

# Generative AI with LMs & Andrew NG Course

what is a LM capable of? basically, they are "next word generators"  
so they can do any kind of question answering, summarizing & Entity extraction  
augmenting LMs with connected APIs.

- LMs are distinguished by their "parameter size" (E.g. 10M (parameters), 2B (parameters))

Larger LMs: more conceptual understanding → best for complicated a multi-purpose

Smaller LMs: Less " " but faster → focused tasks

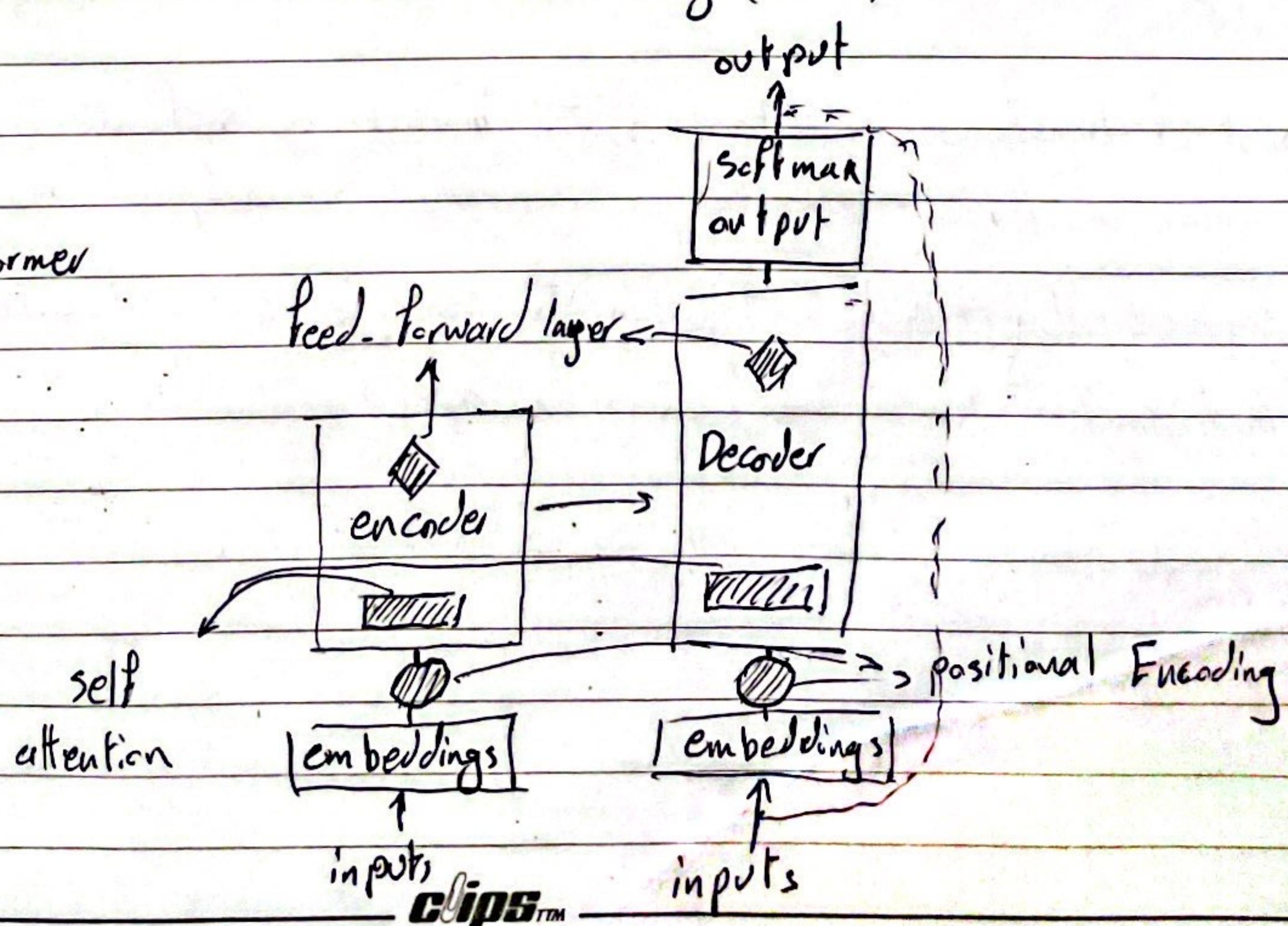
Before the Transformers (structure used in LMs and latest GenAI), we used RNNs  
for next word prediction.

