

MAE 688: Machine Learning for Mechanical Engineers

Lecture 6-Transformer



Dr. Bing Dong

Director, Built Environment Science and Technology (BEST) Lab

Professor

Mechanical and Aerospace Engineering

Syracuse University

Transformer

2017

Attention Is All You Need

Ashish Vaswani*
Google Brain
avaswani@google.com

Noam Shazeer*
Google Brain
noam@google.com

Niki Parmar*
Google Research
nikip@google.com

Jakob Uszkoreit*
Google Research
usz@google.com

Llion Jones*
Google Research
llion@google.com

Aidan N. Gomez* †
University of Toronto
aidan@cs.toronto.edu

Lukasz Kaiser*
Google Brain
lukaszkaizer@google.com



Transformer



BERT 340M



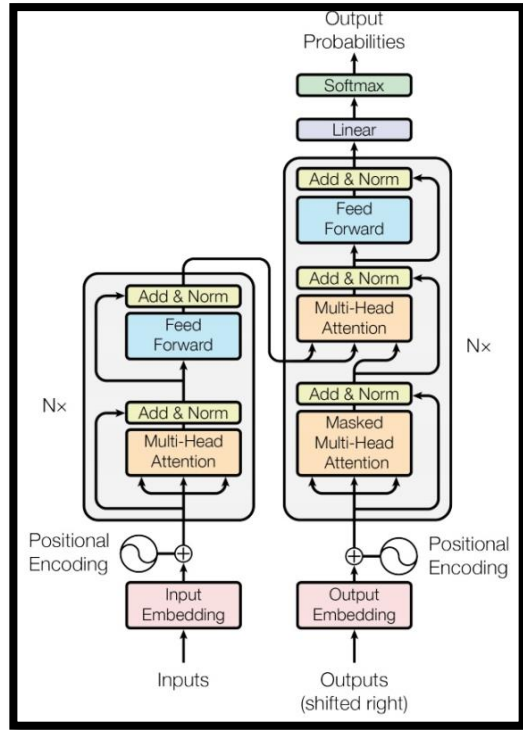
GPT-2 (1.5B)



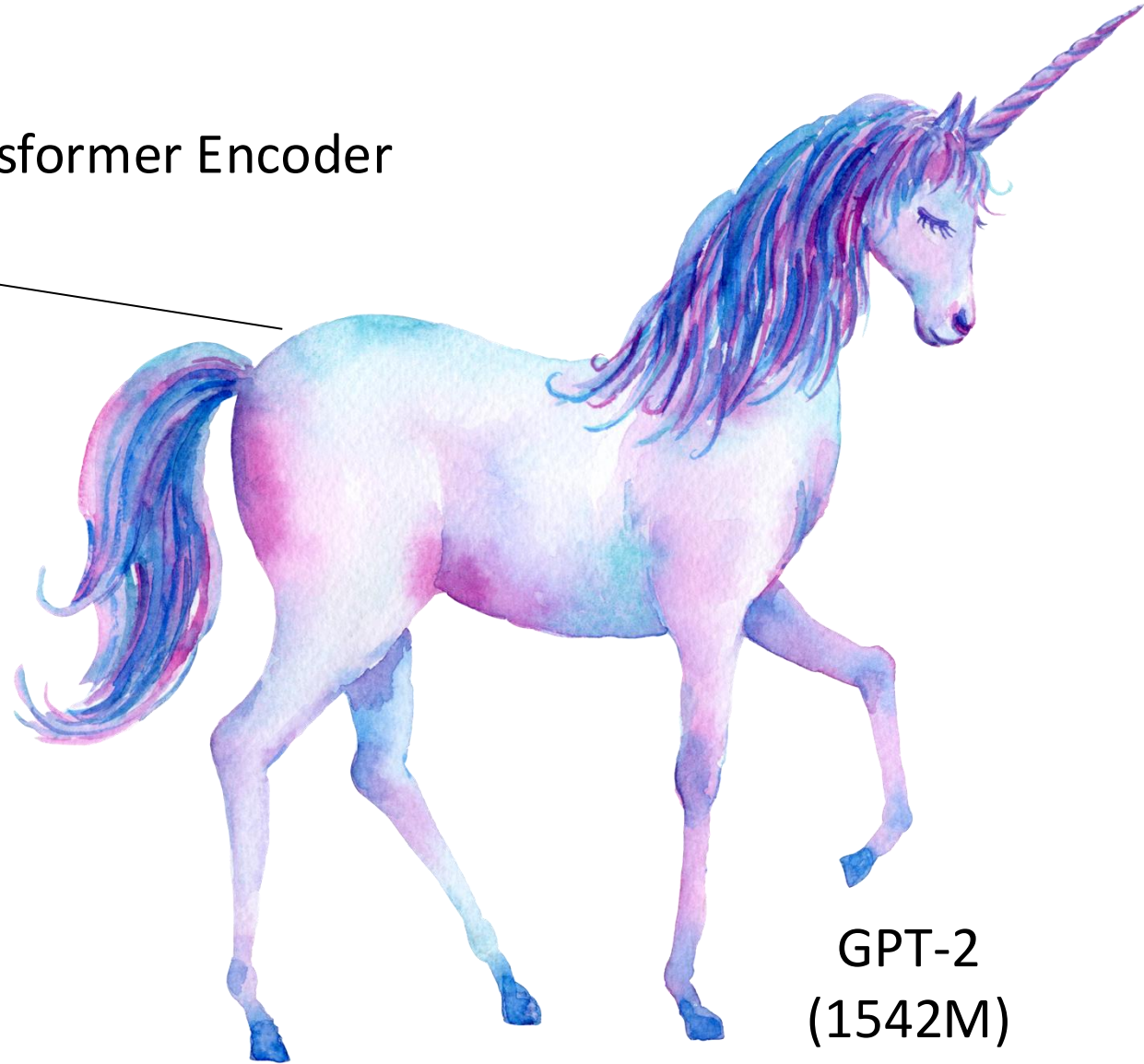
GPT-3 (175B)

GPT-4 (2T)

Generative Pre-Training (GPT)



Transformer Encoder



BERT
(340M)

ELMO
(94M)



GPT-2
(1542M)

Source of image: <https://huaban.com/pins/1714071707/>

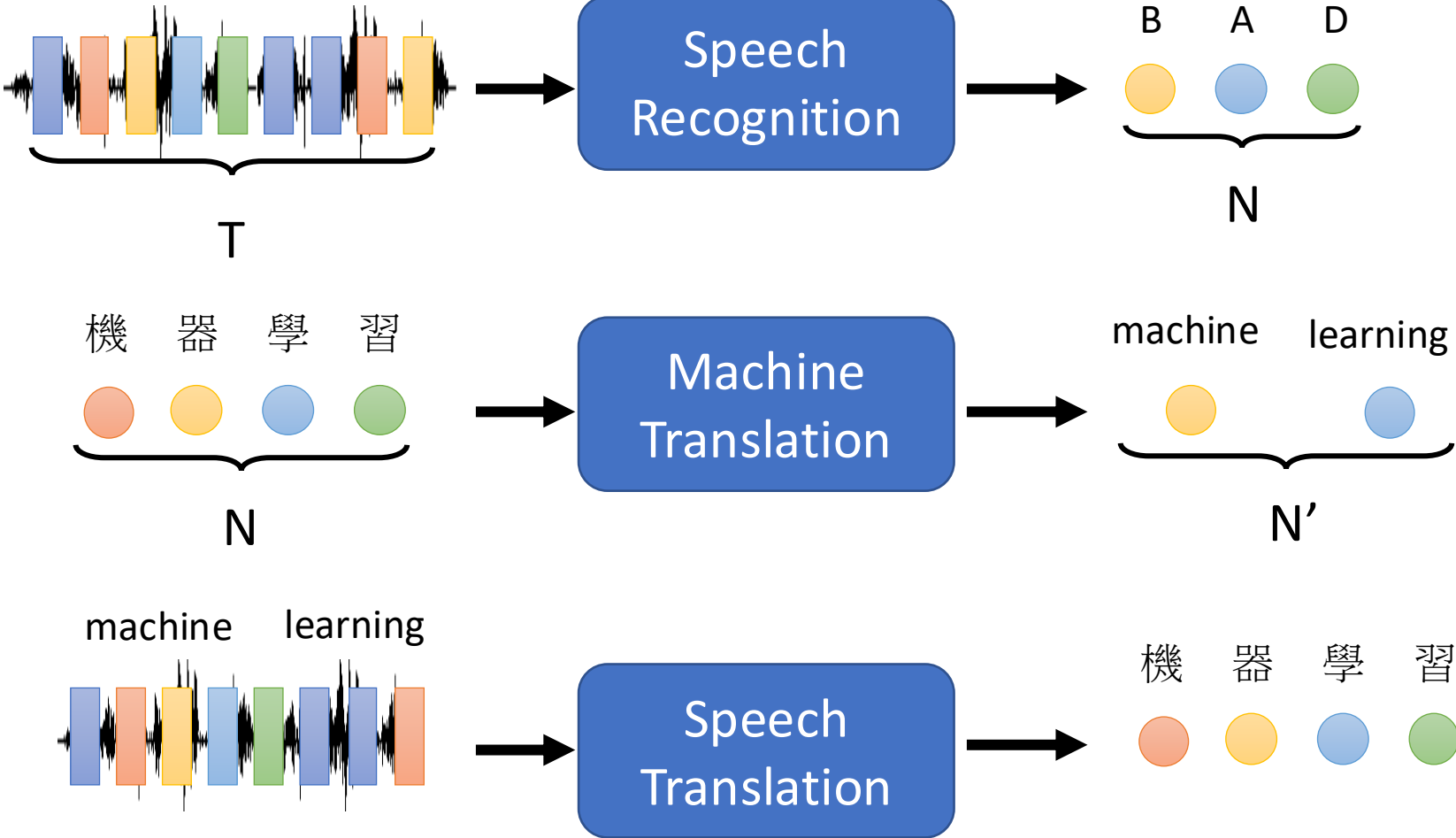
Why Transformer?

RNNs: slow sequential processing

LSTMs/GRUs: limited long-range memory (LLM has the same problem!!)

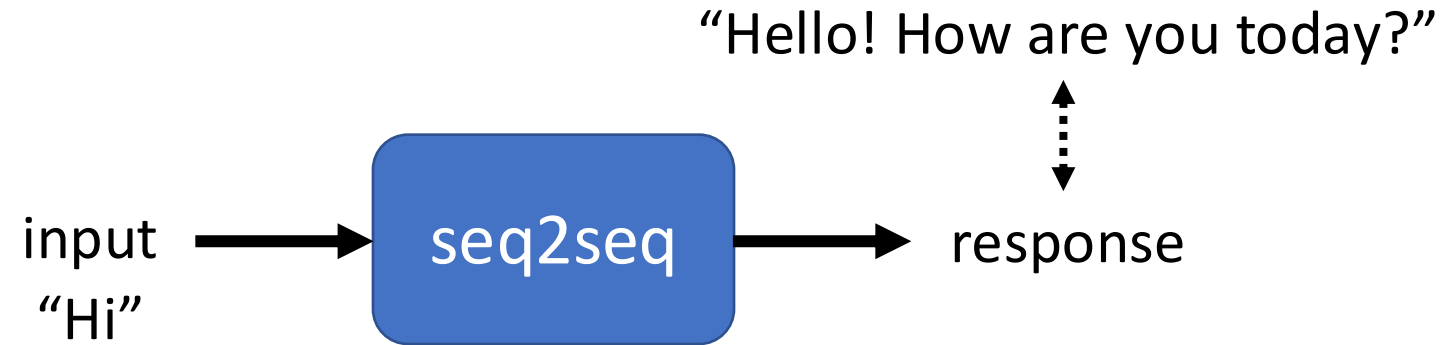
Sequence-to-sequence (Seq2seq)

The output length is determined by model.



Language without text

Seq2seq for Chatbot



Training data:

[PERSON 1:] Hi
[PERSON 2:] Hello ! How are you today ?
[PERSON 1:] I am good thank you , how are you.
[PERSON 2:] Great, thanks ! My children and I were just about to watch Game of Thrones.
[PERSON 1:] Nice ! How old are your children?
[PERSON 2:] I have four that range in age from 10 to 21. You?
[PERSON 1:] I do not have children at the moment.
[PERSON 2:] That just means you get to keep all the popcorn for yourself.
[PERSON 1:] And Cheetos at the moment!
[PERSON 2:] Good choice. Do you watch Game of Thrones?
[PERSON 1:] No, I do not have much time for TV.
[PERSON 2:] I usually spend my time painting: but, I love the show.

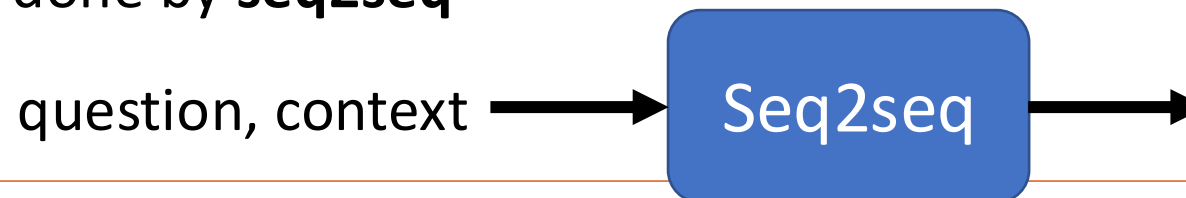
Most Natural Language Processing applications ...

Question Answering (QA)

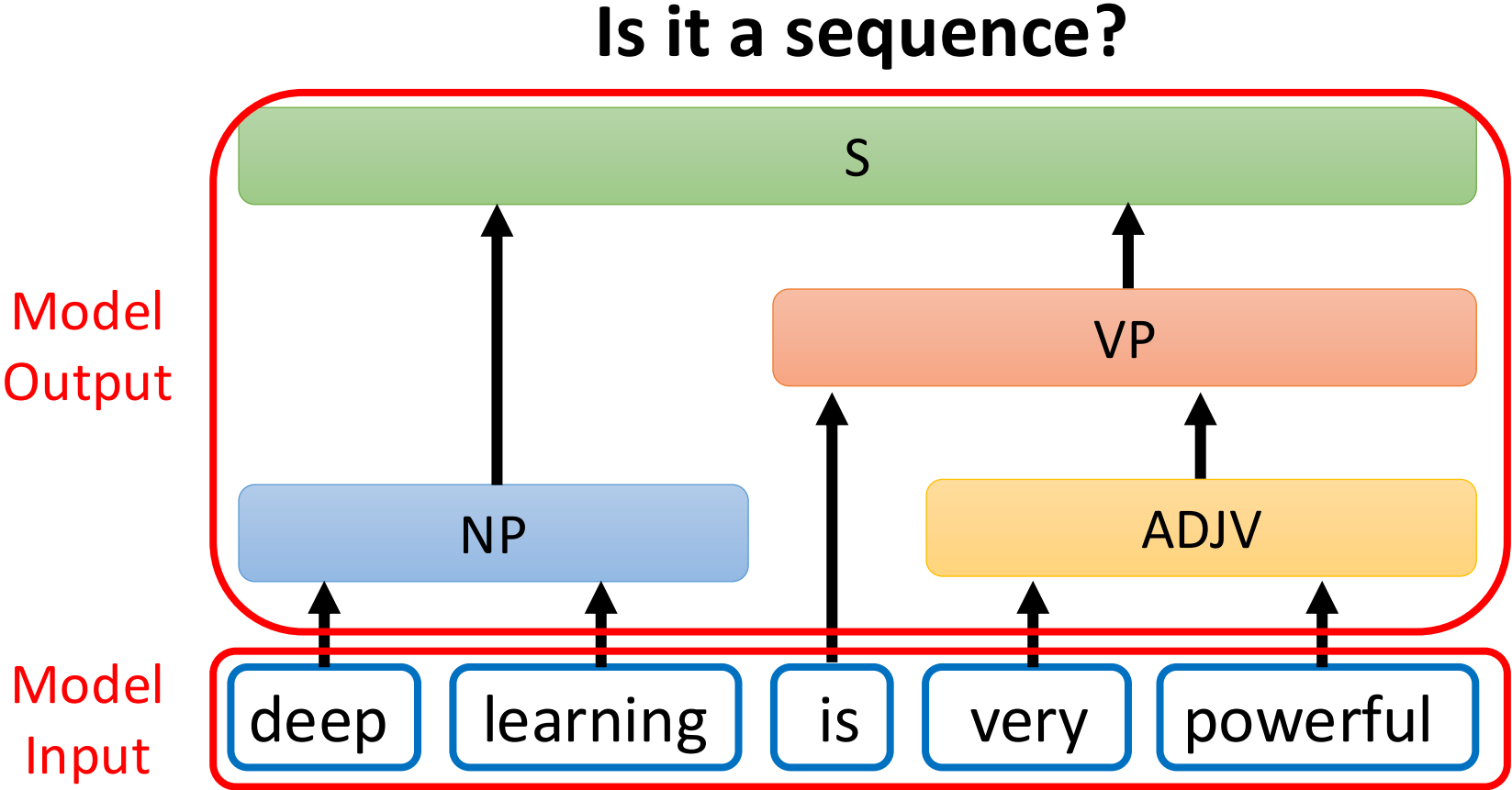
<u>Question</u>	<u>Context</u>	<u>Answer</u>
What is a major importance of Southern California in relation to California and the US?	...Southern California is a major economic center for the state of California and the US....	major economic center
What is the translation from English to German?	Most of the planet is ocean water.	Der Großteil der Erde ist Meerwasser
What is the summary?	Harry Potter star Daniel Radcliffe gains access to a reported £320 million fortune ...	Harry Potter star Daniel Radcliffe gets £320M fortune ...
Hypothesis: Product and geography are what make cream skimming work. Entailment , neutral, or contradiction?	Premise: Conceptually cream skimming has two basic dimensions – product and geography.	Entailment
Is this sentence positive or negative? (sentiment analysis)	A stirring, funny and finally transporting re-imagining of Beauty and the Beast and 1930s horror film.	positive



QA can be done by seq2seq



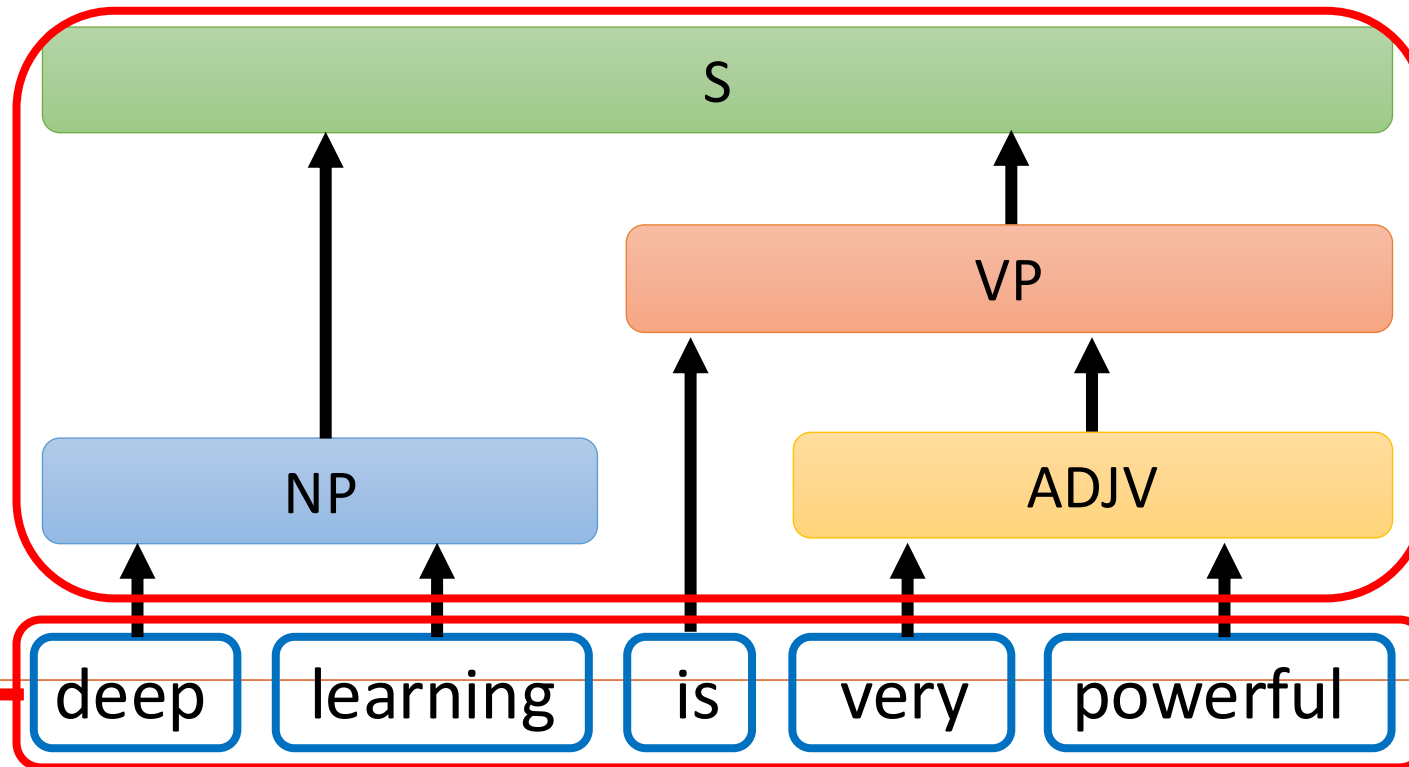
Seq2seq for Syntactic Parsing



Seq2seq for Syntactic Parsing

(S (NP deep learning) (VP is
(ADJV very powerful)))

Seq2seq!



Seq2seq for Syntactic Parsing

(S (NP deep learning) (VP is
(ADJV very powerful)))

Grammar as a Foreign Language

Oriol Vinyals*
Google
vinyals@google.com

Lukasz Kaiser*
Google
lukaszkaizer@google.com

Terry Koo
Google
terrykoo@google.com

Slav Petrov
Google
slav@google.com

Ilya Sutskever
Google
ilyasu@google.com

Geoffrey Hinton
Google
geoffhinton@google.com

<https://arxiv.org/abs/1412.7449>

deep

learning

is

very

powerful

Seq2seq for Multi-Label Classification

An object can belong to multiple classes.



