

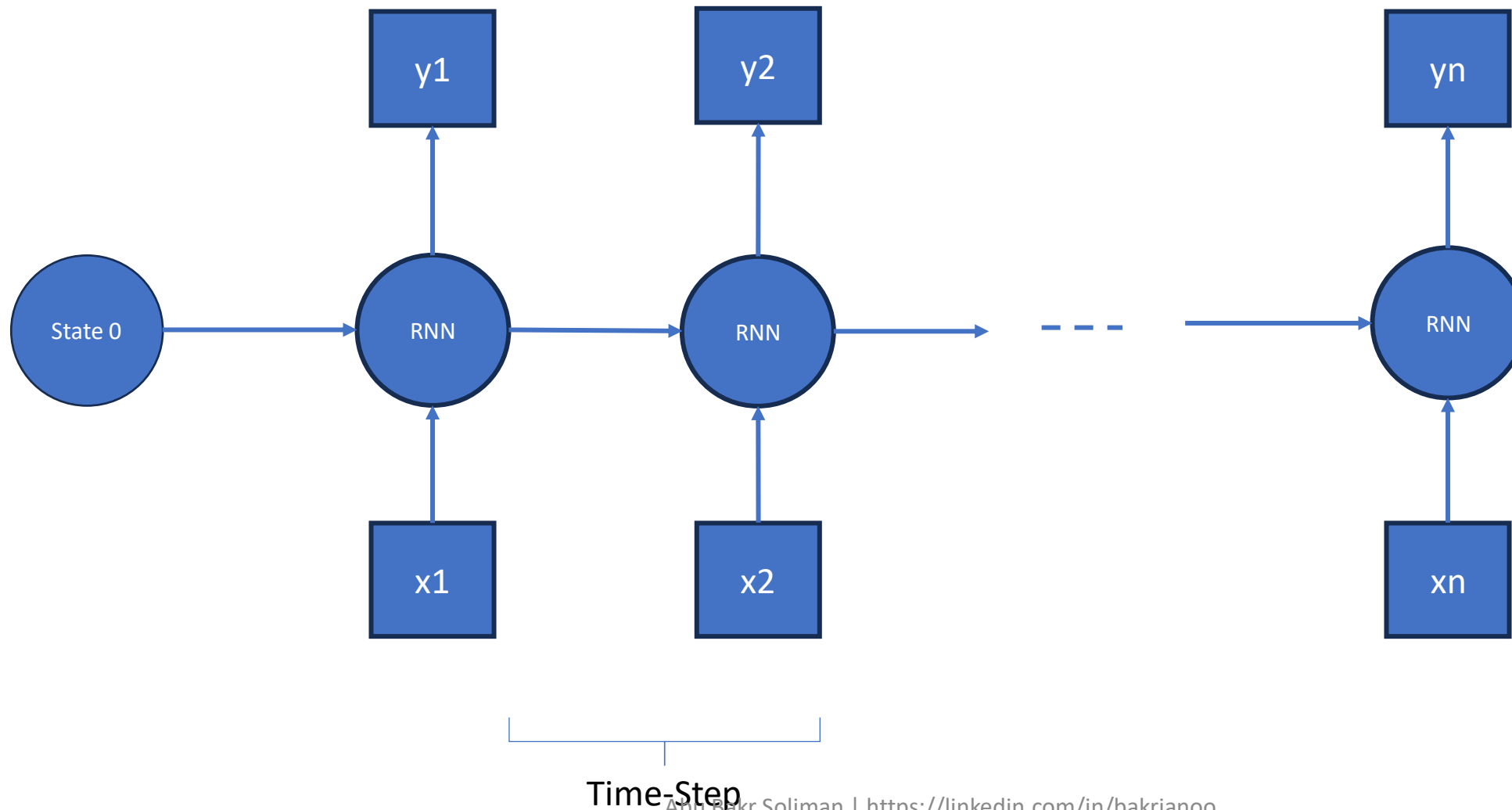
Transformers

Abu Bakr Soliman

Sequence to Sequence (Seq2Seq)

Input	x1	x2	..	xn
Output	y1	y2	..	yn

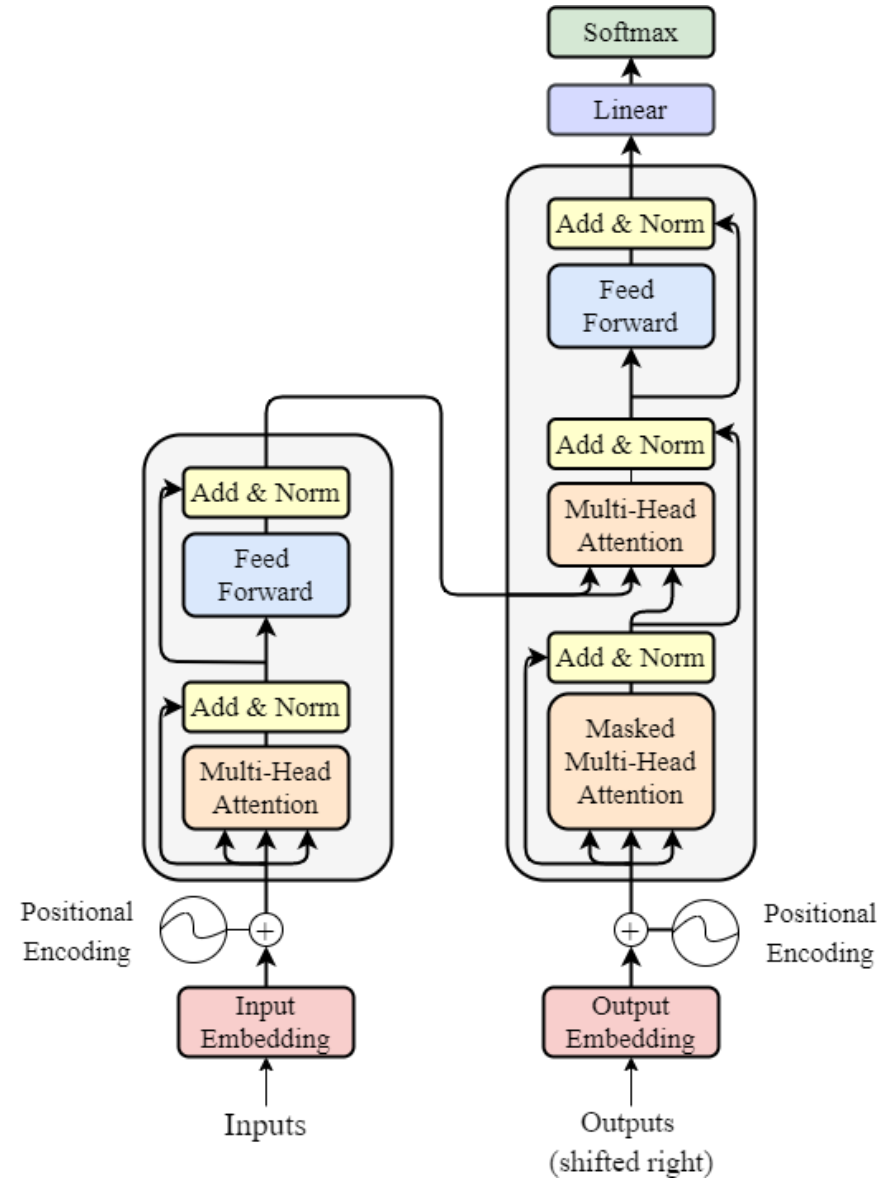
Recurrent Neural Network (RNN)



RNN Problems

- Long Sequences = Slow Computations
- Vanishing and Exploding Gradients
- Vanished Memory

