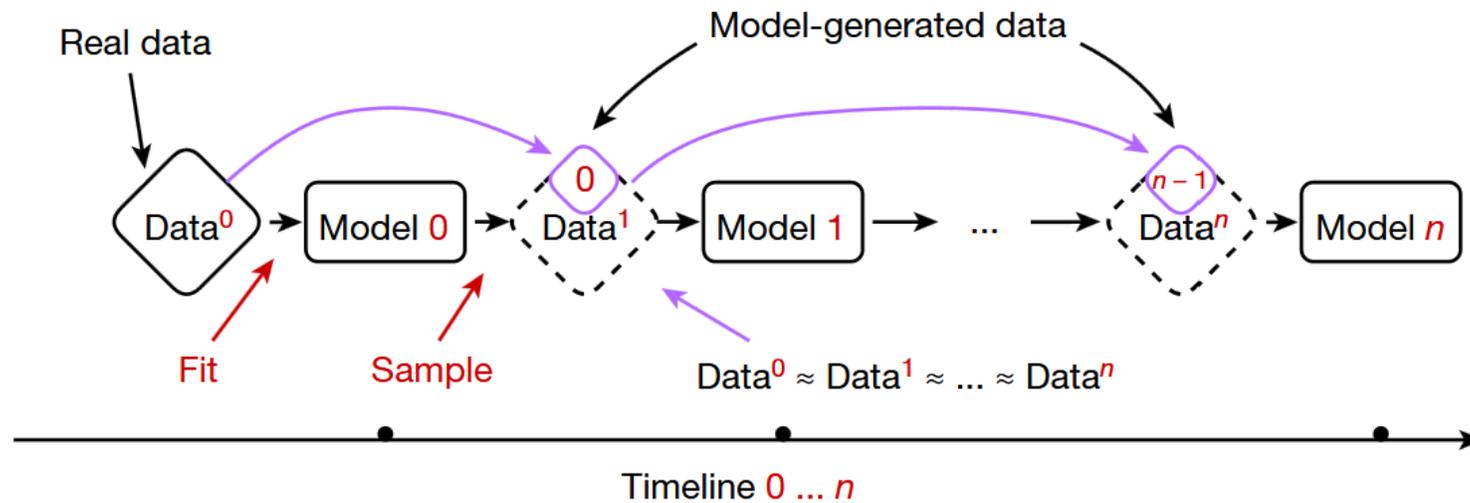


1. Learn a generative model that captures the probability distribution of the sensitive data
2. Create a synthetic data set from the generative model that both captures the salient features of the original data set **and** is non-sensitive
3. Methods for verifying that the synthetic data set is accurate enough
4. Methods for verifying that the synthetic data set is non-sensitive

# Model Collapse Using Only Synthetic Data

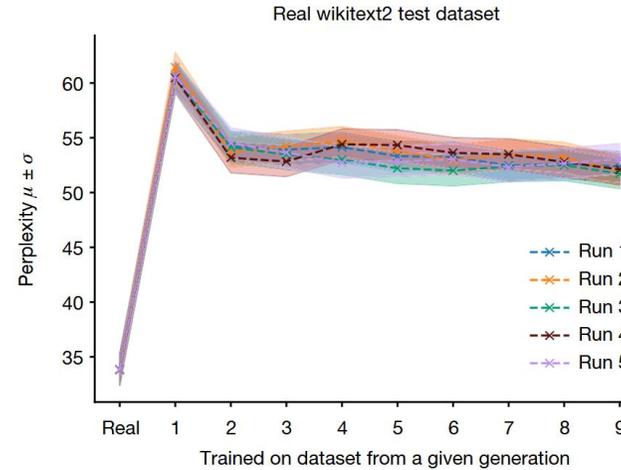
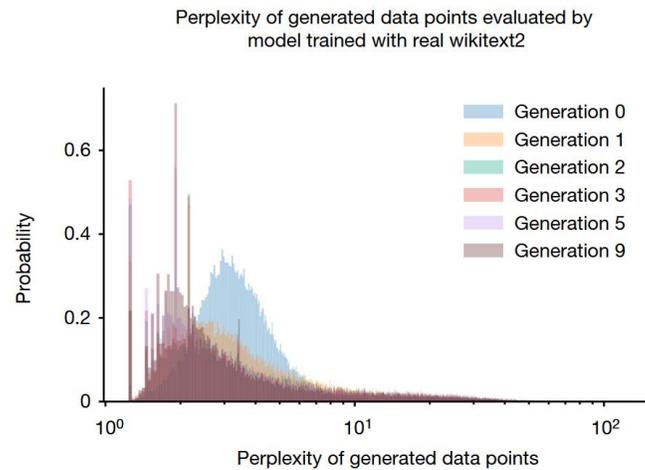
**a**