Class 1
Class 3



Class 1



Class 3
Class 9
Class 17



Class 10



Class 9



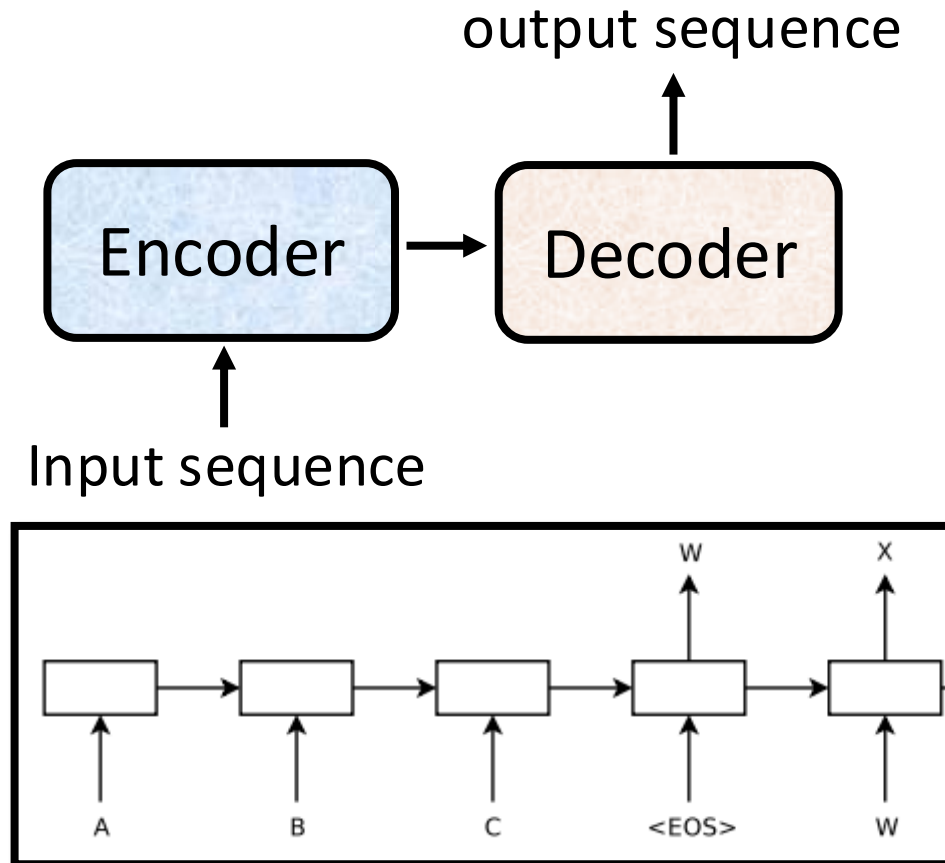
Class 7



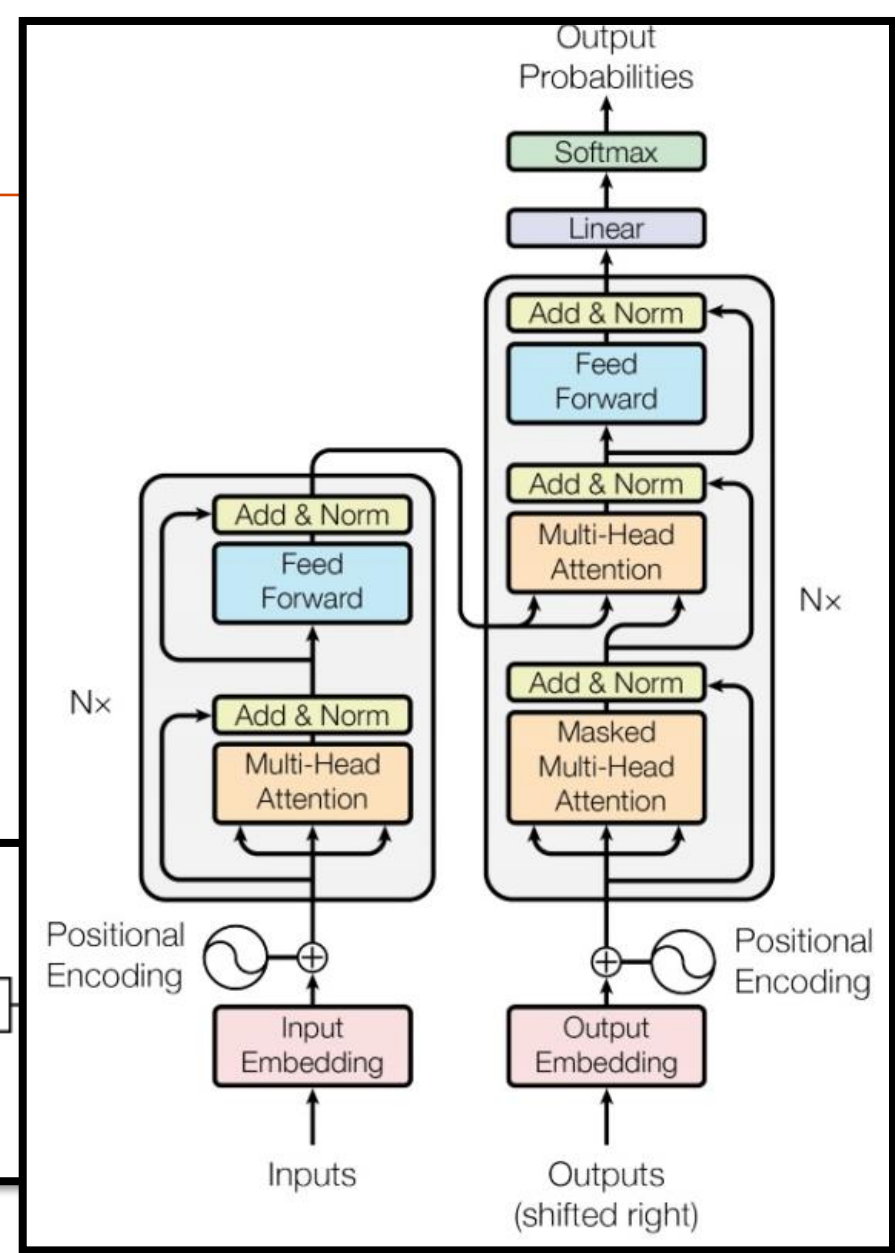
Class 13

<https://arxiv.org/abs/1909.03434>
<https://arxiv.org/abs/1707.05495>

Seq2seq

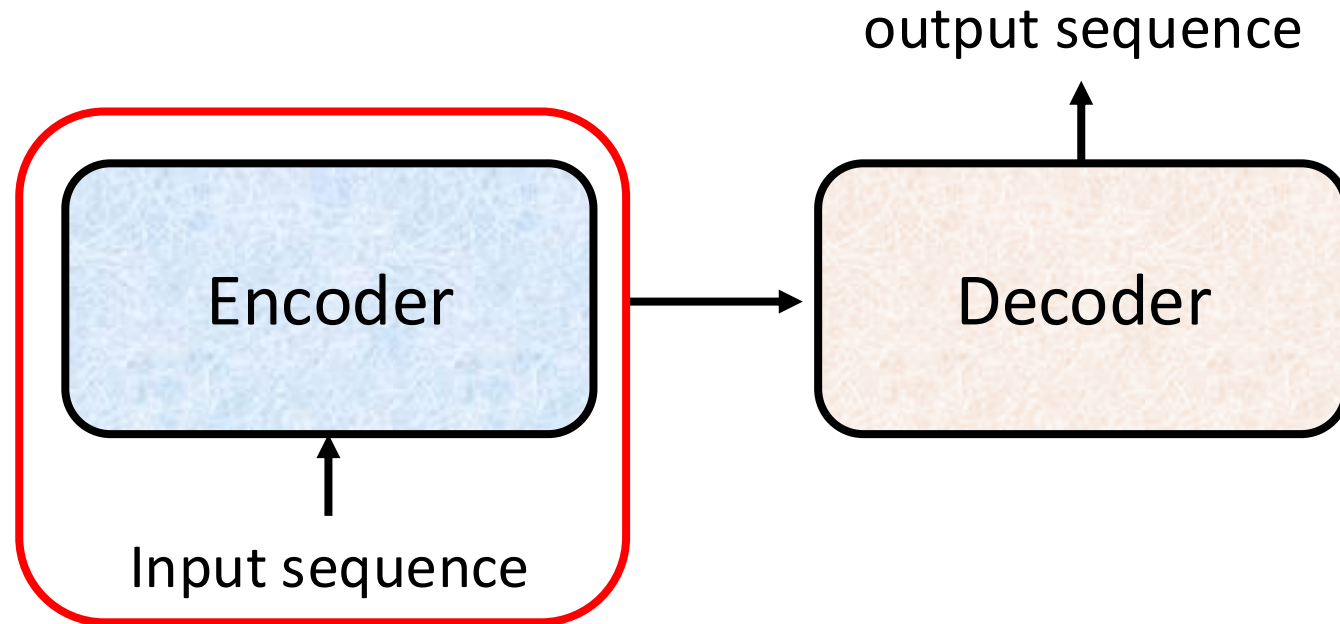


Sequence to Sequence Learning with Neural Networks



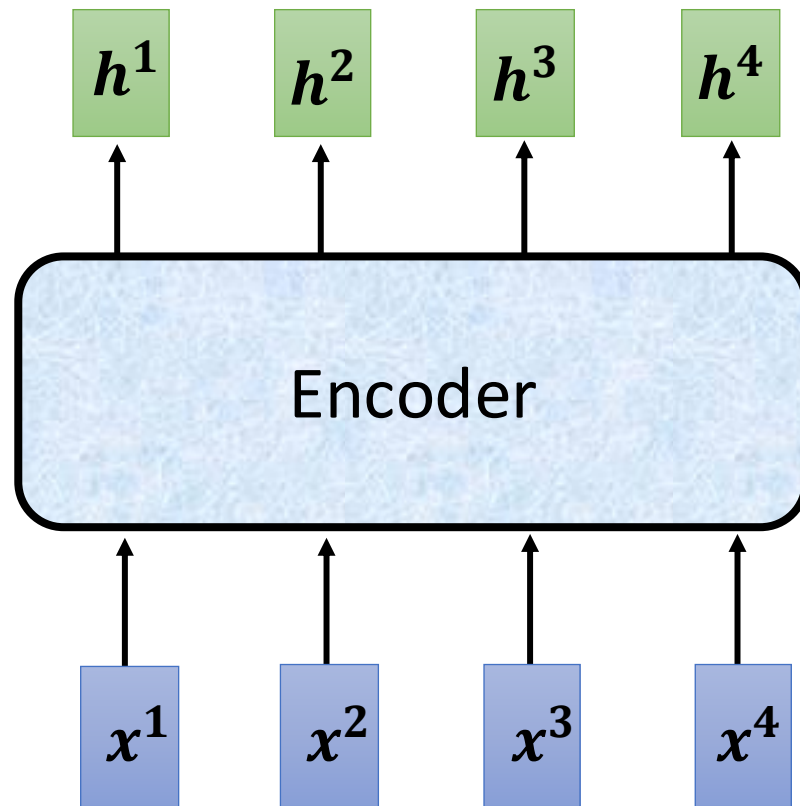
Transformer

Encoder

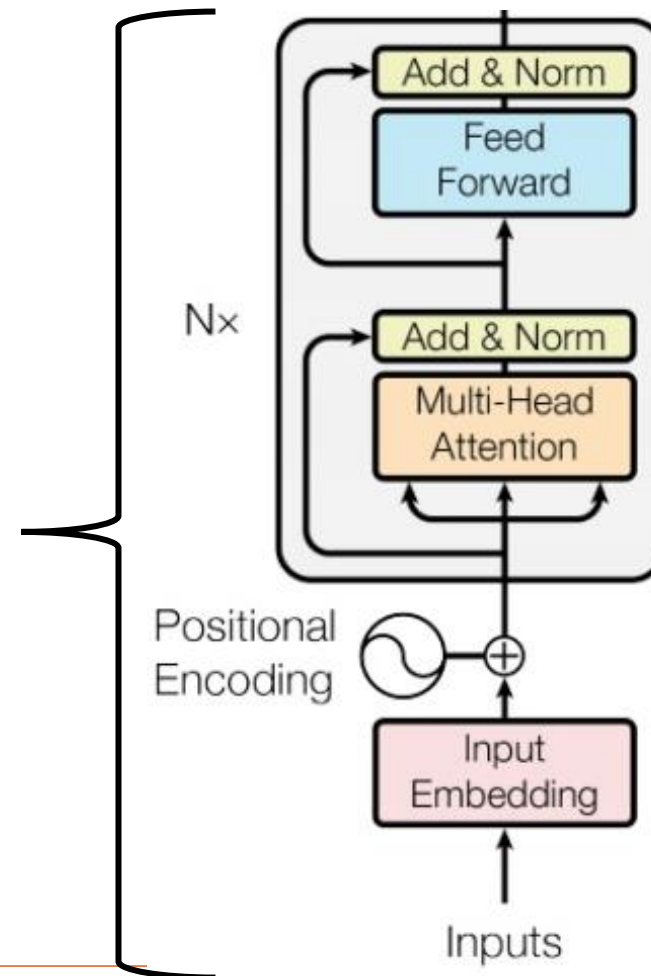


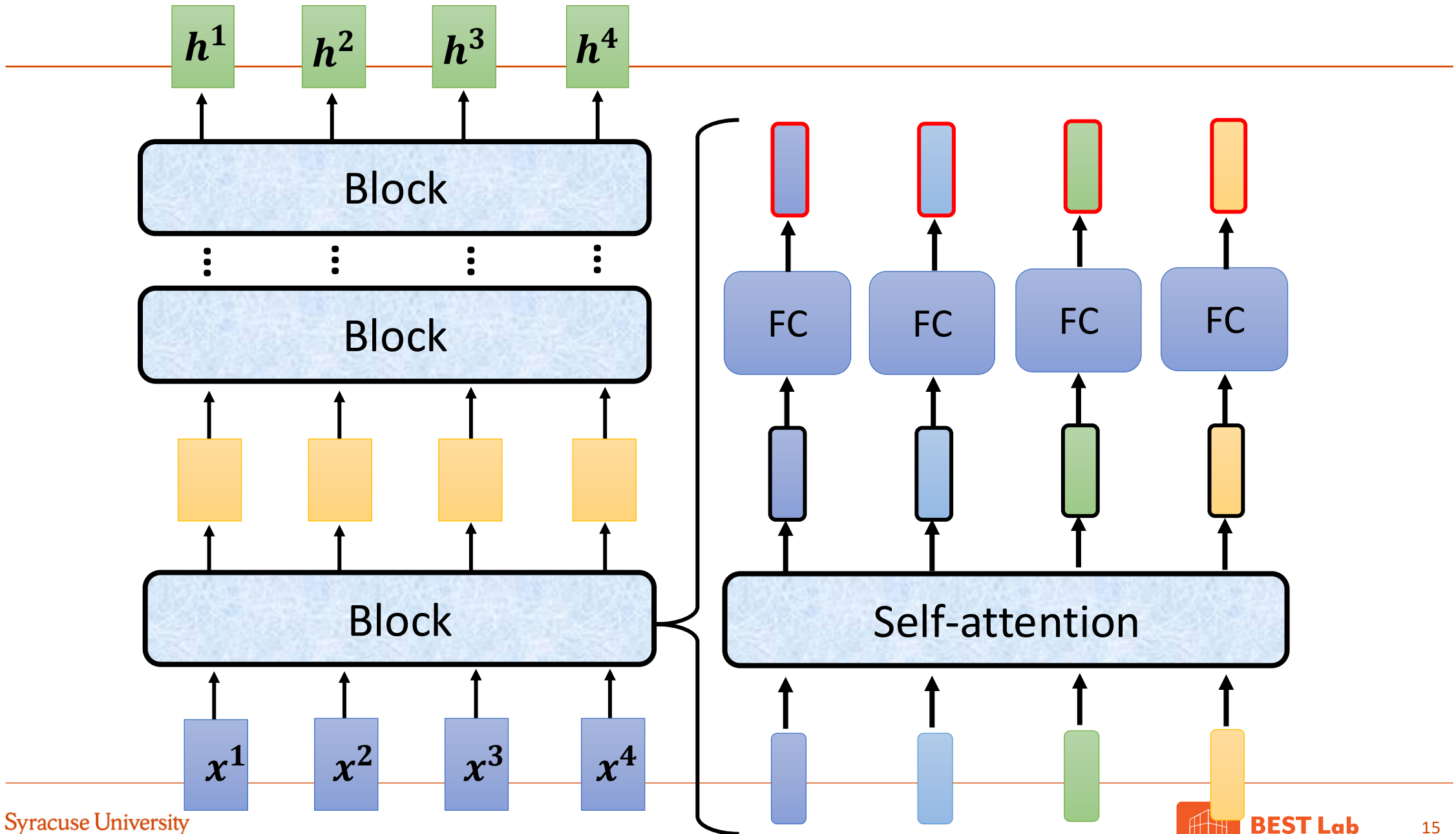
Encoder

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

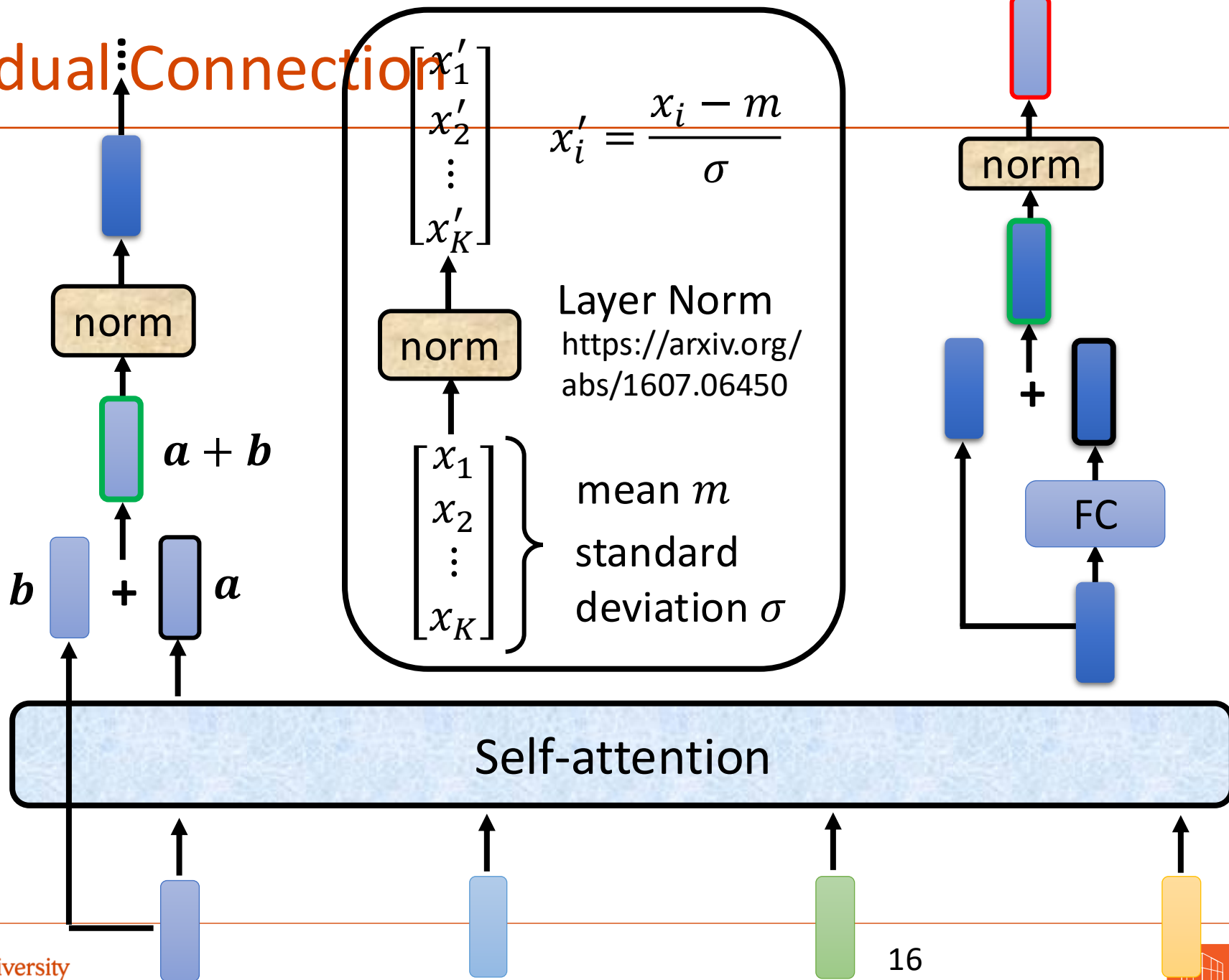


Transformer's Encoder

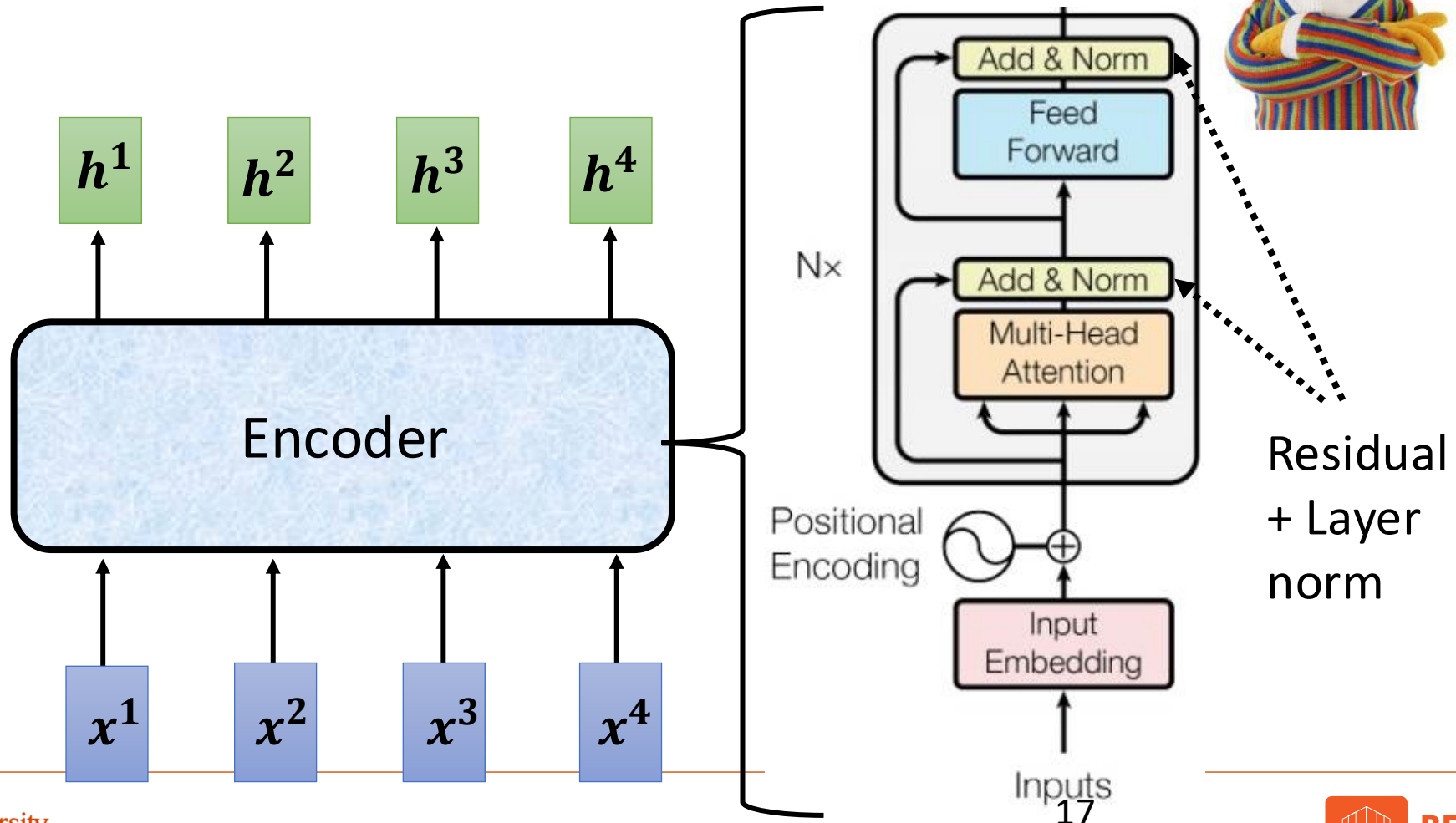




