

Transformer Model

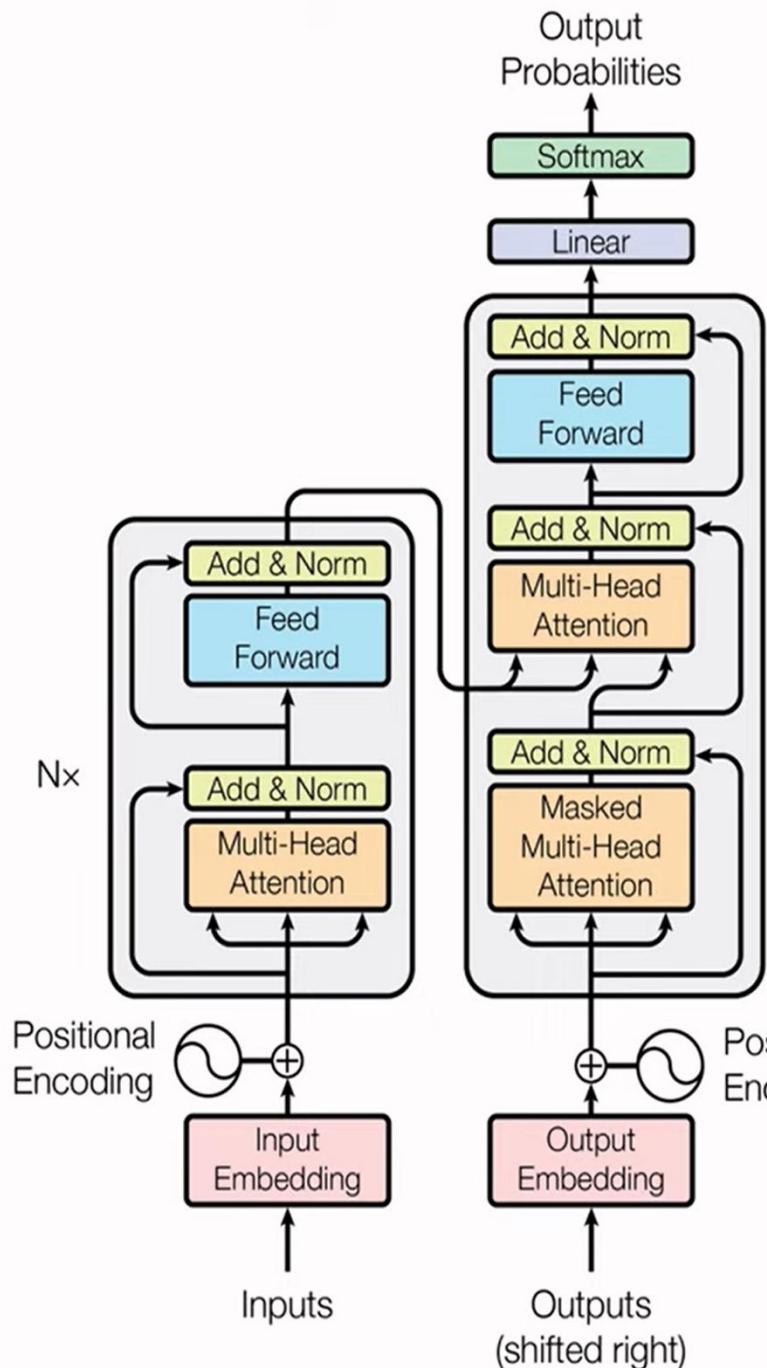
Architecture

Advantages and Disadvantages

Performance Comparison

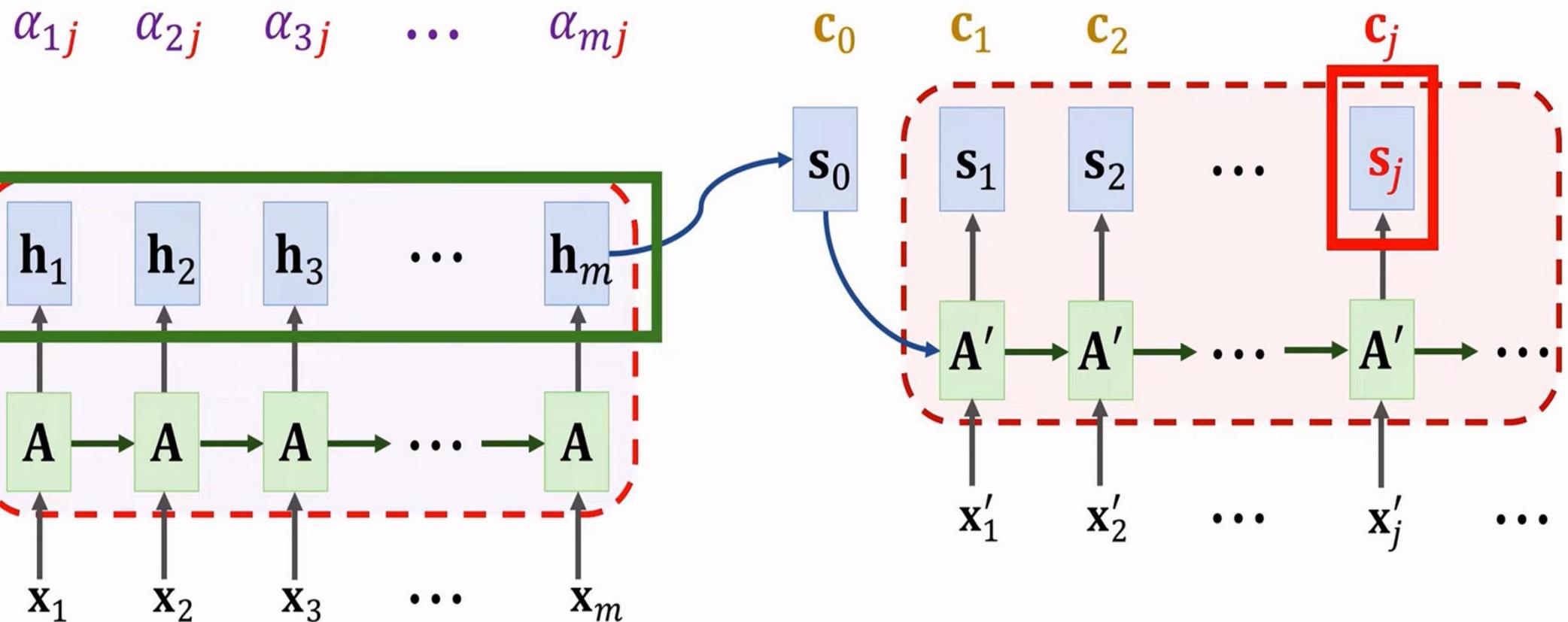
Transformer Model

- Transformer is a Seq2Seq model.
- Transformer is not RNN.
- Purely based on attention and dense layers.
- Higher accuracy than RNNs on large datasets.