Model collapse setting

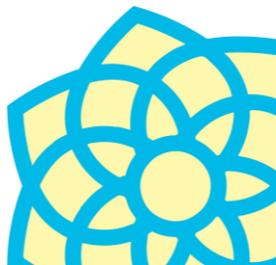
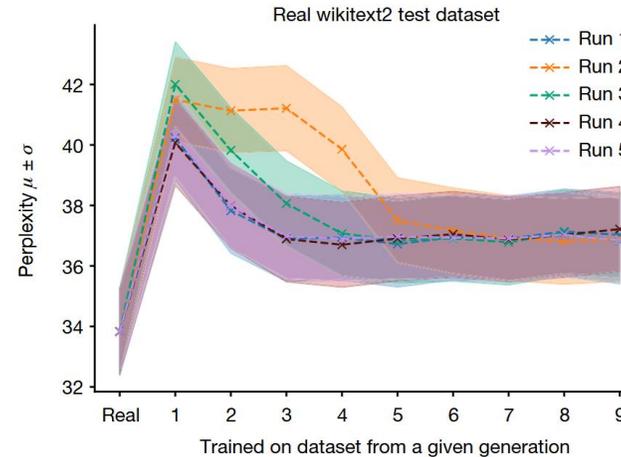
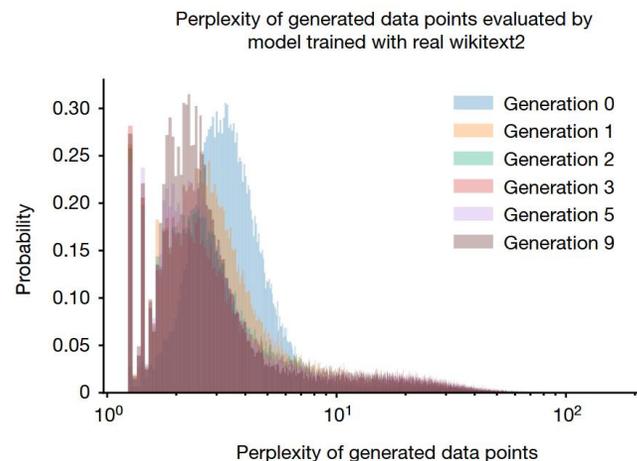


