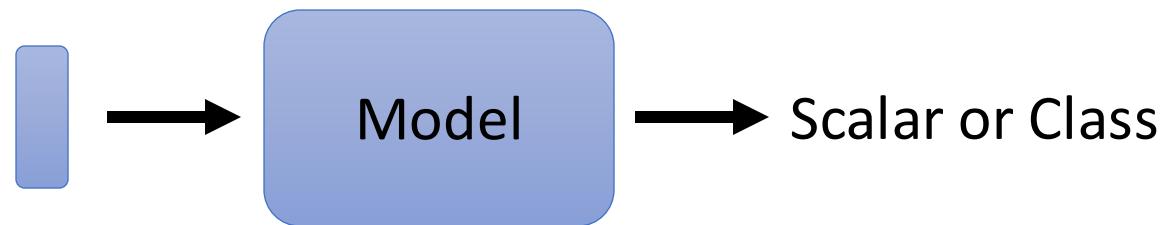


Transformer

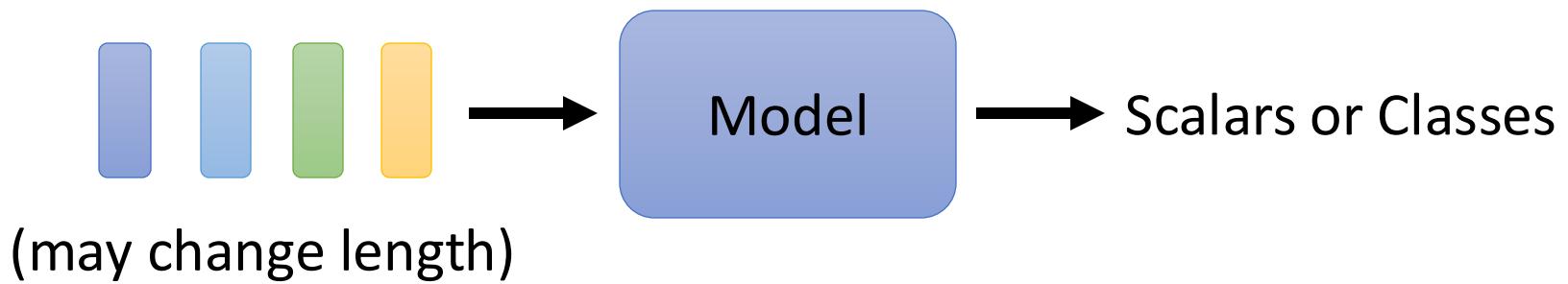
Slides credit: Hung-Yi Lee

Sophisticated Input

- Input is a **vector**

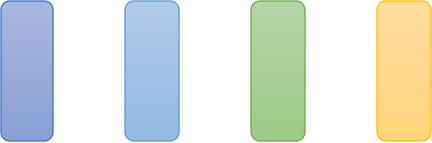


- Input is a **set of vectors**



Vector Set as Input

this is a cat



The word "this" is represented by a blue rectangle.
The word "is" is represented by a light blue rectangle.
The word "a" is represented by a green rectangle.
The word "cat" is represented by an orange rectangle.

One-hot Encoding

apple = [1 0 0 0 0]

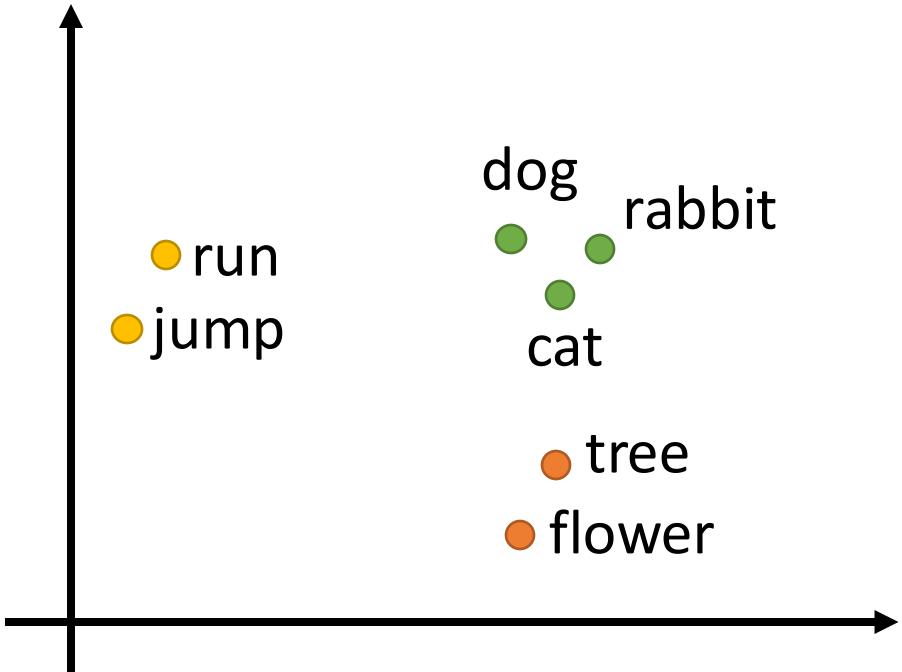
bag = [0 1 0 0 0]

cat = [0 0 1 0 0]

dog = [0 0 0 1 0]

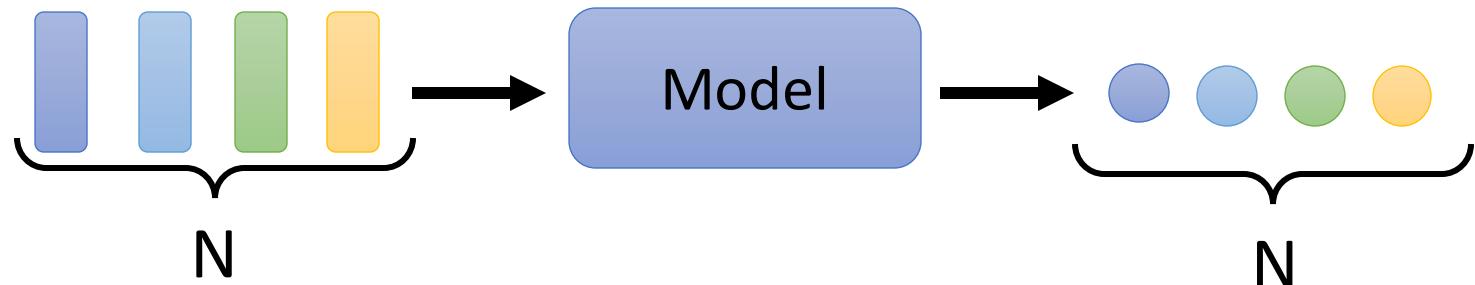
elephant = [0 0 0 0 1]

Word Embedding

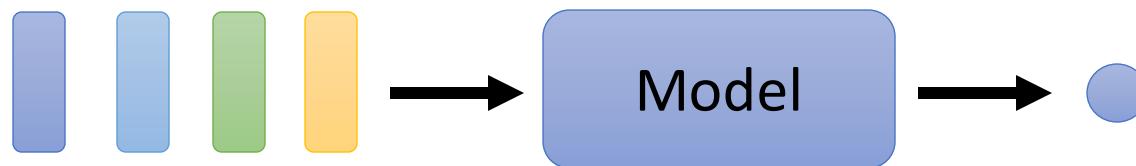


What is the output?

- Each vector has a label.

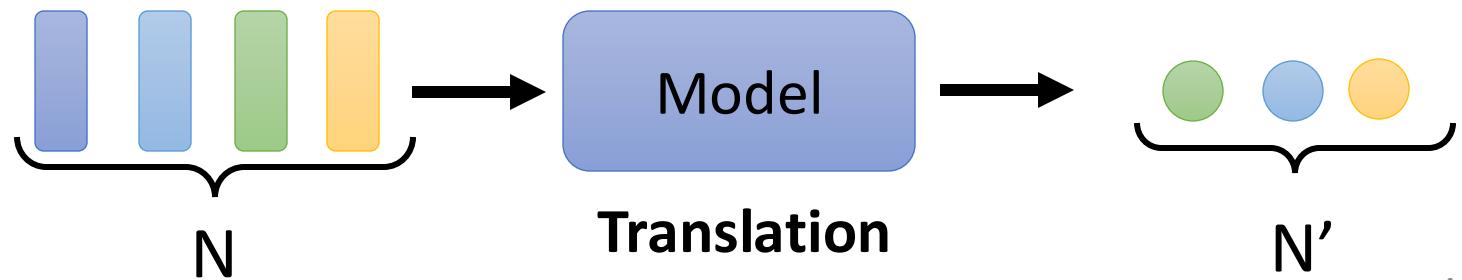


- The whole sequence has a label.



- Model decides the number of labels itself.

seq2seq



Sequence Labeling

Is it possible to consider the context?

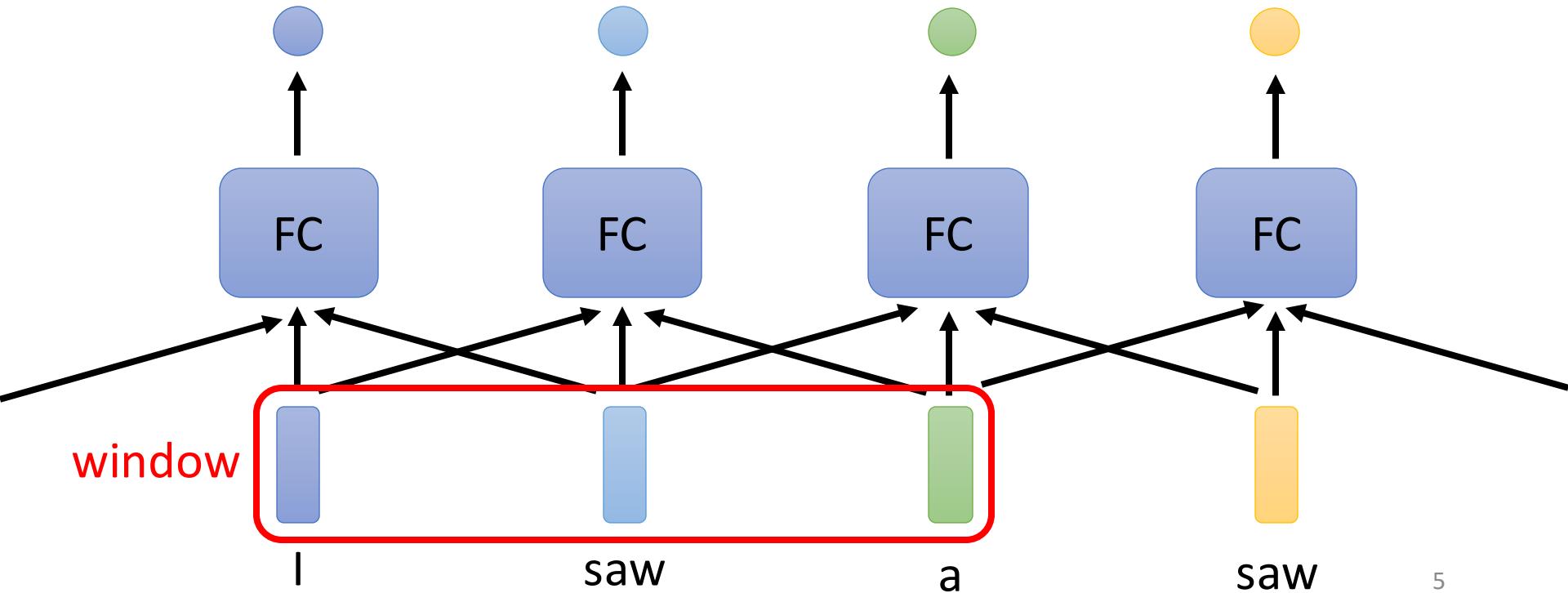


Fully-connected

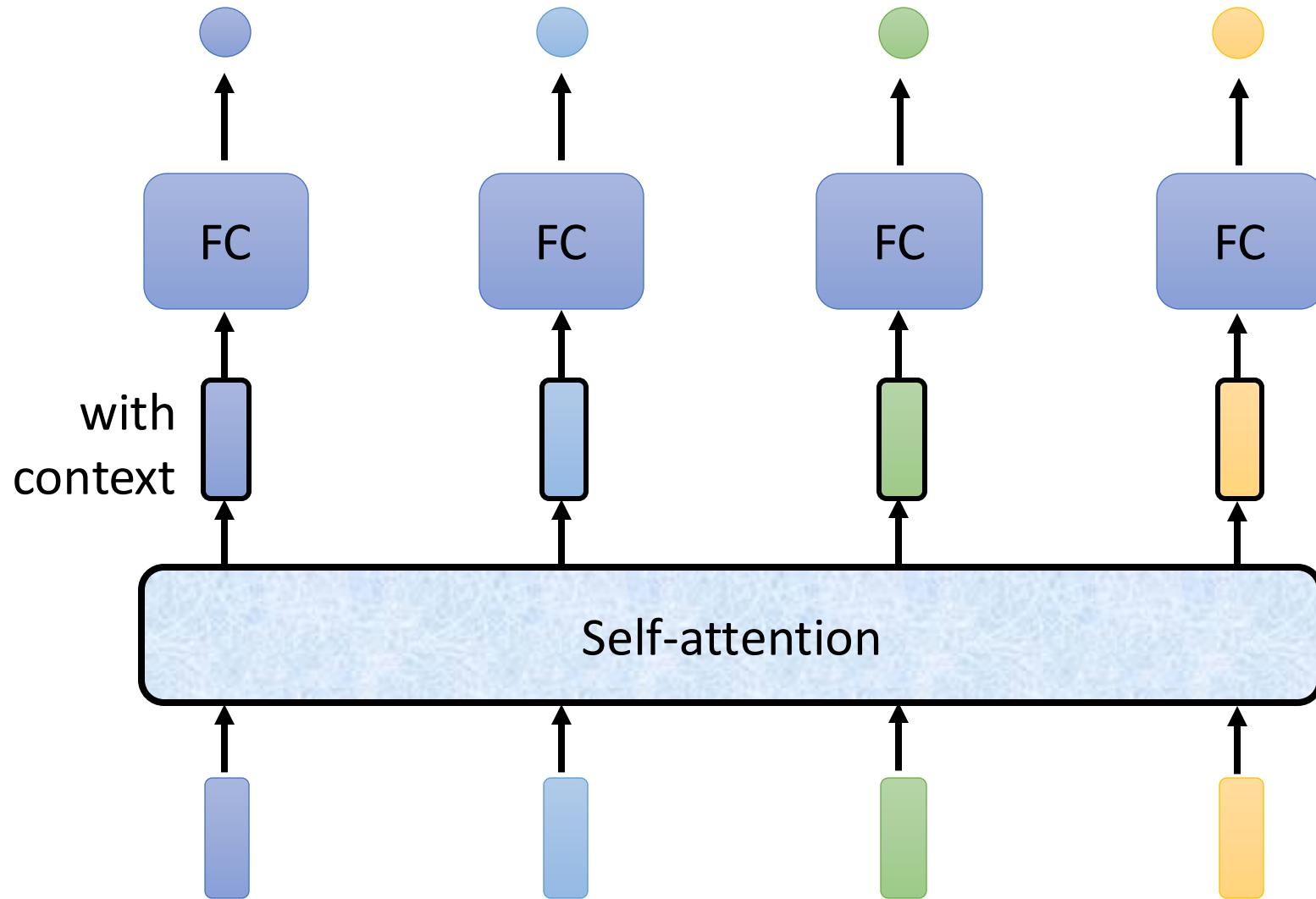
FC can consider the neighbor

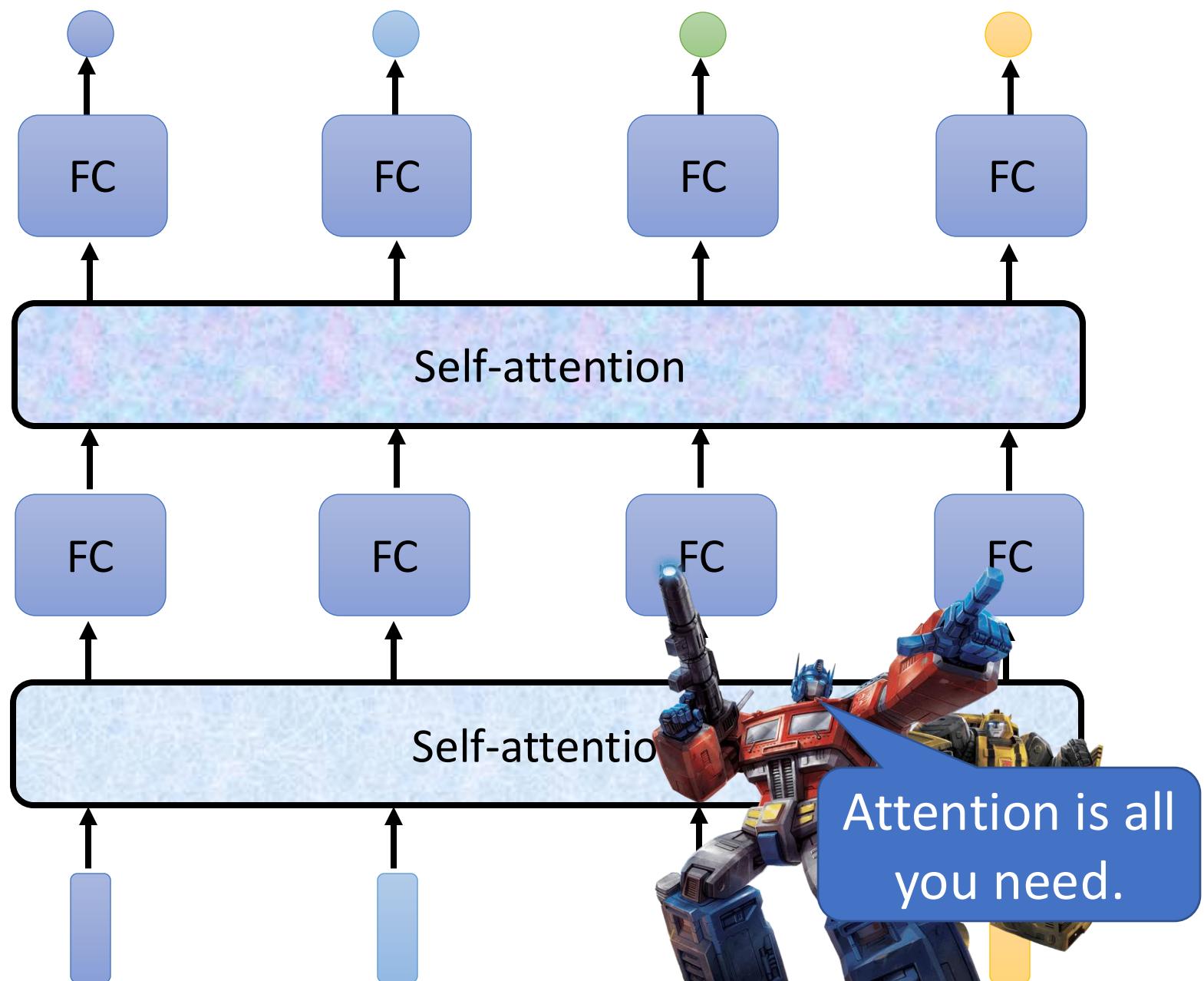
How to consider the whole sequence?

a window covers the whole sequence?



