



## Tokens-to-Token ViT: Training Vision Transformers from Scratch on ImageNet



## 预备知识

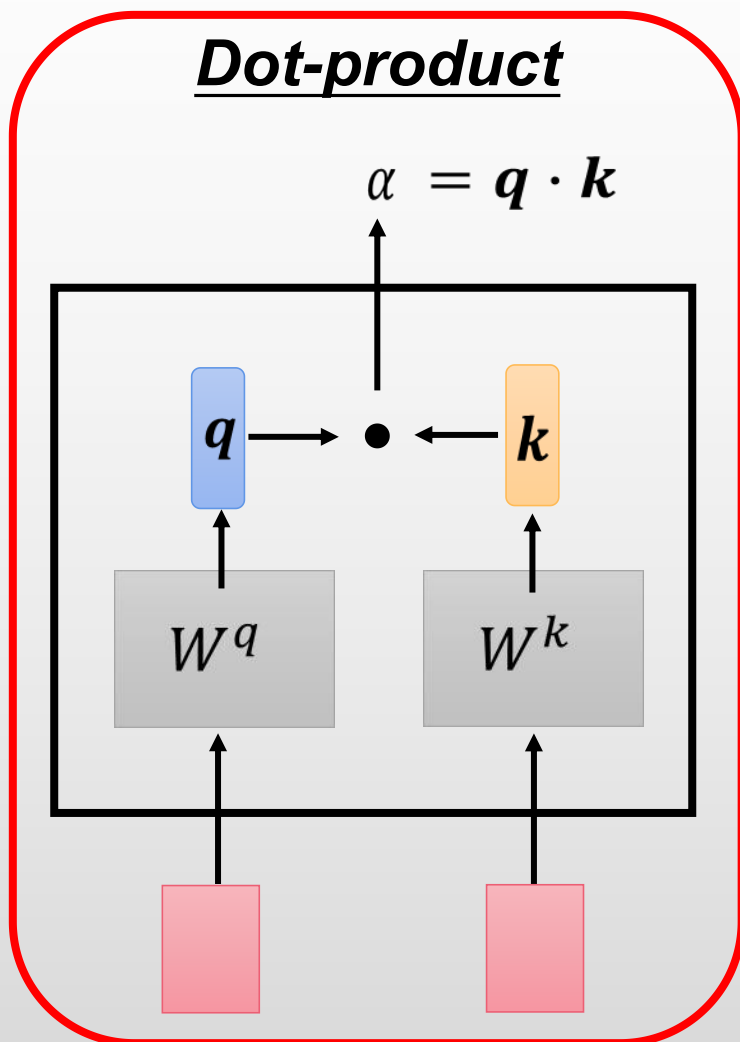
1.self-attention

2.transformer

3ViT



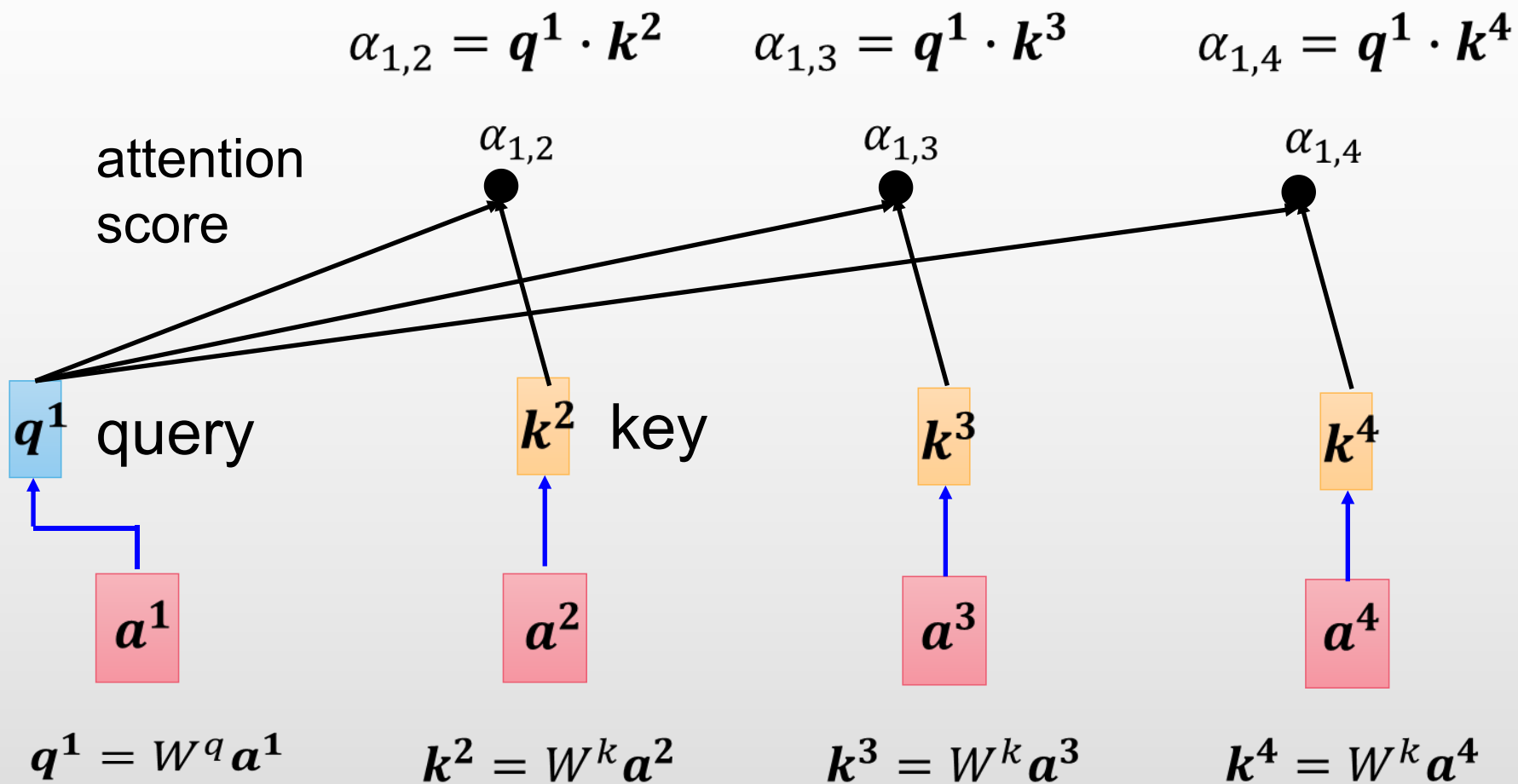
## SELF-ATTENTION



- 2005年，Bahdanau等人在论文《Neural Machine Translation by Jointly Learning to Align and Translate》
- Google 机器翻译团队在NIPS 2017上发表的《Attention is All You Need》
- GoogleMind 2014年发表《Recurrent Models of Visual Attention》



