

# Introduction to Transformers

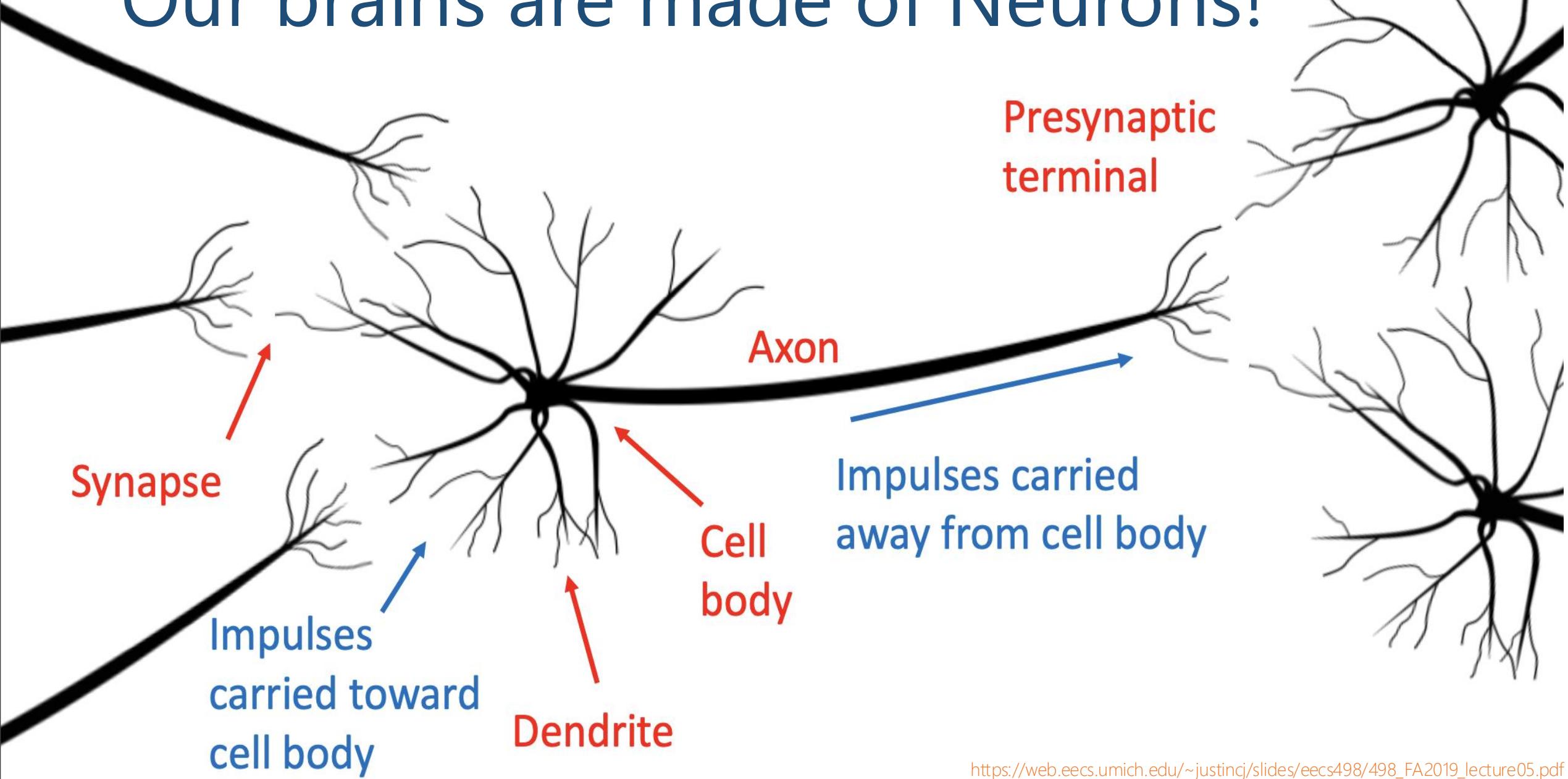
Chirag Agarwal

Assistant Professor

School of Data Science

University of Virginia

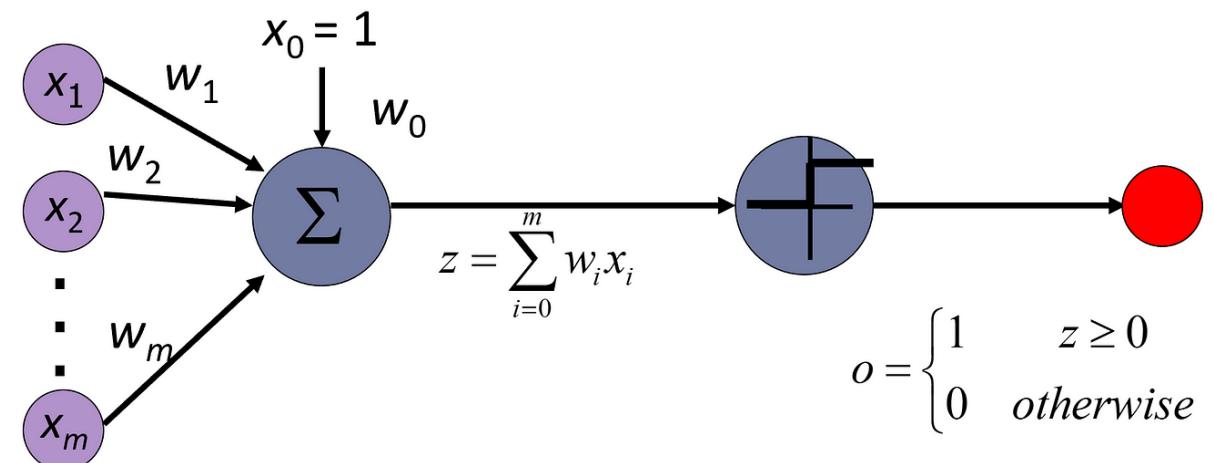
# Our brains are made of Neurons!



# Biological Neuron

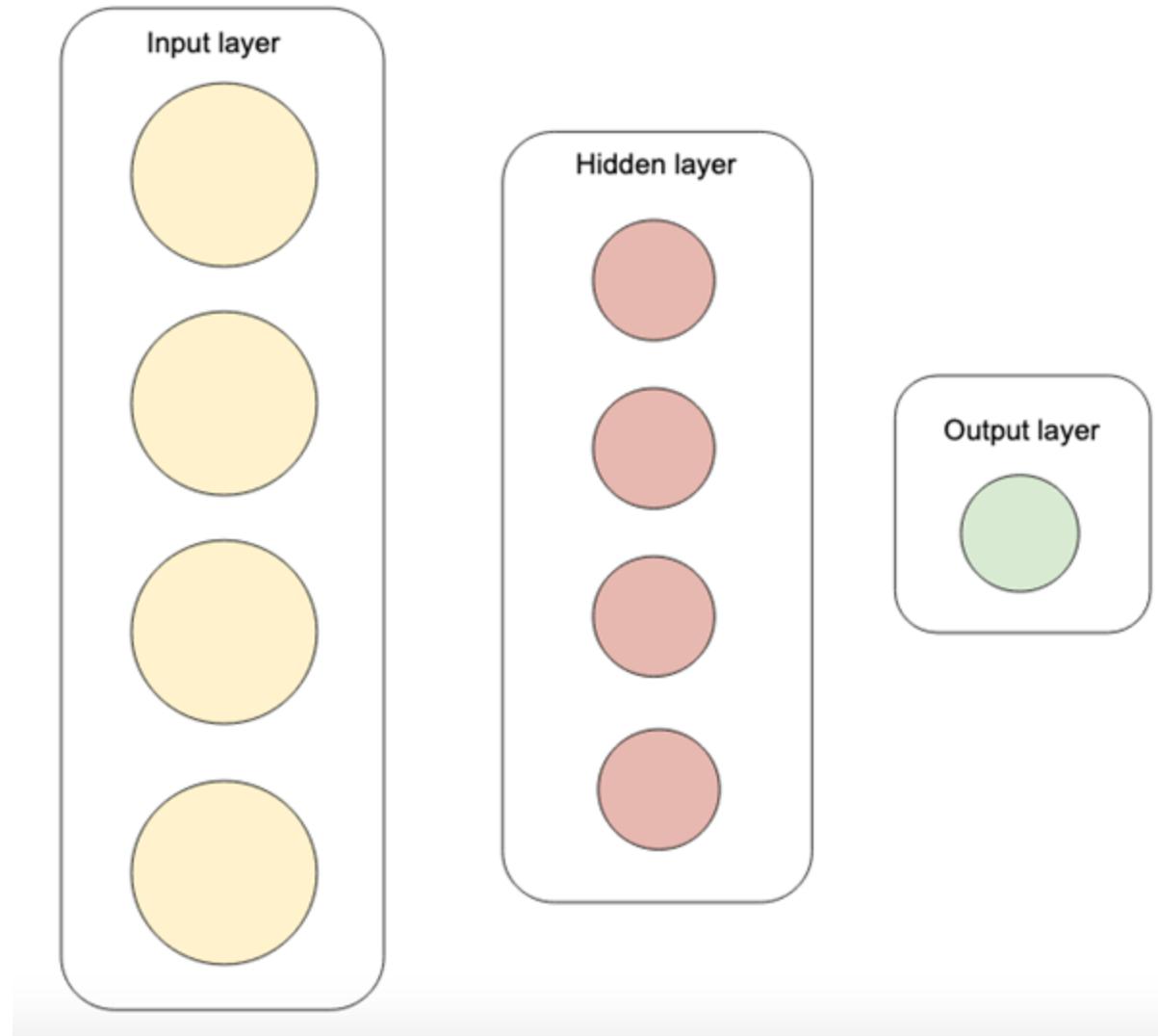


# Artificial Neuron

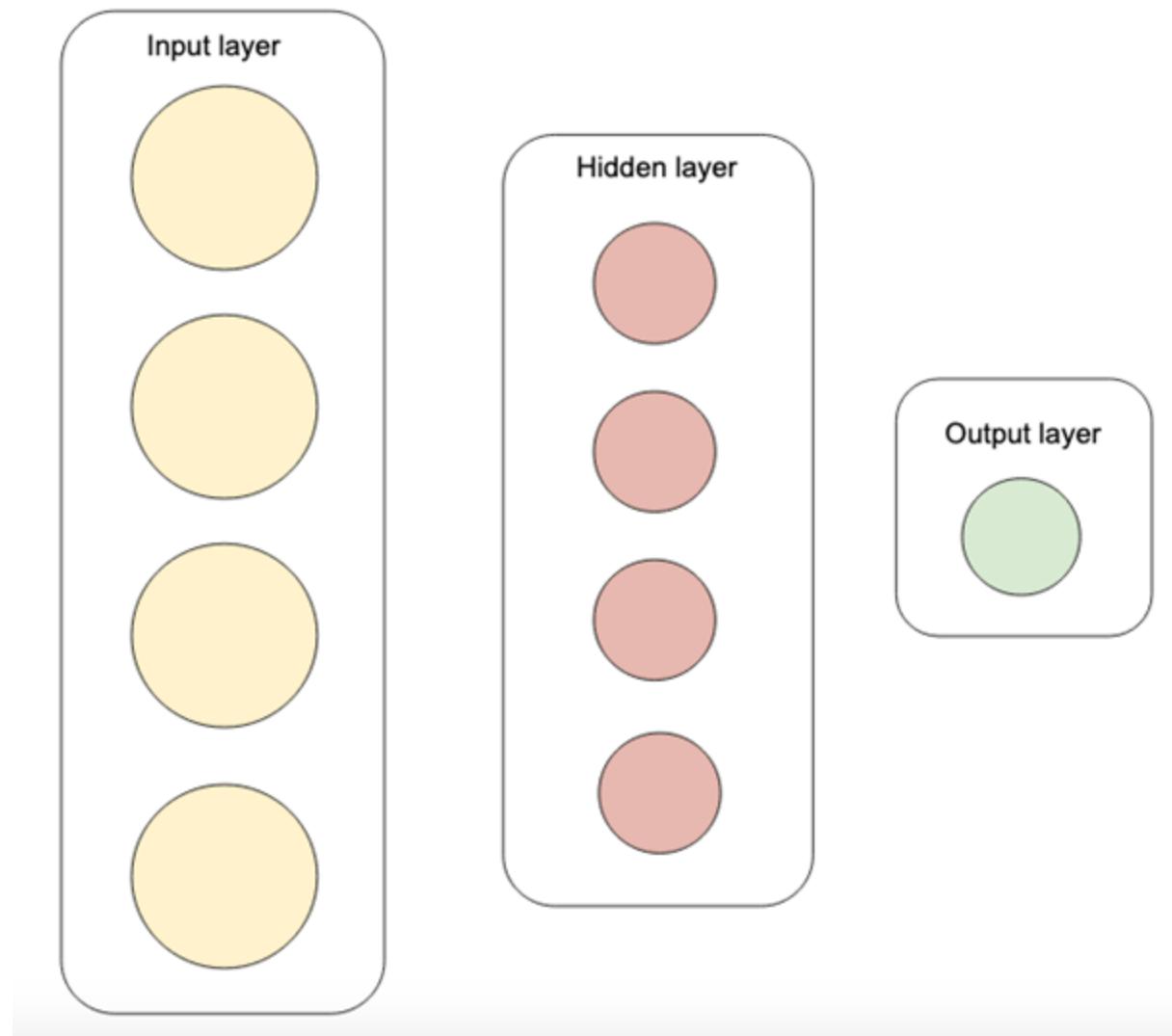


**Let's play  
a Game!**

- We have 9 players arranged in three layers
- There is an input layer, a hidden layer, and an output layer.

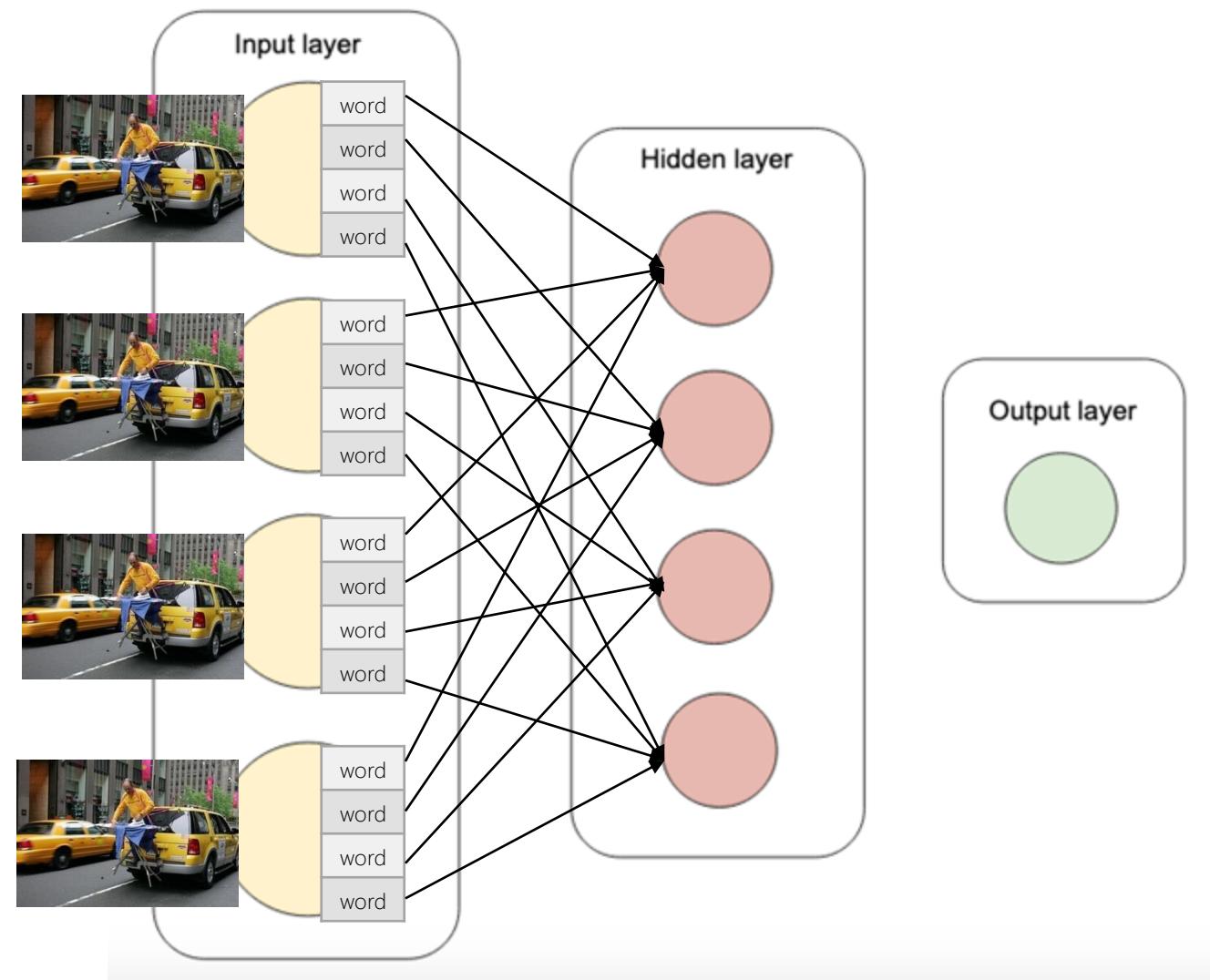


# Step 1:



# Step 1:

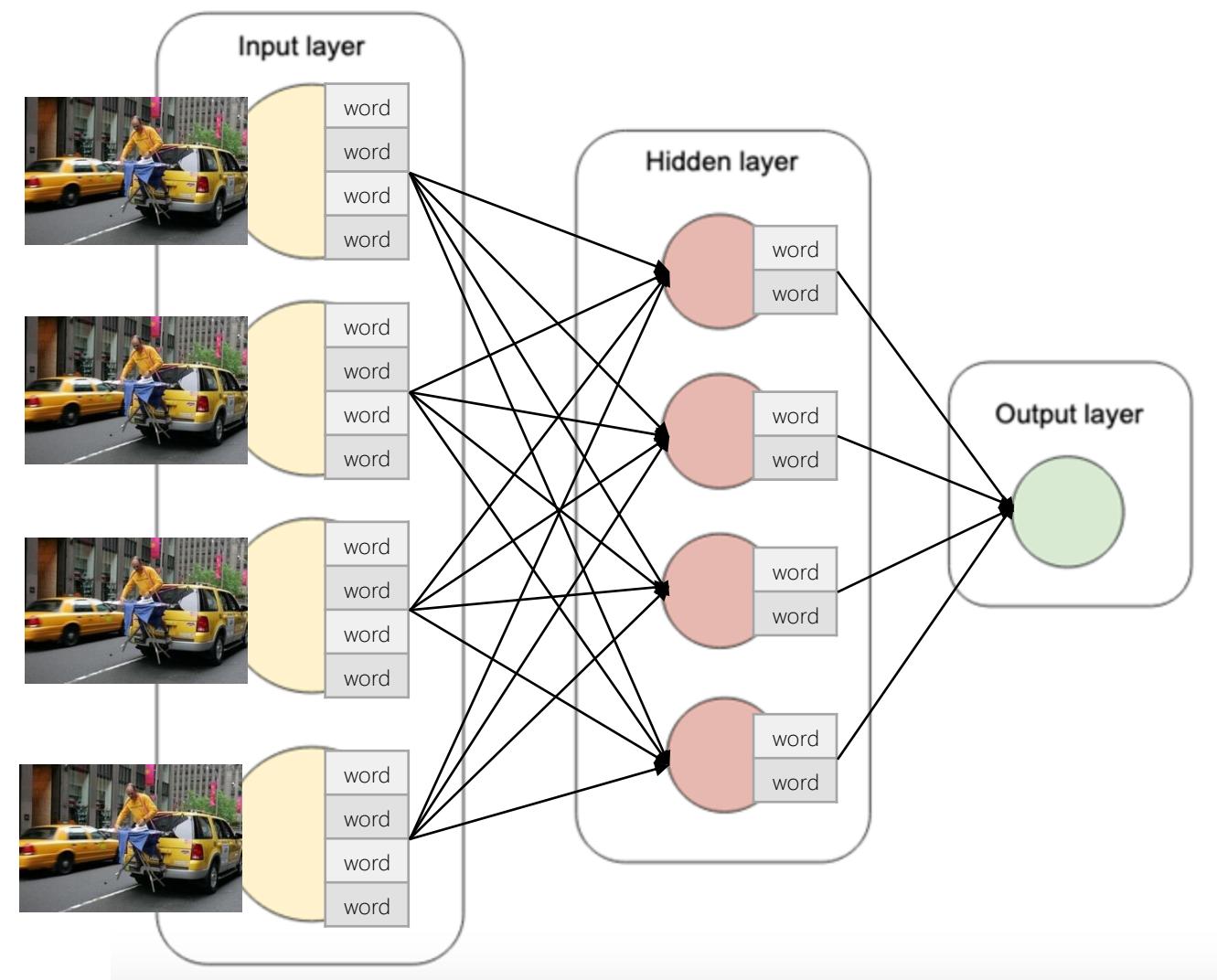
- players in the Input layer receives the image and write down 1 word on each slot (4 slots each)
- send 1 word to each of the 4 players in the hidden layer



# Step 2:

Players in the Hidden layer:

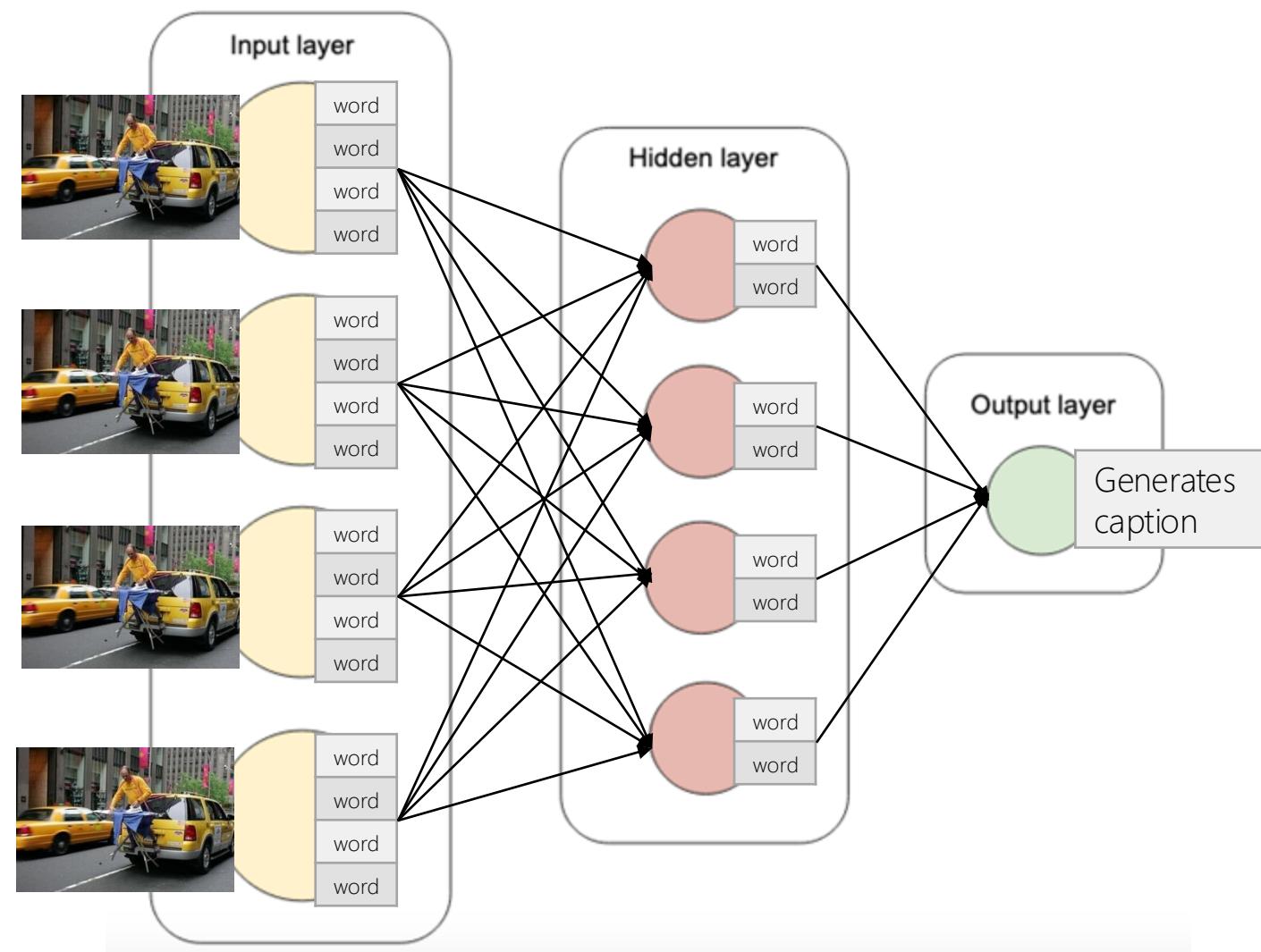
- receive 4 words each
- select 2 words to pass to the output node



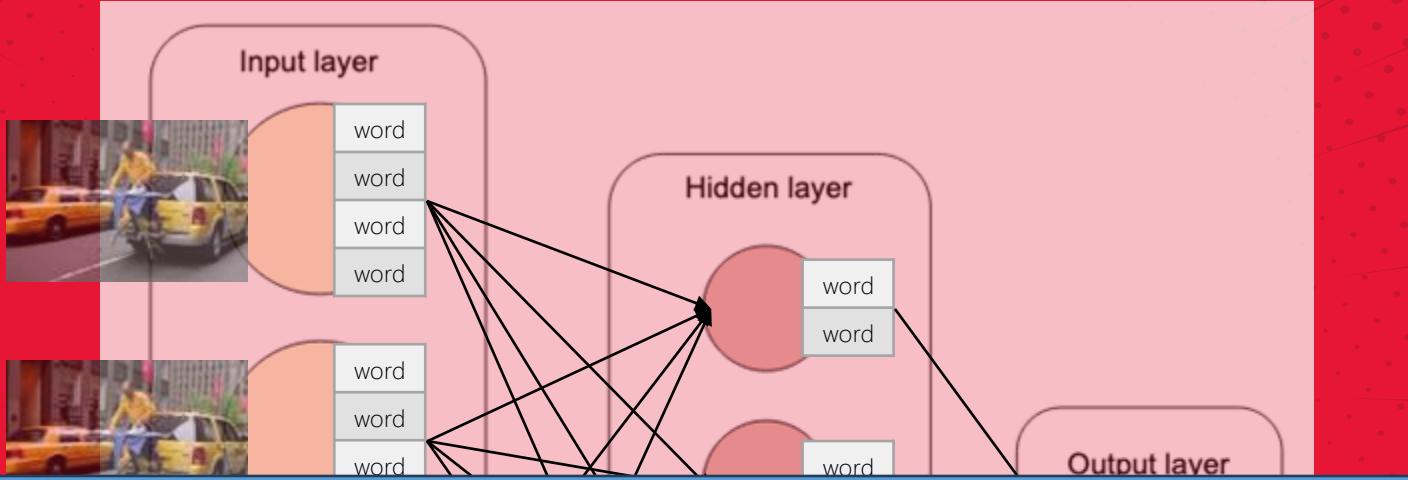
# Step 3:

The output node:

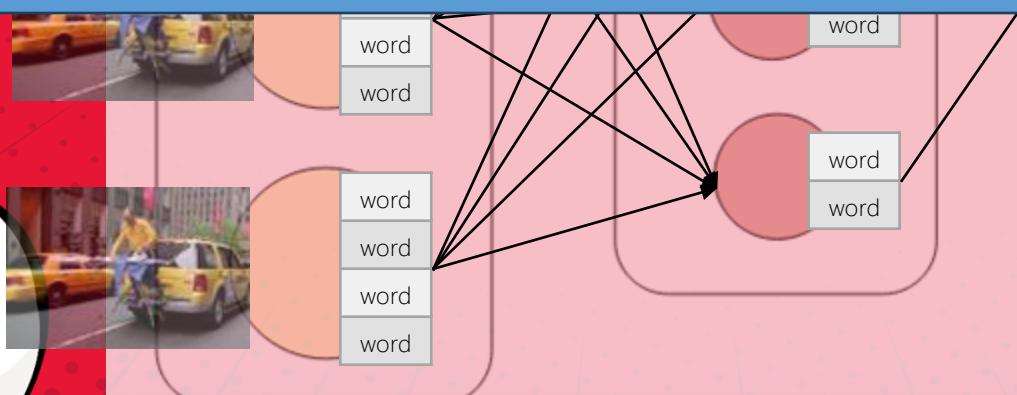
- receive 8 different words
- Picks 4 words to make a caption for an image with some filler words



# You just learned what is Feedforward!!



Feeding Forward!



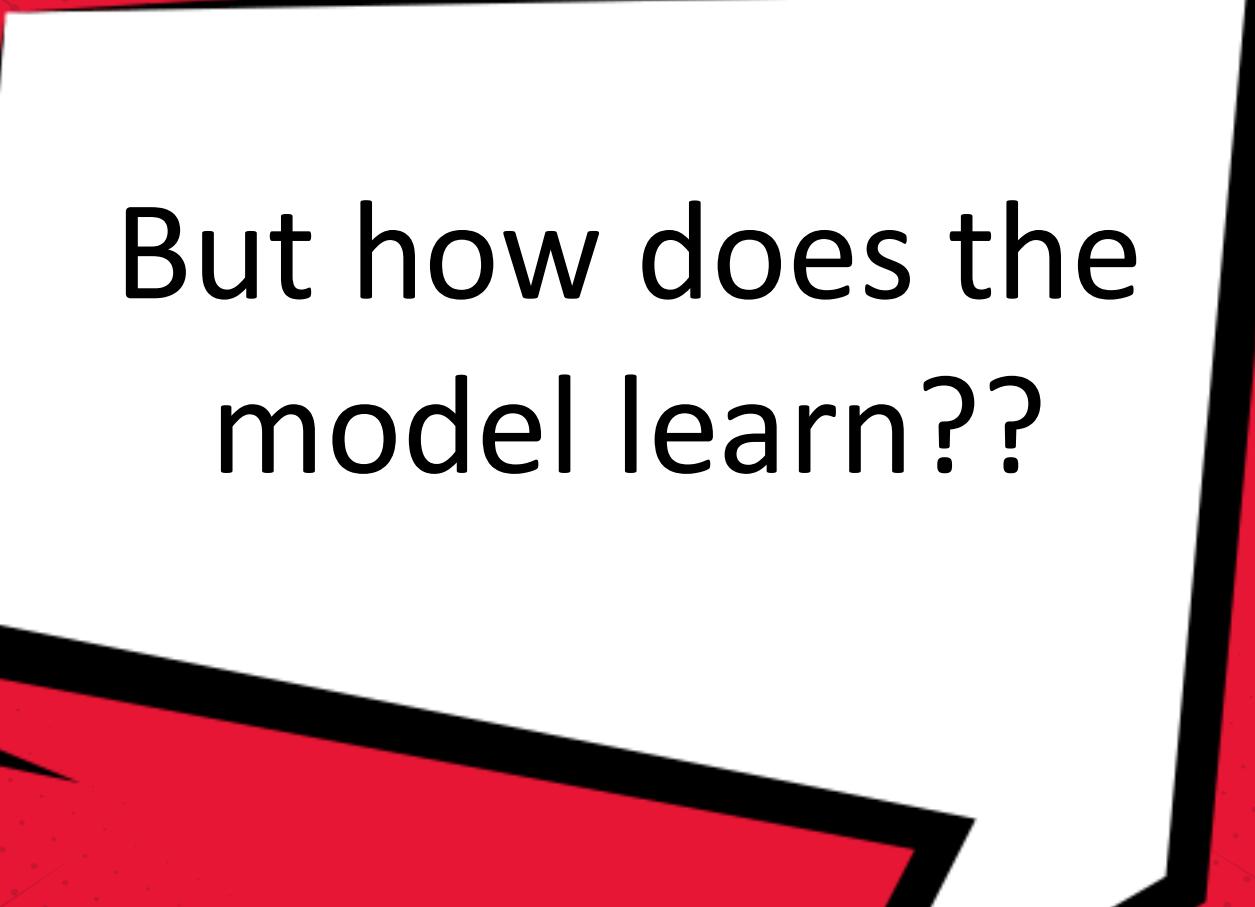
## Step 4: Evaluate

Original  
Caption

Wrinkles?  
Not in  
NYC!

Predicted  
Caption

Sprinkles?  
Got in  
NYC!

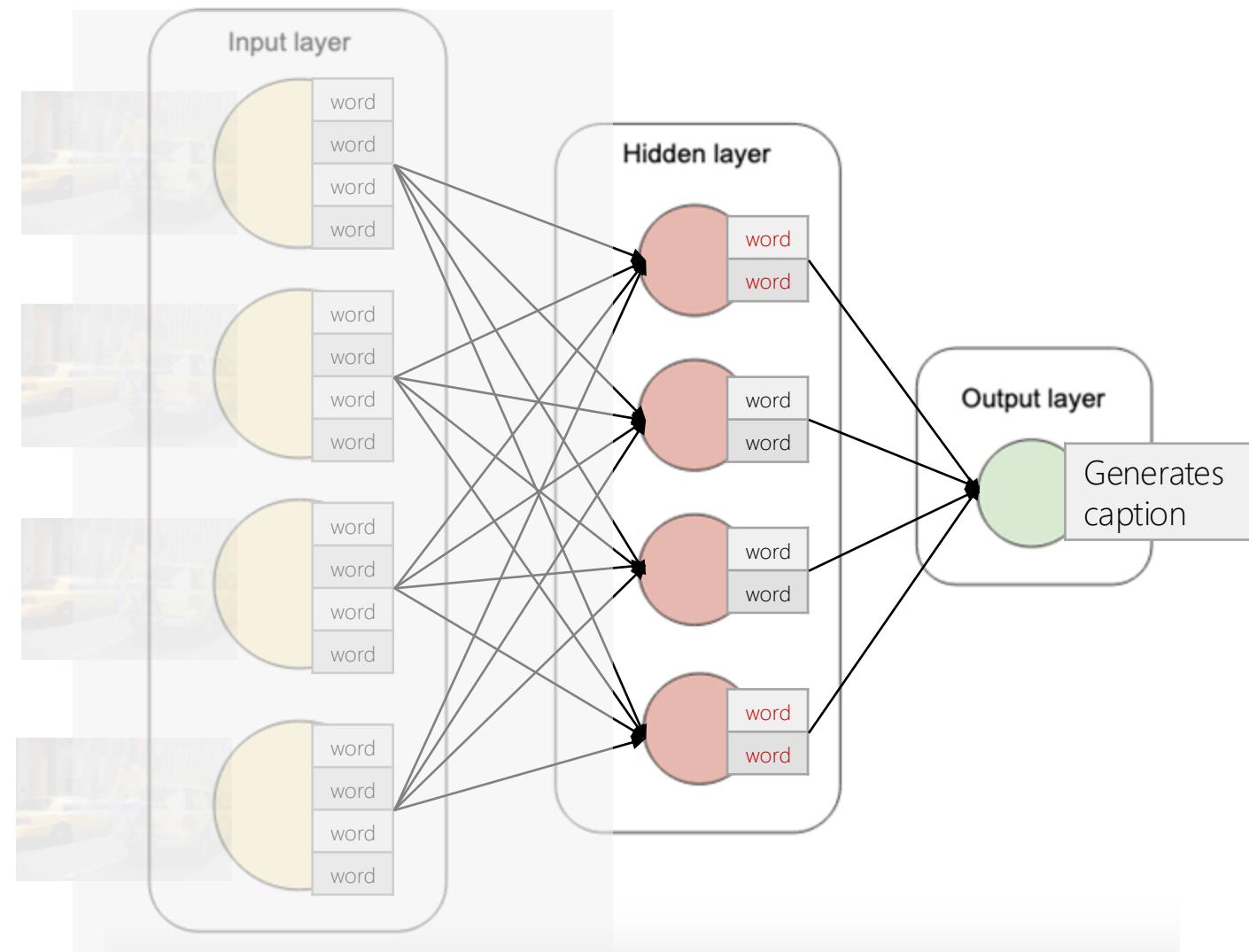


But how does the  
model learn??

# Step 5:

The output node:

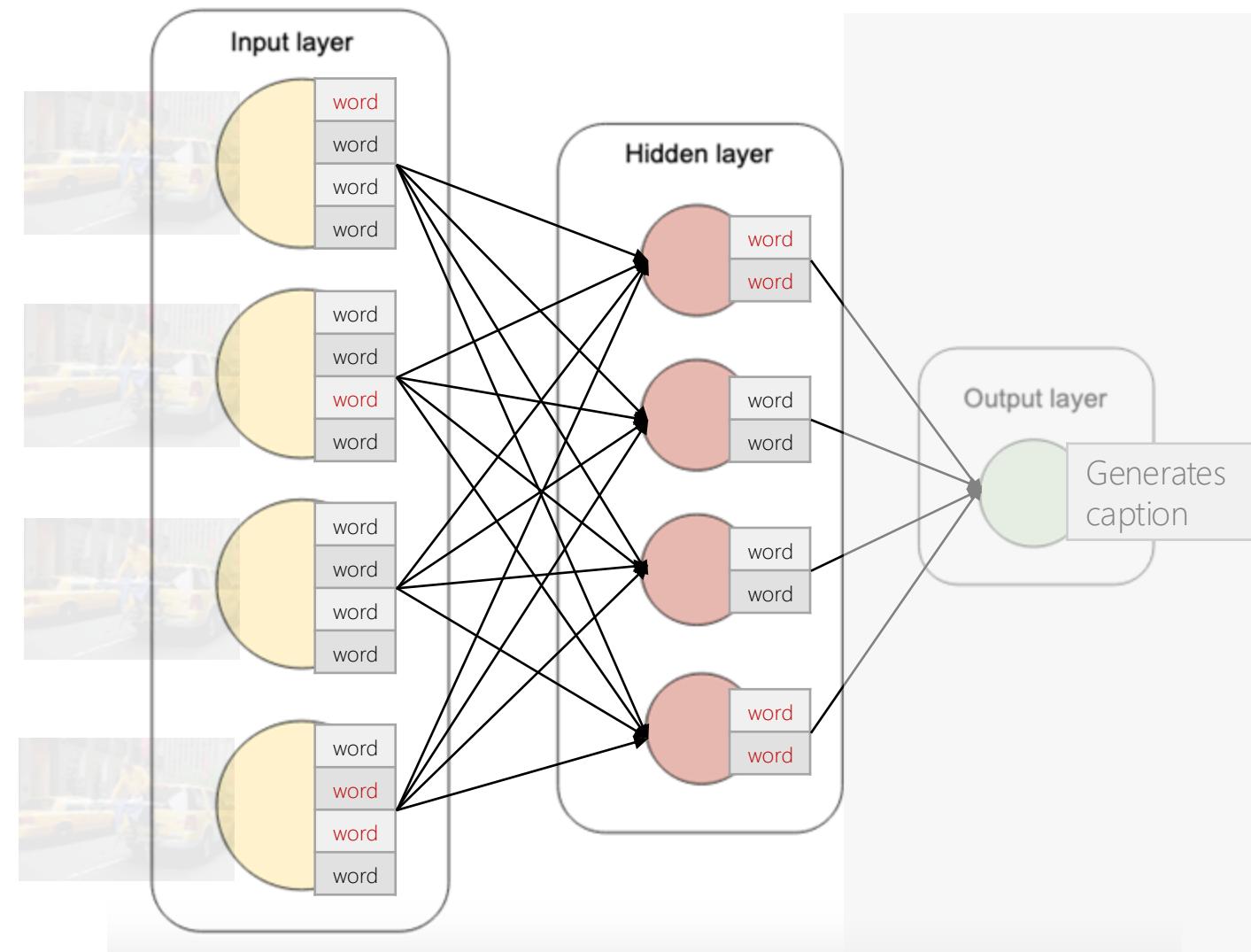
- Identifies the senders' words that were used in the original caption, i.e., the "correct" words
- Identifies which links gave "better" information



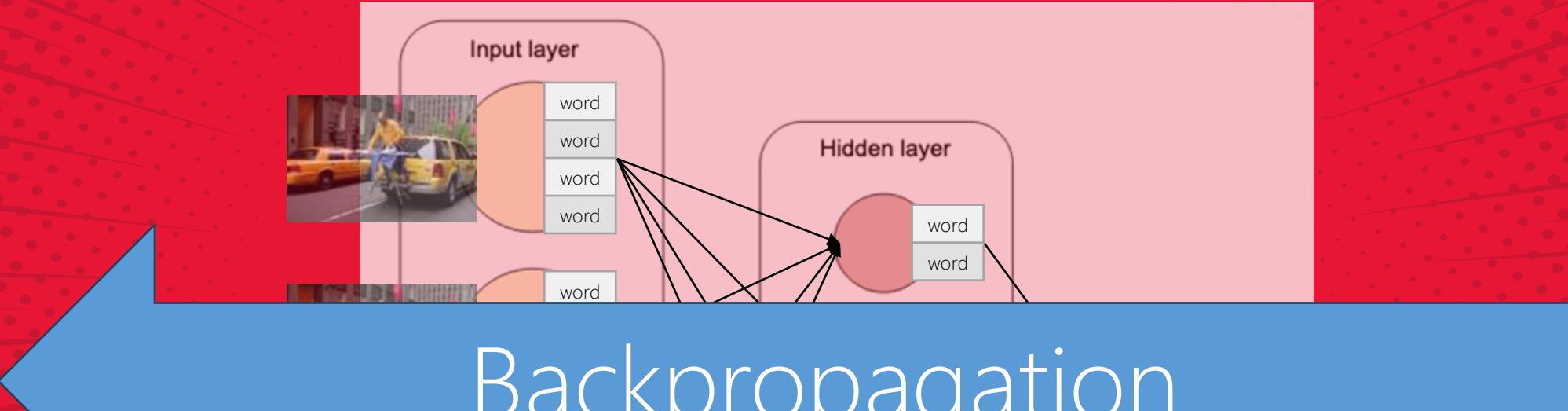
# Step 6:

The hidden layer nodes:

- Circles the senders' words that were in the original caption. (these are "correct")
- Identifies which links gave "better" information



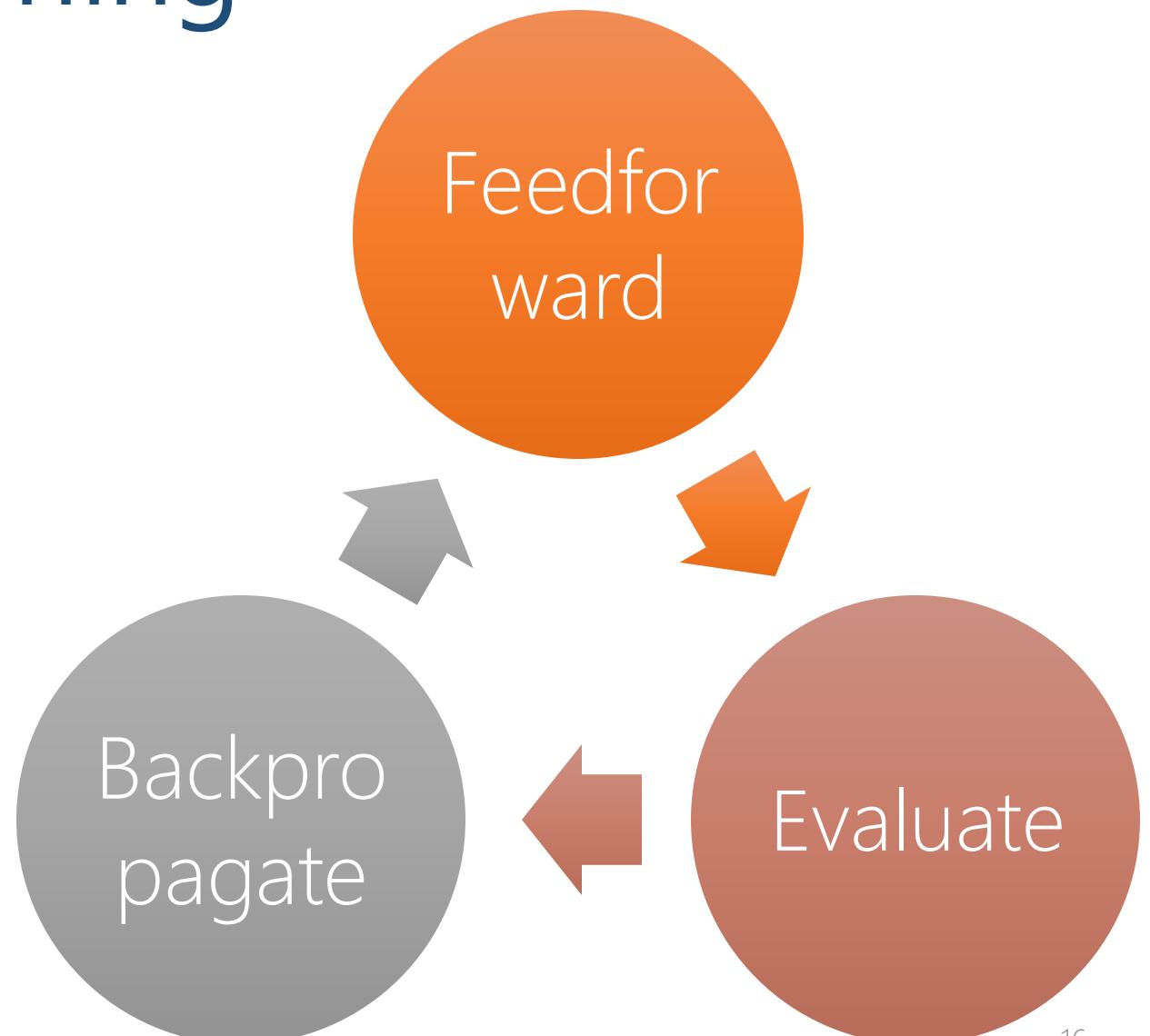
This is essentially Backpropagation!!