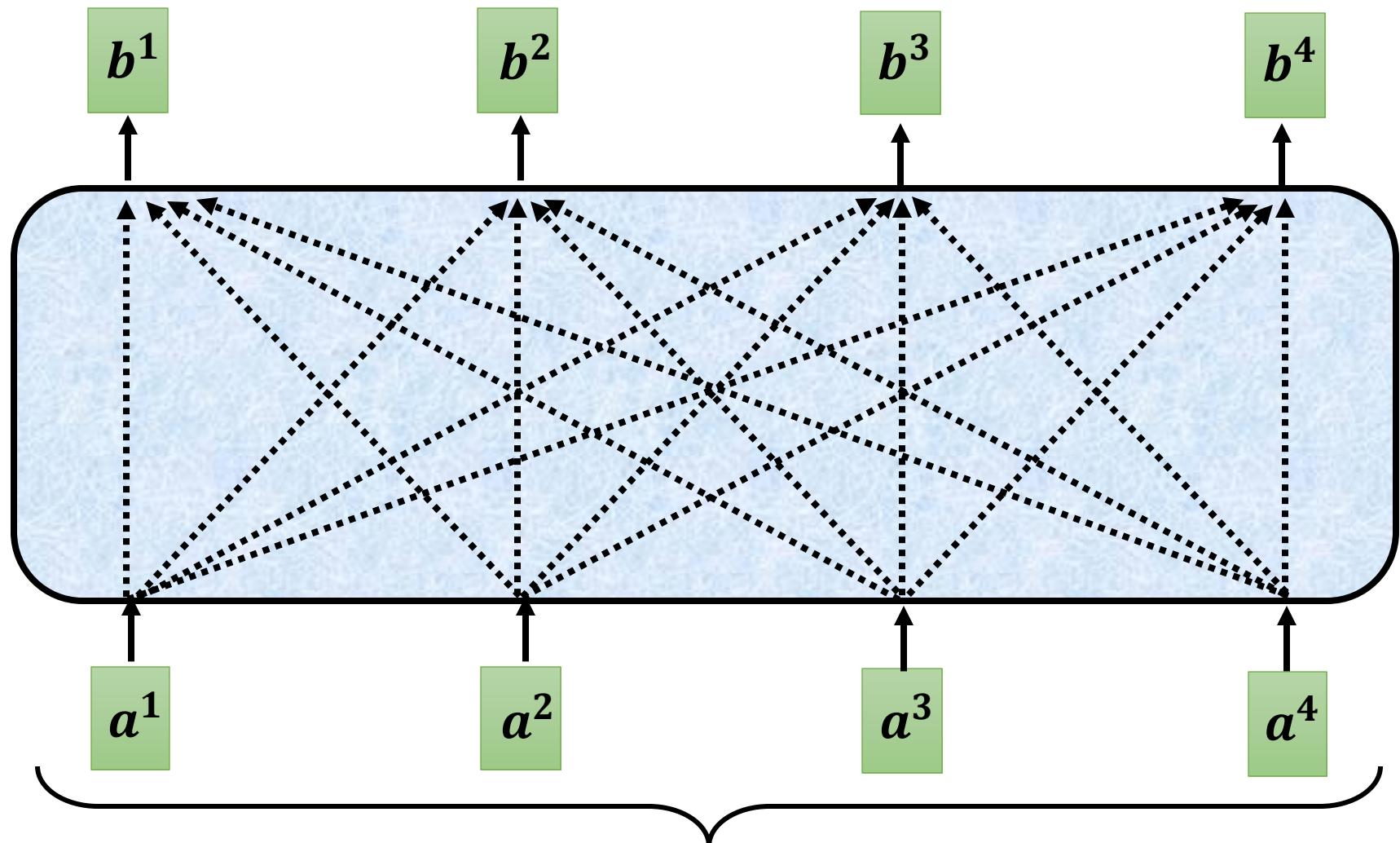
Self-attention





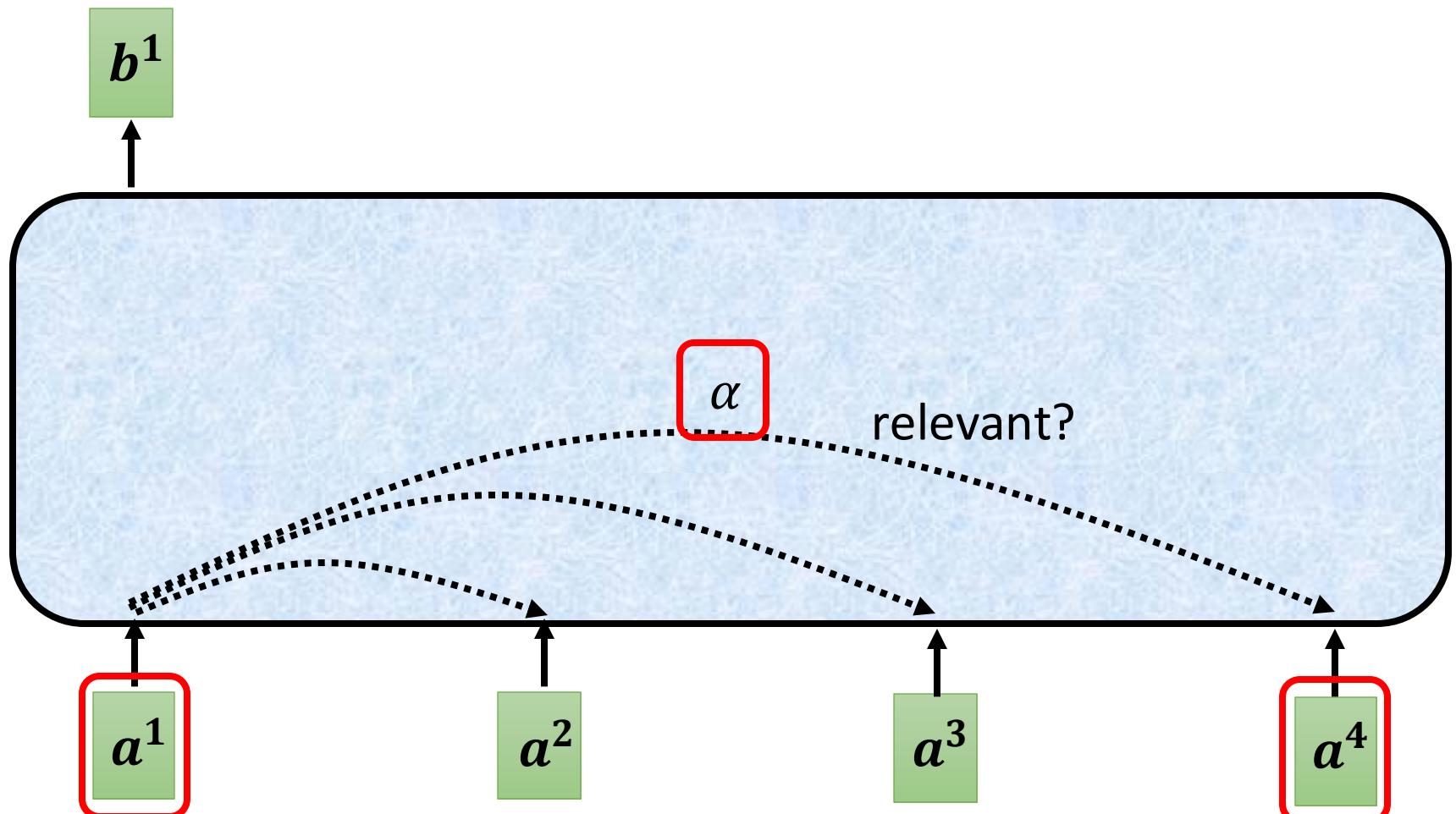
<https://arxiv.org/abs/1706.03762>

Self-attention



Can be either **input** or a **hidden layer**

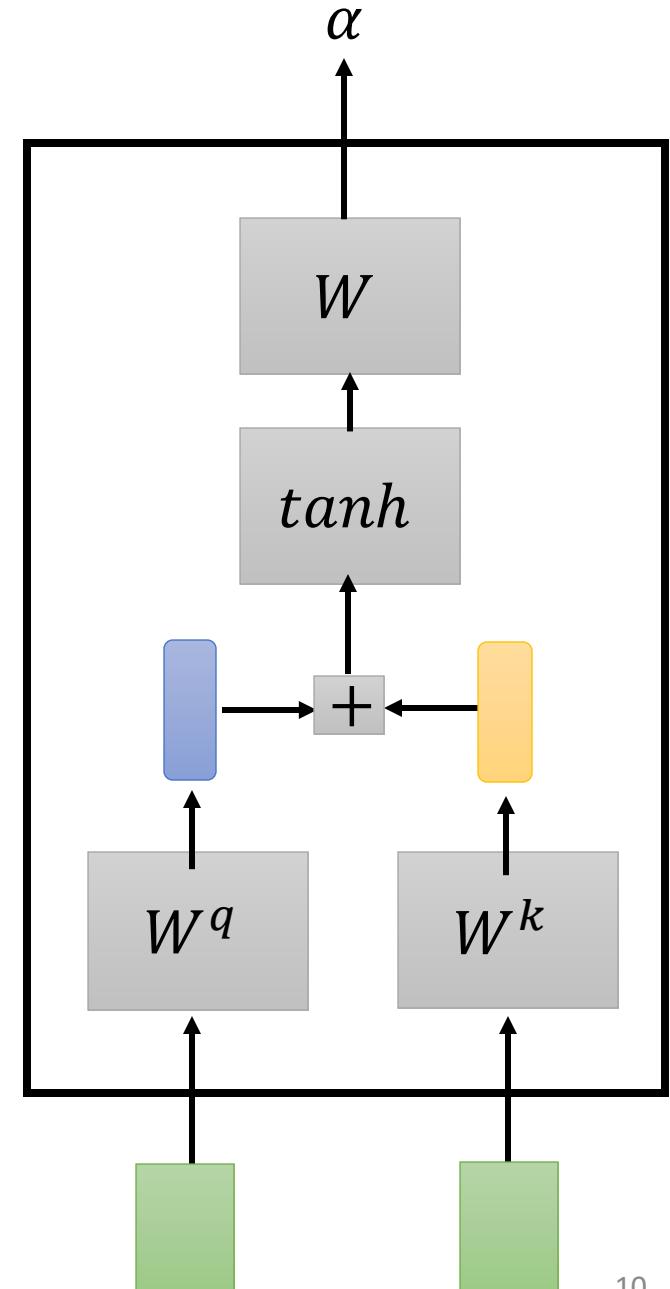
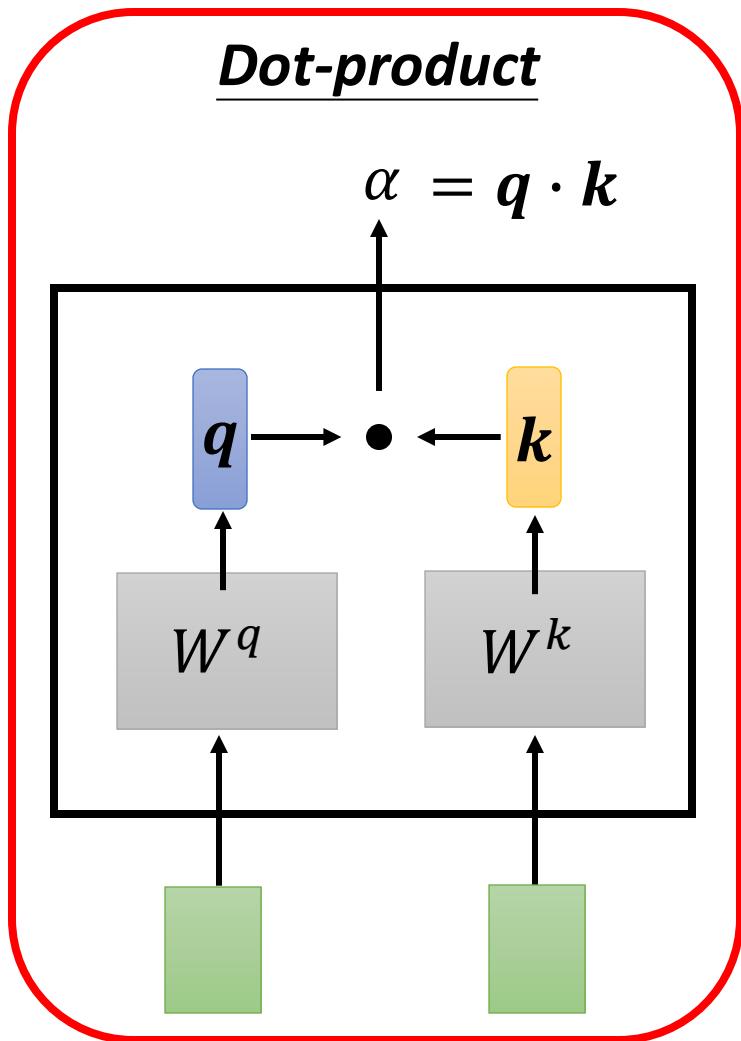
Self-attention