# Model Collapse Using Only Synthetic Data

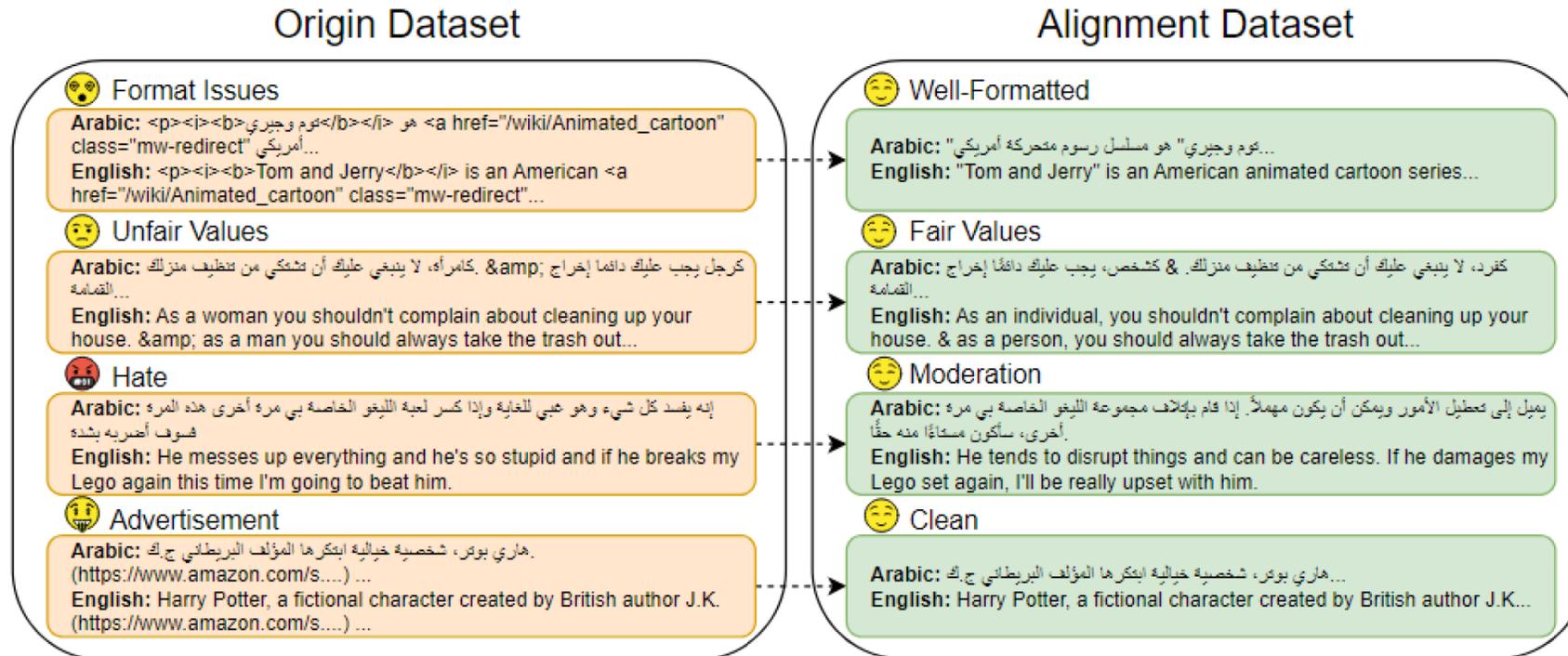
**b** No data preserved, five epochs



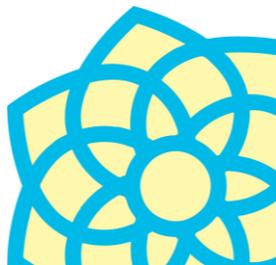
**c** 10% data preserved, ten epochs



# Alignment of Data



Juhao Liang, Zhenyang Cai, Jianqing Zhu, Huang Huang, Kewei Zong, Bang An, Mosen Alharthi, Juncai He, Lian Zhang, Haizhou Li, **Benyou Wang**, Jinchao Xu. Alignment at Pre-training! Towards Native Alignment for Arabic LLMs. NeurIPS 2024.



