

DS 6051 Decoding Large Language Models

Introduction to Transformers

Chirag Agarwal

Assistant Professor

School of Data Science

University of Virginia

Recap

- Introduction to Tokenizers
- Byte-Pair Encoding
- Problem with Tokenizers
- Extending Tokenizers

The Early Ages

- Prior to transformers

- Seq2Seq models - Sequence to Sequence
- Recurrent Neural Networks
- Long-Short Term Memory - LSTM
- Gated Recurrent Unit

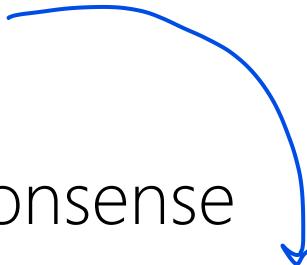


Problems:

- ✗ Context: Cannot understand
- ✗ Long term dependencies

The Current Age

- Transformers - *Solve the probs above*



- Logical, Mathematical and Commonsense Reasoning
- Alignment using RLHF
↳ align/match LLM to human principles (fairness/robustness)
- Safety testing and jailbreaks
- Ethical and Fair models



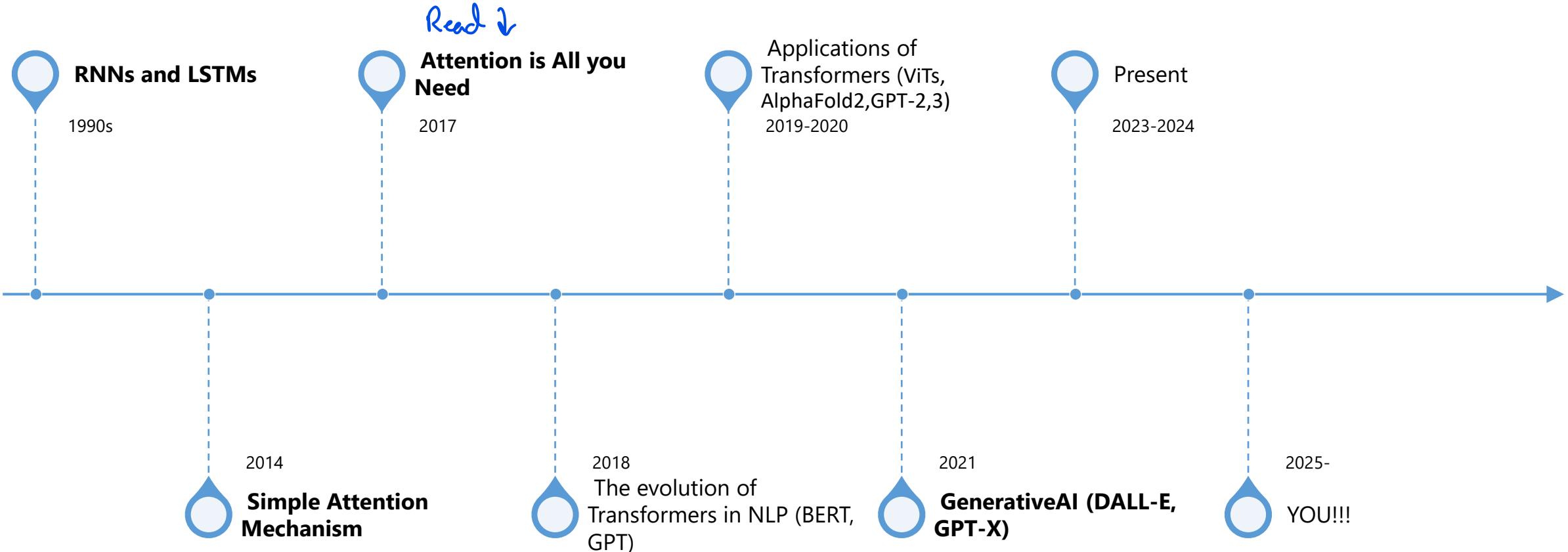
Context



Long term dependencies



Evolution of Attention



The Binding Model!!



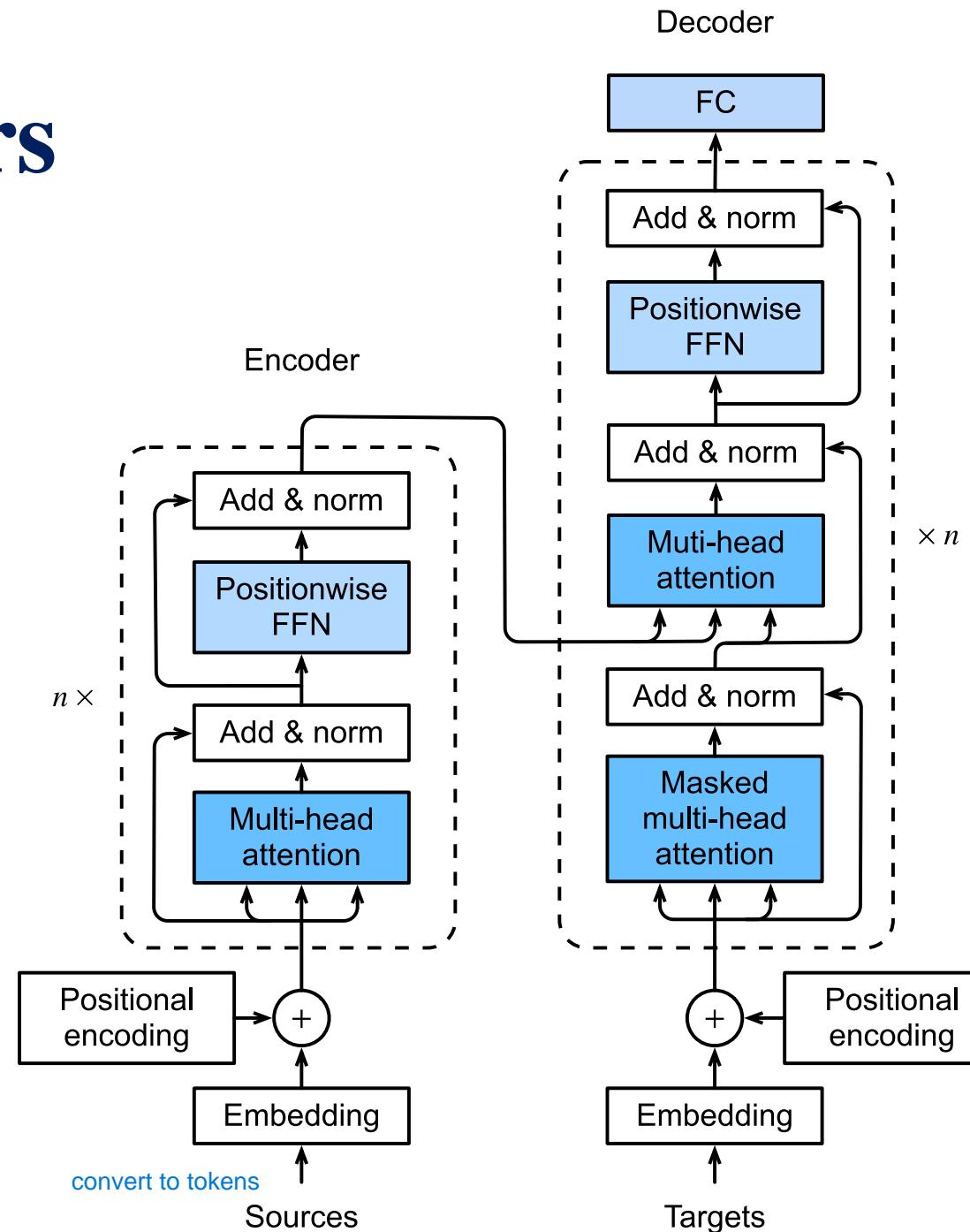
Any data you can convert to a sequence you can use transformers on it!!!

↳ tokenize



Transformers

Attention Is All You Need
Image →



Input Processing

1st task is to tokenize

- Input
- Embedding Layer
- Position Embeddings

Let's consider a Dialogue Completer

Input Dialogue

It is our choices, Harry, that show what we truly are,

If you want to know what a man's like,
take a good look at

It matters not what someone is born,

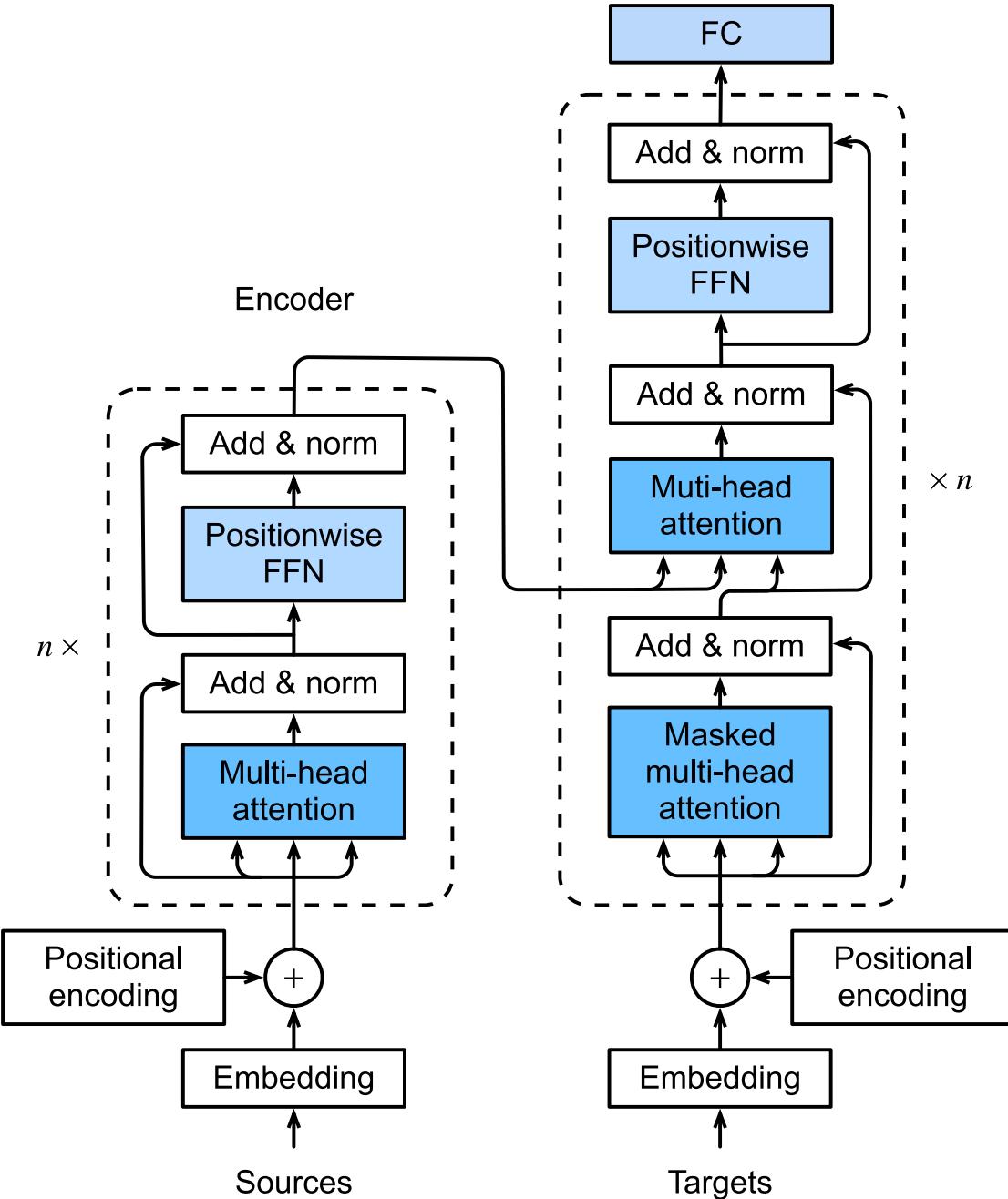
Dialogue Completion

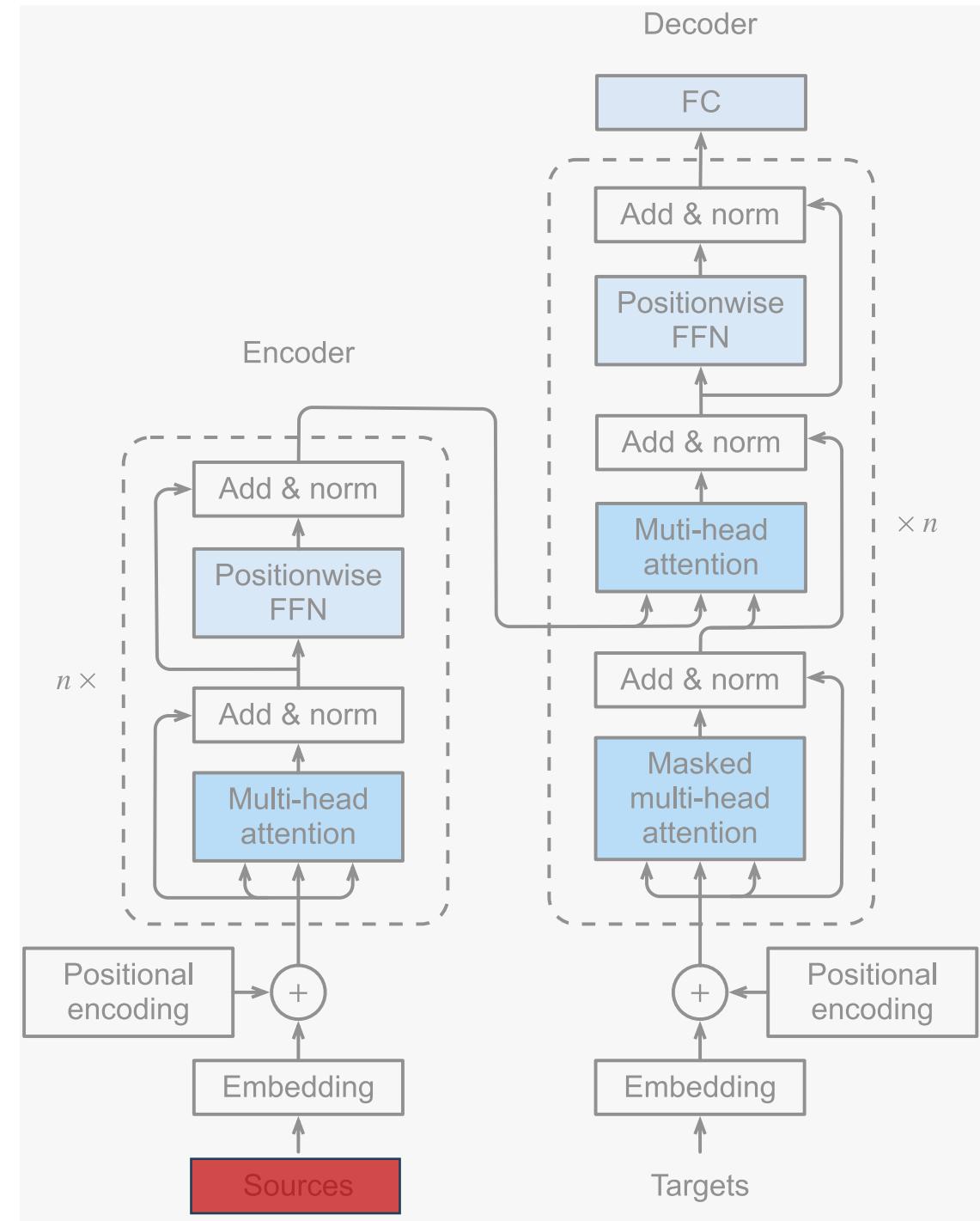
<start> far more than our abilities
<end> *Prompts*

<start> how he treats his inferiors, not his equals <end>

<start> but what they grow to be <end>

Decoder



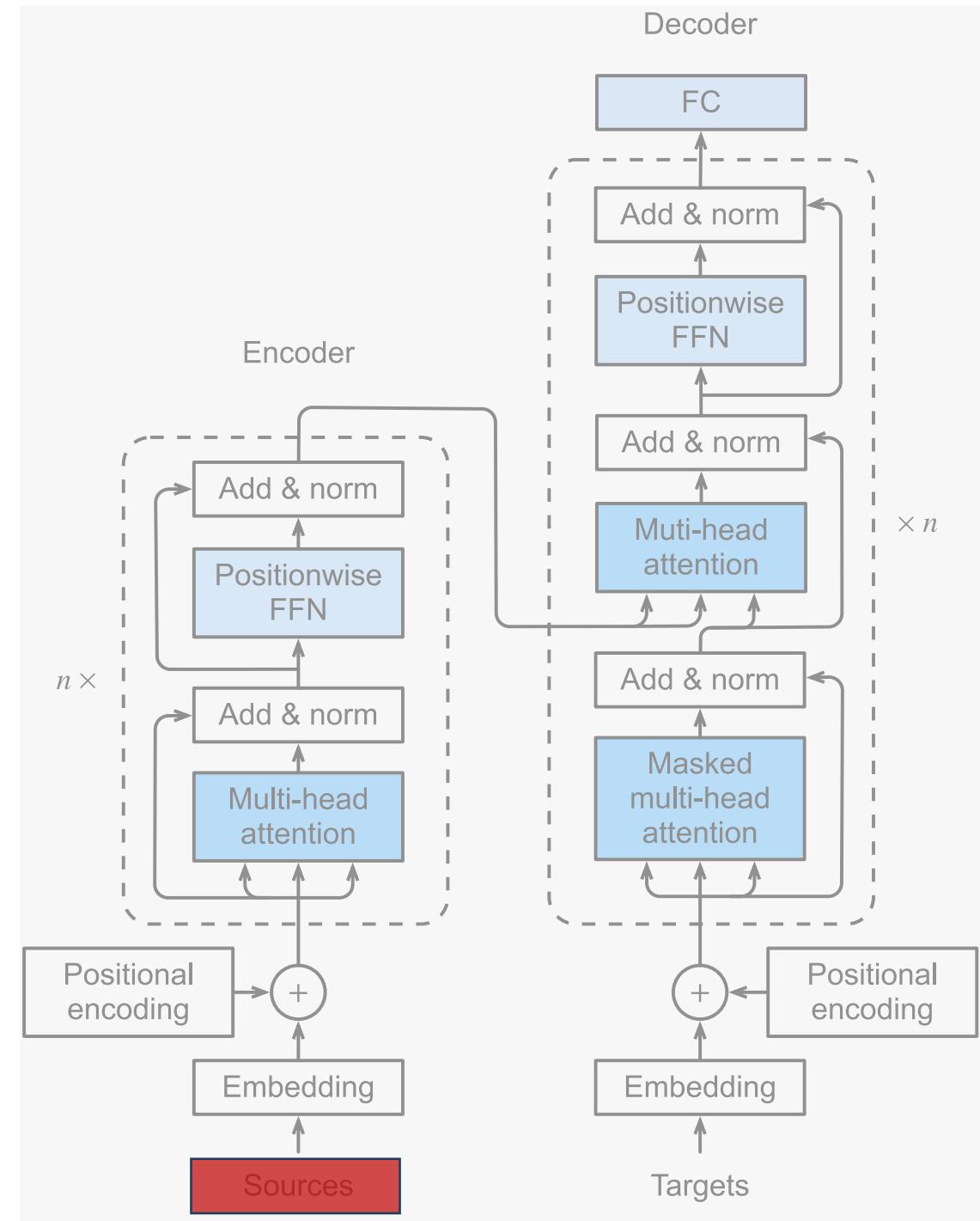


Text

Input

It matters not what someone is born,

↑
ignore



- First LLM is given a vocab to recognize

Input

Text

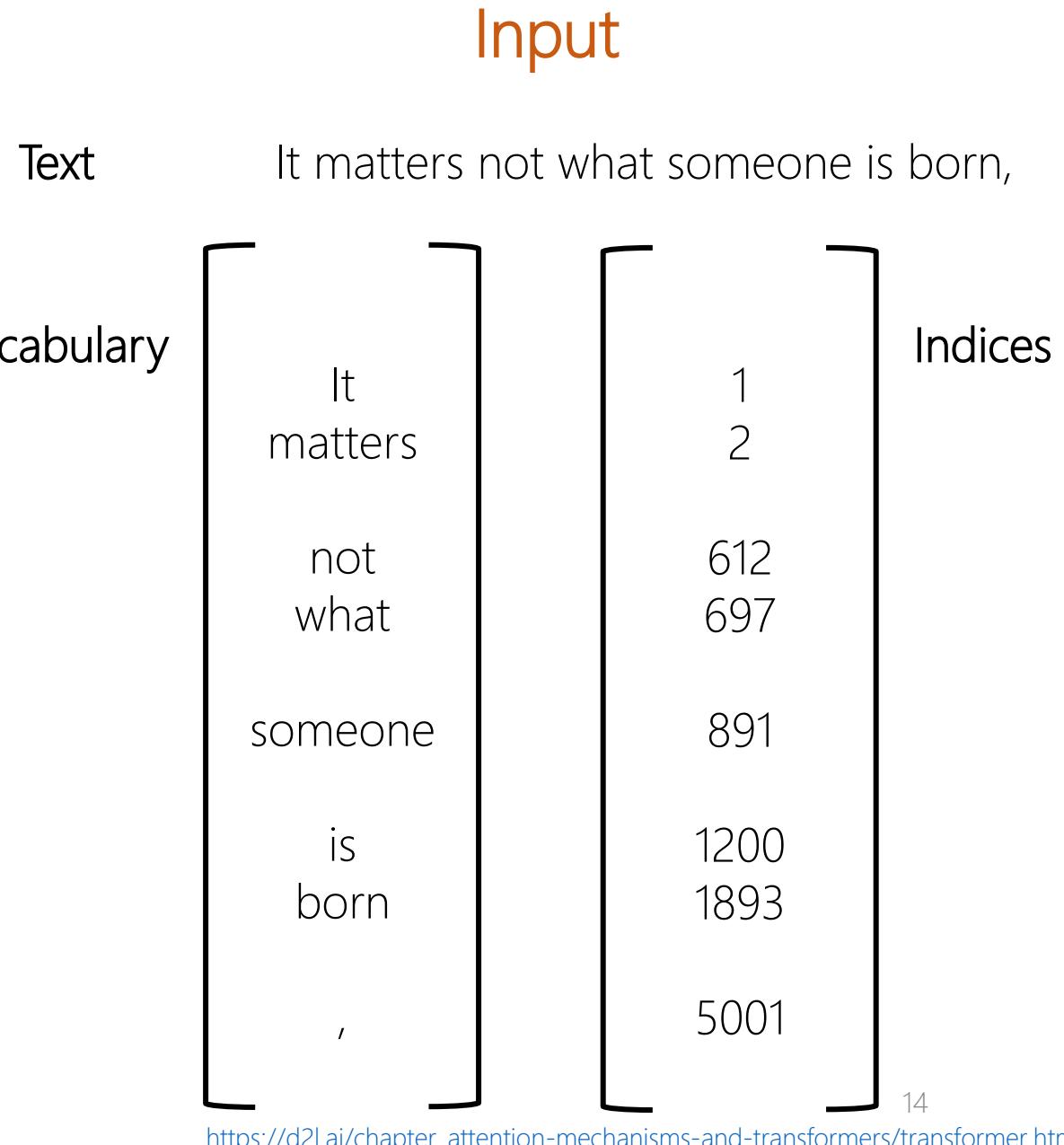
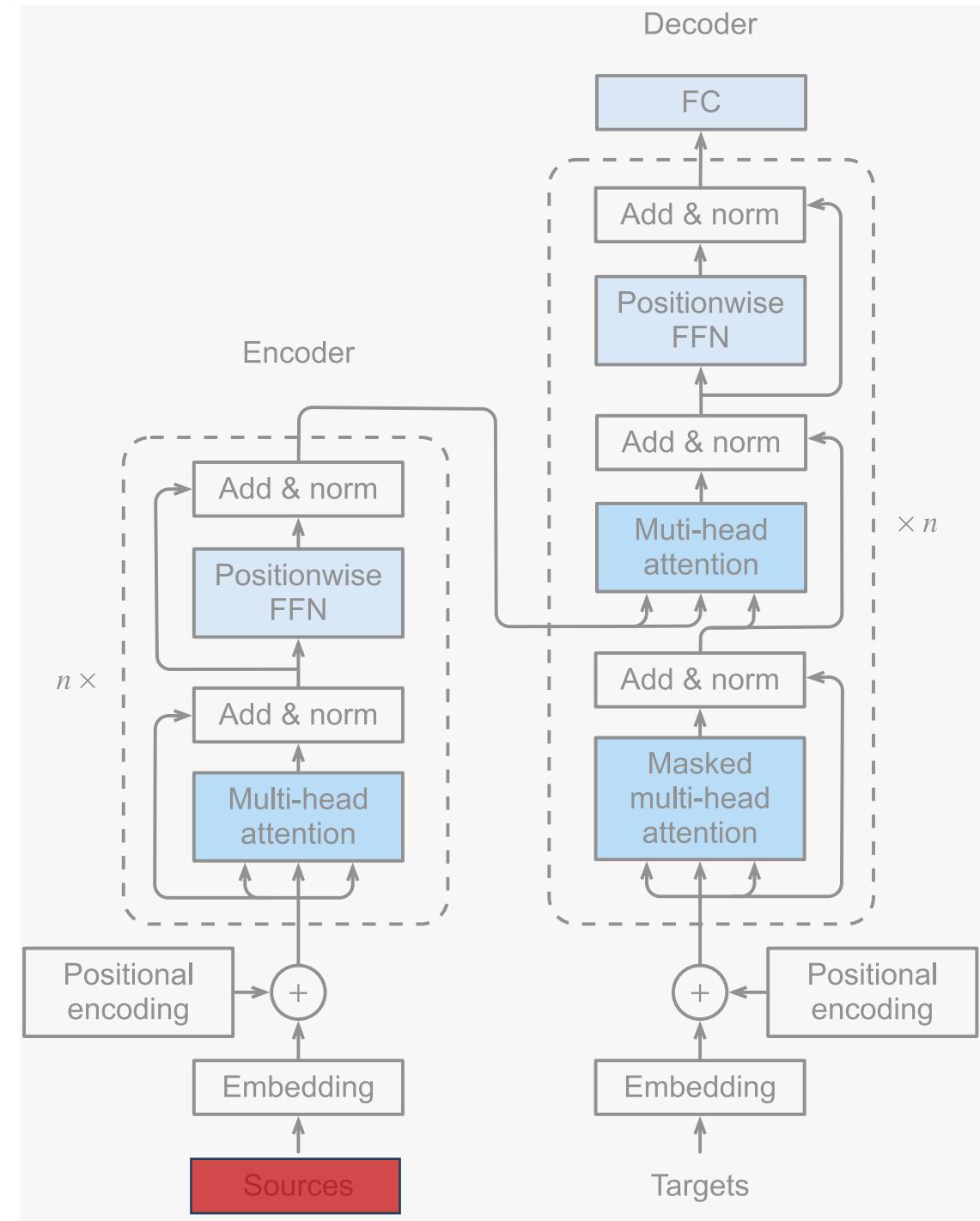
It matters not what someone is born,

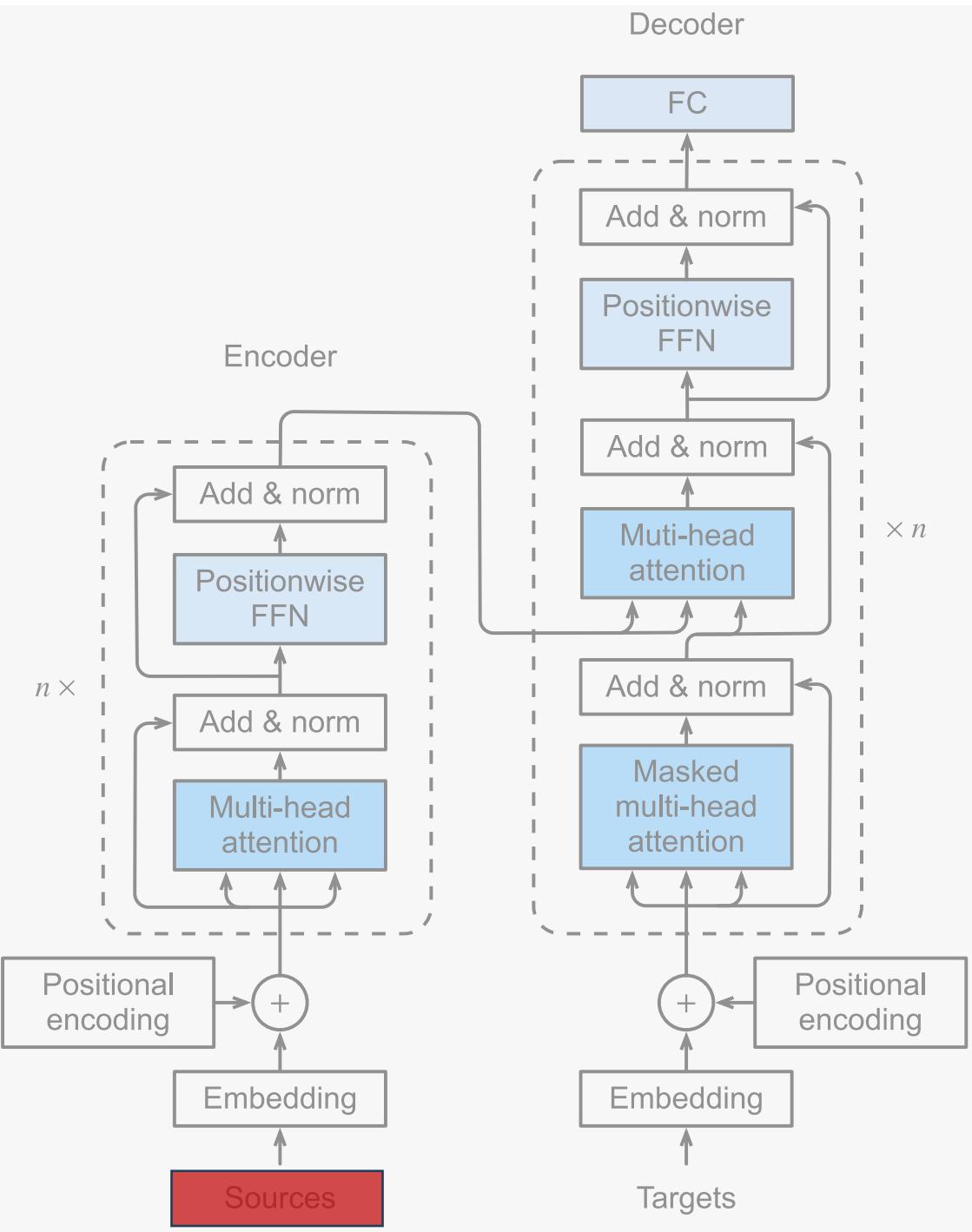
Vocabulary

...	Harry
0	It
1	matters
2	Ron
501	not
612	what
697	...
...	someone
891	<YouKnowWho>
1111	is
1200	born
1893	Hermione
3019	,
5001	

Indices

of those words in the vocab, retrieve the tokens of words we input





Input

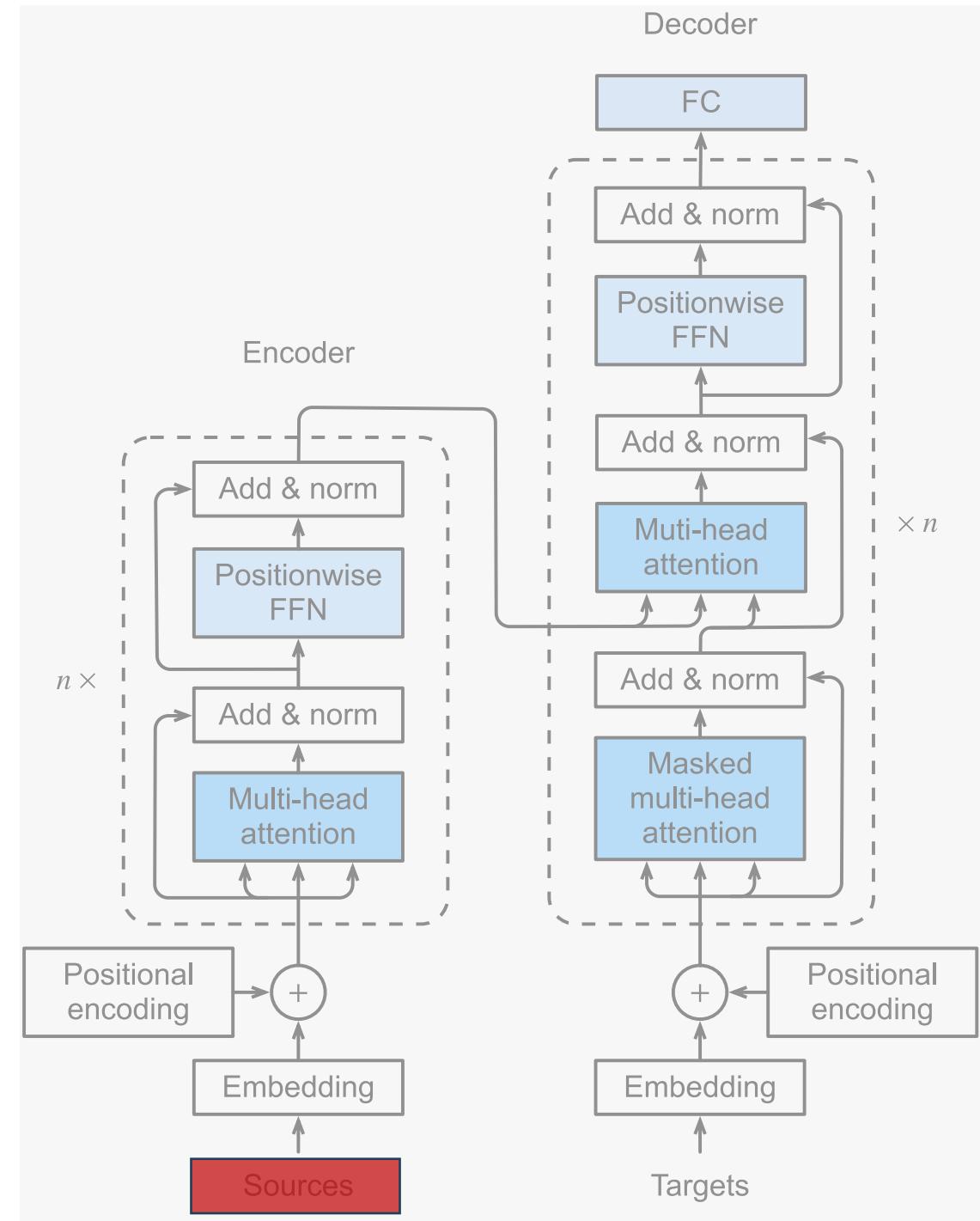
Text

It matters not what someone is born,

Vocabulary
Indices

x_0	x_1	x_2	x_3	x_4	x_5	x_6	x_7
1	2	612	697	891	1200	1893	5001

Now that is the input sample of indices, exact location of those words in the vocab it has
This is passed to the embedding layer



Input

Text

It matters not what someone is born,

Vocabulary
Indices

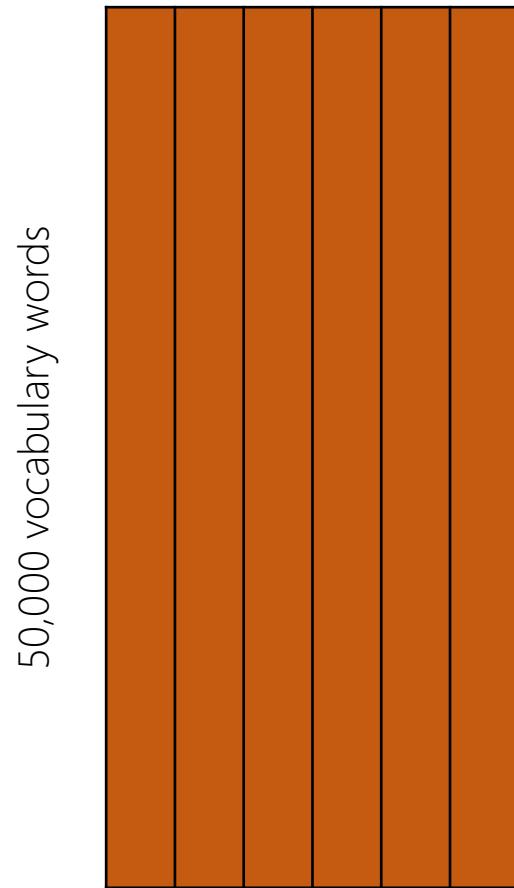
V is vocab amount
D dimensions are what it is scaled down to

embedding layer 50K X 128
but we don't have to send the
whole 50K we can send a subset of just
relevant to the input

Input to be
passed to the
embedding layer

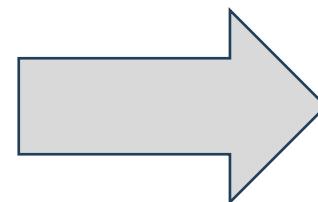
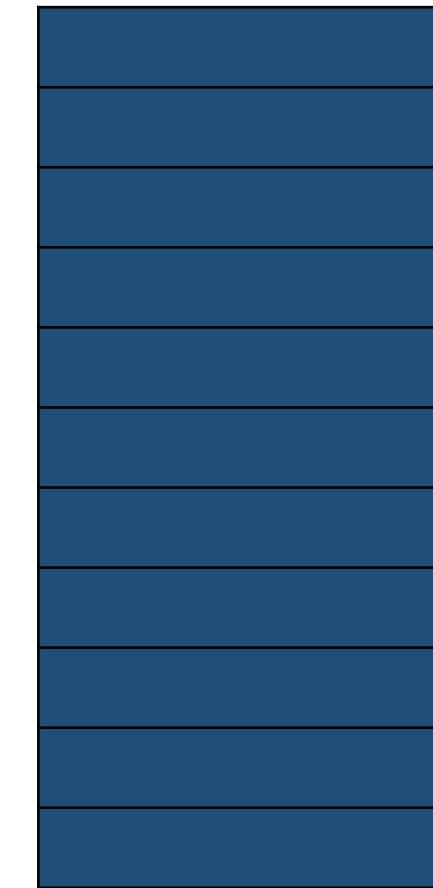
Hidden Layer Weight Matrix

