

# Transformer Model

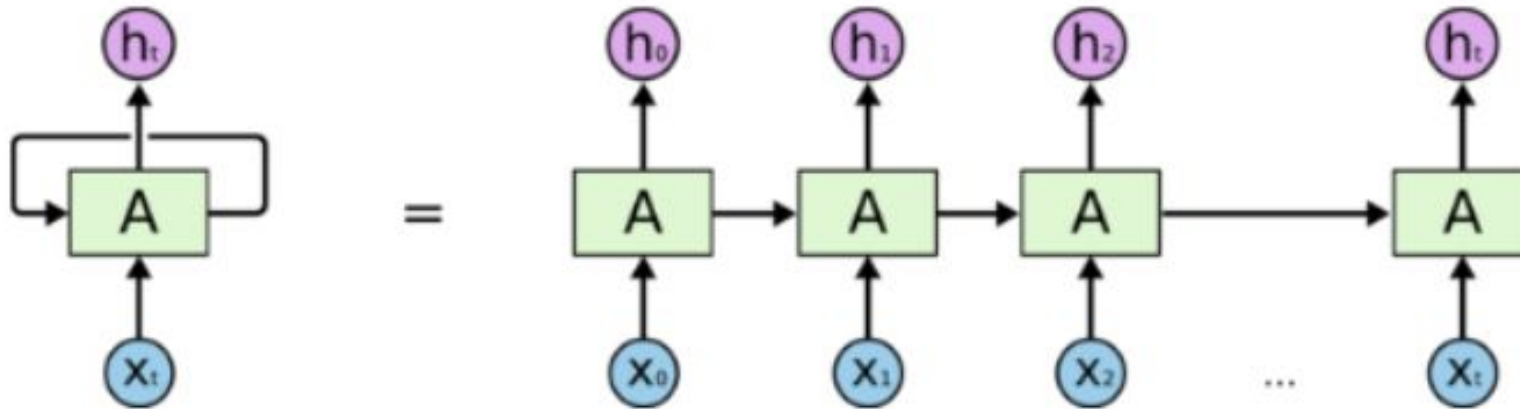
Architecture

Advantages and Disadvantages

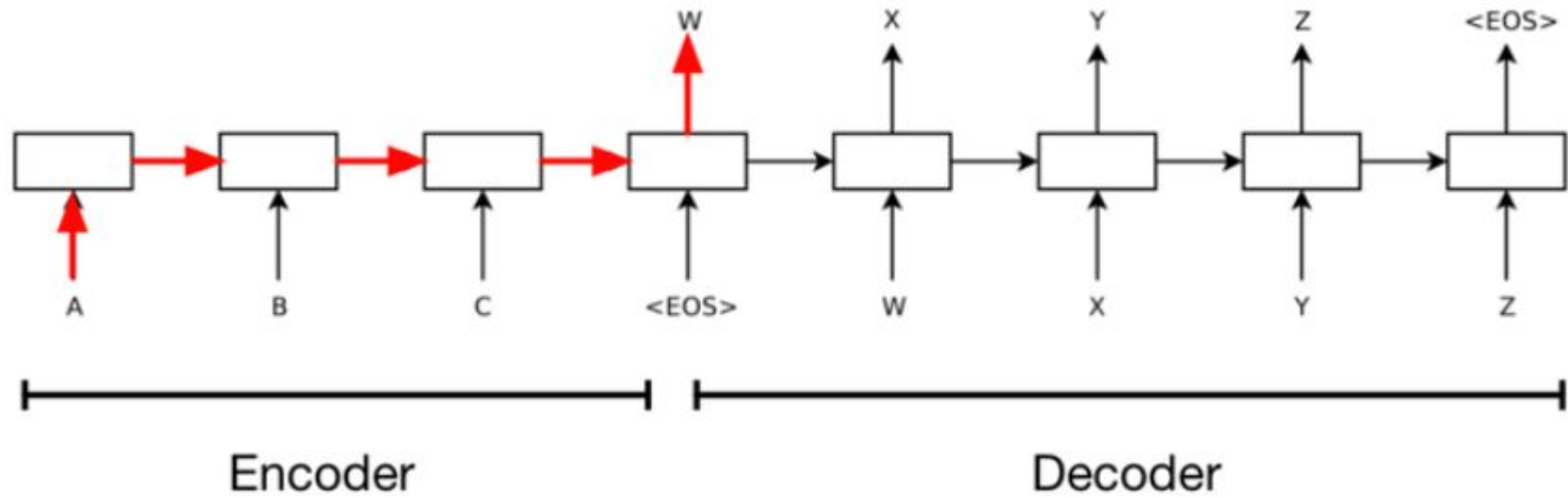
Performance Comparison

# RNN's

- Can handle time series data
- Suffers from vanishing / exploding gradient problem
- Solution? LSTM

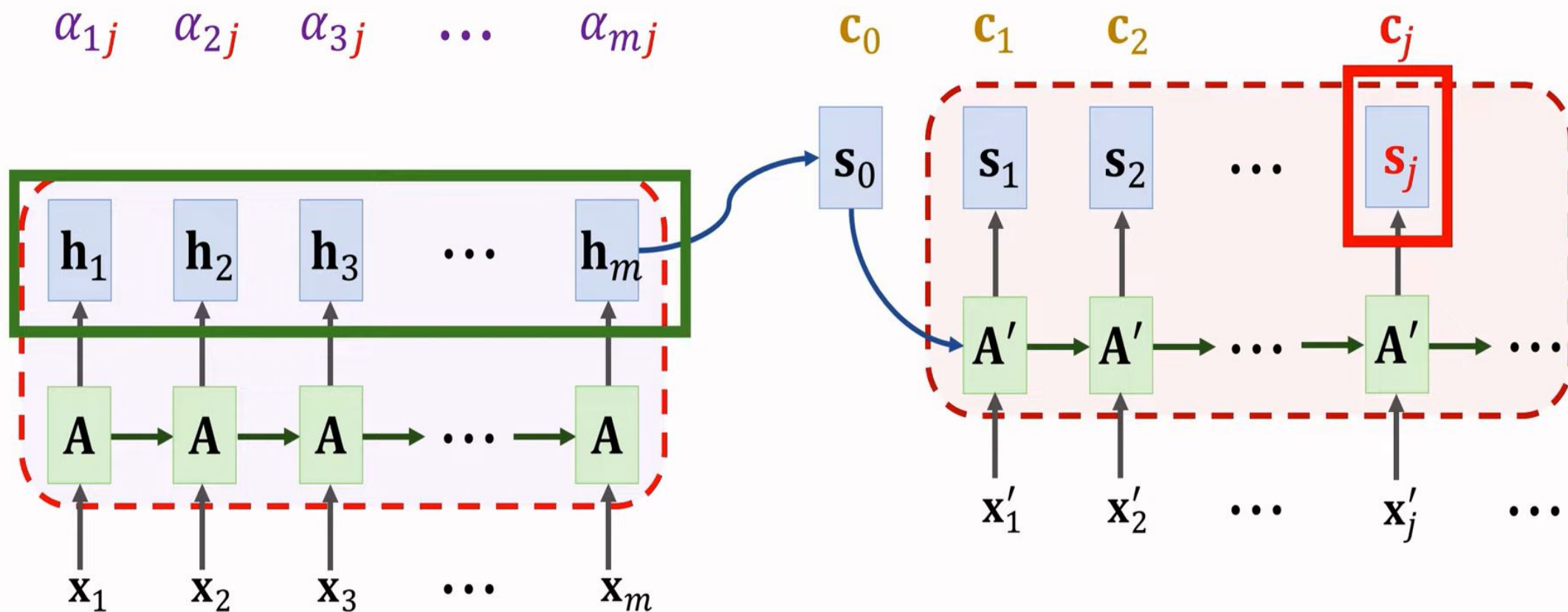


# Encoder & Decoder



# Attention for Seq2Seq Model

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



# Attention for Seq2Seq Model

Query:  $\mathbf{q}_{:j} = \mathbf{W}_Q \mathbf{s}_j$ ,      Key:  $\mathbf{k}_{:i} = \mathbf{W}_K \mathbf{h}_i$ ,      Value:  $\mathbf{v}_{:i} = \mathbf{W}_V \mathbf{h}_i$ .

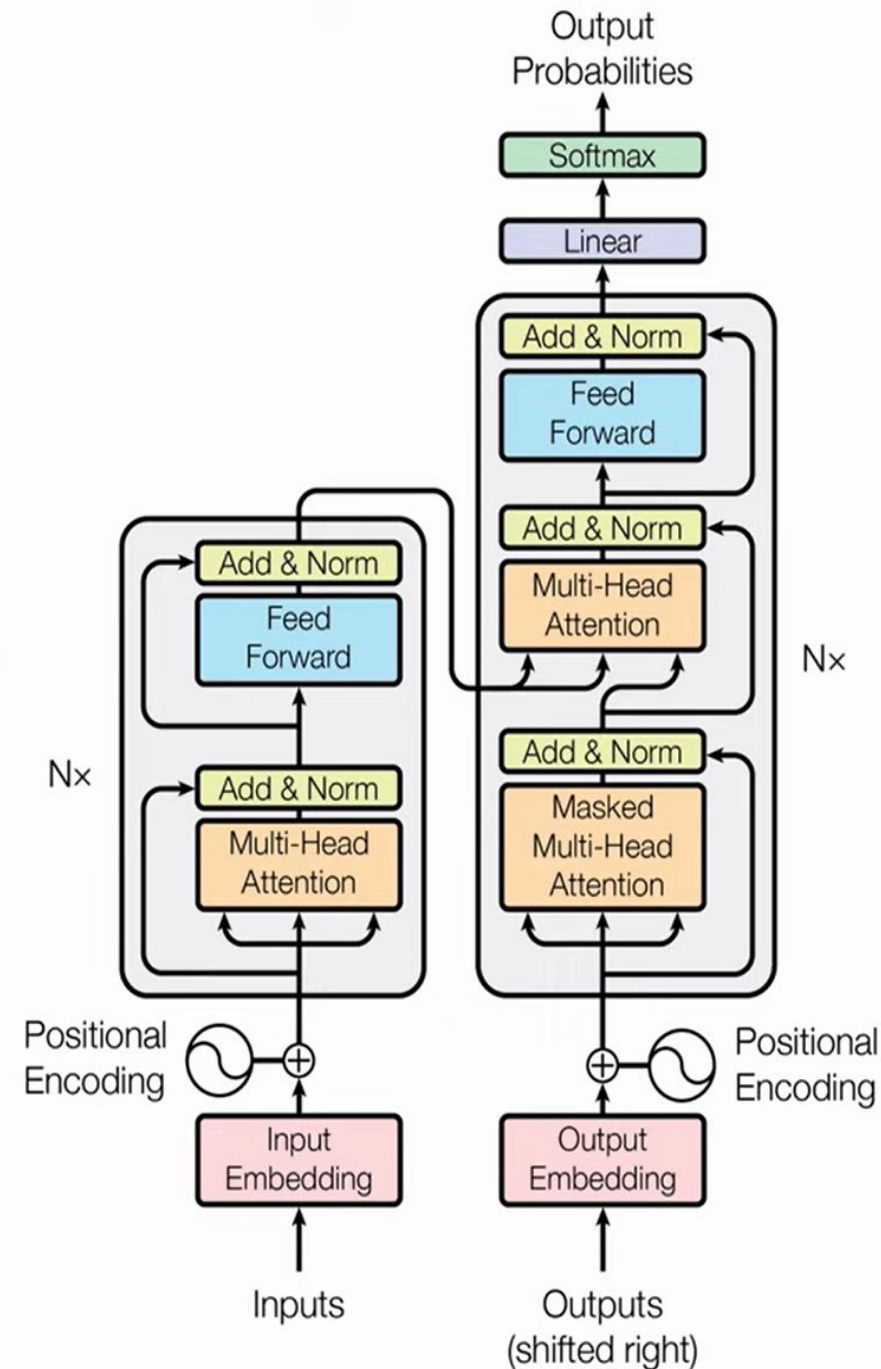
Weights:  $\alpha_{:j} = \text{Softmax}(\mathbf{K}^T \mathbf{q}_{:j}) \in \mathbb{R}^m$ .

