



澳門大學
UNIVERSIDADE DE MACAU
UNIVERSITY OF MACAU



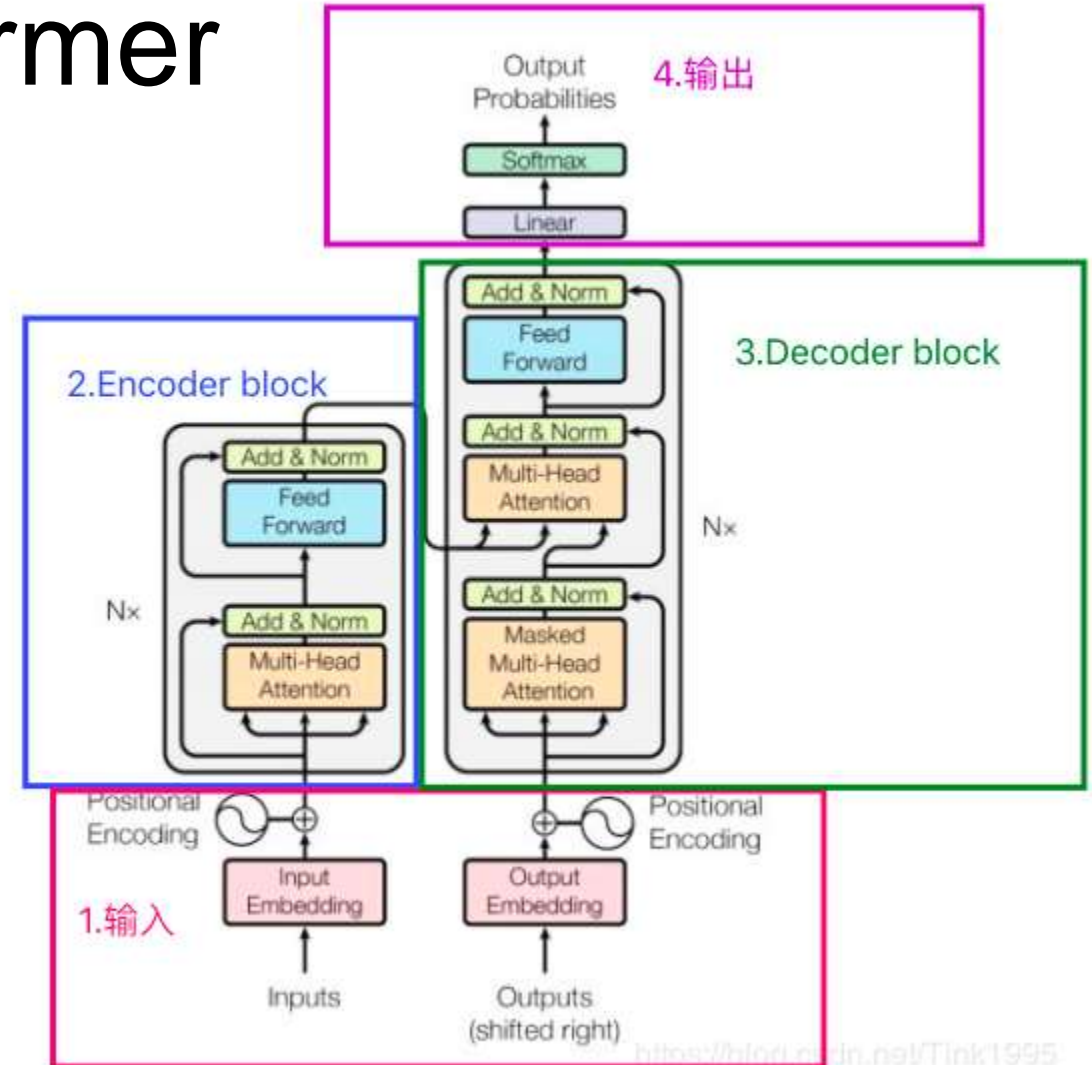
Transformer

Yu Jiening
Yu Jiening@umac.mo

um 澳大

Architecture of transformer

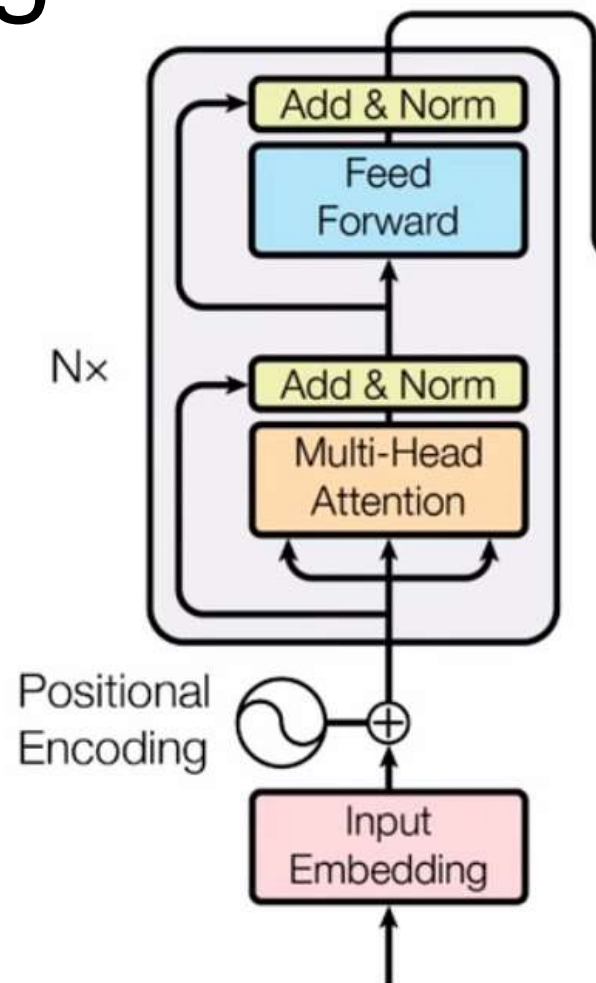