## SELF-ATTENTION

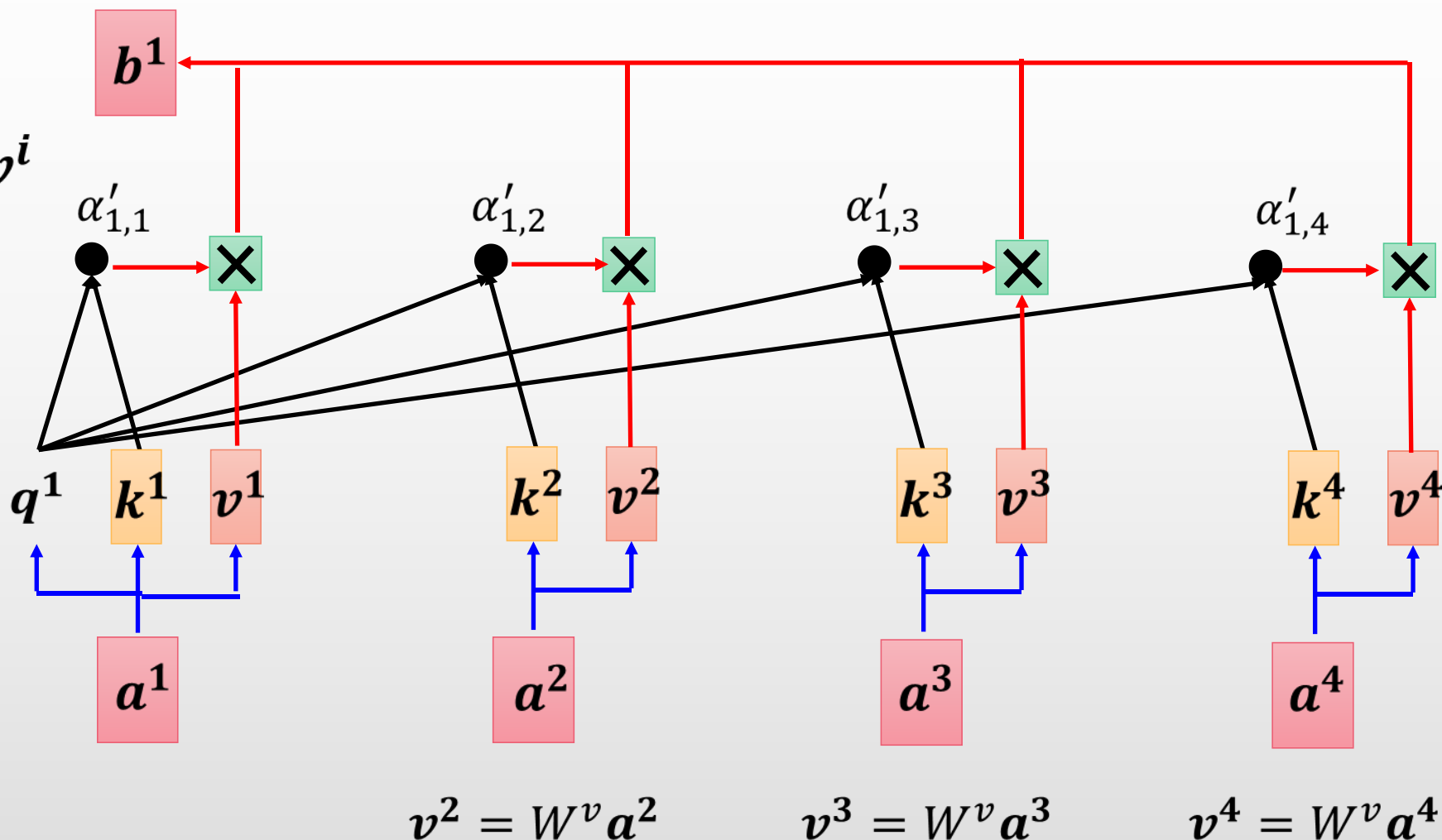




# Tokens-to-Token ViT: Training Vision Transformers from Scratch on ImageNet

李玉光

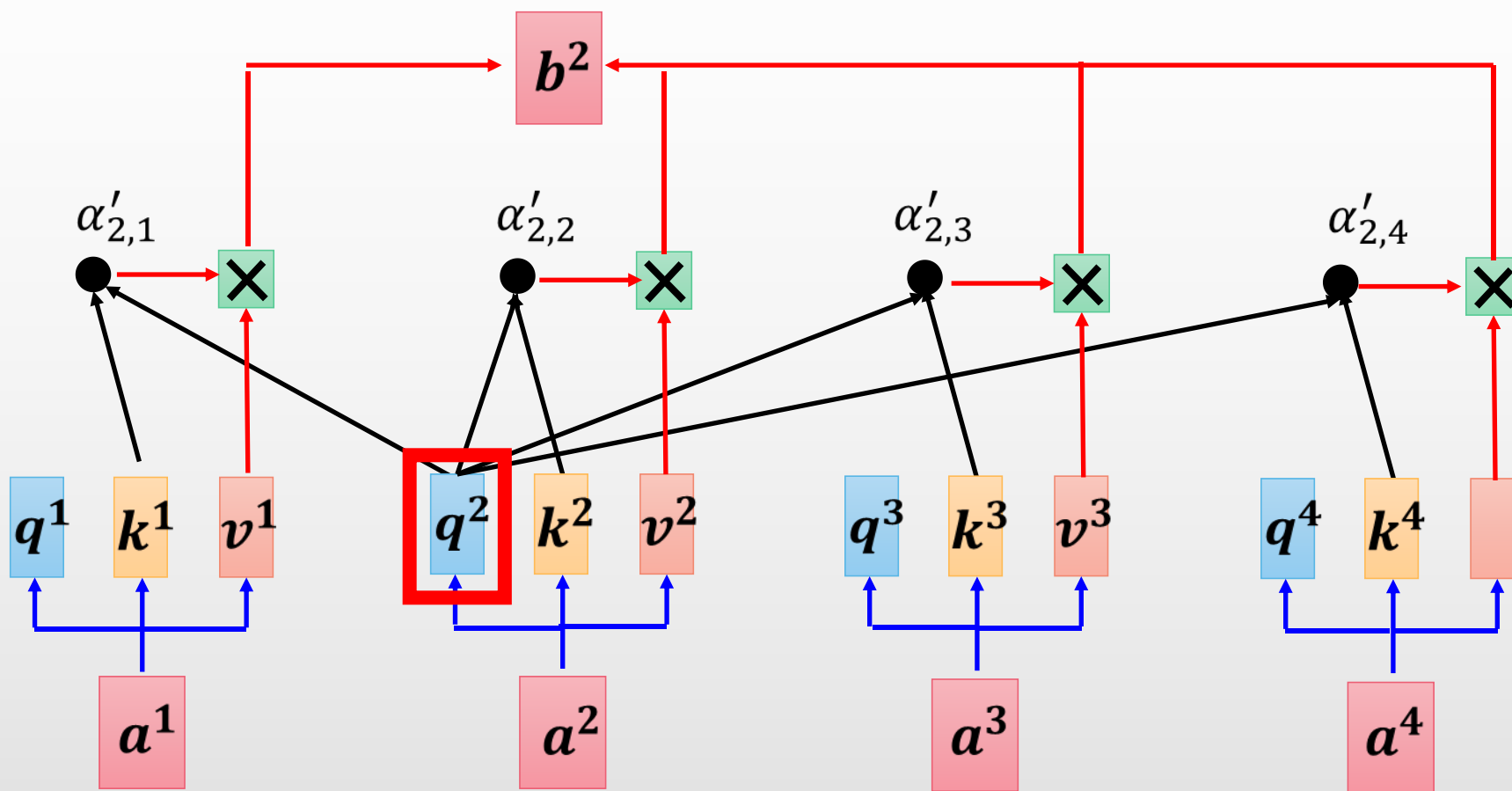
$$b^1 = \sum_i \alpha'_{1,i} v^i$$





# Tokens-to-Token ViT: Training Vision Transformers from Scratch on ImageNet

李玉光





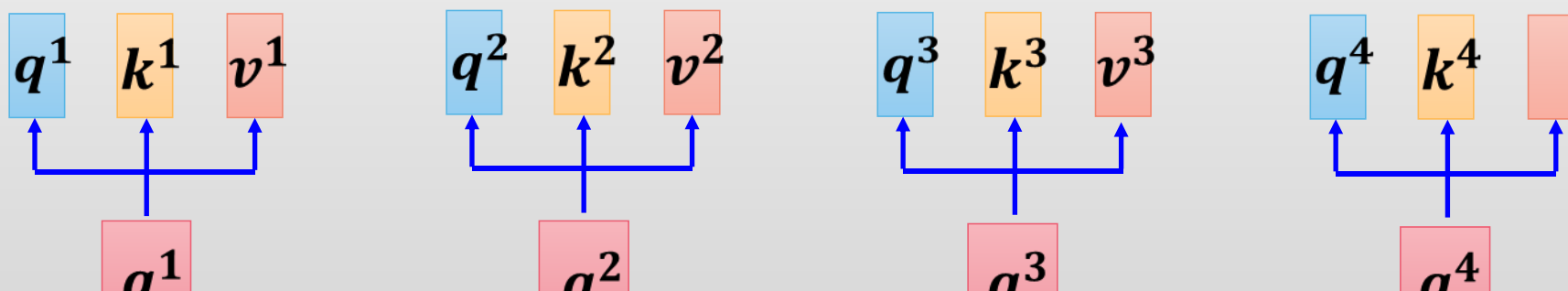
# Tokens-to-Token ViT: Training Vision Transformers from Scratch on ImageNet

李玉光

$$q^i = W^q a^i \quad \begin{array}{c} \boxed{q^1} \boxed{q^2} \boxed{q^3} \boxed{q^4} \\ Q \end{array} = \begin{array}{c} \boxed{W^q} \quad \boxed{a^1} \boxed{a^2} \boxed{a^3} \boxed{a^4} \\ I \end{array}$$

$$k^i = W^k a^i \quad \begin{array}{c} \boxed{k^1} \boxed{k^2} \boxed{k^3} \boxed{k^4} \\ K \end{array} = \begin{array}{c} \boxed{W^k} \quad \boxed{a^1} \boxed{a^2} \boxed{a^3} \boxed{a^4} \\ I \end{array}$$

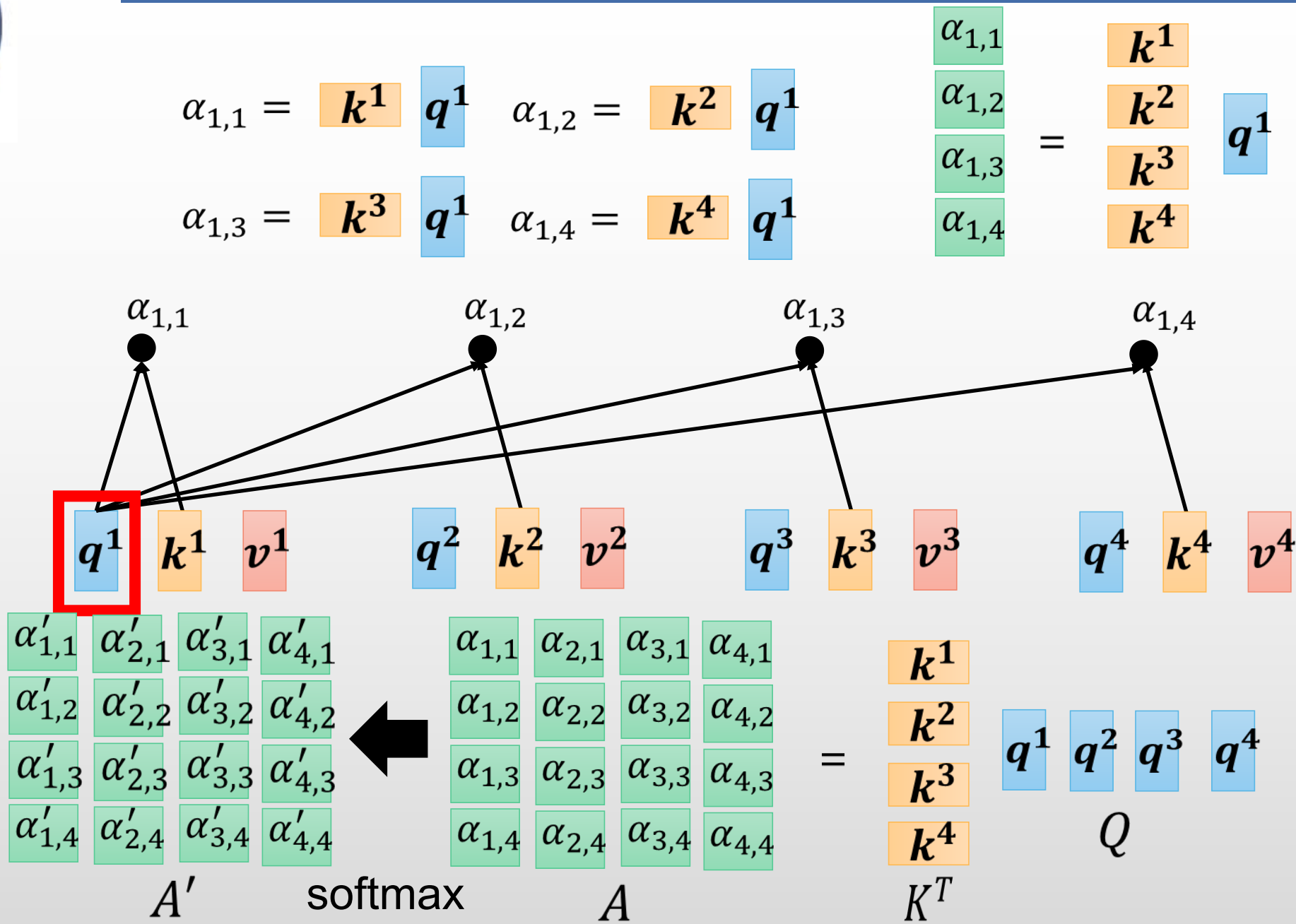
$$v^i = W^v a^i \quad \begin{array}{c} \boxed{v^1} \boxed{v^2} \boxed{v^3} \boxed{v^4} \\ V \end{array} = \begin{array}{c} \boxed{W^v} \quad \boxed{a^1} \boxed{a^2} \boxed{a^3} \boxed{a^4} \\ I \end{array}$$