Context vector:  $\mathbf{c}_j = \alpha_{1j} \mathbf{v}_{:1} + \dots + \alpha_{mj} \mathbf{v}_{:m}$ .

- Query:  $\mathbf{q}_{:j} = \mathbf{W}_Q \mathbf{s}_j$ . (To match others.)
- Key:  $\mathbf{k}_{:i} = \mathbf{W}_K \mathbf{h}_i$ . (To be matched.)
- Value:  $\mathbf{v}_{:i} = \mathbf{W}_V \mathbf{h}_i$ . (To be weighted averaged.)

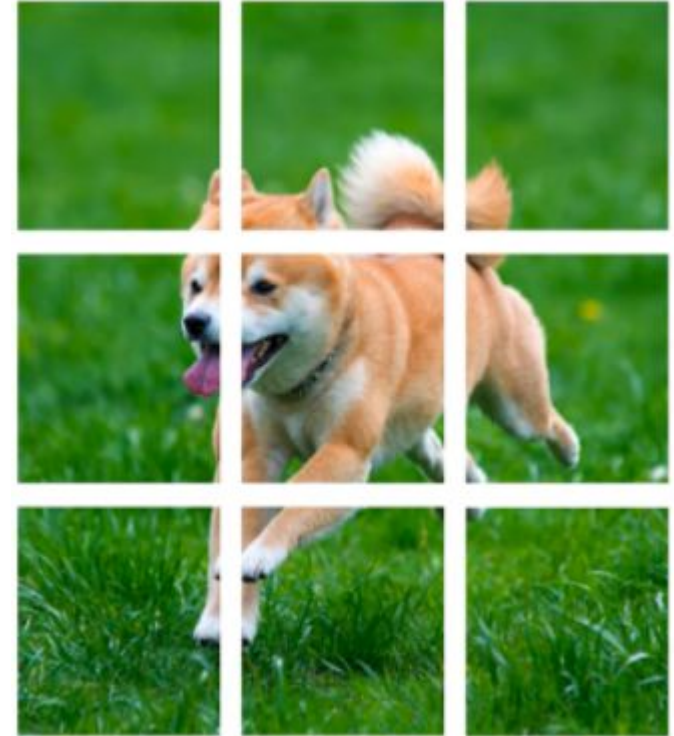
# Transformer Model

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



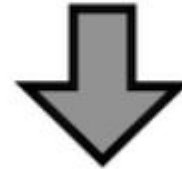
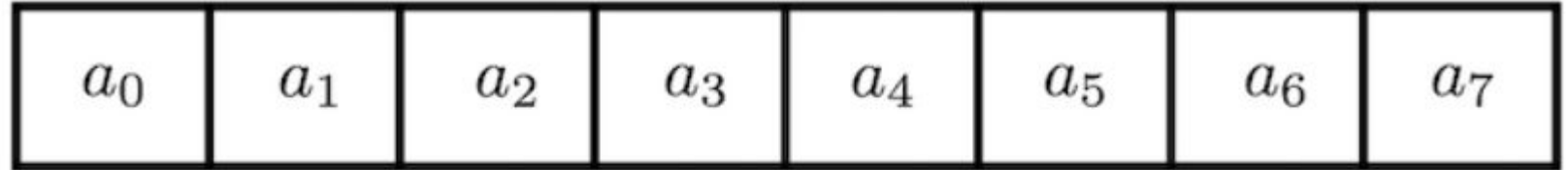
# Split Image into Patches

- May also be overlapped
- Number of patches can be calculated:
- $N = H * W / P * P$  ( $P$  = Patch size)



# Positional encoding

- Just count



- Consider 500 length sequence ?
- want encoding to be a number in the range  $[-1,1]$



- Normalize the above encoding

$a_0$	$a_1$	$a_2$	$a_3$	$a_4$	$a_5$	$a_6$	$a_7$
-------	-------	-------	-------	-------	-------	-------	-------



0	0.14	0.29	0.43	0.57	0.71	0.86	1
---	------	------	------	------	------	------	---

- Consider variable sequence length ?
- want encoding to be consistent throughout

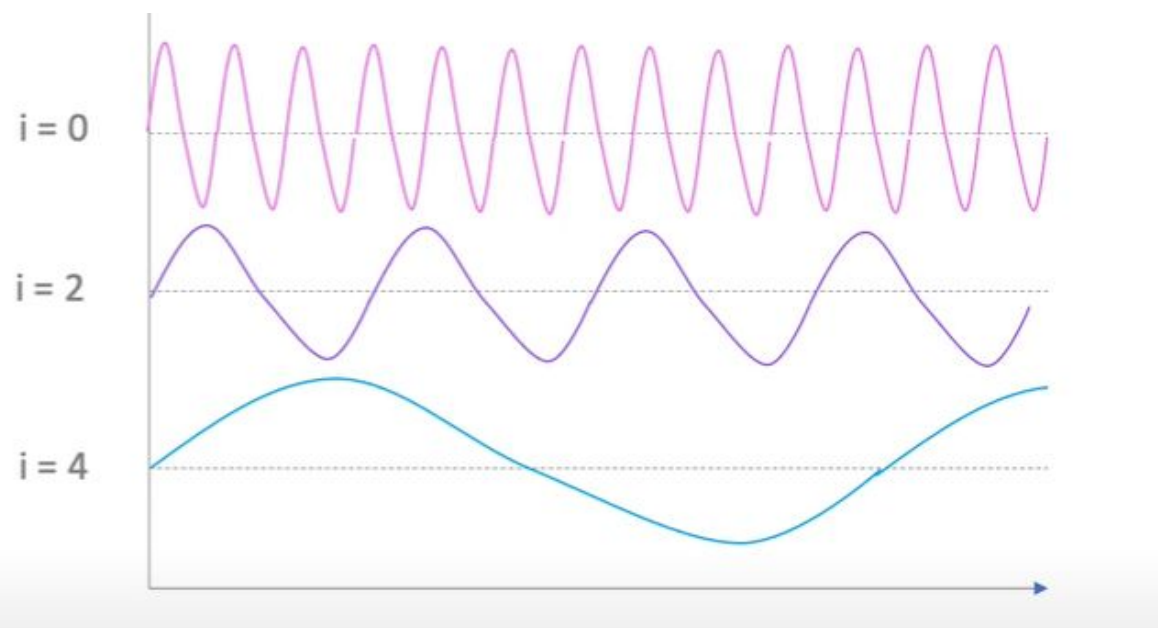
# Using frequencies

- Used sine and cosine curves
- Range => [-1,1] and same output on same input
- $d$  = embedding dimension,  $i$  = entries of positional encoding vector

$$\vec{p}_t^{(i)} = f(t)^{(i)} := \begin{cases} \sin(\omega_k \cdot t), & \text{if } i = 2k \\ \cos(\omega_k \cdot t), & \text{if } i = 2k + 1 \end{cases}$$

