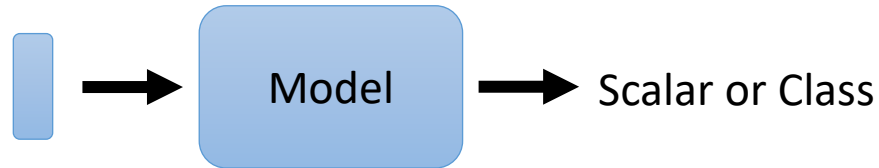


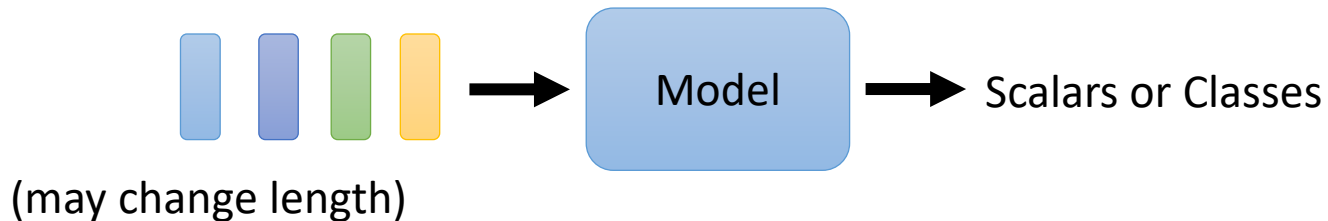
Self-Attention, Transformers, BERT, GPT

Sophisticated Input

- Input is a **vector**




- Input is a **set of vectors**



Vector Set as Input

this is a cat



A 4-dimensional embedding

cat =>

| | | | |
|-----|------|------|-----|
| 1.2 | -0.1 | 4.3 | 3.2 |
| 0.4 | 2.5 | -0.9 | 0.5 |
| 2.1 | 0.3 | 0.1 | 0.4 |

...

...

One-hot Encoding

apple = [1 0 0 0 0]

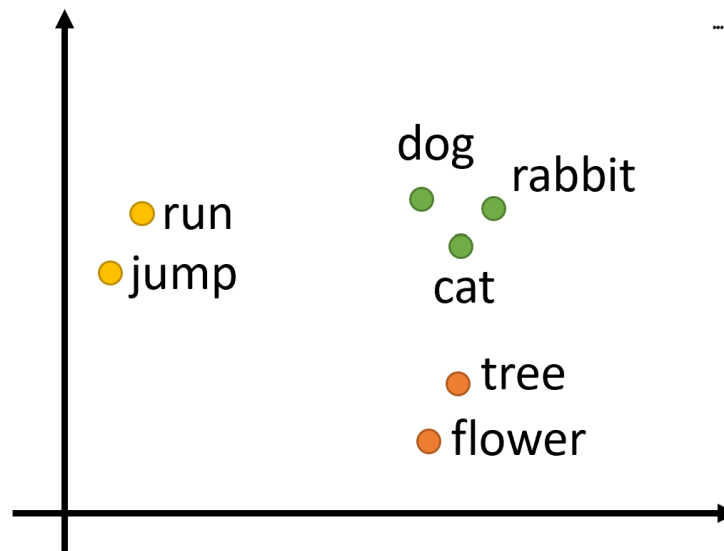
bag = [0 1 0 0 0]

cat = [0 0 1 0 0]

dog = [0 0 0 1 0]

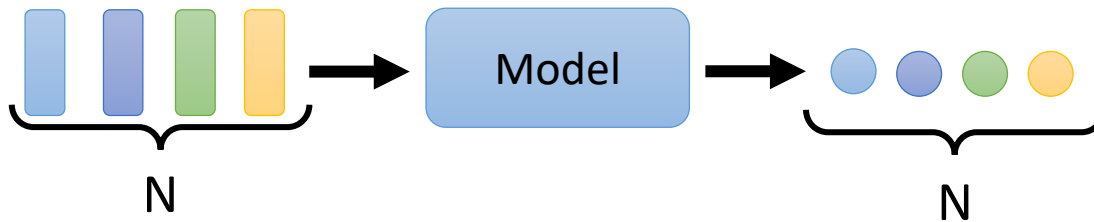
elephant = [0 0 0 0 1]

Word Embedding

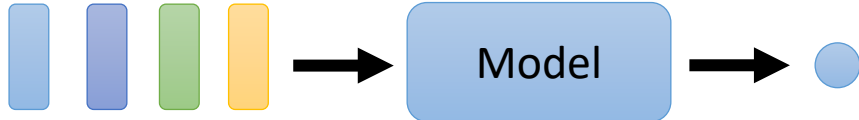


What is the output?

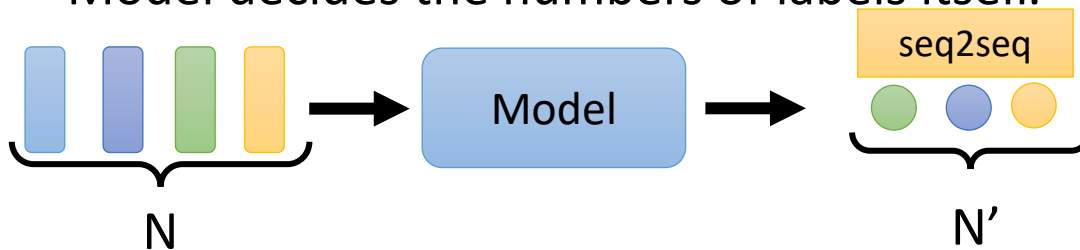
- Each vector has a label.



- The whole sequence has a label.



- Model decides the numbers of labels itself.

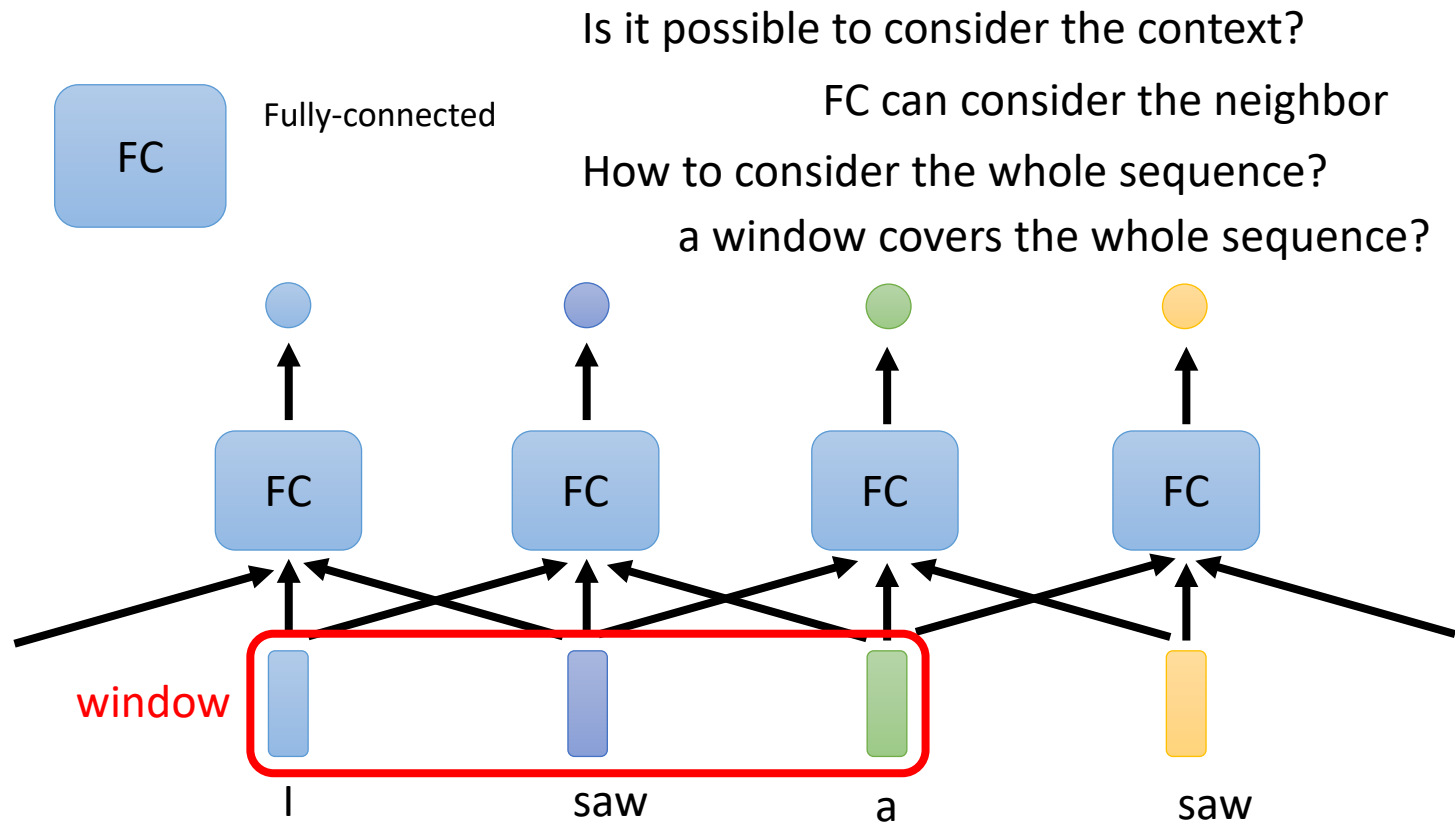


I saw a saw
↓ ↓ ↓ ↓
N V DET N

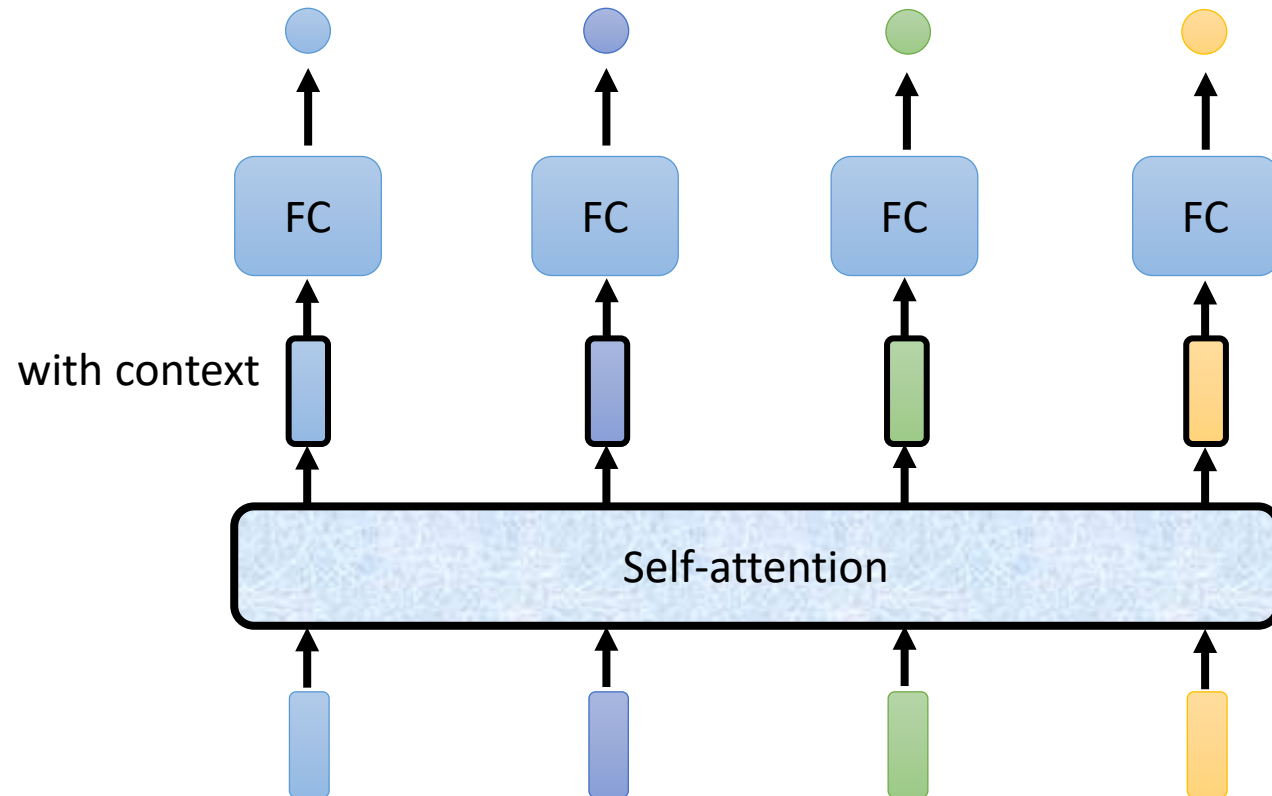
POS tagging

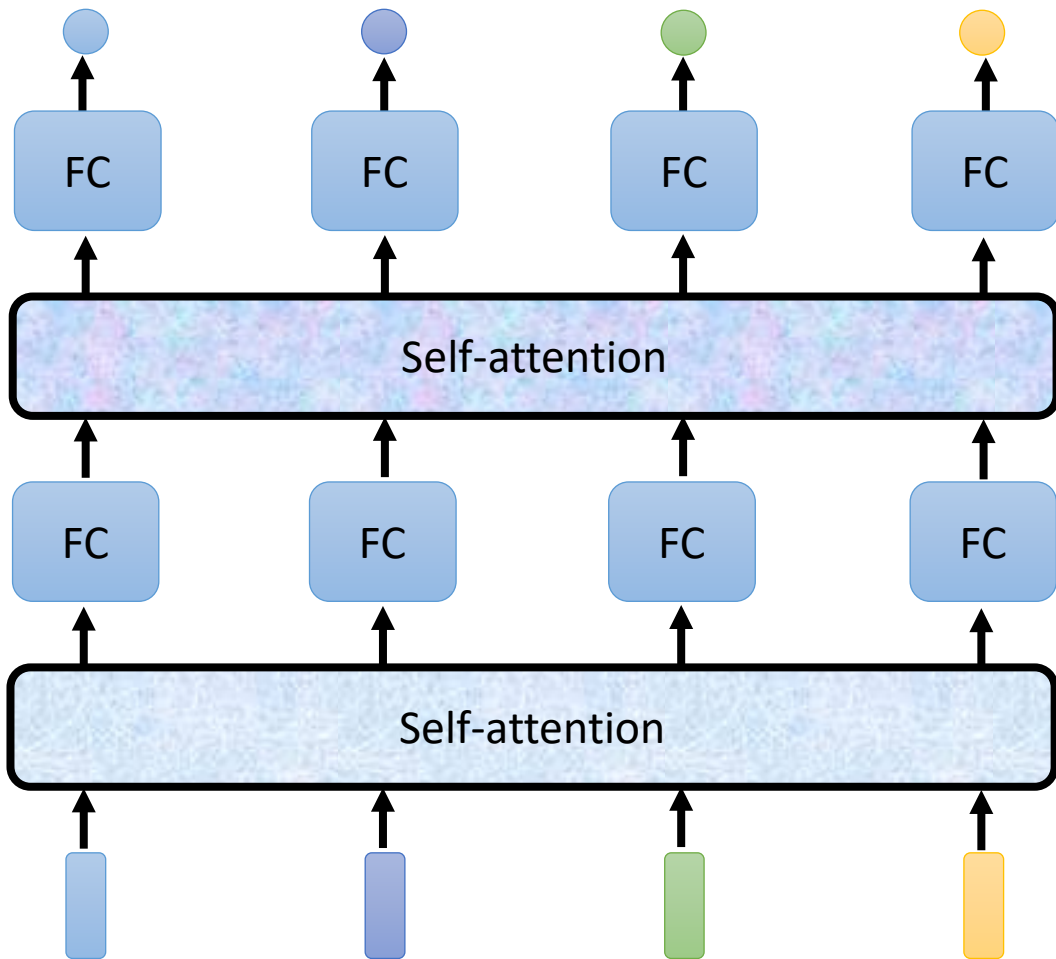
this is good
Sentiment analysis ↓
positive