# Tokens-to-Token ViT: Training Vision Transformers from Scratch on ImageNet

李玉光

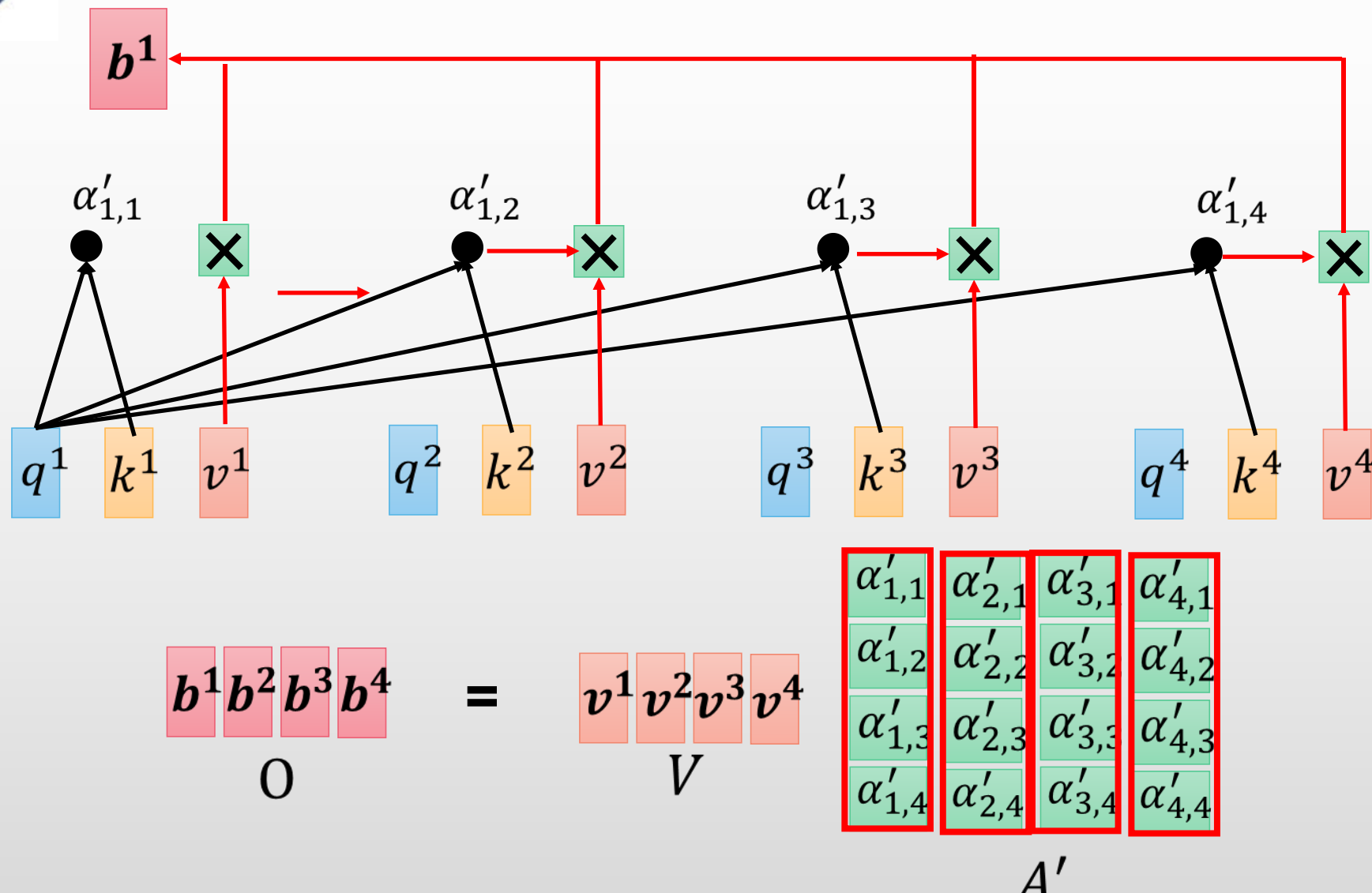






# Tokens-to-Token ViT: Training Vision Transformers from Scratch on ImageNet

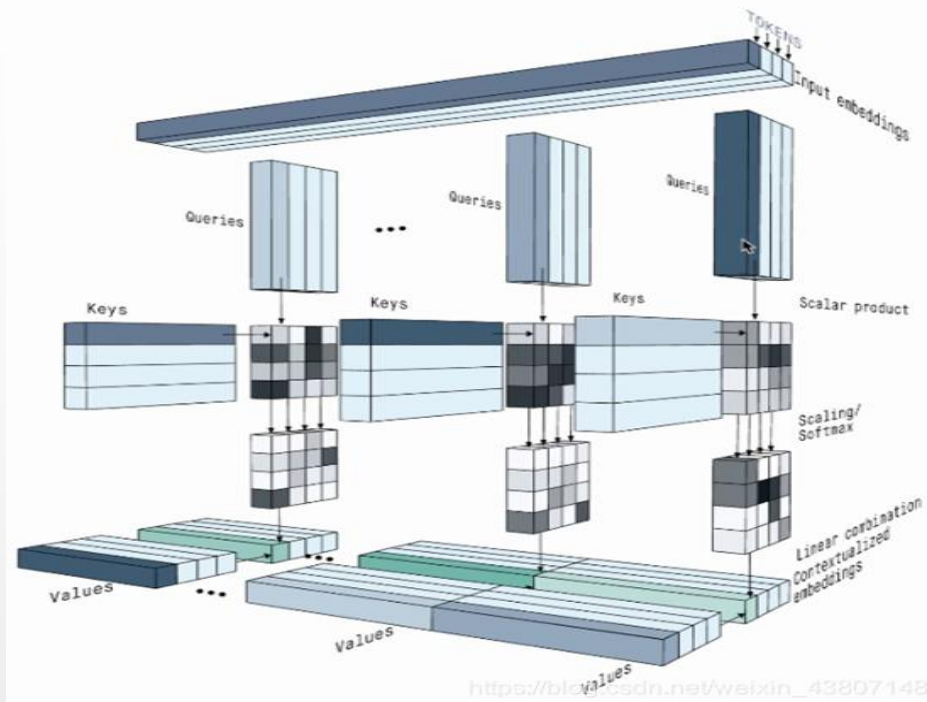
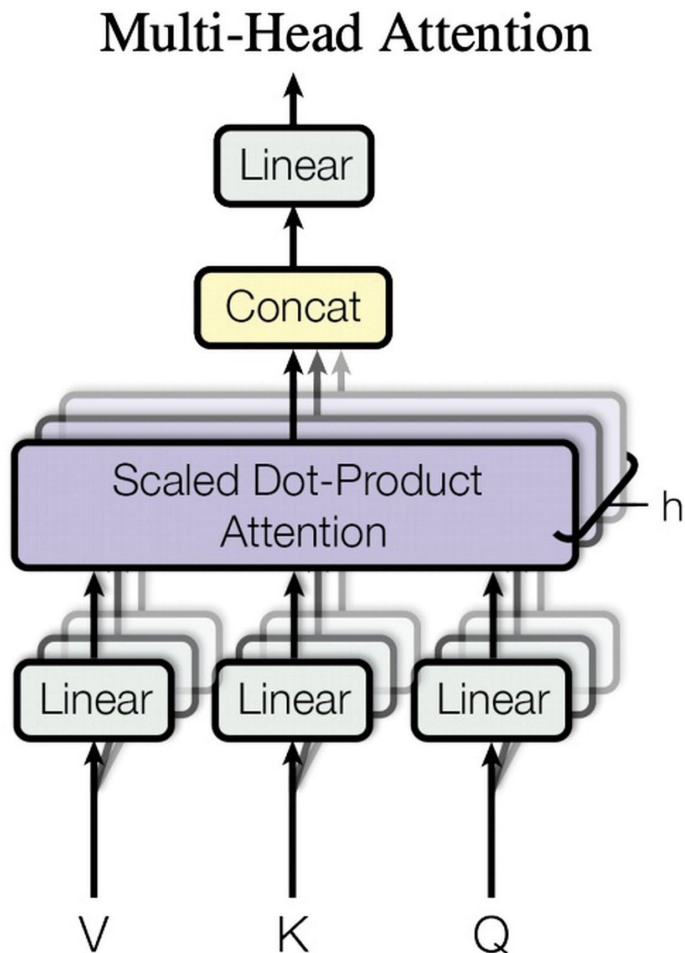
李玉光





# Tokens-to-Token ViT: Training Vision Transformers from Scratch on ImageNet

李玉光



$$Q_i = QW_i^Q, K_i = KW_i^K, V_i = VW_i^V, i = 1, \dots, 8$$

$$head_i = \text{Attention}(Q_i, K_i, V_i), i = 1, \dots, 8$$

$$\text{MultiHead}(Q, K, V) = \text{Concat}(head_1, \dots, head_8)W^O$$

这里，我们假设

$$Q, K, V \in R^{512}, W_i^Q, W_i^K, W_i^V \in R^{512 \times 64}, W^O \in R^{512 \times 512}, head_i \in R^{64}$$



# Tokens-to-Token ViT: Training Vision Transformers from Scratch on ImageNet

李玉光

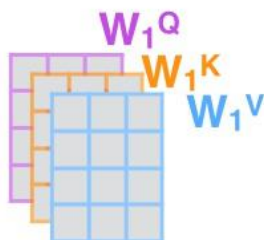
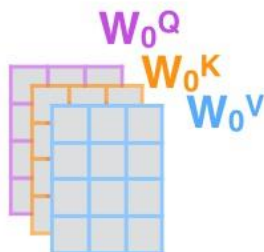
1) This is our input sentence\*

Thinking  
Machines

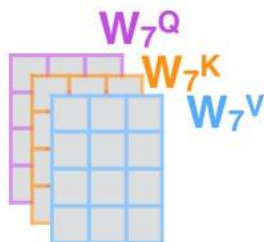
2) We embed each word\*



3) Split into 8 heads. We multiply  $X$  or  $R$  with weight matrices



...



