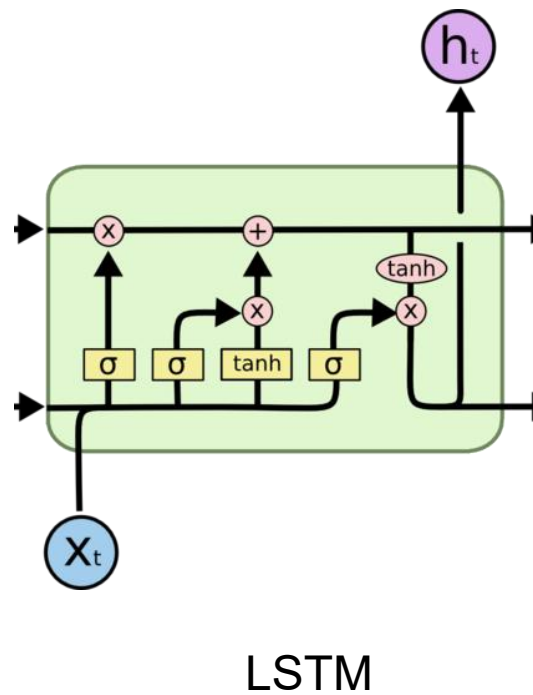
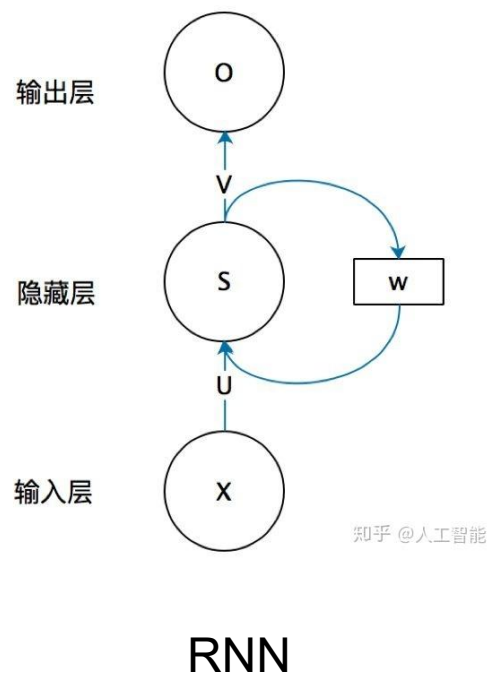


Transformer

Reporter: Bowen Xu

Background

Traditional translation model

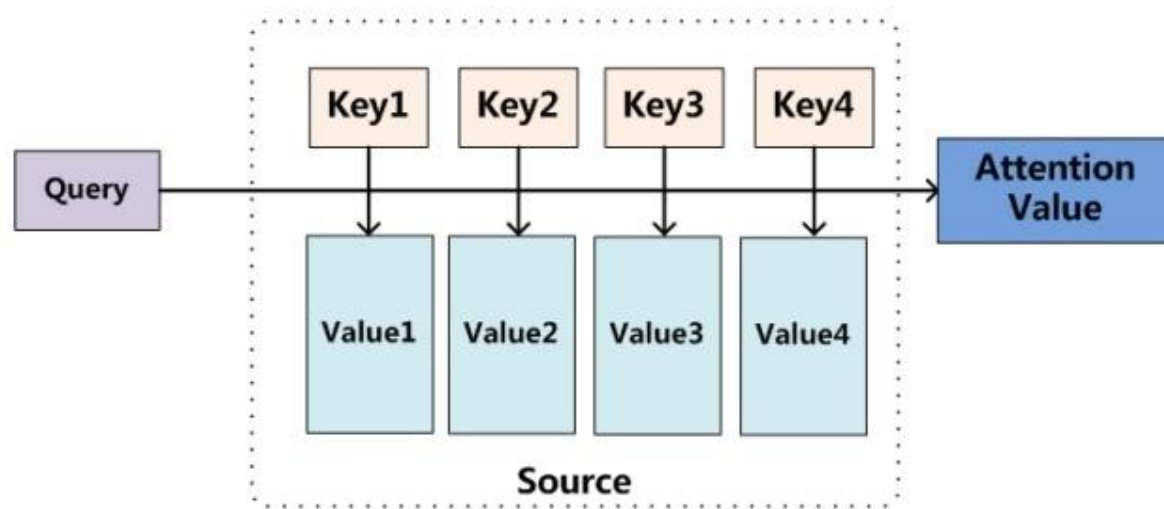


Drawbacks:

- Sequential nature precludes **parallelization**
- **Information loss** in long term