Attention is All You Need



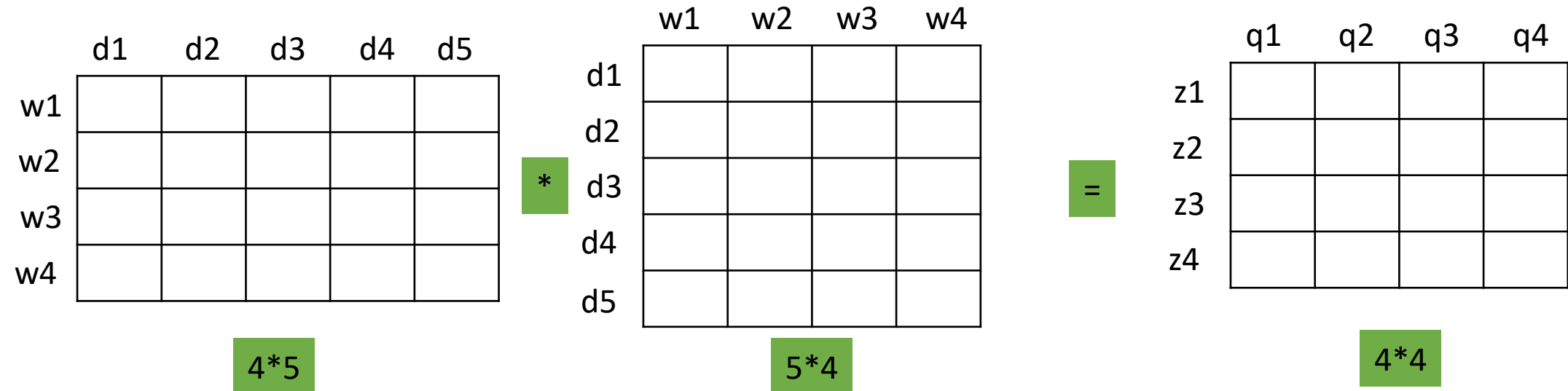
Vectors

Patient ID
10200
3600

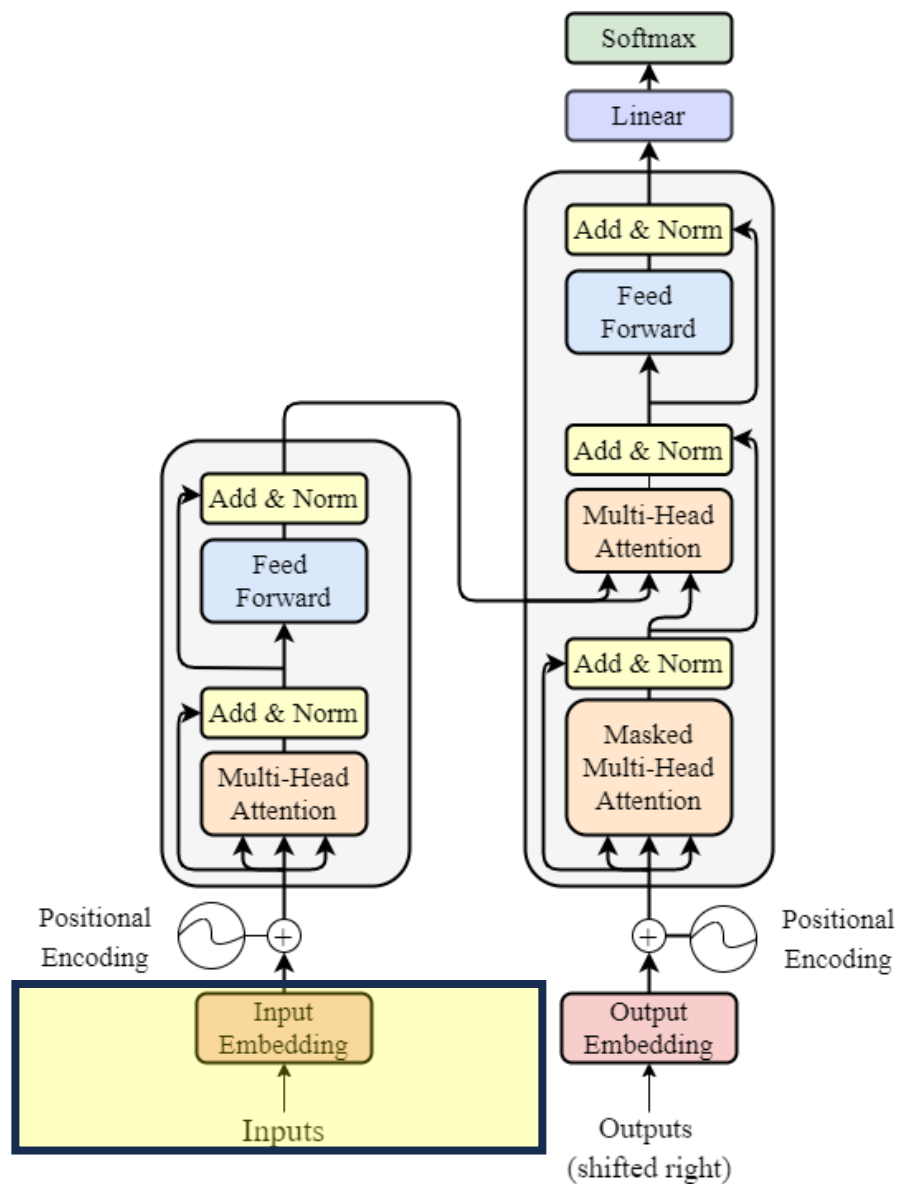
Vectors

Patient ID	Height (CM)	Weight (KG)	Systolic (mmHg)	Diastolic (mmHg)
10200	165	71	120	80
3600	180	92	110	85