4) Calculate attention using the resulting  $Q/K/V$  matrices



...



5) Concatenate the resulting  $Z$  matrices, then multiply with weight matrix  $W^O$  to produce the output of the layer



...



\* In all encoders other than #0, we don't need embedding. We start directly with the output of the encoder right below this one

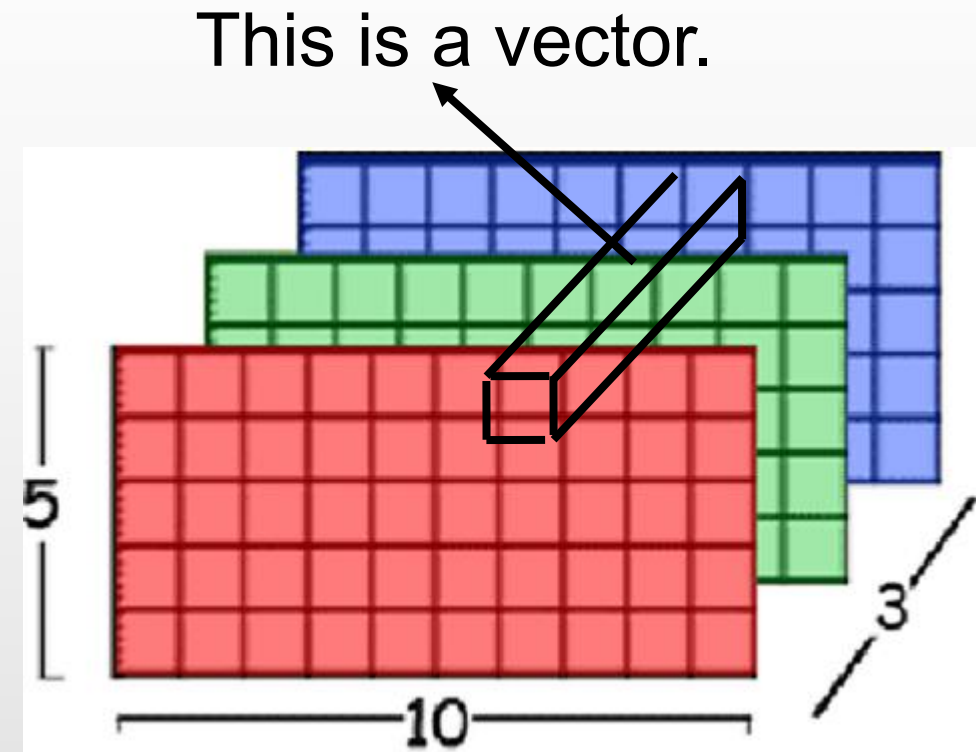
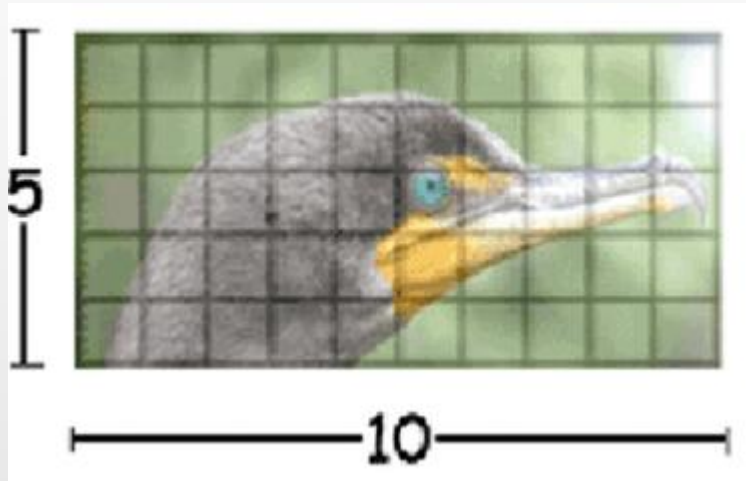


知乎 @随时学丫



## Self-attention for Image

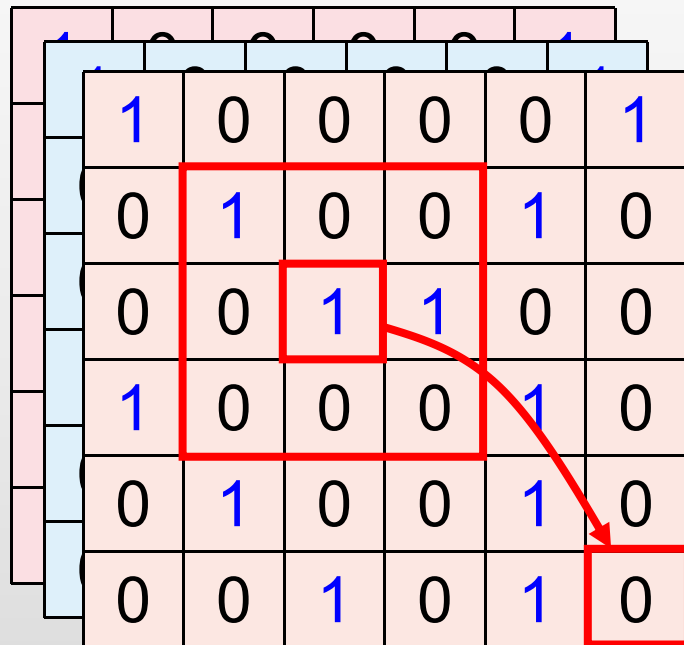
An **image** can also be considered as a **vector set**.



Source of image: [https://www.researchgate.net/figure/Color-image-representation-and-RGB-matrix\\_fig15\\_282798184](https://www.researchgate.net/figure/Color-image-representation-and-RGB-matrix_fig15_282798184)



## Self-attention v.s. CNN



CNN: self-attention that can only attends in a receptive field

➤ CNN is simplified self-attention.

Self-attention: CNN with learnable receptive field

➤ Self-attention is the complex version of CNN.