Introduction

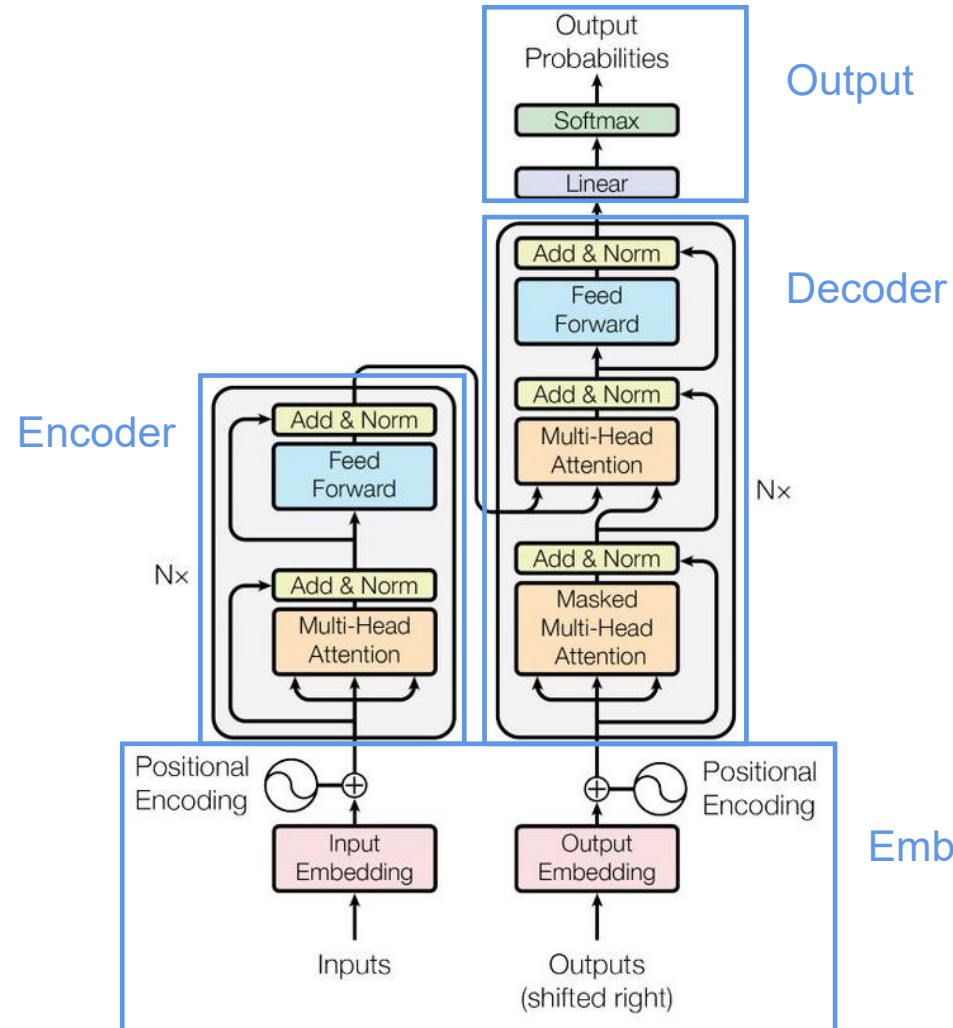
Based entirely on Attention Mechanism



Advantages:

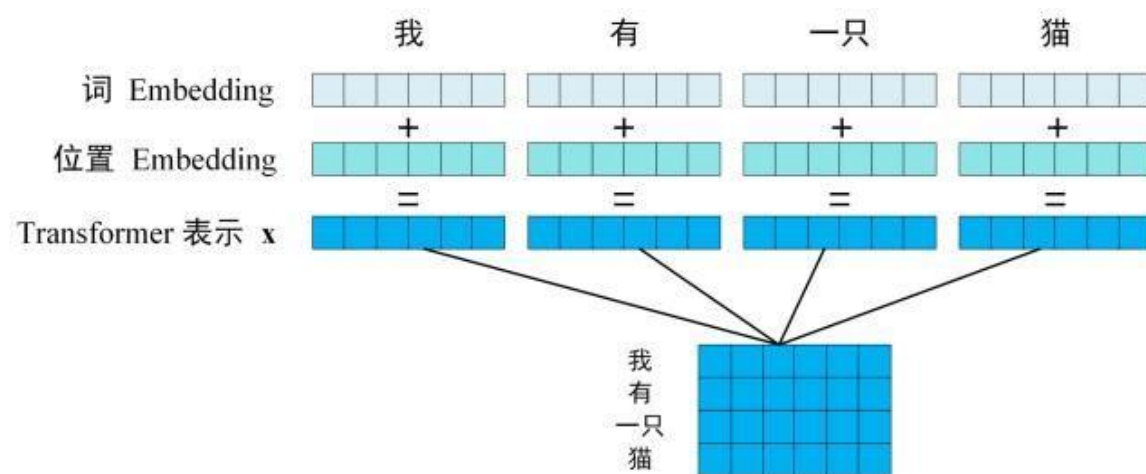
- Allowing for **parallelization**
- Allowing modeling of dependencies without regard to distance

Model Architecture



- Embedding input
- **Encoder**
- **Decoder**
- Output

Embedding Input



- Word Embedding: word2vec, GloVe, one-hot, etc.
- Positional Embedding:

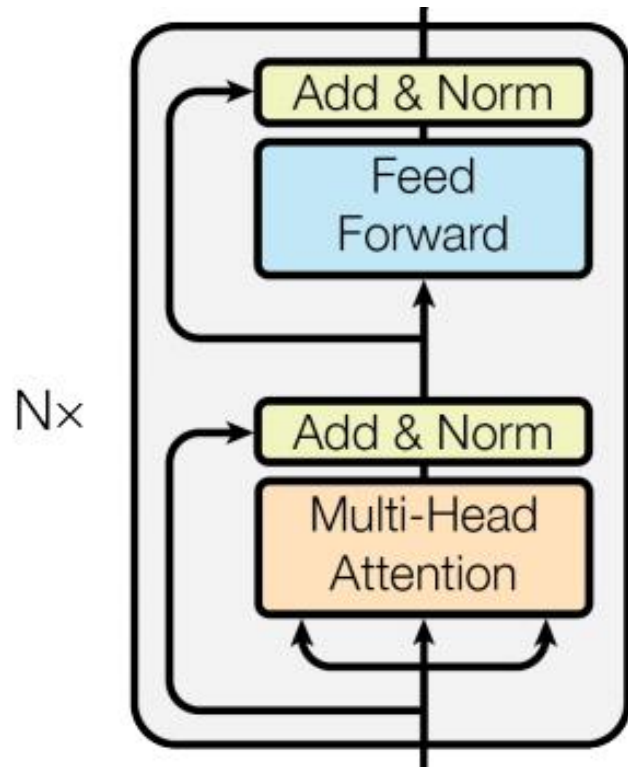
$$PE(pos, 2i) = \sin\left(\frac{pos}{10000^{\frac{2i}{d_{model}}}}\right)$$

$$PE(pos, 2i + 1) = \cos\left(\frac{pos}{10000^{\frac{2i}{d_{model}}}}\right)$$

- d_{model} : dimension of input
- pos: position of words

Encoder

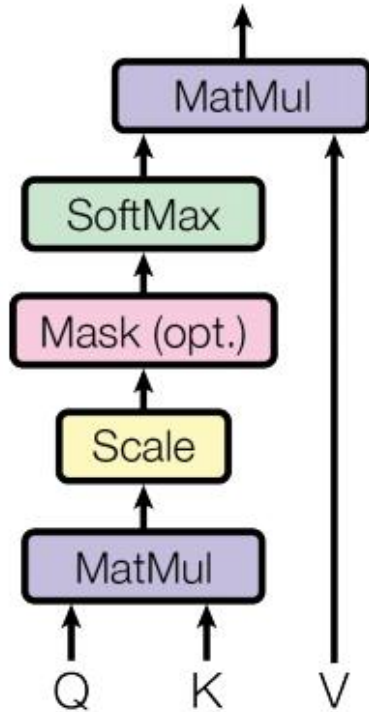
Encoder has $N=6$ identical layers



Structure of one layer:

- **Multi-Head Attention**
- Feed Forward
- Add & Norm

Self-Attention Mechanism



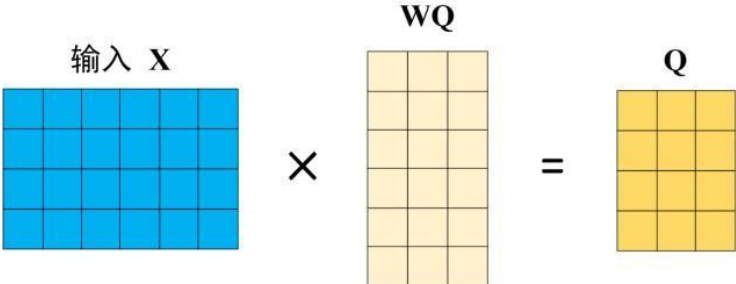
Self-Attention(without mask):

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

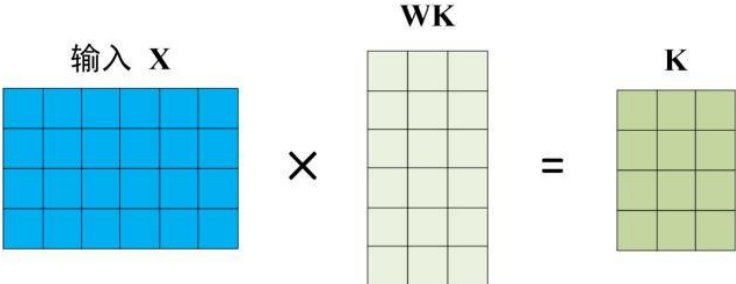
- d_k : Number of columns (dimension) of Q and K