where

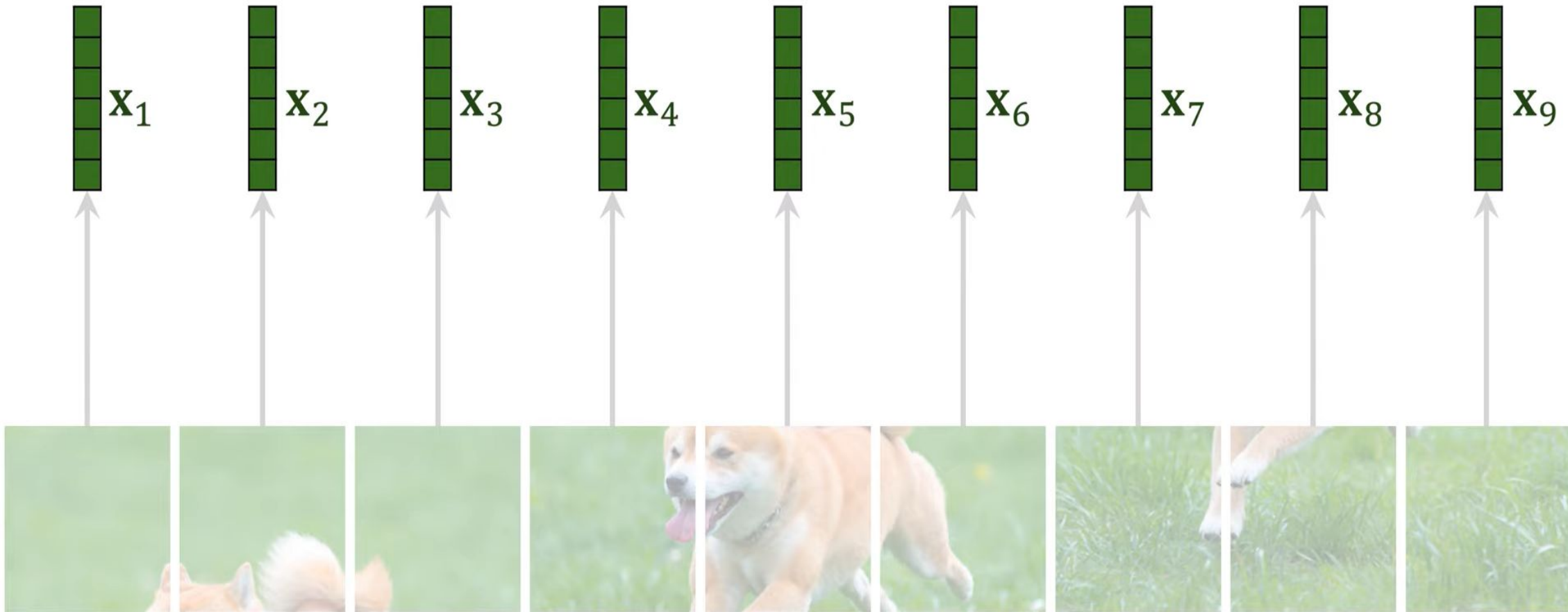
$$\omega_k = \frac{1}{10000^{2k/d}}$$

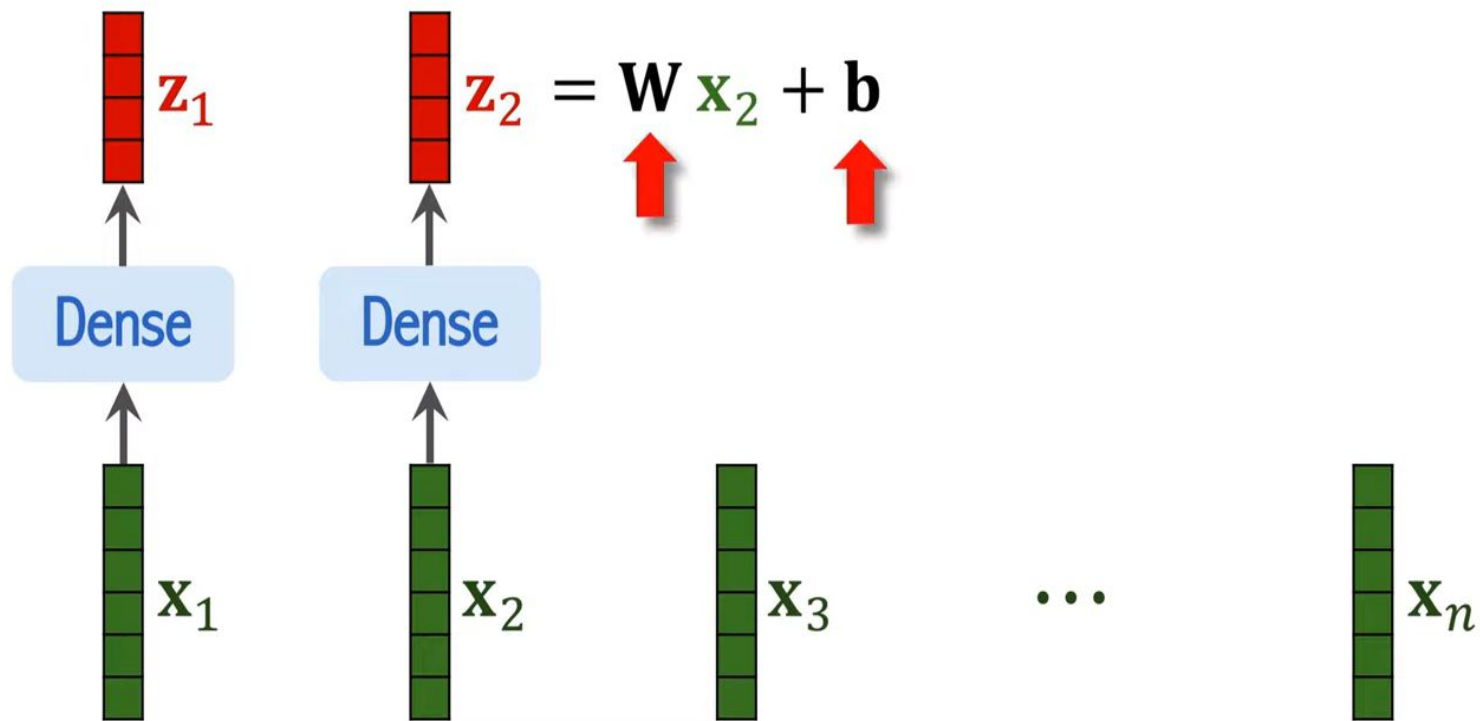


$$\vec{p}_t = \begin{bmatrix} \sin(\omega_1.t) \\ \cos(\omega_1.t) \\ \\ \sin(\omega_2.t) \\ \cos(\omega_2.t) \\ \\ \vdots \\ \\ \sin(\omega_{d/2}.t) \\ \cos(\omega_{d/2}.t) \end{bmatrix}_{d \times 1}$$