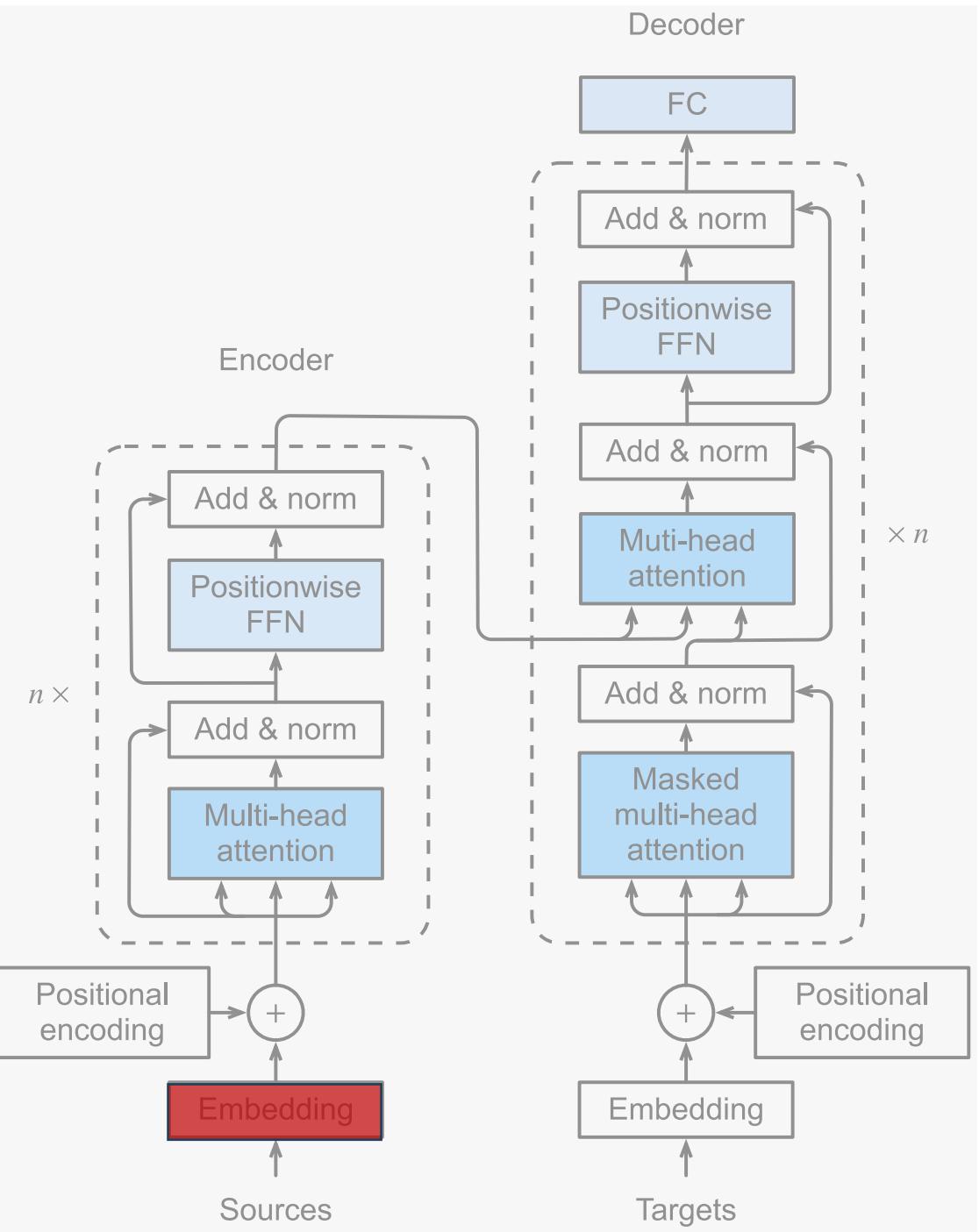
100 hidden layer neurons



Word Embedding Dictionary

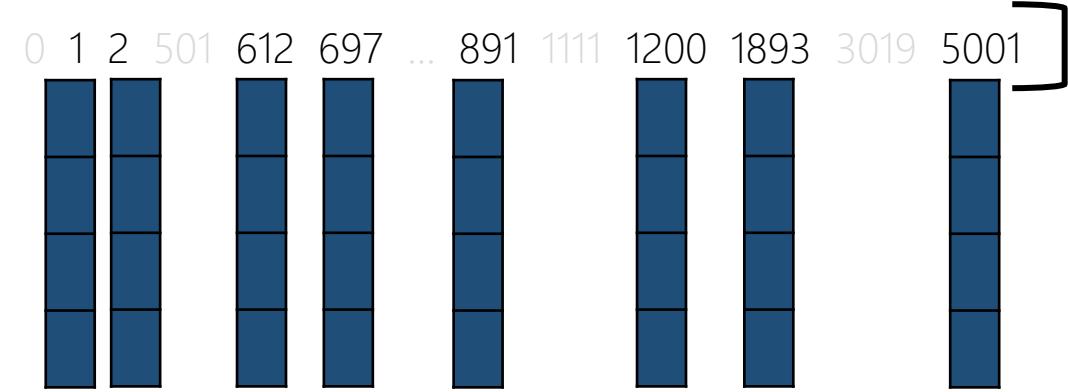
100 embedding features





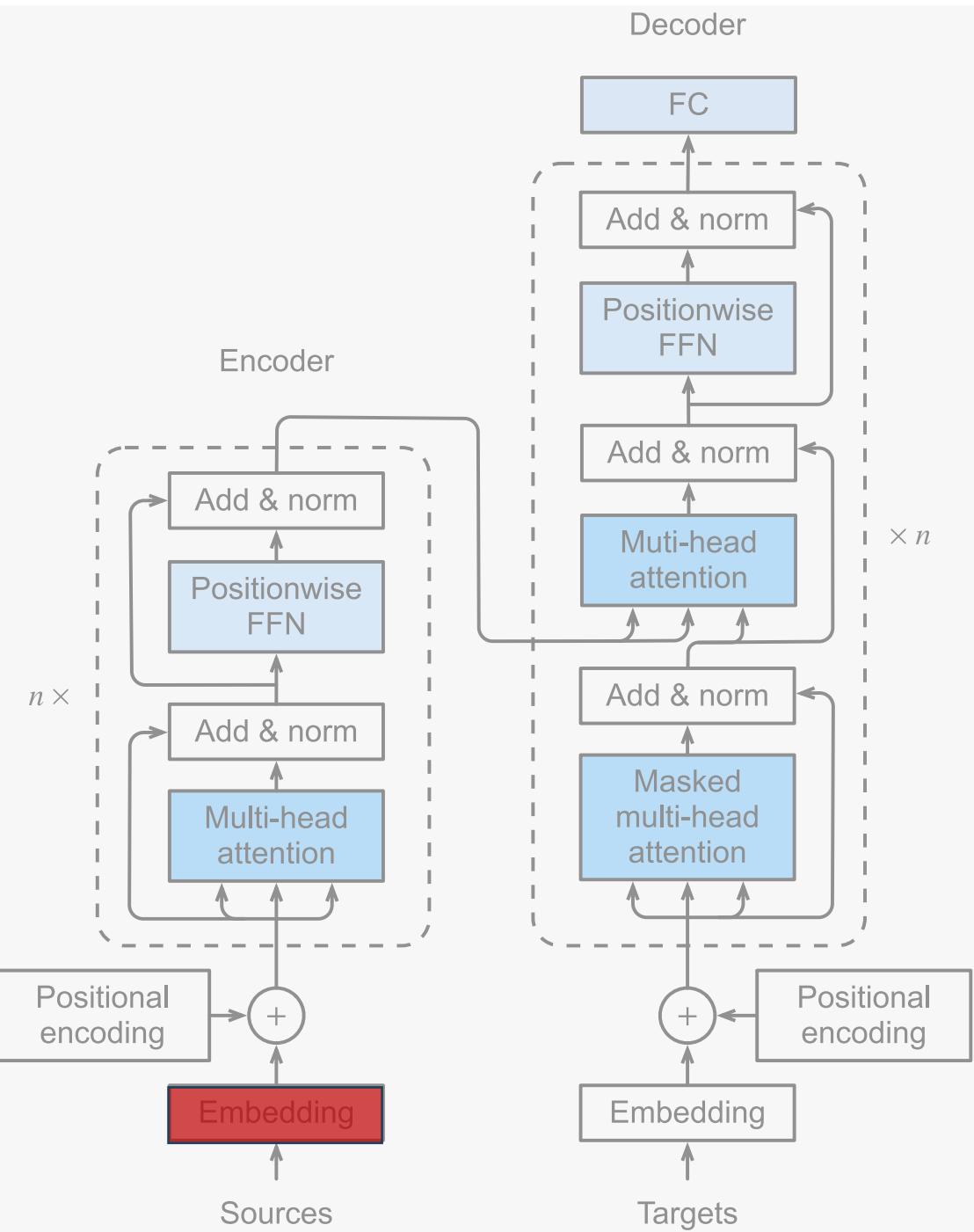
Input Embeddings

Vocabulary
Indices



Input now 7x4

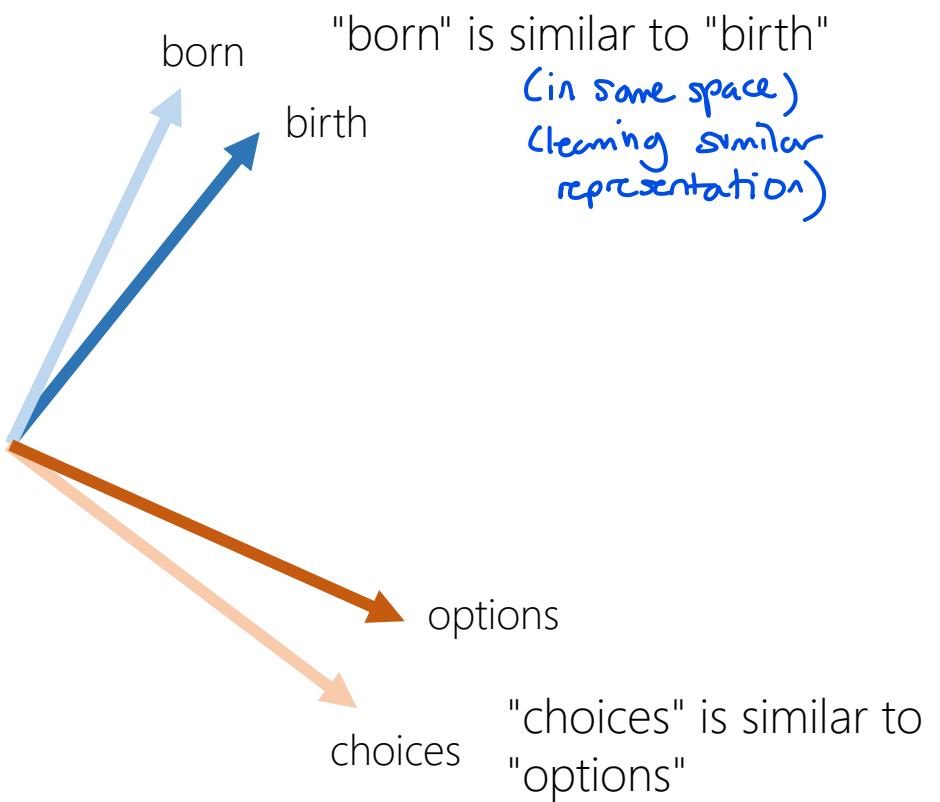
Input Embeddings

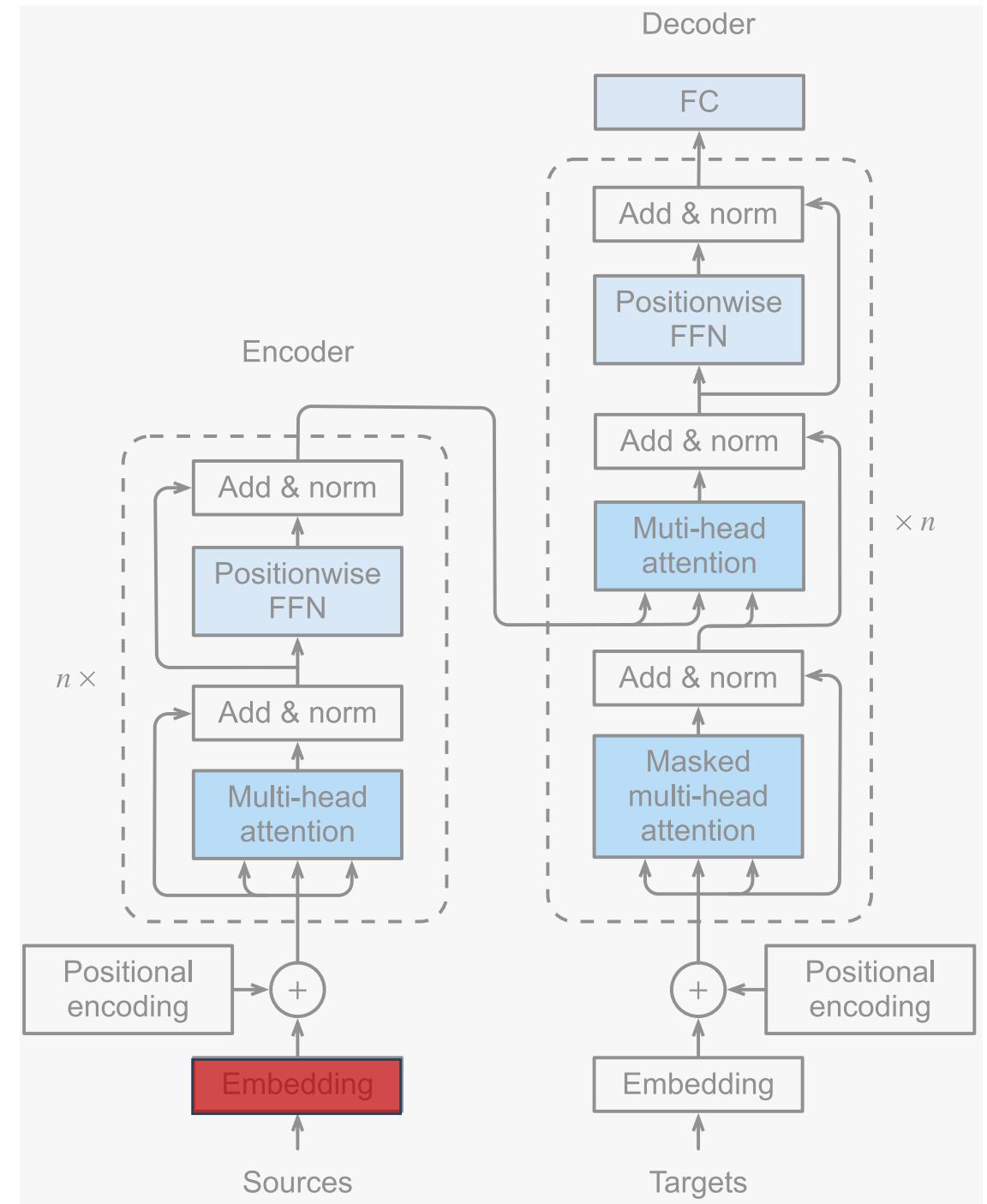


Vocabulary
Indices

... 0 1 2 501 612 697 ... 891 1111 1200 1893 3019 5001

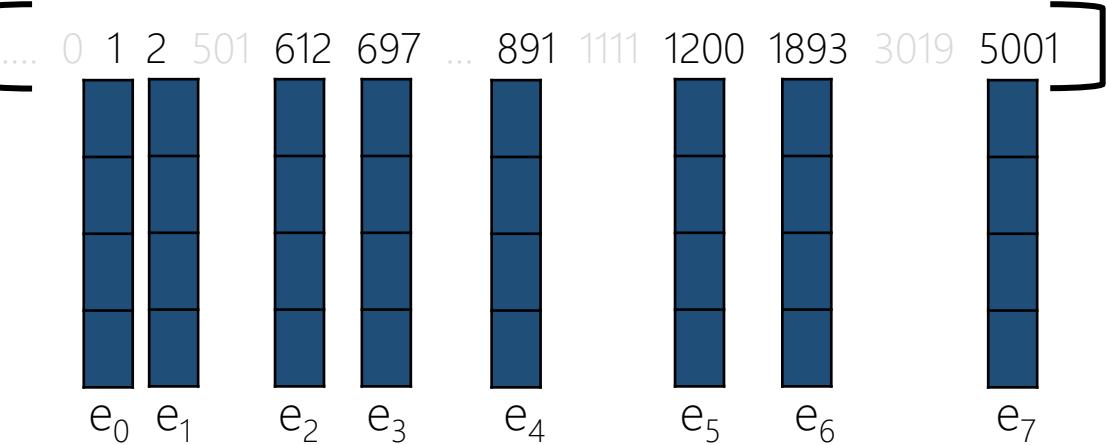
born

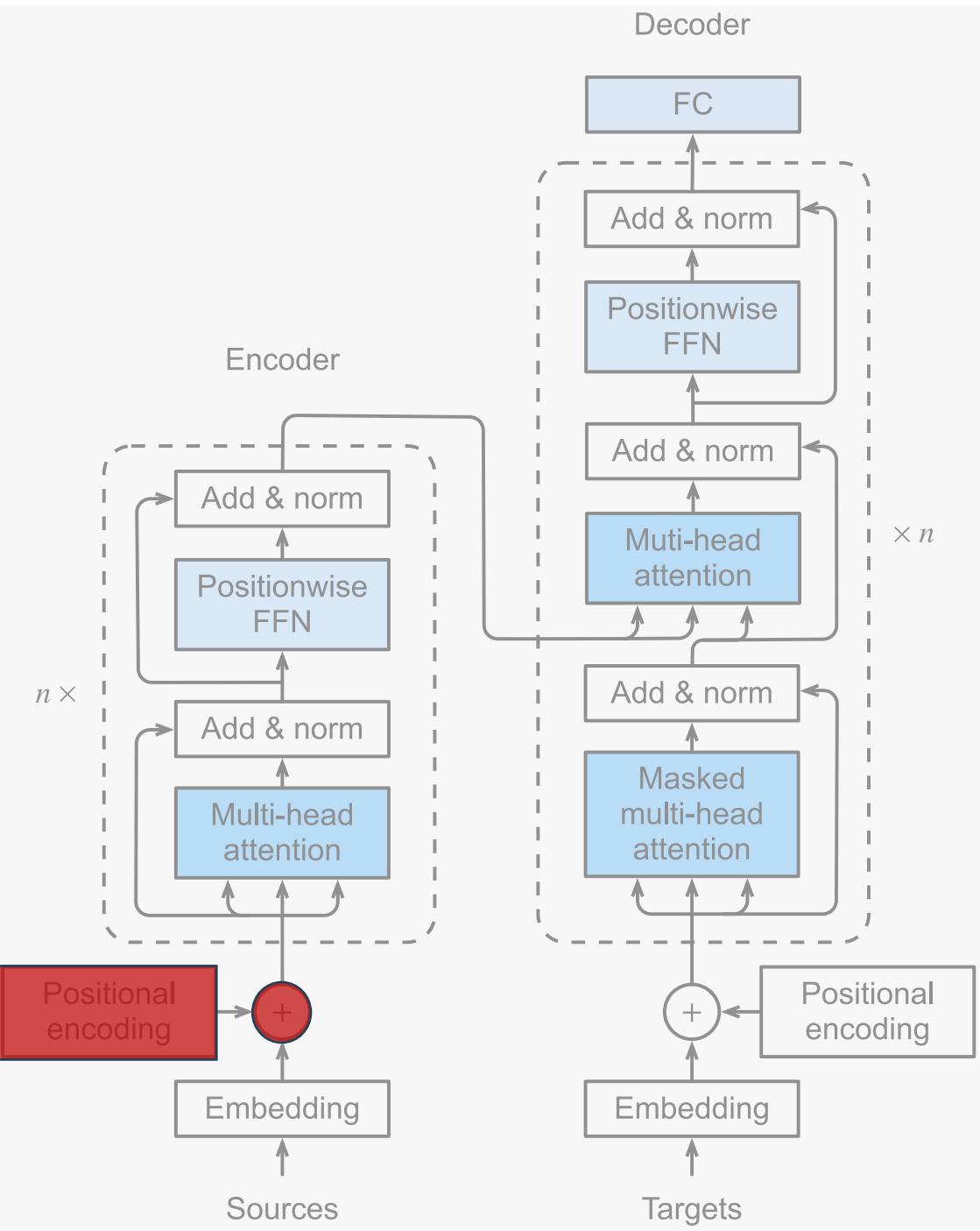




Input Embeddings

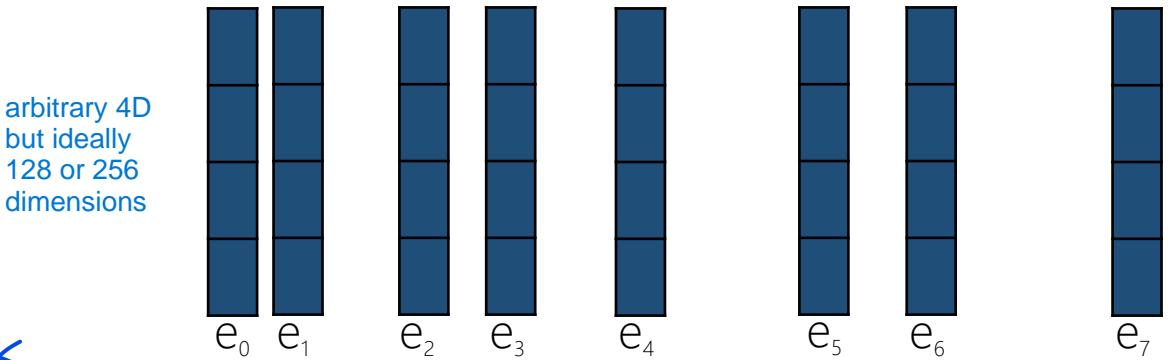
Vocabulary
Indices





Why do we need Positional Embedding?

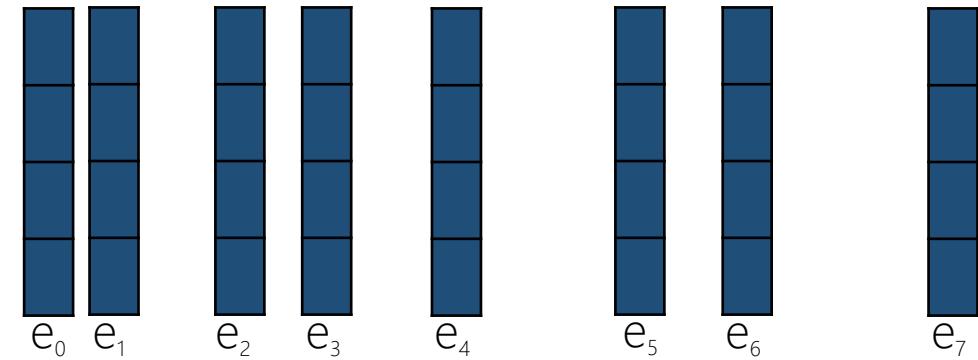
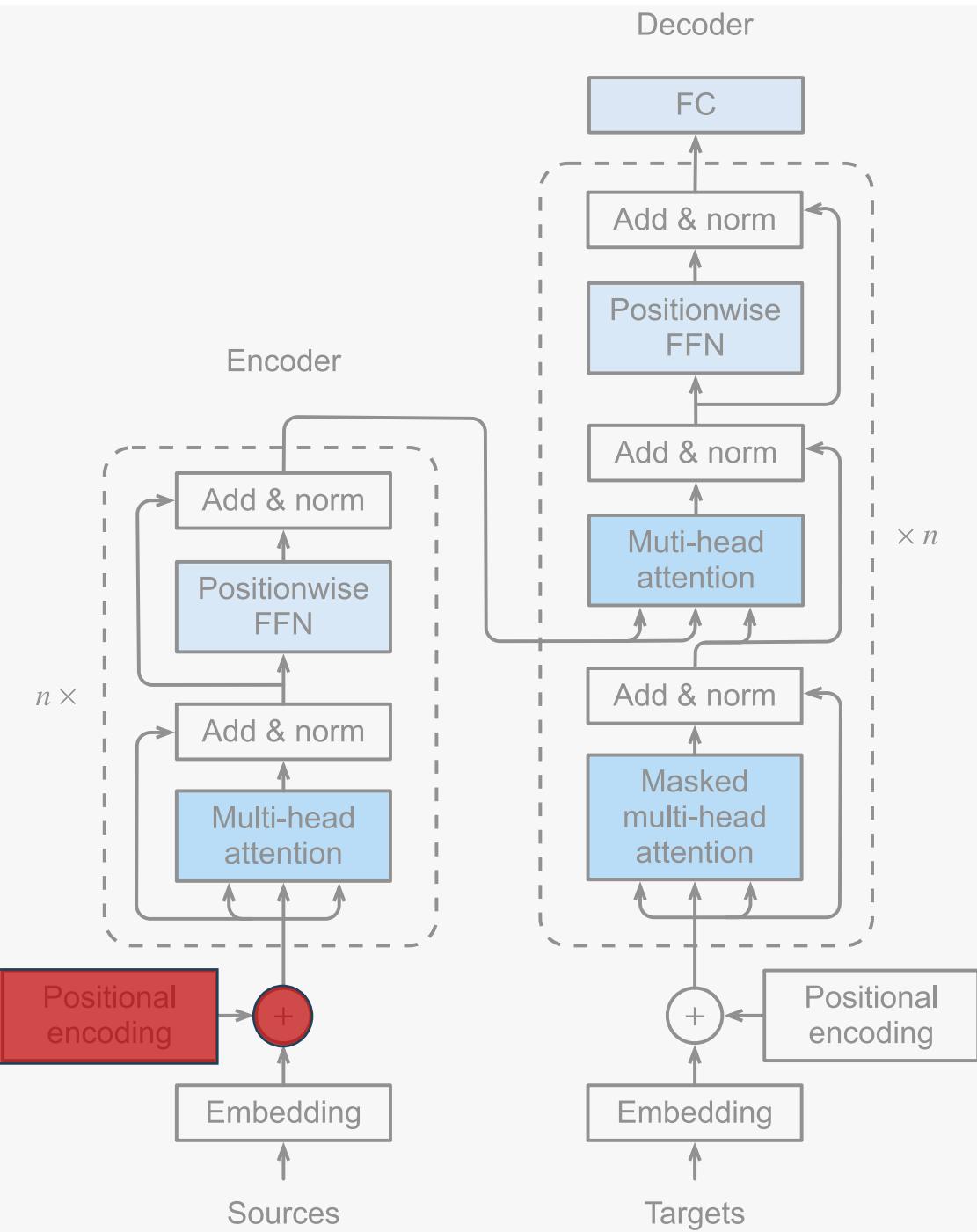
we don't send sequentially, we process in parallel, might miss context/position of words close by



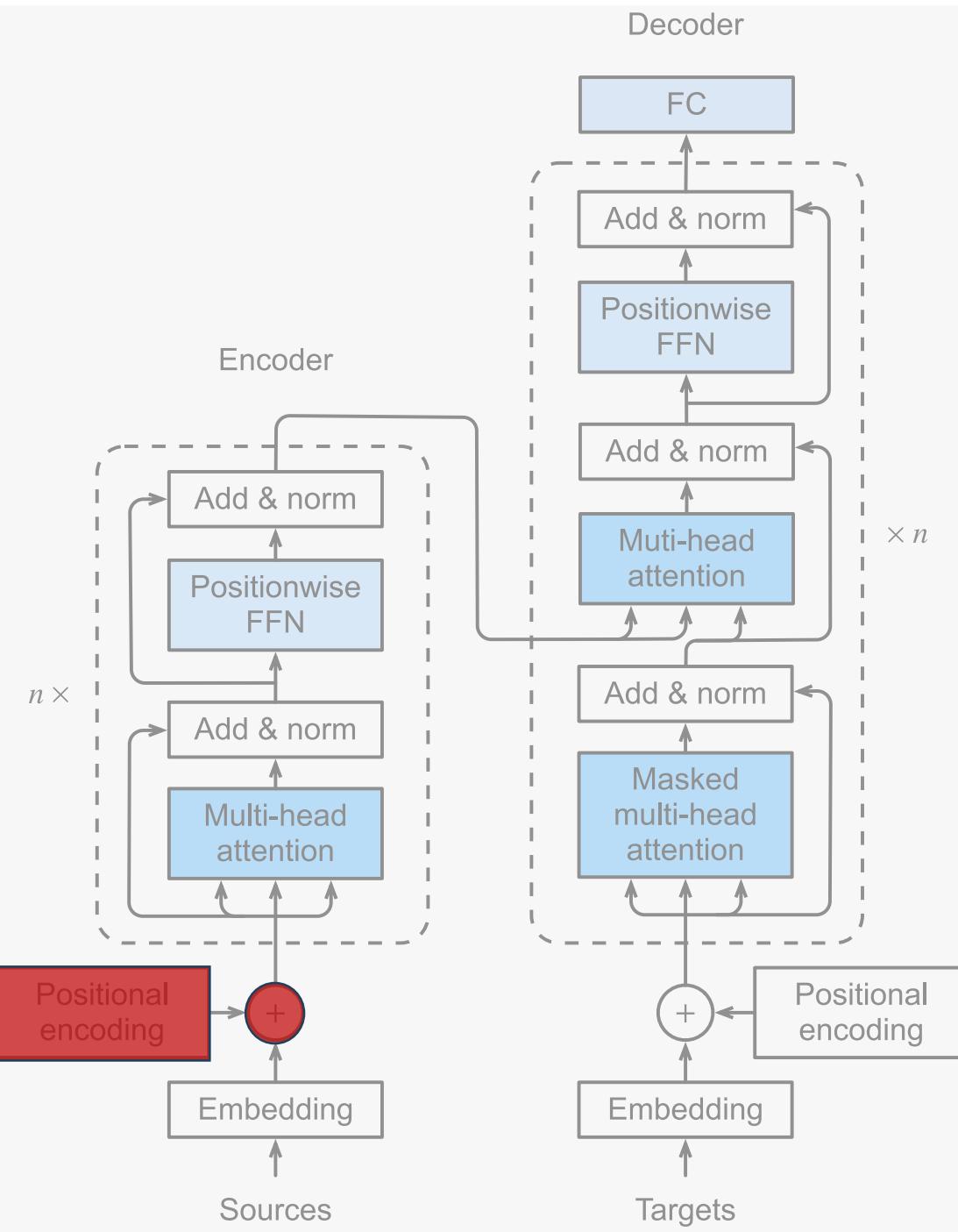
Positional encoding informs the transformer of the order of the words to not lose that info

Recurrent Neural Network
or
Long-Short Term Memory

Why do we need Positional Embedding?



Transformers



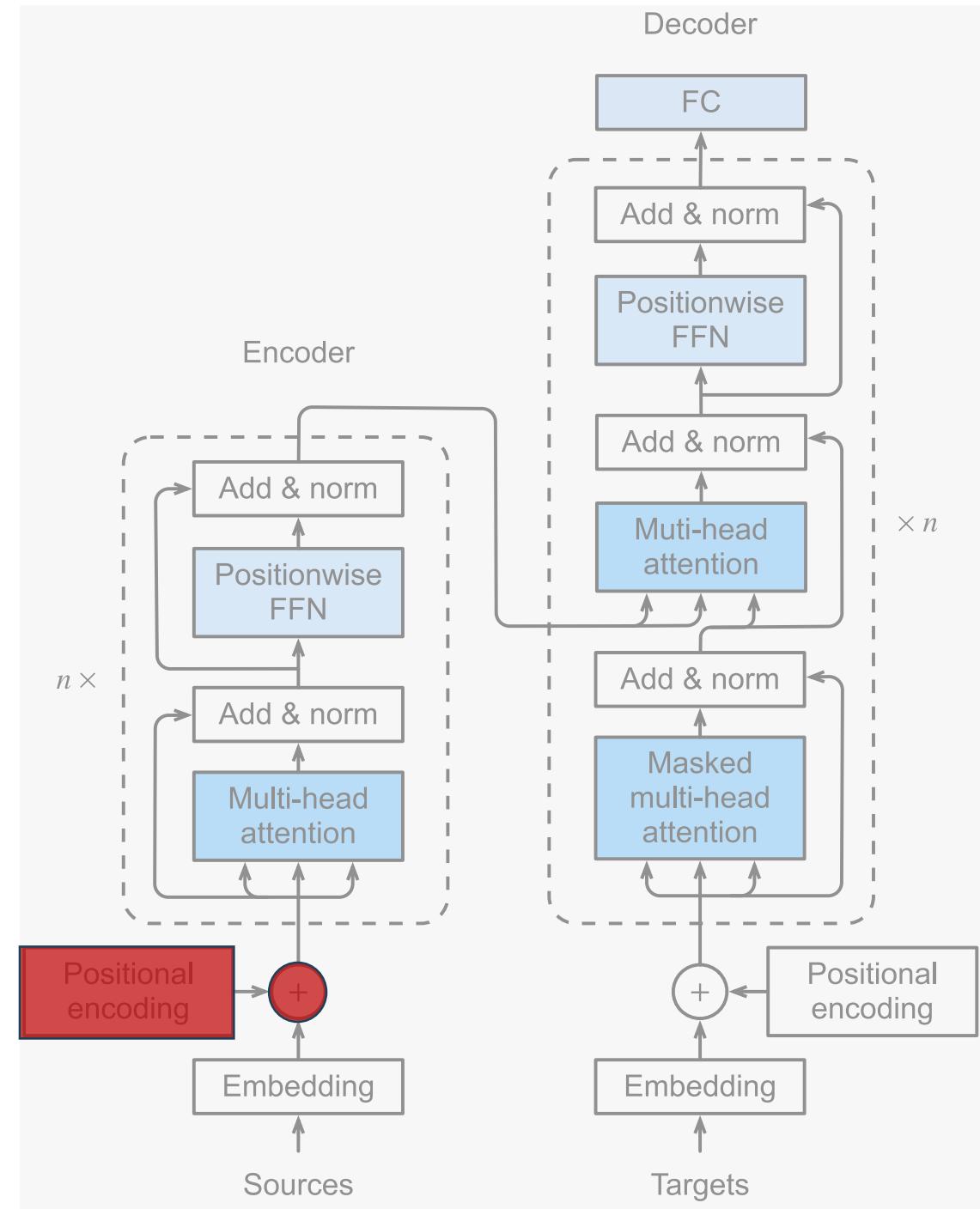
Why do we need Positional Embedding?

Importance of position

Shows why we need the word order to noted

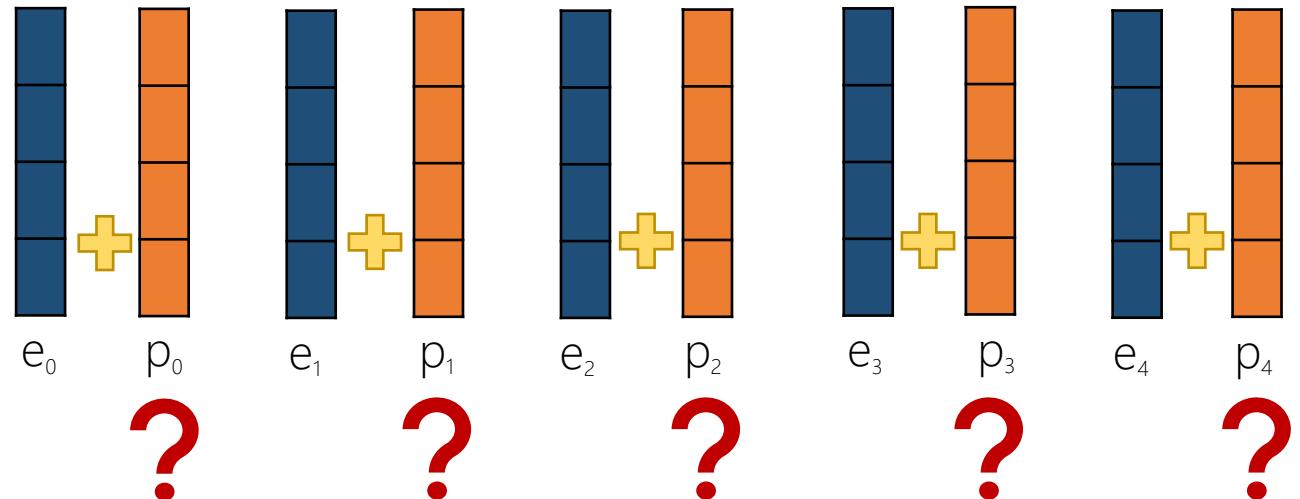
England invented Football, which is now played in 211 countries, including the USA, Canada, and Brazil.

Brazil invented Football, which is now played in 211 countries, including the USA, Canada, and England.

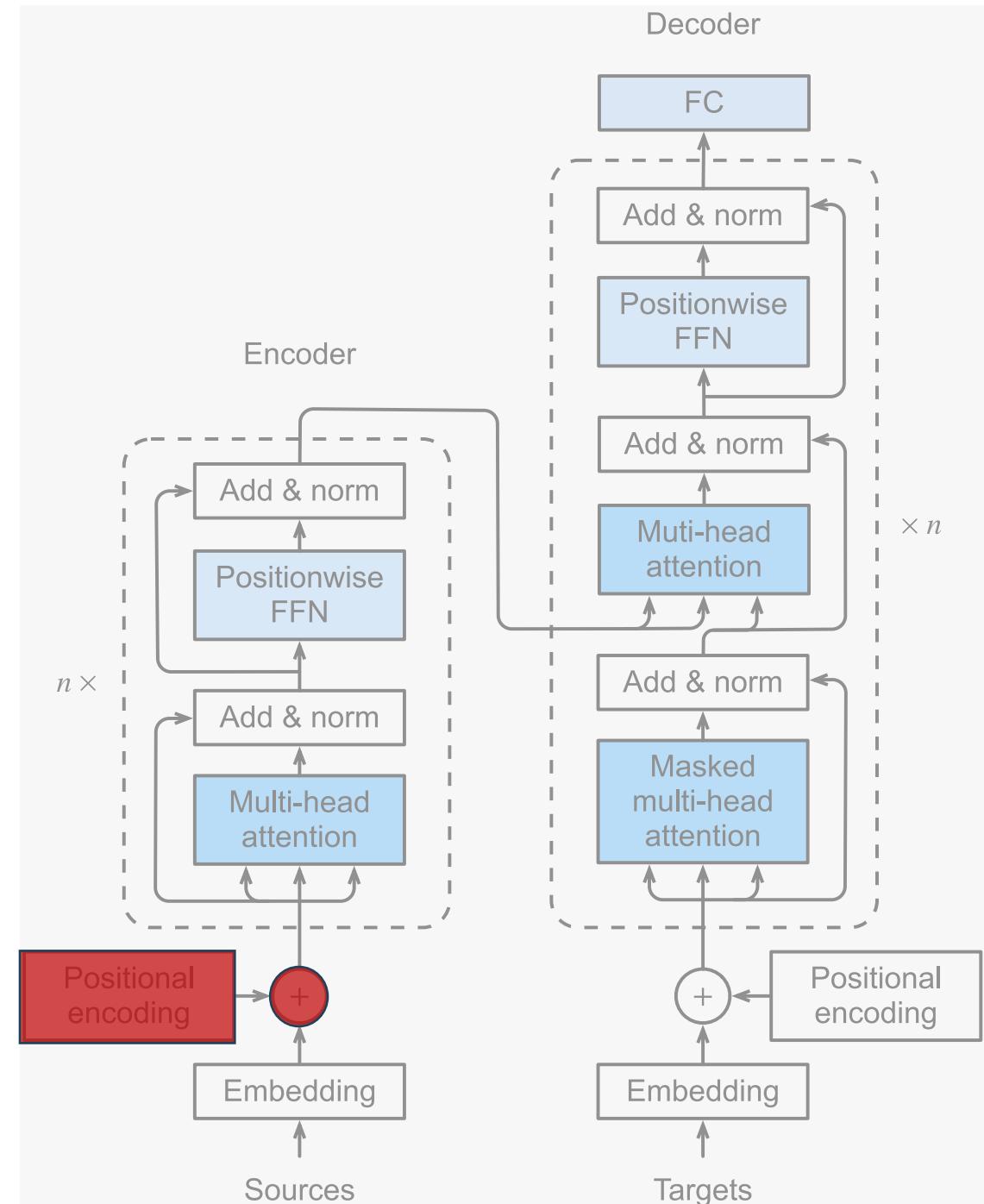


What to use as Positional Embedding?

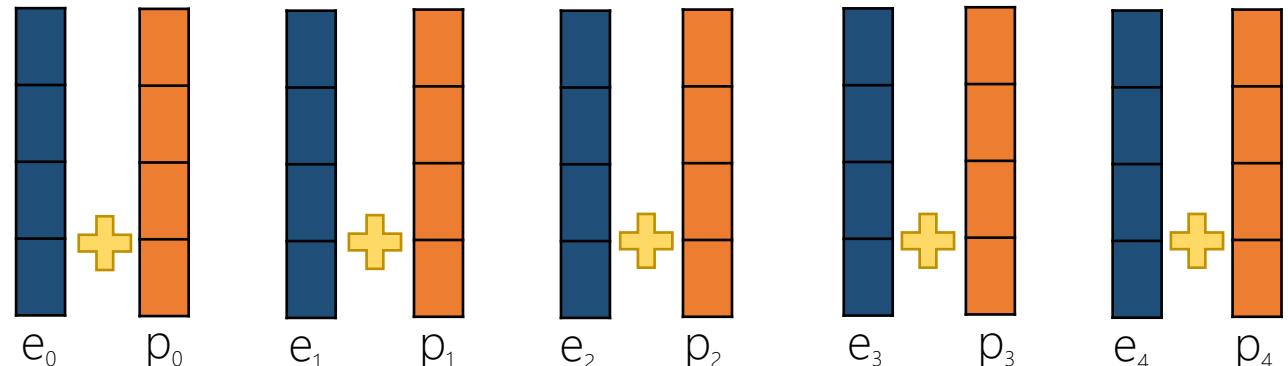
Cannot just use a vector of positions because like some prompts are 1000 words and this is too much for the LLM to process, so we normalize and note position more smartly



tries to
give more
position info



Frequencies for Positional Embedding



This normalizes position

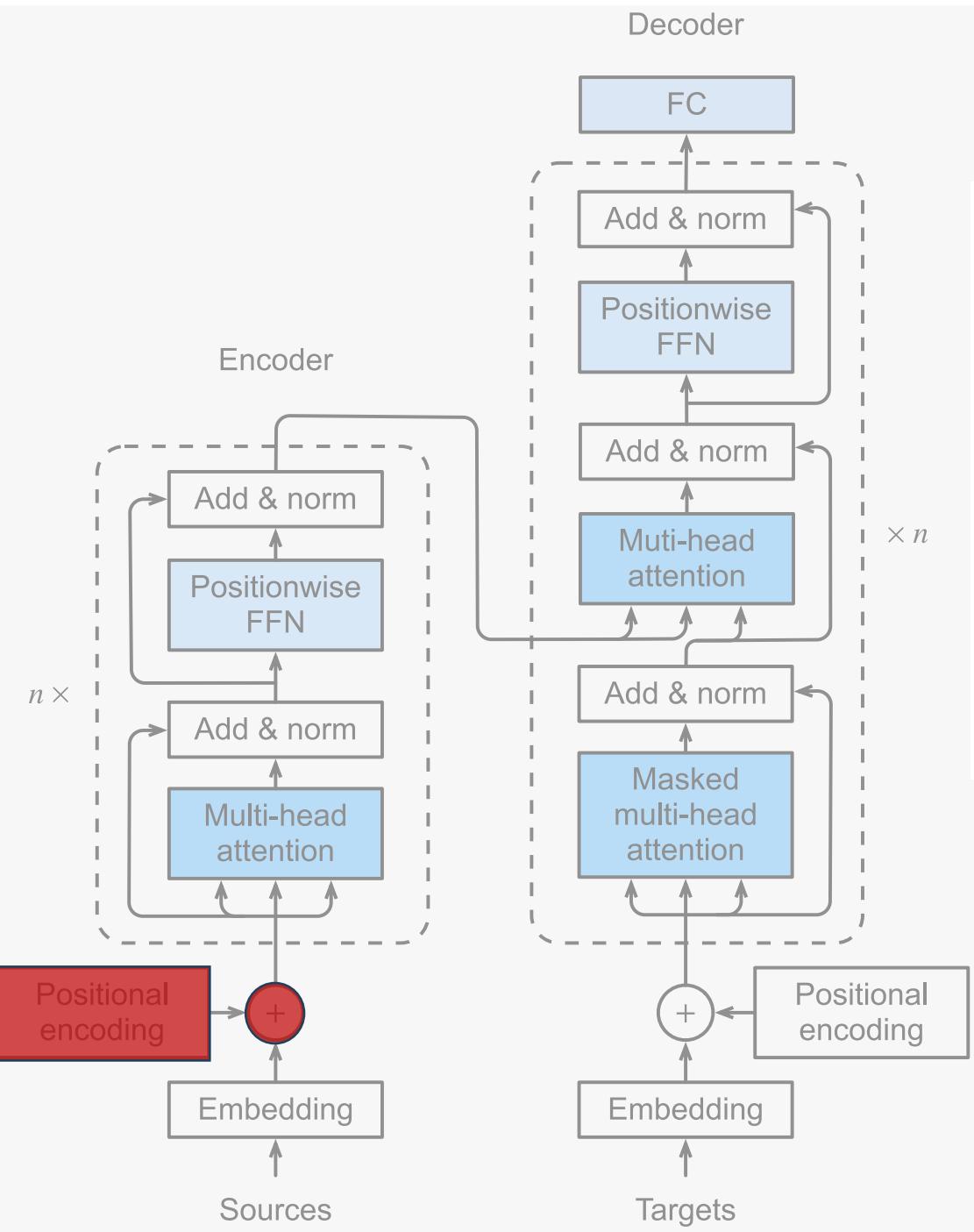


$$PE_{(pos,2i)} = \sin(pos/10000^{2i/d_{\text{model}}})$$

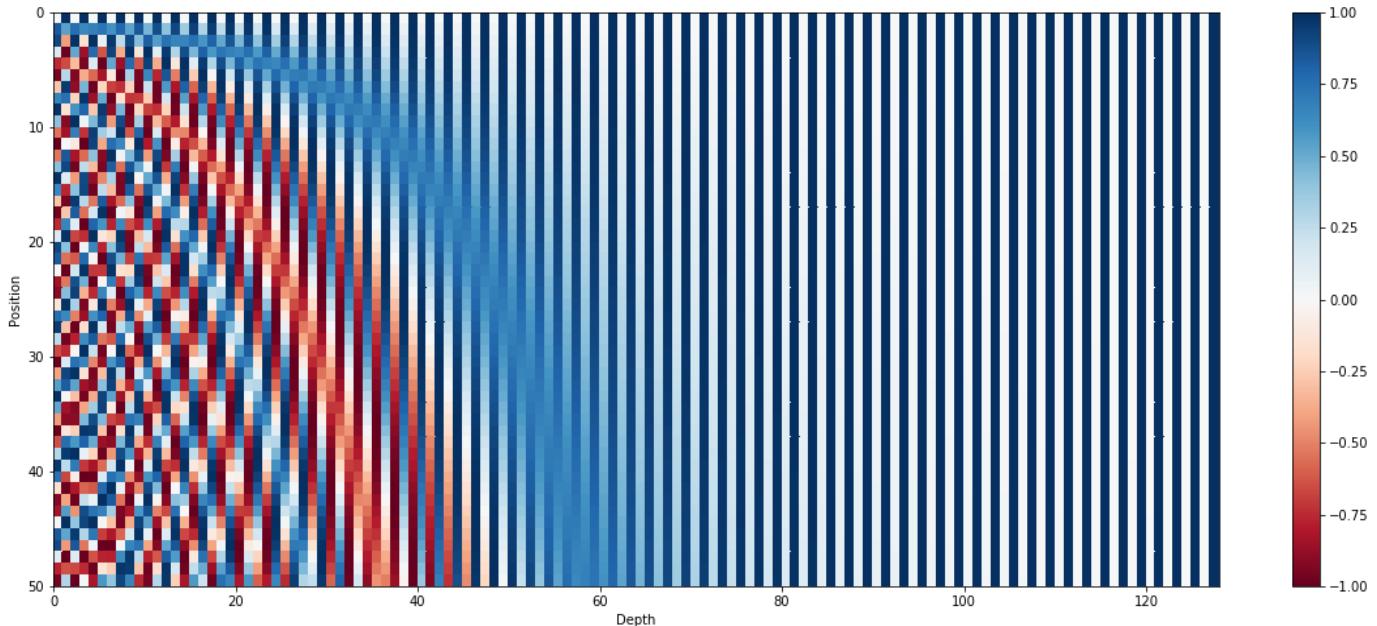
↑ location of word in sentence
i-th location in dim

$$PE_{(pos,2i+1)} = \cos(pos/10000^{2i/d_{\text{model}}})$$

↑ dimension- 4 here



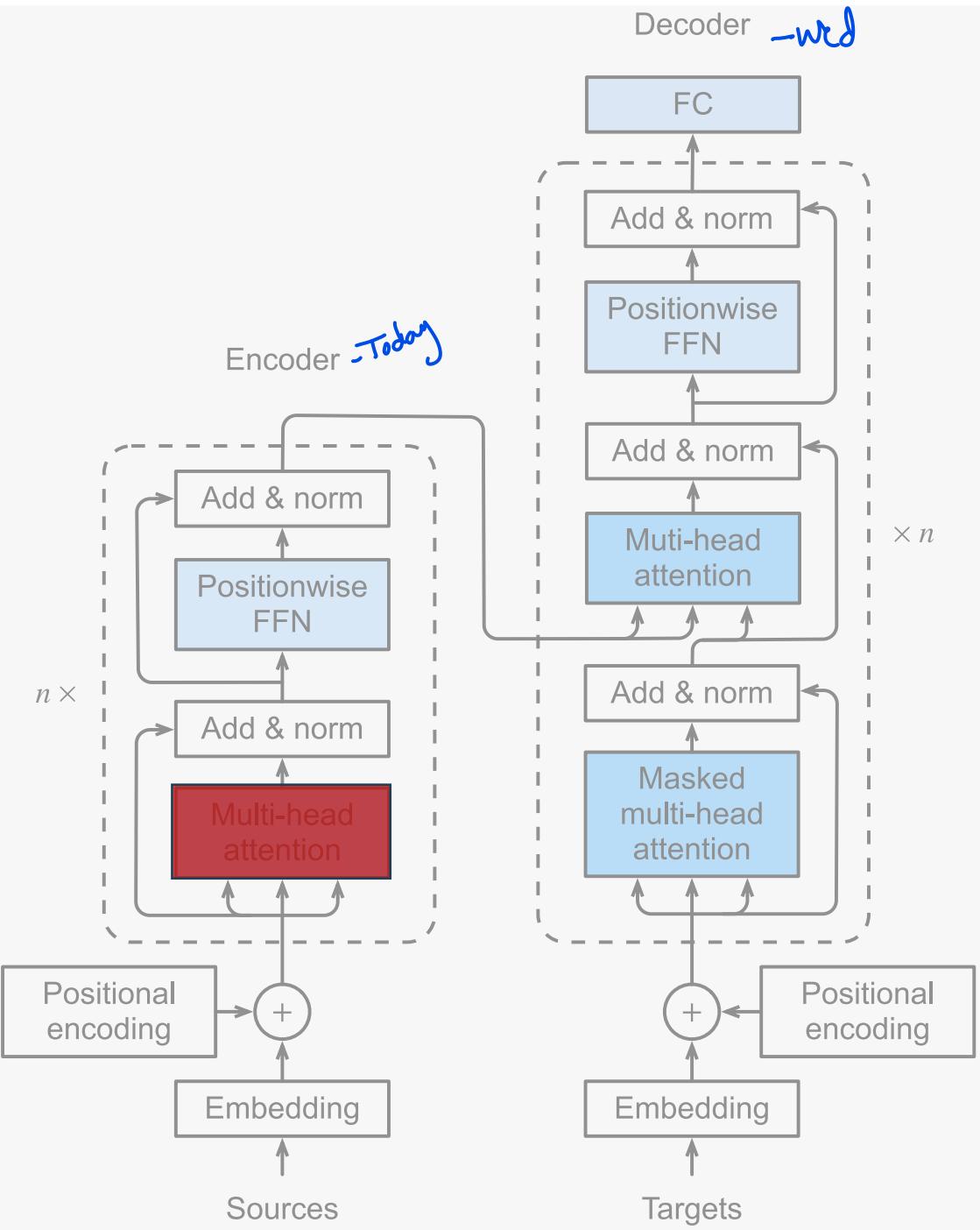
Frequencies for Positional Embedding



$$PE_{(pos,2i)} = \sin(pos/10000^{2i/d_{\text{model}}})$$

$$PE_{(pos,2i+1)} = \cos(pos/10000^{2i/d_{\text{model}}})$$

Multi-Head Attention



But before that what is attention
and why do we need it?

See video

Psychology lessons: Selective Attention

The act of focusing on a particular object for a period of time while simultaneously ignoring irrelevant information that is also occurring



Spotlight Effect!

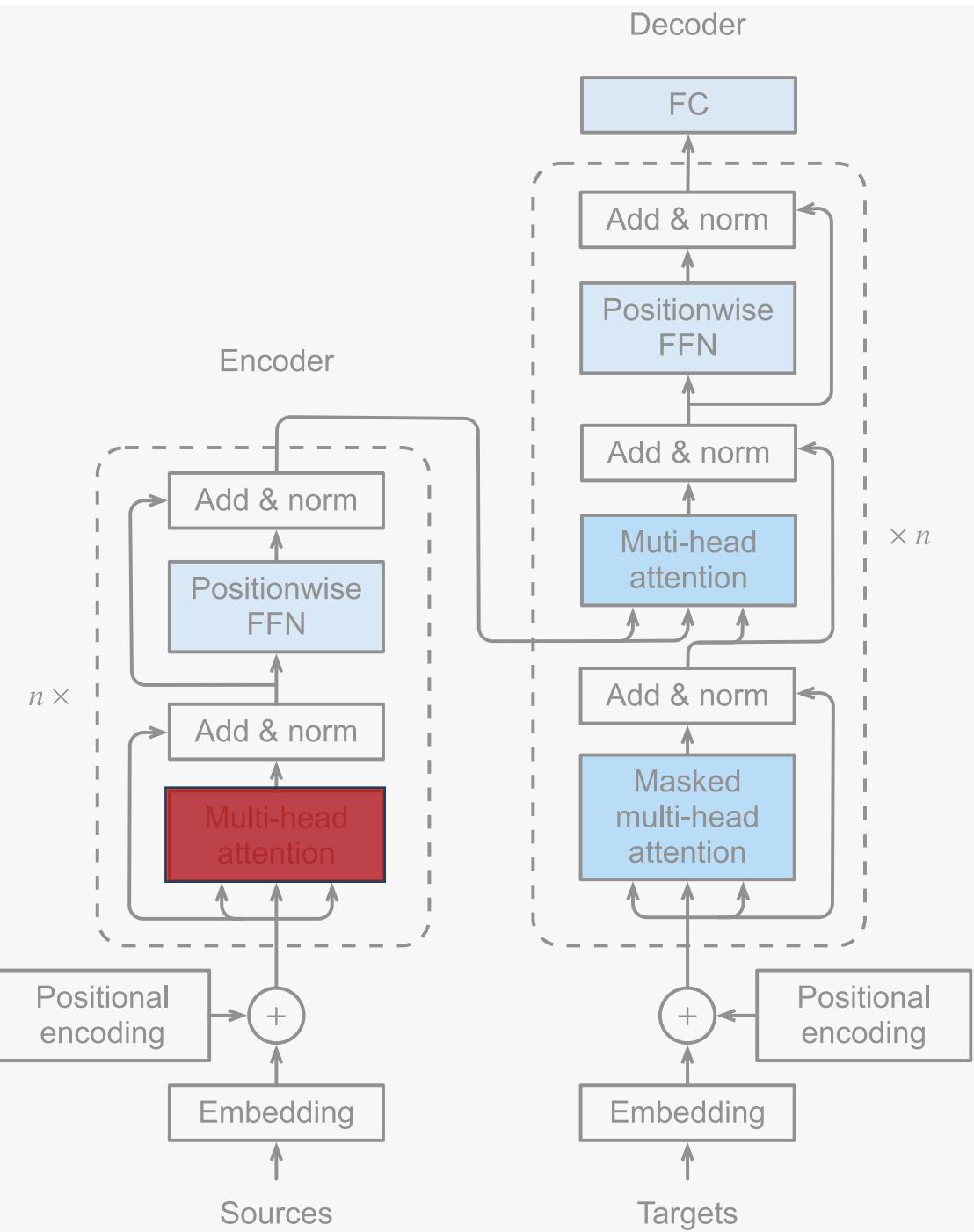


We miss the gorilla running thru while counting, its the spotlight effect/ selective attention

An aerial photograph of a small, flat island in the middle of a vast, clear blue ocean under a bright, cloudy sky. A large, semi-transparent green oval is overlaid on the image, centered over the island. The island has a few small buildings and some vegetation. The water is a vibrant turquoise color.

we focus on the island first

Attention

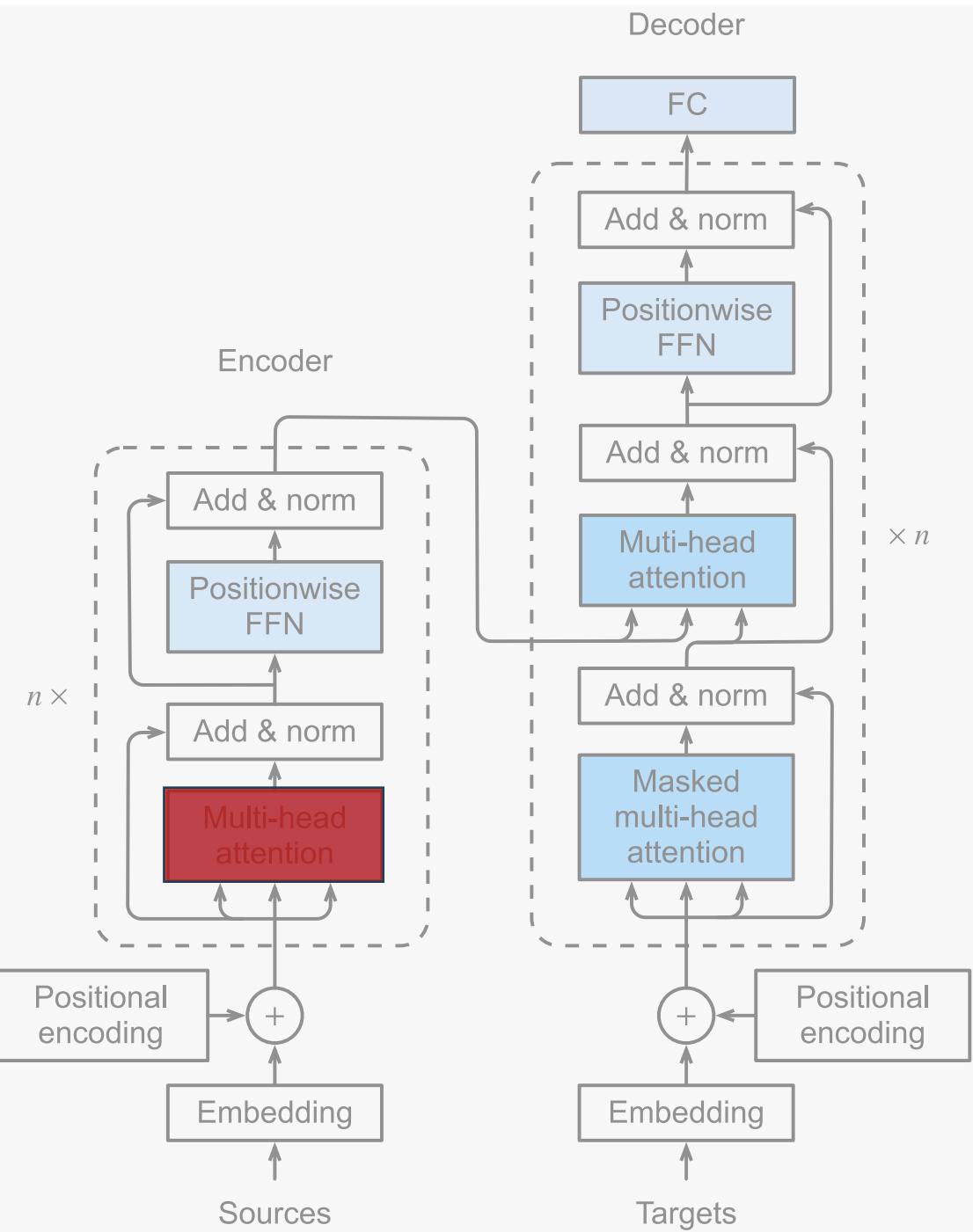


Who is this Harry Potter character?