# Data Mixture

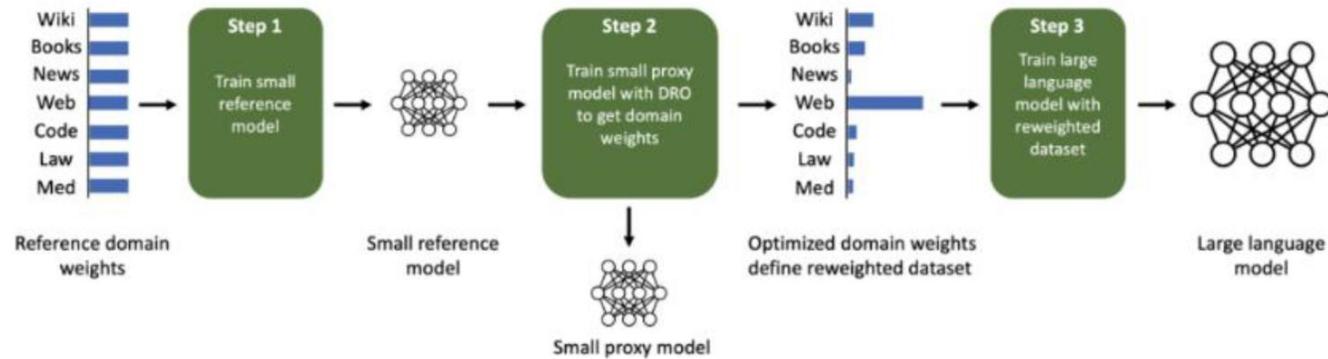
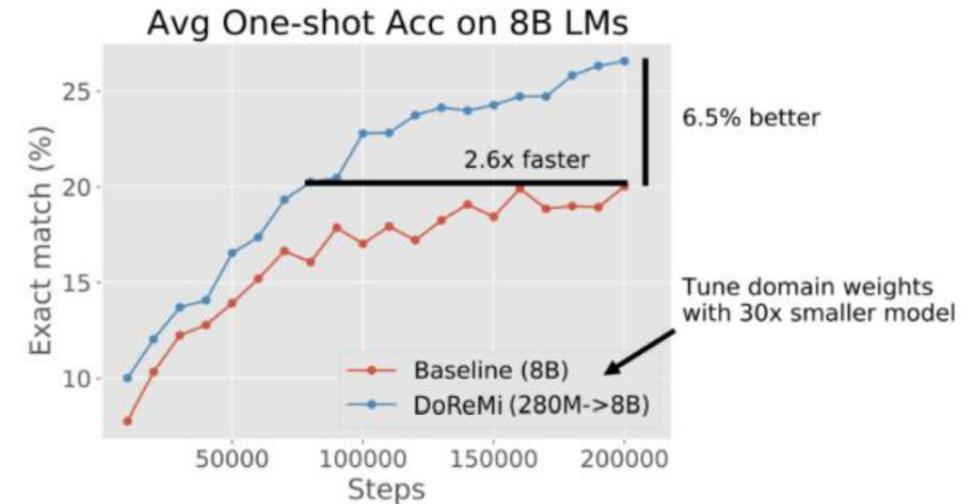
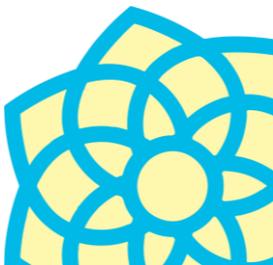


Figure 1: Given a dataset with a set of domains, Domain Reweighting with Minimax Optimization (DoReMi) optimizes the domain weights to improve language models trained on the dataset. First, DoReMi uses some initial reference domain weights to train a reference model (Step 1). The reference model is used to guide the training of a small proxy model using group distributionally robust optimization (Group DRO) over domains (Nemirovski et al., 2009, Oren et al., 2019, Sagawa et al., 2020), which we adapt to output domain weights instead of a robust model (Step 2). We then use the tuned domain weights to train a large model (Step 3).



Xie, S.M., Pham, H., Dong, X., Du, N., Liu, H., Lu, Y., Liang, P., Le, Q.V., Ma, T. and Yu, A.W., 2023. DoReMi: Optimizing Data Mixtures Speeds Up Language Model Pretraining. arXiv preprint arXiv:2305.10429.



# Data Order

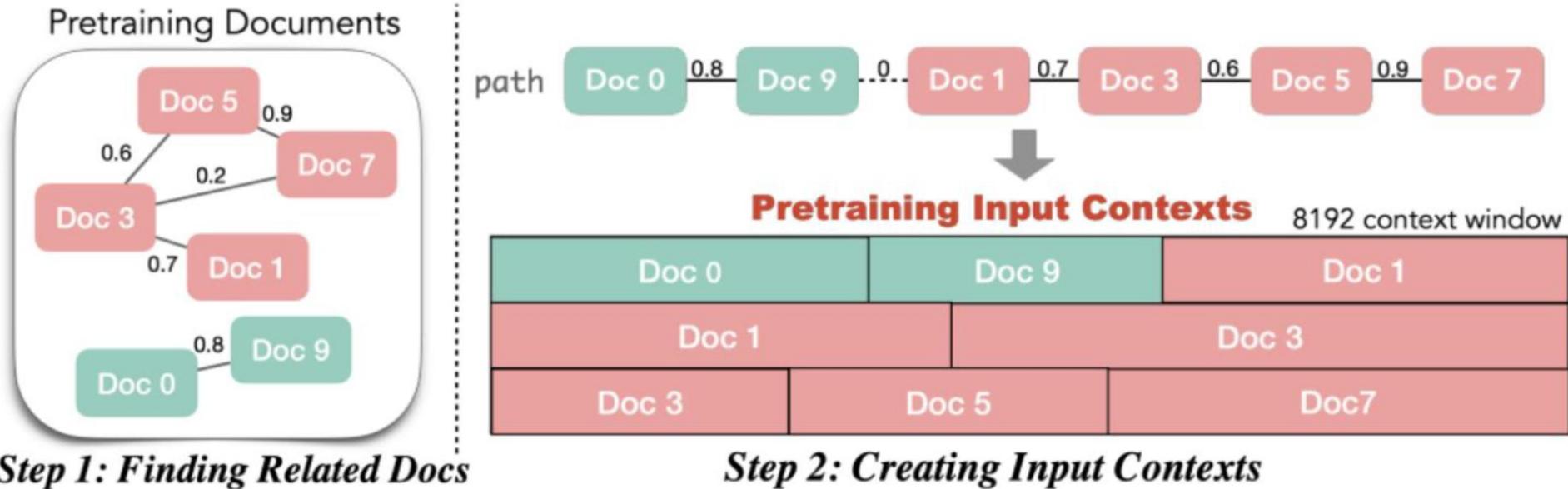
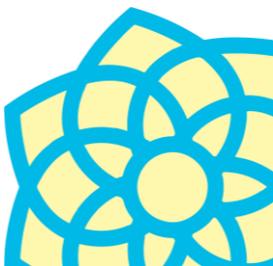