Residual Connection

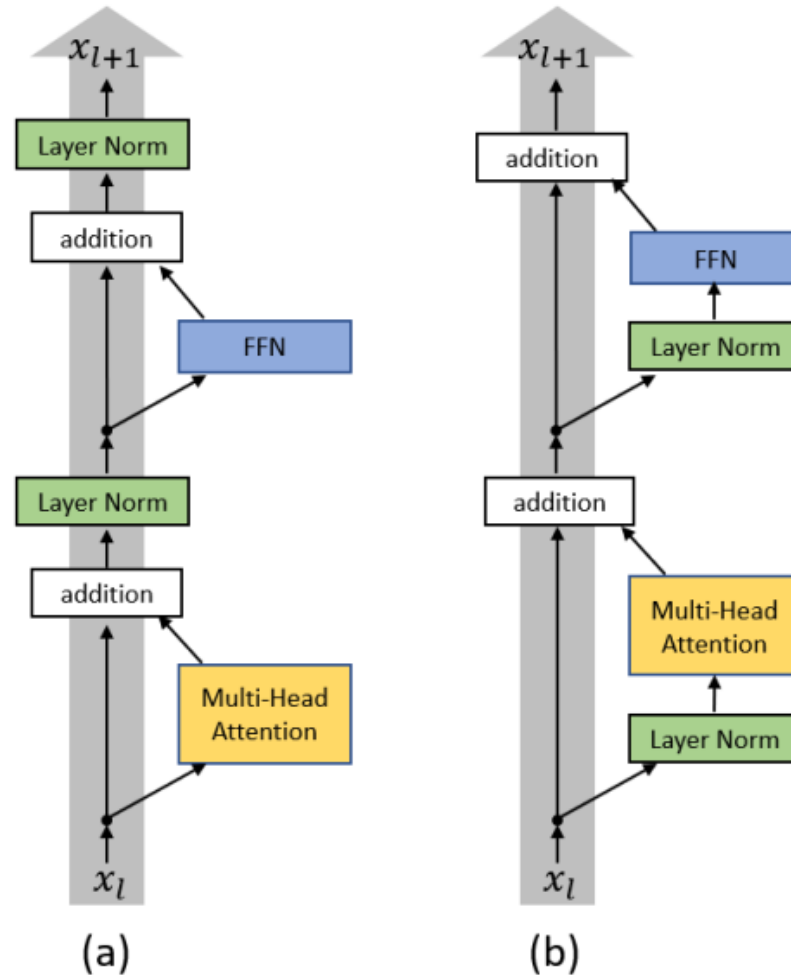


I use the **same** network architecture as **transformer encoder**.

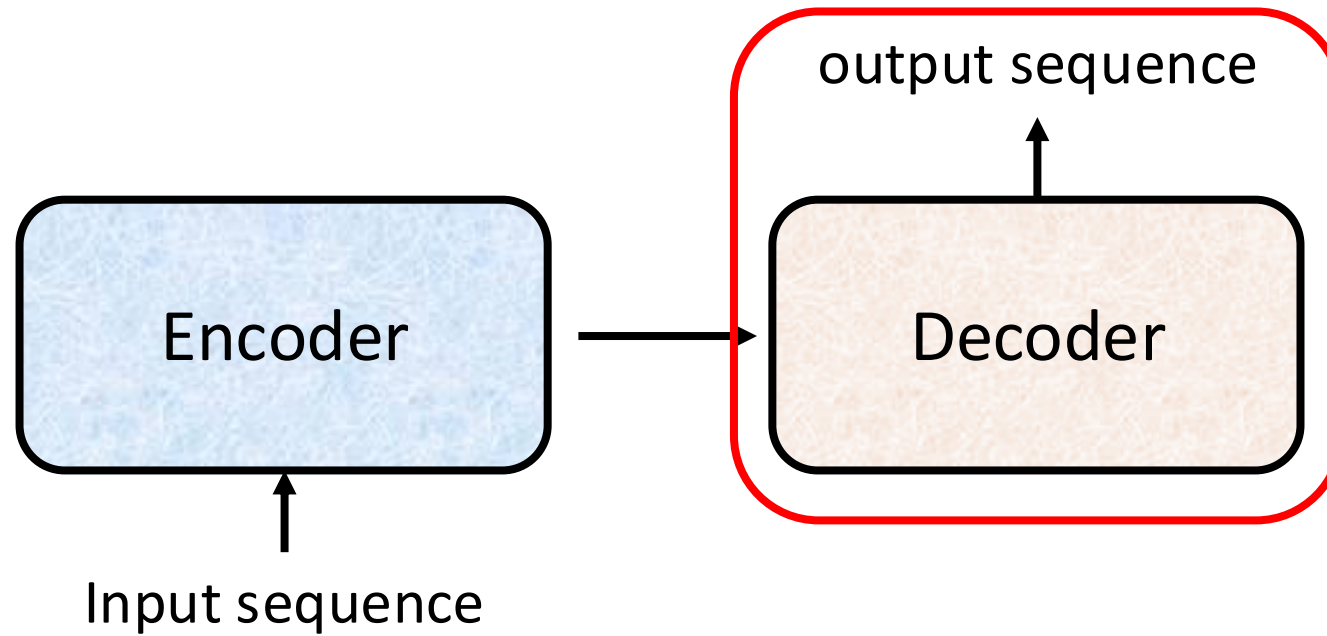


To learn more

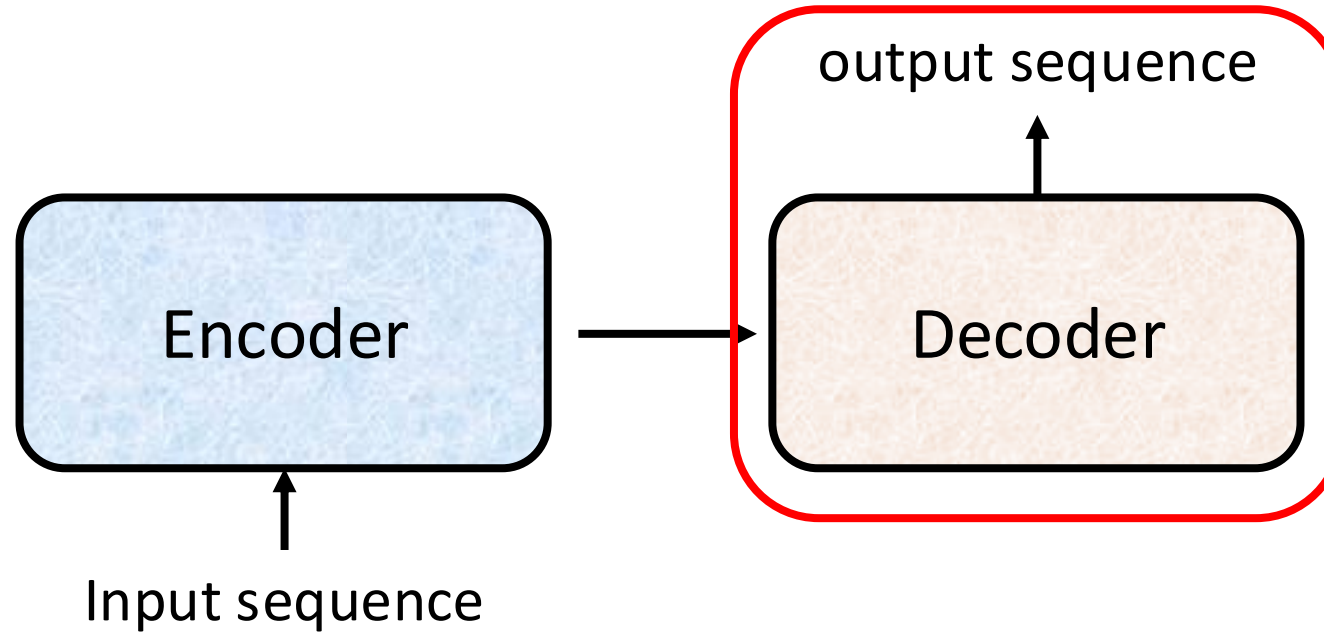
- On Layer Normalization in the Transformer Architecture
- <https://arxiv.org/abs/2002.04745>
- PowerNorm: Rethinking Batch Normalization in Transformers
- <https://arxiv.org/abs/2003.07845>



Decoder

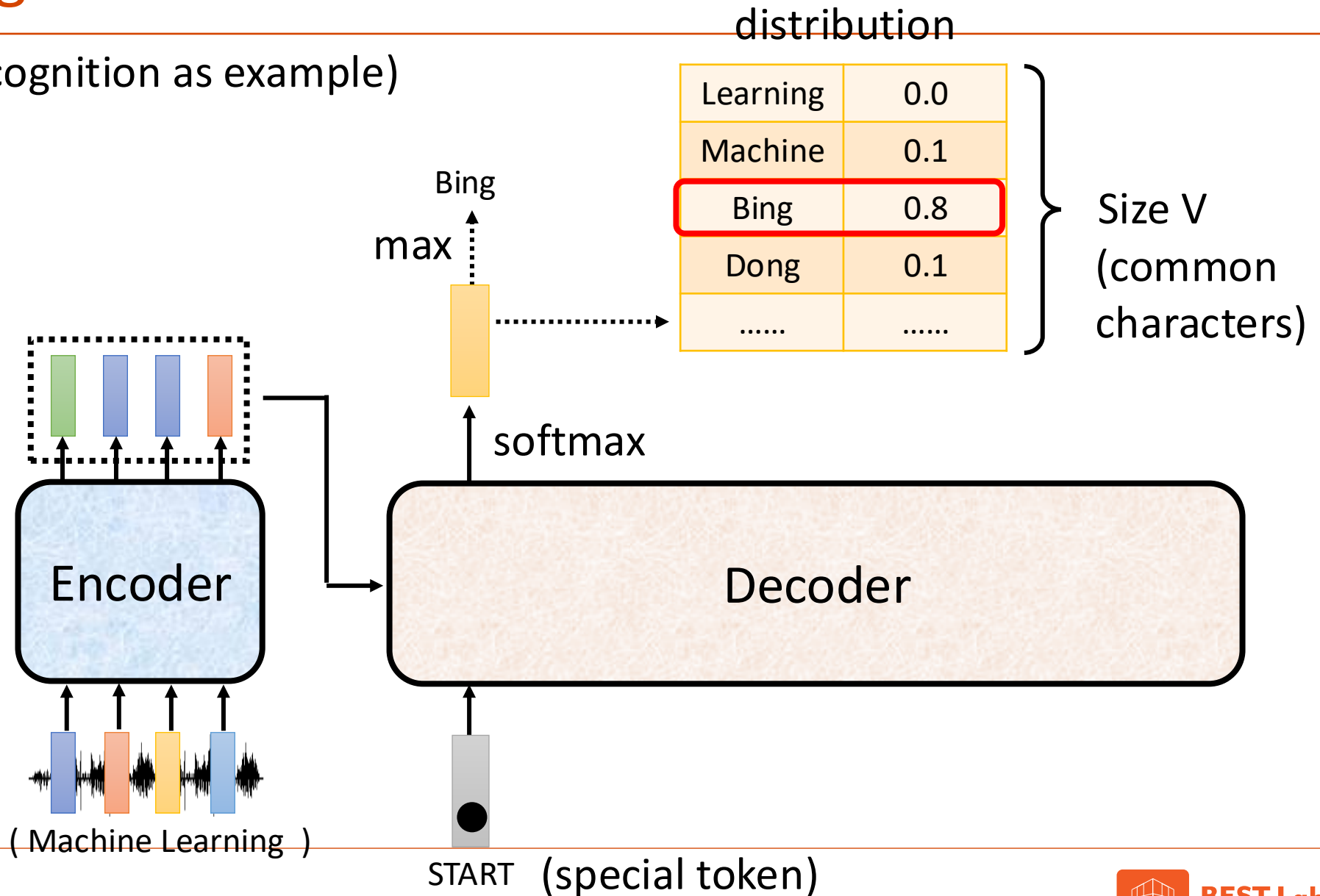


Decoder – Autoregressive (AT)

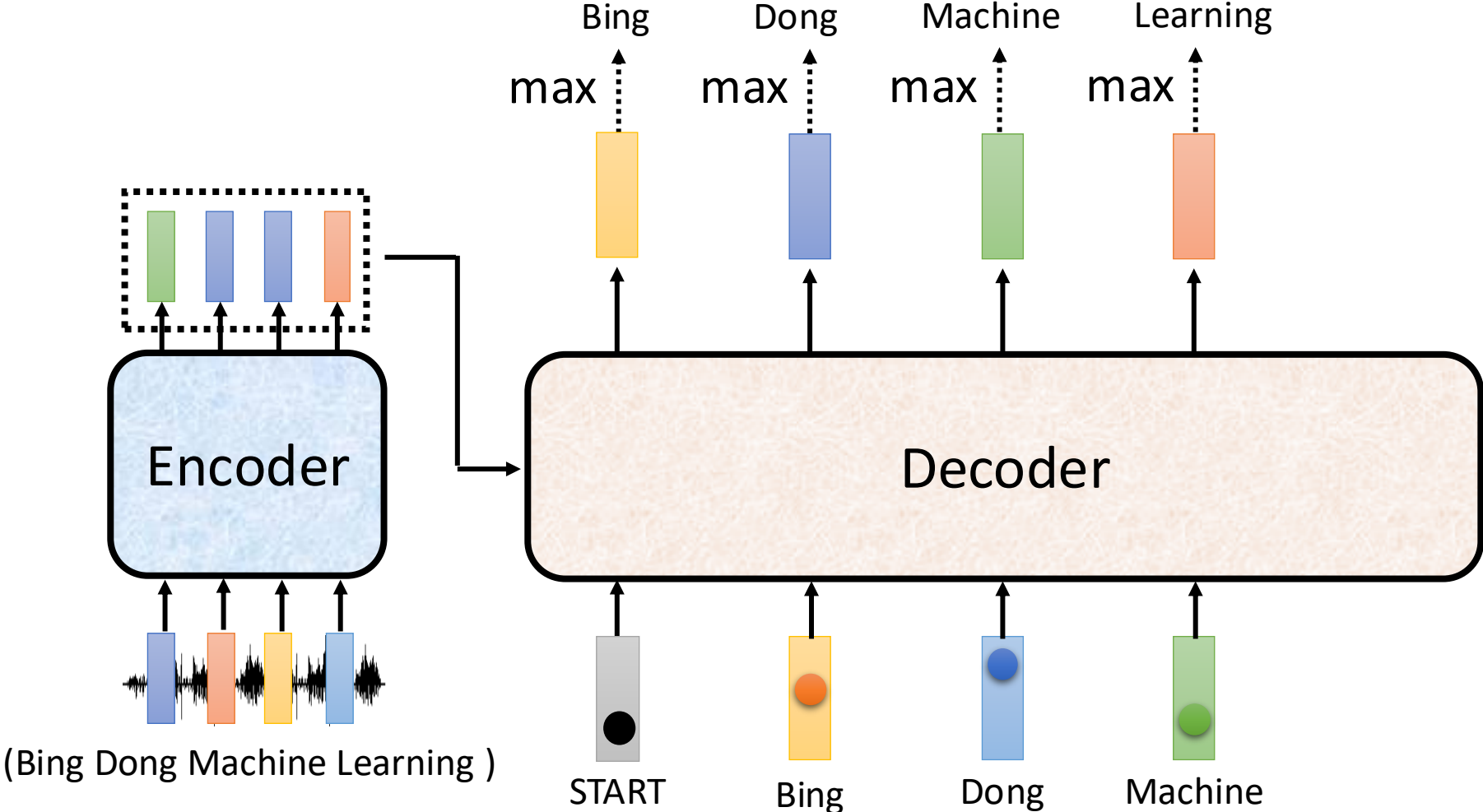


Autoregressive

(Speech Recognition as example)

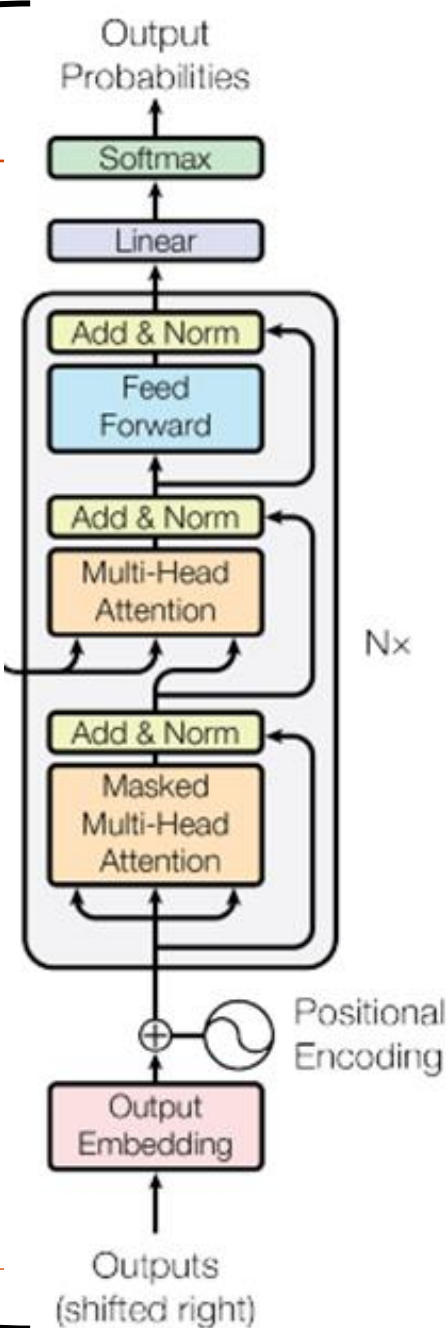
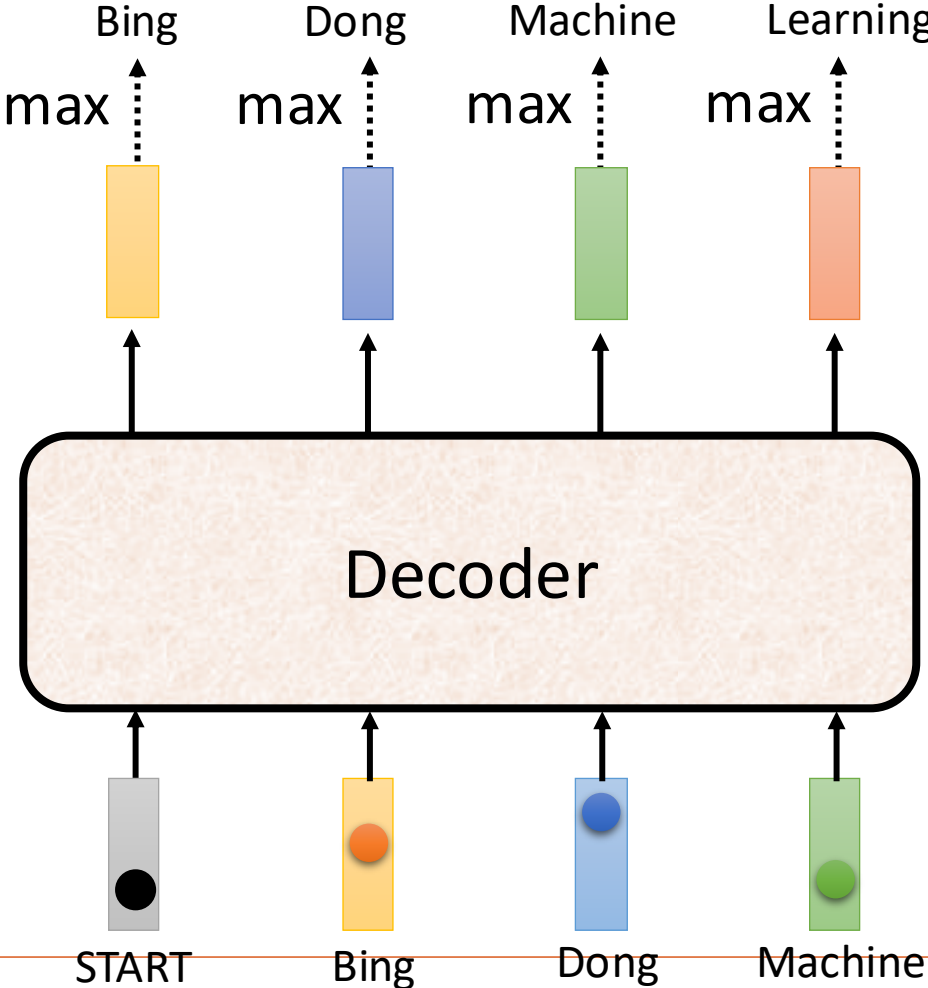