- Consists of four parts: input, output, encoder, decoder.



Positional Encoding

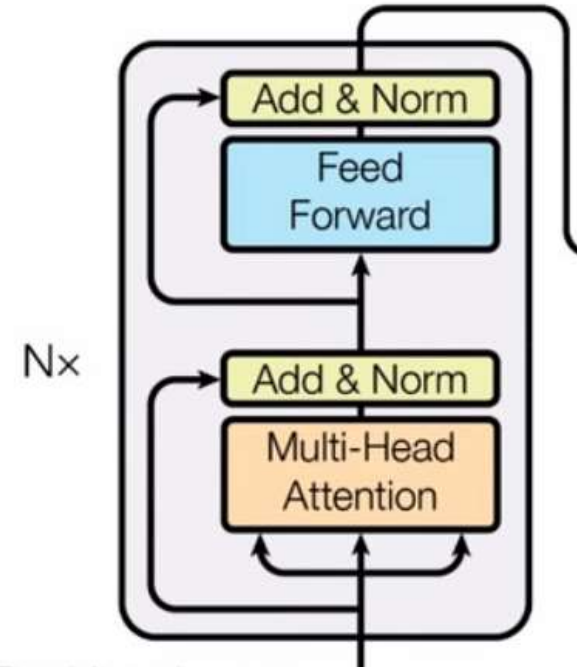
$$PE(pos, 2i) = \sin\left(\frac{pos}{10000^{\frac{2i}{d_{model}}}}\right)$$