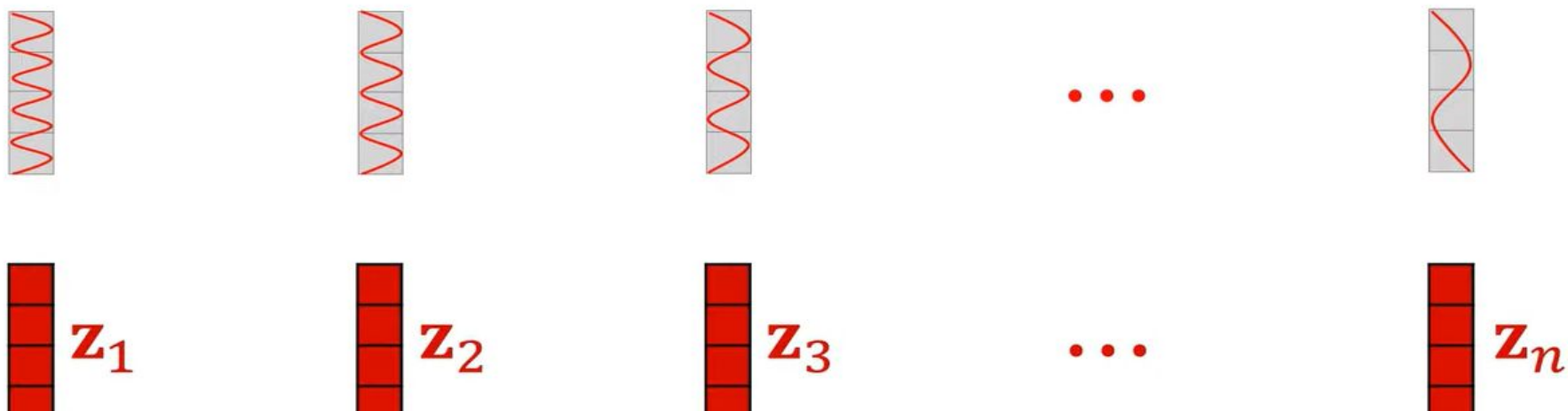
# Vectorization

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

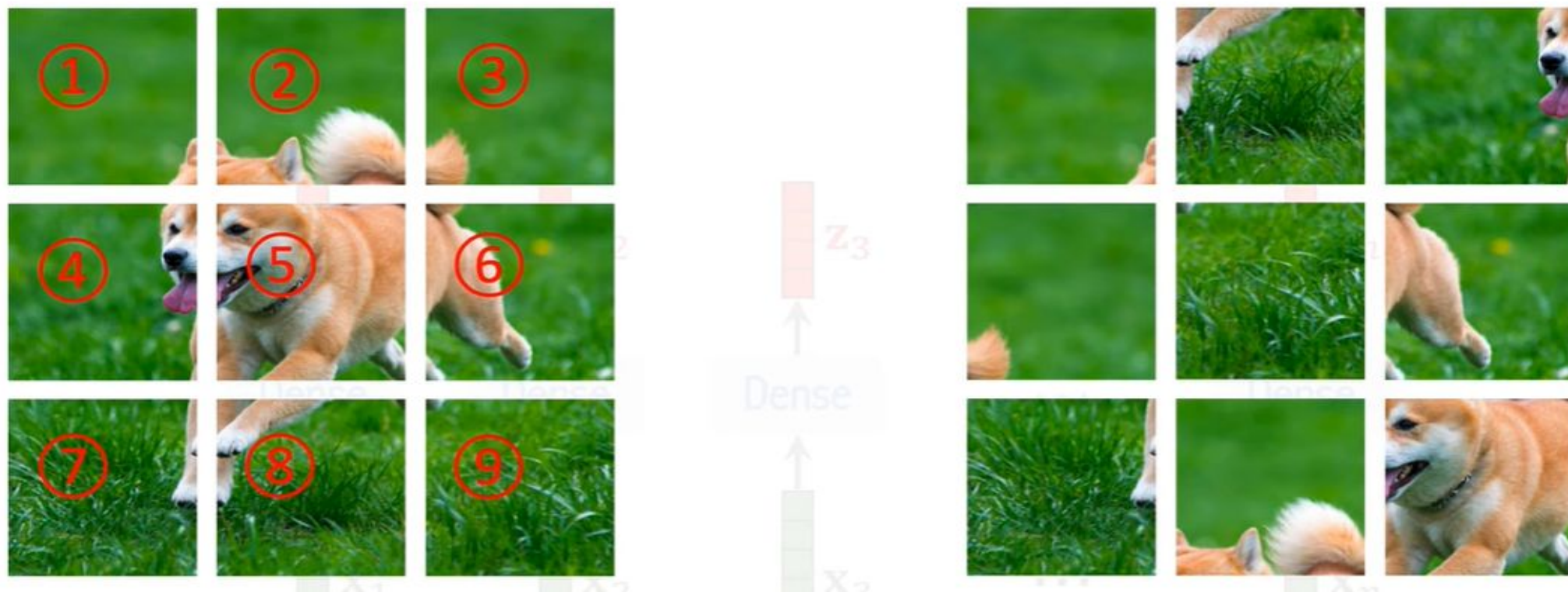




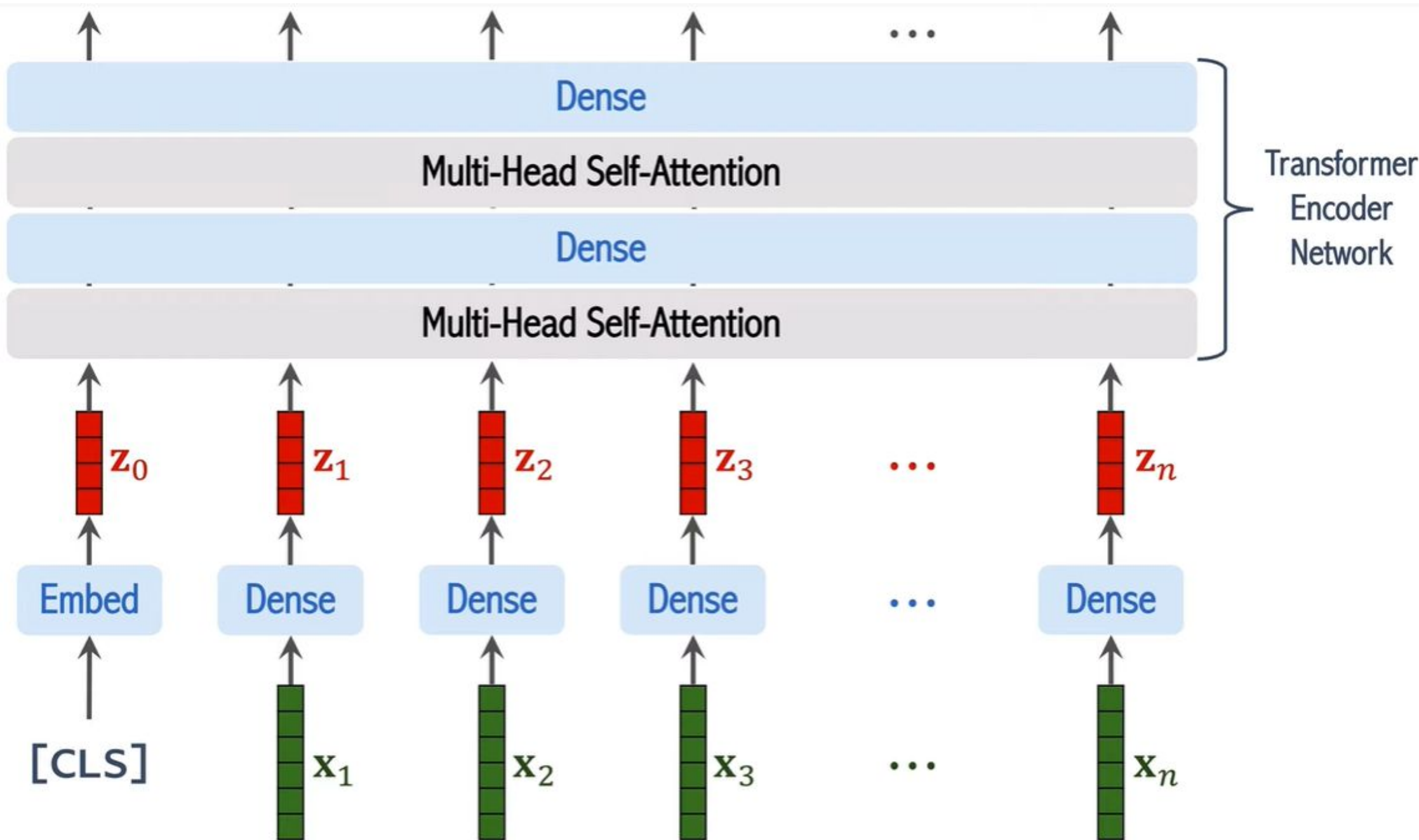
Add positional encoding vectors to  $z_1, z_2, \dots, z_n$ .

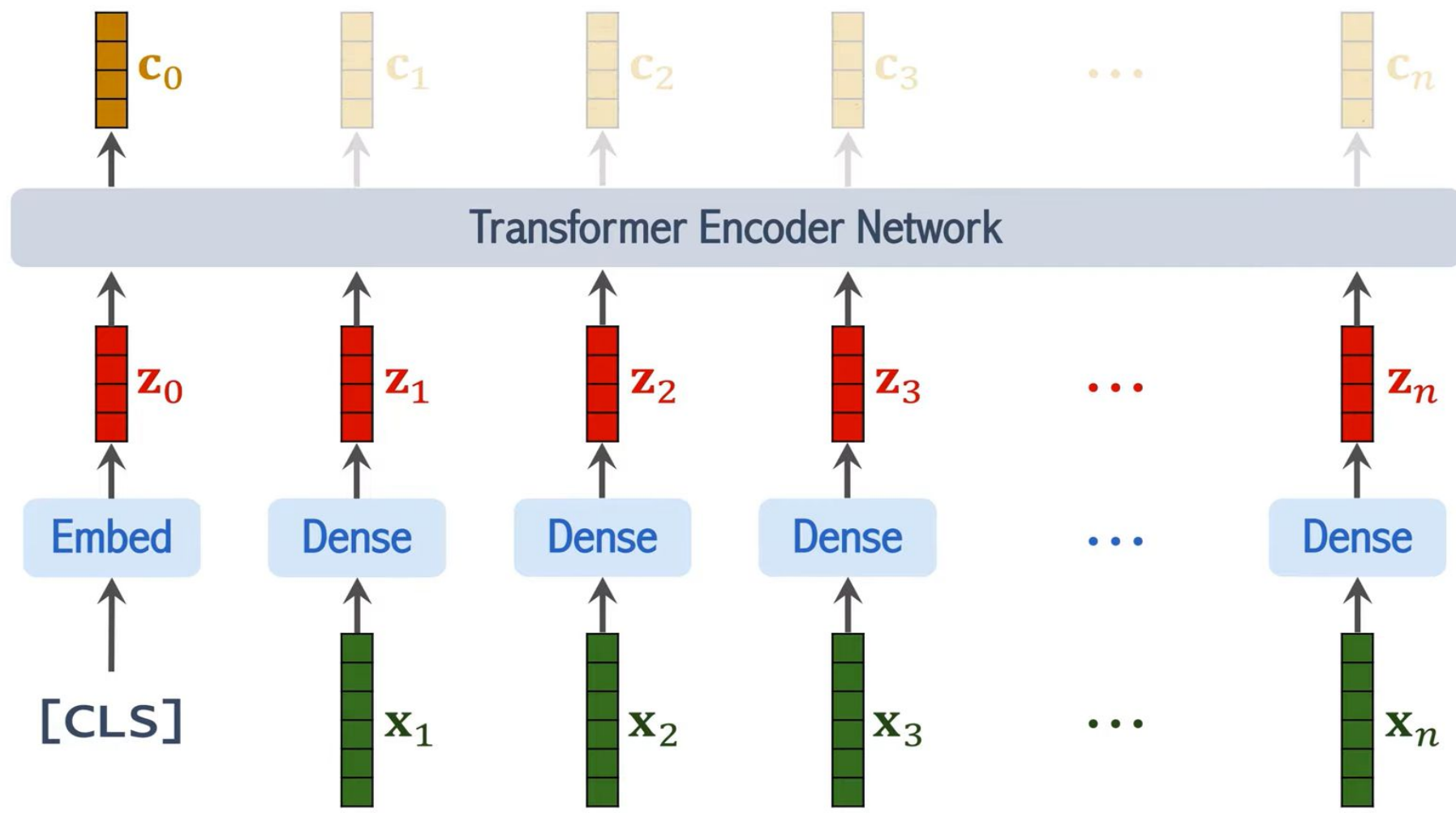
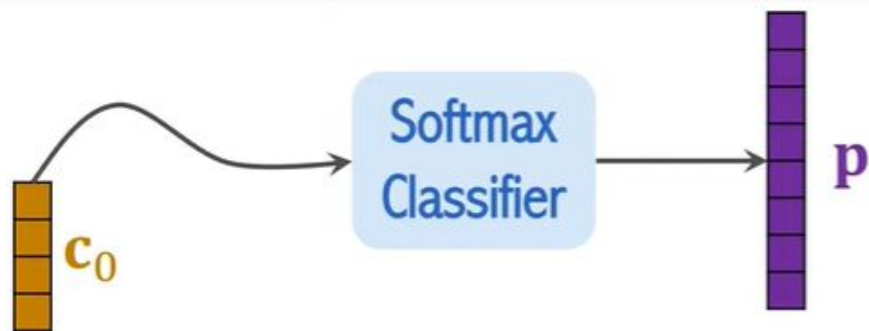


Add positional encoding vectors to  $\mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_n$ . (Why?)