Autoregressive

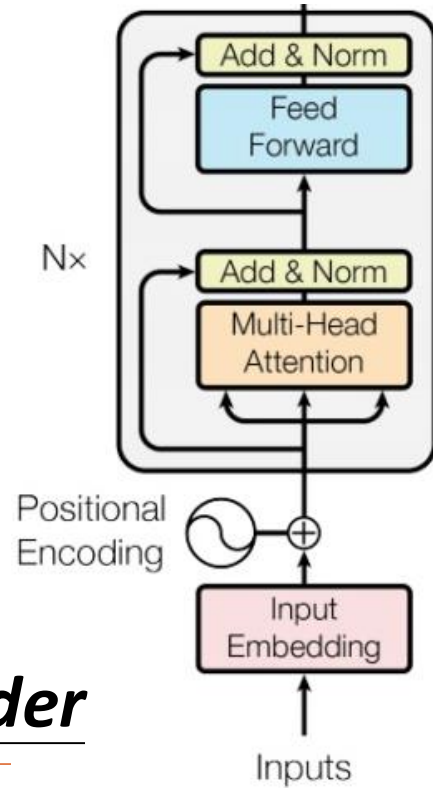


Autoregressive

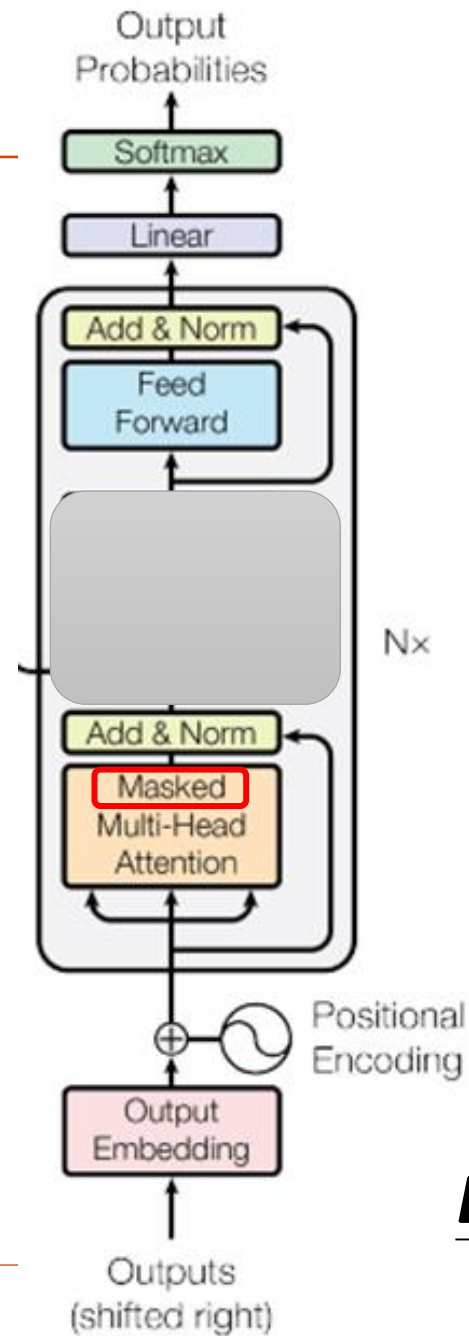
ignore the input from the encoder here 😊



Encoder

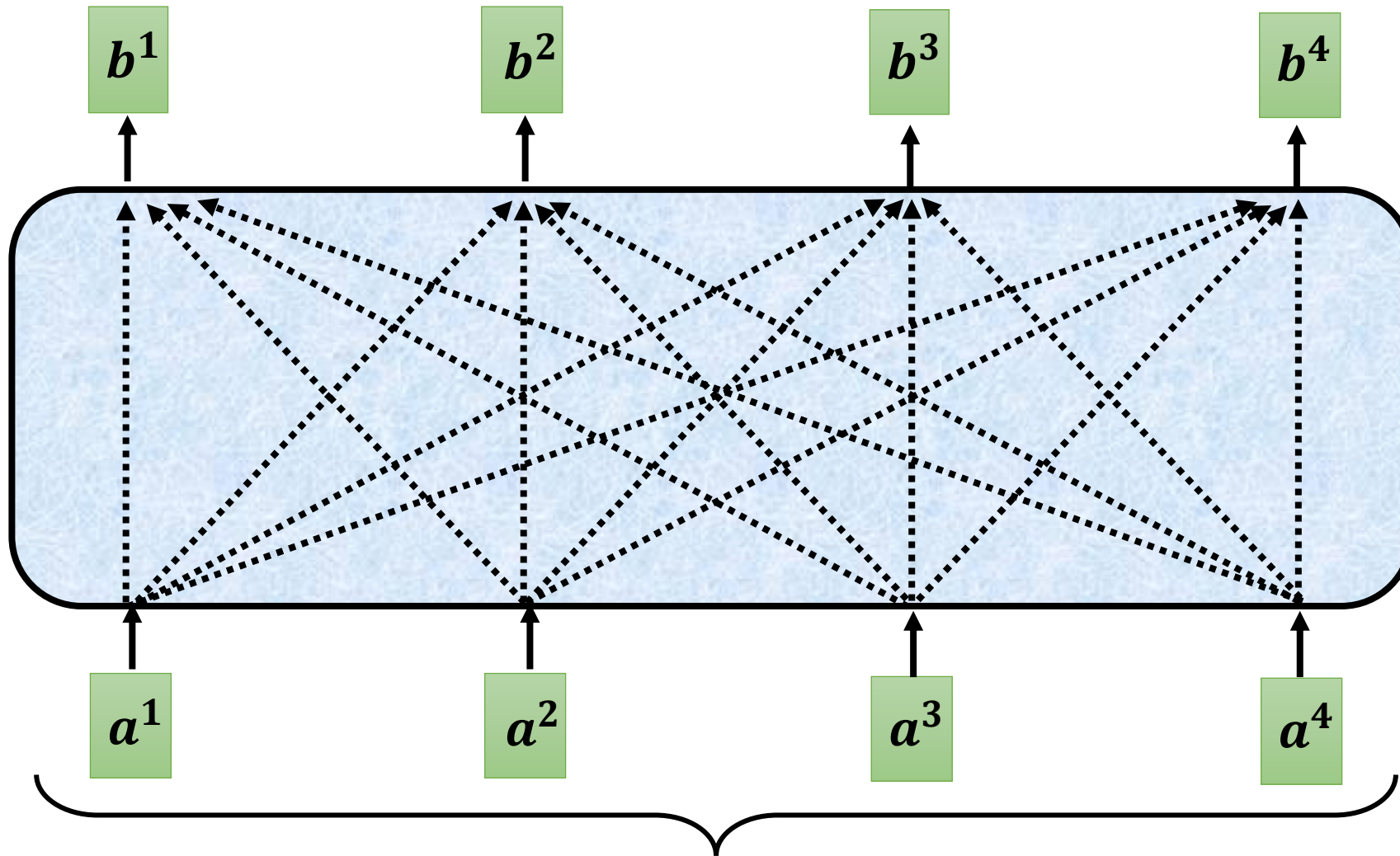


Decoder



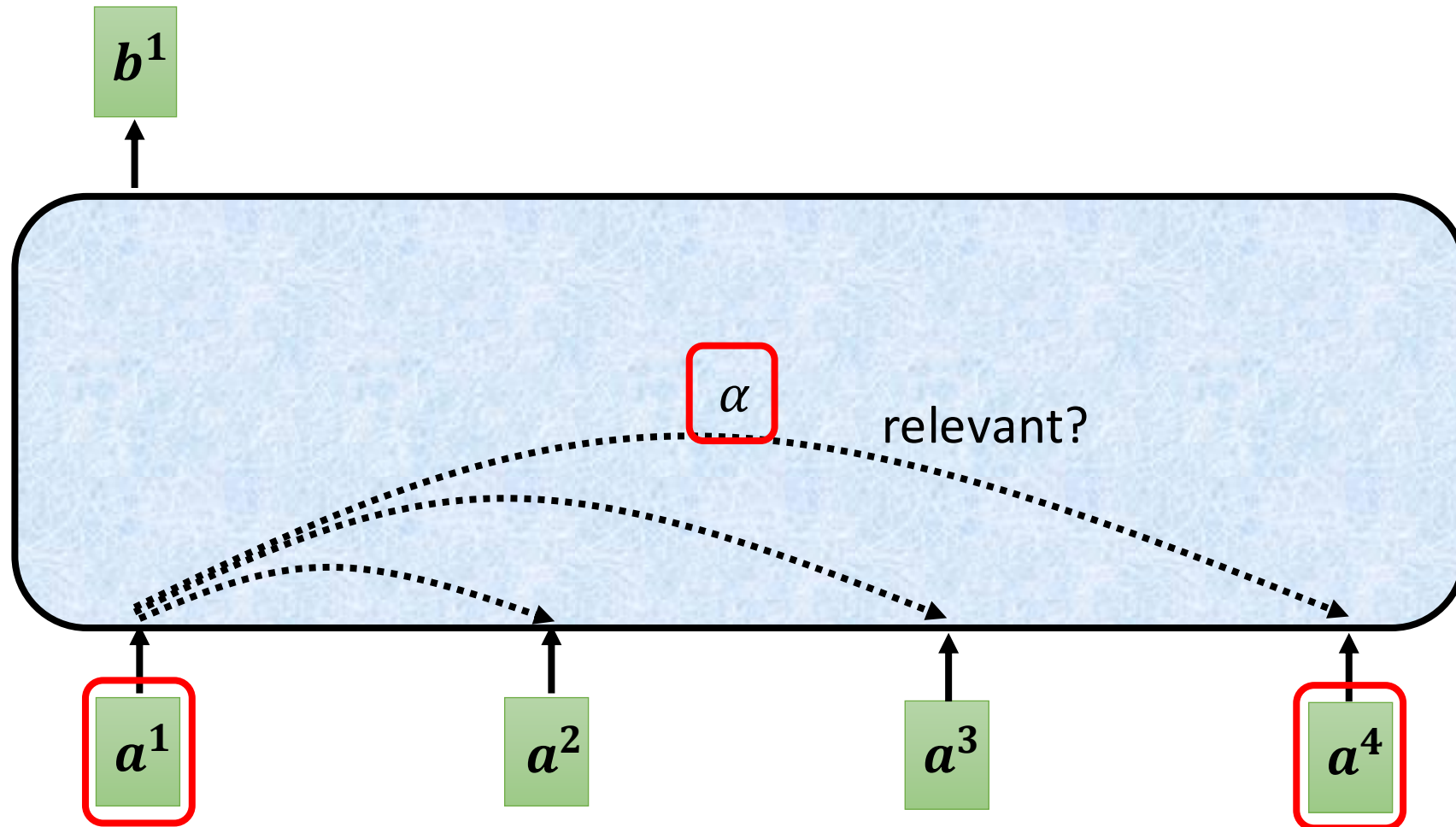
What's inside Self-attention layer?

Self-attention



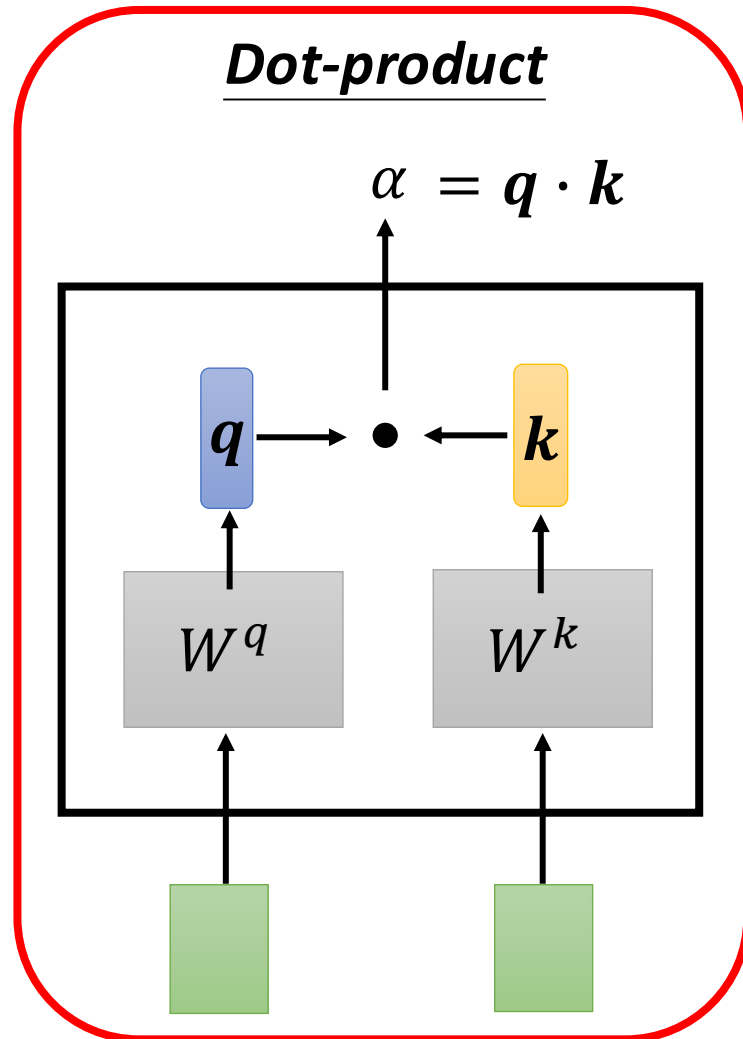
Can be either **input** or a **hidden layer**