Matrices Operations



Attention is All You Need



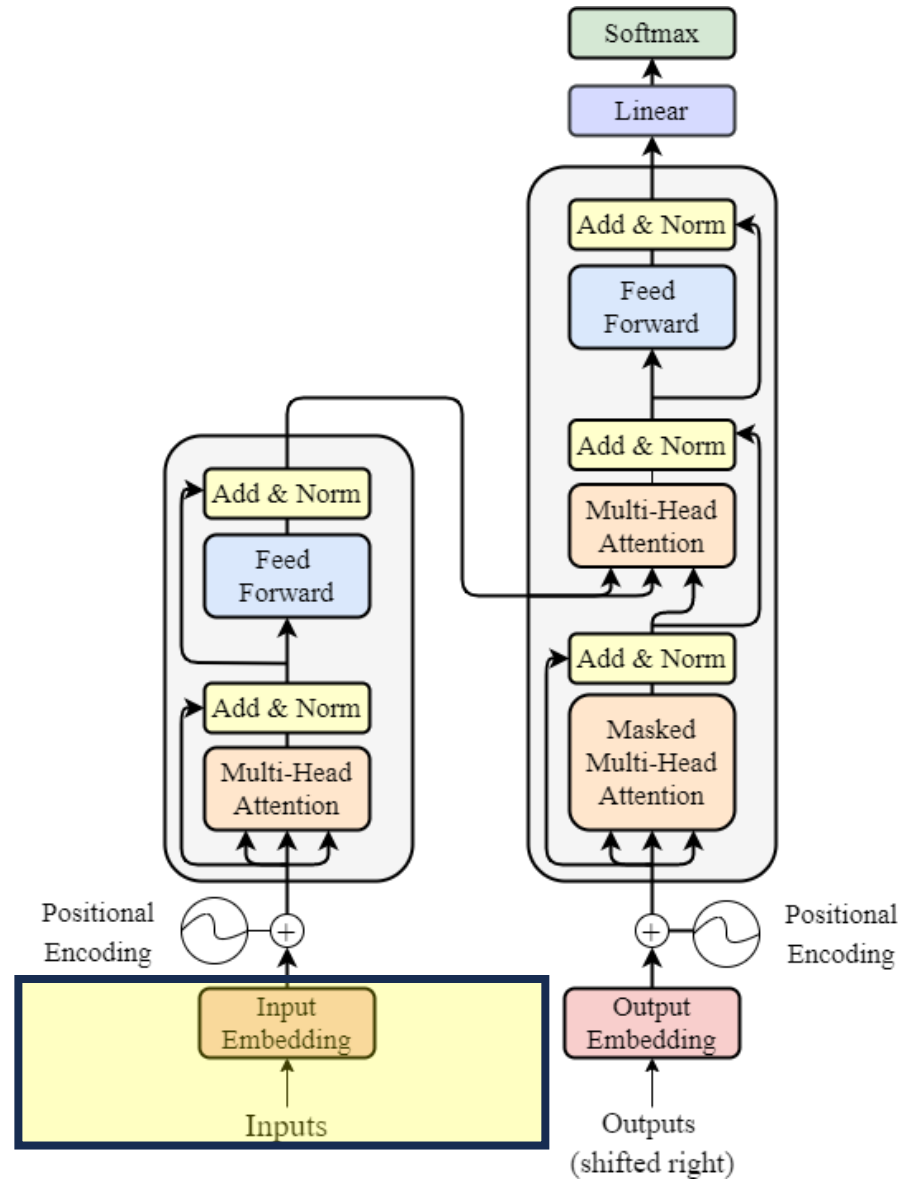
Tokenizer

Token	ID
Hello	1
Go	2
ed	3
..	..
red	32000

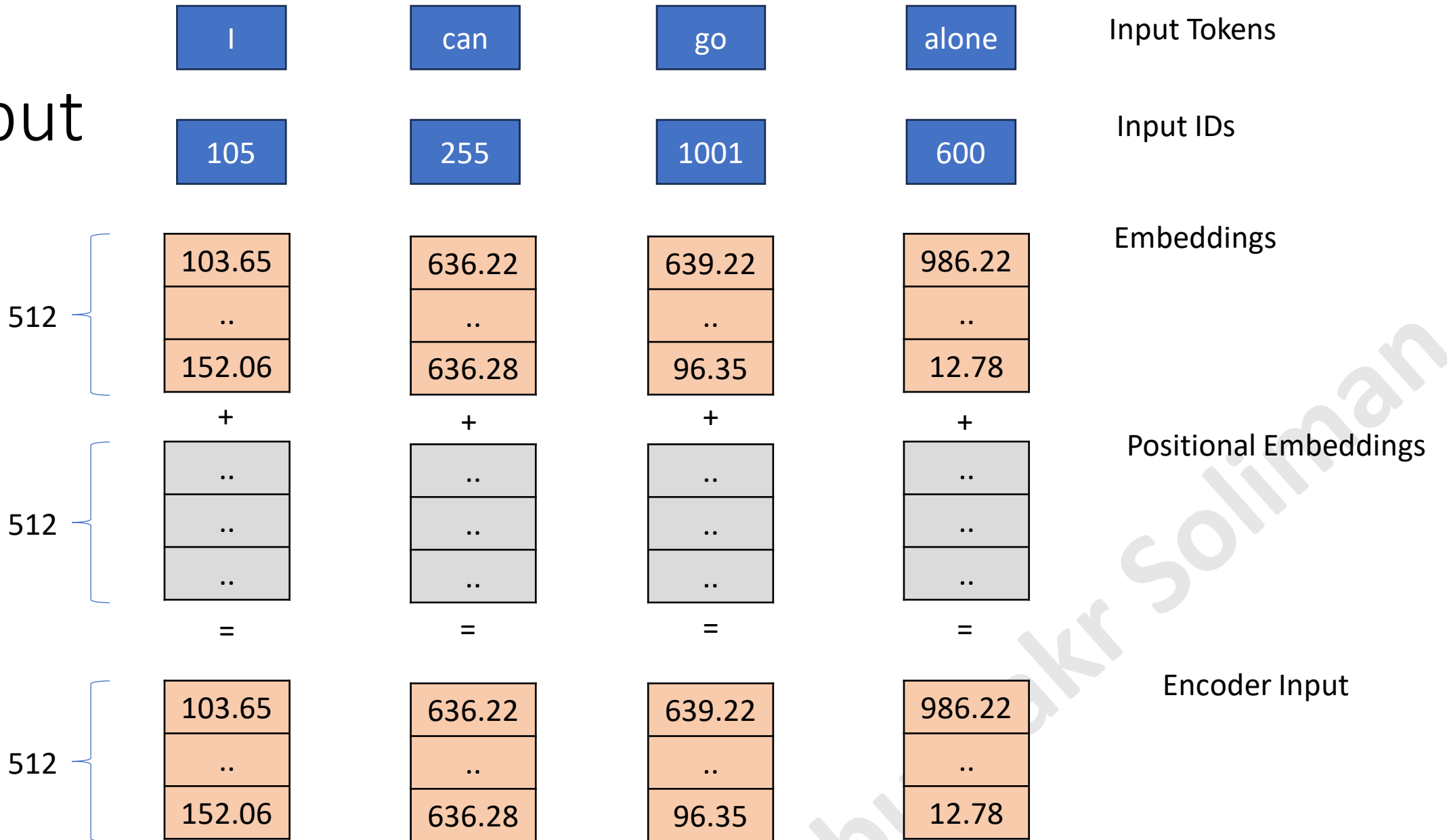
Input

I can go alone

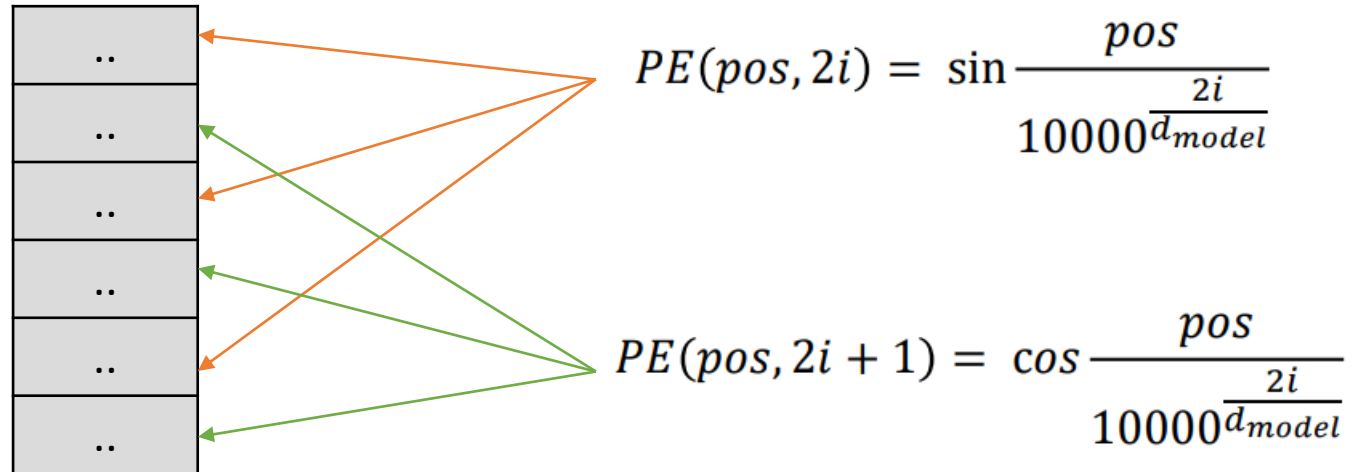
	I	can	go	alone	Input Tokens
	105	255	1001	600	Input IDs
512	<div>103.65</div> <div>633.01</div> <div>25.33</div> <div>..</div> <div>152.06</div>	<div>636.22</div> <div>2.01</div> <div>96.25</div> <div>..</div> <div>636.28</div>	<div>639.22</div> <div>9.36</div> <div>78.22</div> <div>..</div> <div>96.35</div>	<div>986.22</div> <div>7.22</div> <div>9.36</div> <div>..</div> <div>12.78</div>	Embeddings

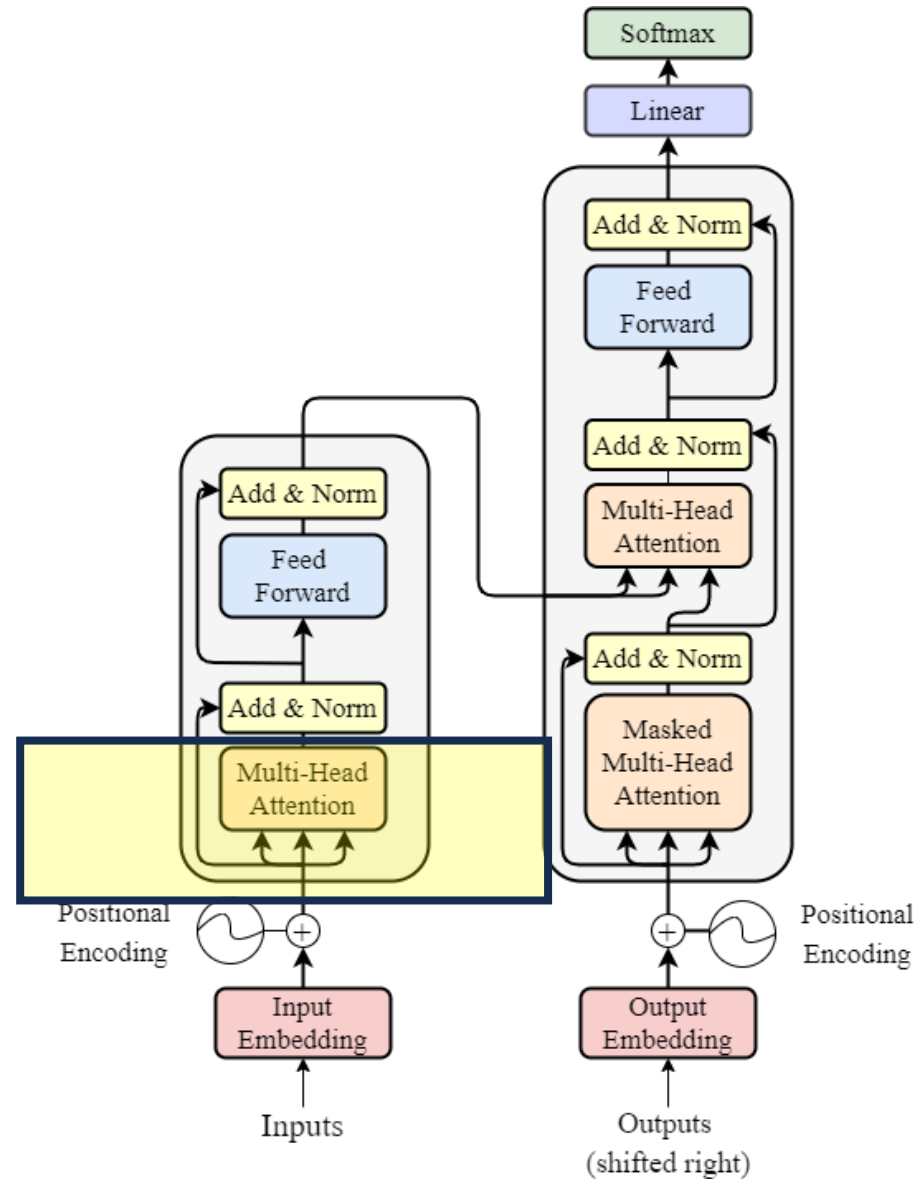


Input



I	can	go	alone	Input Tokens
105	255	1001	600	Input IDs

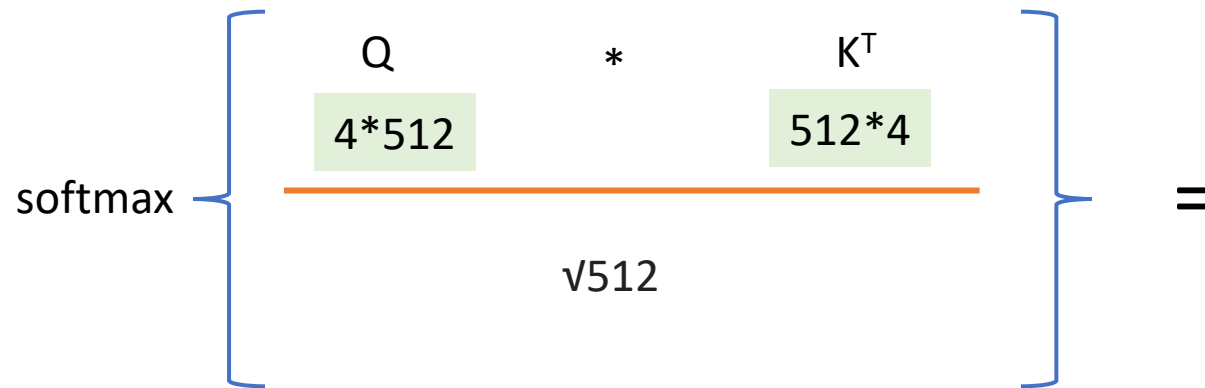




Self-Attention

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

$d_k = d_{\text{model}}$	512
seq	4



	I	can	go	alone
I	0.7	0.2	0.1	0.1
can	0.3	0.5	0.1	0.1
go	0.1	0.3	0.4	0.2
alone	0.05	0.05	0.1	0.8

4*4

Self-Attention

$$\text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)$$

	I	can	go	alone
I	0.7	0.2	0.1	0.1
can	0.3	0.5	0.1	0.1
go	0.1	0.3	0.4	0.2
alone	0.05	0.05	0.1	0.8