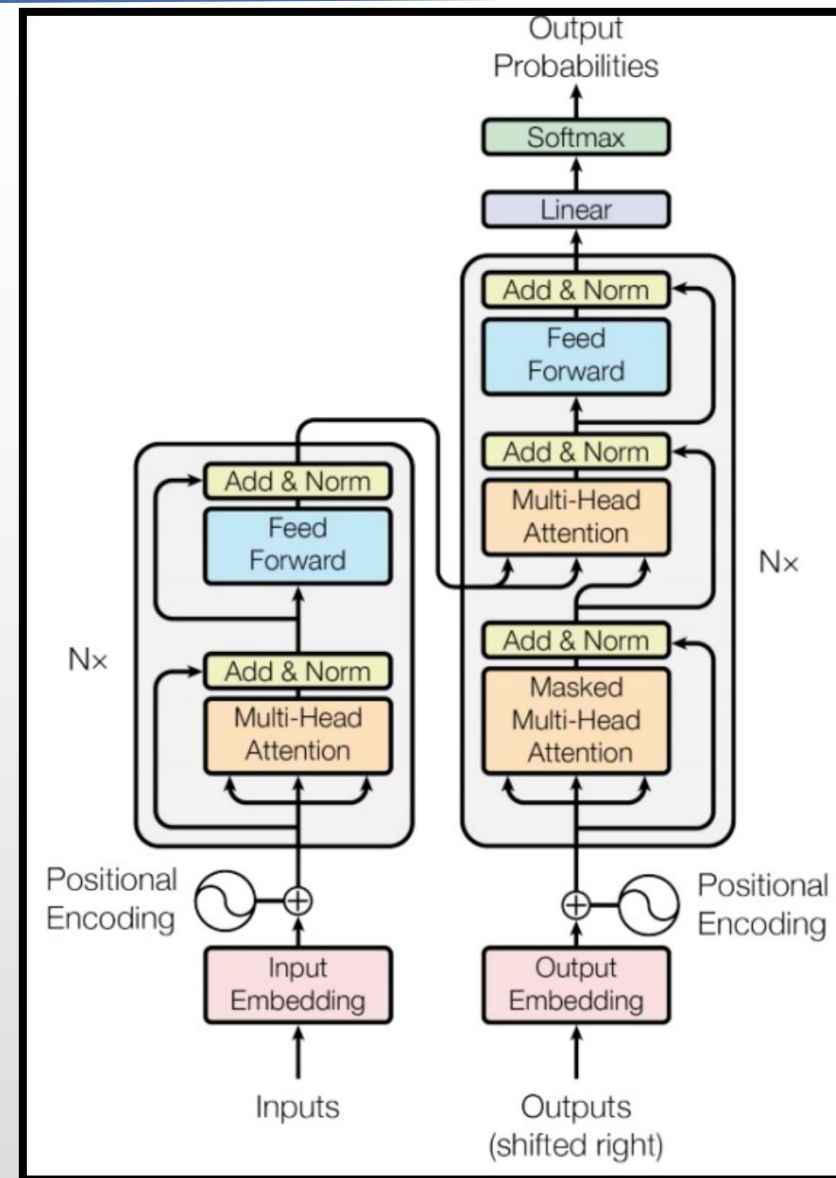


# Tokens-to-Token ViT: Training Vision Transformers from Scratch on ImageNet

李玉光

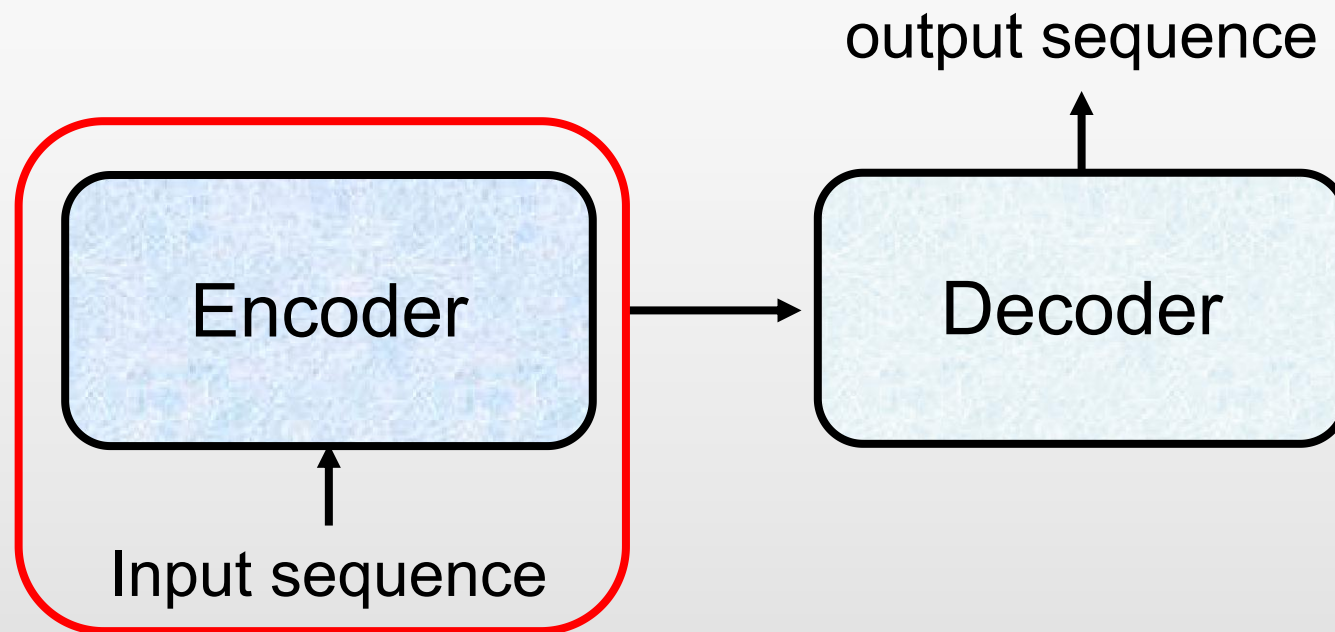
## TRANSFORMER

Transformer is the first transduction model relying entirely on self-attention to compute representations of its input and output without using sequence aligned RNNs or convolution.





# Encoder

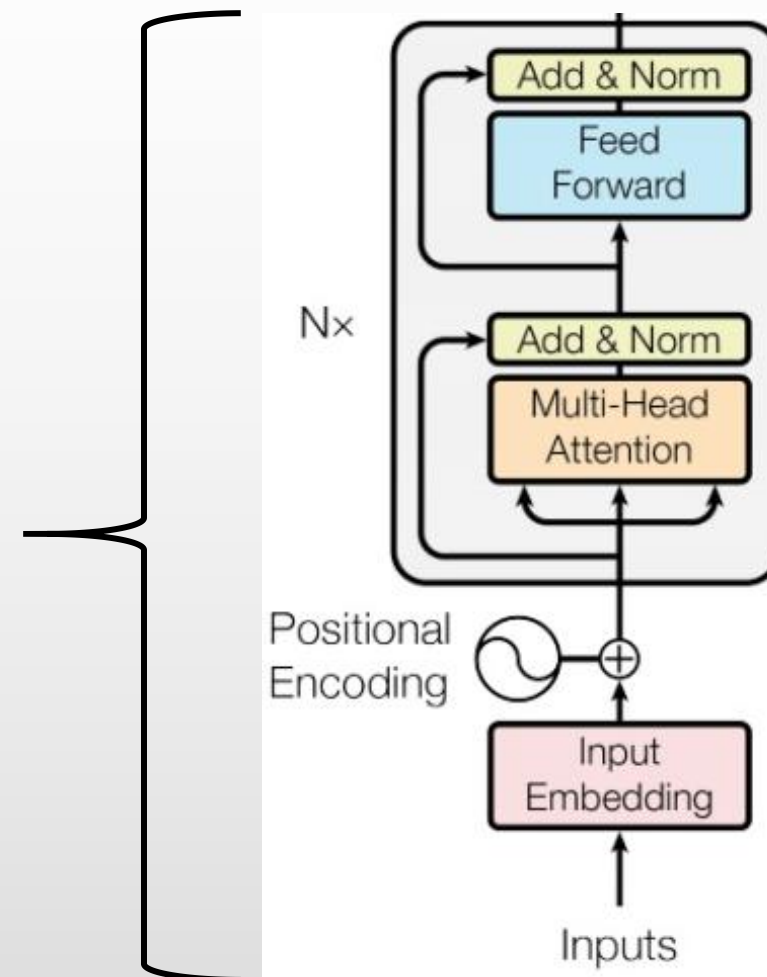
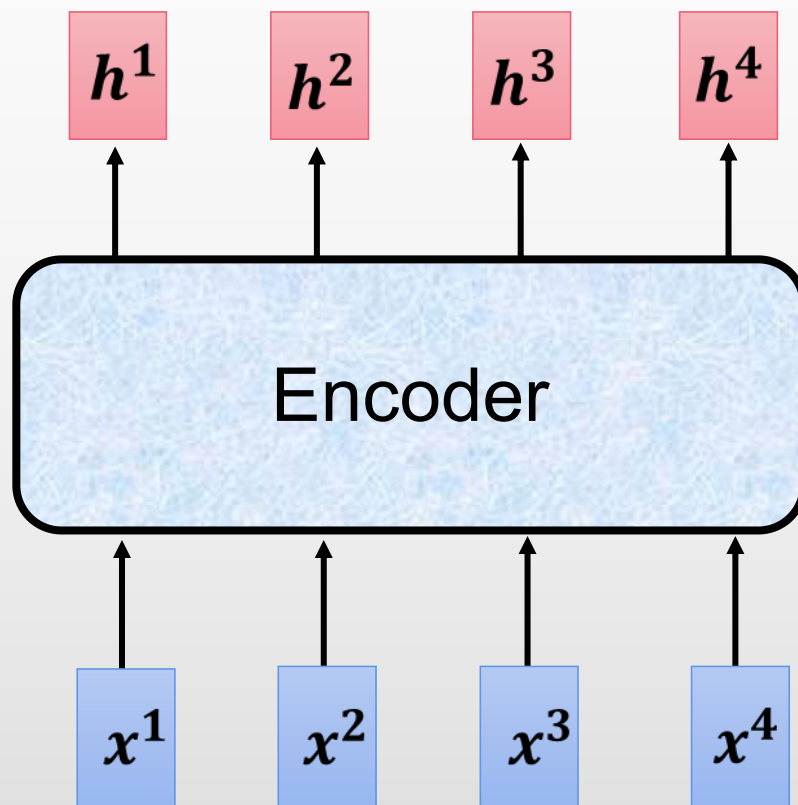