Backpropagation

# Neural Network learning

1. Each neuron finds out which words are right/wrong
2. adjusts its weights for how it picks words so the network can do better next time



# Back to Transformers

# The Early Ages

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

✗ Context

✗ Long term dependencies



# The Current Age

- Transformers
- Logical, Mathematical and Commonsense Reasoning
- Alignment using RLHF
- Safety testing and jailbreaks
- Ethical and Fair models

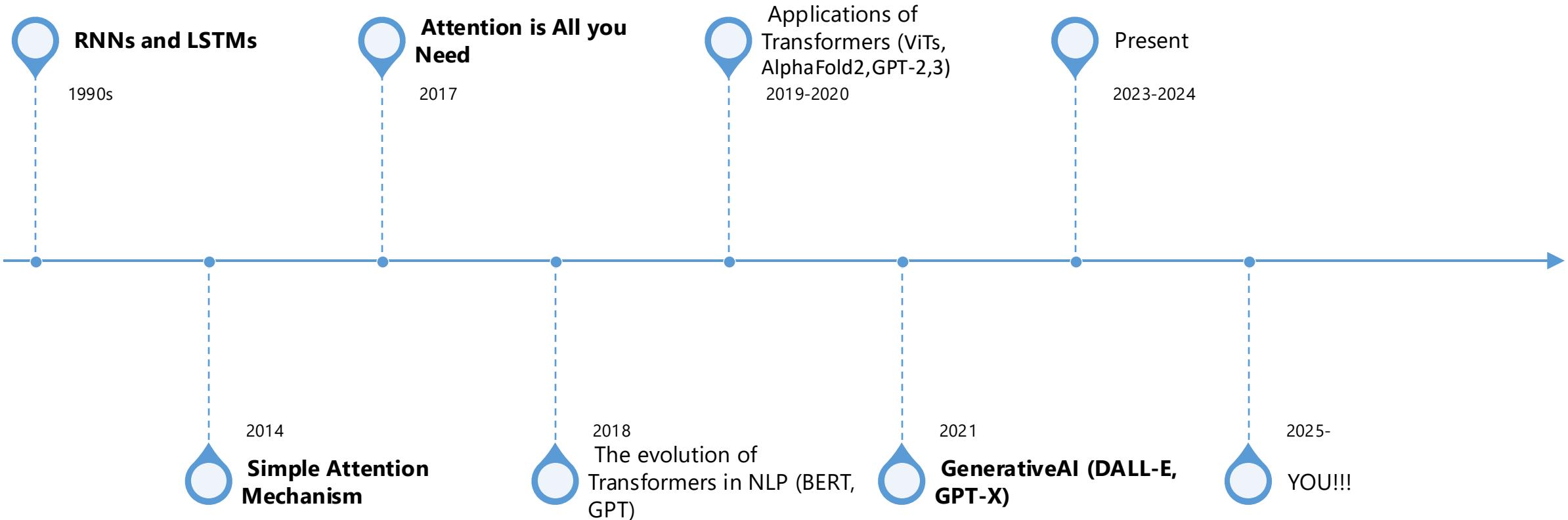


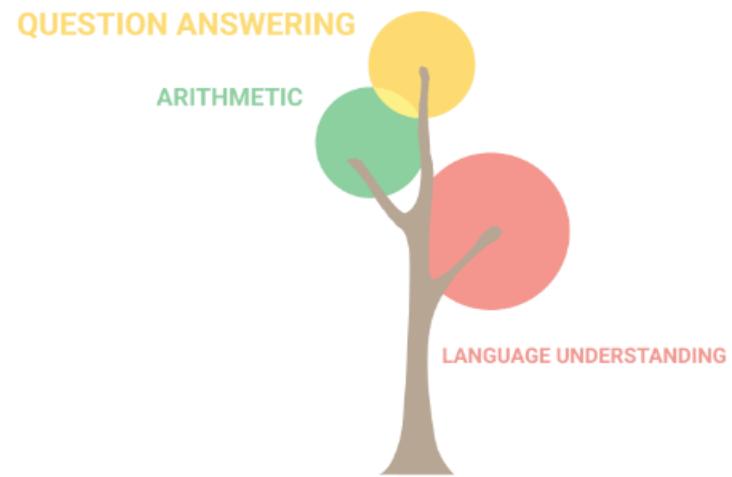
Context



Long term dependencies

# Evolution of Attention



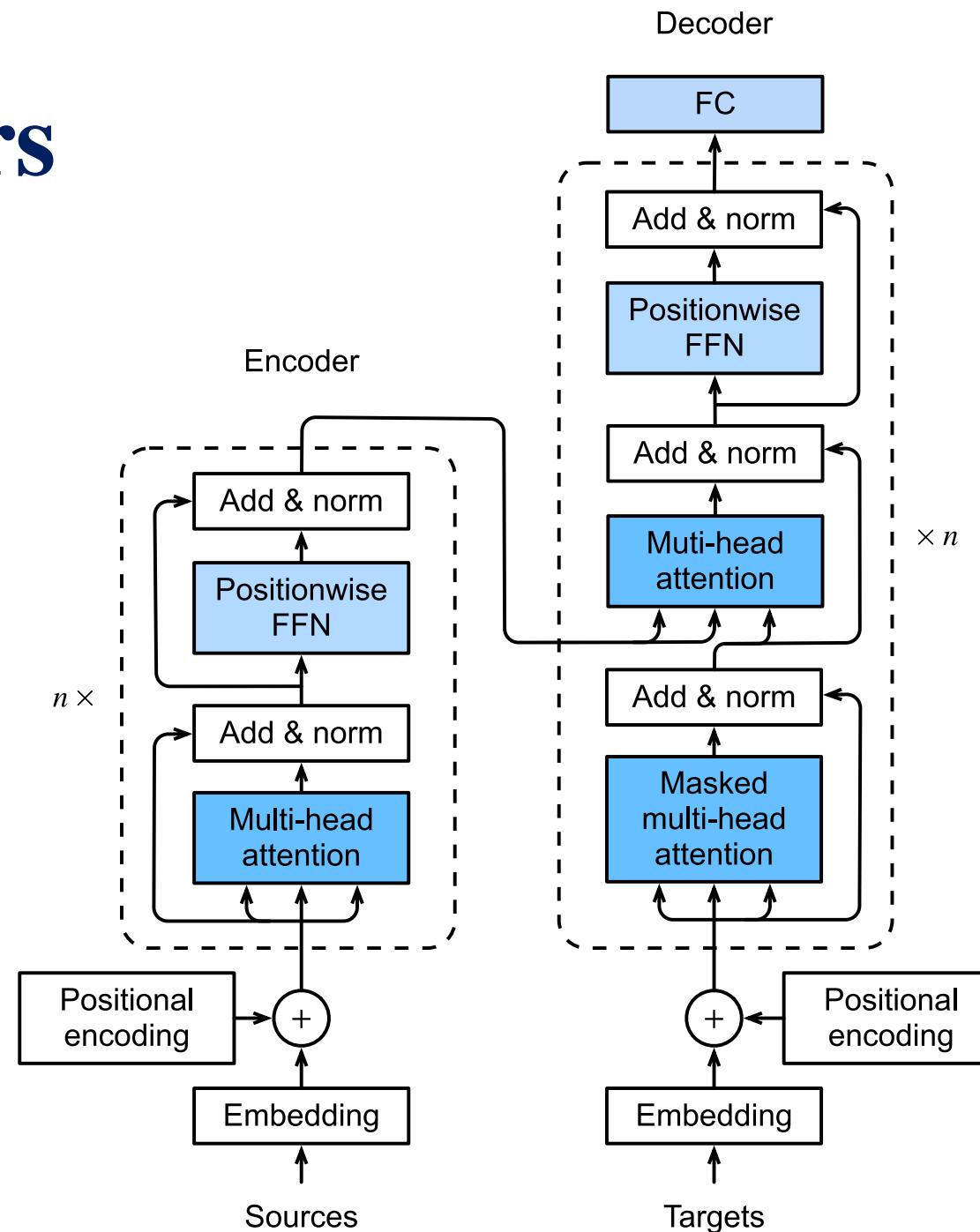


8 billion parameters

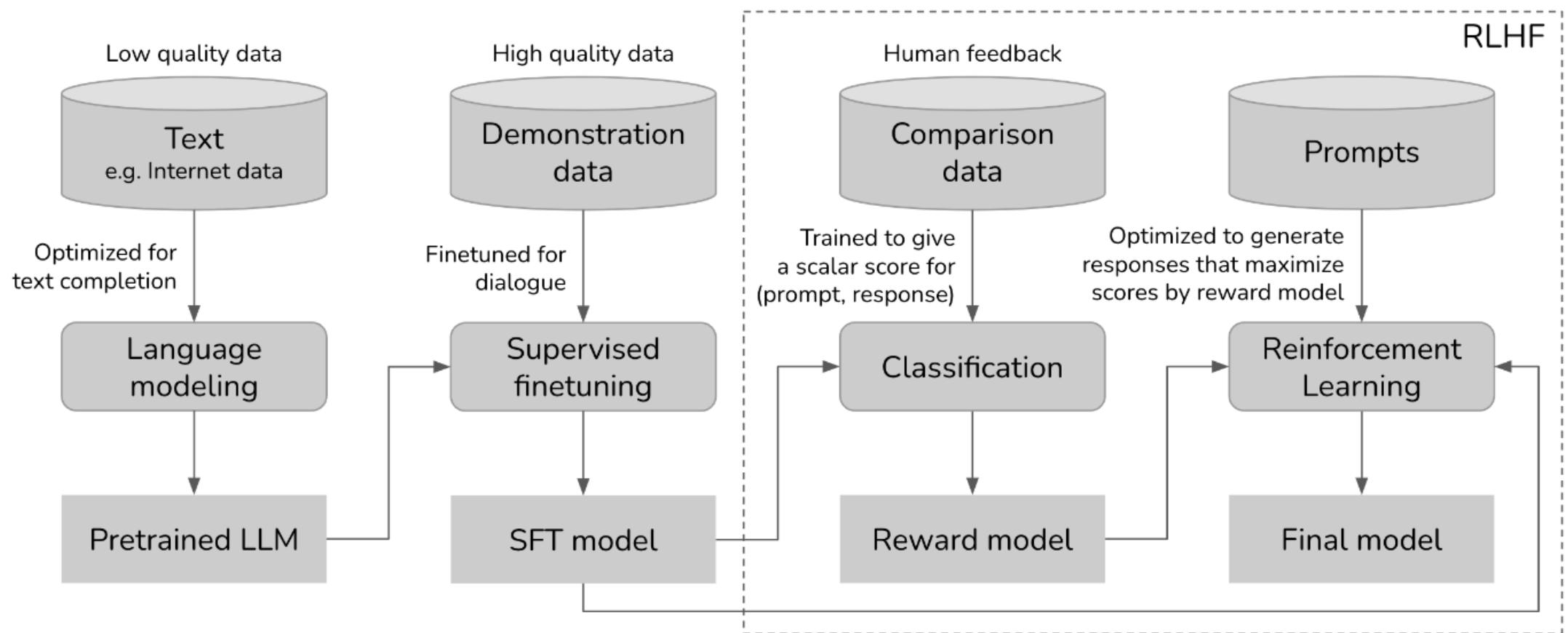




# Transformers



# Training Pipeline of Large Language Models



# Input Processing

- 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

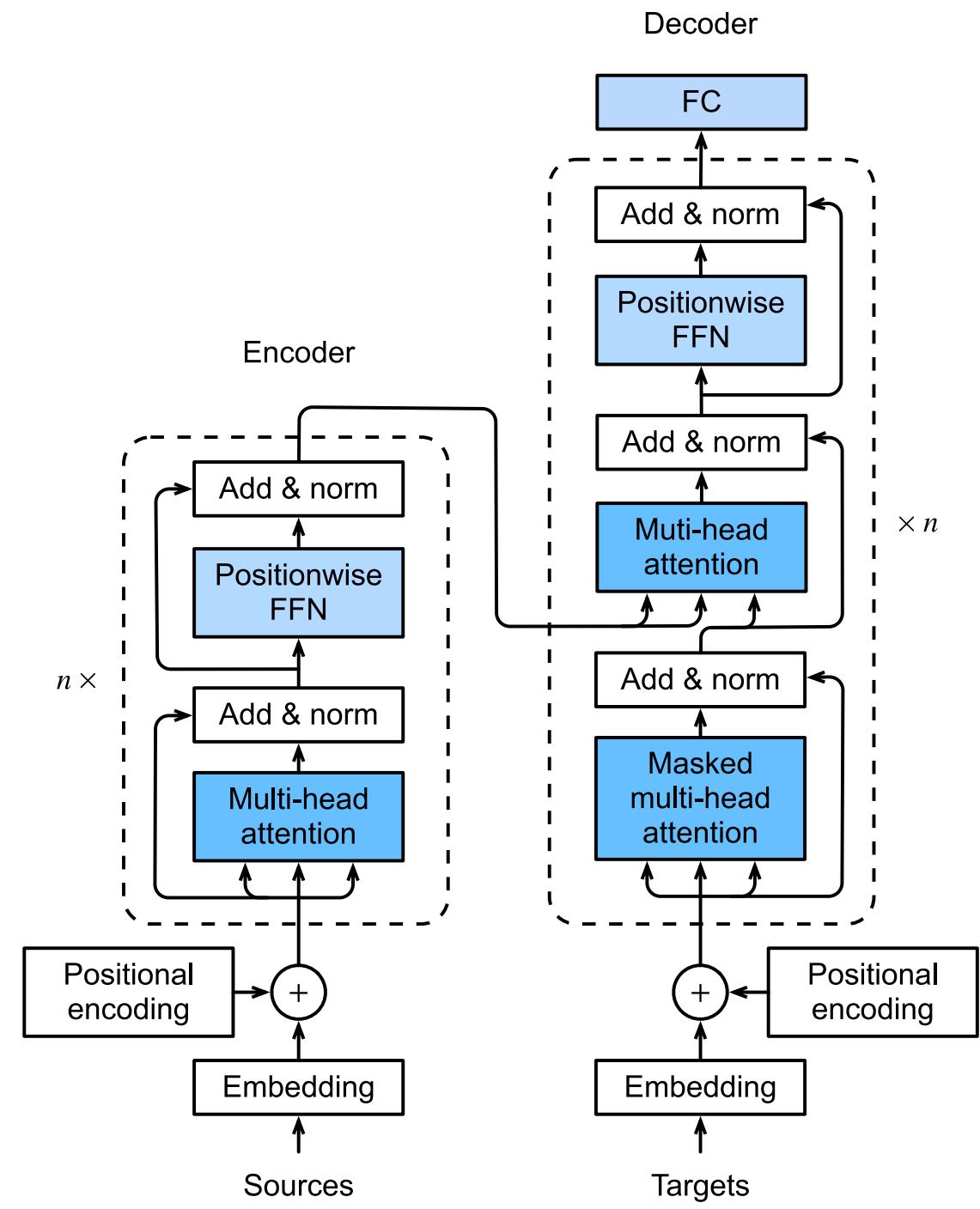
It matters not what someone is born,

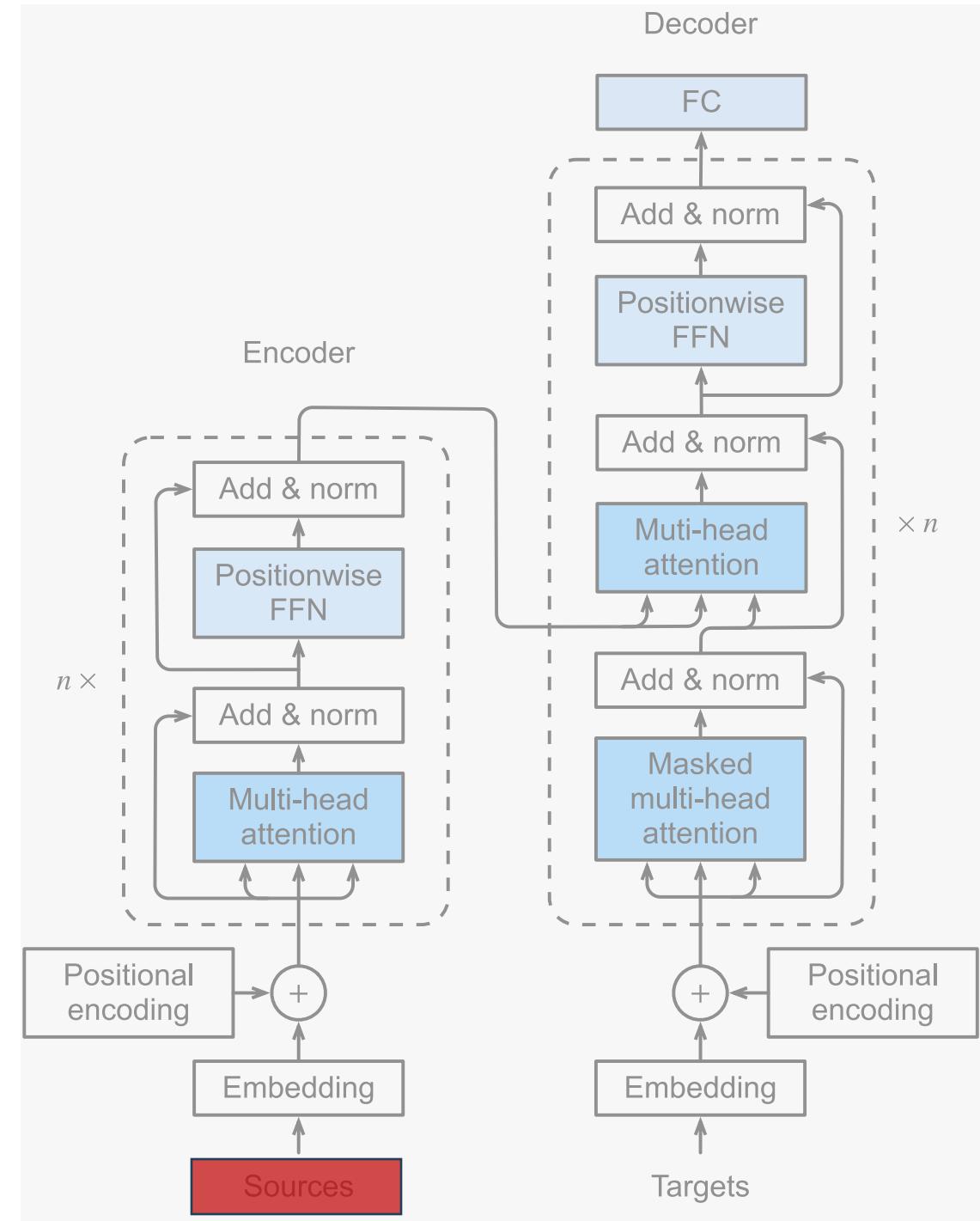
## Dialogue Completion

<start> far more than our abilities  
<end>

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

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

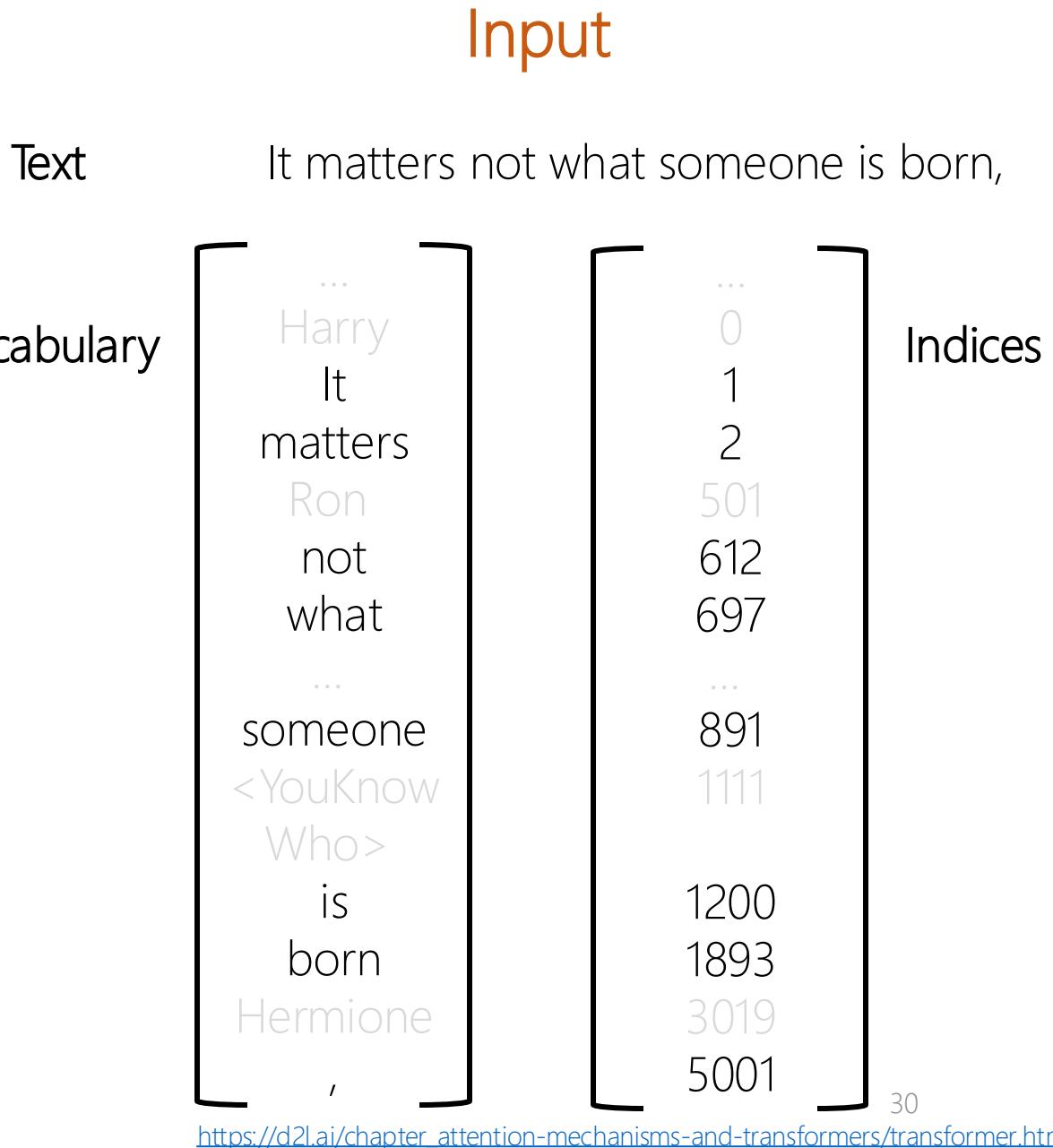
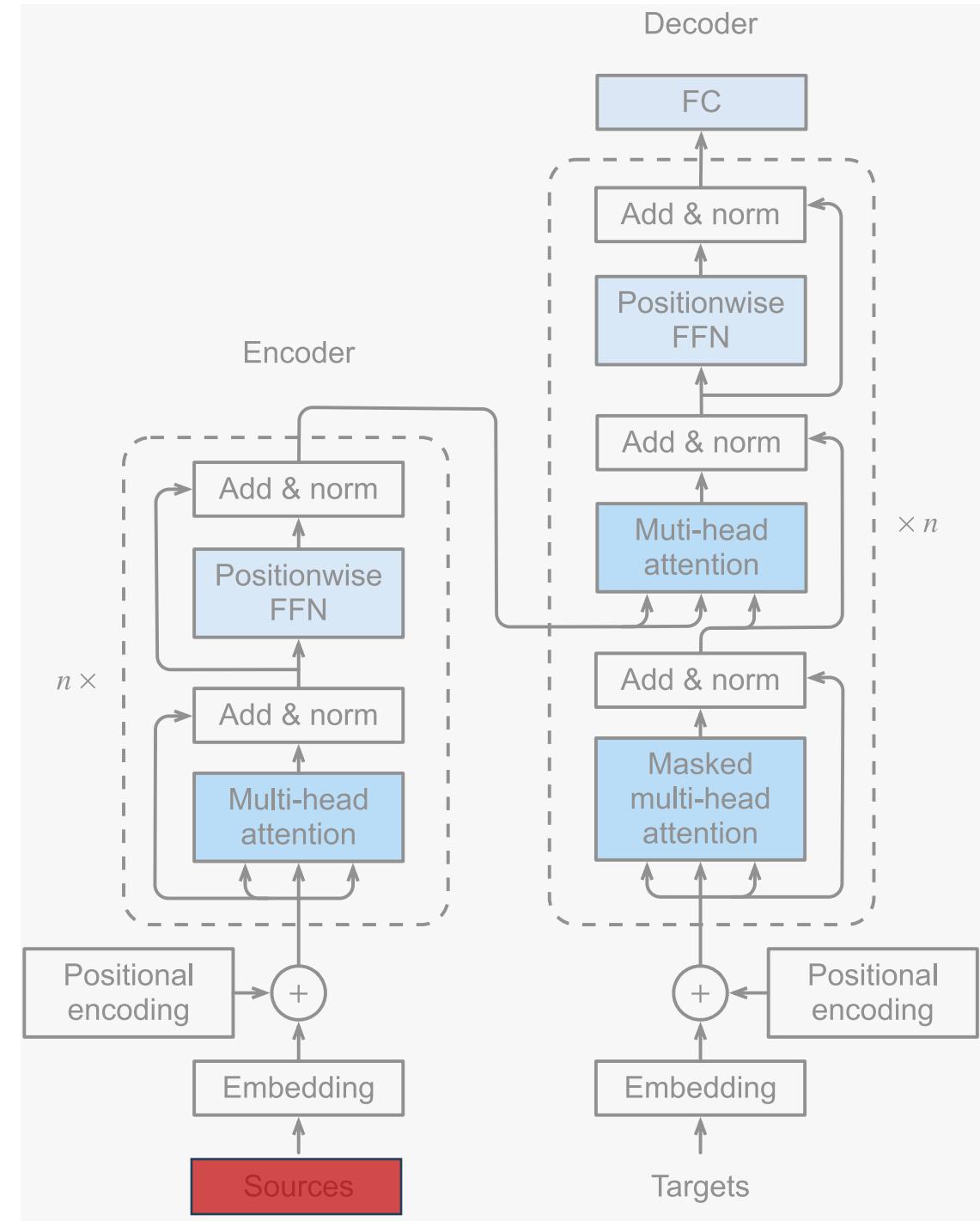


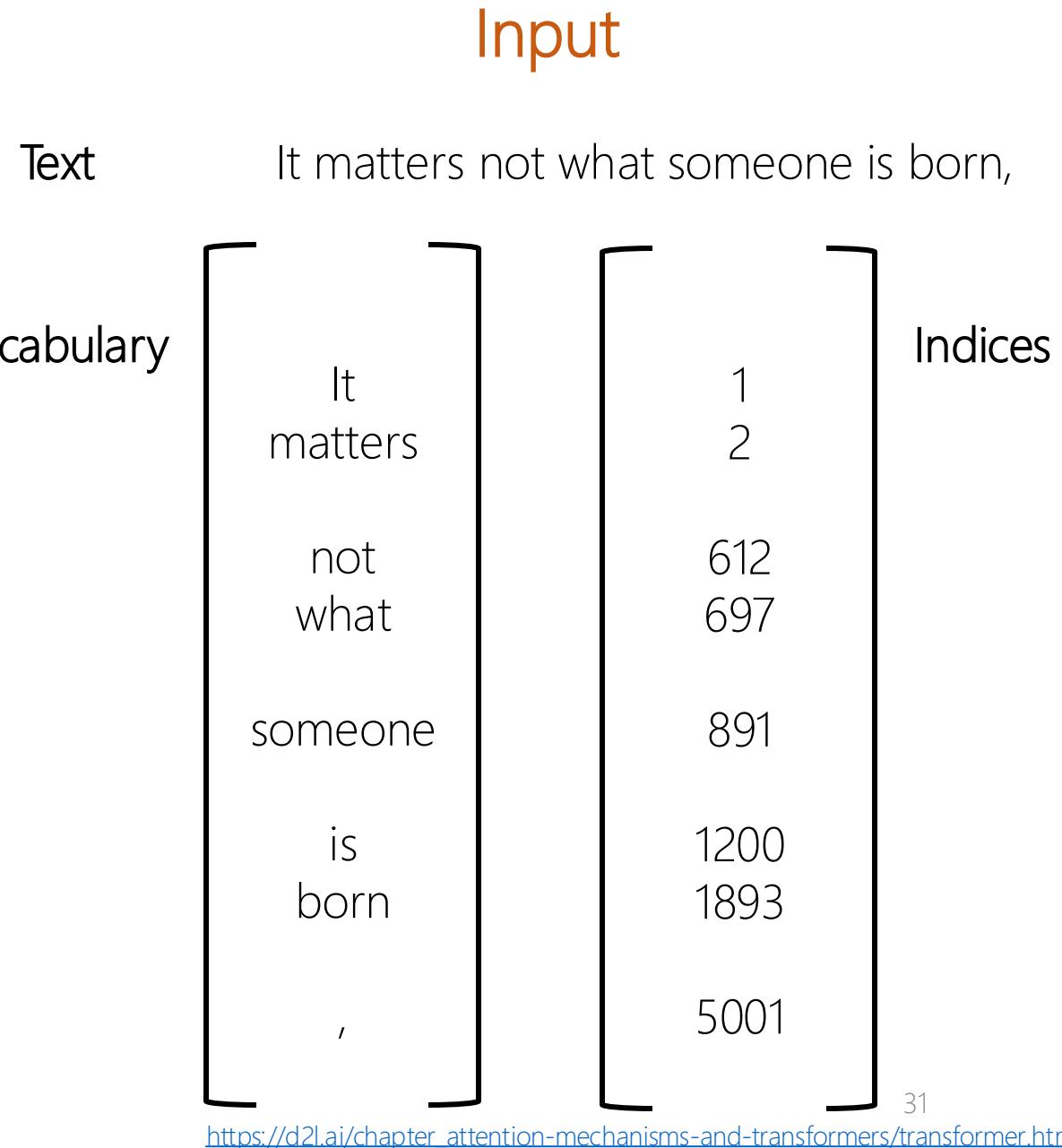
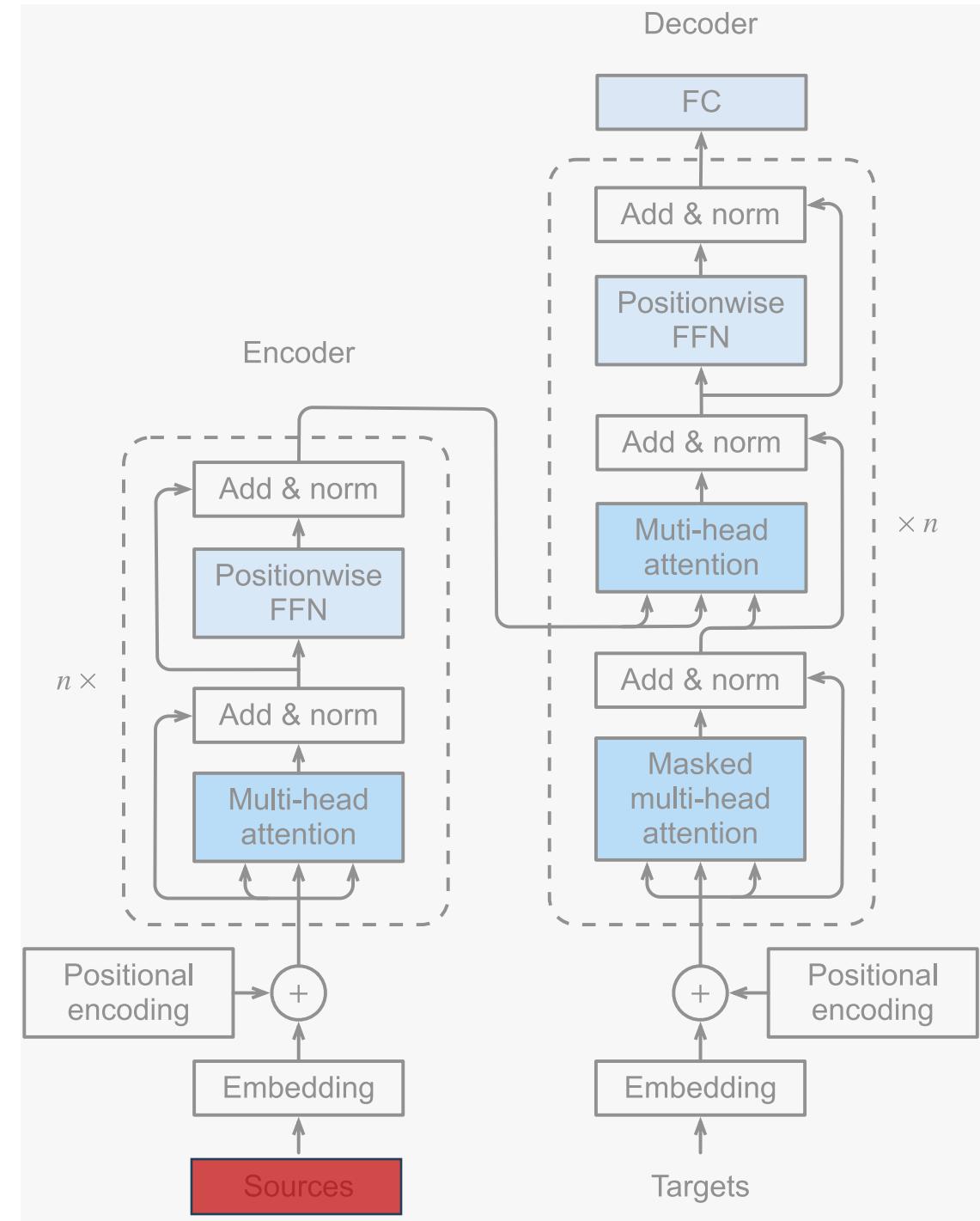


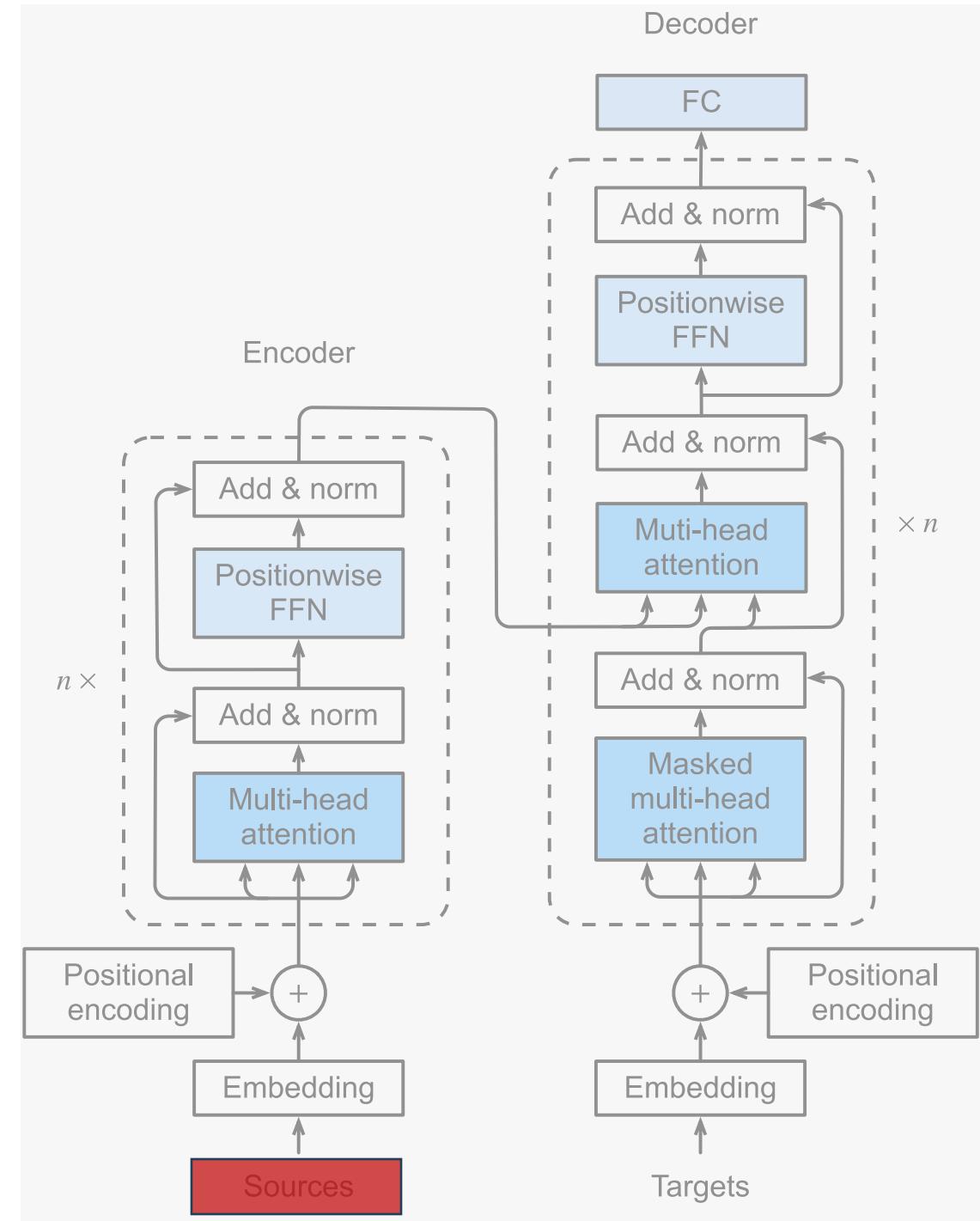
Text

Input

It matters not what someone is born,







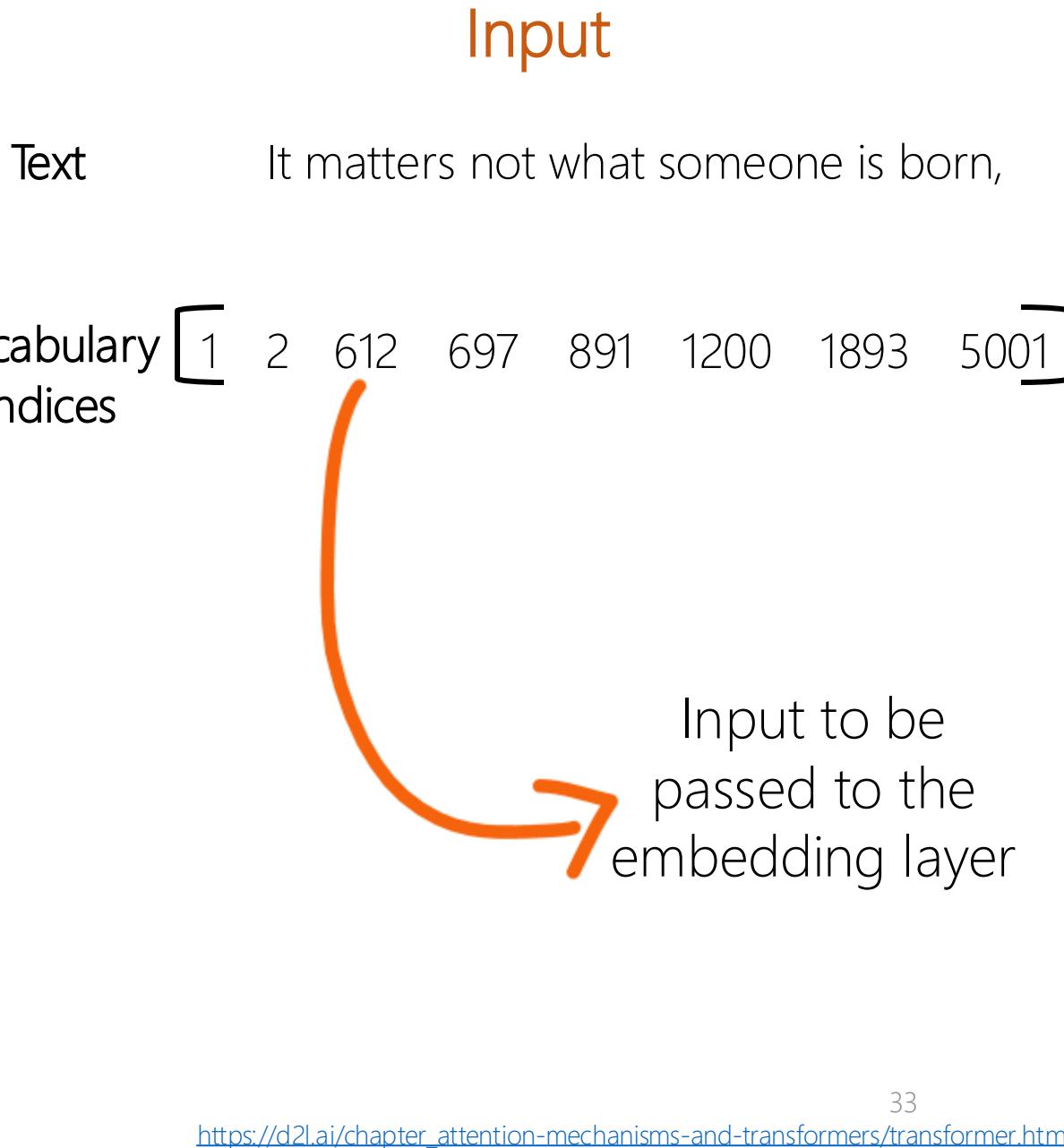
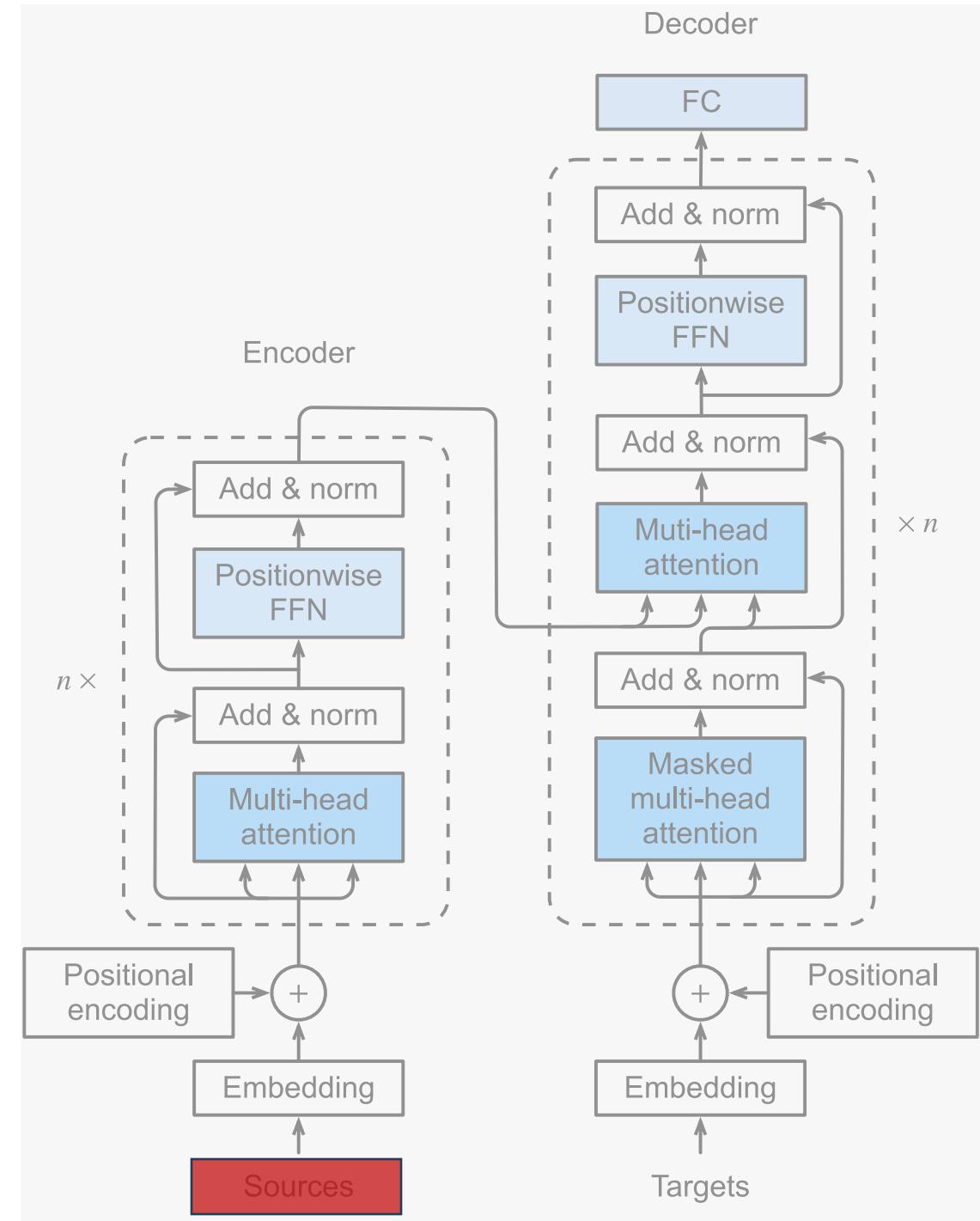
## Input

Text

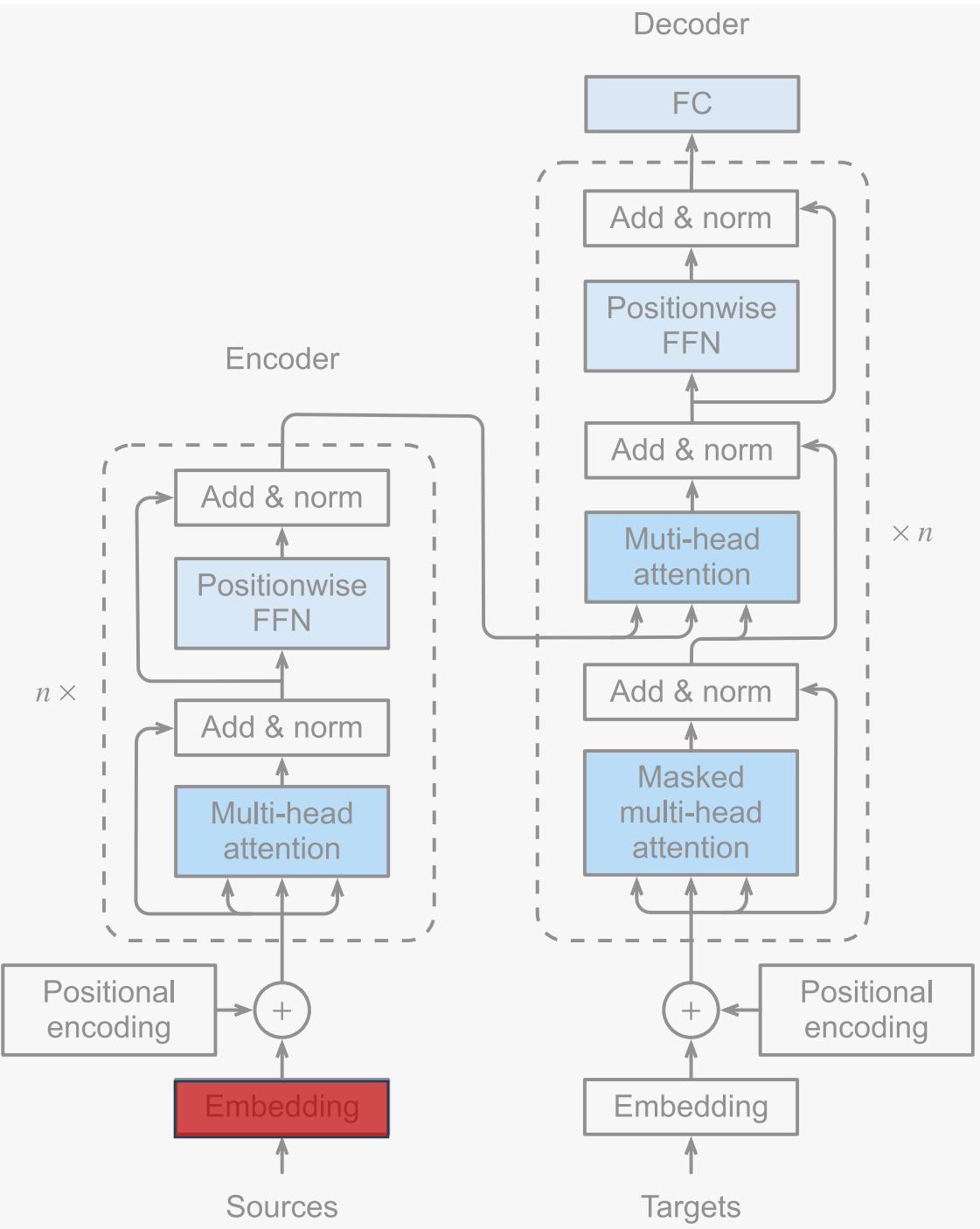
It matters not what someone is born,

Vocabulary  
Indices

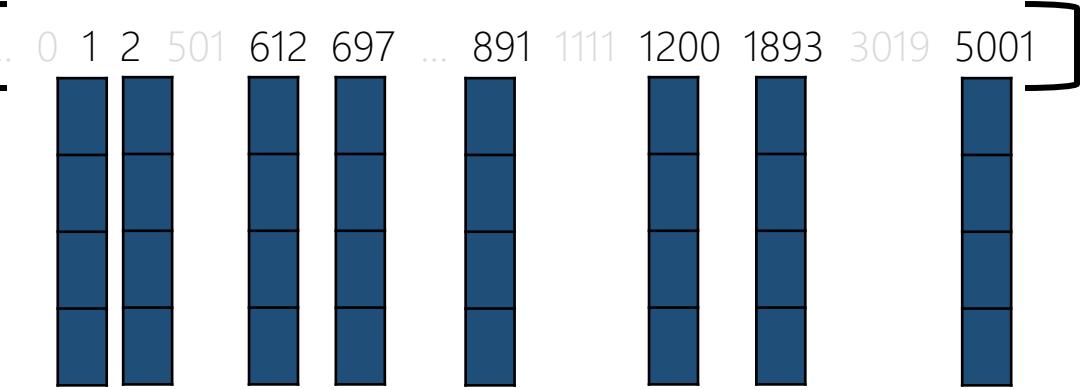
$x_0$	$x_1$	$x_2$	$x_3$	$x_4$	$x_5$	$x_6$	$x_7$
1	2	612	697	891	1200	1893	5001

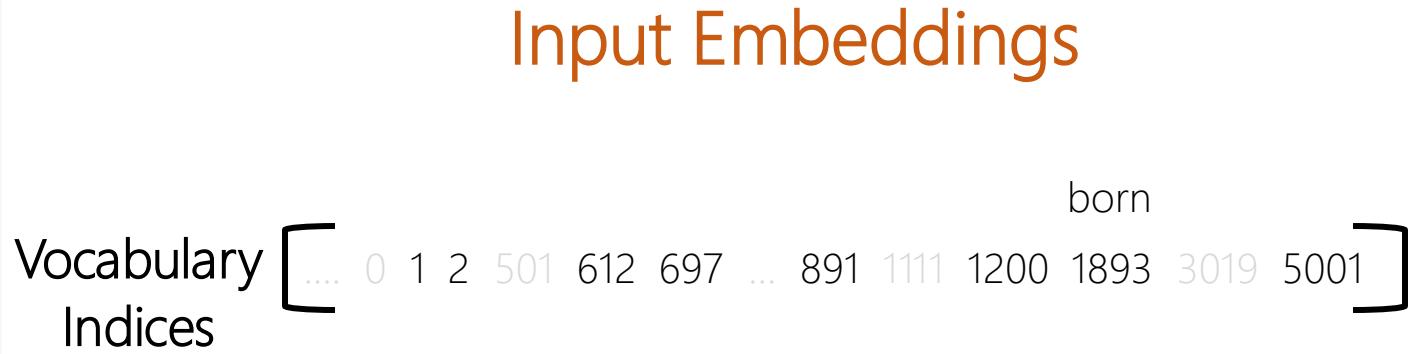
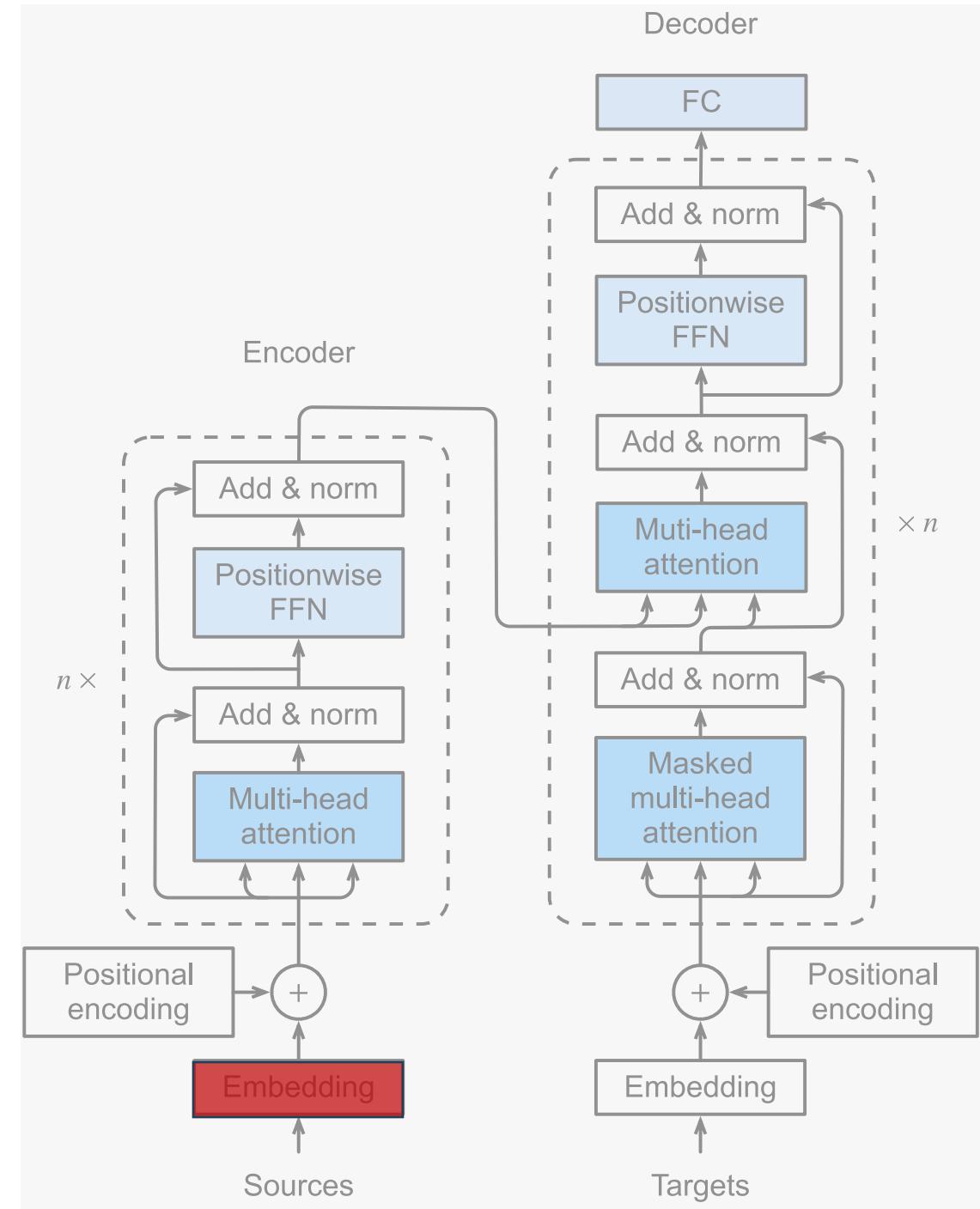