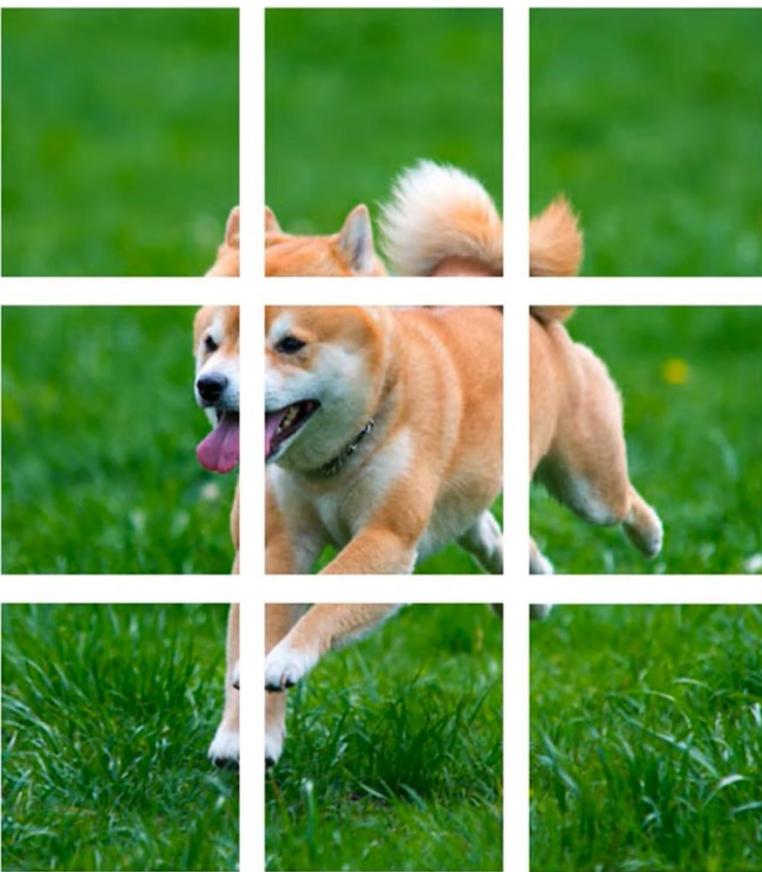


Attention for Seq2Seq Model

Weights: $\alpha_{ij} = \text{align}(\mathbf{h}_i, \mathbf{s}_j)$.



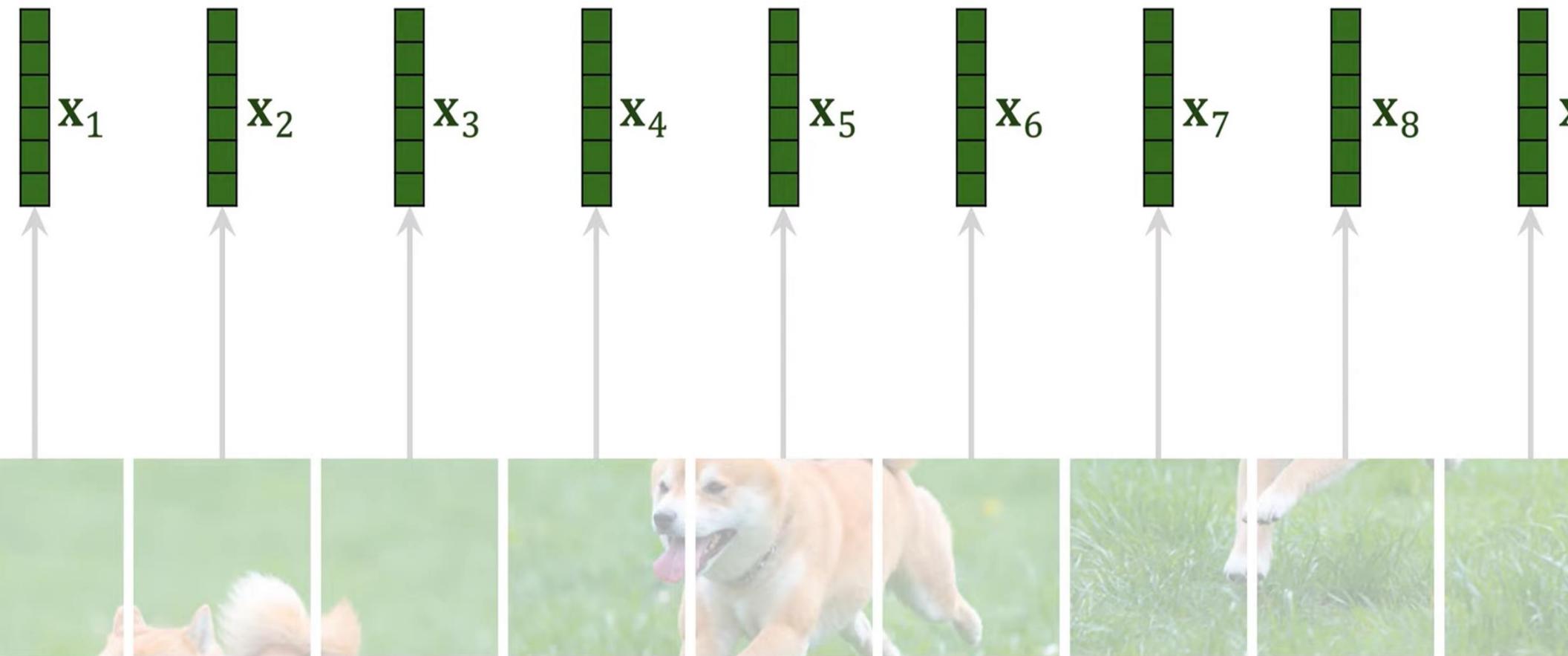
Split Image into Patches

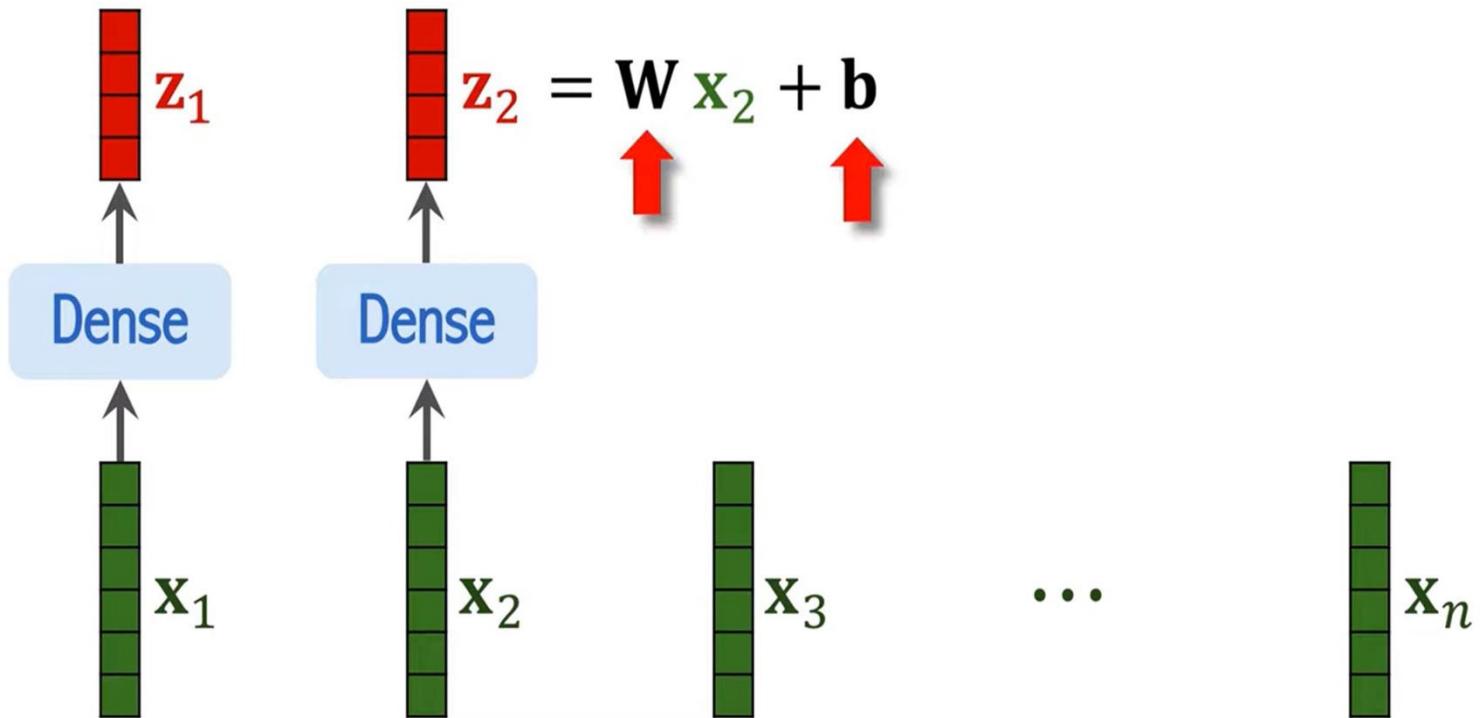


- Here, the patches do not overlap.