Self-Attention Input

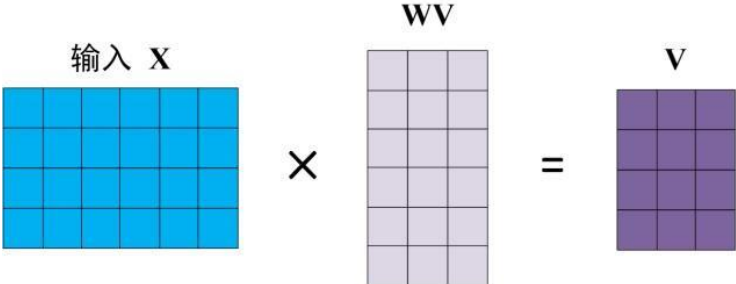
输入 X


$$\begin{matrix} \text{输入 X} \\ \times \\ \text{WQ} \end{matrix} = \text{Q}$$

输入 X

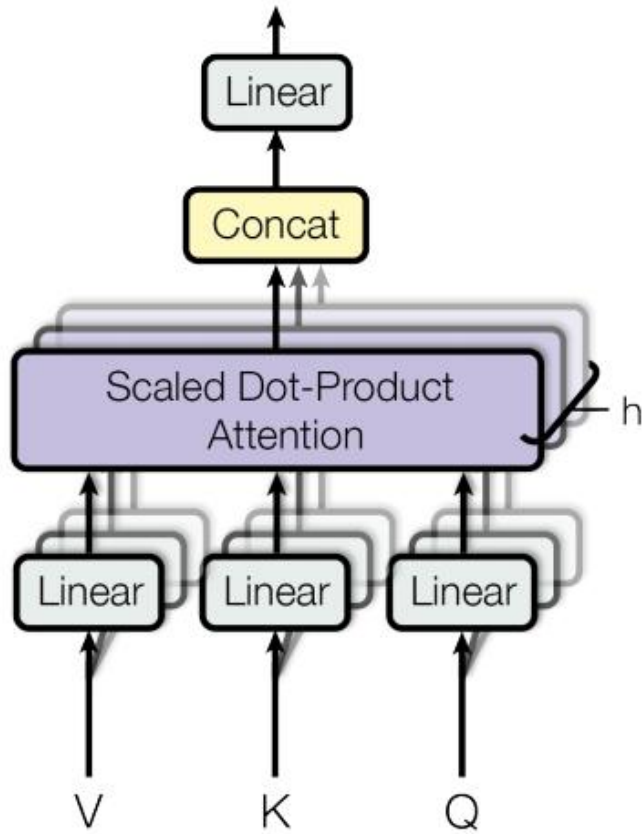

$$\begin{matrix} \text{输入 X} \\ \times \\ \text{WK} \end{matrix} = \text{K}$$

输入 X


$$\begin{matrix} \text{输入 X} \\ \times \\ \text{WV} \end{matrix} = \text{V}$$

- X: Input of each layer of encoder
- W^Q, W^K, W^V : Linear transformation matrix

Multi-Head Attention

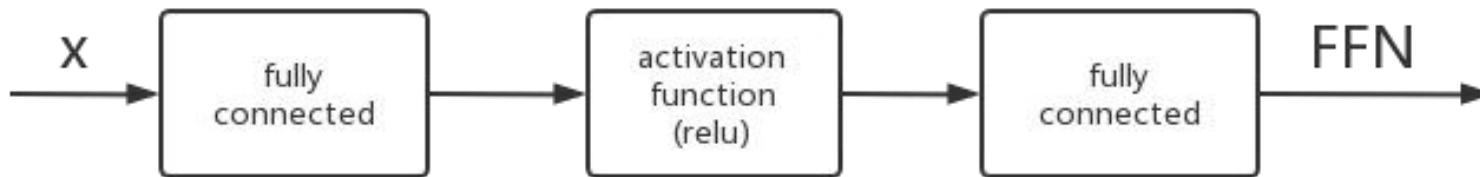


$$\text{MultiHead}(Q, K, V) = \text{Concat}(\text{head}_1, \dots, \text{head}_h)W^O$$

- where $\text{head}_i = \text{Attention}(QW_i^Q, KW_i^K, VW_i^V)$

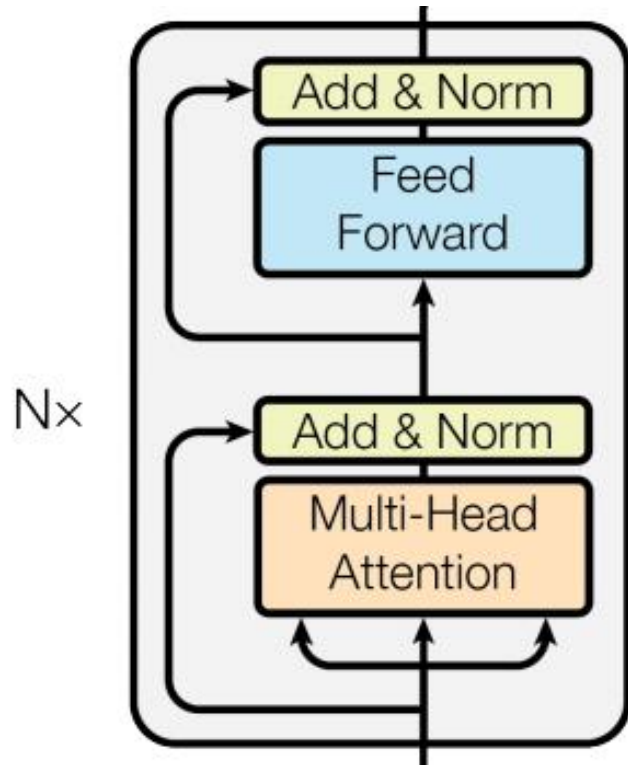
Multi-Head attention ensures the **parallelism**