Find the relevant vectors in a sequence

Self-attention

Additive

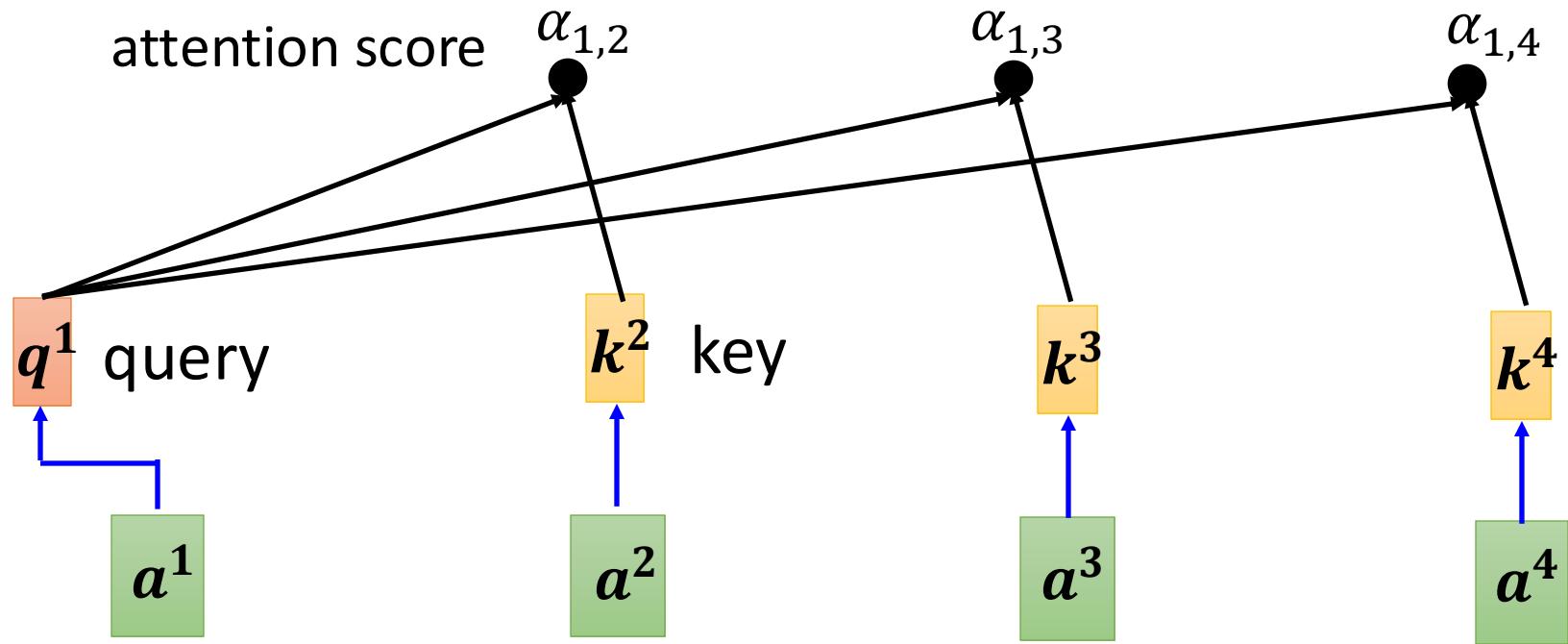


Self-attention

$$\alpha_{1,2} = q^1 \cdot k^2$$

$$\alpha_{1,3} = q^1 \cdot k^3$$

$$\alpha_{1,4} = q^1 \cdot k^4$$



$$q^1 = W^q a^1$$

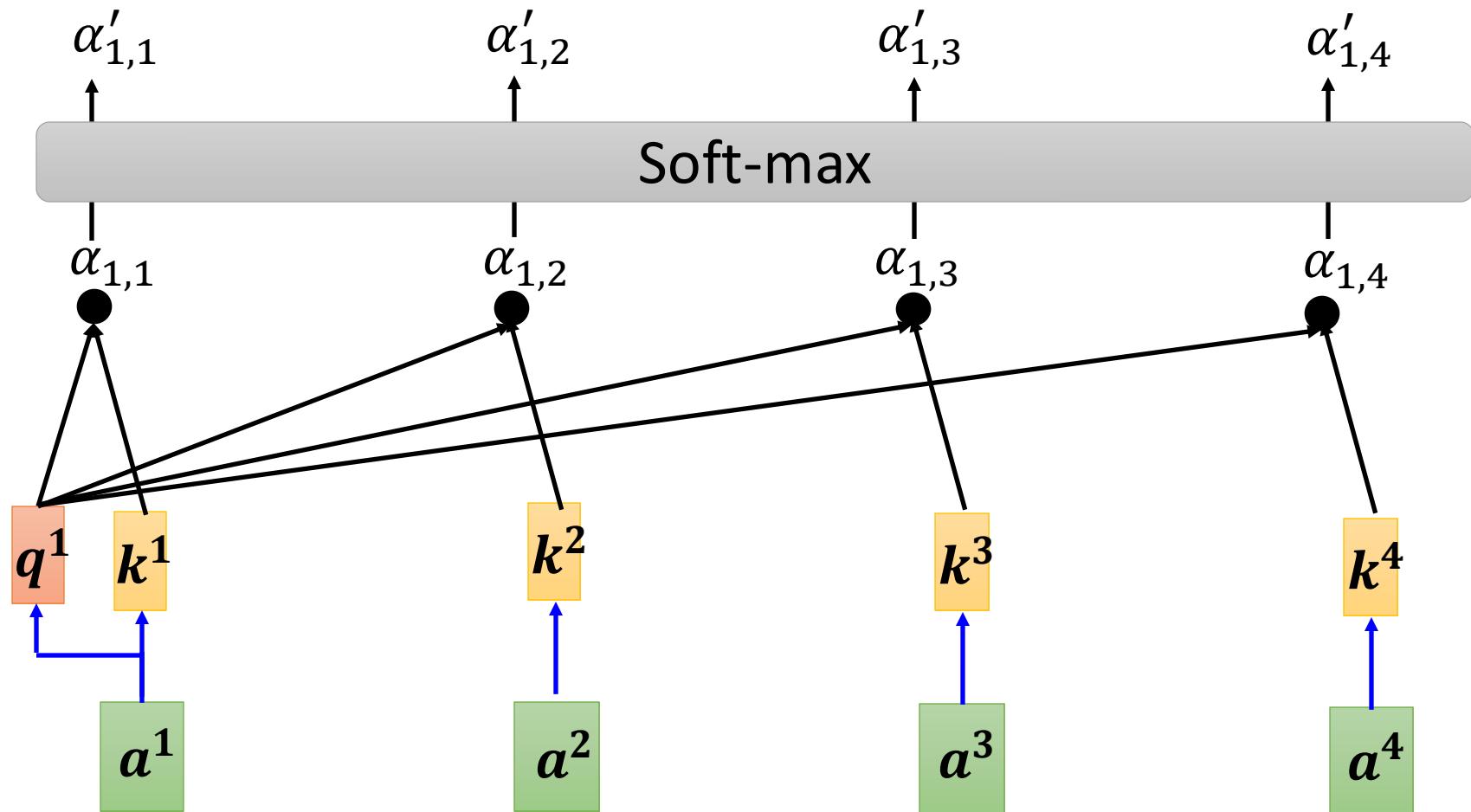
$$k^2 = W^k a^2$$

$$k^3 = W^k a^3$$

$$k^4 = W^k a^4$$

Self-attention

$$\alpha'_{1,i} = \exp(\alpha_{1,i}) / \sum_j \exp(\alpha_{1,j})$$



$$q^1 = W^q a^1$$

$$k^2 = W^k a^2$$

$$k^3 = W^k a^3$$

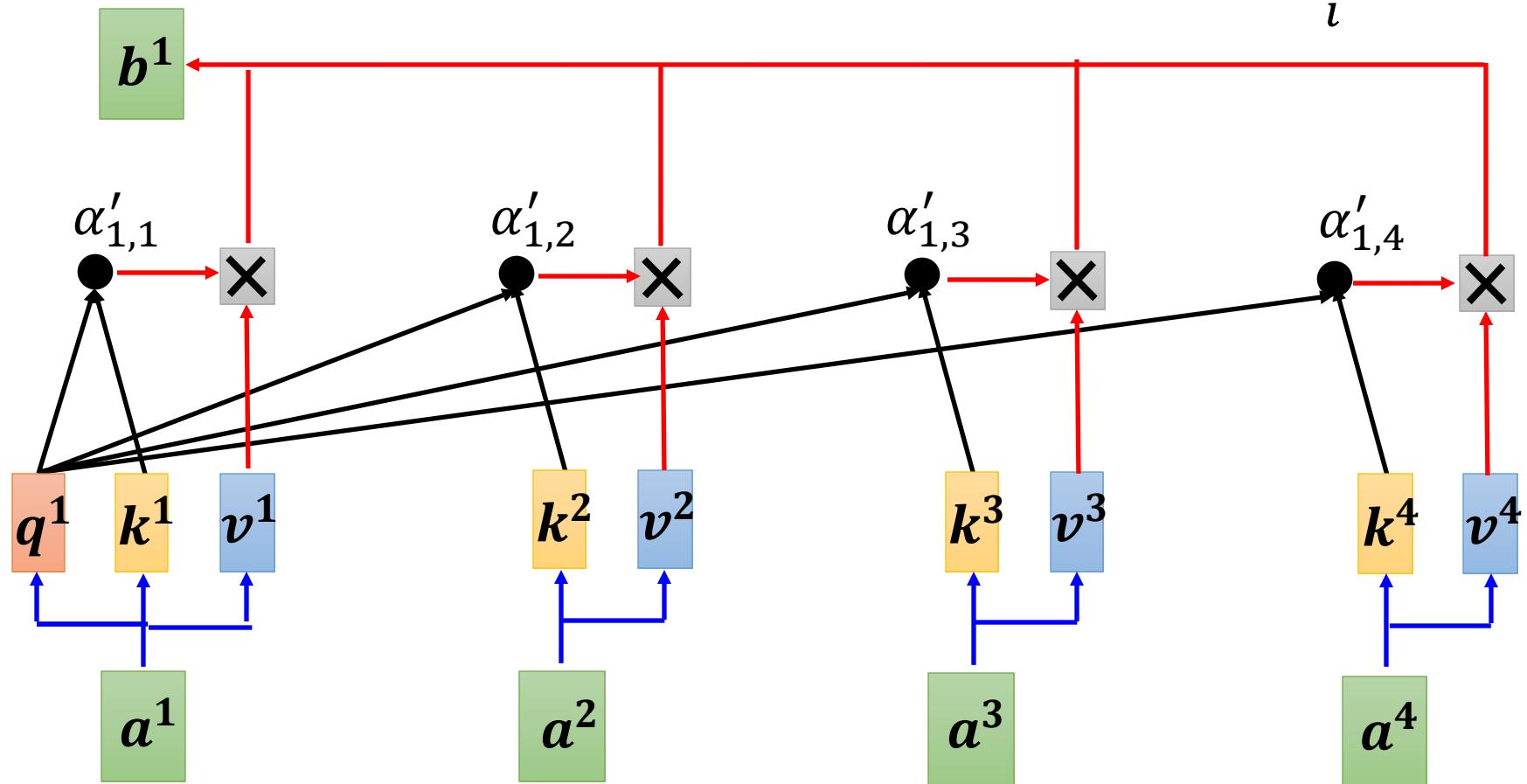
$$k^4 = W^k a^4$$

$$k^1 = W^k a^1$$

Self-attention

Extract information based
on attention scores

$$b^1 = \sum_i \alpha'_{1,i} v^i$$



$$v^1 = W^v a^1$$

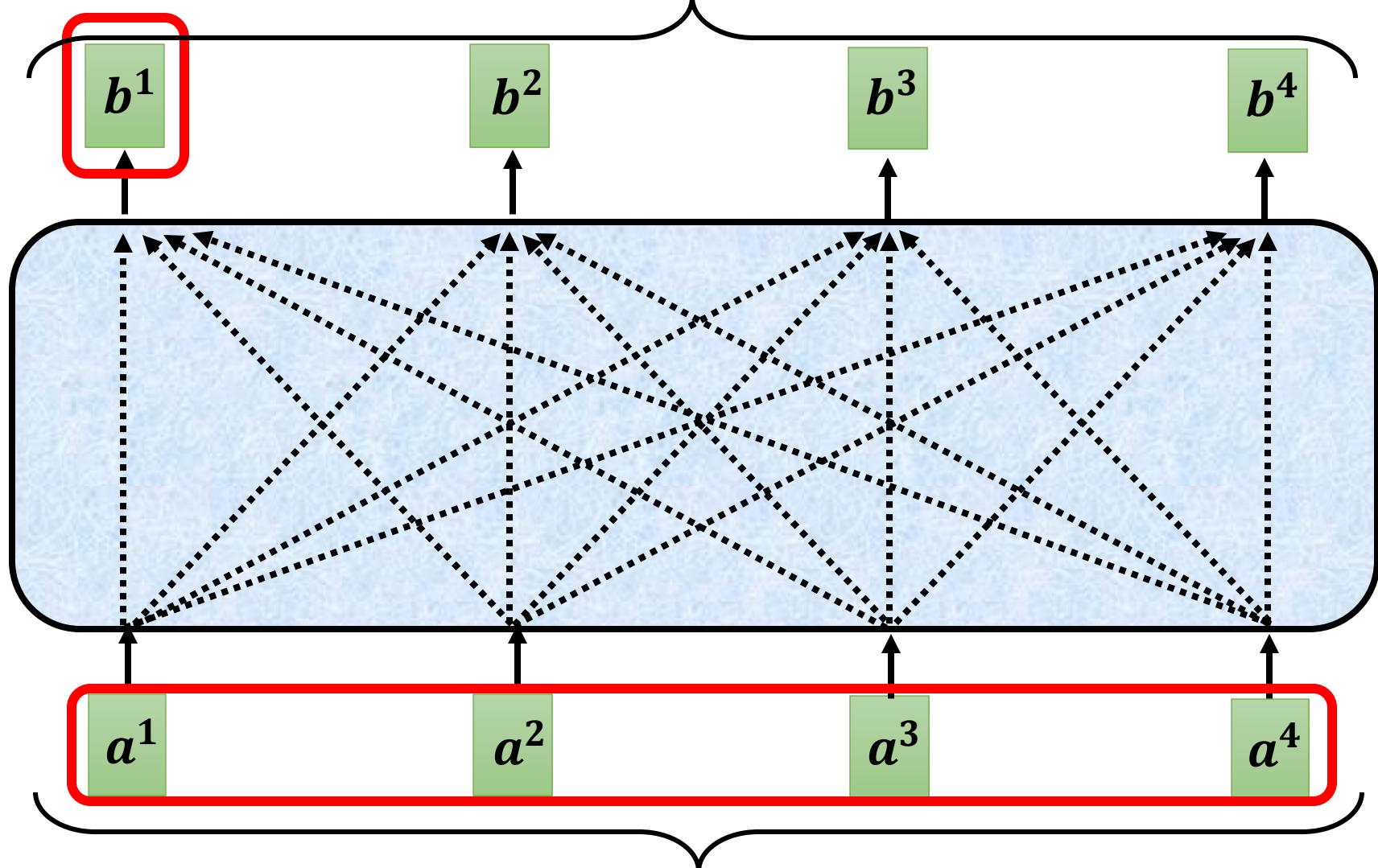
$$v^2 = W^v a^2$$

$$v^3 = W^v a^3$$

$$v^4 = W^v a^4$$

Self-attention

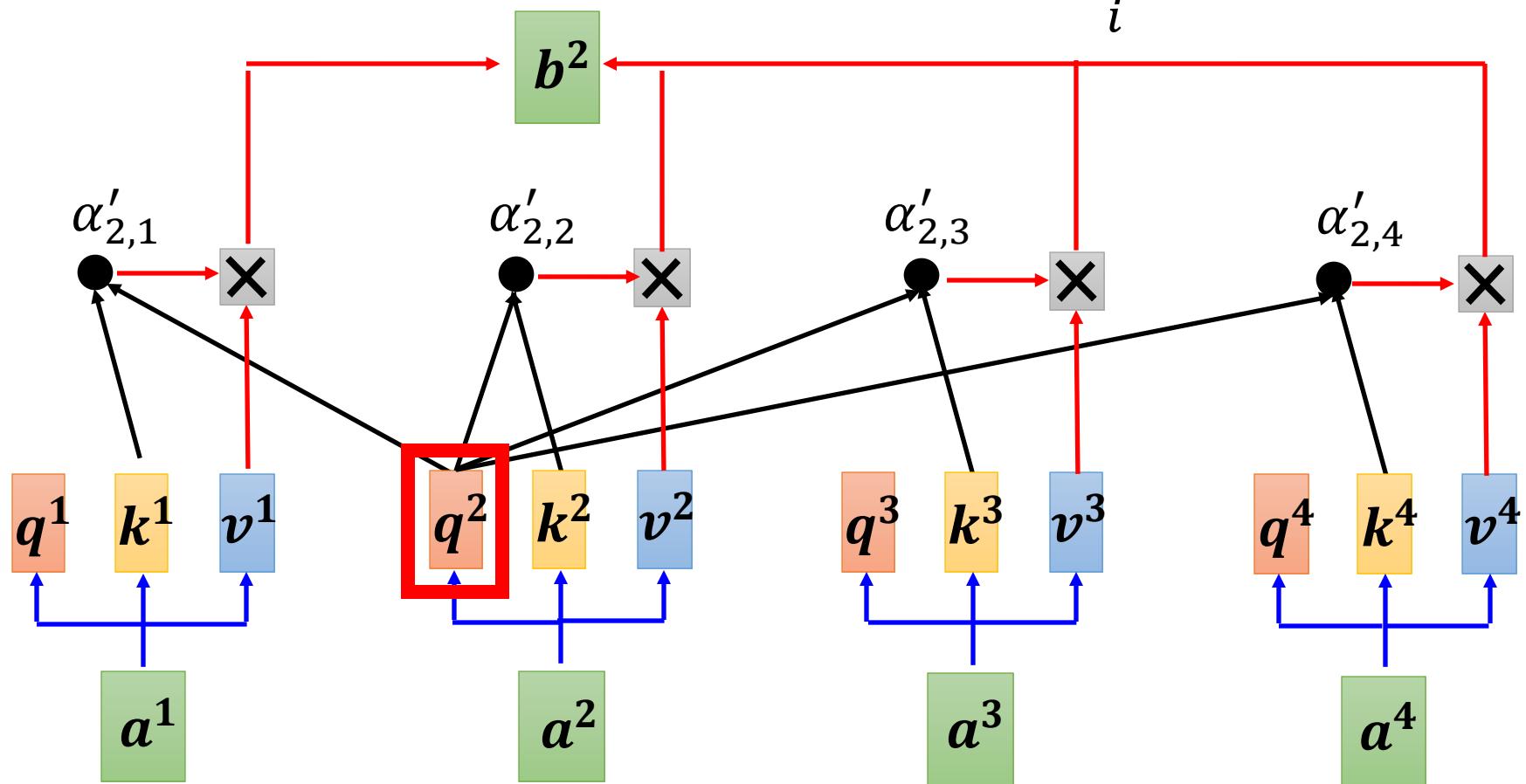
parallel



Can be either **input** or a **hidden layer**

Self-attention

$$b^2 = \sum_i \alpha'_{2,i} v^i$$



Self-attention

$$q^i = W^q a^i$$

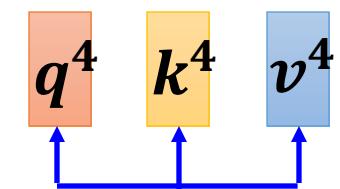
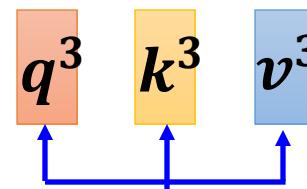
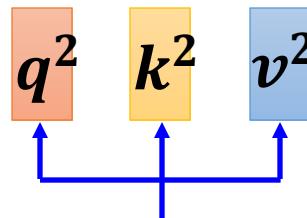
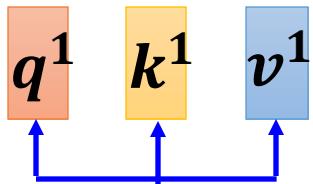
$$\begin{matrix} q^1 & q^2 & q^3 & q^4 \end{matrix} = \begin{matrix} W^q \\ Q \end{matrix} \quad \begin{matrix} a^1 & a^2 & a^3 & a^4 \end{matrix} = \begin{matrix} W^q \\ I \end{matrix}$$

$$k^i = W^k a^i$$

$$\begin{matrix} k^1 & k^2 & k^3 & k^4 \end{matrix} = \begin{matrix} W^k \\ K \end{matrix} \quad \begin{matrix} a^1 & a^2 & a^3 & a^4 \end{matrix} = \begin{matrix} W^k \\ I \end{matrix}$$

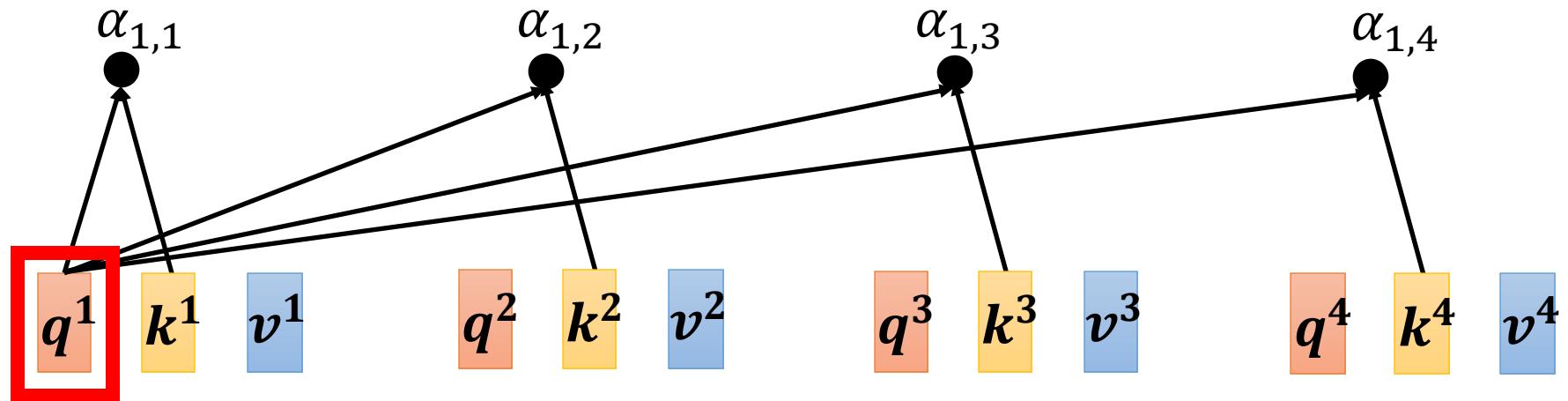
$$v^i = W^v a^i$$

$$\begin{matrix} v^1 & v^2 & v^3 & v^4 \end{matrix} = \begin{matrix} W^v \\ V \end{matrix} \quad \begin{matrix} a^1 & a^2 & a^3 & a^4 \end{matrix} = \begin{matrix} W^v \\ I \end{matrix}$$



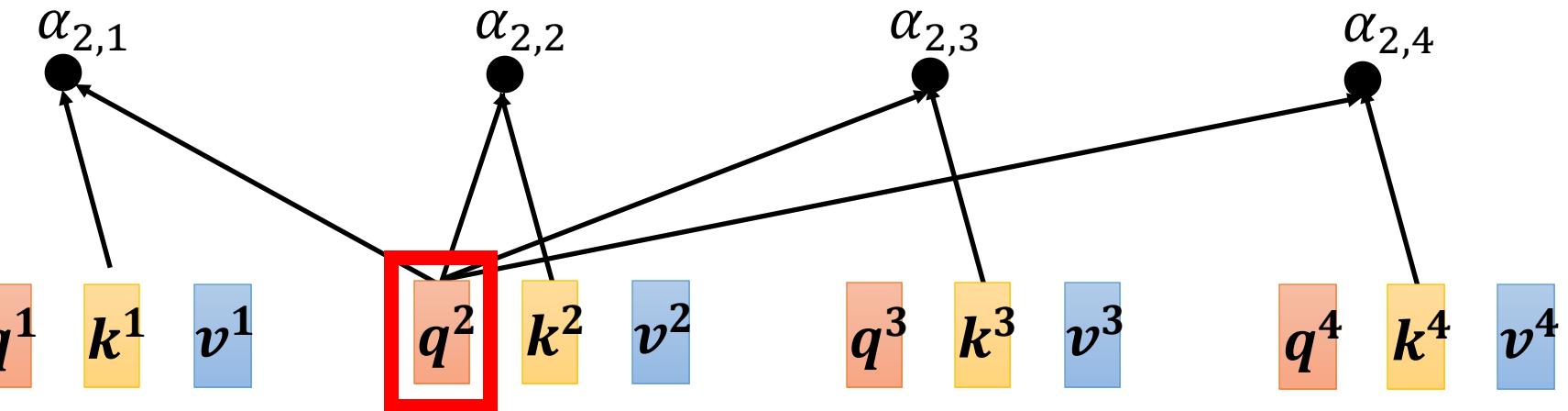
Self-attention

$$\begin{array}{ll} \alpha_{1,1} = \boxed{k^1} \quad \boxed{q^1} & \alpha_{1,2} = \boxed{k^2} \quad \boxed{q^1} \\ \alpha_{1,3} = \boxed{k^3} \quad \boxed{q^1} & \alpha_{1,4} = \boxed{k^4} \quad \boxed{q^1} \end{array} \quad \begin{array}{l} \alpha_{1,1} \\ \alpha_{1,2} \\ \alpha_{1,3} \\ \alpha_{1,4} \end{array} = \begin{array}{c} \boxed{k^1} \\ \boxed{k^2} \\ \boxed{k^3} \\ \boxed{k^4} \end{array} \quad \boxed{q^1}$$