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



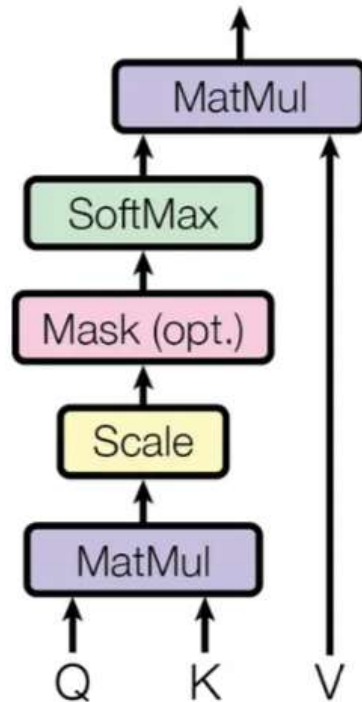
Encoder

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

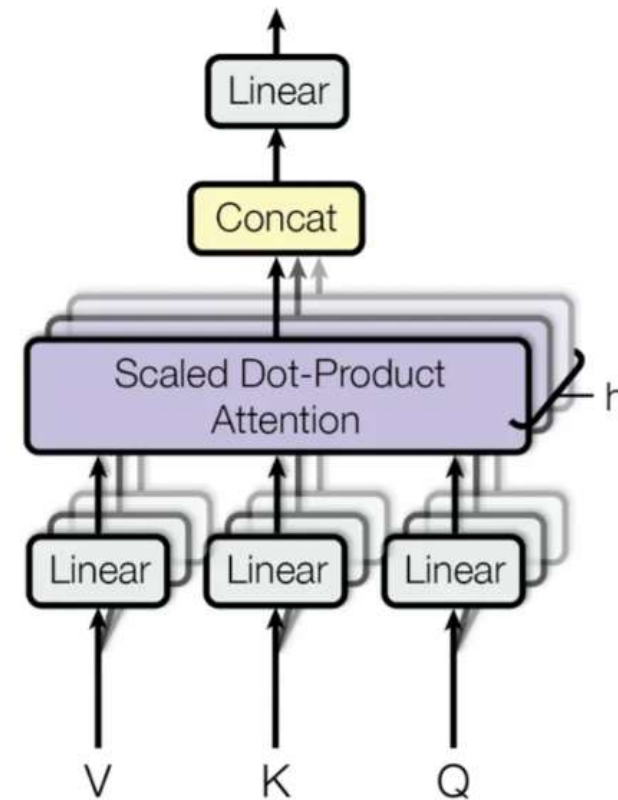


Multi-Head Attention

Scaled Dot-Product Attention



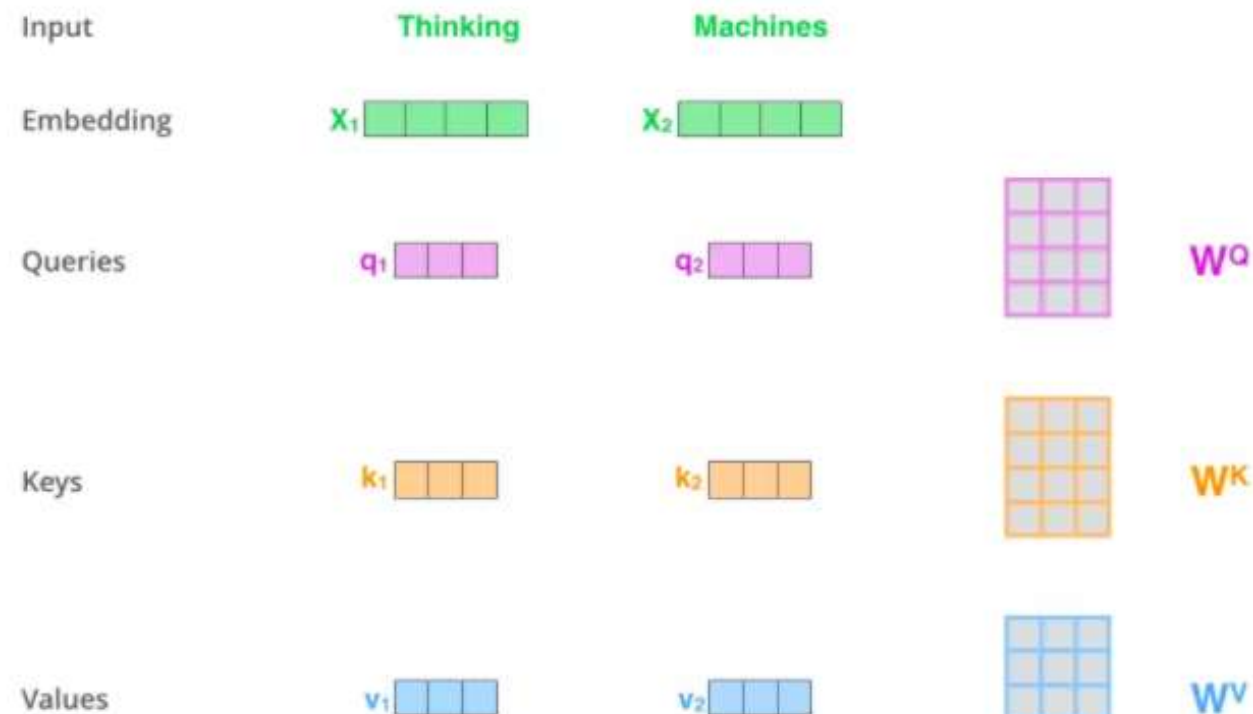