Feed Forward



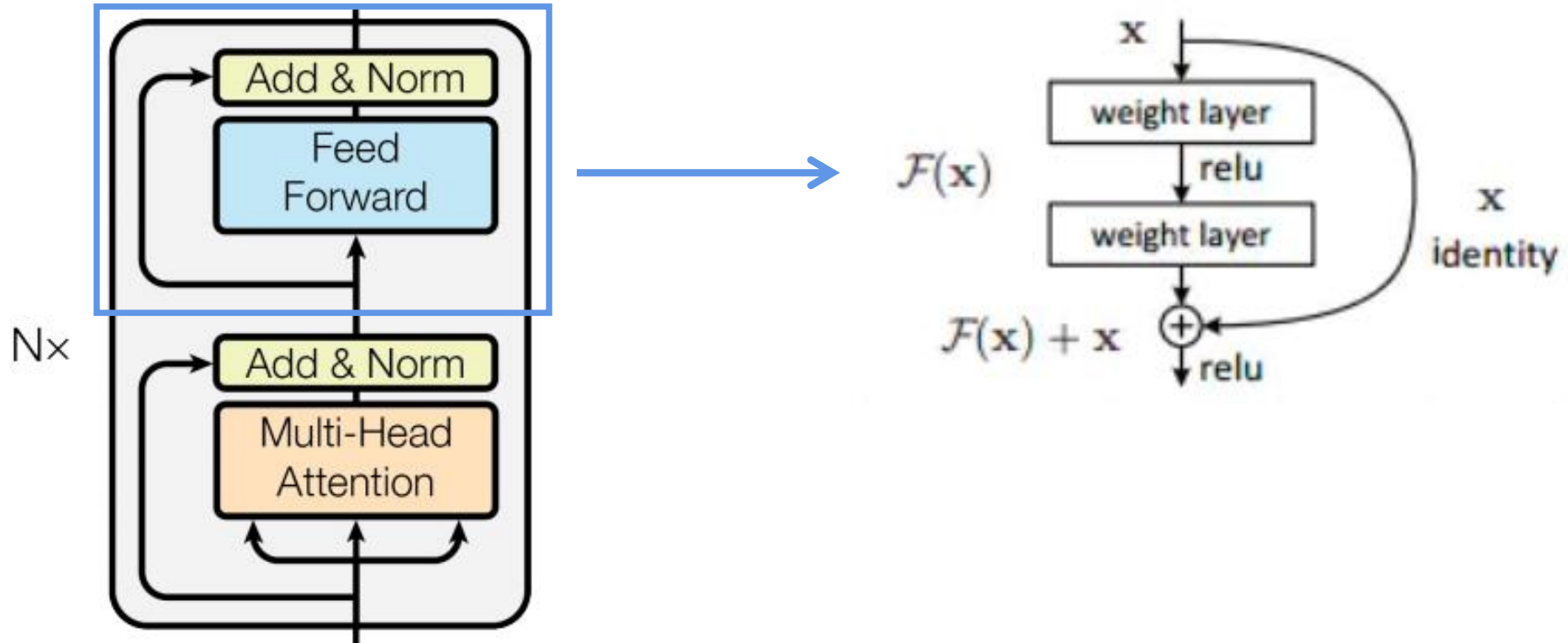
$$FFN(x) = \max(0, xW_1 + b_1)W_2 + b_2$$