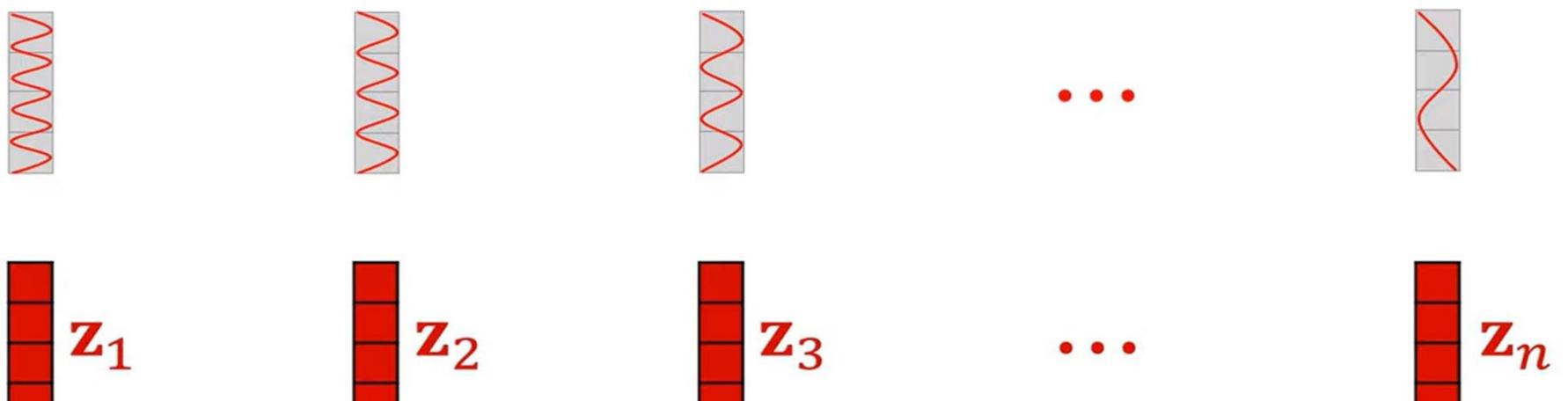
Vectorization

If the patches are $d_1 \times d_2 \times d_3$ tensors, then the vectors are $d_1 d_2 d_3 \times 1$.