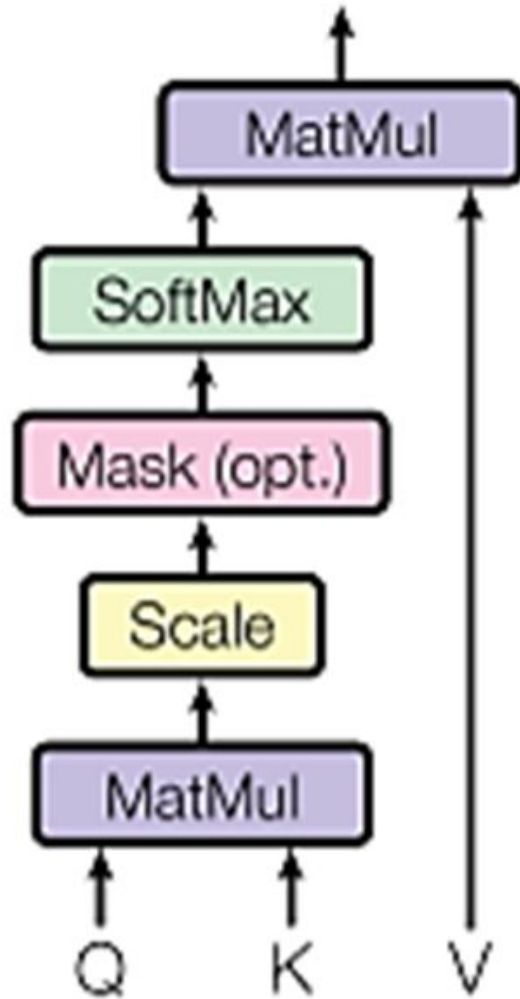








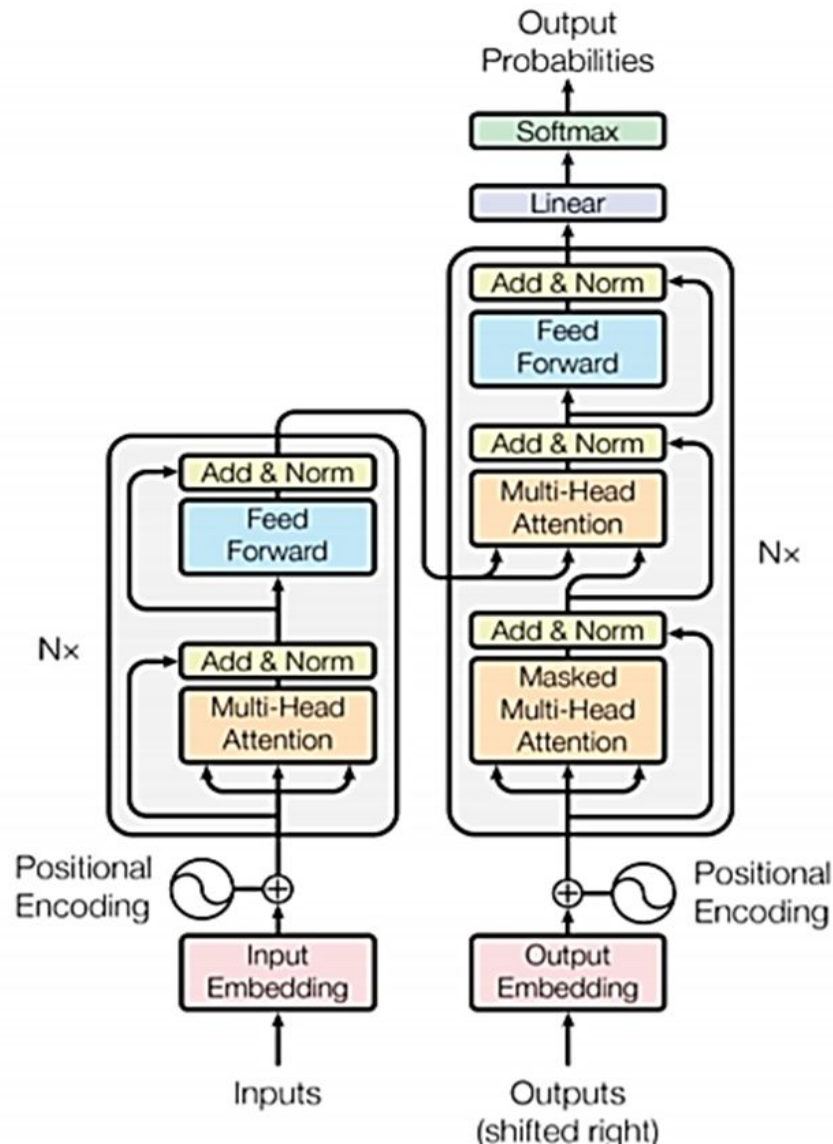
# Self Attention Mechanism



- Multiplying each of the encoder input vectors with three weights matrices ( $W(Q)$ ,  $W(K)$ ,  $W(V)$ ) that we trained during the training process.
- Multiply the Query vector of the current input with the key vectors from other inputs.
- Divide the score by square root of dimensions of the key vector ( $d_k$ ).
- Softmax function on all self-attention scores we calculated wrt the query word
- We multiply the value vector.
- Sum up the weighted value vectors.

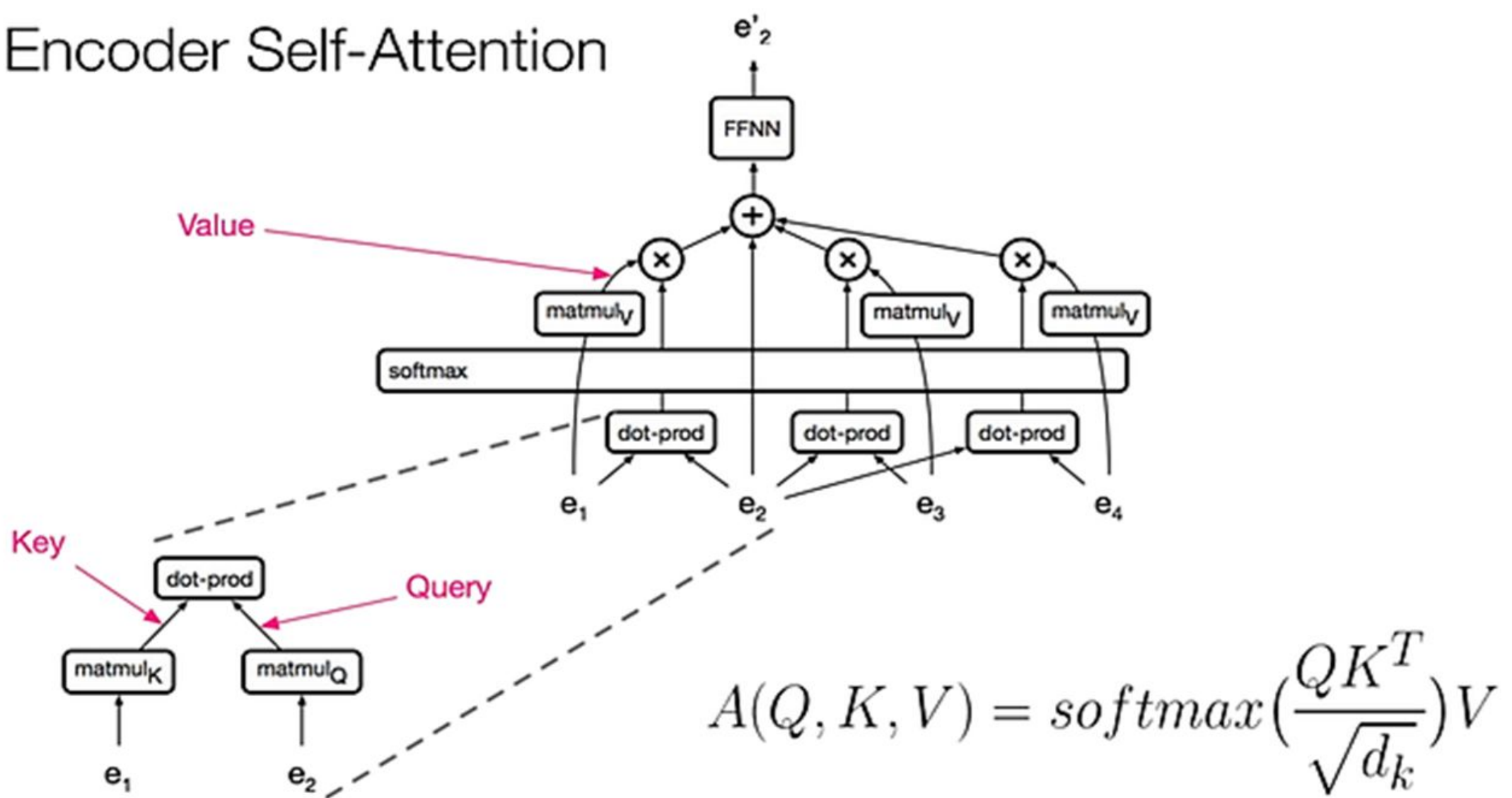
$$Attention(Q, K, V) = softmax\left(\frac{QK^T}{\sqrt{d_k}}\right) V$$

# Encoder Decoder Attention



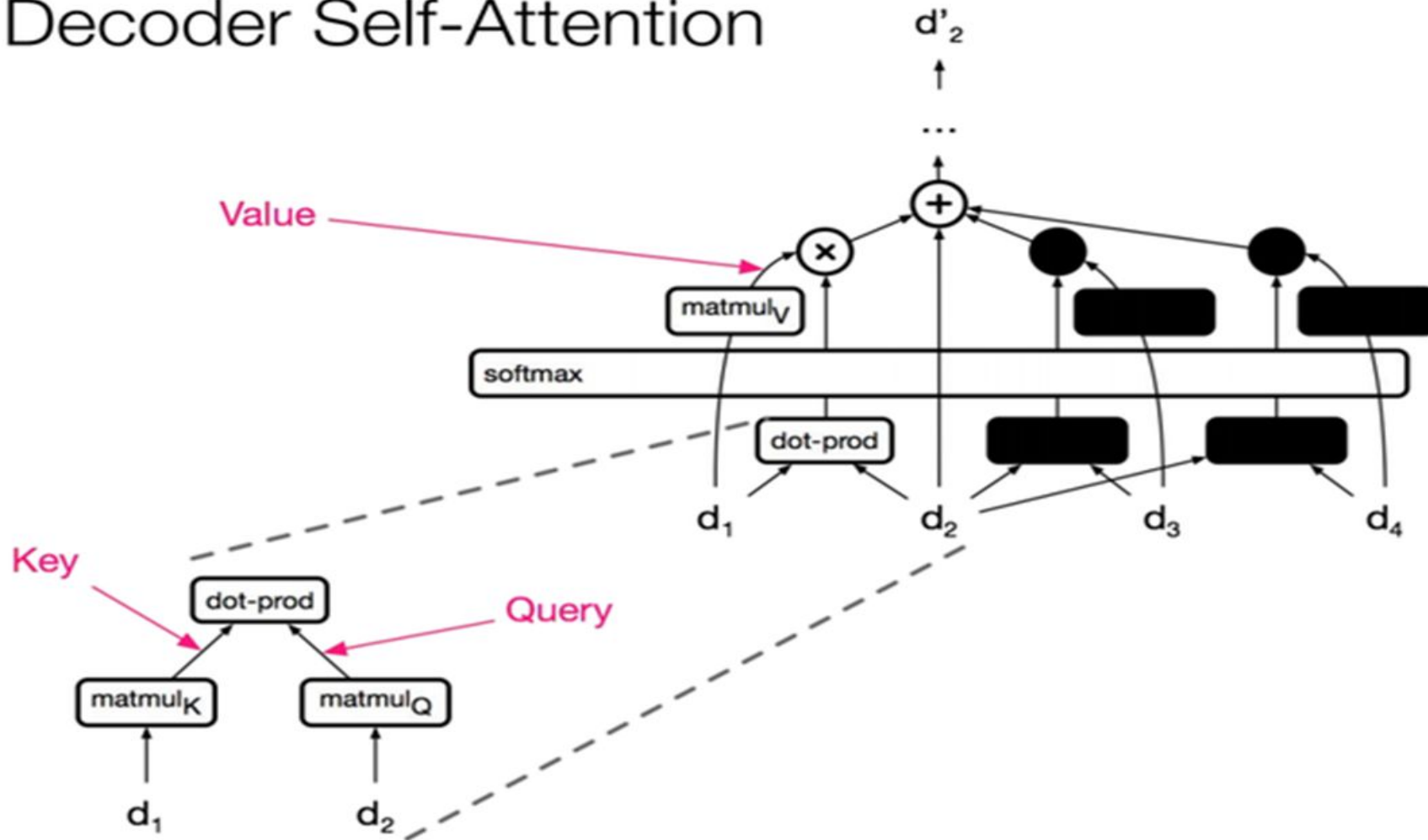
In encoder-decoder attention layer the queries come from the previous decoder layer while the keys and values come from the encoder output. This allows each position in the decoder to give attention to all the positions of the input sequence.

