# Create Input Target Pairs

o Last step before we create
vector embeddings is to creat input target pairs

* what do these input target pairs look like?

Best sample

D is target, before is input

iter 1 LLMs [learn]

2 LLMs learn [To]

3 LLMs learn to [predict] ← Target to predict

4 LLMs learn to predict [one]

5 ~

6 ~

7 ~

← LLMs cannot access
words past the target

(8 LLMs learn to predict one word at a [time]

* Given a test sample

 ↳ Extract input blocks as subsamples that serve as <u>input to LLM</u>

 ↳ LLMs prediction task during training is to predict the next
word that follows the input block

 ↳ During training, we mask at all words that are past the target

* Data loader text fetches input-output target pairs using a sliding
window approach

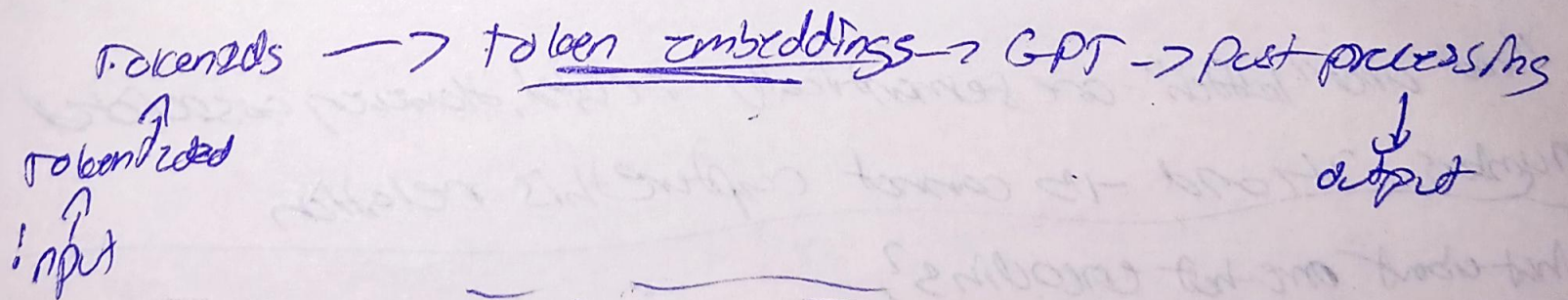test: "In the heart of the city stood the old library, a relic

from"

Tensor input: $x = [[$ "In", "the", "heart", "of" $]$,
    2
   3
  4
           2)

output
tensor   $y = [[$ "the", "heart", "of", "city" $]$
      N    2    3    4

\# To implement efficient data loader, we collect inputs in a tensor $x$ where each row represents one input context. The second tensor $y$ contains corresponding prediction targets (next word), which are created by shifting input by one position.

## Token Embeddings

Tokenizds $\longrightarrow$ token embeddings $\rightarrow$ GPT $\rightarrow$ Post processing

Tokenized

↑

input

dropout

---

① conceptual understanding of why token embeddings are needed

② Demo

③ How are token embeddings created for LLMs