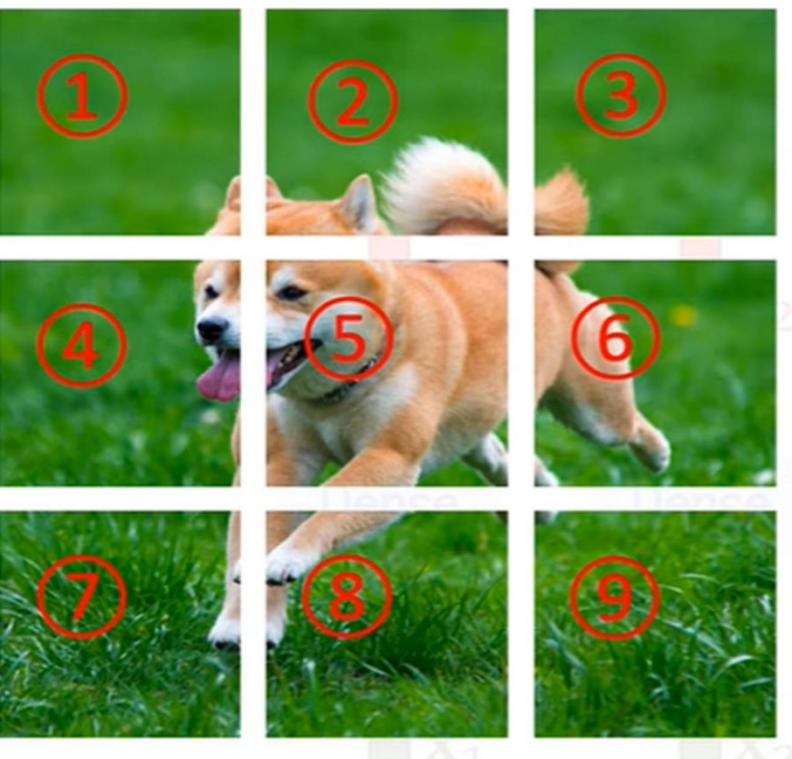


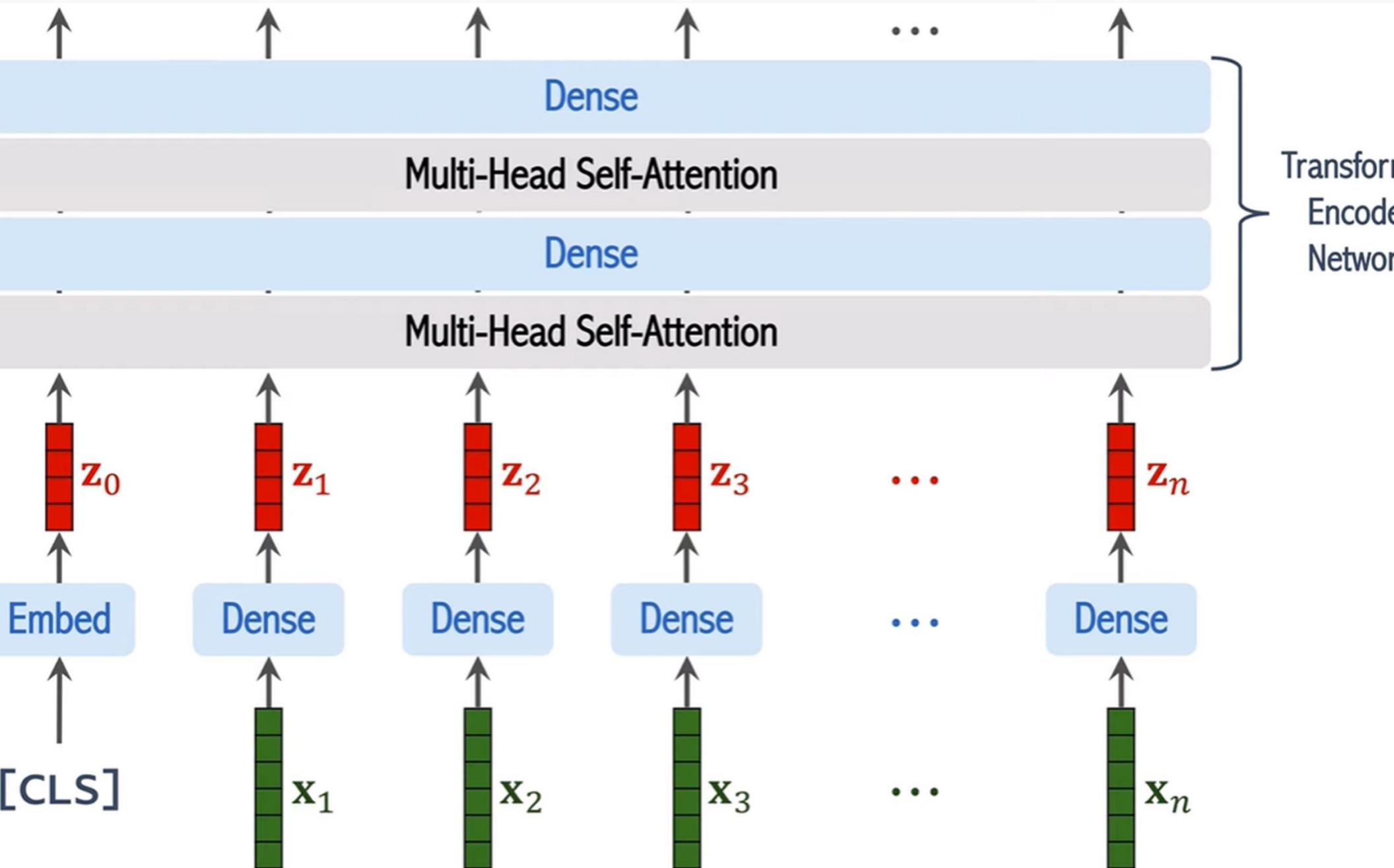


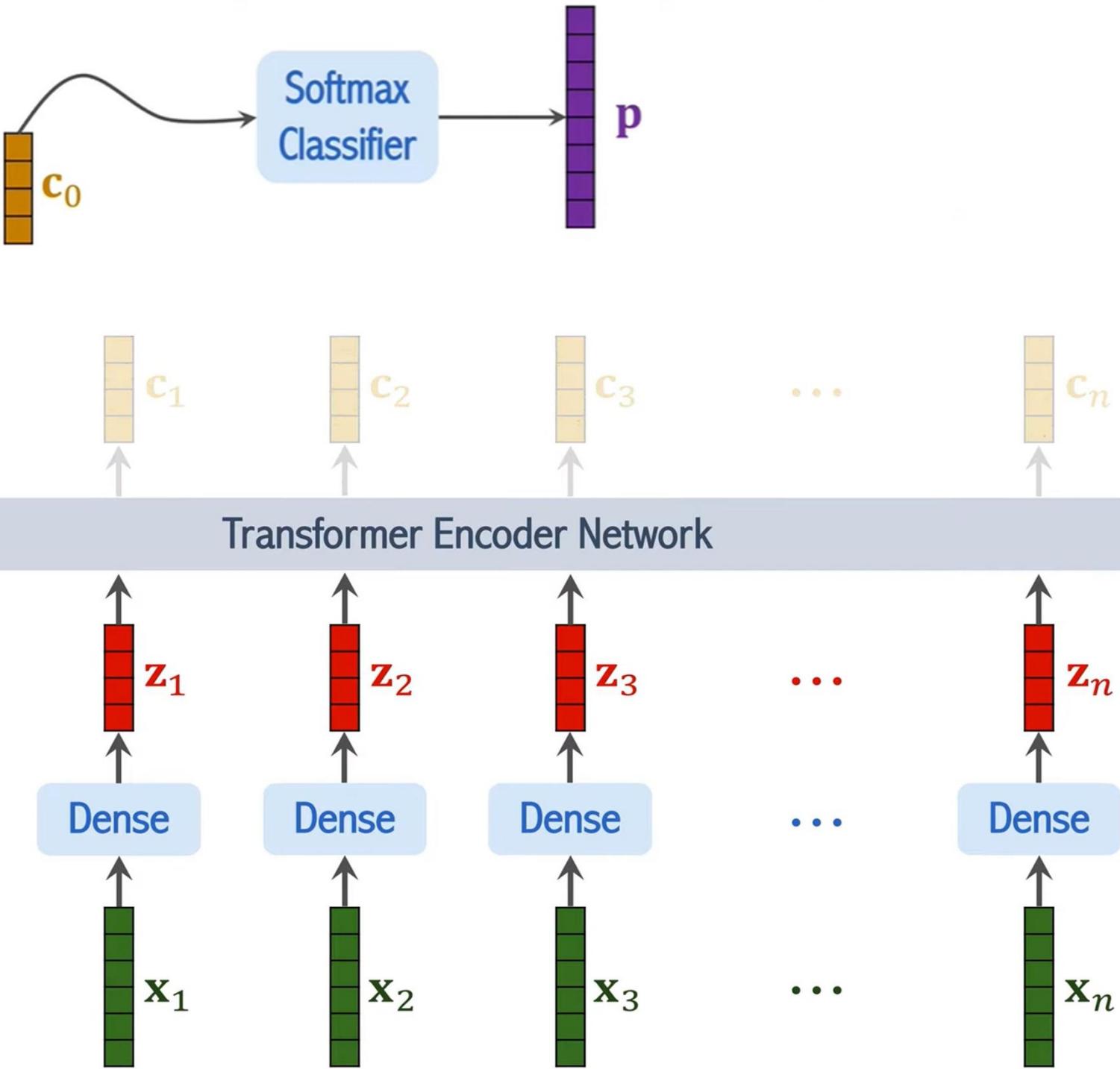
Add positional encoding vectors to $\mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_n$.



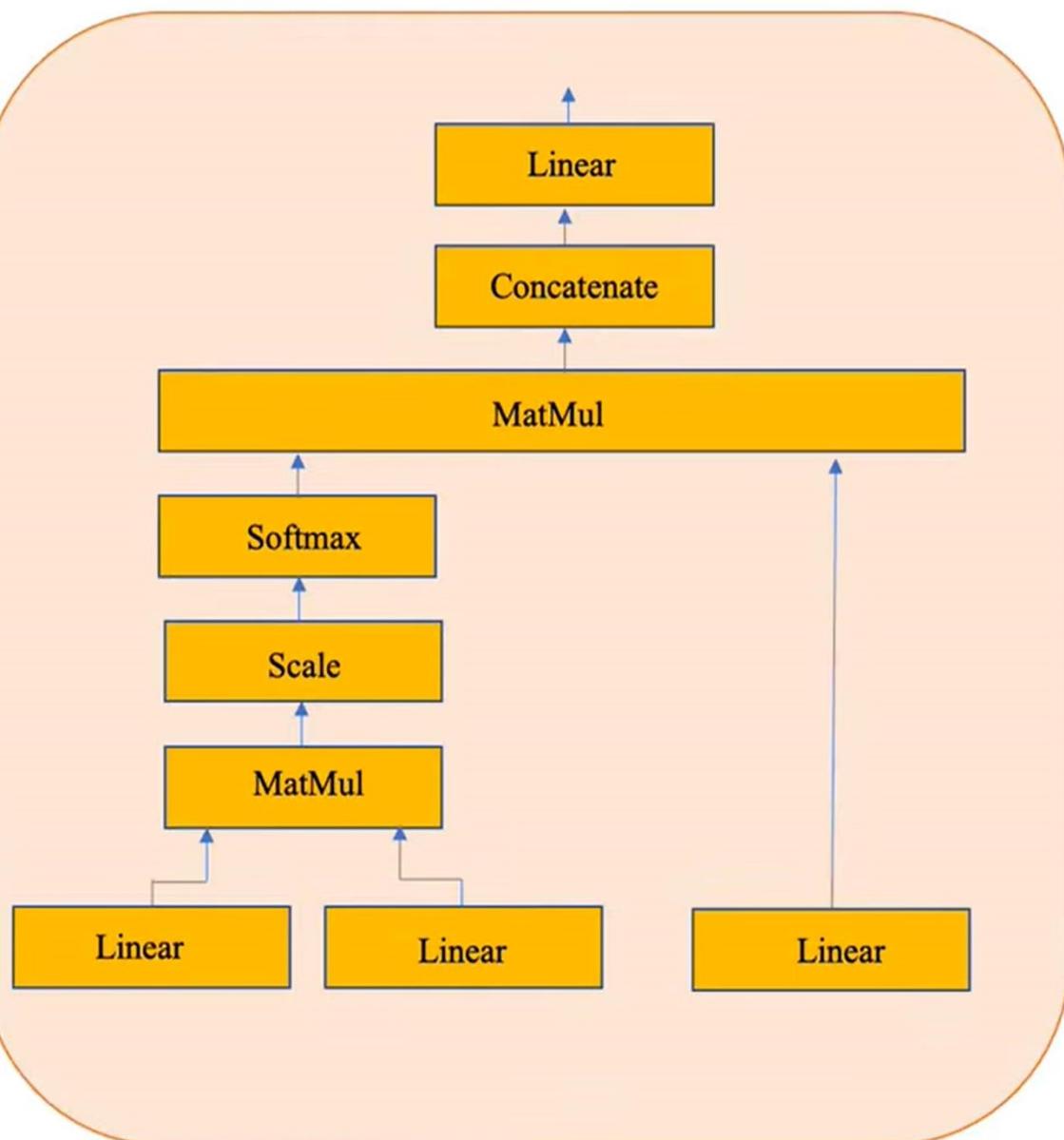
Add positional encoding vectors to $\mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_n$. (Why?)