# Encoder Self-Attention



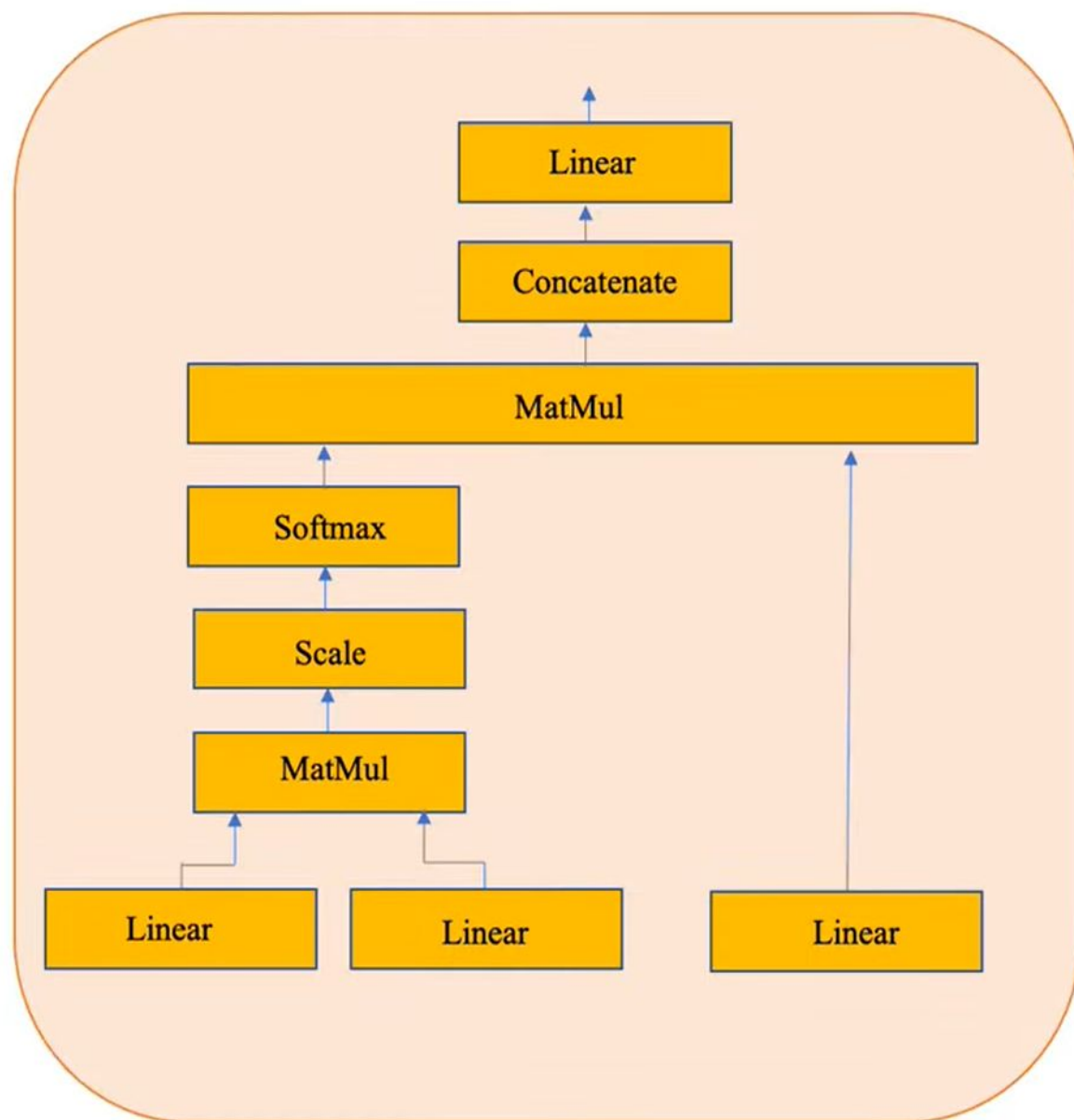
Encoder Self-Attention layer receives key, value, and query input from the output of the previous encoder layer. Each position in the encoder can get attention score from every position in the previous encoder layer.

# Decoder Self-Attention



Decoder Self-Attention is similar to self-attention in encoder where all queries, keys, and values come from the previous layer. The self-attention decoder allows each position to attend each position up to and including that position. The future values are masked with  $(-\text{Inf})$ . This is known as masked-self attention.

## Multi-Head Attention



## Linear Layer

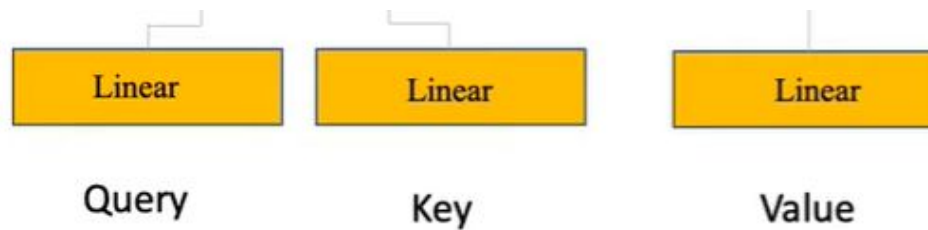
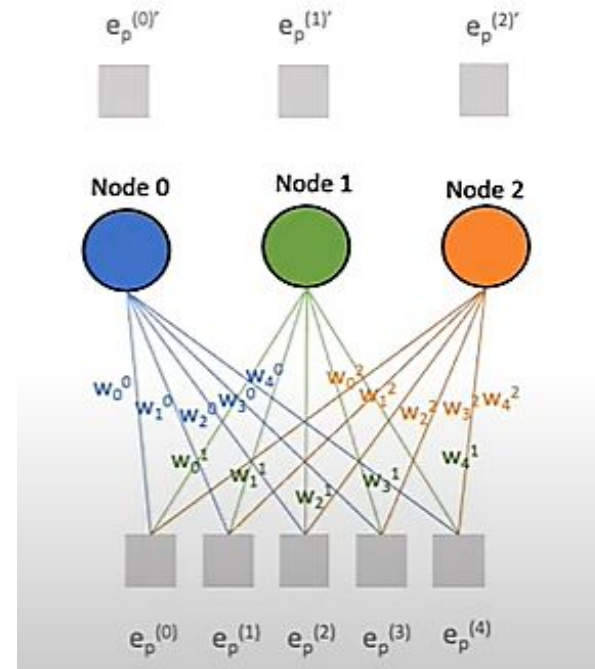