Now if you two don't mind, I'm going to bed before either of you come up with another clever idea to get us killed - or worse, expelled.

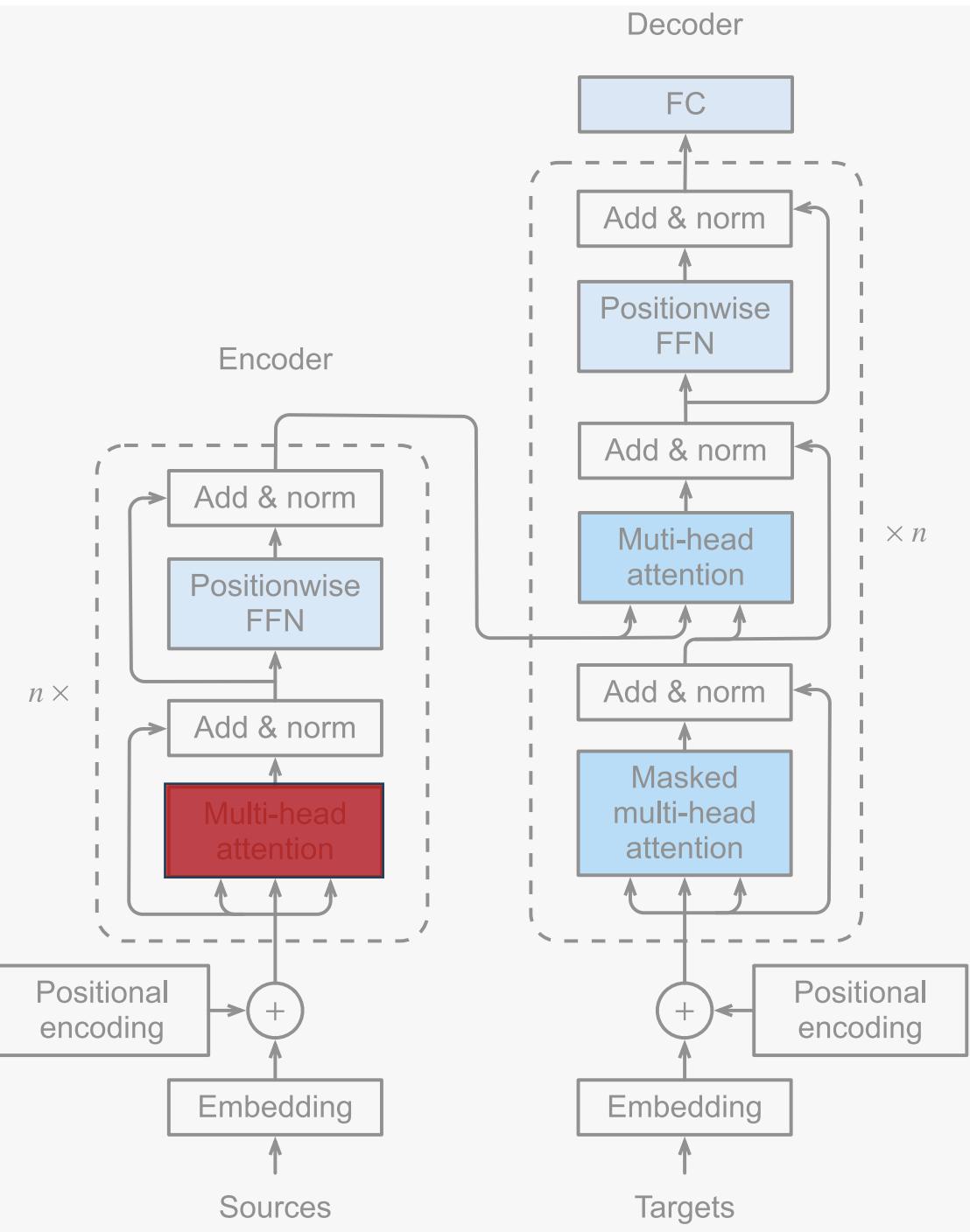
Focus on important words in the sentence

Attention



Who is this Harry Potter character?

Now if you two don't mind, I'm going to bed before either of you come up with another clever idea to get us killed - or worse, expelled.

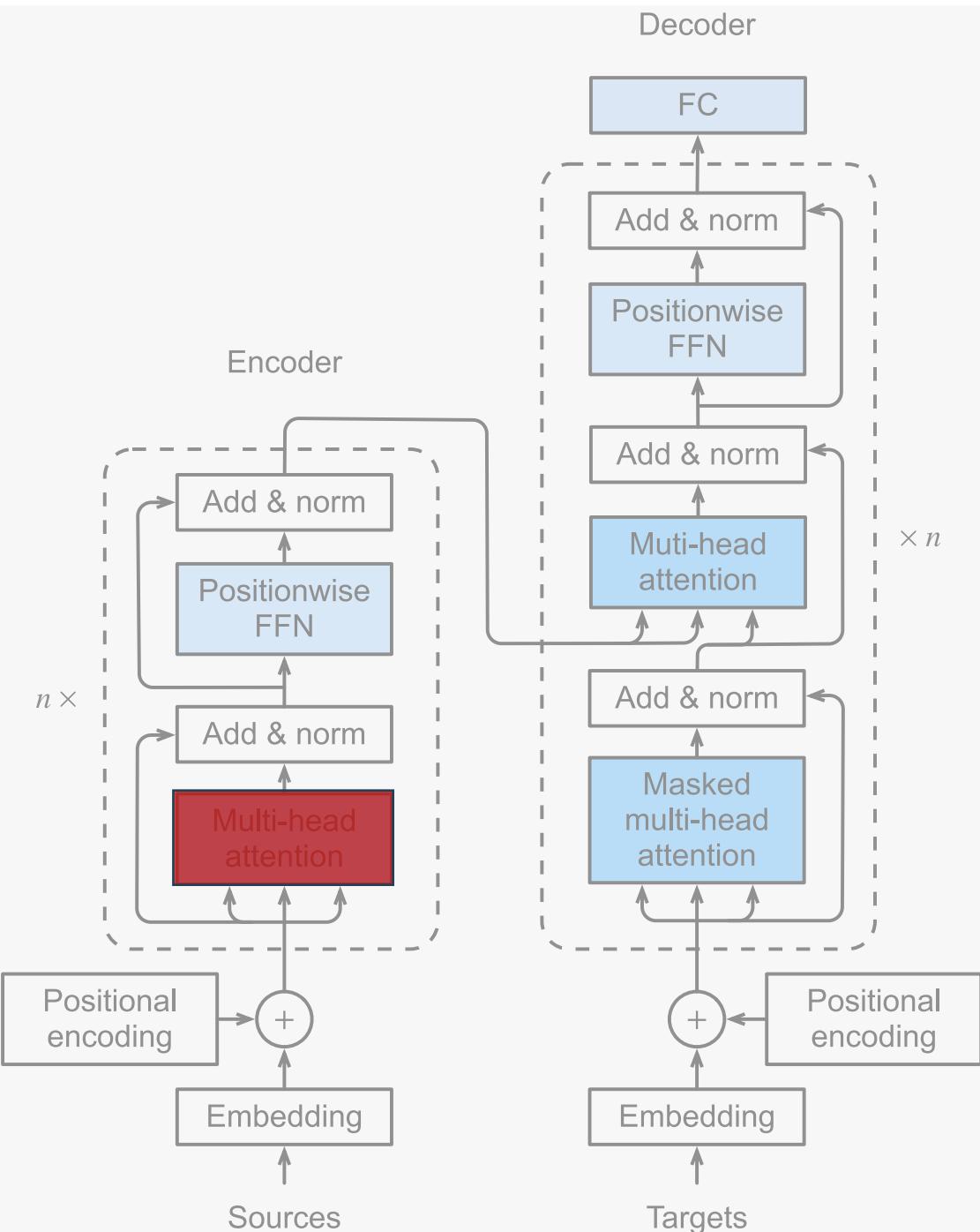


Who is this Harry Potter character?

Now if you two don't mind, I'm going to bed before either of you come up with another clever idea to get us killed - or worse, expelled.

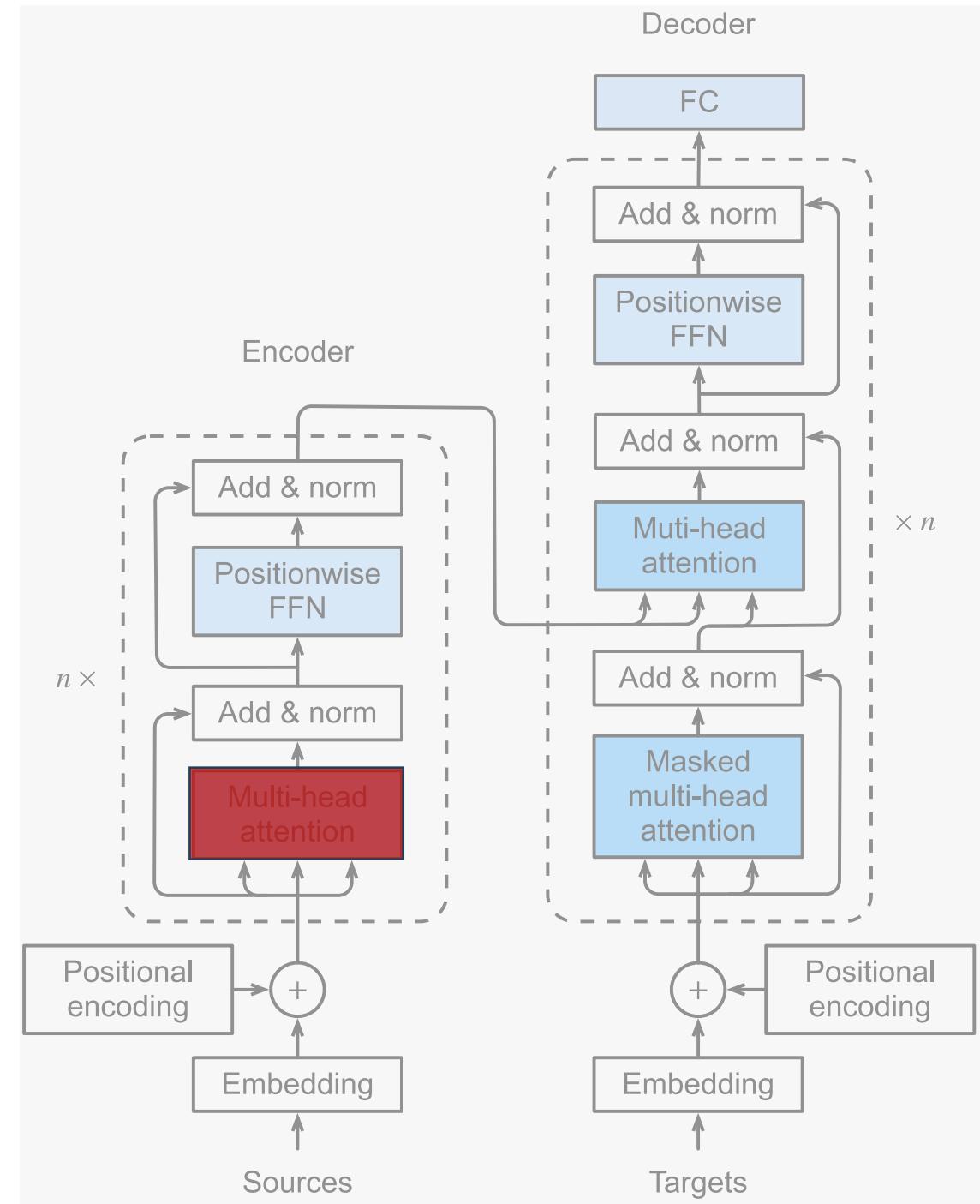
Self-Attention

vs cross attention (later)



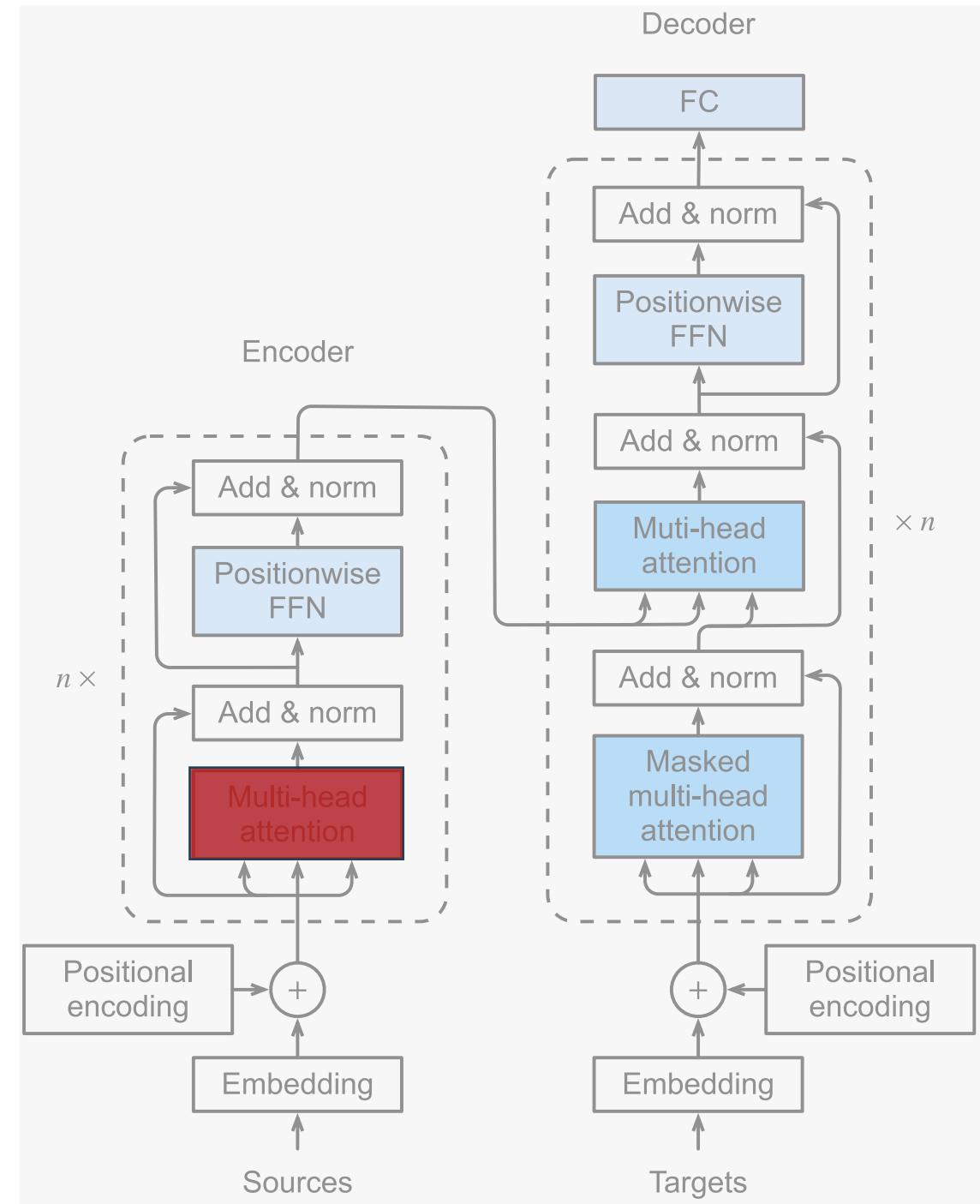
The dog began to **bark** loudly when it saw someone approaching the tree with rough **bark**.

Every context must be understood. bark vs bark, where dog is doing, etc

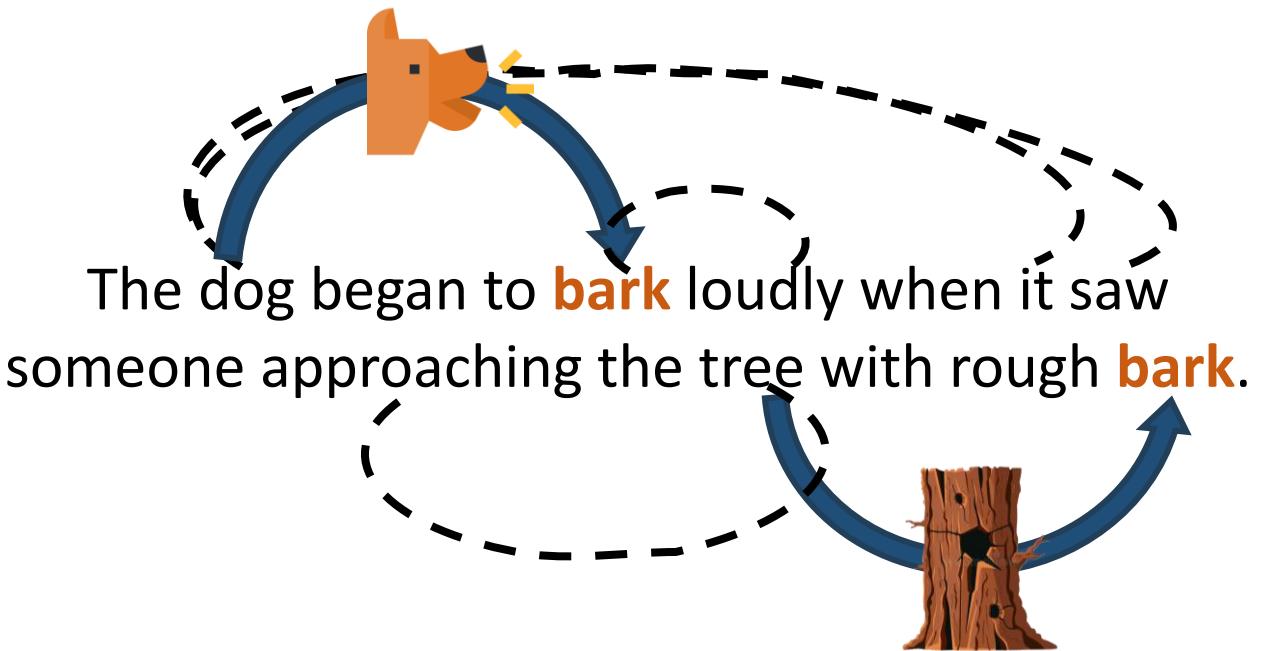


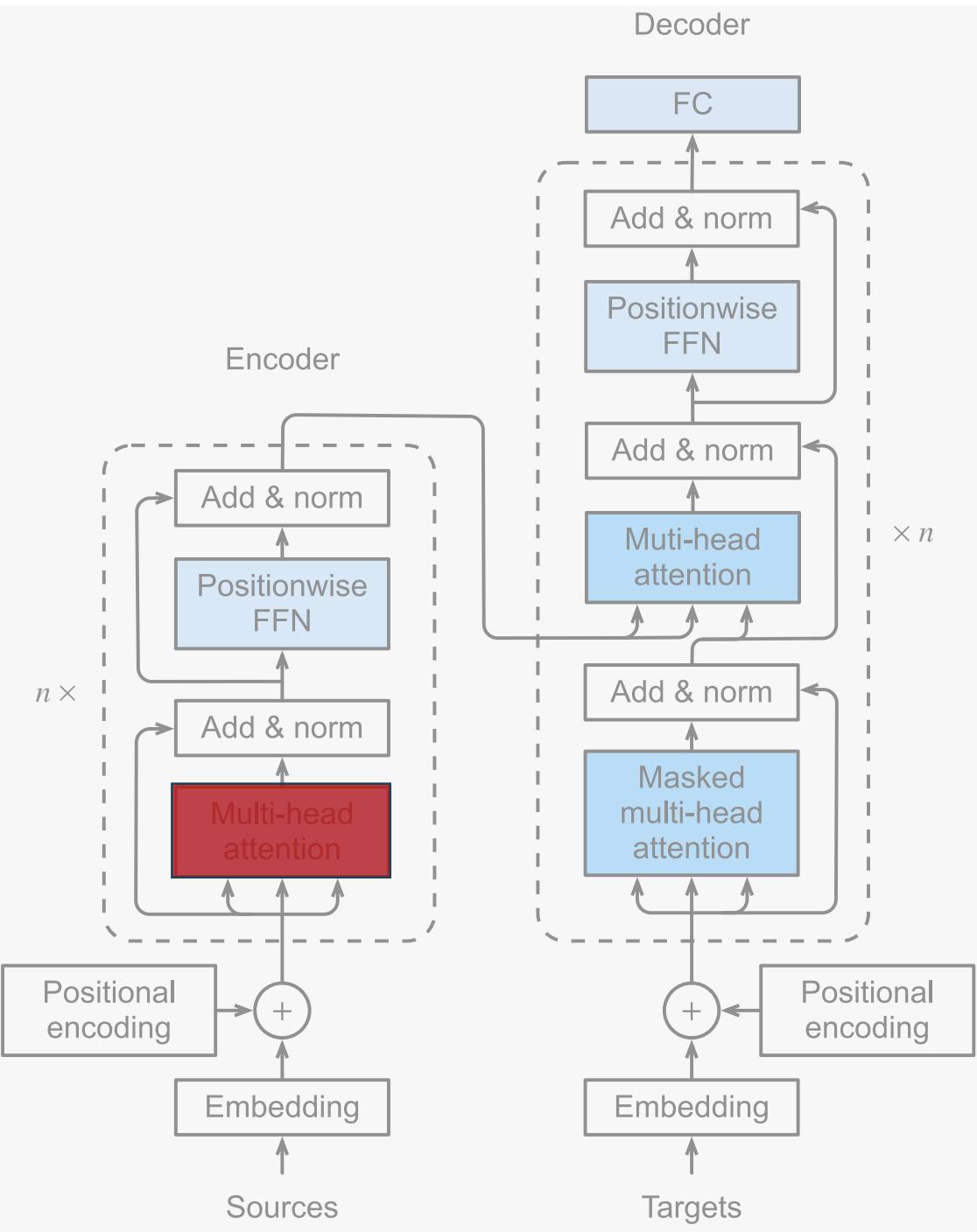
Self-Attention





Self-Attention

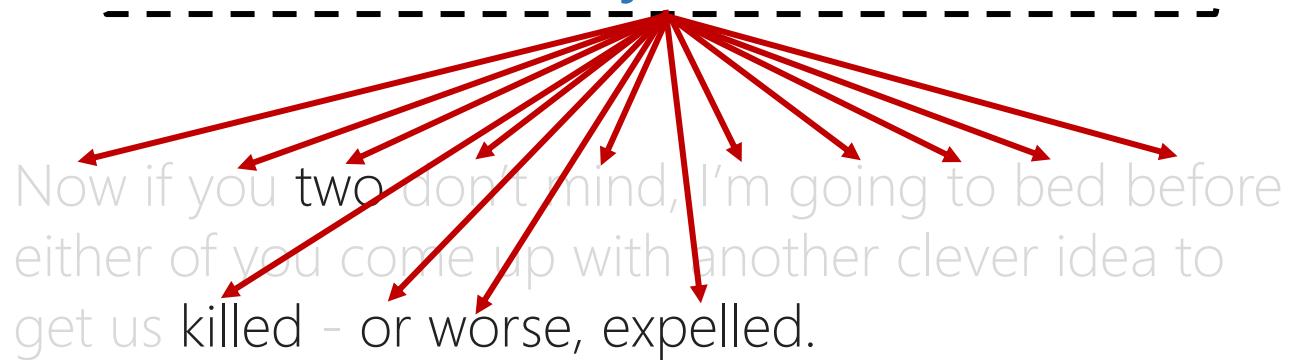




Simple Attention

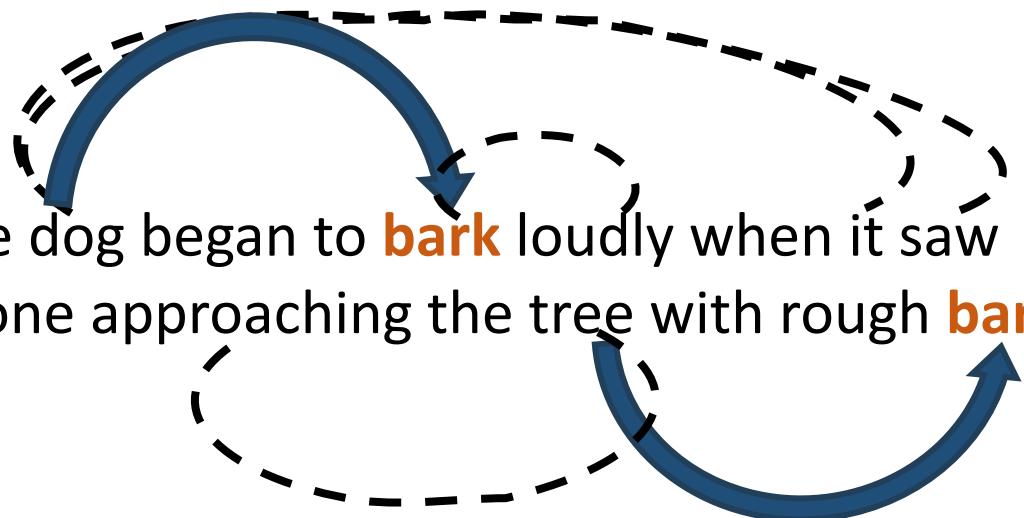
it attends to every word then the few needed to answer the question

Who is this Harry Potter character?

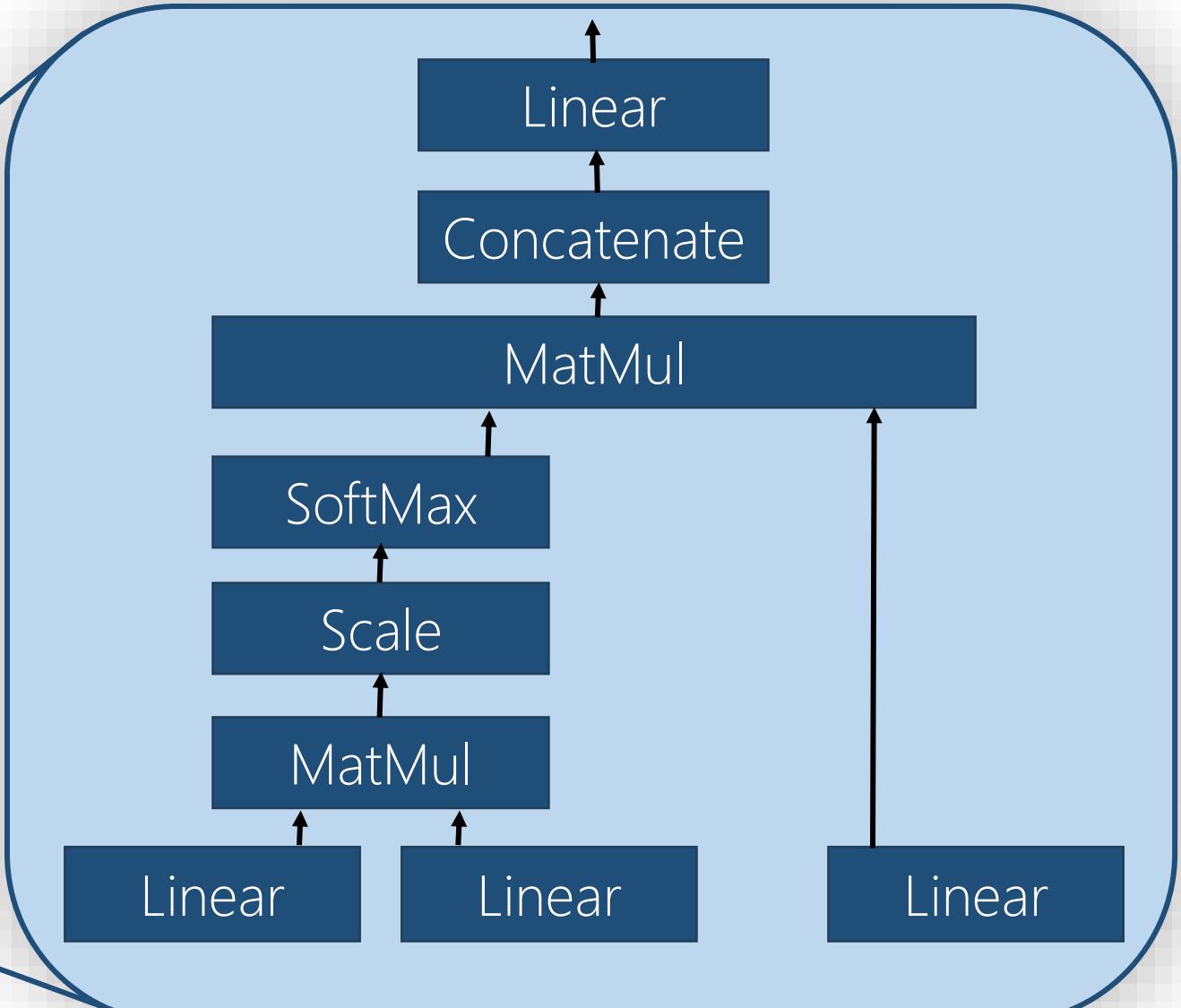
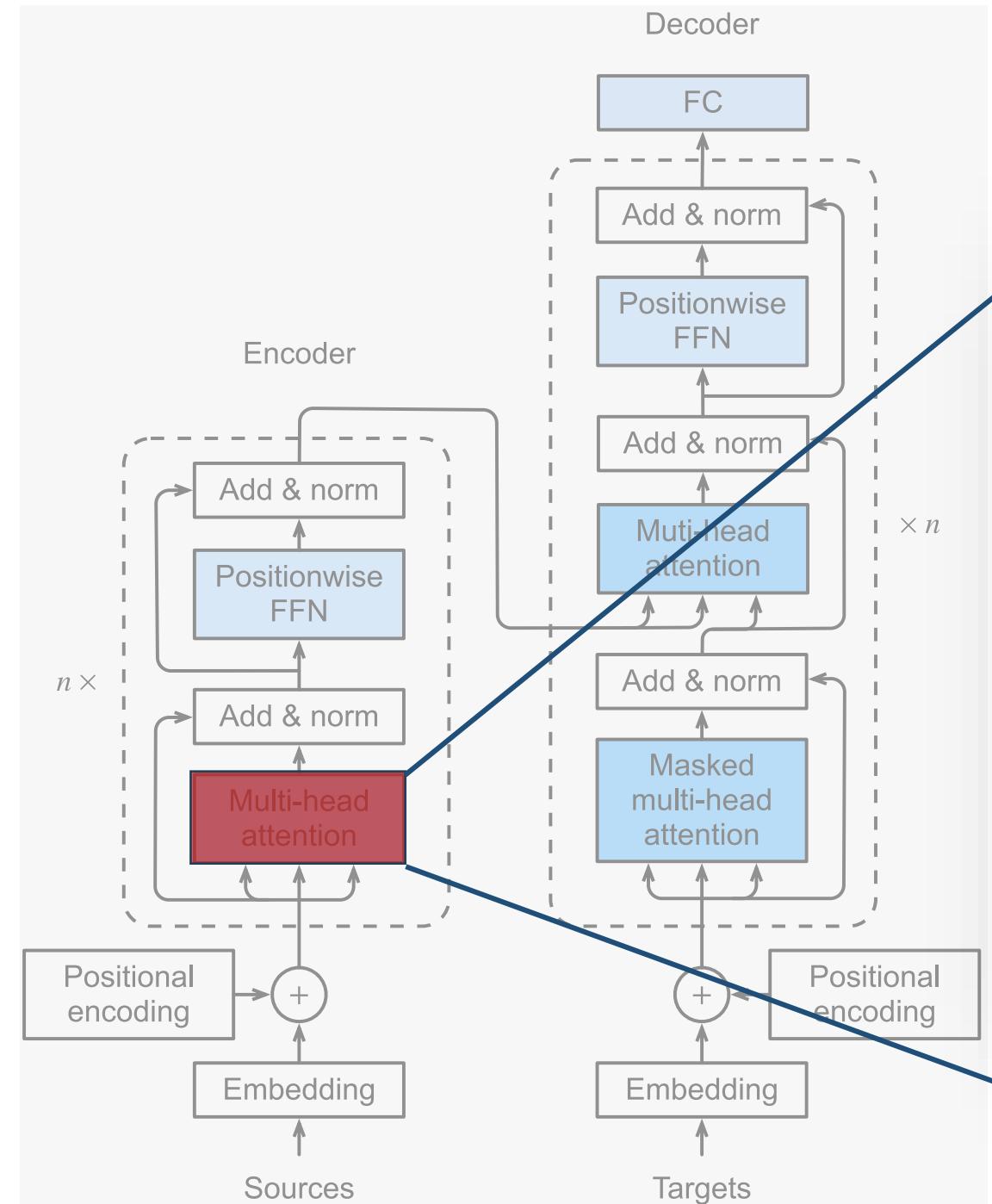


Self-Attention

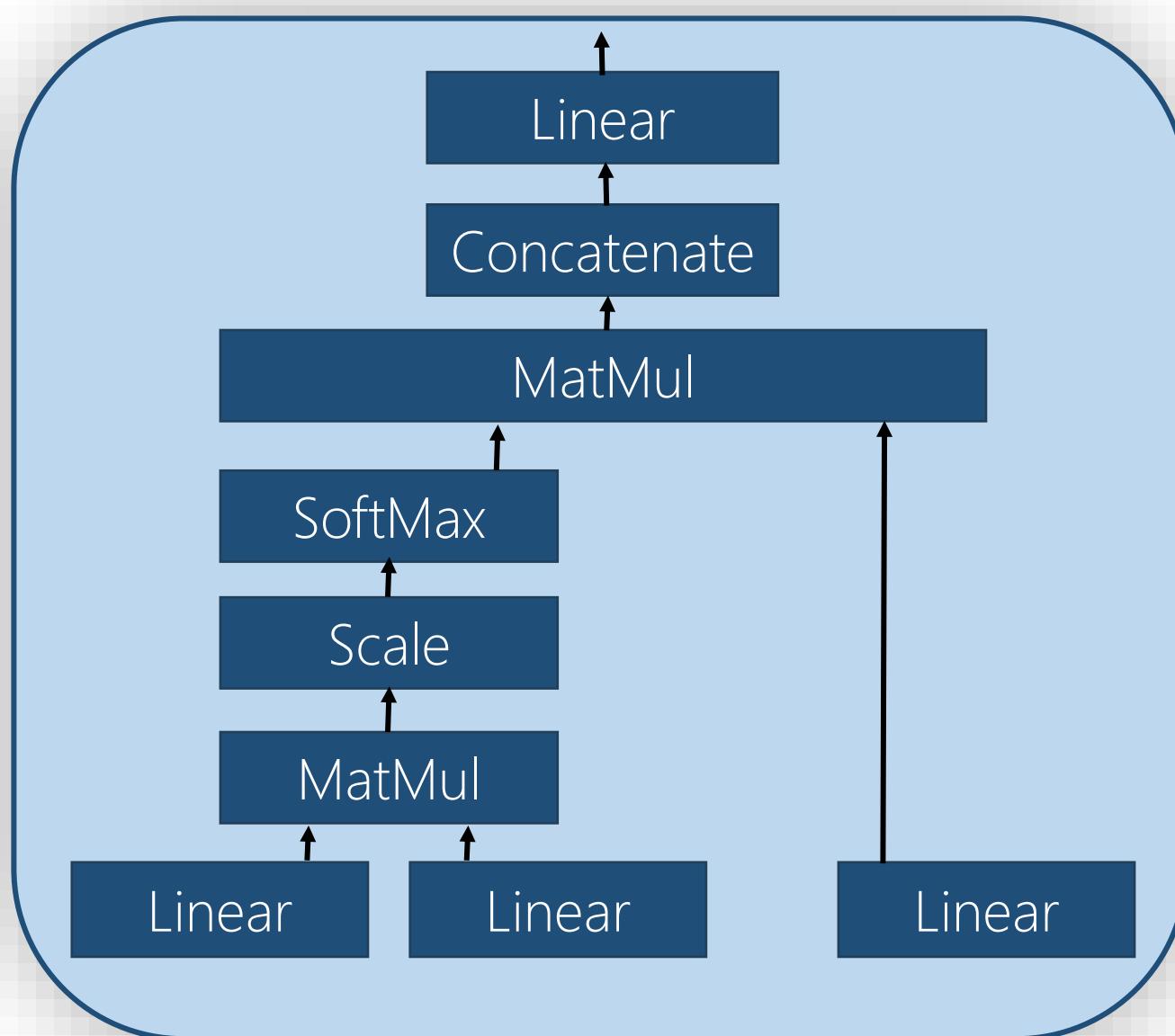
The dog began to bark loudly when it saw someone approaching the tree with rough bark.



Multi-Head Attention



Multi-Head Attention



linear transformation

concat diff vectors

dot product

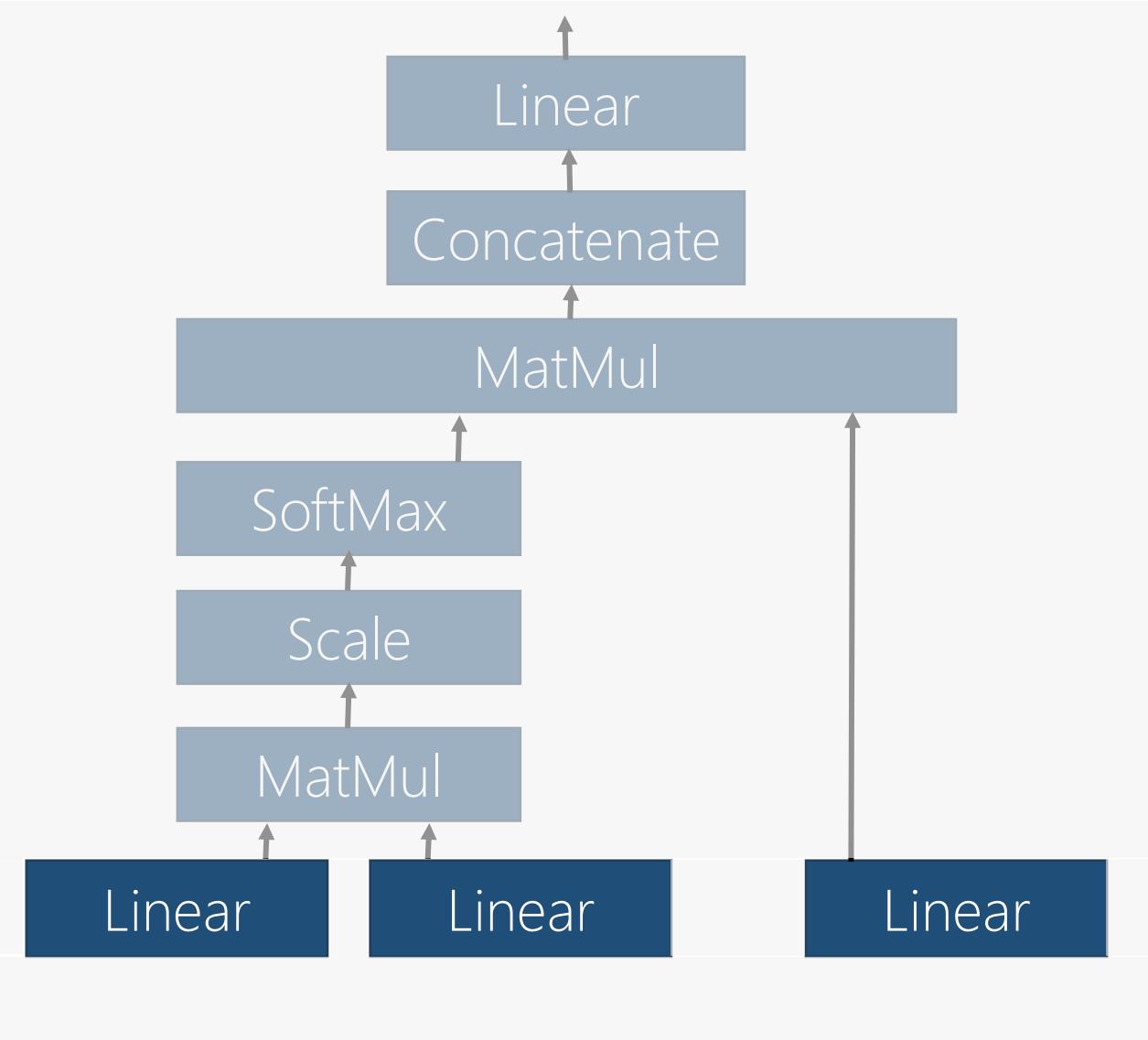
normalizing to 0-1

scaling

mat mul is the dot product

linear transform of the input

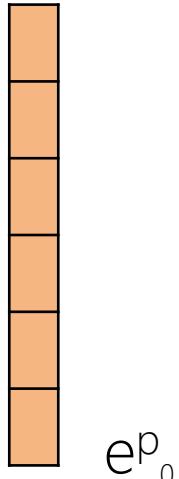
Multi-Head Attention



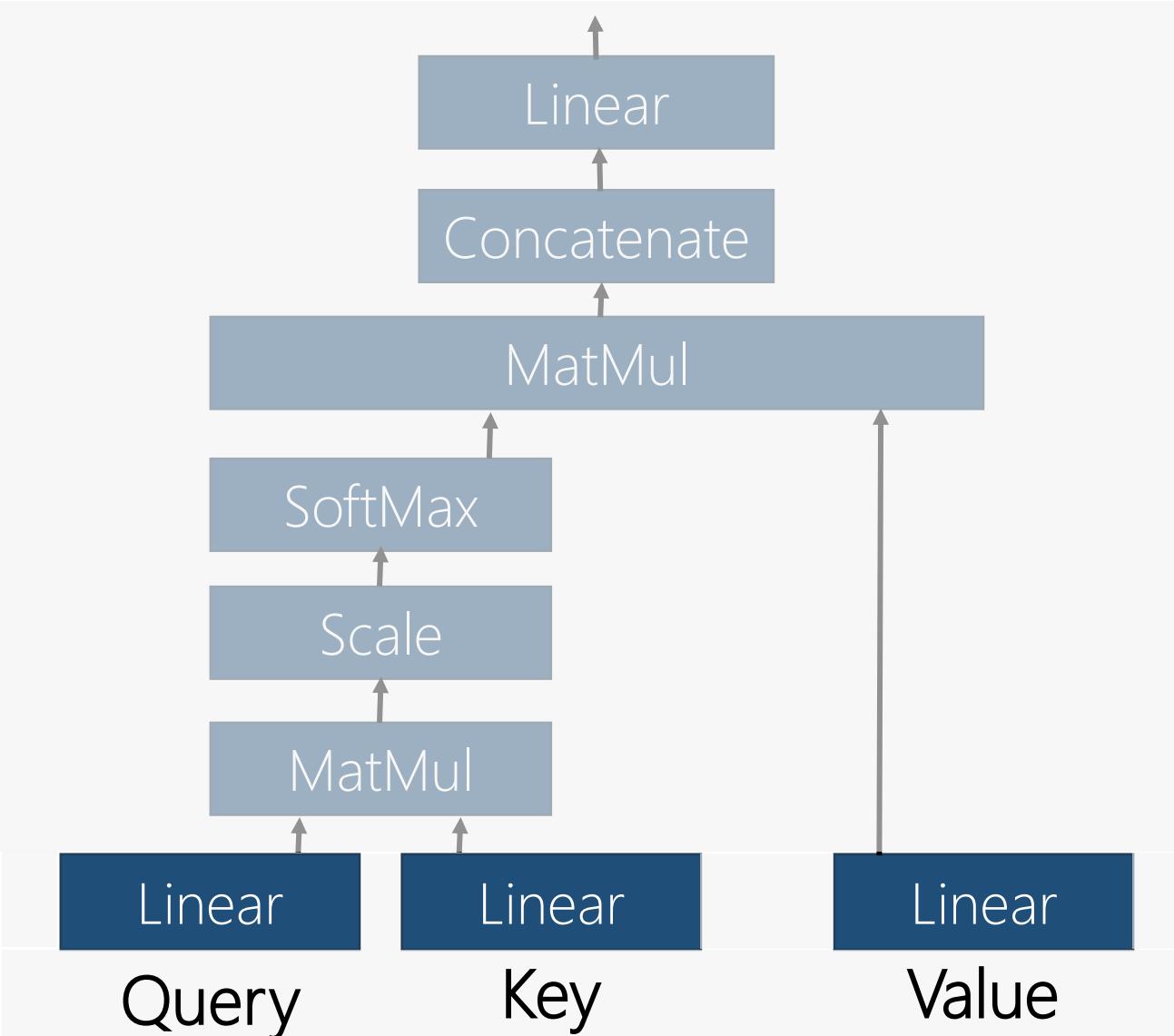
What does a linear layer do?