# Tokens-to-Token ViT: Training Vision Transformers from Scratch on ImageNet

李玉光

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

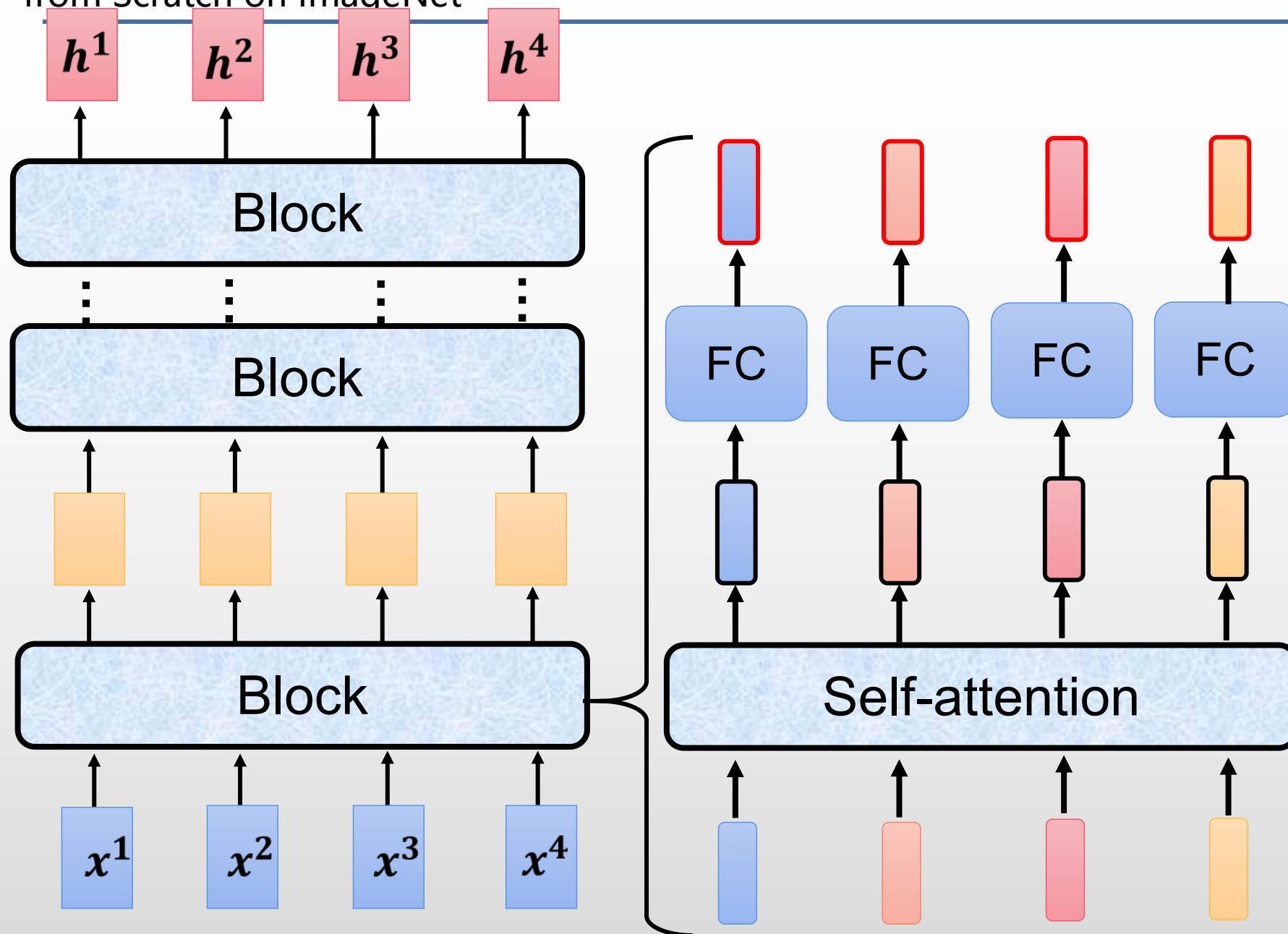






# Tokens-to-Token ViT: Training Vision Transformers from Scratch on ImageNet

李玉光

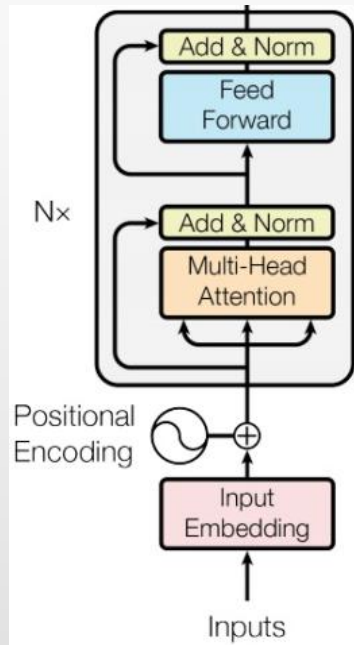




# Tokens-to-Token ViT: Training Vision Transformers from Scratch on ImageNet

李玉光

## POSITIONAL ENCODING



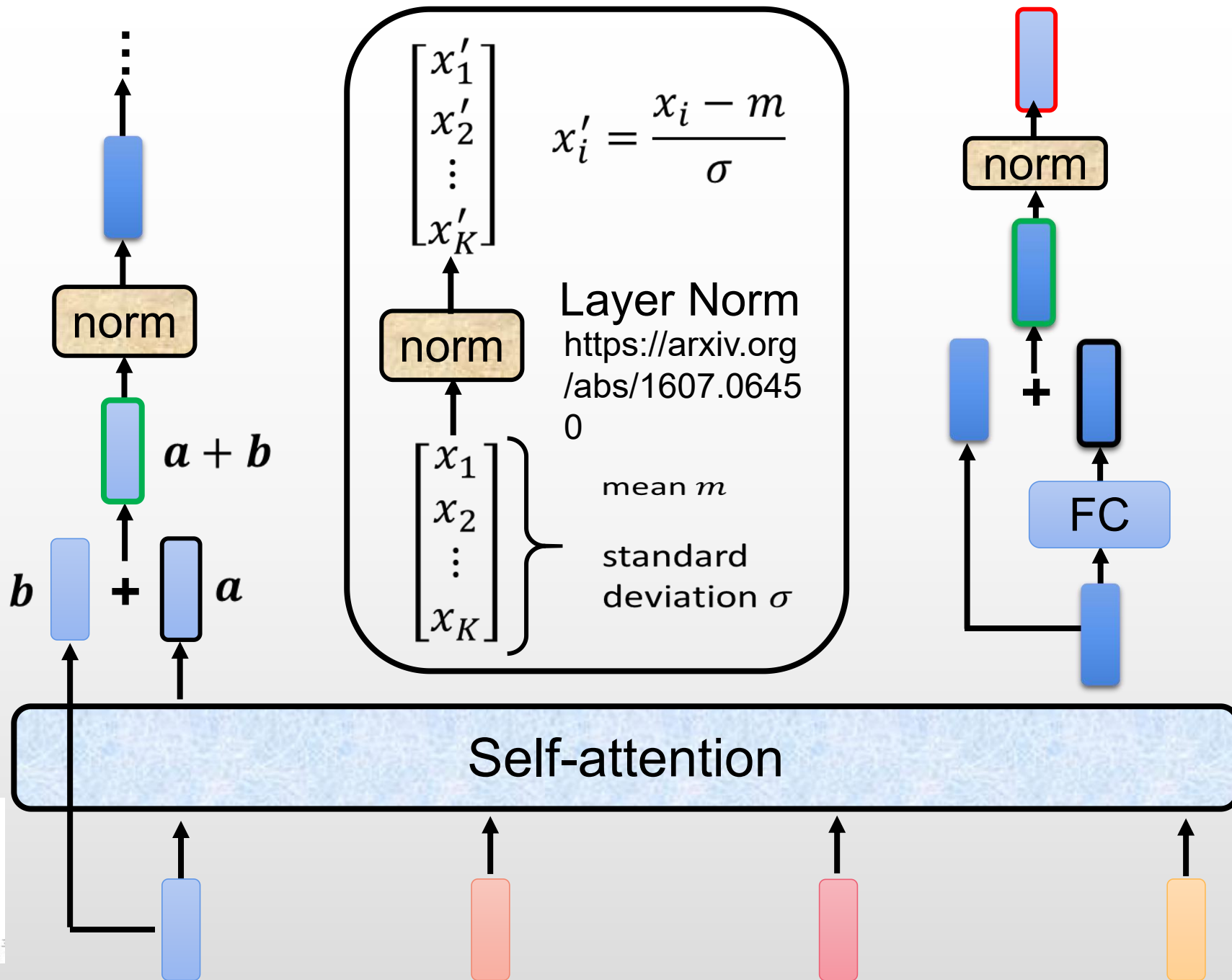
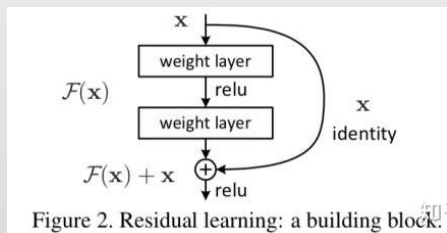
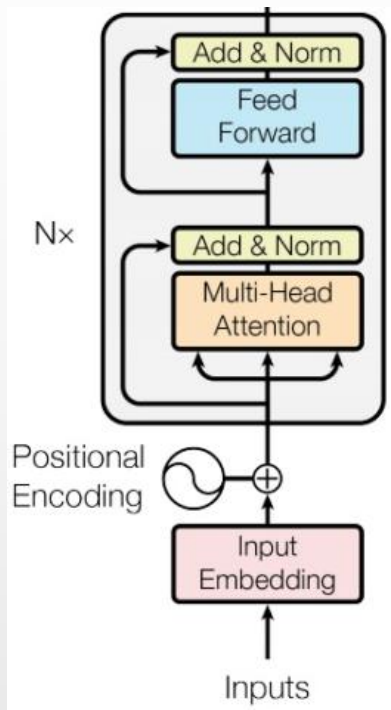
$$PE = pos = 0, 1, 2, \dots, T - 1$$

$$PE = pos / (T - 1)$$

$$PE(pos) = \sin\left(\frac{pos}{\alpha}\right)$$

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

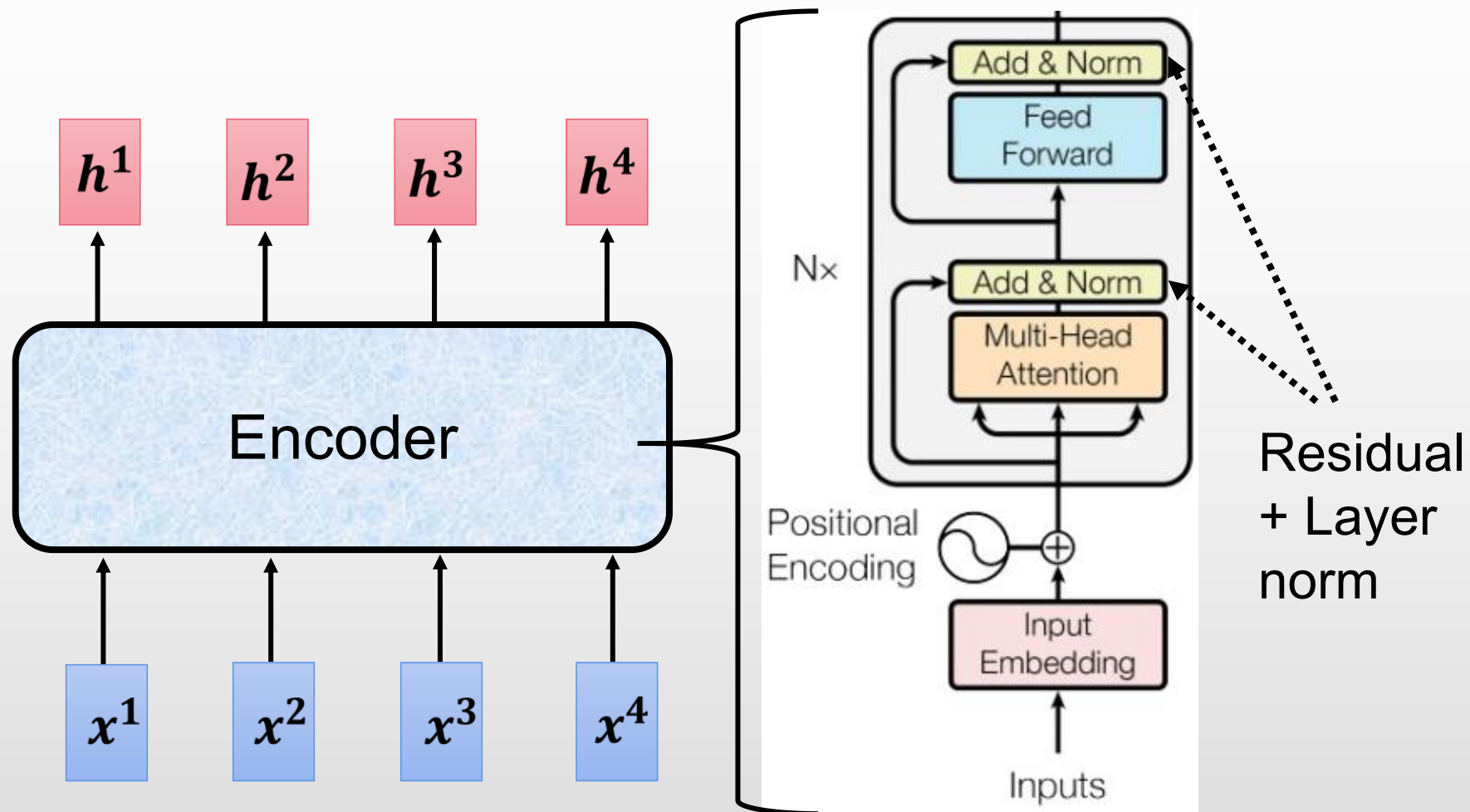
$$PE(pos, 2i + 1) = \cos\left(\frac{pos}{10000^{2i/d_{\text{model}}}}\right)$$