i) Mapping inputs onto the outputs

ii) Changing matrix/vector dimensions

## Linear Layer

$$W_0^0 = 0.6$$

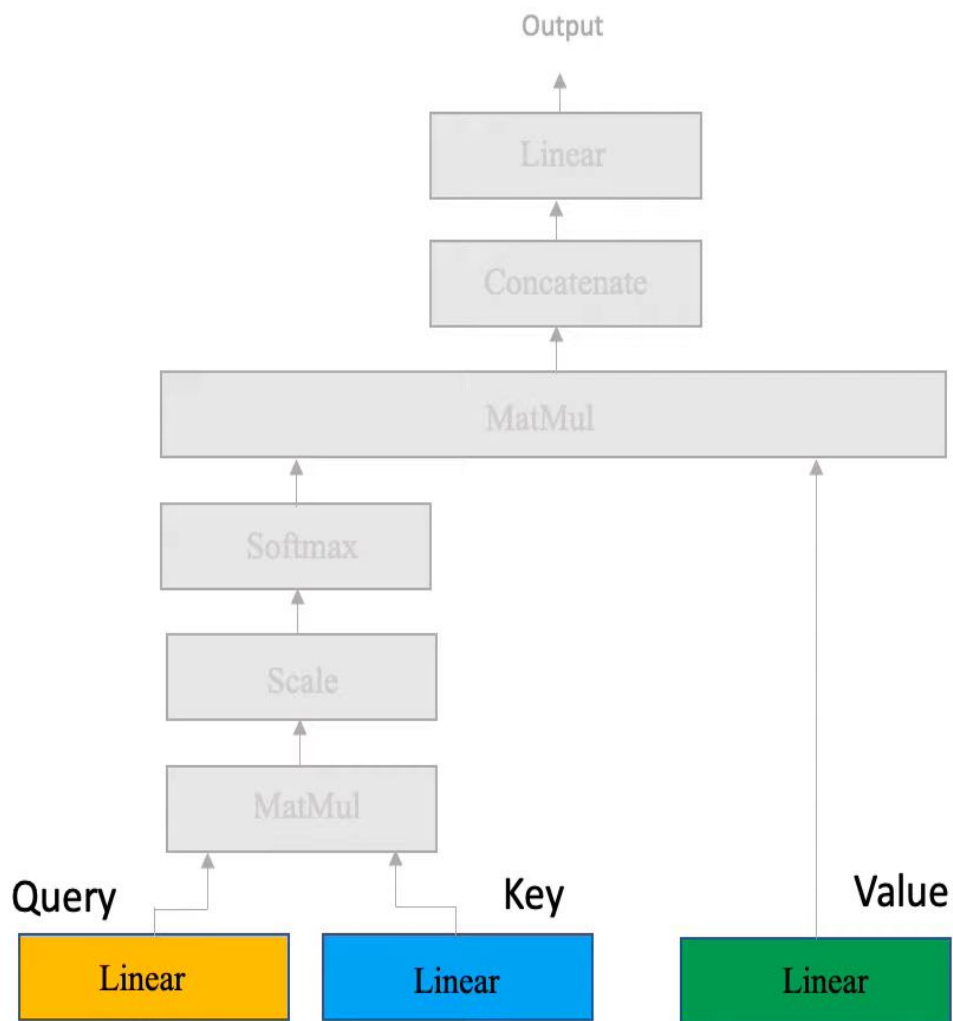


# Cosine Similarity

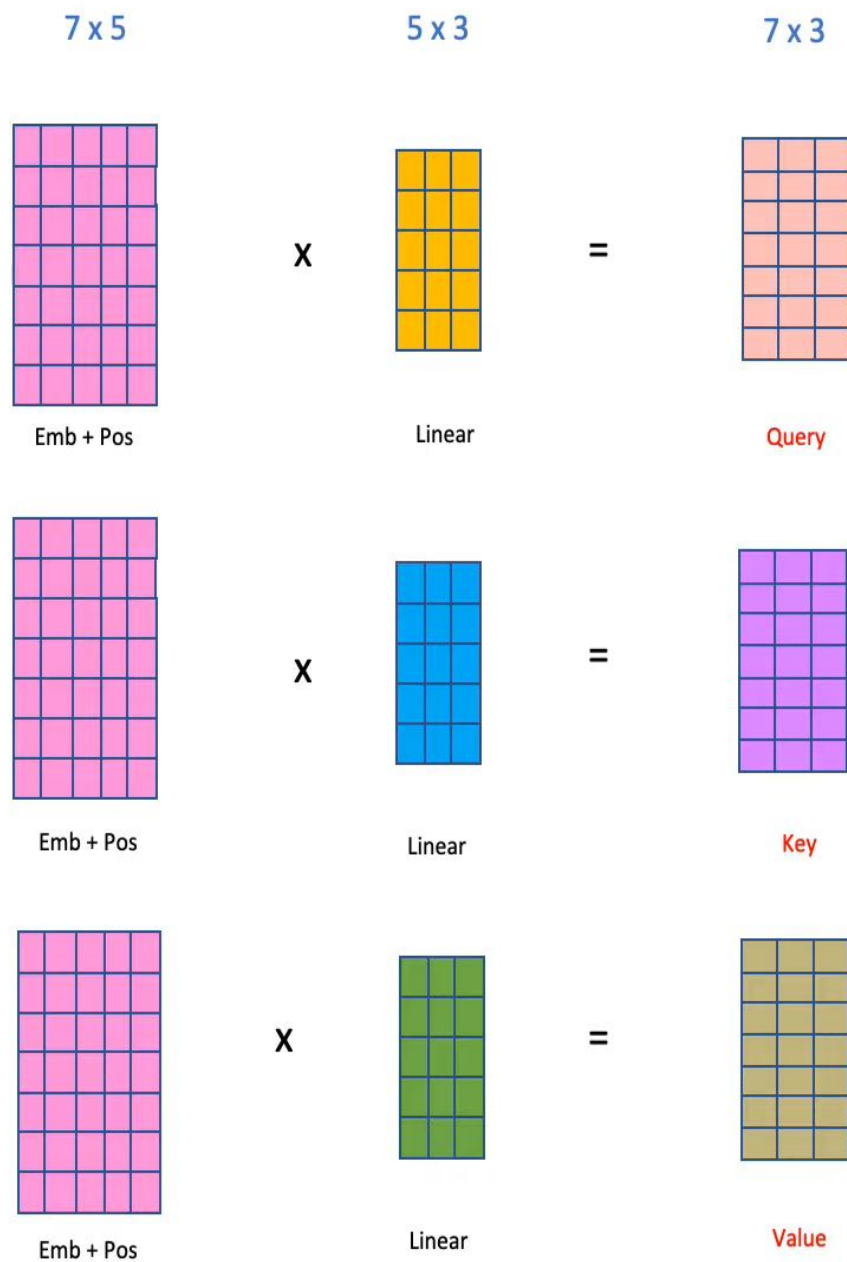
$$\text{Cos}(A, B) = \frac{A \cdot B}{|A| |B|}$$

## Similarity b/w Matrices

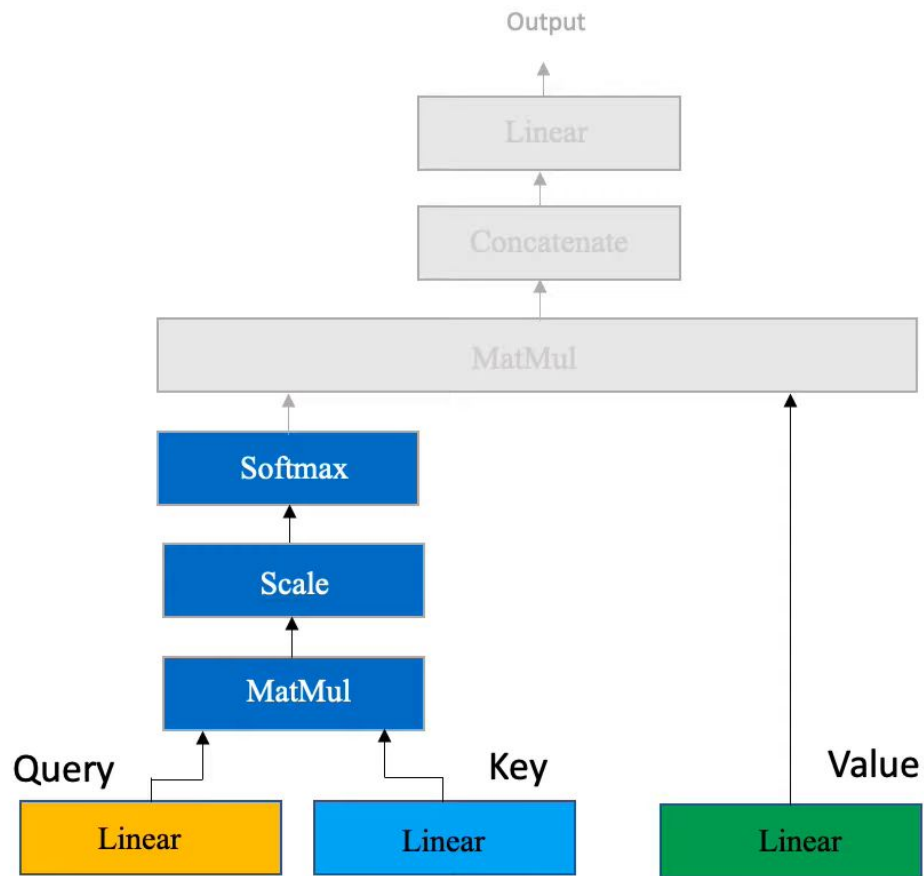
$$\text{similarity}(A, B) = \frac{A \cdot B^T}{\text{scaling}}$$