Input dims could be 128, 256 etc
It tries to breakdown to smaller dimension

Linear



Multi-Head Attention



Analogous to Data Retrieval/Search

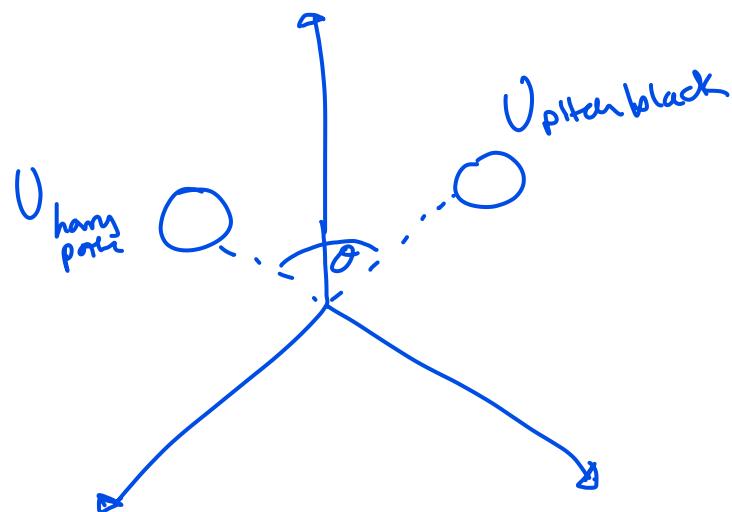
Query, key, value-foundational elements of attention

Query: info looking for
Key: where in database (of what the video is about in youtube for ex)
Value: what info can I share

inputting a 7×4 matrix for the query and key: 7= # words, 4 is the dim of the words

Similarity between Query and Key?

input the image I took



- $\cos(\theta) = 1$

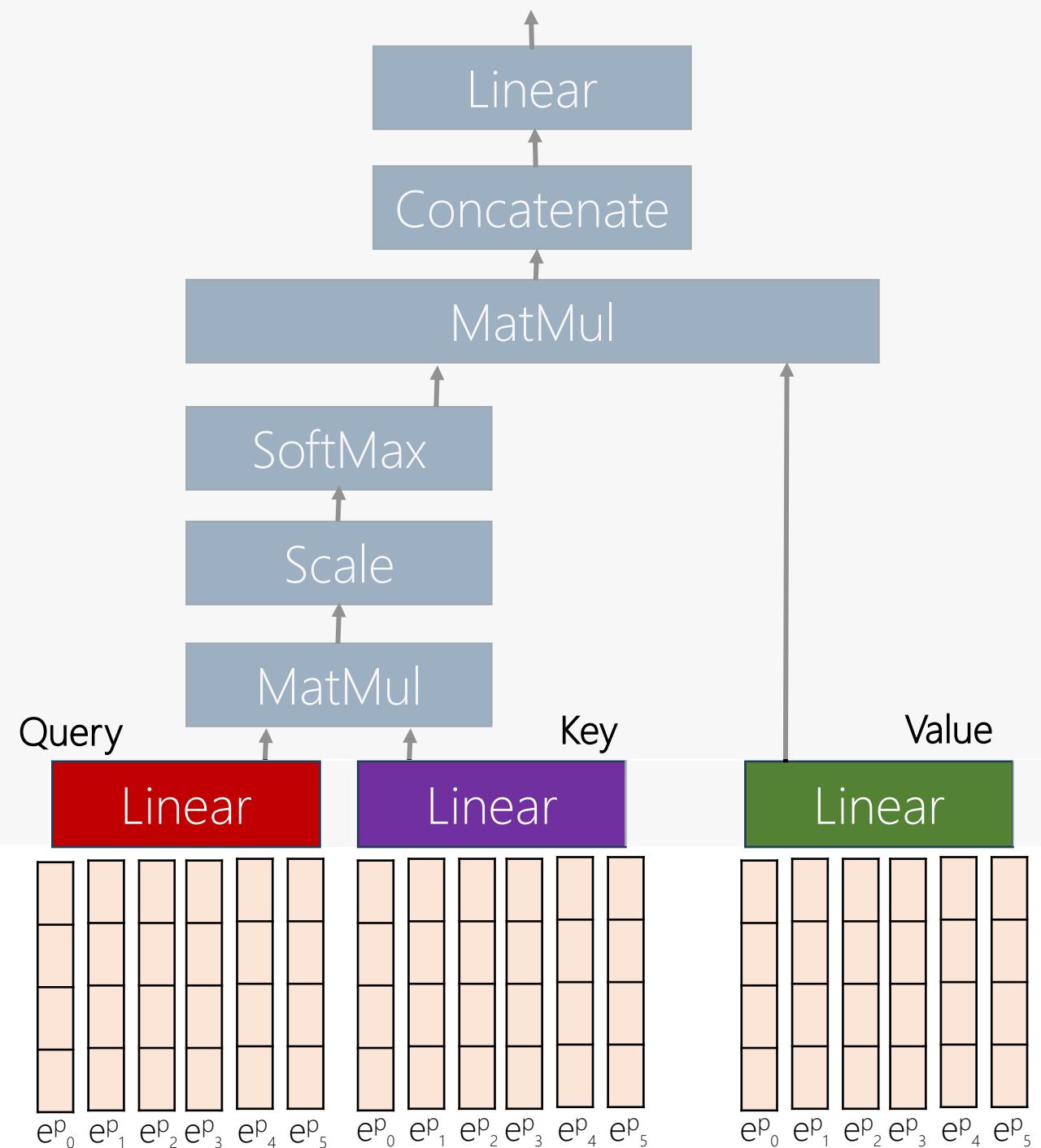
$\theta = 0 \rightarrow A$ and B point in same direction

- $\cos(\theta) = -1$

$\theta = 180 \rightarrow A$ and b pt in opp directions

- $\cos(\theta) = 0$

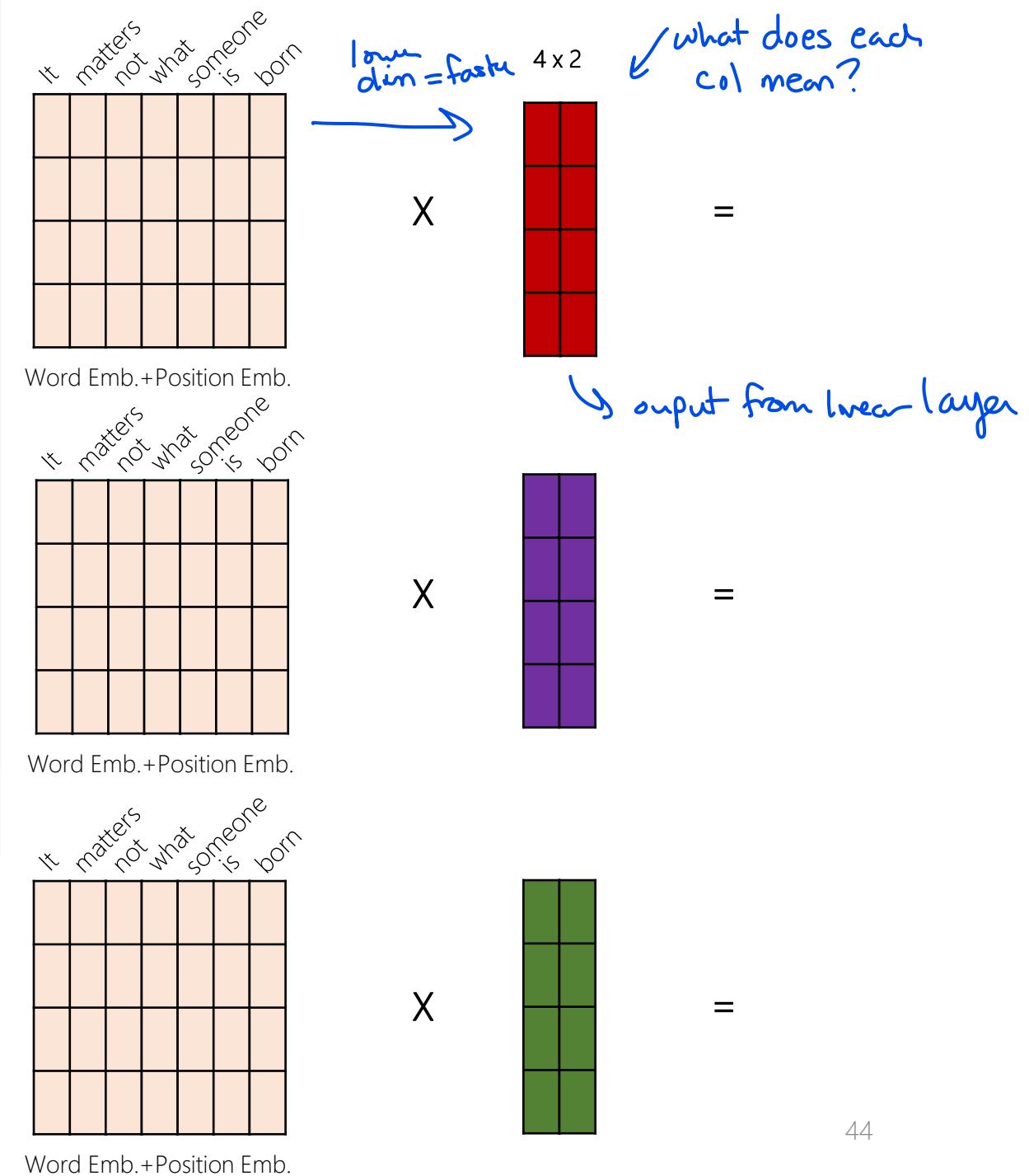
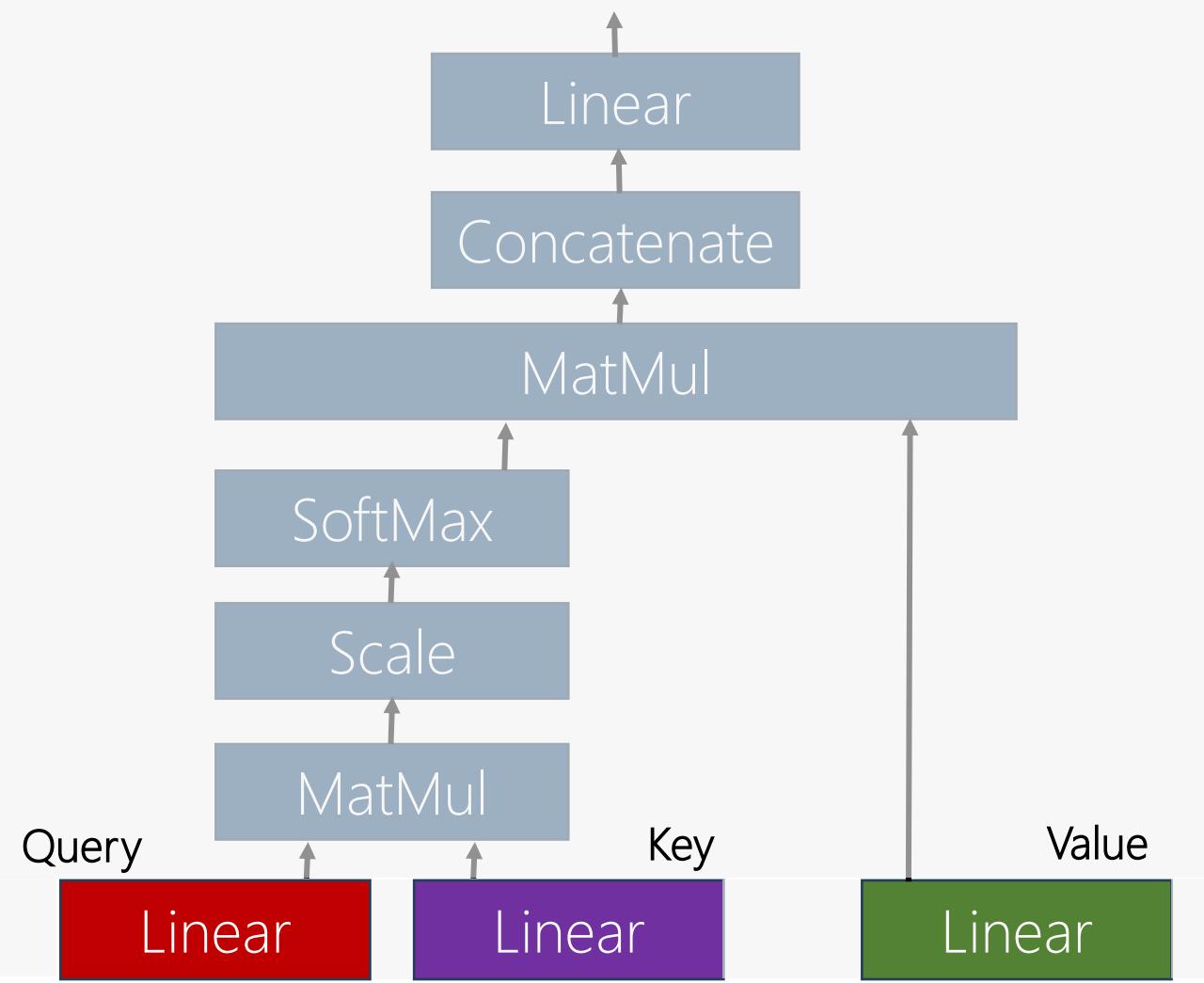
$\theta = 90 \rightarrow A$ and B $\hat{=}$ 90°
(Orthogonal)

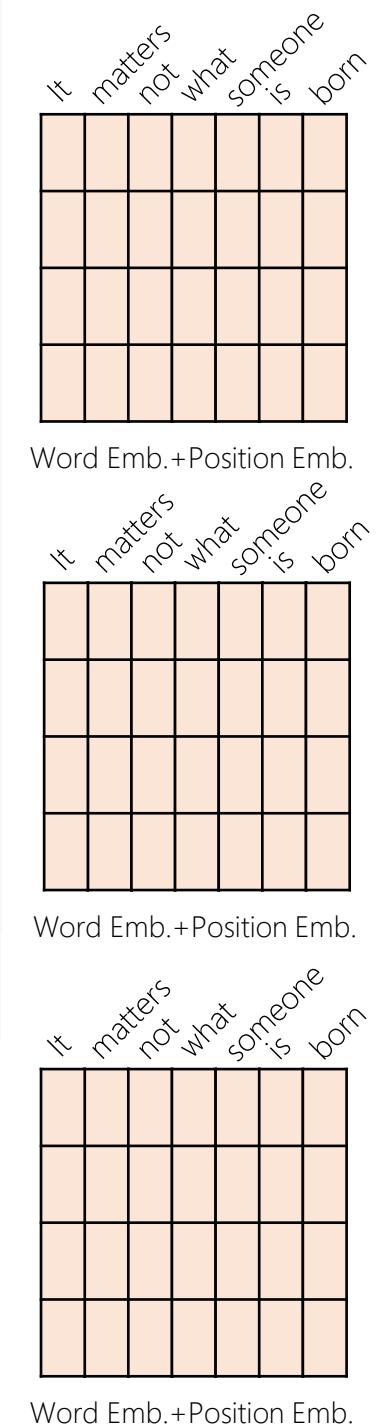
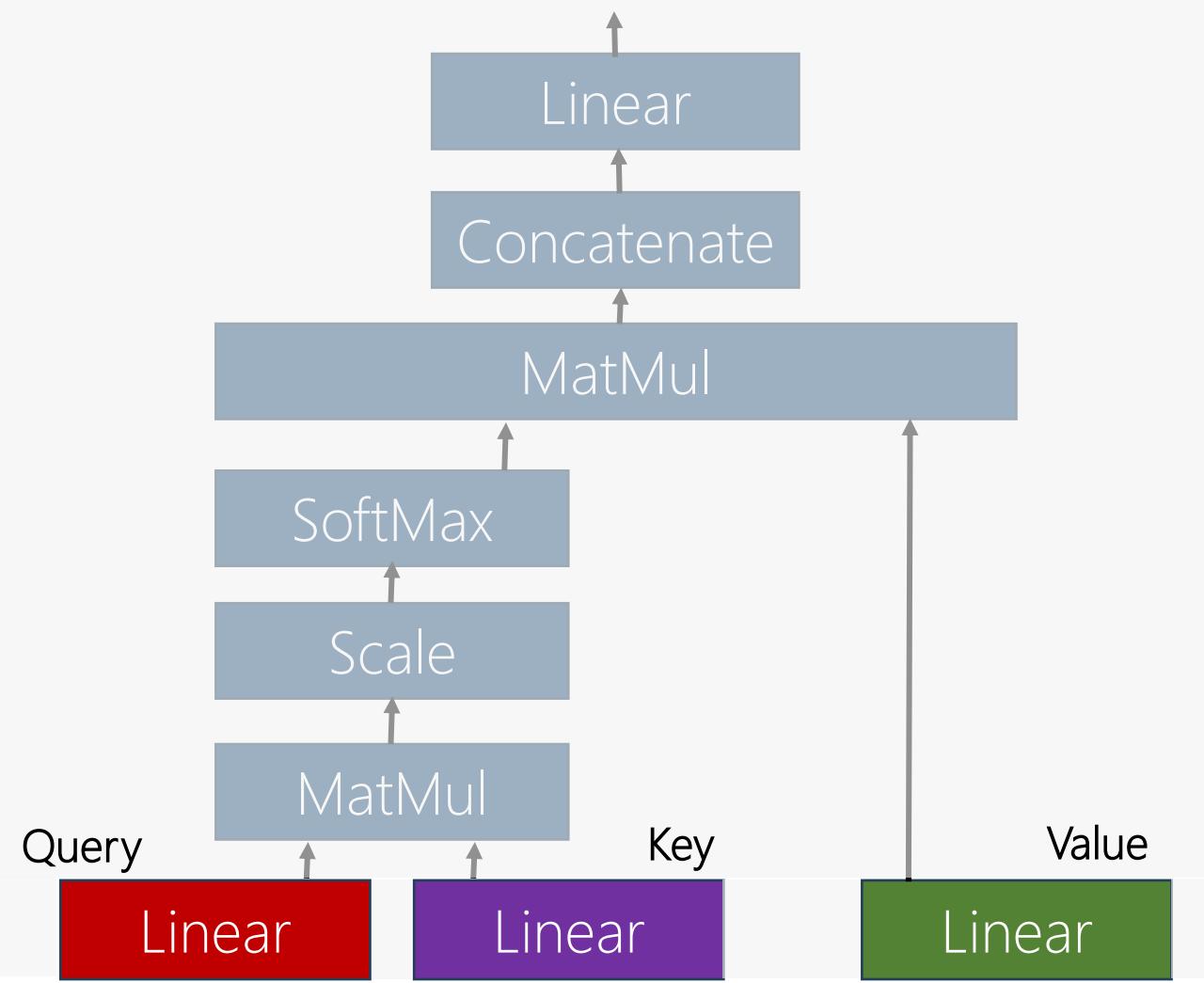


Multi-Head Attention

Analogous to Data Retrieval/Search

Not sending one word at a time

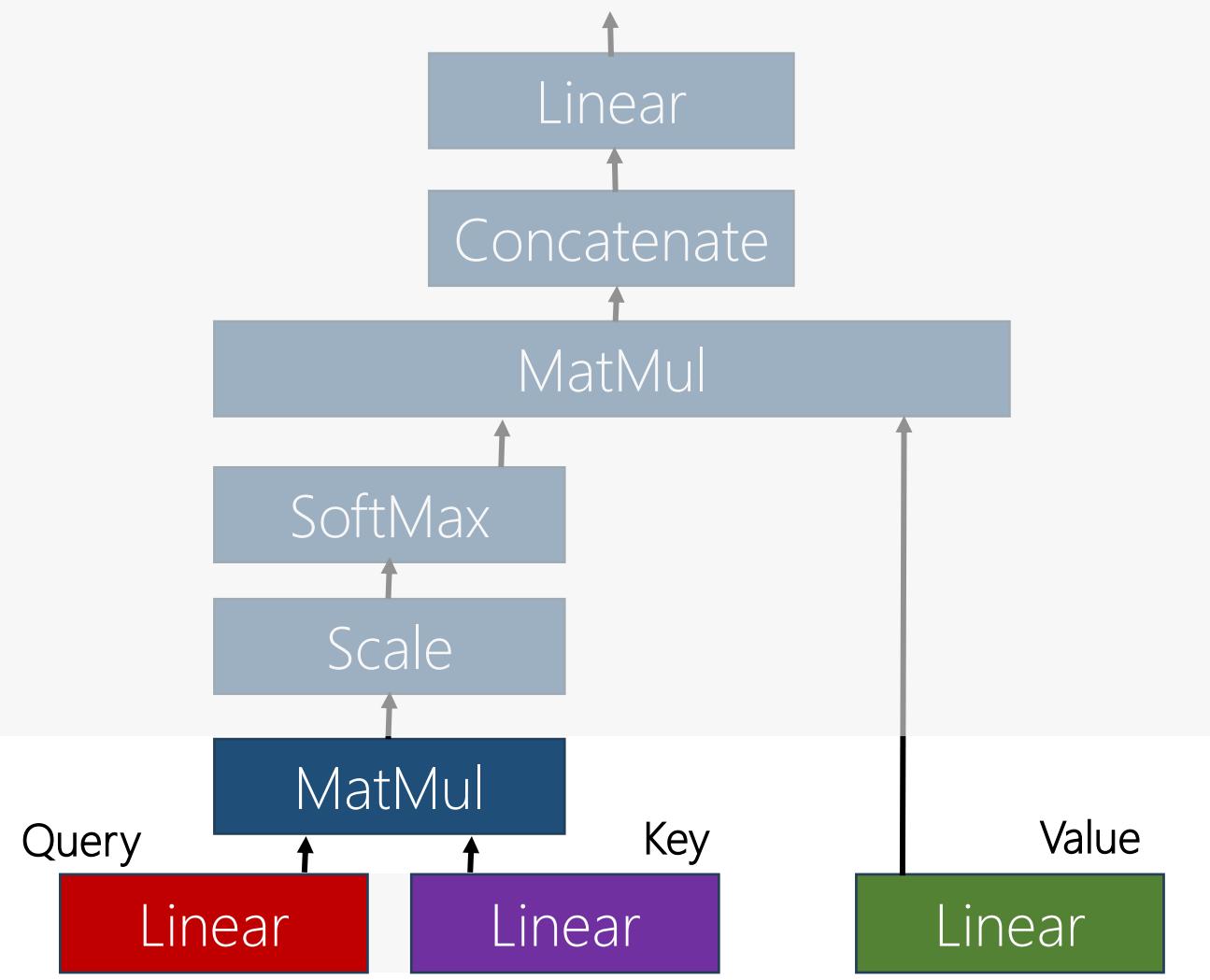




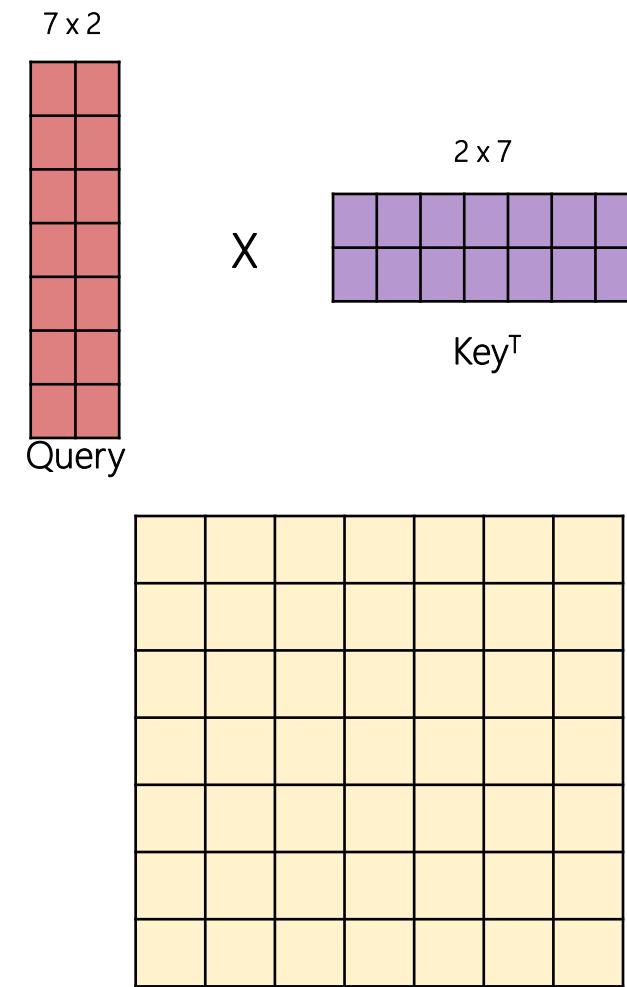
$$\begin{array}{c}
 4 \times 2 \\
 \times \\
 = \\
 \text{Query} \\
 \text{Key} \\
 \text{Value}
 \end{array}$$

The diagram shows the computation of Query, Key, and Value matrices from Word Embedding + Position Embedding matrices. The input matrix is labeled "Word Emb.+Position Emb." and contains the words "It", "matters", "not", "what", "someone", and "is", "born". The matrix is 6 rows by 8 columns. The output consists of three separate matrices:

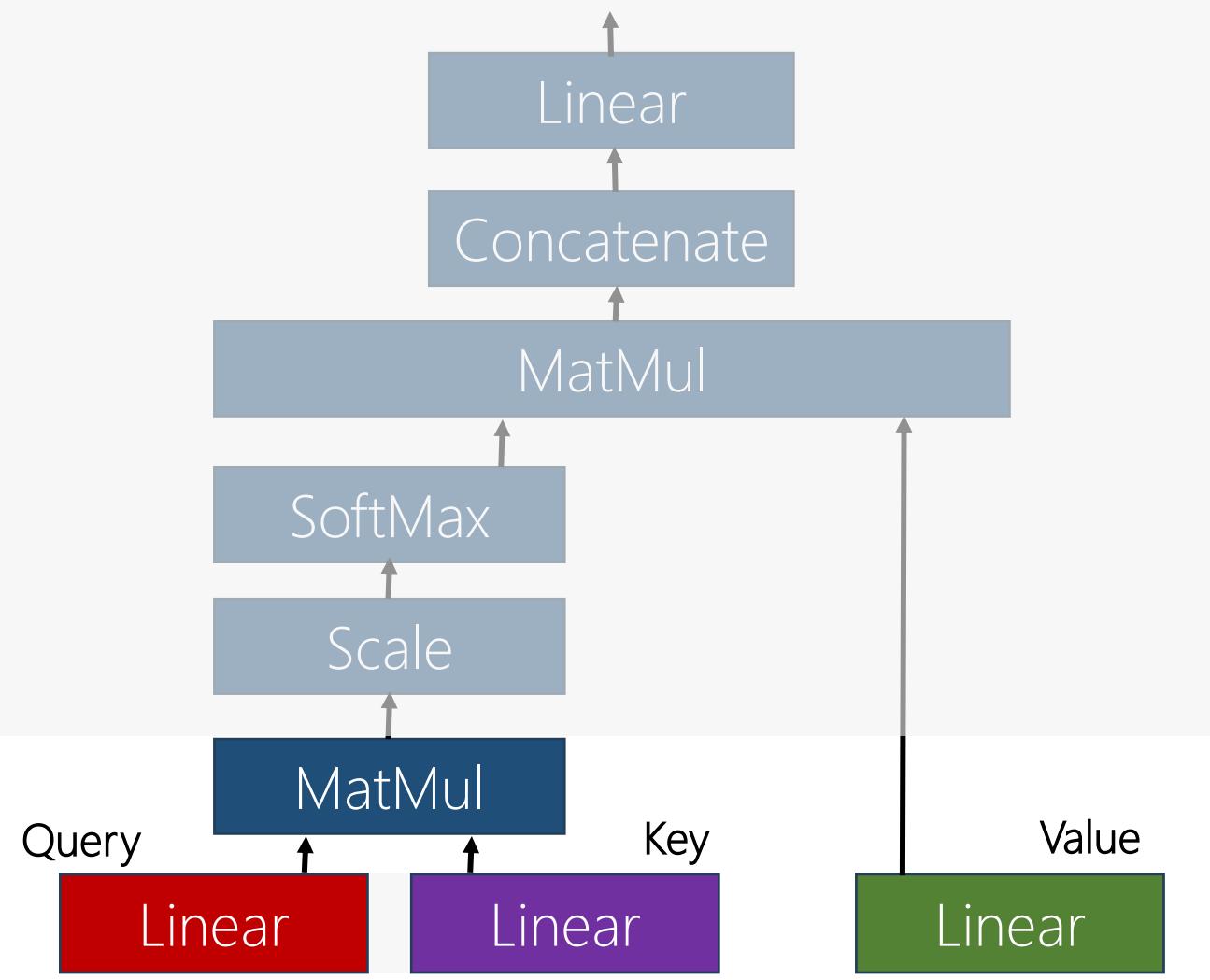
- Query**: A 6x2 matrix where each row is a 2x4 matrix from the input. The first column is red and the second column is blue.
- Key**: A 6x3 matrix where each row is a 3x3 matrix from the input. The first column is red, the second is blue, and the third is green.
- Value**: A 6x4 matrix where each row is a 4x4 matrix from the input. The first column is red, the second is blue, the third is green, and the fourth is orange.



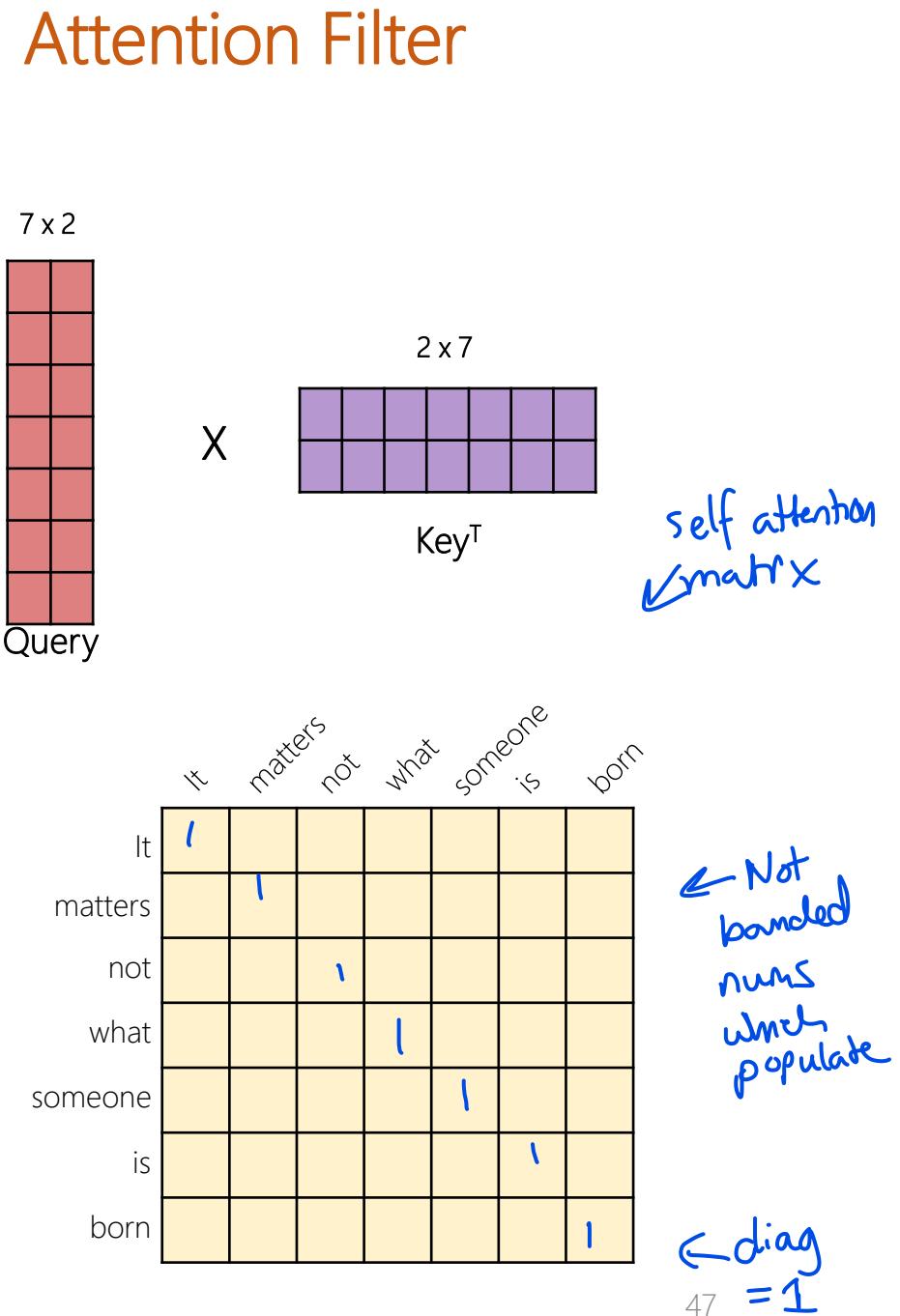
Attention Filter

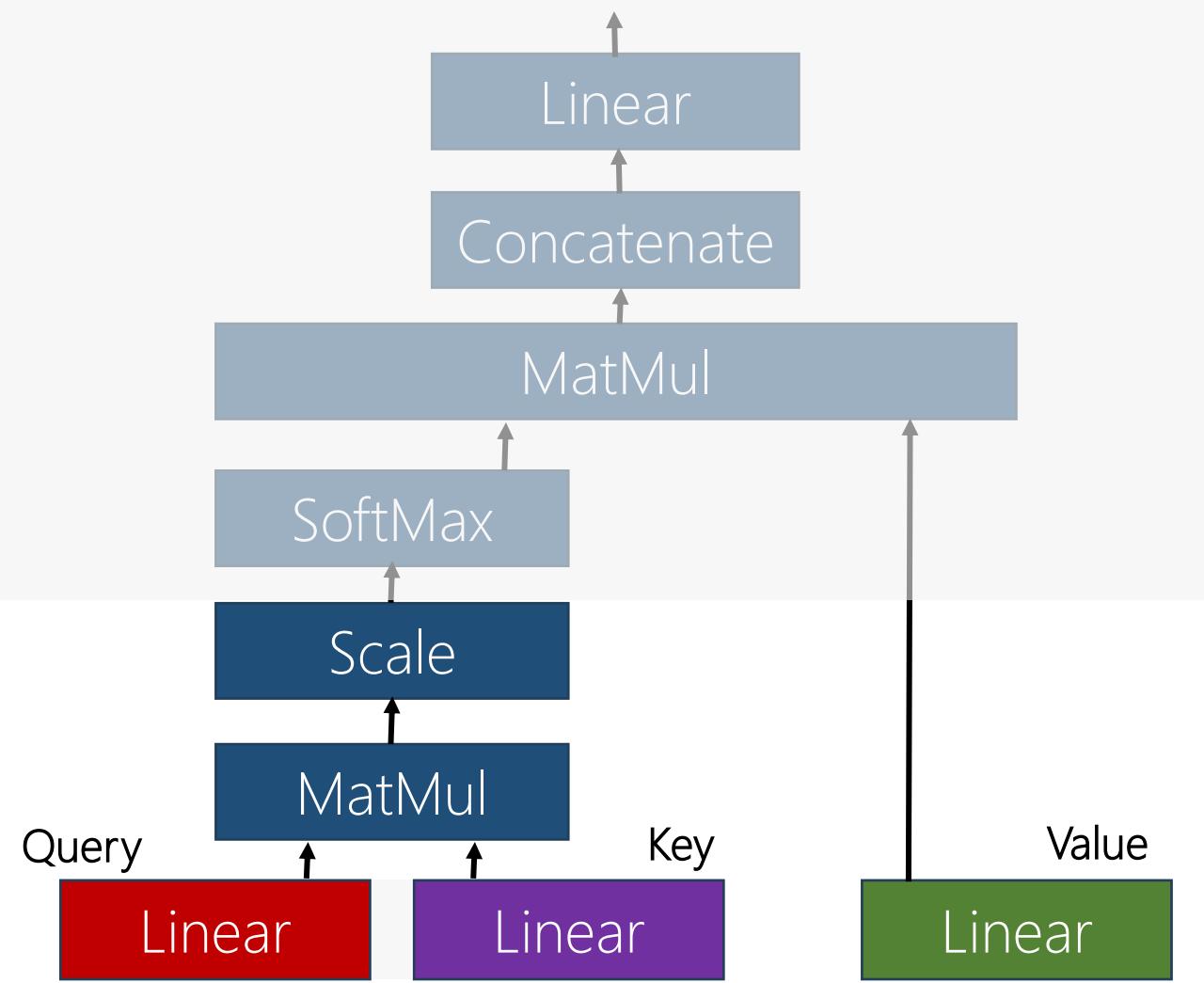


← self
attention
matrix



Q: What do
the vals
rep?



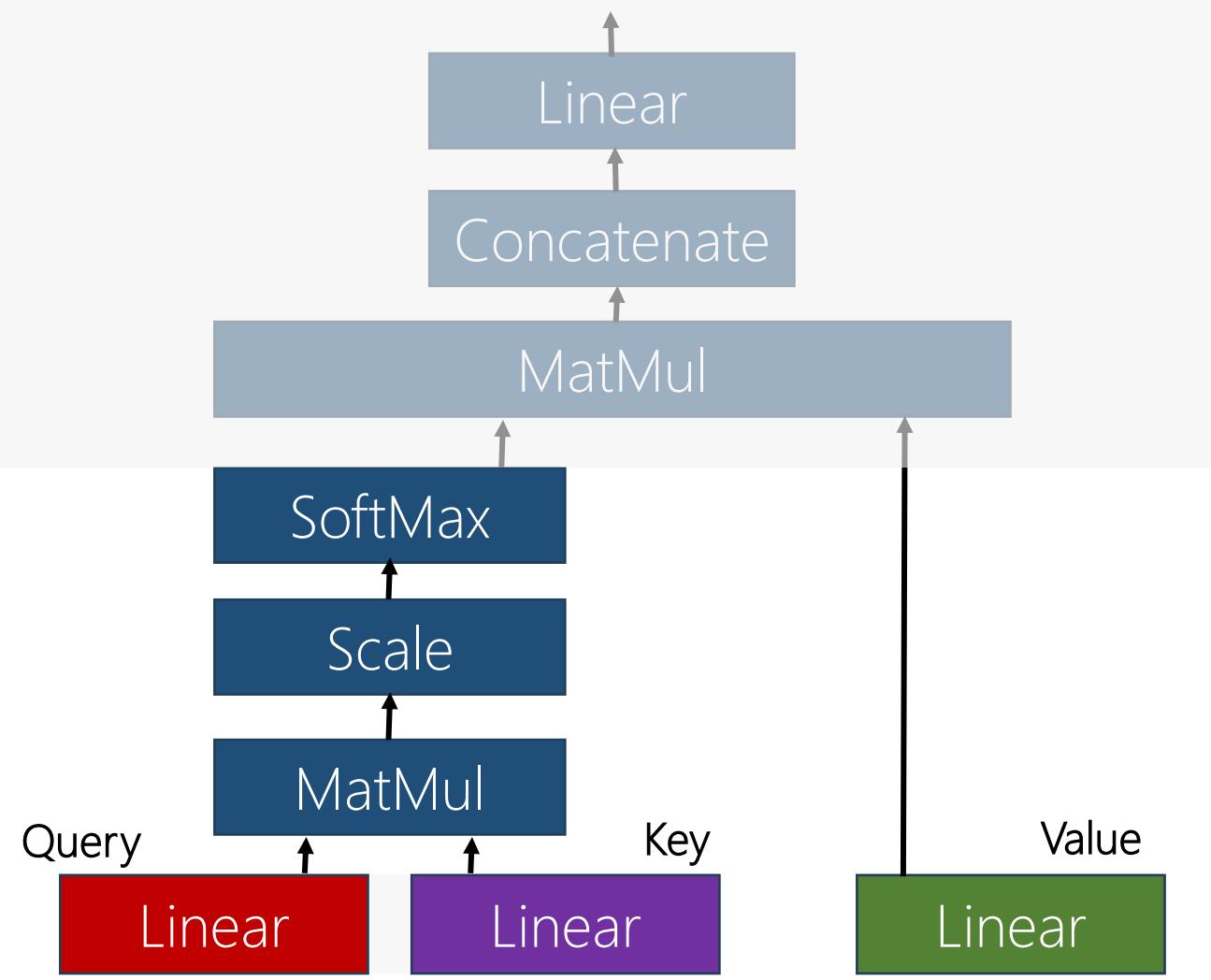


Attention Filter

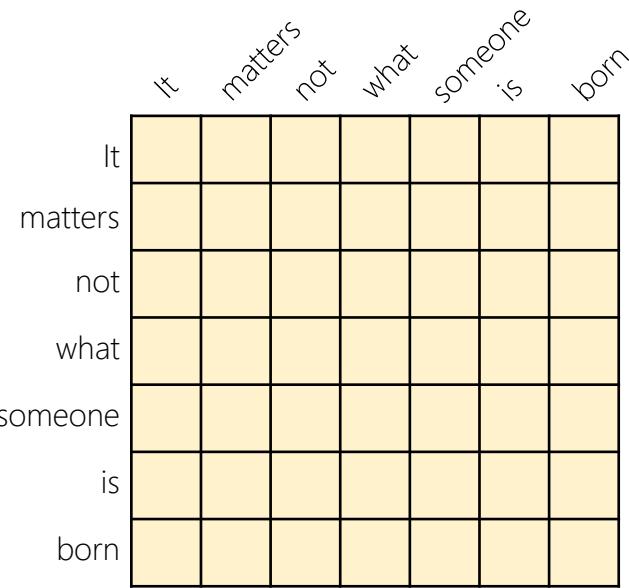
	It	matters	not	what	someone	is	born
It	1.0	0.0	0.0	0.0	0.0	0.0	0.0
matters	0.0	1.0	0.0	0.0	0.0	0.0	0.0
not	0.0	0.0	1.0	0.0	0.0	0.0	0.0
what	0.0	0.0	0.0	1.0	0.0	0.0	0.0
someone	0.0	0.0	0.0	0.0	1.0	0.0	0.0
is	0.0	0.0	0.0	0.0	0.0	1.0	0.0
born	0.0	0.0	0.0	0.0	0.0	0.0	1.0

Scale using dimension! *of the vector*

(still not in a closed space of $0 \rightarrow 1$)



Attention Filter

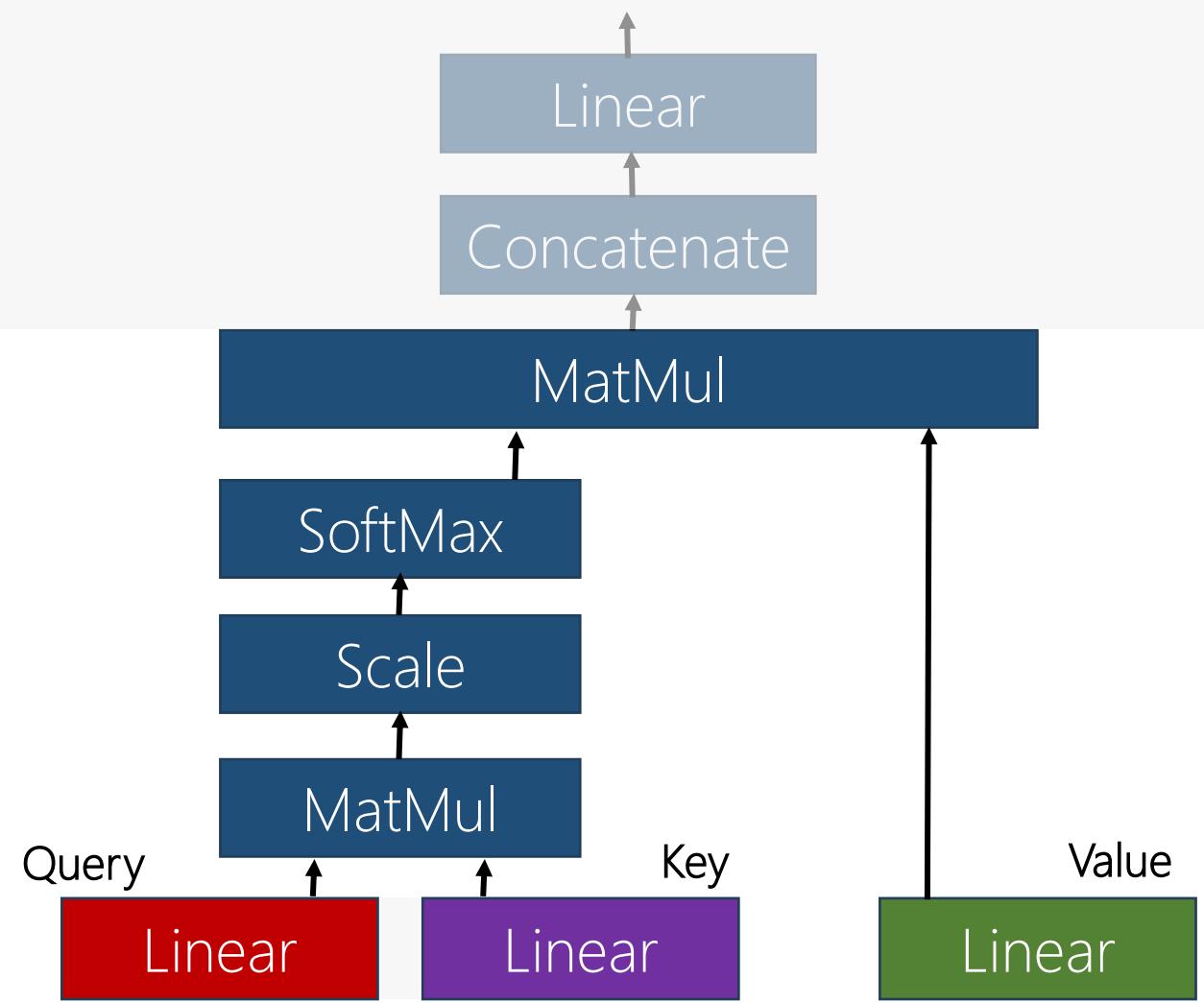


Normalize using SoftMax!