# Input Embeddings

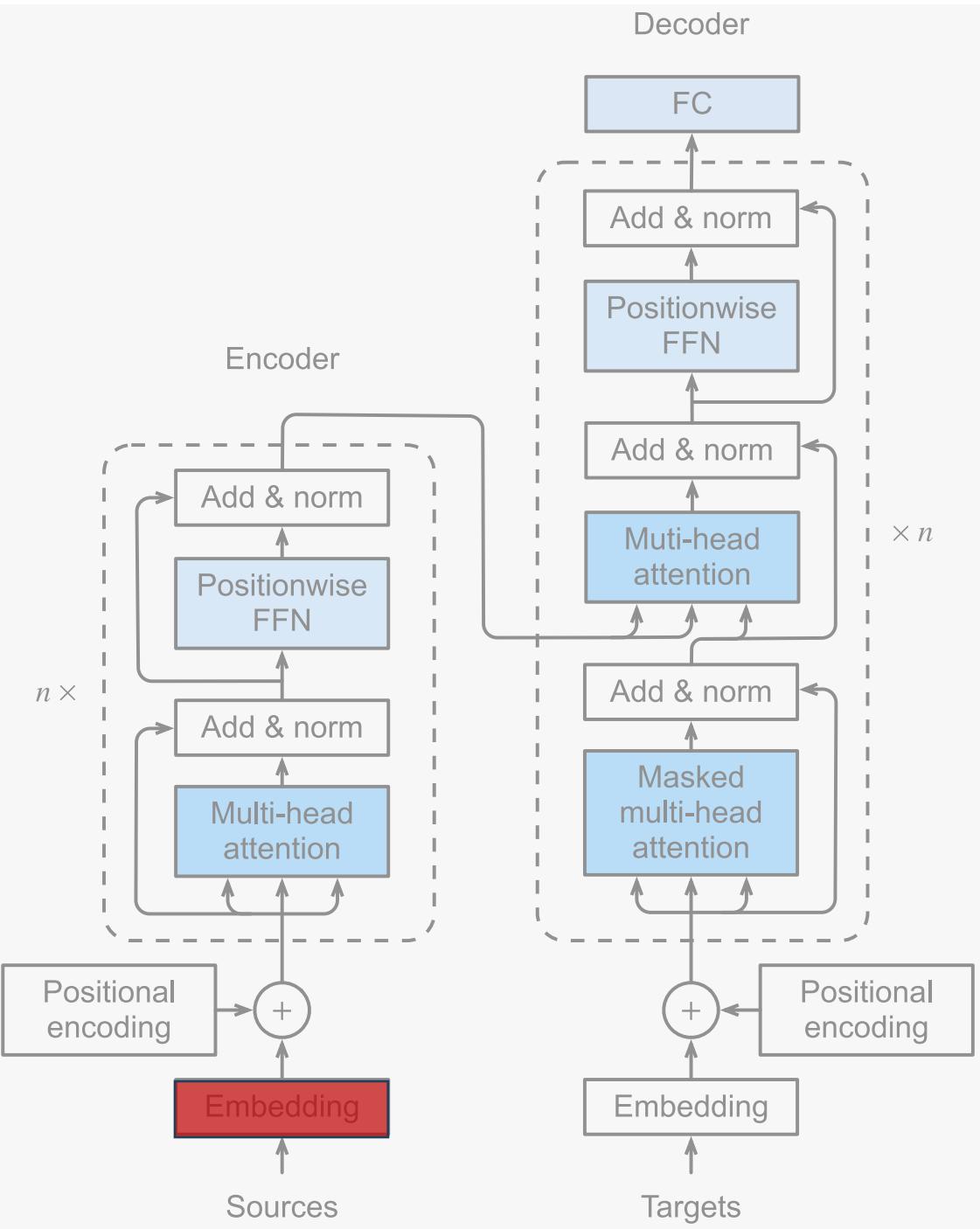


Vocabulary  
Indices

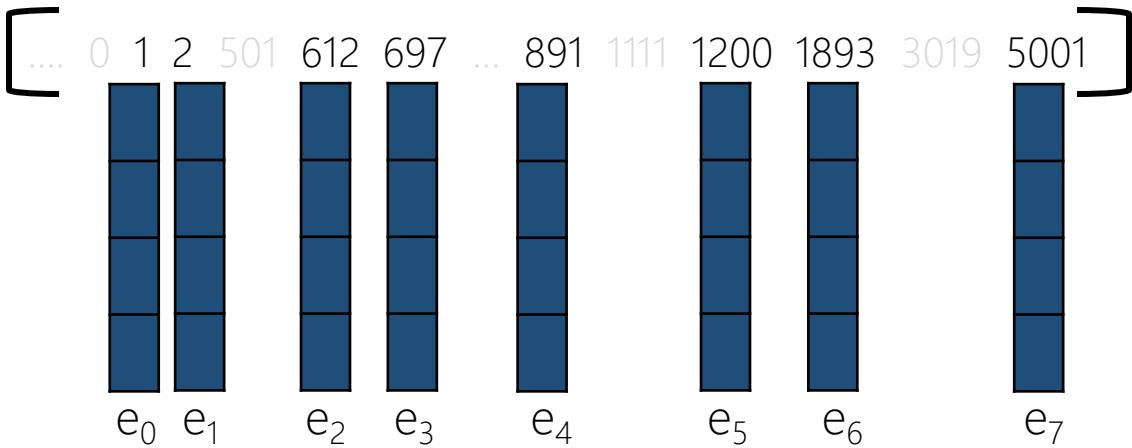


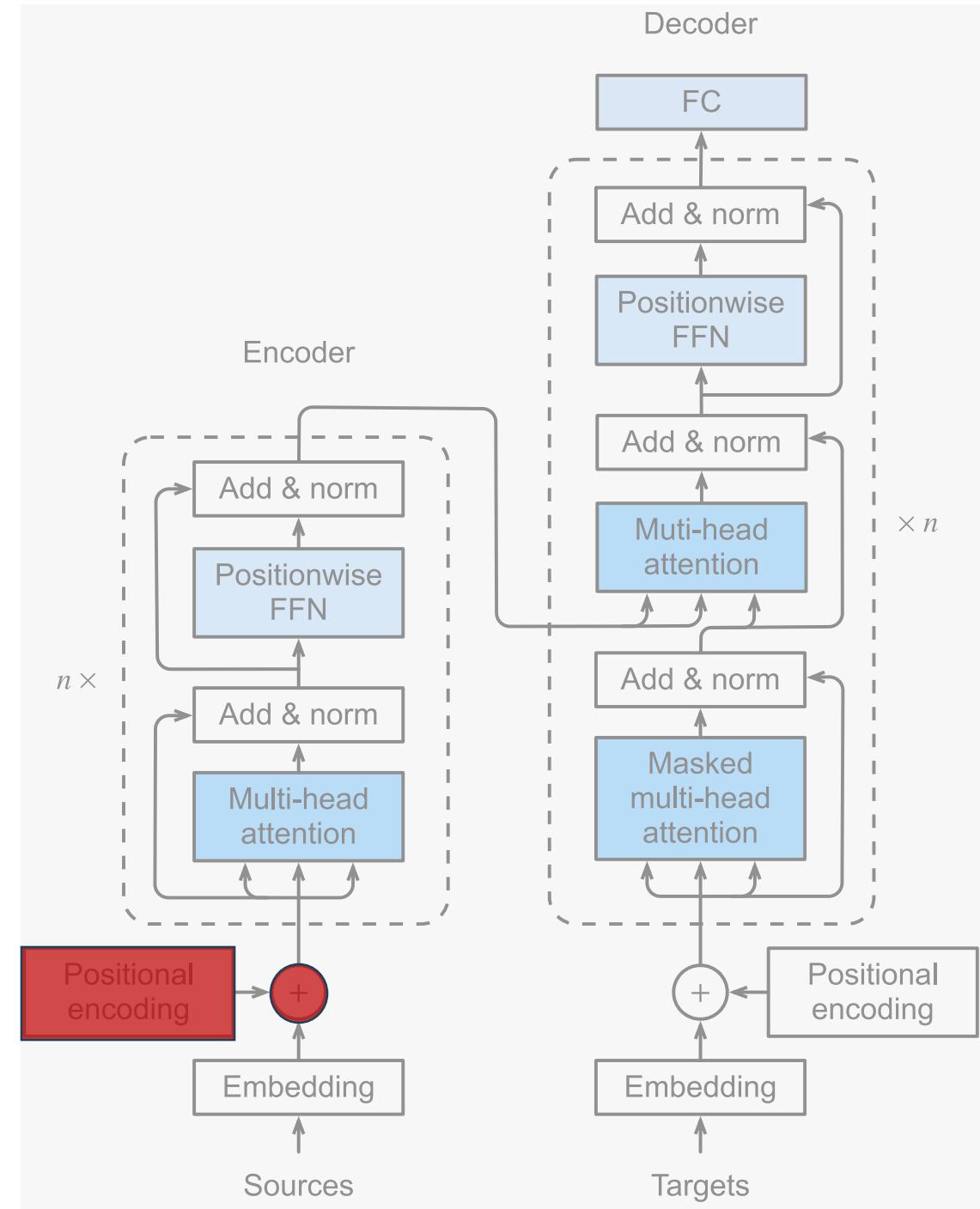


# Input Embeddings

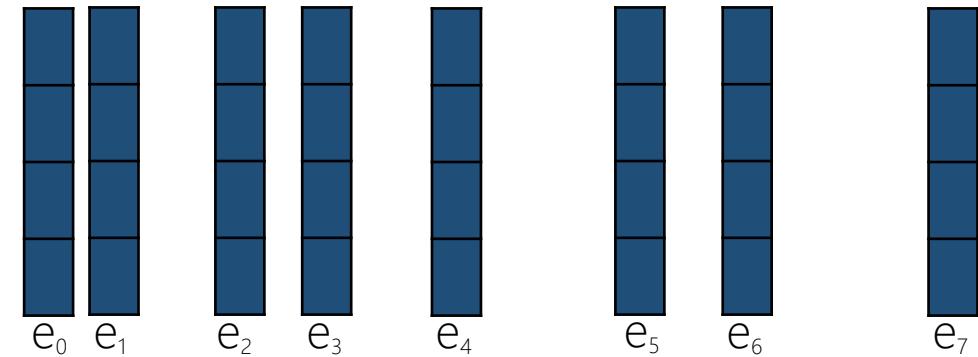


Vocabulary  
Indices

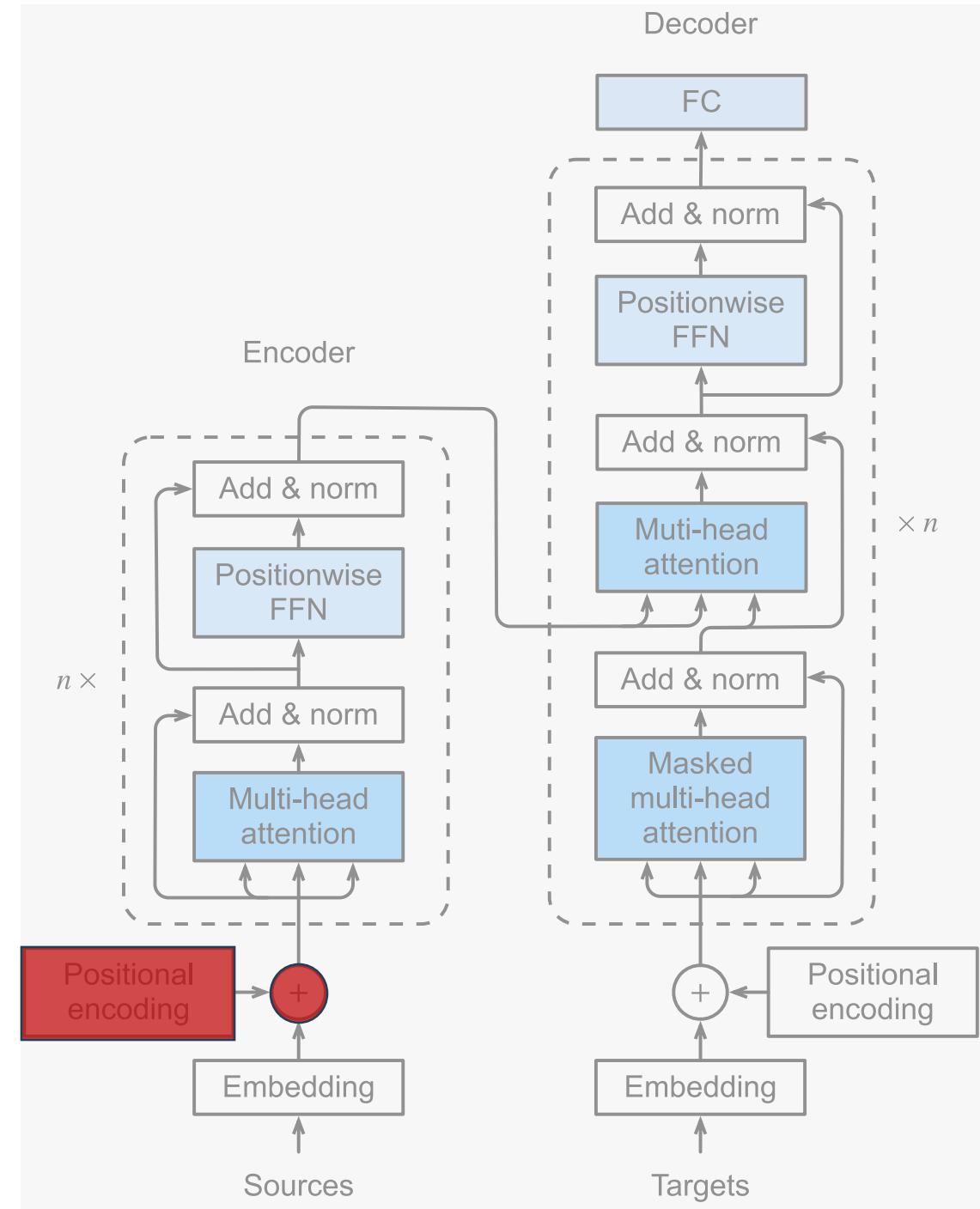




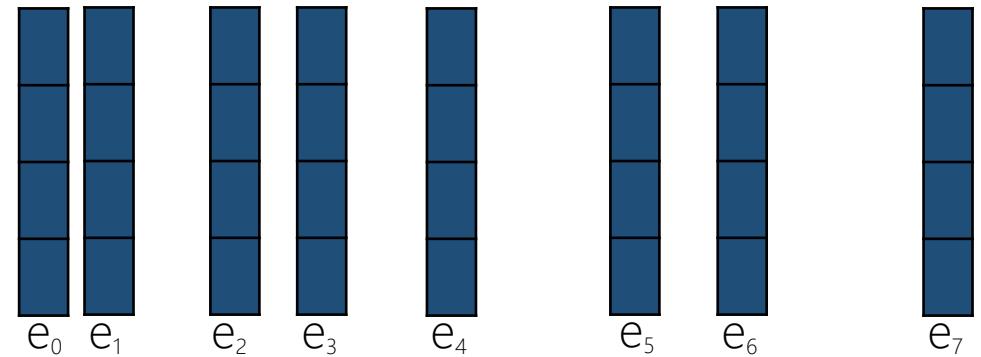
## Why do we need Positional Embedding?



Recurrent Neural Network  
or  
Long-Short Term Memory

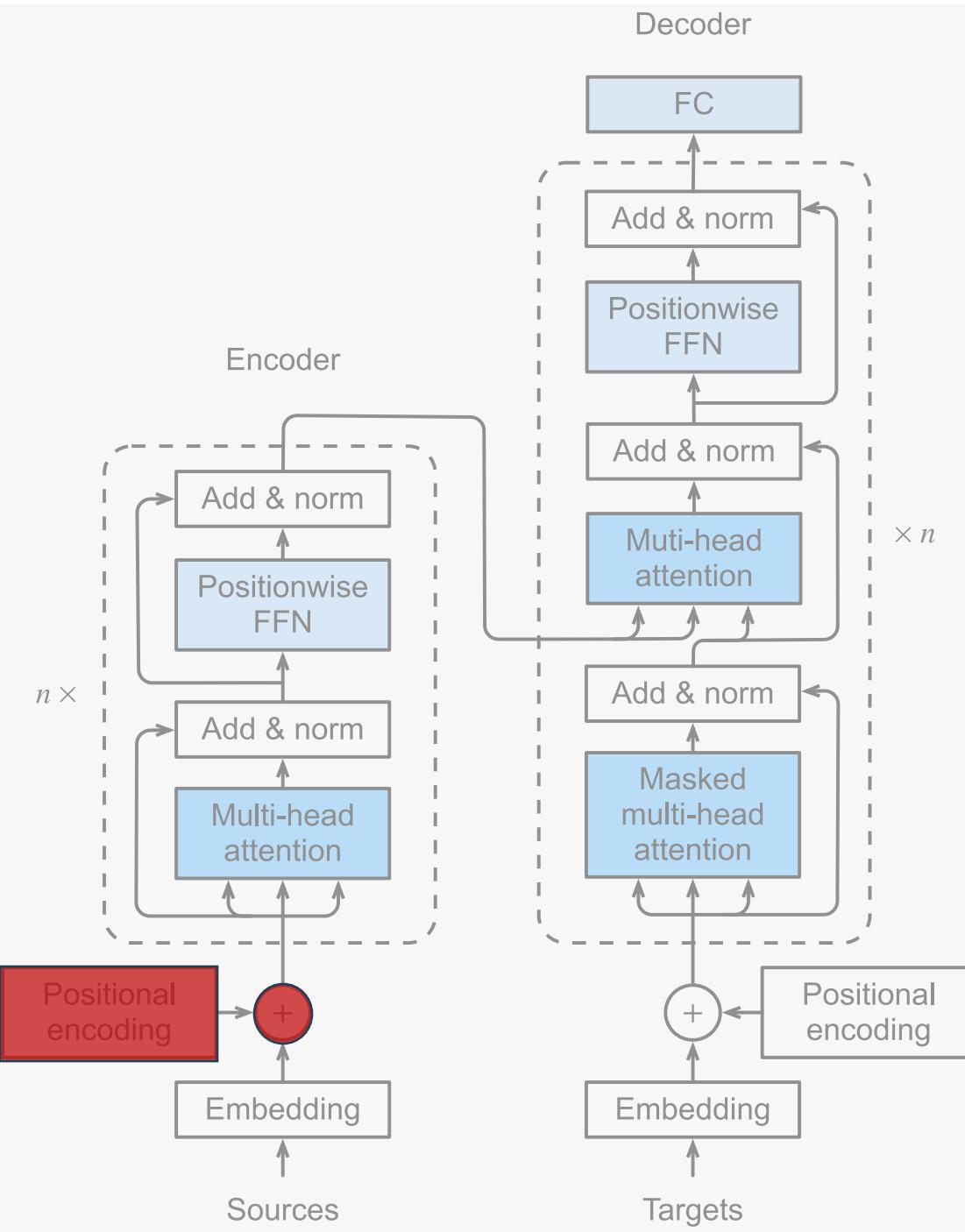


# Why do we need Positional Embedding?



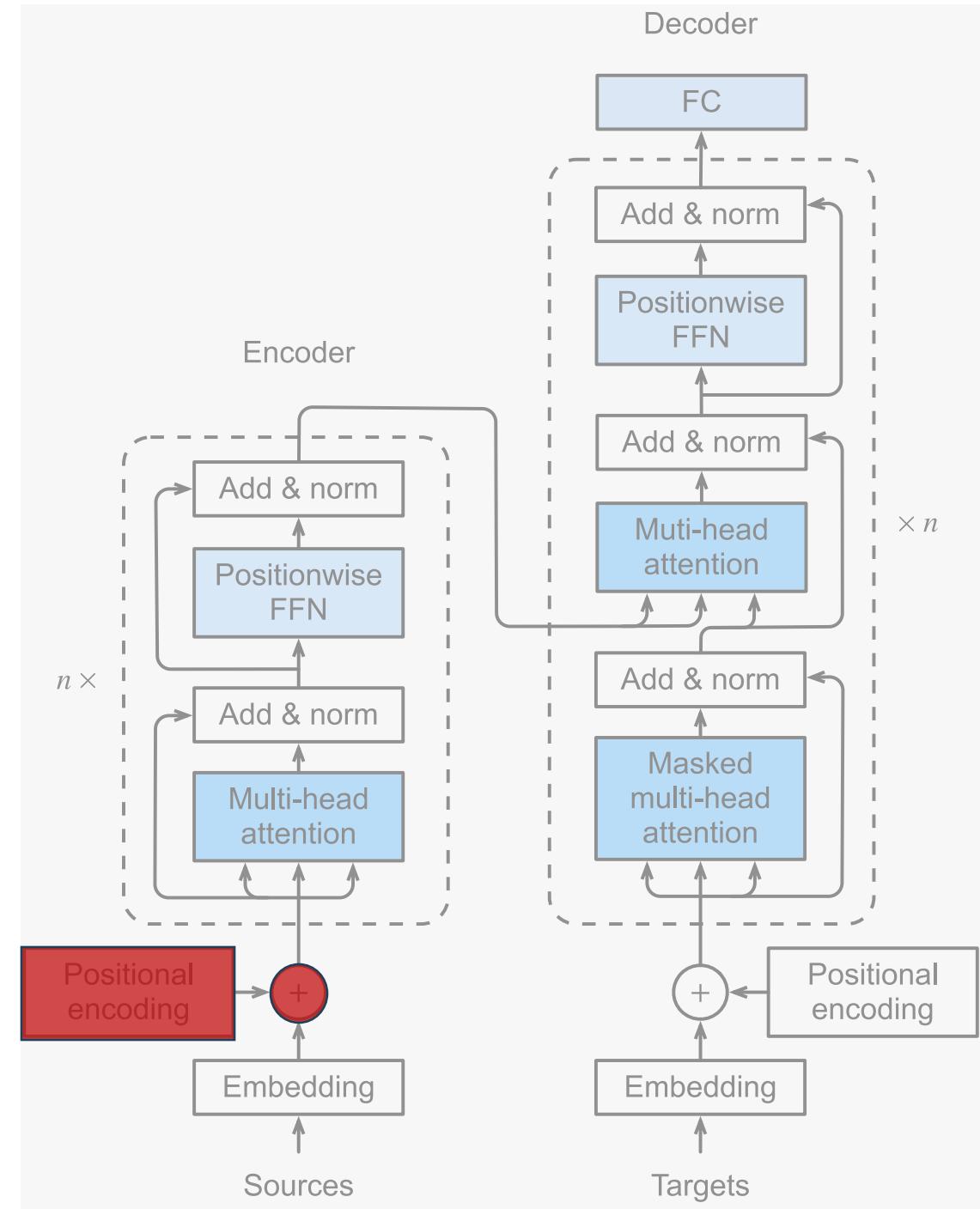
Transformers

# Why do we need Positional Embedding?

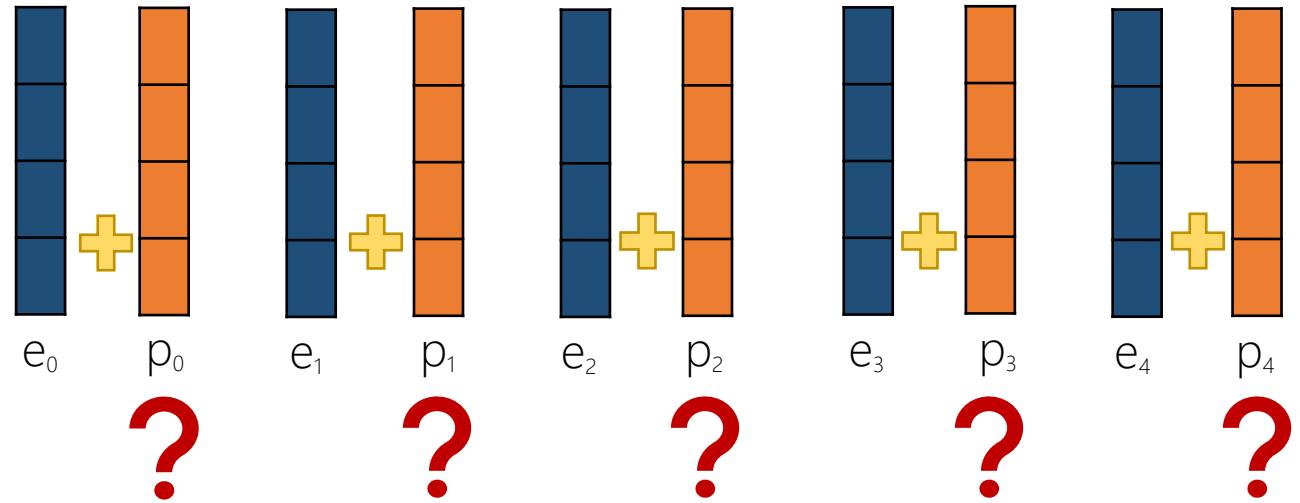


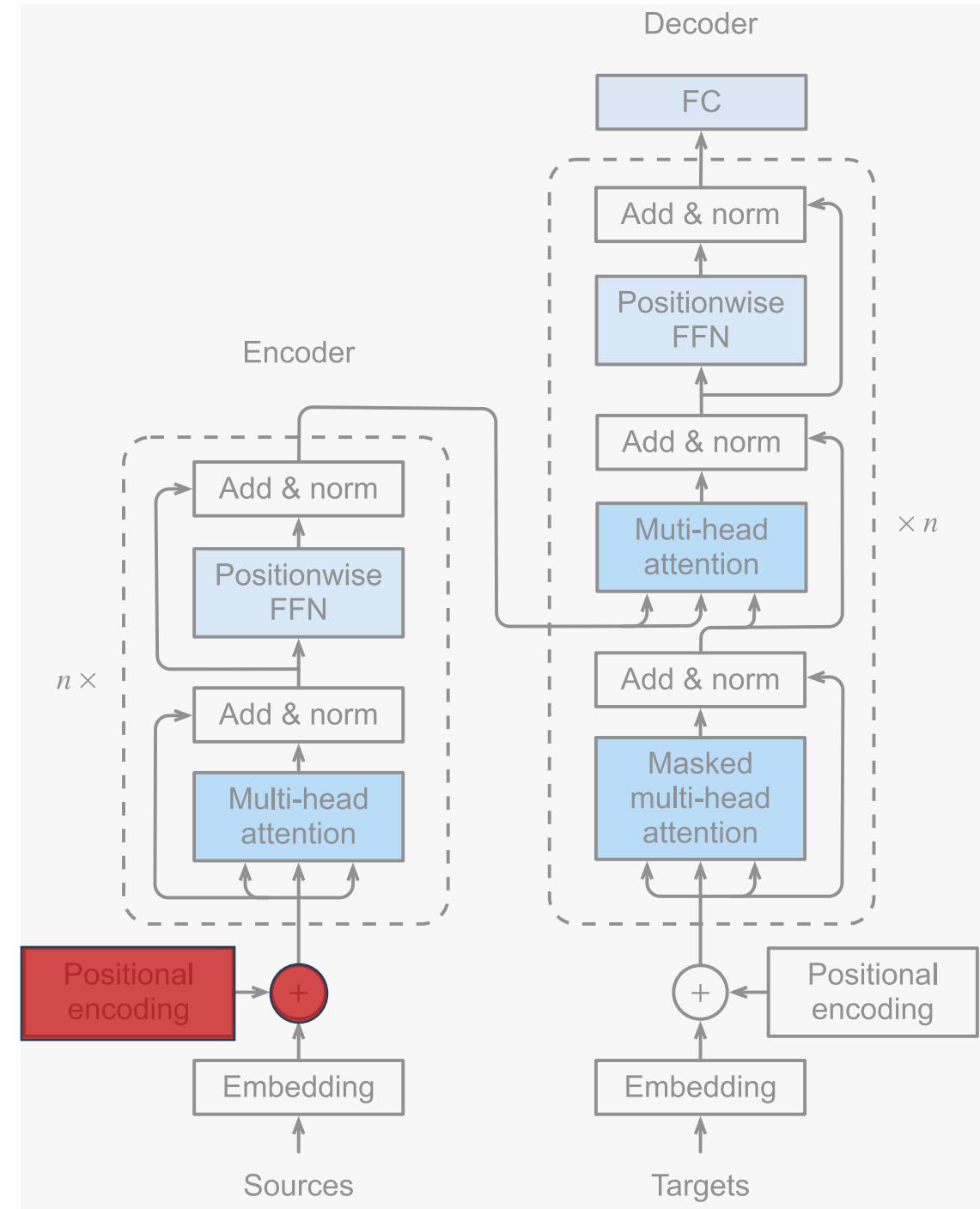
England invented Football, which is now played in 211 countries, including the USA, Canada, and Brazil.

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

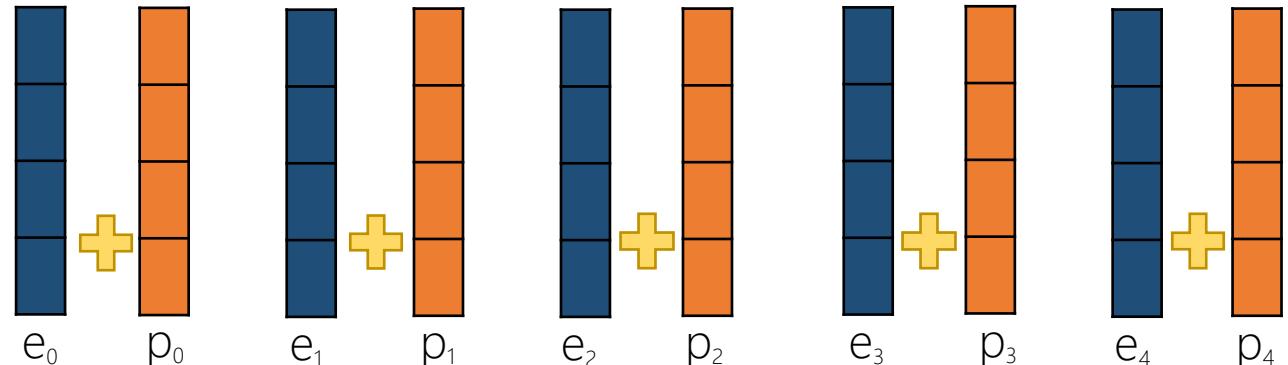


# What to use as Positional Embedding?



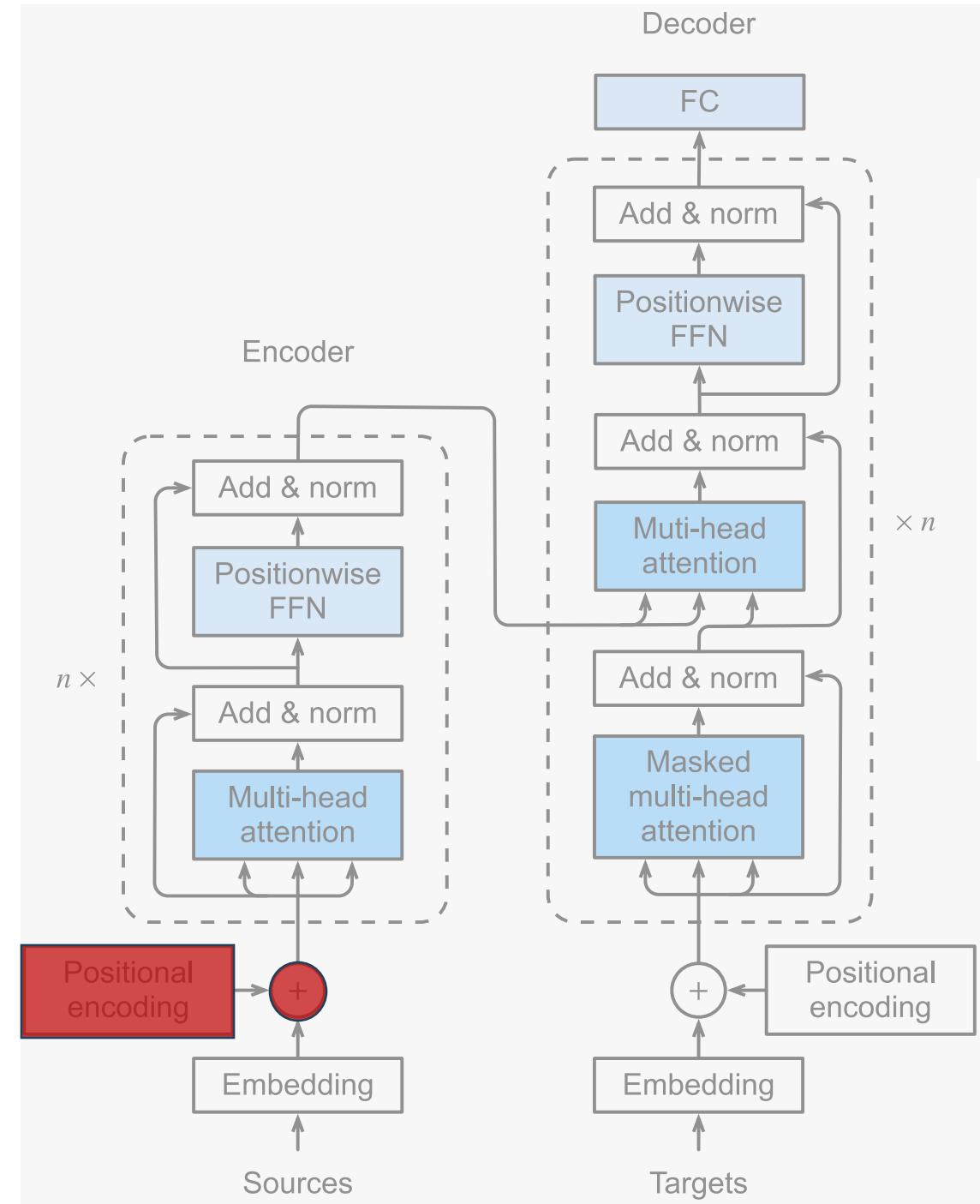


## Frequencies for Positional Embedding

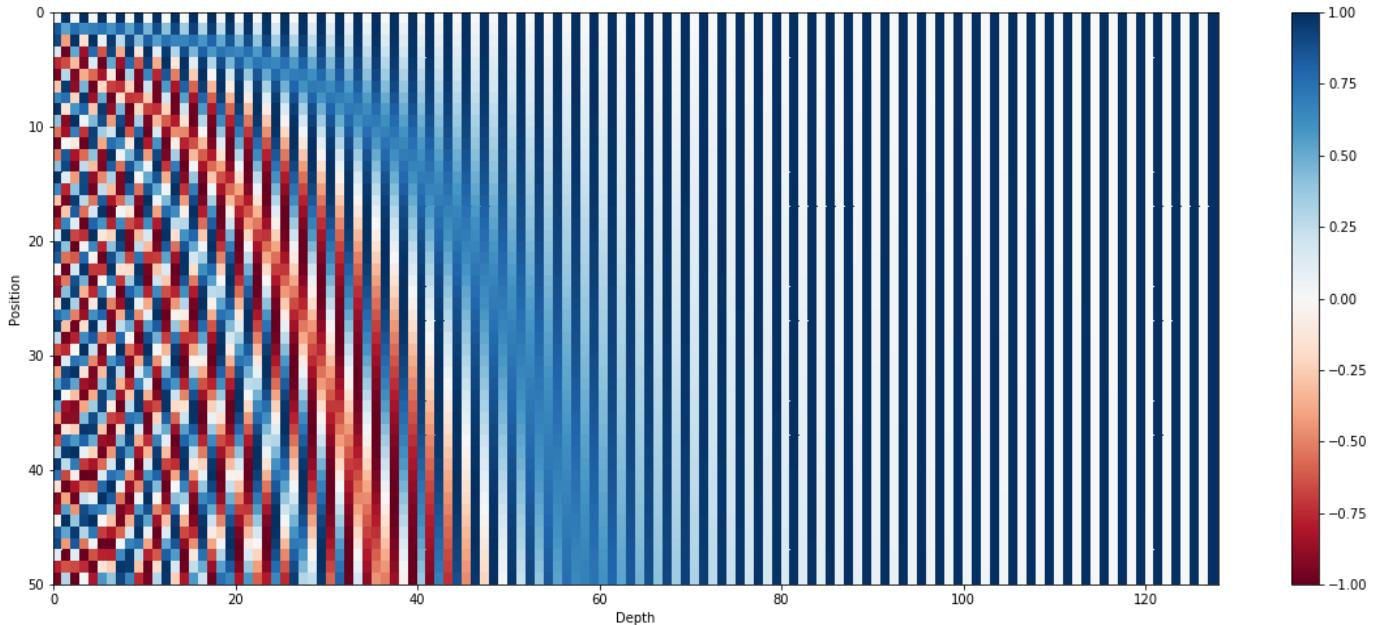


$$PE_{(pos,2i)} = \sin(pos/10000^{2i/d_{\text{model}}})$$

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



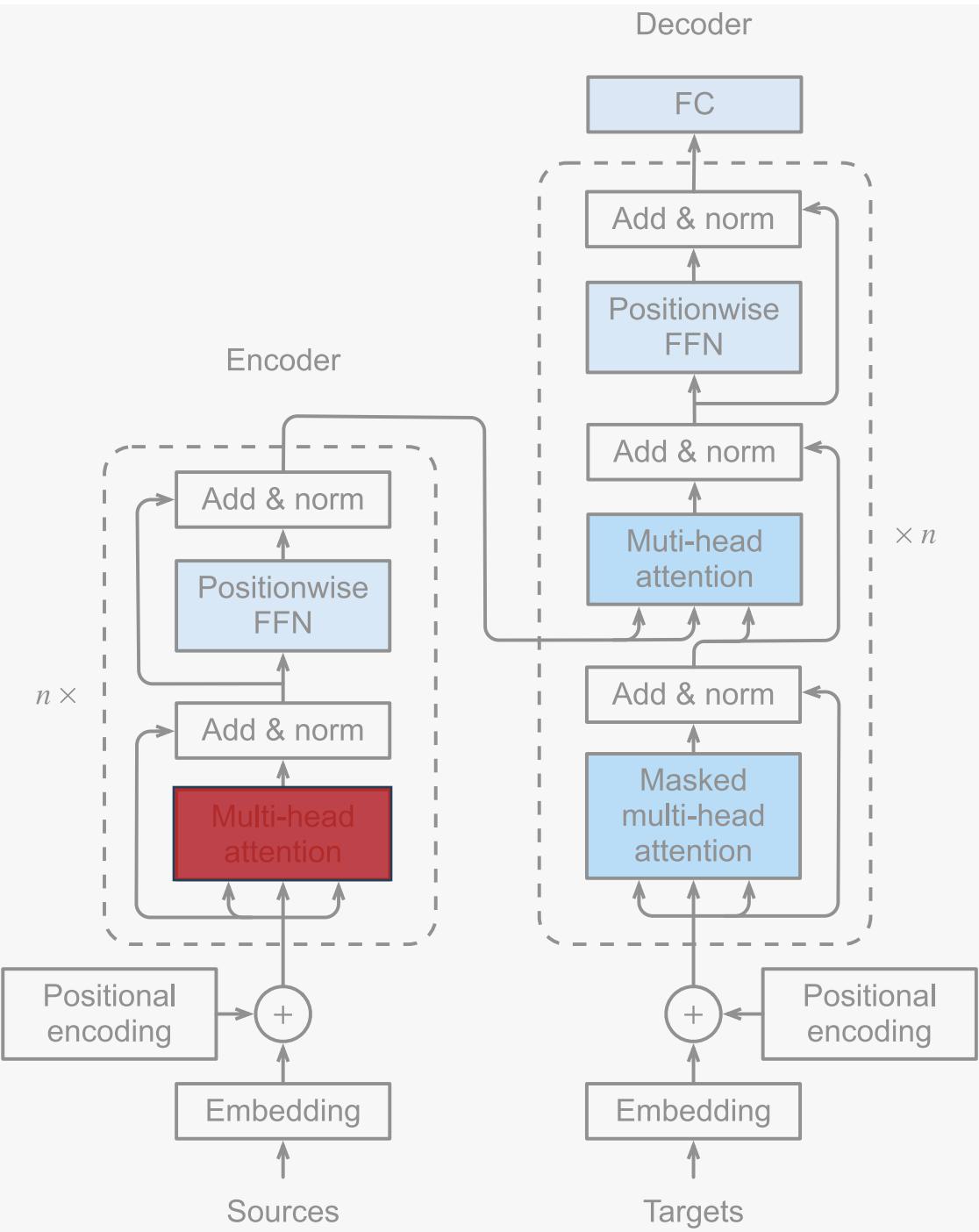
## Frequencies for Positional Embedding



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

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

# Multi-Head Attention



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

# Spotlight Effect!



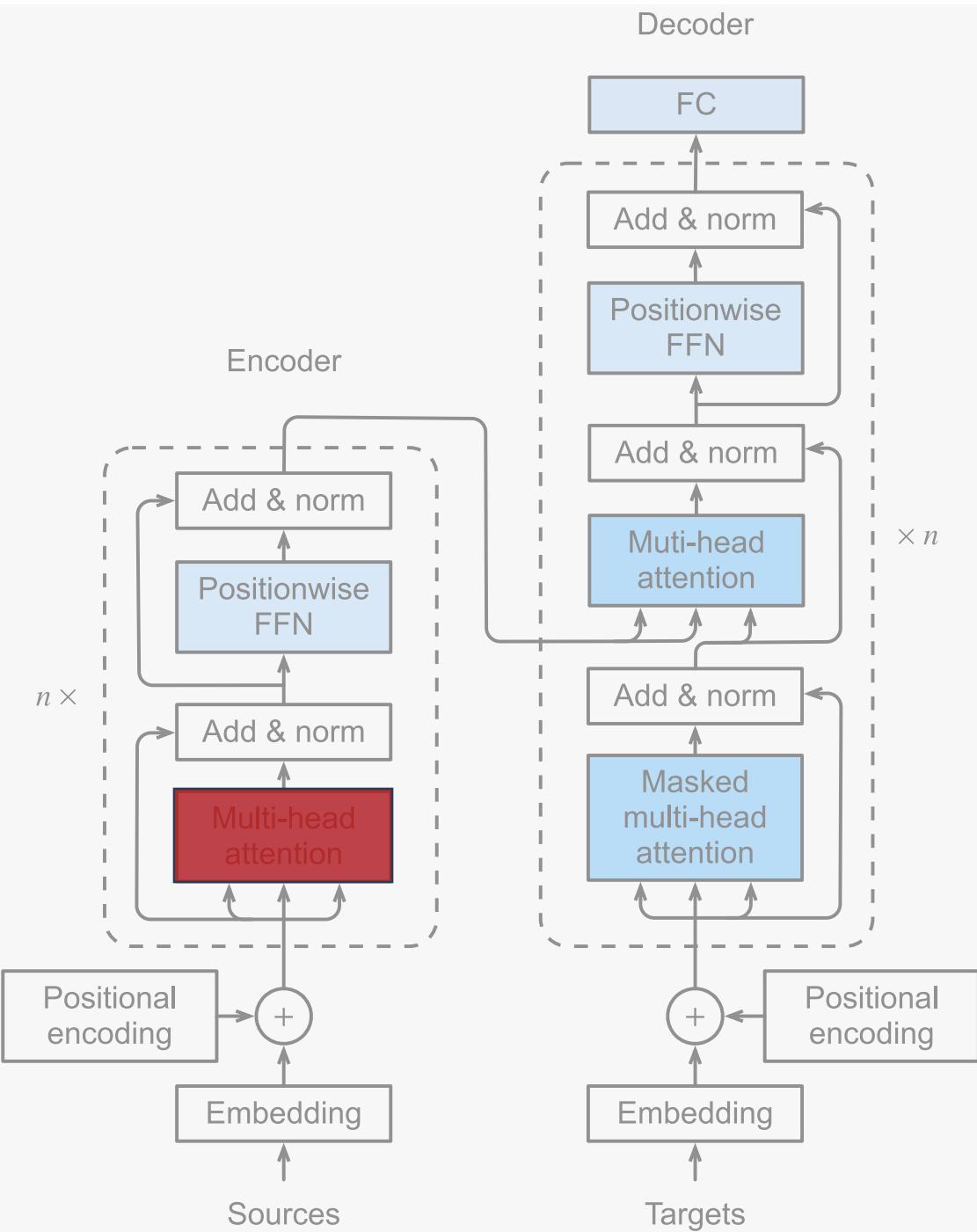
# Psychology lessons: Selective Attention

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





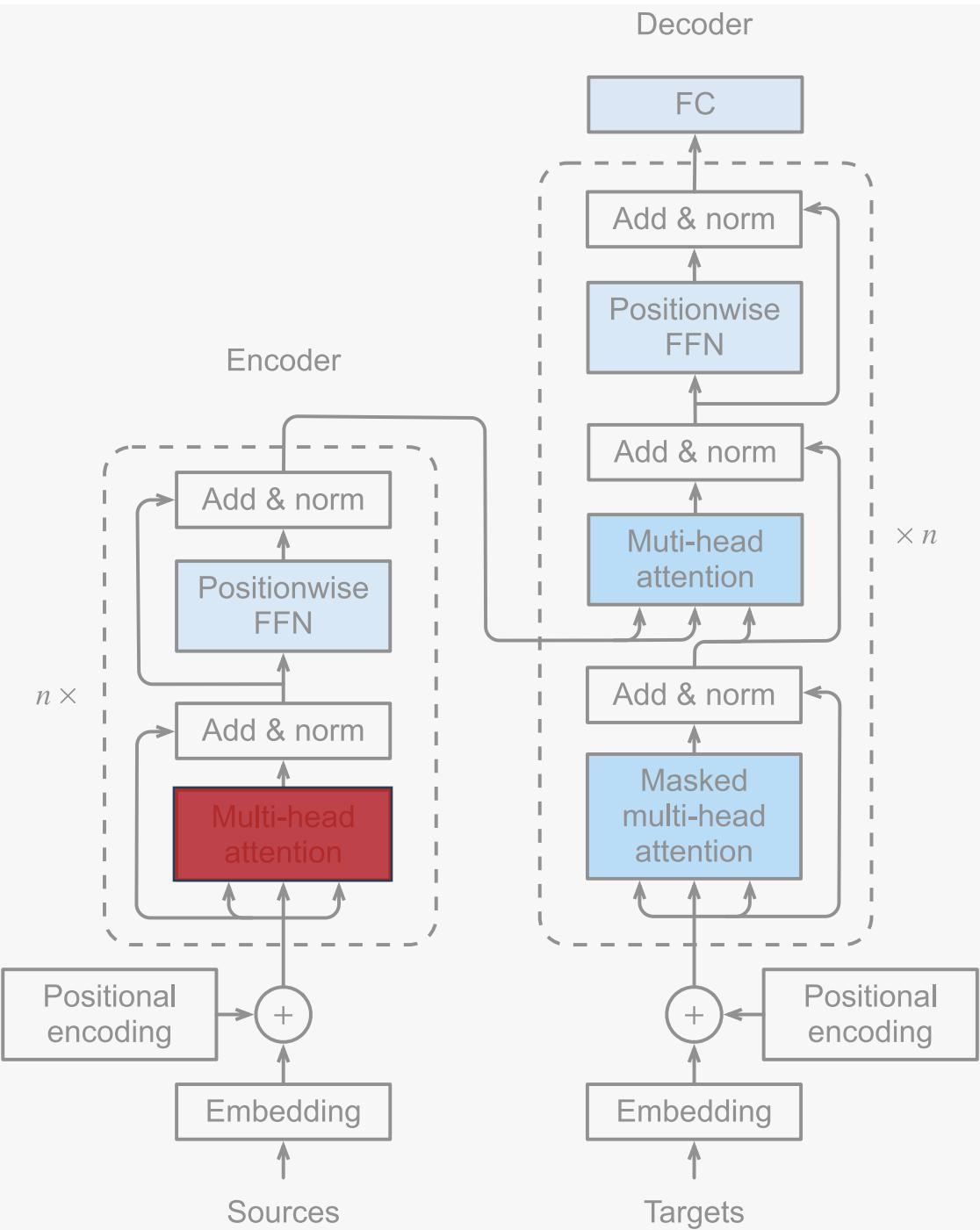
# Attention



Who is this Harry Potter character?

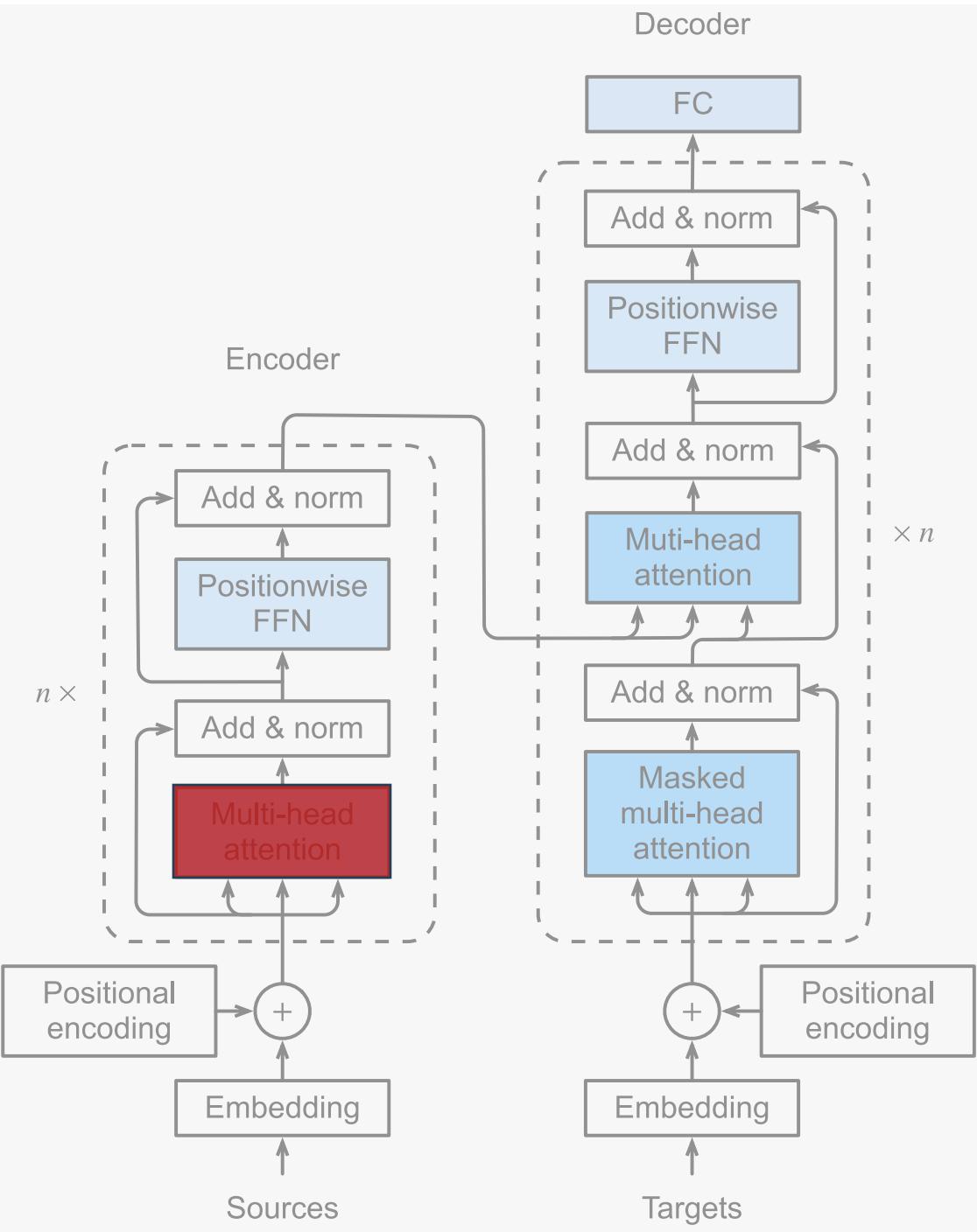
Now if you two don't mind, I'm going to bed before either of you come up with another clever idea to get us killed - or worse, expelled.

# Attention



Who is this Harry Potter character?

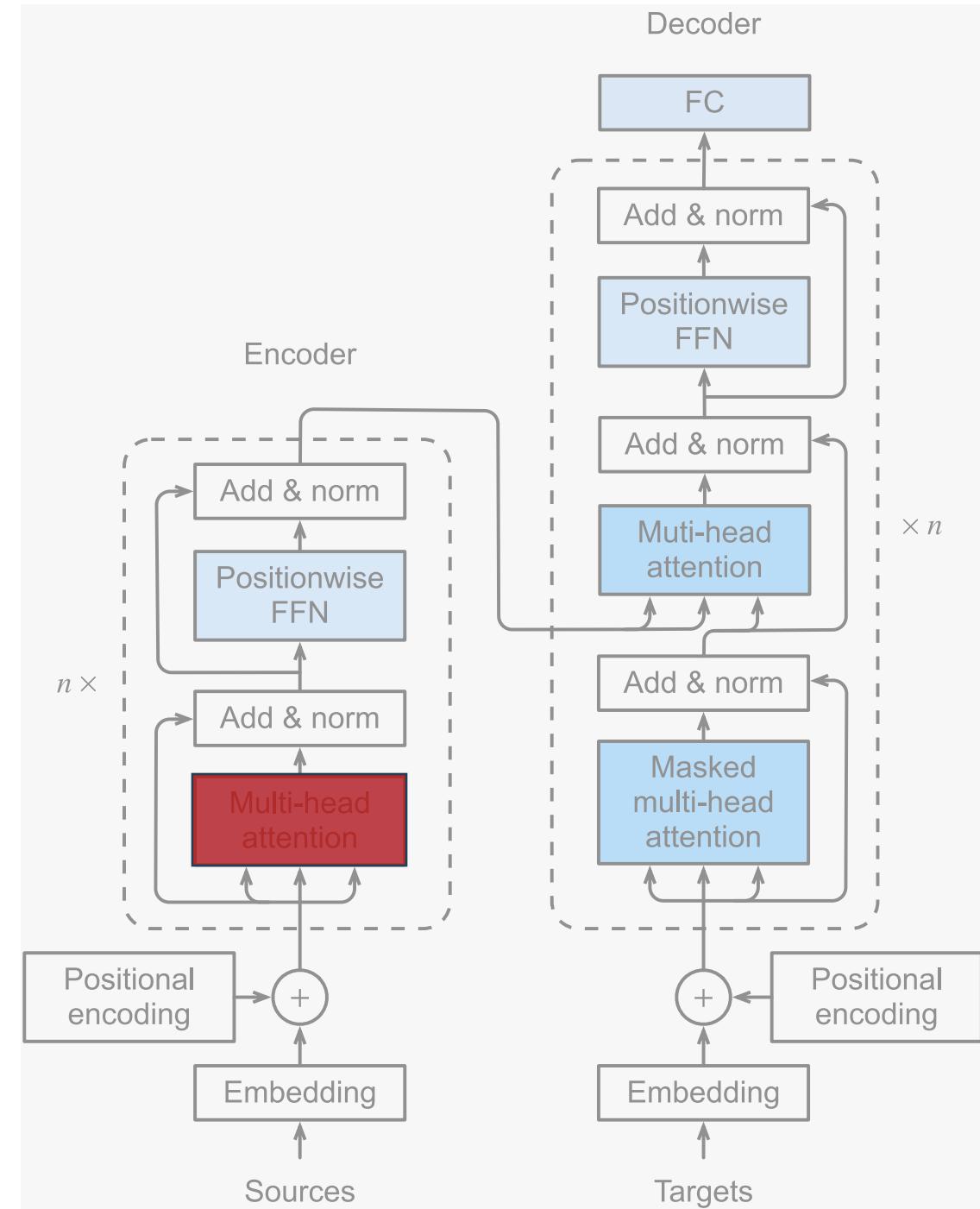
Now if you two don't mind, I'm going to bed before either of you come up with another clever idea to get us killed - or worse, expelled.