what we lacked? { - memory: longer sentences couldn't be analyzed  
{ - multi-meaning words in sentence  
{ - two/multiple sided sentences } solved them  
} using Transformer

But how? the transformers use "self attention": a vector, containing relation of every  
word in a sentence (token)  
word (or even a whole text), called "embedding vector"

\* input size of an LM is based on its size of embeddings (vector)  
(max)

structure of a transformer



let's break down some of its parts:

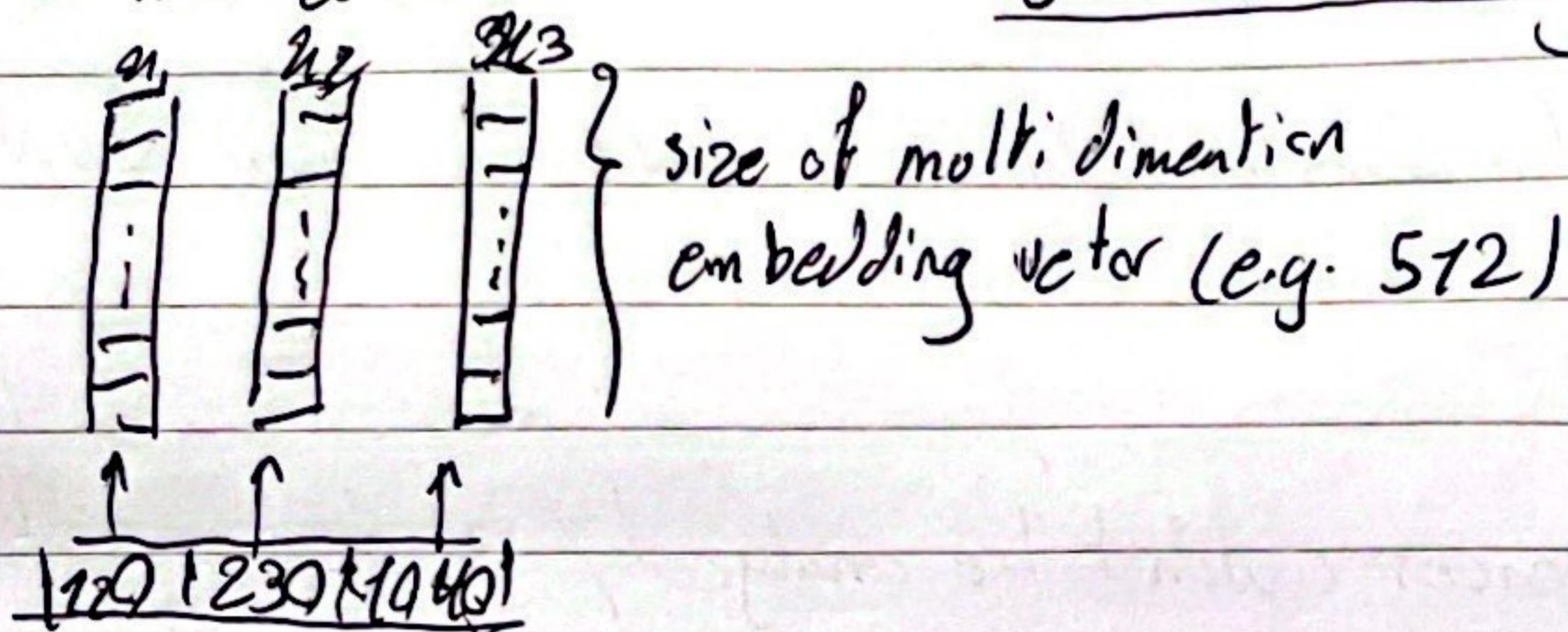
① every input gets tokenized, each word being represented as a unique ID, based on dictionary in use: the teacher thought the