Multi-Head Attention



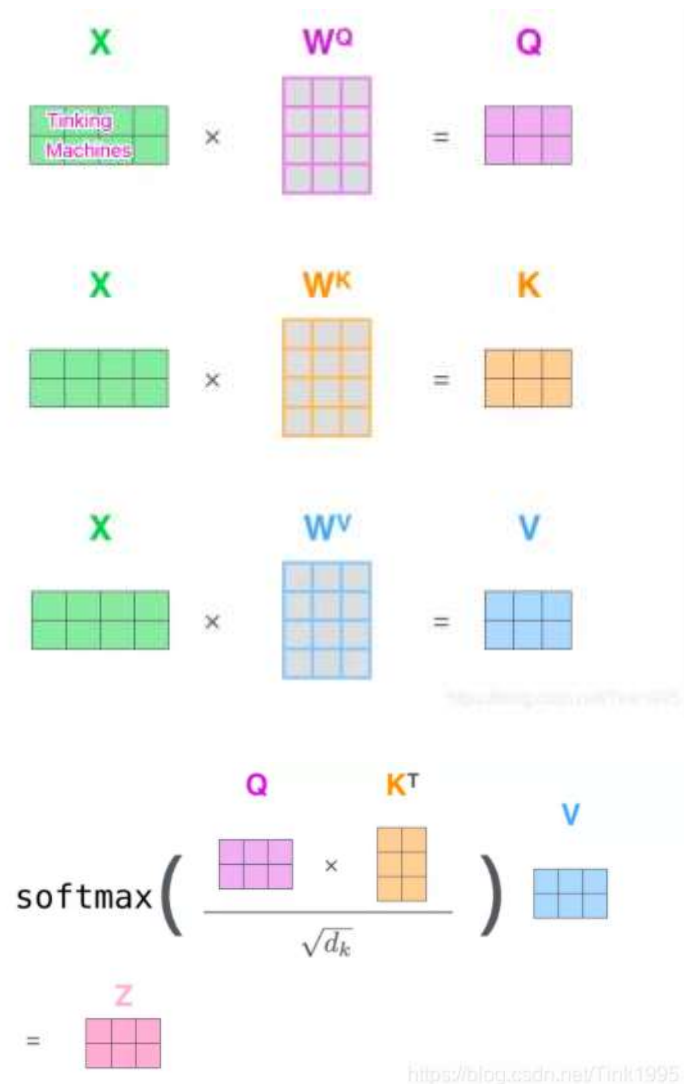
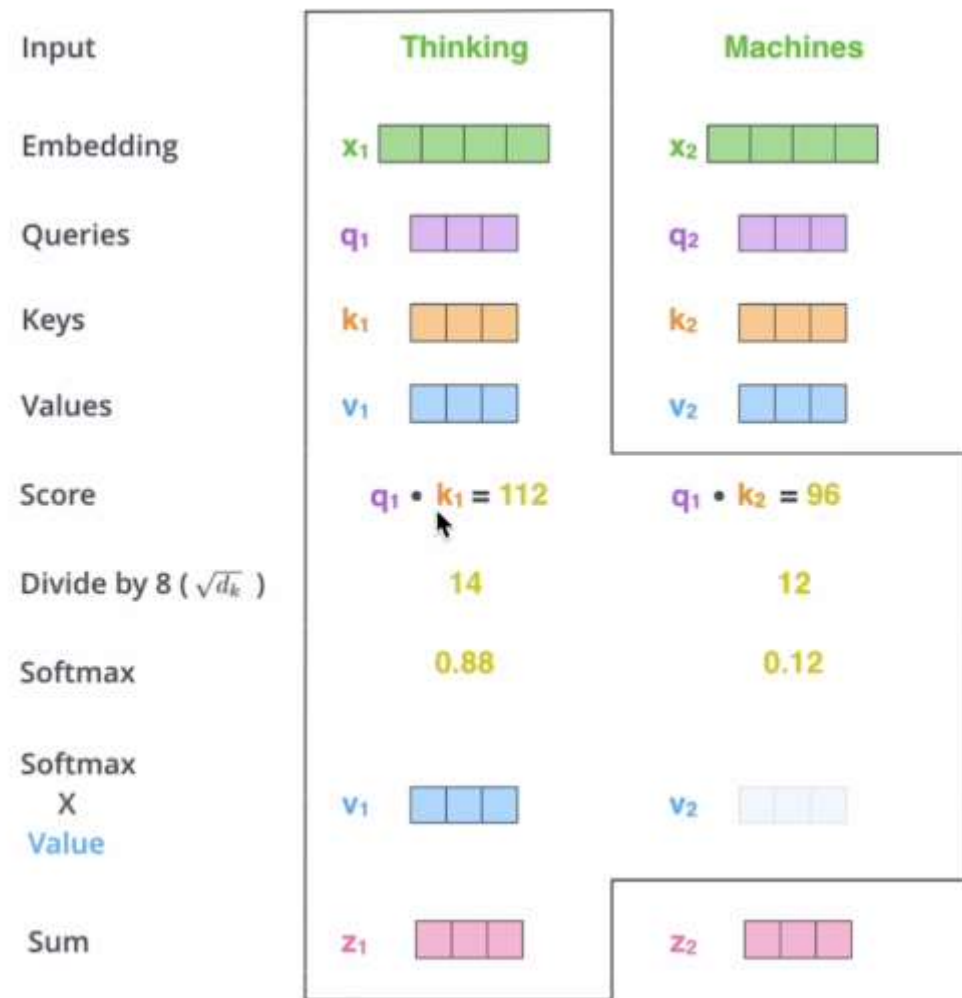
$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

$$\begin{aligned} & \text{MultiHead}(Q, K, V) \\ &= \text{Concat}(\text{head1}, \dots, \text{headh})W_0 \end{aligned}$$