Multi-Head Attention

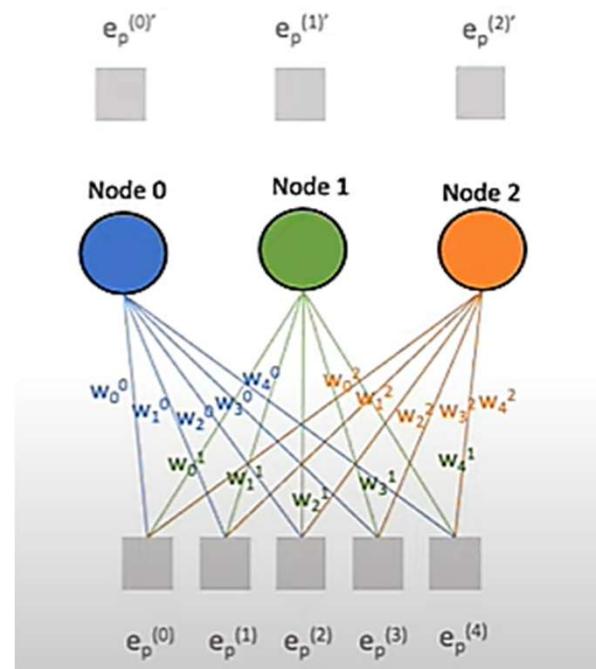


Linear Layer



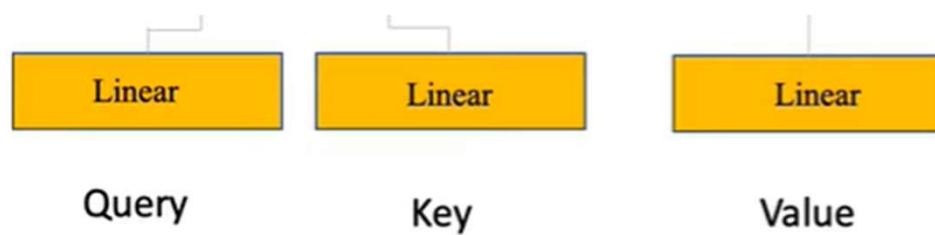
Linear Layer

$$W_0^0 = 0.6$$



) Mapping inputs onto the outputs

i) Changing matrix/vector dimensions

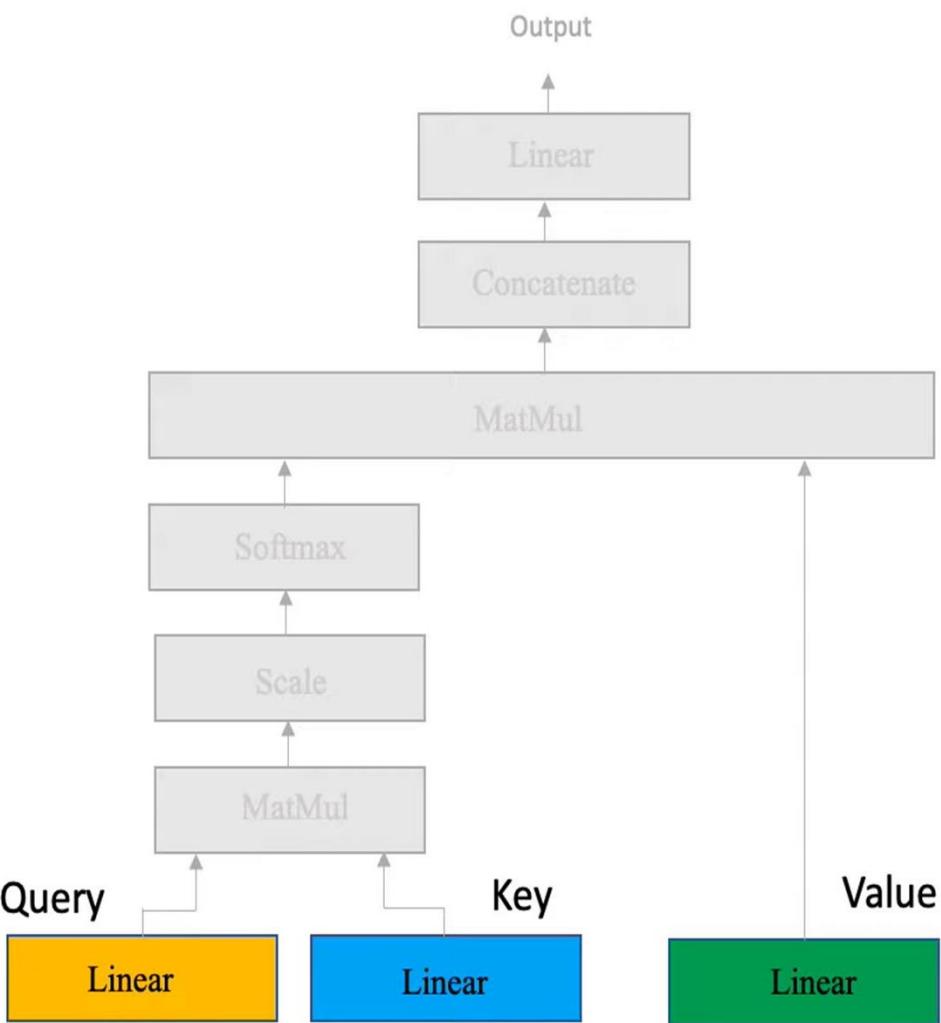


Cosine Similarity

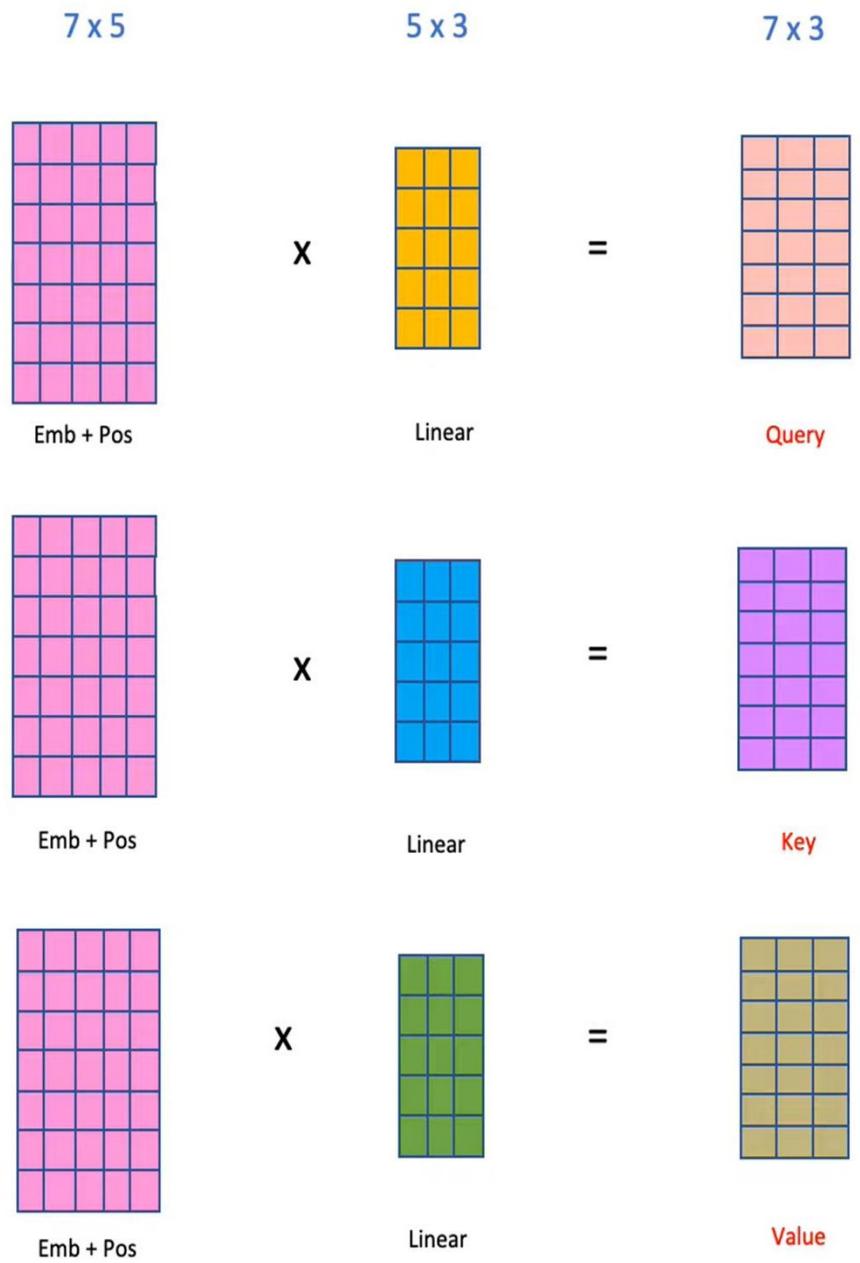
$$\text{Cos}(A, B) = \frac{A \cdot B}{|A| |B|}$$

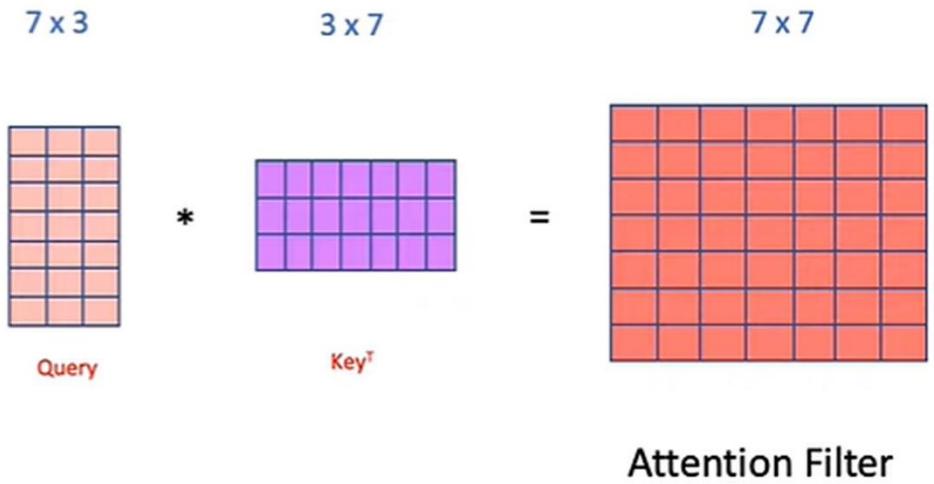