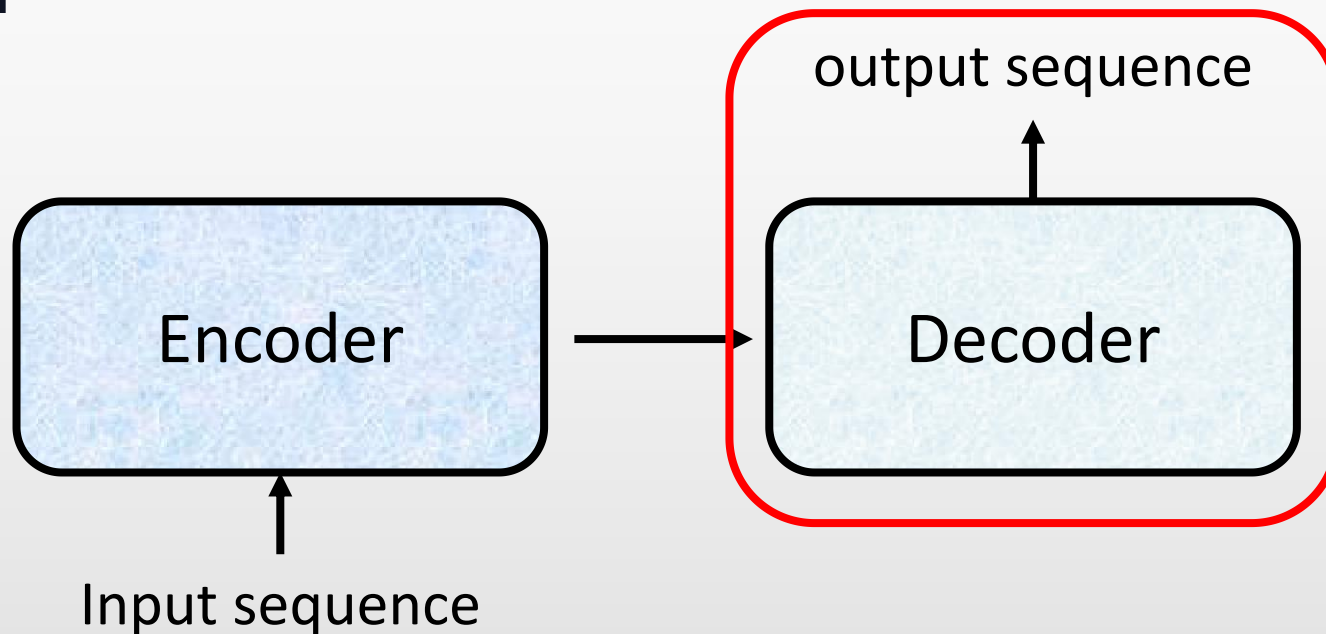
# Tokens-to-Token ViT: Training Vision Transformers from Scratch on ImageNet

李玉光





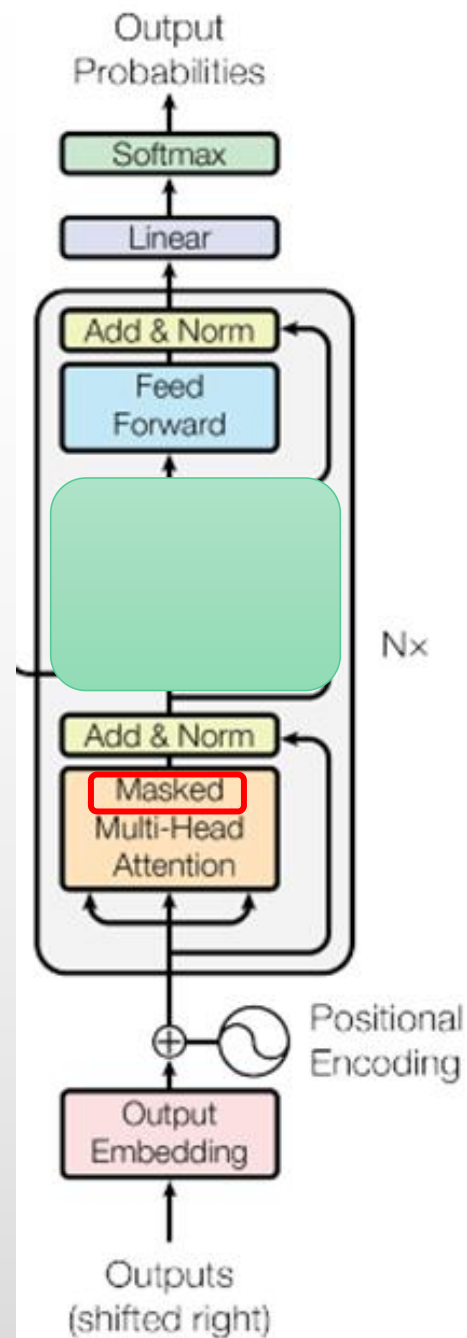
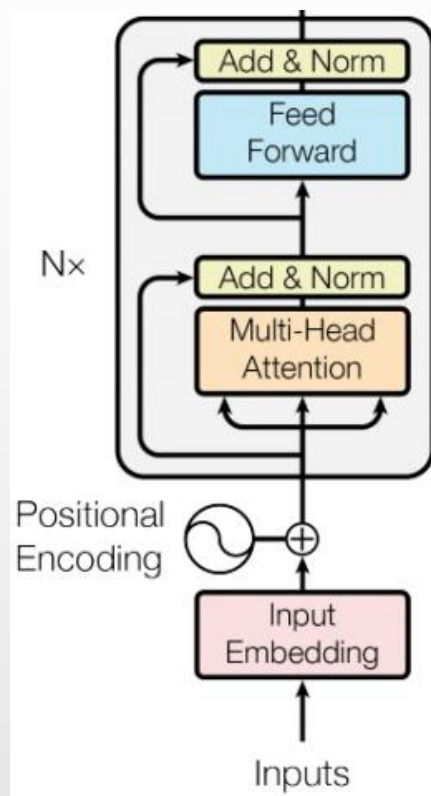
# Decoder