4*4

*

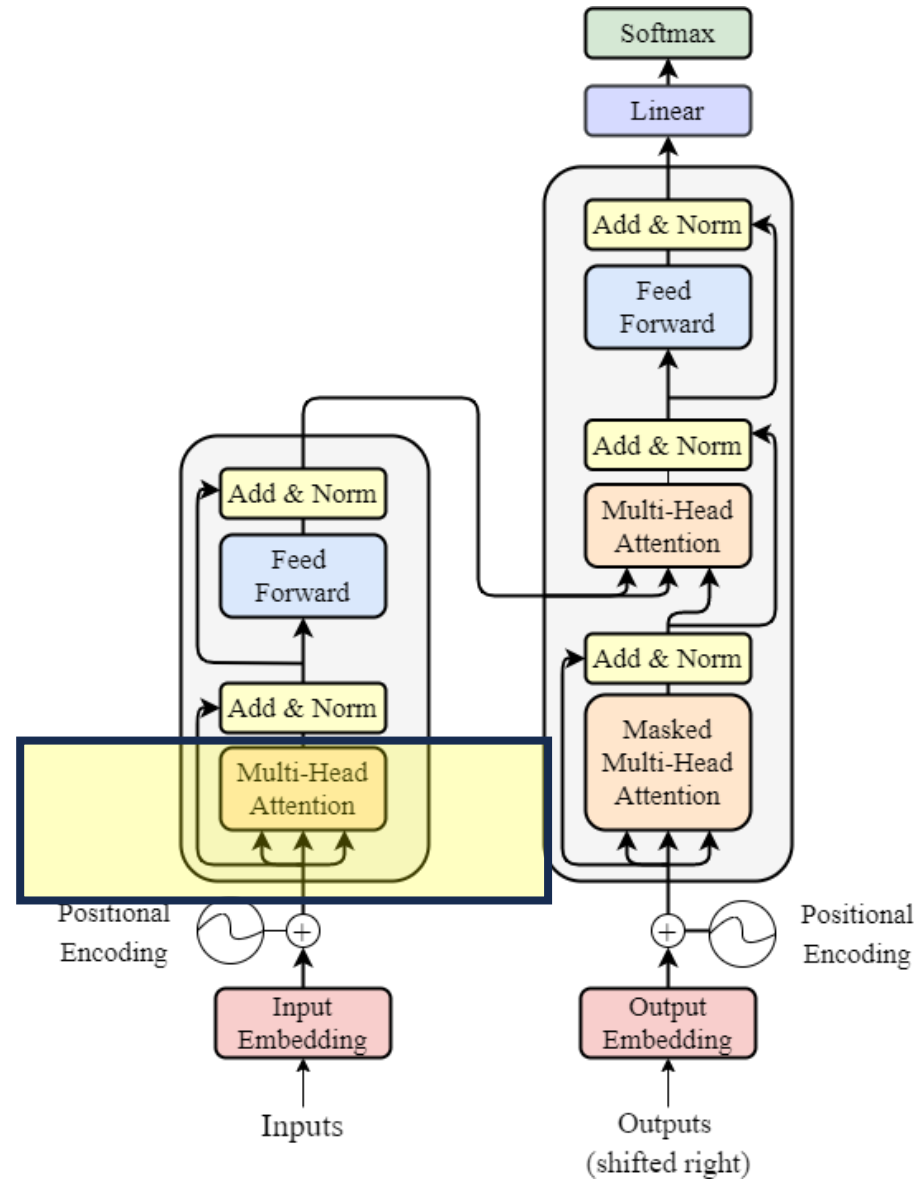
V

=

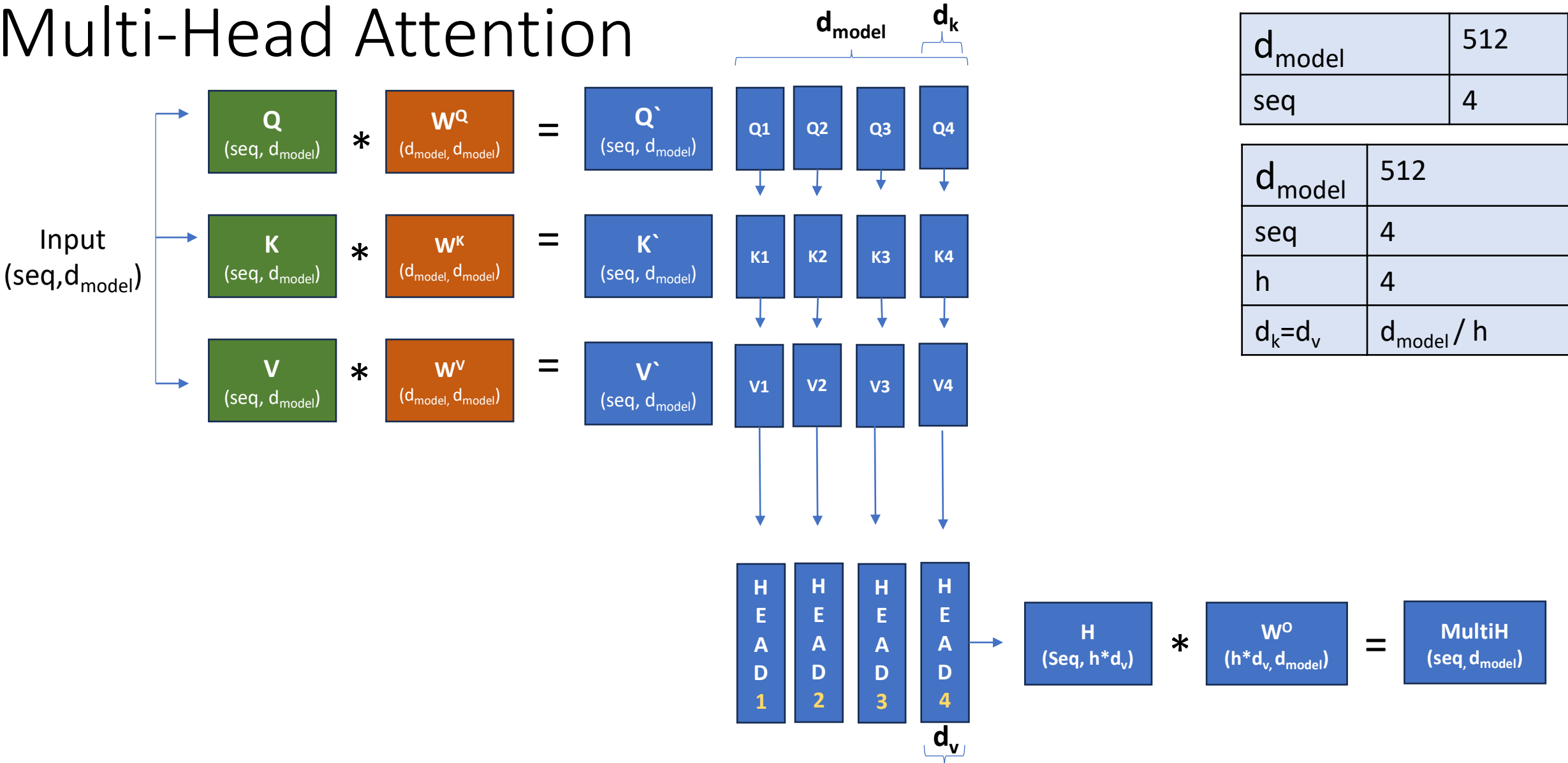
4*512

Attention(Q, K, V)