Self-attention

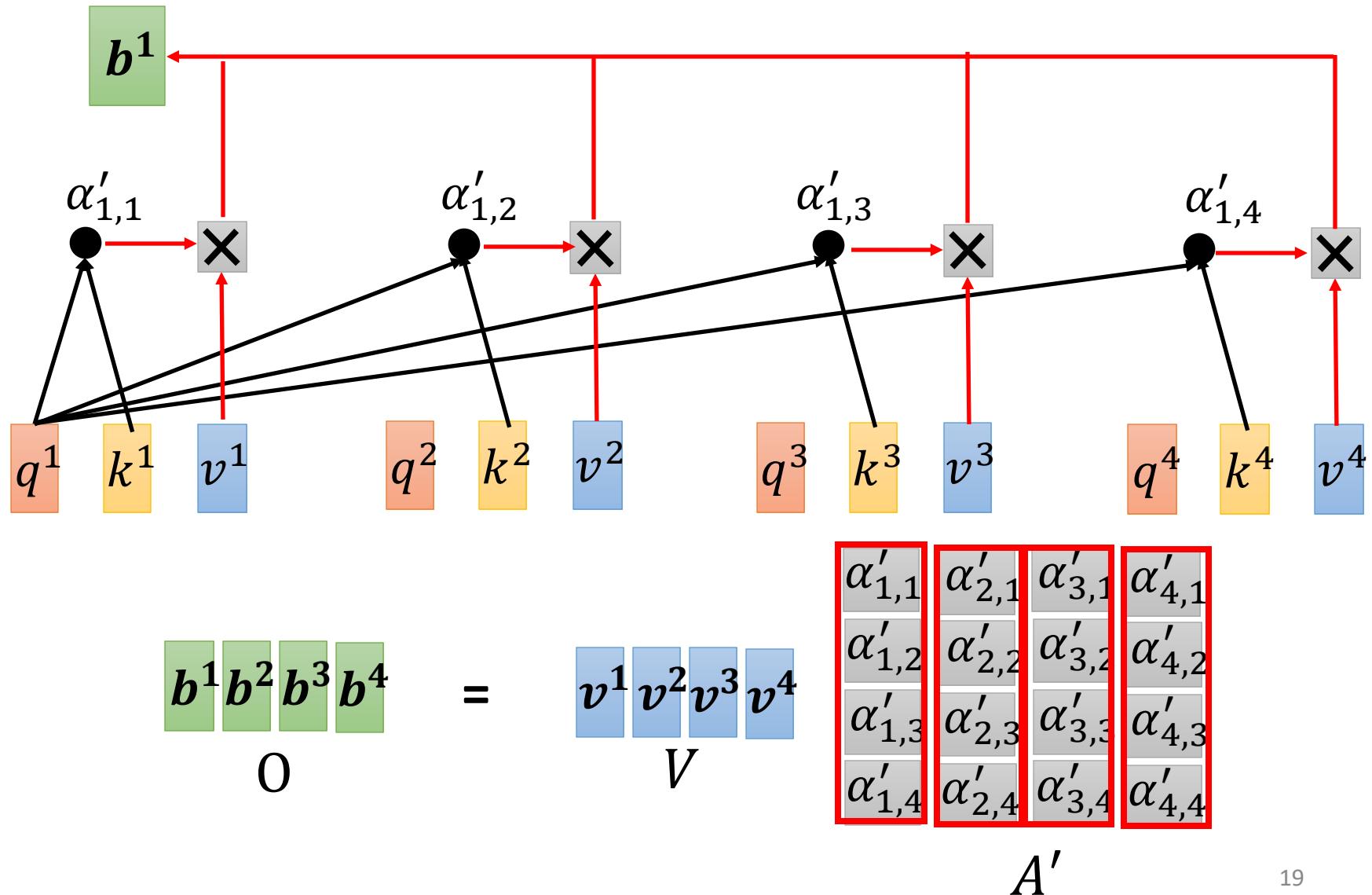
$$\begin{array}{ll} \alpha_{1,1} = k^1 q^1 & \alpha_{1,2} = k^2 q^1 \\ & \alpha_{1,3} = k^3 q^1 \quad \alpha_{1,4} = k^4 q^1 \end{array} \quad \begin{array}{l} \alpha_{1,1} \\ \alpha_{1,2} \\ \alpha_{1,3} \\ \alpha_{1,4} \end{array} = \begin{array}{c} k^1 \\ k^2 \\ k^3 \\ k^4 \end{array} \quad q^1$$



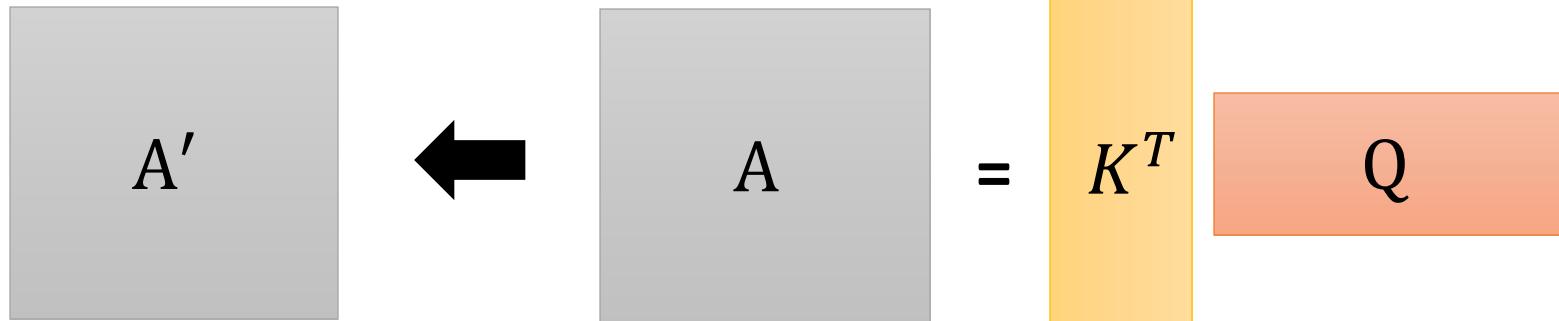
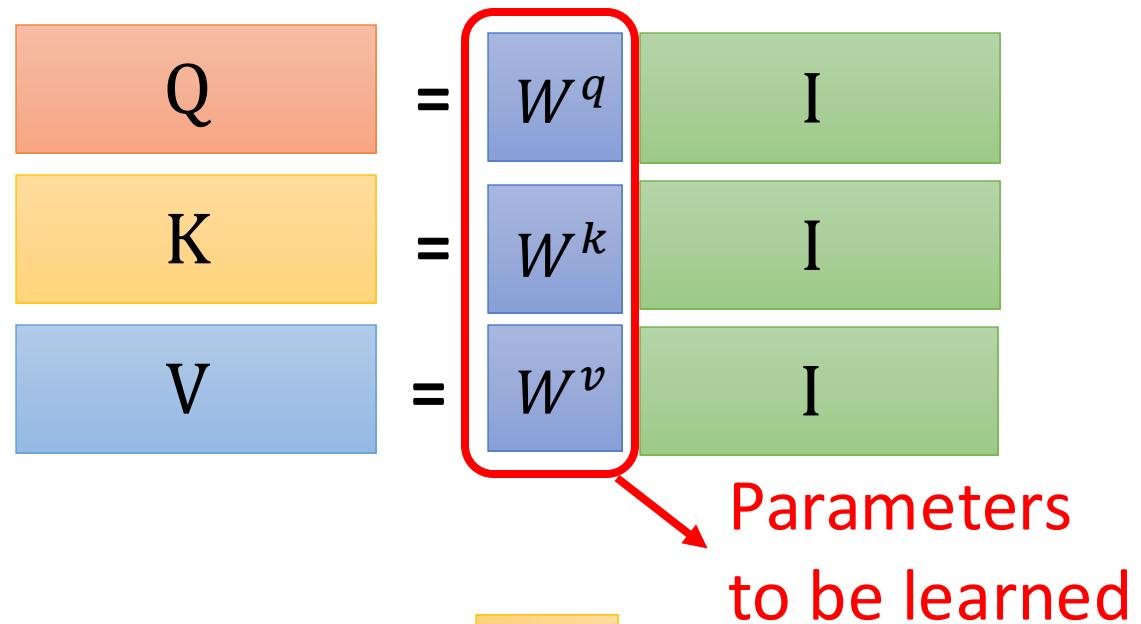
$$\begin{array}{cccc} \alpha'_{1,1} & \alpha'_{2,1} & \alpha'_{3,1} & \alpha'_{4,1} \\ \alpha'_{1,2} & \alpha'_{2,2} & \alpha'_{3,2} & \alpha'_{4,2} \\ \alpha'_{1,3} & \alpha'_{2,3} & \alpha'_{3,3} & \alpha'_{4,3} \\ \alpha'_{1,4} & \alpha'_{2,4} & \alpha'_{3,4} & \alpha'_{4,4} \end{array} \xleftarrow{\text{softmax}} \begin{array}{cccc} \alpha_{1,1} & \alpha_{2,1} & \alpha_{3,1} & \alpha_{4,1} \\ \alpha_{1,2} & \alpha_{2,2} & \alpha_{3,2} & \alpha_{4,2} \\ \alpha_{1,3} & \alpha_{2,3} & \alpha_{3,3} & \alpha_{4,3} \\ \alpha_{1,4} & \alpha_{2,4} & \alpha_{3,4} & \alpha_{4,4} \end{array} = \begin{array}{c} k^1 \\ k^2 \\ k^3 \\ k^4 \end{array} \quad Q$$

$$A' \quad \text{softmax} \quad A \quad K^T$$

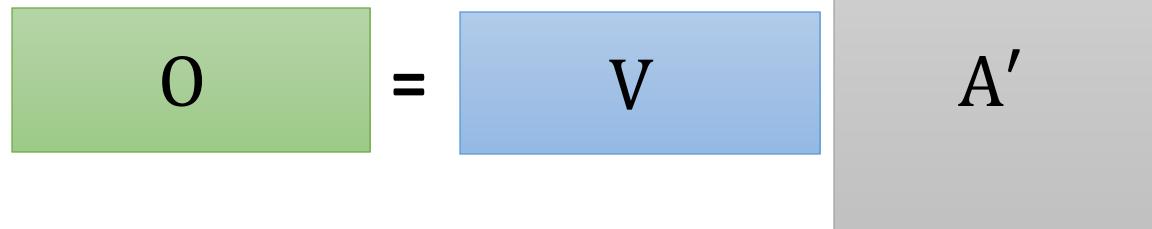
Self-attention



Self-attention



Attention Matrix



Many applications ...



Transformer

<https://arxiv.org/abs/1706.03762>



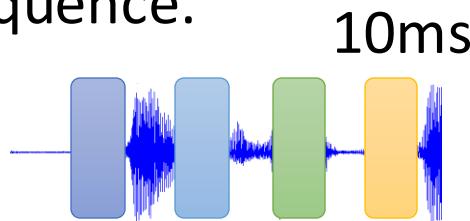
BERT

<https://arxiv.org/abs/1810.04805>

Widely used in Natural Language Processing (NLP)!

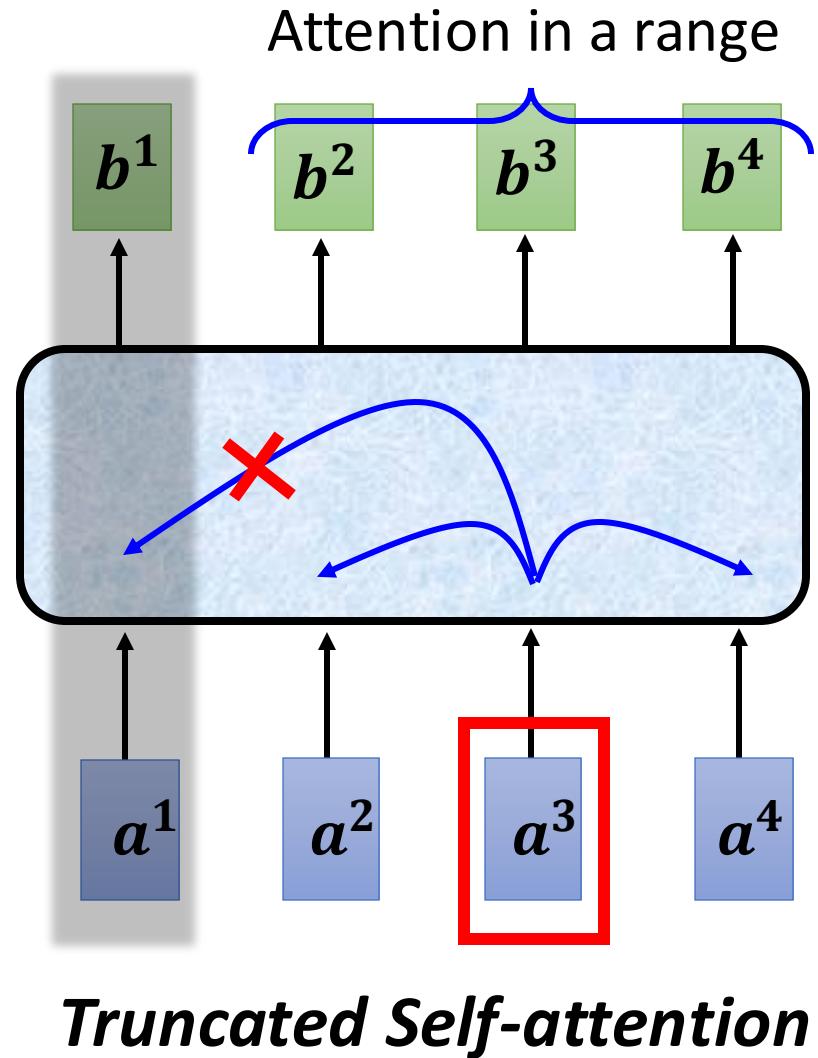
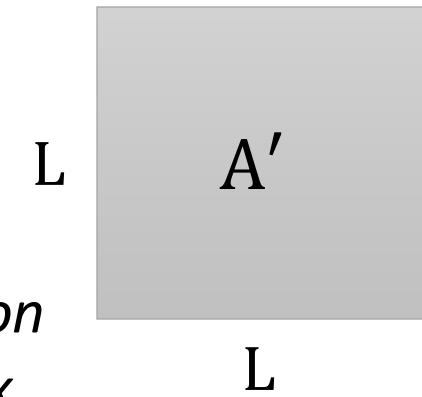
Self-attention for Speech

Speech is a very long vector sequence.



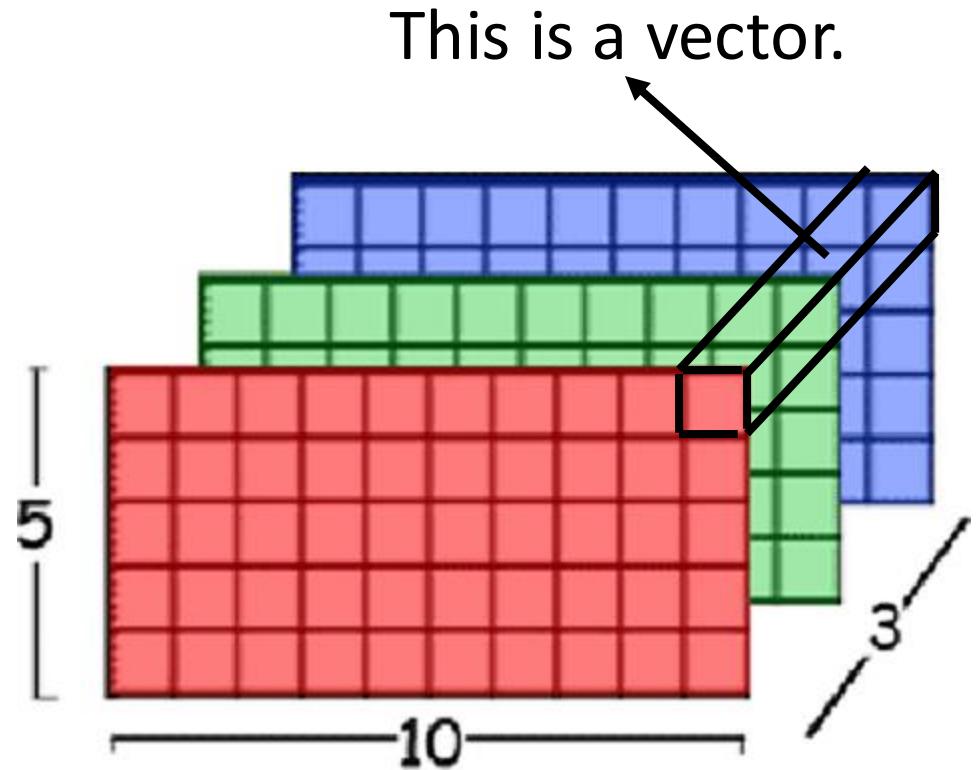
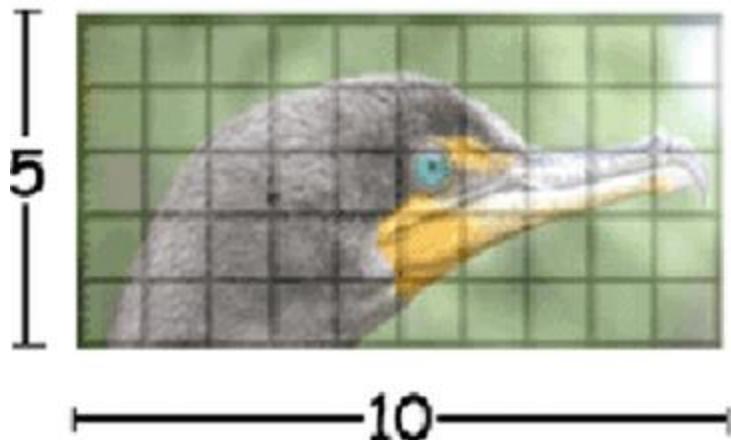
If input sequence is length L

*Attention
Matrix*



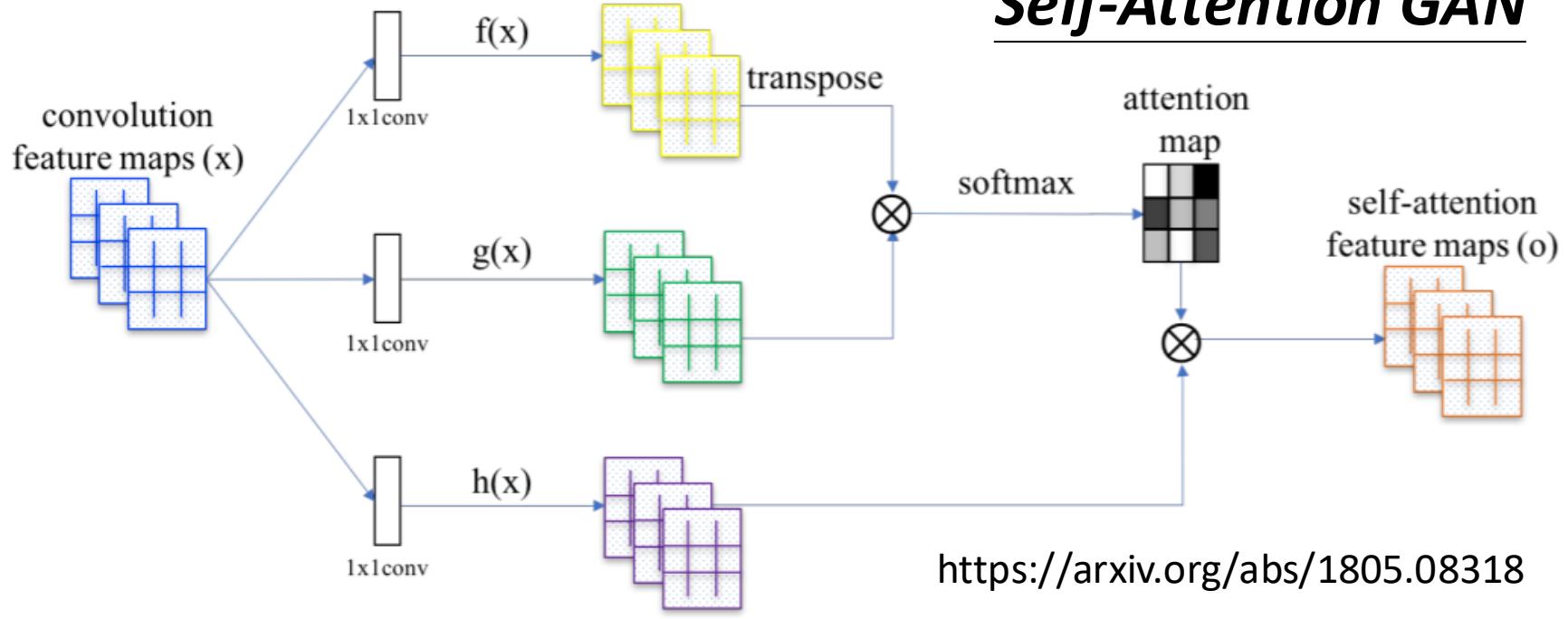
Self-attention for Image

An **image** can also be considered as a **vector set**.