李玉光

## Encoder



## Decoder

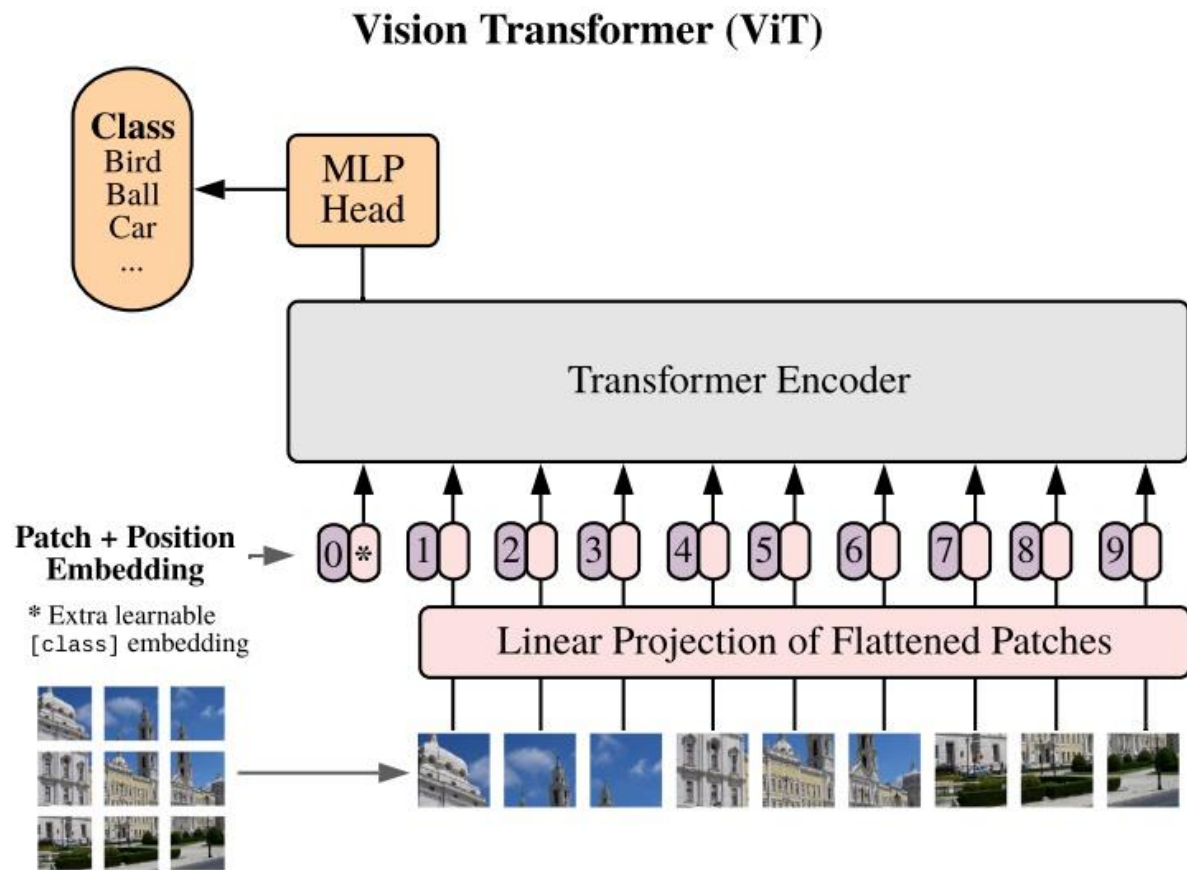




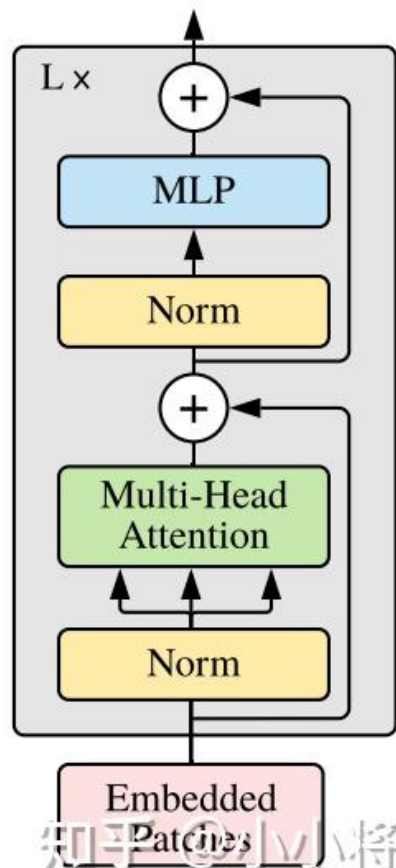
VIT

# Tokens-to-Token ViT: Training Vision Transformers from Scratch on ImageNet

李玉光



## Transformer Encoder





# Tokens-to-Token ViT: Training Vision Transformers from Scratch on ImageNet

---

李玉光

Tokens-to-Token ViT: Training Vision Transformers  
from Scratch on ImageNet

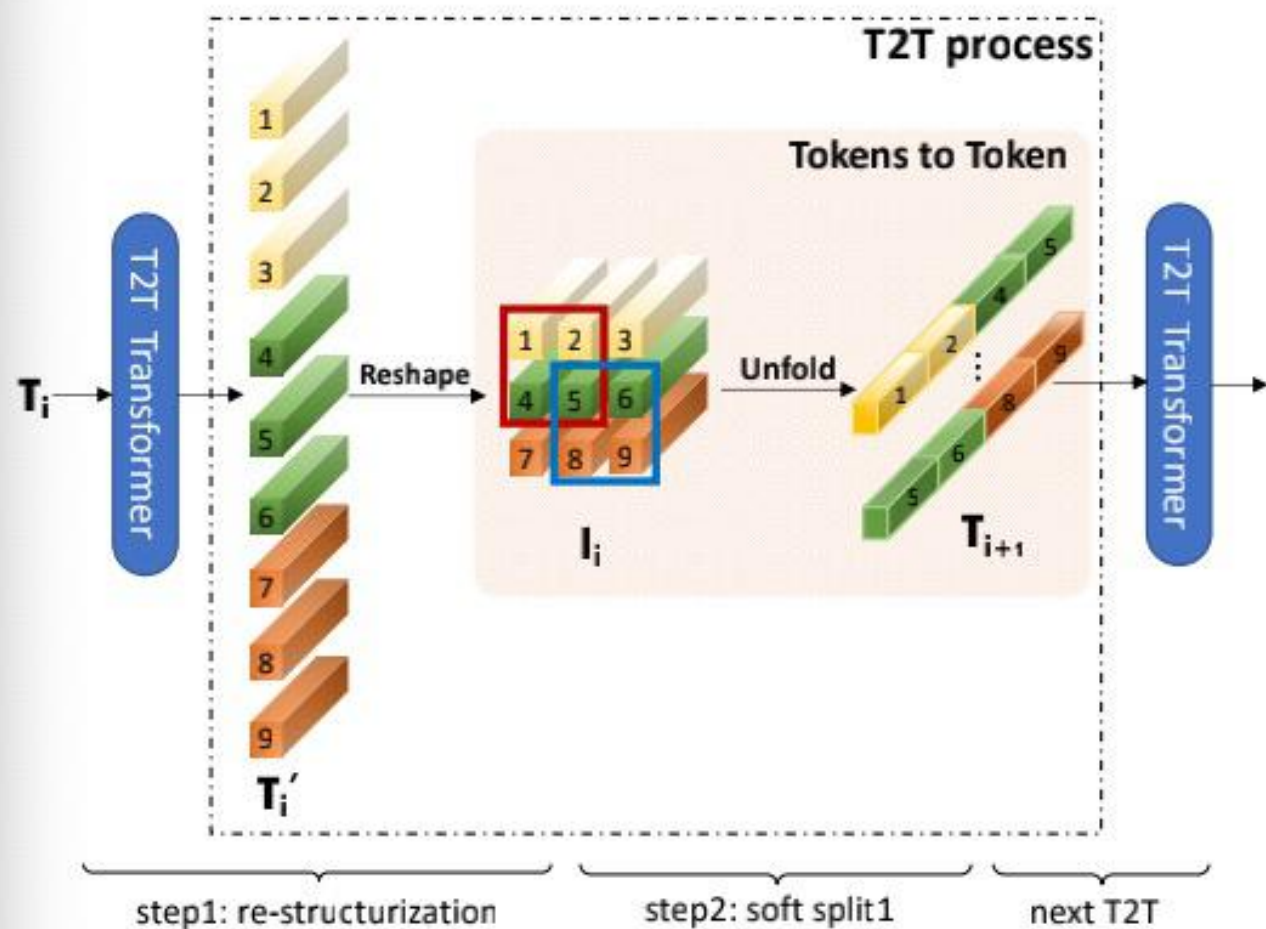




# Tokens-to-Token ViT: Training Vision Transformers from Scratch on ImageNet

李玉光

## T2T-ViT



### 1 ) Restructurization

$$T' = \text{MLP}(\text{MSA}(T))$$

$$I = \text{reshape}(T')$$

### 2 ) Soft Split

$$T_{i+1} = SS(I_i), i = 1, \dots, (n - 1)$$



## Backbone

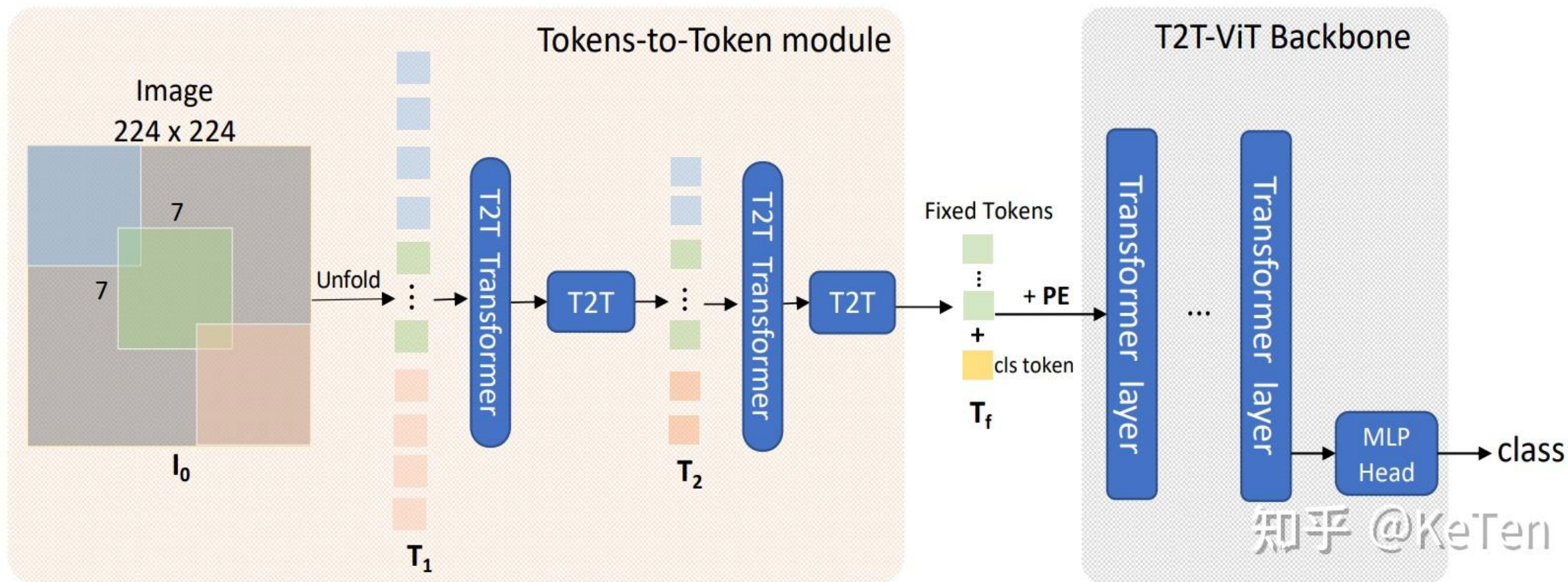
- Dense Connection , 类似于DenseNet ;
- Deep-narrow vs shallow-wide结构 , 类似于Wide-ResNet一文的讨论 ;
- Channel Attention , 类似SENet ;
- More Split Head , 类似ResNeXt ;
- Ghost操作 , 类似GhostNet。

结论 : Deep-Narrow结构可以在通道层面通过减少通道维度减少冗余 , 可以通过提升深度提升特征丰富性。



# Tokens-to-Token ViT: Training Vision Transformers from Scratch on ImageNet

李玉光





# Tokens-to-Token ViT: Training Vision Transformers from Scratch on ImageNet

---

李玉光

感谢大家！