4*512

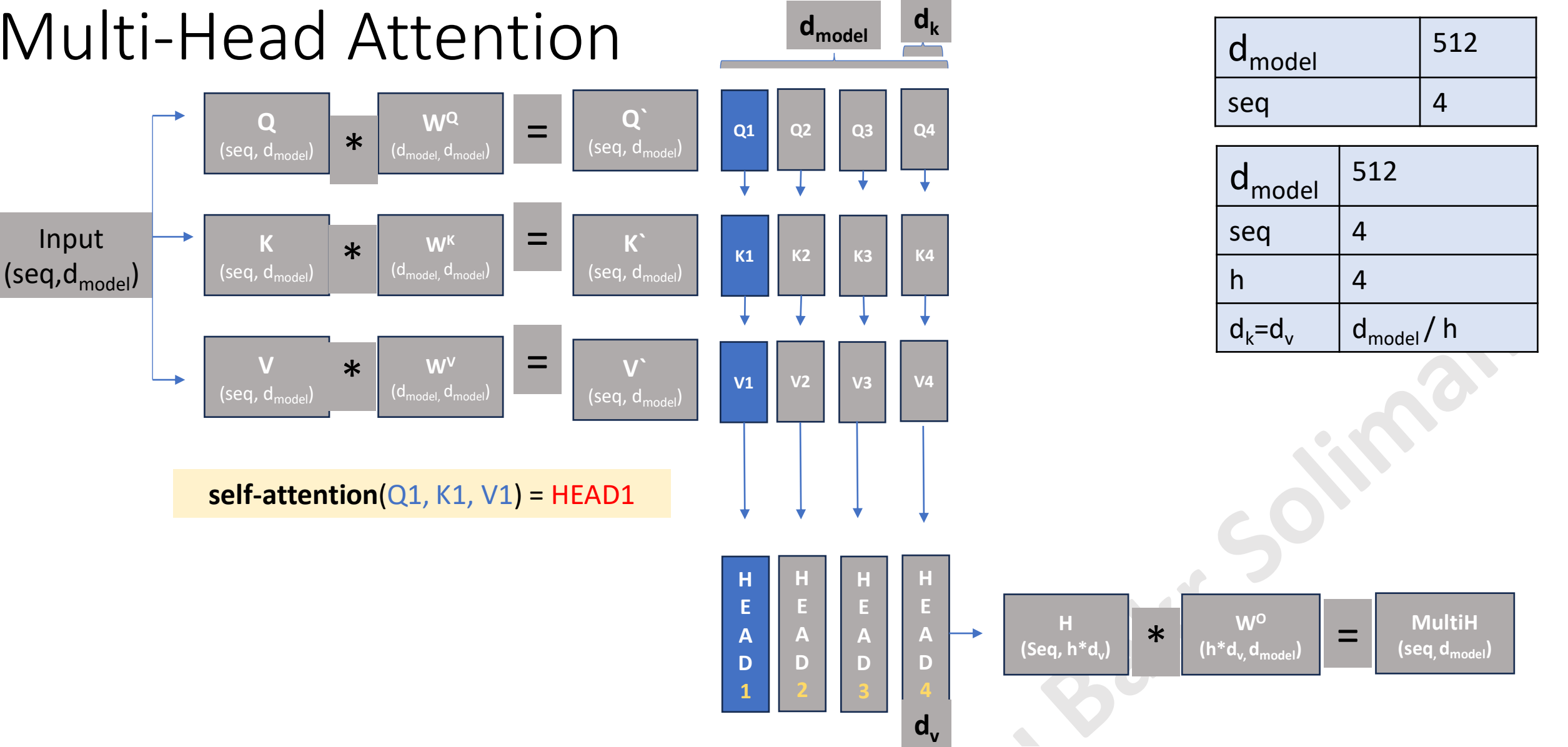


Multi-Head Attention



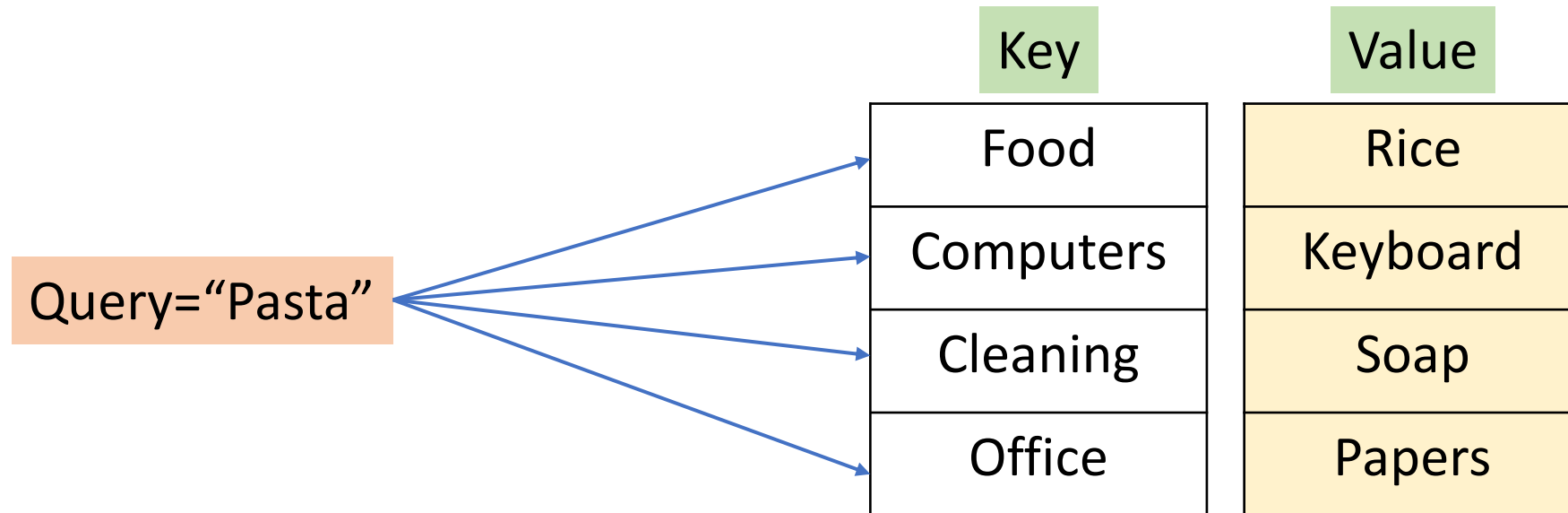
$MultiHead(Q, K, V) = Concat(head_1 \dots head_h)W^O$

Multi-Head Attention

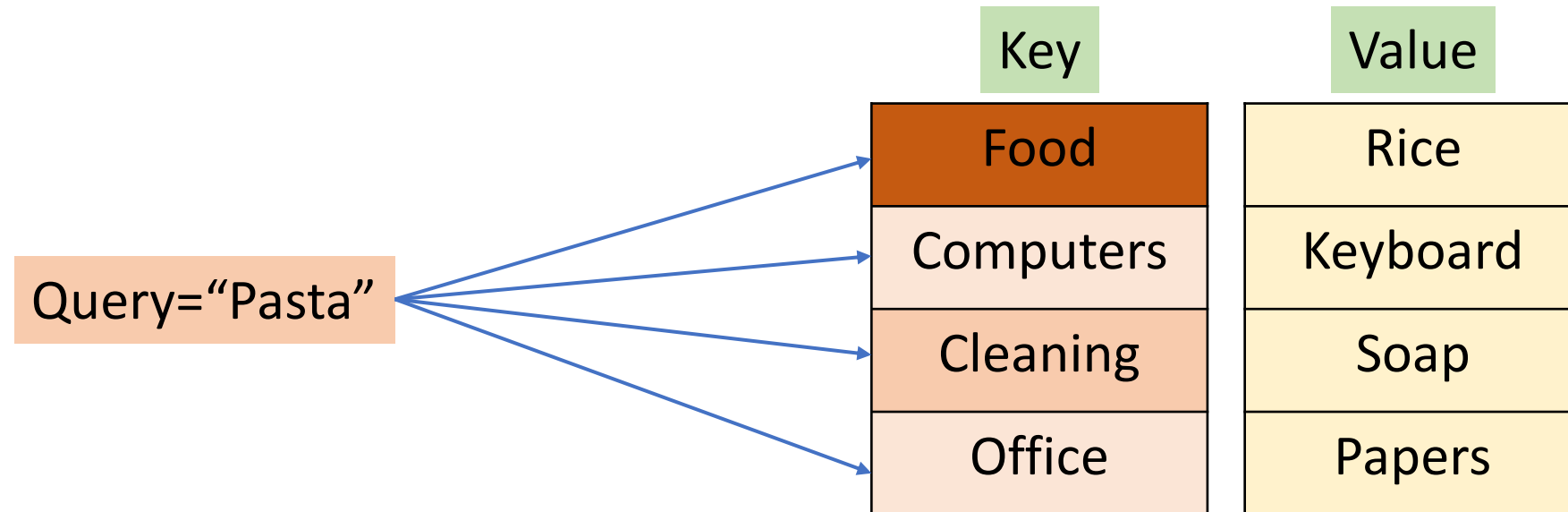


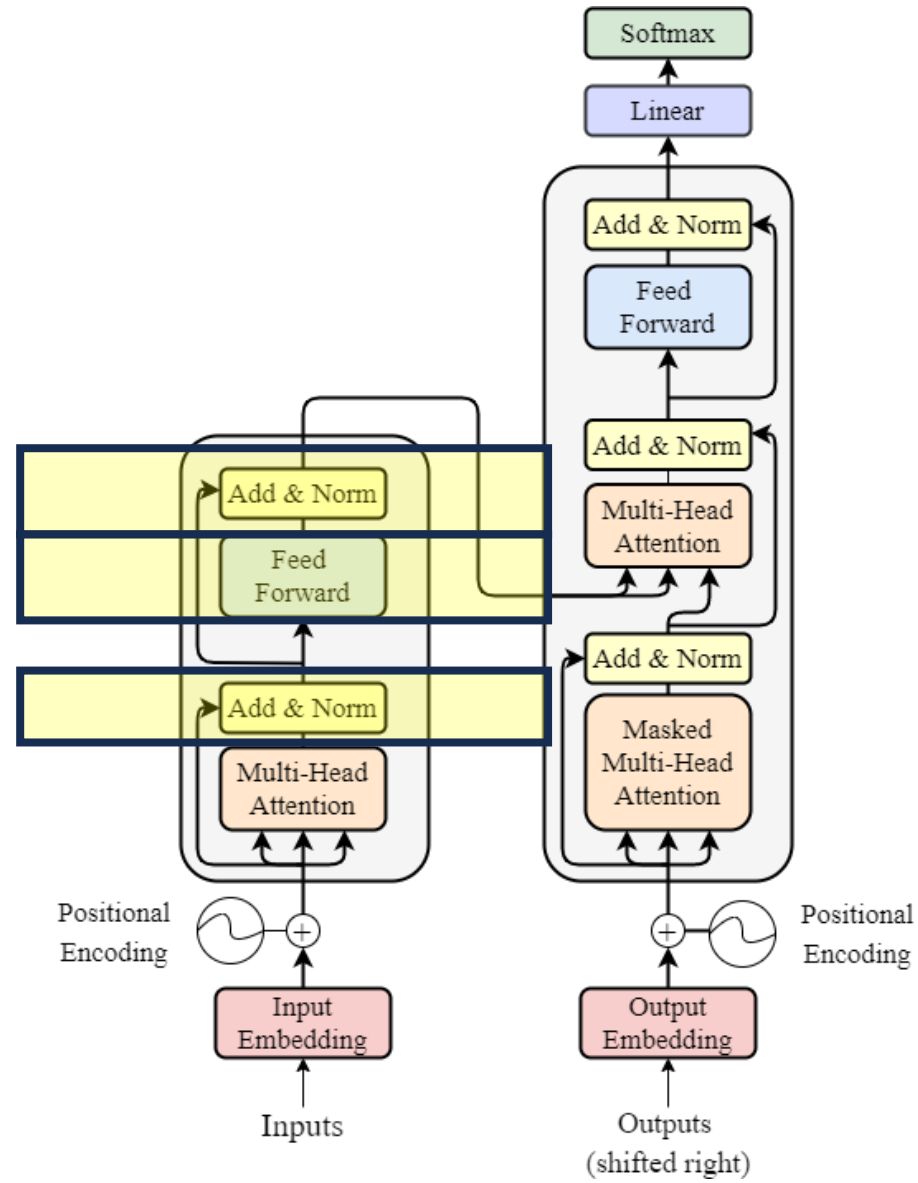
$MultiHead(Q, K, V) = Concat(head_1 \dots head_h)W^O$