To bounded vals of 0 → 1

1 = high attention (high similarity)

0 = not attended (not related at all)



This is 1
attention heads

Weighted Value

weights for how close each word
is to each other
 7×7

Normalized Attention Filter

vector representing the words
numerically

X

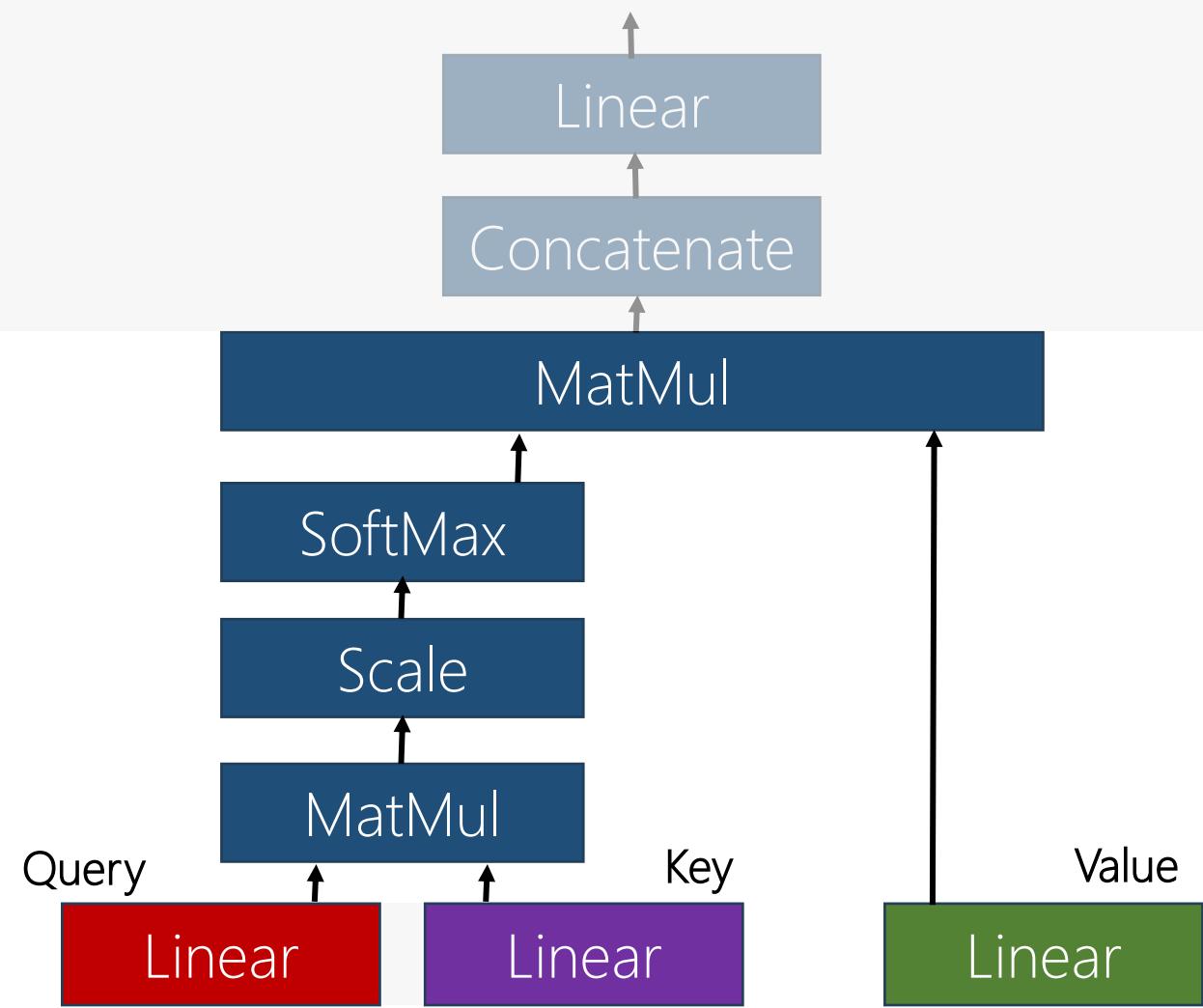
$$\text{Attention}(Q, K, V) = \sigma \left(\frac{Q \cdot K^T}{\sqrt{d_k}} \right) \cdot V$$

Annotations for the attention formula:

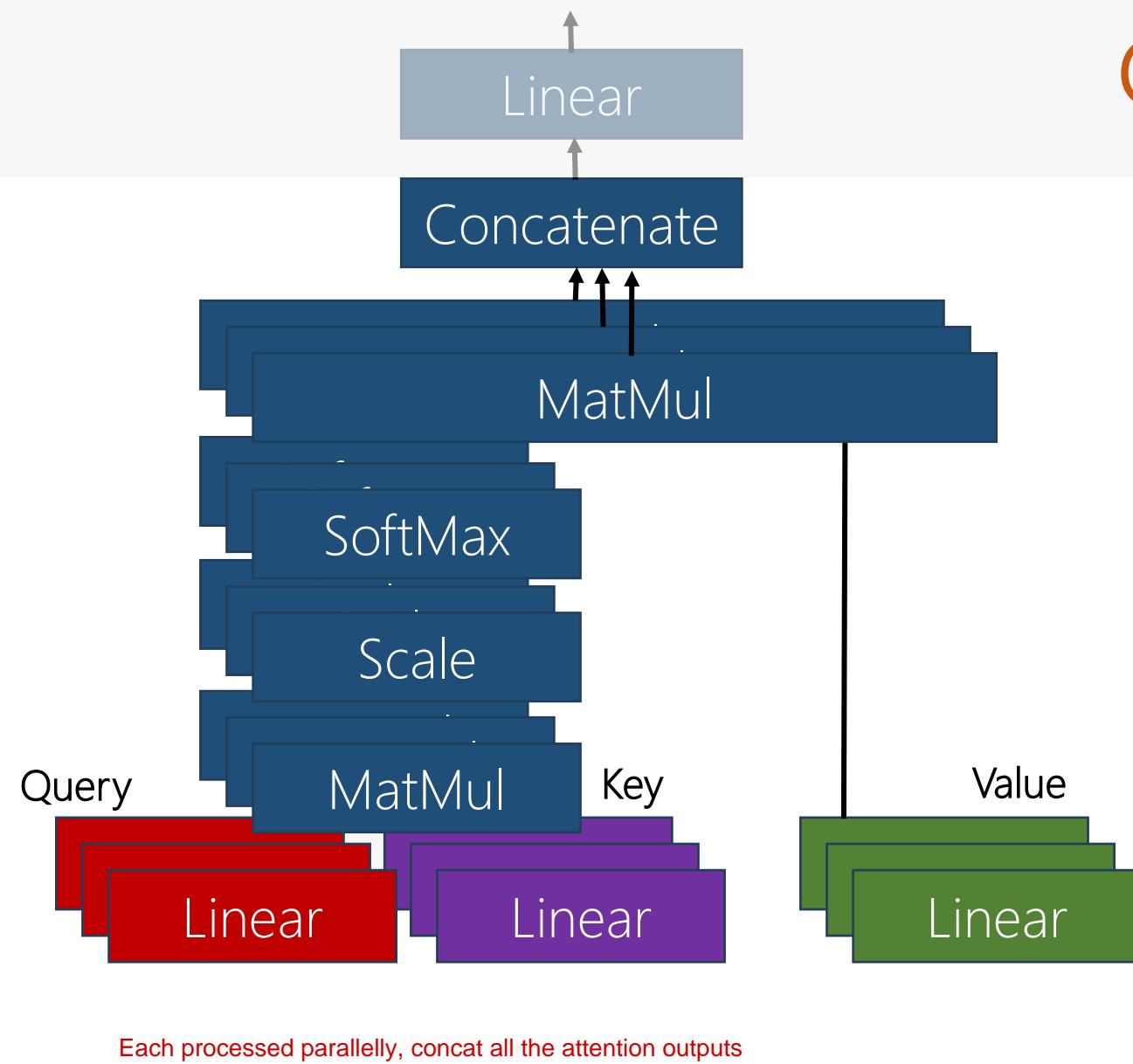
- 7x2 matrix: Q and K (blue arrows)
- softmax val: σ (blue arrow)
- numeration: $\sqrt{d_k}$ (blue arrow)
- val vector: V (blue arrow)
- output is weighted words rep

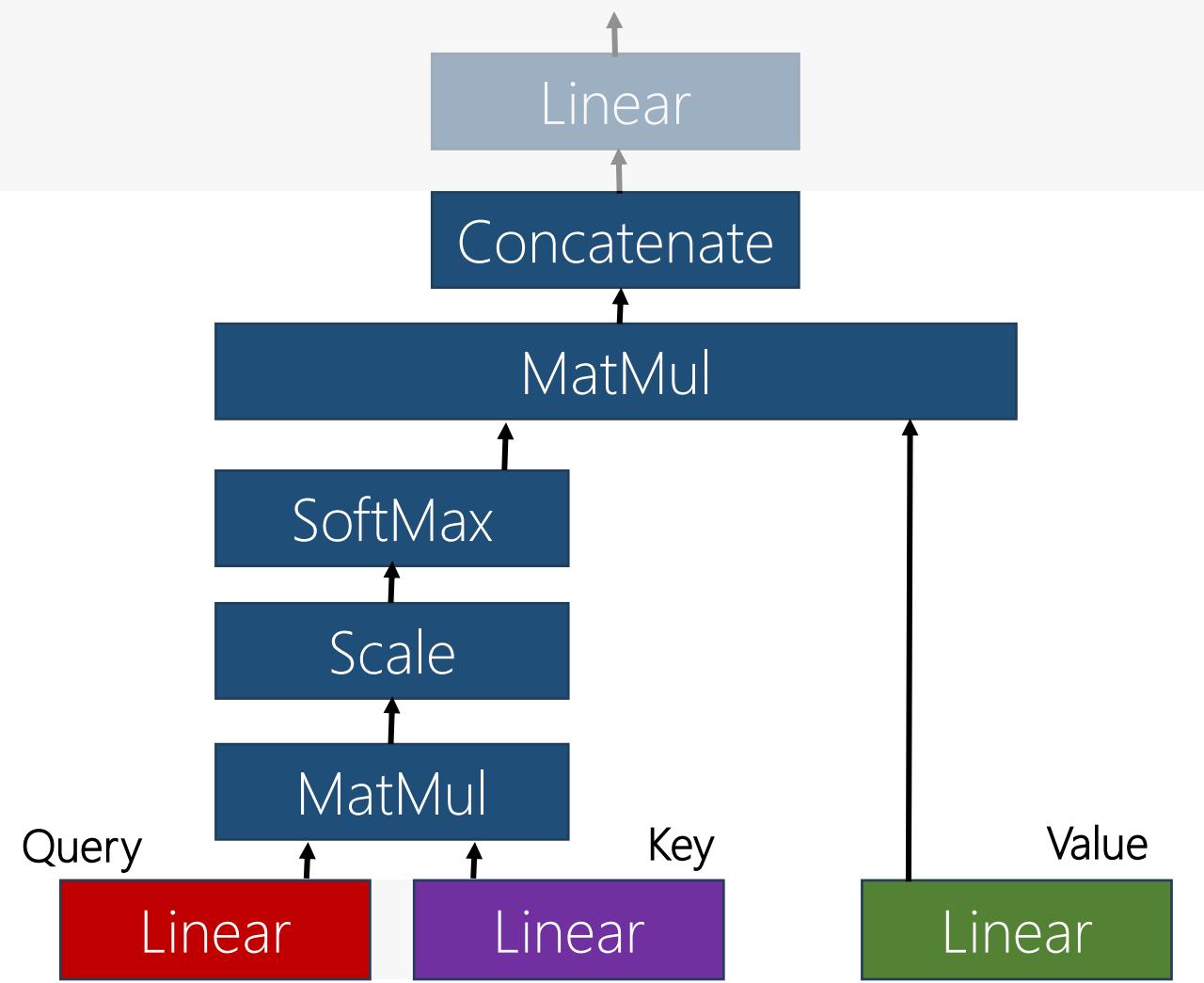
What's to Concatenate??

We do the above multiple times to get multiple attention heads to give more learning capabilities

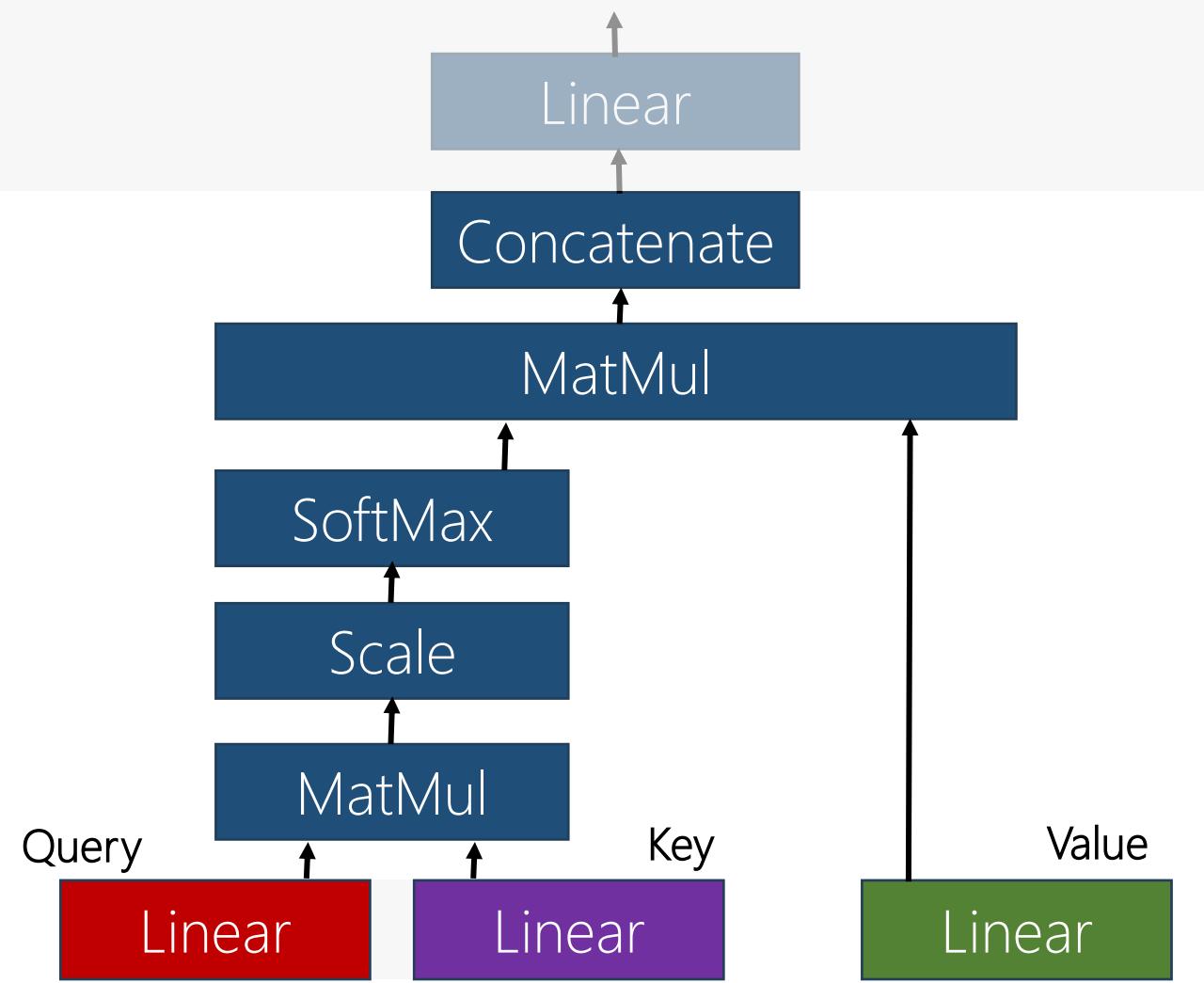


Concatenate Outputs from all Attention Heads

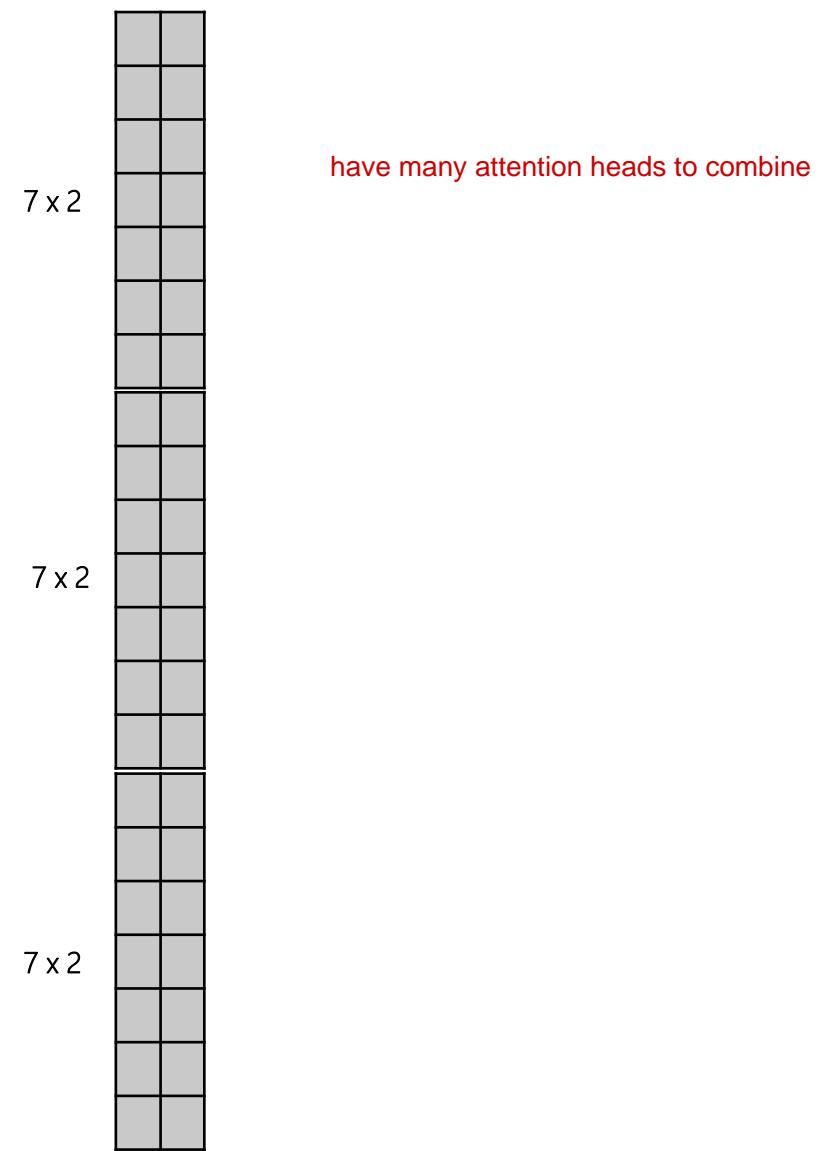




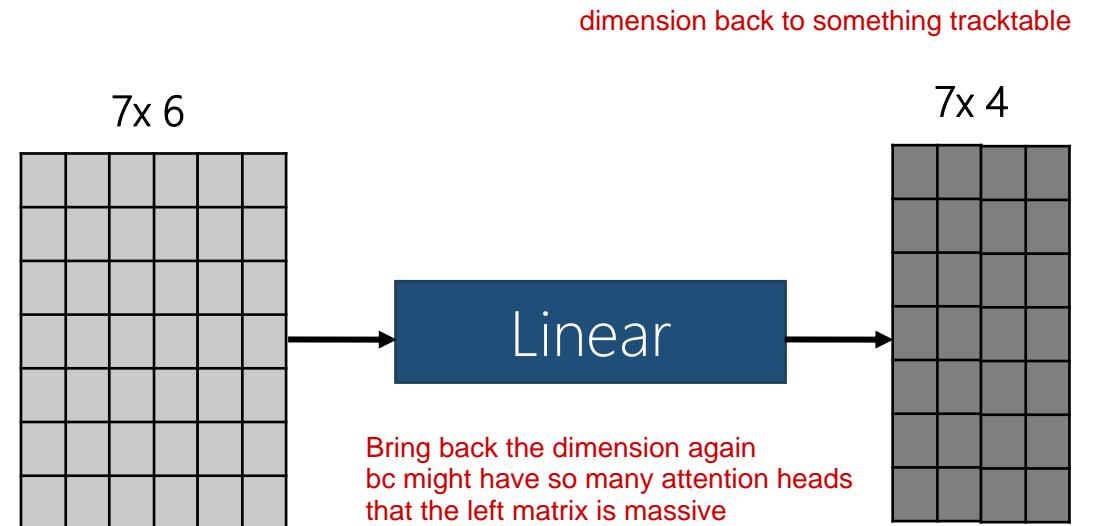
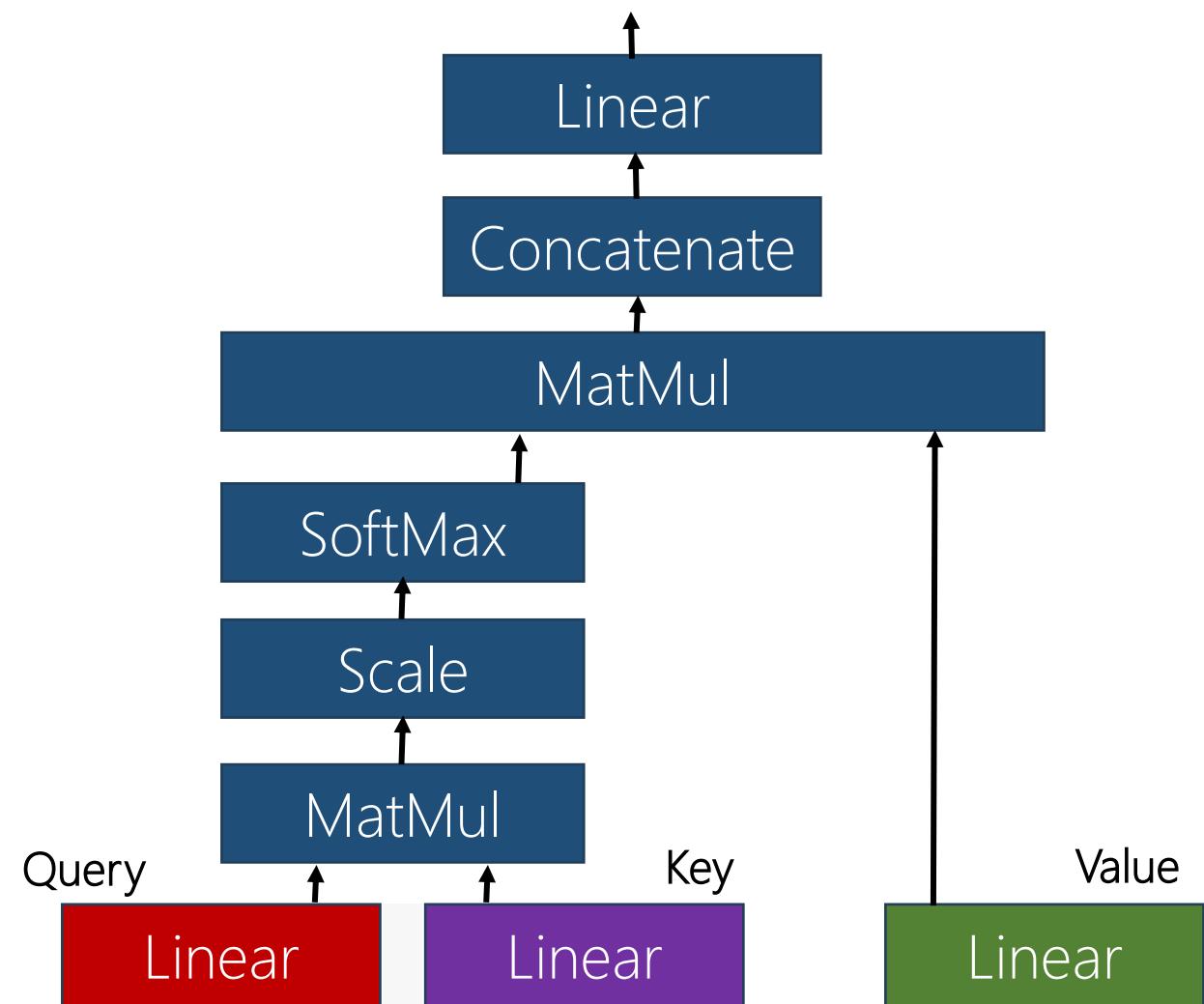
Another Linear layer??
Why? Why? Why?



Another Linear layer??
Why? Why? Why?

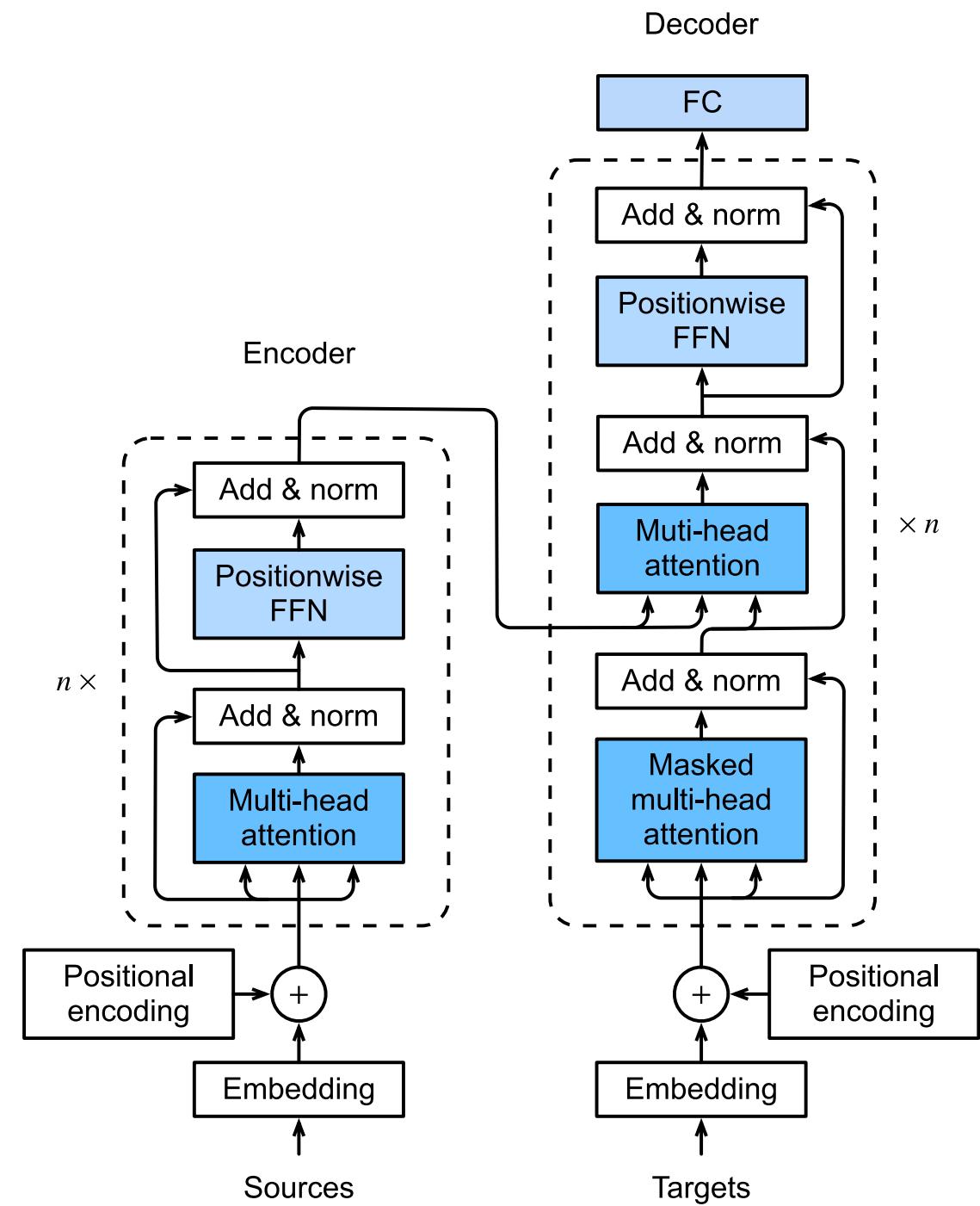


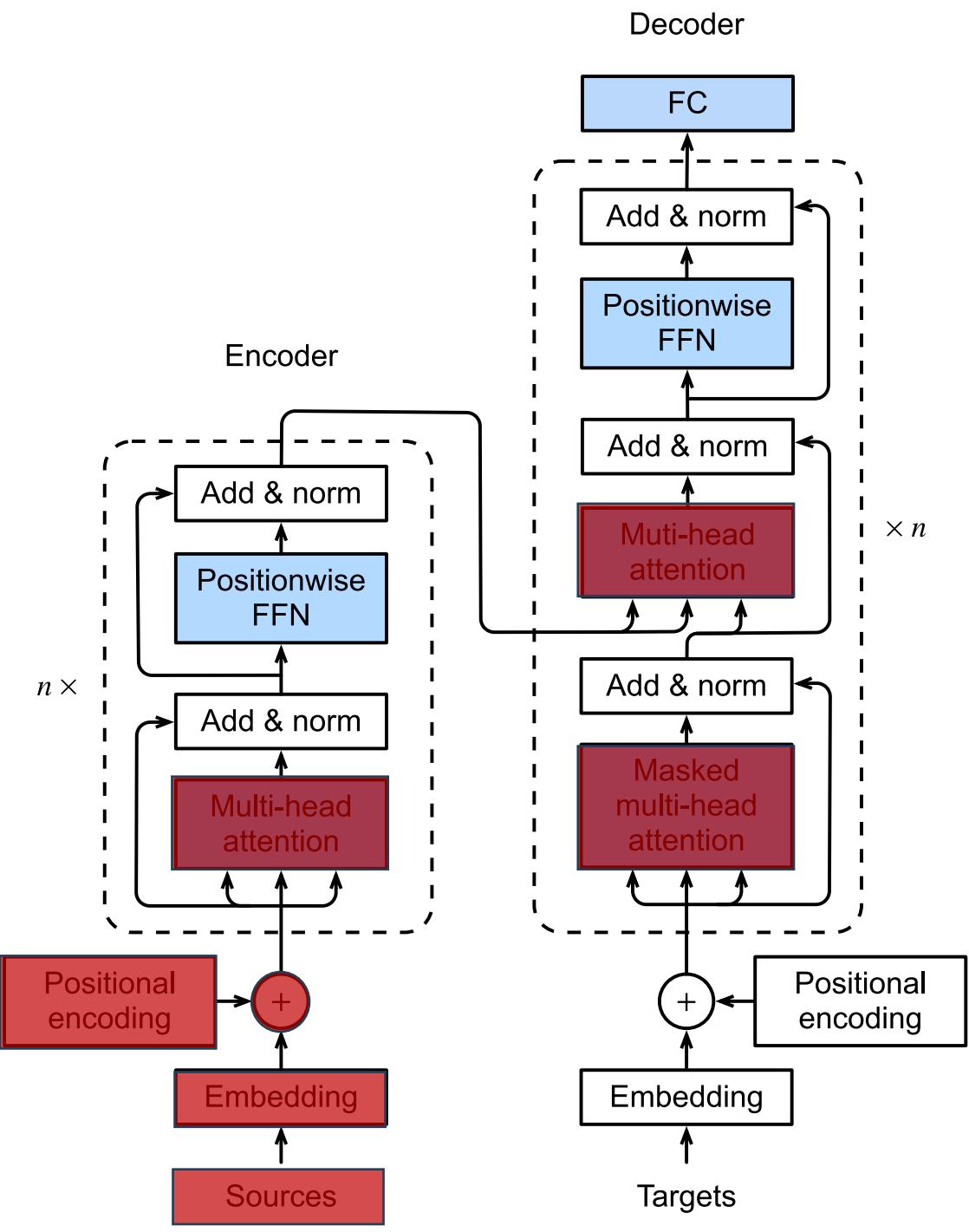
Another Linear layer?? Why? Why? Why?



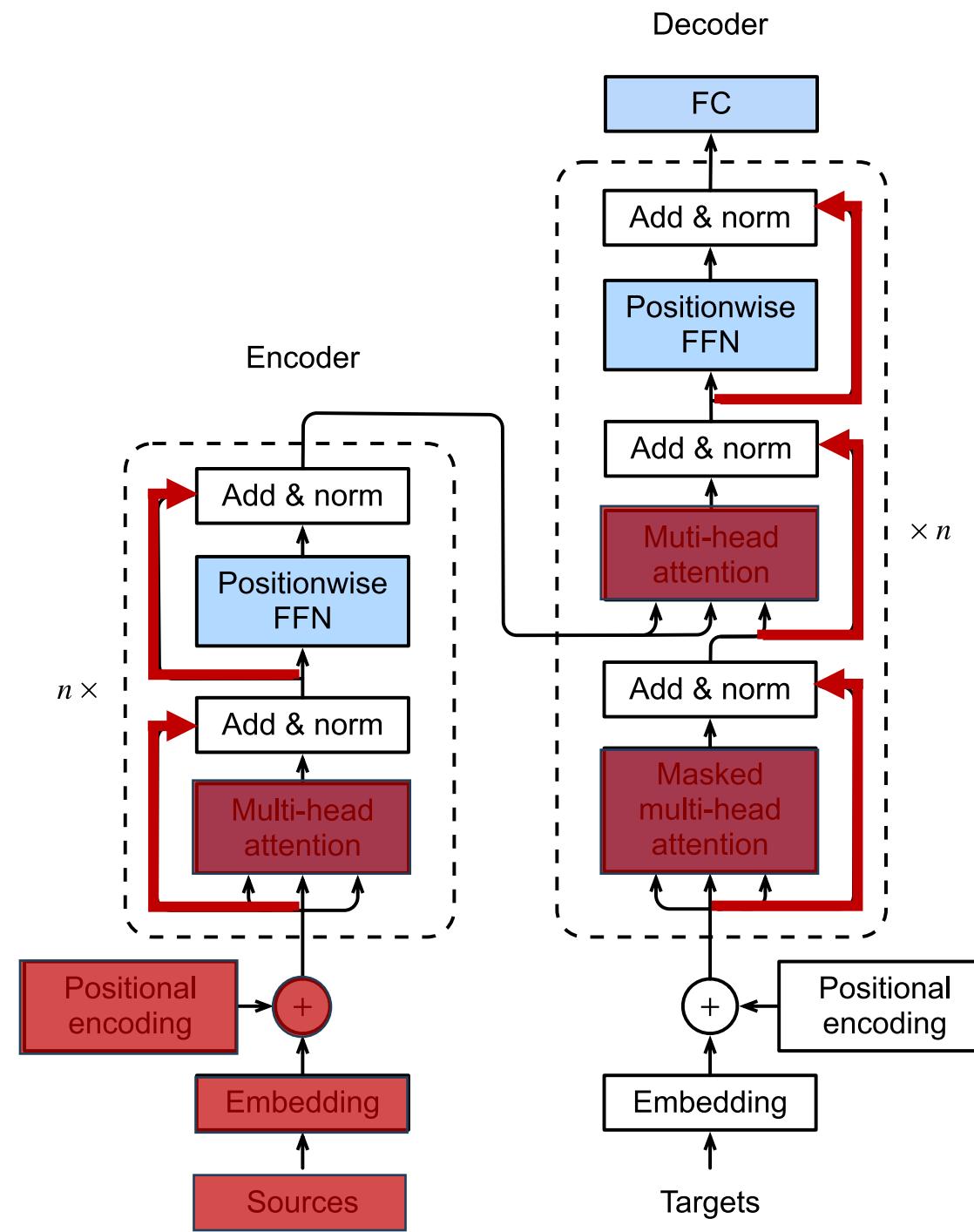
Continue here next time.....

What about the rest of the Transformer?



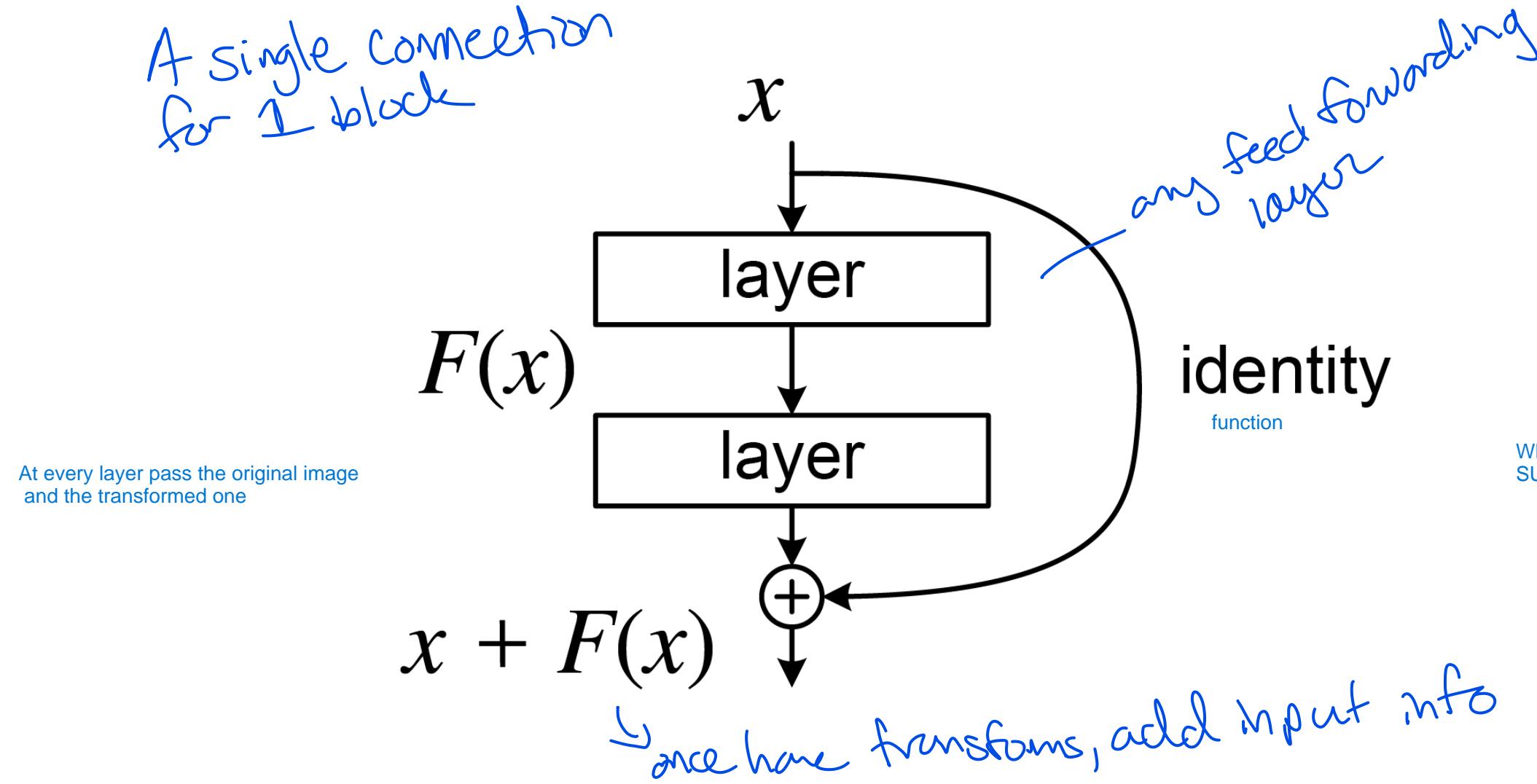


Residual Connection



Let's go down the memory lane...

Residual neural network



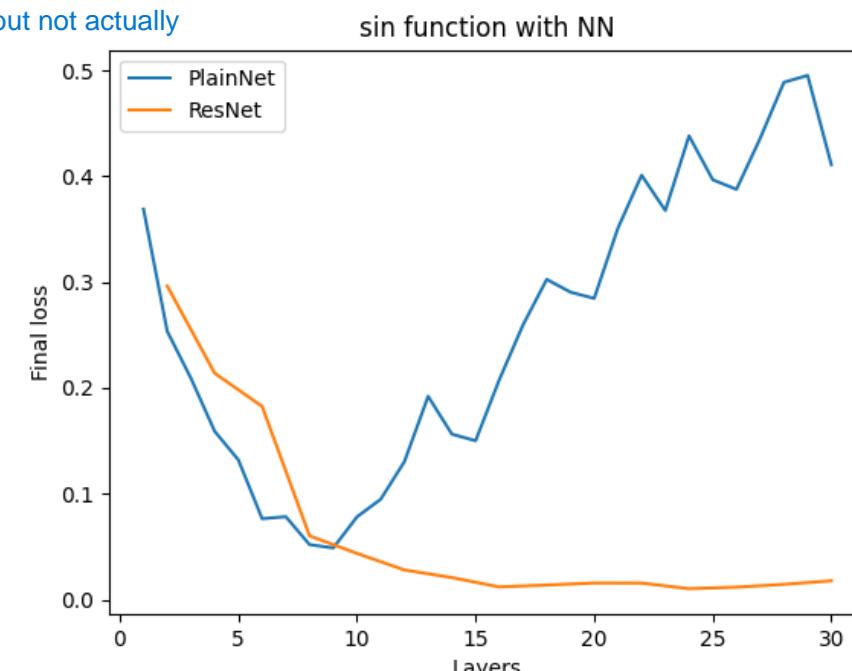
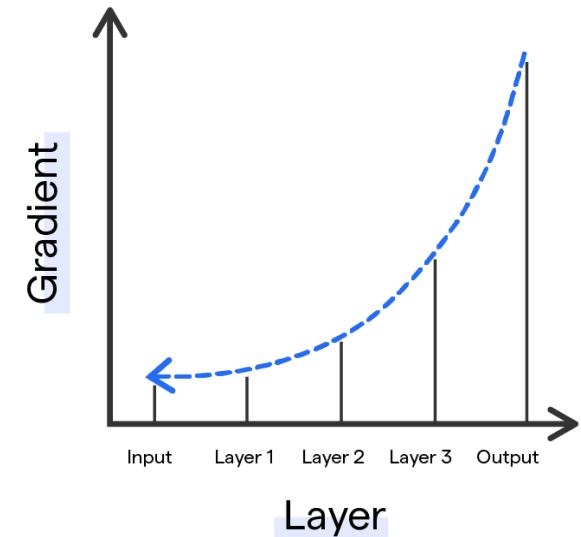
Remember backprop trying to learn weights: it's trying to learn $F(x)-x$ as the residual

WHAT MATH WERE WE SUPPOSED TO DO???????

Residual neural network

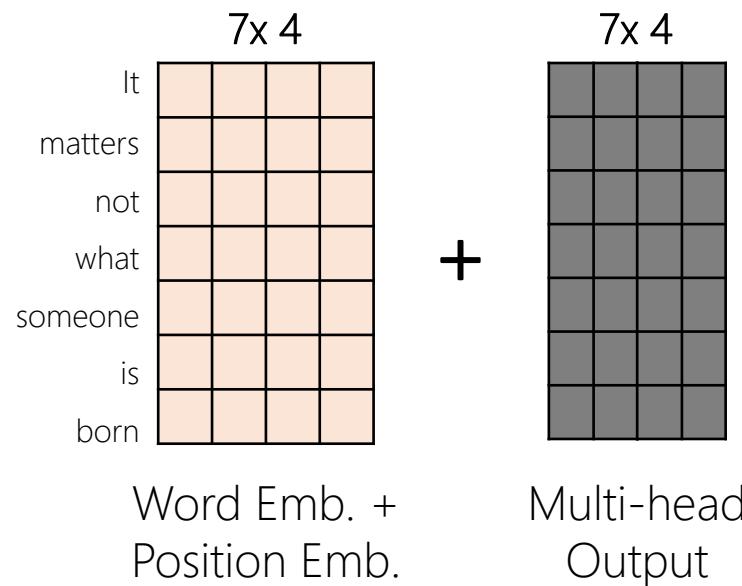
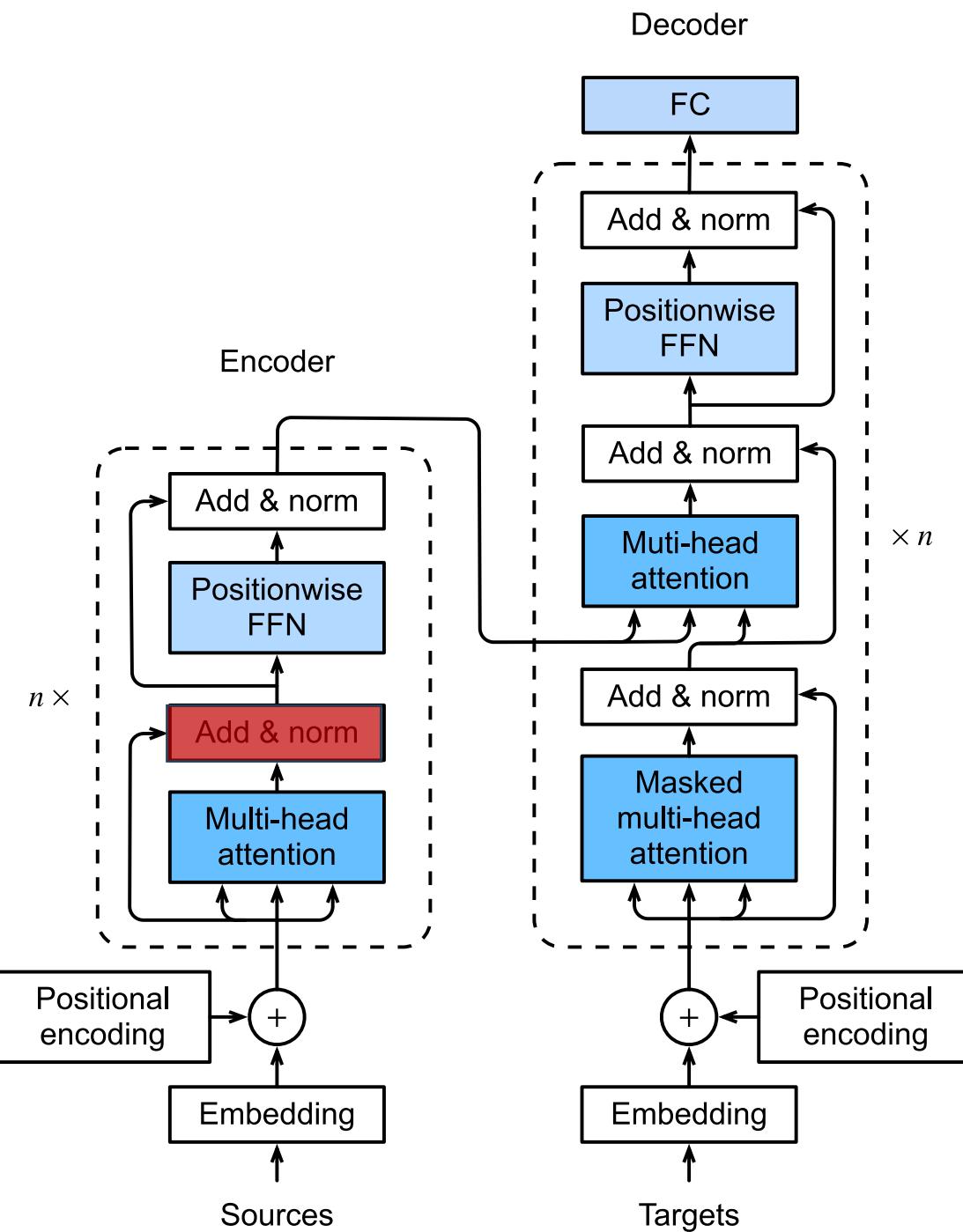
If train a network with more and more layers it starts diverging instead of converging

- **Vanishing Gradient:** As we backpropagate the gradients back through the network from the output layer towards the input the repeated multiplication of these small derivative values leads to increasingly smaller gradients could be close to zero but not actually
- **Degradation problem:** As a network deepens, the accuracy can start to decline even when training properly, which is not simply overfitting



Add & Norm

What are we adding?



Ppl now do layer normalization not batch normalization to fix covariance shifts

Pytorch has a layer norm function

Add & Norm

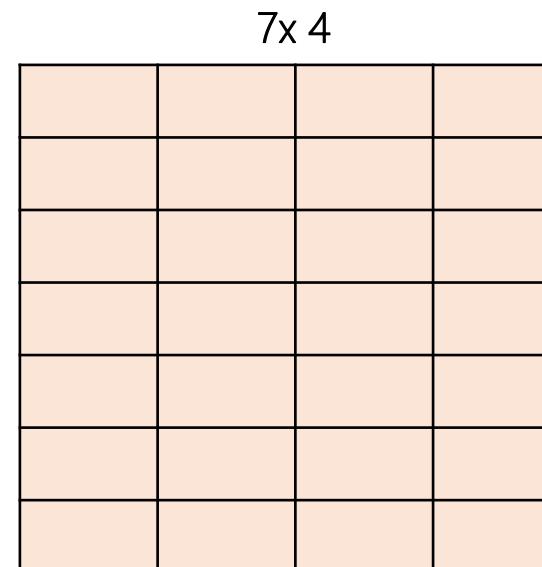
The diagram illustrates the Transformer architecture, divided into two main sections: the Encoder and the Decoder.