Add & Norm

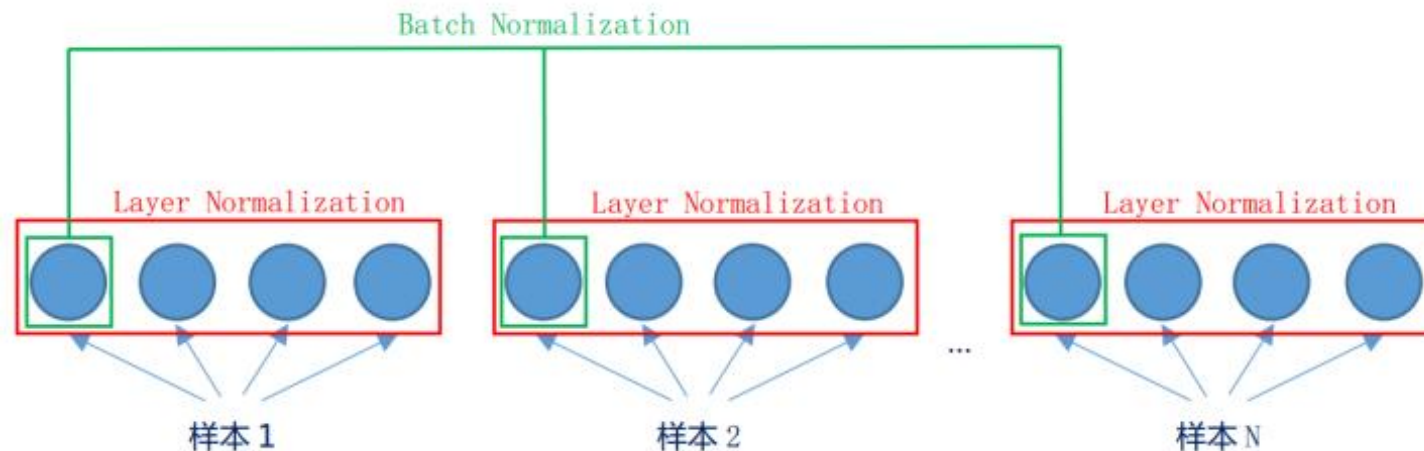


- Add: Residual connection
- Norm: Layer normalization

Residual Connection



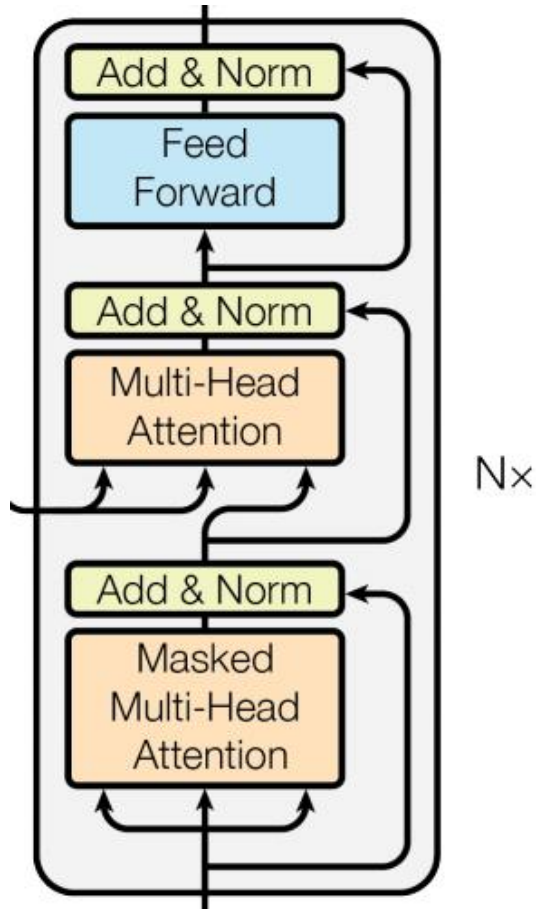
Layer Normalization



- **Layer norm:** Normalization **between different dimensions** in the same sample
- Batch norm: Normalization between different samples in the same dimension

Decoder

Decoder has $N=6$ identical layers



Structure of one layer:

- Masked Multi-Head Attention
- Multi-Head Attention
- Feed Forward
- Add & Norm

Mask Operation

0						
1						
2						
3						
4						

输入矩阵 X



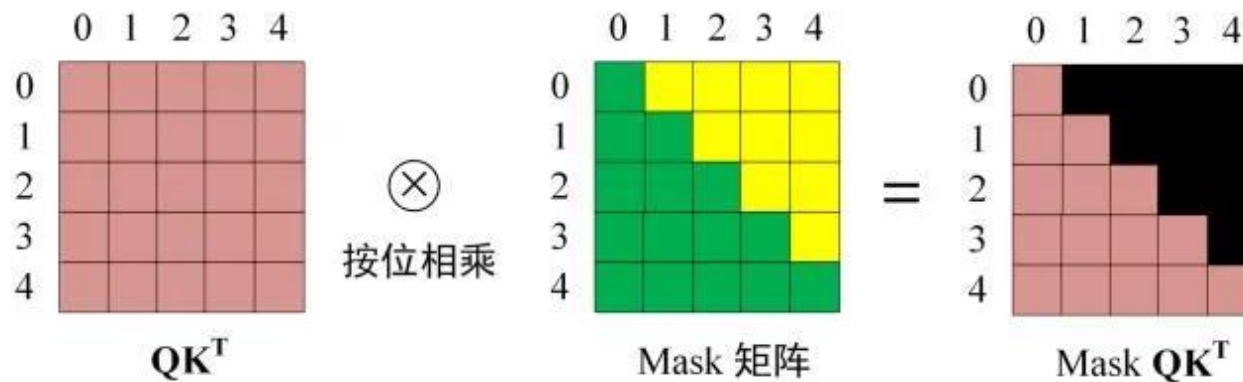
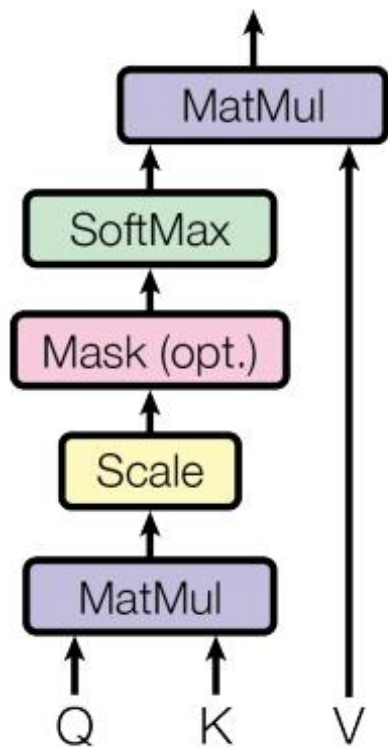
	0	1	2	3	4
0	不遮挡	遮挡	遮挡	遮挡	遮挡
1	不遮挡	不遮挡	遮挡	遮挡	遮挡
2	不遮挡	不遮挡	不遮挡	遮挡	遮挡
3	不遮挡	不遮挡	不遮挡	不遮挡	遮挡
4	不遮挡	不遮挡	不遮挡	不遮挡	不遮挡