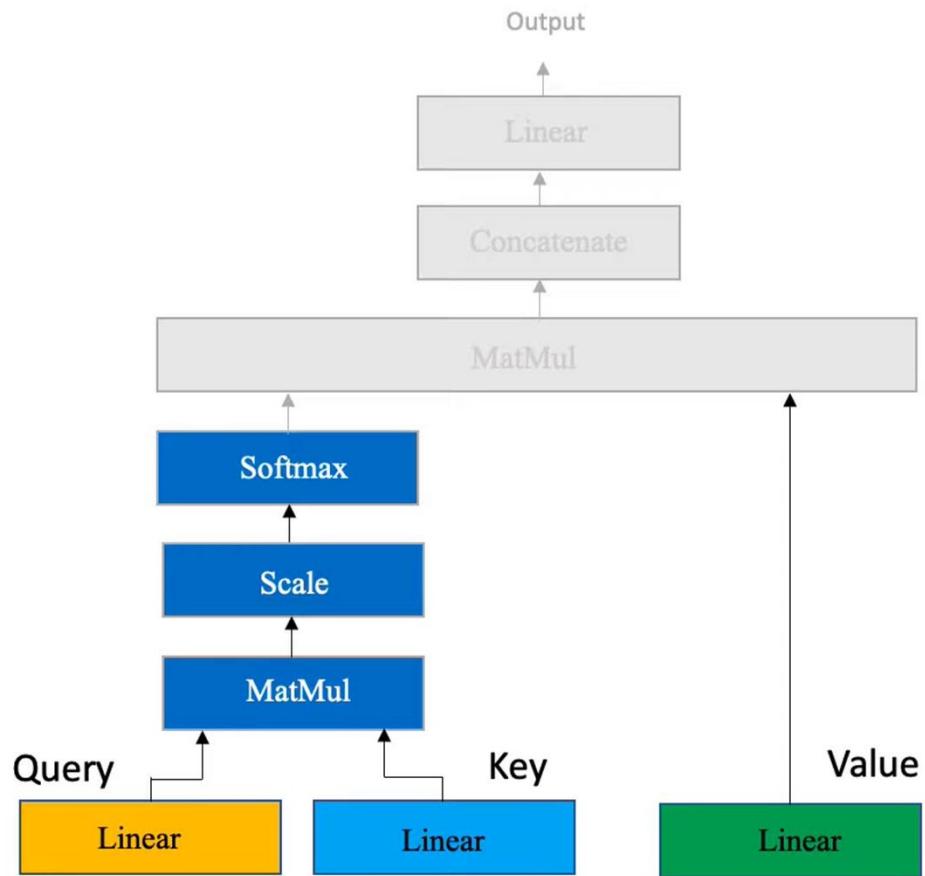
Similarity b/w Matrices

$$\text{similarity}(A, B) = \frac{A \cdot B^T}{\text{scaling}}$$

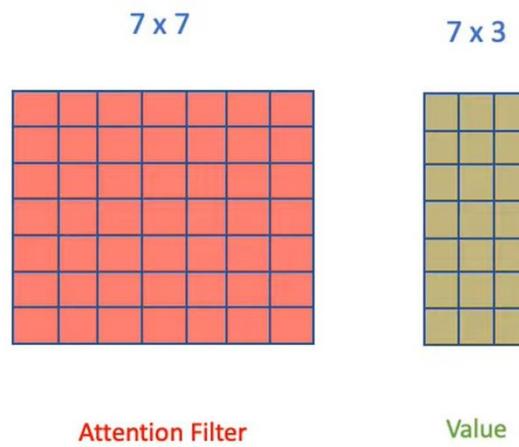


Multi-Head Attention





Attention Filter



Attention Filter

$$\text{softmax}\left(\frac{Q \cdot K^T}{\sqrt{d_k}}\right) V$$

Intuition

Attention Filter



*

Original Image



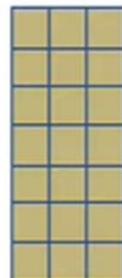
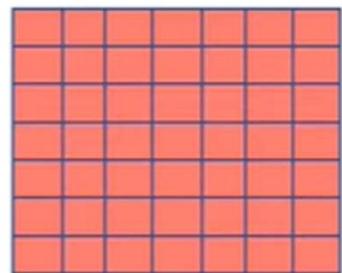
Filtered Image