Encoder: The Encoder processes "Sources" (represented by a dashed box) through a series of layers. Each layer consists of a "Multi-head attention" block (blue), followed by an "Add & norm" block (red). The output of the final "Add & norm" block is passed through a "Positionwise FFN" (blue) and another "Add & norm" block (red). A "Positional encoding" block (white) is added to the input before the first layer. The final output of the Encoder is passed to the Decoder.

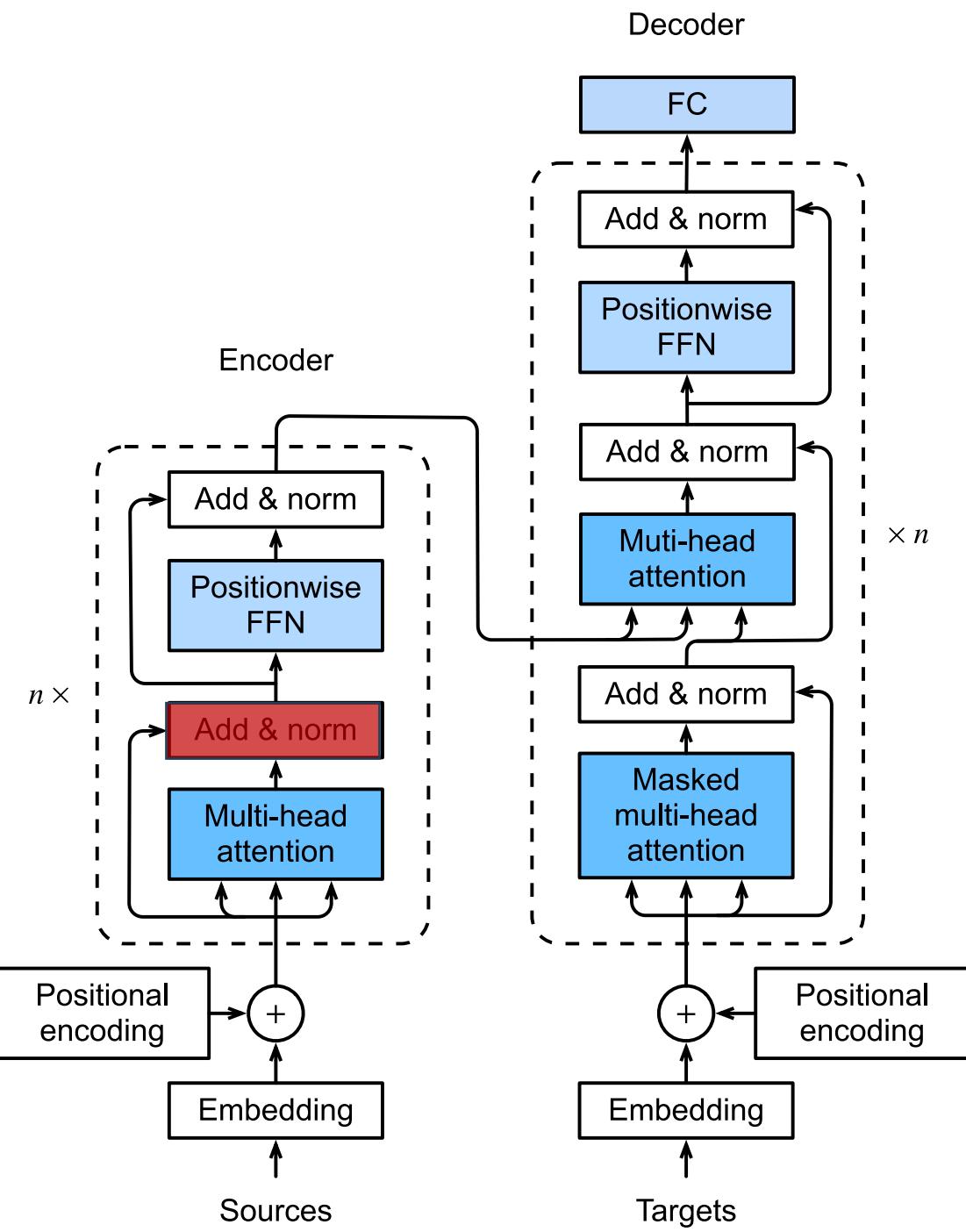
Decoder: The Decoder processes "Targets" (represented by a dashed box) through a series of layers. Each layer consists of a "Masked multi-head attention" block (blue), followed by an "Add & norm" block (blue). The output of the final "Add & norm" block is passed through a "Multi-head attention" block (blue), another "Add & norm" block (blue), a "Positionwise FFN" (blue), and a final "Add & norm" block (blue). A "Positional encoding" block (white) is added to the input before the first layer. The final output of the Decoder is processed by an "FC" (Fully Connected) layer.

Solves internal covariance shift inside the model representation

LayerNorm is used to stabilize the training process and addresses the internal covariate shift (ICS) problem, where the distribution of activations within a layer changes during training, making it difficult for the network to learn effectively.



Add & Norm

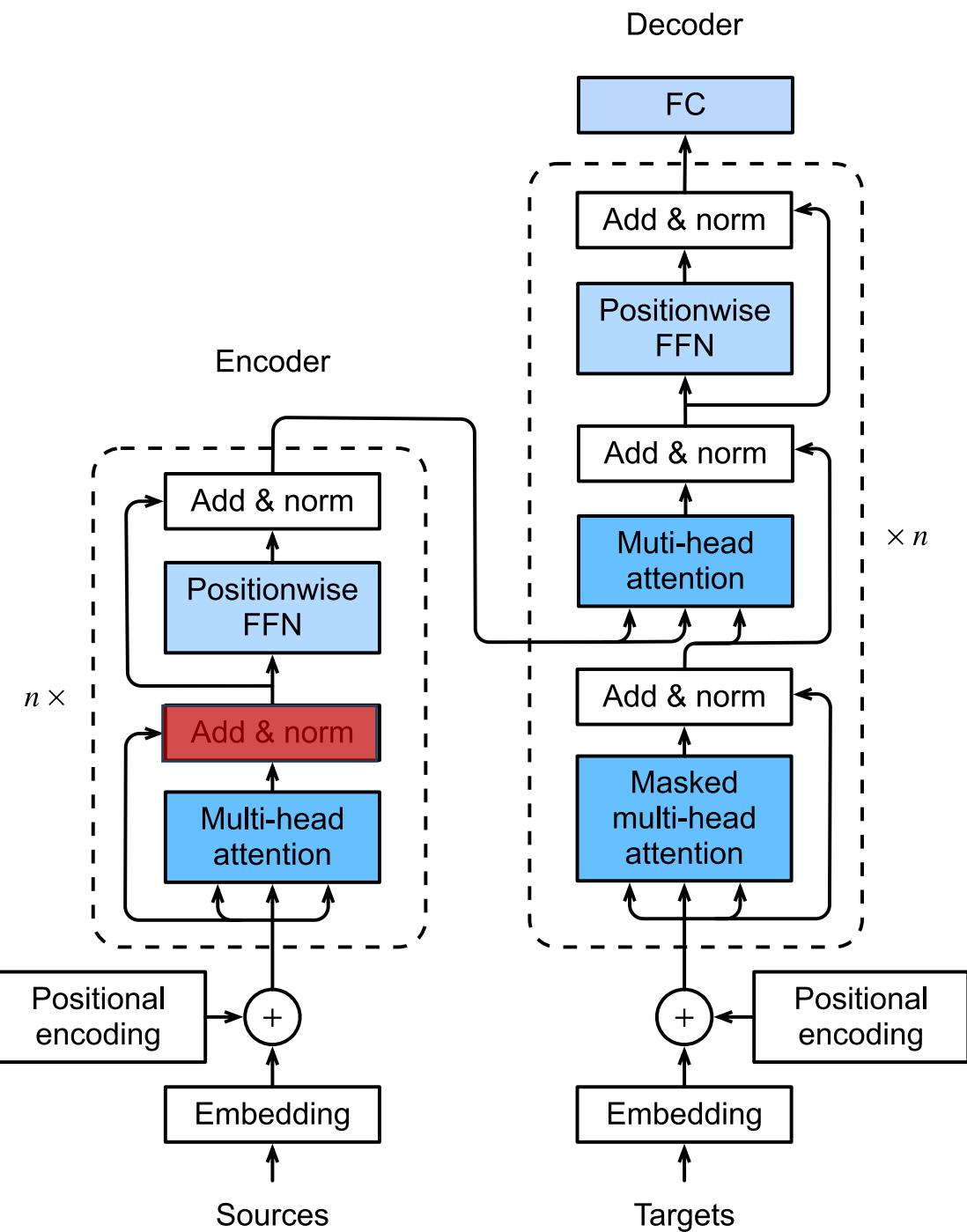


Attention output:

7x 4

It	1.56	2.12	0.91	2.87
matters	0.45	1.23	2.76	0.67
not	2.03	0.58	1.41	1.29
what	0.92	2.31	0.14	2.55
someone	1.80	0.61	2.98	1.52
is	2.67	0.33	1.99	1.74
born	0.48	2.40	1.68	0.29

Add & Norm



add the layers and normalize

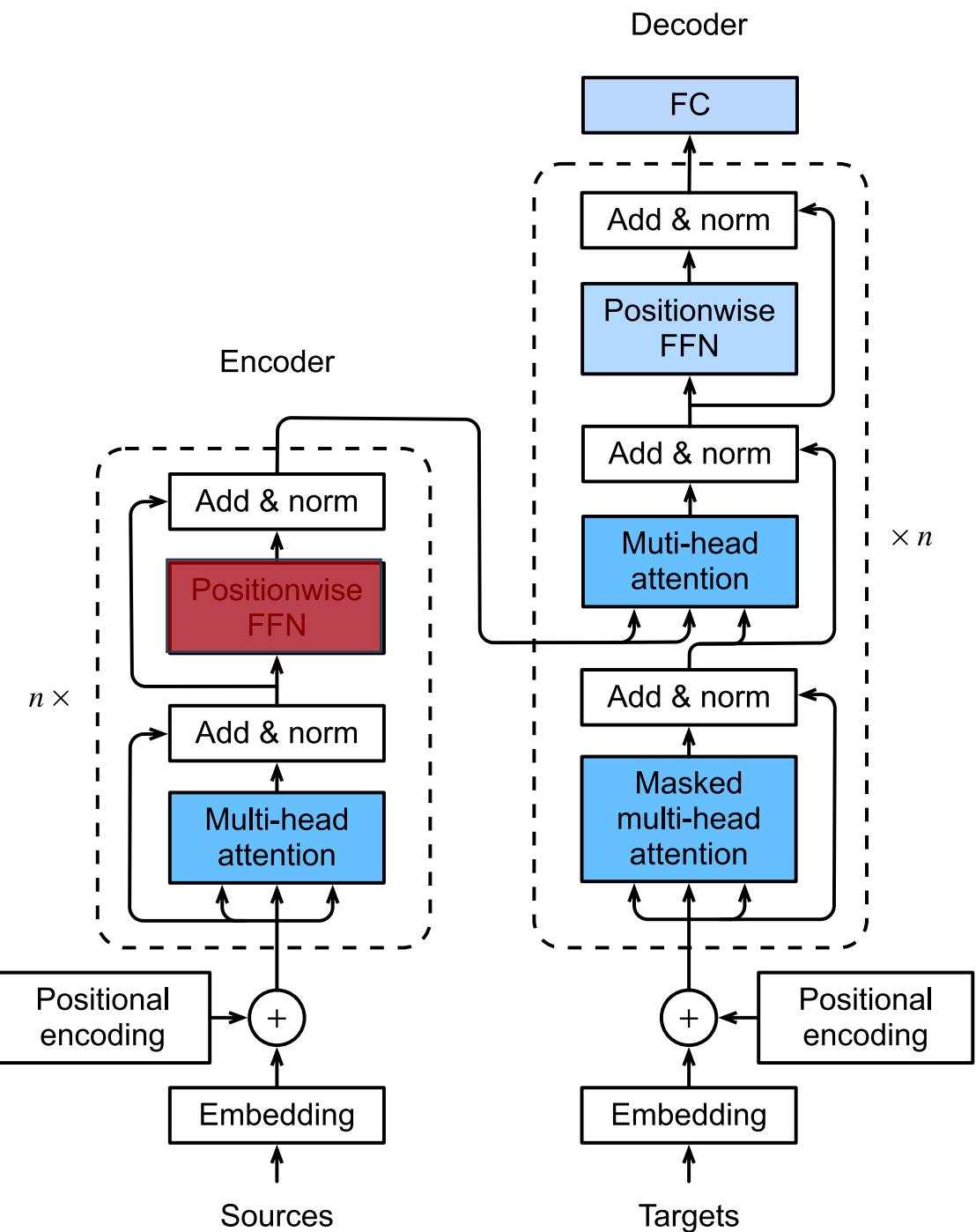


- calc mean and stdv row wise, want normalization of representation of each word
- 7 is num words, 4 is embedding dimension

	Mean (μ)	Stdev. (σ)
It	1.56	2.12
matters	0.45	0.91
not	2.03	2.76
what	0.58	0.67
someone	1.41	1.29
is	0.14	2.55
born	1.80	1.52
	2.67	1.74
	0.48	1.68
	2.40	0.29

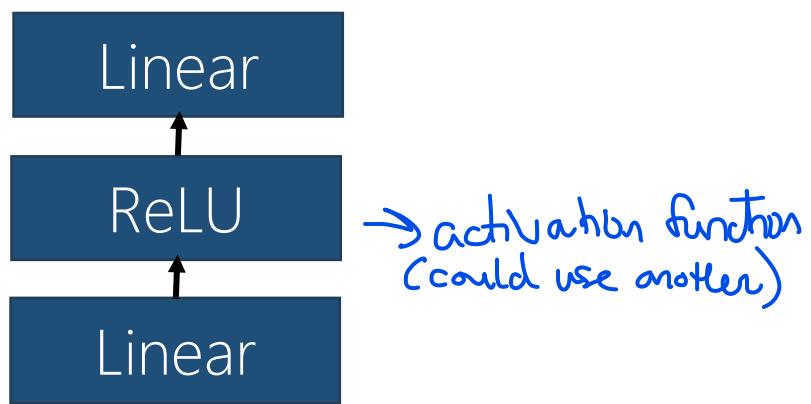
$$x'_i = \frac{x_i - \mu_i}{\sqrt{\sigma_i^2 + \epsilon}}$$

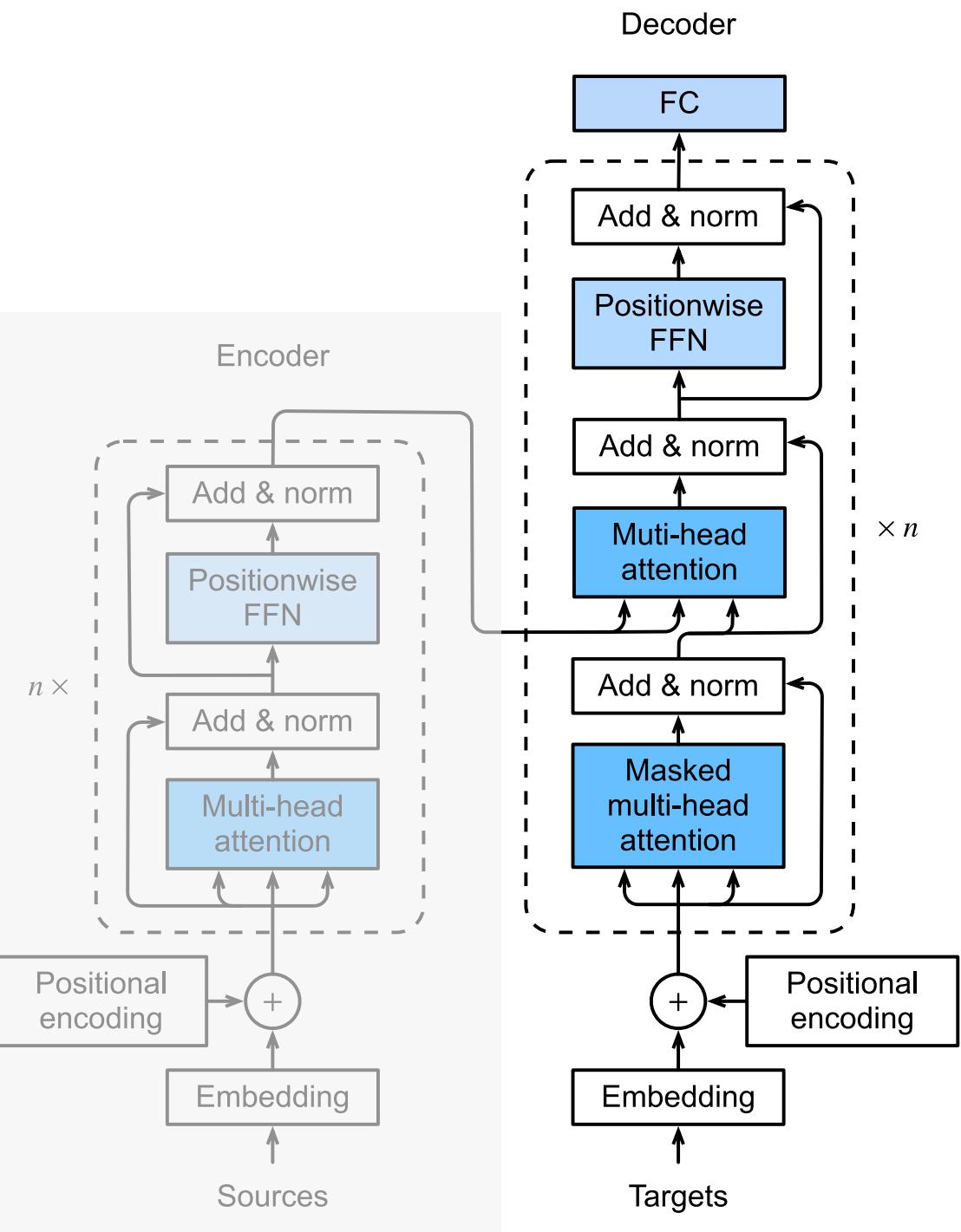
normalize
each
value
in the
attention
matrix ?



A bunch of Linear is all we need!

Now use
Gelu

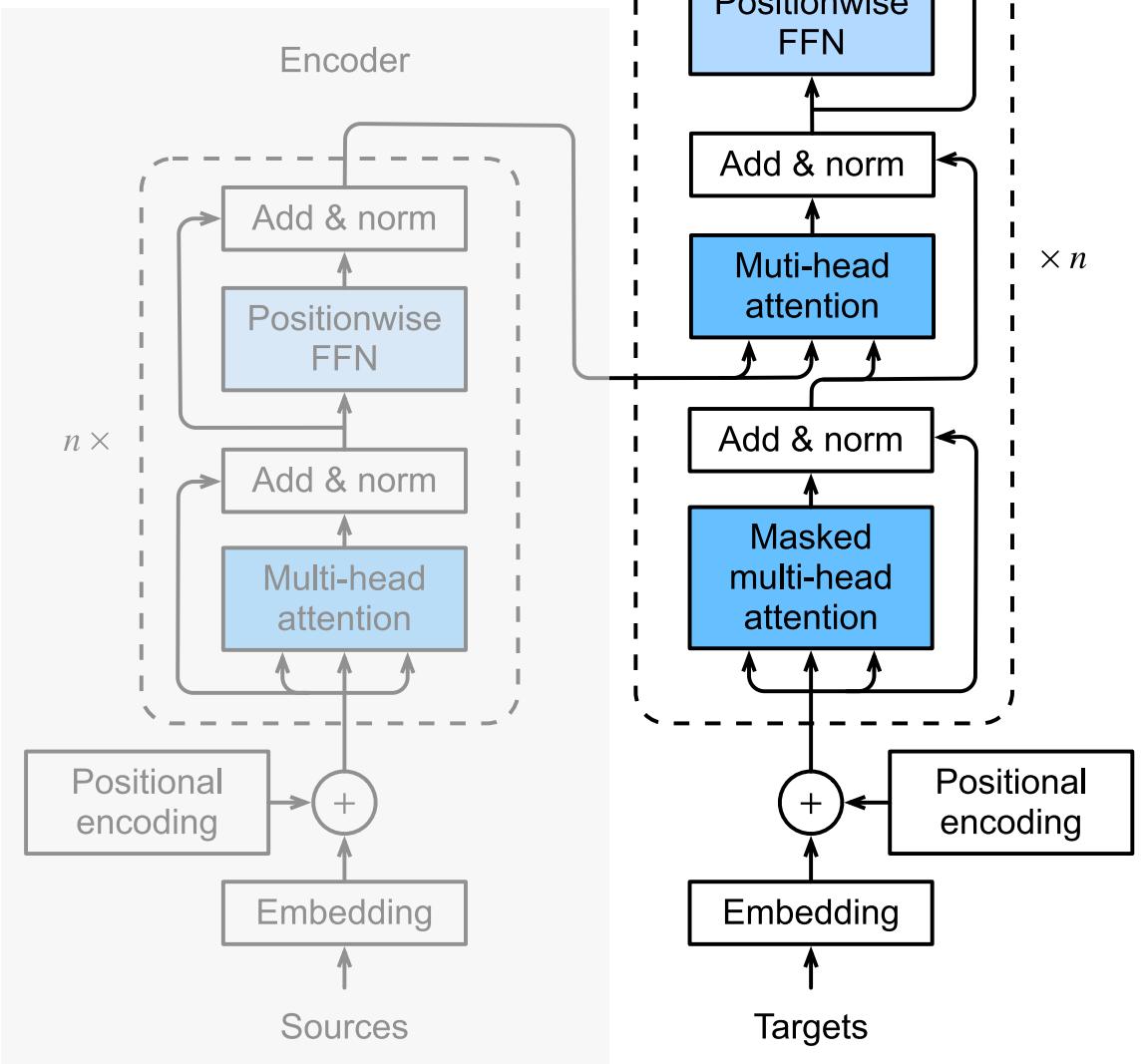


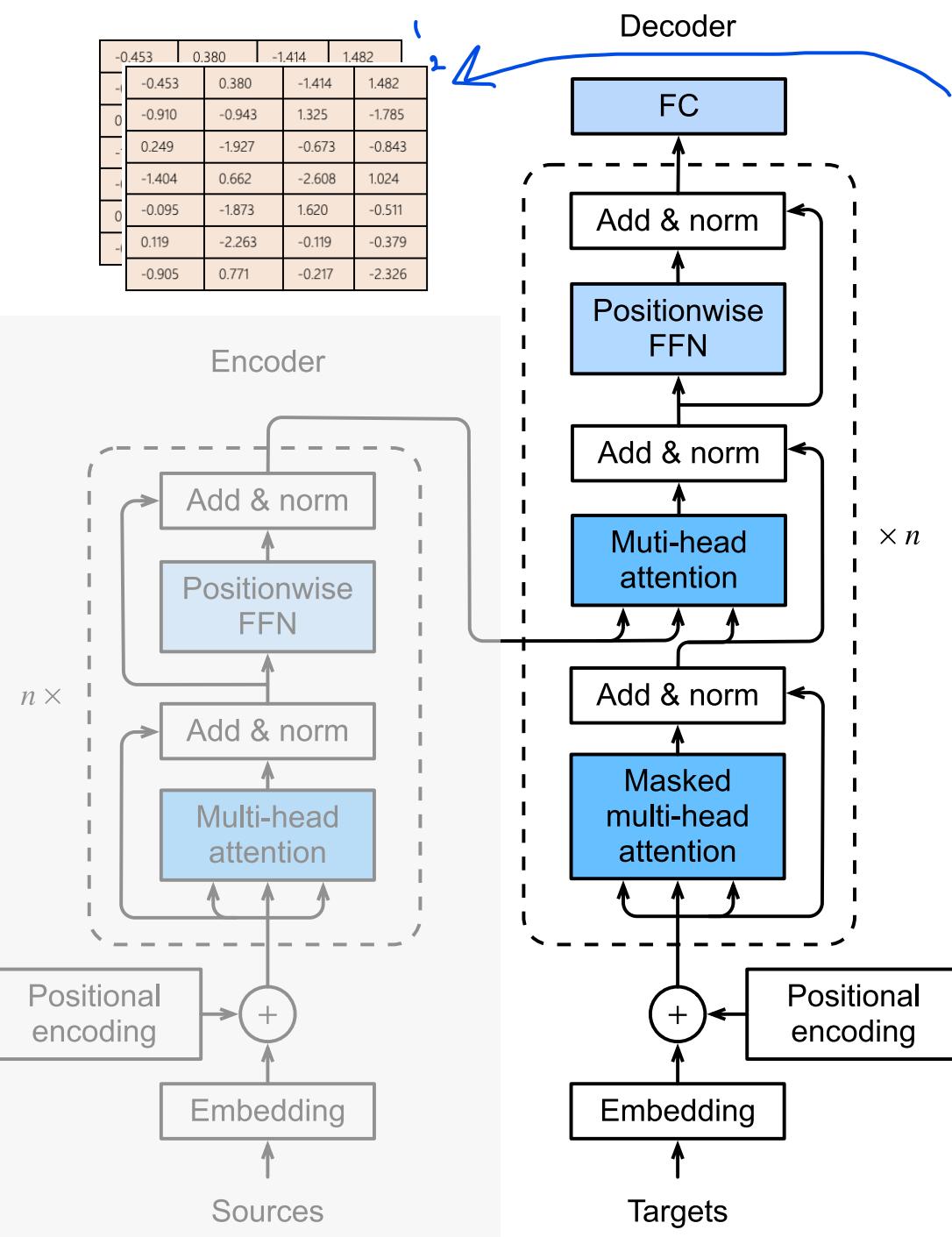


That covers our Encoder part of the Transformer!



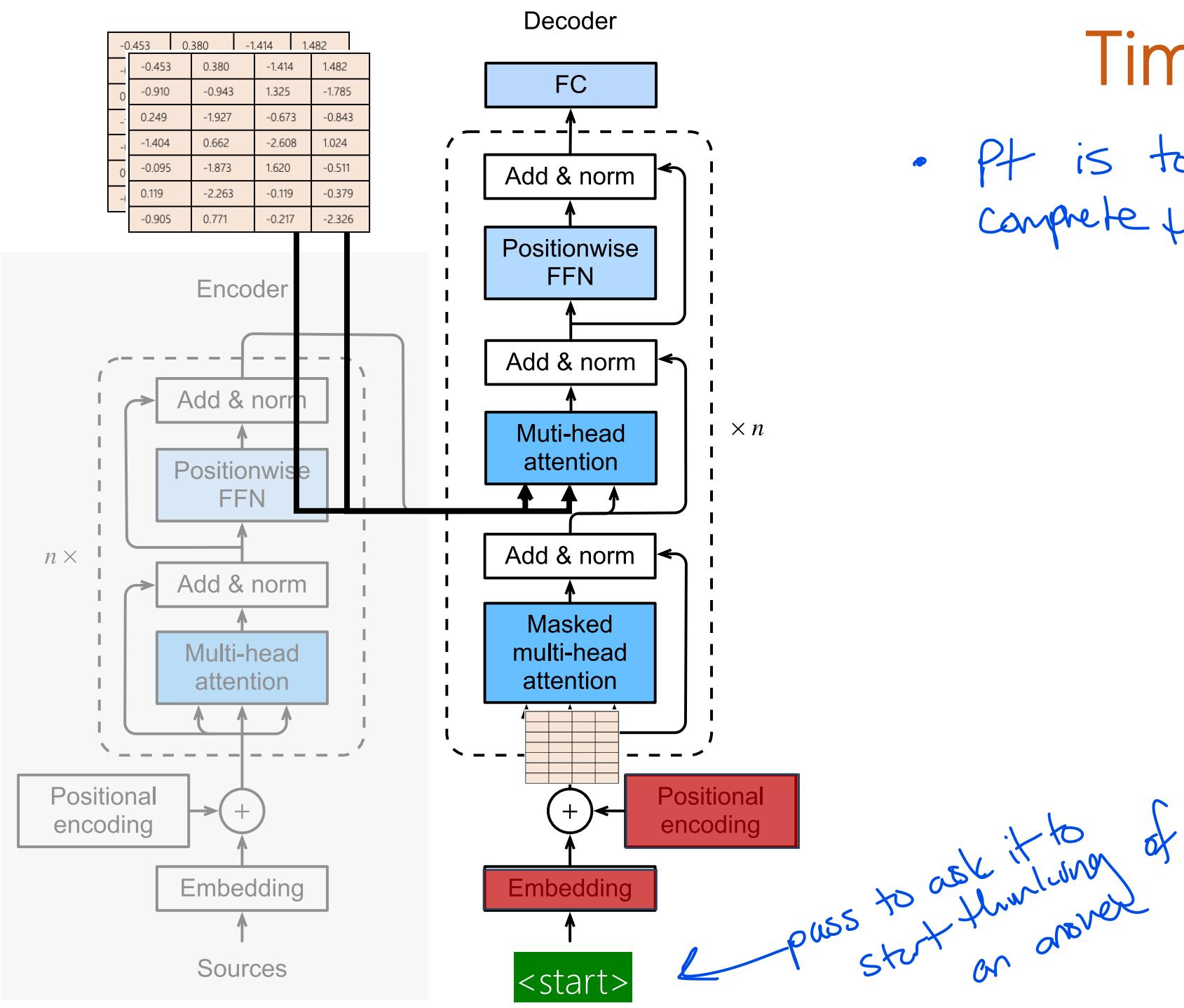
lt	-0.453	0.380	-1.414	1.482
matters	-0.910	-0.943	1.325	-1.785
not	0.249	-1.927	-0.673	-0.843
what	-1.404	0.662	-2.608	1.024
someone	-0.095	-1.873	1.620	-0.511
is	0.119	-2.263	-0.119	-0.379
born	-0.905	0.771	-0.217	-2.326



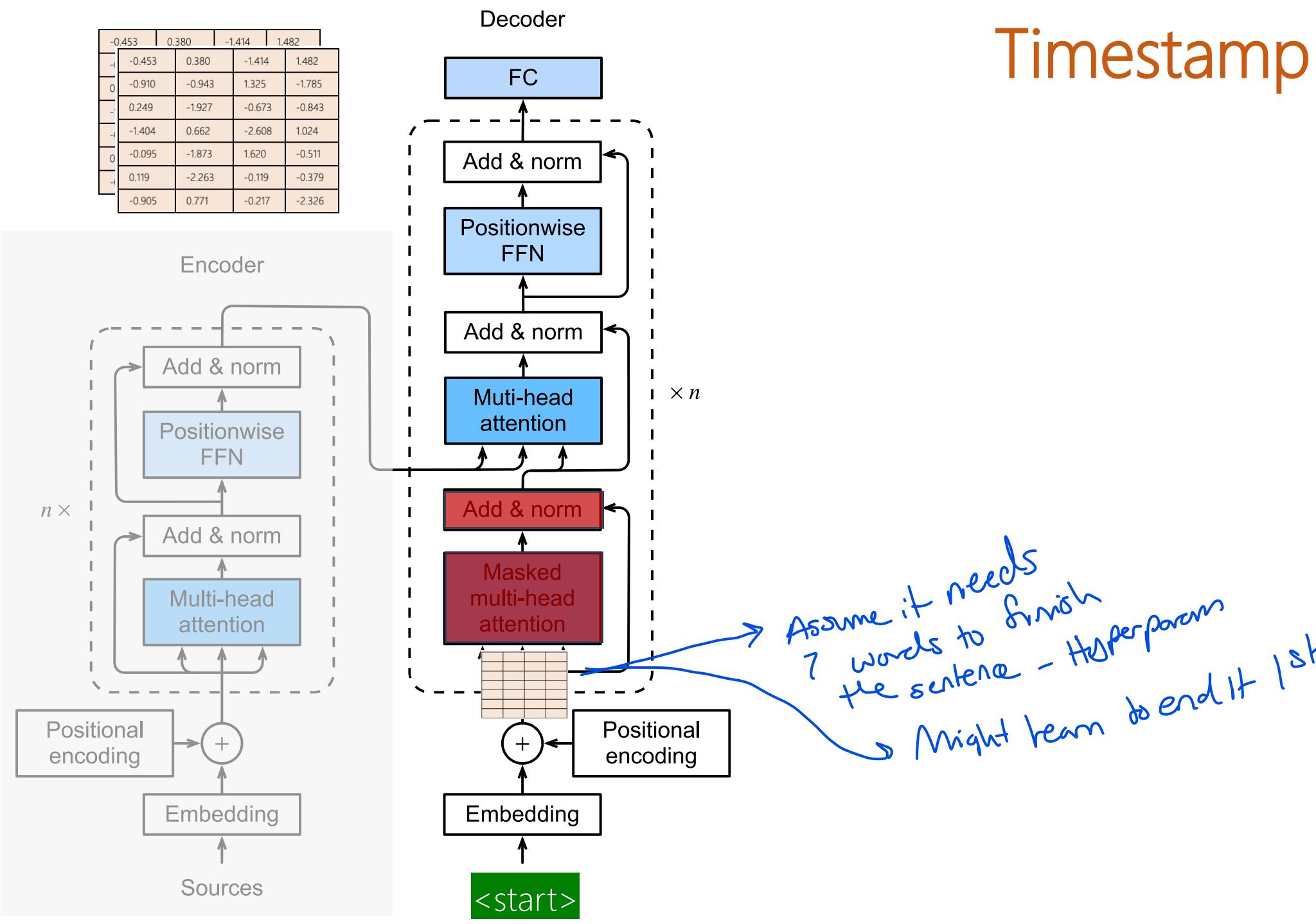


Timestamp = 1

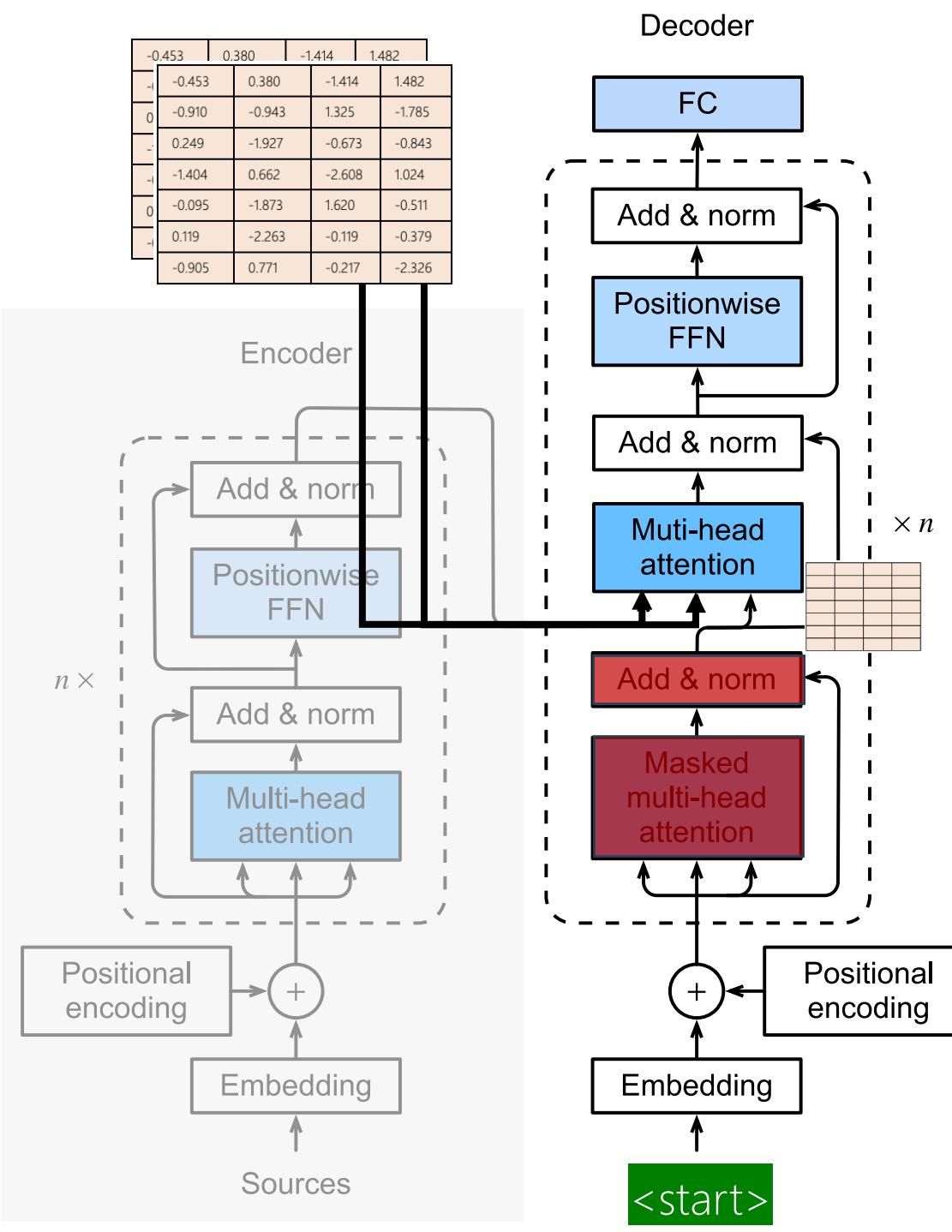
- PT is to ask the model to complete the phrase

