

Positional Encoding

In embedding layer, some tokens DO gets mapped to the same vector representation, regardless of where the tokens DO'S positioned in the input sentence

* The cat sat on the mat

* On the mat the cat sat

$\begin{bmatrix} 5 \\ 1 \\ 5 \end{bmatrix}$

from embeddings layer,

Percent vector embedding approach has no idea of position of 'cat'

Some tokens DO'S result in the same embedding vectors

$\begin{bmatrix} -123 & 456 & 789 \\ 101 & 212 & 345 \\ -123 & 456 & 789 \end{bmatrix}$

③ It is helpful to inject additional position info to LLM

④ There are 2 types of positional embeddings

↓
absolute

For each position in input sequence, unique embeddings is added to the token's embedding to convey its exact location

The positional vectors have same dimension as the original token embeddings (for easy addition matrix purposes)

↓
relative

The emphasis is on relative position or distance between tokens. Model learns the relationships in terms of 'how far apart' rather than at which exact position

↓
Advantage: model can generalize better to sequence of varying lengths, even if it has not seen such lengths during training

⑤ Both types of positional encodings enable models to understand the order and relationship between tokens, ensuring more accurate and context aware predictions

⑥ The choice between the 2 depends on specific application and nature of data being processed

Absolute: Suitable when fixed order of tokens is crucial, such as sequence generation

Relative: Suitable for tasks like language modelling over long sequences, where the same phrase can appear in different parts of the sequence

⑦ OpenAI models (GPT) use absolute positional embeddings that are optimized during training process. This optimization is part of the model training itself.

⑧ Get input data from data loader

Inputs = $\begin{bmatrix} \dots \end{bmatrix}$ batch size = 8
 \rightarrow 8 test samples w/

context length = 4

4 tokens size

vector Dims,

Inputs = $\begin{bmatrix} 10 & 8 & 20 & 21 \\ 1 & 3 & 4 & 5 \\ \vdots & \vdots & \vdots & \vdots \\ 9 & 8 & 16 & 23 \end{bmatrix}$

Token 21's \rightarrow $\begin{bmatrix} \text{vector} \\ \text{size} \end{bmatrix}$
 50 256

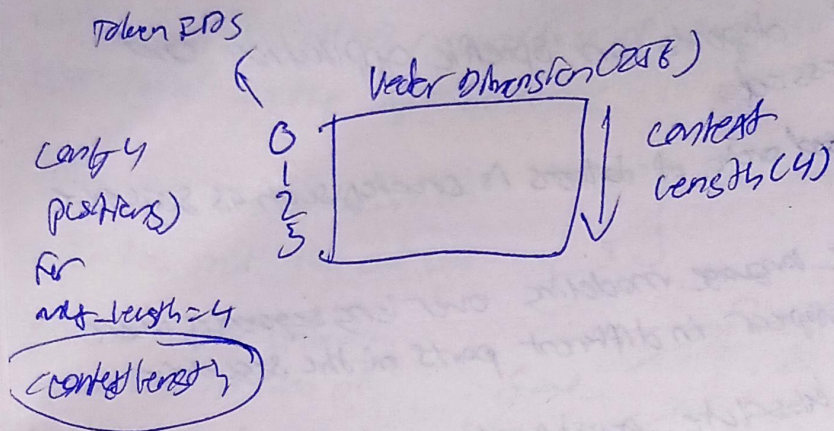
One embedding vector of size 256 length is generated for each token in input

$\begin{bmatrix} (256) & (256) & (256) & (256) \\ \vdots & \vdots & \vdots & \vdots \end{bmatrix}$

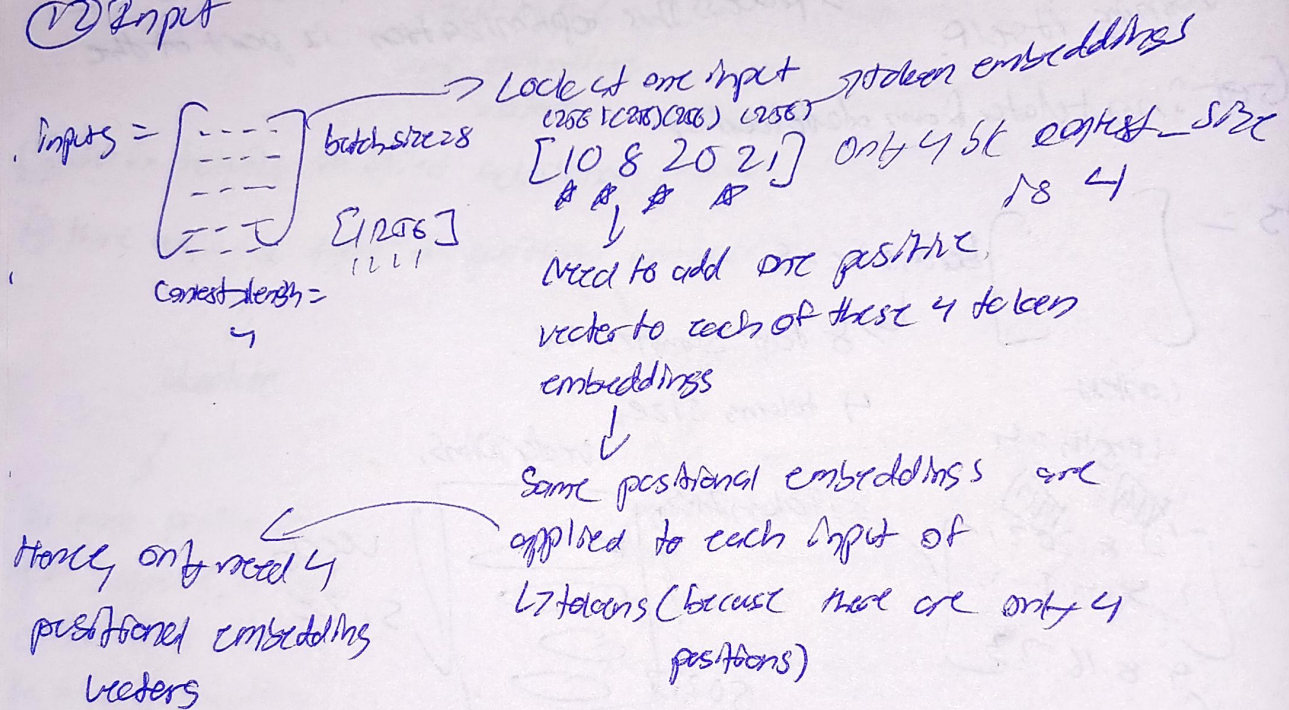
$8 \times 4 \times 256$

for each of these 8x4 values, we have a 256 dimension tensor

① Create another embedding after
encoding



② Input



only encode on FOUR positions, since context is 4.

③ Generate 4 positional embedding vectors from positional embedding matrix

⑬ Generate 4 positional embedding vectors from the positional embedding matrix

$[24, 256]$

⑭ Add the token embeddings + position embeddings

$$\begin{bmatrix} (256) & (256) & (256) & (256) \\ \vdots & & & \\ (256) & (256) & (256) & (256) \end{bmatrix} + \begin{bmatrix} (256) & (256) & (256) & (256) \end{bmatrix}$$

$8 \times 4 \times 256$ 4×256

(token embeddings)

Per position, each w/ 256 vector

(position embeddings)

= Input Embeddings
 $(8 \times 4 \times 256)$

broadcasts 4×256 to all 8 samples in batch to support the matrix addition