Figure 2: **Illustration of IN-CONTEXT PRETRAINING.** IN-CONTEXT PRETRAINING first finds related documents at scale to create a document graph (§2.1) and then builds pretraining input contexts by traversing the document graph (§2.2). Along the path, documents are concatenated into a sequence and subsequently divided to form fixed-sized input contexts (e.g., 8192 token length).

Shi, W., Min, S., Lomeli, M., Zhou, C., Li, M., Lin, V., Smith, N.A., Zettlemoyer, L., Yih, S. and Lewis, M., 2023. In-Context Pretraining: Language Modeling Beyond Document Boundaries. arXiv preprint arXiv:2310.10638.



# Data Masking

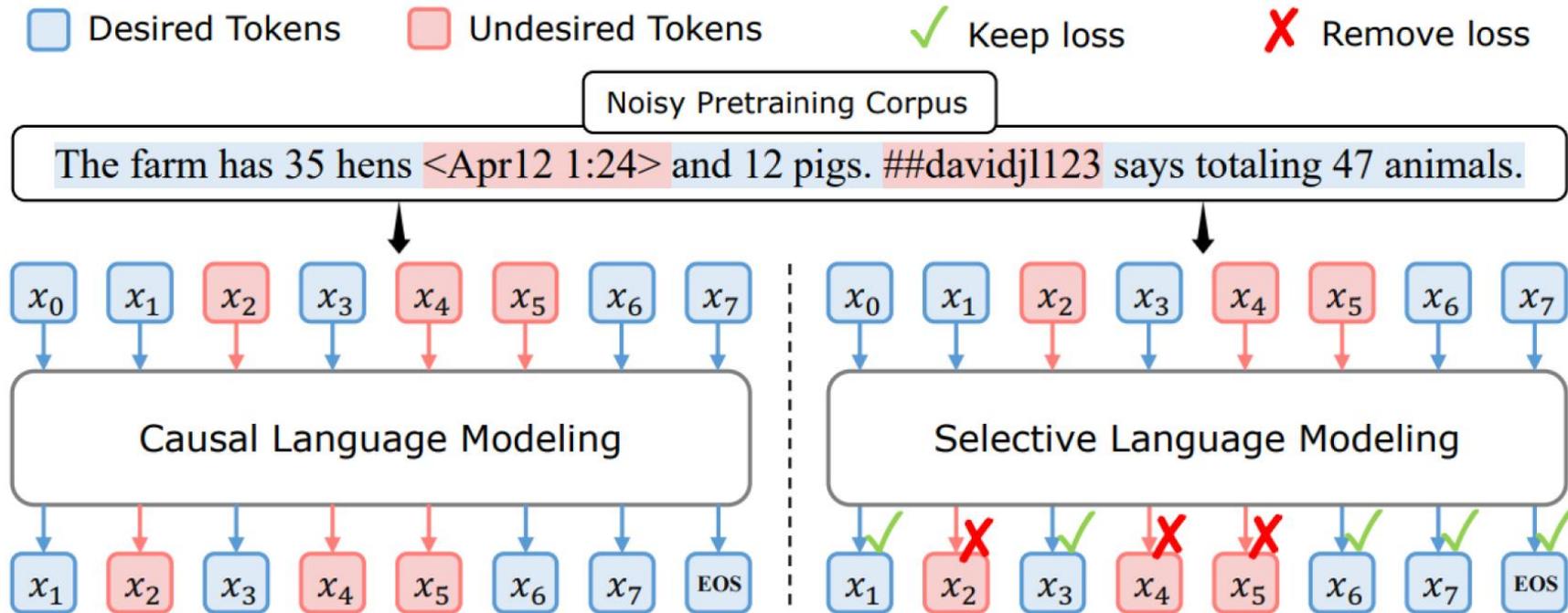
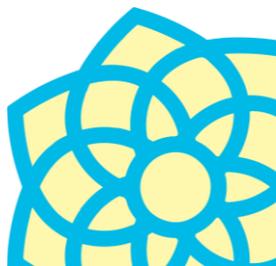
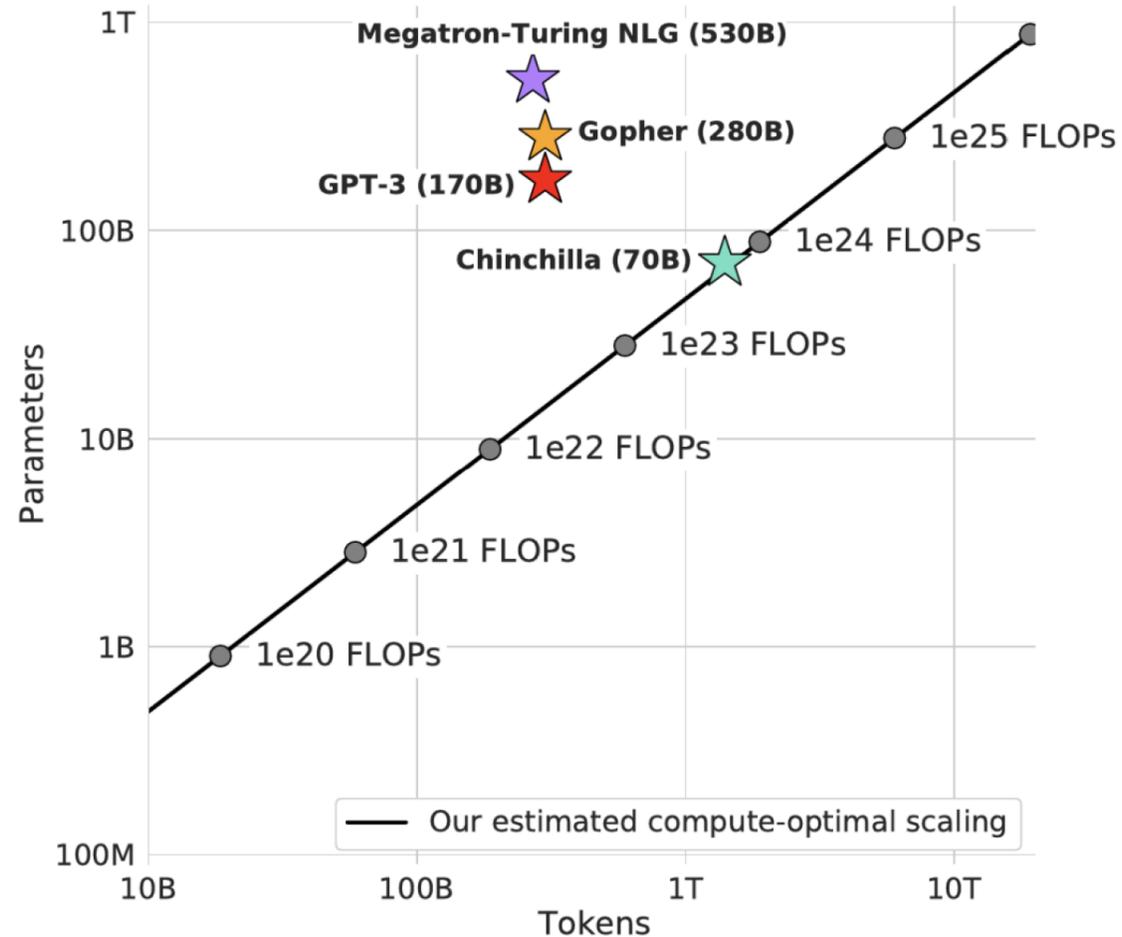


Figure 2: **Upper:** Even an extensively filtered pretraining corpus contains token-level noise. **Left:** Previous Causal Language Modeling (CLM) trains on all tokens. **Right:** Our proposed Selective Language Modeling (SLM) selectively applies loss on those useful and clean tokens.

