

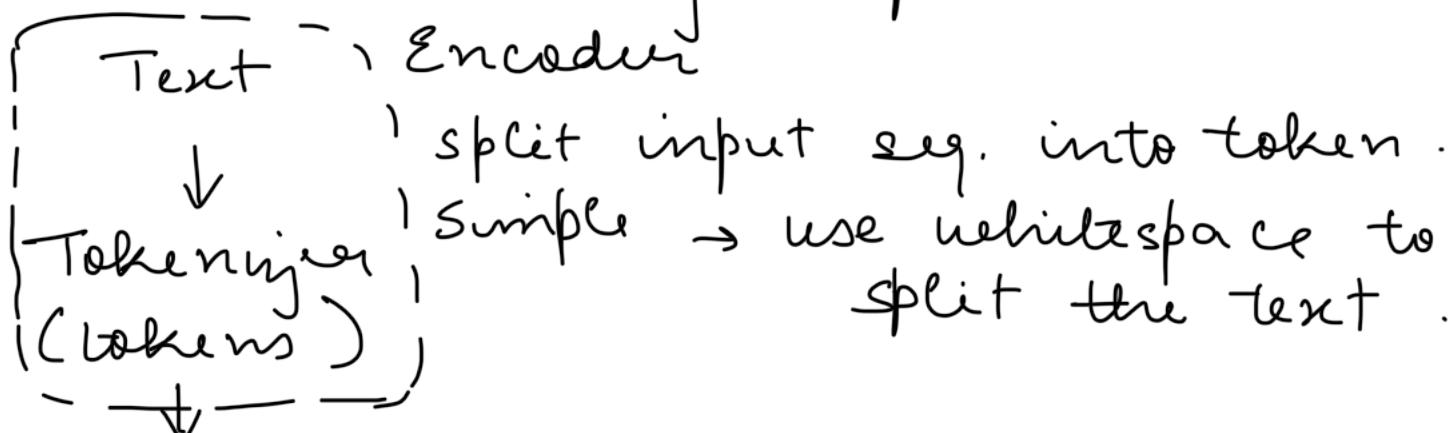
# DLP Week-2 Notes

## LI Tokenizer

- Transformer is a simple encoder-decoder model with attention mechanism at its core.

Input       $\xrightarrow{\text{Task}}$       Output  
 $\text{len}(\text{seq. of words}) \neq \text{len}(\text{seq. of words})$

- Abstract: Tokenizer Pipeline



Token ID      Each token gets a unique ID.

↓  
 Embeddings      Each ID gets a unique embedding vector

↓  
 LLM      Model predicts token IDs.

↓  
 ID to token  
 ↓  
 Token to words  
 Decoder

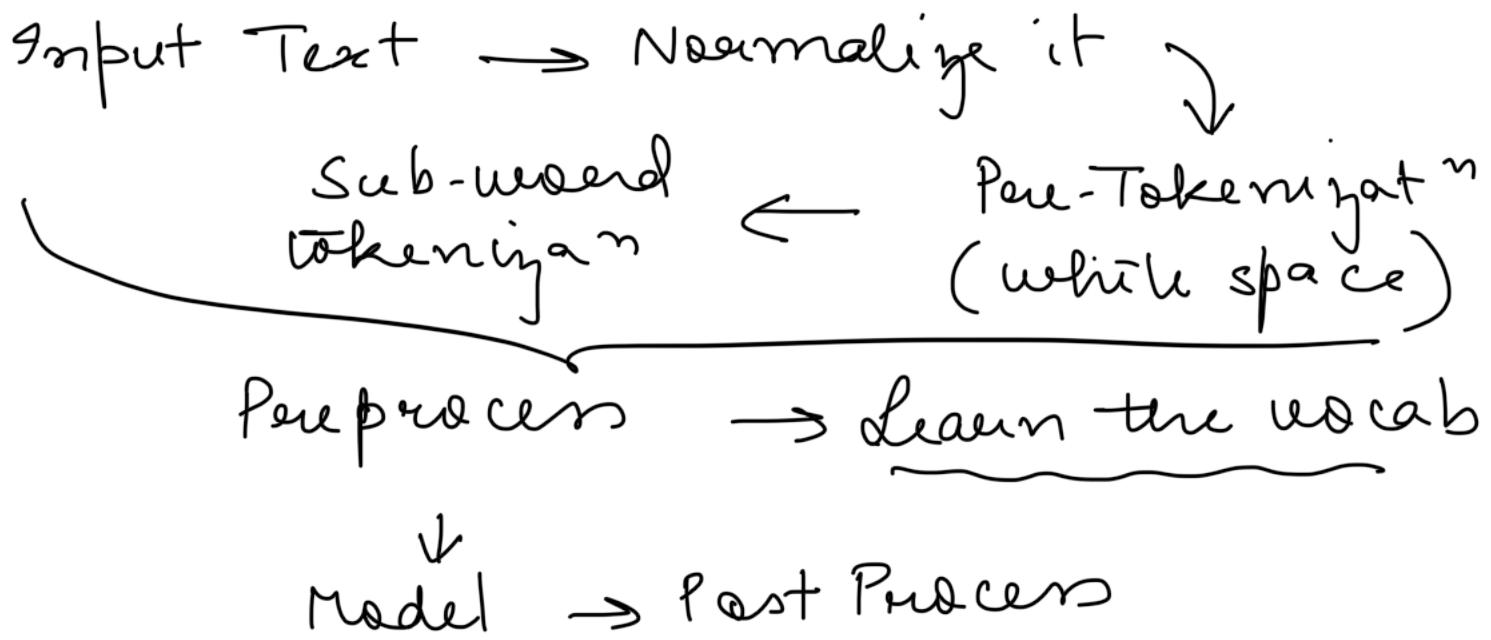
None decoder takes IDs and forms token. In the post-processing step, the tokens are finally converted to words.