Taken IDs ↗ "Tokenizer"  
[120 | 290 | 1040 | 931]

use the same tokenizer for both training and generation

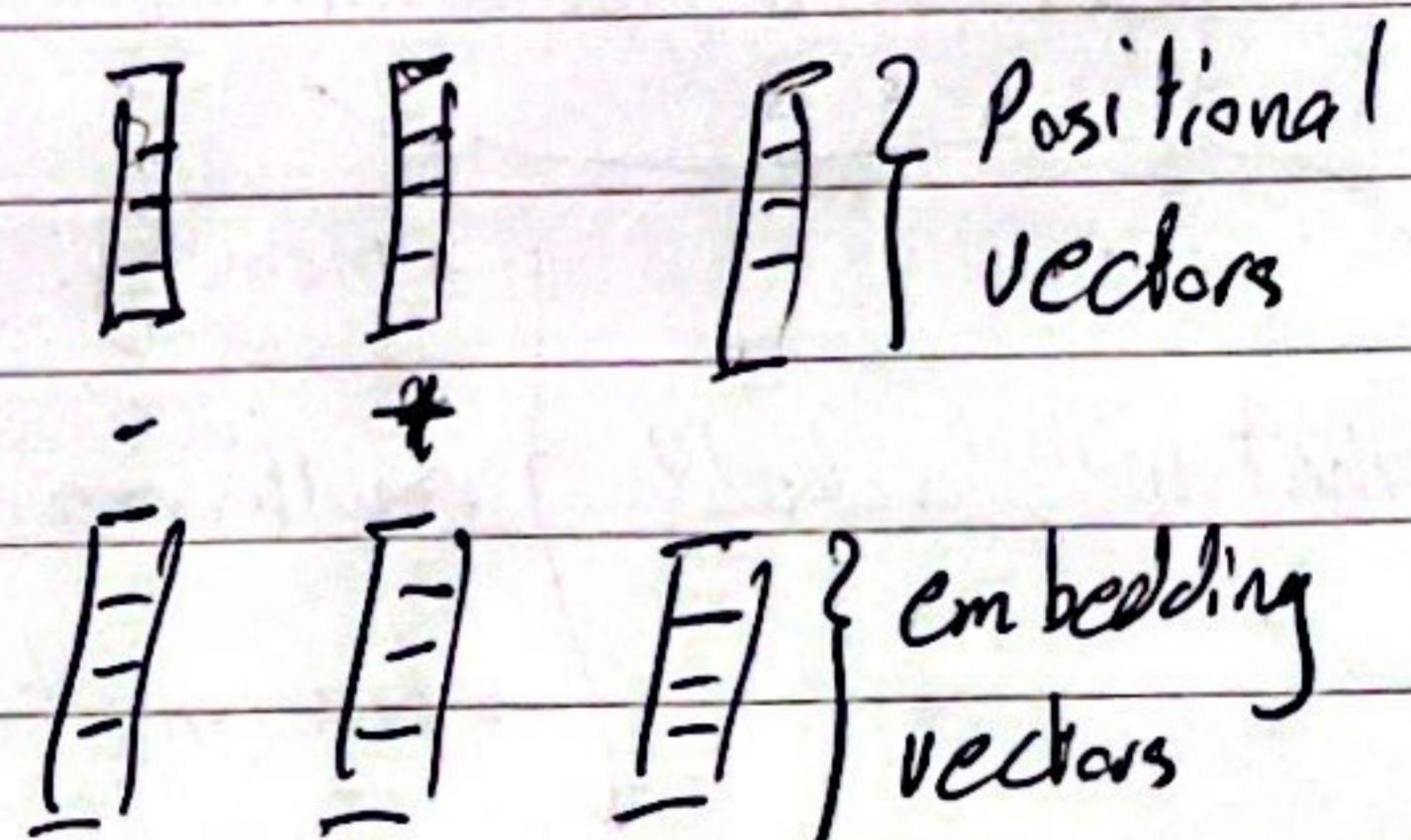
② Taken ID vector gets passed to a (learnable) embedding layer.

so that: each taken ID → high dimensional vector



represents the connection between words, etc.