Self-attention

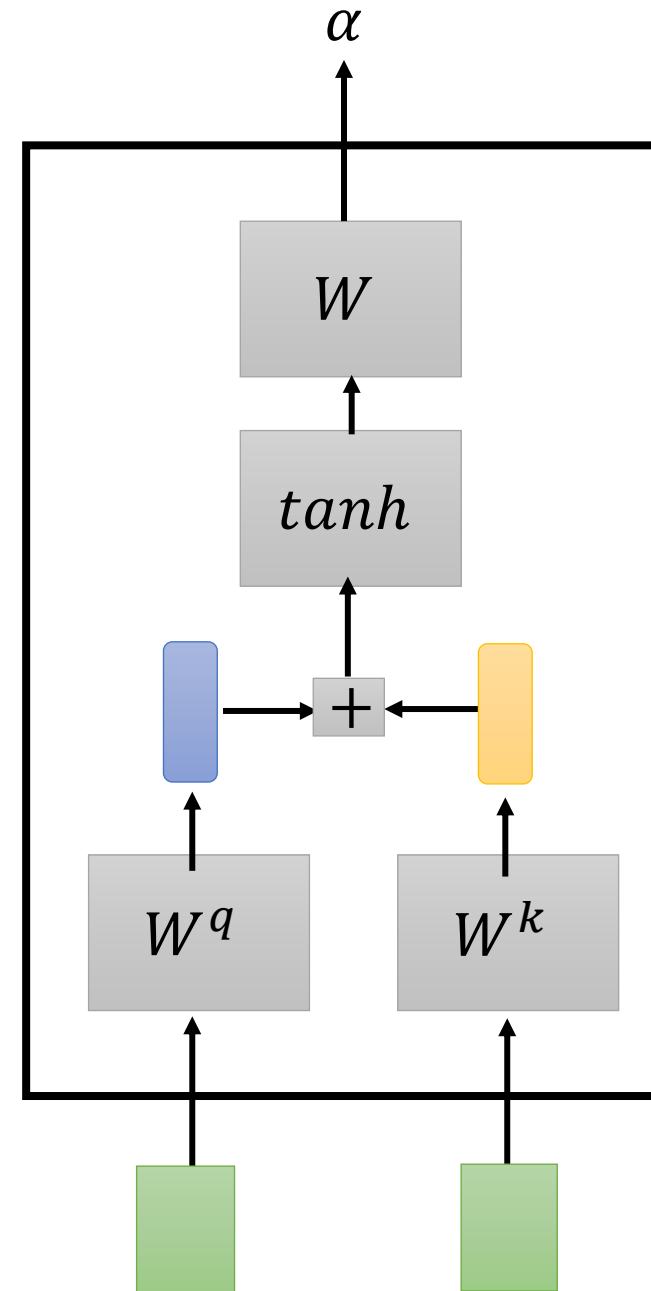


Find the relevant vectors in a sequence

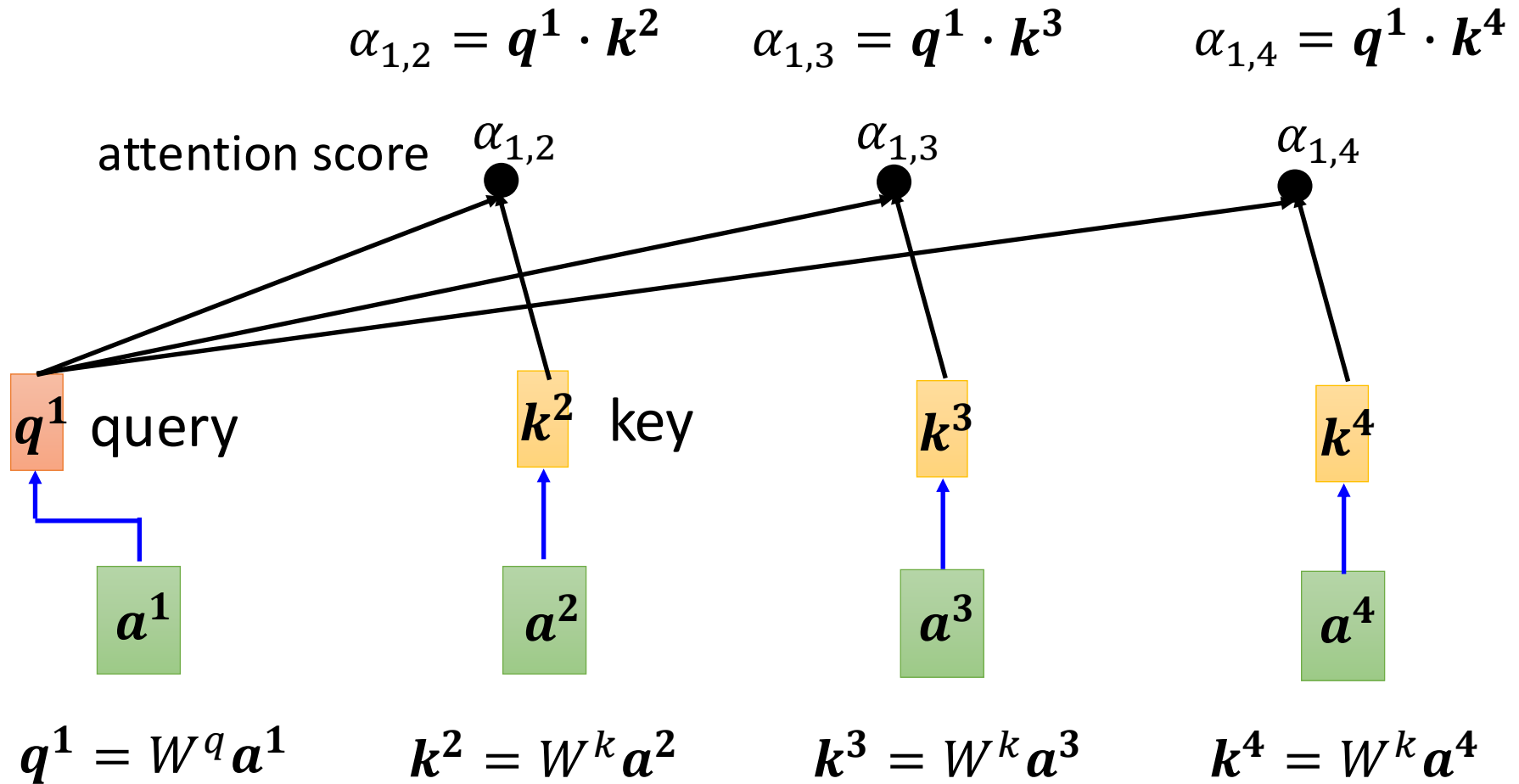
Self-attention



Additive

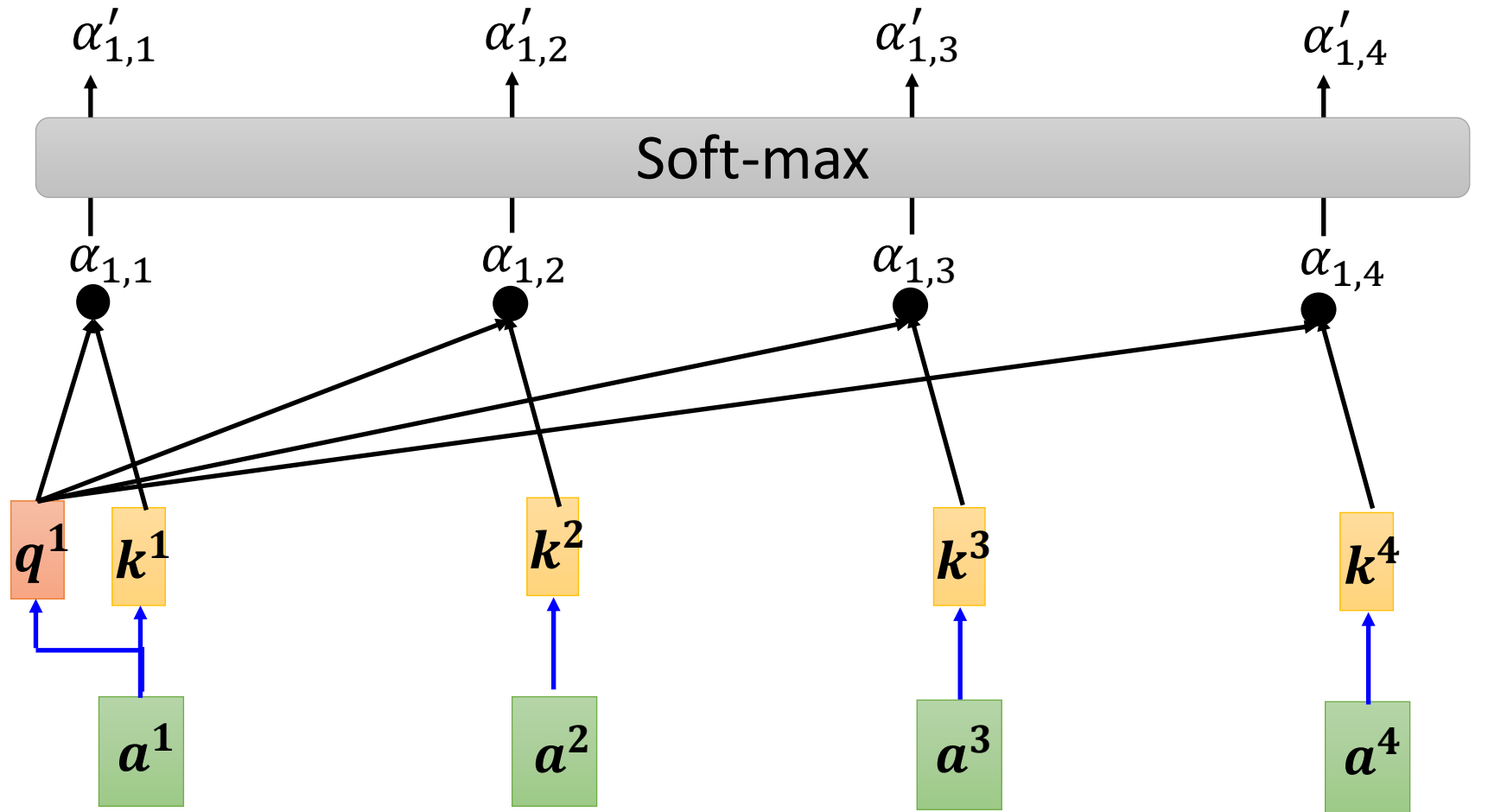


Self-attention



Self-attention

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



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

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

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

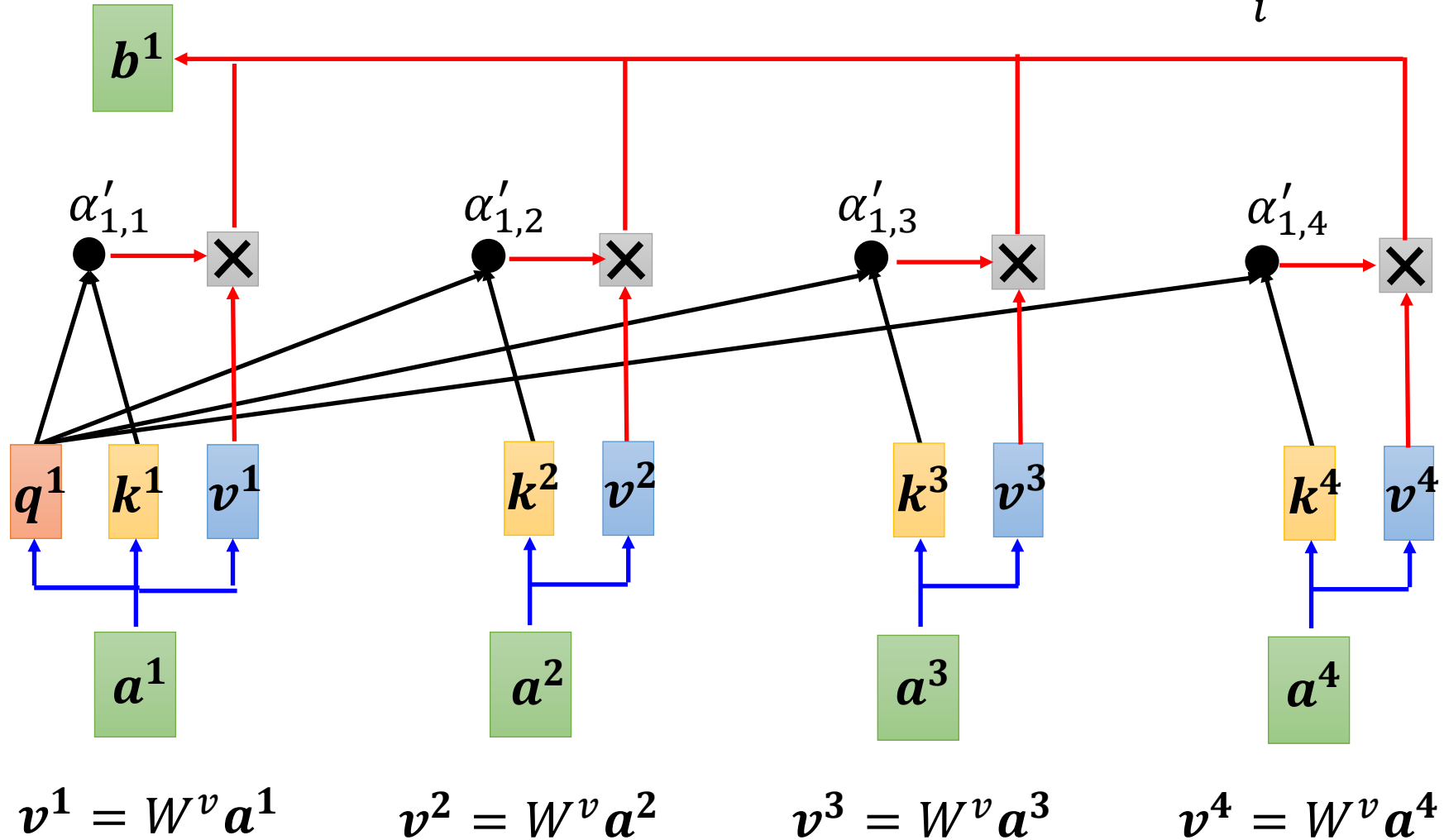
$$k^3 = W^k a^3$$

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

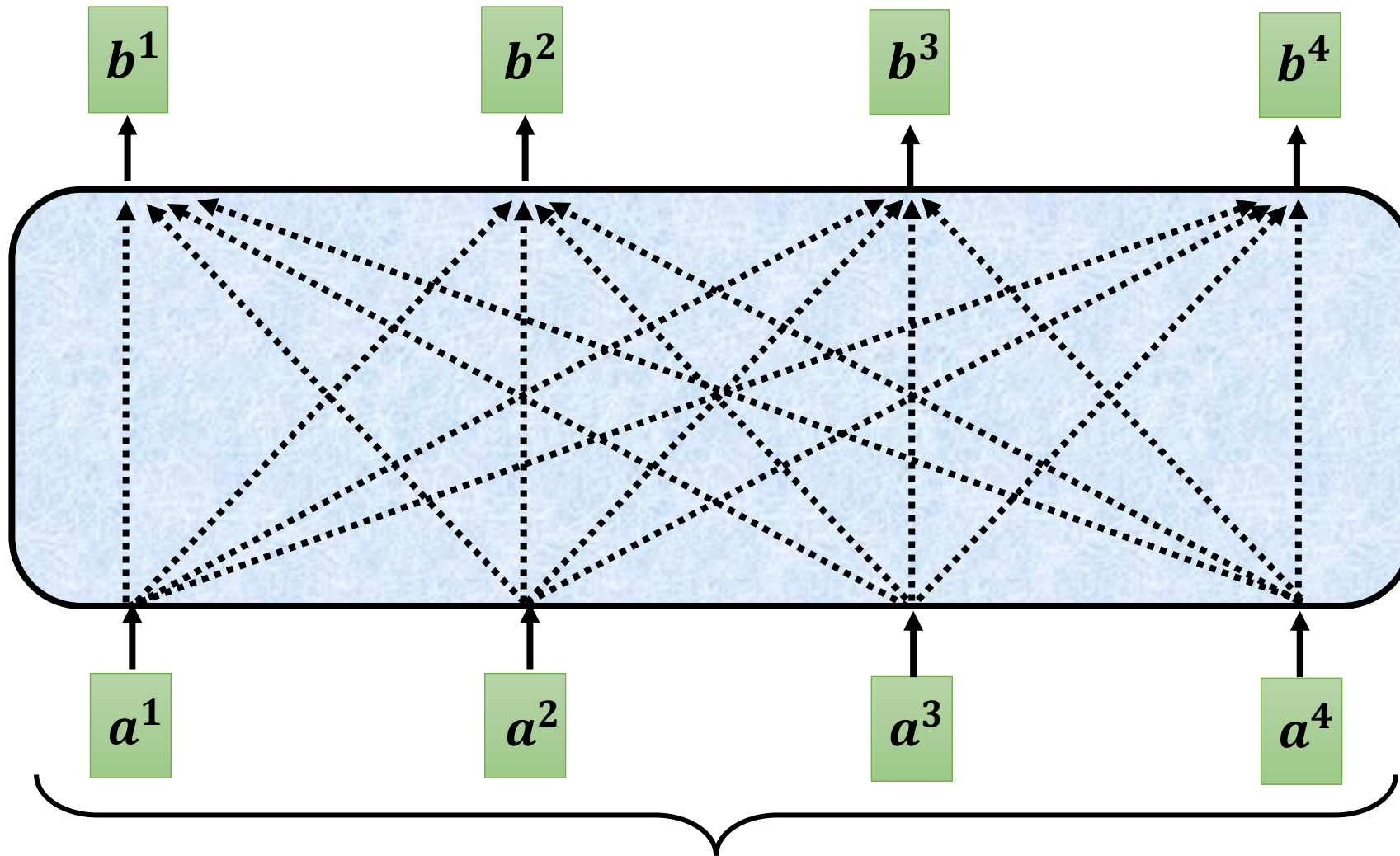
Self-attention

Extract information based on attention scores

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



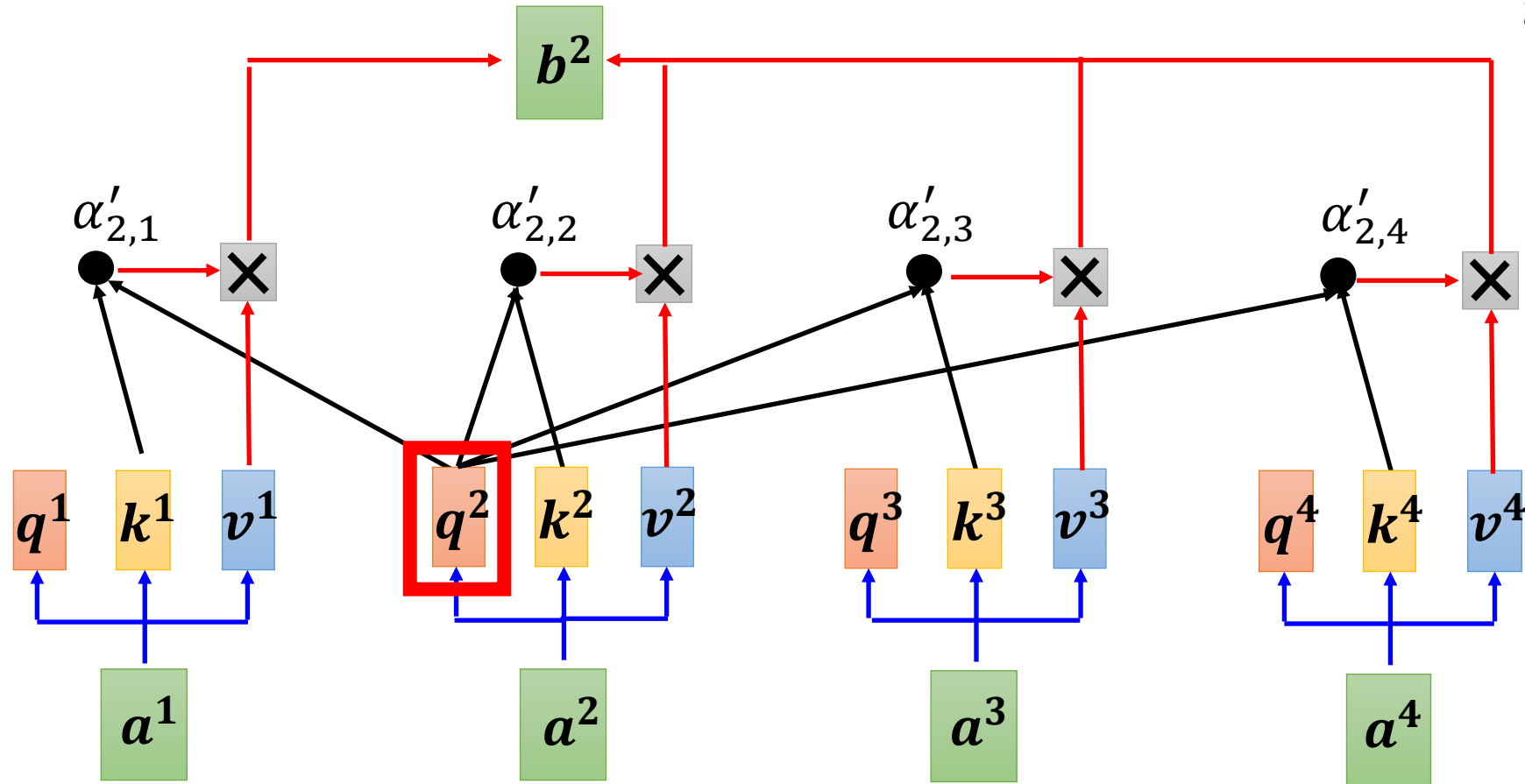
Self-attention



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