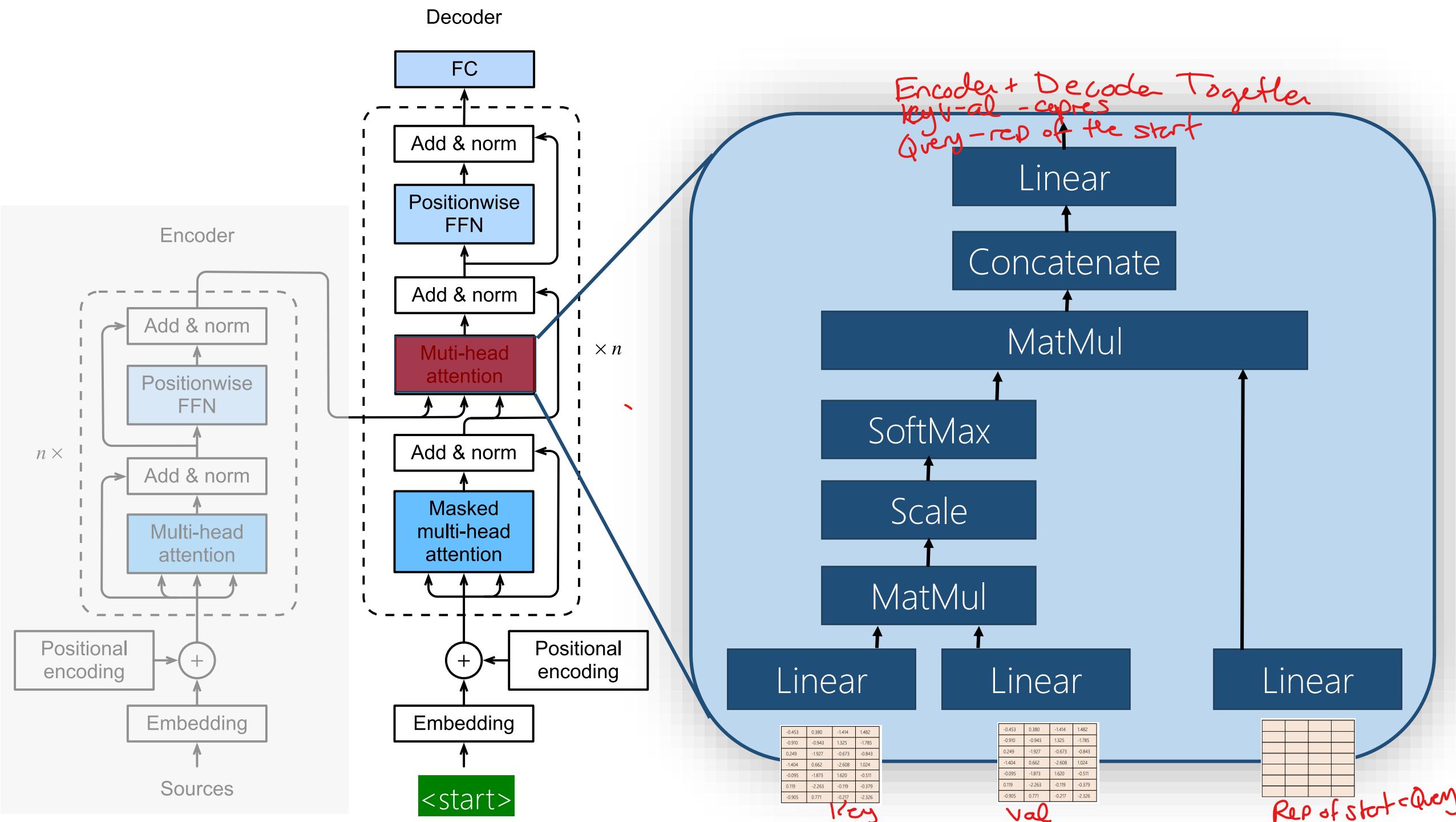


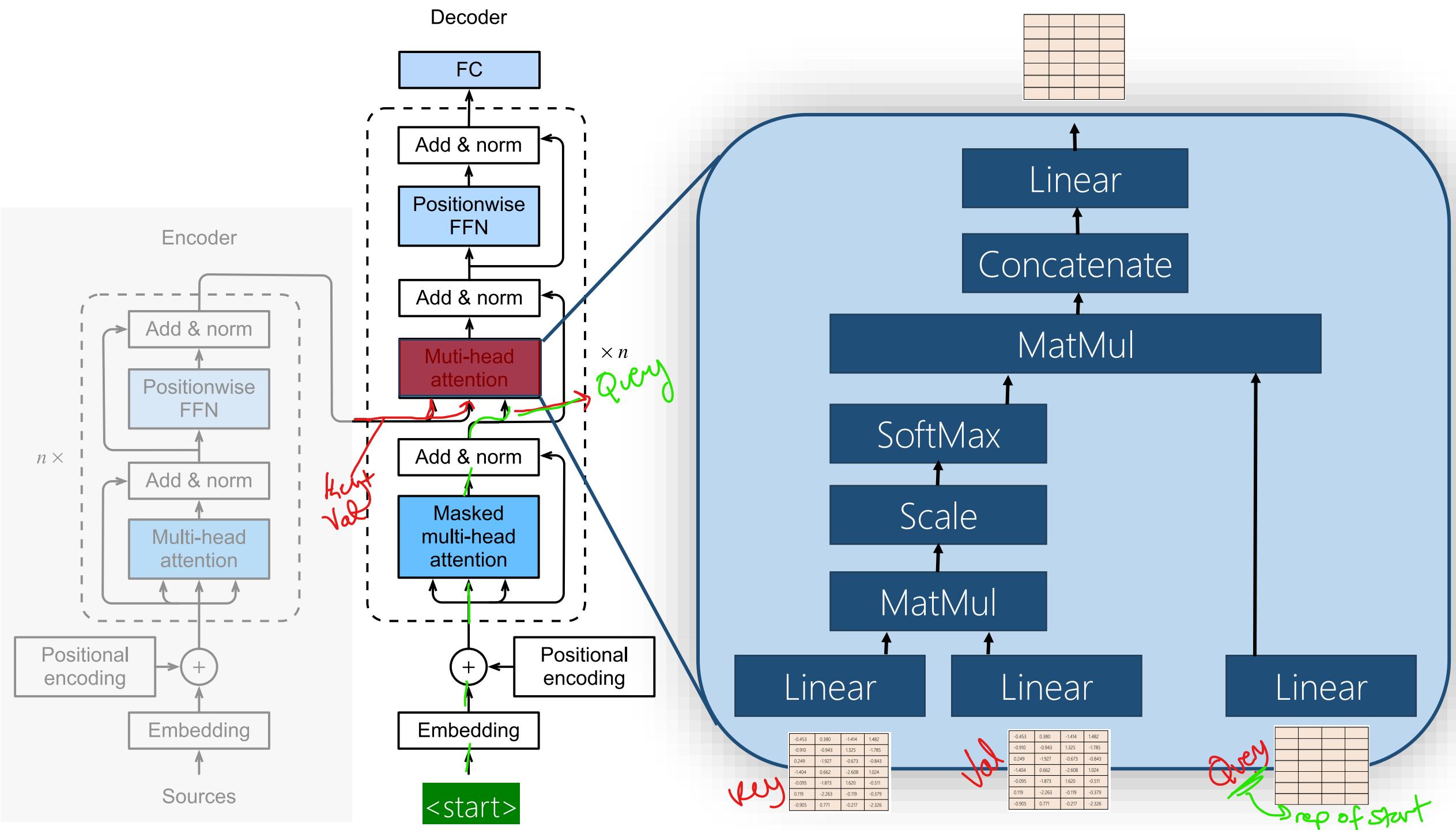
Timestamp = 1

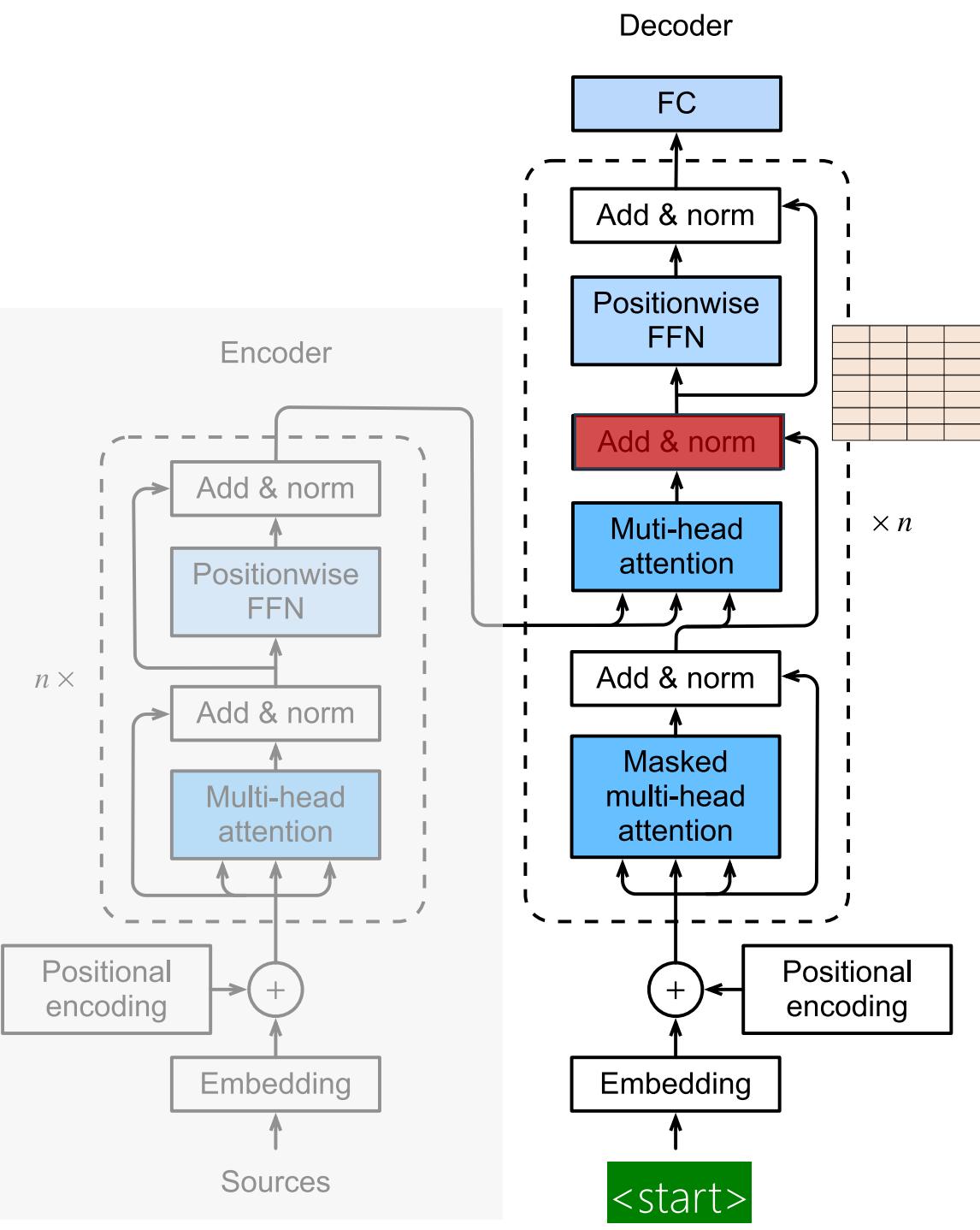


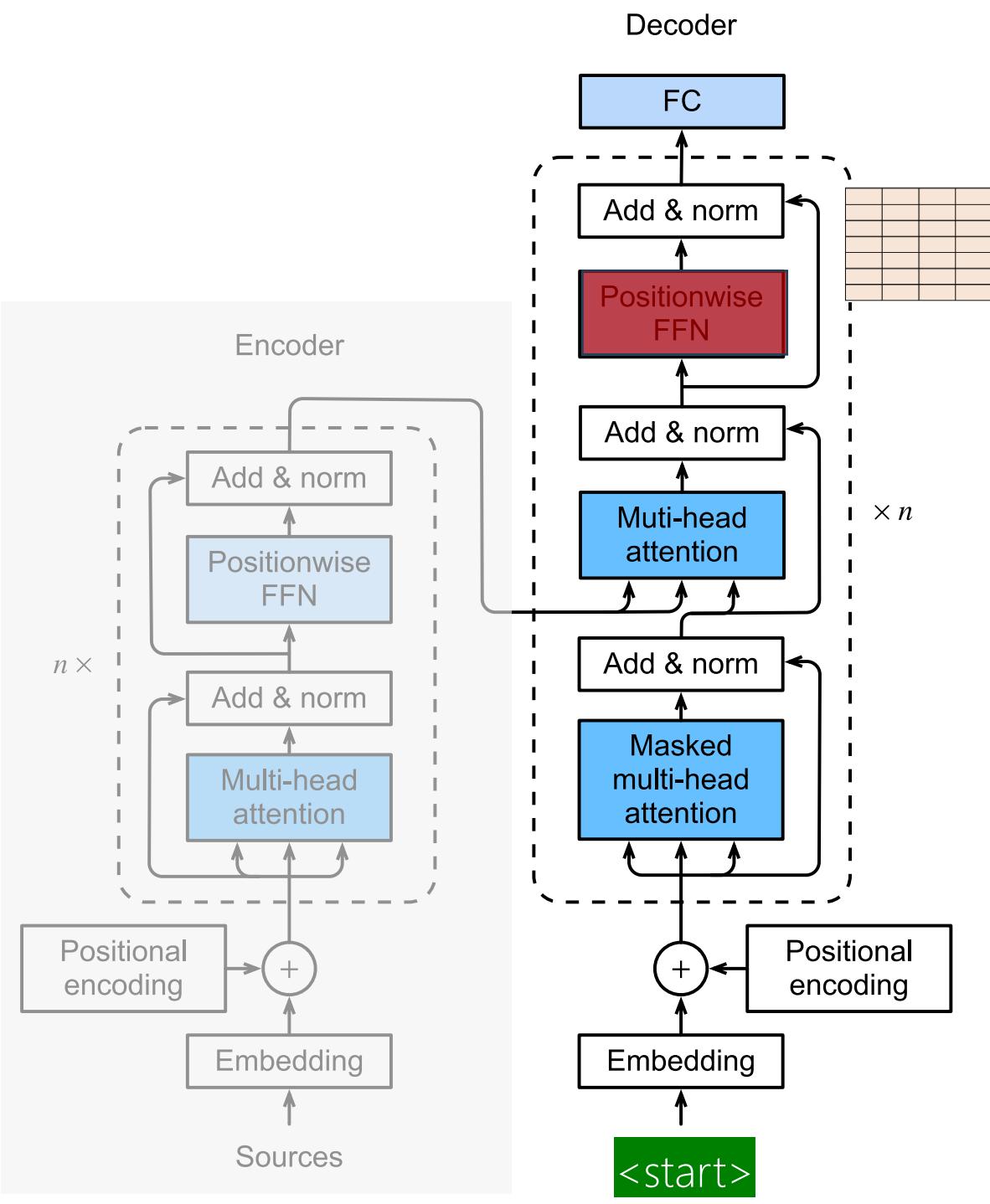
Timestamp = 1

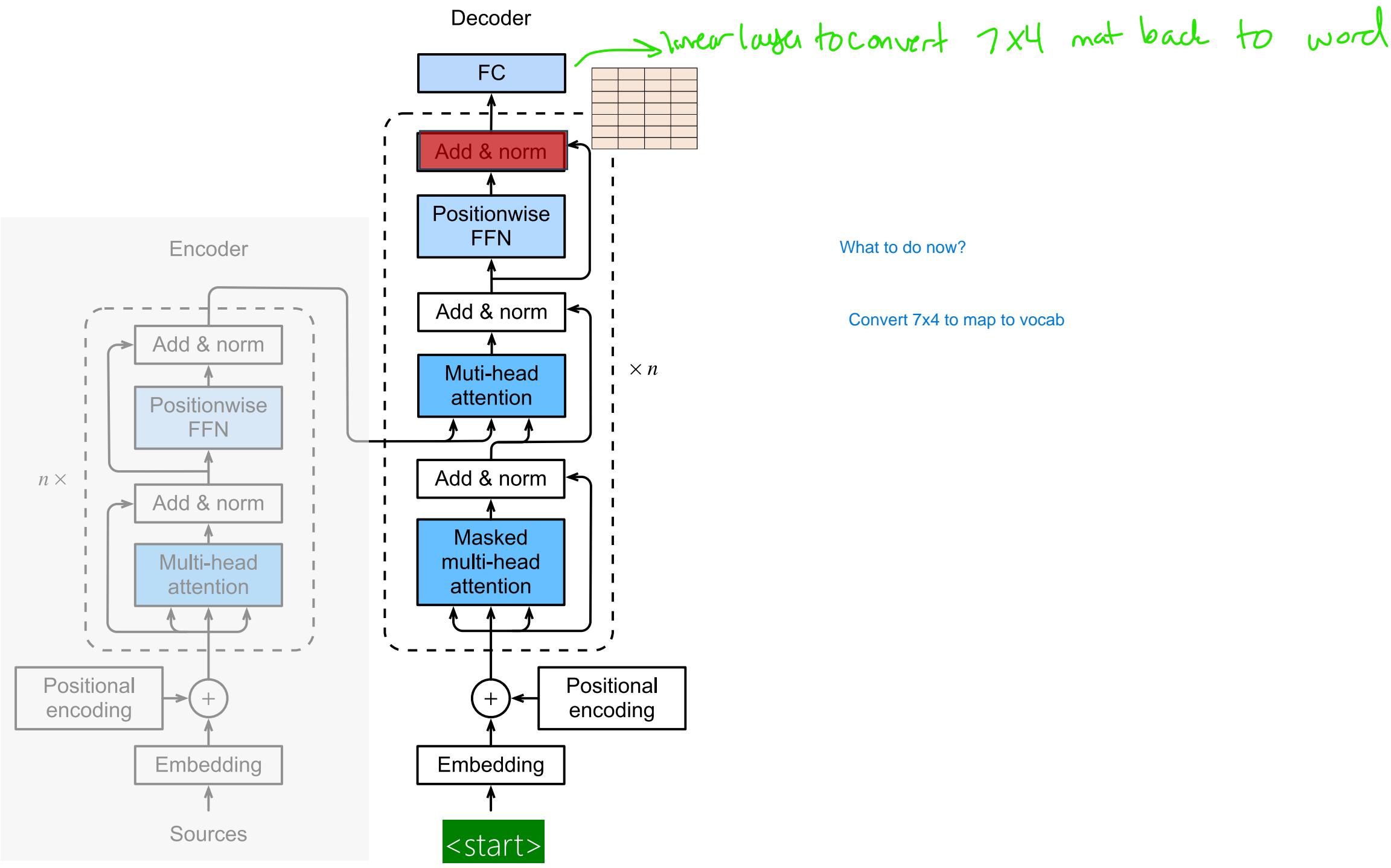


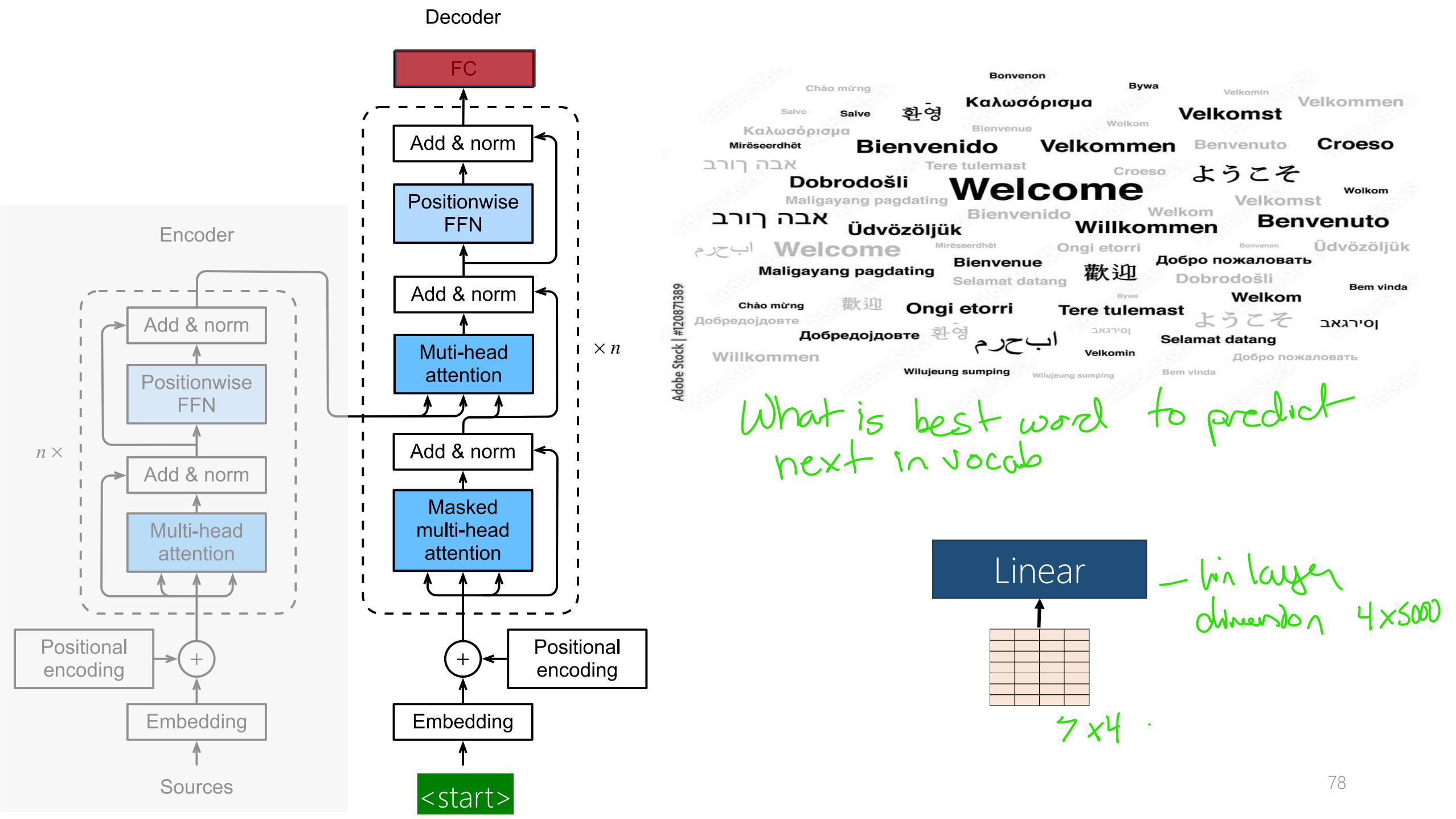


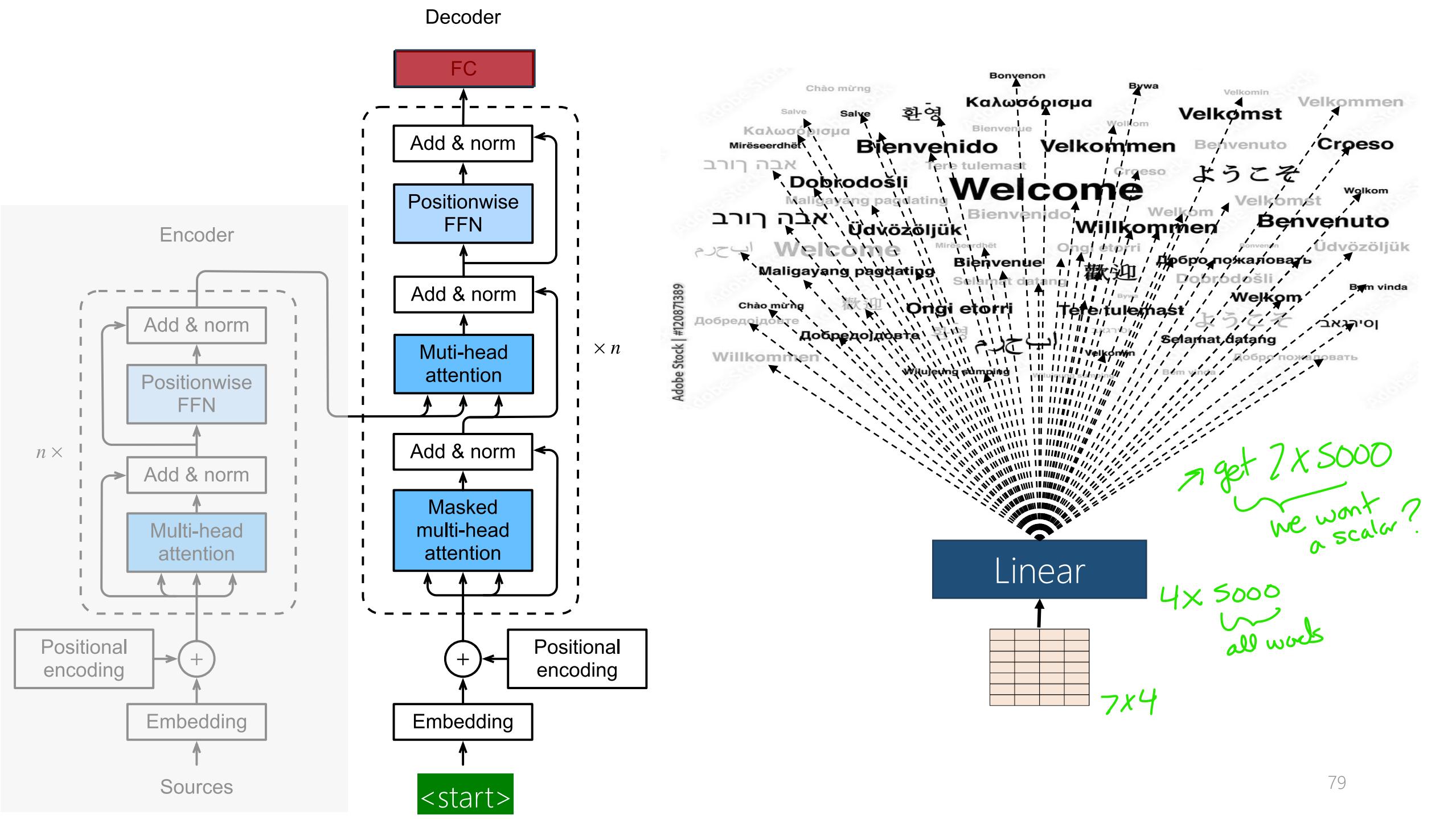


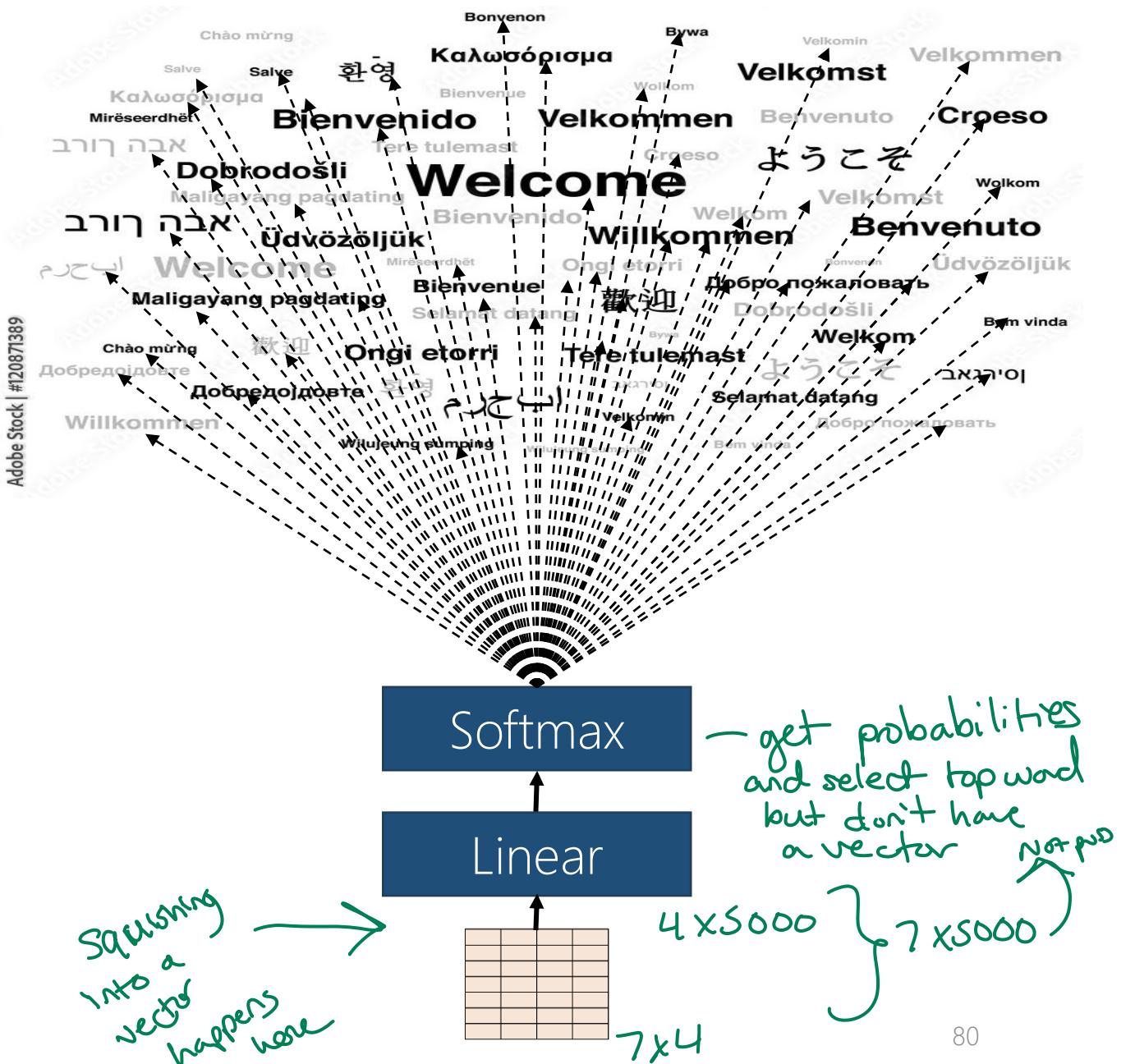
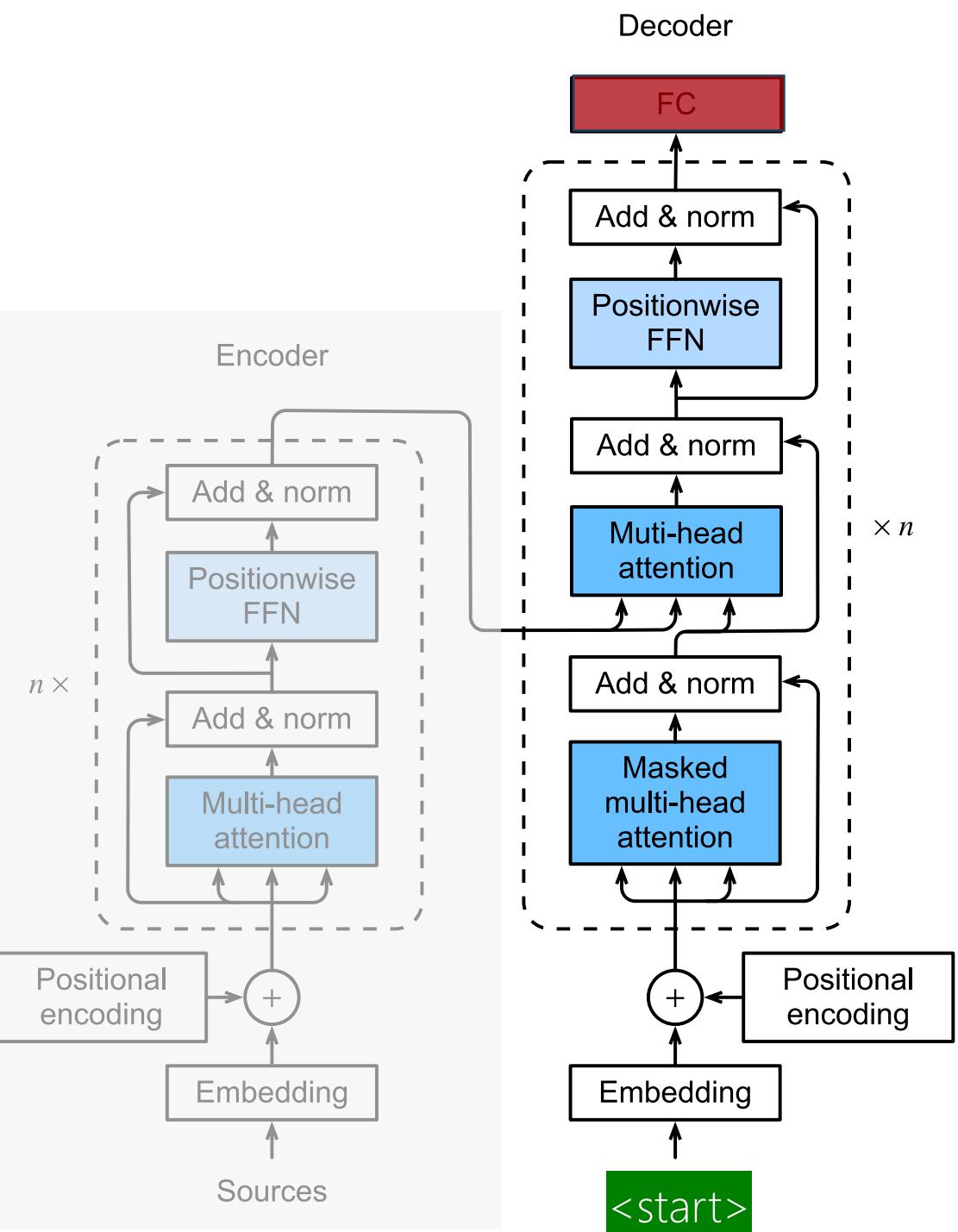






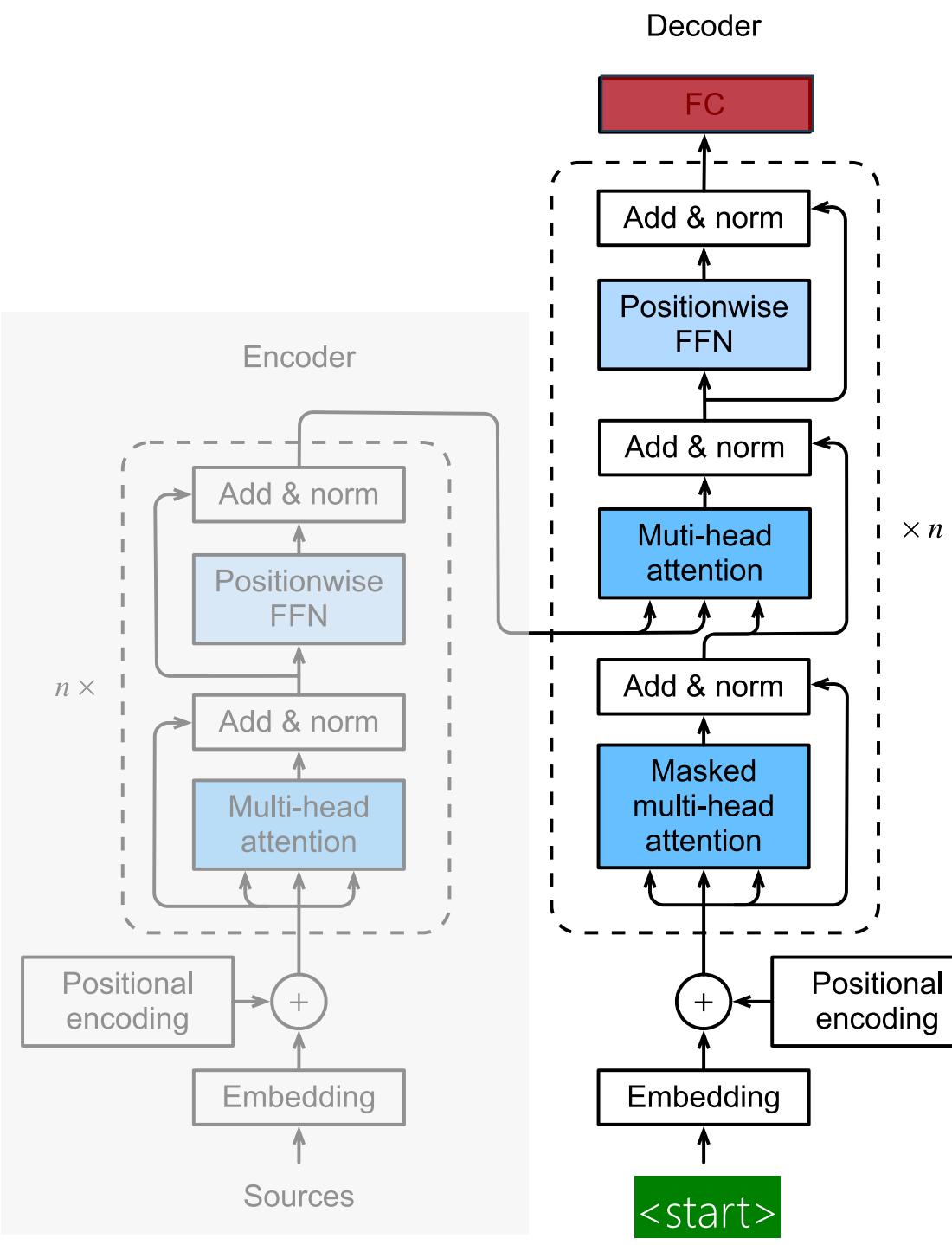






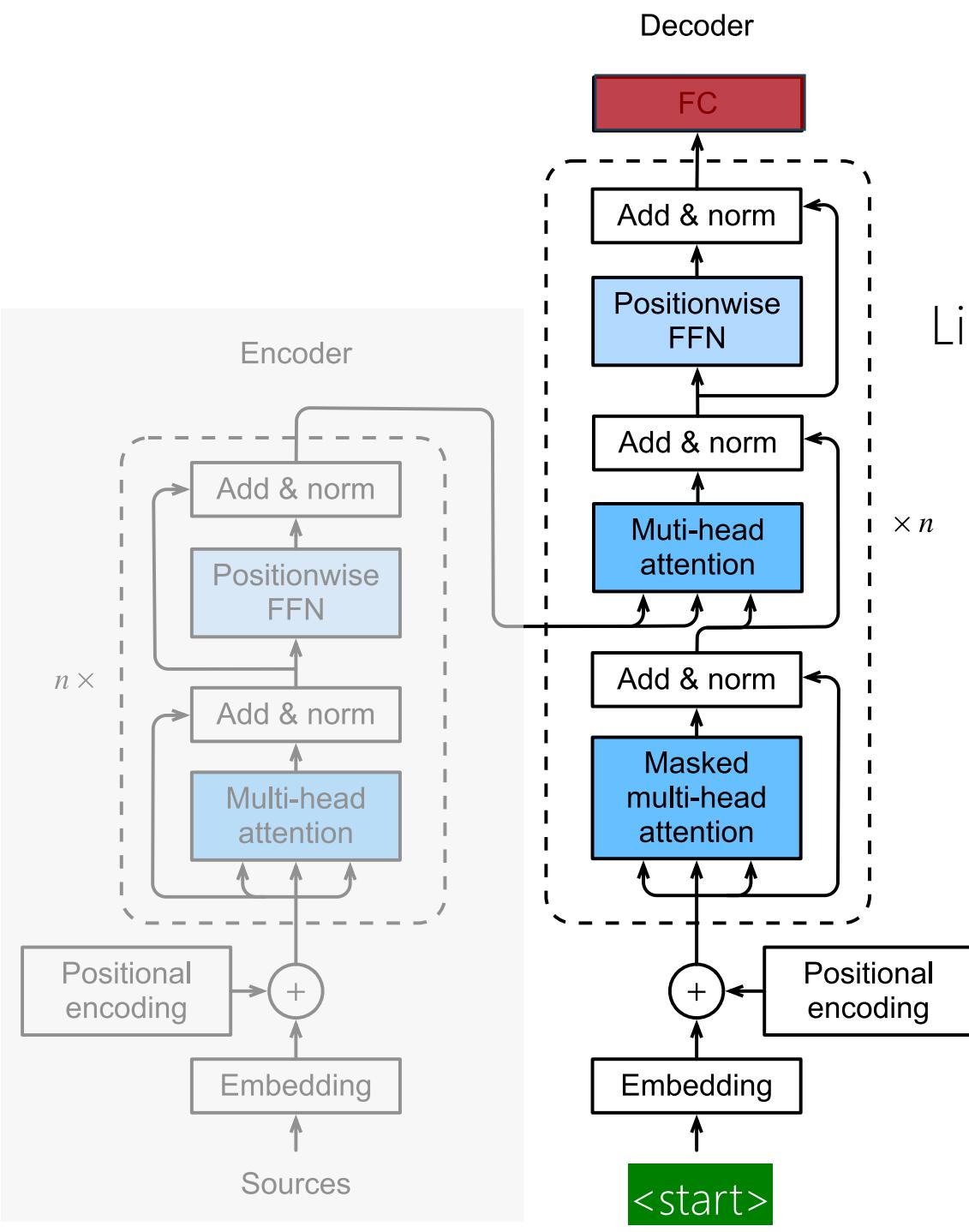
→ get probabilities
and select top word
but don't have
a vector NOT POSS

Squashing
into a
vector
happens
here



Wait a minute!!! Is this even possible?

No! Need a vector,
have a 7×500 matrix

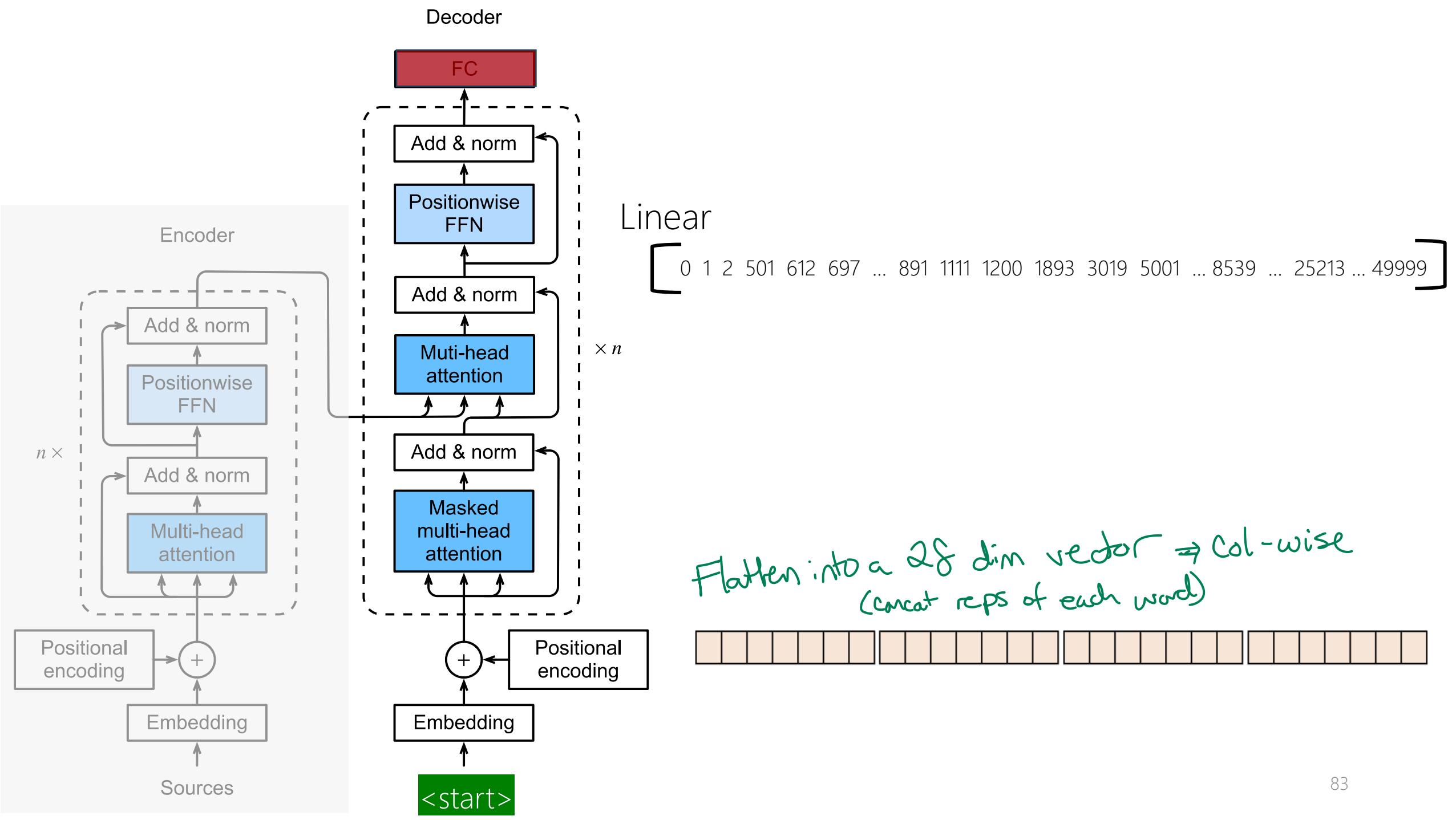


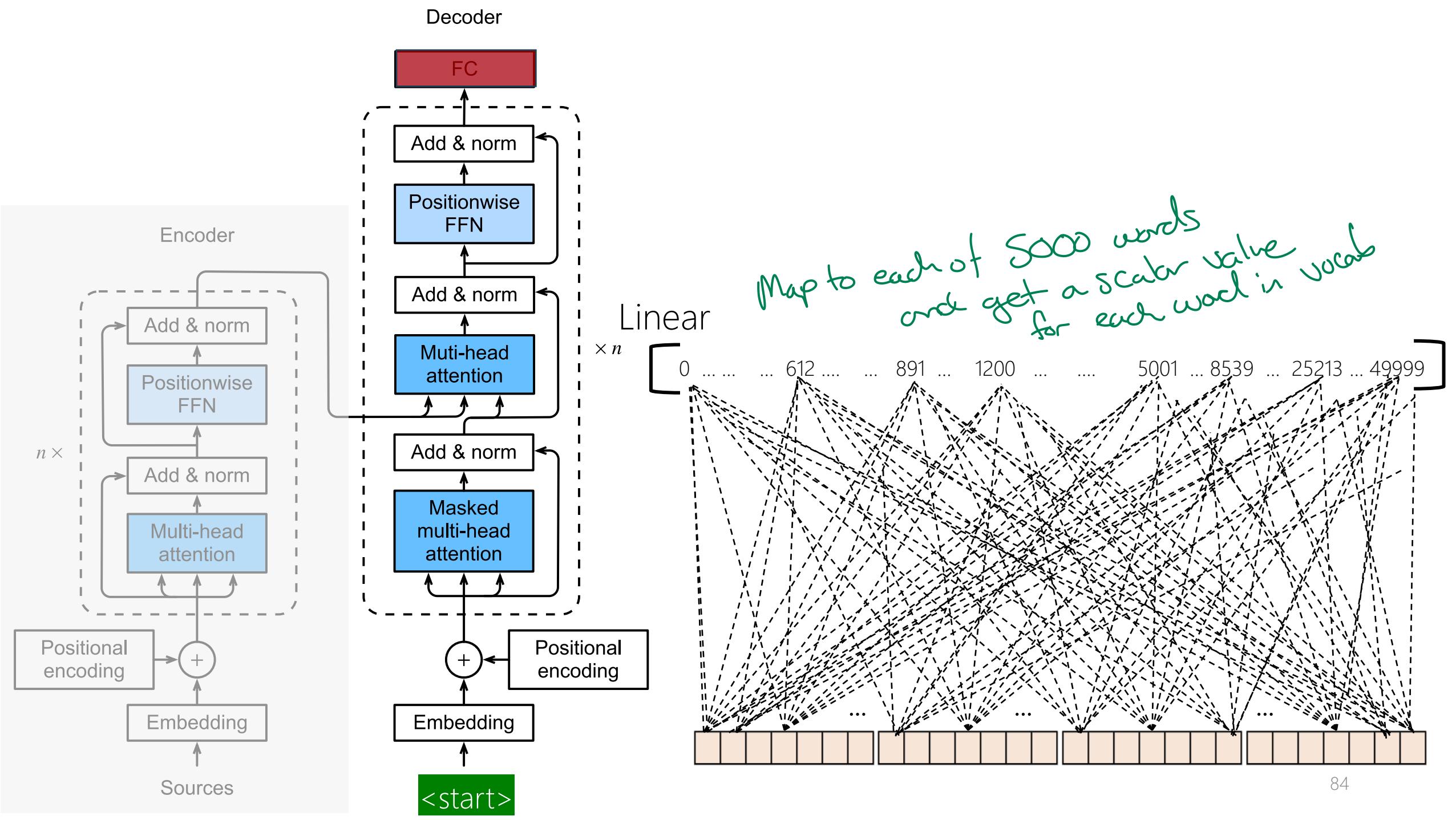
Linear

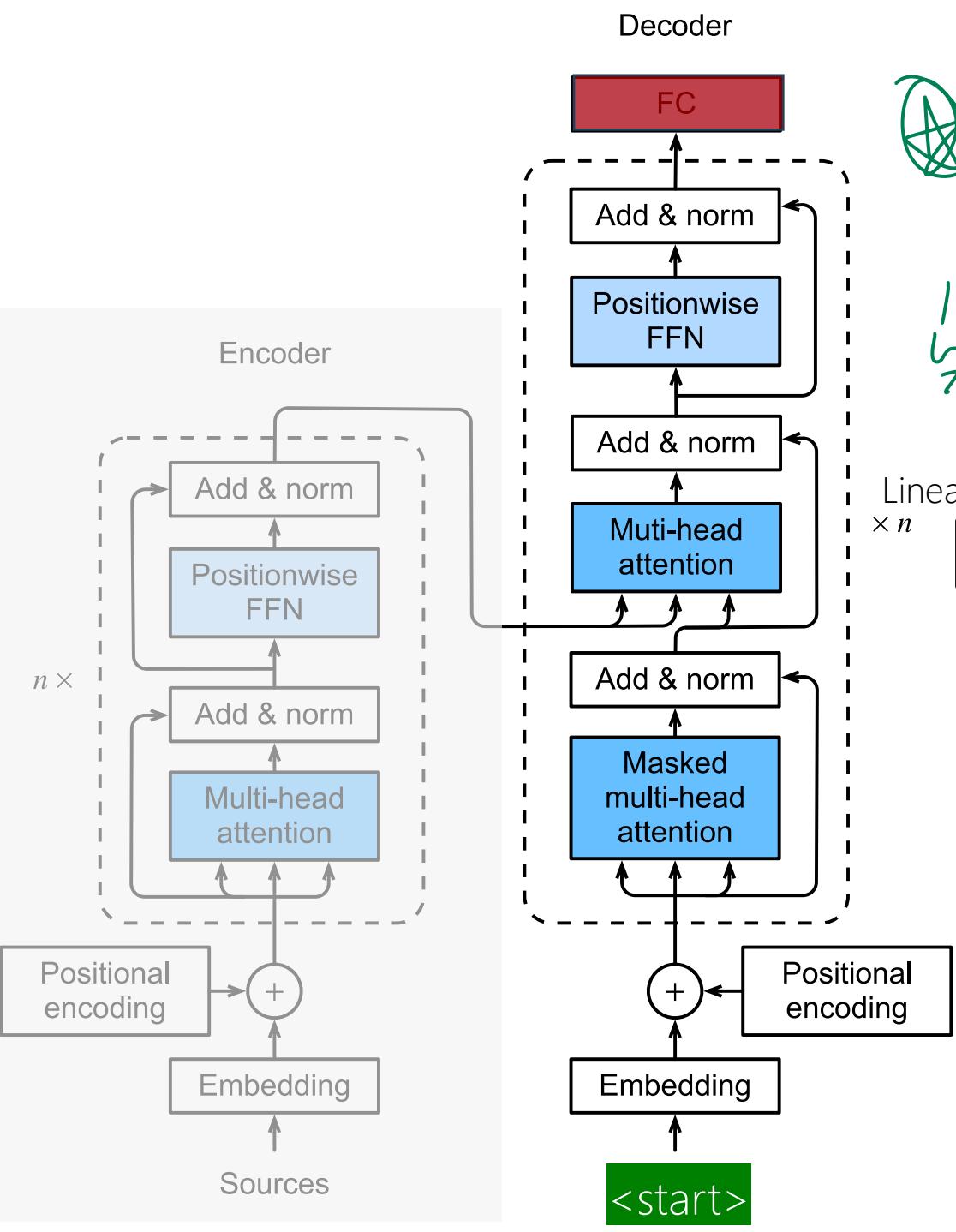
$\times n$

7x4

213 ... 49999
↑
5000
vocab



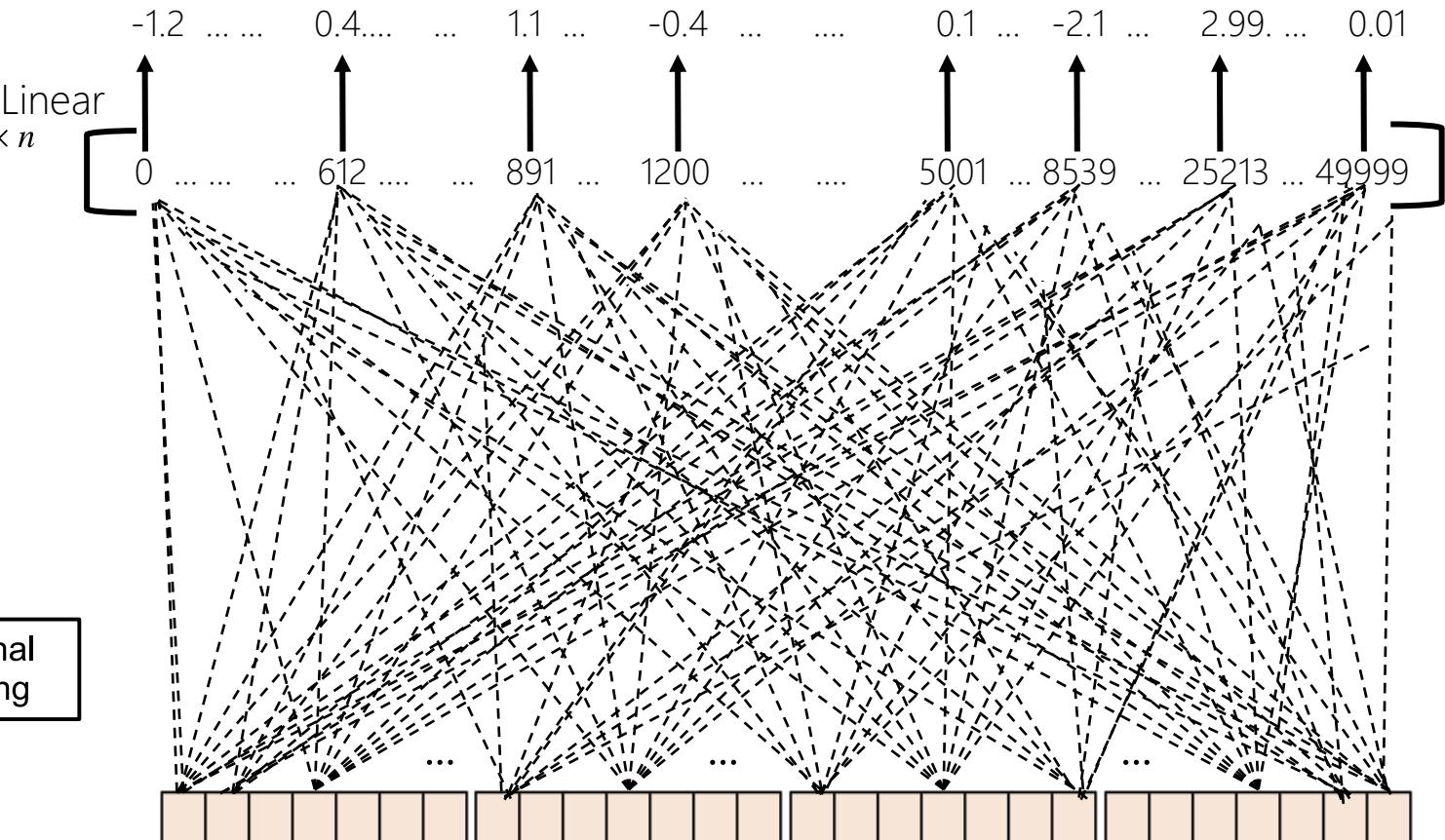


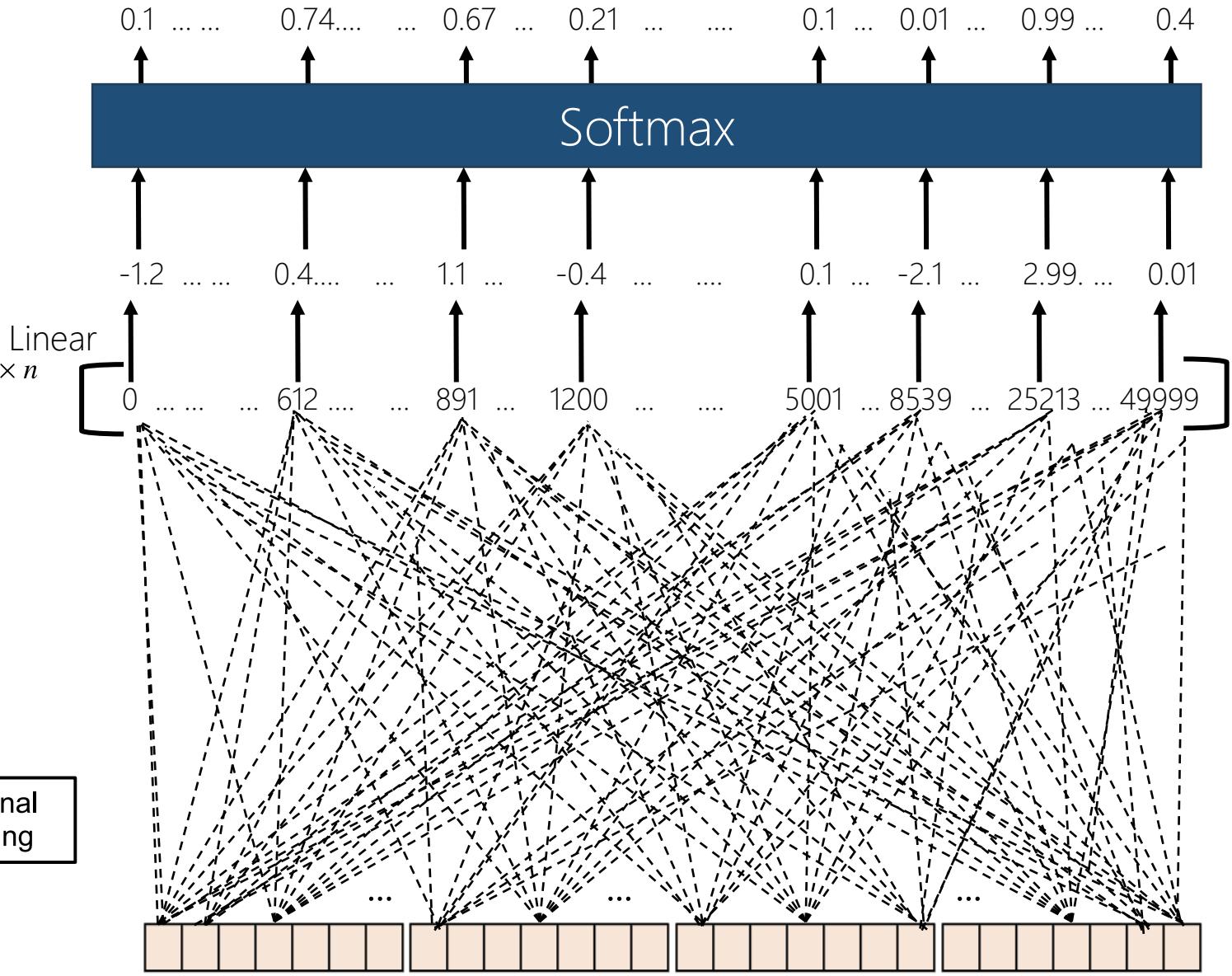
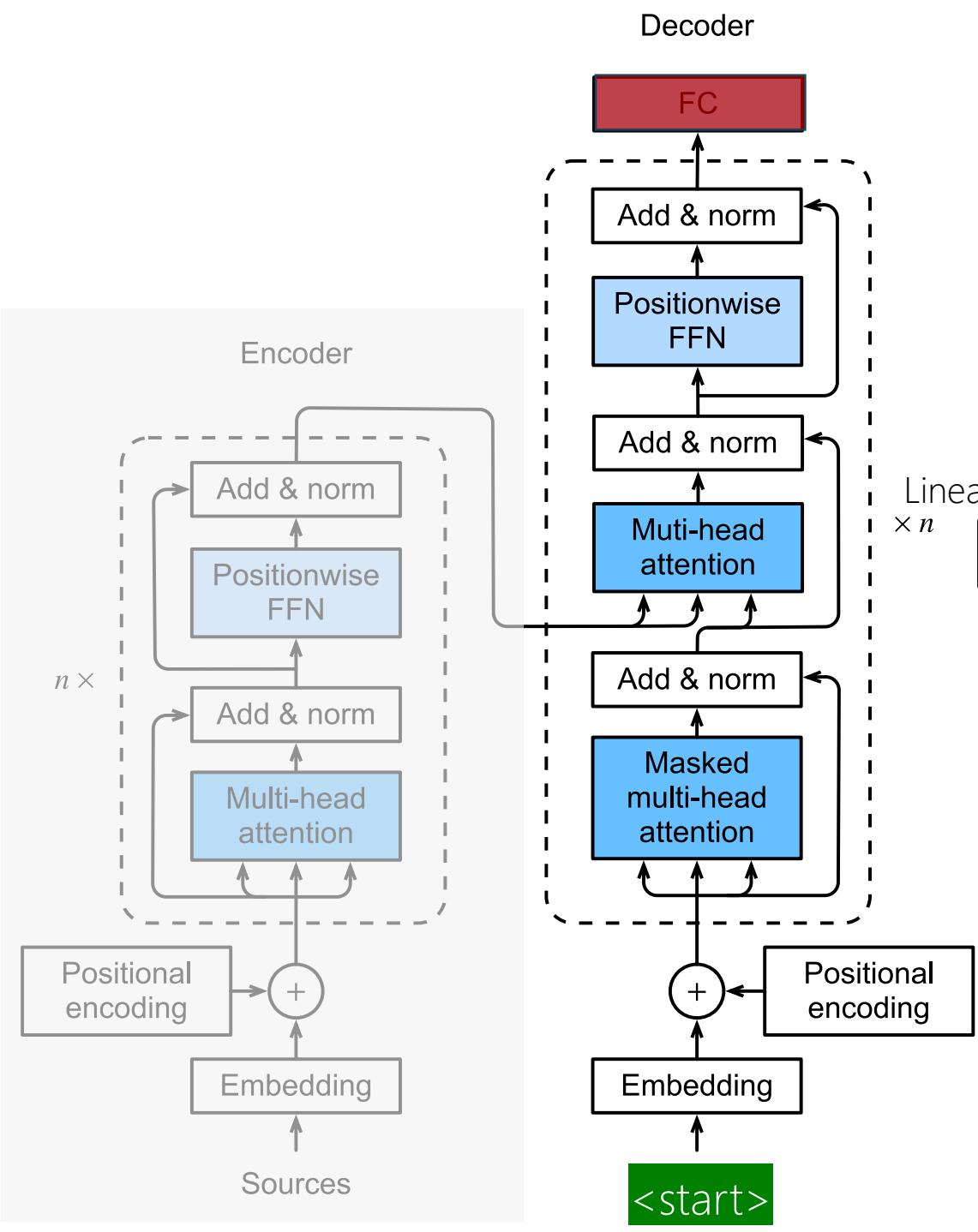