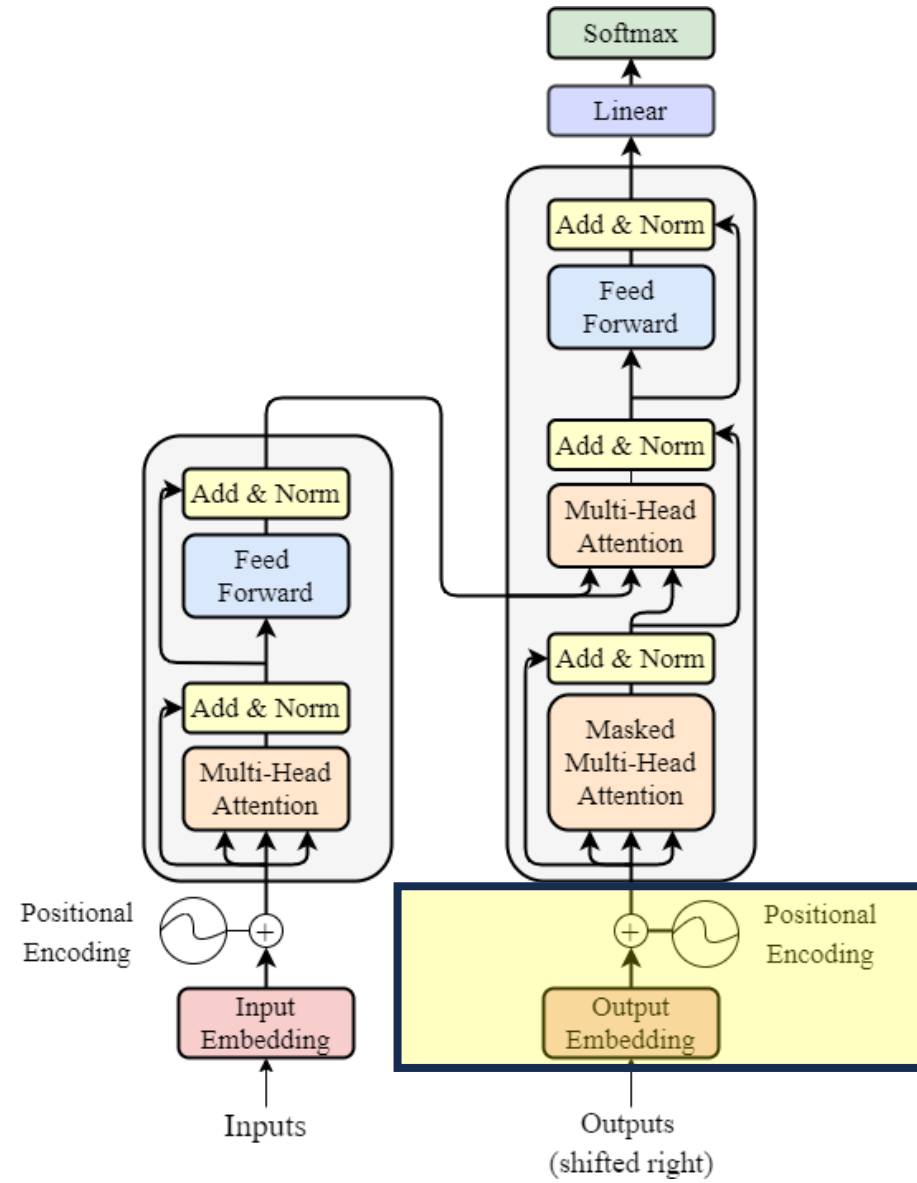
Query, Key & Value

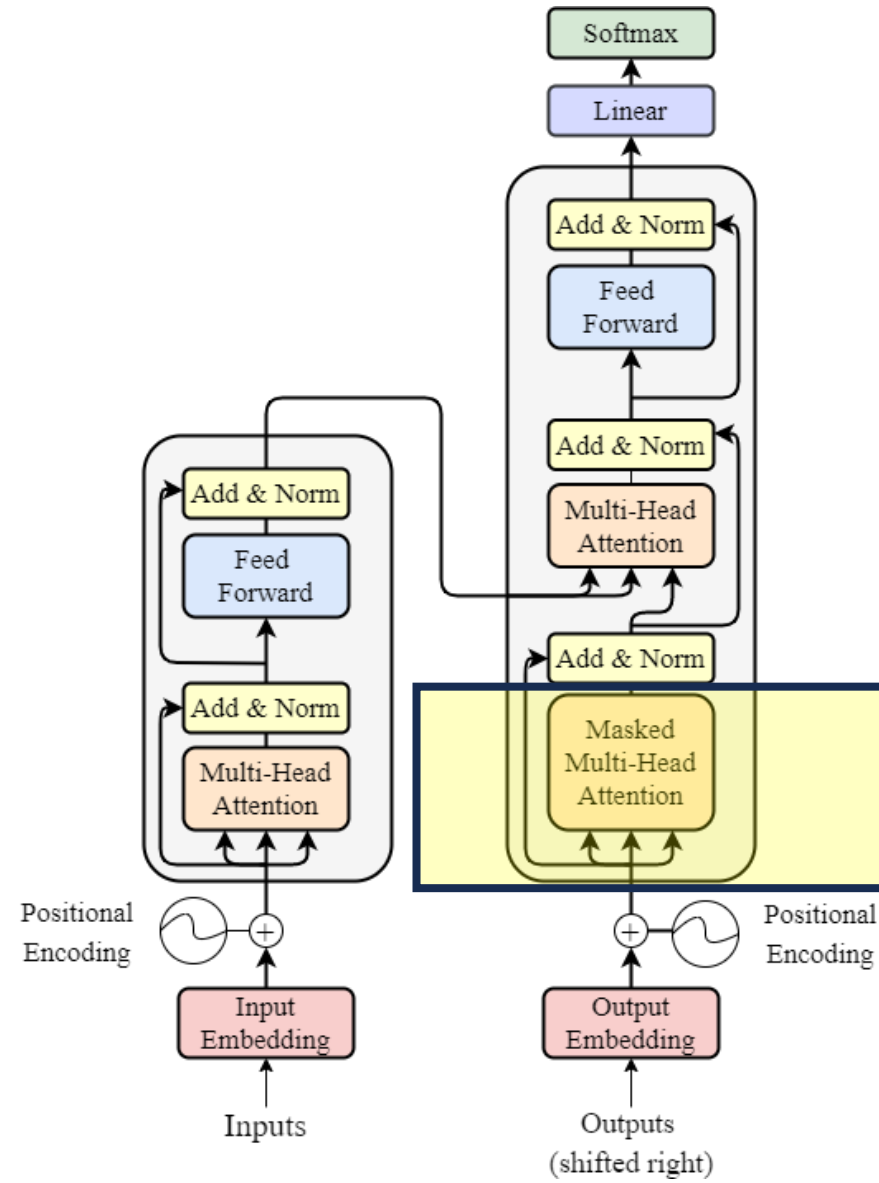


Query, Key & Value









Masked Attention

softmax {
$$\frac{Q * K^T}{\sqrt{512}}$$
 }

	I	can	go	alone
I	0.7	0.2	0.1	0.1
can	0.3	0.5	0.1	0.1
go	0.1	0.3	0.4	0.2
alone	0.05	0.05	0.1	0.8

4*4

Masked Attention

Causal Model: The model must not be able to see the **future** words

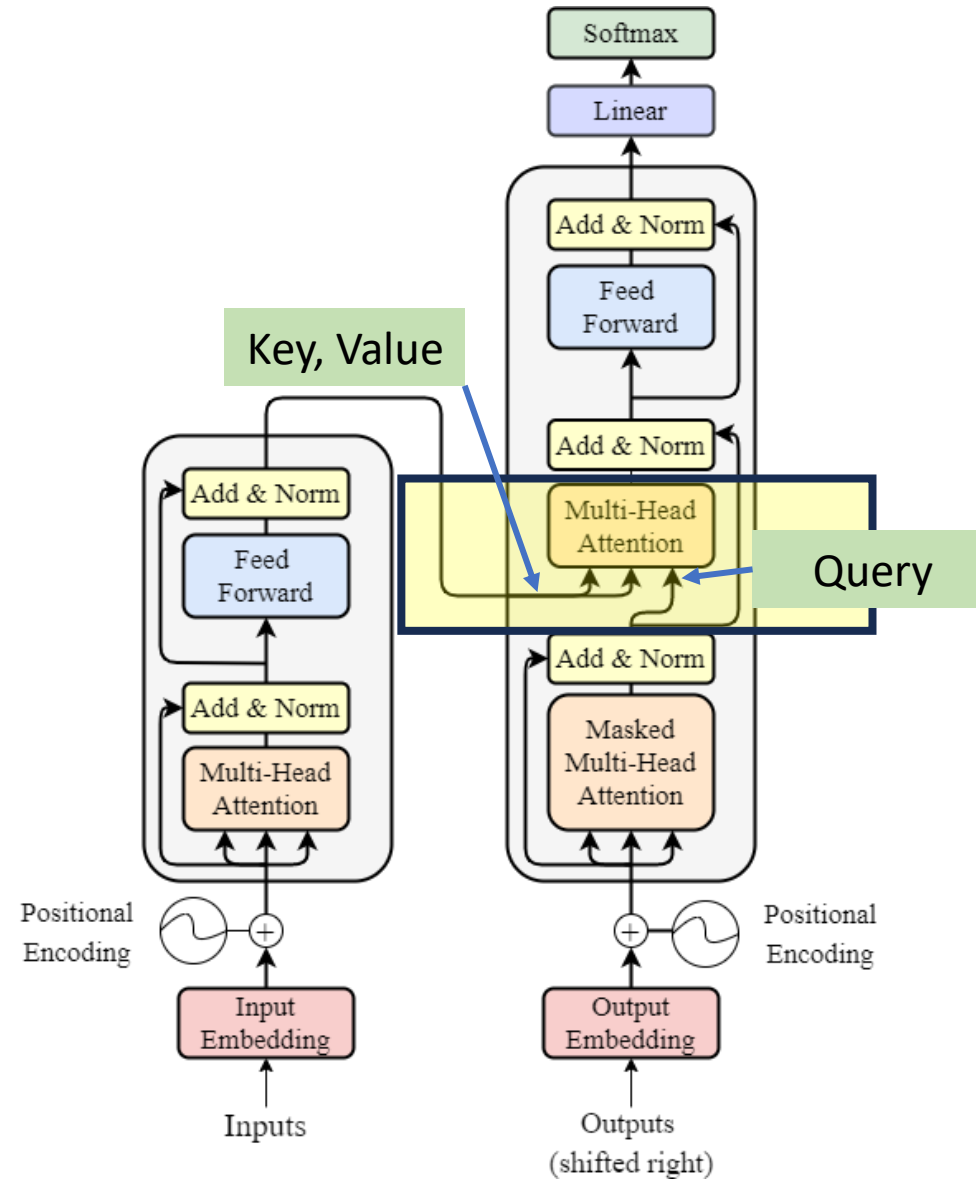
$$\text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)$$

$$\left(\frac{QK^T}{\sqrt{d_k}}\right)$$

	I	can	go	alone
I	63.3	1.2	2.6	7.2
can	3.25	0.3	1.2	2.1
go	12.0	11.9	52.9	2.9
alone	1.6	63.1	14.2	101.3



