

- * When working w/ multiple test cases, we add $\langle \text{end of text} \rangle$ tokens between these tests
- * These $\langle \text{end of text} \rangle$ tokens act as markers, signaling the start of end of a particular segment
- * This leads to more effective processing and understanding by the LLM

Byte pair encoding

- * ~~But~~ words themselves are broken into subwords, and these are the tokens
- (Ch)(as)(ed), one possible subword breakdown used to train GPT-2, GPT-3, and OG Model for Chat GPT

Tokenization Algorithms

① word based

② sub-word based

③ character based

My hobby is playing cricket

↓
[My, 'hobby', 'is', 'playing', 'cricket']

o Problem: what to do w/ out of vocab (OOV) words, diff meaning of similar words (how to not)

My hobby is playing cricket
↓
['a', 'x', 'h', ...]

↓
very small vocabulary, even largest has fixed # of characters (English 256)

problem ←
o meaning associated w/ words is completely lost. Longer then individual words

Character based tokenization
 8 tokens

ex:
 "boy" should not be split
 "boys" shall be split into
 "boy" and "s"

② Sub-word based tokenization
 Rule 1: Do not split frequently used words into smaller subwords

Rule 2: Split rare words into smaller, meaningful subwords

① The subword splitting helps model learn that different words w/ some root word as "token" are similar in

② Helps model learn that like "tokens" and "tokenizing" are similar in

"tokenization" and
 "modernization" are made
 up of diff. root words
 but have same suffix
 "ization" and are used
 in syntactic situations.

Byte Pair Encoding (BPE) is a subword
 tokenization algorithm
 * BPE also: most common pair of consecutive bytes
 of data is replaced w/ a byte that does not
 occur in data.

ex: Original data = "acacacacacac"

① byte pair 'ac' occurs the most, we will replace w/ Z as Z does not
 occur in the data (ac is a pair)

② compressed data: ZabcdZabac

③ Next common byte pair is 'ab', we will replace this by Y

④ ZYdZYc
 only byte pair left. Appears only once, so we do not
 encode it.

wdwac is
 new word

⑤ How is BPE used for LLMs?

(a) BPE ensures that most common
 words in vocab are represented
 as a single token, while rare
 words are broken down into two
 or more subword tokens

(b) let us take this example:

* { "old": 8, "older": 3, "finest": 9, "onest": 43 }

* Preprocessing: Need to add end token "<eos>"
 at the end of each word.

{ "old <eos>": 8, "older <eos>": 3, "finest <eos>": 9, "onest <eos>": 43 }

* Now split words into
 characters and count their
 frequency

#	Token	freq
1	<eos>	23
2	a	14
3	c	10
4	d	5
5	e	3
6	f	2
7	n	9
8	s	13
9	t	13
10	w	4

MEXO

PAGE

Next step! BPPB algorithm is to look for most frequent pairing
 (e) \rightarrow Merge them and perform the same iteration again and again until we reach the token limit or iteration limit

Iteration 1: Start w/ second most common token 'e', most common to the pair starting w/ 'e' is "es".

Intable, subtract the 'e' count and the 's' count by 13, Add 'es' to token list and assign it 13 frequency. 'es' \rightarrow 13, 'e' \rightarrow 0, 's' \rightarrow 0

Iteration 2: Merge tokens "es" and "t" as they have appeared 13 times in our dataset

Intable, subtract the 'es' and 't' count by 13, Add 'est' to token list and assign it 13 frequency. 'est' \rightarrow 13, 'es' \rightarrow 0, 't' \rightarrow 0

Iteration 3

(8) we see that "est <w>" has appeared 13 times,

so Intable, subtract the 'est' and '<w>' by 13, Add 'est <w>' to token list and assign it 13 frequency. 'est' \rightarrow 0, '<w>' \rightarrow 0, 'est <w>' \rightarrow 13
 helps also understand difference between estimate and highest

Iteration 4 "o" and "l" has appeared 10 times

so Intable, subtract the 'o' and 'l' by 10, Add 'ol' to token list and assign it 10 frequency, 'o' \rightarrow 4, 'l' \rightarrow 4, 'ol' \rightarrow 10

Iteration 5 'ol' and 'd' has appeared 10 times

so Intable, subtract 'ol' and 'd' by 10...

"Does not make sense to merge 'f' 'i' to 'fi' because does not appear much, we merged to cover b/c 'ol' and 'est' were present in lots of words,

Final table:

List of 11 tokens will serve as vocabulary.

Stopping criteria can either be token count, or # of iterations.

#	Token	Freq
1	<w>	10
2	o	4
3	l	4
4	e	3
5	r	3
6	f	9
7	i	9
8	n	9
9	w	4
10	est <w>	13
11	old	10