

Tokenization in Transformer :From Text to Tokens

Author: Shih-Shinh Huang

Date: July 11, 2024



Outline

- Introduction
 - About Transformer
 - About Tokenization
- Tokenization Strategy
 - Character Tokenization
 - Word Tokenization
 - Subword Tokenization





Introduction

- About Transformer
 - The Transformer neural network is a novel and widely used architecture in NLP.
 - Text classification
 - Word labeling
 - Question-Answering

A. Vaswani, "Attention Is All You Need," *NIPS*, 2017.**NLP**: natural language processing





Introduction

- About Transformer
 - Most data in NLP is human-readable <u>text strings</u>

"Tokenization is a core task of NLP!"

• The transformer models can only receive a list of **integer numbers** as input.

$$[5, 14, \dots, 4, 1] \longrightarrow \boxed{\text{Transformer}_{Model}}$$

"Tokenication is a to use transformer, models, for N Transformer core task of NLP!" Model

convert

Conv



Introduction

- About Tokenization
 - Tokenization is a step in NLP to convert <u>text</u> <u>strings</u> to a list of <u>integer numbers</u>
 - break down text strings into meaningful units (tokens)
 - encode tokens into transformer-readable integers (<u>token IDs</u>).

tokens





- Character Tokenization
 - split text strings into a list of <u>characters</u> and consider each character as a token.



token list





- Character Tokenization
 - transform the tokens into a list of integers

character-based vocabulary

$a \rightarrow 1$	$A \rightarrow 27$	$0 \rightarrow 53$	$\rightarrow 63$
$b \rightarrow 2$	$B \rightarrow 28$	$1 \rightarrow 54$	$! \rightarrow 64$
$c \rightarrow 3$	$C \rightarrow 29$	$2 \rightarrow 55$. ightarrow 65
$z \rightarrow 26$	$Z \rightarrow 52$	$9 \rightarrow 62$	

map the character in our vocabulary to a unique integer vocabulary size of English language is small (≈ 256)



- Character Tokenization
 - transform the tokens into a list of integers







- Character Tokenization: Discussion
 - Advantages:
 - The vocabulary is easily identified and its size is small.
 - There is no **<u>out of vocabulary</u>**(OOV) tokens
 - Drawbacks: The token sequences is very long and with less meaning and linguistic structures.

'Tokenizationa ts i a toona taska of NIP?'' t a s k of NLP!

sequence length = 35





- Word Tokenization
 - split text strings into a list of **words** by using whitespace and consider each word as a token.

tokenization	is	a	core	task	of	NLP!
--------------	----	---	------	------	----	------

token list





- Word Tokenization
 - transform the tokens into a list of integers

word-based vocabulary



map the word in our vocabulary to a unique integer





- Word Tokenization
 - transform the tokens into a list of integers



encode each token with an integer by mapping





- Word Tokenization: Discussion
 - Advantages: a token as a word in a sentence has a lot of contextual and semantic information
 - Drawbacks:
 - the vocabulary size is extremely large
 - some tokens will be UNKNOWN (OOV tokens) that lose information

UNKNOWN	s a	core	task	of	UNKNOWN
---------	-----	------	------	----	---------





- Subword Tokenization
 - The main idea is to address the issues faced by character and word tokenizations.
 - split the rarely used words into smaller subwords
 - \Rightarrow reduce vocabulary size and alleviate OOV problem
 - keep the frequently used words as unique entities
 - \Rightarrow preserve contextual and semantic information





Subword Tokenization





- Subword Tokenization:
 - WordPiece: used by BERT and DistilBERT

M. Schuster, et. al. "Japanese and Korean Voice Search," ICASSP, 2012.

• Unigram: used by XLNet and ALBERT

T. Kudo, "Subword Regularization: Improving Neural Network Translation Models with Multiple Subword Candidates," *ACL*, 2018.

• Bye-Pair Encoding (BPE): used by GPT-2

R. Sennrich, *et. al.* "Neural Machine Translation of Rare Words with Subword Units," *ACL*, 2016.