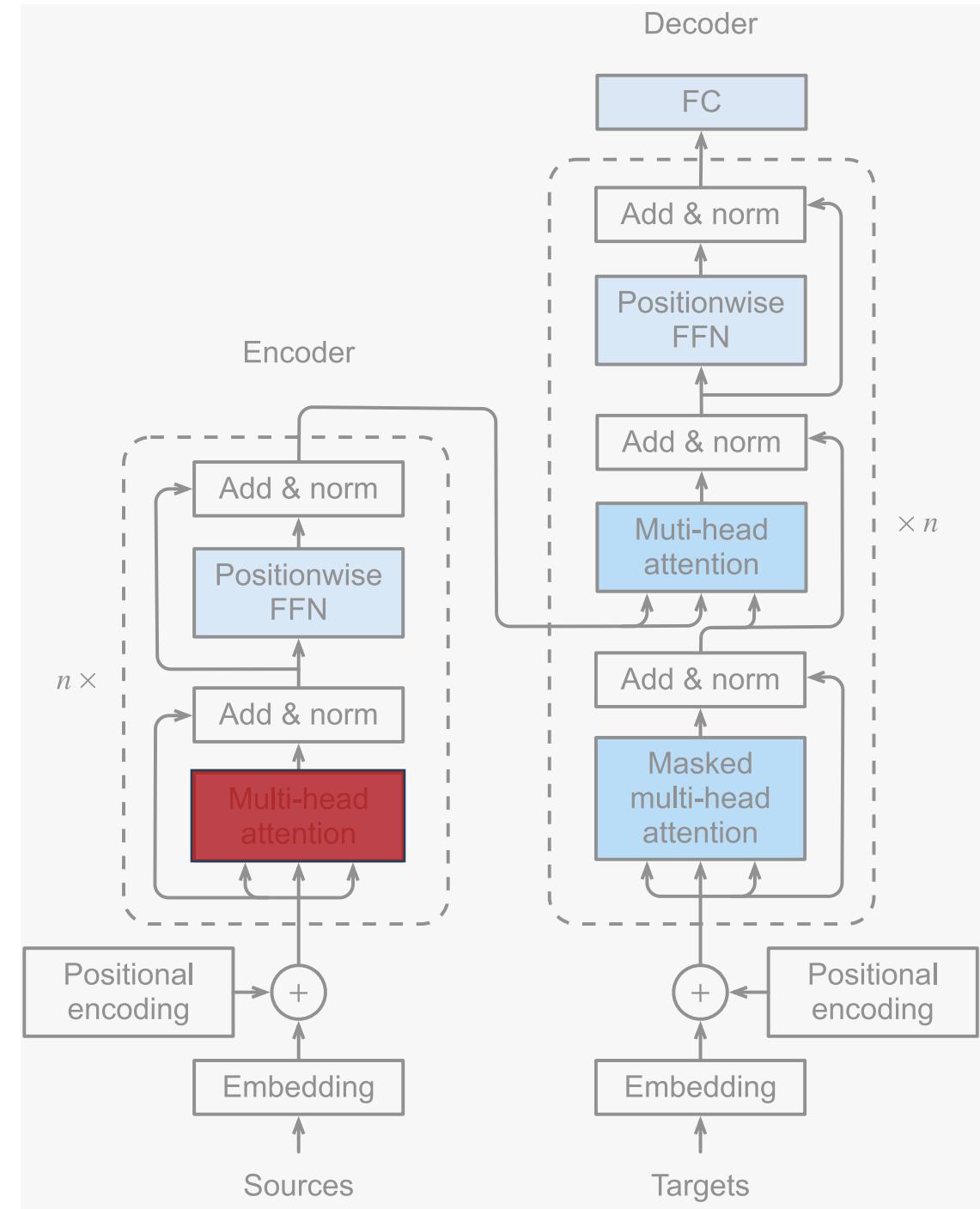
## Who is this Harry Potter character?

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

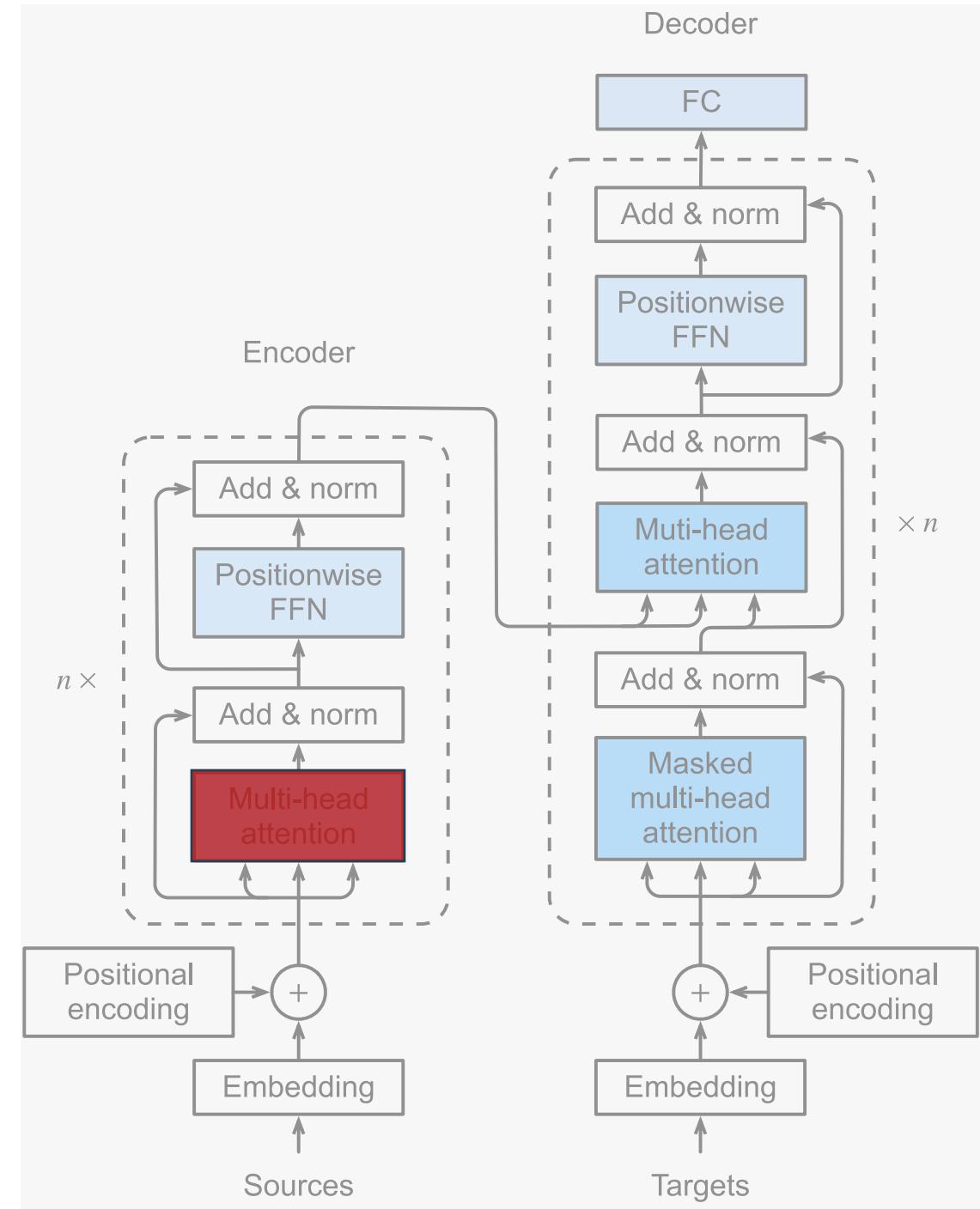


The dog began to **bark** loudly when it saw  
someone approaching the tree with rough **bark**.

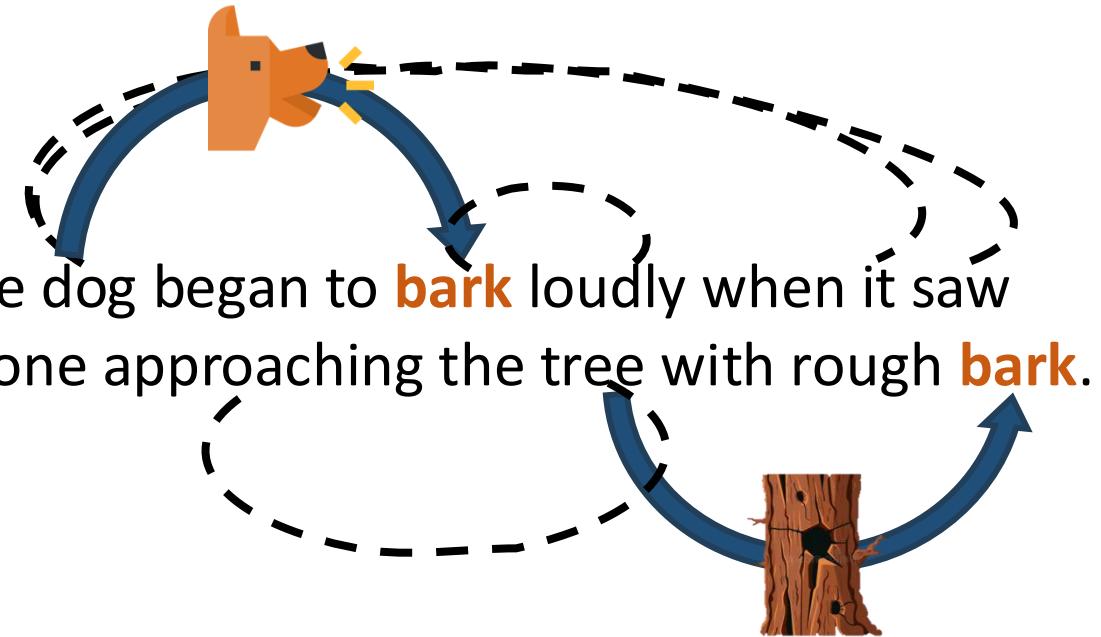


## Self-Attention





## Self-Attention



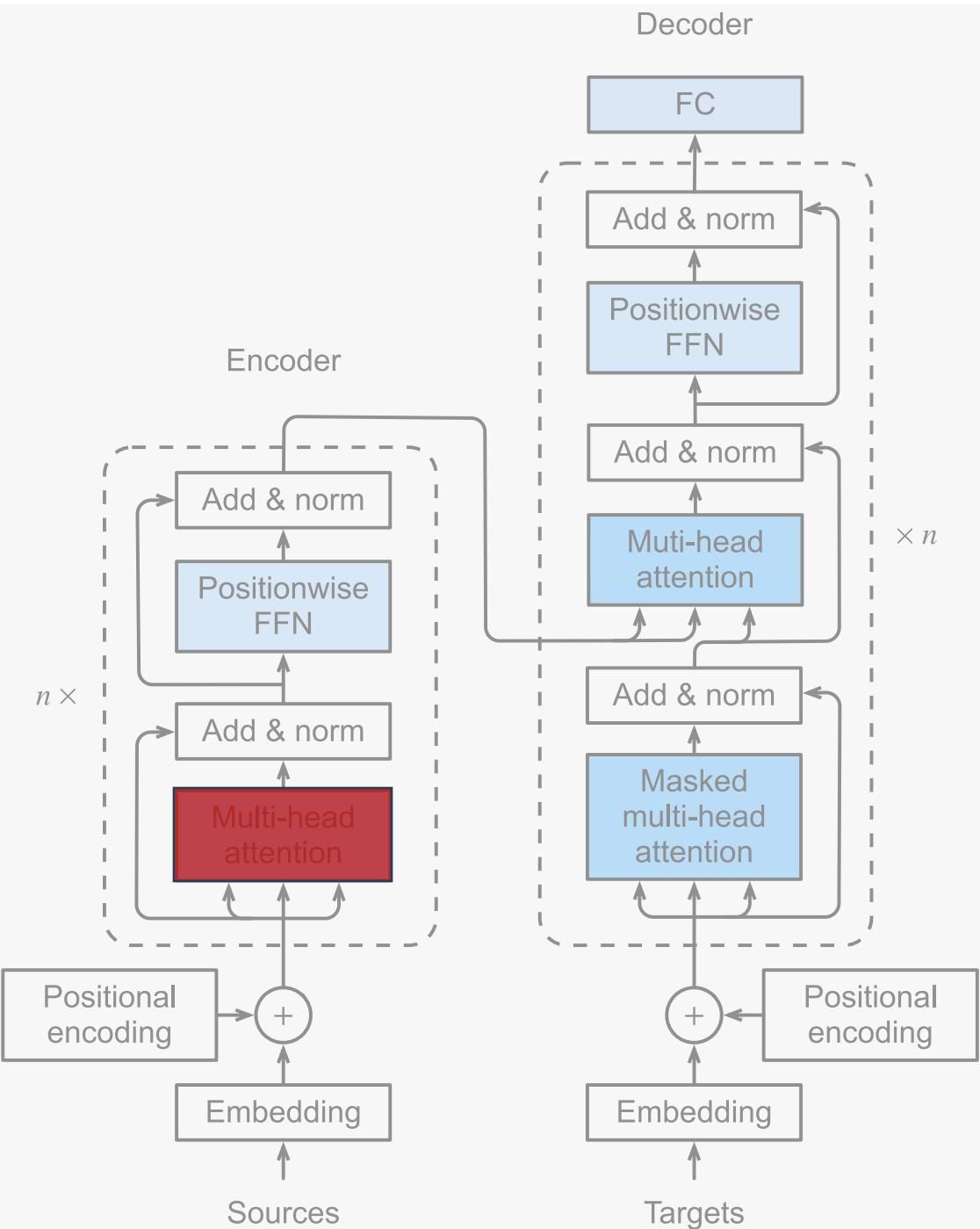
# Simple Attention

Who is this Harry Potter character?

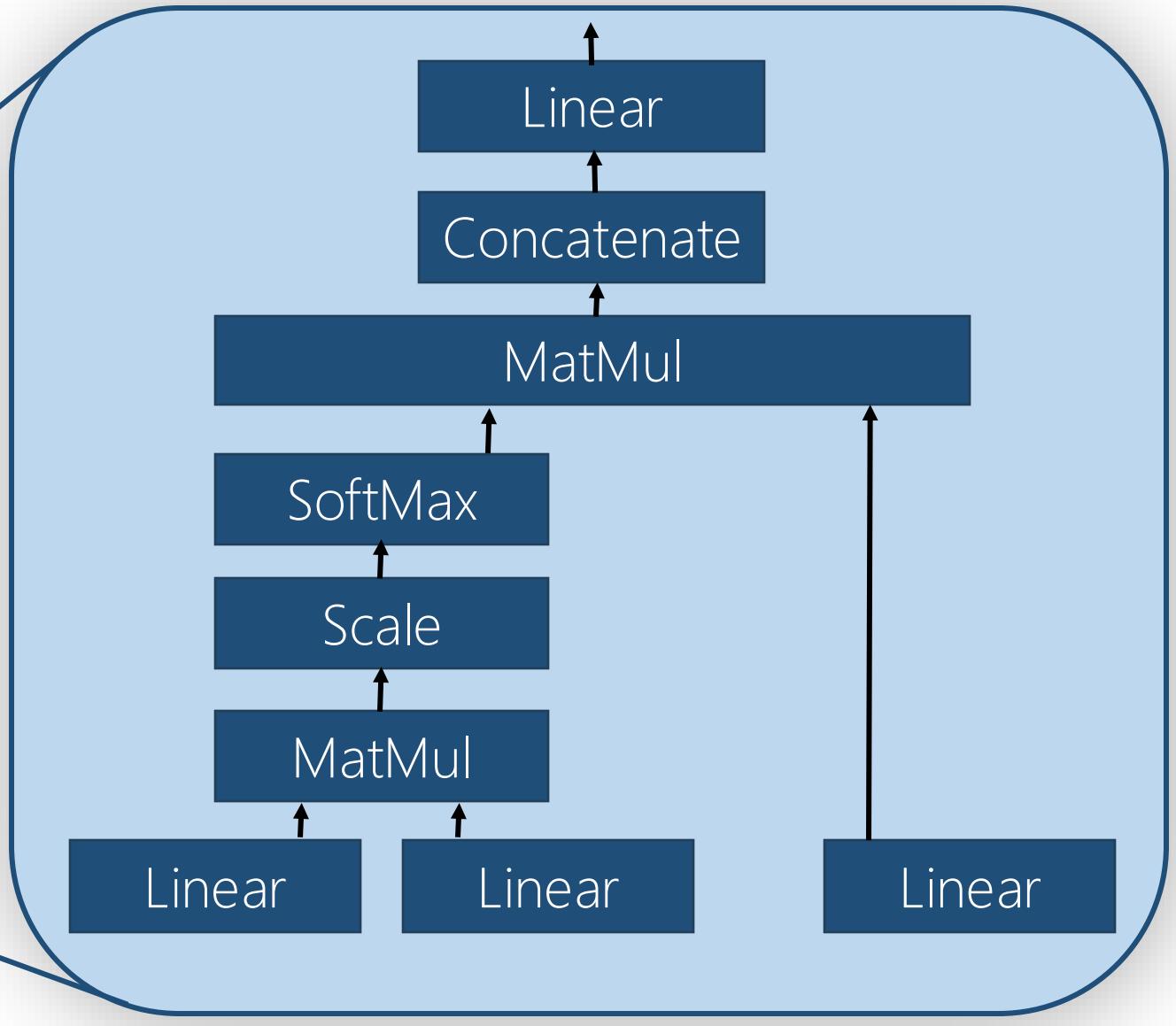
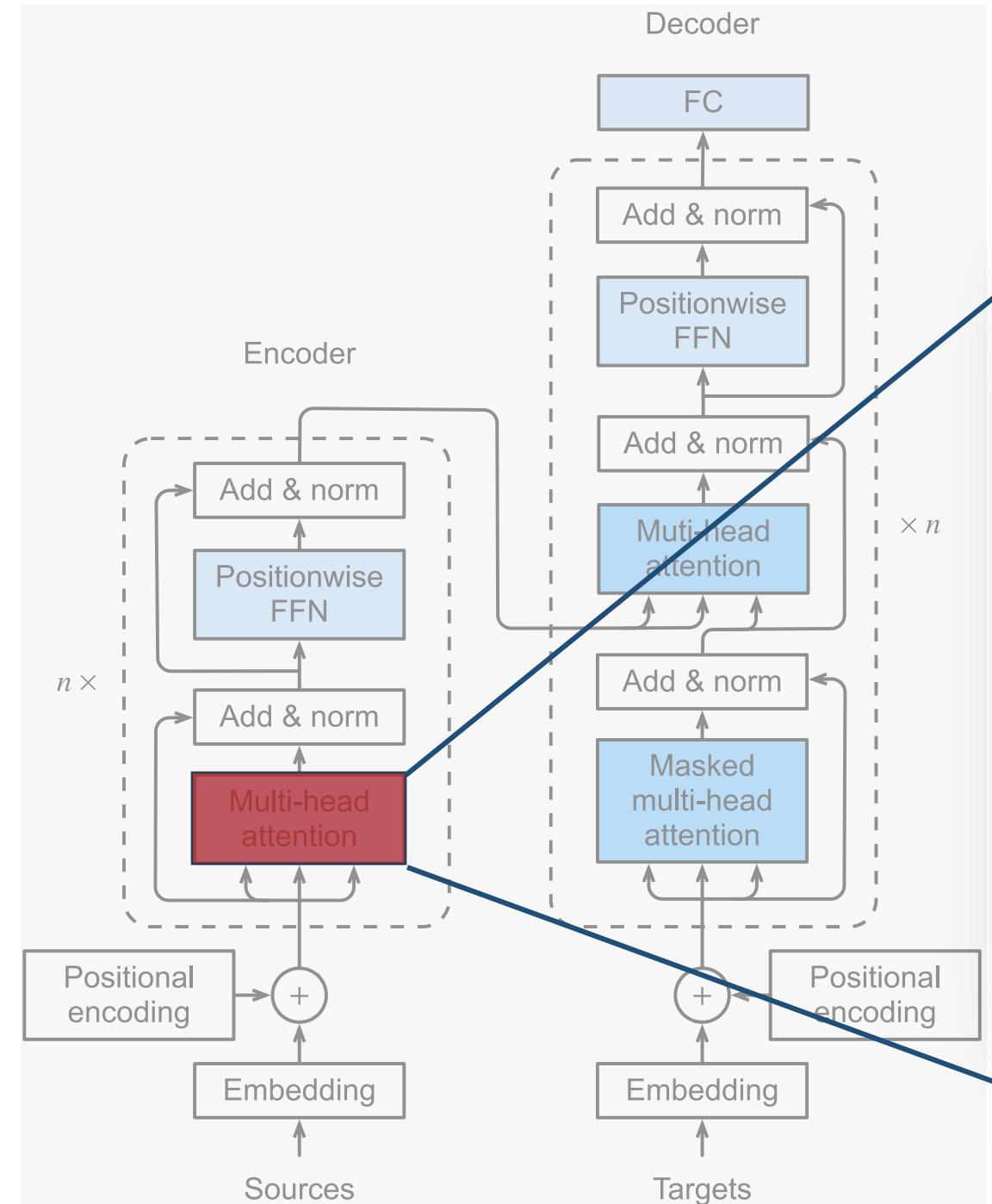
Now if you two don't mind, I'm going to bed before either of you come up with another clever idea to get us killed - or worse, expelled.

## Self-Attention

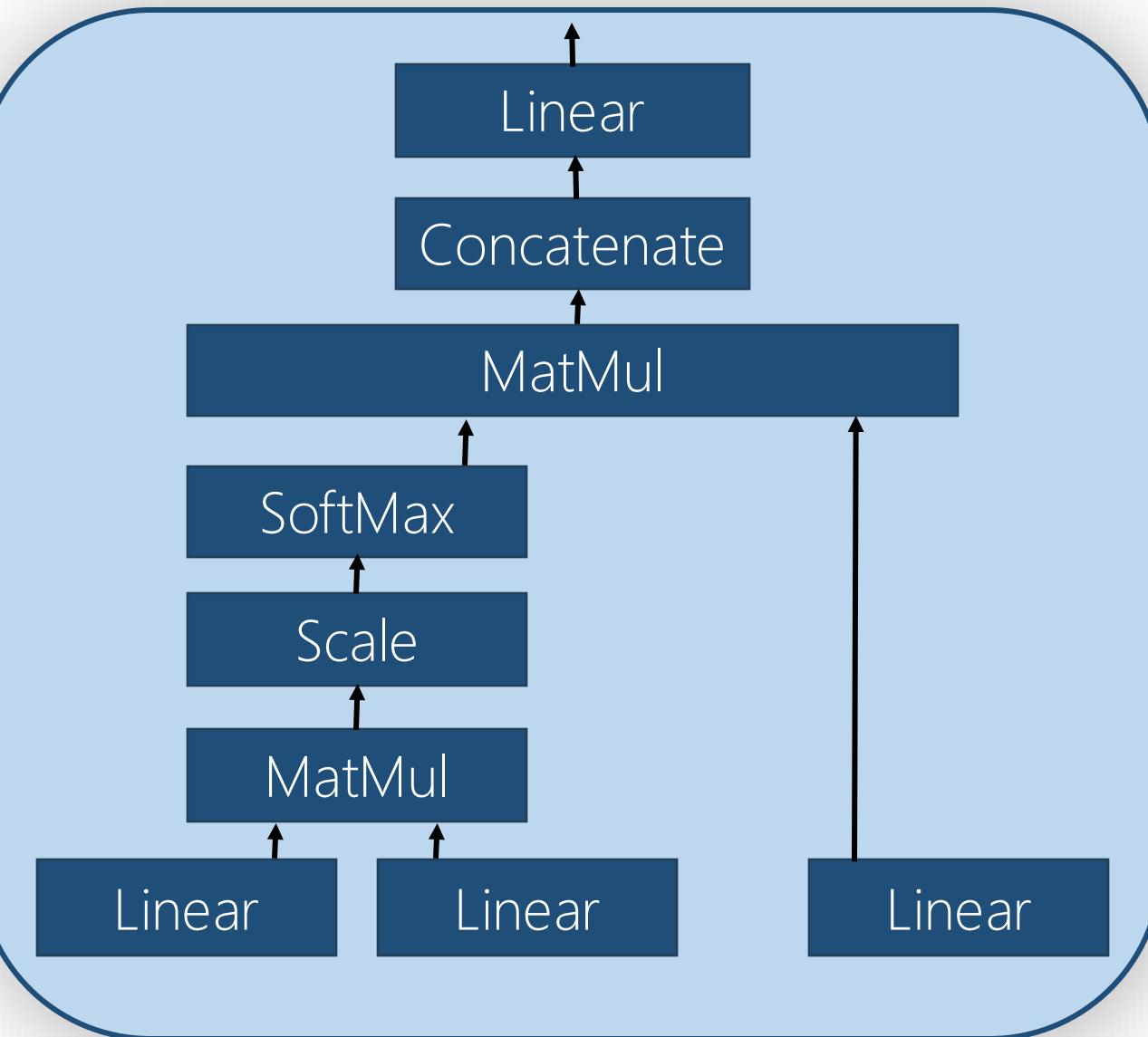
The dog began to bark loudly when it saw someone approaching the tree with rough bark.



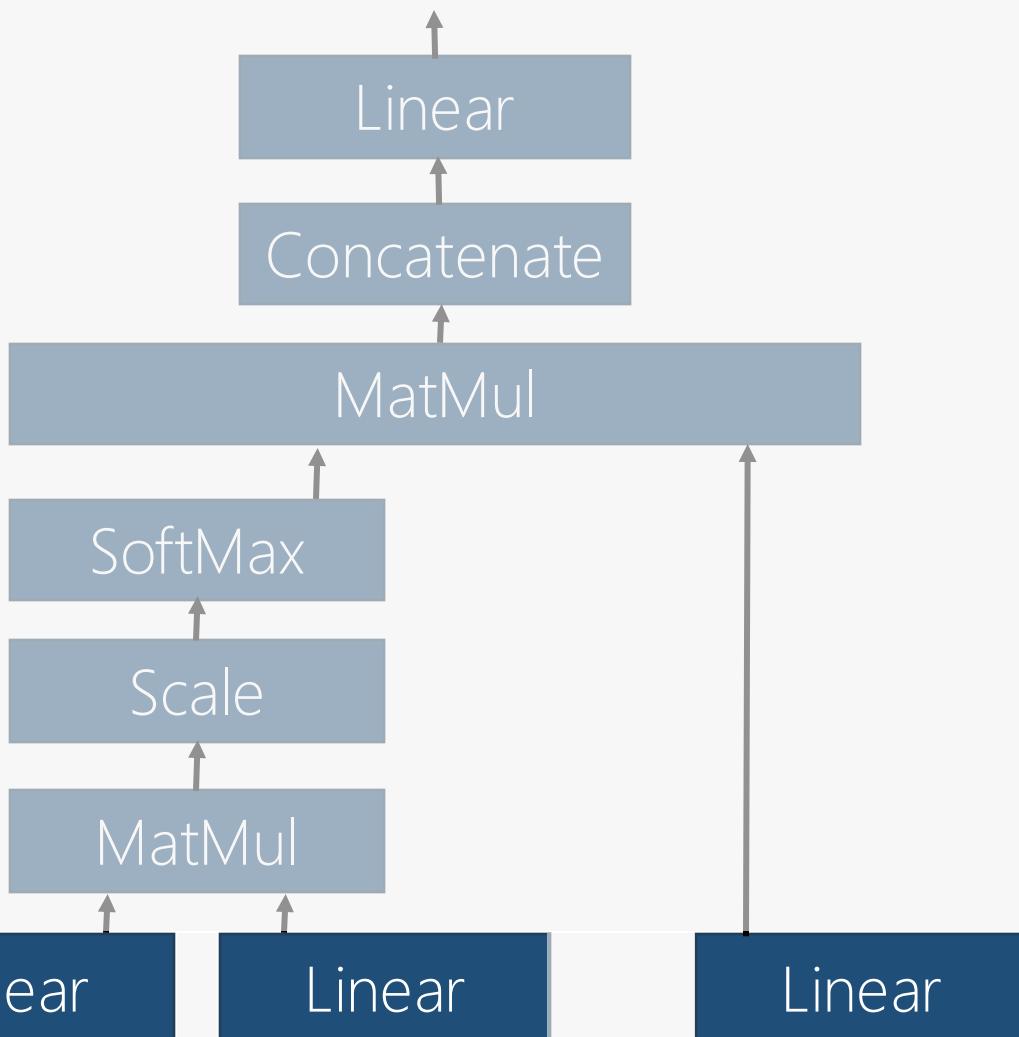
# Multi-Head Attention



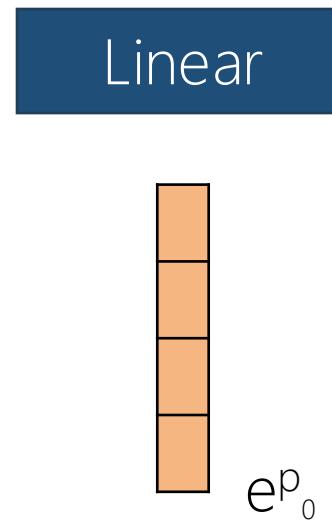
# Multi-Head Attention



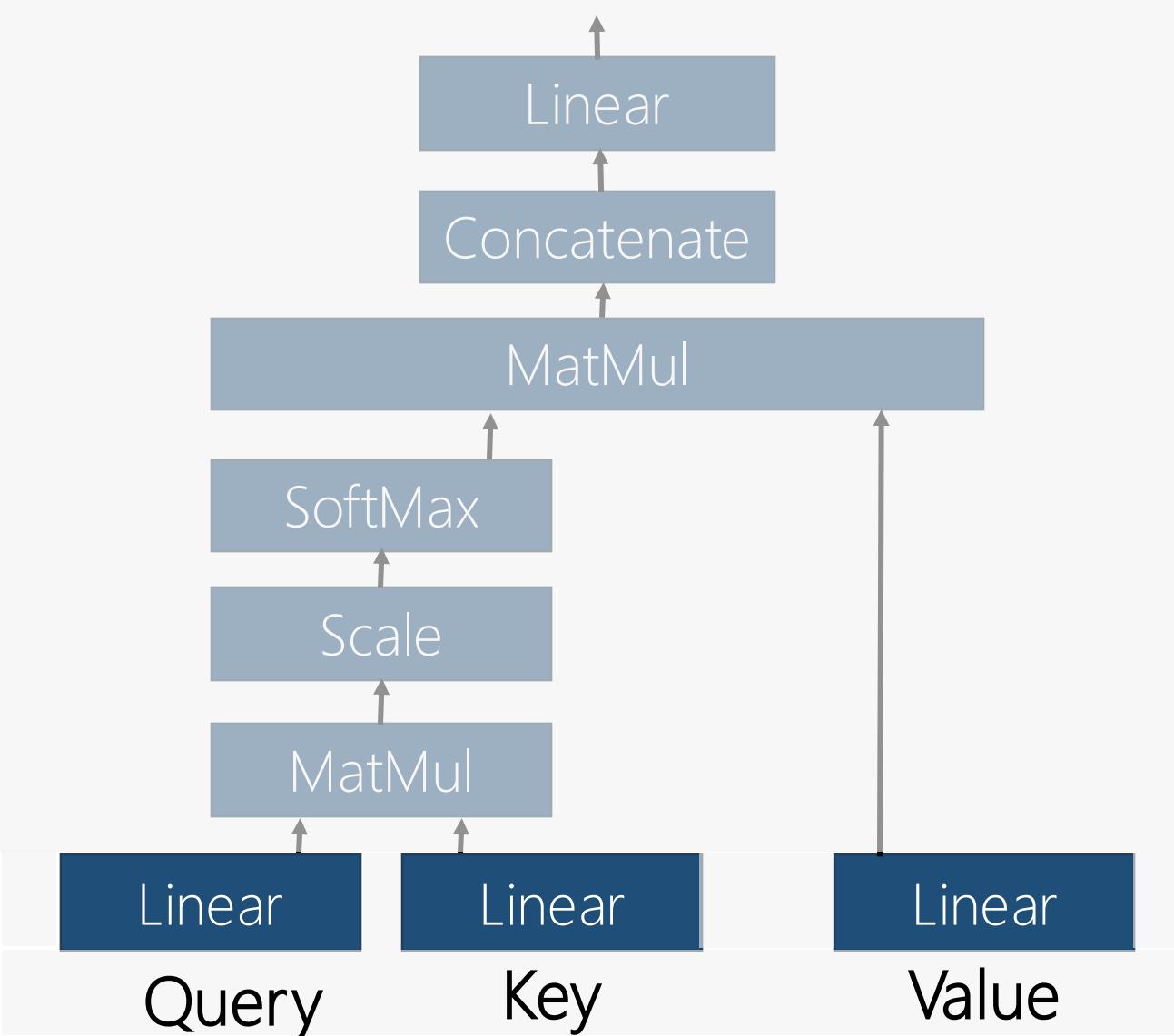
# Multi-Head Attention



What does a linear layer do?

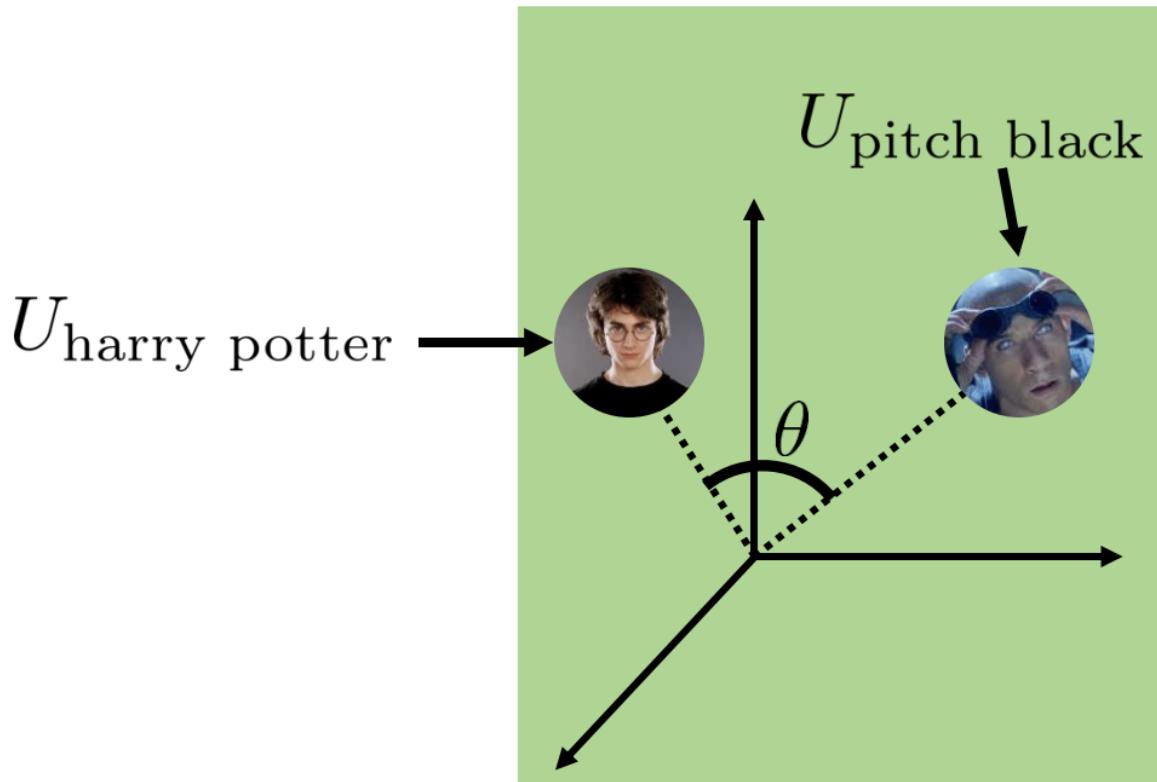


# Multi-Head Attention



Analogous to Data Retrieval/Search

# Similarity between Query and Key?



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

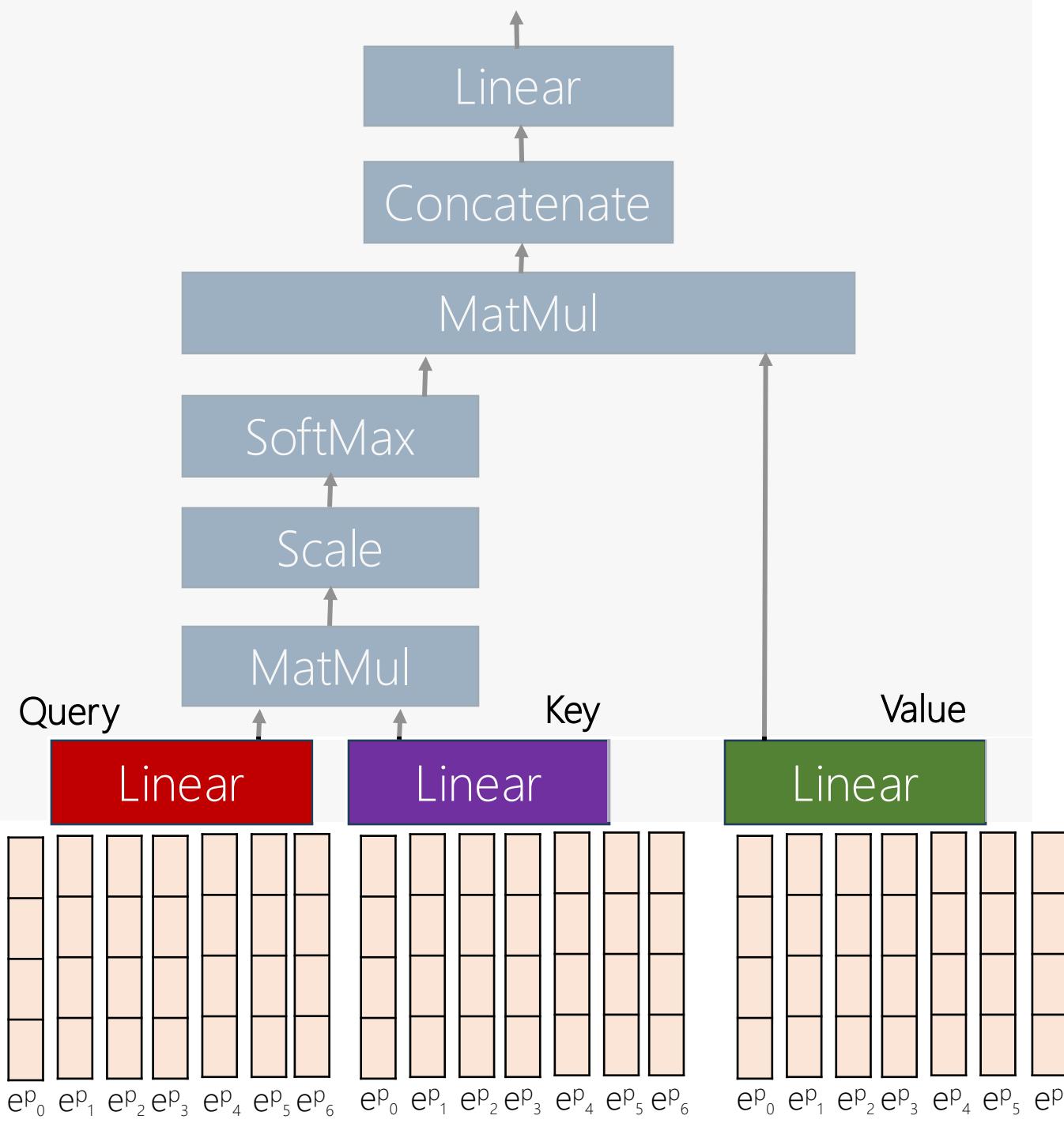
(theta = 0) →  $A$  and  $B$  point in exactly the same direction

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

(theta = 180) →  $A$  and  $B$  point in opposite directions (won't actually happen for 0/1 vectors)

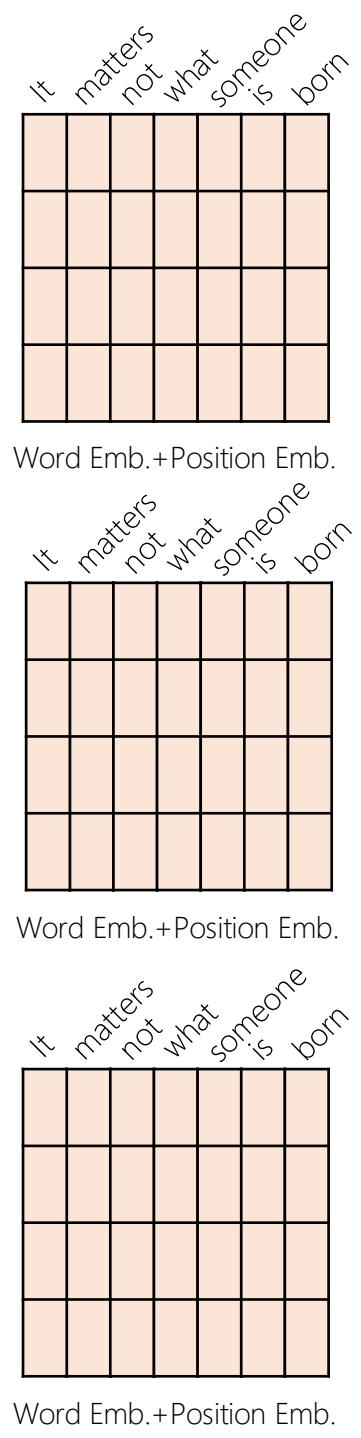
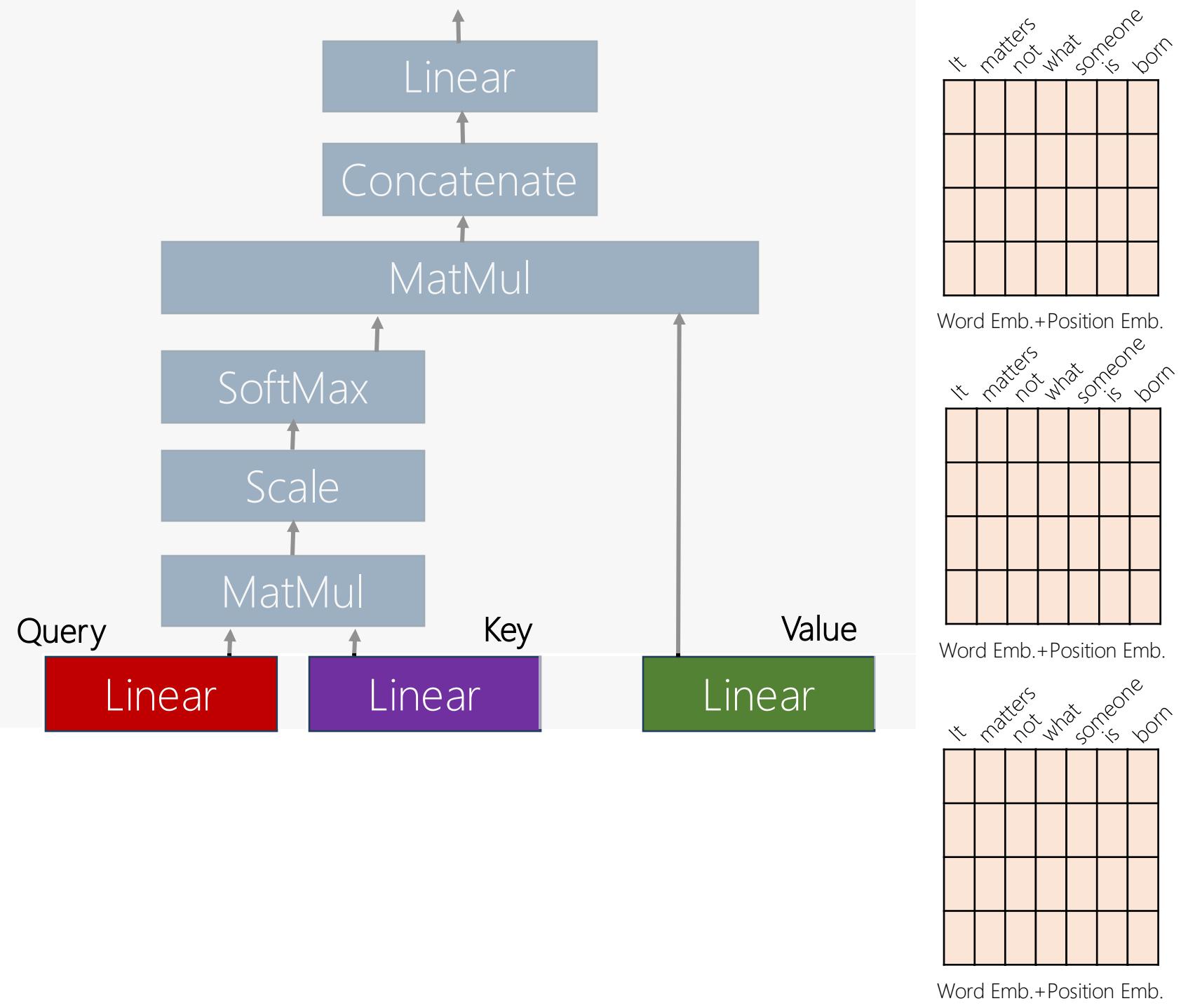
$$\cos(\theta) = 0$$

(theta = 90) →  $A$  and  $B$  are orthogonal

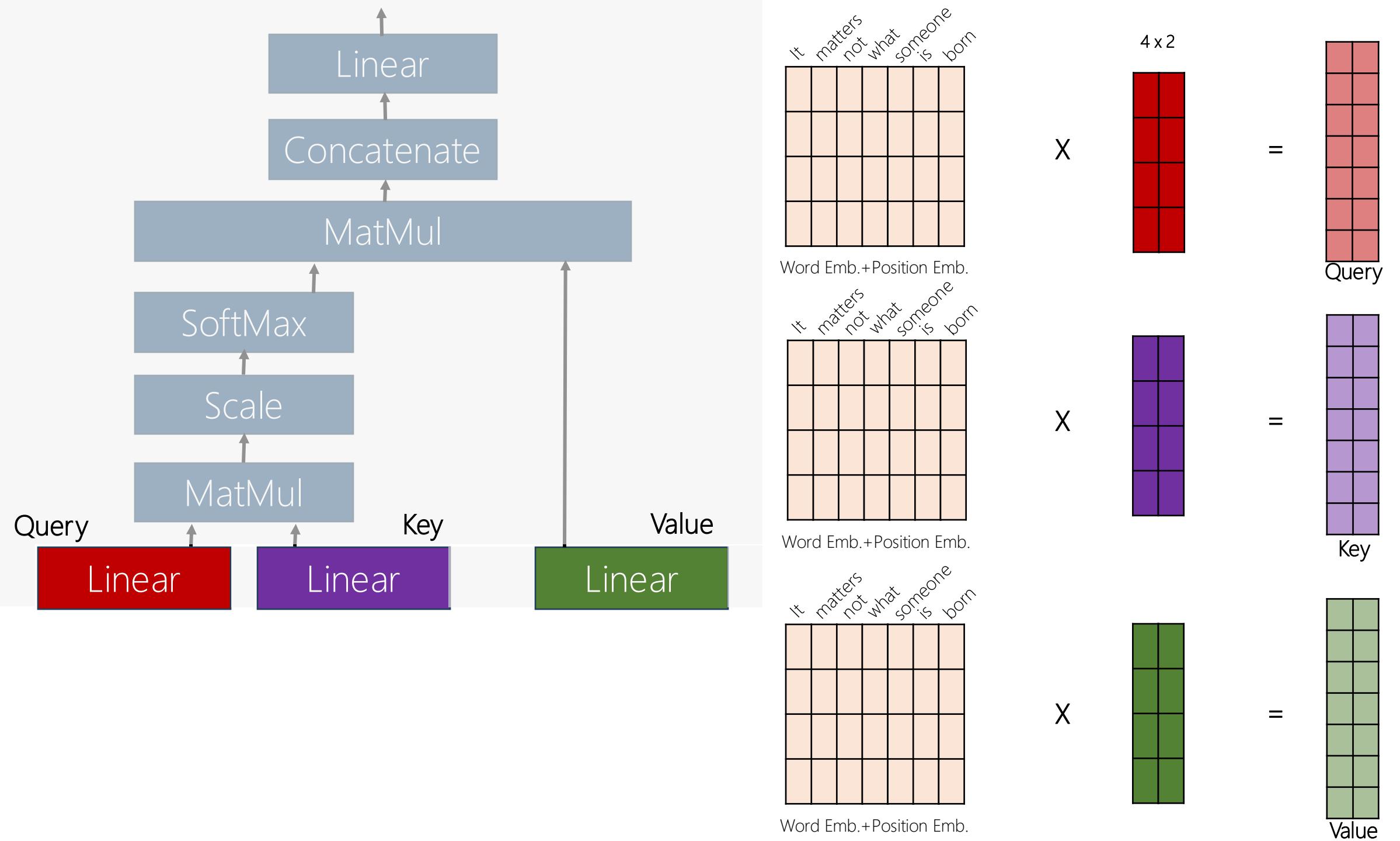


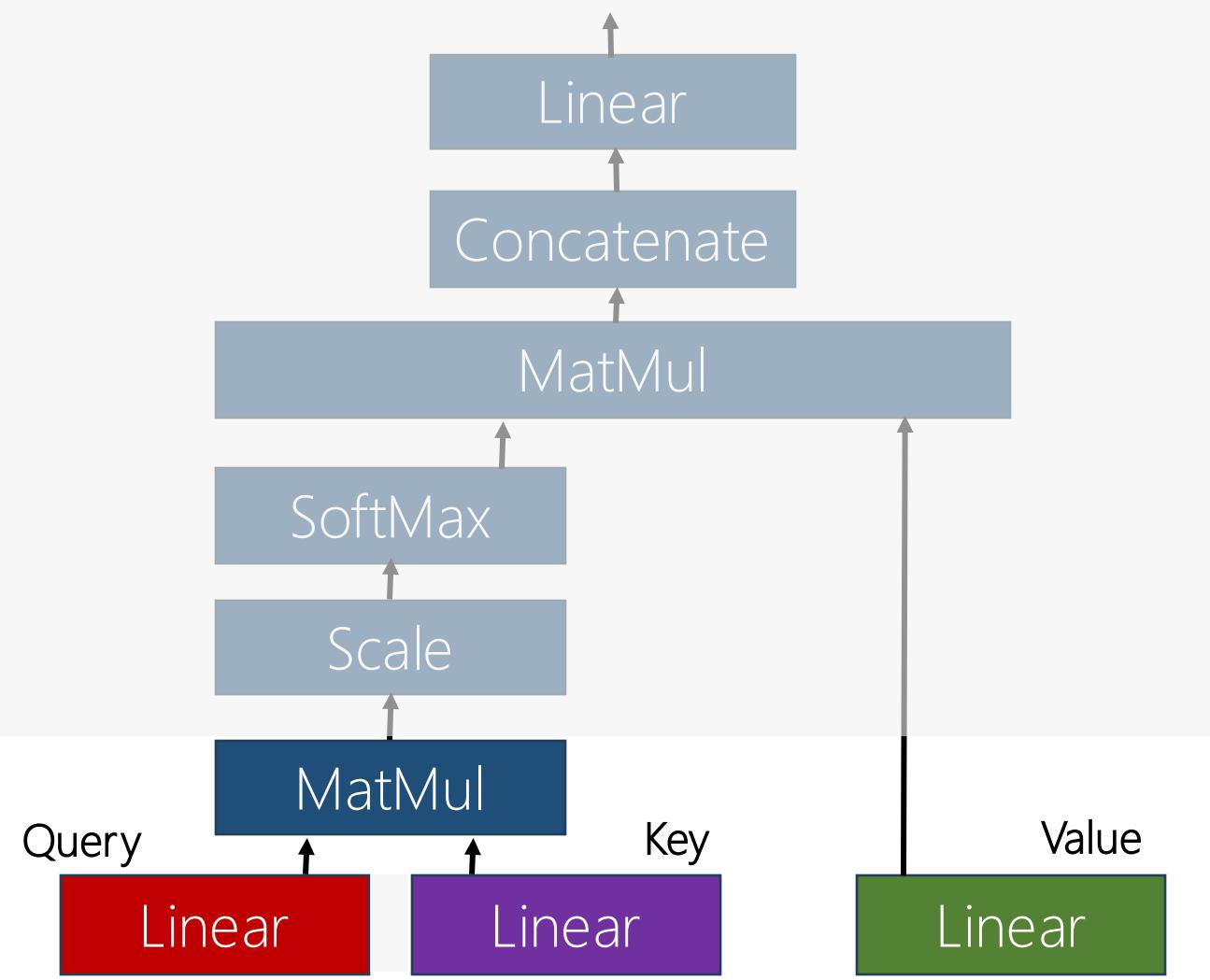
## Multi-Head Attention

Analogous to Data Retrieval/Search

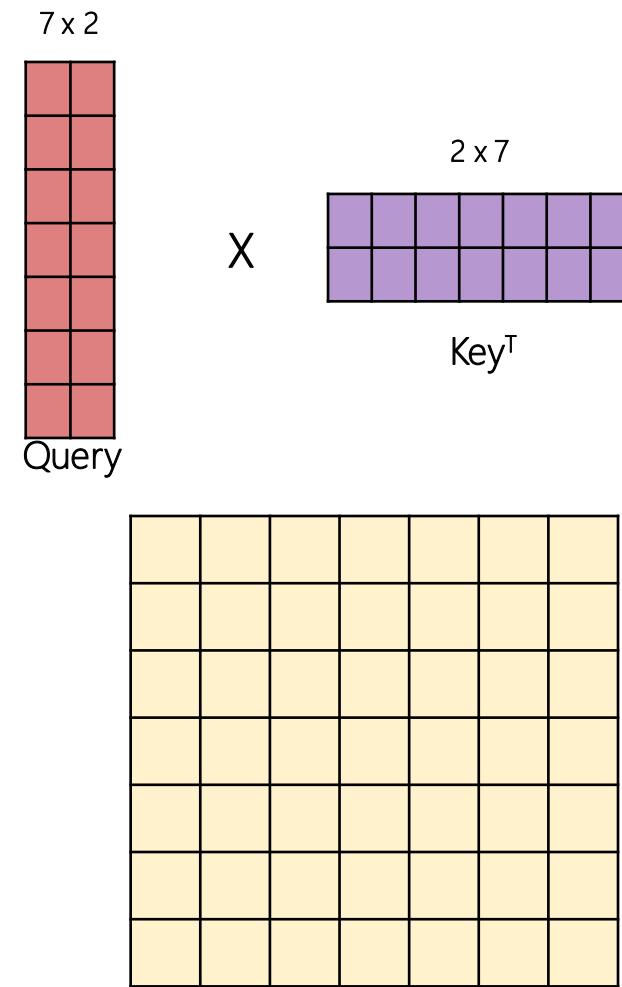


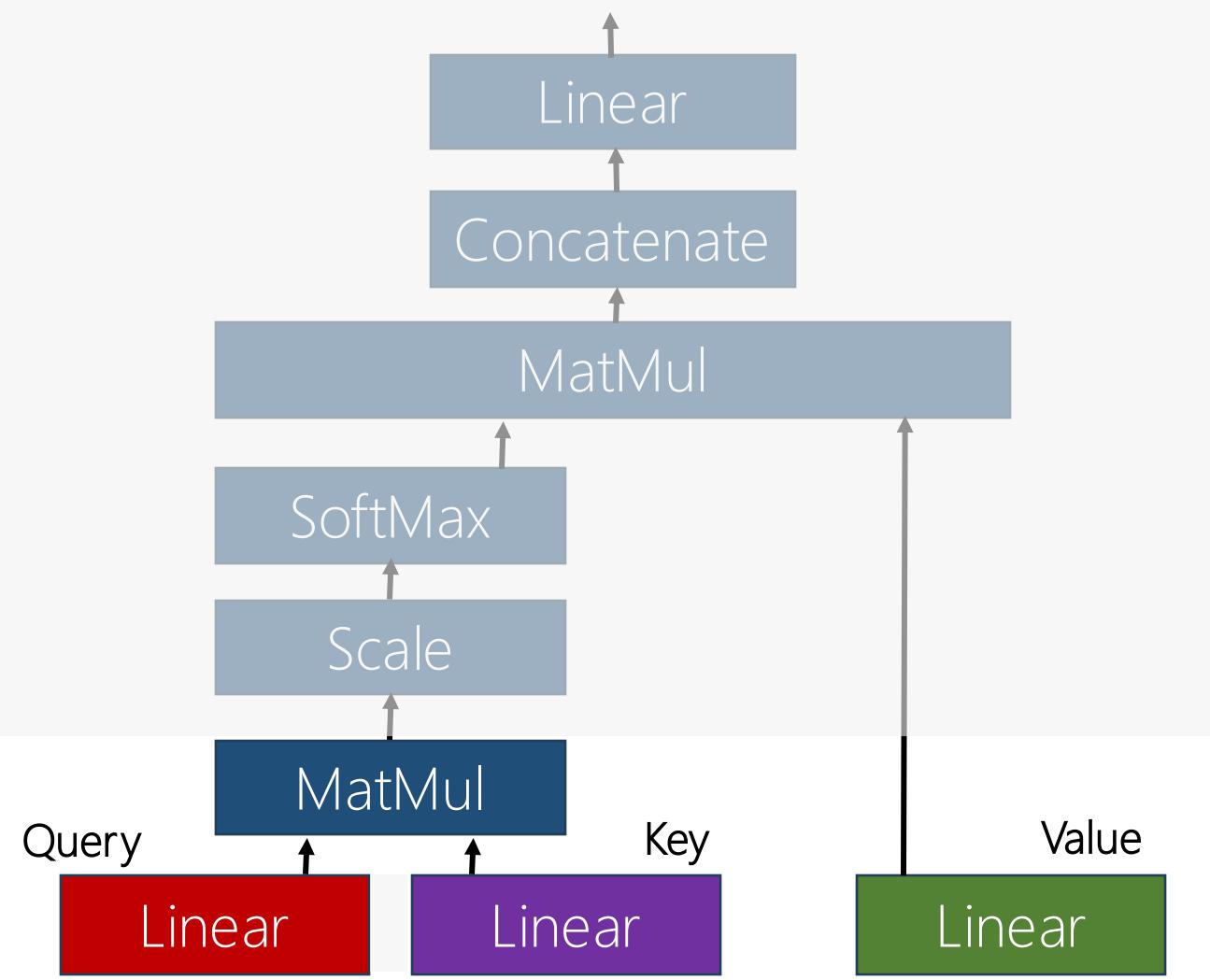
$$\begin{array}{c}
 4 \times 2 \\
 \times \quad = \\
 \textcolor{red}{\boxed{\text{Red}}} \quad \textcolor{purple}{\boxed{\text{Purple}}} \quad \textcolor{green}{\boxed{\text{Green}}}
 \end{array}$$