Sequence Labeling



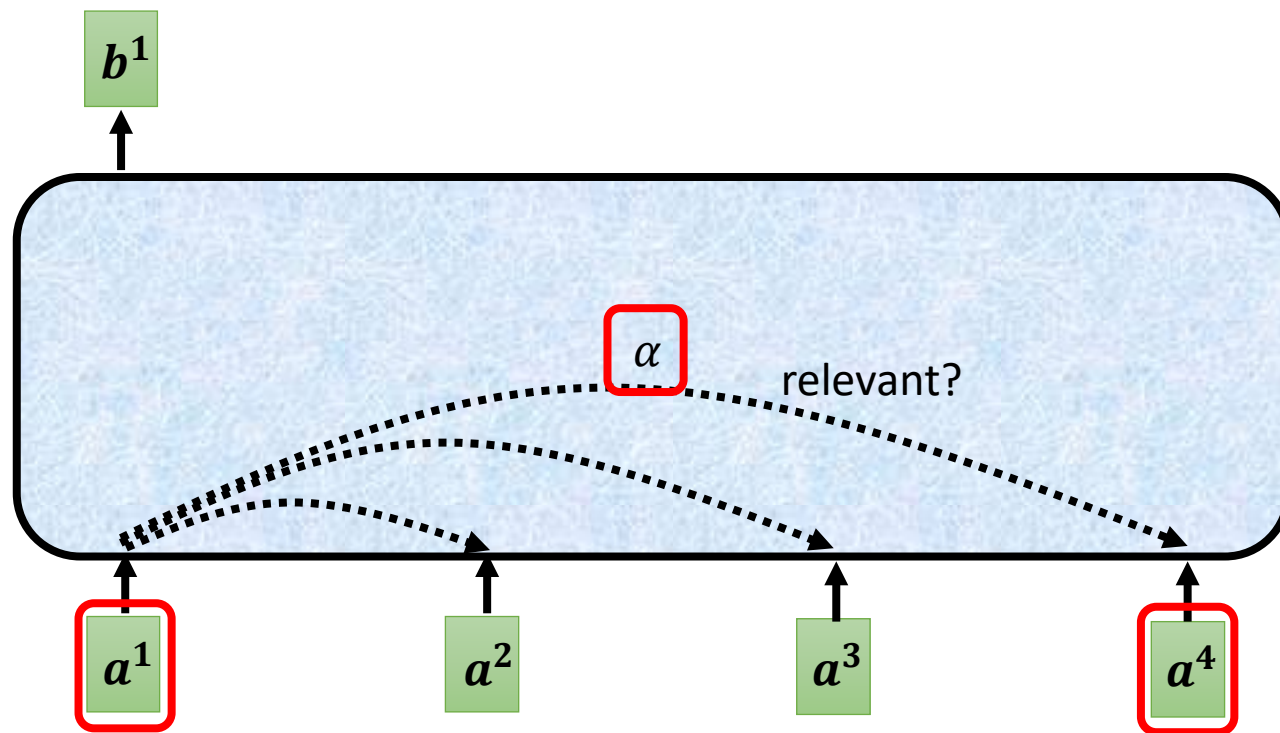
Self-attention





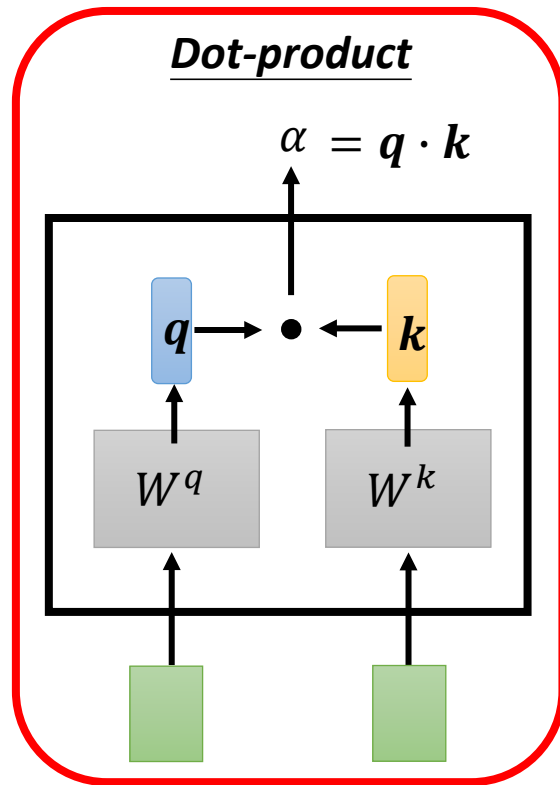
<https://arxiv.org/abs/1706.03762>

Self-attention

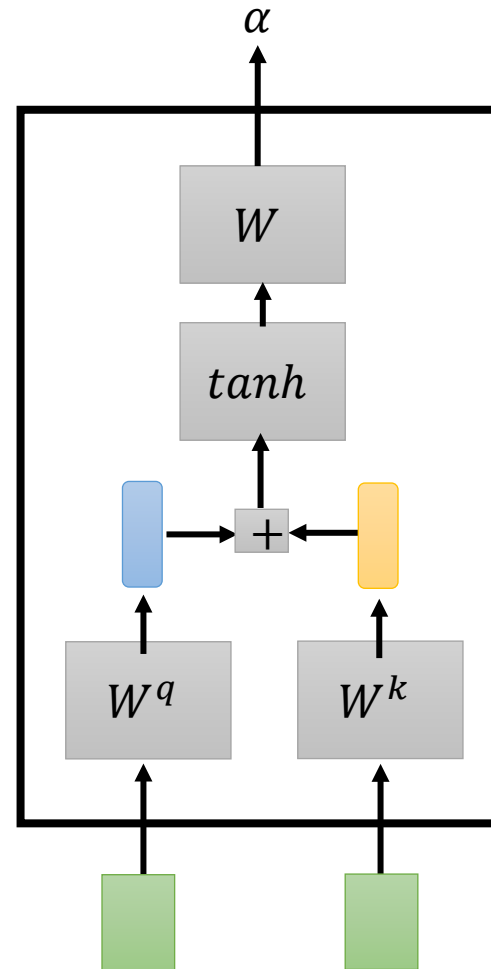


Find the relevant vectors in a sequence

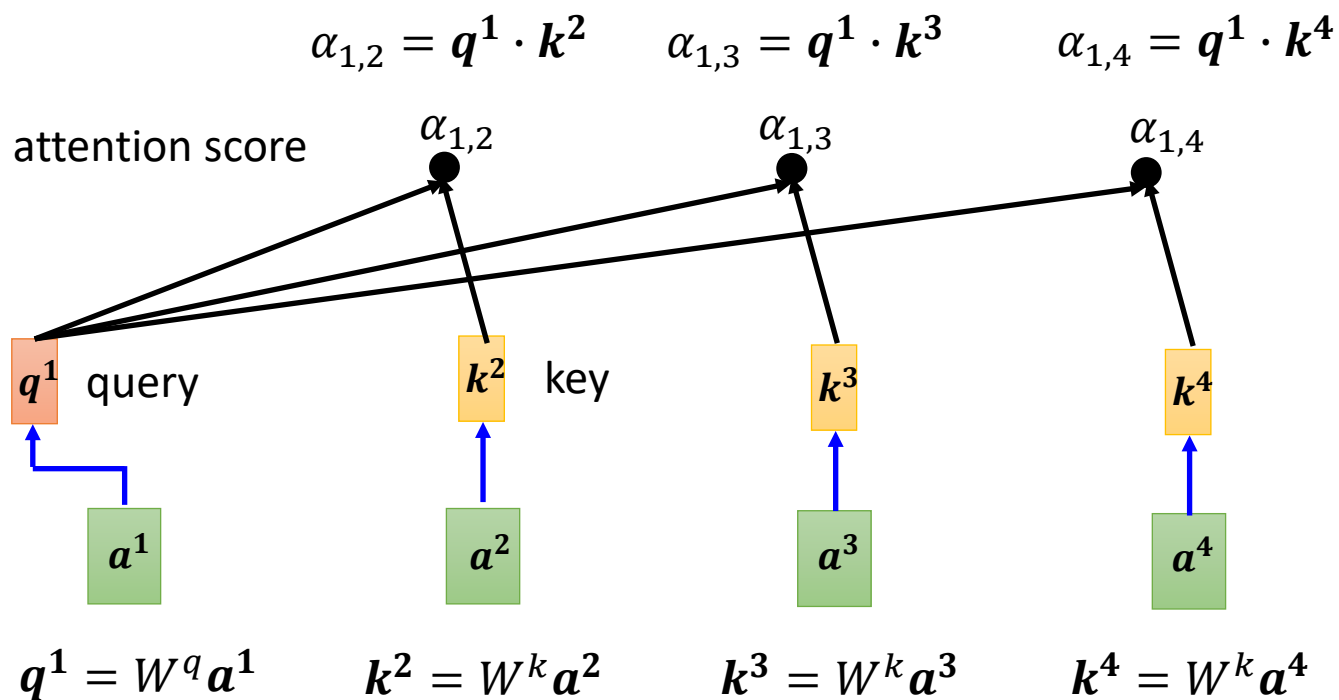
Self-attention



Additive

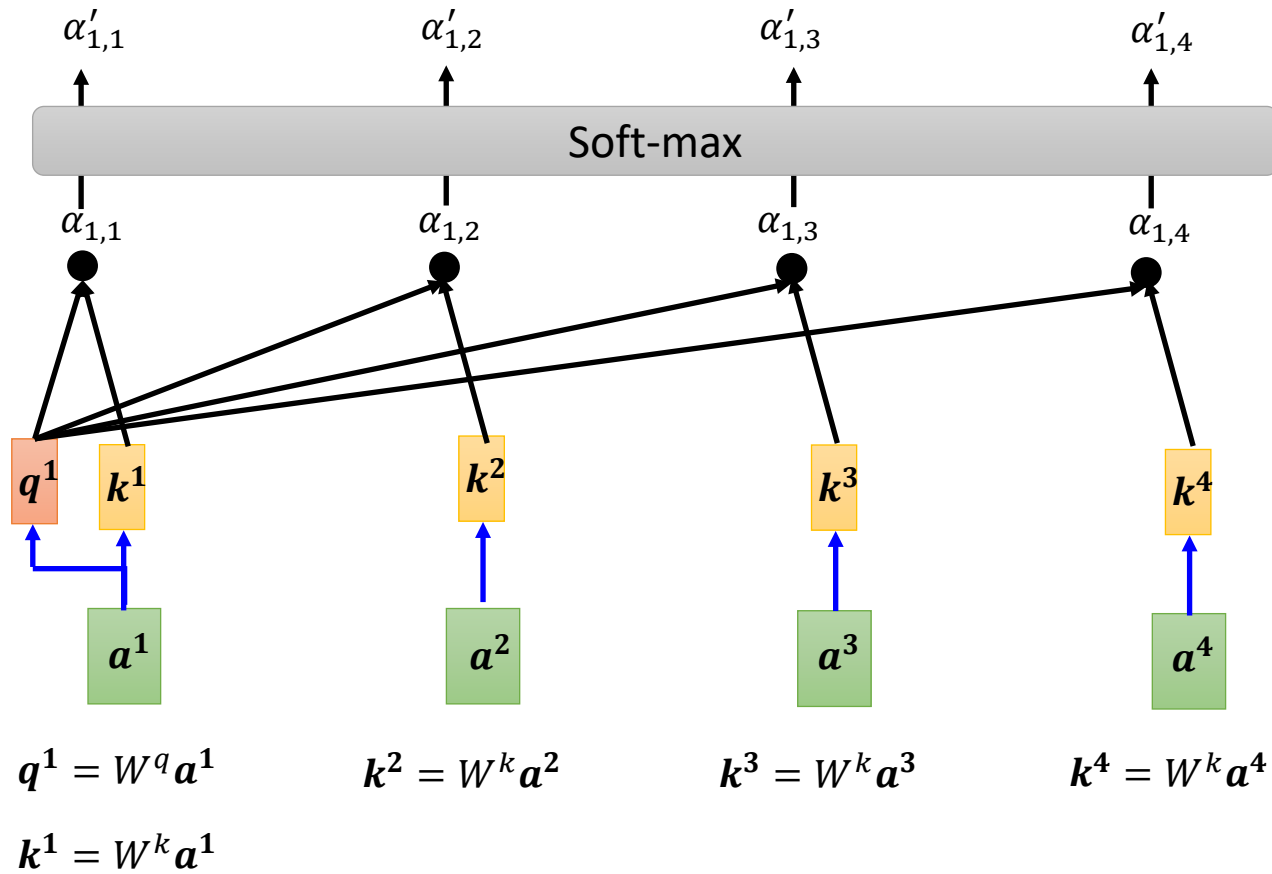


Self-attention



Self-attention

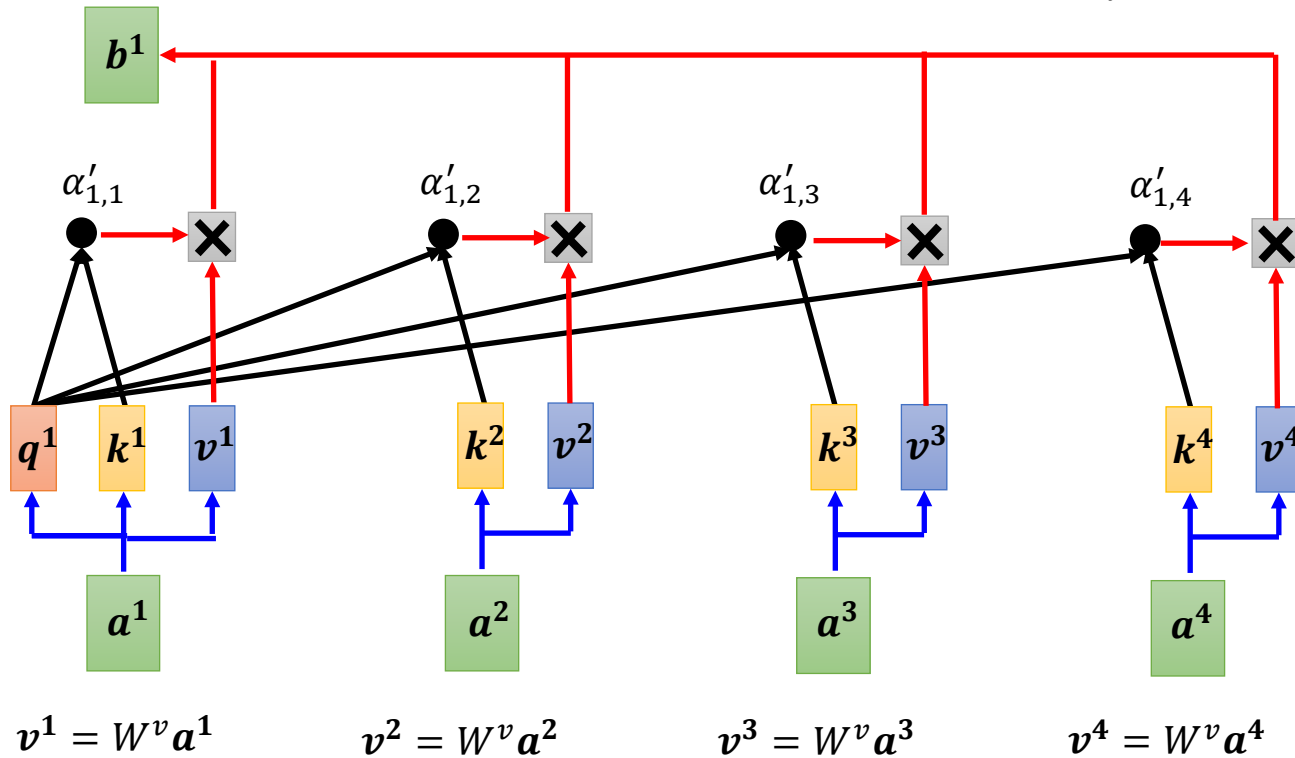
$$\alpha'_{1,i} = \exp(\alpha_{1,i}) / \sum_j \exp(\alpha_{1,j})$$