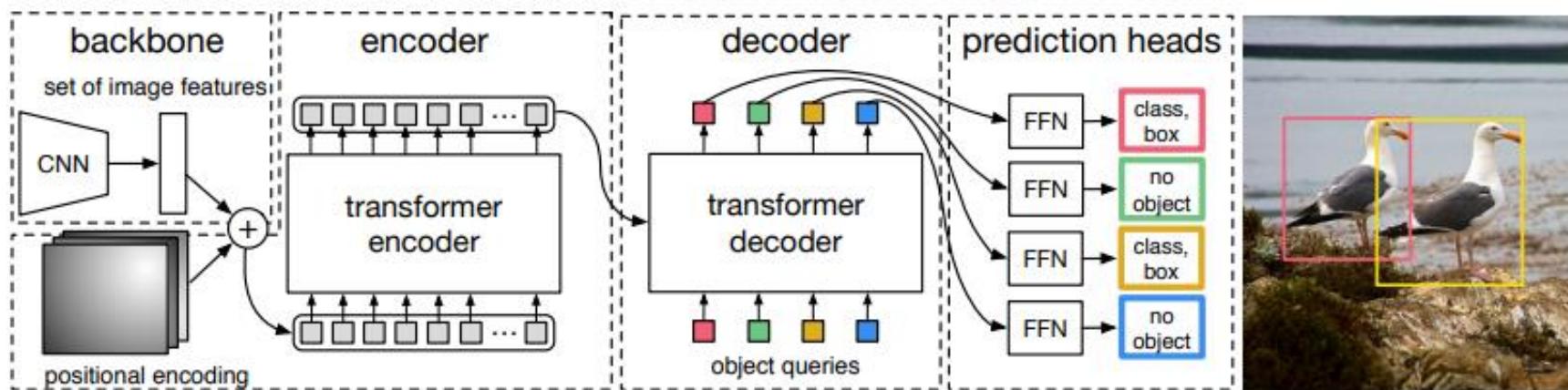


Source of image: https://www.researchgate.net/figure/Color-image-representation-and-RGB-matrix_fig15_282798184

Self-Attention GAN



DEtection Transformer (DETR)

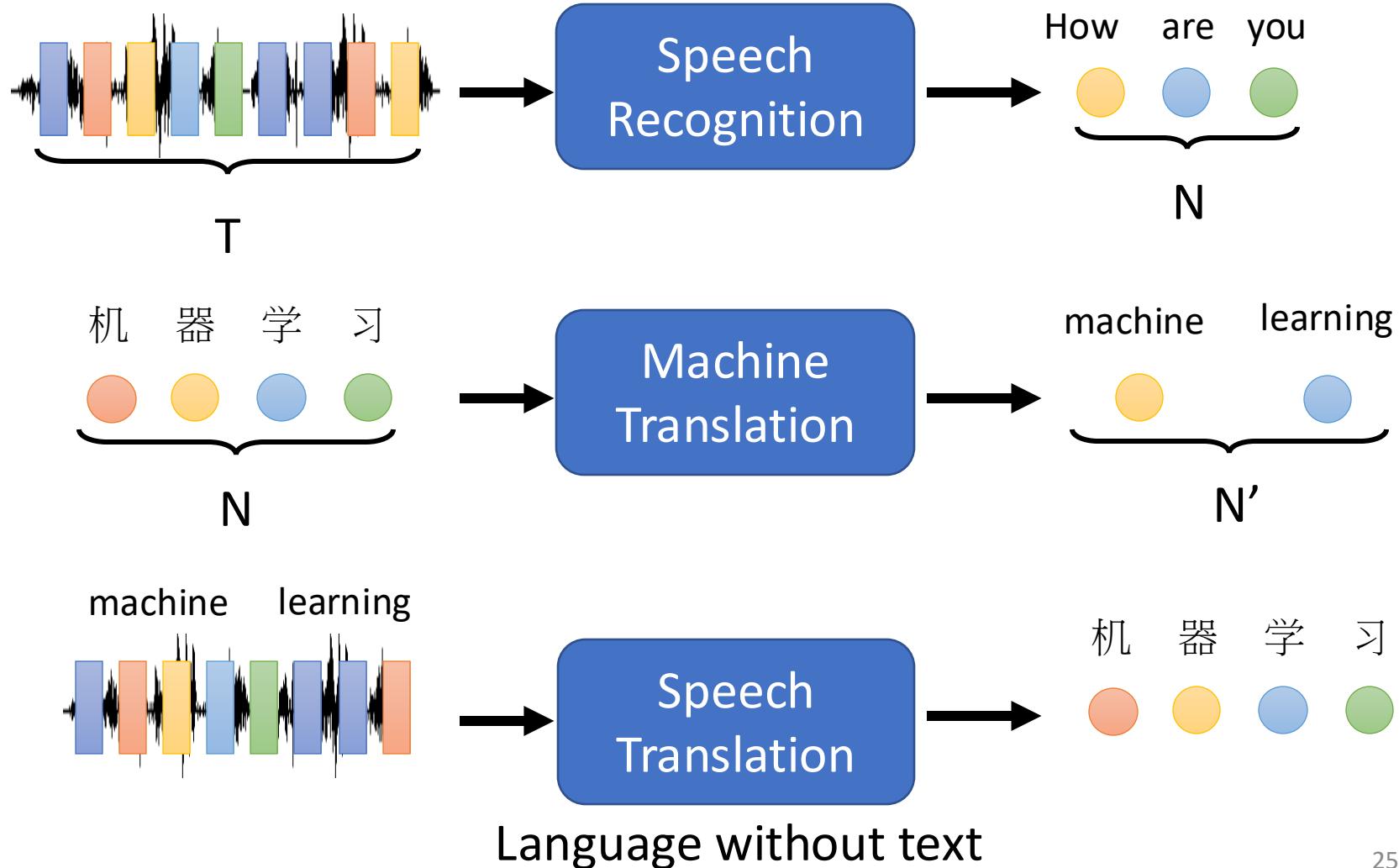


<https://arxiv.org/abs/2005.12872>

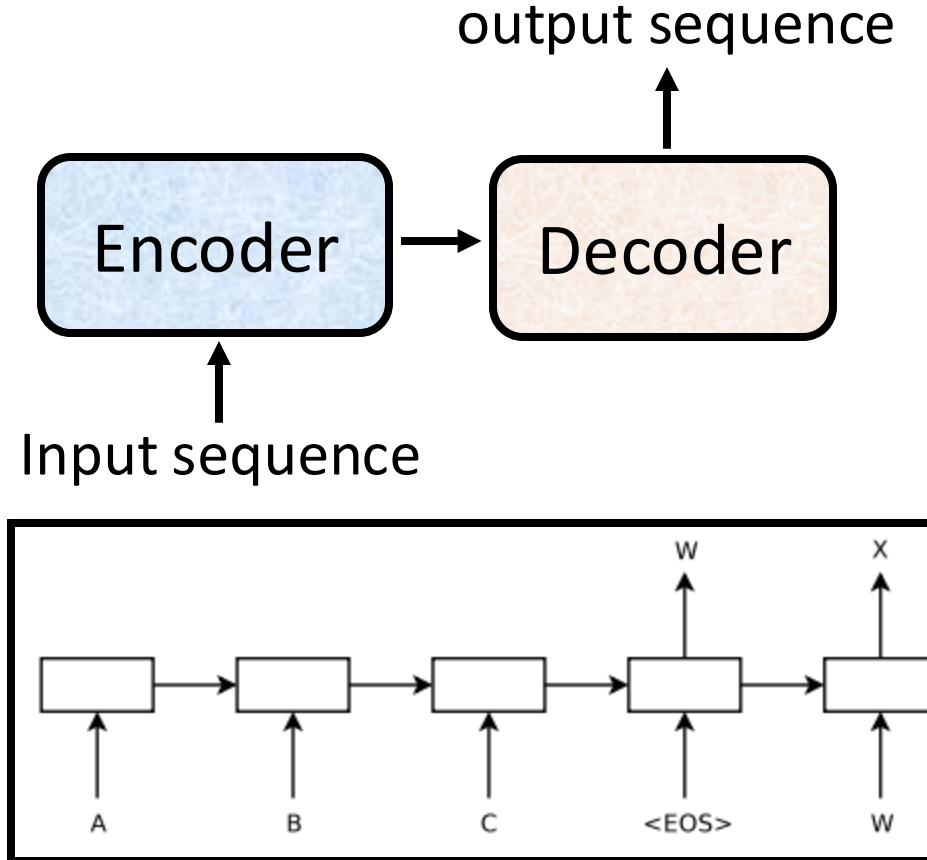
Sequence-to-sequence (Seq2seq)

Input a sequence, output a sequence

The output length is determined by model.