	I	can	go	alone
I	63.3	$-\infty$	$-\infty$	$-\infty$
can	3.25	96.1	$-\infty$	$-\infty$
go	12.0	11.9	52.9	$-\infty$
alone	1.6	63.1	14.2	101.3



Training

I	can	go	alone
---	-----	----	-------



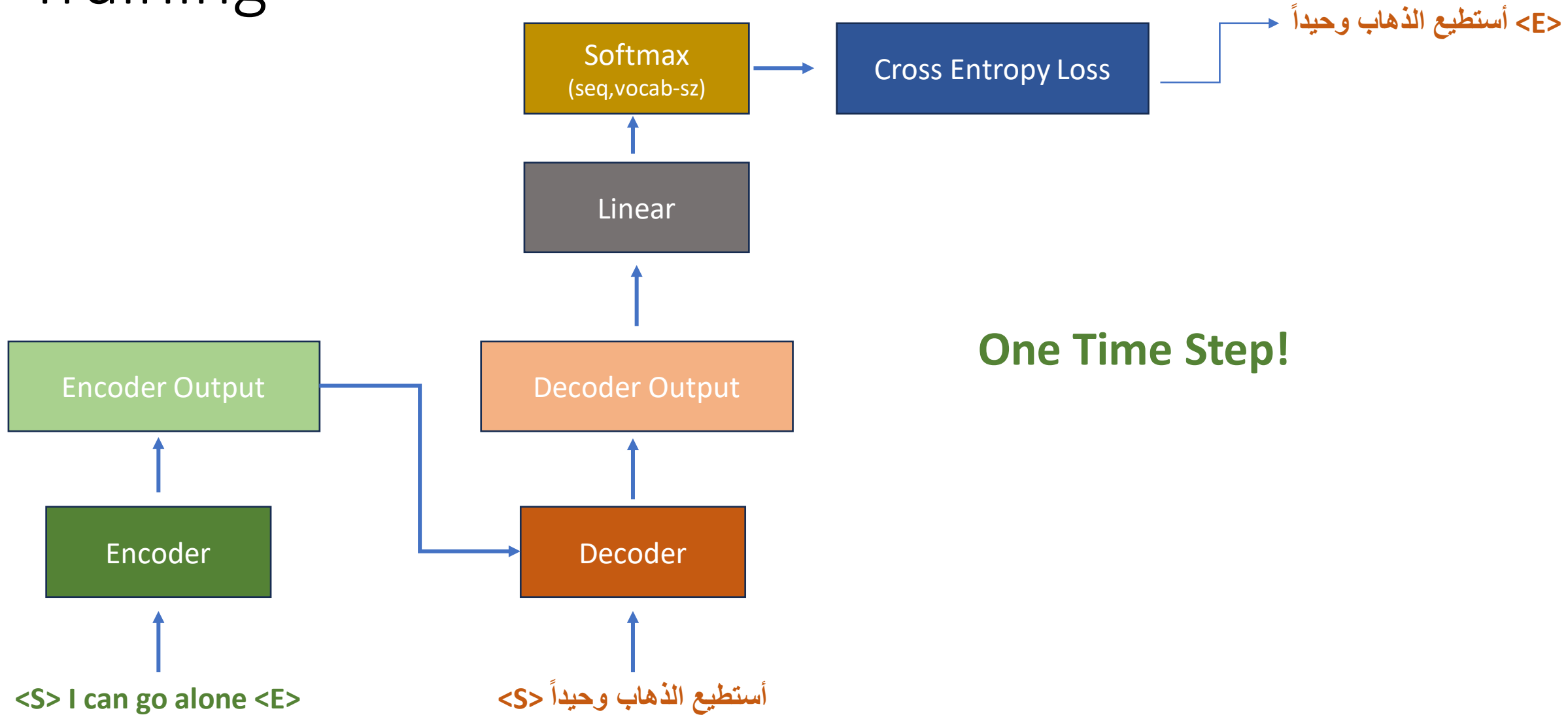
أستطيع	الذهاب	وحيداً
--------	--------	--------