Self-attention

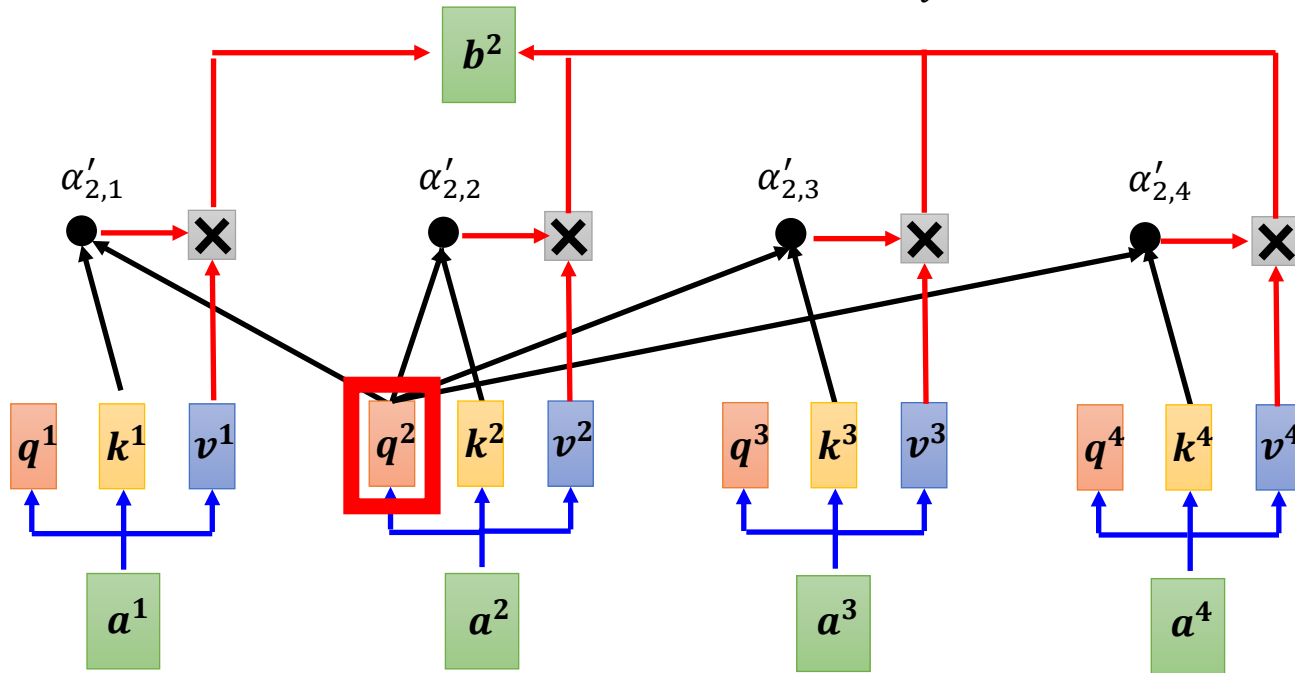
Extract information based on
attention scores

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

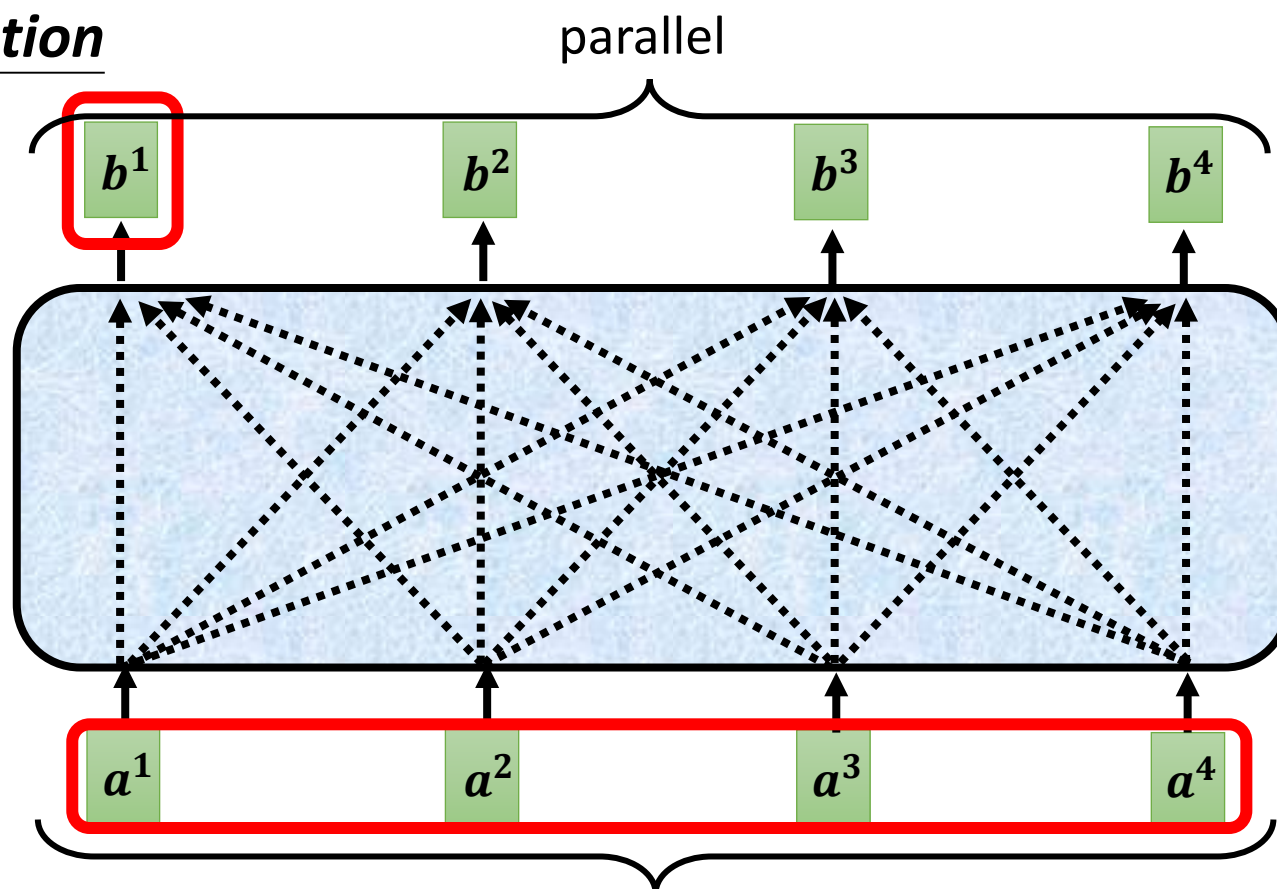


Self-attention

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

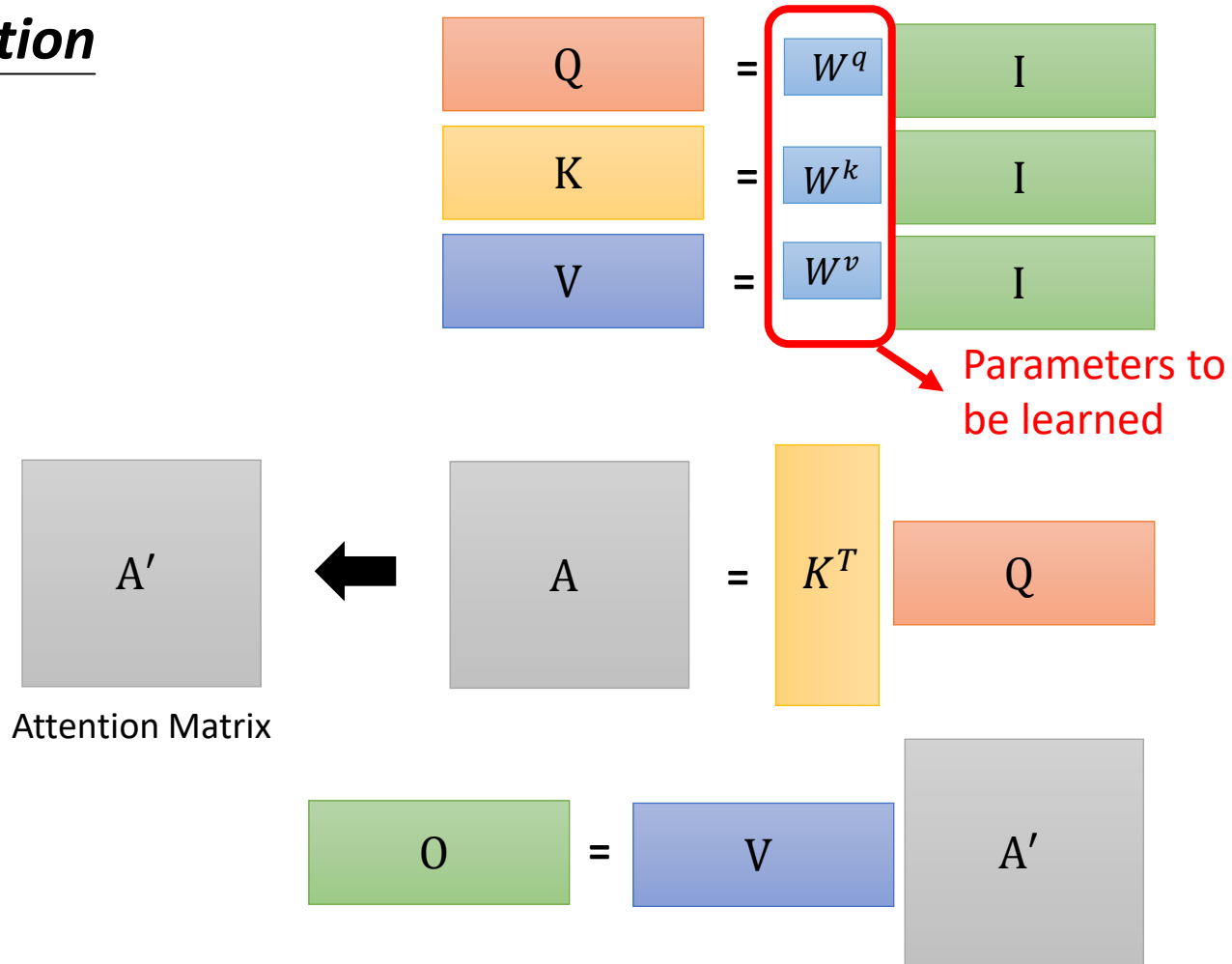


Self-attention



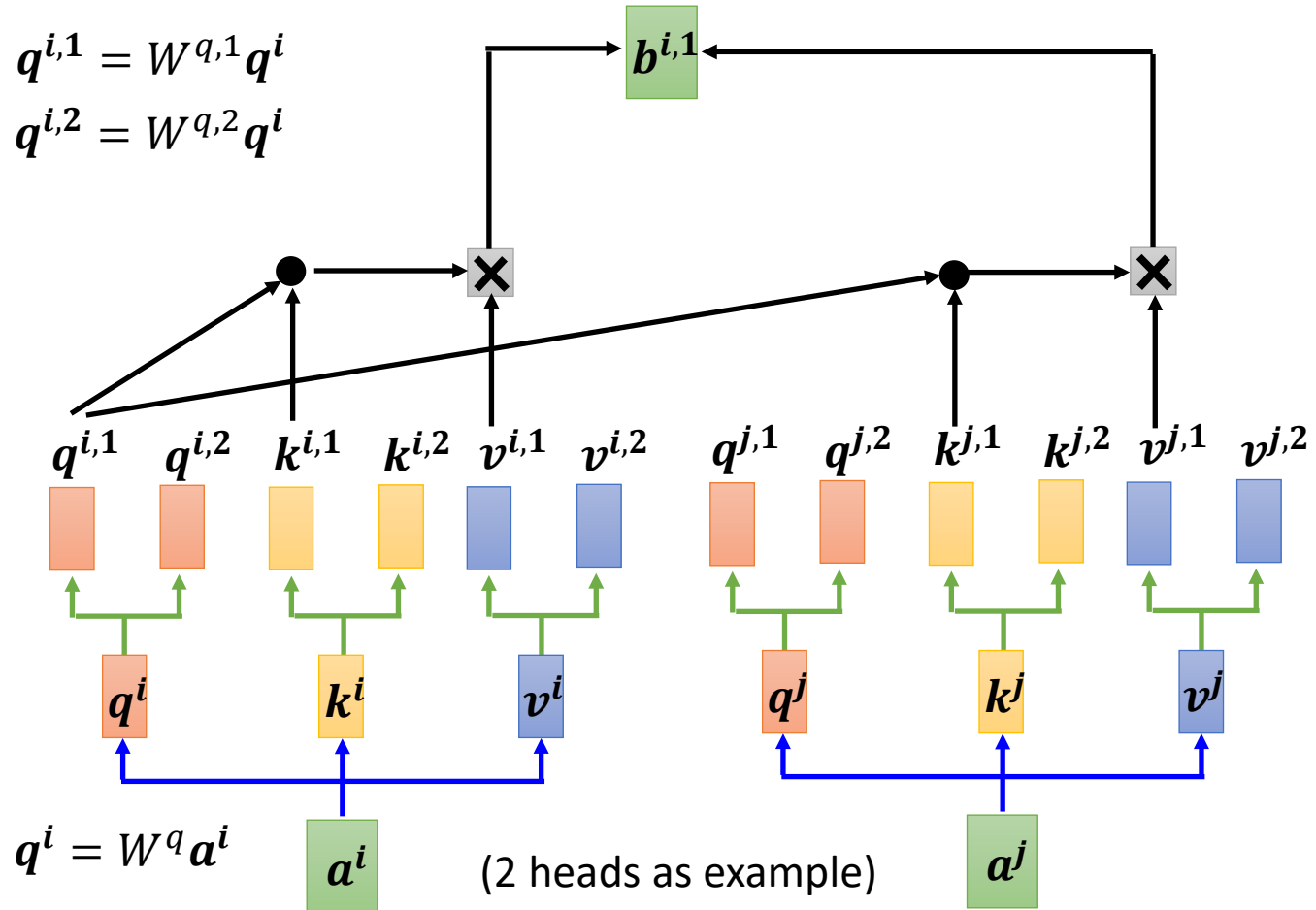
Can be either **input** or a **hidden layer**