③ Add positional encodings → to preserve the word order

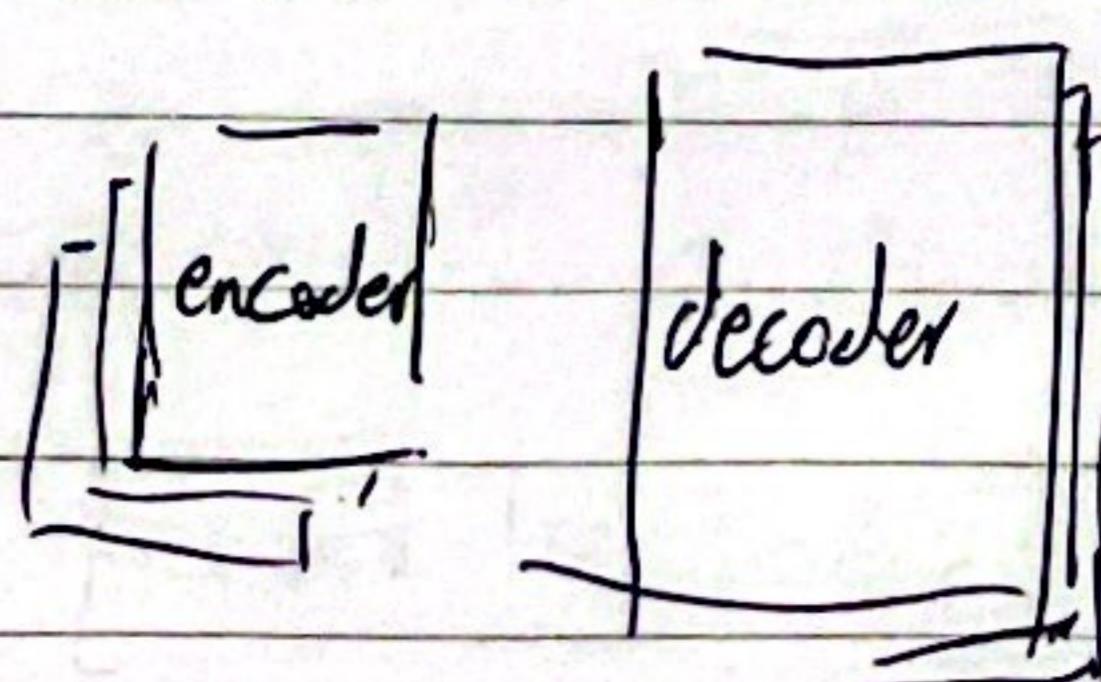


④ Feed the resulting vector (embed + positional) into a self attention layer.

→ which learns to better capture contextual relation between input words

The self attention weights learnt during training, helps the model to capture the relation between each word, with any other word in the input  
(and the importance of, hence, "attention")

In most of transformer architectures, the "multi-headed self-attention" is used, meaning that there are parallel self-attention layers, which are independently learning the importance of each encoding (embed + positional) in the sentence/text. (12 ~ 100 number of multi-heads)  
in fact, each attention head, learns an aspect aspect of language (people, name, rhyme, ...)