<S>	I	can	go	alone	<E>
-----	---	-----	----	-------	-----



<S>	أستطيع	الذهاب	وحيداً	<E>
-----	--------	--------	--------	-----

Training



Inference

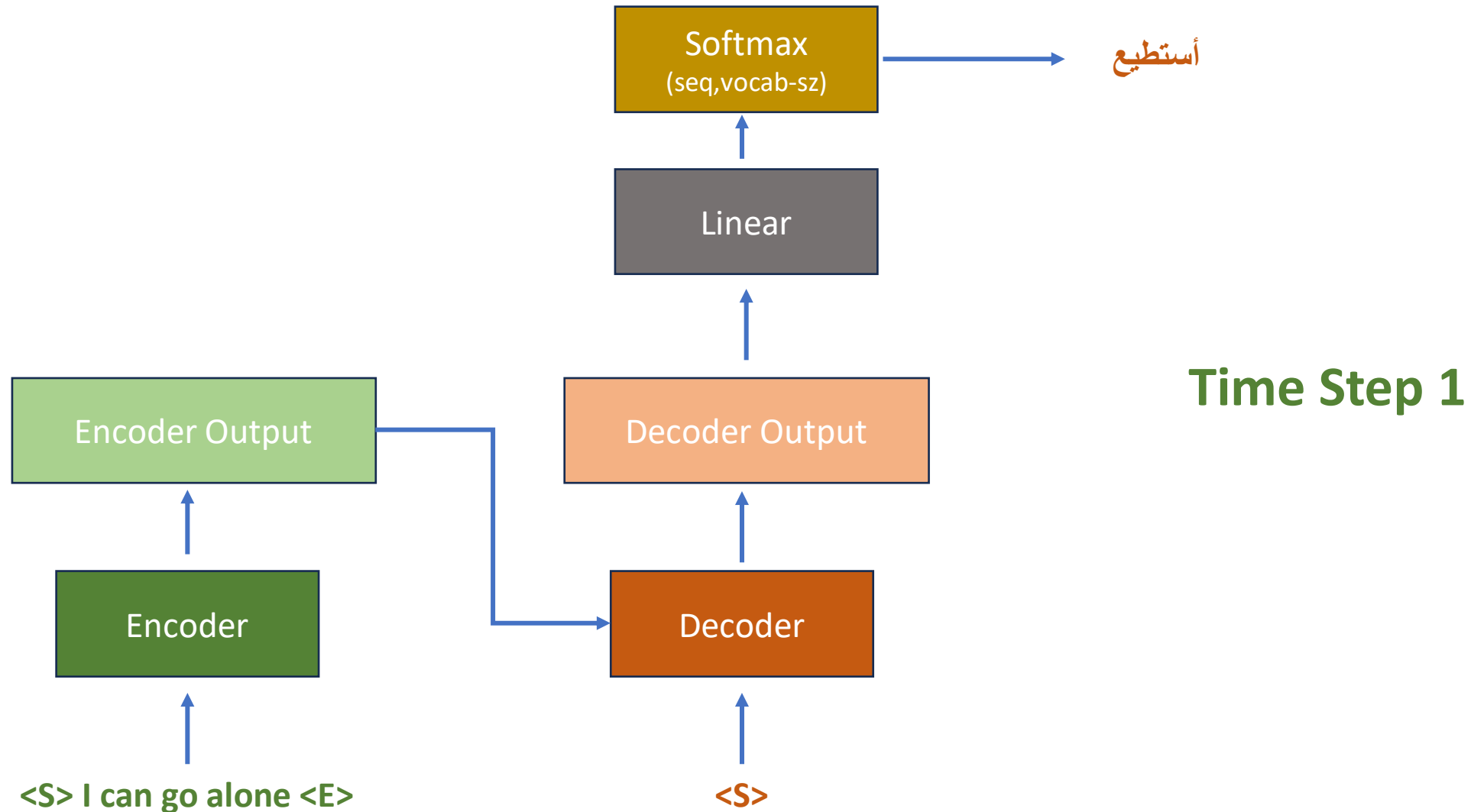
I	can	go	alone
---	-----	----	-------

<S>	I	can	go	alone	<E>
-----	---	-----	----	-------	-----

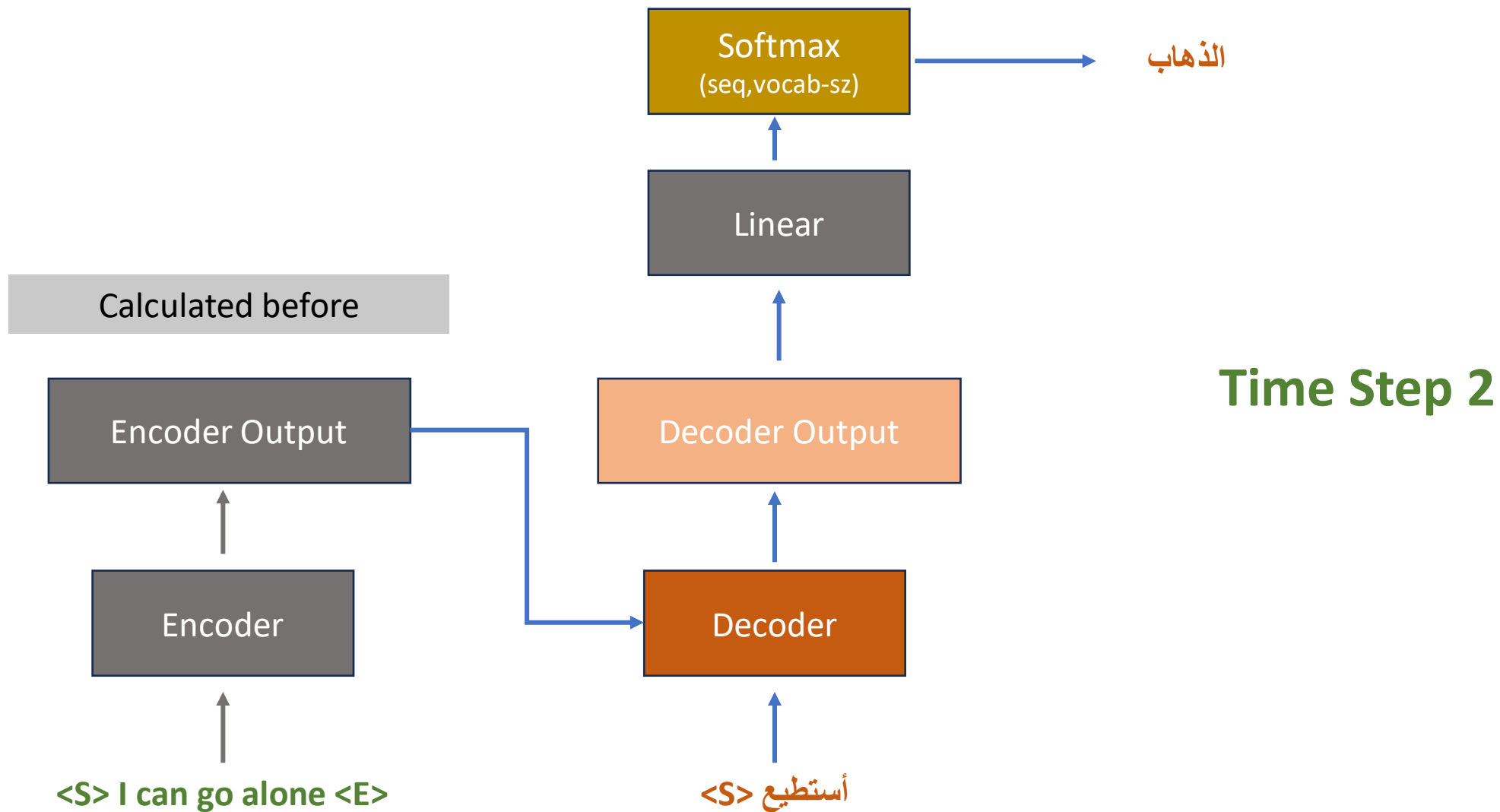


<S>

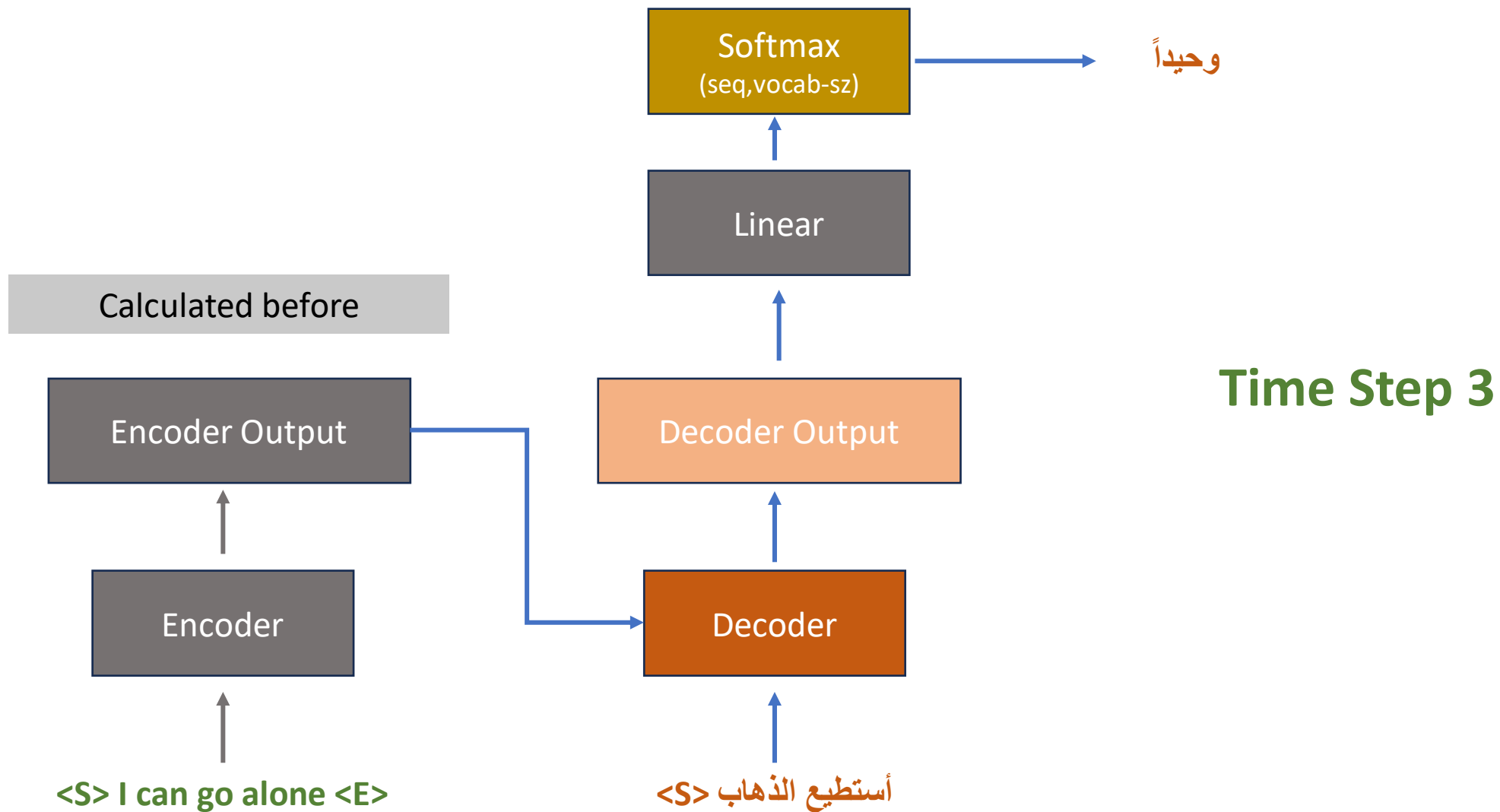
Inference



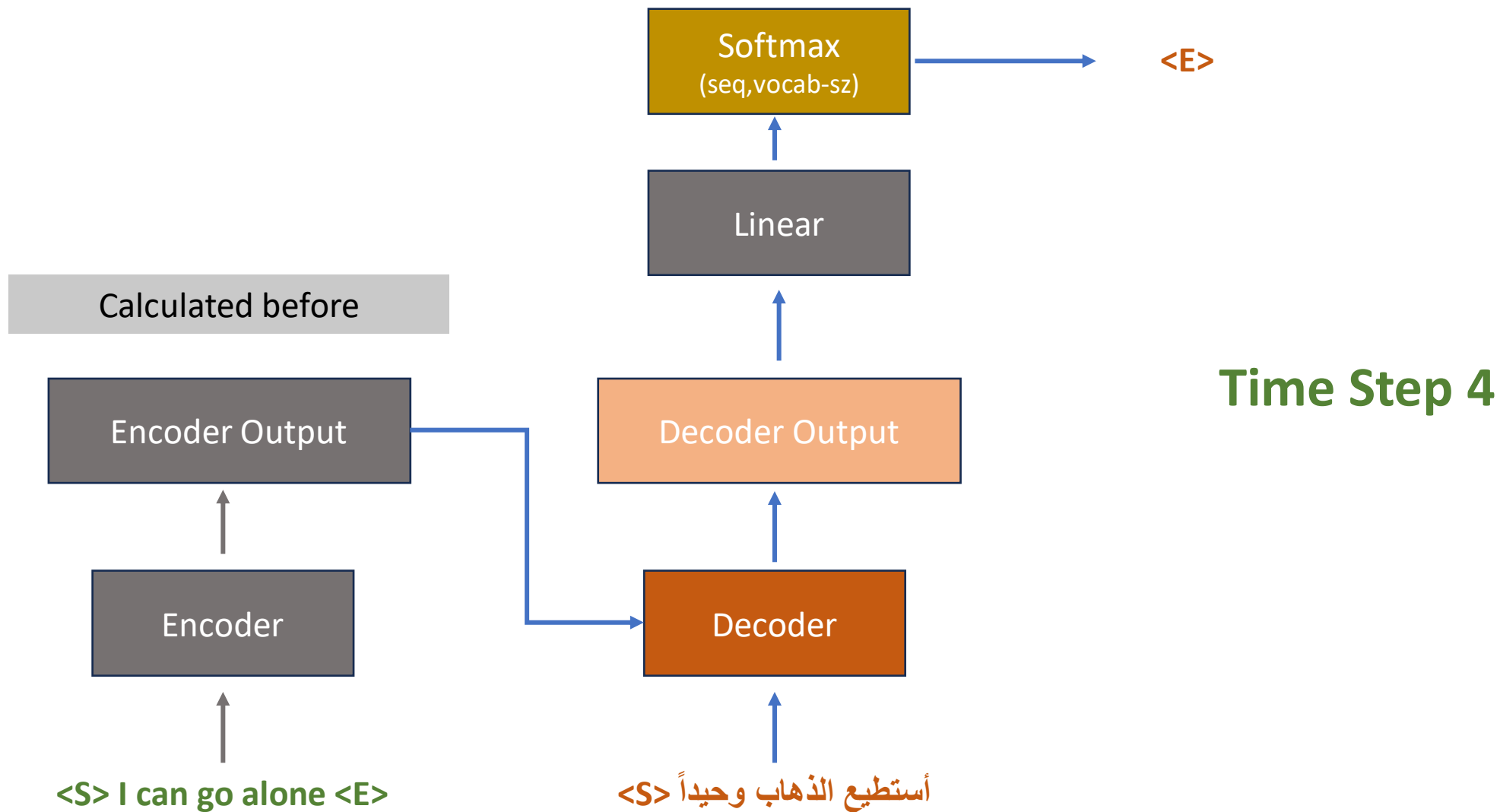
Inference



Inference



Inference



How to generate text: using different decoding methods for language generation with Transformers

Published March 1, 2020

[Update on GitHub](#)



[patrickvonplaten](#)
[Patrick von Platen](#)

 [Open in Colab](#)



Follow Me



Abu Bakr Soliman, MSc

Developing REAL AI solutions and strategies. Follow me to know more.

<https://www.linkedin.com/in/bakrianoo/>