# Lecture - 7 : Tokenization

## Data preparation and sampling

→ we will look at "Data preparation and sampling" in this lecture.

1. How do you prepare input text for training LLMS?

Step-1: (a.) Splitting text into individual word and subwords & token

Step-2: (b.) convert token into token IDs.

Step-3: (c.) Encode token IDs into vector representation.

⌄

: Output text :

| Post processing steps |

| GPT |

☐☐☐ ☐☐☐ ☐☐☐ — Token embedding
(step-2)

| 4013 | | 201 | | 302 | | 1134 | — Token IDs
(step-3)

| This | | is | | an | | example | — Tokenized text
(step-1)

| This is an example | — Input text

## 2. Let's us look at Step-1:
## Tokenizing Text

(a.) Dataset used: "The verdict" by Edith Wharton.

(b.) Download and load in python.

(c.) Tokenize the short story.
↳ use python's regular expression liberary.

(d.) convert Token into Token IDs

---

compute training ⟶ Tokenized
   dataset              text

"The quick brown
fox jumps the lazy
dog"

| The | quick |

| brown | fox | --

Each unique token
is mapped to an
unique integer
Called token ID.

usually sorted ↓
alphabetically
↑
vocabulary
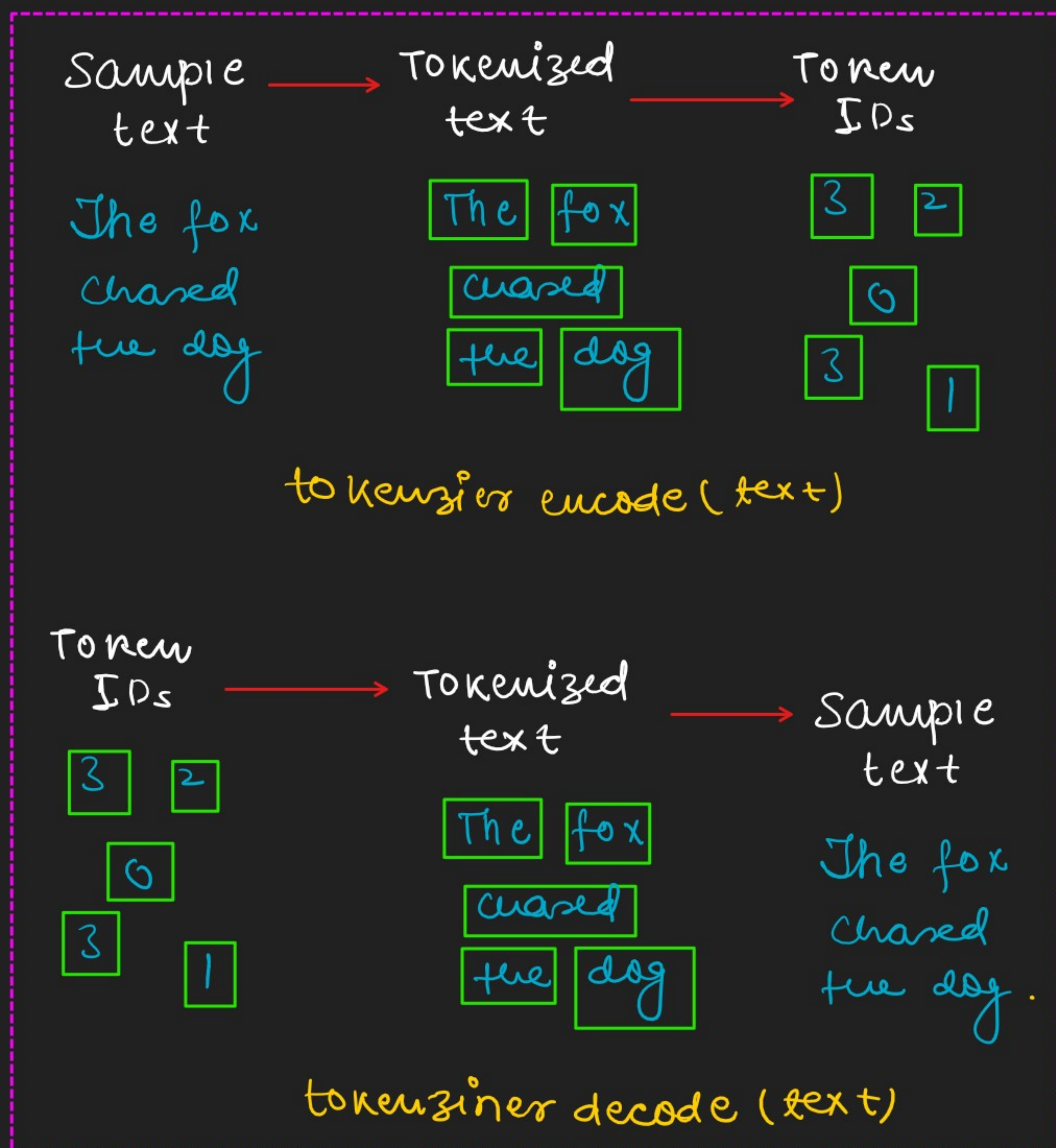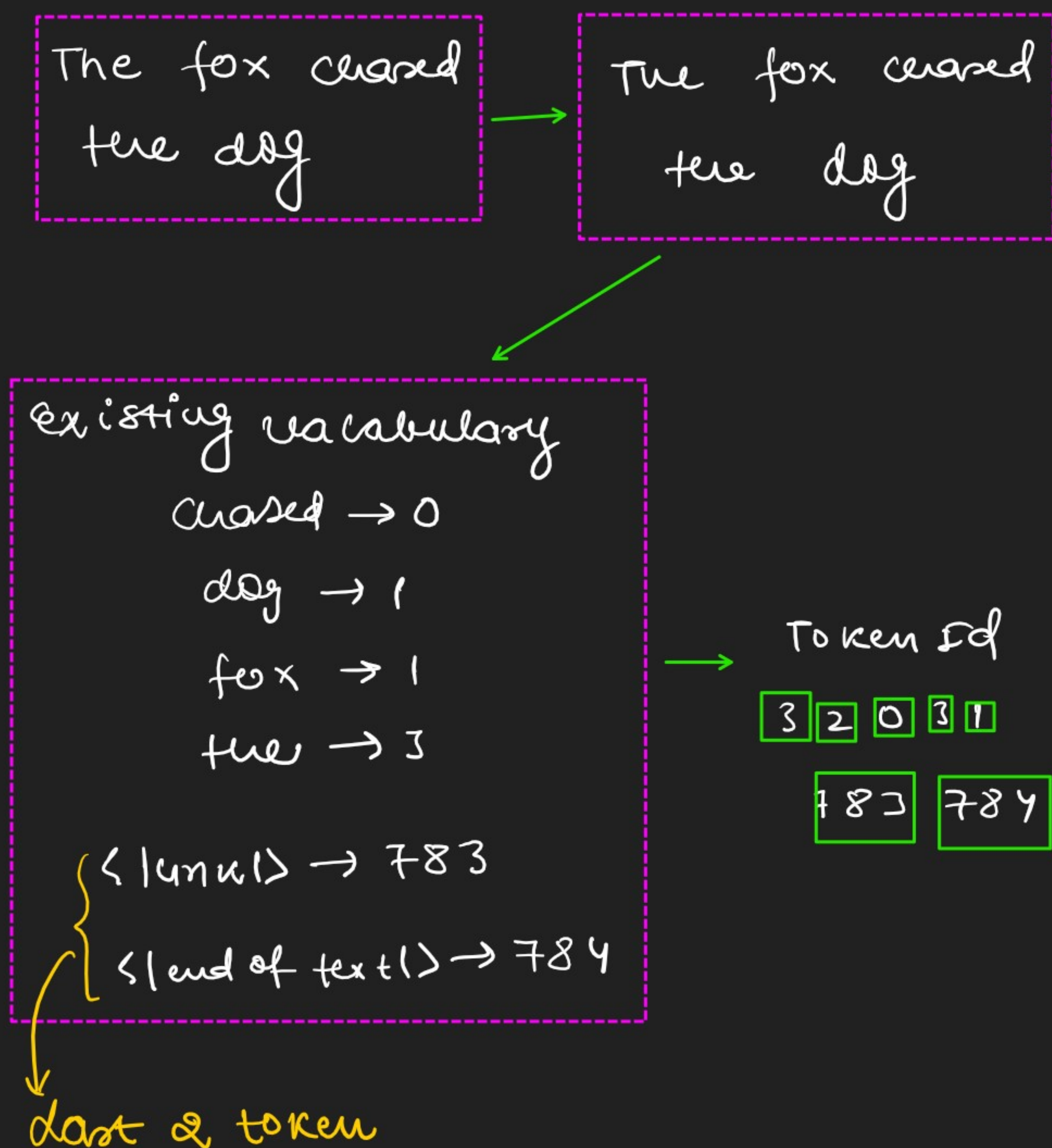
| brown | → 0
| dog | → 1
| fox | → 2
| jump | → 3
| lazy | → 4

↳ token
   Id

↓
unique
tokens

# (e.) Implement the token class in python.

Sample text → Tokenized text → Token IDs

The fox
chased
the dog

| The | fox |
| Chased |
| the | dog |

| 3 | 2 |
| 0 |
| 3 |
| 1 |

tokenizer encode (text)

Token IDs → Tokenized text → Sample text

| 3 | 2 |
| 0 |
| 3 | 1 |

| The | fox |
| Chased |
| the | dog |

The fox
chased
the dog.

tokenizer decode (text)

# (f.) Adding special text token

The fox chased the dog → The fox chased the dog

existing vocabulary
  chased → 0
  dog → 1
  fox → 1
  the → 3

  <|unk|> → 783
  <|end of text|> → 784

Token Id

| 3 | 2 | 0 | 3 | 1 |
| 783 | 784 |

last 2 token

# more details on <|end of text|>

* when working with multiple text sources, we add <|end of text|> token b/w these text.

* these <|end of text|> token acts as markers, signgling the start of end a particular segment.