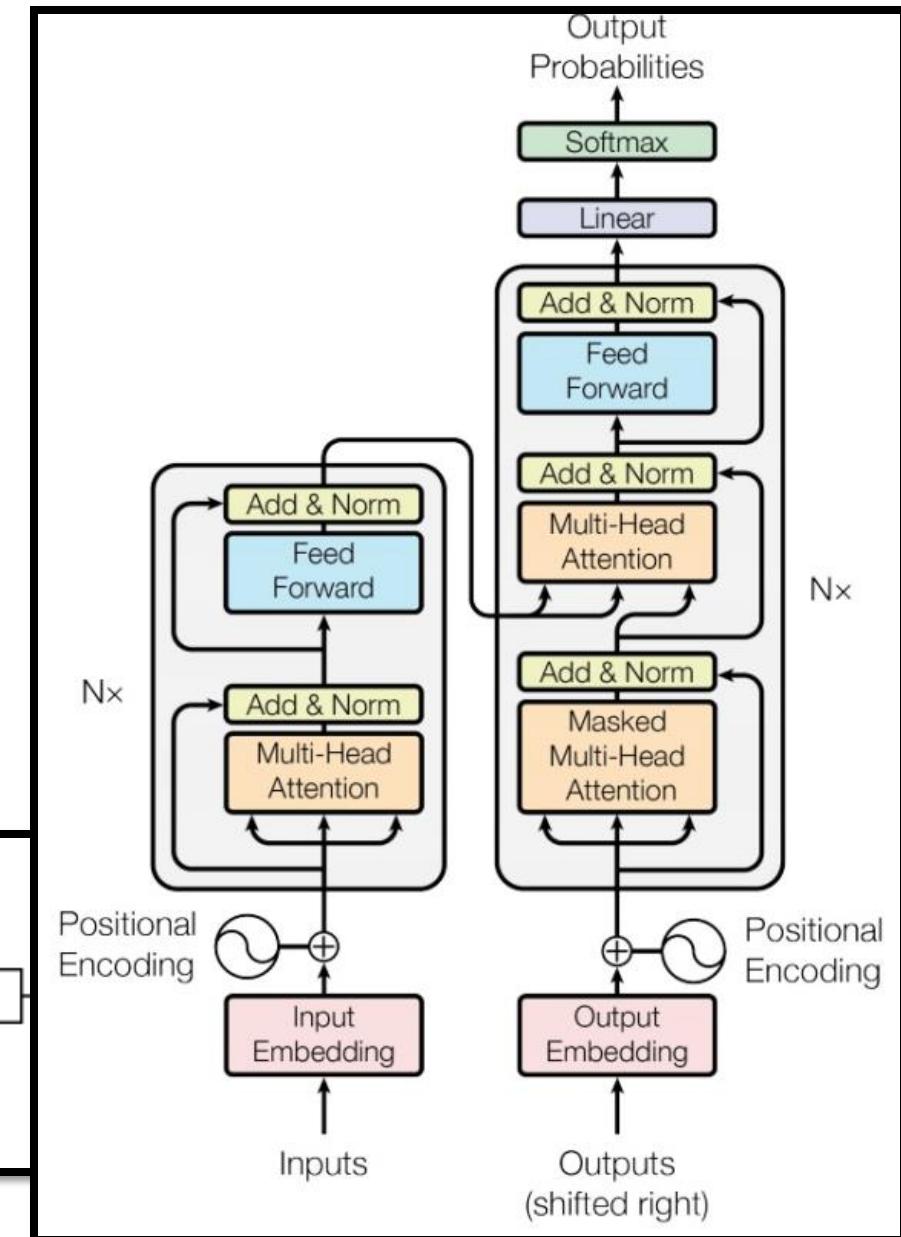


Seq2seq



Sequence to Sequence Learning with
Neural Networks

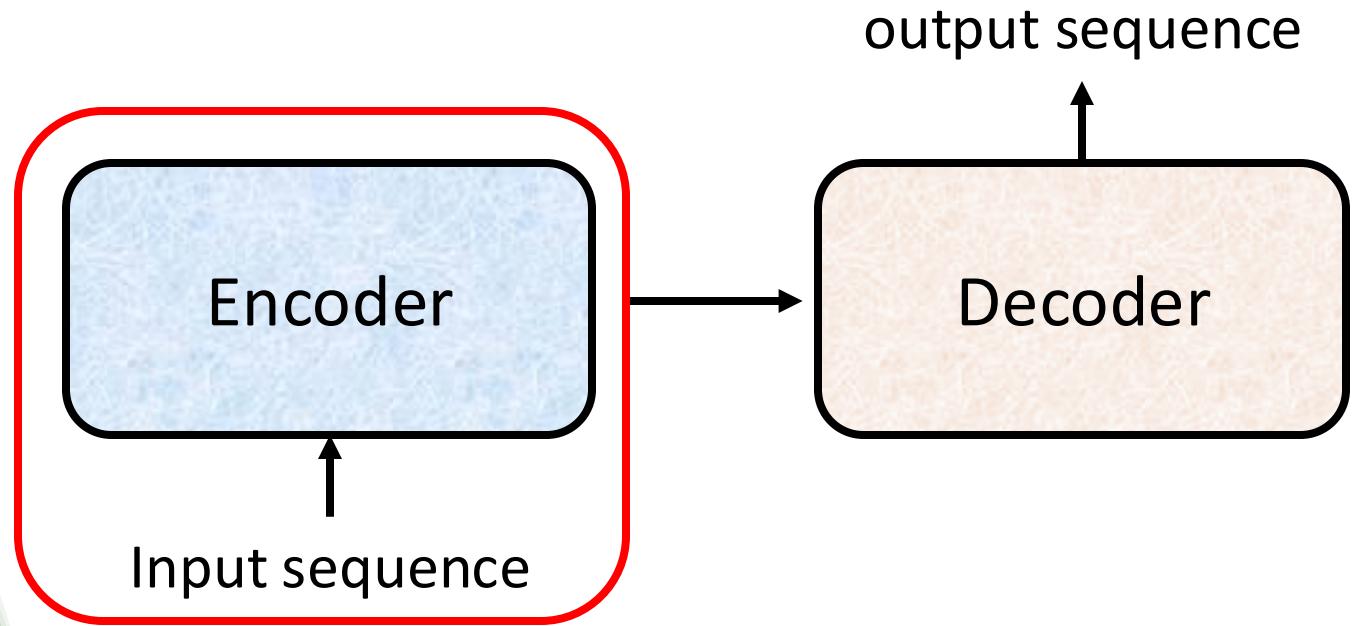
<https://arxiv.org/abs/1409.3215>



Transformer

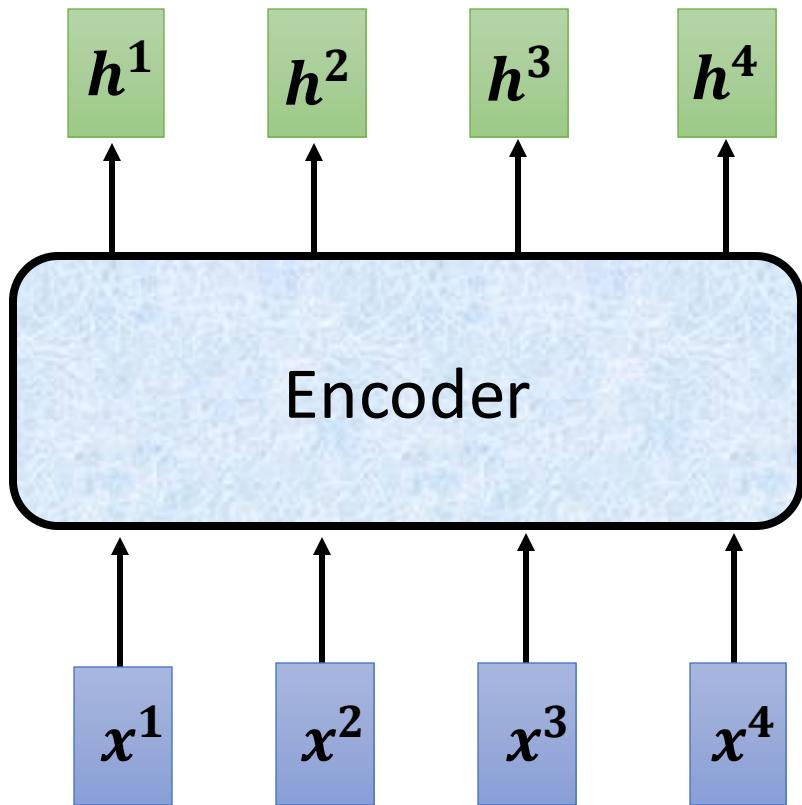
<https://arxiv.org/abs/1706.03762>

Encoder

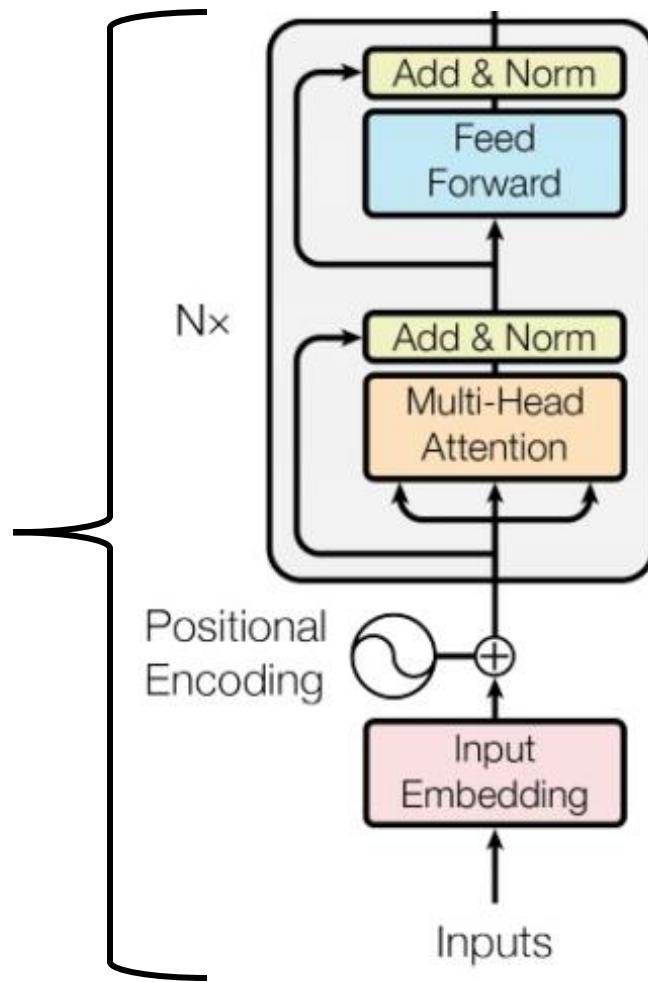


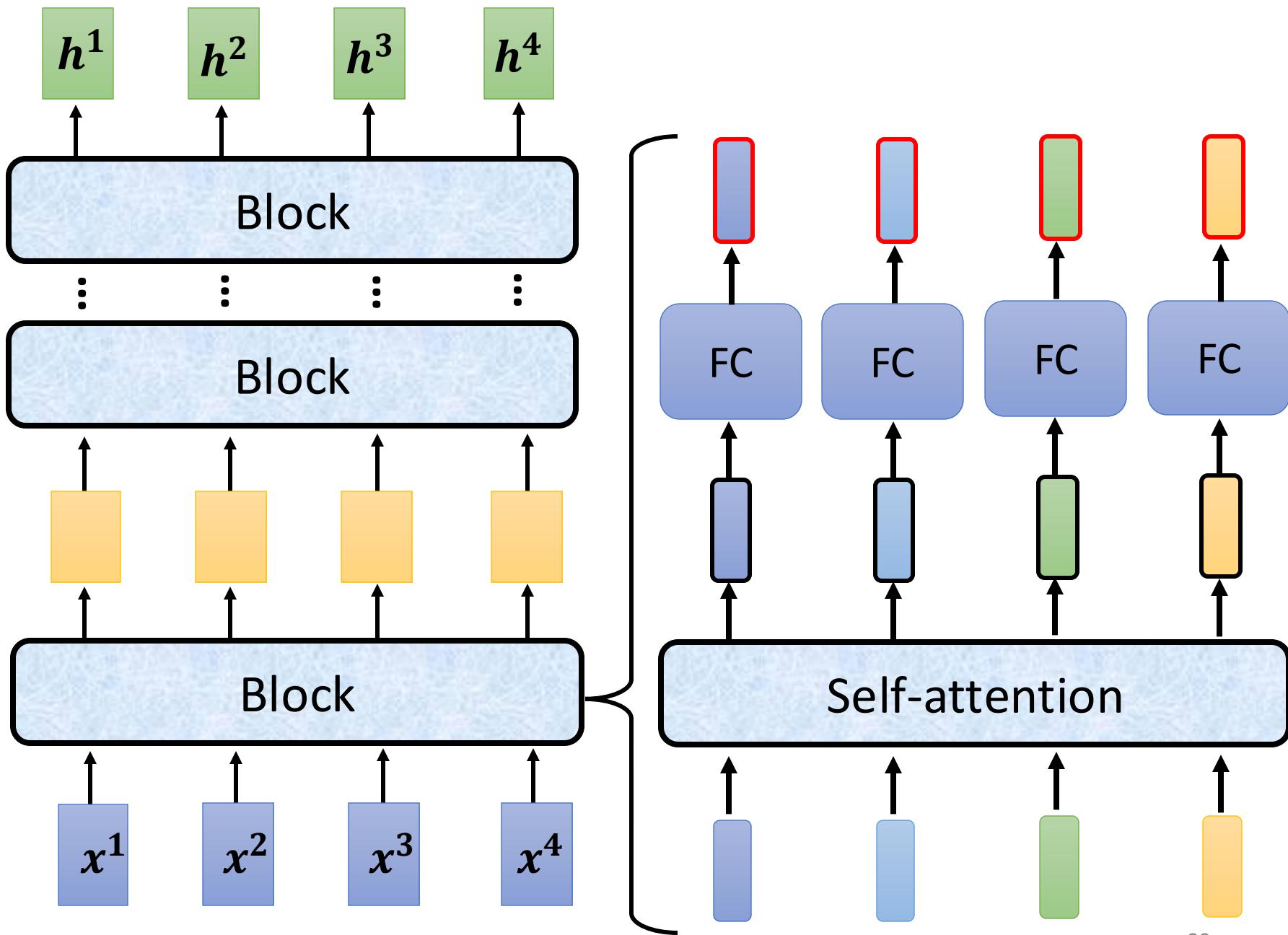
Encoder

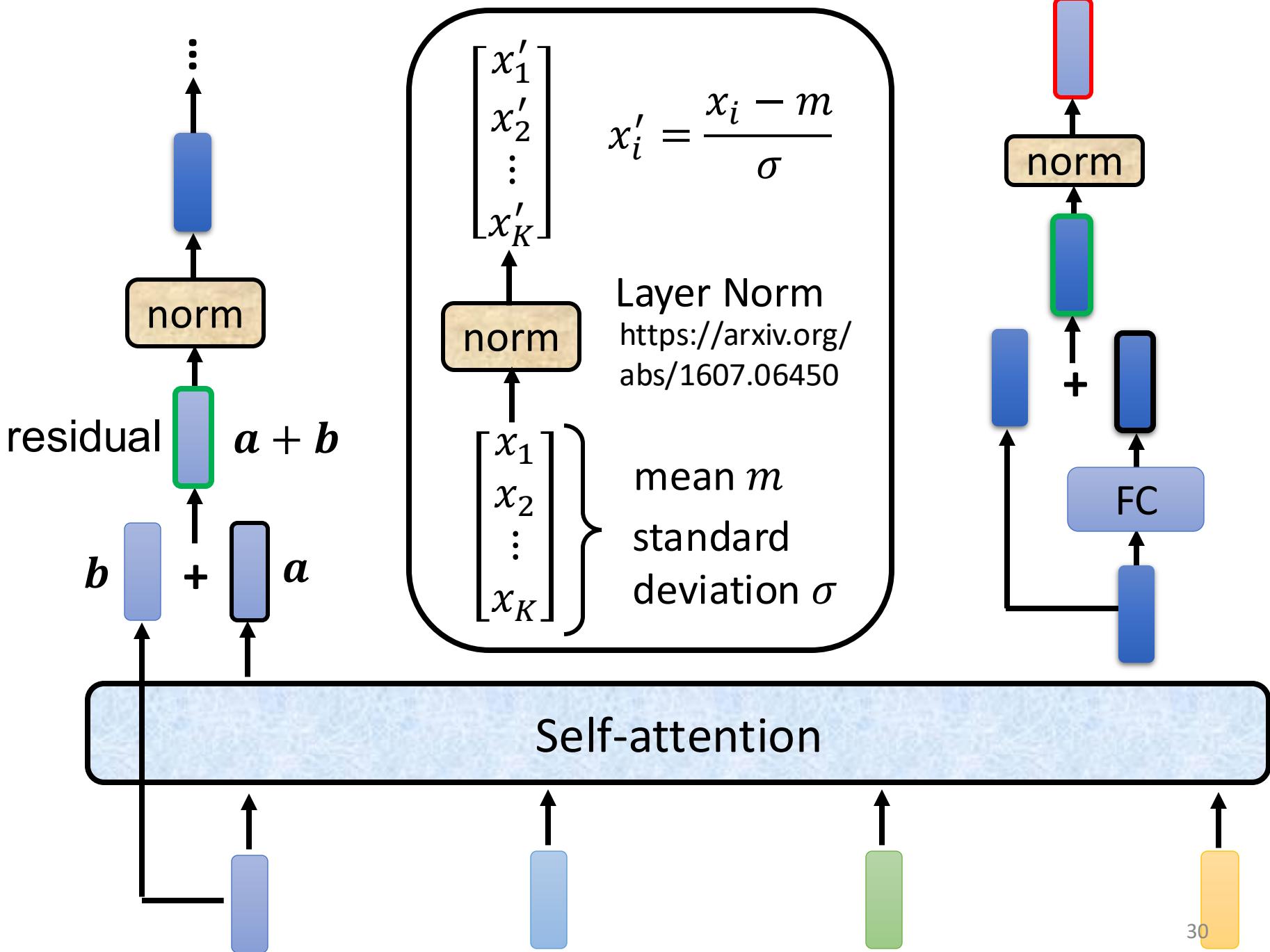
You can use **RNN** or **CNN**.

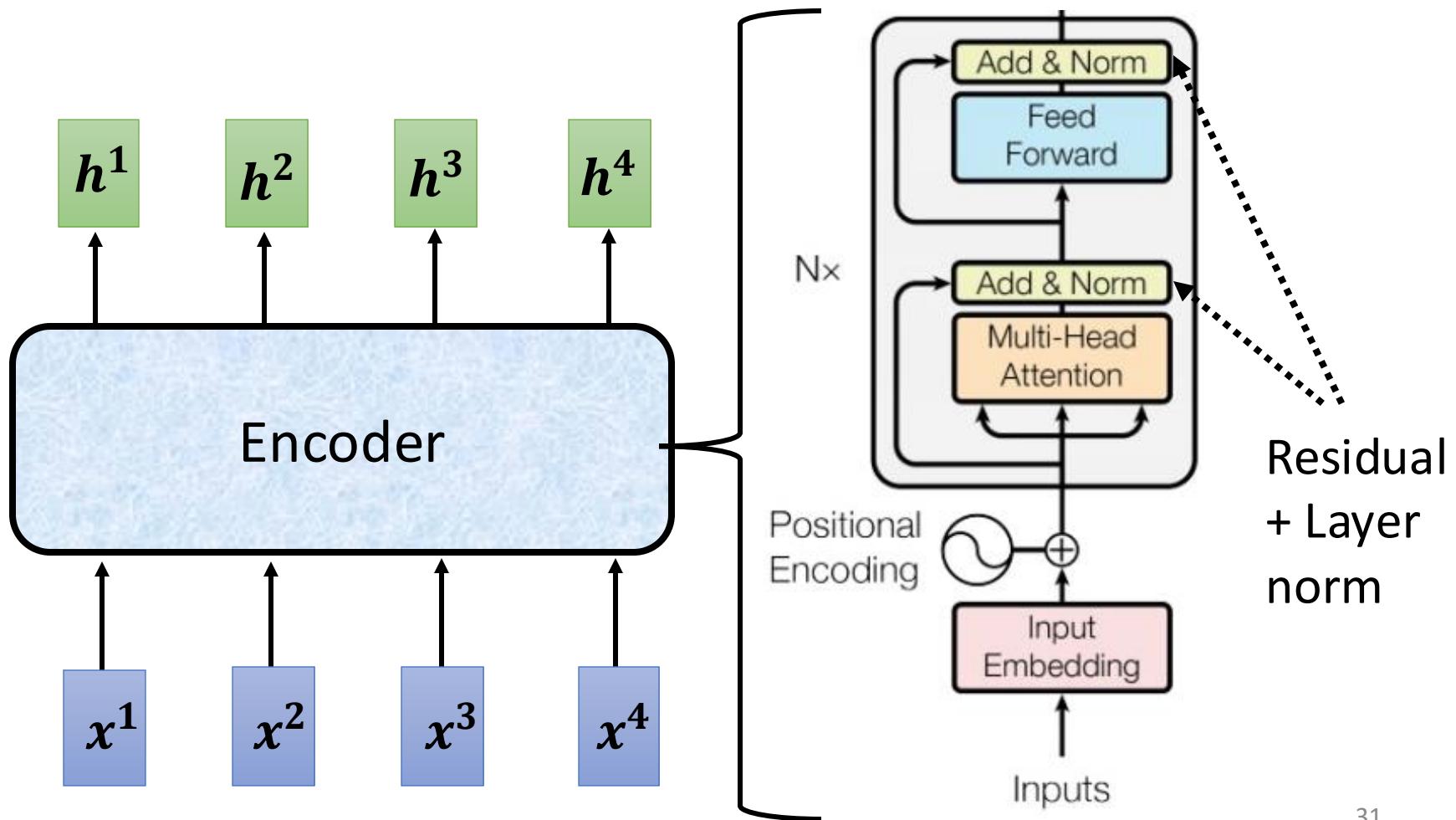


Transformer's Encoder



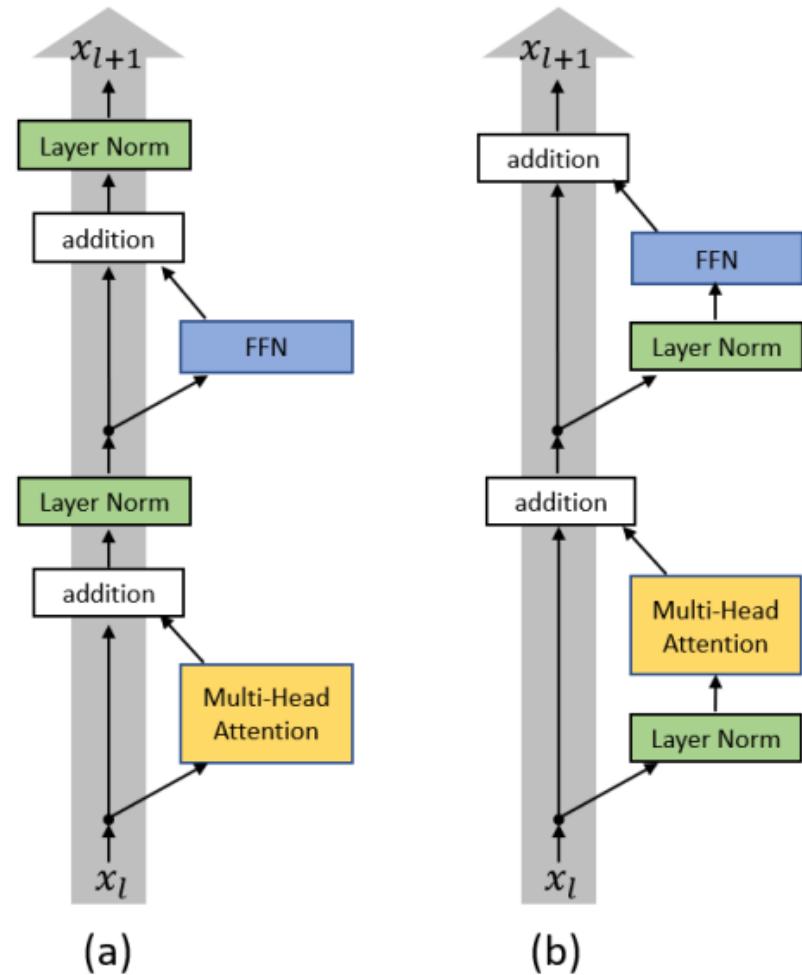




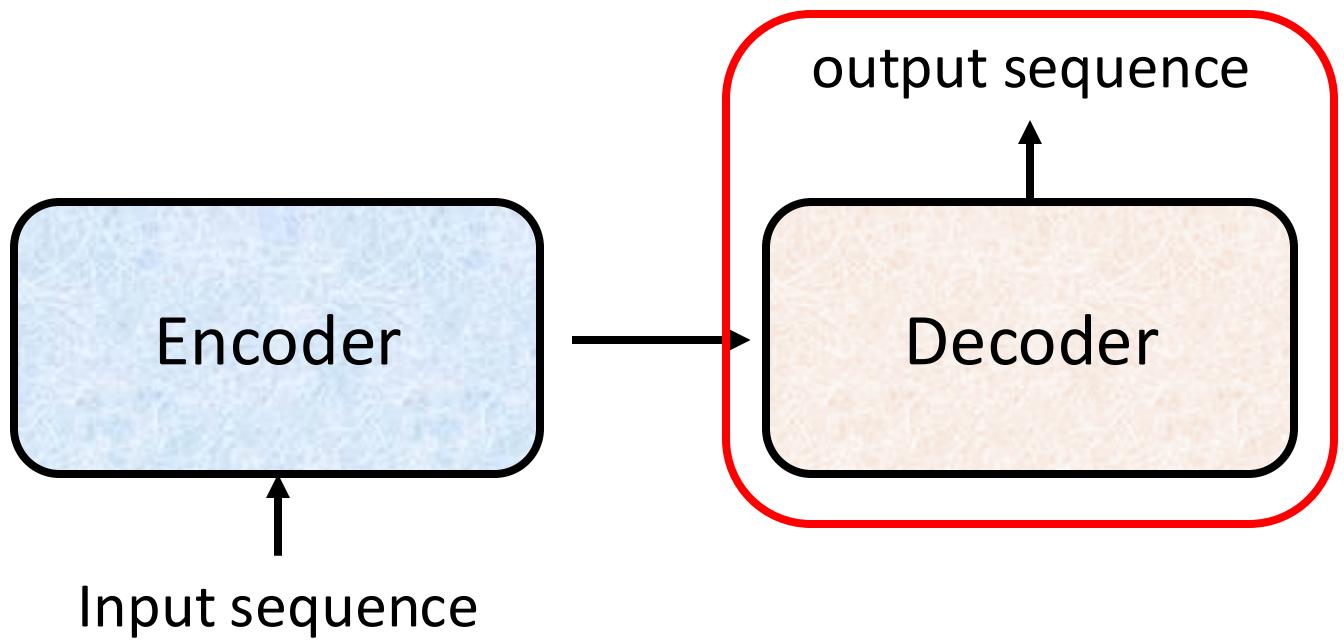


To learn more

- On Layer Normalization in the Transformer Architecture
- <https://arxiv.org/abs/2002.04745>
- PowerNorm: Rethinking Batch Normalization in Transformers
- <https://arxiv.org/abs/2003.07845>