Self-attention



Multi-head Self-attention

Different types of relevance

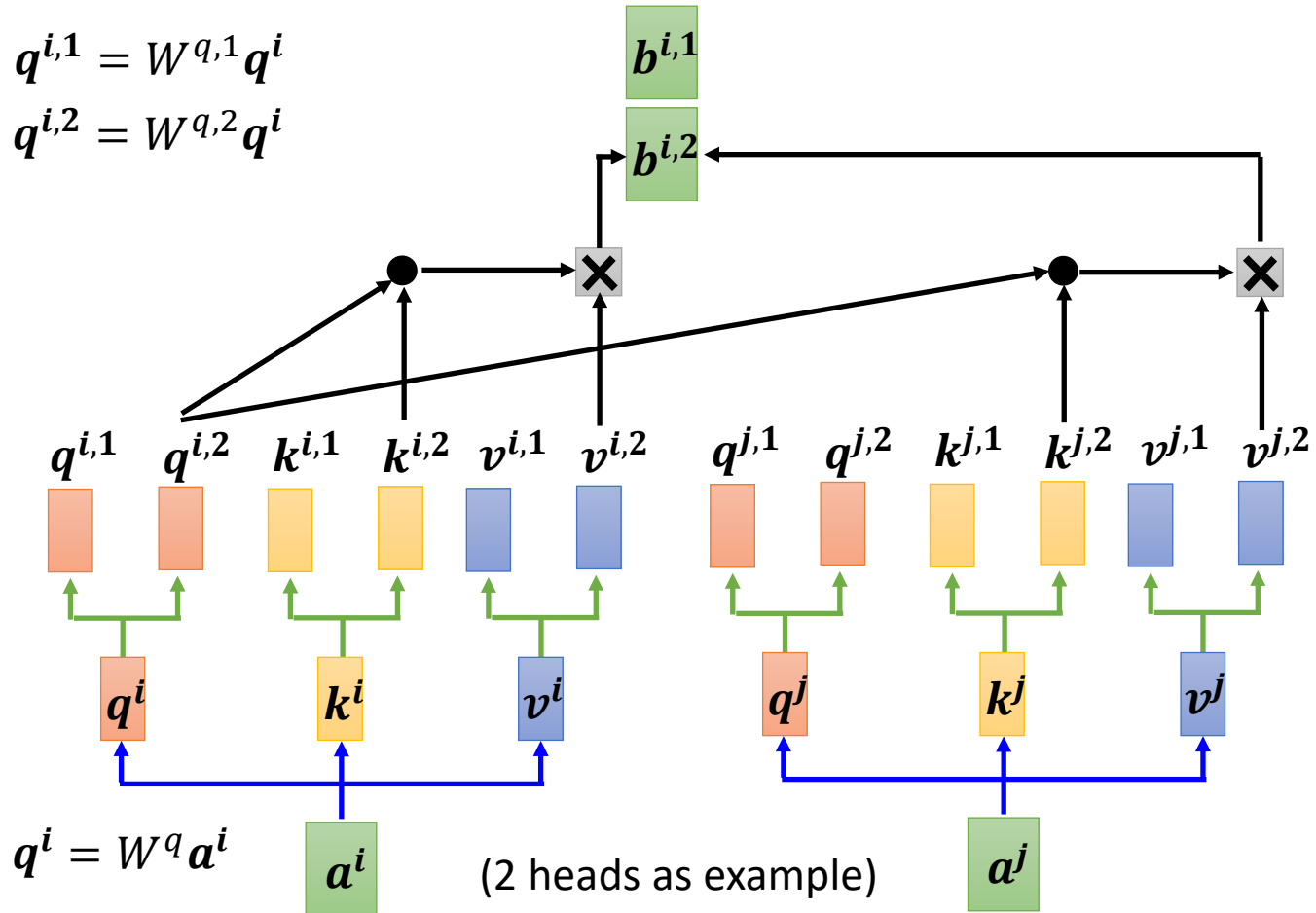


Multi-head Self-attention

Different types of relevance

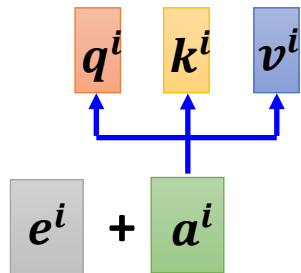
$$q^{i,1} = W^{q,1} q^i$$

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

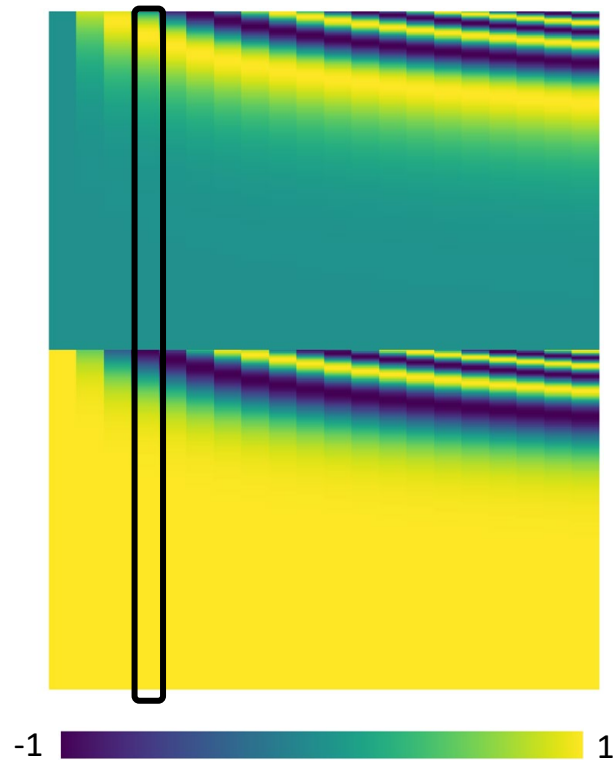


Positional Encoding

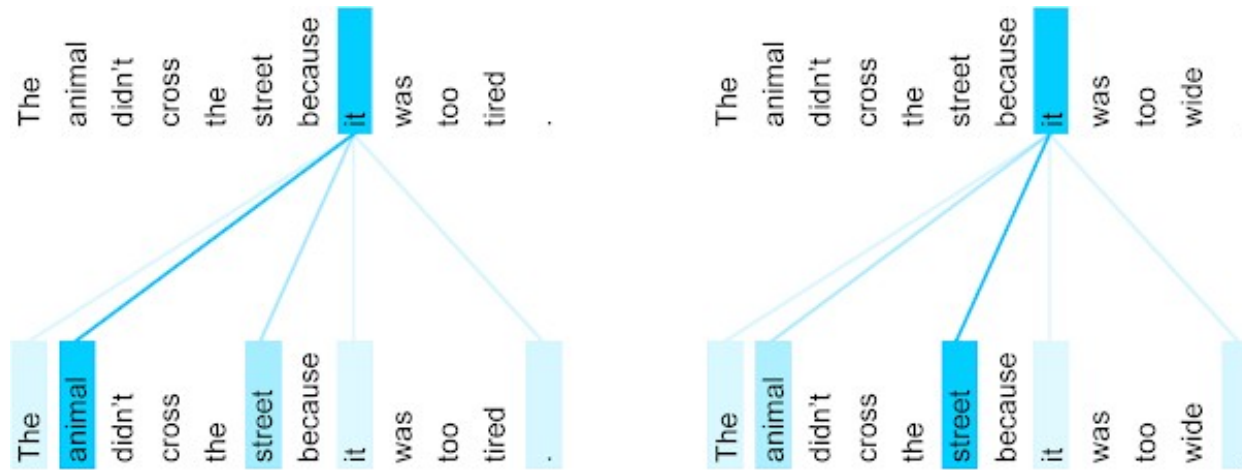
- No position information in self-attention.
- Each position has a unique positional vector e^i
- **hand-crafted**
- **learned from data**



Each column represents a positional vector e^i



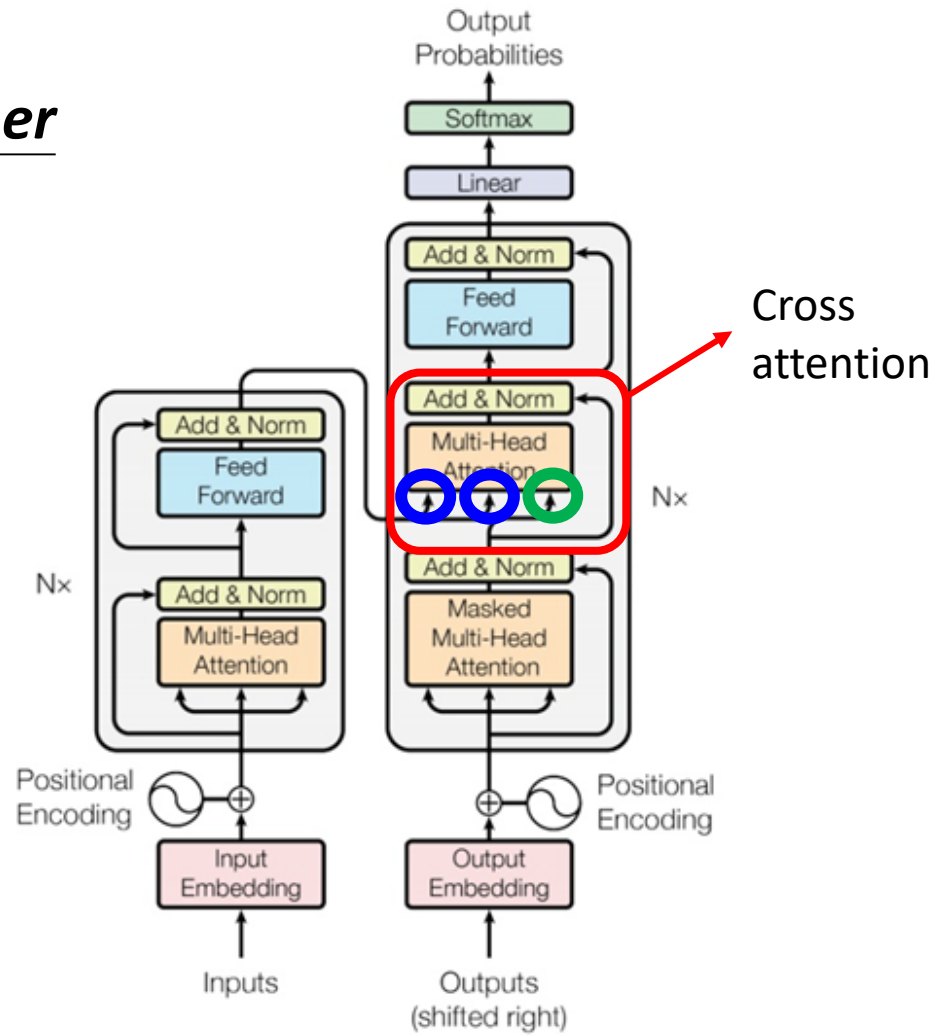
Attention Visualization



The encoder self-attention distribution for the word "it" from the 5th to the 6th layer of a Transformer trained on English to French translation (one of eight attention heads).

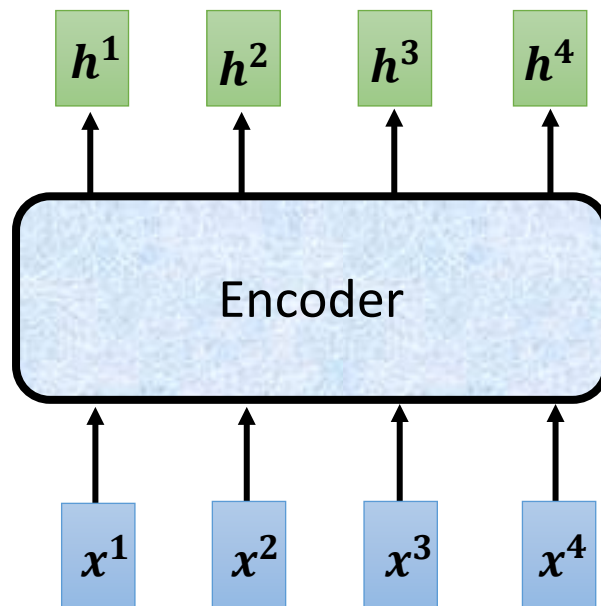
<https://ai.googleblog.com/2017/08/transformer-novel-neural-network.html>