7×7

7×3

7×3



*



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

Attention Filter

Original Value

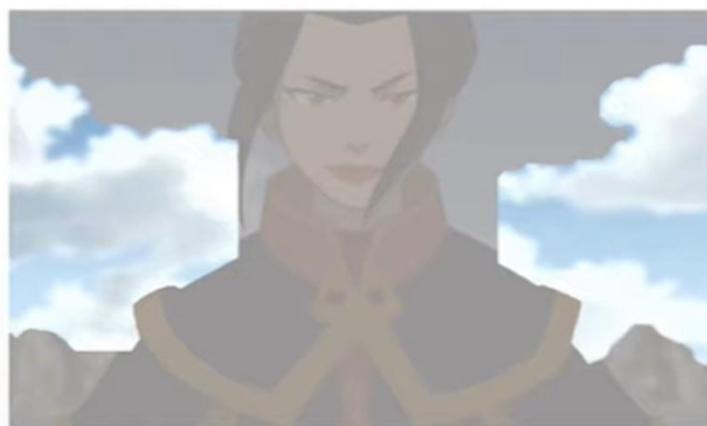
Filtered Value

Intuition for Multi-head Attention

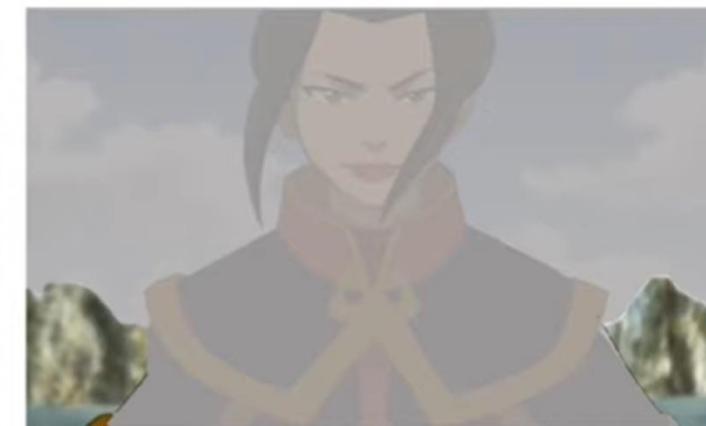
Attention Filter 1



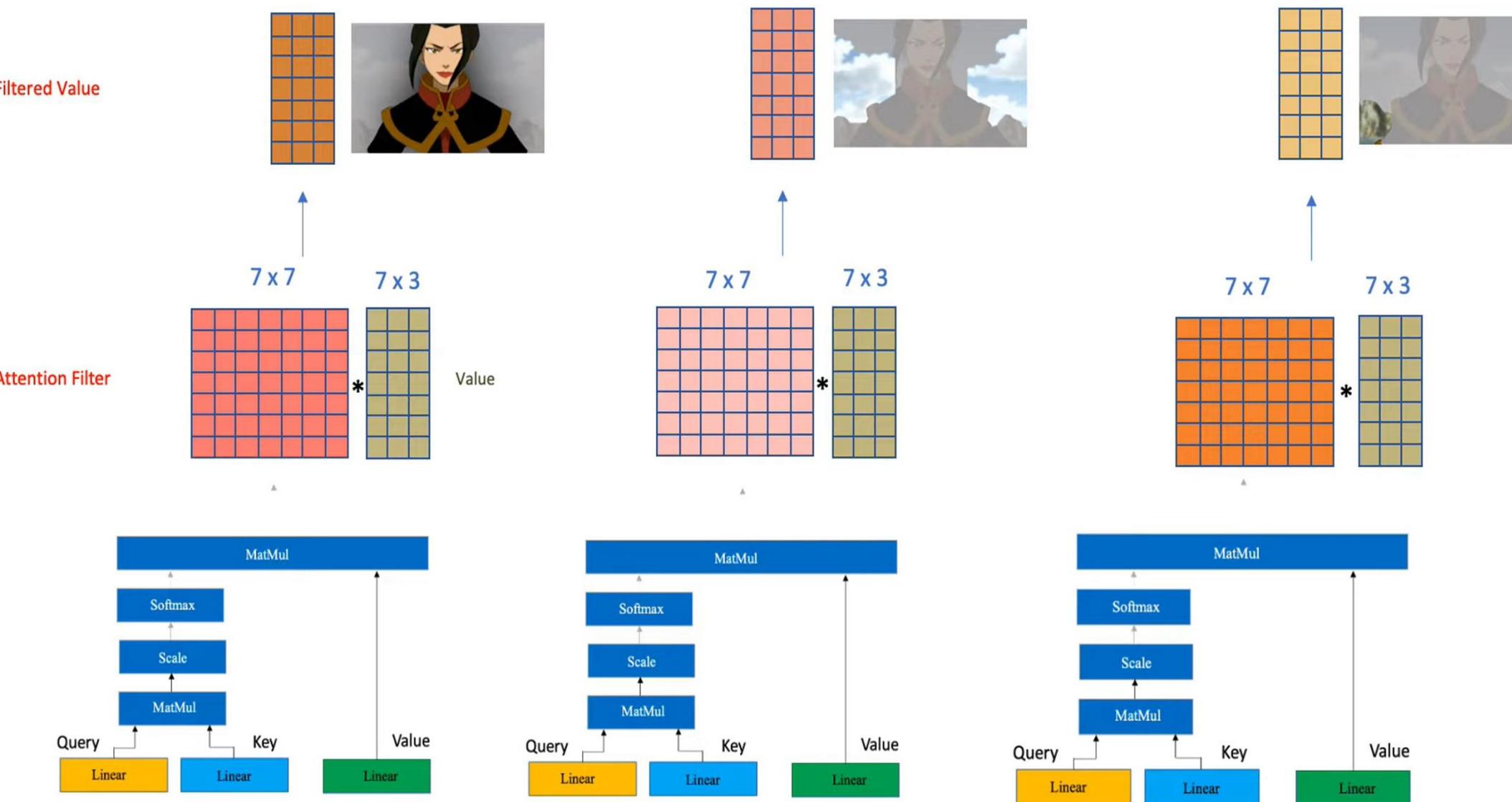
Attention Filter 2



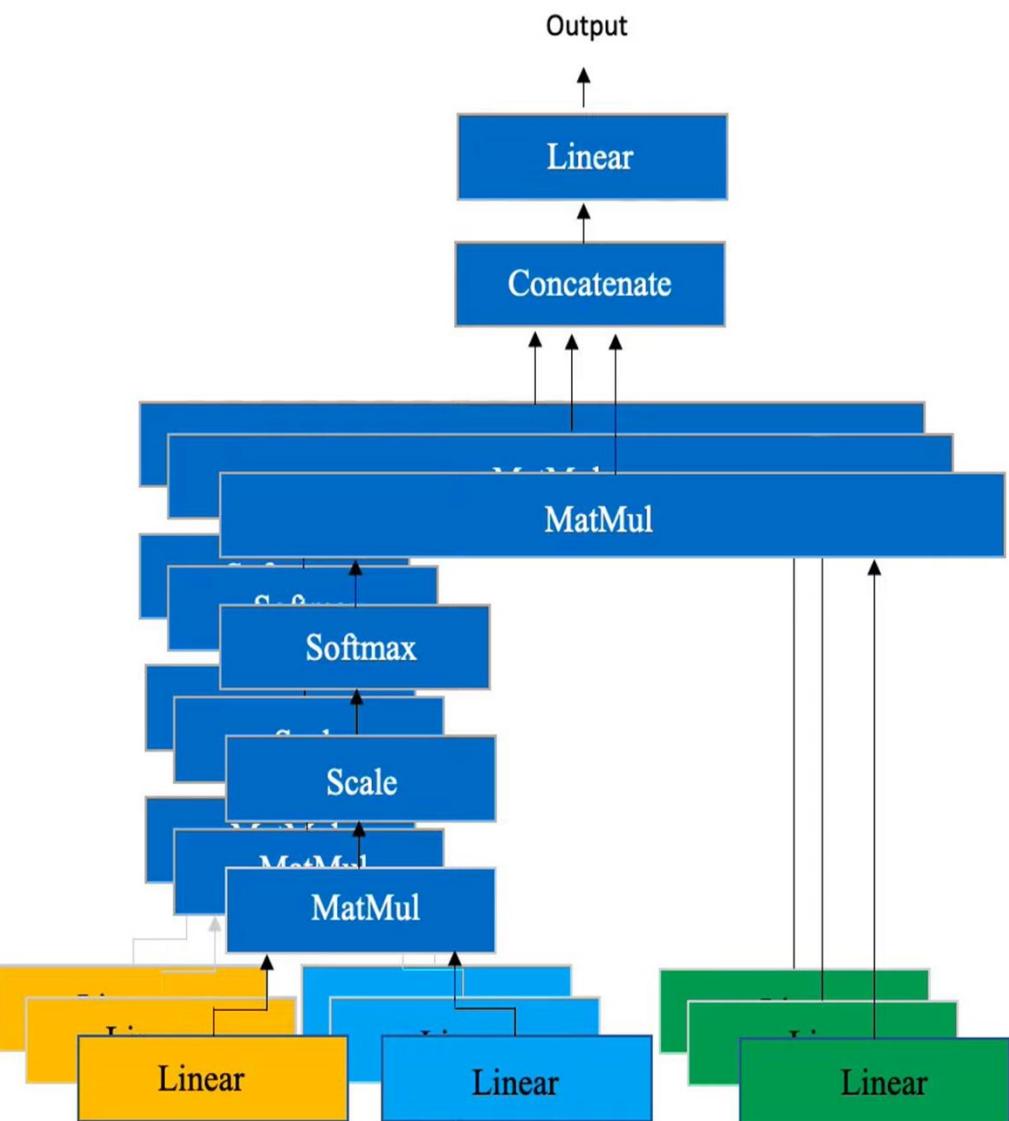
Attention Filter 3



Multi-Head Attention



Multi-Head Attention



21×3



Z

Linear

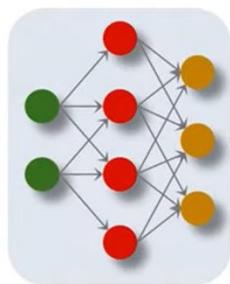
7×5



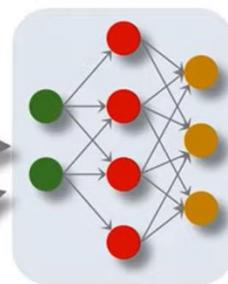
Output of
Multihead
Attention

Training of Transformer Model

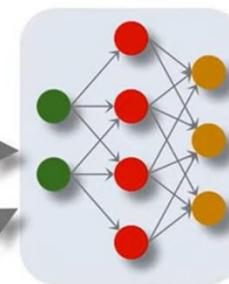