Mask 矩阵

Sequence Mask:

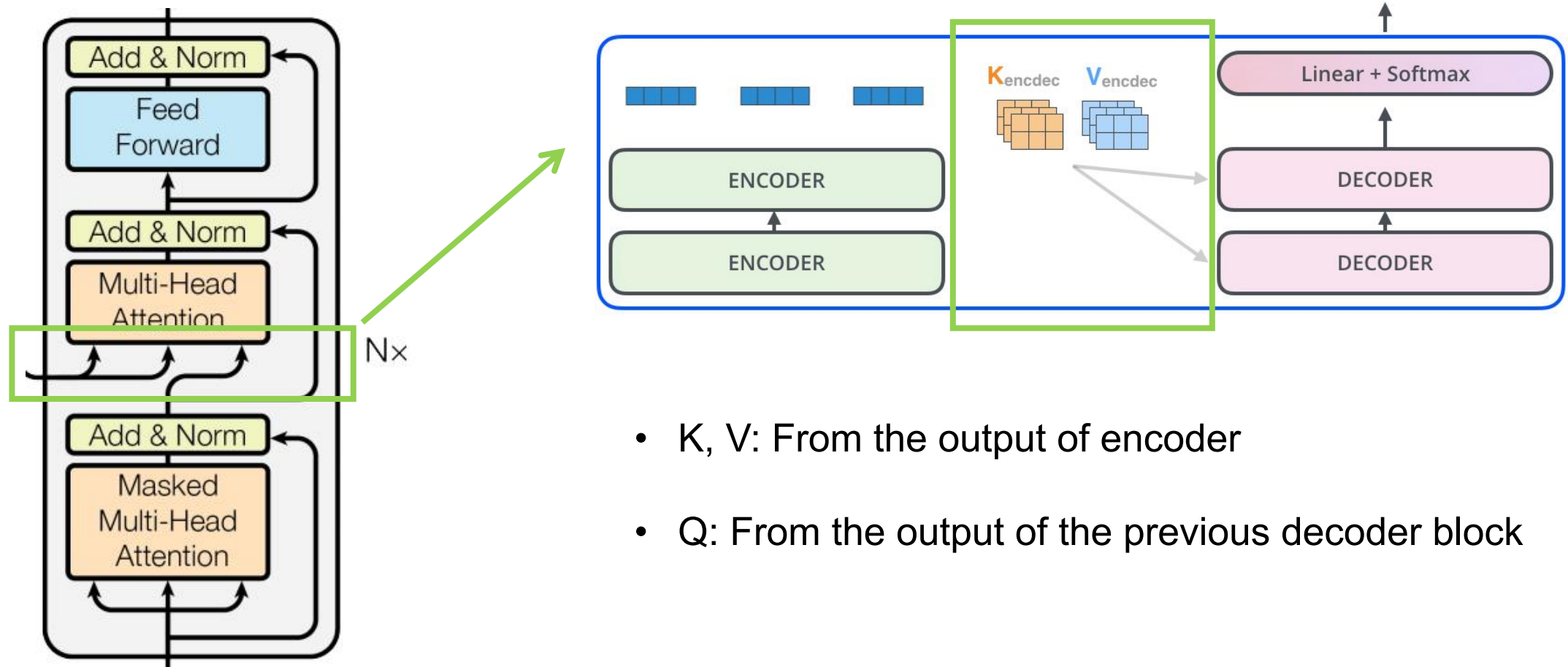
- Why: The decoder cannot see the future information.
- How: Setting masked position to $-\infty$.