Transformer

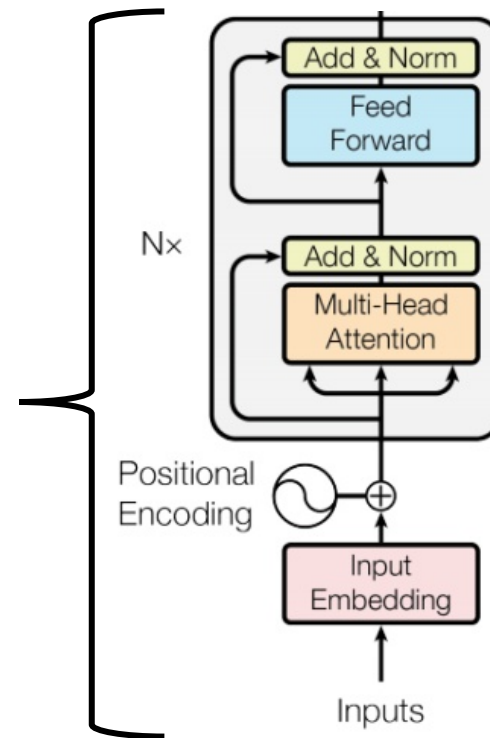


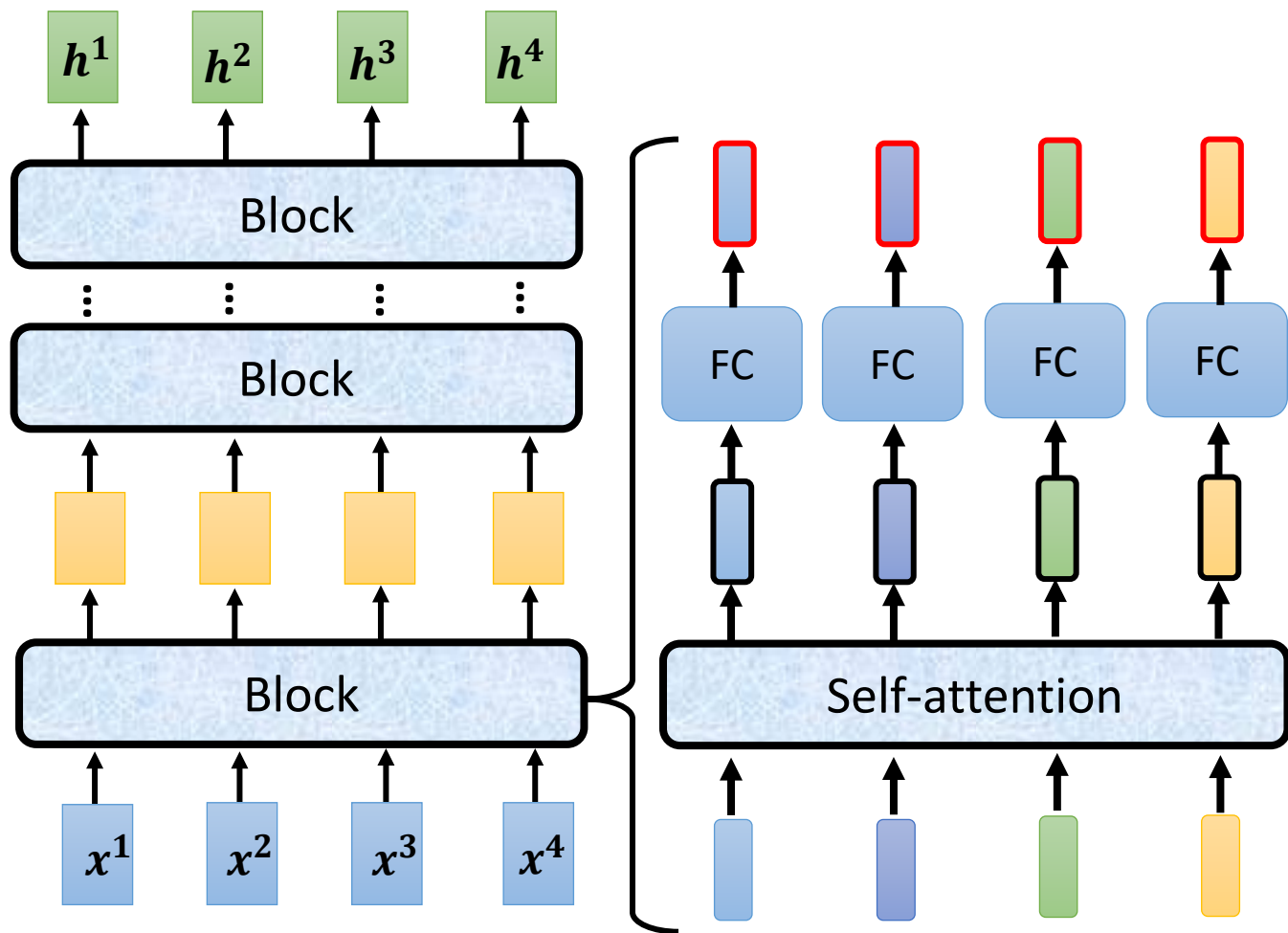
Encoder

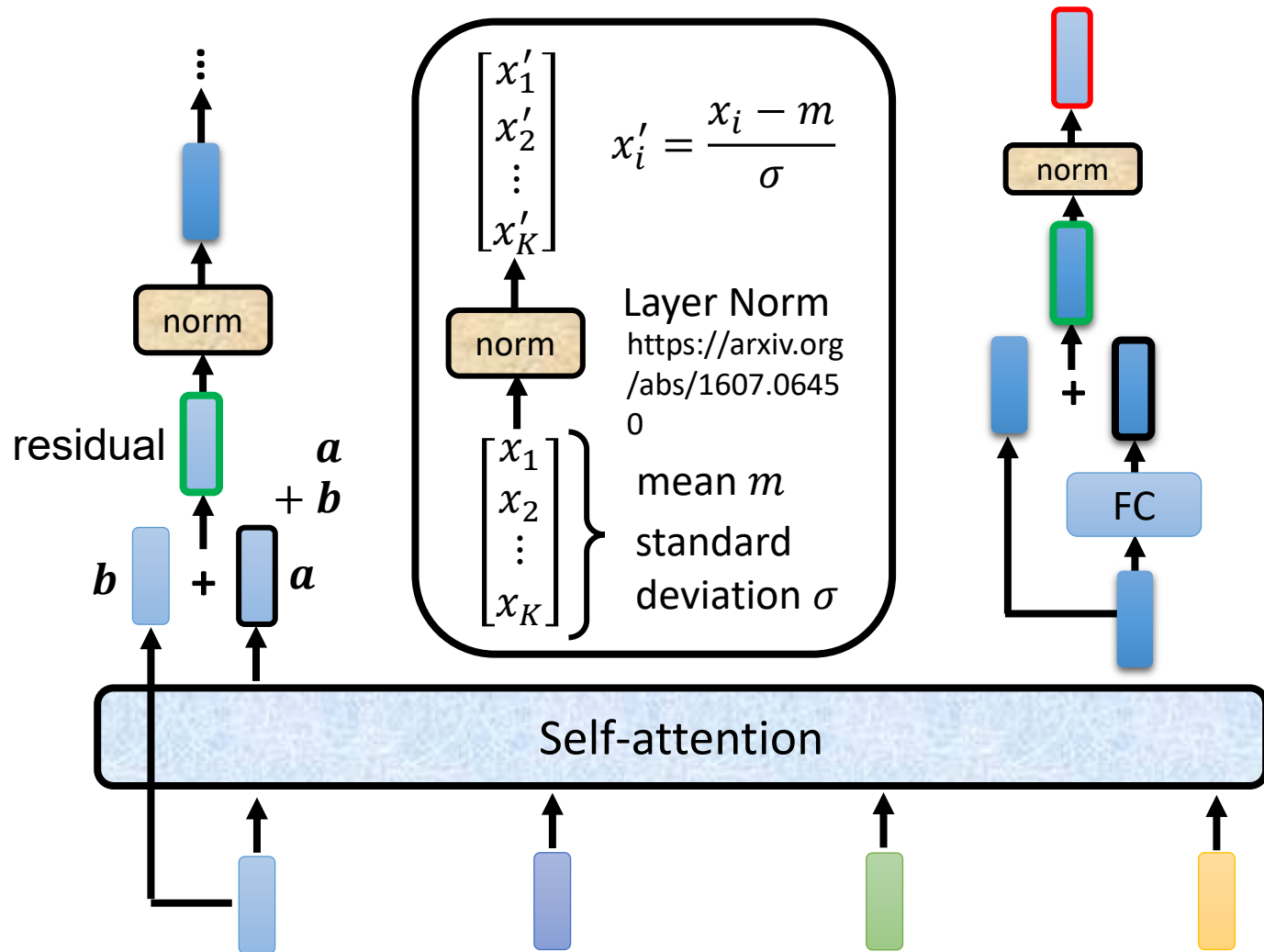
You can use **RNN** or **CNN**.



Transformer's Encoder

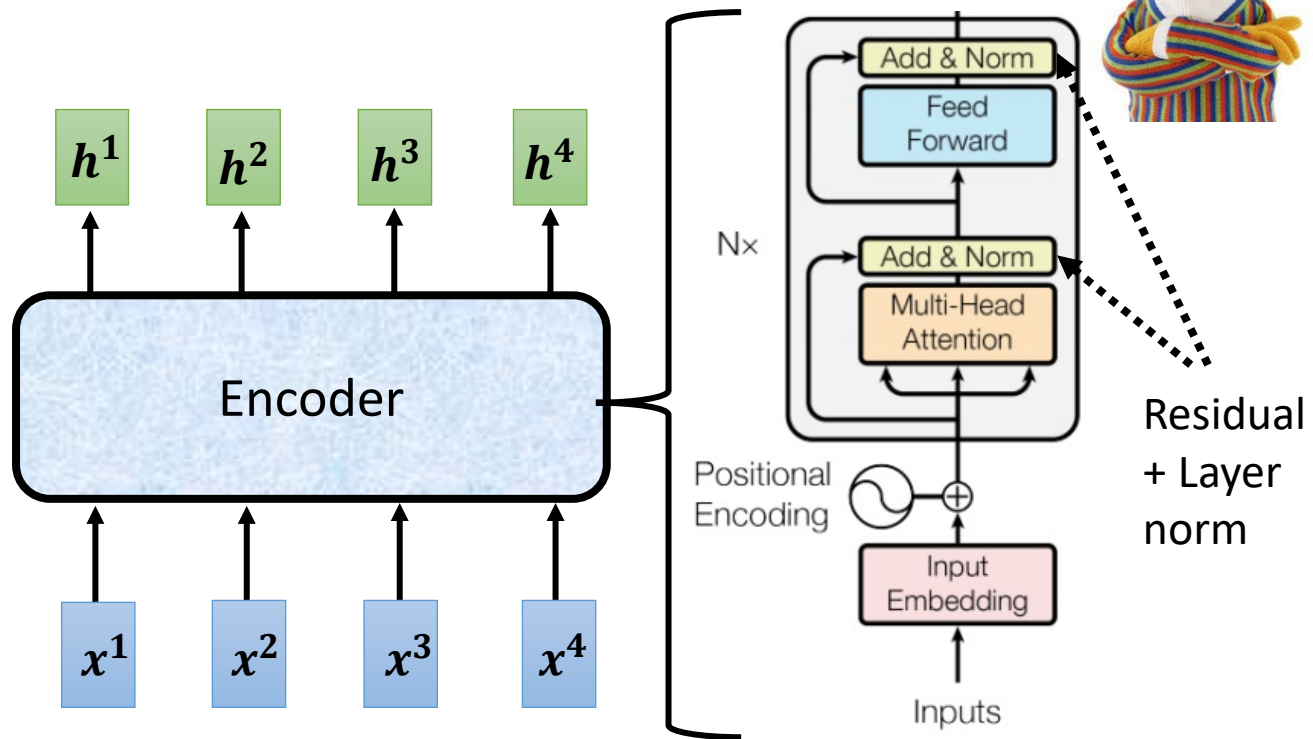






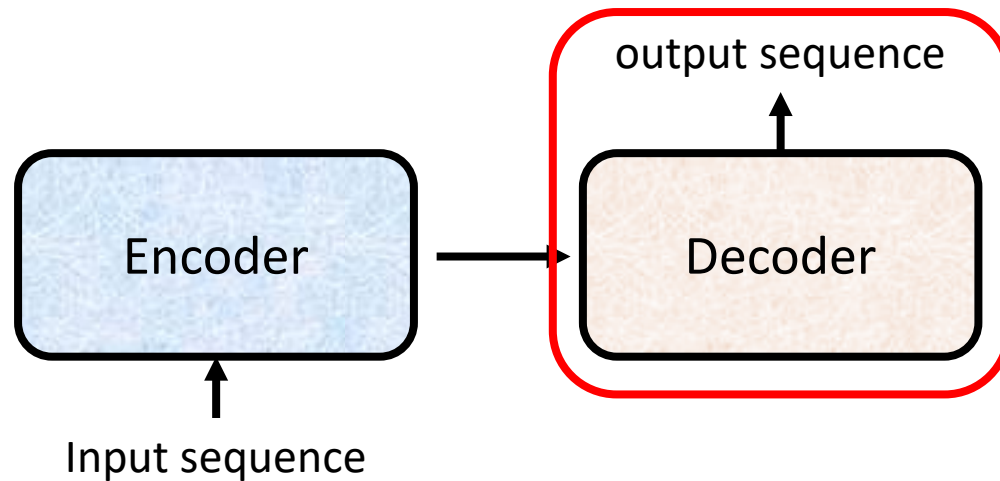
BERT

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



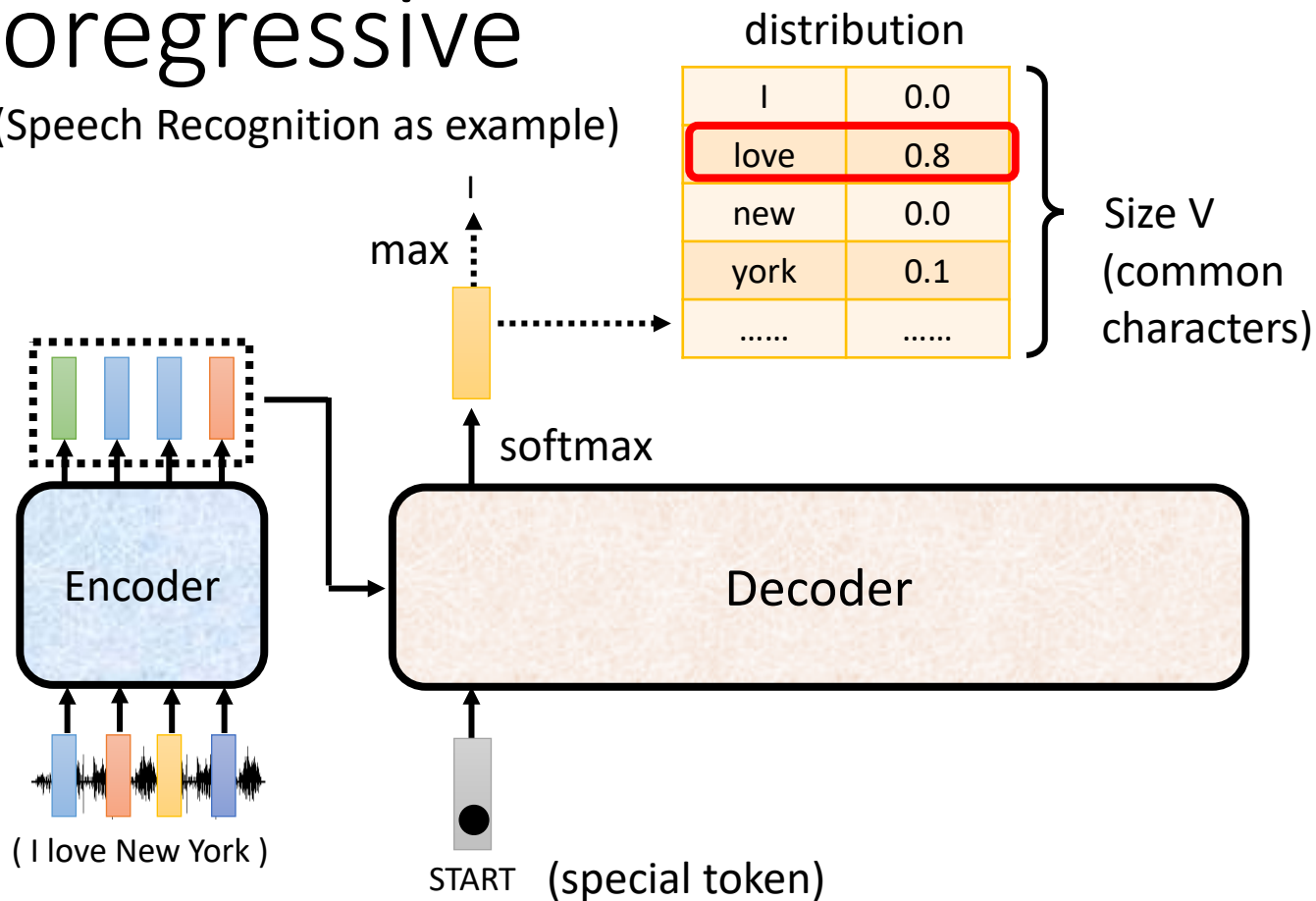
Decoder

- Autoregressive (AT)

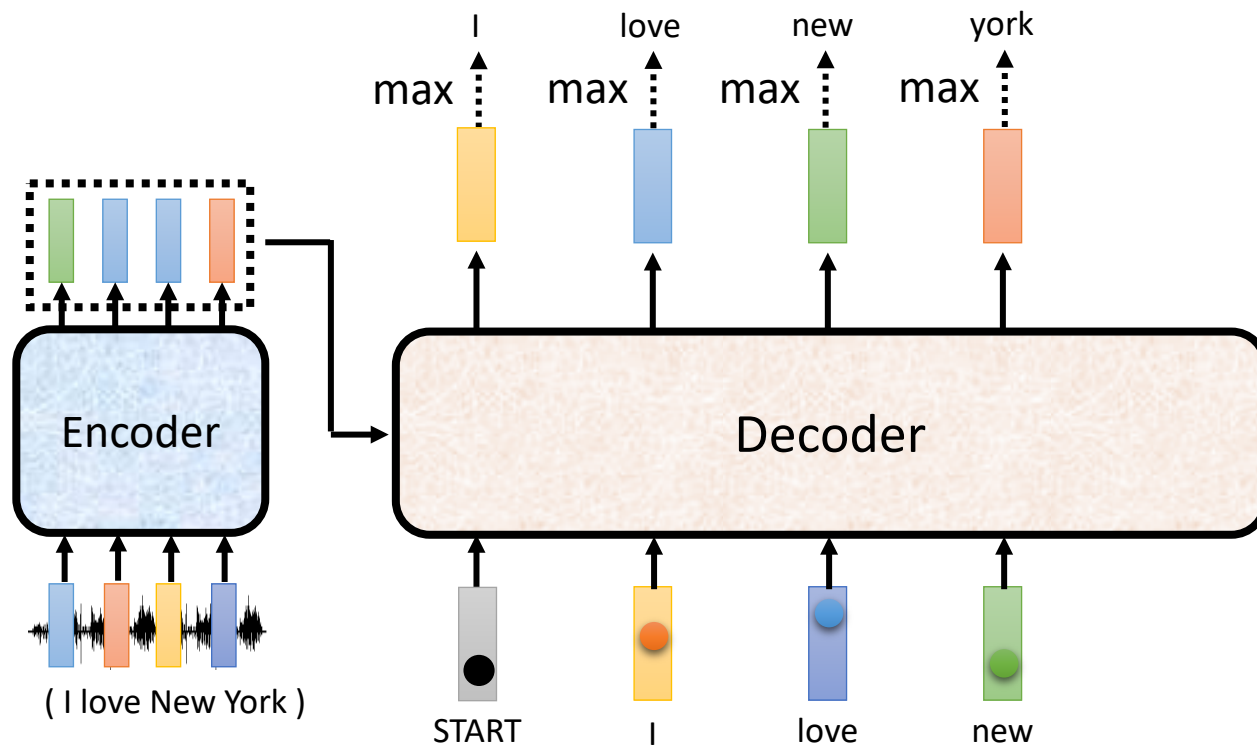


Autoregressive

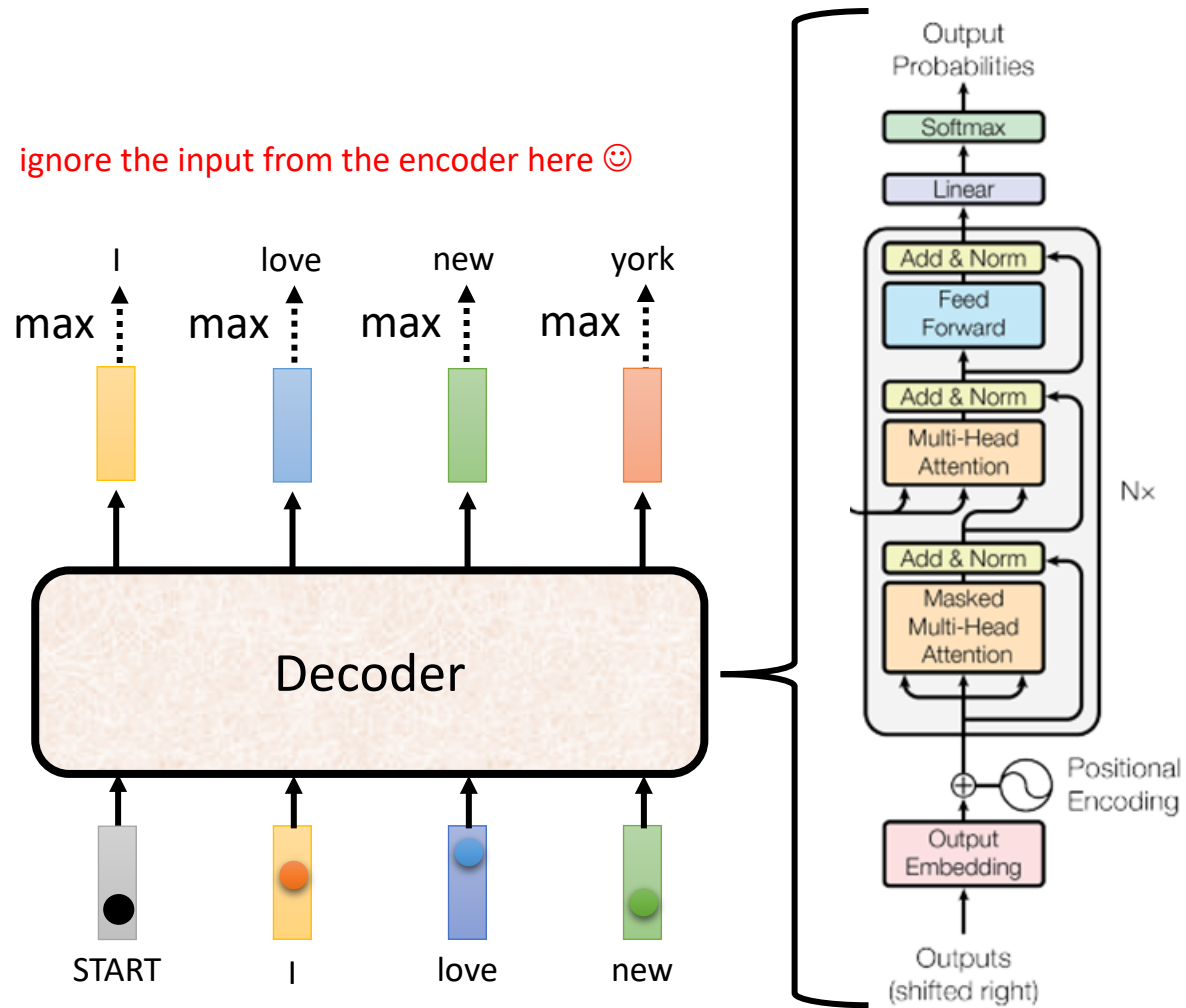
(Speech Recognition as example)