Randomly
Initialized



Pretrained



Fine-tuned



Test
Accuracy

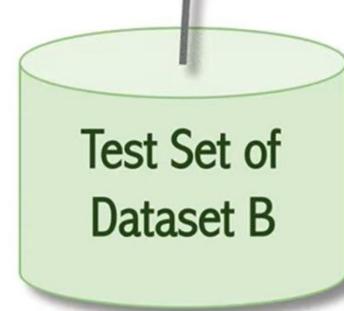
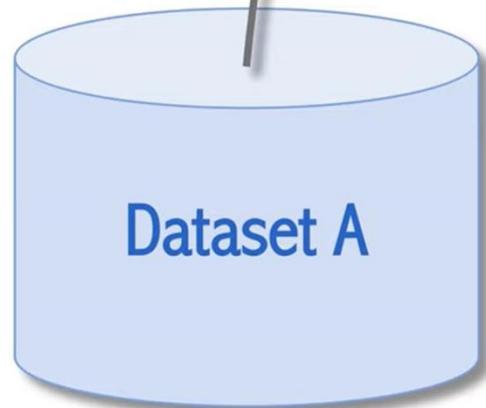
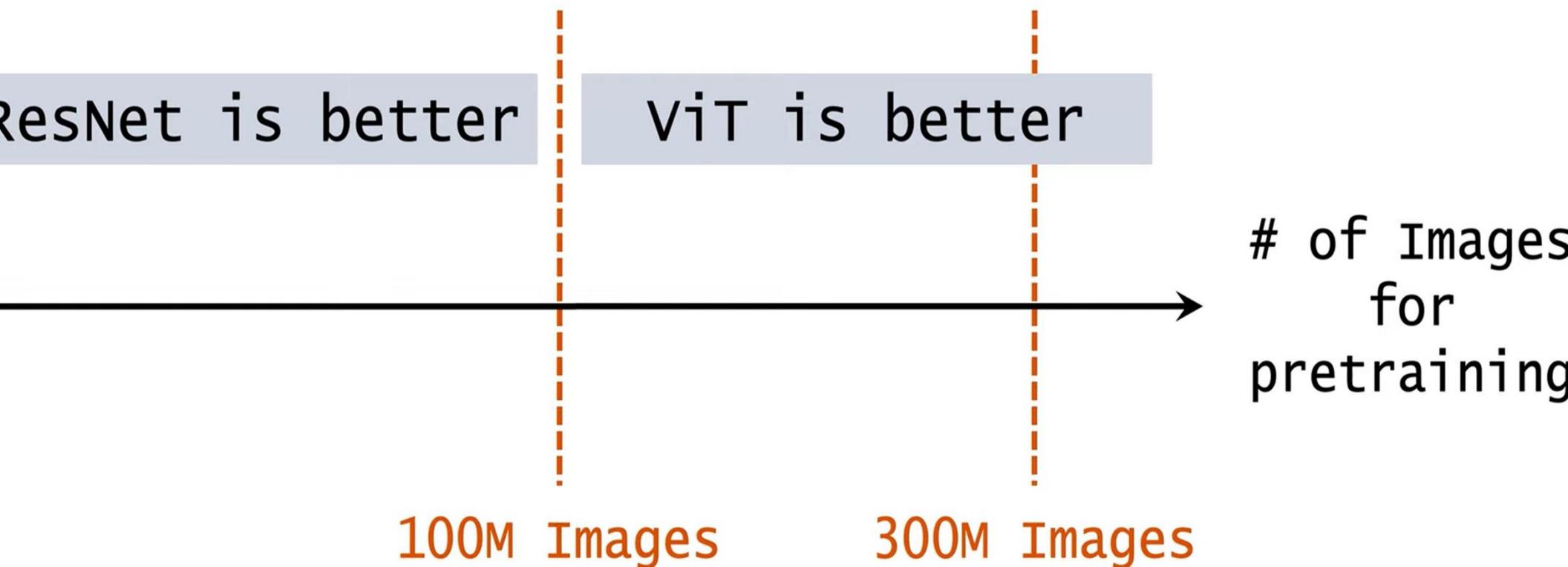


Image Classification Accuracies

- Pretrain the model on **Dataset A**, fine-tune the model on **Dataset B**, and evaluate the model on **Dataset B**.
- Pretrained on **ImageNet (small)**, ViT is slightly **worse** than ResNet.
- Pretrained on **ImageNet-21K (medium)**, ViT is **comparable** to ResNet.
- Pretrained on **JFT (large)**, ViT is slightly **better** than ResNet.

Image Classification Accuracies



THANK YOU