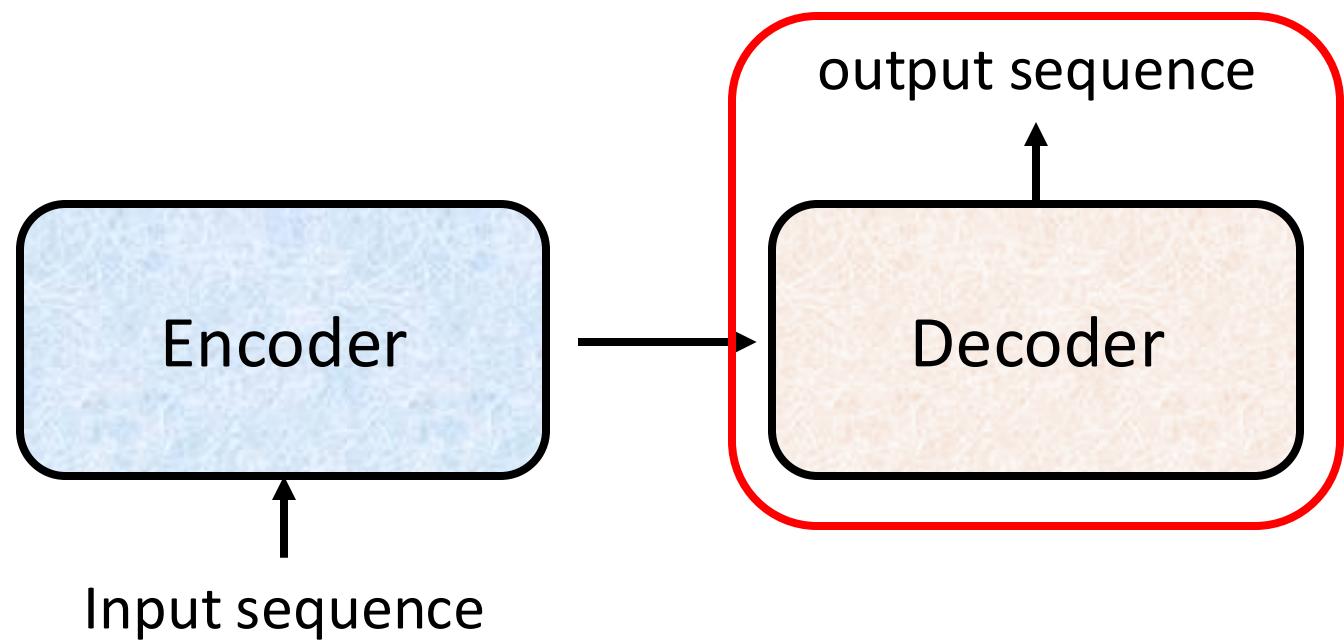


Decoder

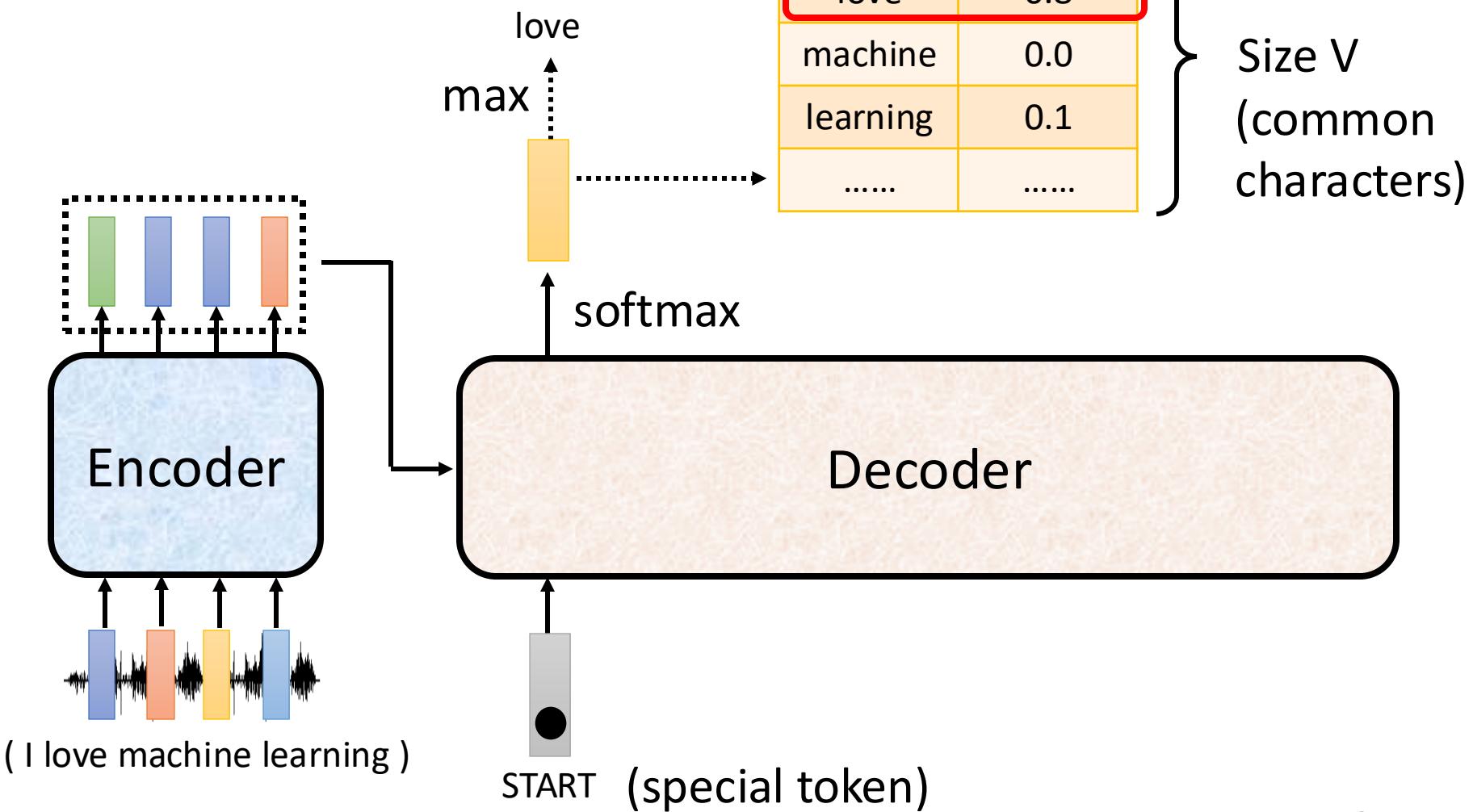


Decoder

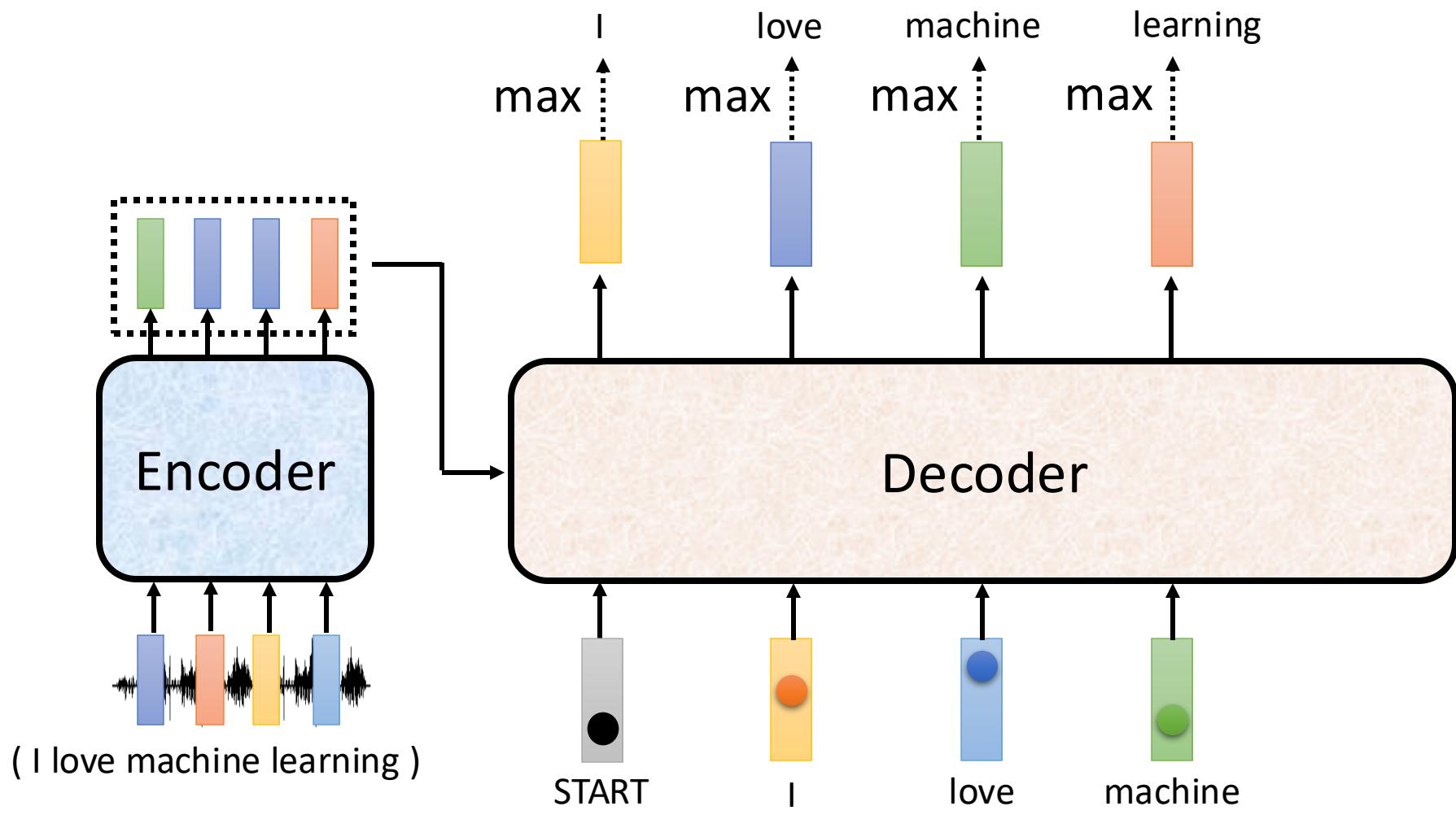
– Autoregressive (AT)

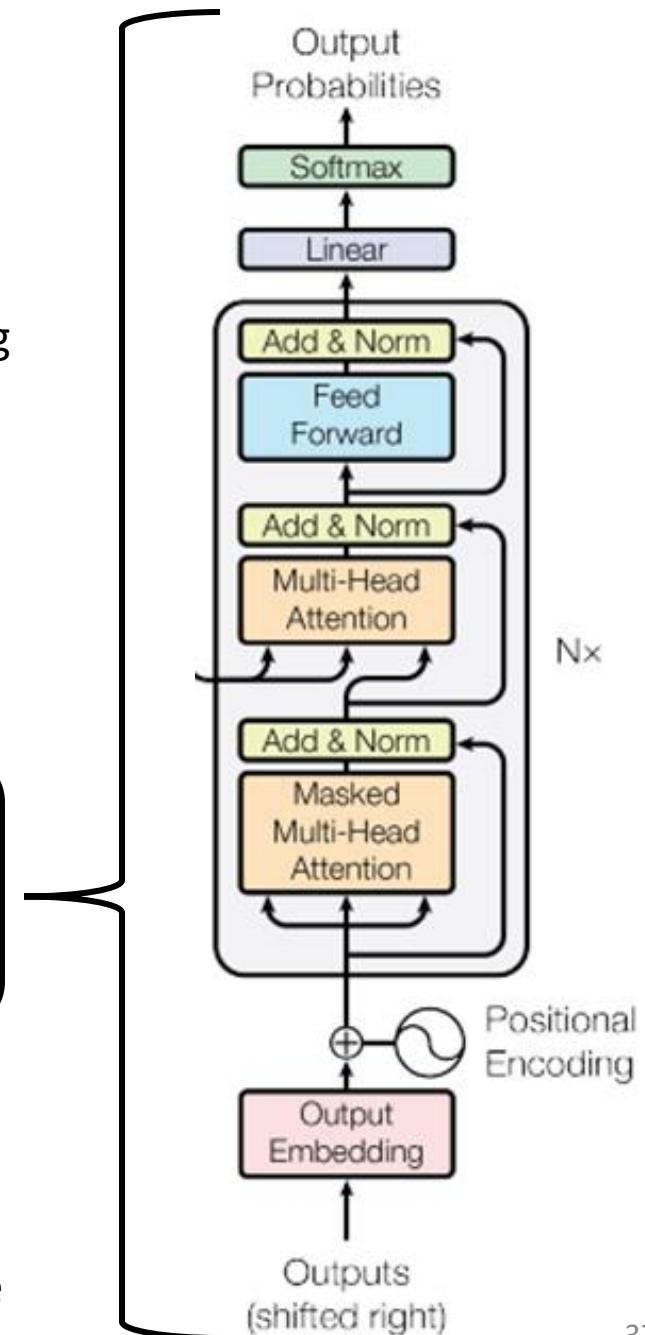
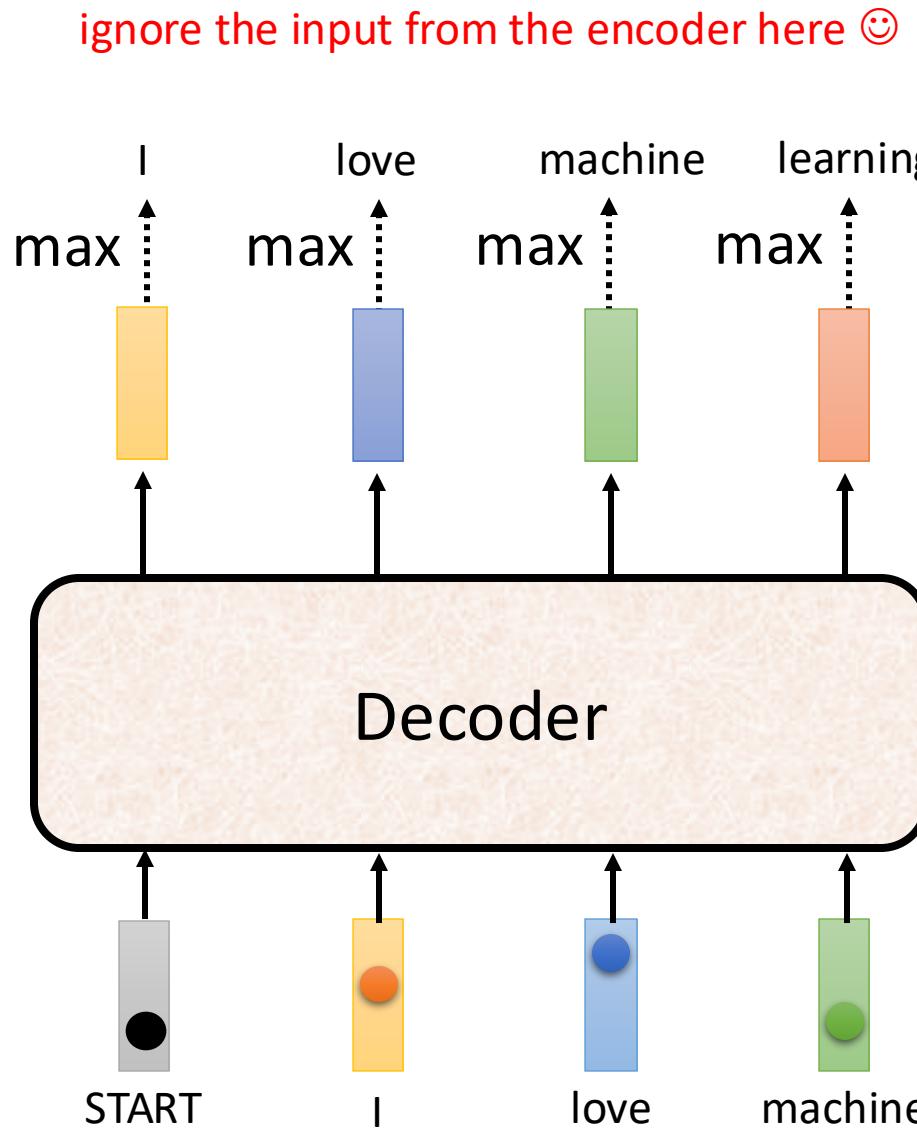


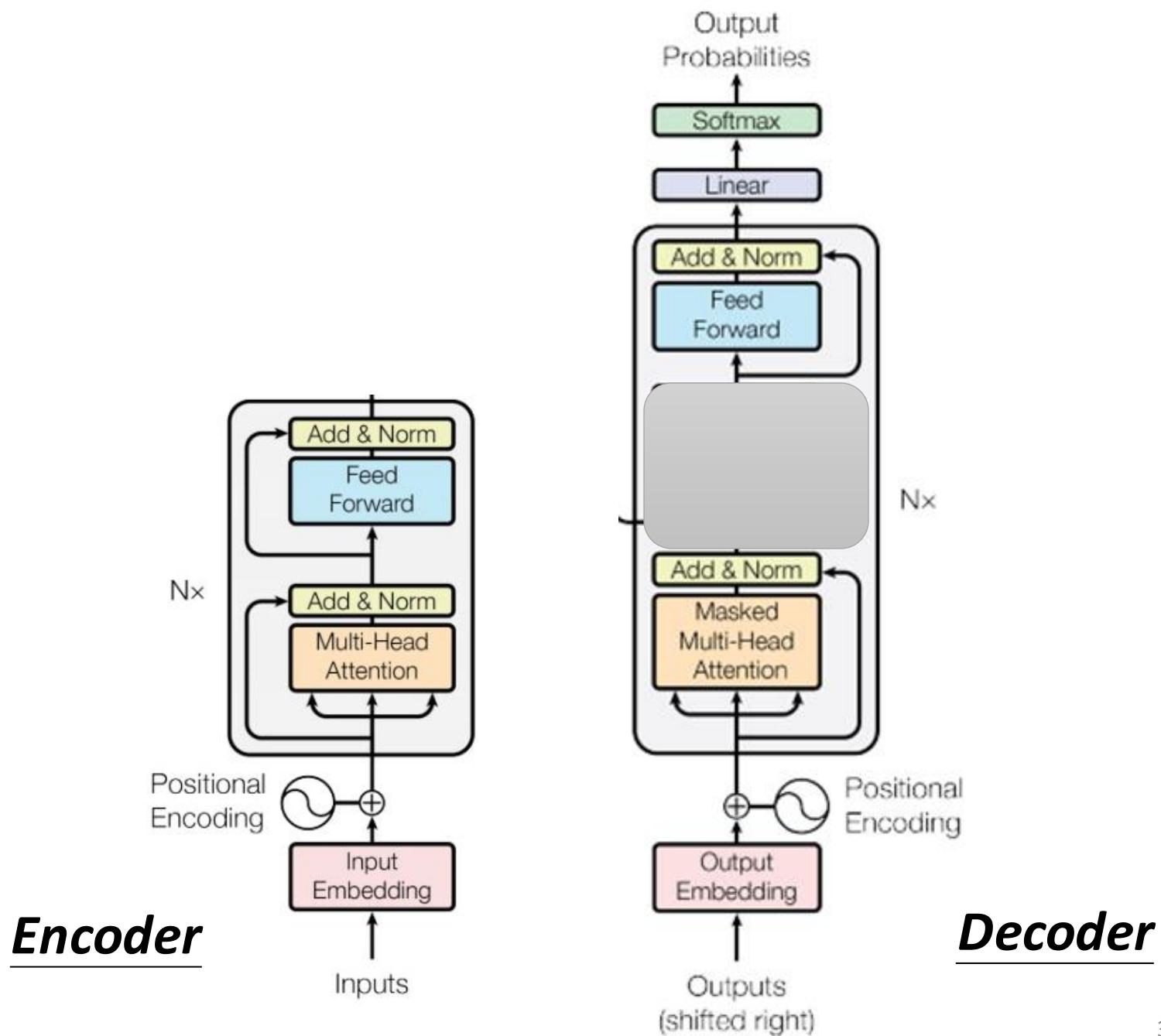
Autoregressive (Speech Recognition as example)