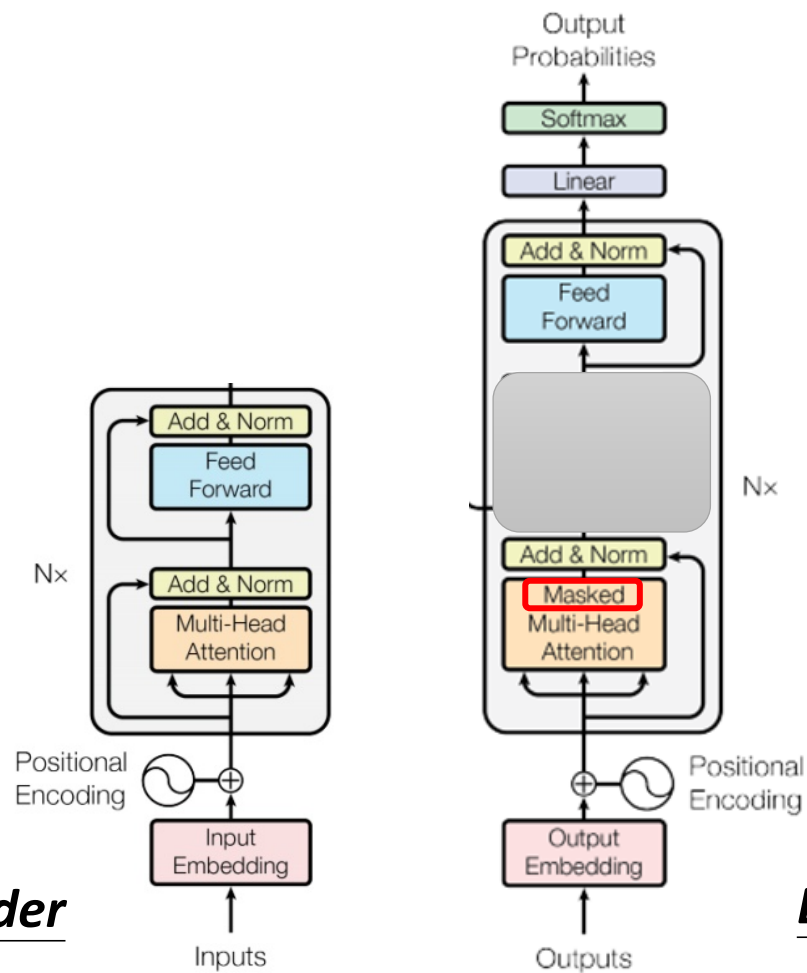
Autoregressive



ignore the input from the encoder here 😊

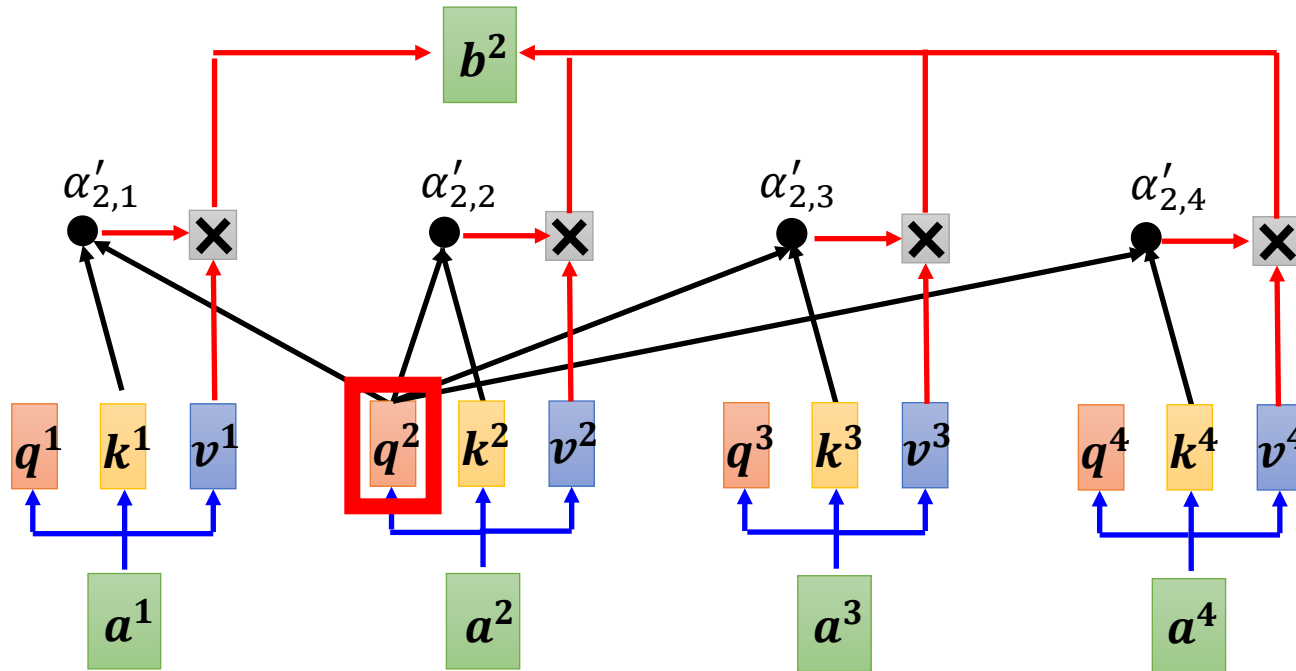


Encoder



Decoder

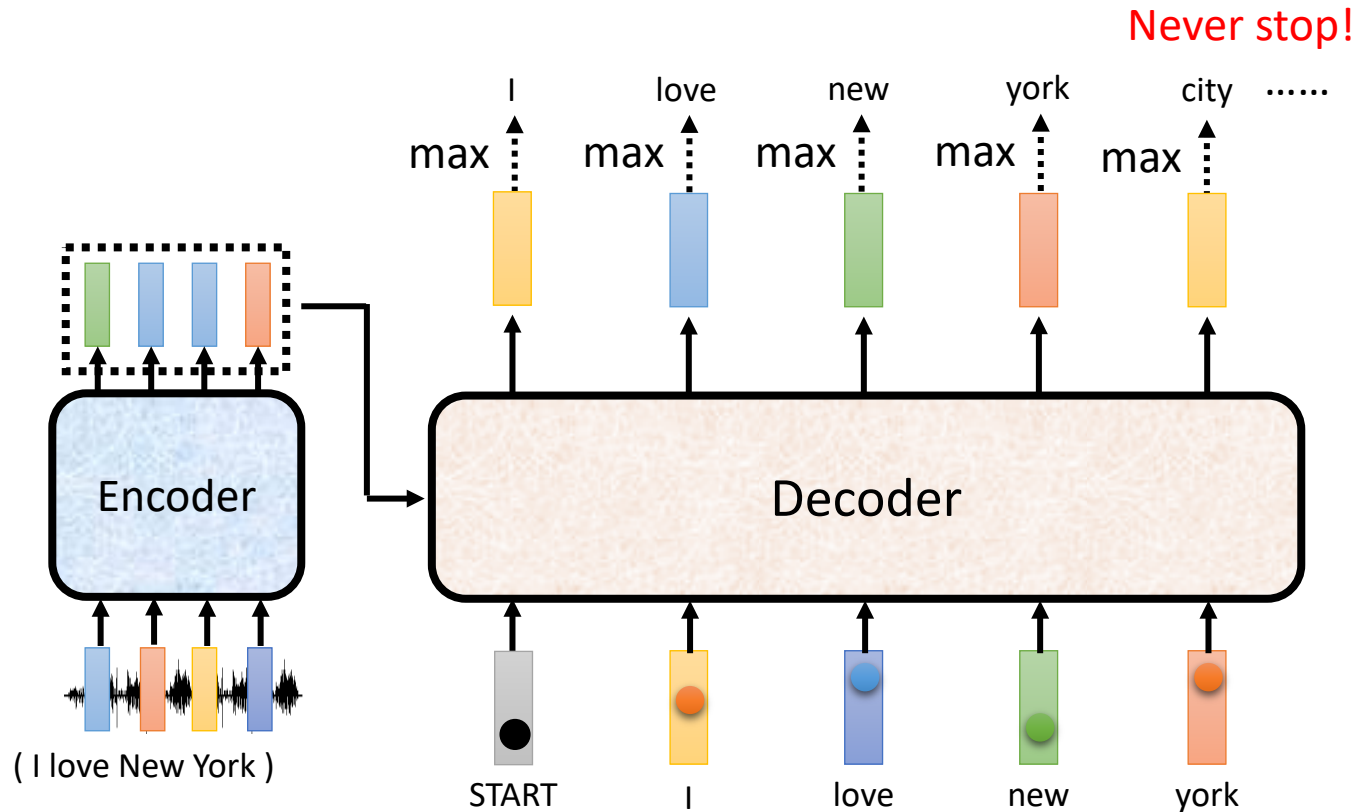
Self-attention ➔ Masked Self-attention



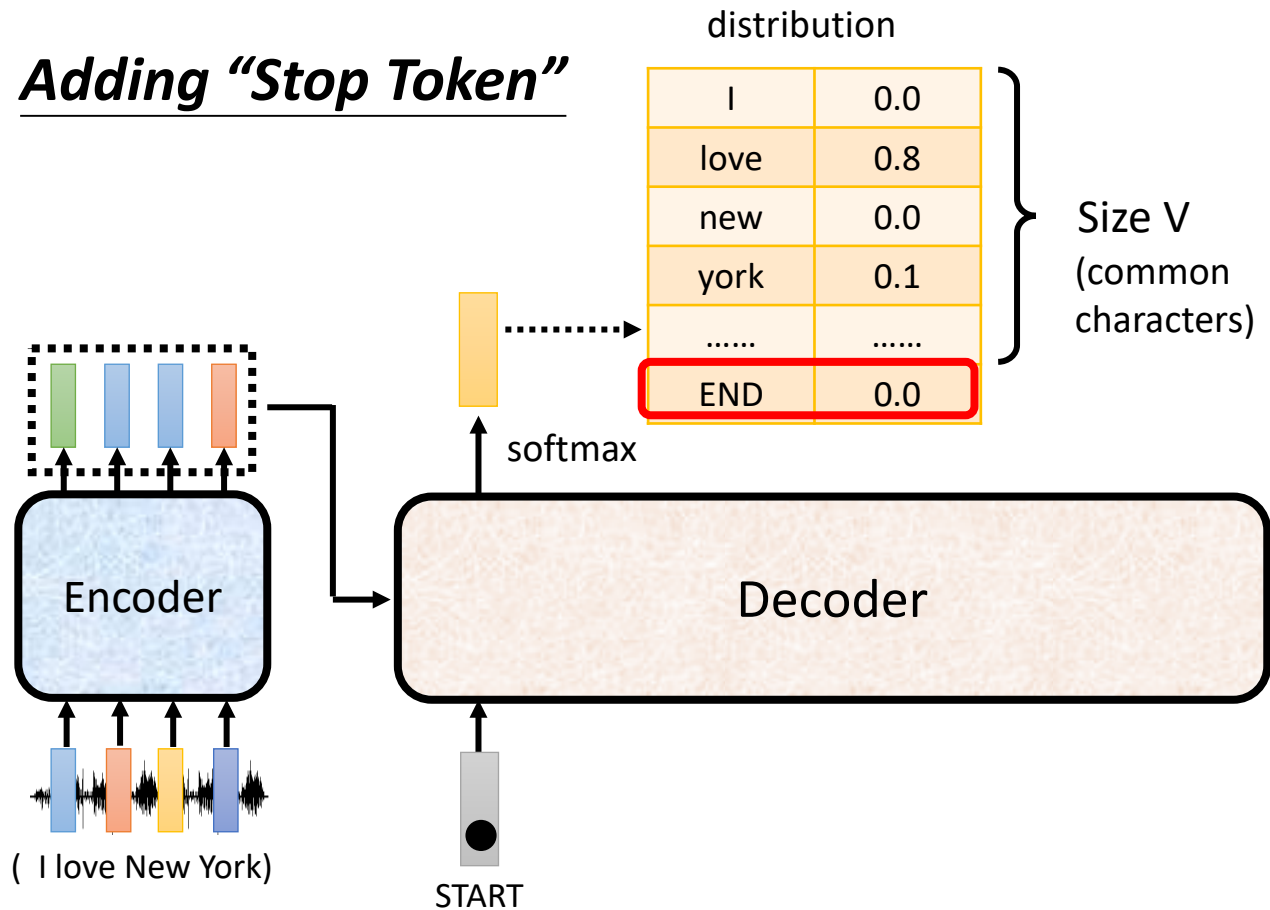
Why masked? Consider how does decoder work

Autoregressive

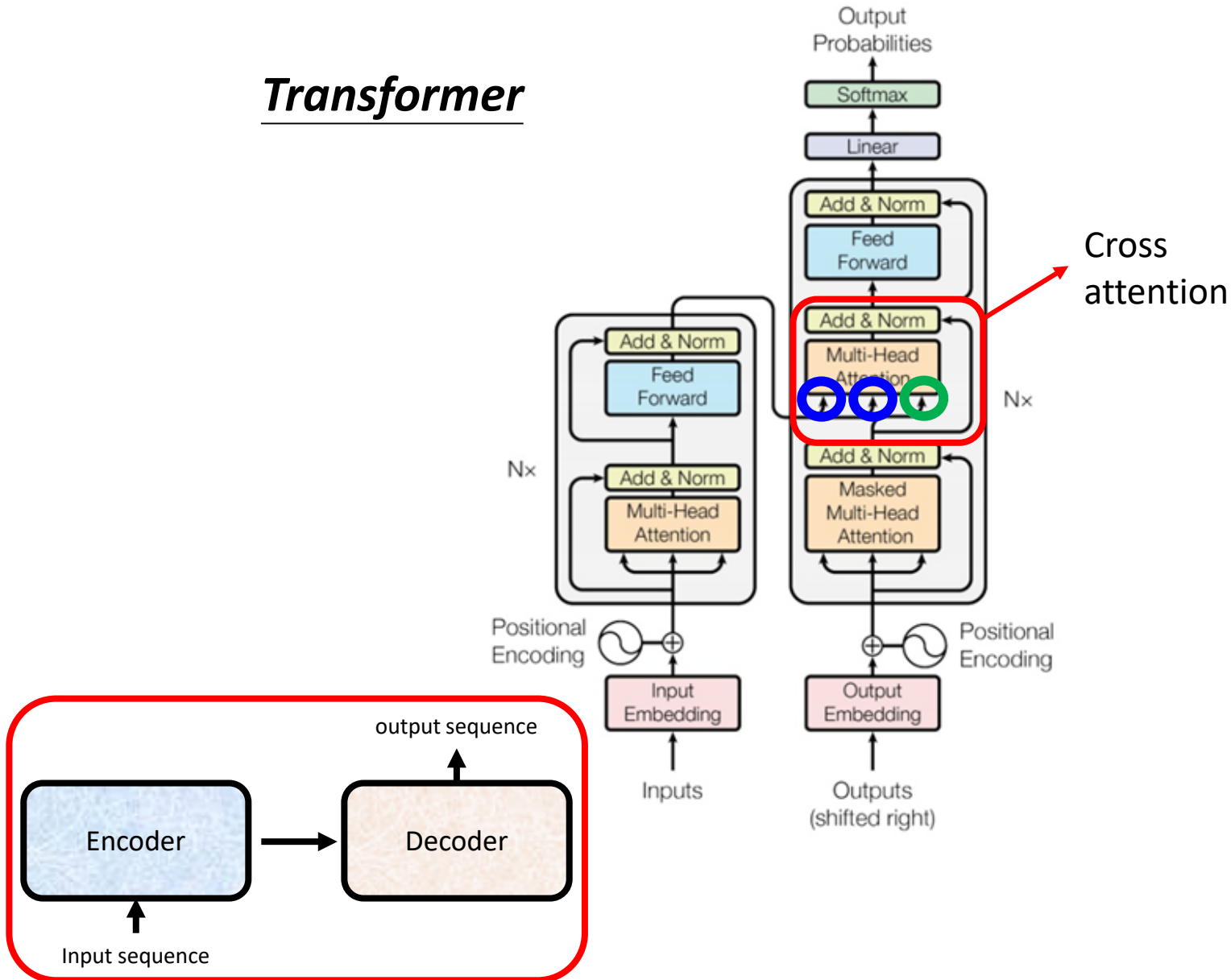
We do not know the correct output length.