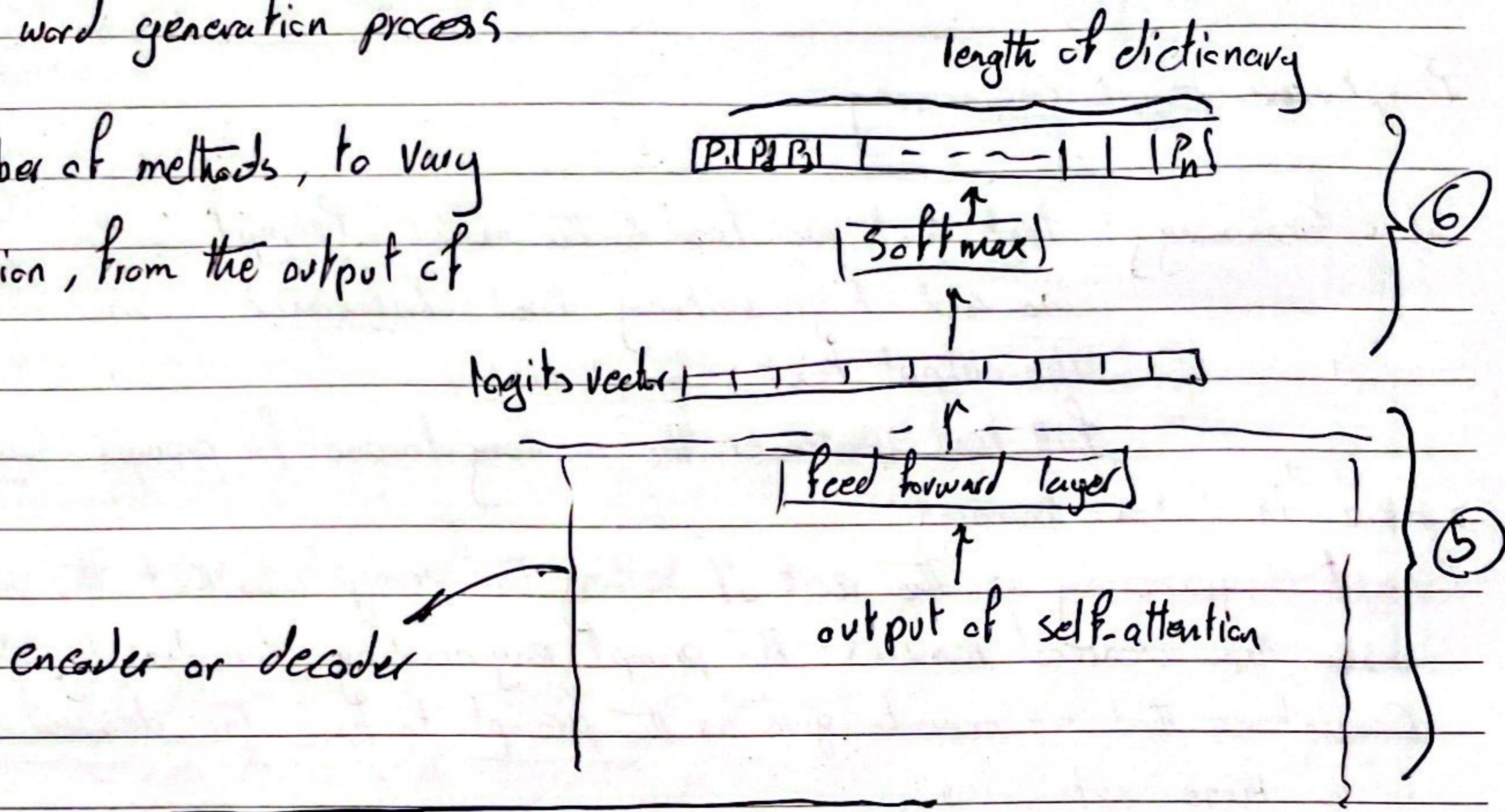
(multi-headed self-attention)

(3)

⑤ after the attention weights have applied to the data, it will be fed to a feed-forward network, and its output is a vector of logits, with each value proportional to the probability score, for each and every token for the tokenization dictionary

⑥ then we can feed it to a softmax layer, outputting a vector with length of dictionary, containing prob. score for each token.  
The highest score, in most cases, is the next word that should be used to start, continue the word generation process

A) there are a number of methods, to vary the final selection, from the output of softmax layer



### Generating text with transformers

So, how do transformers work after the training phase?

A) Path of encoder:

Input → Tokenizer → embedding + positional → self-attention (attention vectors) → f.f. layer → inserted in the decoder for influence on generation  
Parallel steps (because multi-head)

B) Path of decoder

Input → Tokenizer → embed + positional → self-attention → inserted vector → f.f. → softmax → generate a token  
This loop continues as long as the tokens finish

other architectures of transformers;

(seq-to-vec)

① seq-to-seq encoder-only → used for classification tasks (like sentimental analysis) like BERT

② Enc-dec models → tasks such as text generation → BART and T5

③ decoder-only models → can be generalized to most tasks → GPT - BLOOM - Claude - -

## Prompt and prompt engineering

Some terminology: text that we feed to the model: Prompt

often limited size

the act of generating text: Inference

the output text: Completion

full text space or the memory to use for prompt: Context window

Aka. "in context learning"

Prompt engineering is the act of setting the prompt, so that the model gives the user the desired results. The prompt engineering, divides by "how many samples" completions that we need to give as the prompt, to have the desired completion" in to these categories:

Zero-shot inference: just give the prompt, without any samples.

in context learning  
the process of training/  
fine-tuning model using  
e.g. zero-shot -

one-shot

give a single completed example to the model as prompt. | prompt engineering

few-shot

n more than 1 sample

.....