Masked Multi-Head Attention

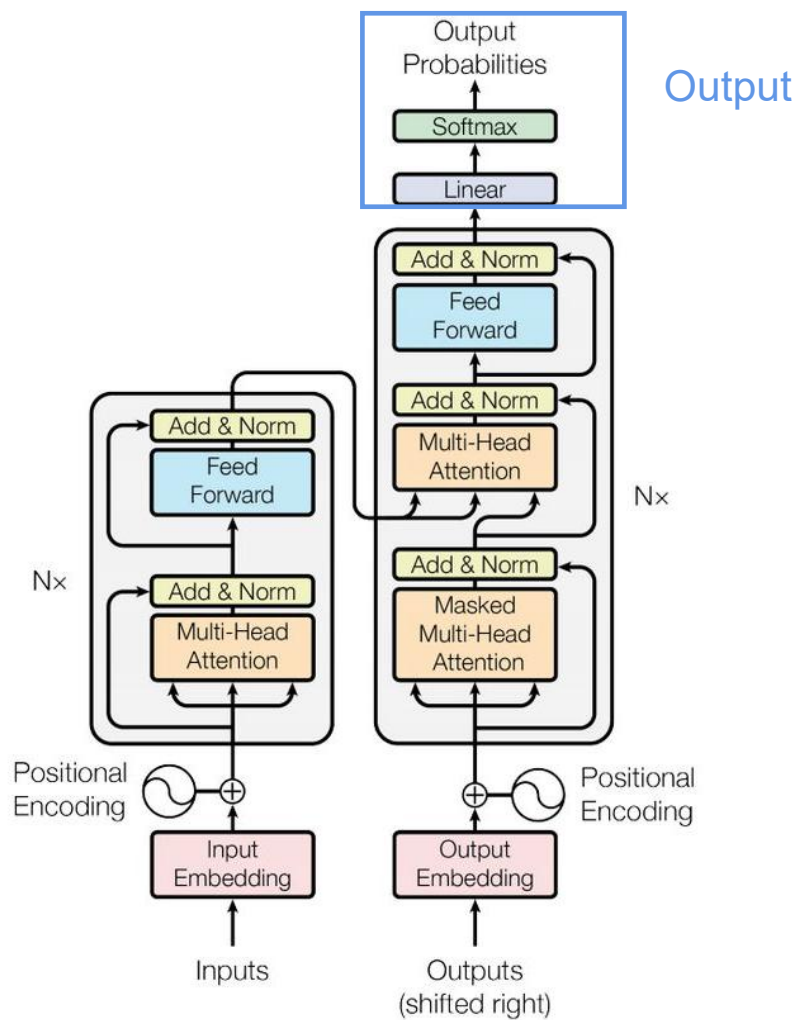


- Mask matrix is a lower triangular matrix (non-0 elements are all 1).

Multi-Head Attention In Decoder



Output



Structure of output block:

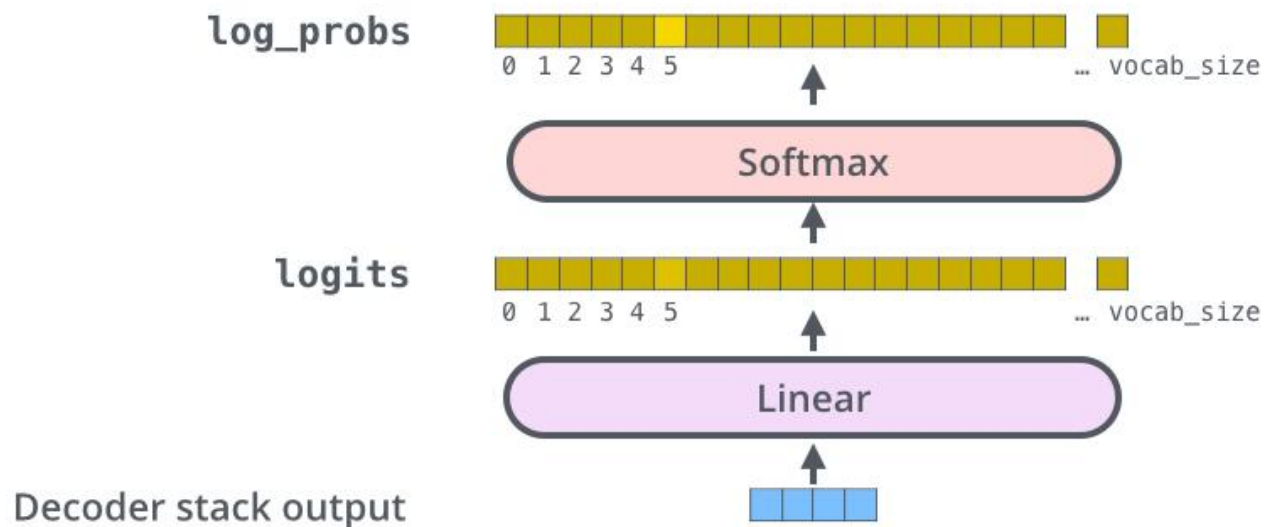
- Linear layer
- Softmax

Output

- How to predict a word:

Which word in our vocabulary
is associated with this index?

Get the index of the cell
with the highest value
(**argmax**)



Conclusion

Advantages:

- Avoid sequence model and can realize **parallelization**.
- Solves the problem of **long distance dependence**.
- Self-attention can produce more **interpretable models**.

Disadvantages:

- Loss of **position information** (Compare with RNN)
- Lost the ability to capture **local features** (Compare with CNN)

Transformer is widely used in Deep Learning fields (NLP, CV, etc.).

*Thank
you*