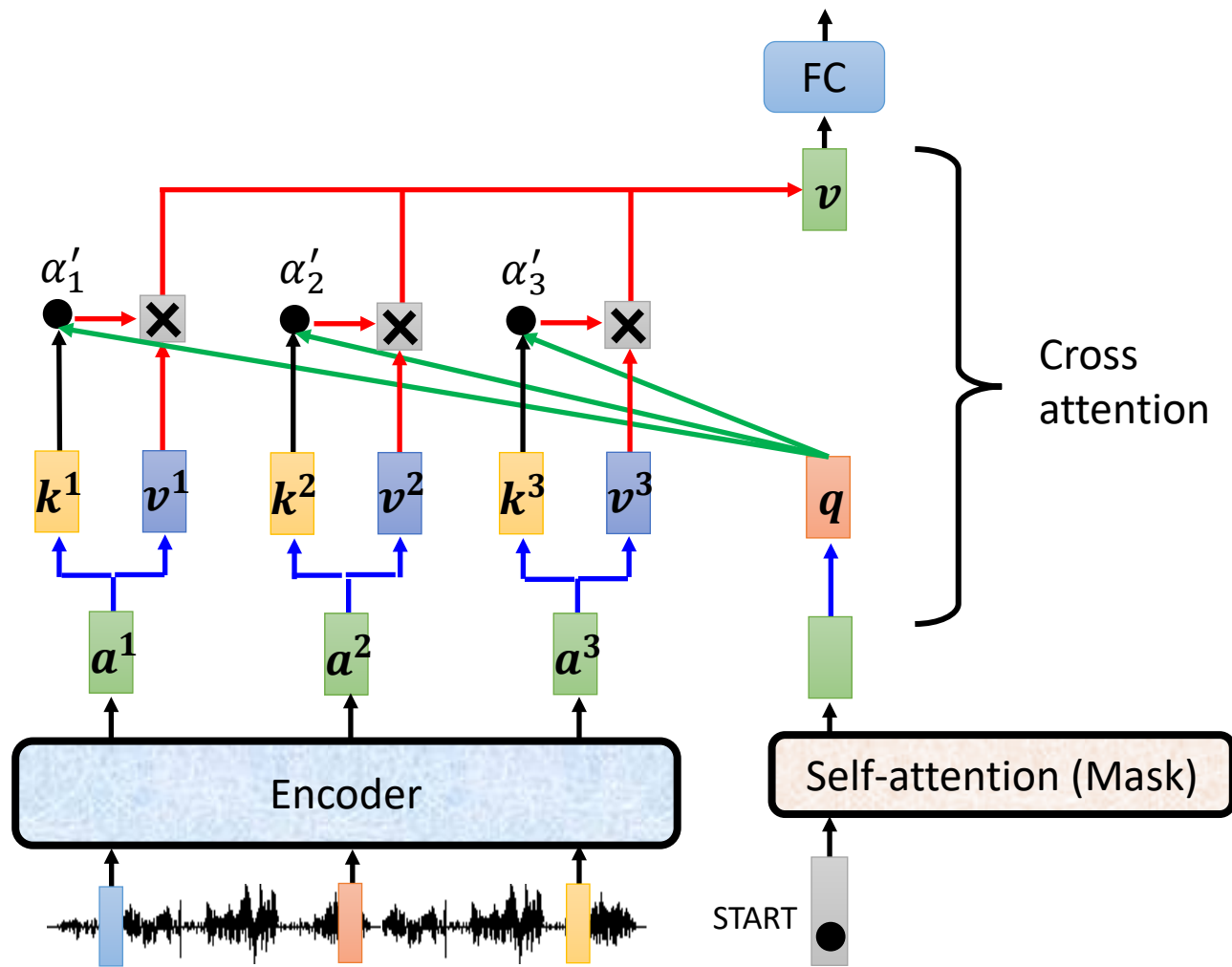


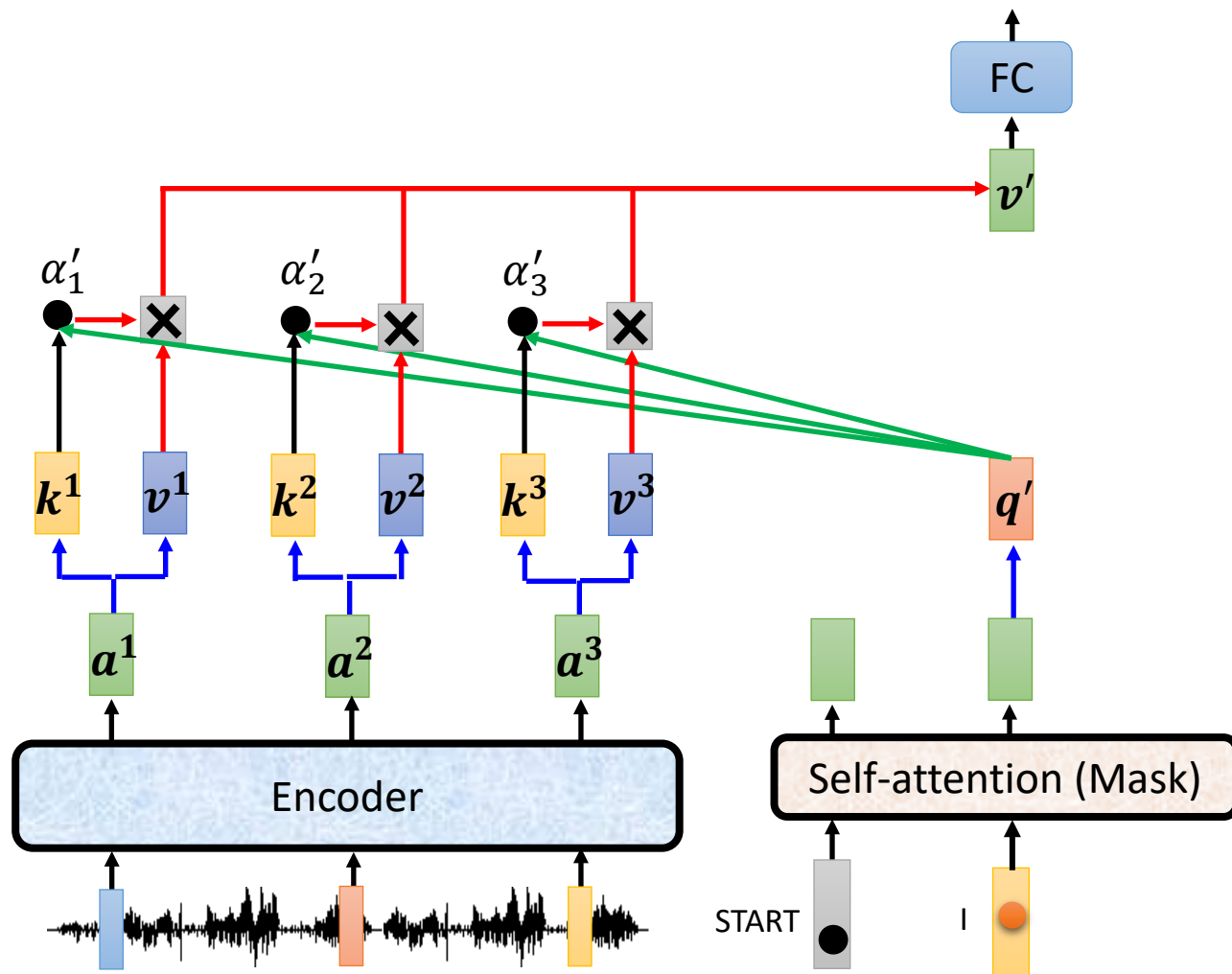
Adding “Stop Token”

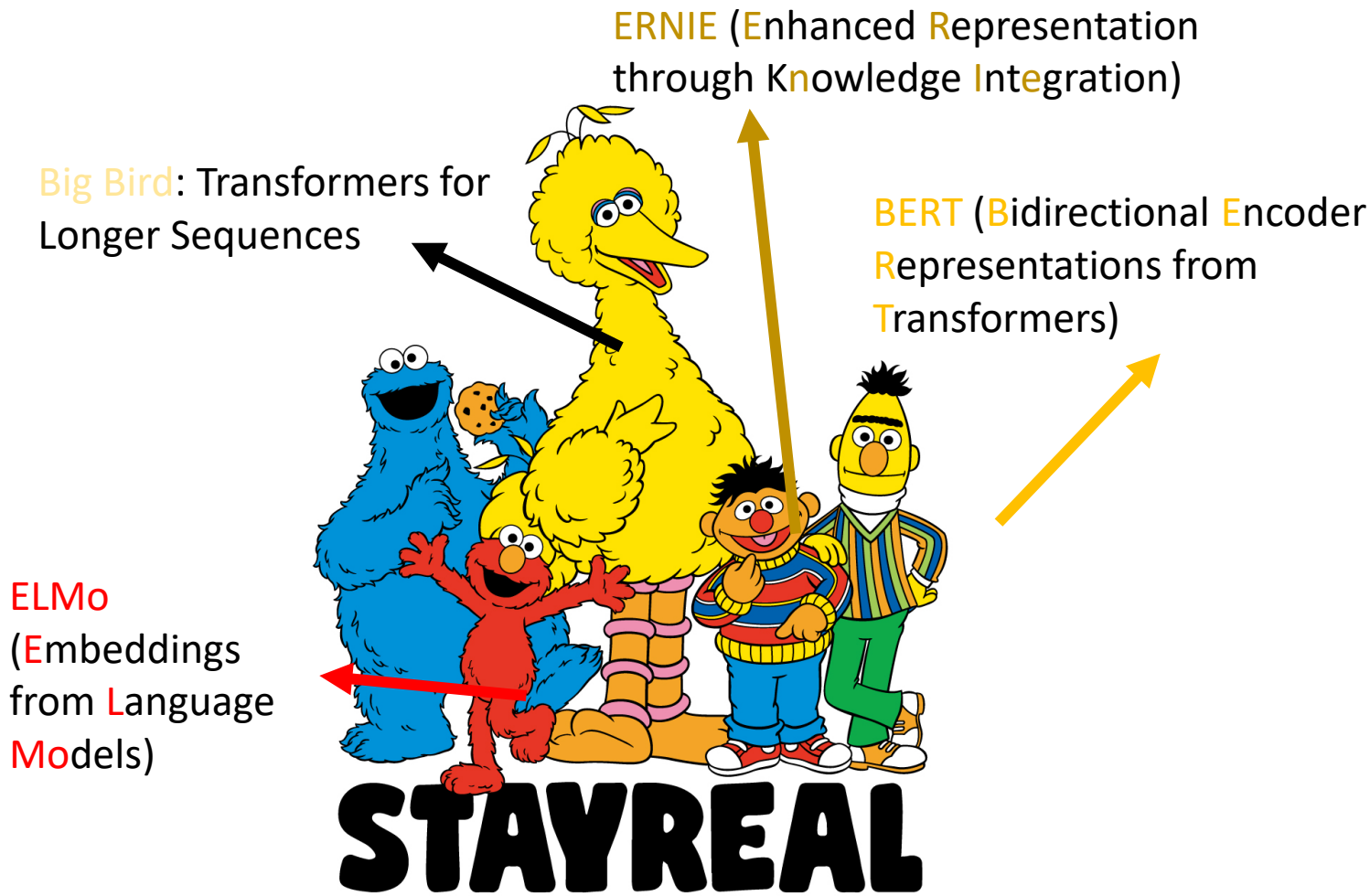


Transformer







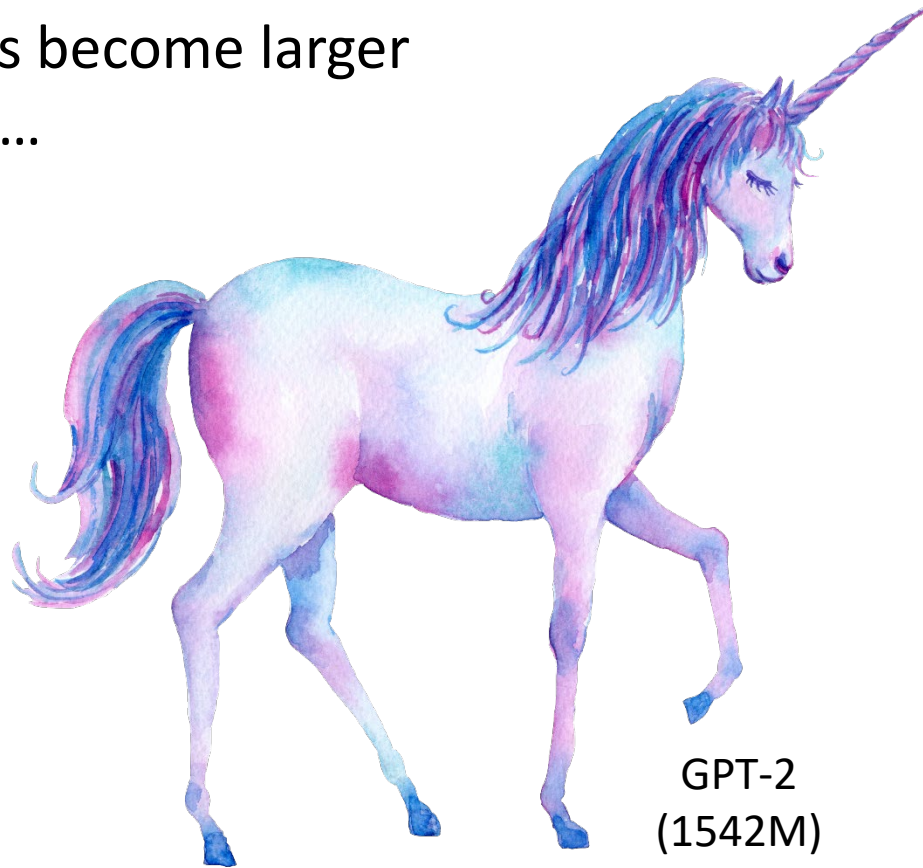


The models become larger
and larger ...