Autoregressive



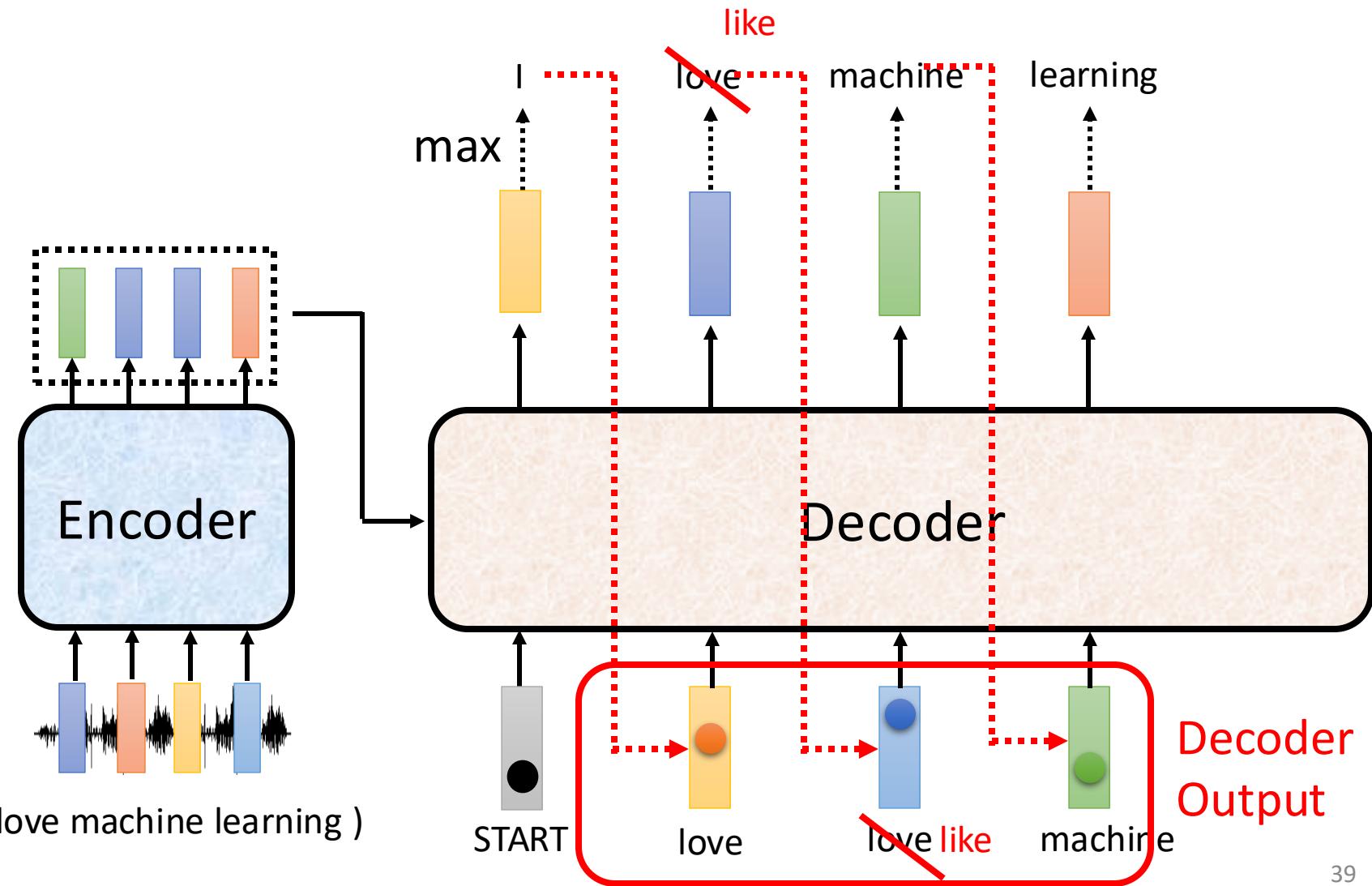




Encoder

Decoder

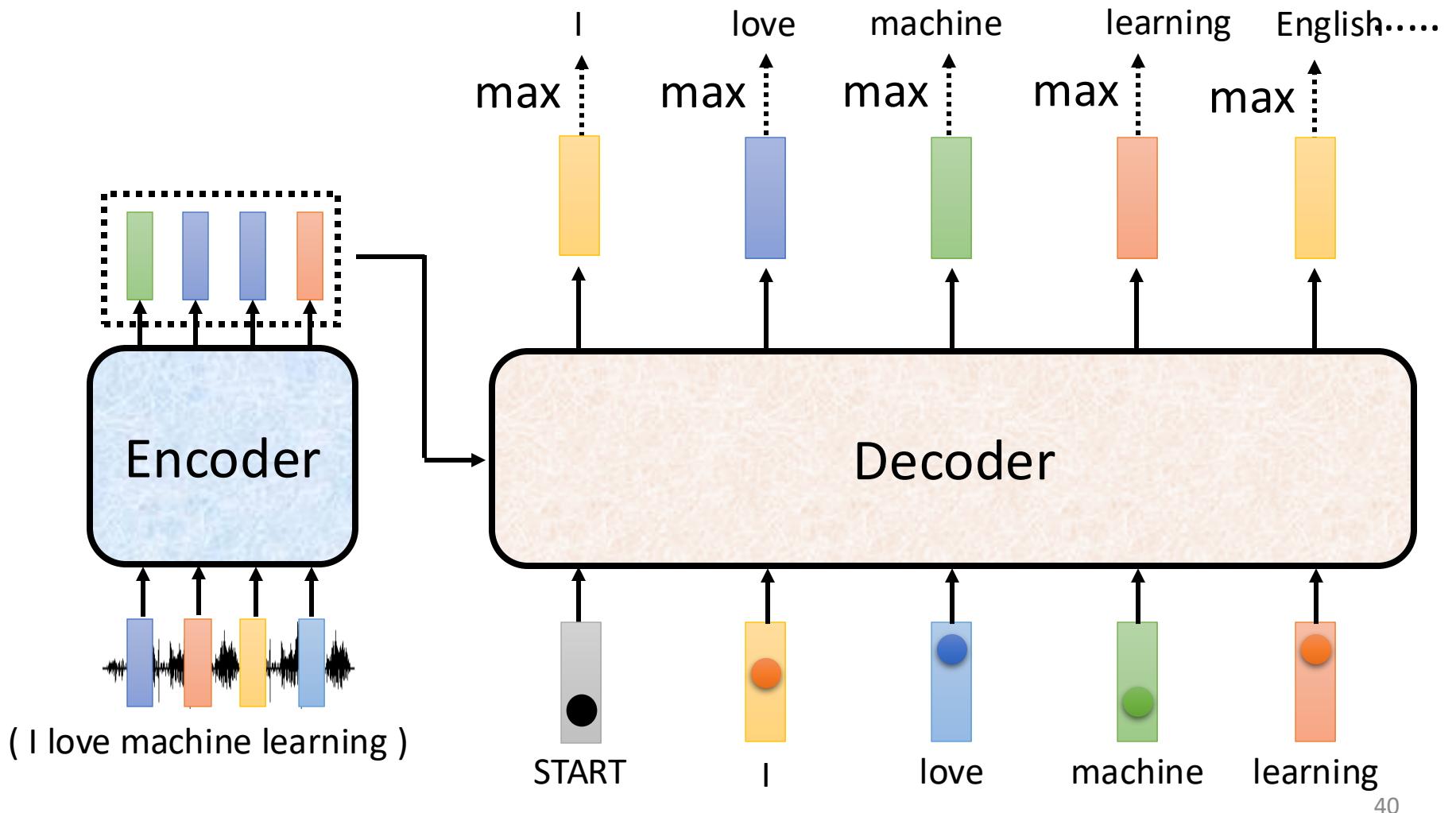
Autoregressive



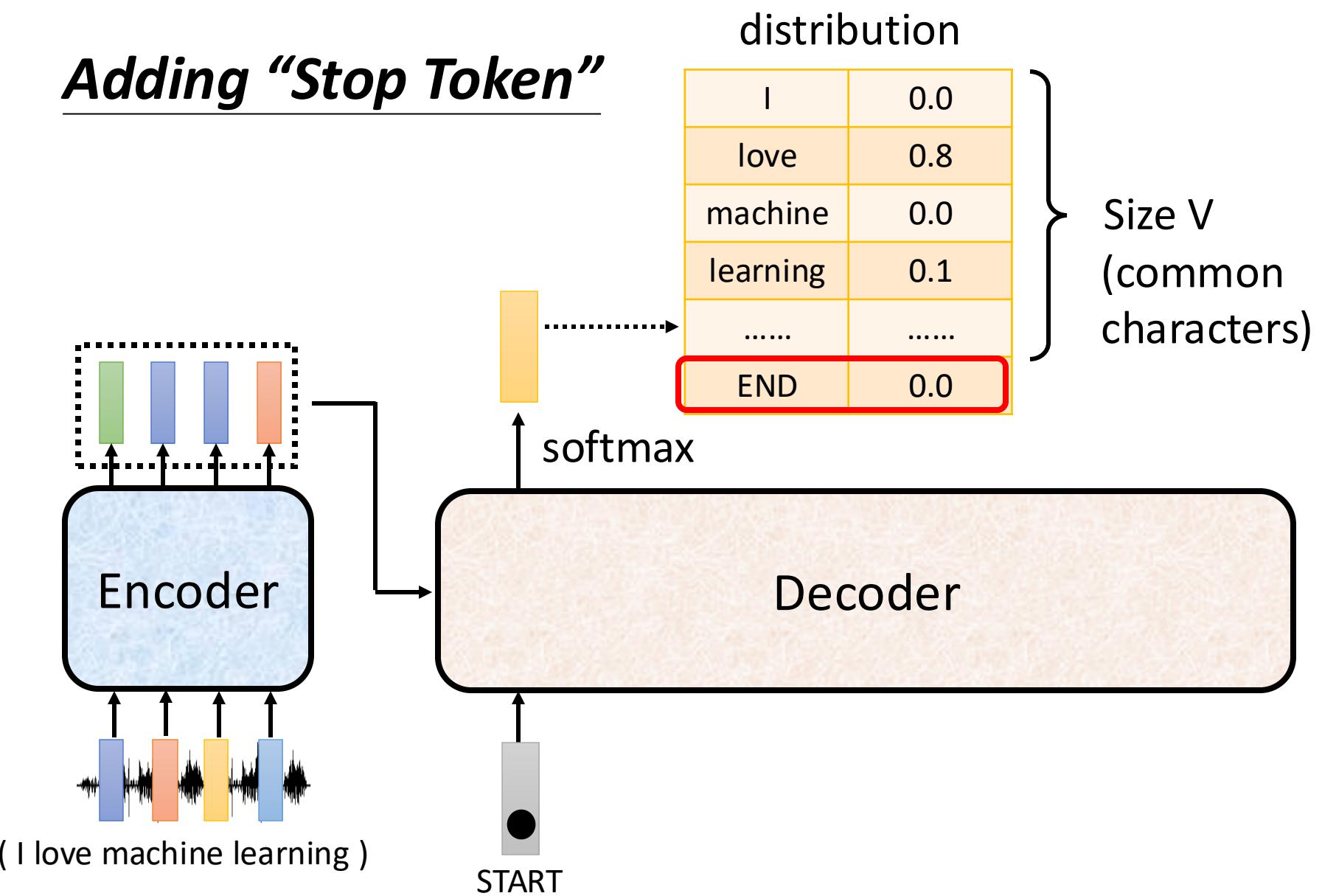
Autoregressive

We do not know the correct output length.

Never stop!

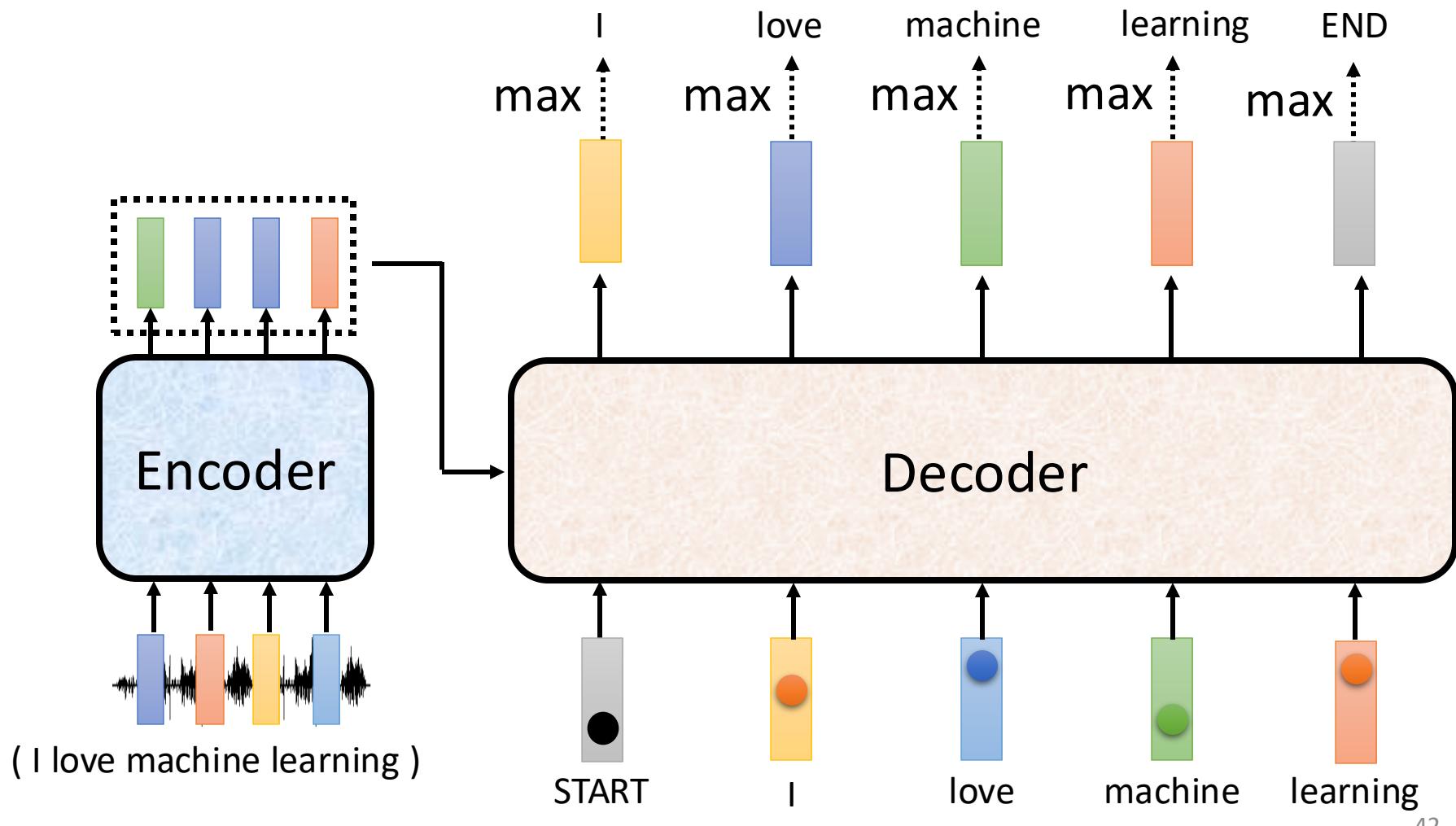


Adding “Stop Token”



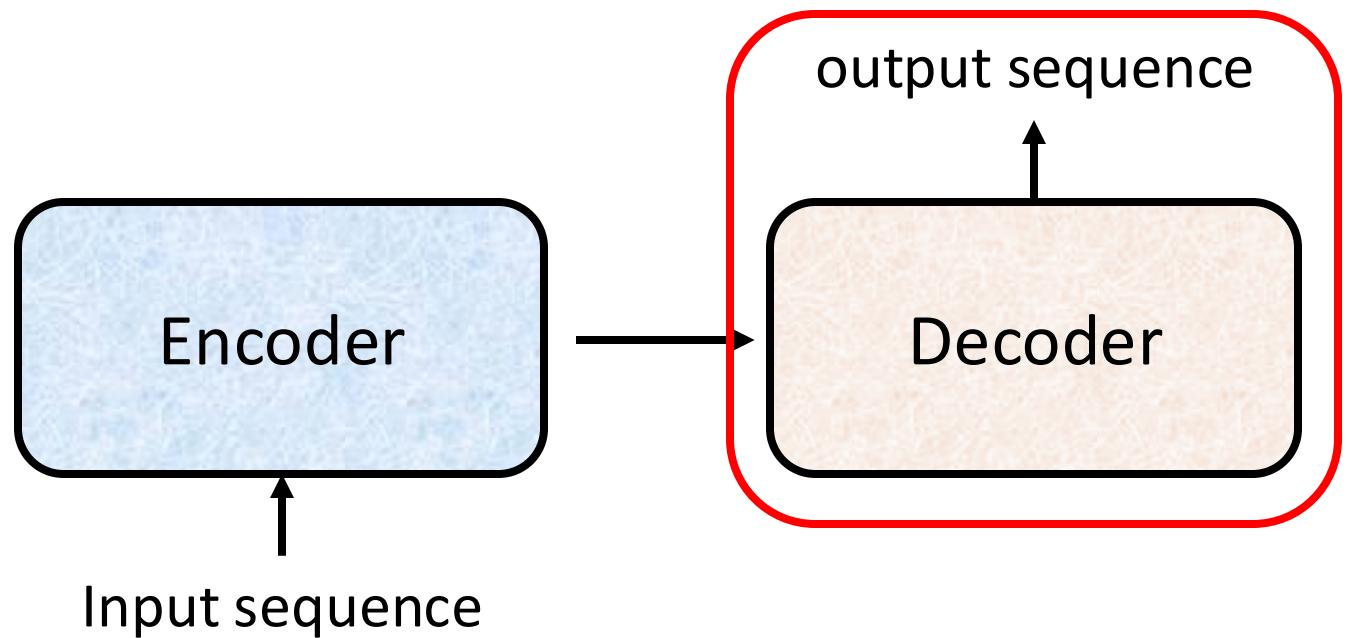
Autoregressive

Stop at here!

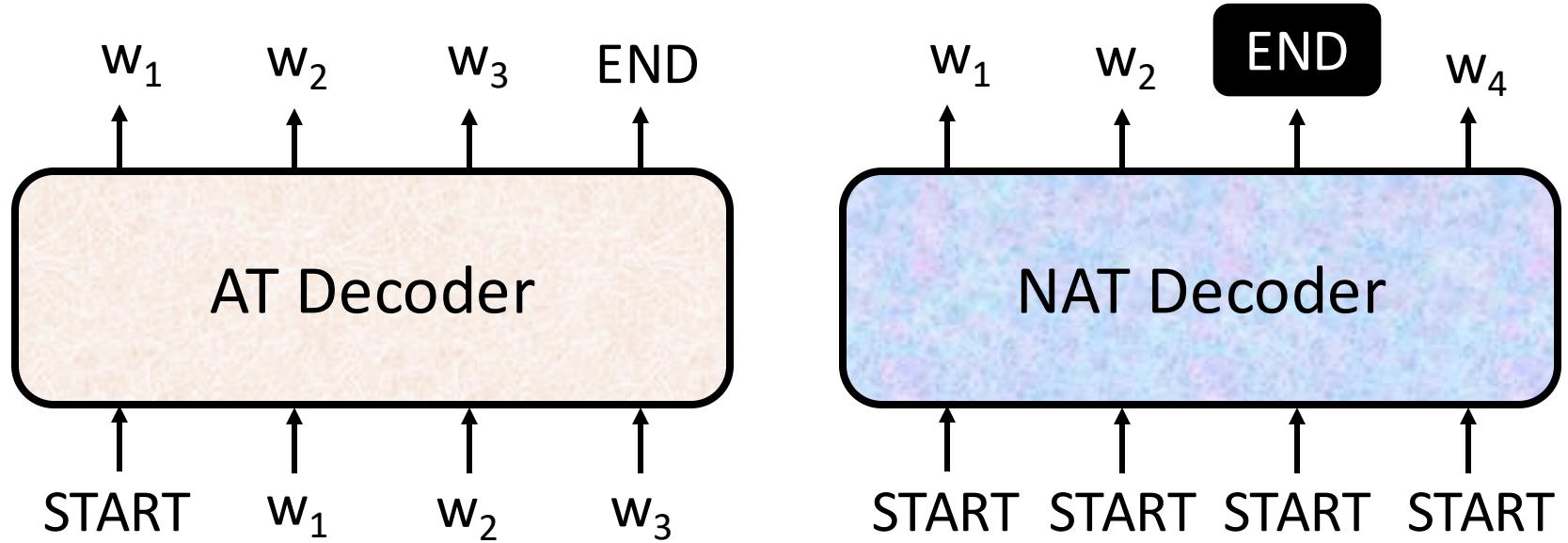


Decoder

– Non-autoregressive (NAT)

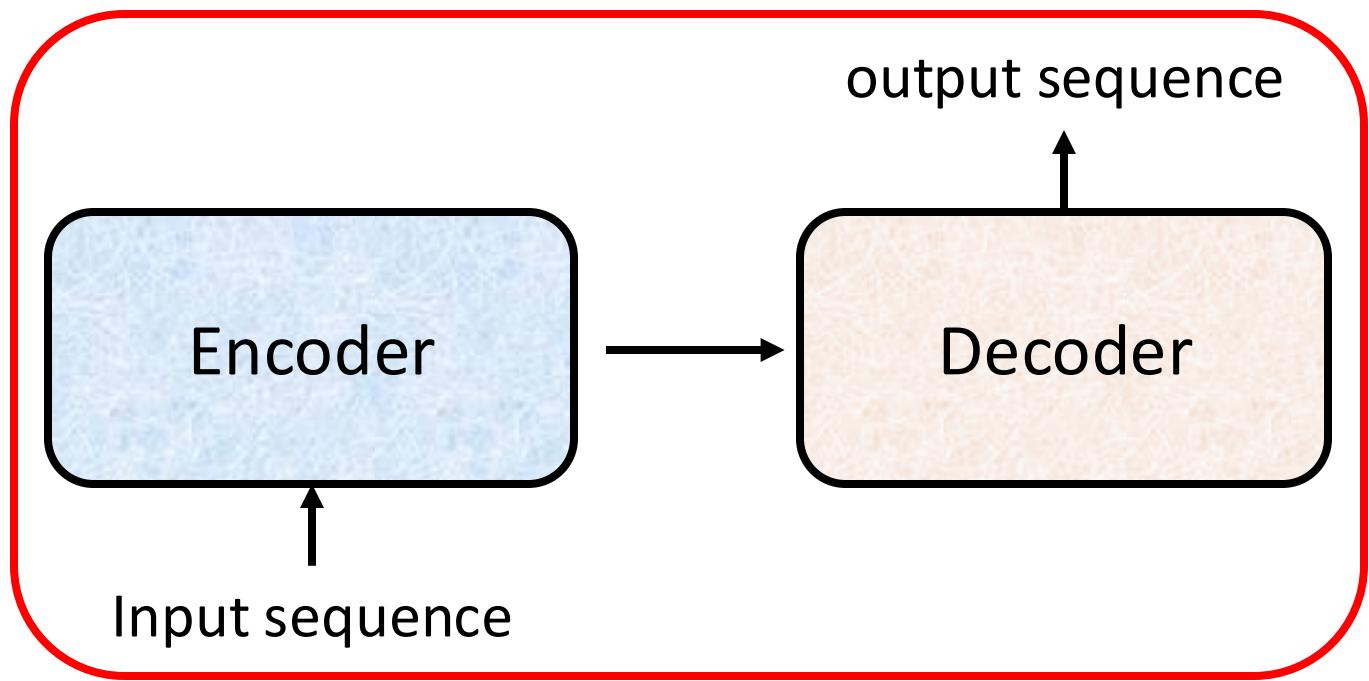


AT v.s. NAT



- How to decide the output length for NAT decoder?
 - Another predictor for output length
 - Output a very long sequence, ignore tokens after END
- Advantage: parallel, more stable generation (e.g., TTS)
- NAT is usually worse than AT (why? **Multi-modality**)

Encoder-Decoder



Concluding Remarks: *Transformer*