RHO-1: Not All Tokens Are What You Need. <https://arxiv.org/pdf/2404.07965>



# Data Scale Matters

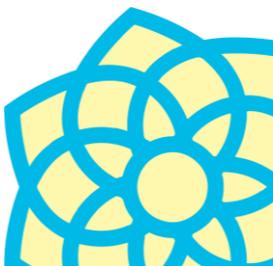


Recent models and its training tokens:

LlaMA-1: 1-1.4 T tokens

LlaMA-2: 2T tokens

Mistral-7B: much more...



# Pre-Training Data Quality Reduces Compute Needs

Recent work finds smaller amounts of higher quality data removes the need for a larger model.

There is increasing evidence that efforts to better curate training corpus, including **deduping, pruning data and investing in synthetic data** can compensate for the need for larger networks and/or improve training dynamics.

	% train examples with dup in train		% valid with dup in train
C4	3.04%	1.59%	4.60%
RealNews	13.63%	1.25%	14.35%
LM1B	4.86%	0.07%	4.92%
Wiki40B	0.39%	0.26%	0.72%

Table 2: The fraction of examples identified by NEARDUP as near-duplicates.

[Lee et al. 2022](#)

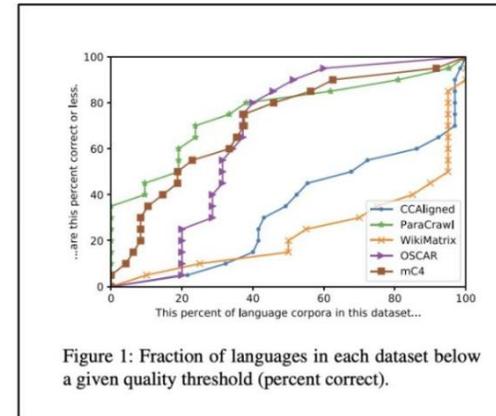


Figure 1: Fraction of languages in each dataset below a given quality threshold (percent correct).

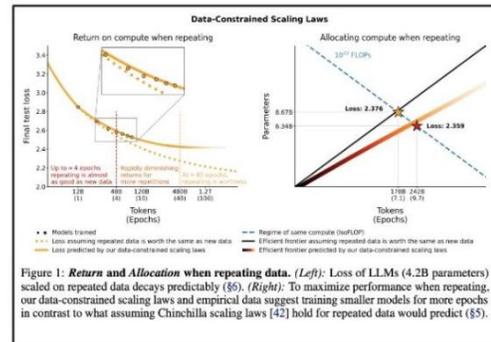
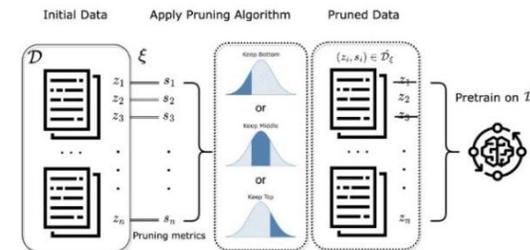


Figure 1: *Return and Allocation when repeating data.* (Left): Loss of LLMs (4.2B parameters) scaled on repeated data decays predictably (§6). (Right): To maximize performance when repeating, our data-constrained scaling laws and empirical data suggest training smaller models for more epochs in contrast to what assuming Chinchilla scaling laws [42] hold for repeated data would predict (§5).