where $\text{head}_i = \text{Attention}(Q_i, K_i, V_i)$

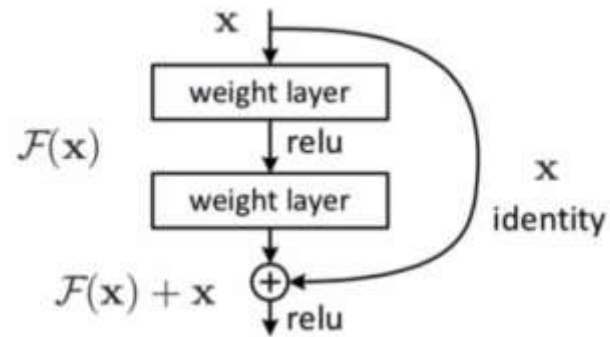


<https://blog.csdn.net/Tink1995>

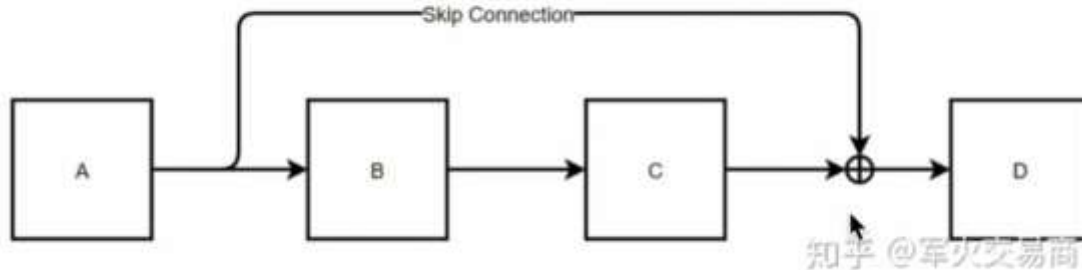


$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

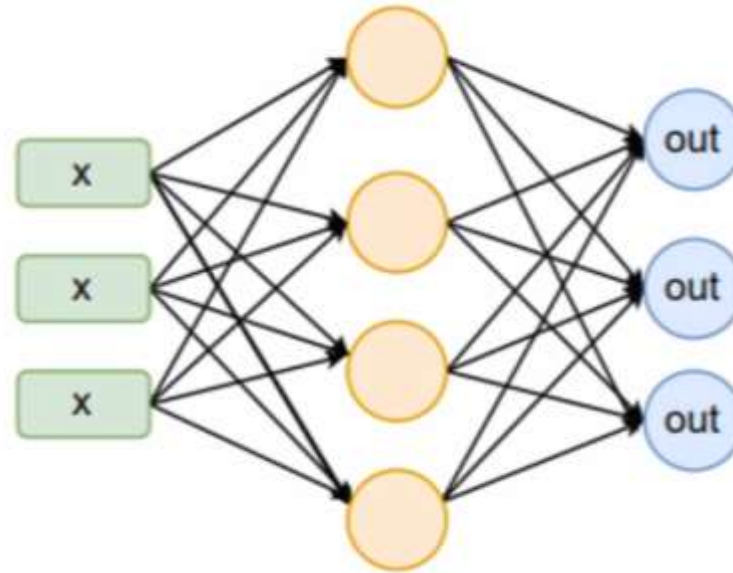
Add & Norm



$$D_{in} = A_{out} + C(B(A_{out}))$$



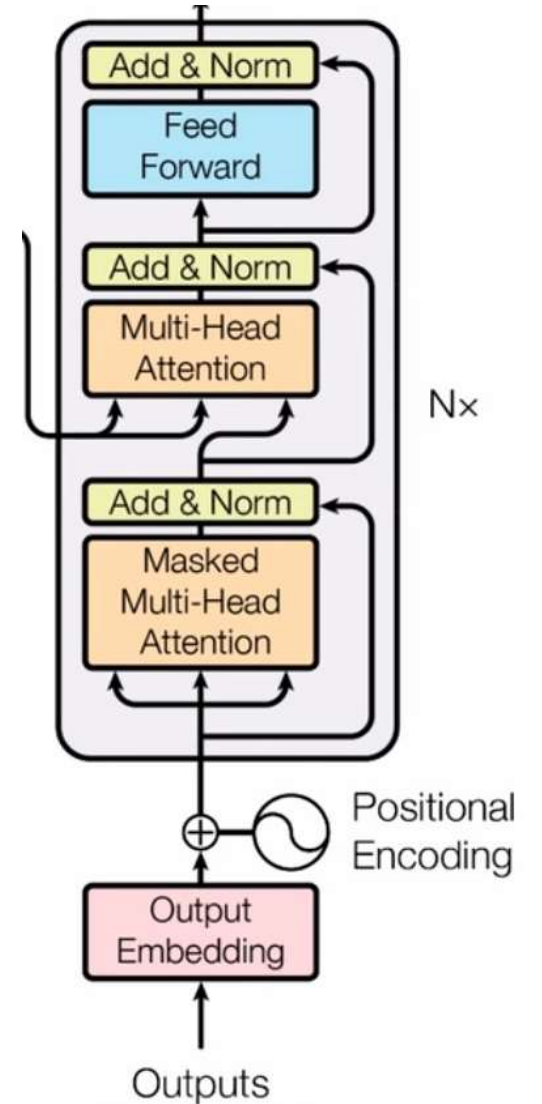
Feed Forward



MLP structure

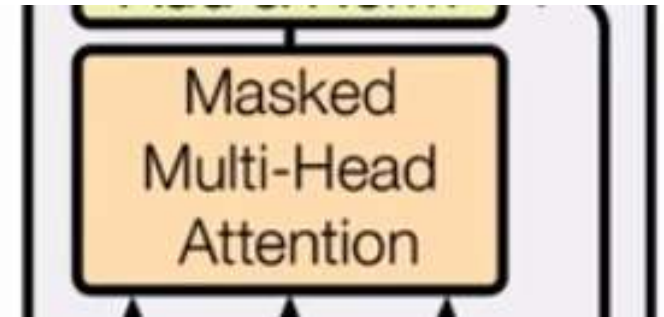
Decoder

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



Masked Multi-Head Attention

- The point of masking is to mask out future information, making the trained model more accurate



Interactive Attention

Calculation of Q, K, V for self-attention:

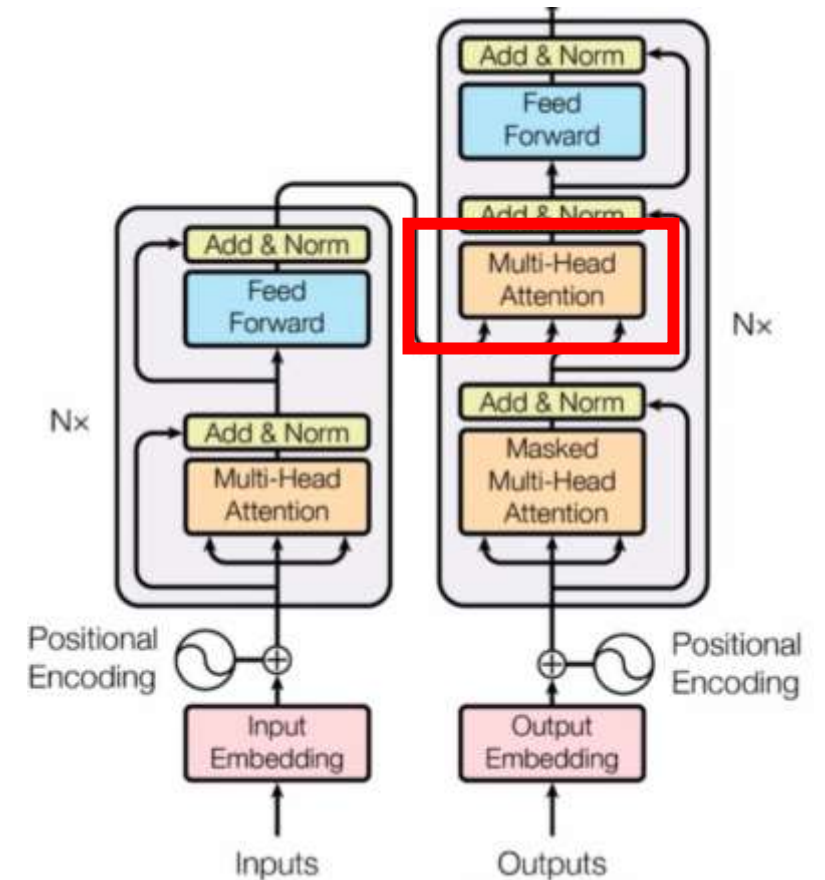
$$Q = W_q X + b_q$$

$$K = W_k X + b_k$$