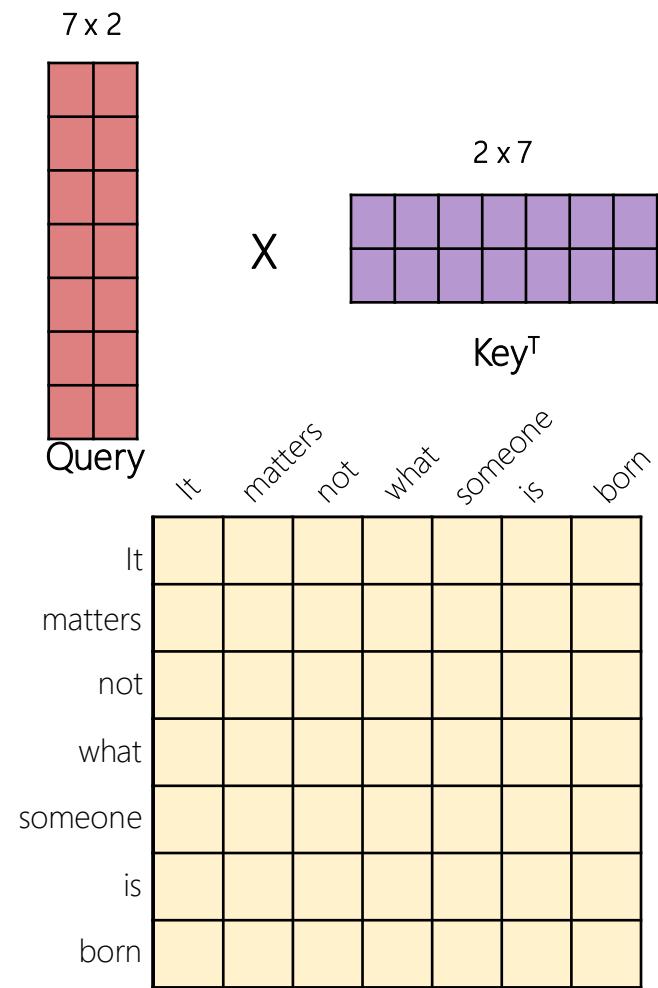


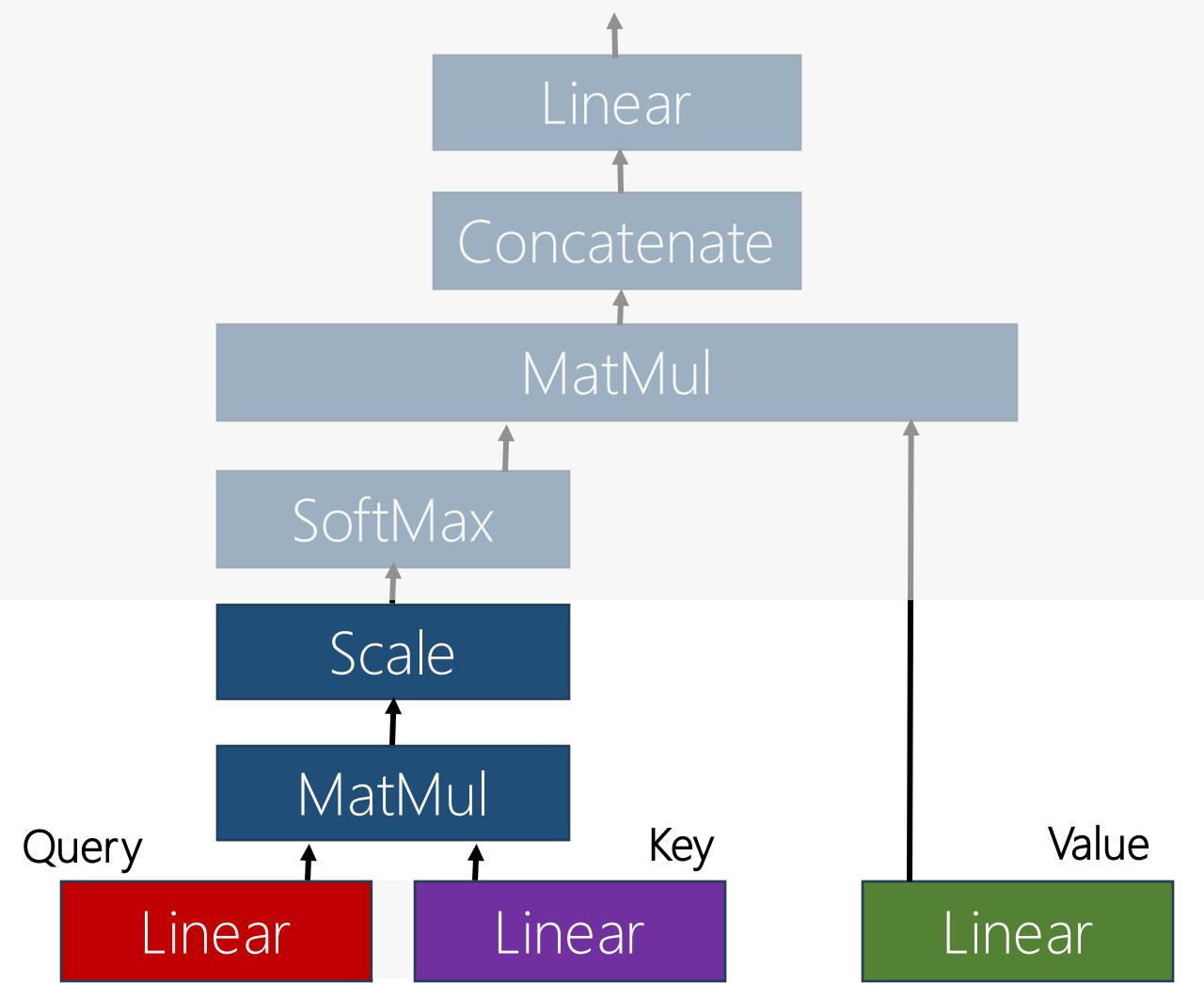
## Attention Filter



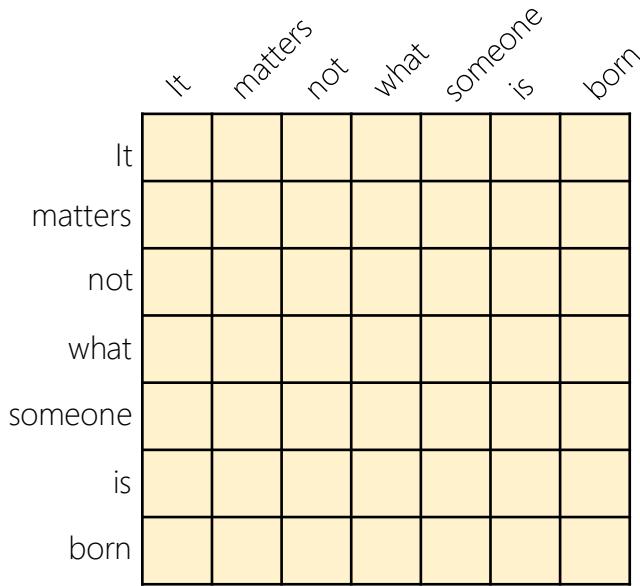


## Attention Filter



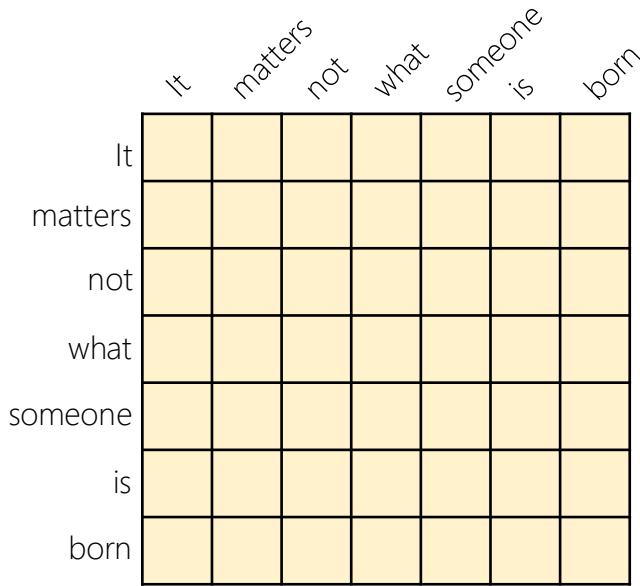
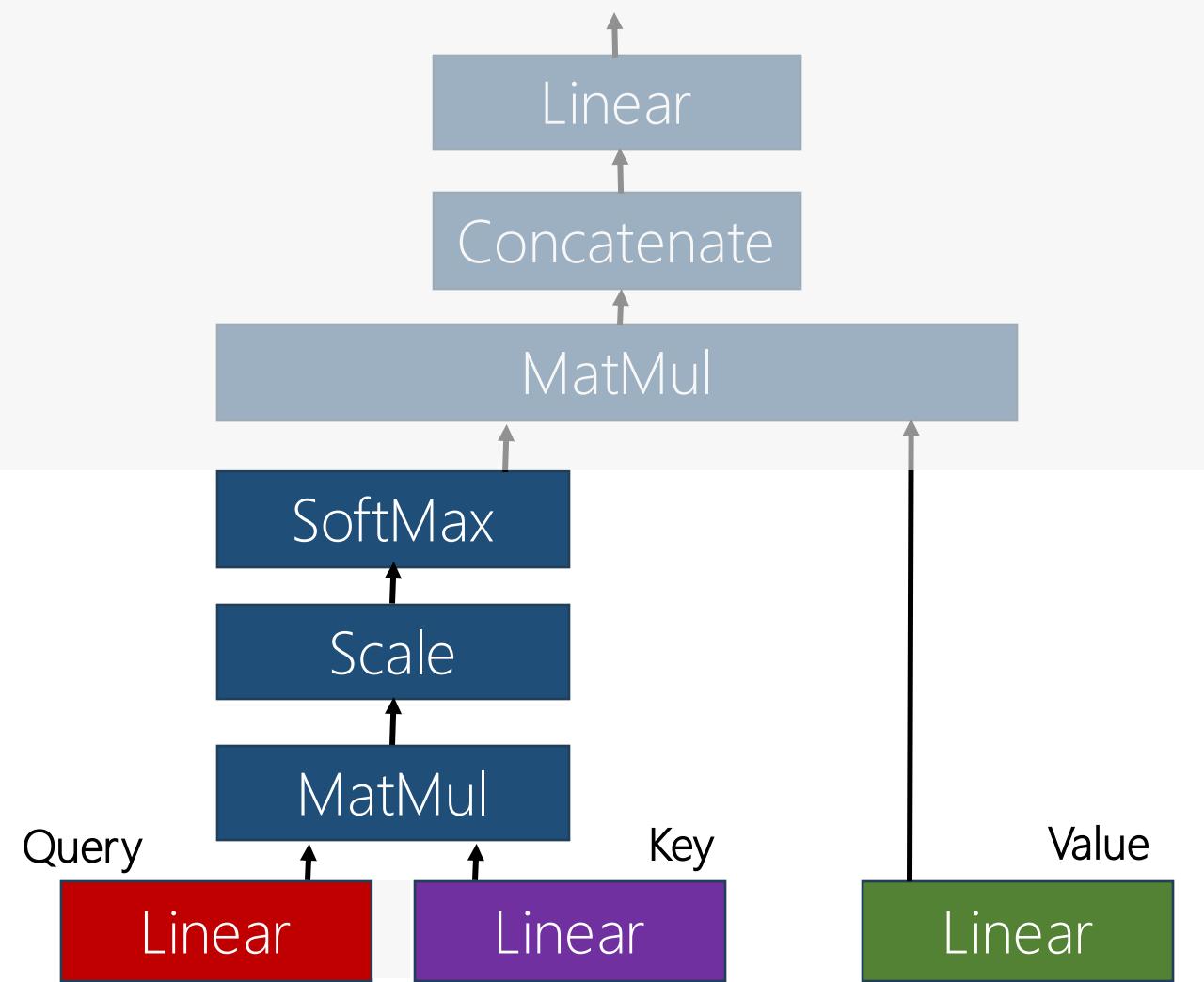


## Attention Filter

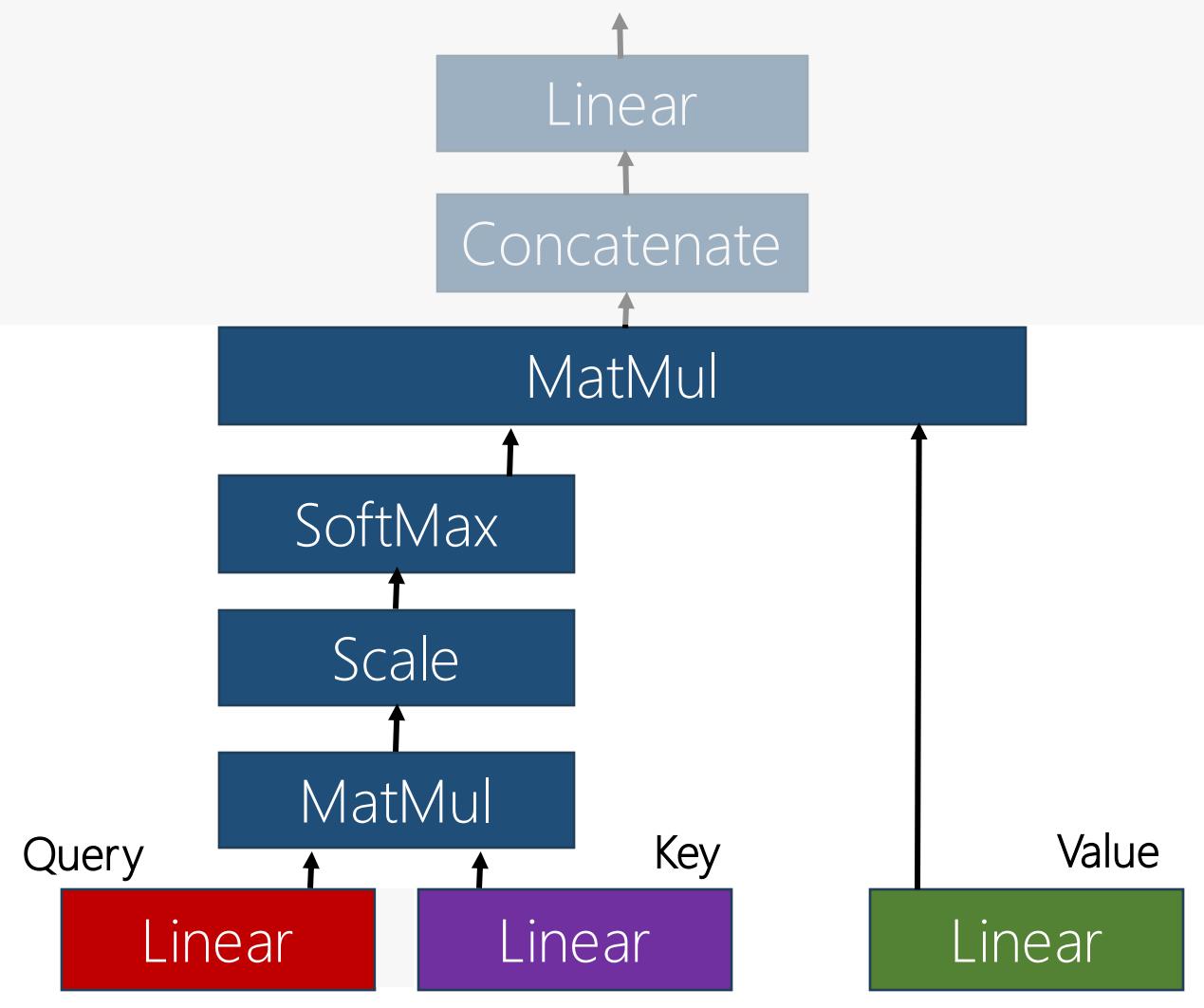


Scale using dimension!

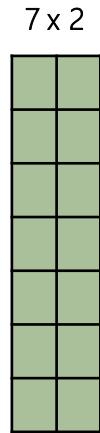
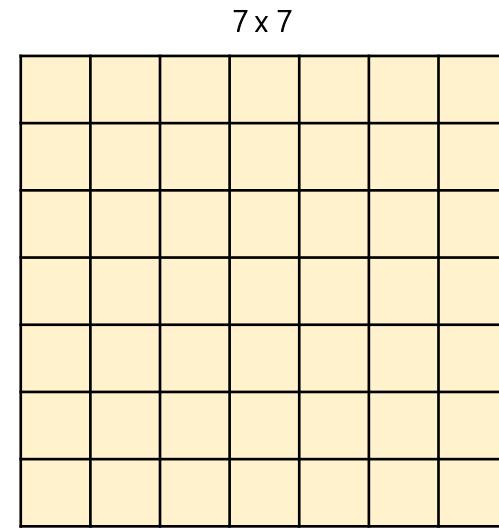
## Attention Filter



Normalize using SoftMax!

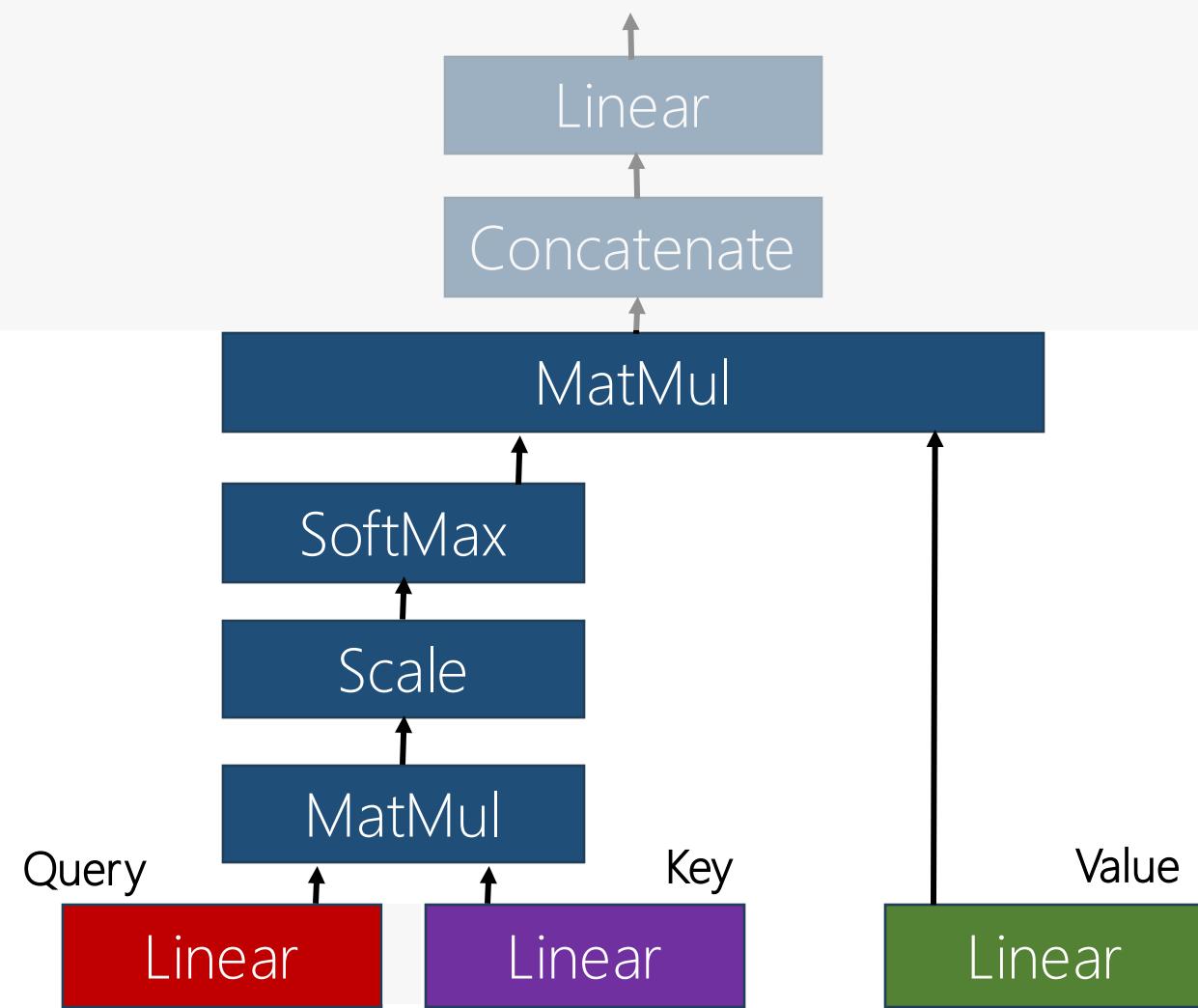


# Weighted Value

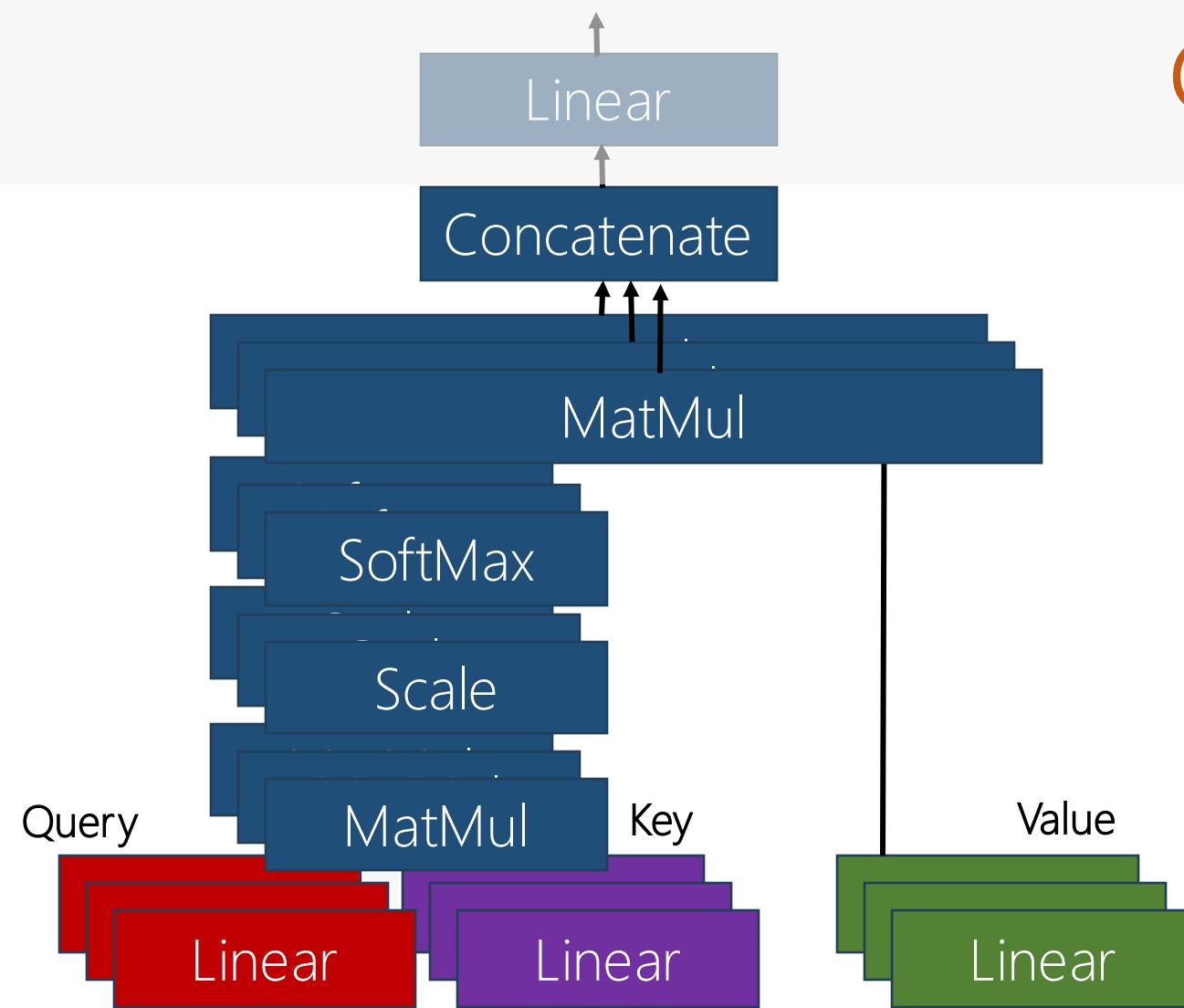


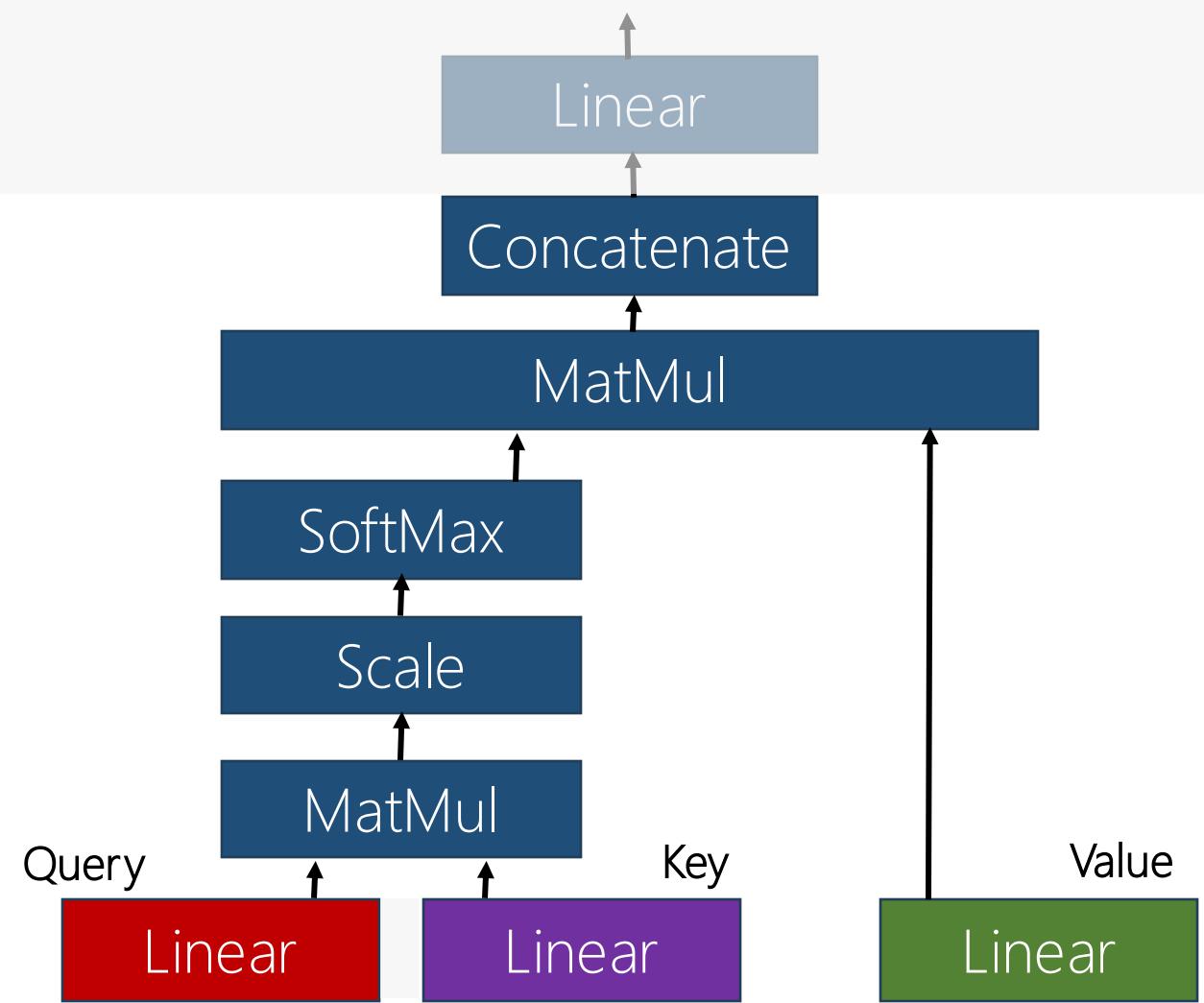
$$\text{Attention}(Q, K, V) = \sigma \left( \frac{Q \cdot K^T}{\sqrt{d_k}} \right) \cdot V$$

# What's to Concatenate??

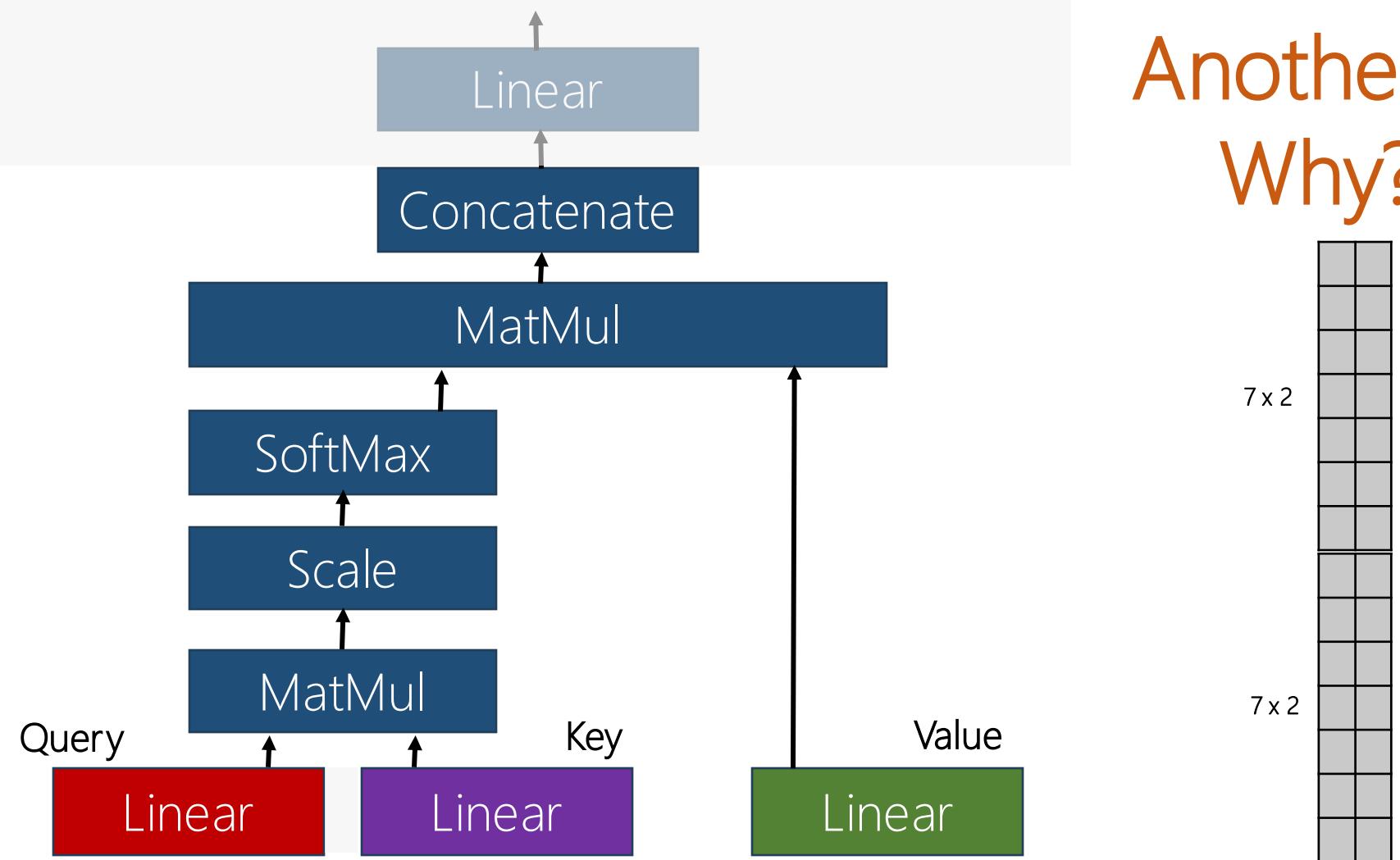


# Concatenate Outputs from all Attention Heads

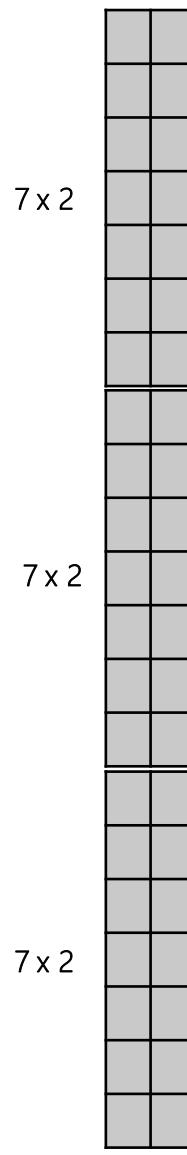




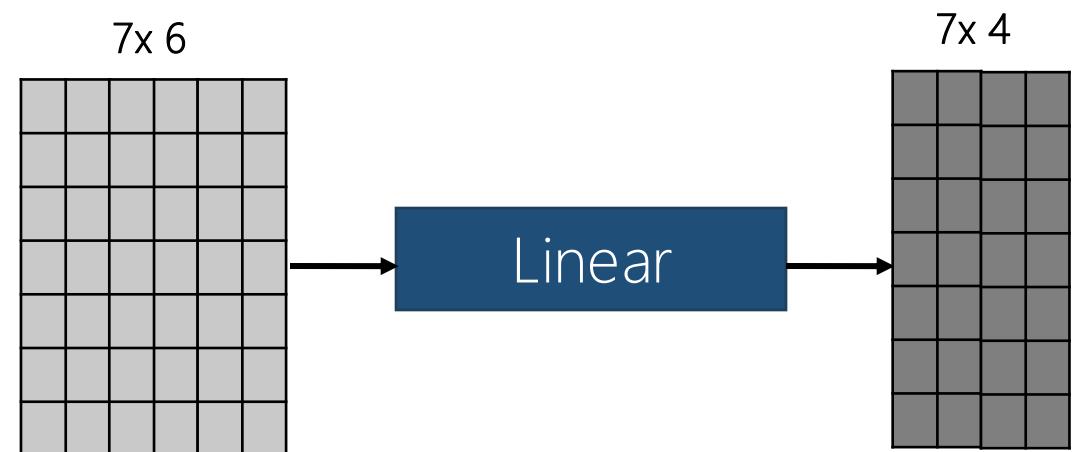
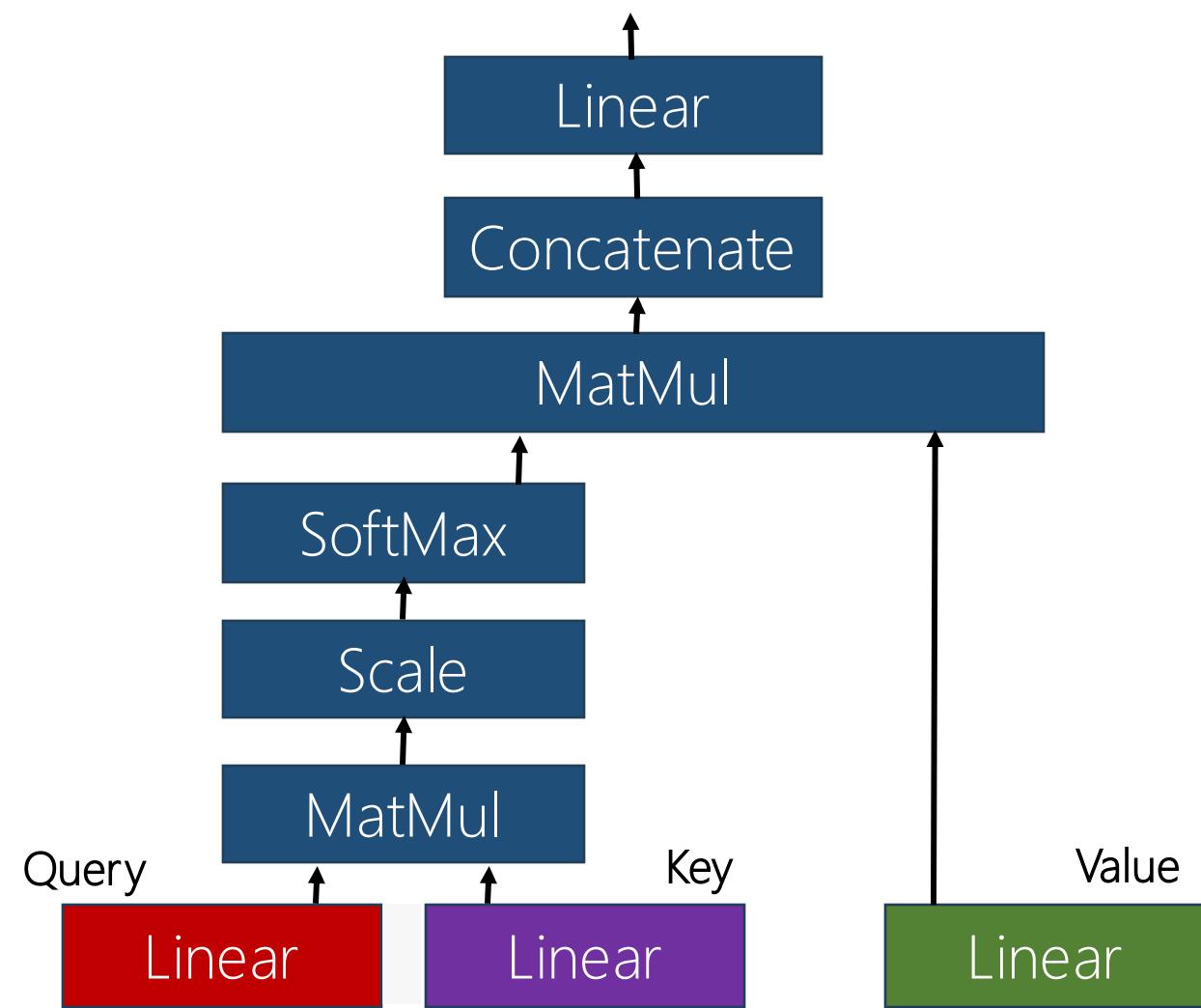
Another Linear layer??  
Why? Why? Why?



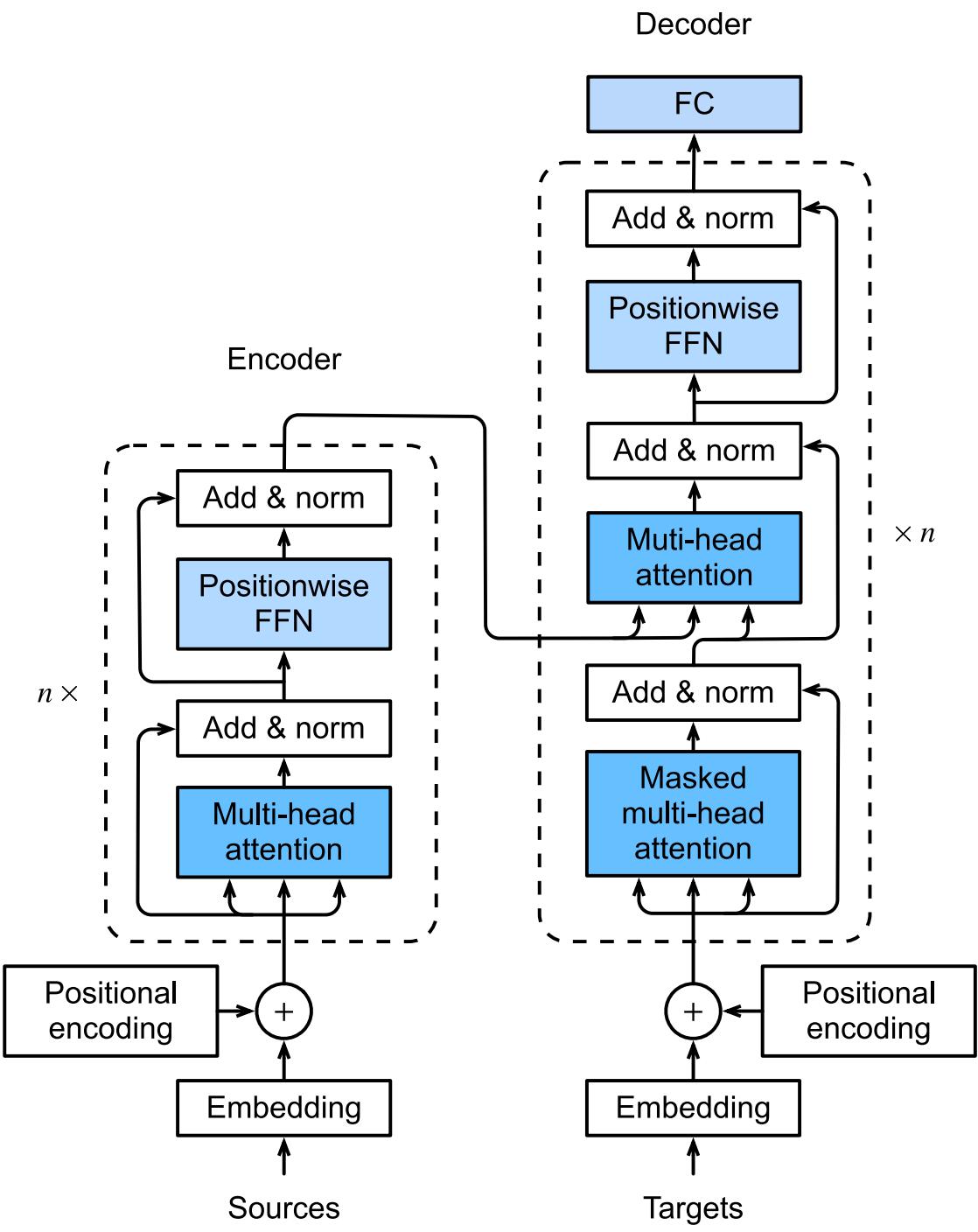
Another Linear layer??  
Why? Why? Why?

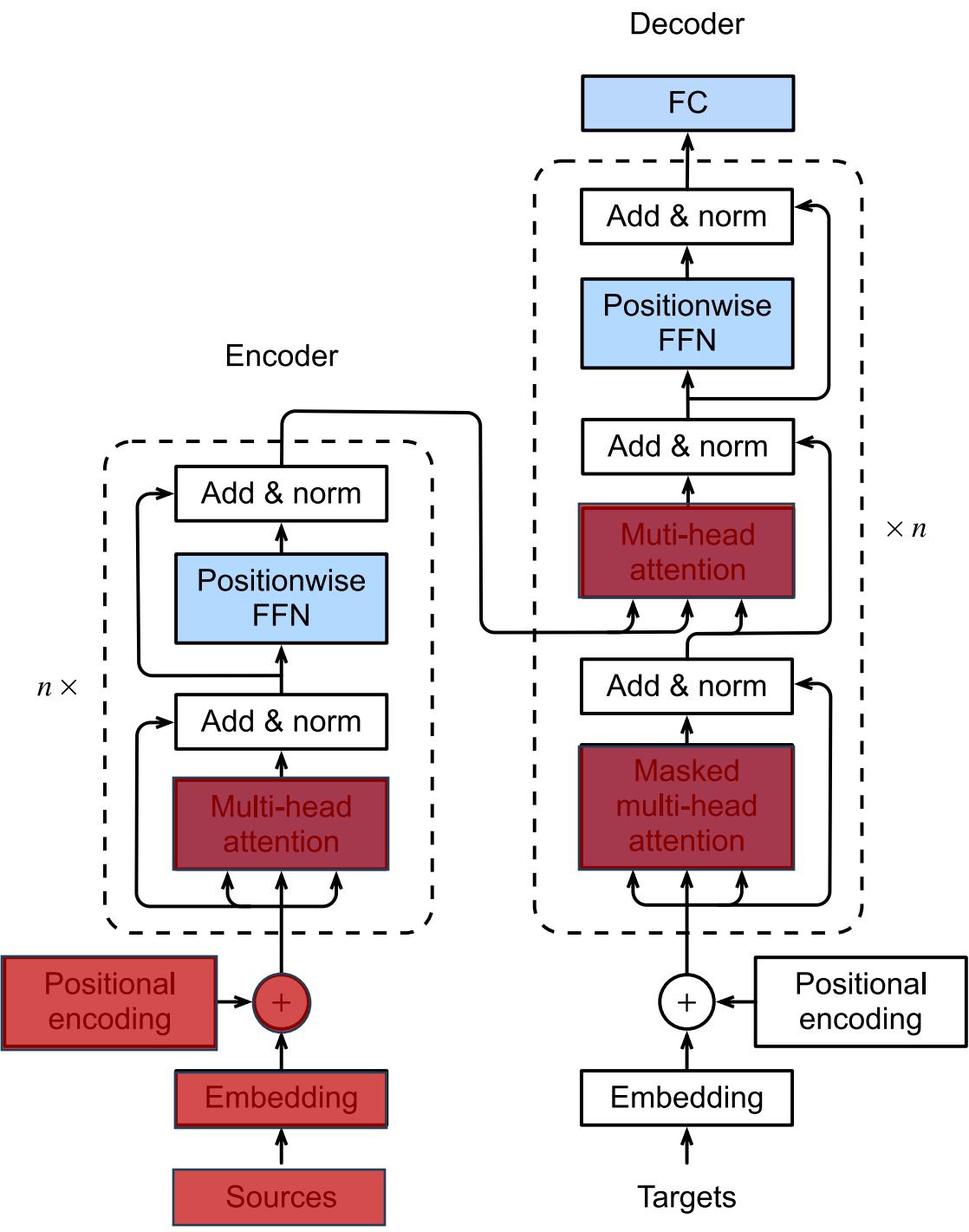


# Another Linear layer?? Why? Why? Why?

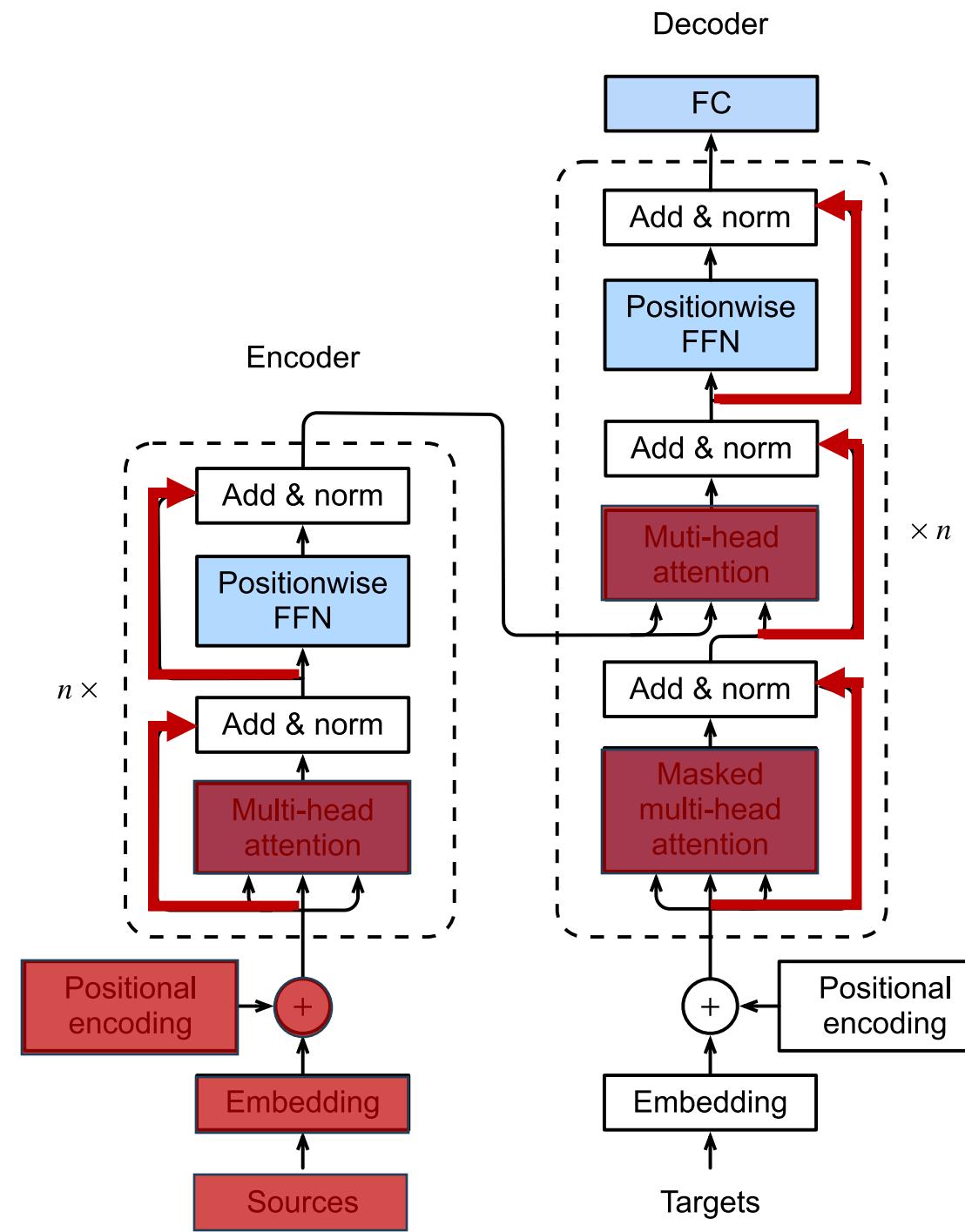


What about the rest of  
the Transformer?