Self-attention → Masked Self-attention

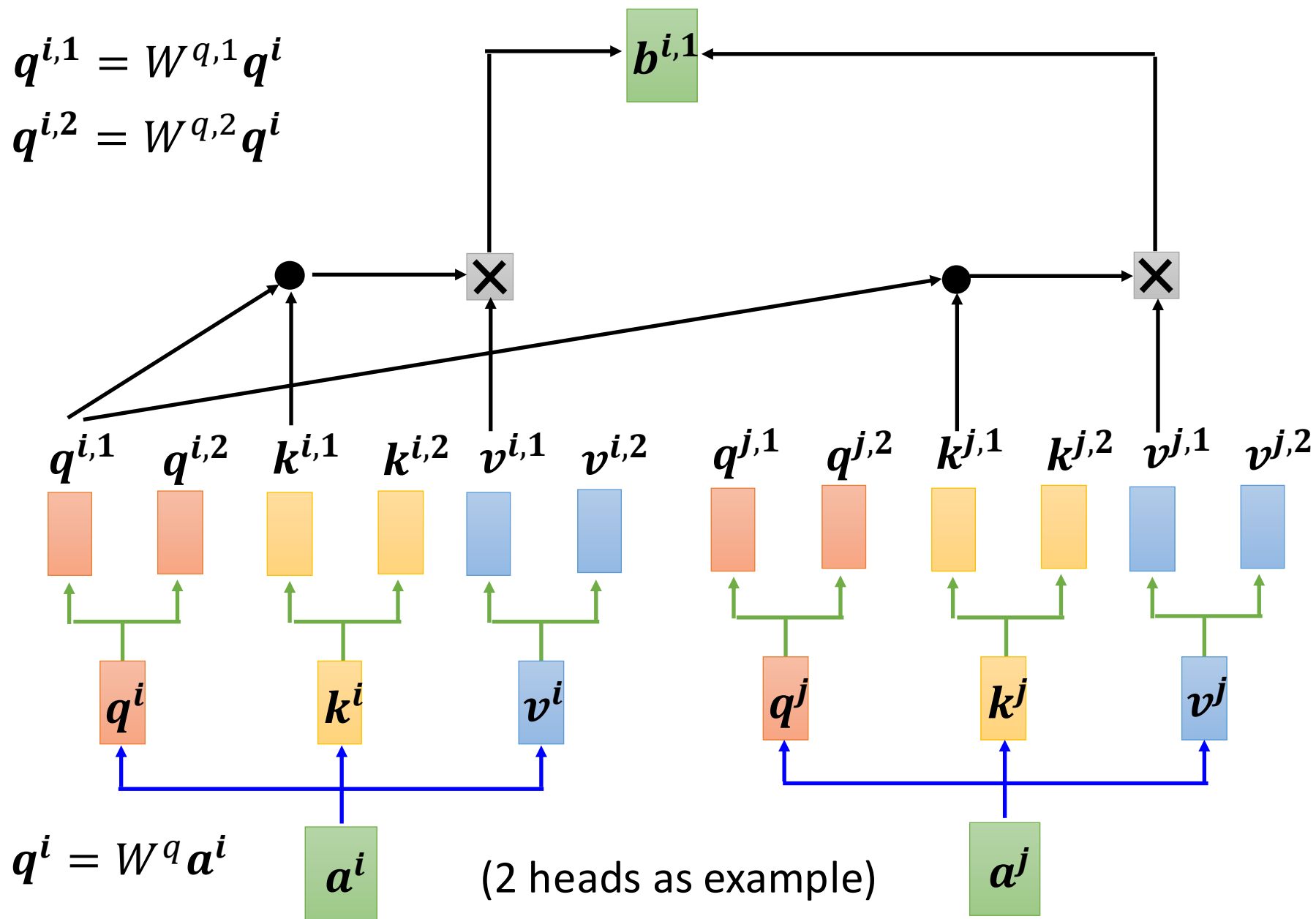
$$b^2 = \sum_i \alpha'_{2,i} v^i$$



Why masked? Consider how does decoder work

Multi-head Self-attention

Different types of relevance

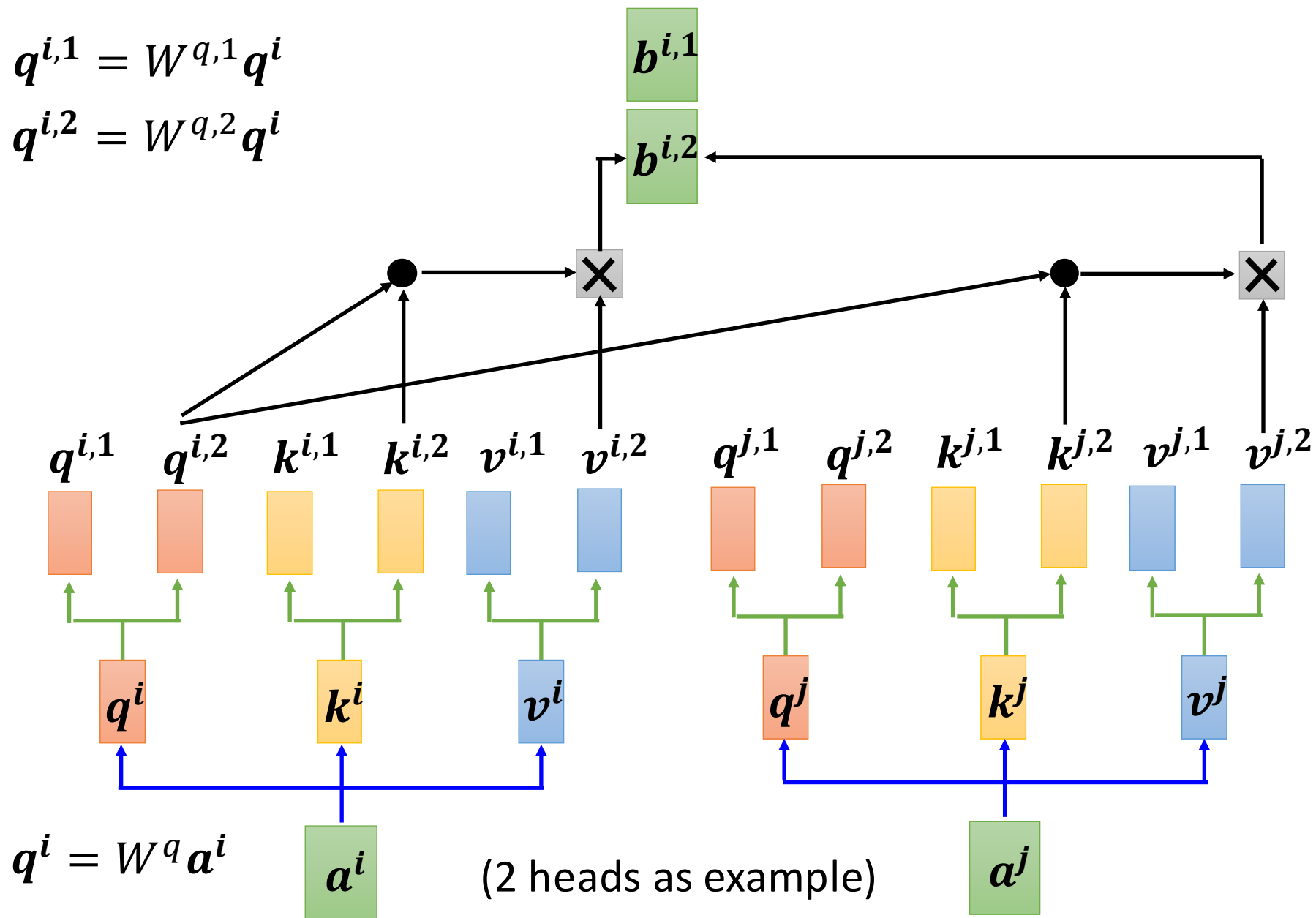


Multi-head Self-attention

Different types of relevance

$$q^{i,1} = W^{q,1} q^i$$

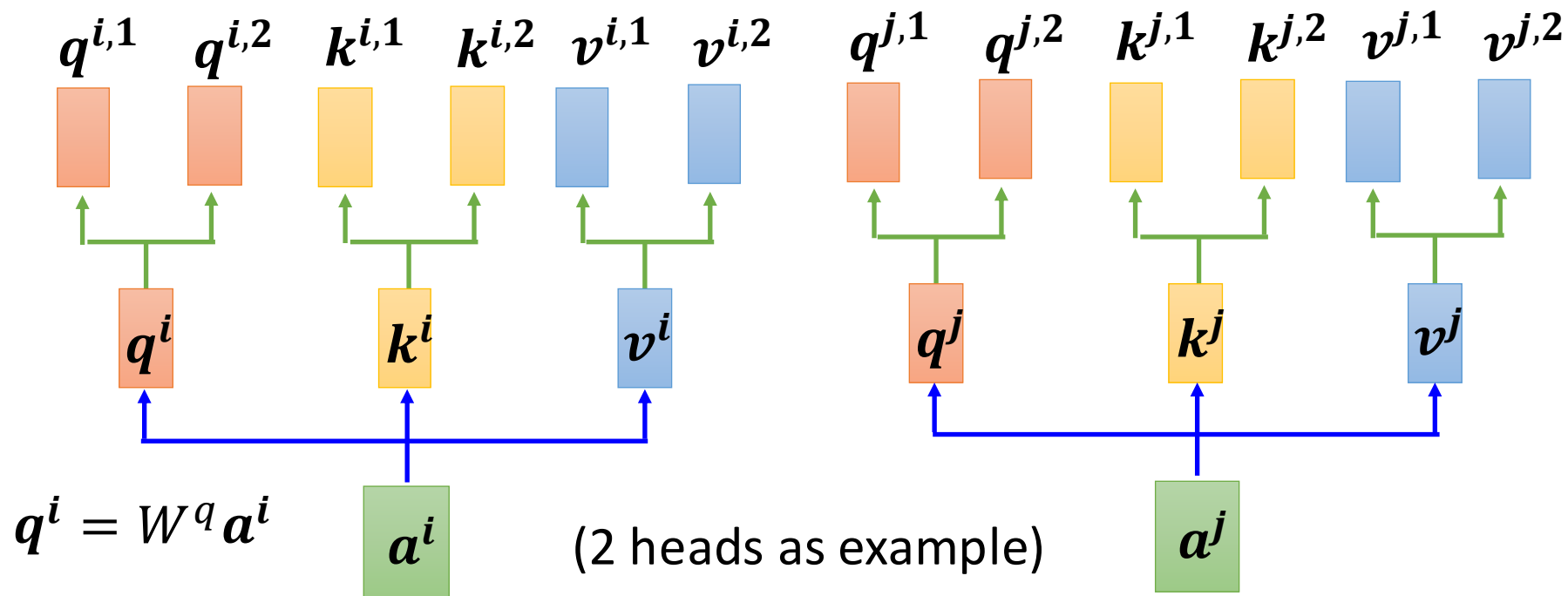
$$q^{i,2} = W^{q,2} q^i$$

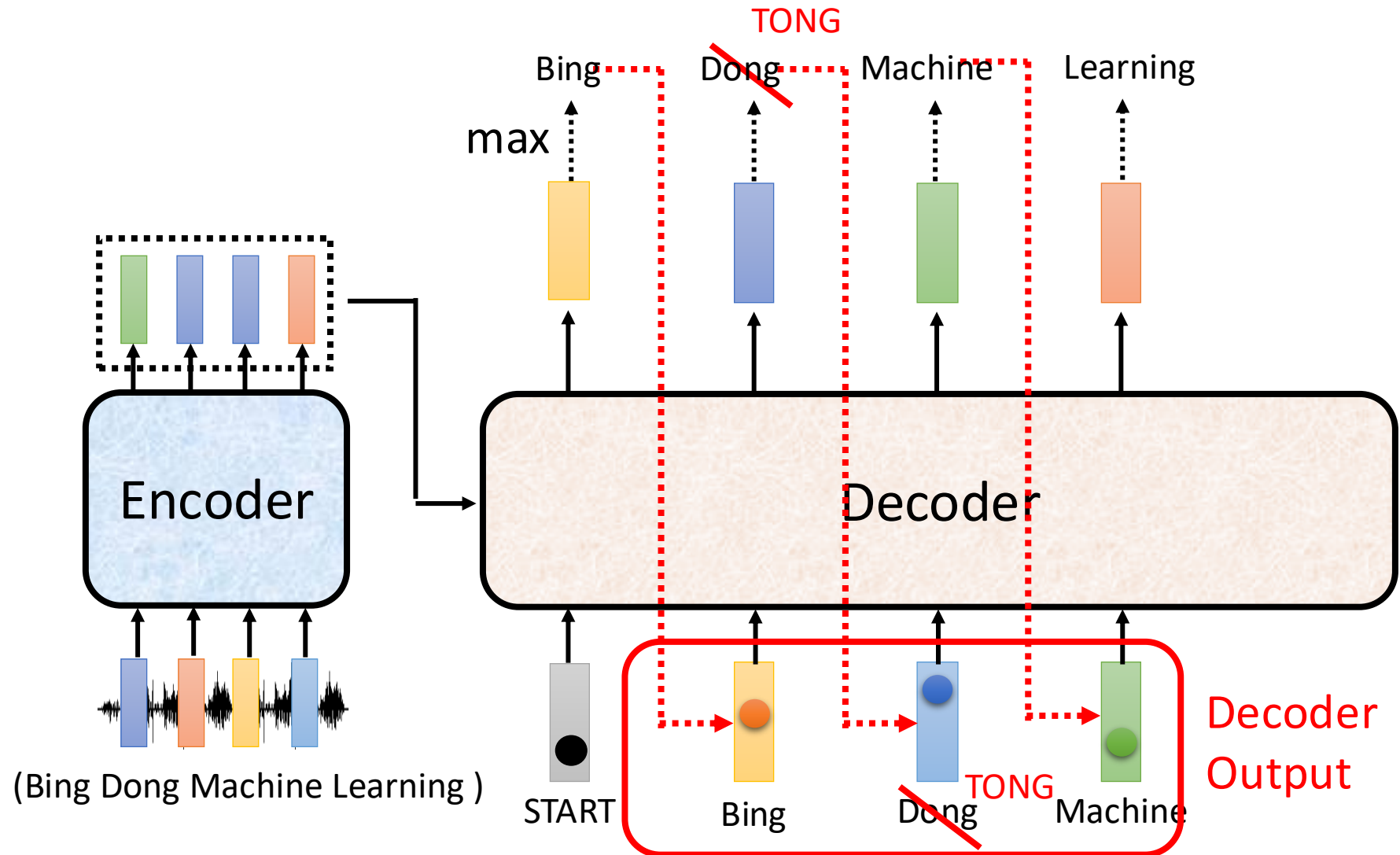


Multi-head Self-attention

Different types of relevance

$$b^i = W^o \begin{bmatrix} b^{i,1} \\ b^{i,2} \end{bmatrix}$$

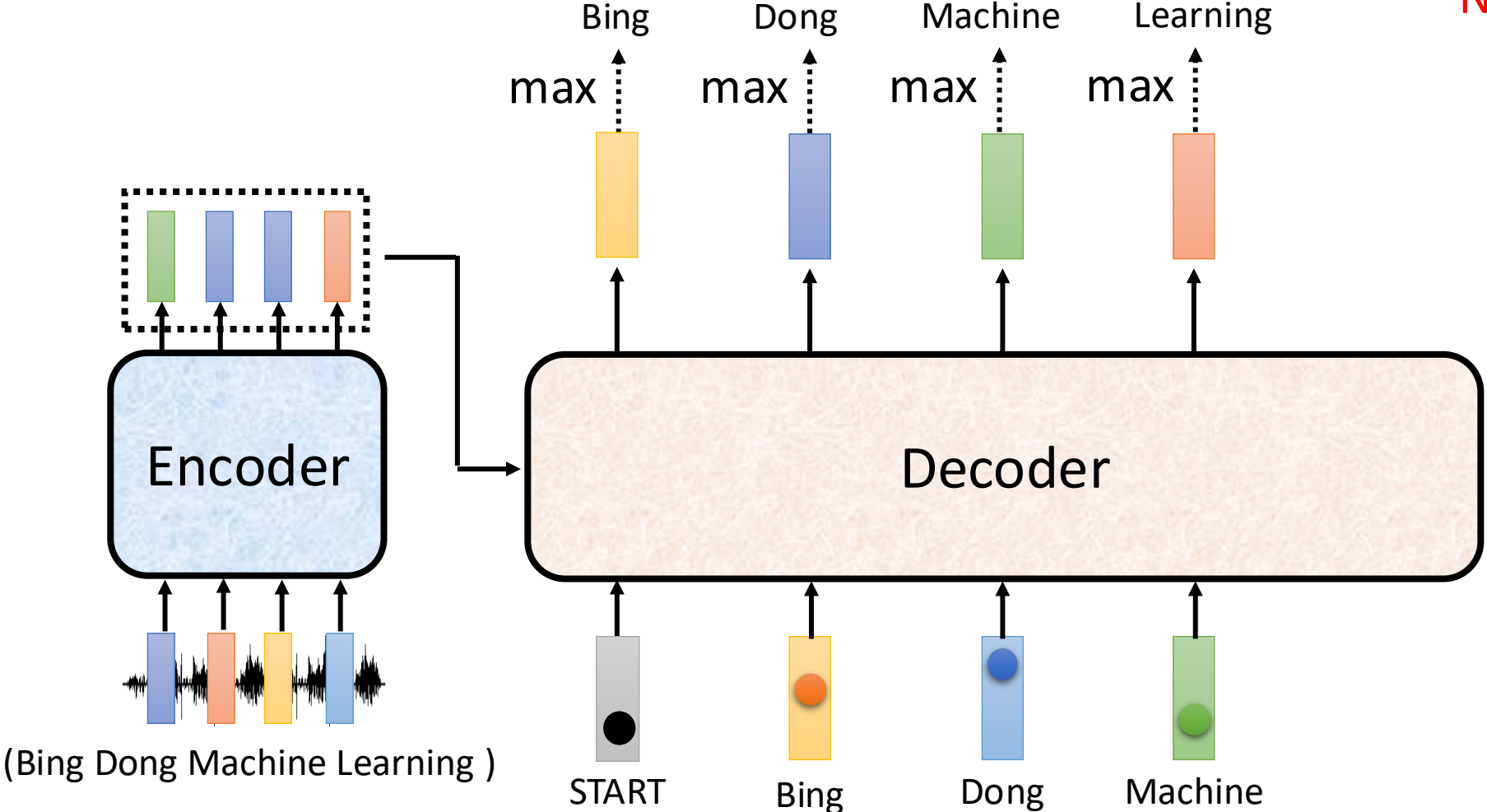




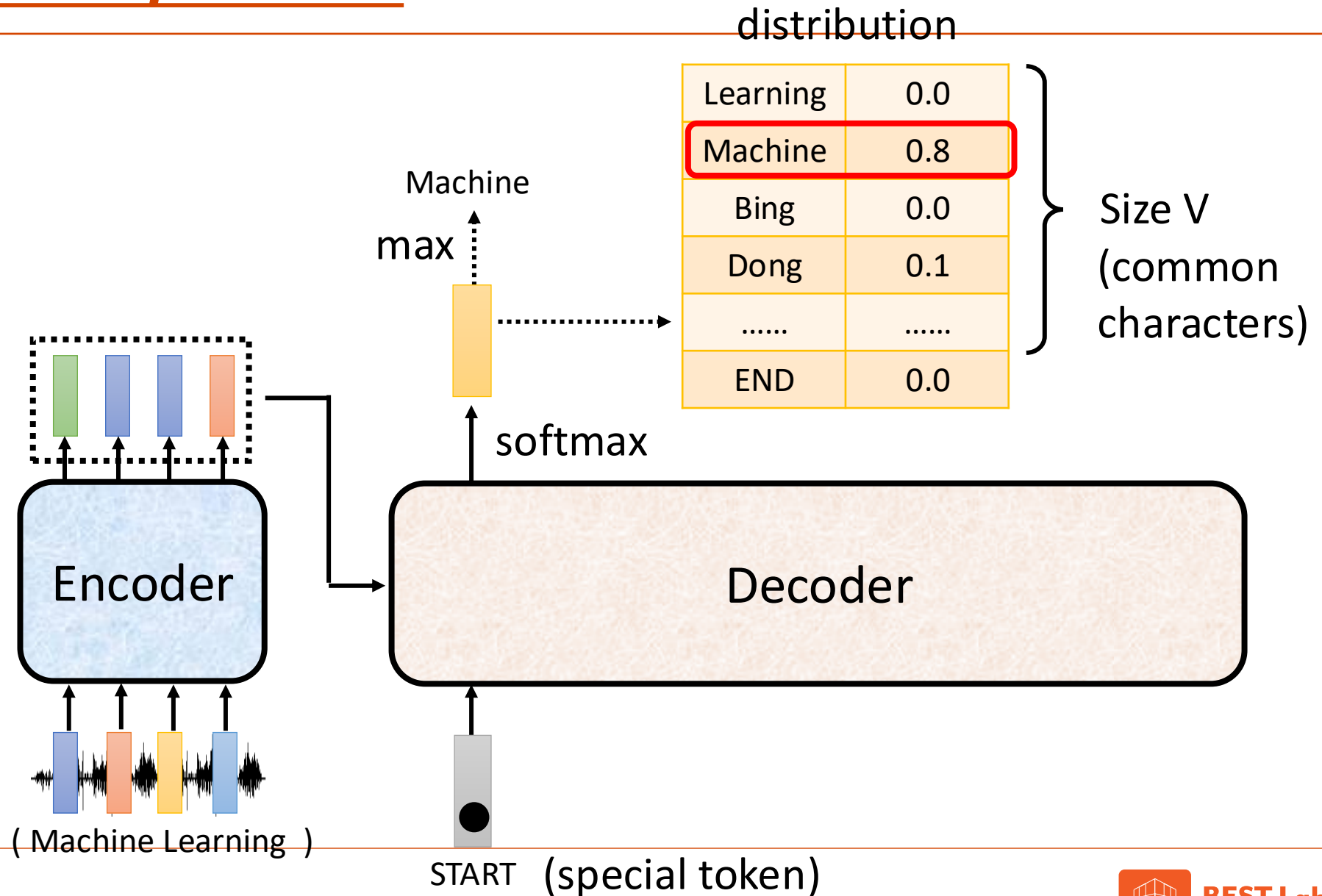
Autoregressive

We do not know the correct output length.

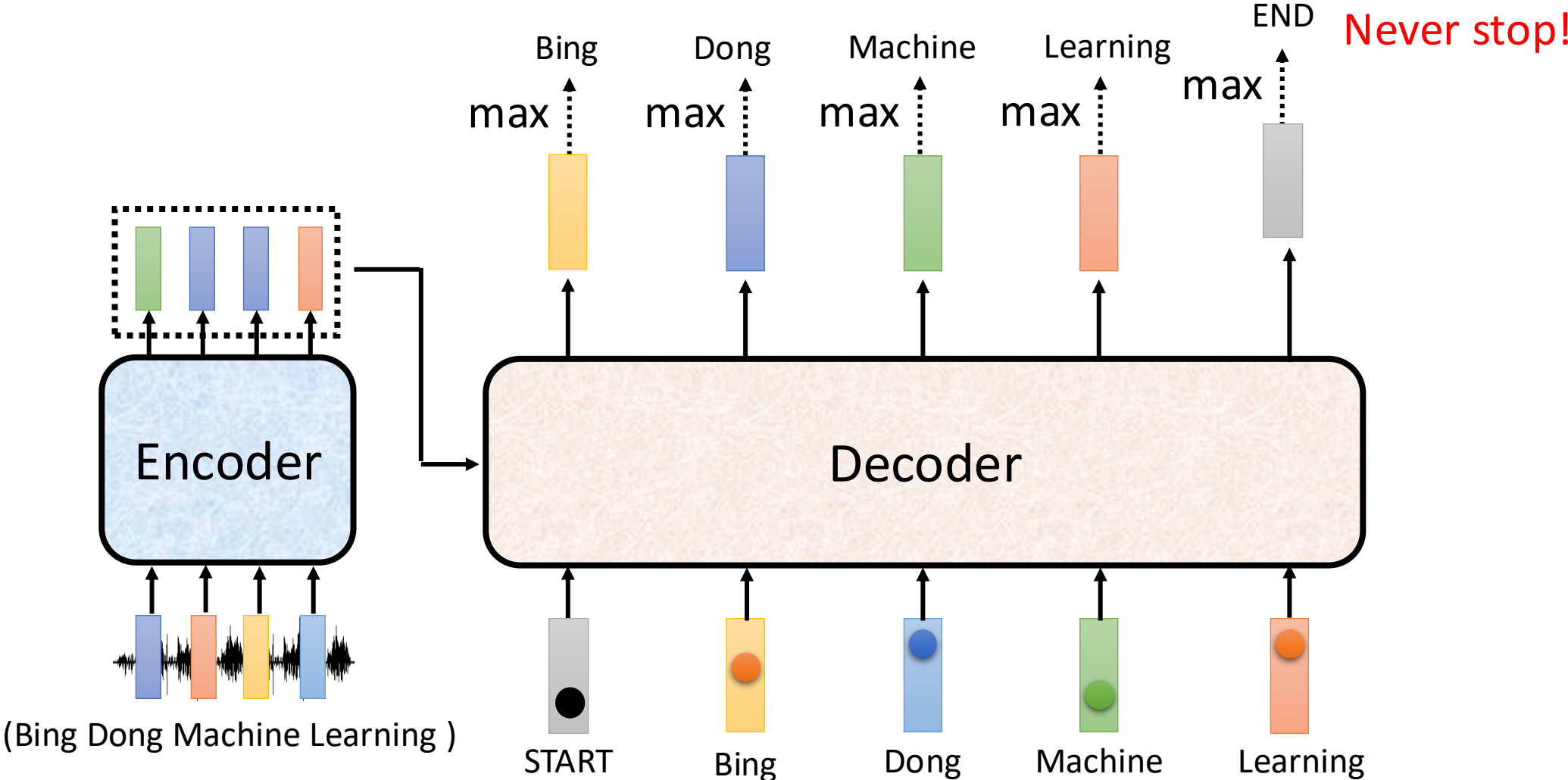
Never stop!



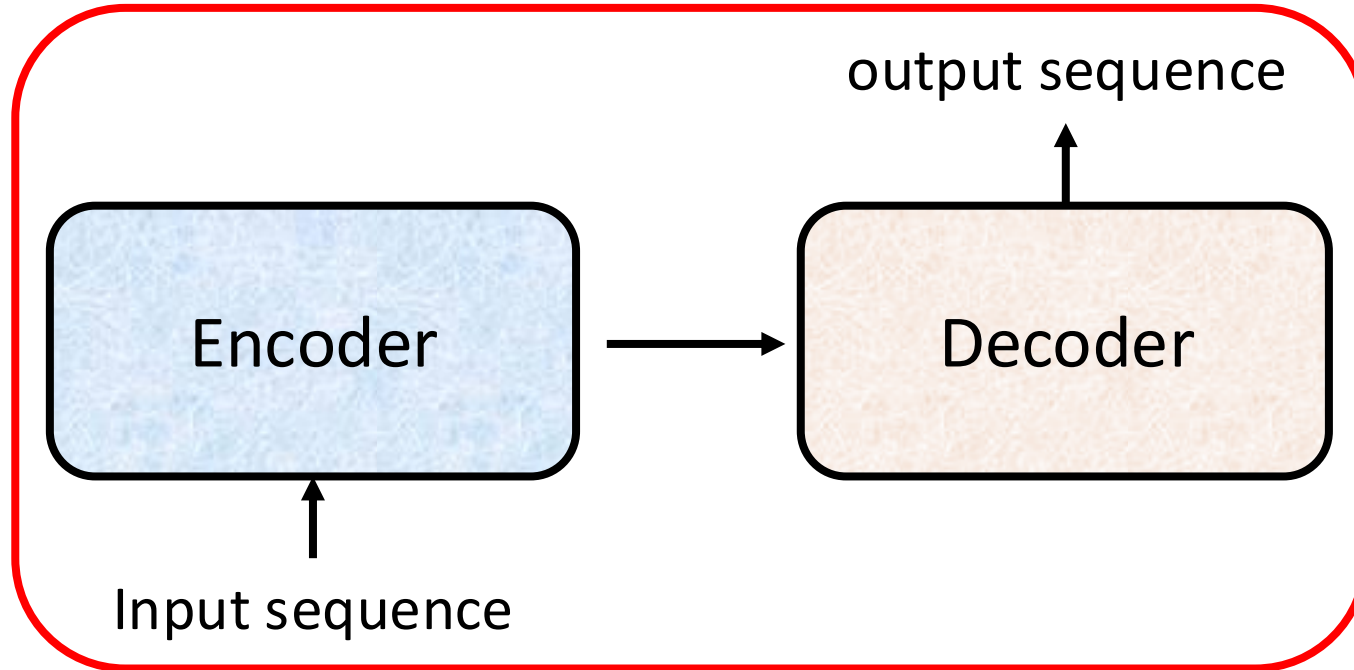
Adding "Stop Token"