[Muennighoff et al. 2023](#)

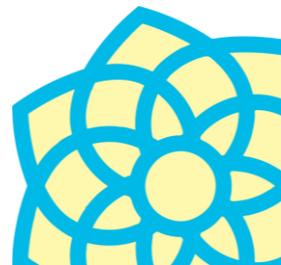
[Kreutzer et al. 2022](#)



[Marion et al. 2023](#)

↪ Cohere For AI

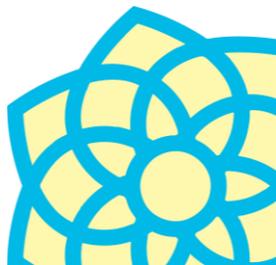
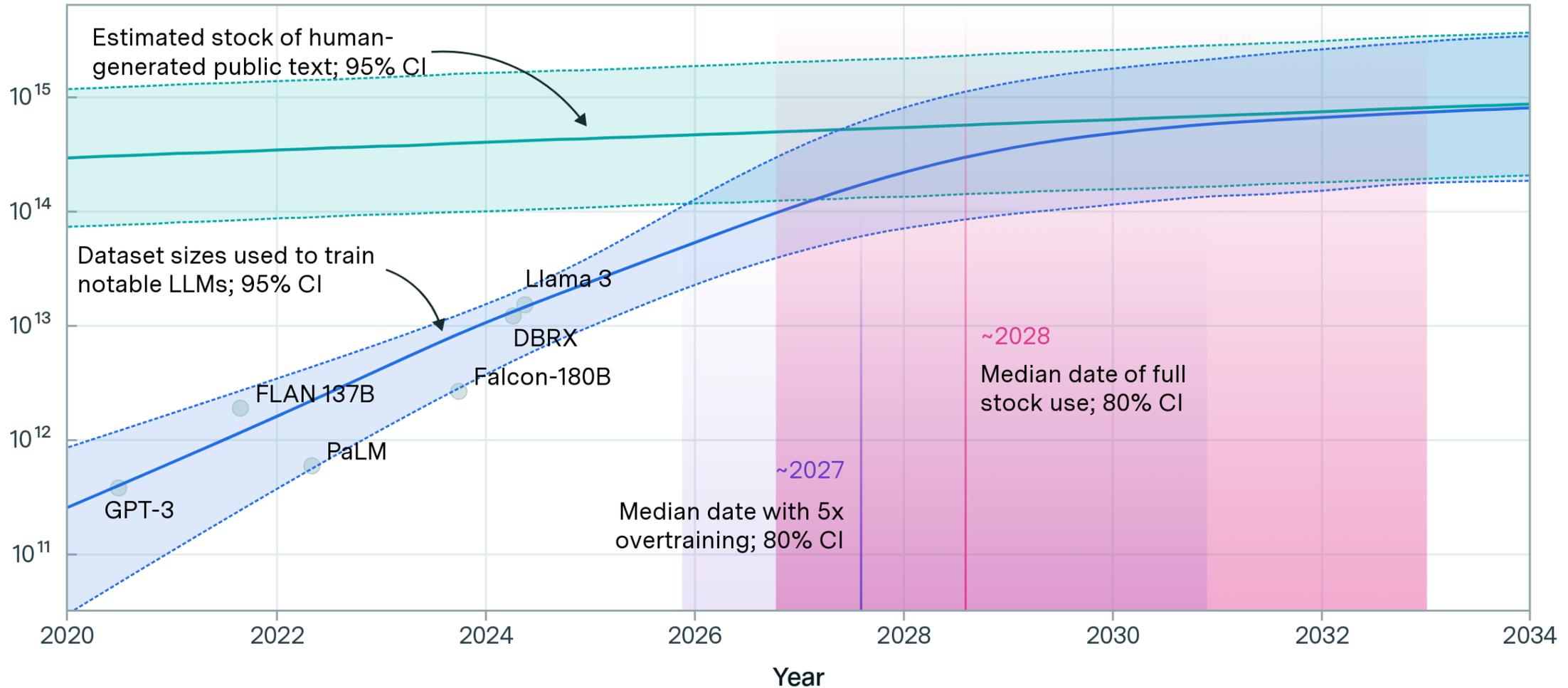
S. Hooker. [On the Limitations of Compute Thresholds as a Governance Strategy](#). 2024.



# Projections of the stock of public text and data usage



Effective stock (number of tokens)



- **LAIM LE3 VT2025:**  
Tokenization  
Data Processing Pipeline

[www.ida.liu.se/~frehe08/llm](http://www.ida.liu.se/~frehe08/llm)