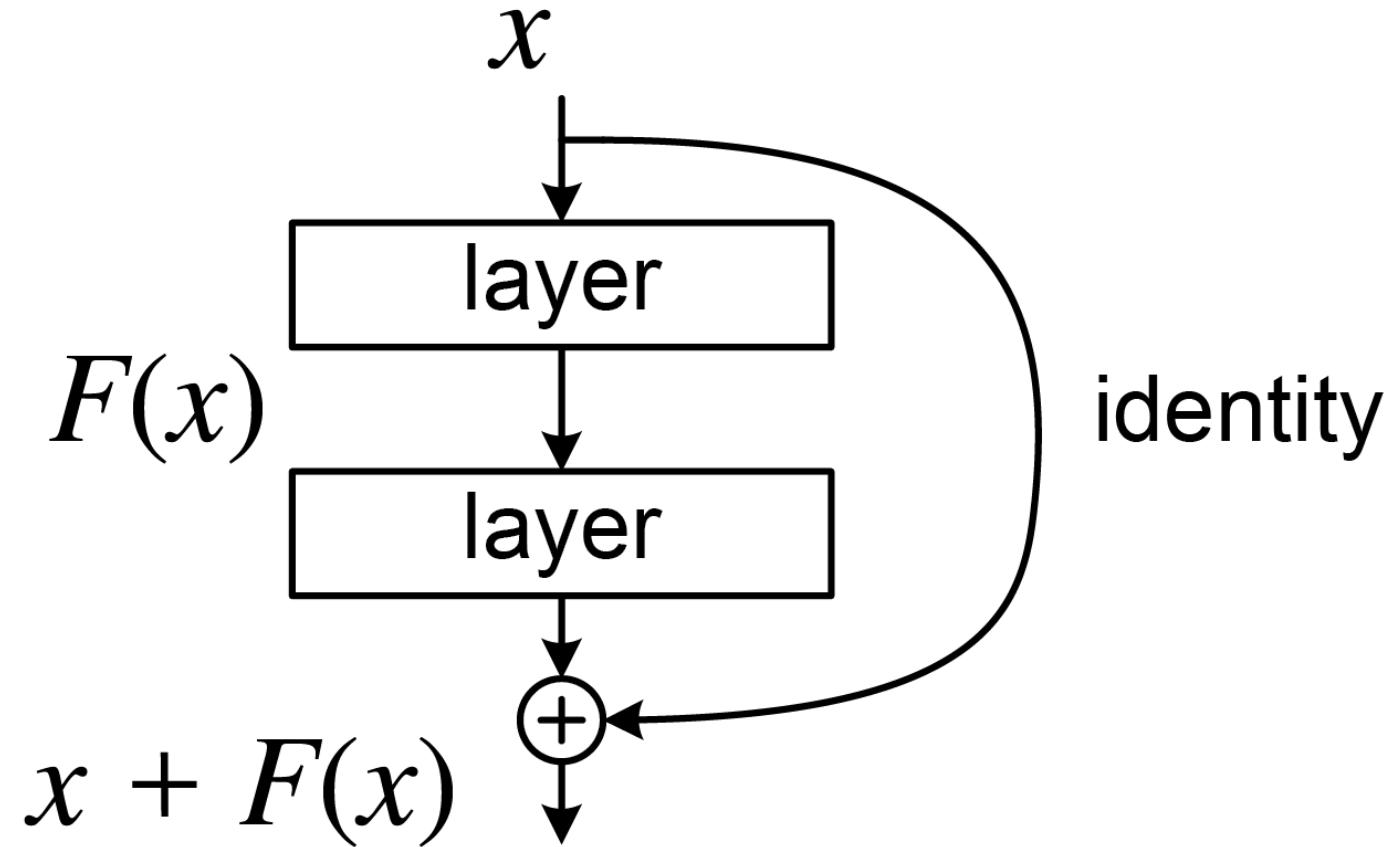


# Residual Connection



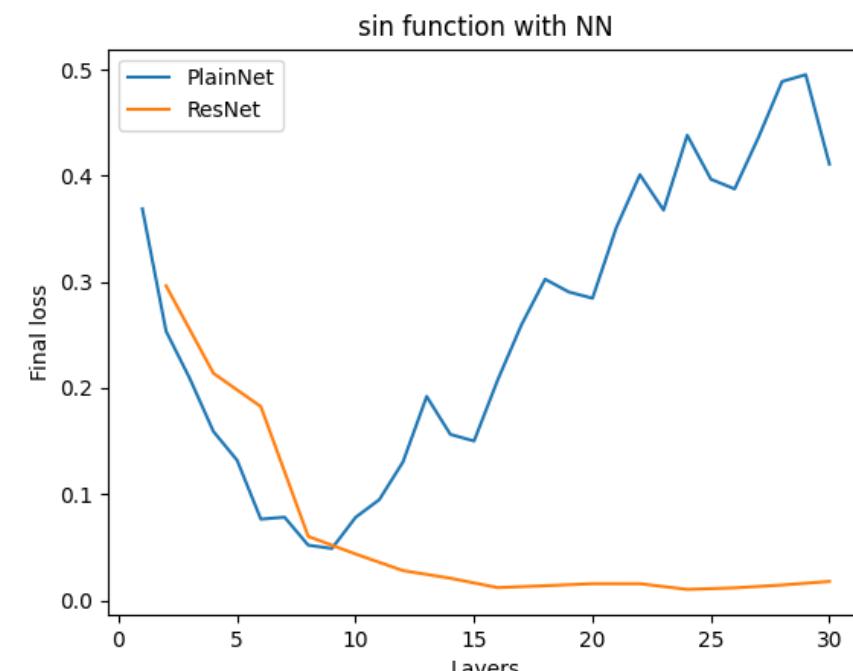
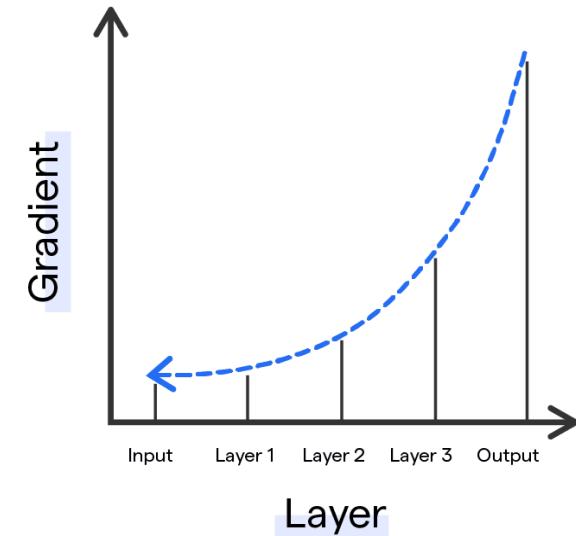
Let's go down the memory lane...

# Residual neural network

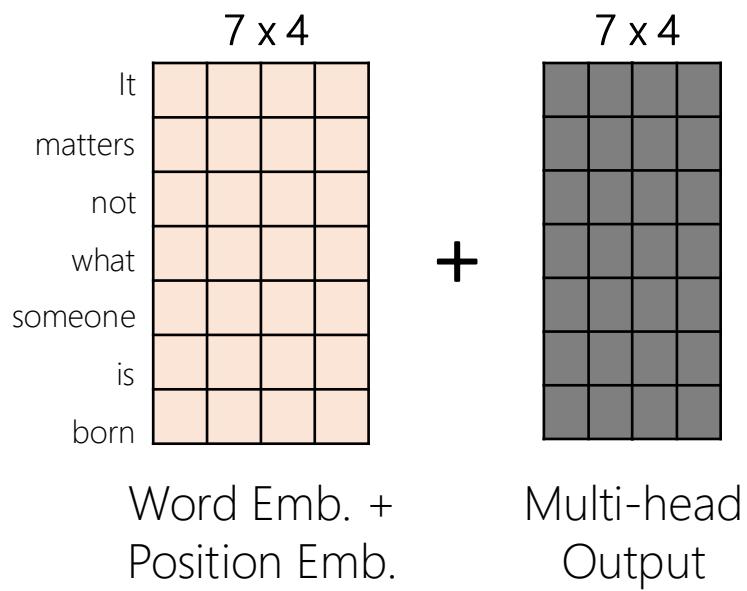
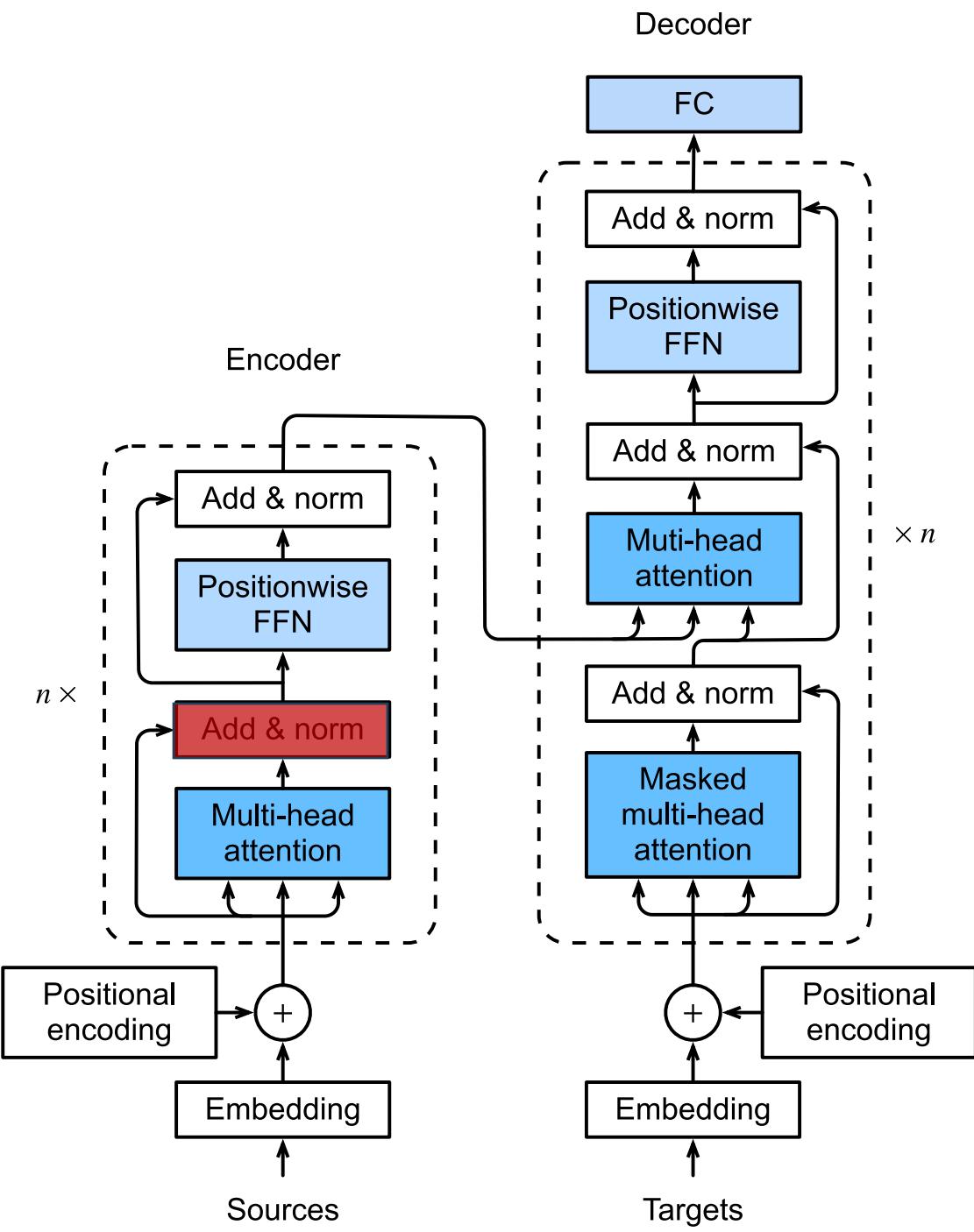


# Residual neural network

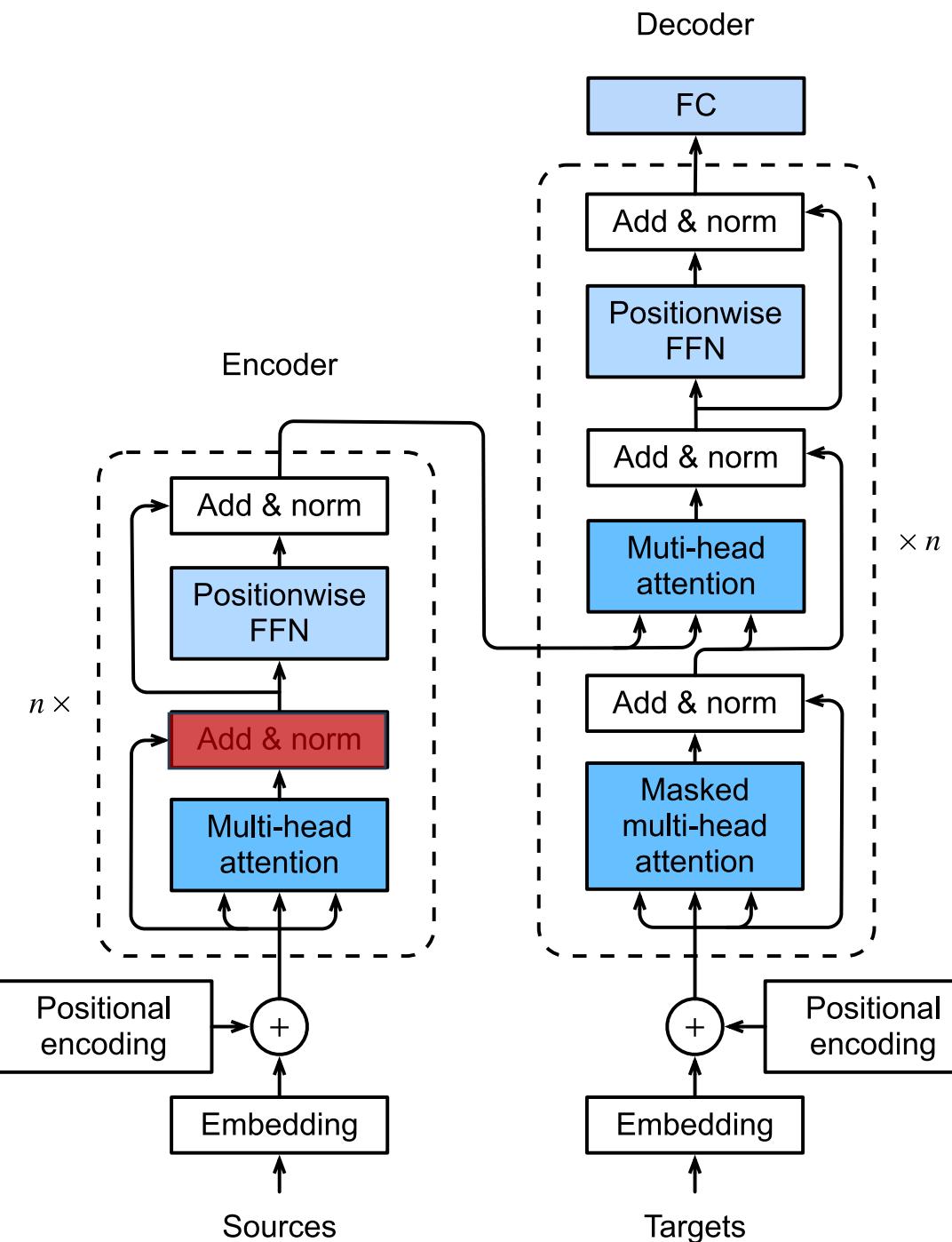
- **Vanishing Gradient:** As we backpropagate the gradients back through the network from the output layer towards the input the repeated multiplication of these small derivative values leads to increasingly smaller gradients
- **Degradation problem:** As a network deepens, the accuracy can start to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



# Add & Norm



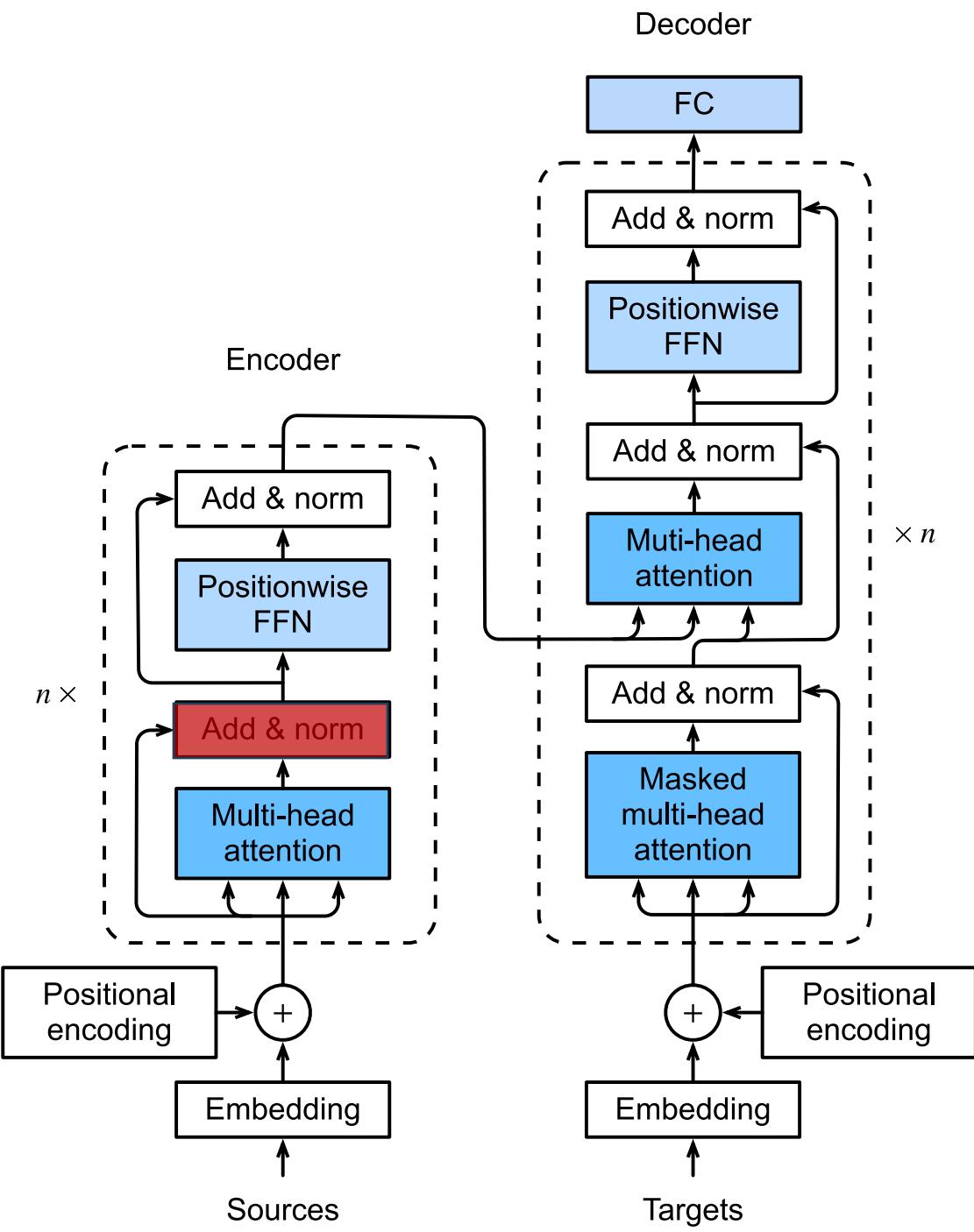
# Add & Norm



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.

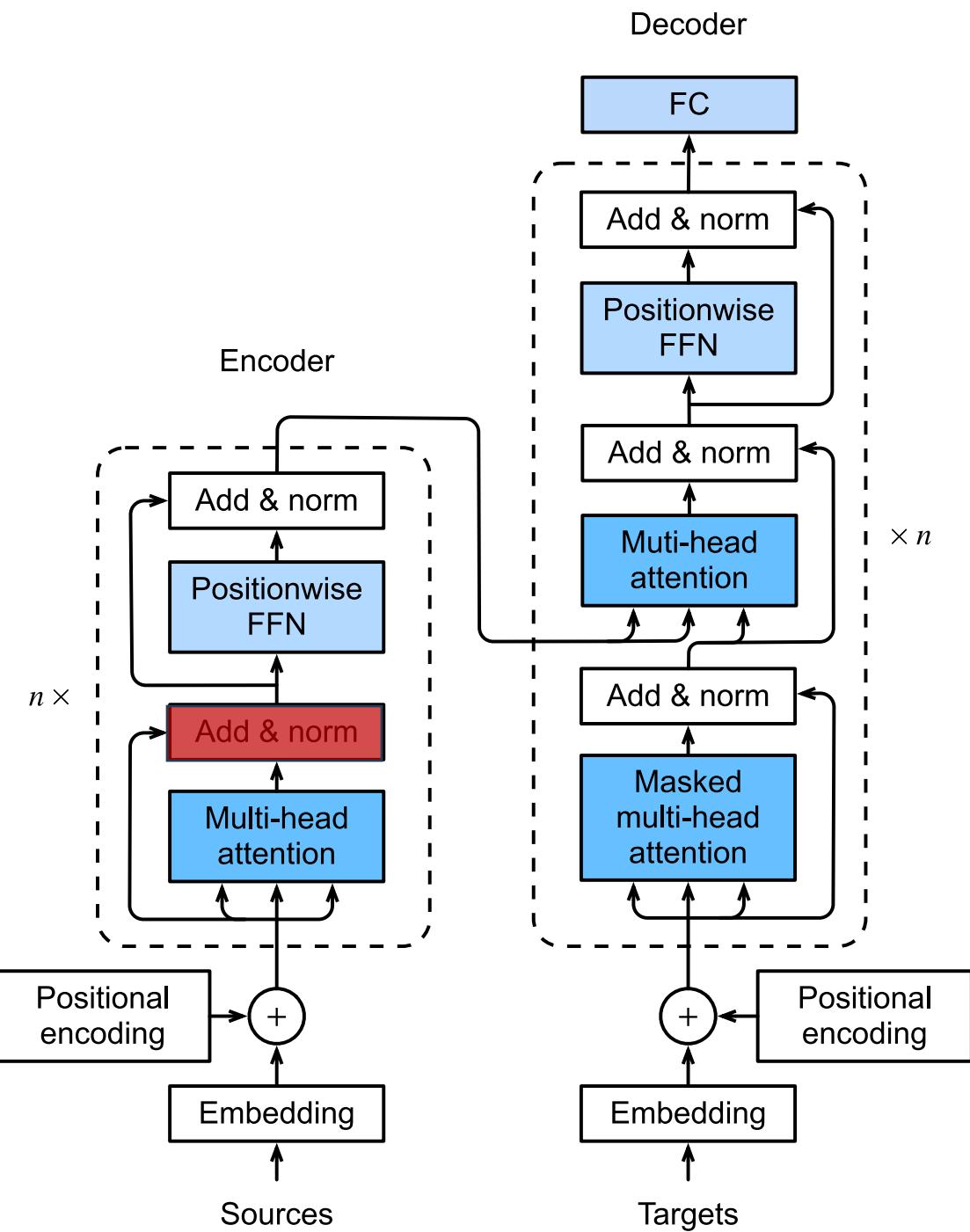
7 x 4			
1.56	2.12	0.91	2.87
0.45	1.23	2.76	0.67
2.03	0.58	1.41	1.29
0.92	2.31	0.14	2.55
1.80	0.61	2.98	1.52
2.67	0.33	1.99	1.74
0.48	2.40	1.68	0.29

# Add & Norm



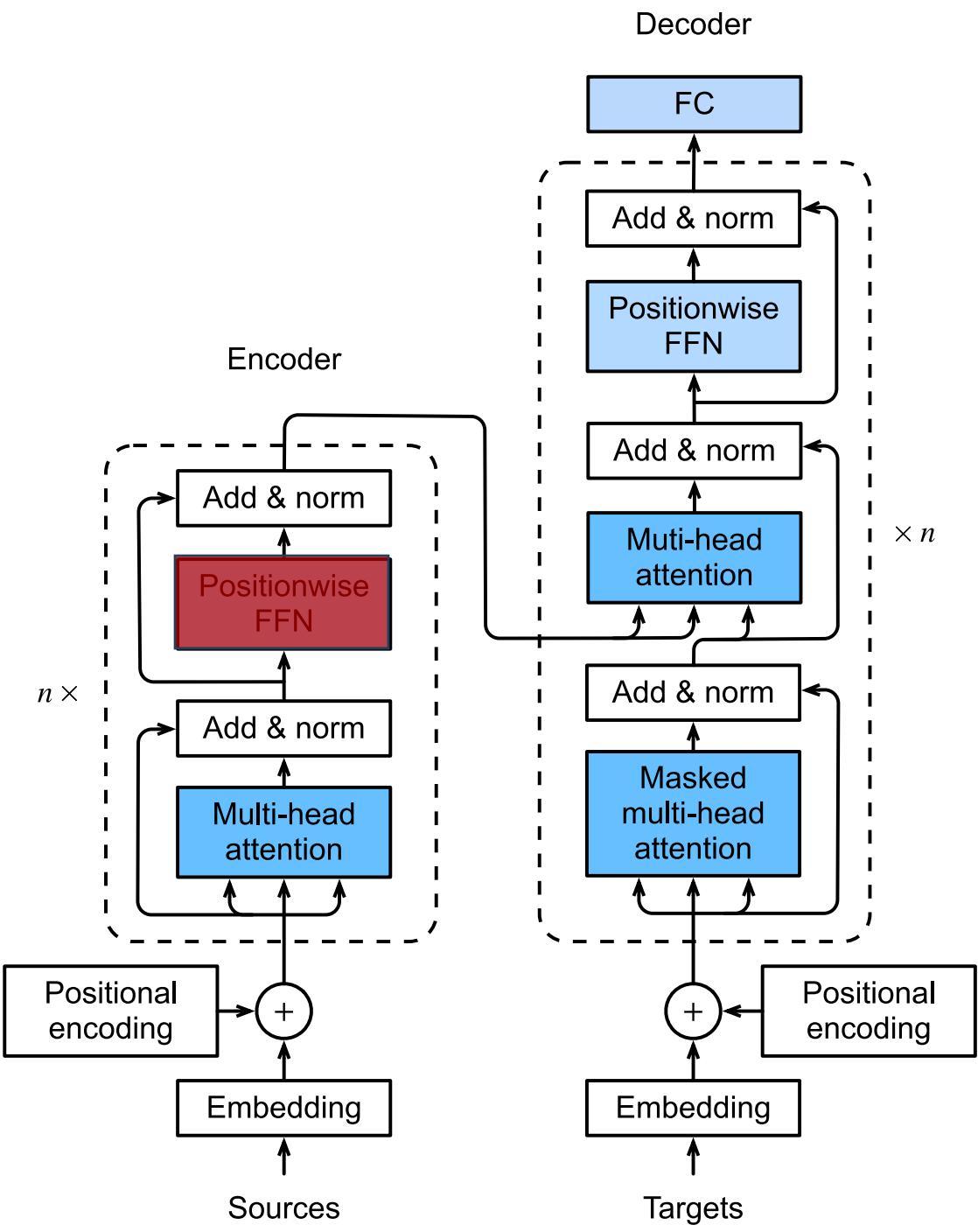
	7 x 4			
It	1.56	2.12	0.91	2.87
matters	0.45	1.23	2.76	0.67
not	2.03	0.58	1.41	1.29
what	0.92	2.31	0.14	2.55
someone	1.80	0.61	2.98	1.52
is	2.67	0.33	1.99	1.74
born	0.48	2.40	1.68	0.29

# Add & Norm

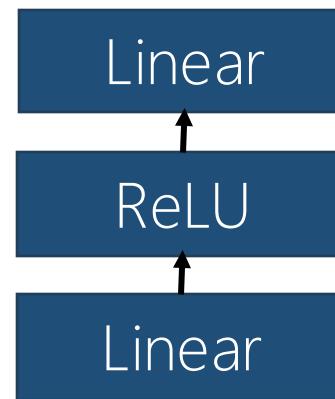


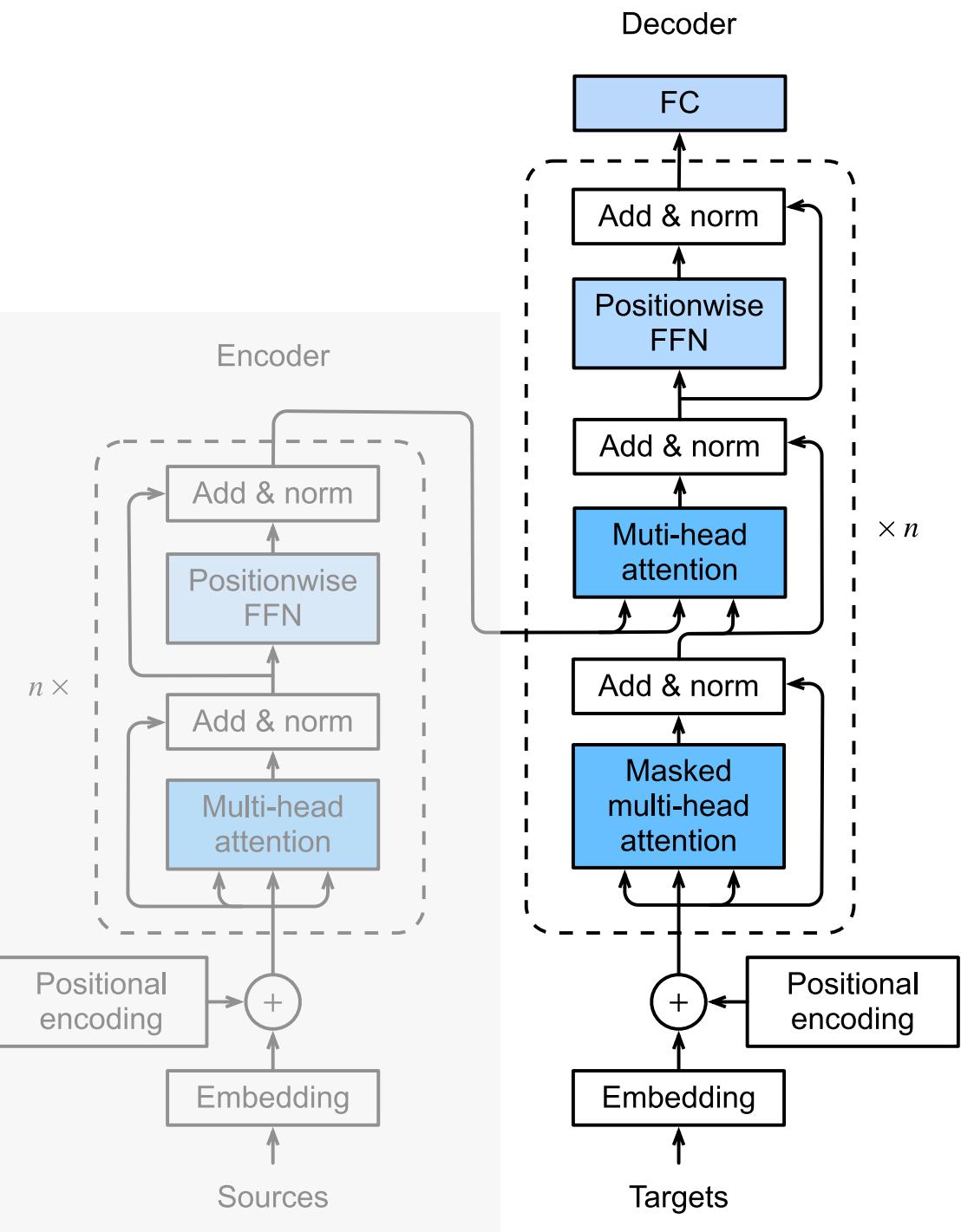
	Mean ( $\mu$ )	Stdev. ( $\sigma$ )
It	1.56	2.12
matters	0.45	2.76
not	2.03	0.58
what	0.92	1.41
someone	1.80	0.14
is	2.67	2.55
born	0.48	1.52
	1.978	0.811
	1.930	0.656
	1.198	0.822

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



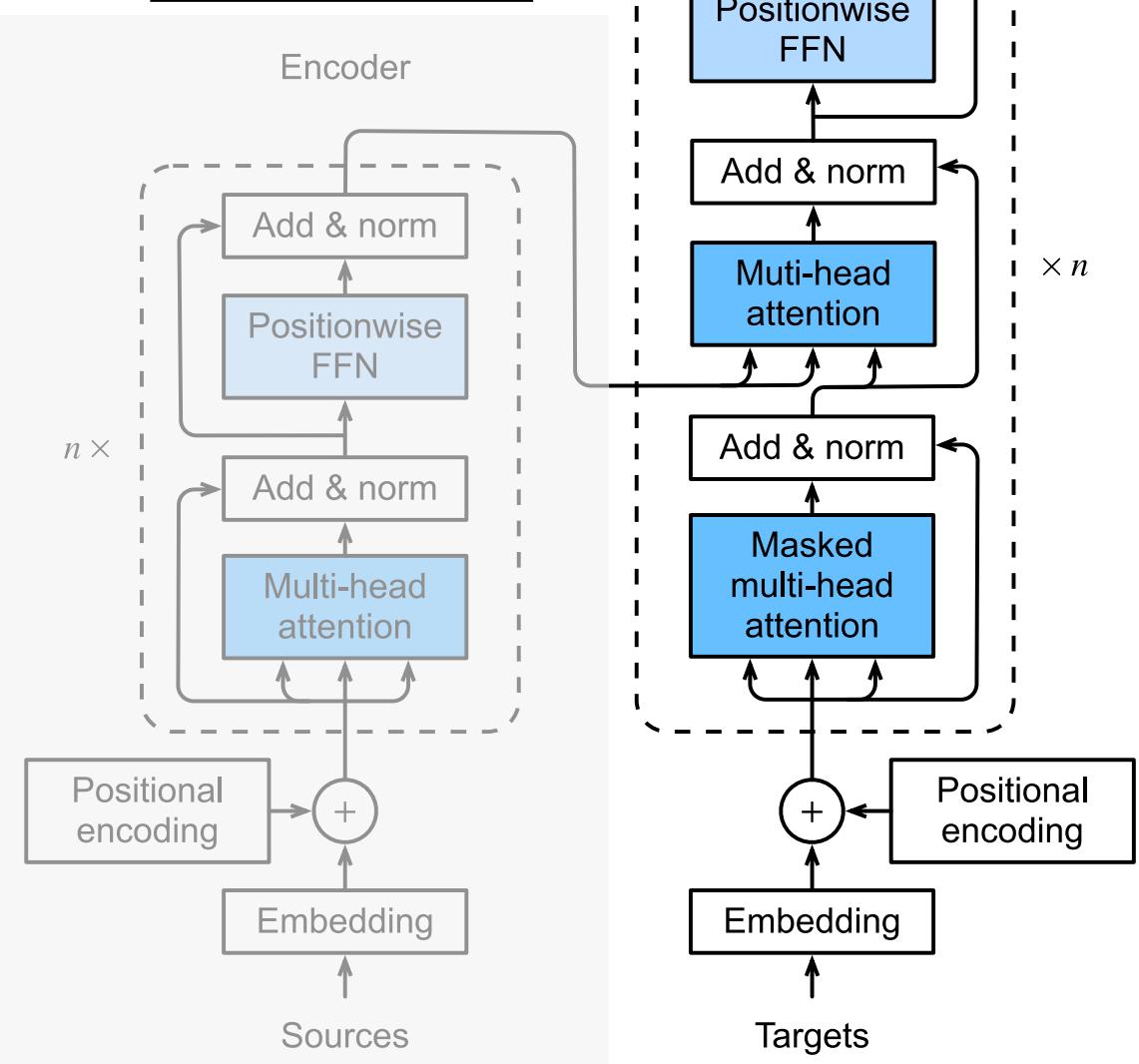
A bunch of Linear is all we need!

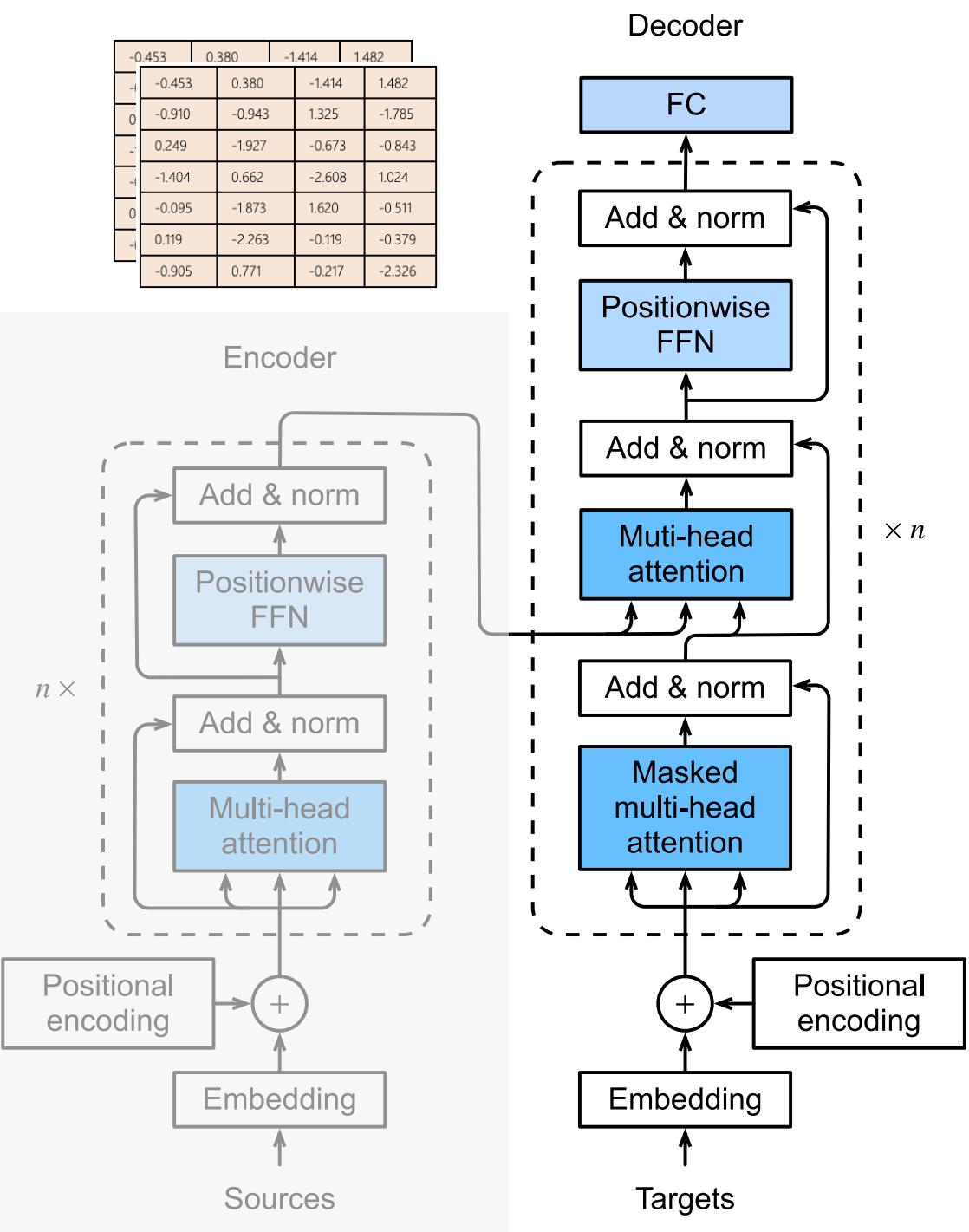




That covers our Encoder part of the Transformer!

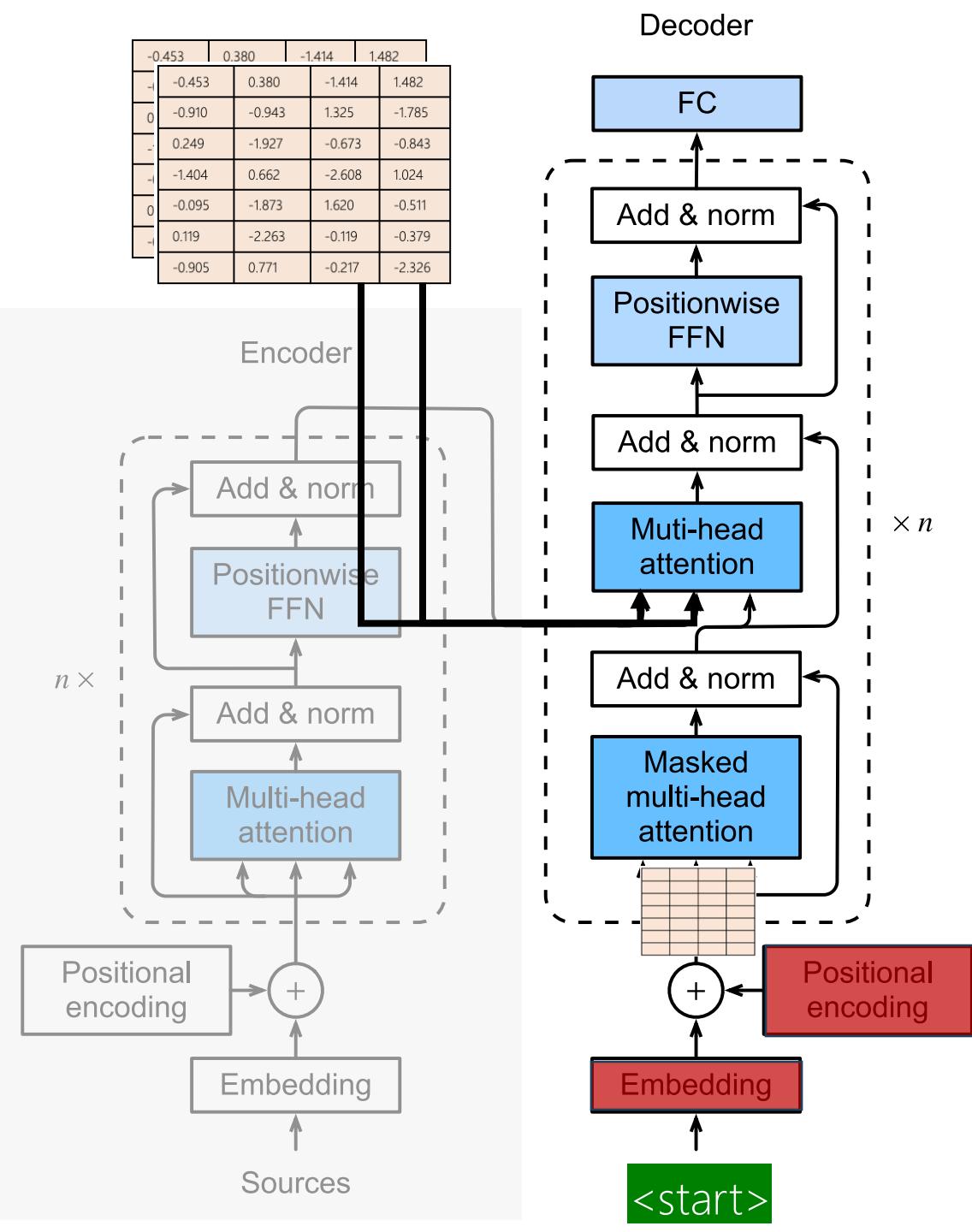
It	-0.453	0.380	-1.414	1.482
matters	-0.910	-0.943	1.325	-1.785
not	0.249	-1.927	-0.673	-0.843
what	-1.404	0.662	-2.608	1.024
someone	-0.095	-1.873	1.620	-0.511
is	0.119	-2.263	-0.119	-0.379
born	-0.905	0.771	-0.217	-2.326