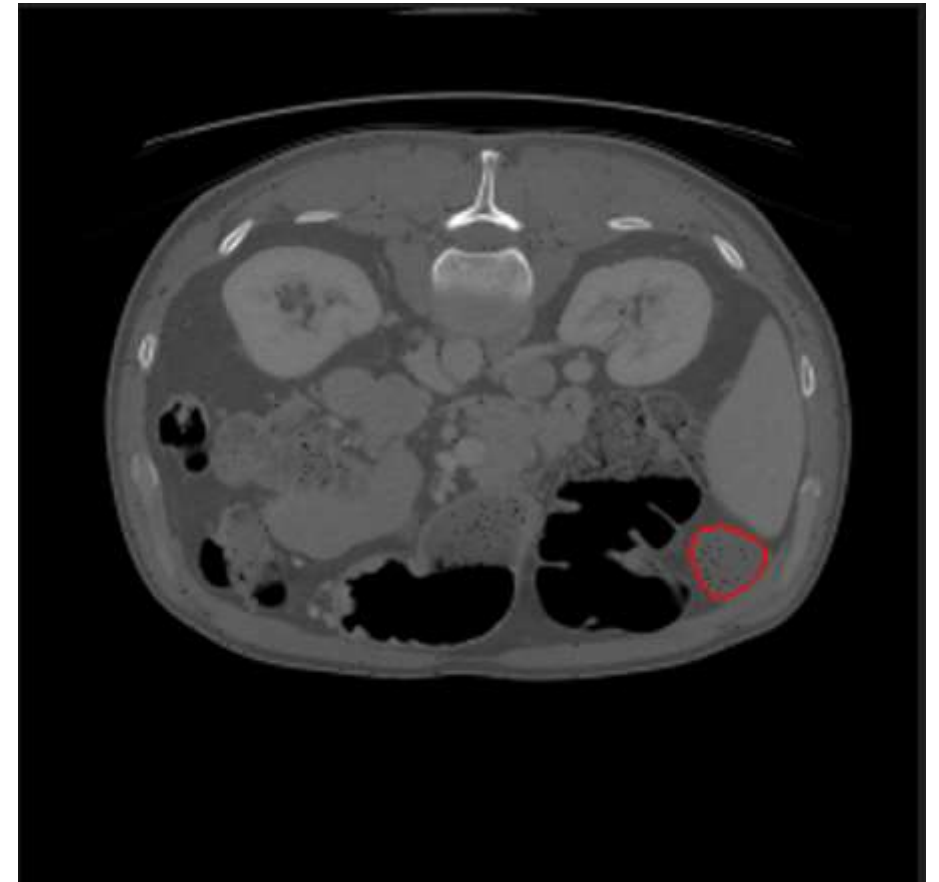
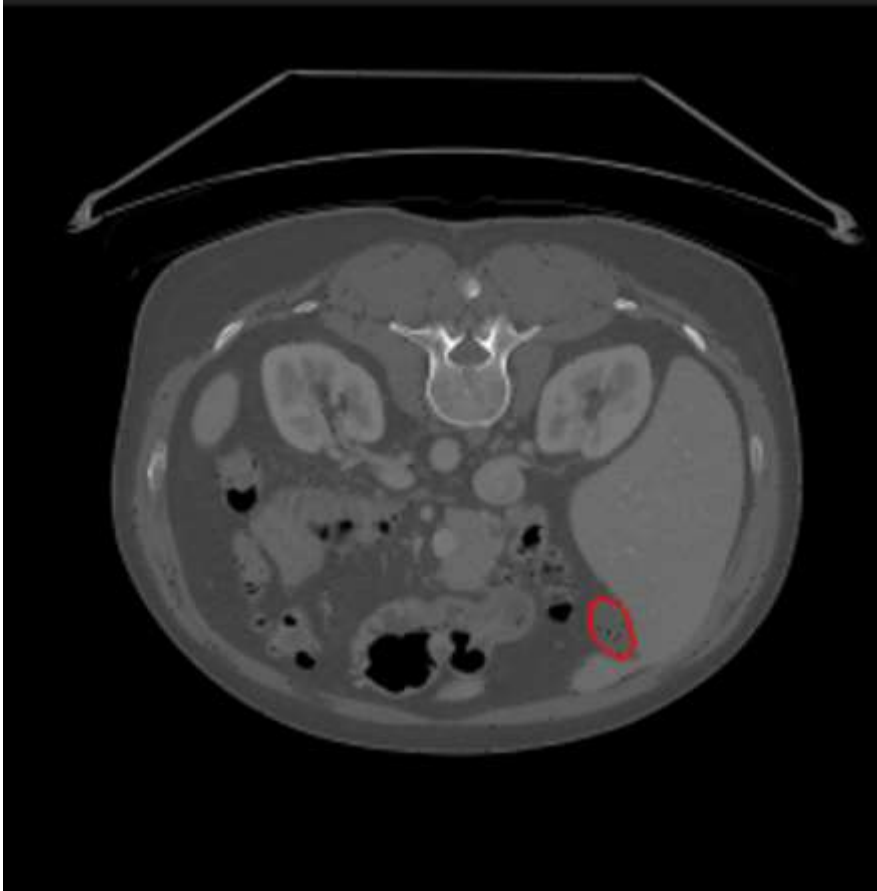
$$V = W_v X + b_v$$