- size of vocab = size of embeddings.

- how to build a vocab.
- Pre-tokenization - splitting text into words
- some sp. tokens → <go>, <cls>, <sep>, <pad>, <unk>  
are also used in BERT training

- L5 HF Tokenizer (Half) NOTE → Few things are needed before going L2, L3, ALY
- if whitespace a good delimiter? Then enjoyed & enjoy are diff.
  - but, we also can't have each characters as tokens. Too large vocab
  - few challenges -
    1. what is optimal size of vocab?
    2. If vocab size is 50k, how to deal with out of vocab words.
    3. How to handle small spelling mistake words?
    4. Open vocab problem. New words.
  - wishlist
    1. moderate size v.
    2. efficiently handle unknown words
    3. be lang. agnostic v size
  - charac. level — 1 - 112  
subword level — 30 - 50k → types word level — > 250k

## L2 Byte Pair Encoding



- Algo →

1. start with a dict., containing words & count
2. append  $\langle w \rangle$  at each word
3. set no. of merges (HPT)
4. create a charac. freq-table (base V)
5. get the freq. count for a pair of characters.
6. merge pairs with max occurrence.

- Obj. → find most frequently occurring byte-pair

word count → gives the vocab size

initial vocab size → charac. based

- final vocab has initial vocab + all merges. The more words broken down into 2 or more subwords.

Fertility (std 2-3) → knowing → 2

word tokenize → tokens

- BPE for non-segmented lang. (Japanese)

### 3 Word-Piece Tokenizer

- BPE → pair of tokens with highest freq.

$$\text{Score} = \frac{\text{count}(\alpha, \beta)}{\text{count}(\alpha) \text{count}(\beta)} \quad ('i', 'n') \over (\cdot), (n)$$

merging occurs on the basis of score of each byte pair  
reverse?

large vocab  $\xrightarrow{N}$  30 - 50 K  
clustering sub-word  
lit. Tokenizer

### 4 Sentence Piece Tokenizer

- any word can have numerous subwords
  - ↪ possibility -  $2^k$  combina<sup>n</sup>s
- follow the merge rule, while merging
- BPE is greedy & deterministic

$$x^k \in \{x_1, x_2, \dots, x_k\}$$

maximizing the likelihood of word  $x$ .

$$P(x) = \prod_{i=1}^n P(x_i)$$

$$\sum_{x \in V} p(x) = 1$$

$$\begin{aligned}
 & P(x_1 = k, \text{now}, \text{ing}) \\
 &= P(k) p(\text{now}) p(\text{ing}) \\
 &= \frac{3}{16} \times \frac{3}{16} \times \frac{7}{16} = \frac{63}{4096} \\
 &\quad x_2 \dots \\
 &\quad x_3 \dots
 \end{aligned}$$

$$x^* = \operatorname{argmax} P(x)$$

$$L = \sum_{S \in I} \log \left( \sum_{x \in S(x^*)} P(x) \right)$$

1. Construct a large seed vocab using BPE
2. E step  $\rightarrow$  estimate the probabs. for every token using freq. tokens.
3. M step  $\rightarrow$  use Viterbi algo. to segment I return optimal segments that maximises the log likelihood.
4. compute likelihood for new segments
5. shrink by removing bottom  $x\%$  of subseeds.  
Eq. whereby

## L5 HF Tokenizer (contd.)

Noemalija " "

PereTokenija " "

Tokenizer " (With V)

Post Processing

from-tokenizers  
insert tokens

normalizer → lowercase, strip accents  
pre-tokenizer → whitespace, BERT, Regex  
Algo. → BPE, subword pieces

Post processor → insert model specific tokens

customization at each step

- after setting the normalizer & pre-tokenizer just call the -train methods.

## L6 Building & Training a Tokenizer

- book corpus lib. dataset (74M, 5GB)
- Note - len(sentences) does not affect the Tokenizer training
- BPE ([UNK])
  - few unknown tokens.
- create a trainer with vocab size & special tokens.
  - continuing\_subword-prefix = "##"
- save it somewhere safe.
  - the early tokens are small characters but last one are like big words.
- total merges  $\leq$  vocab-size
  - slider to view animation.

## L7 Encoder & Decoder

- pass the single sample to 'encode' method of Tokenizer class.
- we can use a visualizer to see how the split works
- atten<sup>n</sup> mask  $\rightarrow$  1 (pay attention to every word)  
↳ BERT models (type id  $\rightarrow$  0, 1)
- when you have a batch of samples, then use [PAD] token to make sure samples do not have irregularity.
- in general, it's fine to increase the length of a seq. to exceed the model's content length.
- hf  $\rightarrow$  h ##f, 😊  $\rightarrow$  [UNK]
- always try to save the tokenizer.
- BERT tokenizer
  - <CLS> at start
  - <SEP> btw. 2 sentences
  - <MASK>spe.
- fast processing step  $\rightarrow$  Template Processing
  - single =
  - pair =
  - special\_tokens = [( ), ( )]
- New, decode,  
remove spe. tokens, merge the #ff thing

## L8 Wrap With Pre-Trained Tokenizer

- PreTrainedTokenizer class.
  - file, unk token, pad token
- ↓
  - you can use this to create a model input.
- Padding is not by default. You have to do it manually.