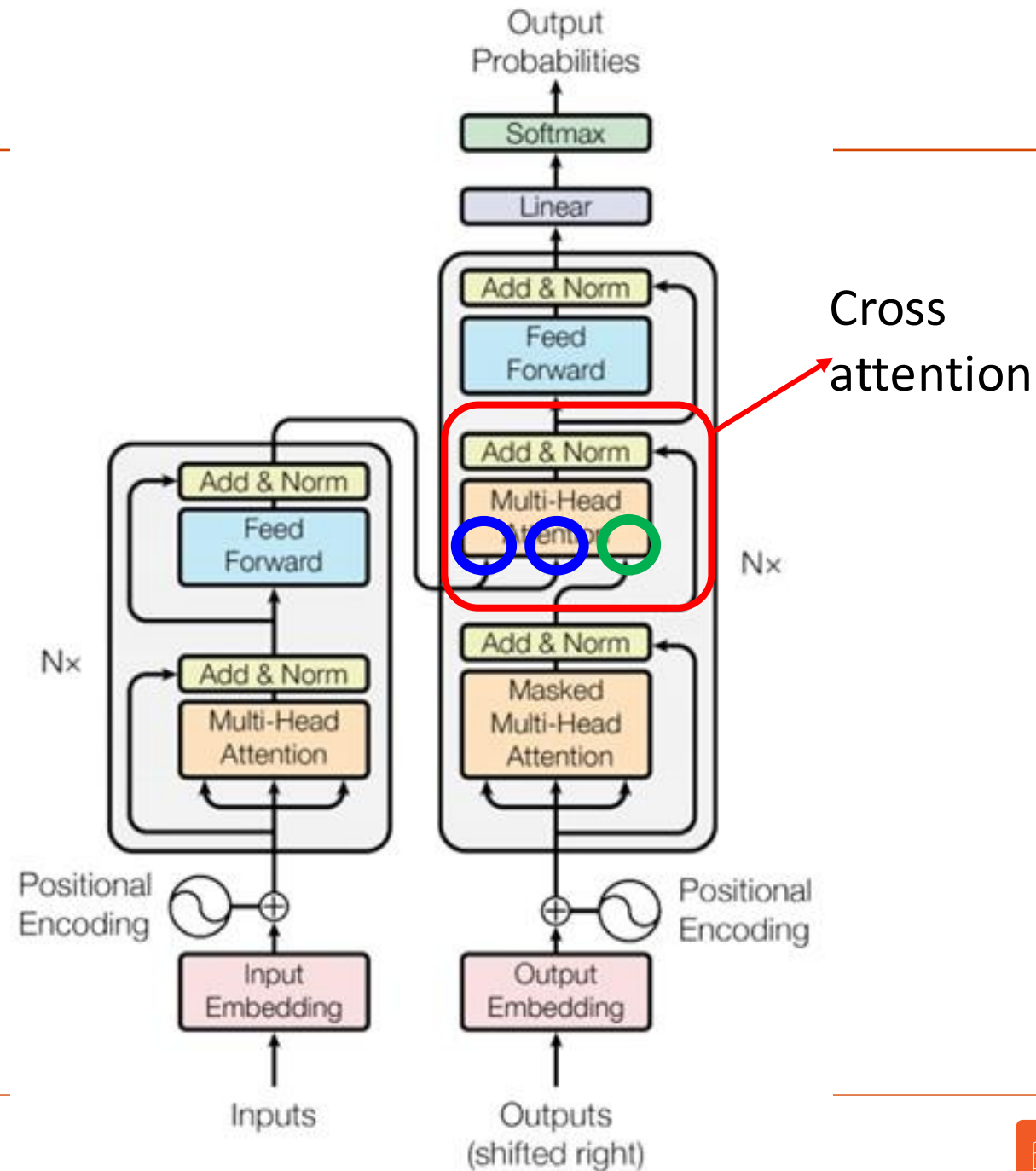
Autoregressive

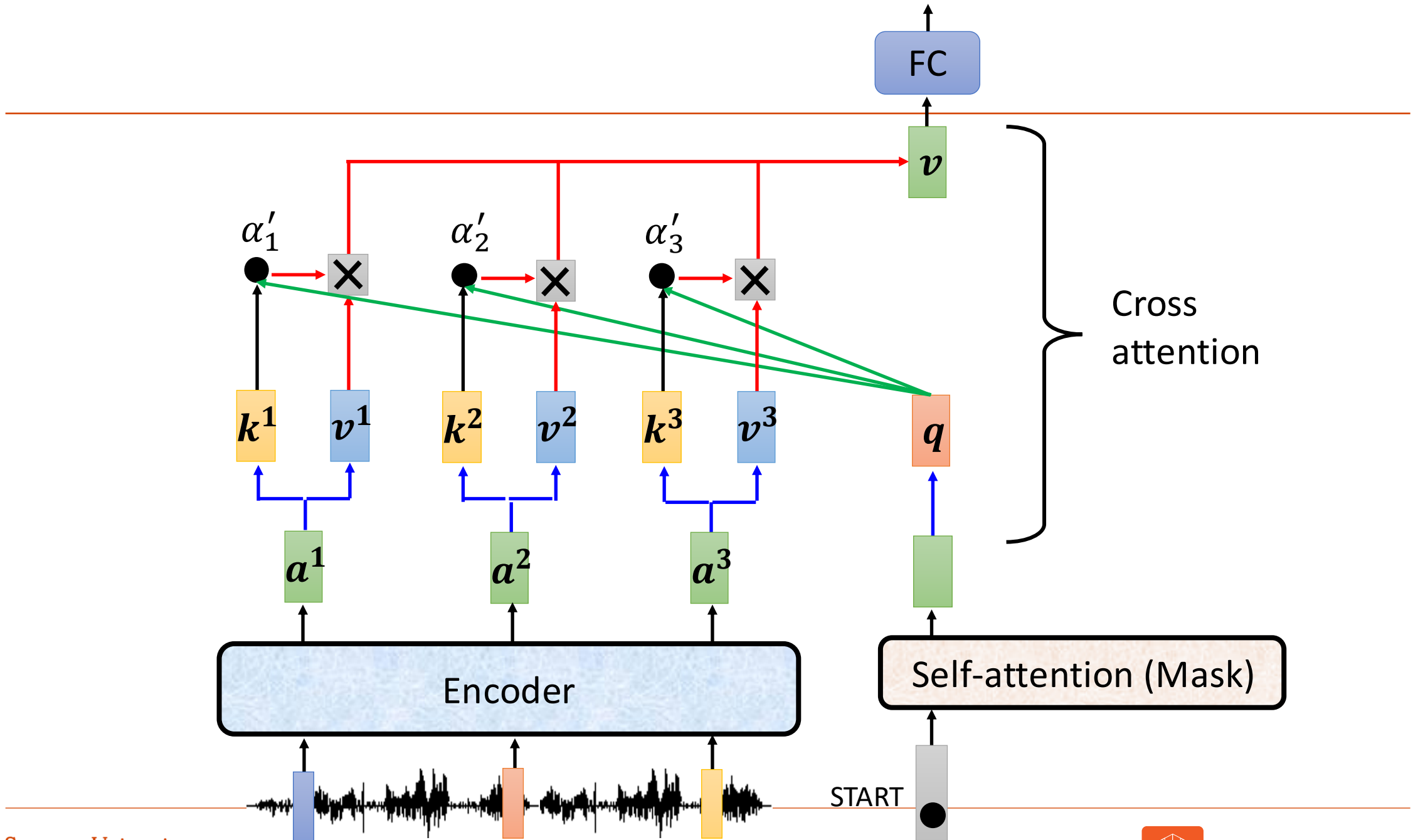


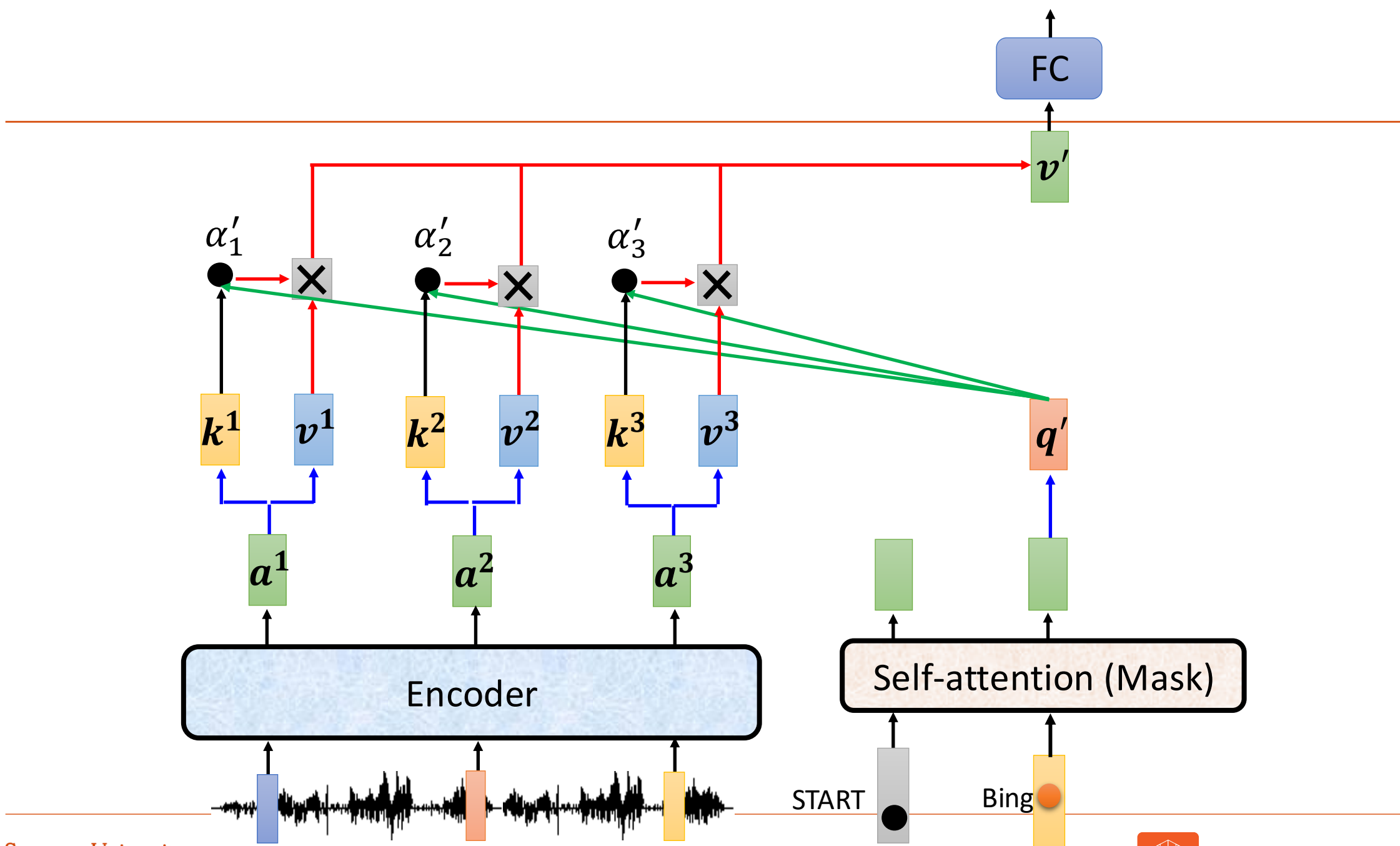
Encoder-Decoder



Transformer



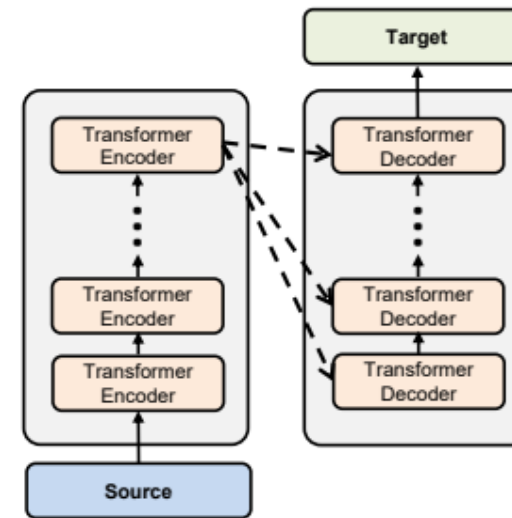




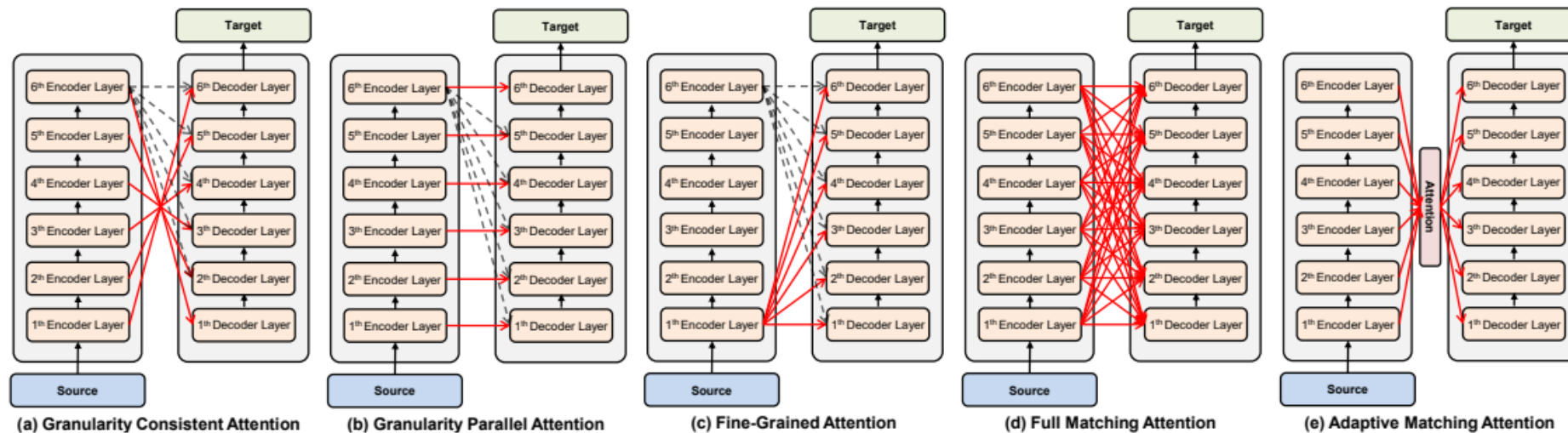
Cross Attention

Source of image:

<https://arxiv.org/abs/2005.08081>



(a) Conventional Transformer



(a) Granularity Consistent Attention

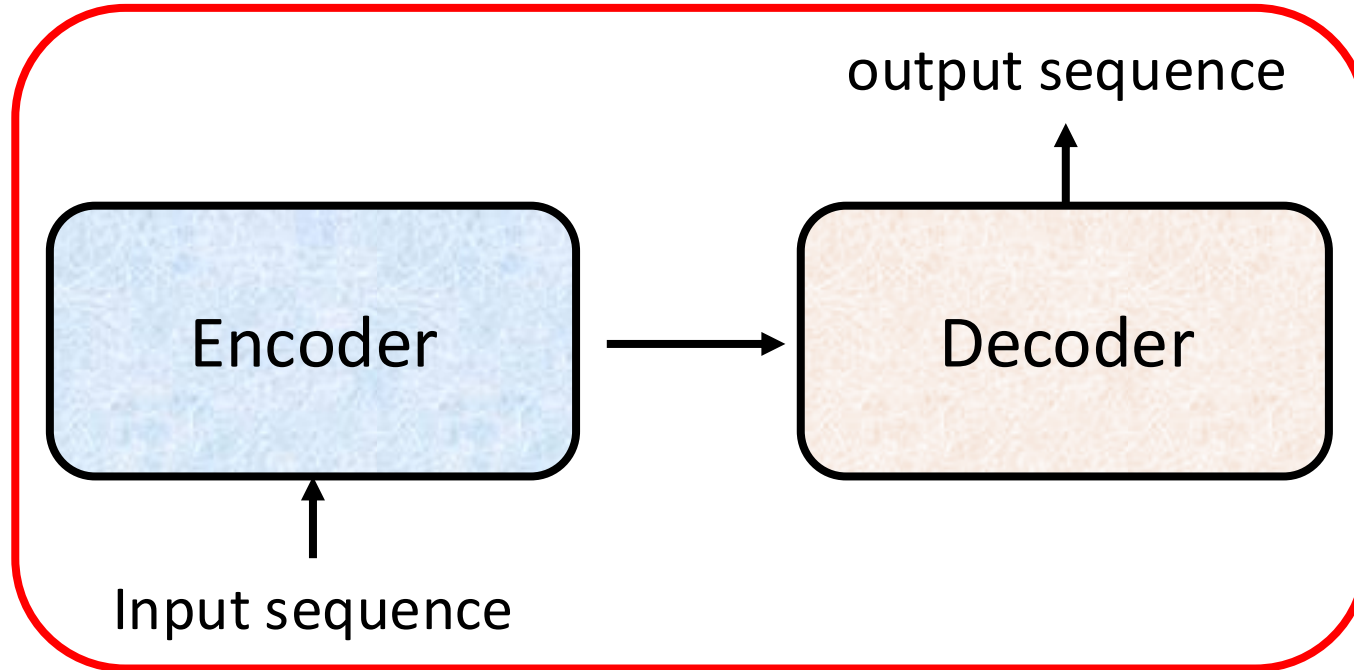
(b) Granularity Parallel Attention

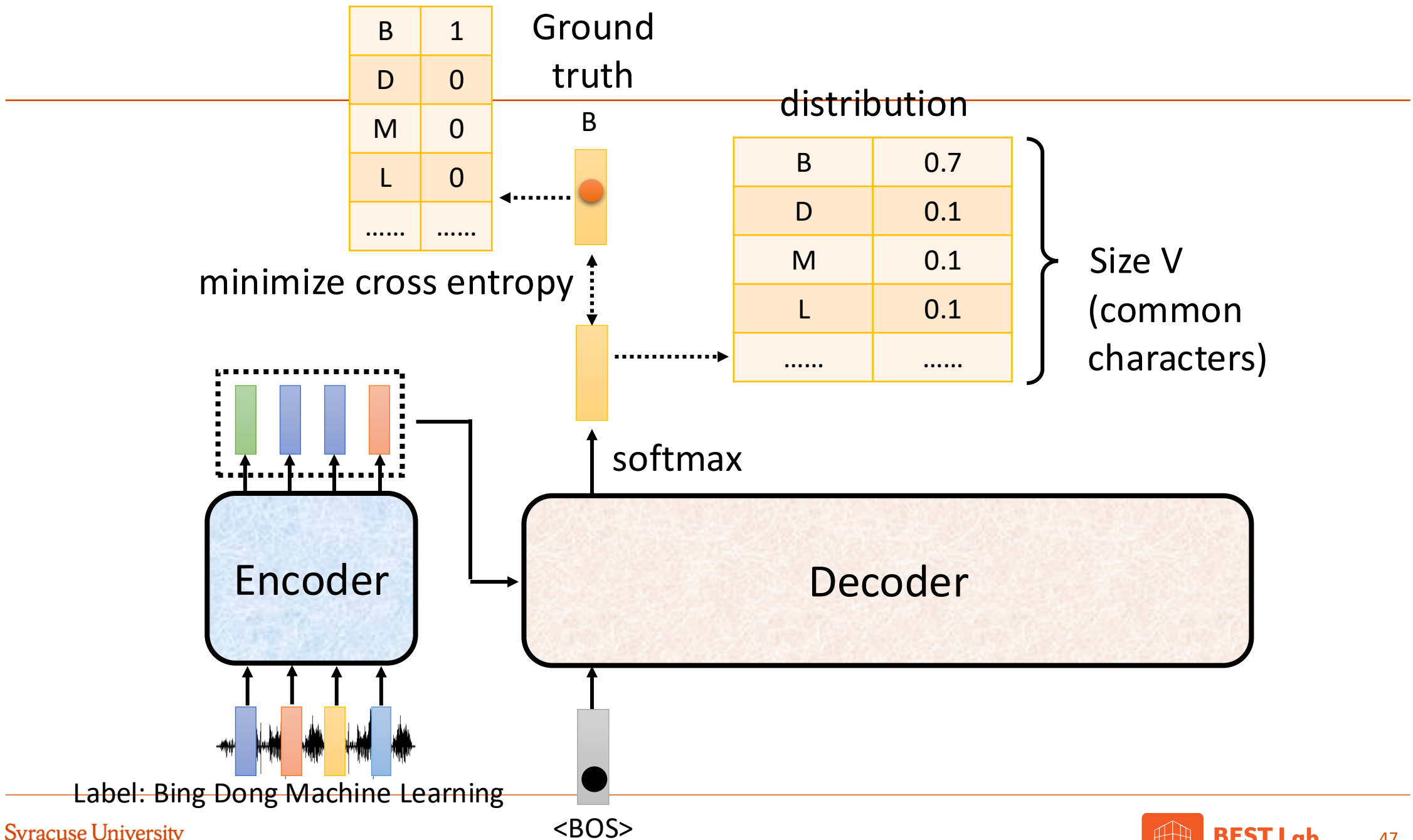
(c) Fine-Grained Attention

(d) Full Matching Attention

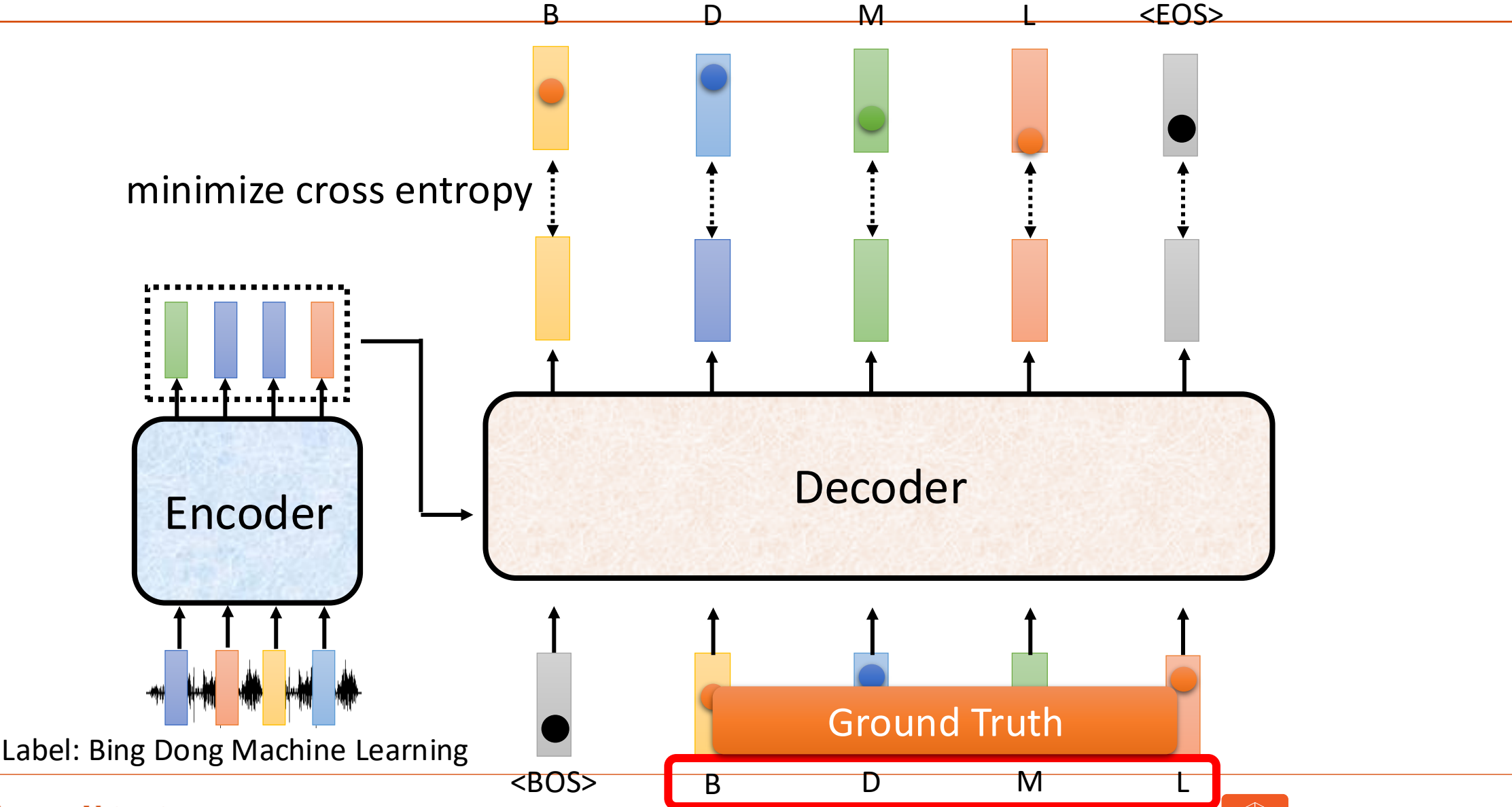
(e) Adaptive Matching Attention

Training





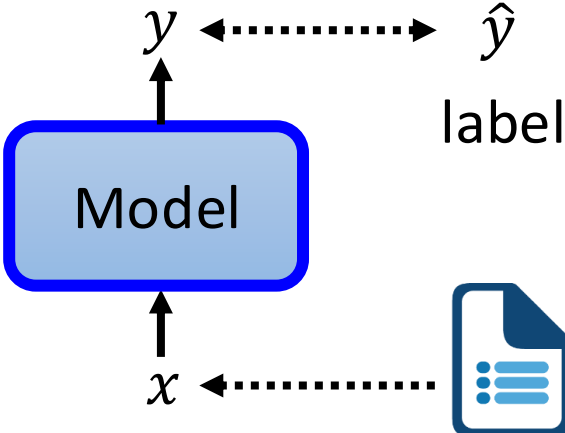
Teacher Forcing: using the ground truth as input.