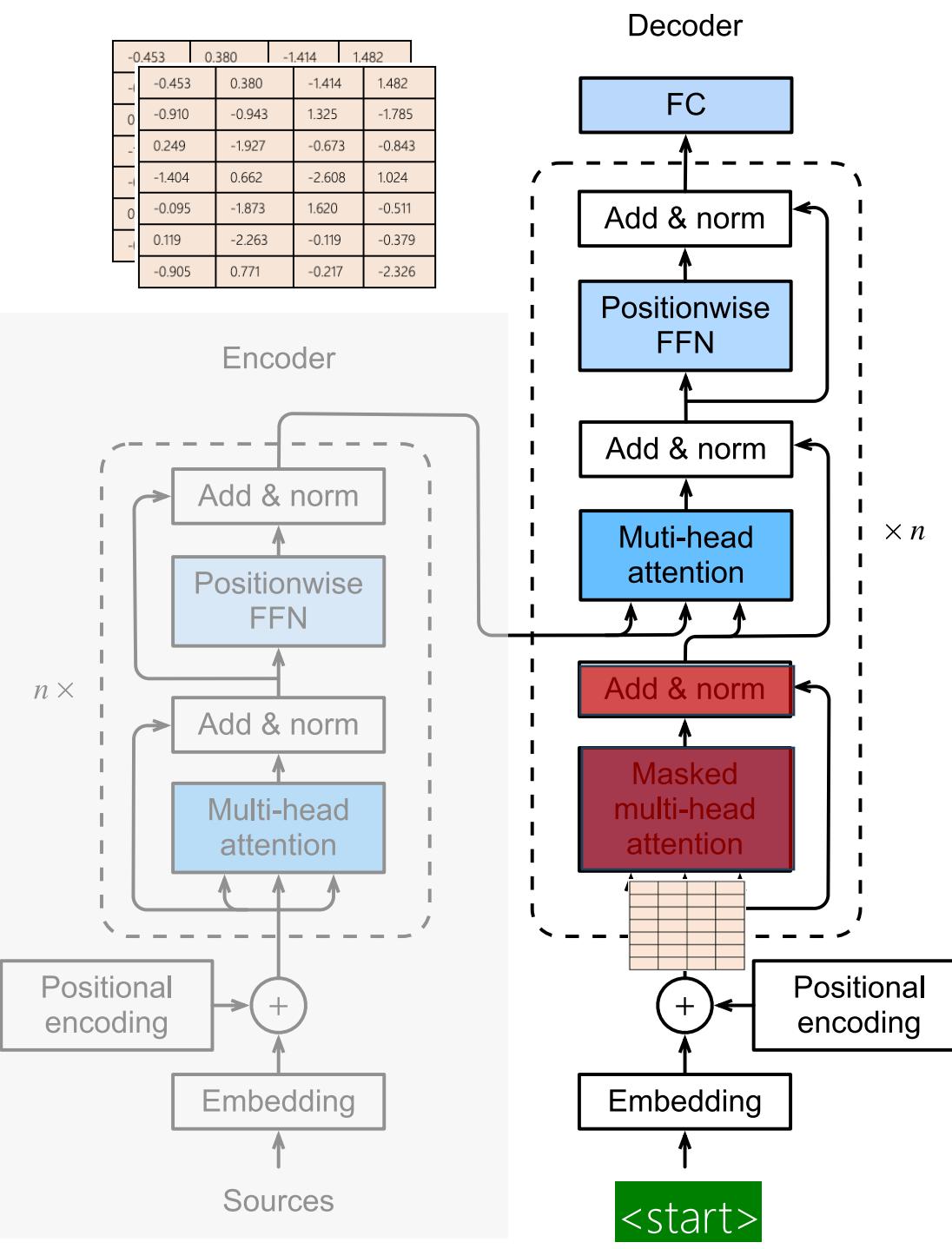


<start>

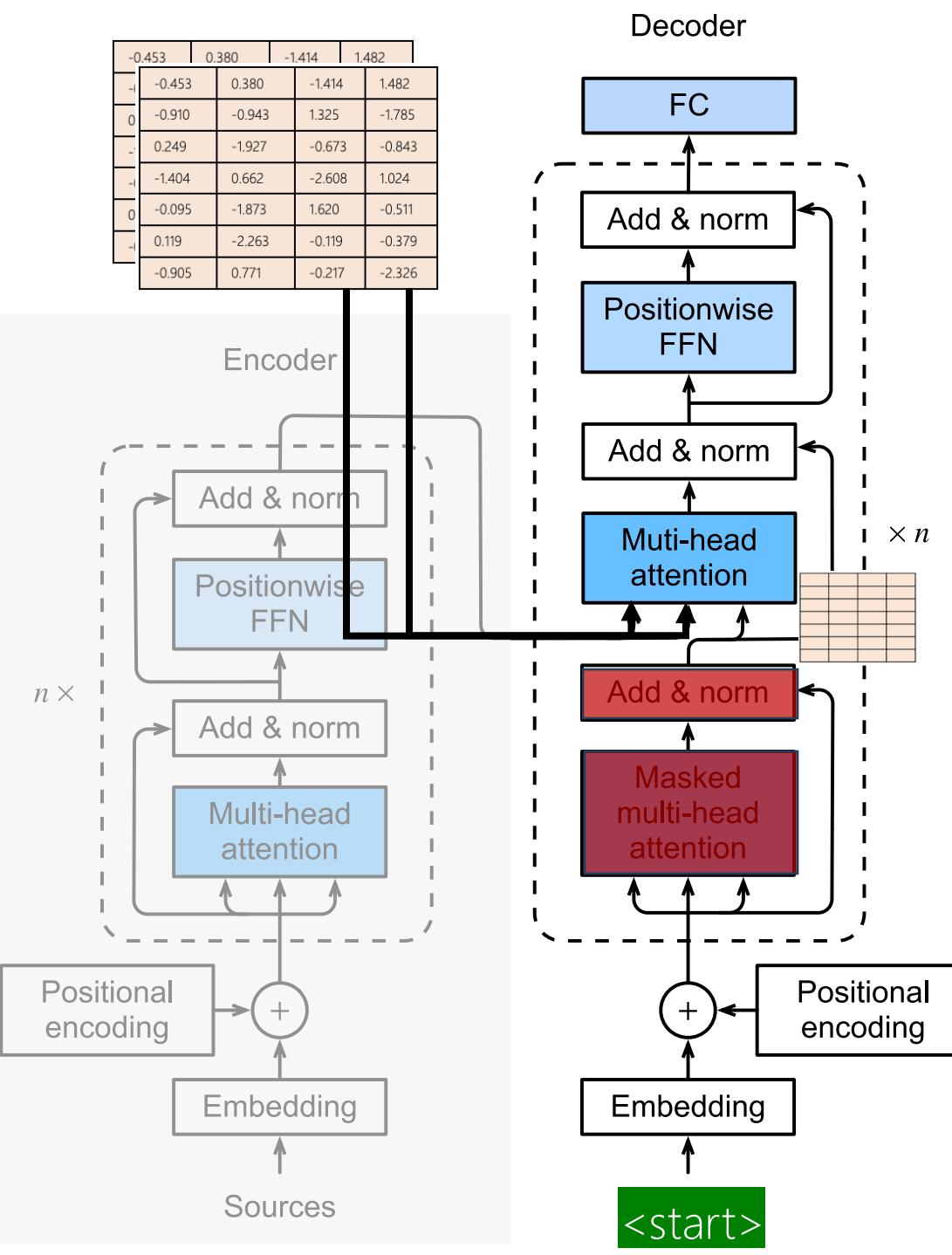
# Timestamp = 1

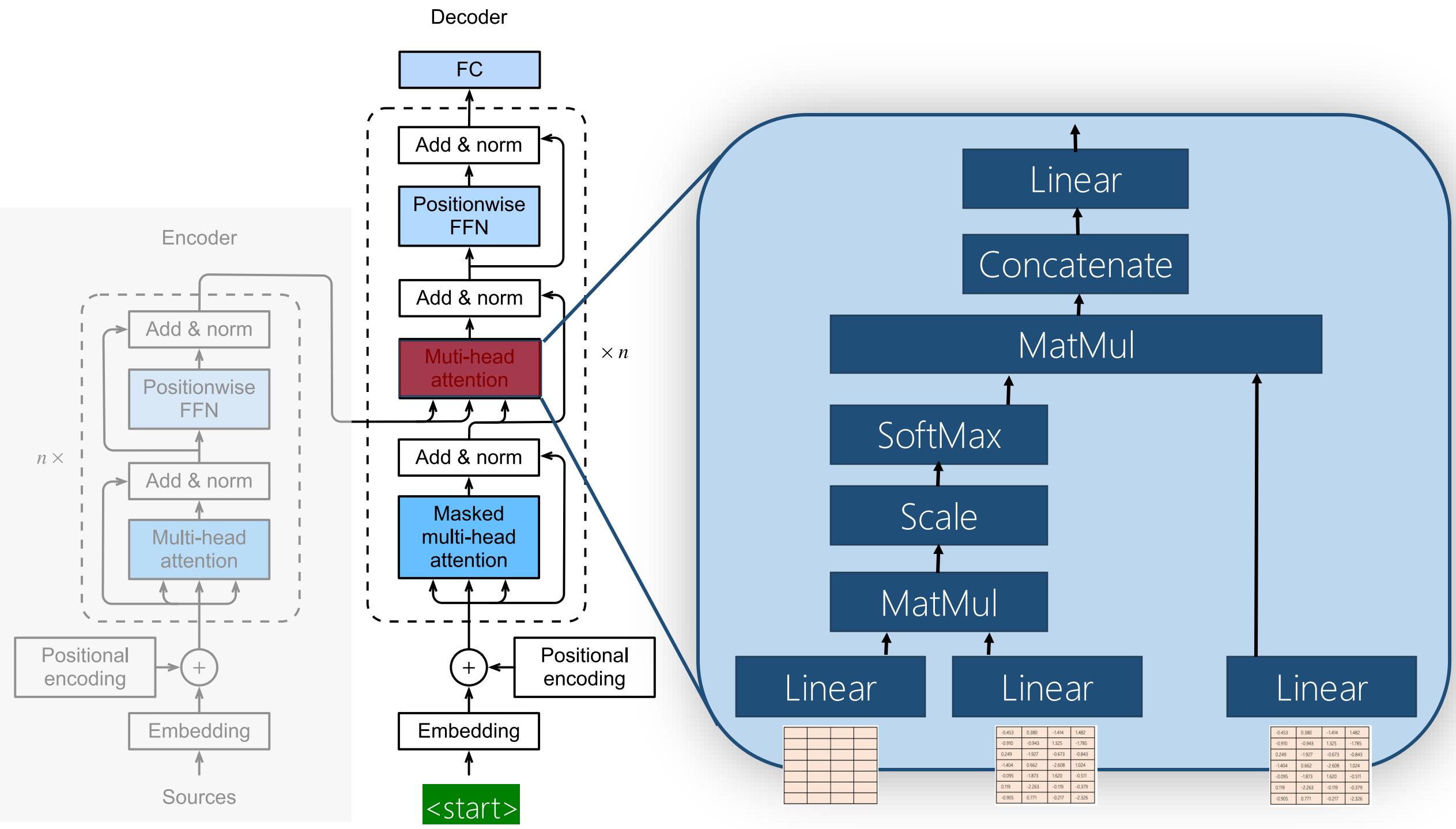


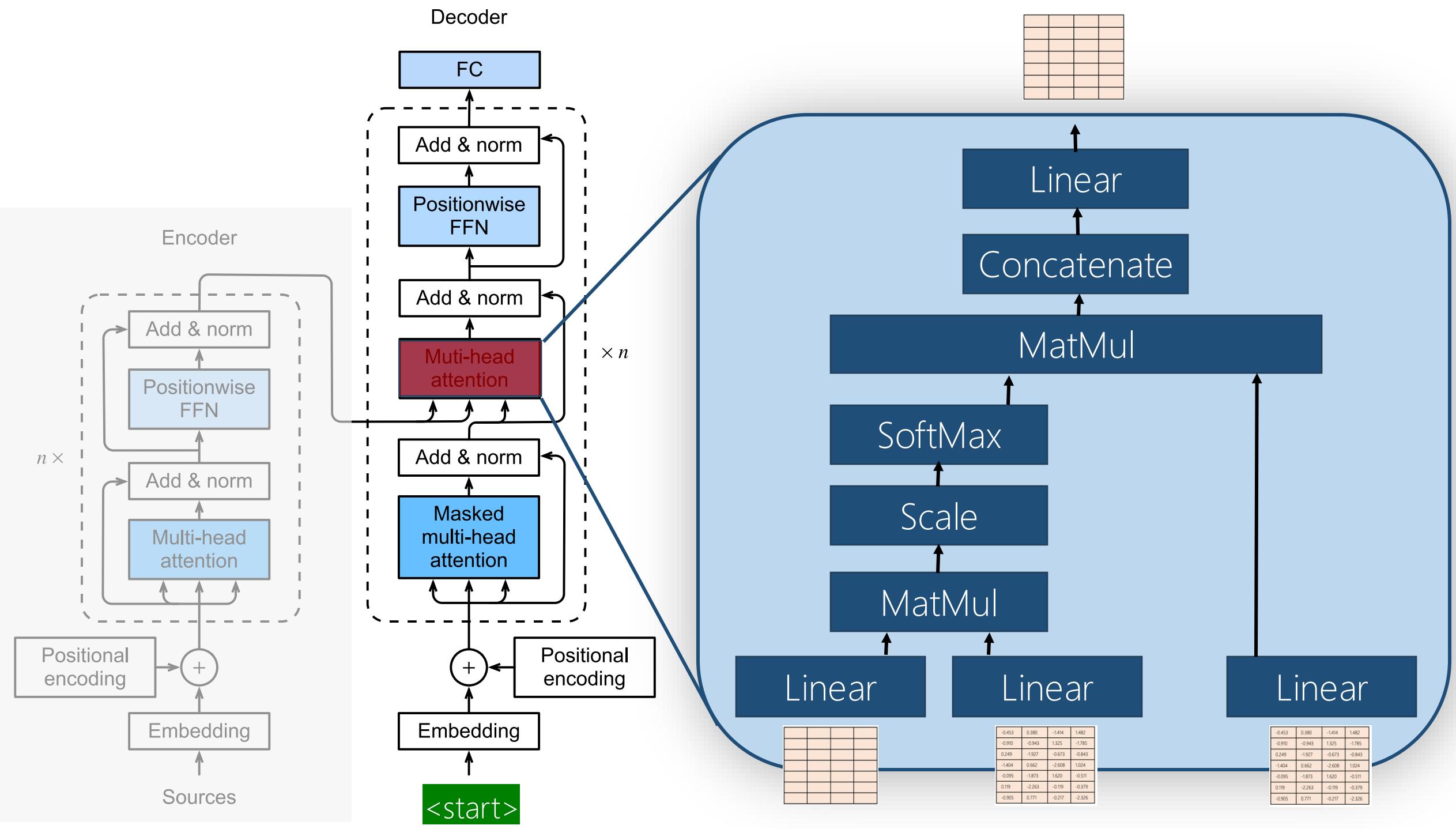
# Timestamp = 1

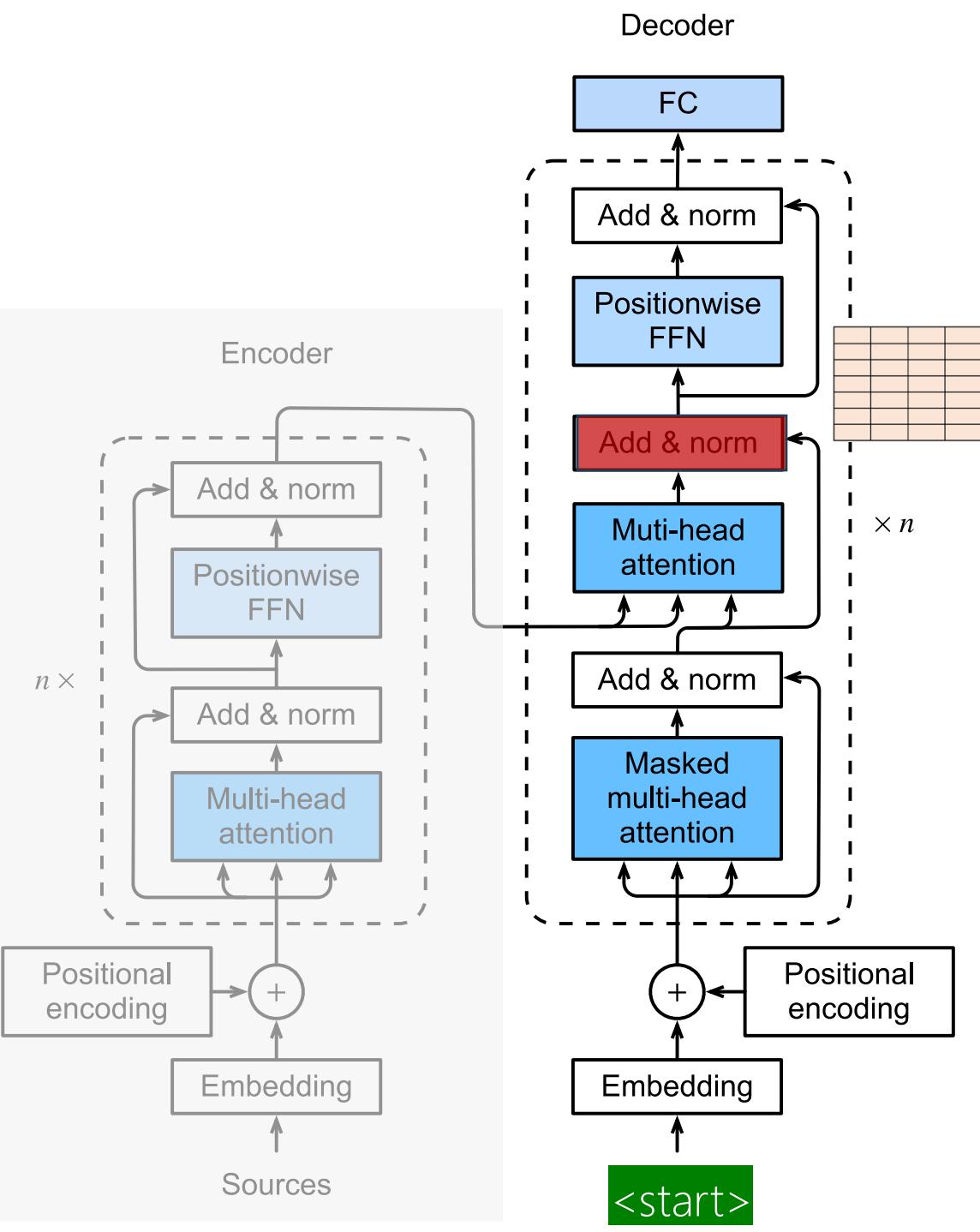


# Timestamp = 1

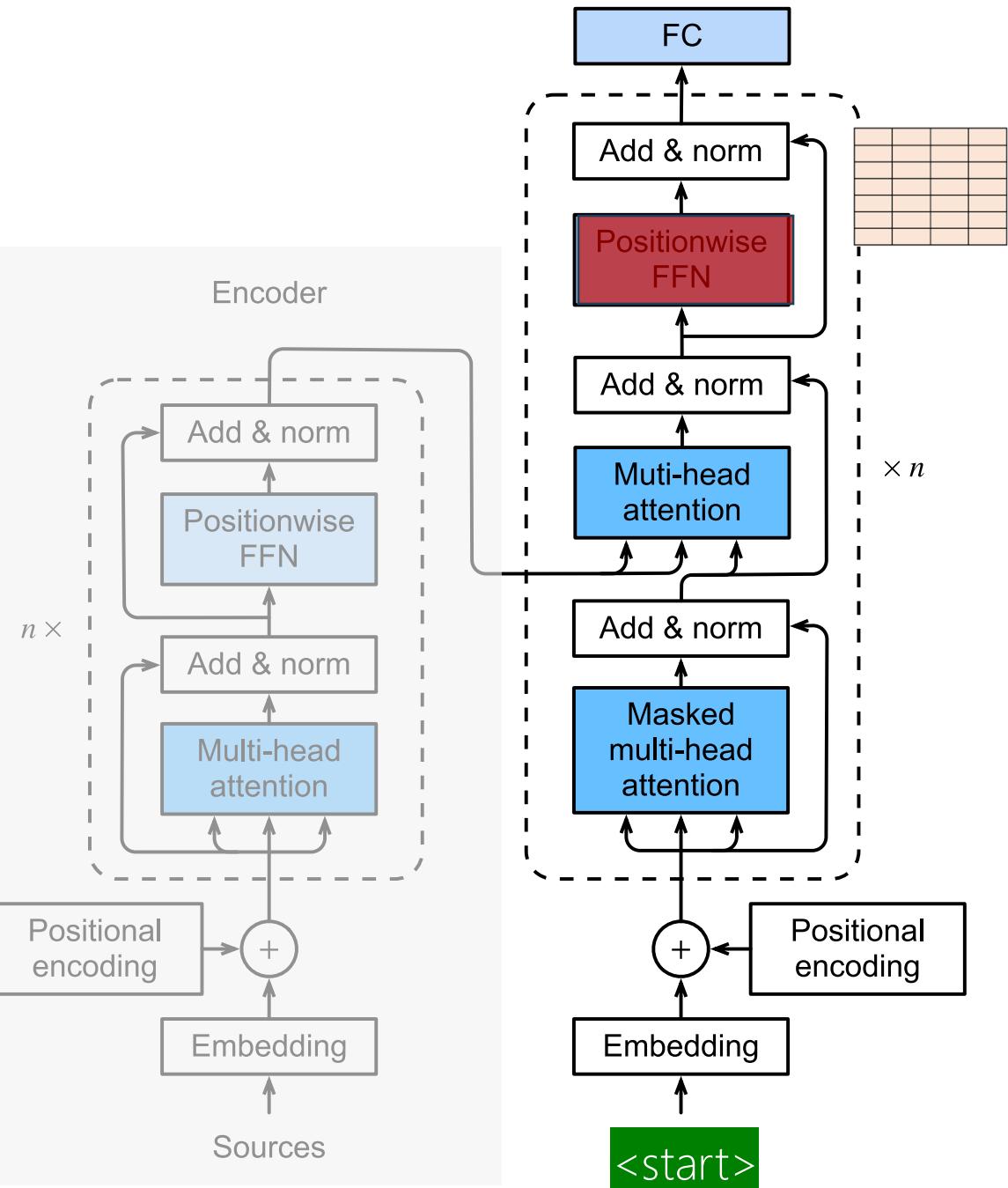


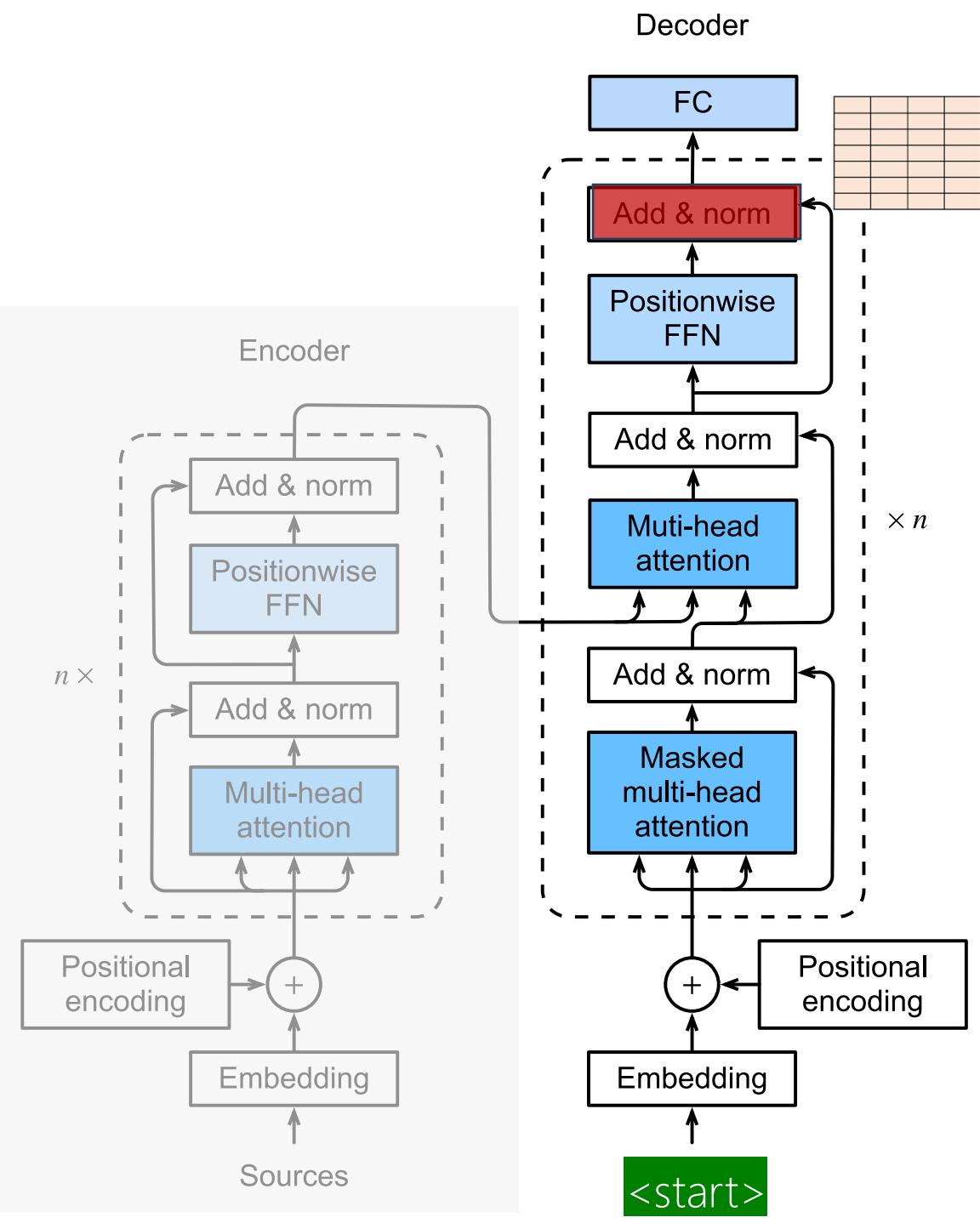






## Decoder

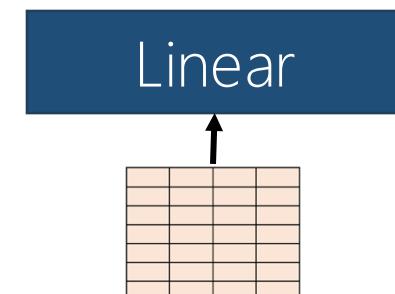
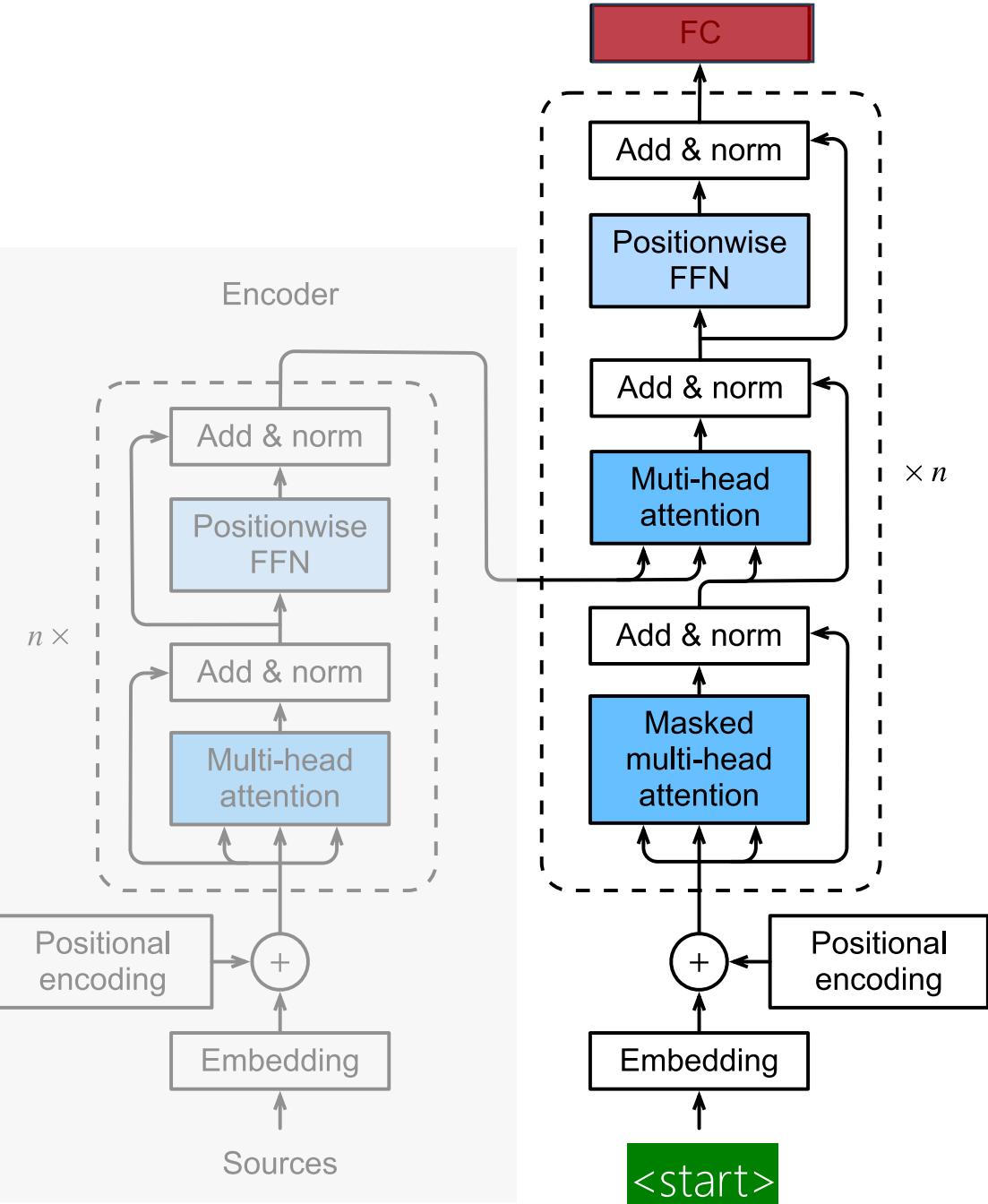


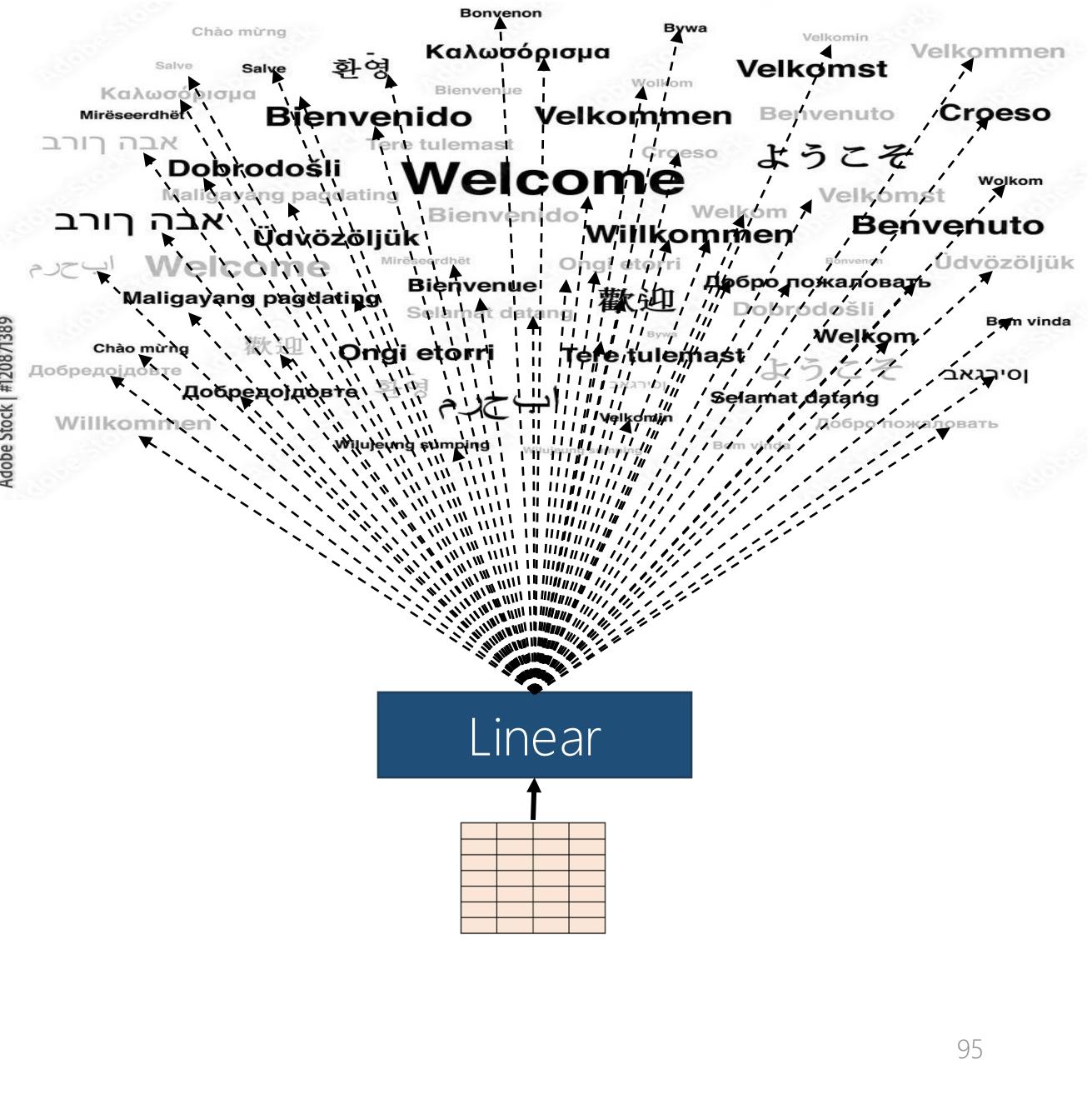
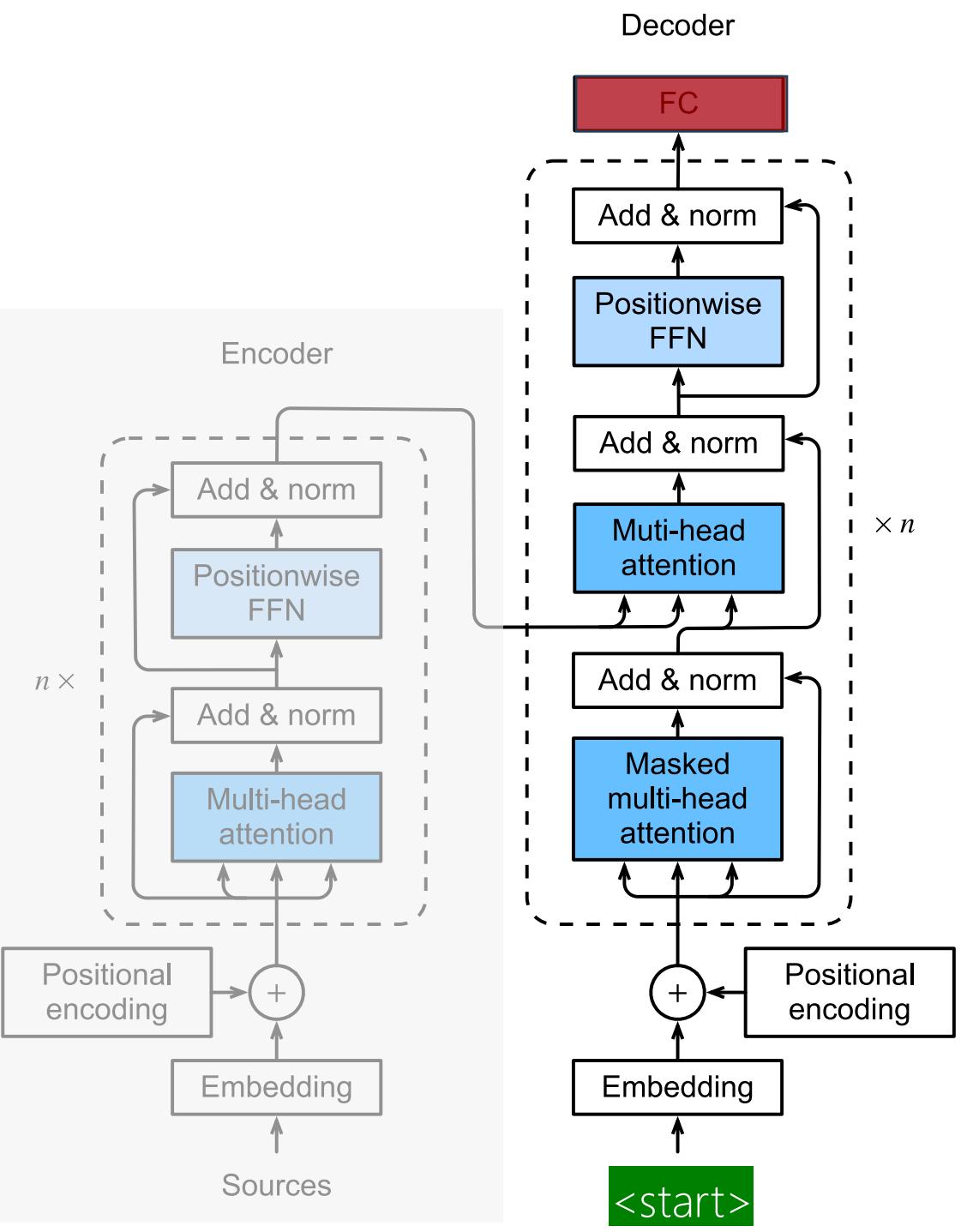


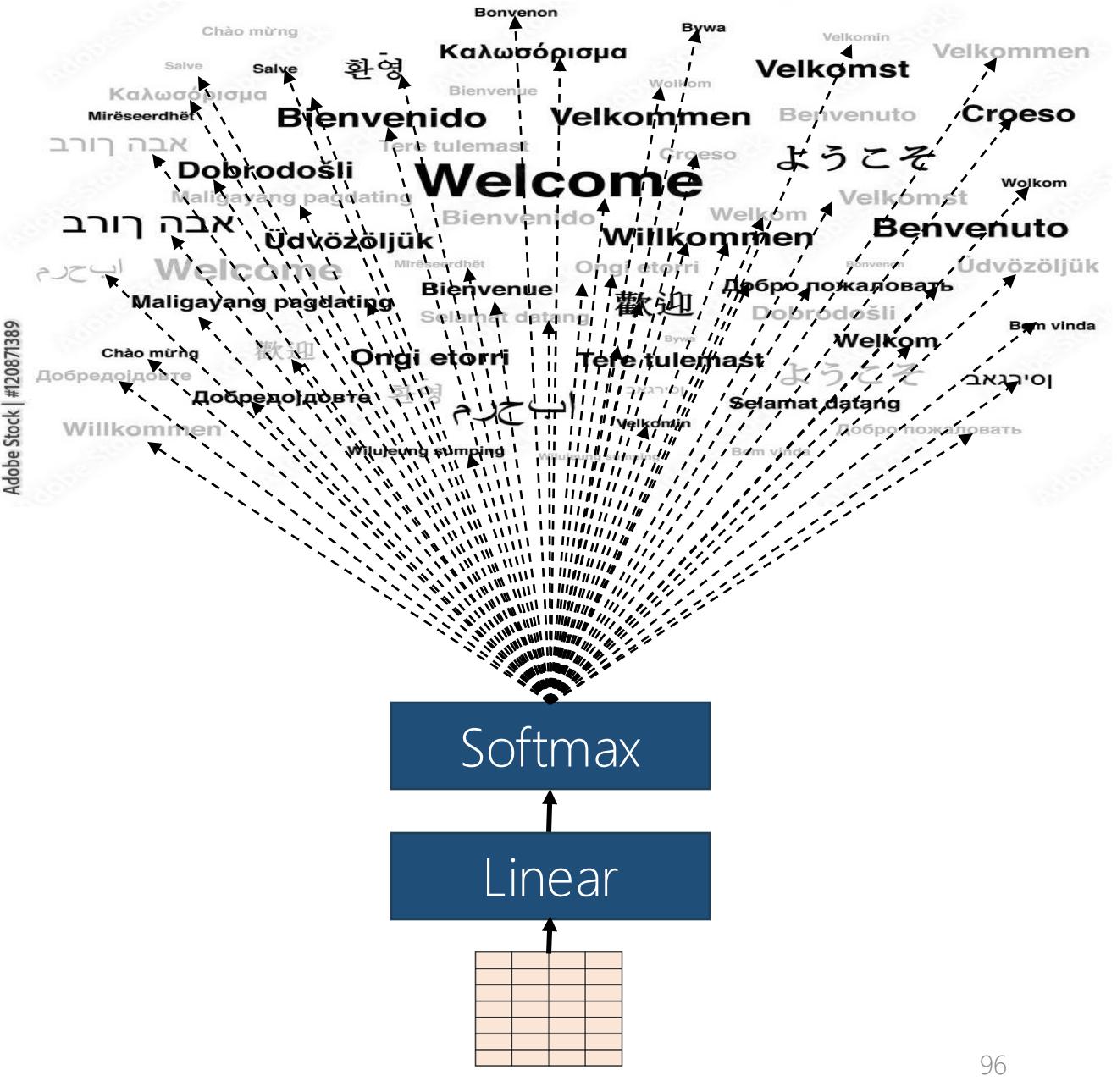
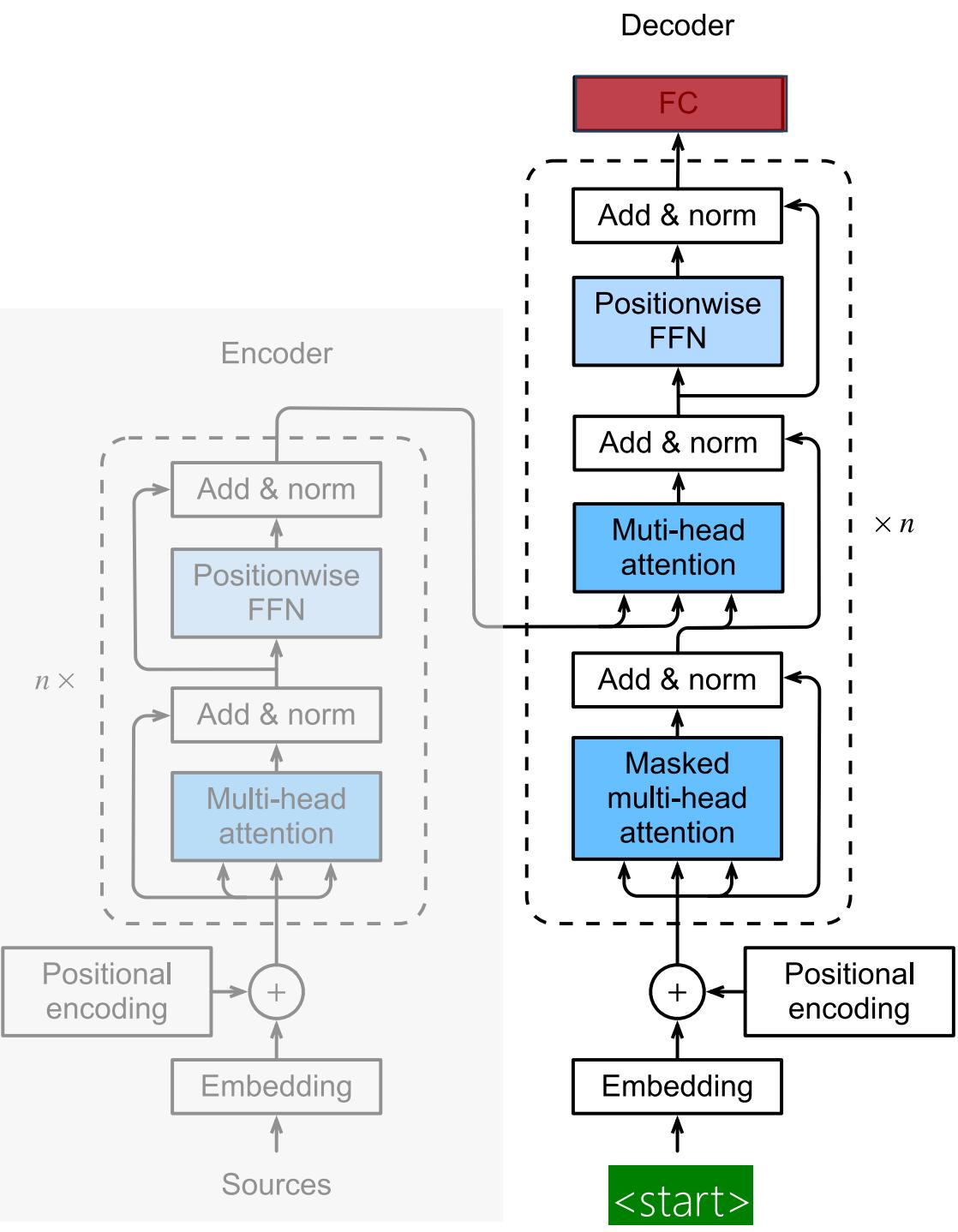
# What to do now?

