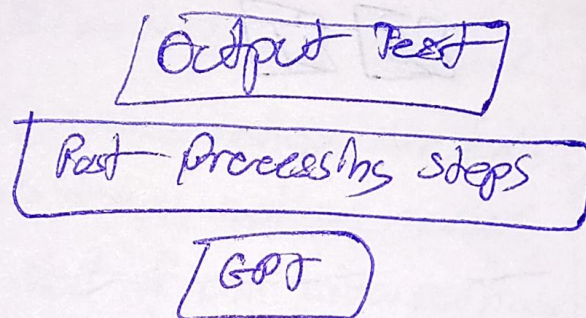


Tokenization depends on the Jupyter notebook
and bleeds for building LLM

① How do we prepare input text for training LLMs?

1. Splitting text into individual word and subword tokens
2. Convert tokens into token IDs
3. Encode token IDs into vector representation

② Let us look at step 1: Tokenizing text



DDD DDD DDD DDD Token Embedding (step 3)
[4013] [201] [302] [1134] Token IDs (step 2)
[This] [is] [an] [example] Tokenized text (step 1)

[This is an example] Input text

- (a) Dataset used: "The Verdict" by Edith Wharton
- (b) Download and load in Python

Complete Training Dataset

→ Tokenized text

→ Vocabulary
(word to index)

"The guide brown fox jumps over the lazy dog"

(The/guide
brown/fox)

brown - 0
dog - 1
fox - 2

Each unique token is mapped to a unique integer called token ID

jumps - 3

over - 4

and - 5

guide - 6

the - 7

Token IDs

unique tokens

(e) Implement the Tokenizer class in Python

Simple Test!

Tokenized text

Token IDs

The fox chased the dog

(The/fox/chased)

23 2 0

the dog

(The/dog)

23 1

Encode

tokenizer

ENCODER

encode(text)

DECODE

Token IDs

23 2 0

23 1

Tokenized text

(The/fox/chased)

(The/dog)

tokenizer decode
(ids)

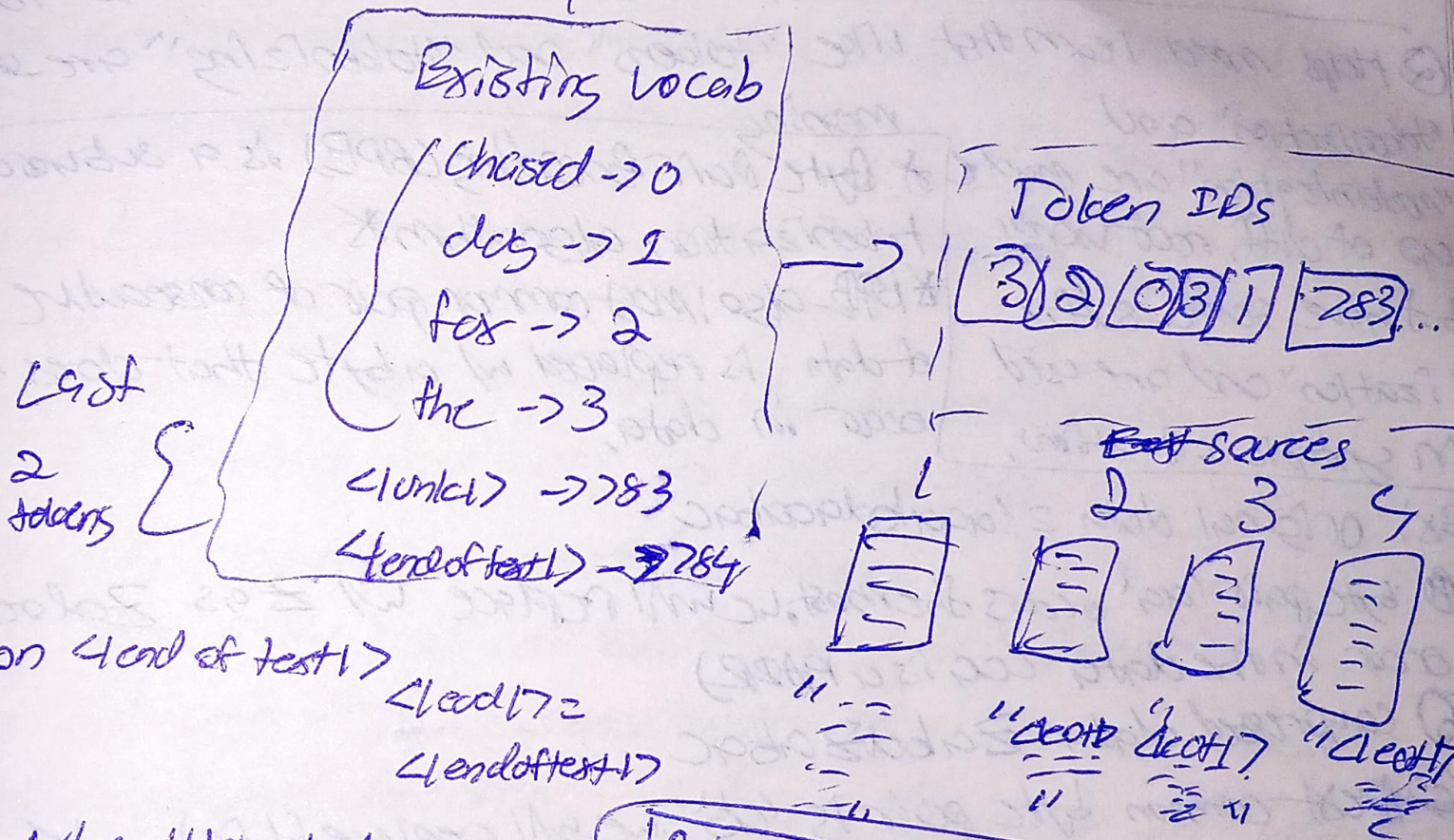
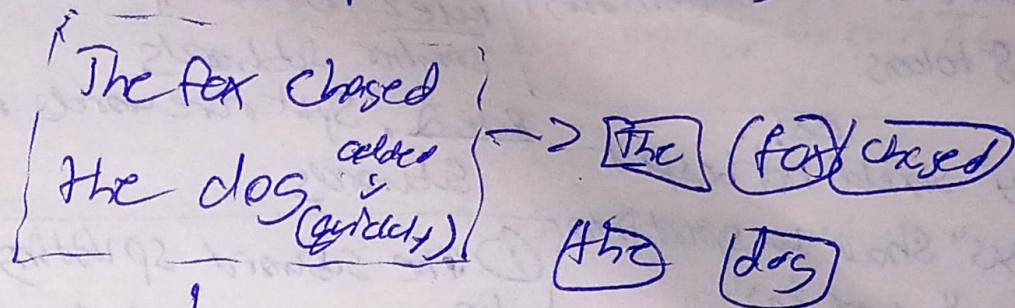
Simple Test

The fox chased the dog

Special Context Tokens

used to handle unknown words
o Modify the vocab and tokenizer we implemented in previous section (v2)

(8) Adding special test tokens



more details on <end of text>

* when working w/ multiple test sources, we add <end of text> tokens between these texts
* these <end of text> tokens act as markers, signaling the start or end of a particular segment

* This leads to more effective processing and understanding by the LLM

Byte pair encoding

* ~~Big~~ 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
GPT-4 Model for Chat GPT