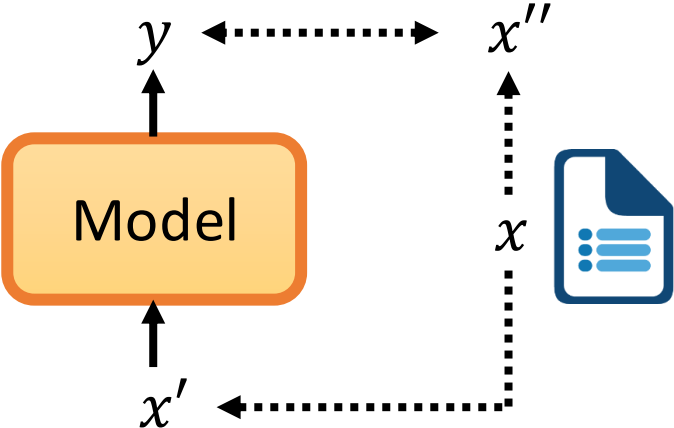
Large Language Model----A Self-Supervised Learned Transformer

Self-Supervised Learning

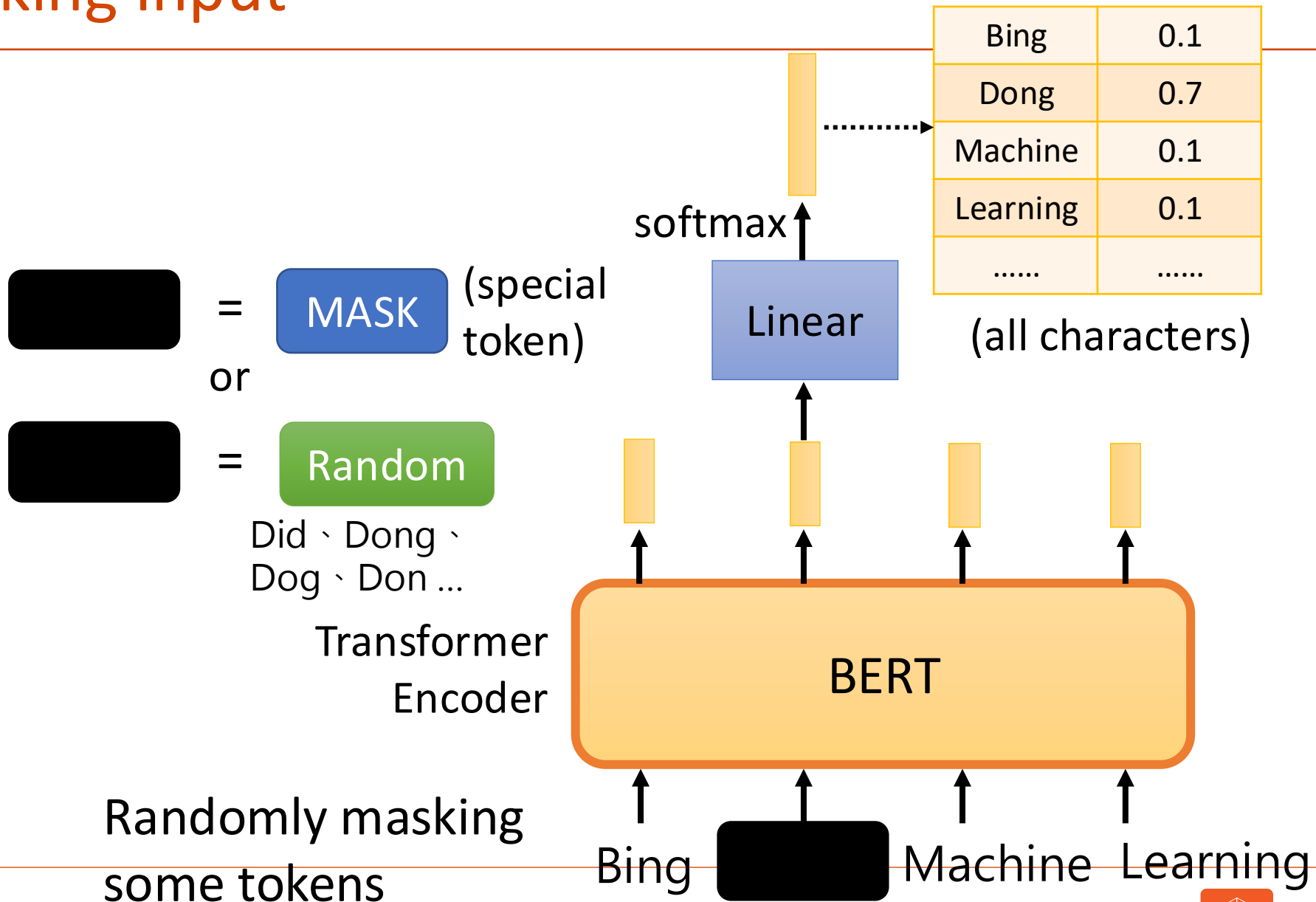
Supervised



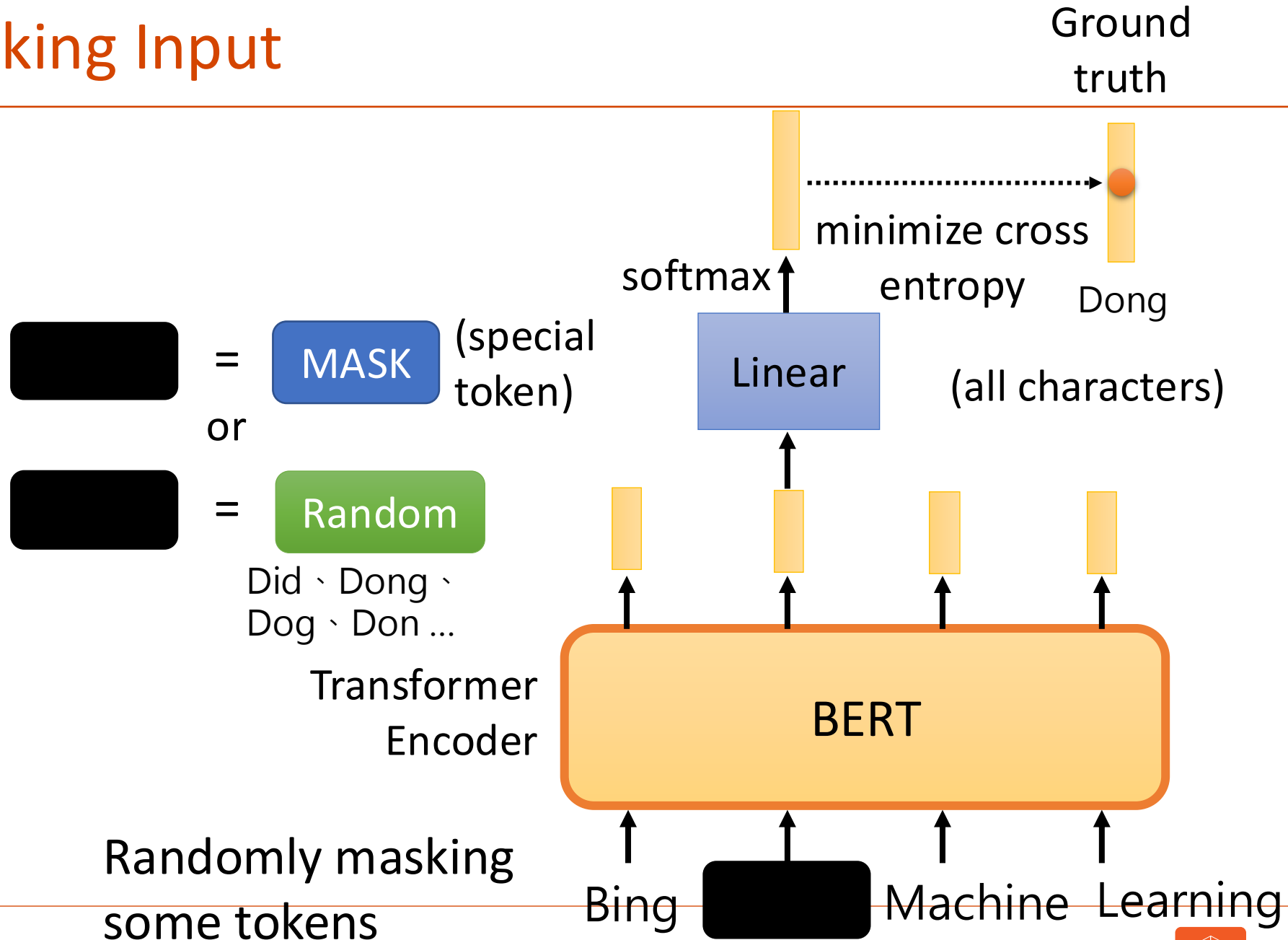
Self-supervised

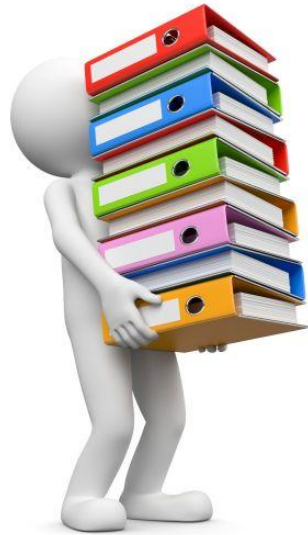


Masking Input

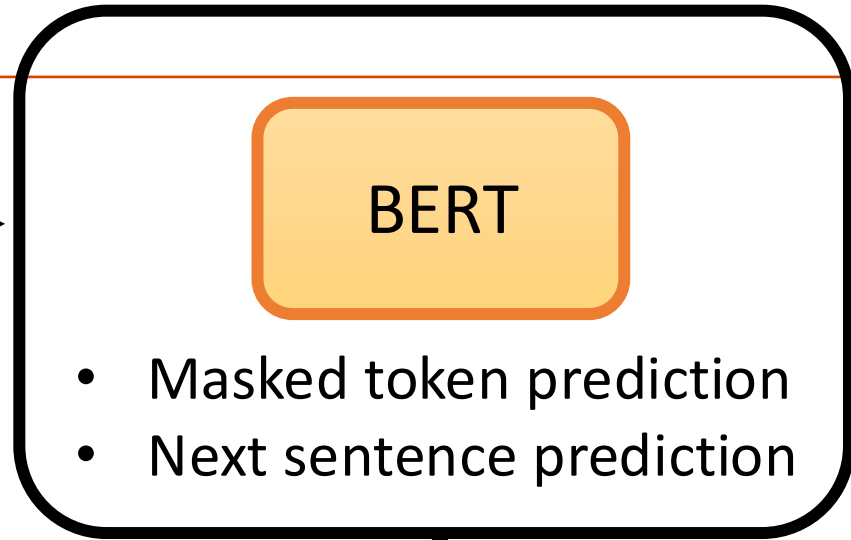


Masking Input

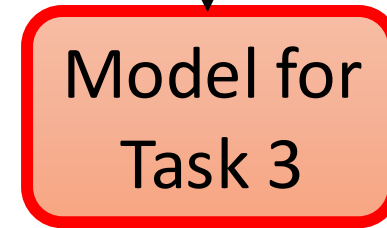
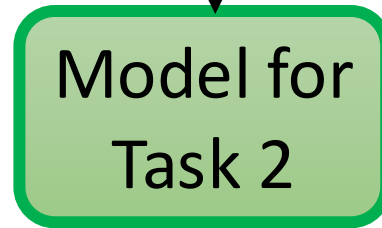
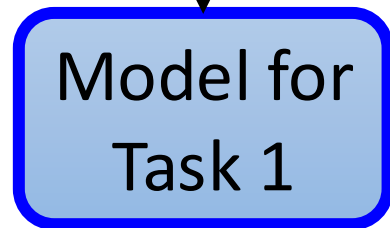




Self-supervised
Learning
Pre-train



Fine-tune



Downstream Tasks

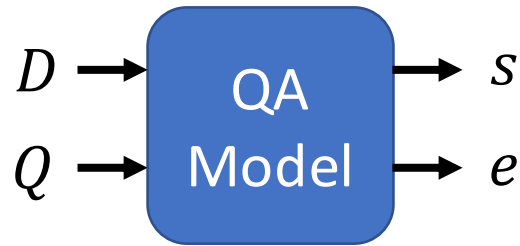
- The tasks we care
- We have a little bit labeled data.

How to use BERT

- Extraction-based Question Answering (QA)

Document: $D = \{d_1, d_2, \dots, d_N\}$

Query: $Q = \{q_1, q_2, \dots, q_M\}$



output: two integers (s, e)

Answer: $A = \{d_s, \dots, d_e\}$

In meteorology, precipitation is any product of the condensation of **17** spheric water vapor that falls under **gravity**. The main forms of precipitation include drizzle, rain, sleet, snow, **graupel** and hail... Precipitation forms as smaller droplets coalesce via collision with other rain drops or ice crystals **within a cloud**. Short, intense periods of rain **77** at **79** are called "showers".

What causes precipitation to fall?

gravity

$s = 17, e = 17$

What is another main form of precipitation besides drizzle, rain, snow, sleet and hail?

graupel

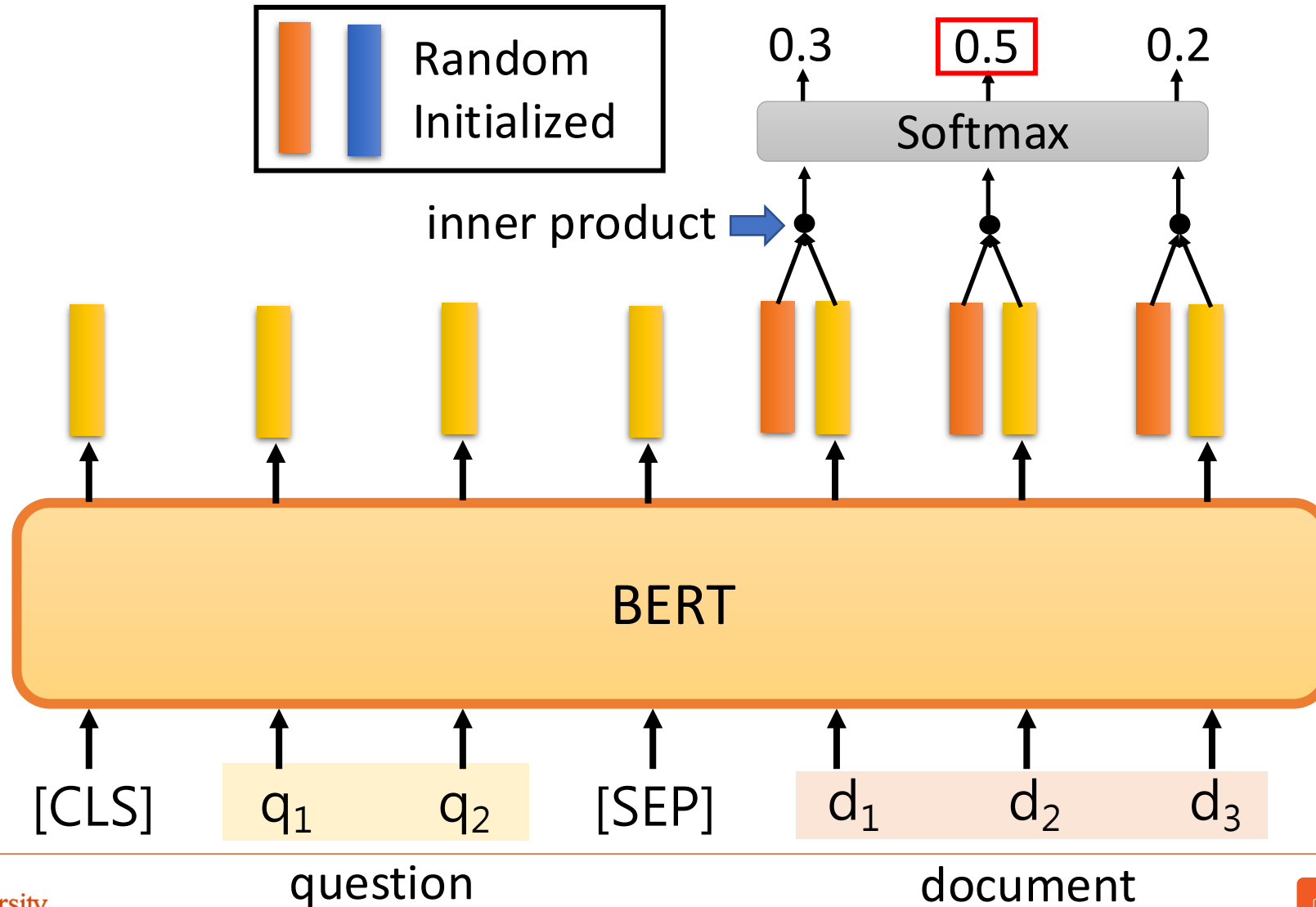
Where do water droplets collide with ice crystals to form precipitation?

within a cloud

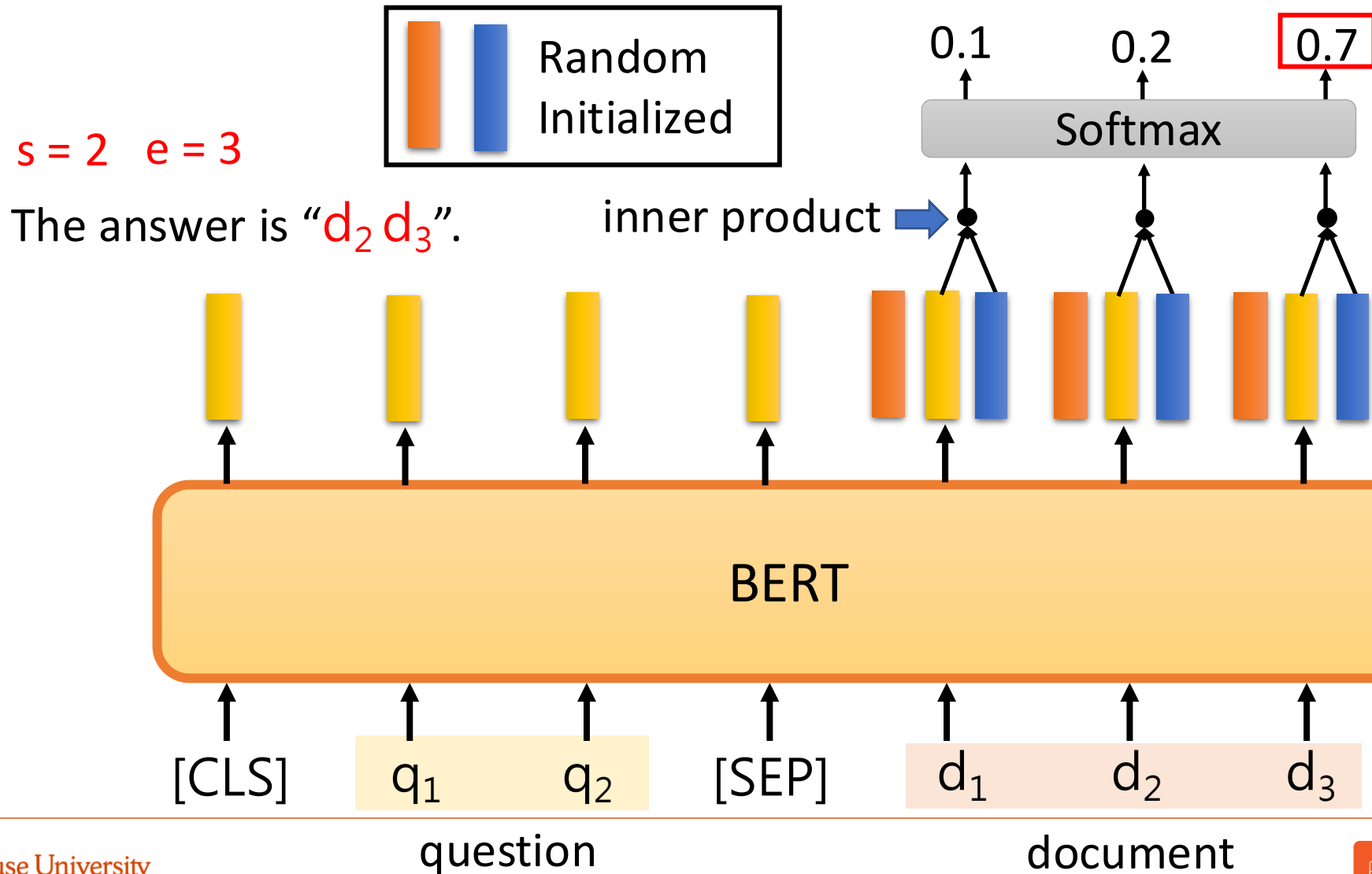
$s = 77, e = 79$

How to use BERT

$s = 2$



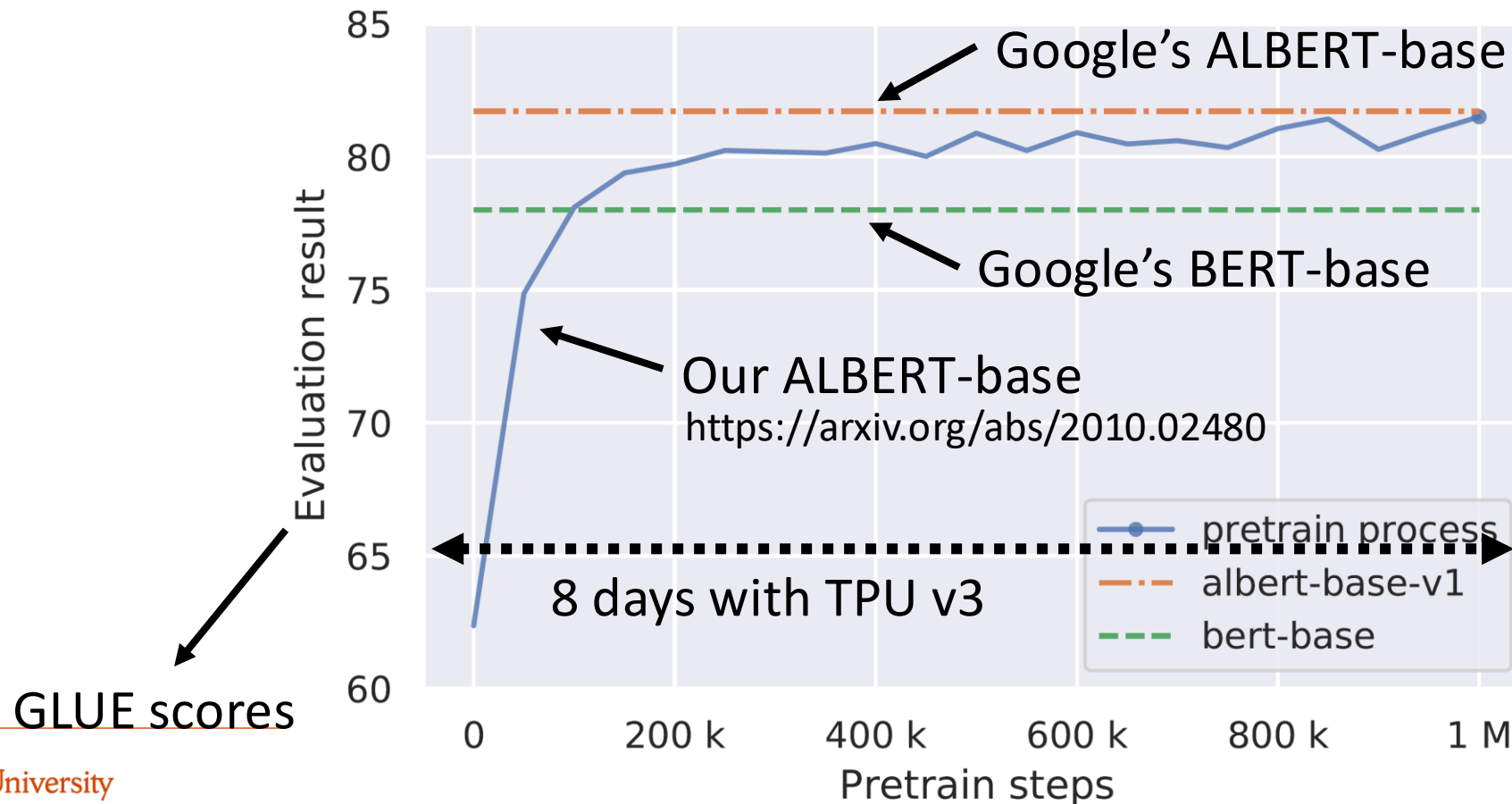
How to use BERT



Training BERT is challenging!

Training data has more than **3 billions** of words.

3000 times of **Harry Potter series**



Transformer

2017

Attention Is All You Need

Ashish Vaswani*
Google Brain
avaswani@google.com

Noam Shazeer*
Google Brain
noam@google.com

Niki Parmar*
Google Research
nikip@google.com

Jakob Uszkoreit*
Google Research
usz@google.com

Llion Jones*
Google Research
llion@google.com

Aidan N. Gomez* †
University of Toronto
aidan@cs.toronto.edu

Lukasz Kaiser*
Google Brain
lukaszkaizer@google.com



Transformer



BERT 340M



GPT-2 (1.5B)



GPT-3 (175B)

GPT-4 (2T)

Any Questions?

bidong@syr.edu