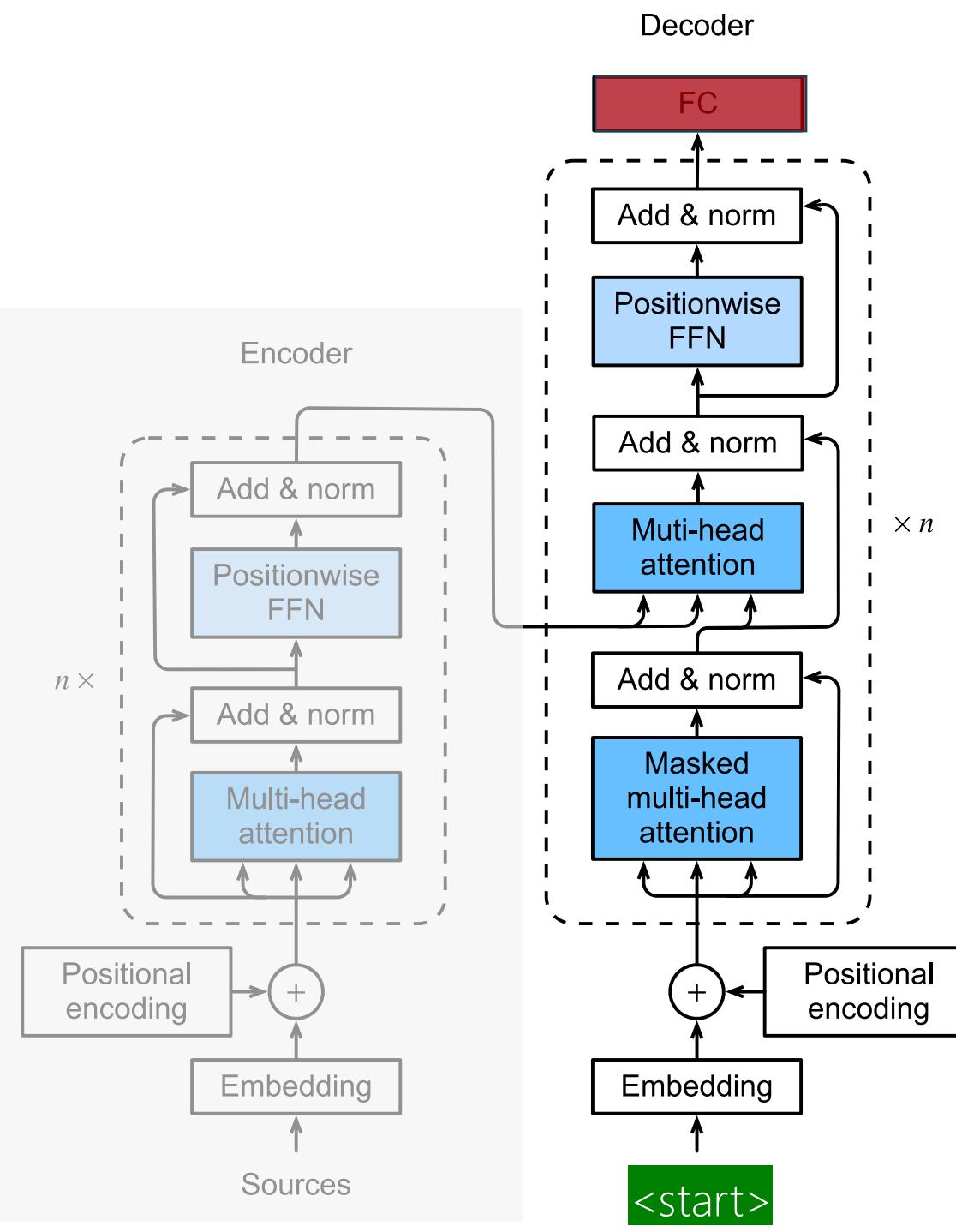


## Decoder

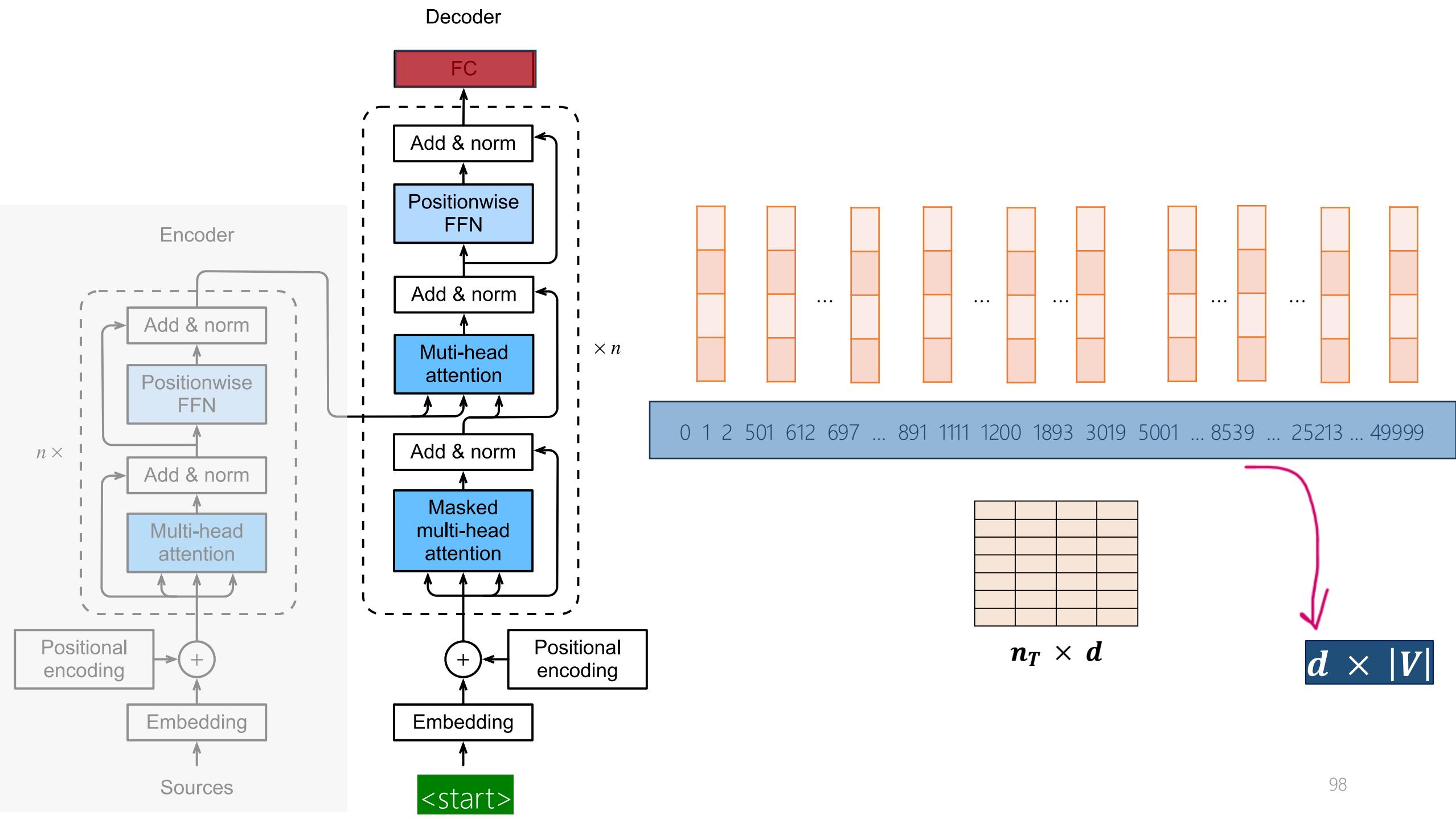


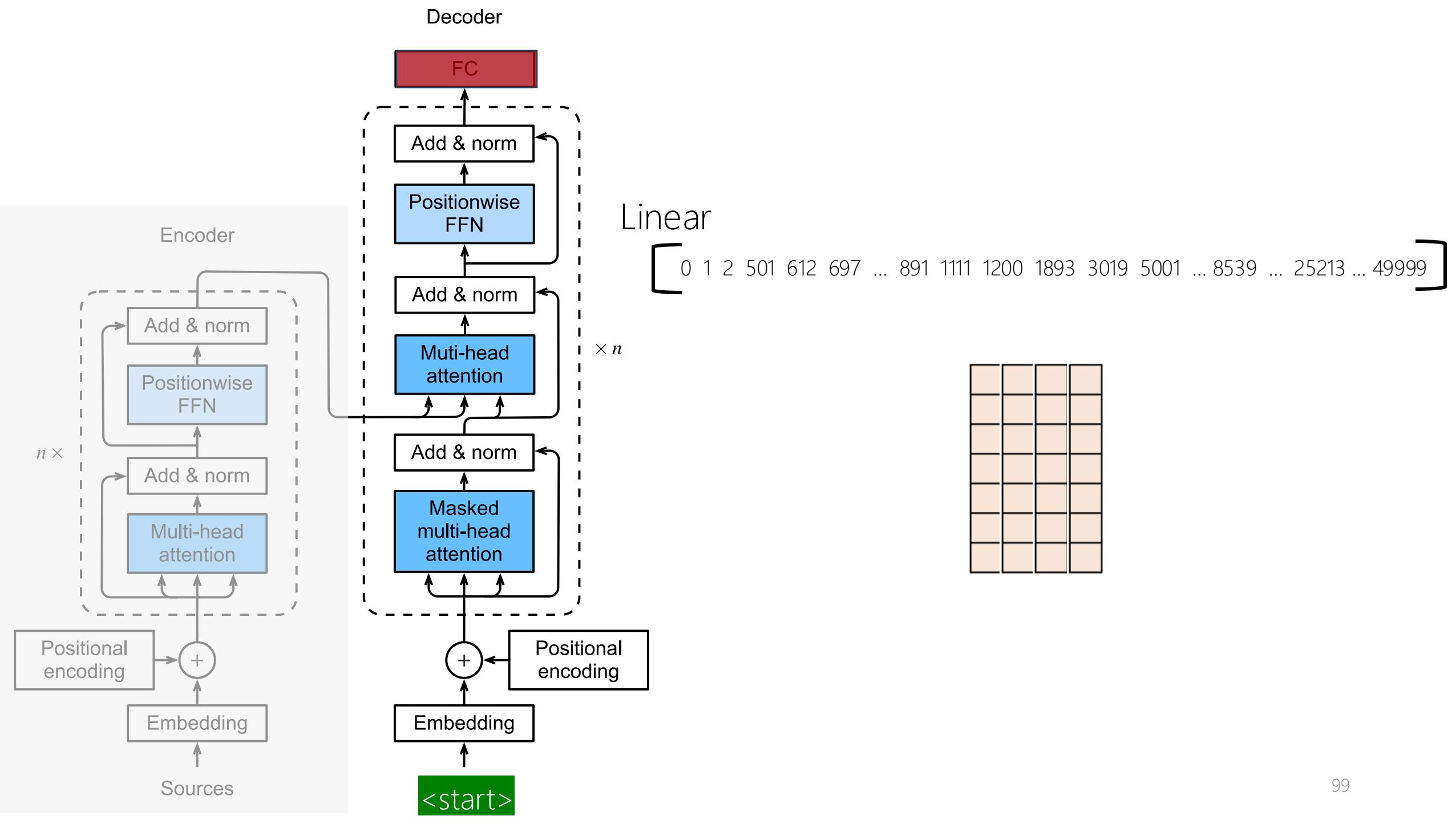


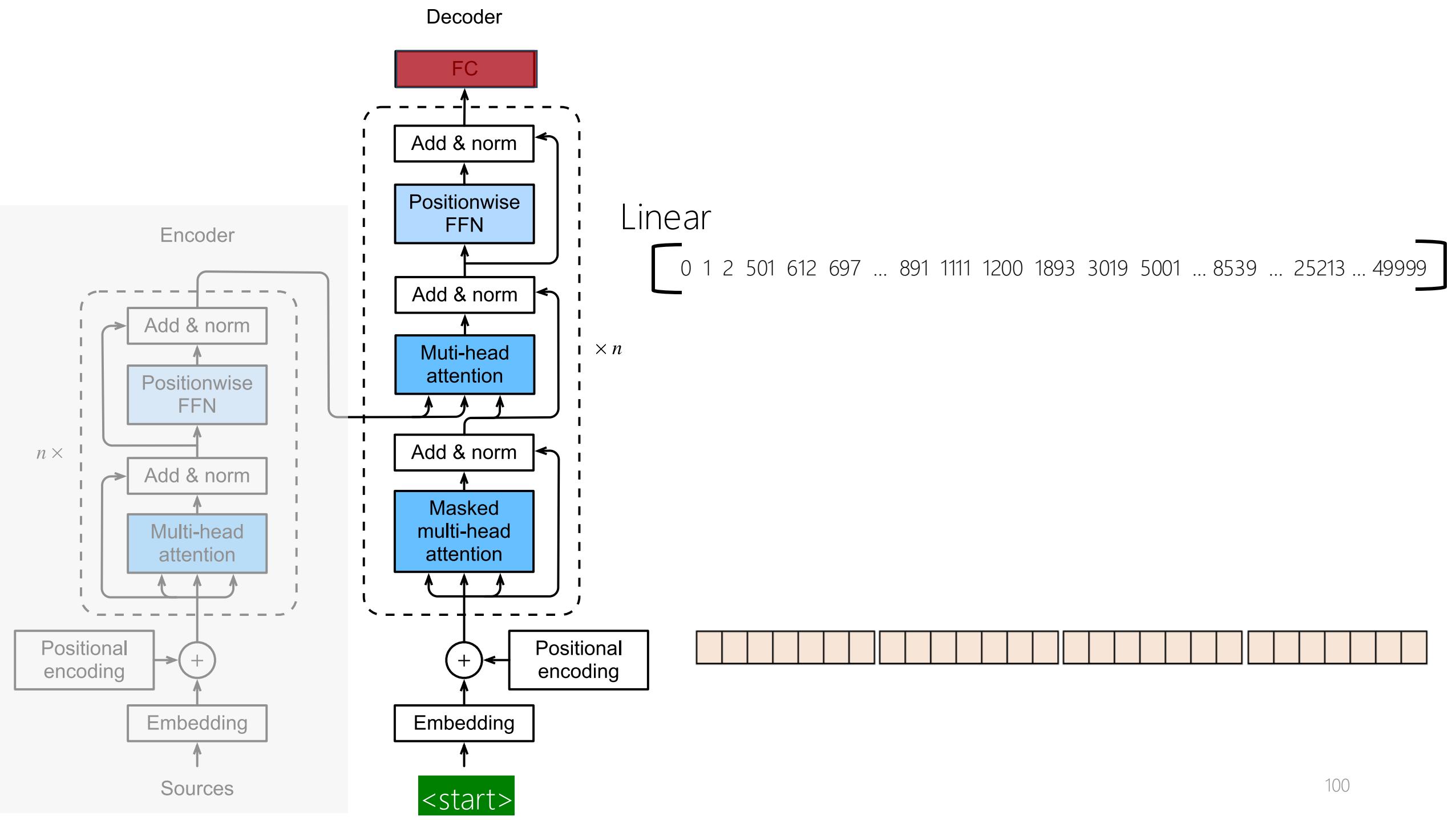


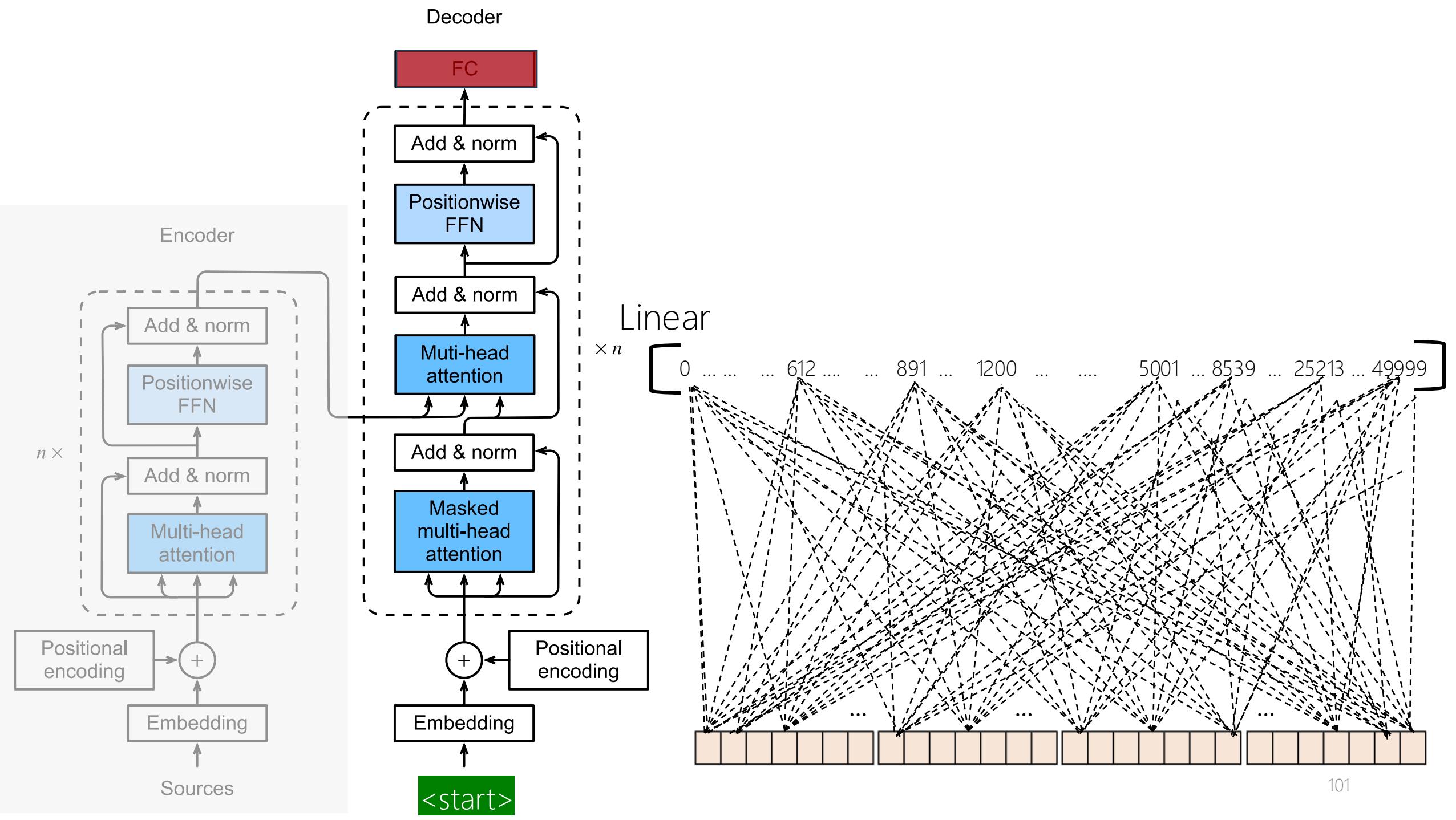


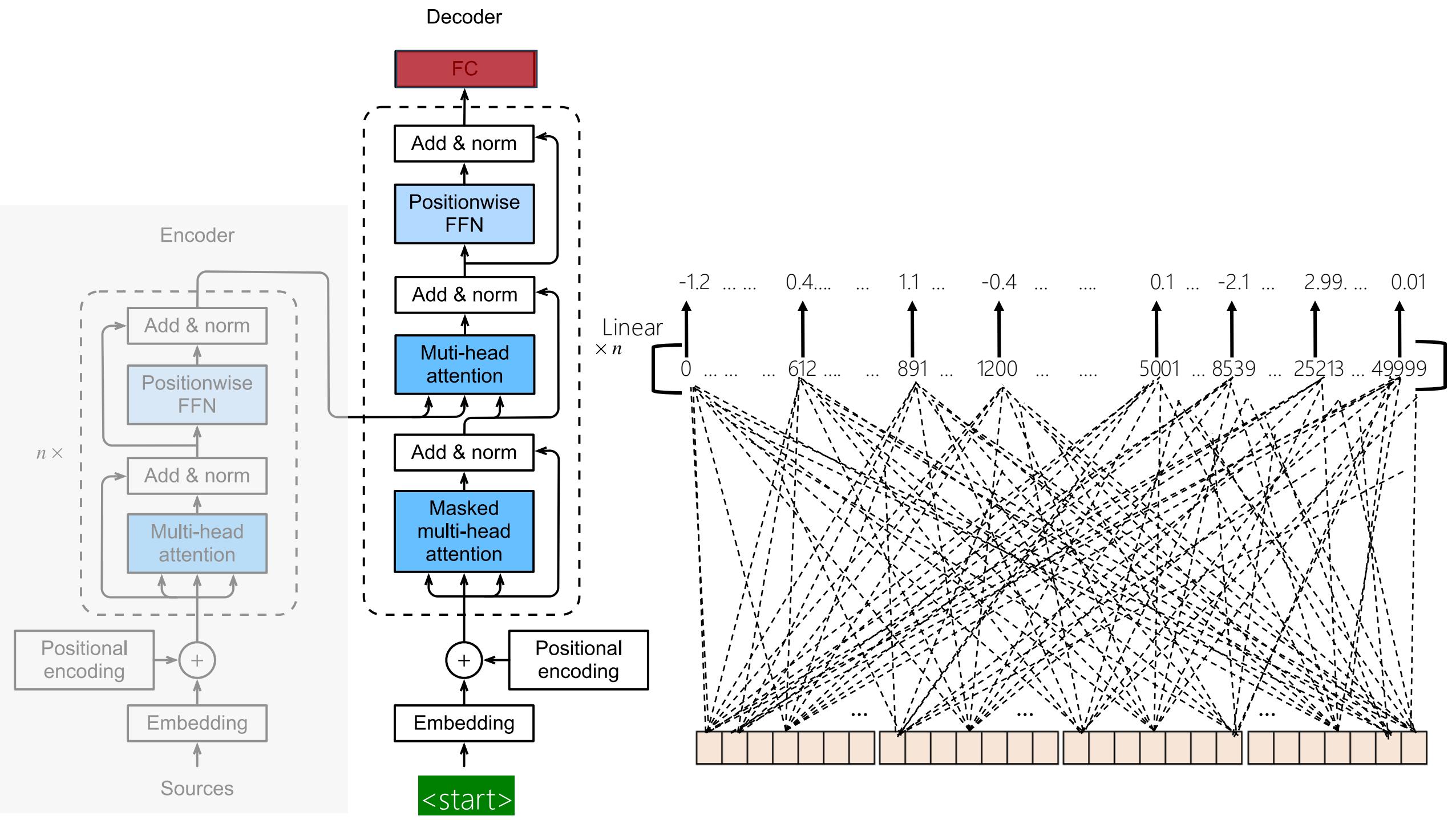
Wait a minute!!! Is this even possible?

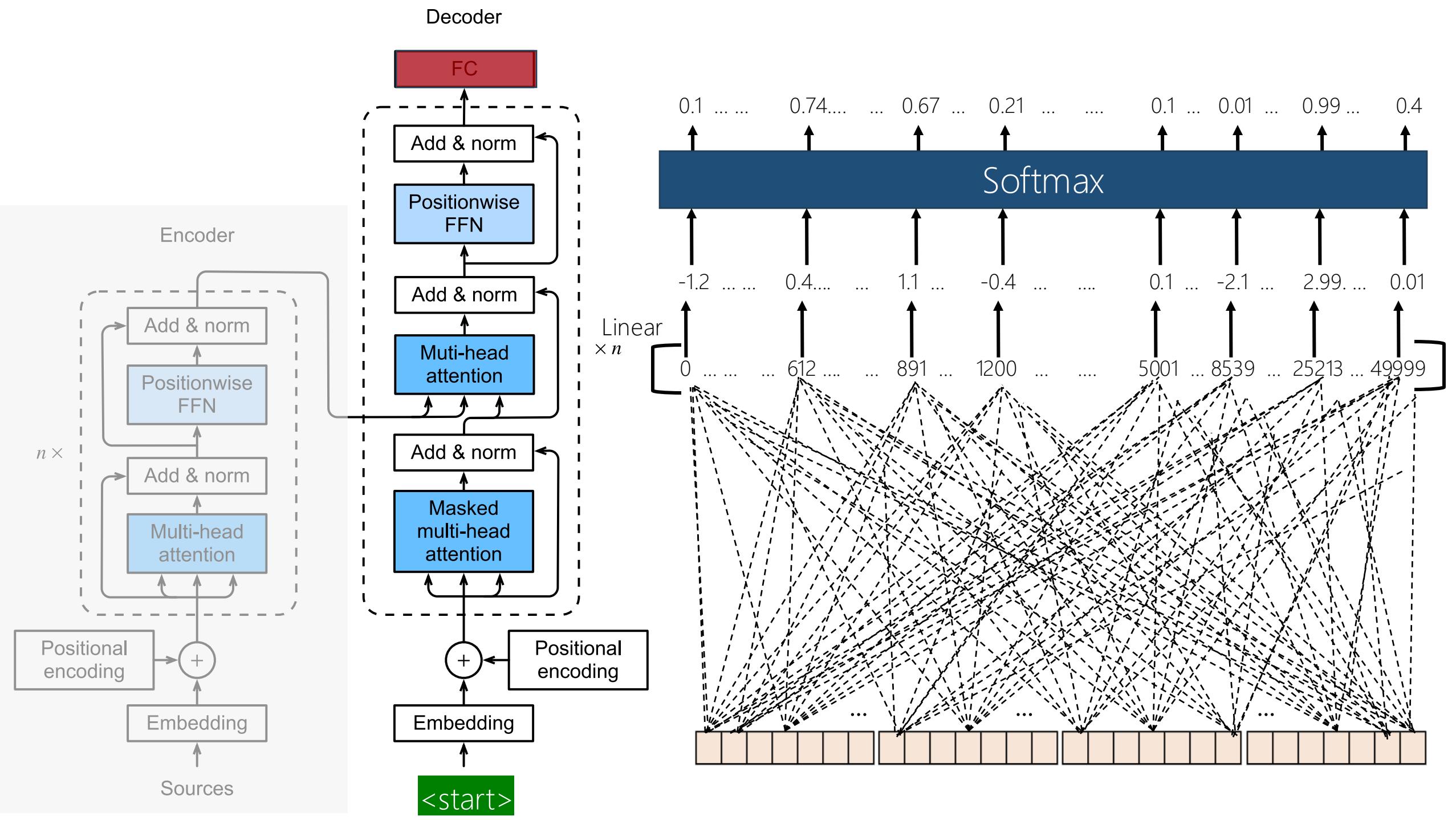


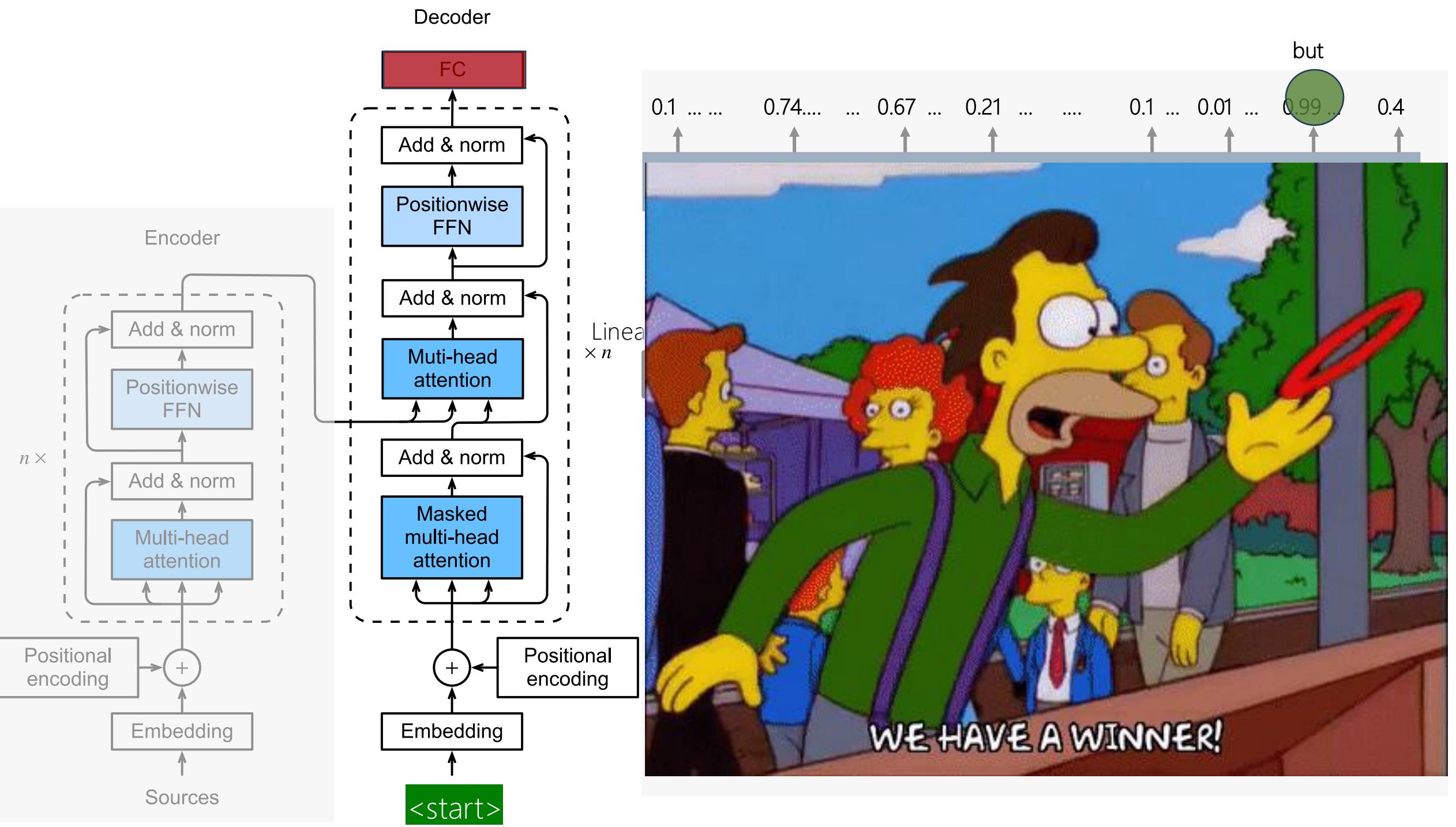




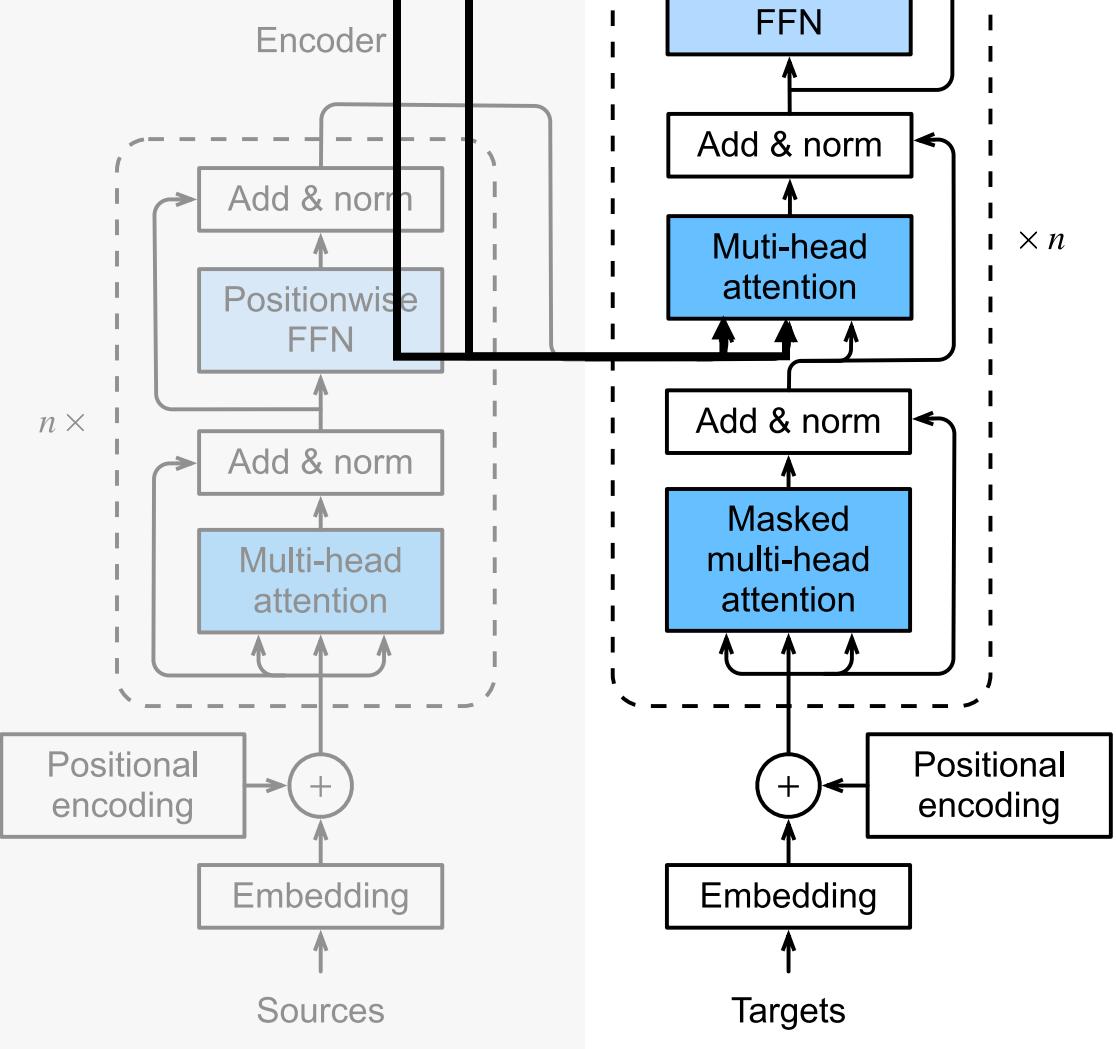




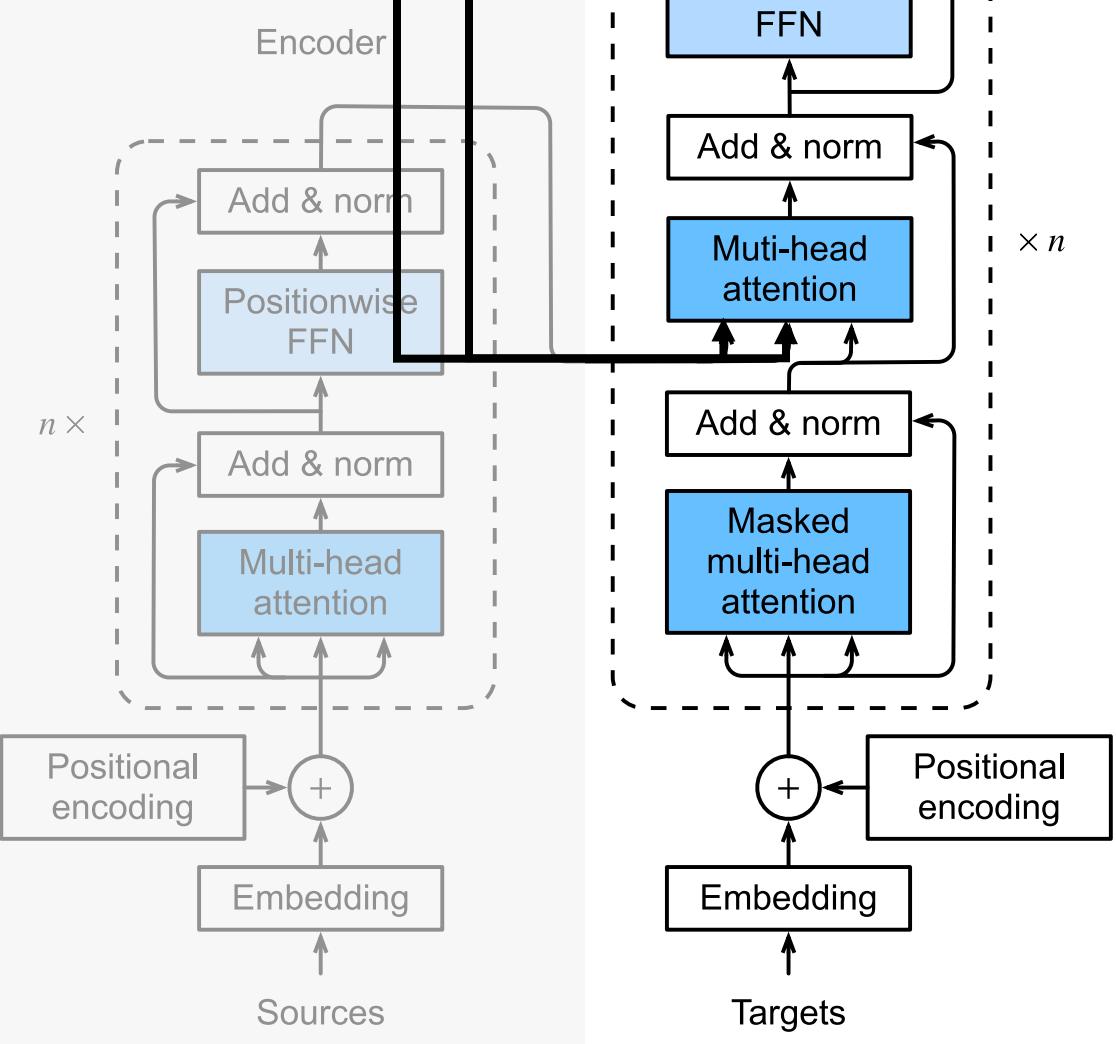




It	-0.453	0.380	-1.414	1.482
matters	-0.453	0.380	-1.414	1.482
not	-0.910	-0.943	1.325	-1.785
what	0.249	-1.927	-0.673	-0.843
someone	-1.404	0.662	-2.608	1.024
is	-0.095	-1.873	1.620	-0.511
born	0.119	-2.263	-0.119	-0.379
	-0.905	0.771	-0.217	-2.326

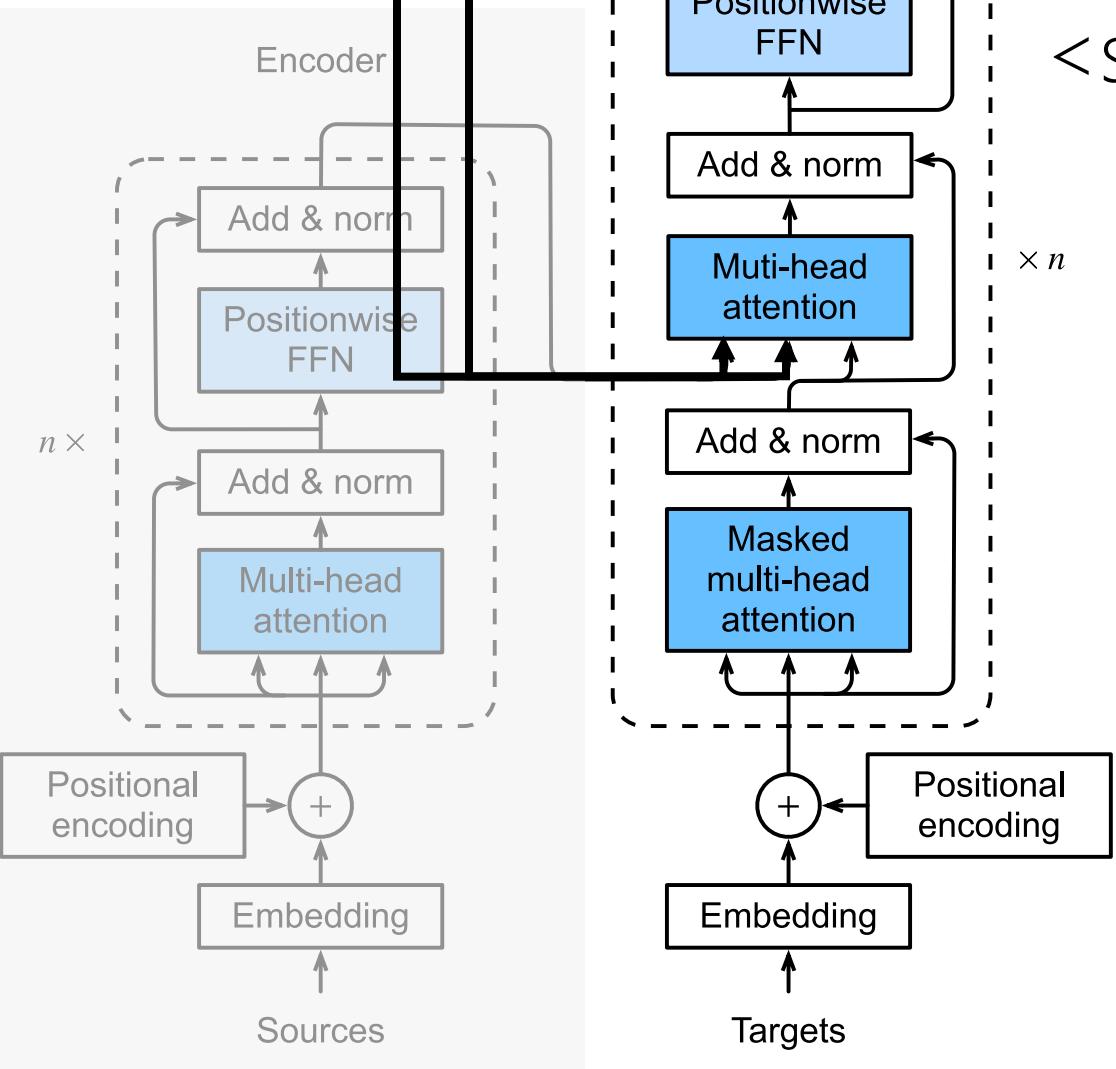


It	-0.453	0.380	-1.414	1.482
matters	-0.453	0.380	-1.414	1.482
not	-0.910	-0.943	1.325	-1.785
what	0.249	-1.927	-0.673	-0.843
someone	-1.404	0.662	-2.608	1.024
is	-0.095	-1.873	1.620	-0.511
born	0.119	-2.263	-0.119	-0.379
	-0.905	0.771	-0.217	-2.326



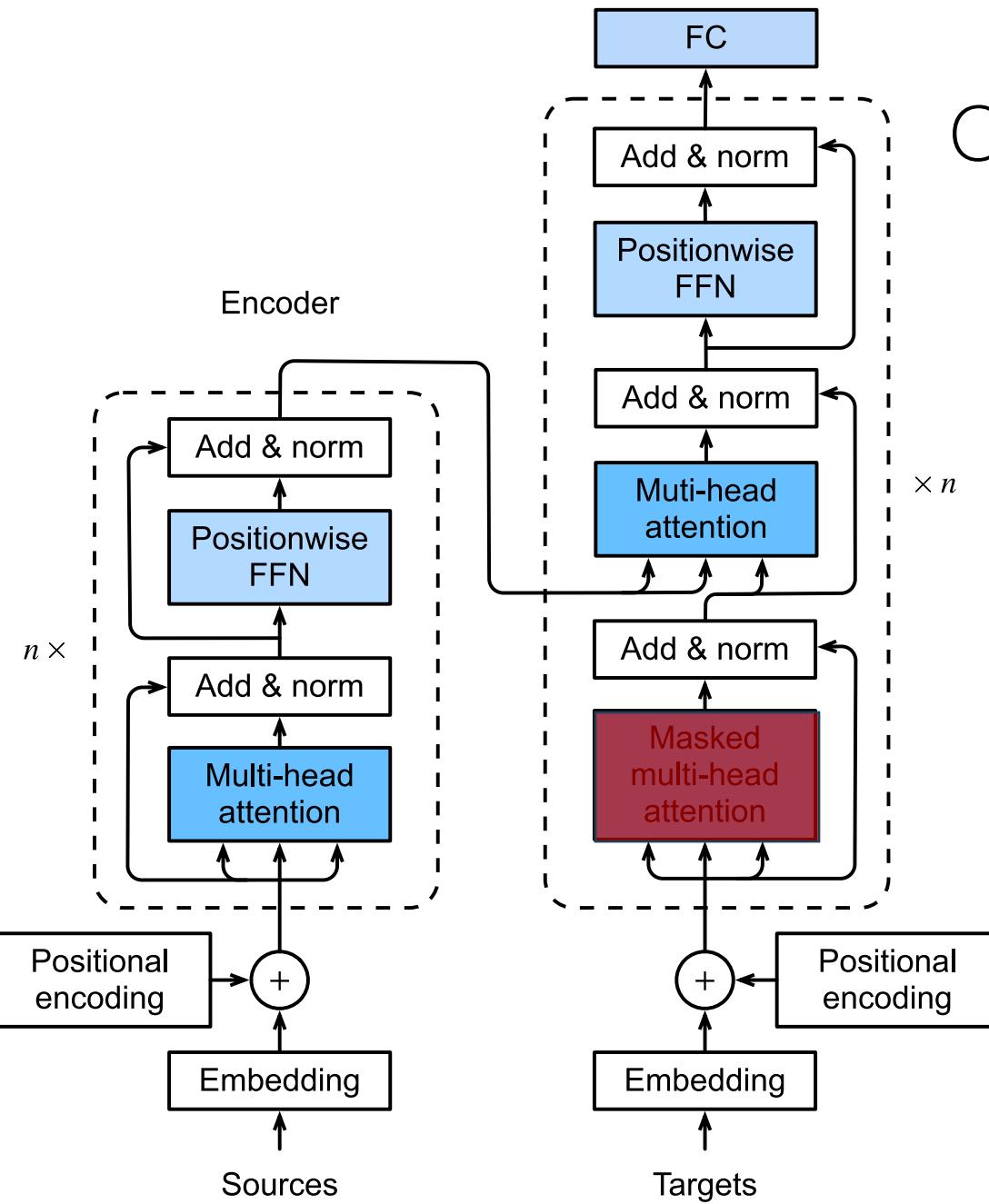
<start> but

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!!!



# 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

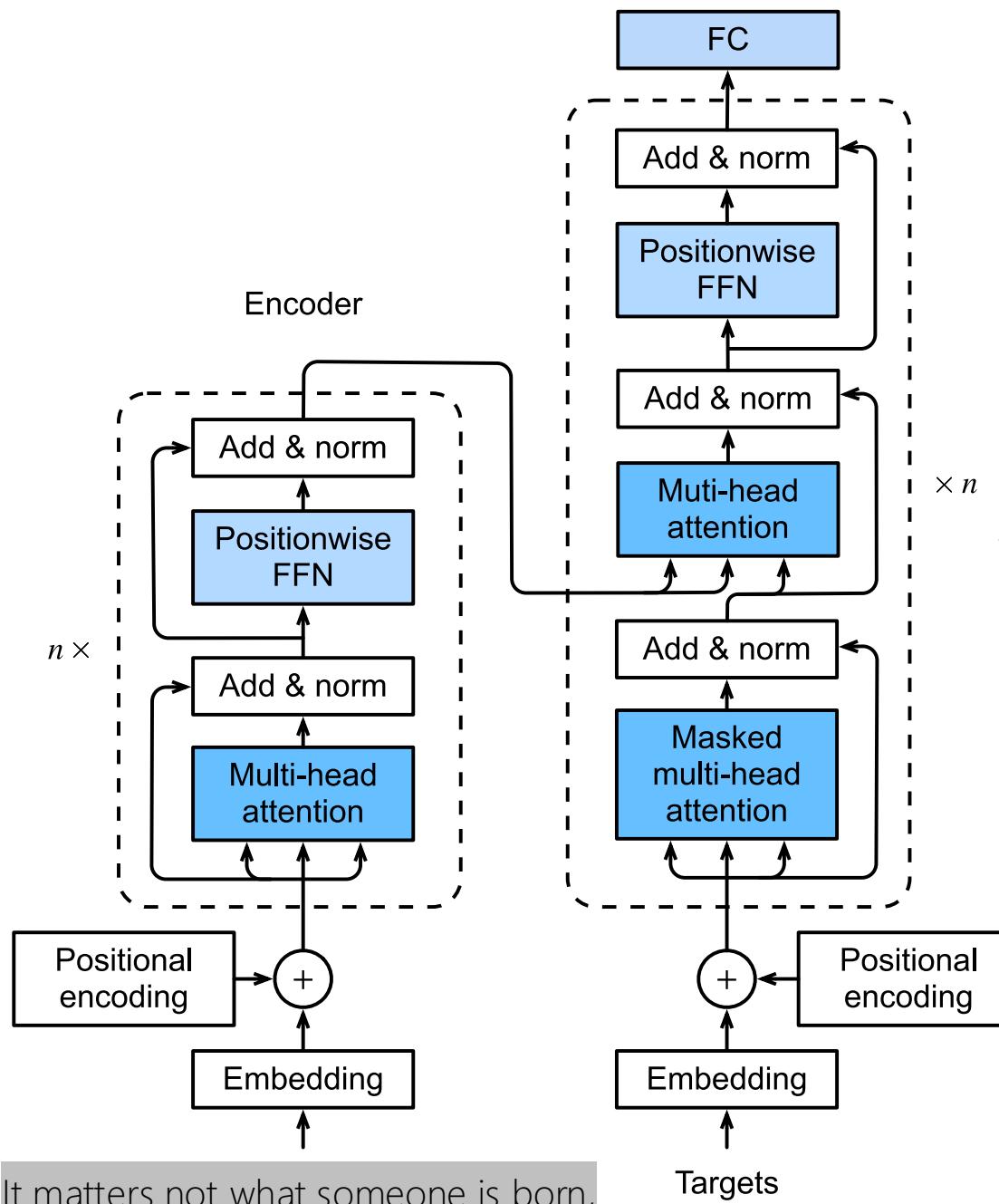
<start> far more than our abilities  
<end>

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

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

# Training

Decoder <start> MASK MASK MASK MASK MASK MASK MASK<end>

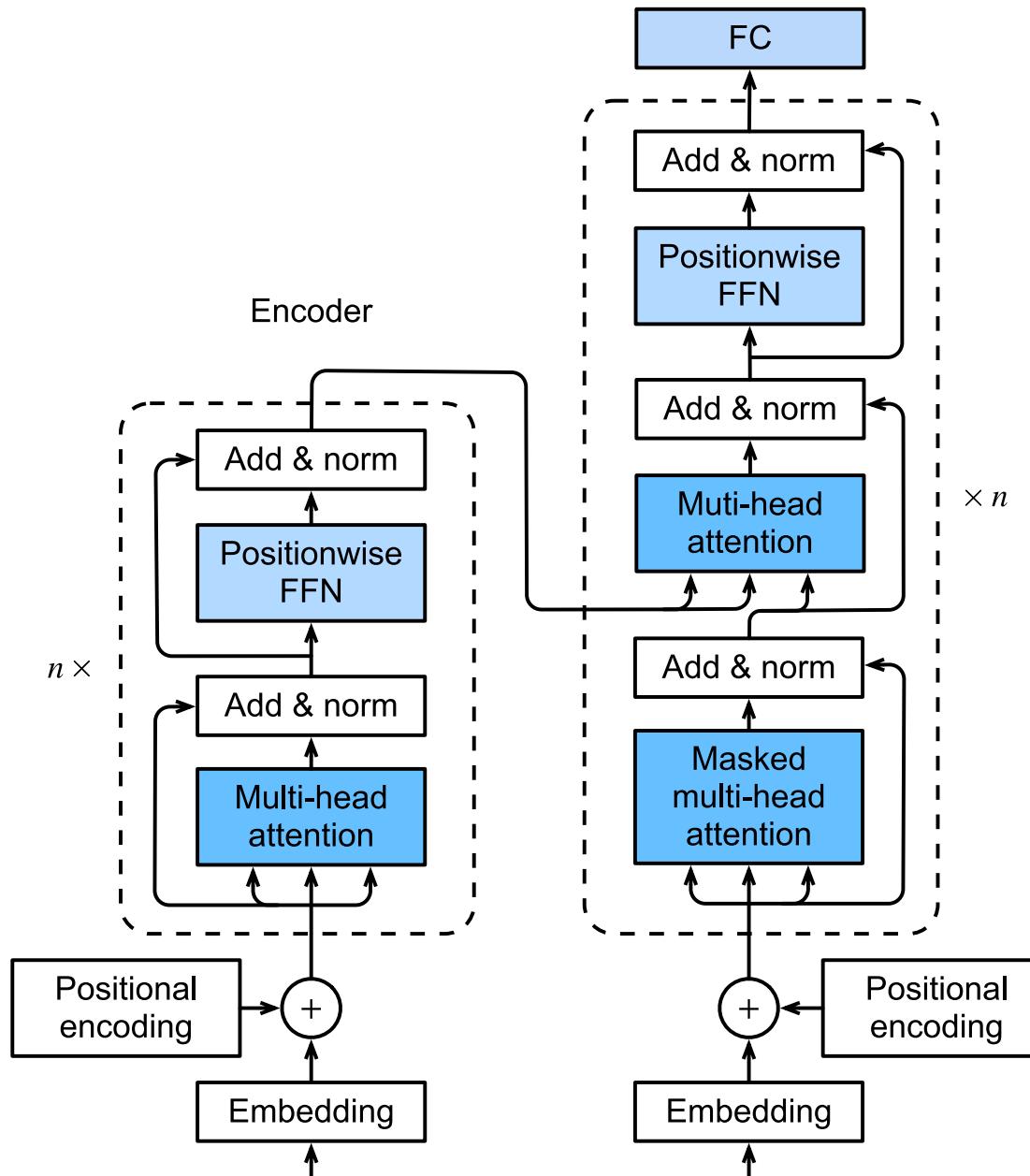


Why MASK? 🤔



# Training

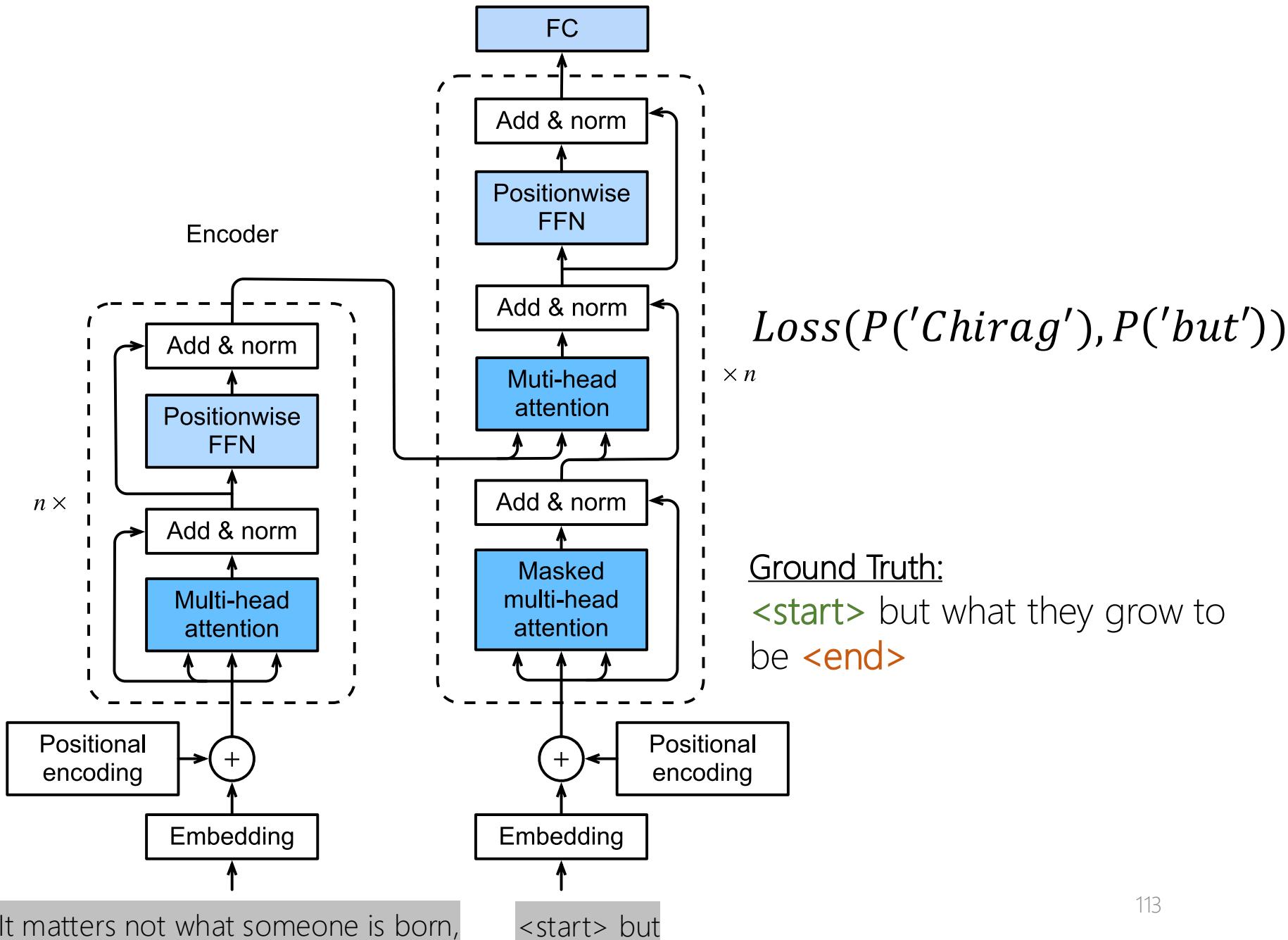
Decoder Chirag MASK MASK MASK MASK MASK MASK <end>



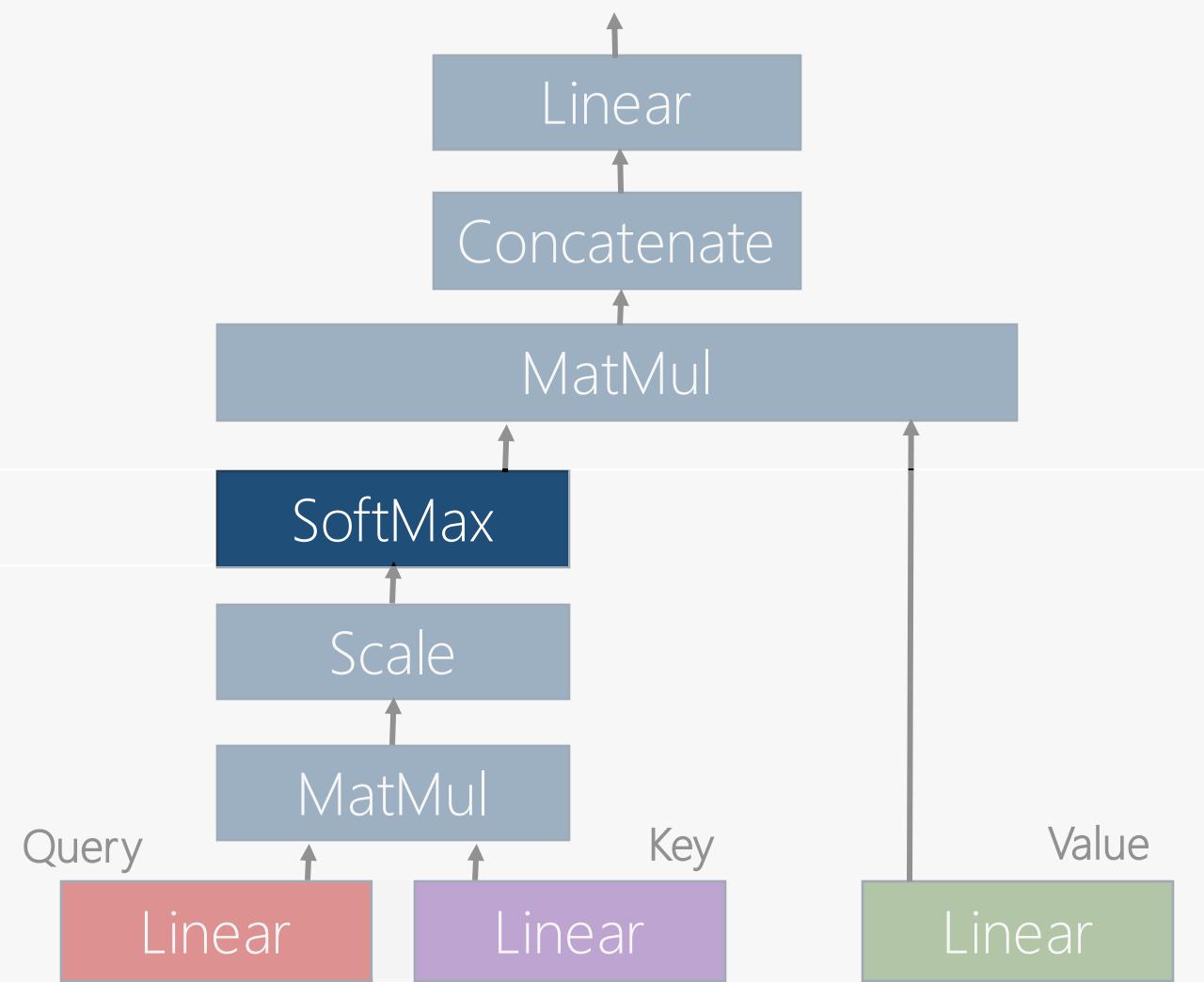
It matters not what someone is born,

<start>

# Training



# Masked Attention



	<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

	<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

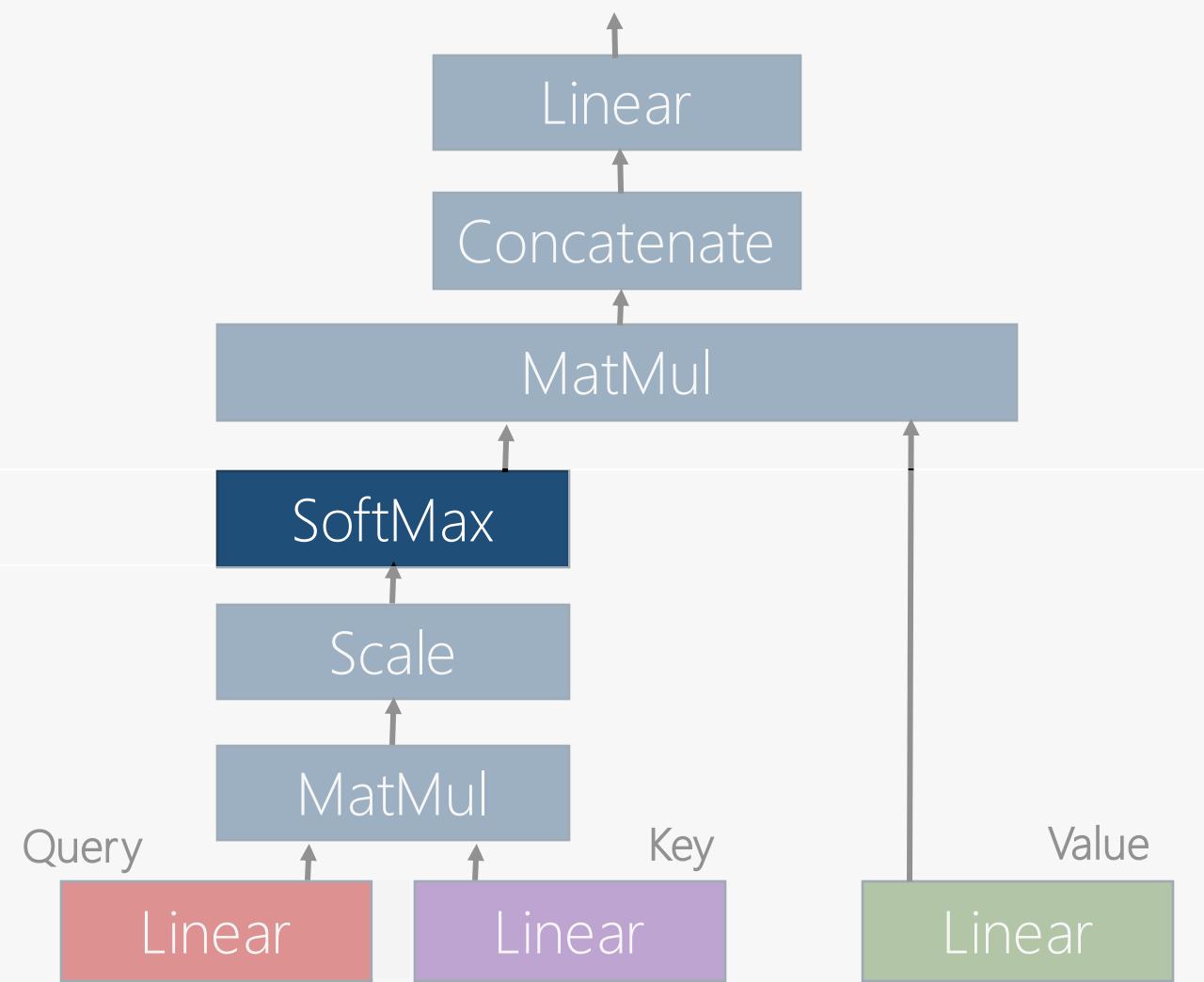
Attention Filter

+

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

Mask Filter

# 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

# Thank you!