ELM
O
(94M)



BERT
(340M)



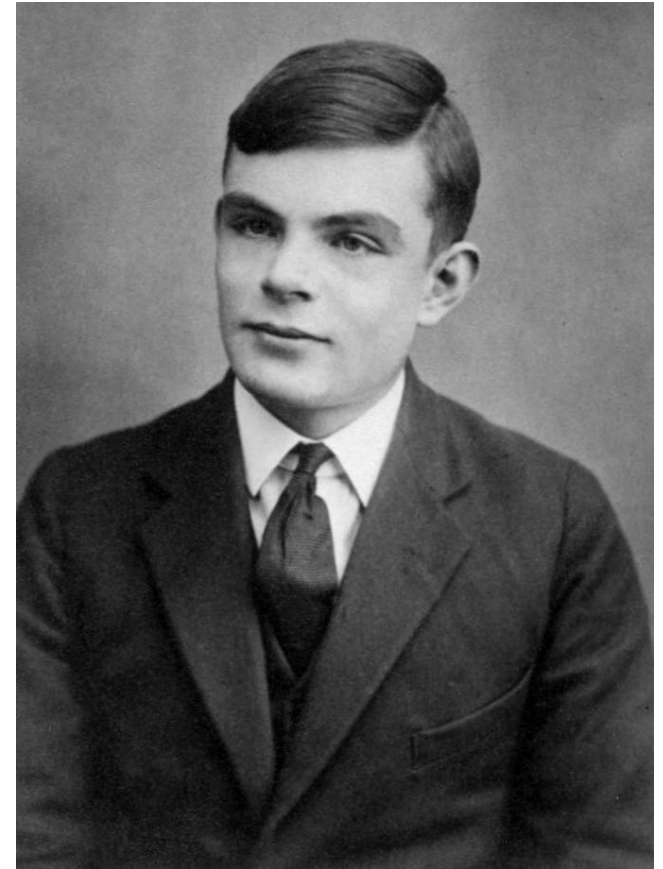
GPT-2
(1542M)

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

The models become larger
and larger ...

GPT-3 is **10** times larger than
Turing NLG.

Turing
NLG (17B)



GPT-2



Megatron (8B)

Outline

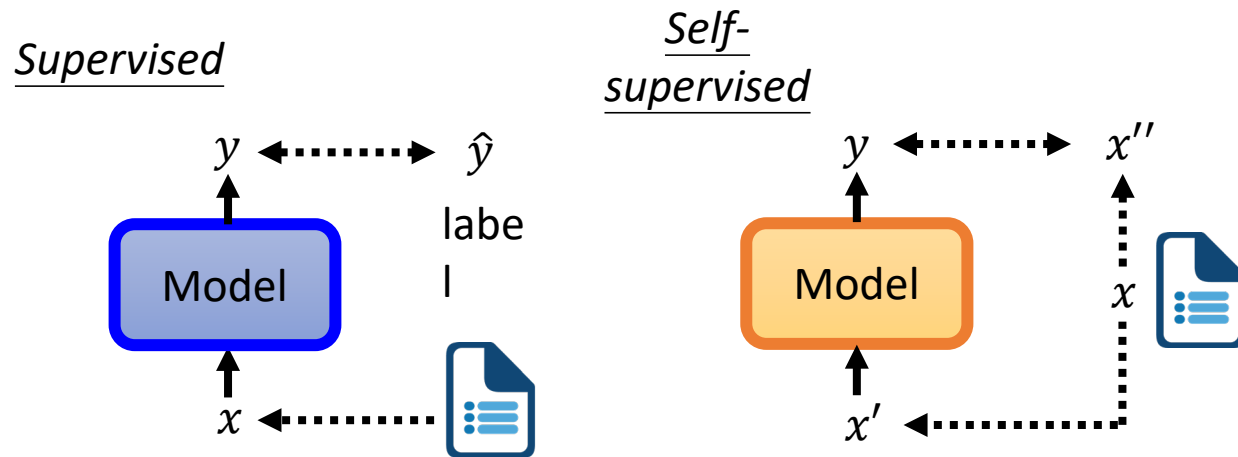


BERT series



GPT series

Self-supervised Learning

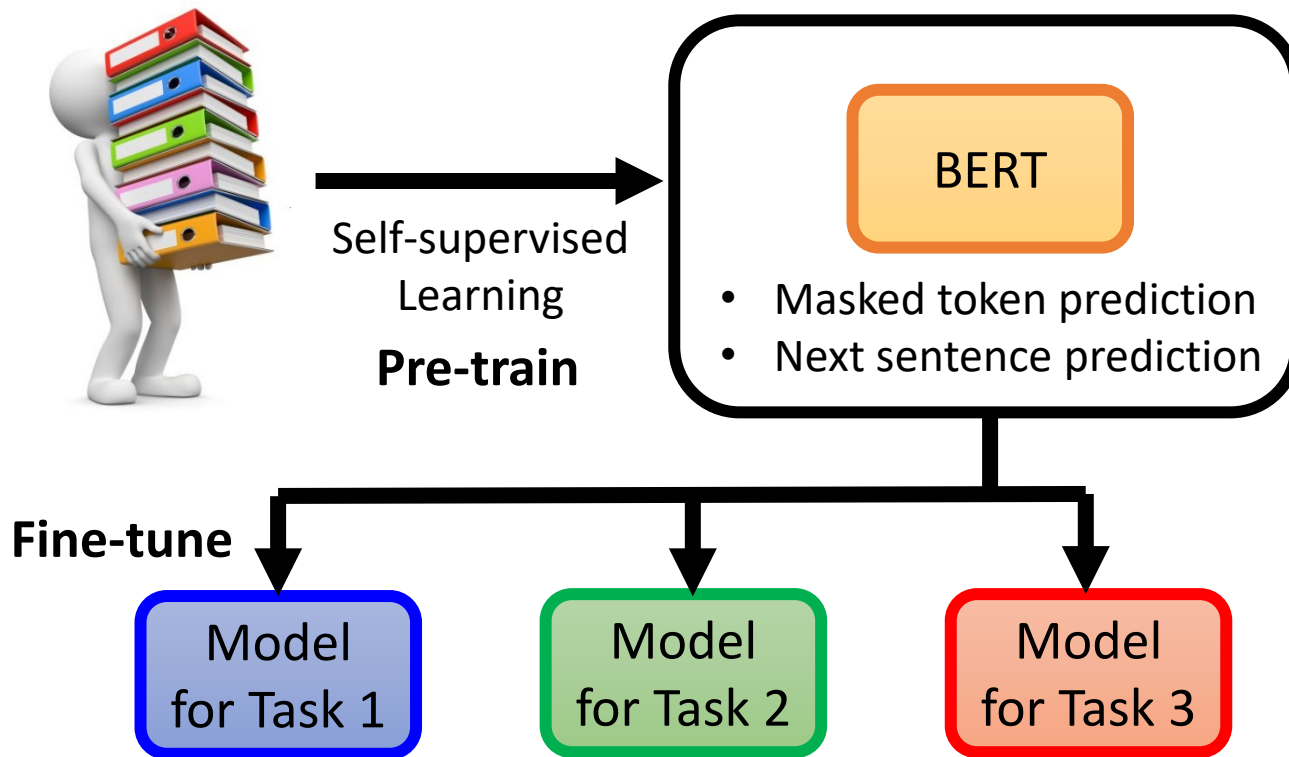


Yann LeCun

2019年4月30日 · 🌐

I now call it "self-supervised learning", because "unsupervised" is both a loaded and confusing term.

In self-supervised learning, the system learns to predict part of its input from other parts of its input. In other words a portion of the input is used as a supervisory signal to a predictor fed with the remaining portion of the input.

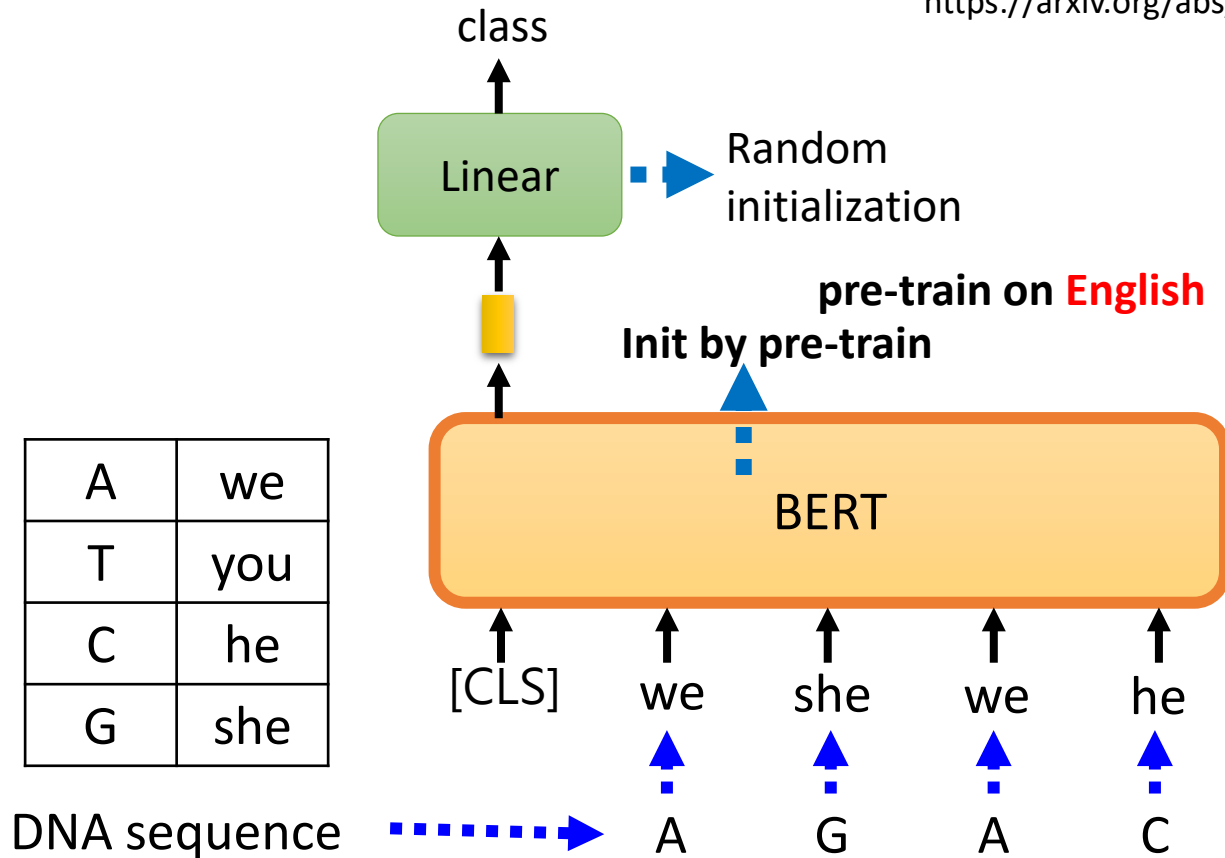


Downstream Tasks

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

Why does BERT work?

<https://arxiv.org/abs/2103.07162>



Outline

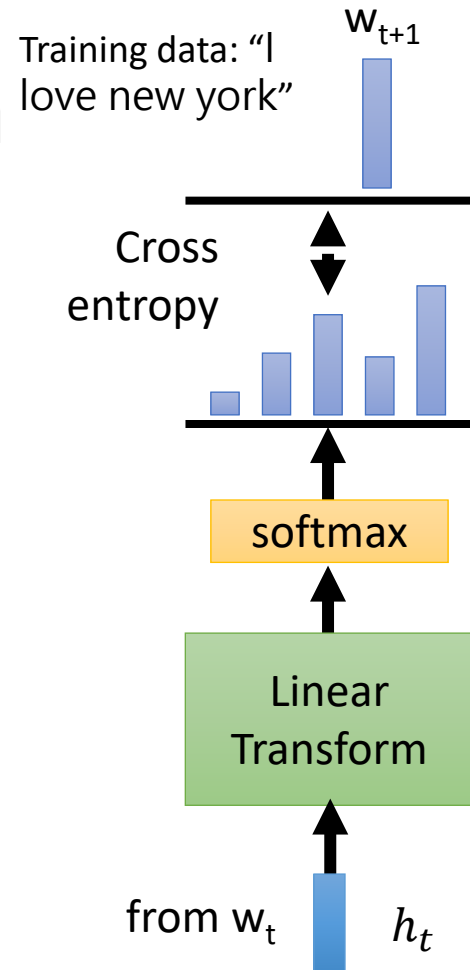
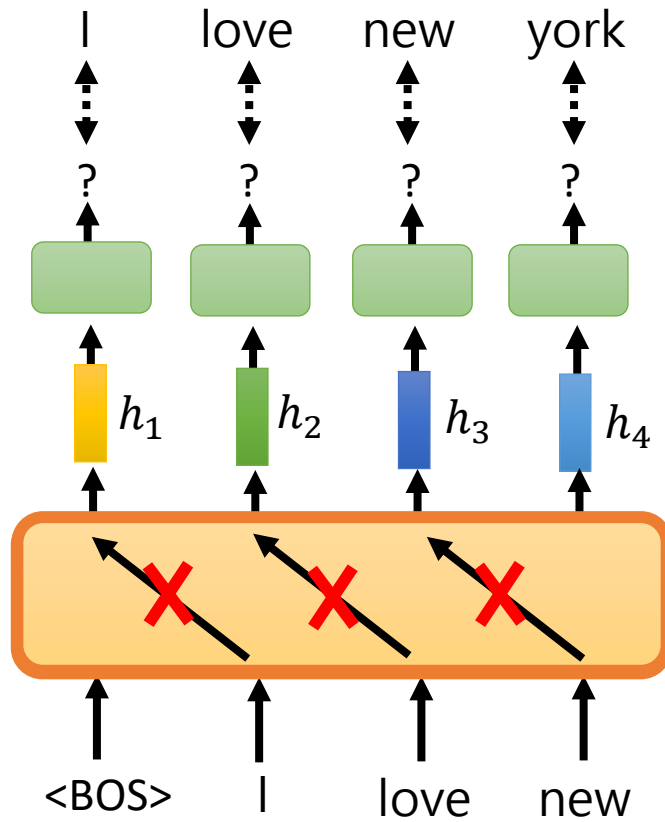


BERT series



GPT series

Predict Next Token



“In-context” Learning

Few-shot Learning

(no gradient descent)

| | | |
|---|--------------------------------|--------------------|
| 1 | Translate English to French: | ← task description |
| 2 | sea otter => loutre de mer | ← examples |
| 3 | peppermint => menthe poivrée | ← |
| 4 | plush girafe => girafe peluche | ← |
| 5 | cheese => | ← prompt |

One-shot Learning

| | | |
|---|------------------------------|--------------------|
| 1 | Translate English to French: | ← task description |
| 2 | sea otter => loutre de mer | ← example |
| 3 | cheese => | ← prompt |

Zero-shot Learning

| | | |
|---|------------------------------|--------------------|
| 1 | Translate English to French: | ← task description |
| 2 | cheese => | ← prompt |

GPT-3

Language Models are Few-Shot Learners

| | | | | |
|---------------------------|-------------------|--------------------|--------------------|----------------|
| Tom B. Brown* | Benjamin Mann* | Nick Ryder* | Melanie Subbiah* | |
| Jared Kaplan [†] | Prafulla Dhariwal | Arvind Neelakantan | Pranav Shyam | Girish Sastry |
| Amanda Askell | Sandhini Agarwal | Ariel Herbert-Voss | Gretchen Krueger | Tom Henighan |
| Rewon Child | Aditya Ramesh | Daniel M. Ziegler | Jeffrey Wu | Clemens Winter |
| Christopher Hesse | Mark Chen | Eric Sigler | Mateusz Litwin | Scott Gray |
| Benjamin Chess | | Jack Clark | Christopher Berner | |
| Sam McCandlish | Alec Radford | Ilya Sutskever | Dario Amodei | |

OpenAI