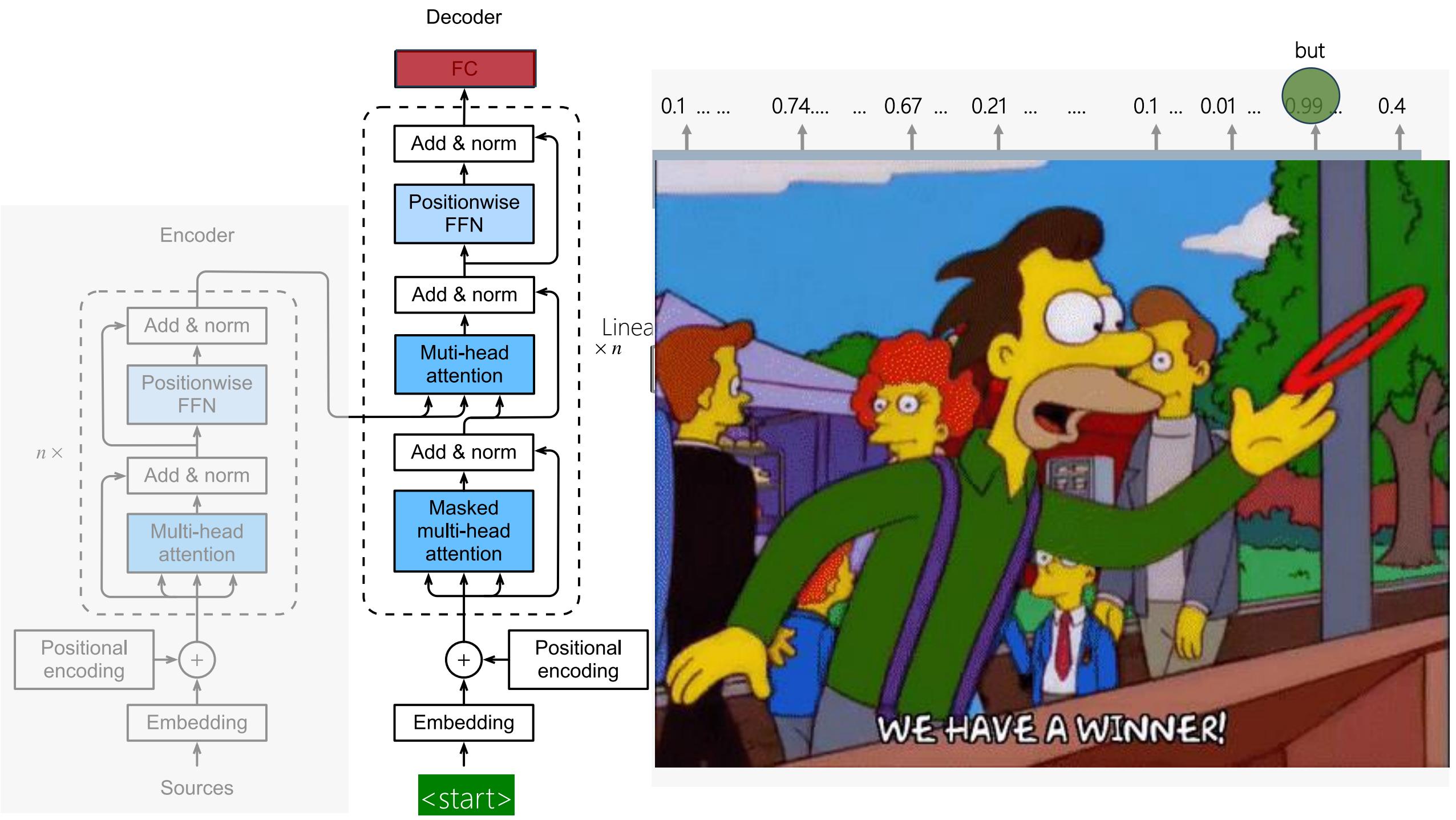
Missed part about what is mult to get scalar

$$1 \times 28 \cdot \underbrace{28 \times 50000}_{7 \times 9 \text{ swivied}} = 1 \times 5000 \text{ vector} - \text{what is this and when?}$$

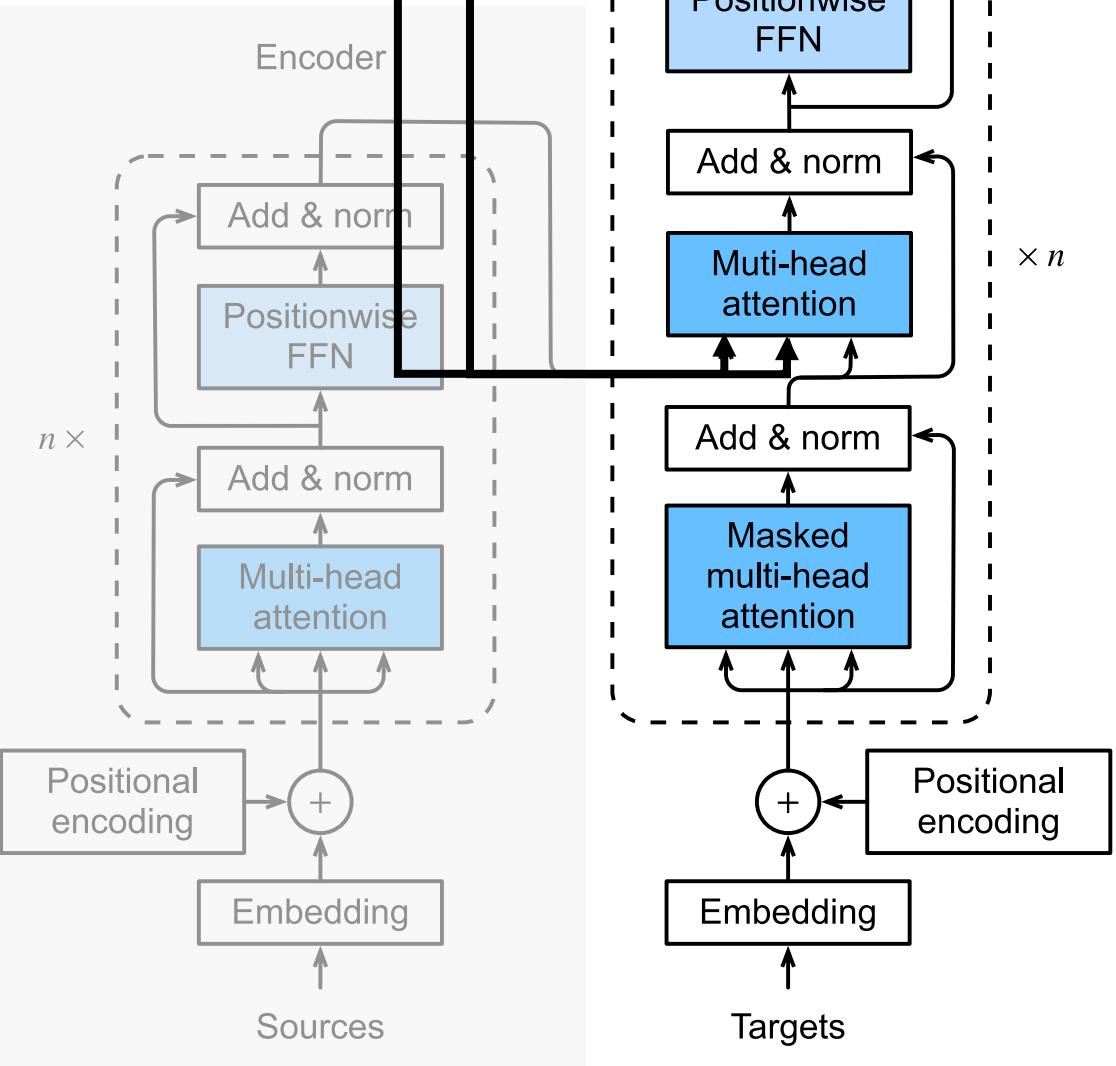
vocab



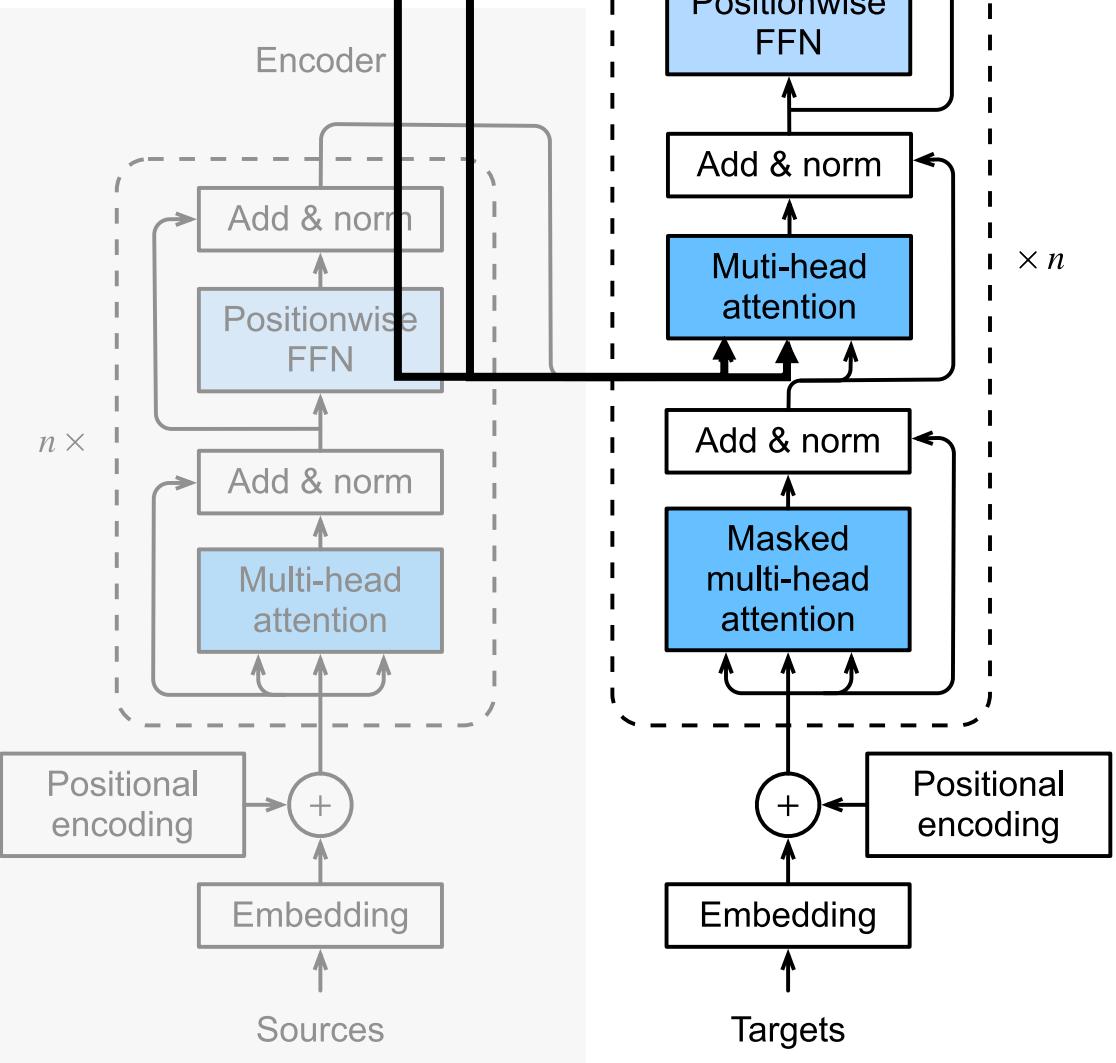




It	-0.453	0.380	-1.414	1.482
matters	-0.453	0.380	-1.414	1.482
not	-0.910	-0.943	1.325	-1.785
what	0.249	-1.927	-0.673	-0.843
someone	-1.404	0.662	-2.608	1.024
is	-0.095	-1.873	1.620	-0.511
born	0.119	-2.263	-0.119	-0.379
	-0.905	0.771	-0.217	-2.326



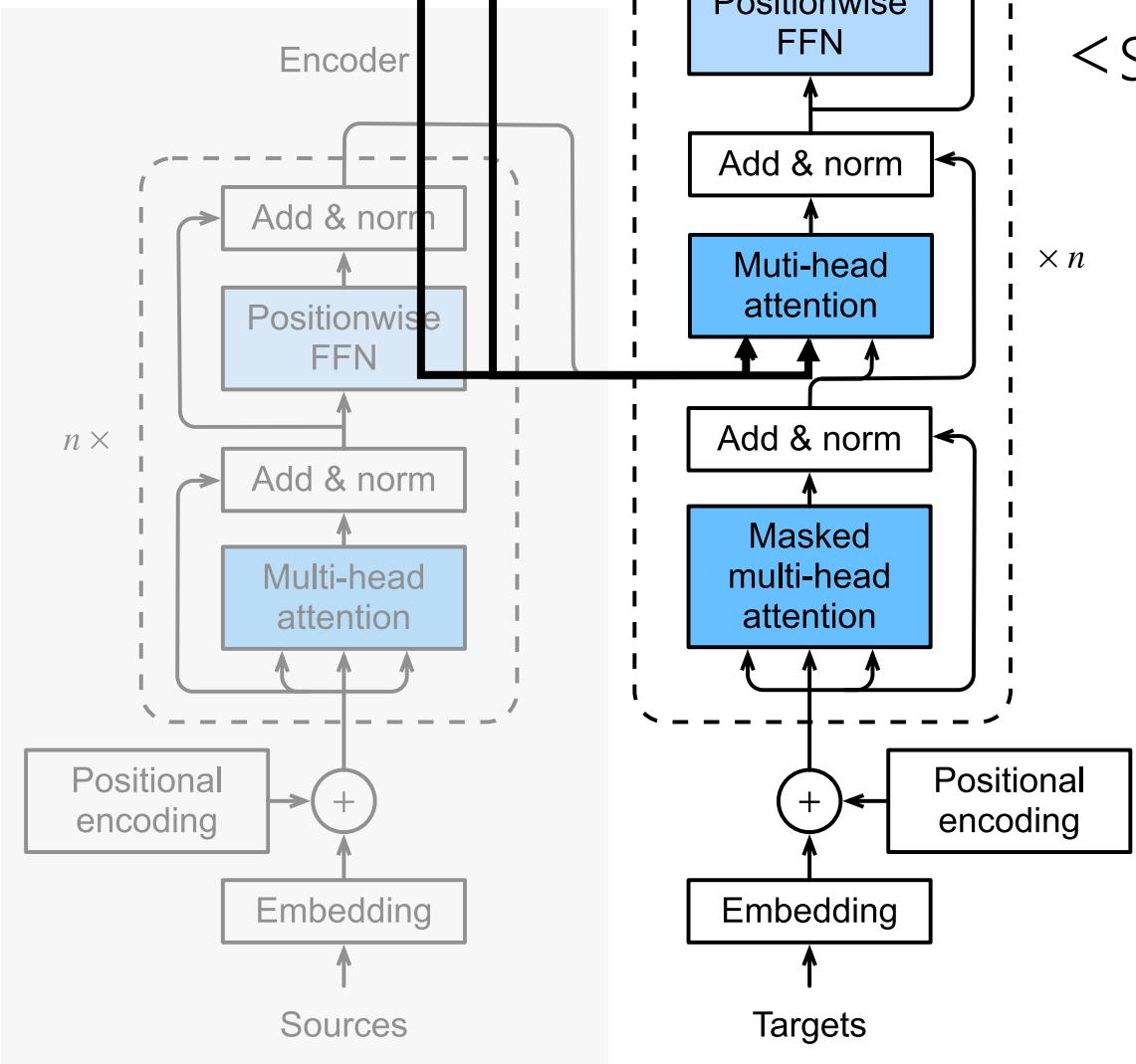
It	-0.453	0.380	-1.414	1.482
matters	-0.453	0.380	-1.414	1.482
not	-0.910	-0.943	1.325	-1.785
what	0.249	-1.927	-0.673	-0.843
someone	-1.404	0.662	-2.608	1.024
is	-0.095	-1.873	1.620	-0.511
born	0.119	-2.263	-0.119	-0.379
	-0.905	0.771	-0.217	-2.326



Targets

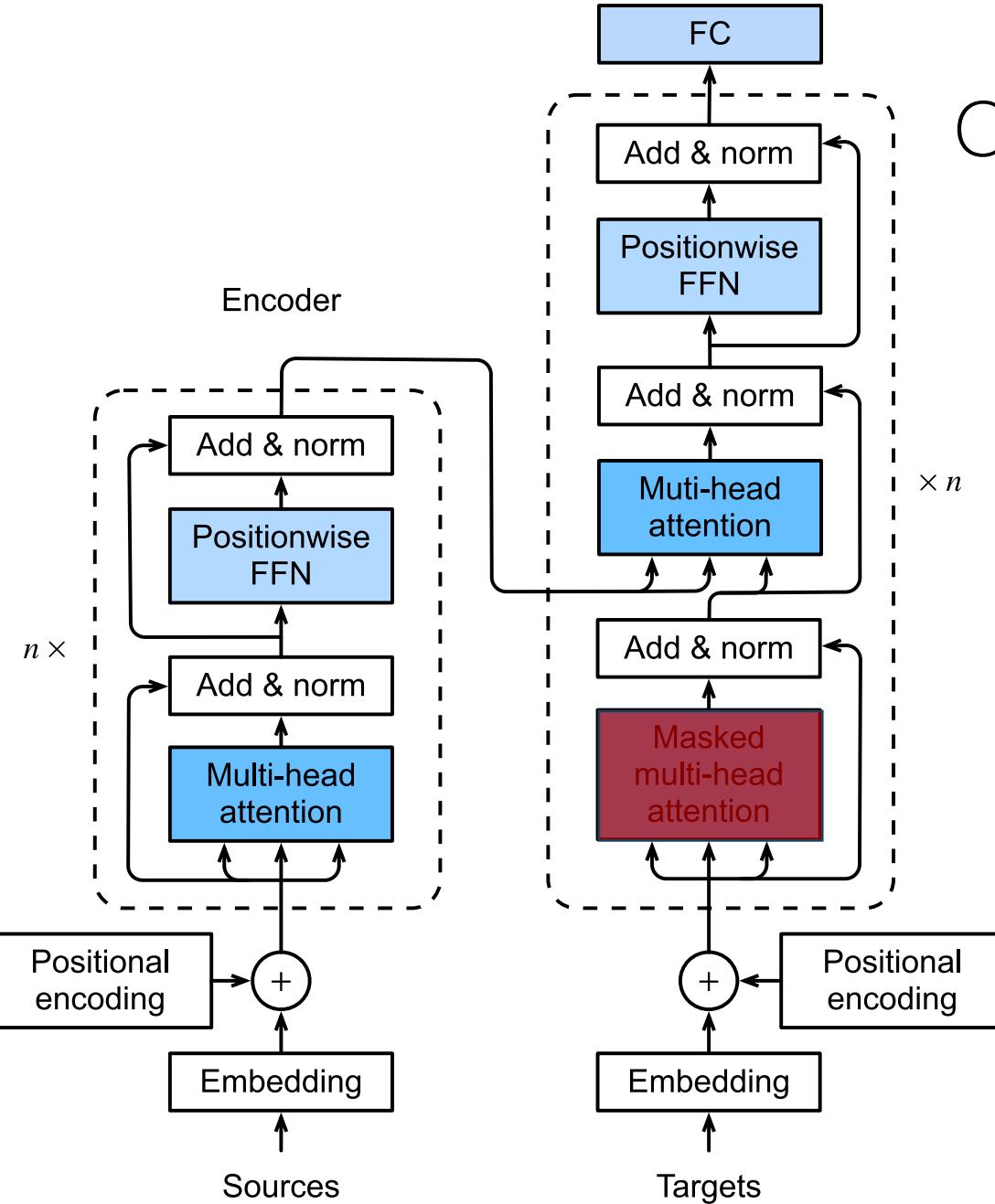
89

It	-0.453	0.380	-1.414	1.482
matters	-0.453	0.380	-1.414	1.482
not	-0.910	-0.943	1.325	-1.785
what	0.249	-1.927	-0.673	-0.843
someone	-1.404	0.662	-2.608	1.024
is	-0.095	-1.873	1.620	-0.511
born	0.119	-2.263	-0.119	-0.379
	-0.905	0.771	-0.217	-2.326



<start> but what they grow to be <end>

Decoder



Oh wait, we never discussed the masked multi-head attention!!!



Not showing the future and trying to influence the past???

The Dialogue Completer Task

Input Dialogue

It is our choices, Harry, that show what
we truly are,

If you want to know what a man's like,
take a good look at

It matters not what someone is born,

Dialogue Completion

?

?

?

The Dialogue Completer Task

Input Dialogue

It is our choices, Harry, that show what we truly are,

If you want to know what a man's like,
take a good look at

It matters not what someone is born,

Dialogue Completion

<start> far more than our abilities
<end>

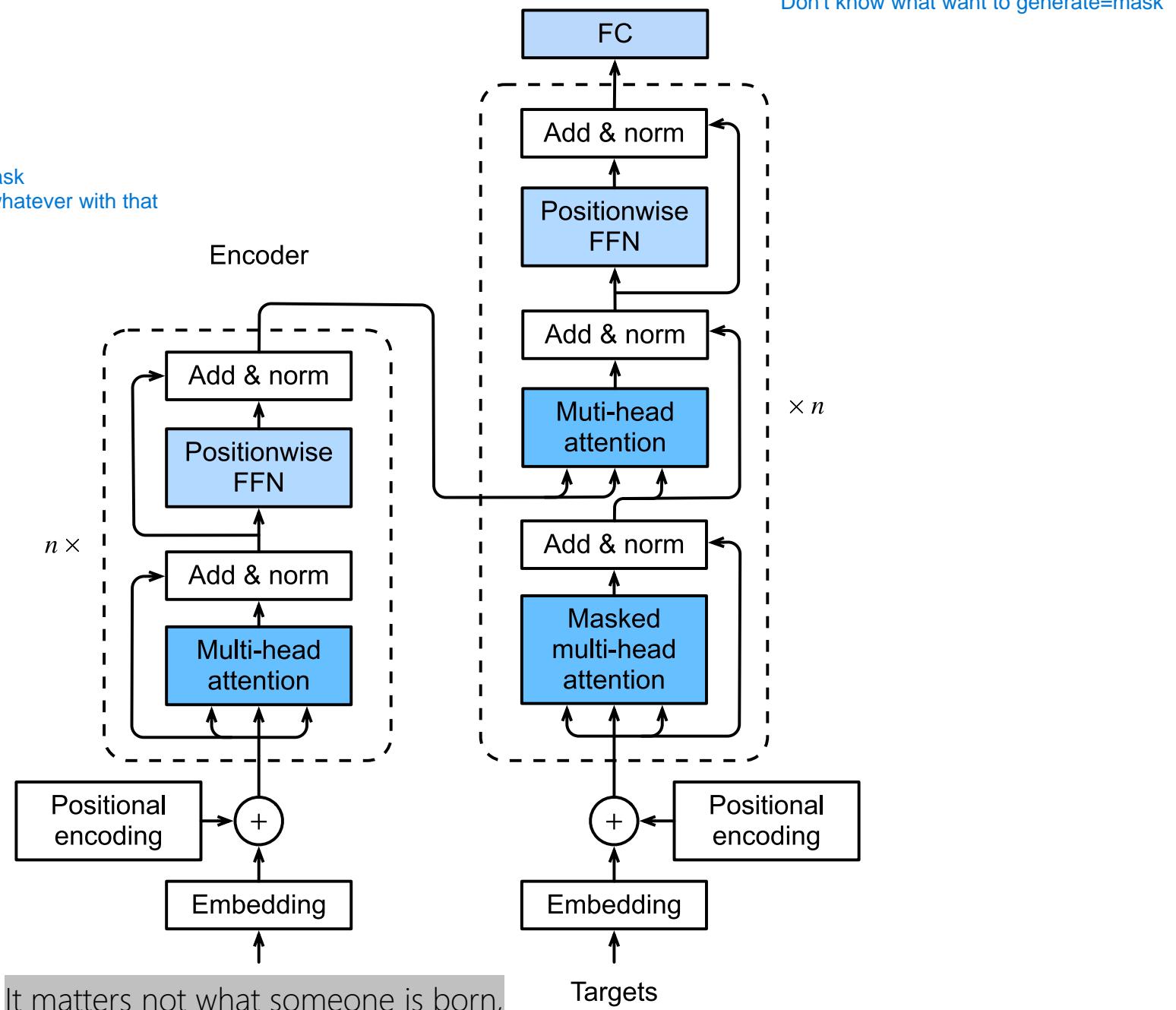
<start> how he treats his inferiors, not his equals <end>

<start> but what they grow to be <end>

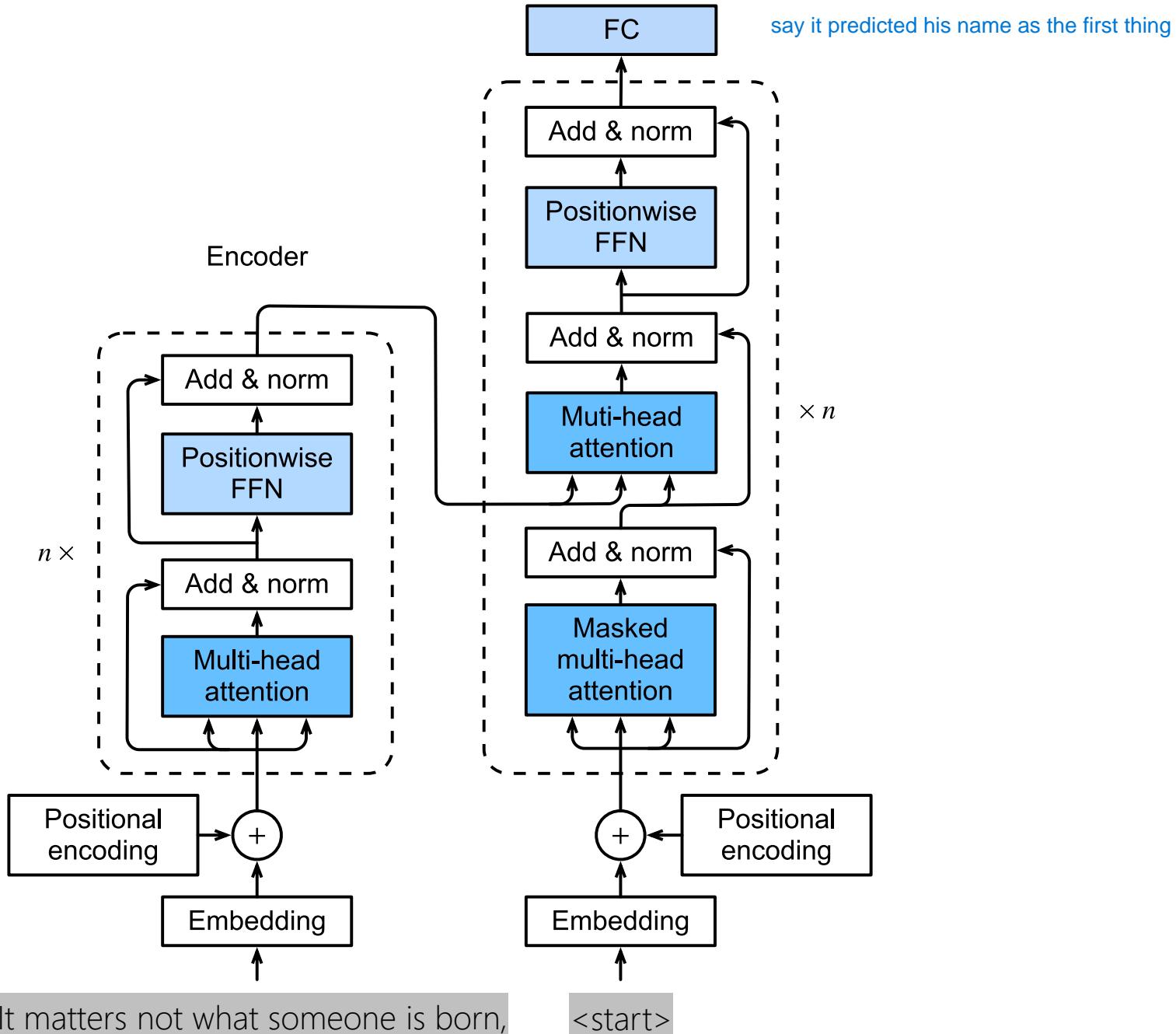
Training

Let the model just learn from the data not the task

Analogy: give toys to a child and child can do whatever with that



Training



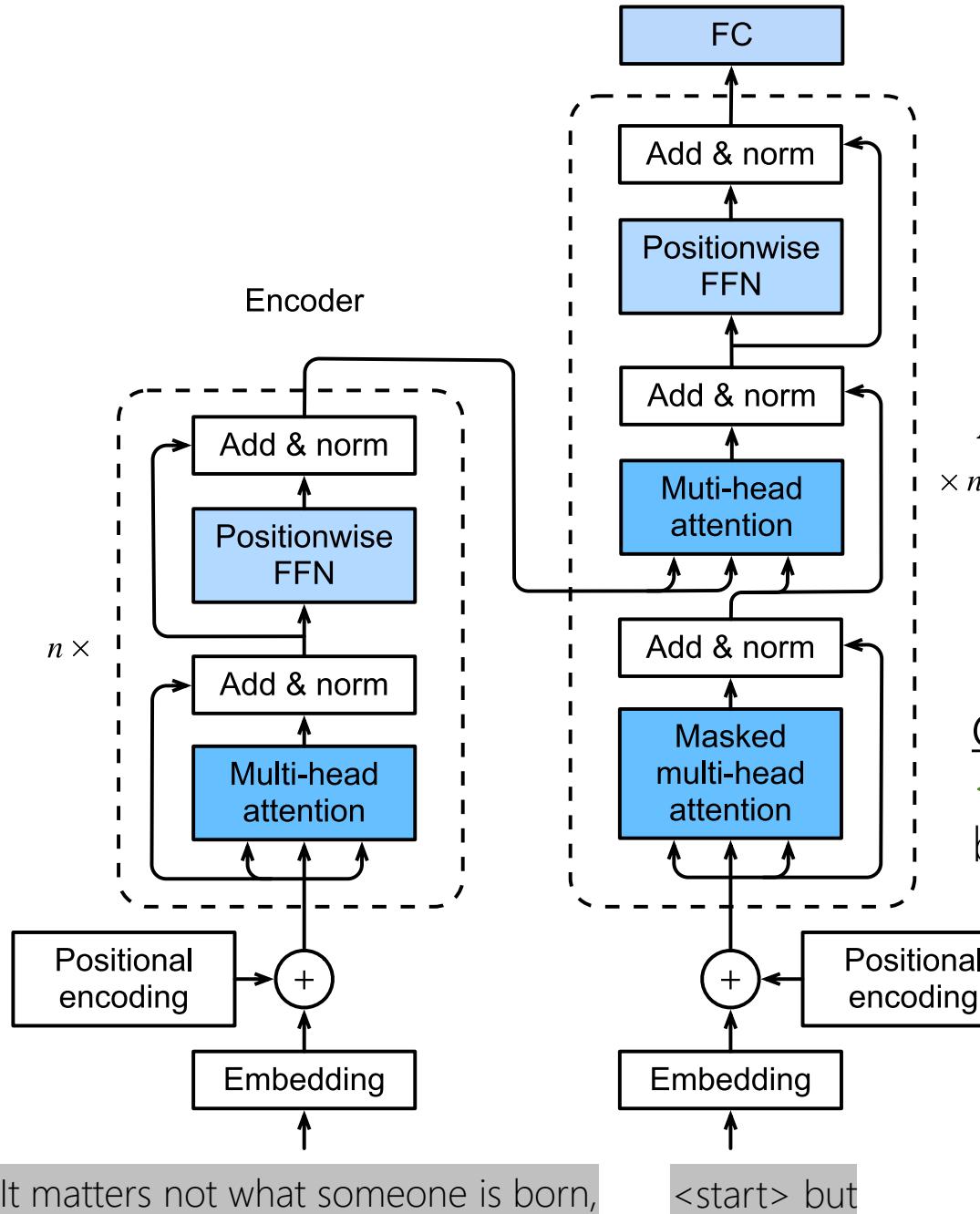
It matters not what someone is born,

<start>

Training

LLM is not generative, it is predictive, just trying to predict the next word

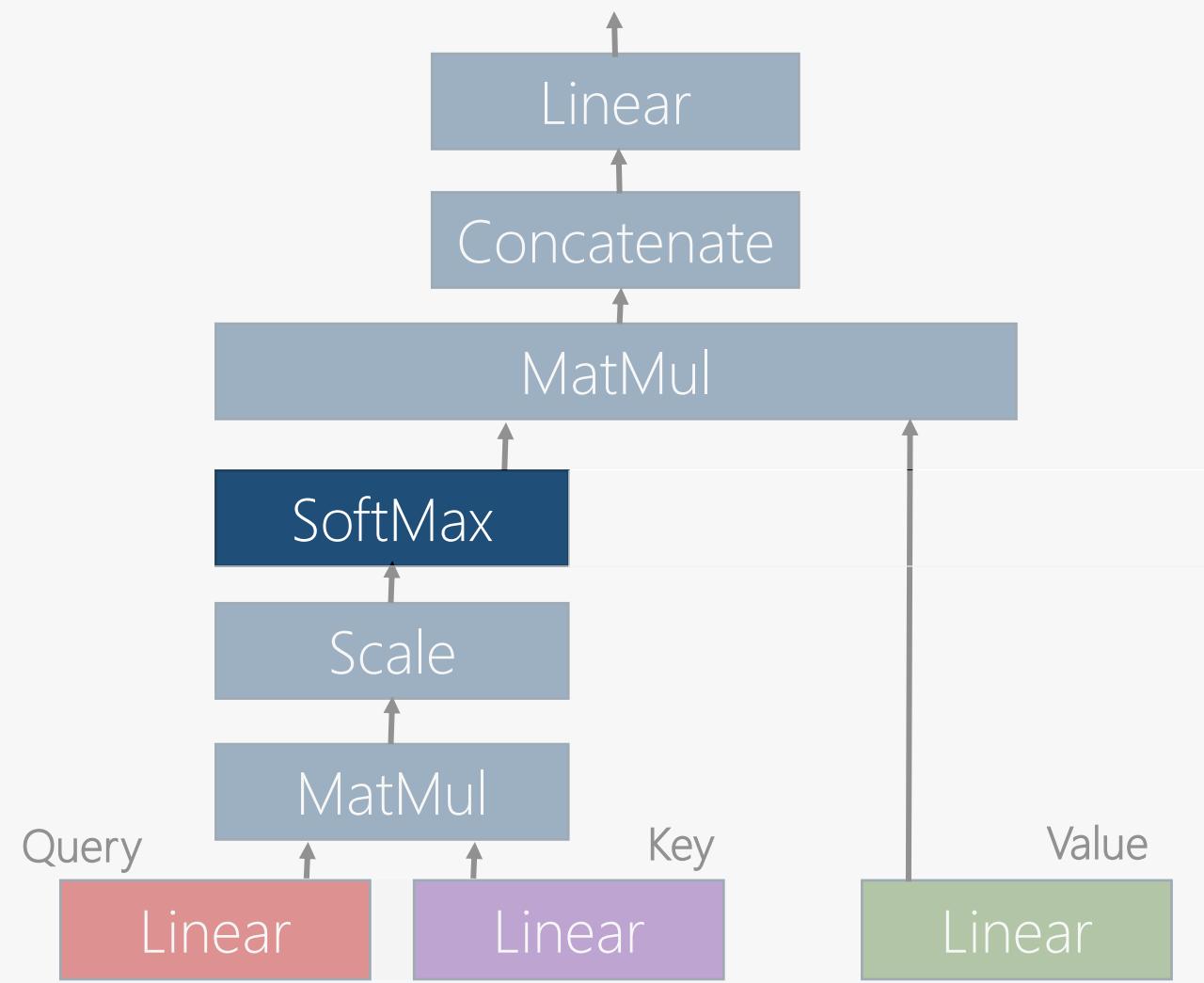
Every word is a class, a unique token



It matters not what someone is born,

<start> but

Masked Attention



Self attention was taking the few same word and

	<start>	but	what	they	grow	to	be	<end>
but	92	35	54	11	39	91	58	7
what	20	21	67	47	13	61	62	3
they	94	54	76	85	39	49	0	58
grow	51	53	72	69	97	46	94	32
to	8	39	22	85	66	95	7	27
be	1	77	5	73	41	20	50	36
<end>	21	90	3	7	92	69	56	97
	91	68	0	56	77	59	81	28

Masked Attention

If just tell the

	<start>	but	what	they	grow	to	be	<end>
<start>	92	35	54	11	39	91	58	7
but	20	21	67	47	13	61	62	3
what	94	54	76	85	39	49	0	58
they	51	53	72	69	97	46	94	32
grow	8	39	22	85	66	95	7	27
to	1	77	5	73	41	20	50	36
be	21	90	3	7	92	69	56	97
<end>	91	68	0	56	77	59	81	28

Attention Filter

soft max includes e^x so soft max of that matrix is 1 if 0 and 0 if -inf

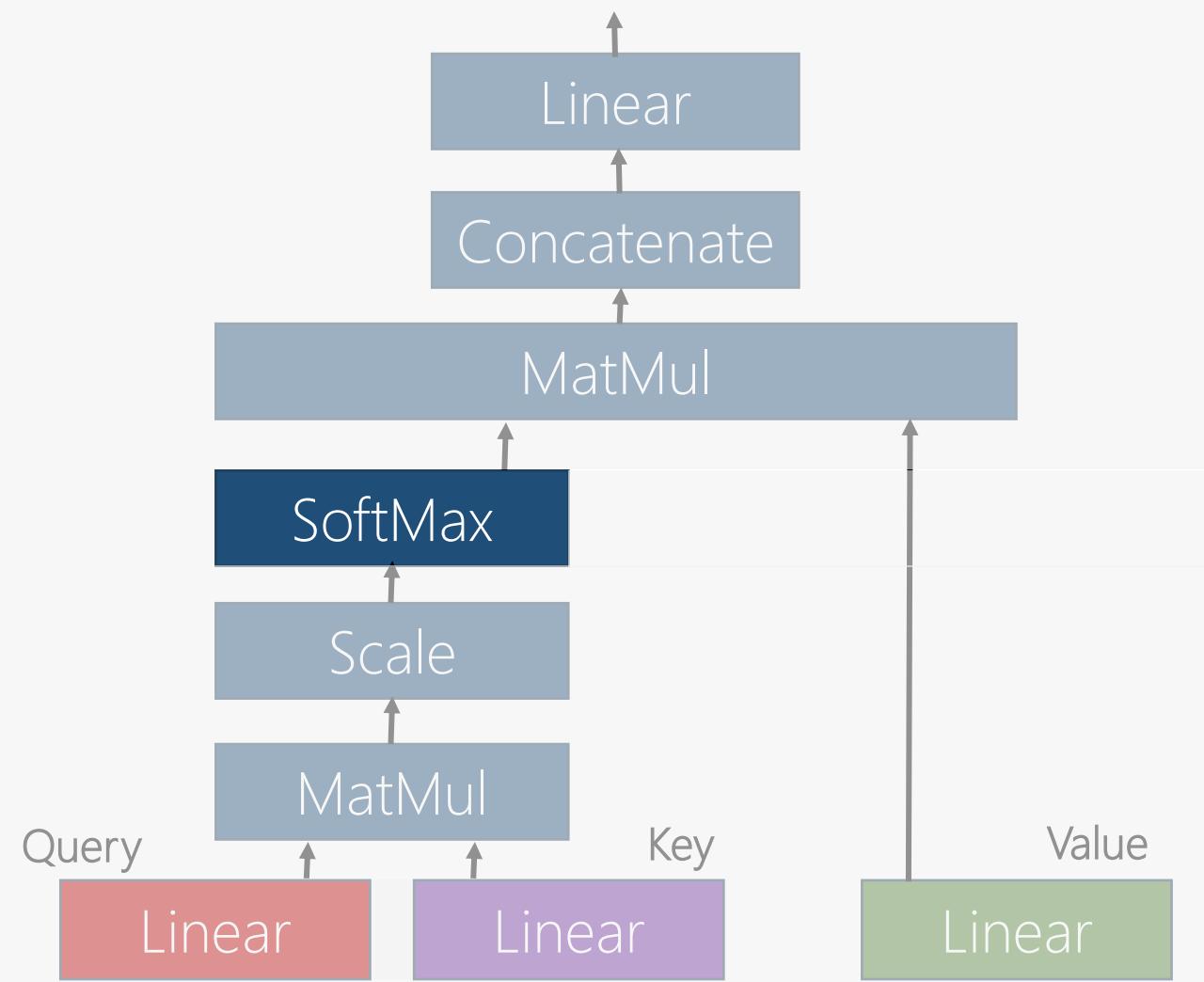
	<start>	but	what	they	grow	to	be	<end>
<start>	0	-inf						
but	0	0	-inf	-inf	-inf	-inf	-inf	-inf
what	0	0	0	-inf	-inf	-inf	-inf	-inf
they	0	0	0	0	-inf	-inf	-inf	-inf
grow	0	0	0	0	0	-inf	-inf	-inf
to	0	0	0	0	0	0	-inf	-inf
be	0	0	0	0	0	0	0	-inf
<end>	0	0	0	0	0	0	0	0

+

Mask Filter

-inf is no information from other words that they want to learn like word but is not getting info except from the start token and itself, not giving it any future info, just itself (or the previous words?????)

Masked Attention



	<start>	but	what*	they	grow	to	be	<end>
<start>	1	0	0	0	0	0	0	0
but	0.27	0.73	0	0	0	0	0	0
what	1	0	0	0	0	0	0	0
they	0	0	0.95	0.05	0	0	0	0
grow	0	0	0	1	0	0	0	0
to	0	0.98	0	0.02	0	0	0	0
be	0	0.12	0	0	0.88	0	0	0
<end>	1	0	0	0	0	0	0	0

Masked Attention Filter

with the soft max,

THE pt is to make sure the training happens correctly

Next Week!!

- Different Training Stages of LLMs