Prompt  
results  
Give the sentimental analysis result:  
"The movie was ~~bad~~ good"

Context window

model

"positive"

Completion

Zero-shot prompting

Give the sentimental - - - - - :

"The movie was ~~bad~~ good"  
result: Positive

Give the - - - - - :

"The movie was nice"  
result: Positive

model

"positive"

oneshot prompting

clips

\* In most cases, the ~~best~~ model's completion won't improve above 5/6 shots. start with zero-shot, end with 5/6 shots if it wasn't satisfying

- \* Larger models (parameter-wise) can easily work with zero-shot prompting. ↳ direction go for another mode
- \* Smaller models may need one or few-shot prompting to give the perfect completion

## Generative Configuration

There are some hyper parameters (non-trainable) variables, that can alter and influence the inference phase:

- 1) Max new tokens: The maximum number of tokens, that can be generated by the transformer. (it's not a hard cap) - model can end the completion with <eos> token before the max new token value
- 2) Sample top K: The output of softmax layer, is a vector containing proba values for each token in the dictionary, to be used as the next token.

To pick the next token from this vector, we can use two methods:

a) greedy: pick the token with highest proba. value → good for short completion  
↳ may end up in repetitive sentences in completion

b) random (-weighted) sampling: select a token using random-weighted strategy across the proba of all tokens. (sample = True)

To control this random sampling, we can use Top k, Top P and temperature.

↳ sample top k, limits the sampling space to the first k number of proba vector

3) Sample top P:  $n_1, n_2, \dots, n_k$  Proba values which their sum does not exceed the value of P

softmax

Cake	0.2
banana	0.09
was	0.12
man	0.015

$k=3$  → select randomly from first 3

$P=0.3$

Cake	0.2
banana	0.09
was	0.02
man	0.015

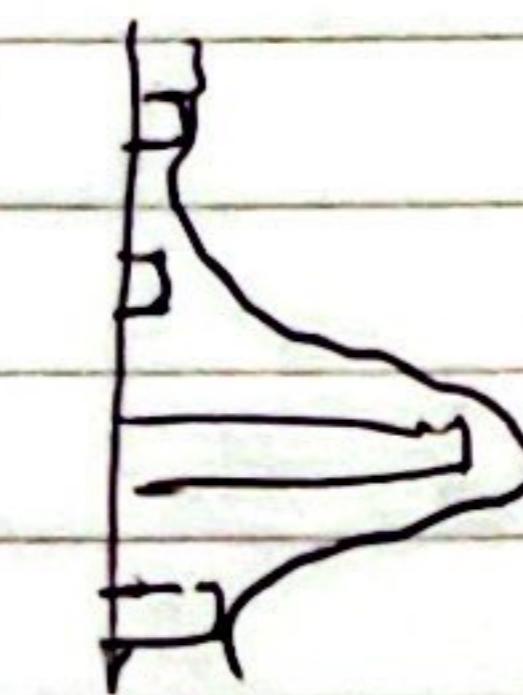
→ select randomly from first value, which their sum is not  $> 0.3$ . (0.29 in this example)

~~A temp=1~~ → default softmax output

4) **temperature**: Directly affecting softmax output, The lower value ( $\approx 0$ ) decrease variability (strongly peaked proba distribution)

Cooler temperature ( $< 1$ )

cake	0.002
banana	0.01
word	0.49
man	0.1



higher temp ( $> 1$ )

cake	0.3
banana	0.2
word	0.5
man	0.1



strongly peaked proba distribution

$\Leftrightarrow$  less randomness

more smooth distribution

$\Rightarrow$  more randomness

lower top k

higher top P ( $\approx 1$ )

higher temp

$\Rightarrow$  more random generation

## Generative AI life cycle of project

most important part of project

choose from existing pretrained models or train your own  
based on the scope of project

start with prompt engineering, then evaluate.

if not satisfied, try fine-tune, v

The generation must align with human feedback

optimize and deploy for the certain application /

use-case

LLMs are prone to inventing / hallucination

$\Rightarrow$  bad at math/reasoning / -

optimize it to overcome those