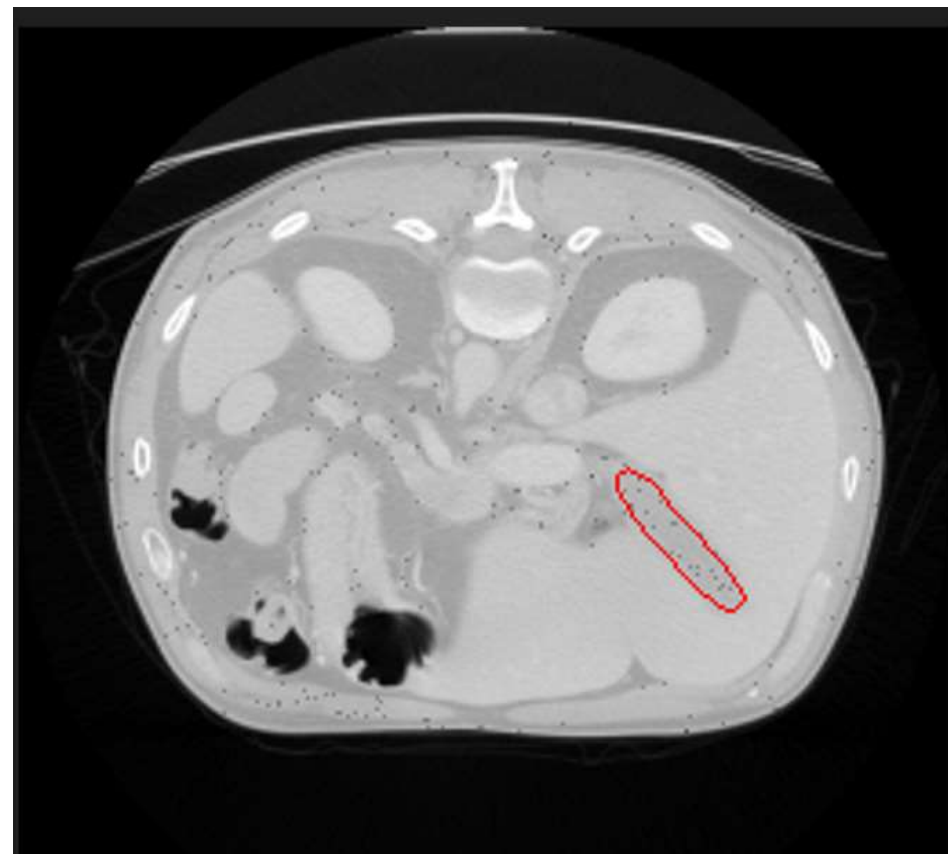
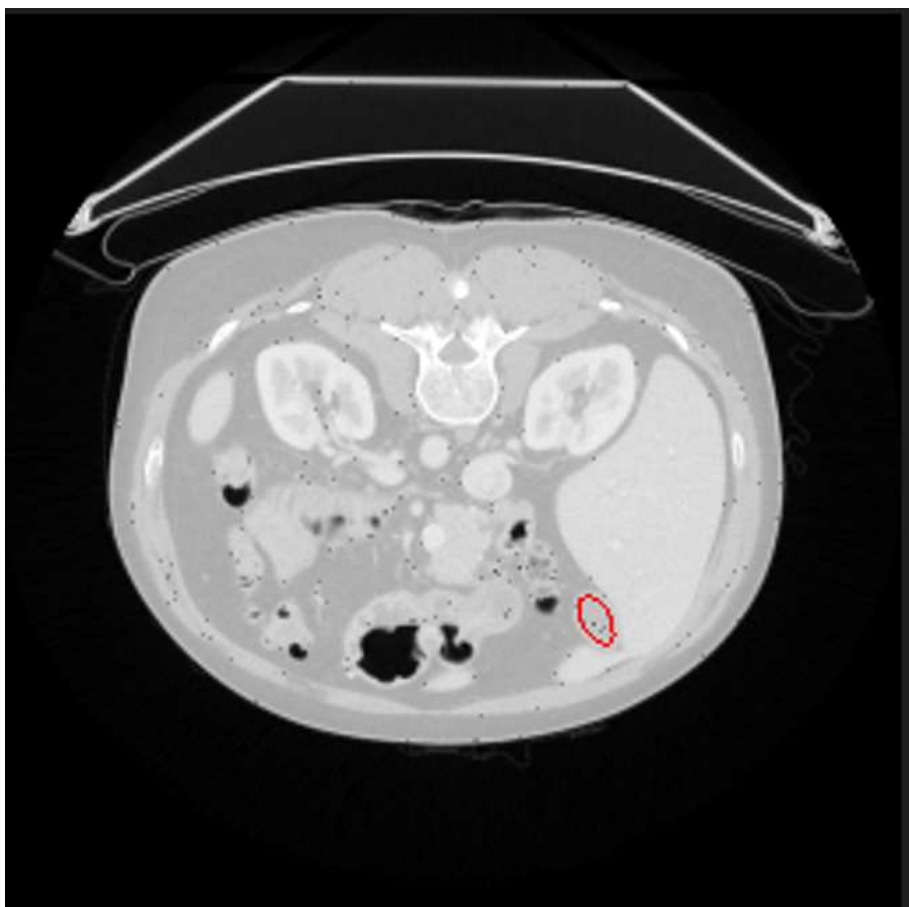
Calculation of Q for interactive attention

$$Q = W_q Out_{encoder} + b_q$$



CT:





Thank You!