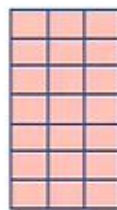
## Multi-Head Attention







7 x 3



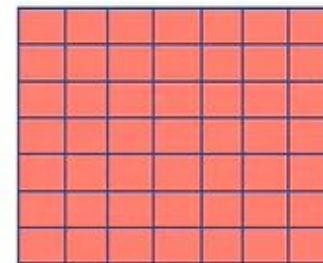
Query

3 x 7



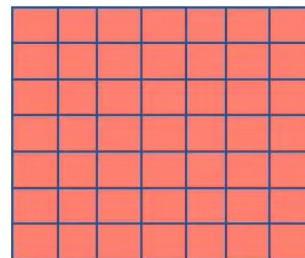
Key<sup>T</sup>

7 x 7



Attention Filter

7 x 7



Attention Filter

7 x 3



Value

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

V

# Intuition



Attention Filter



\*

Original Image

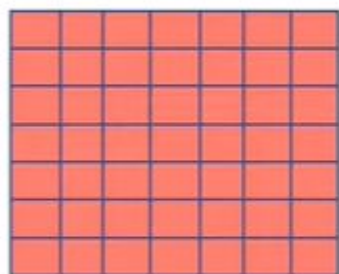


=

Filtered Image

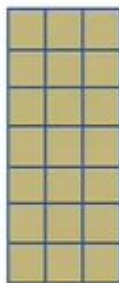


7 x 7



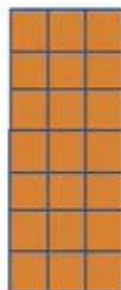
Attention Filter

7 x 3



Original Value

7 x 3



Filtered Value

\*

=

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

## Intuition for Multi-head Attention

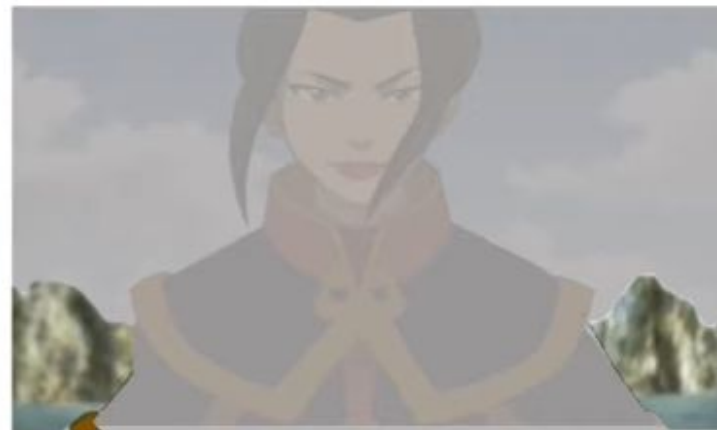
Attention Filter 1



Attention Filter 2



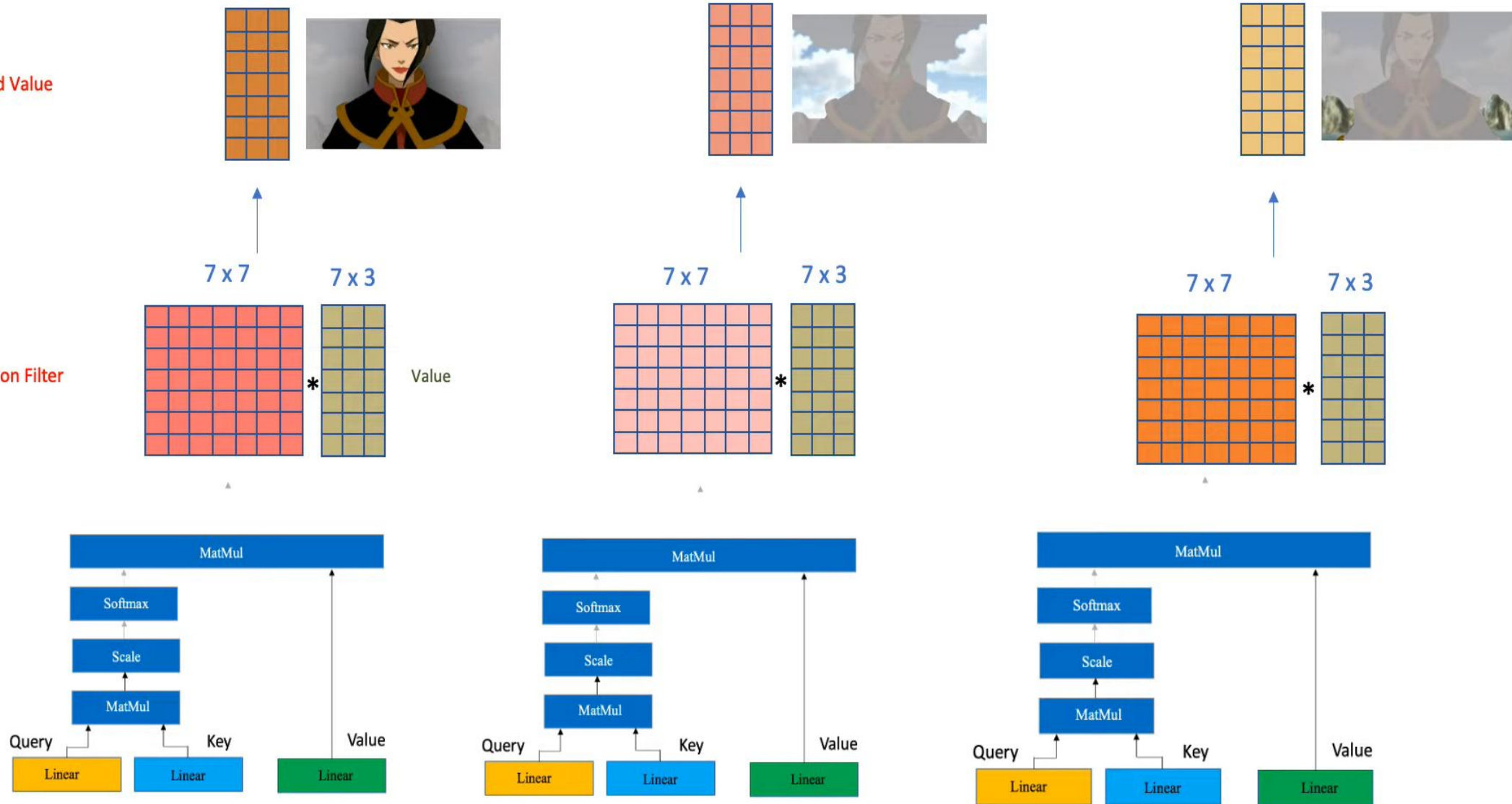
Attention Filter 3



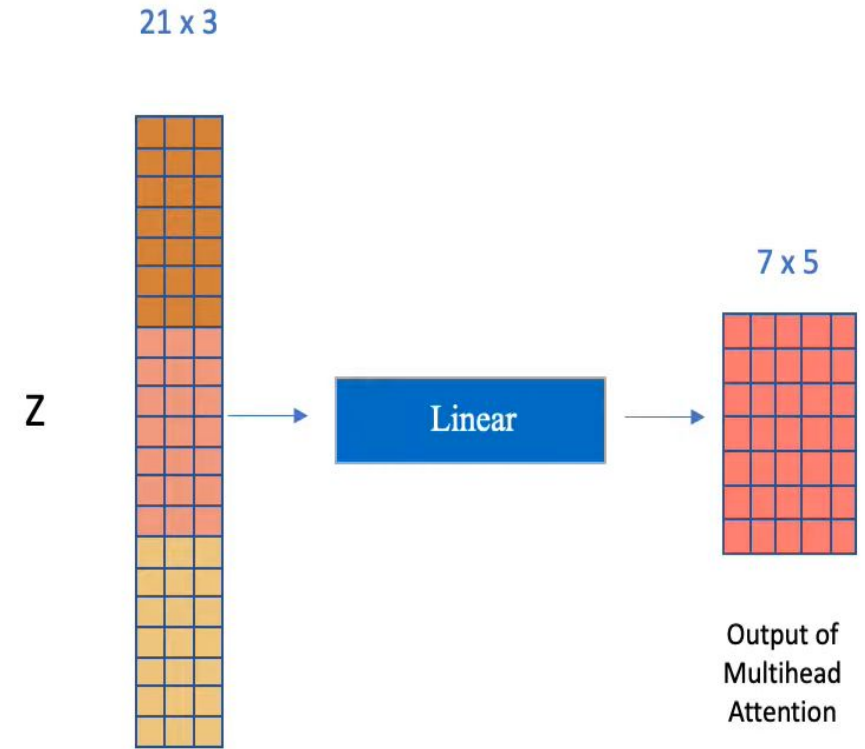
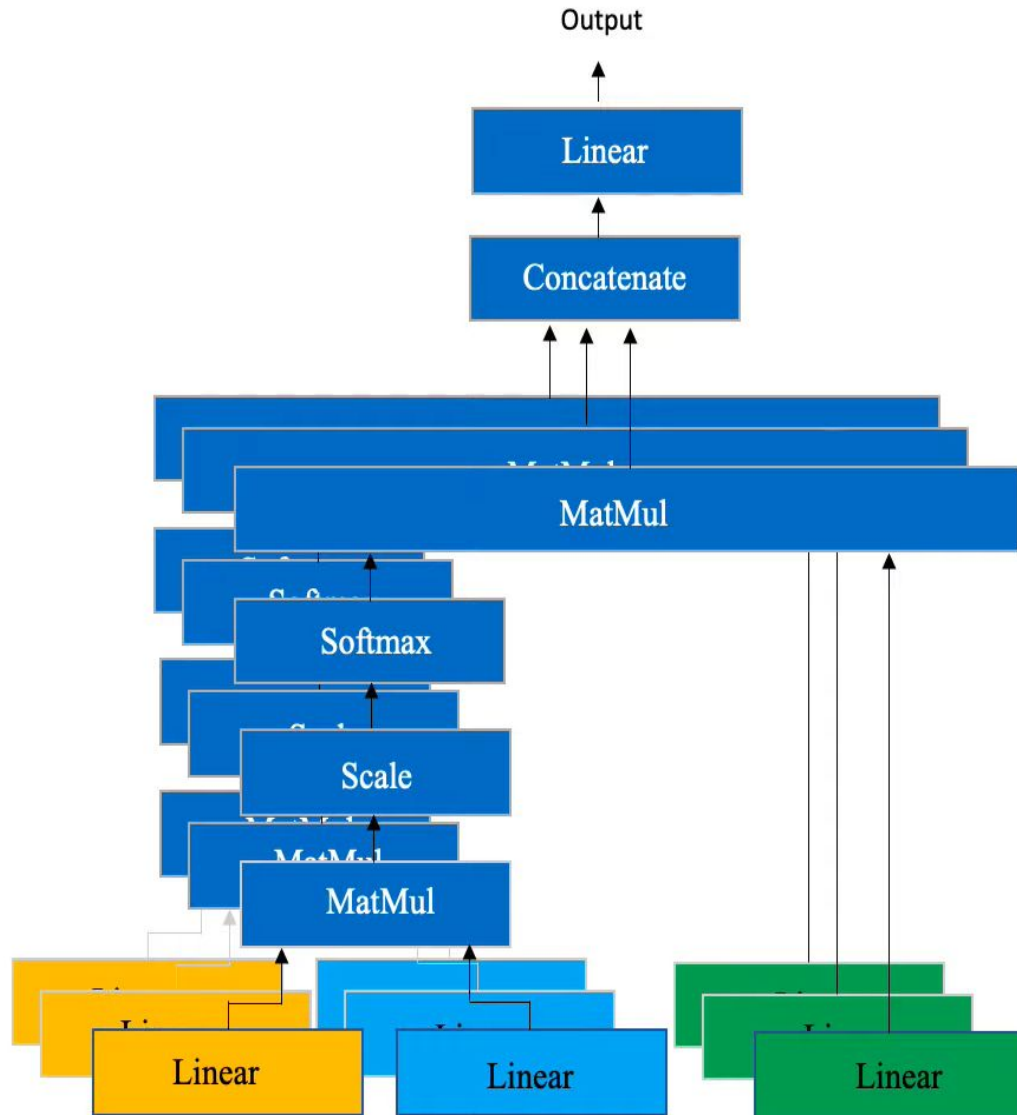
# Multi-Head Attention

Filtered Value

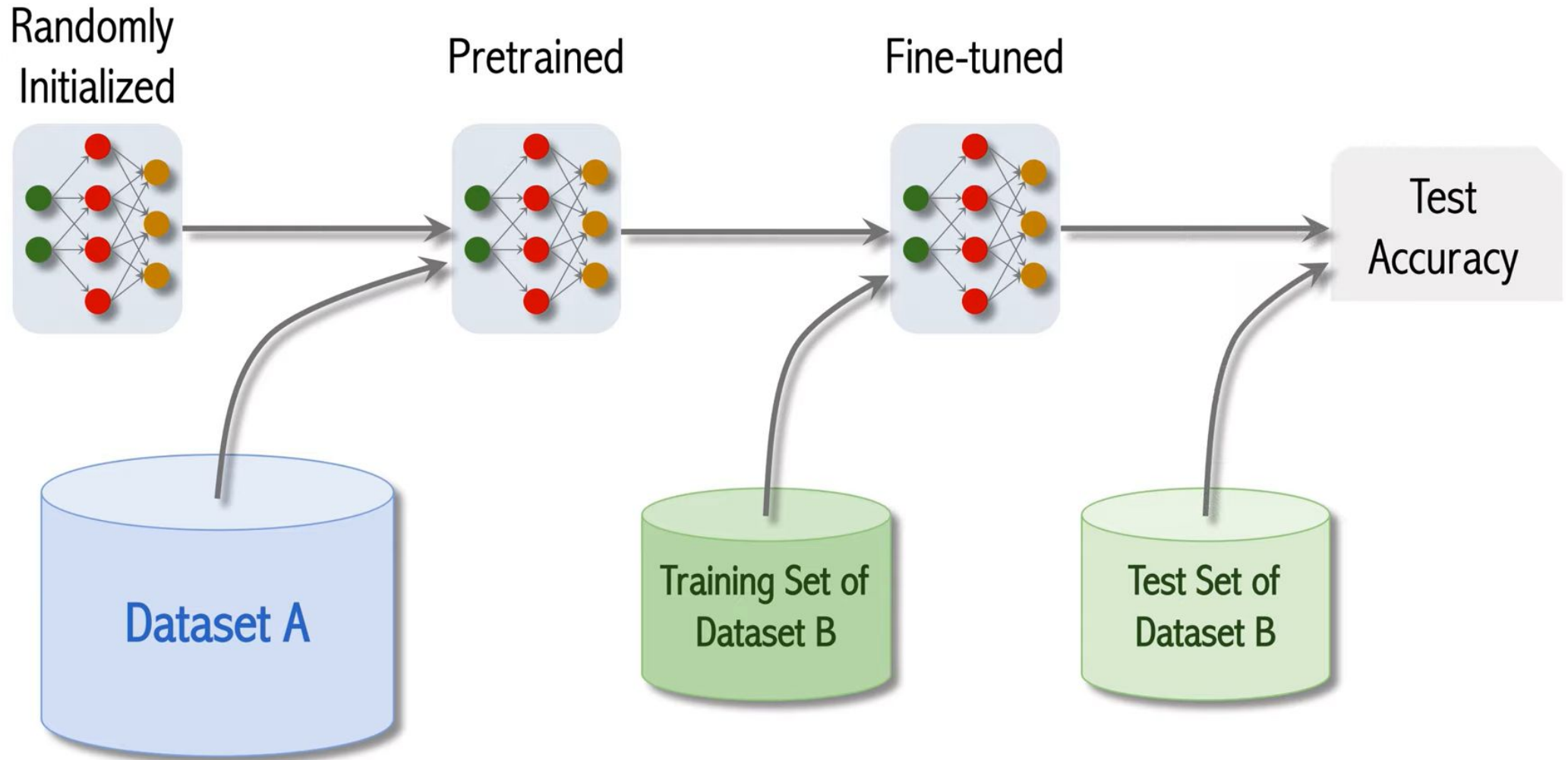
Attention Filter



## Multi-Head Attention



# Training of Transformer Model





# 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

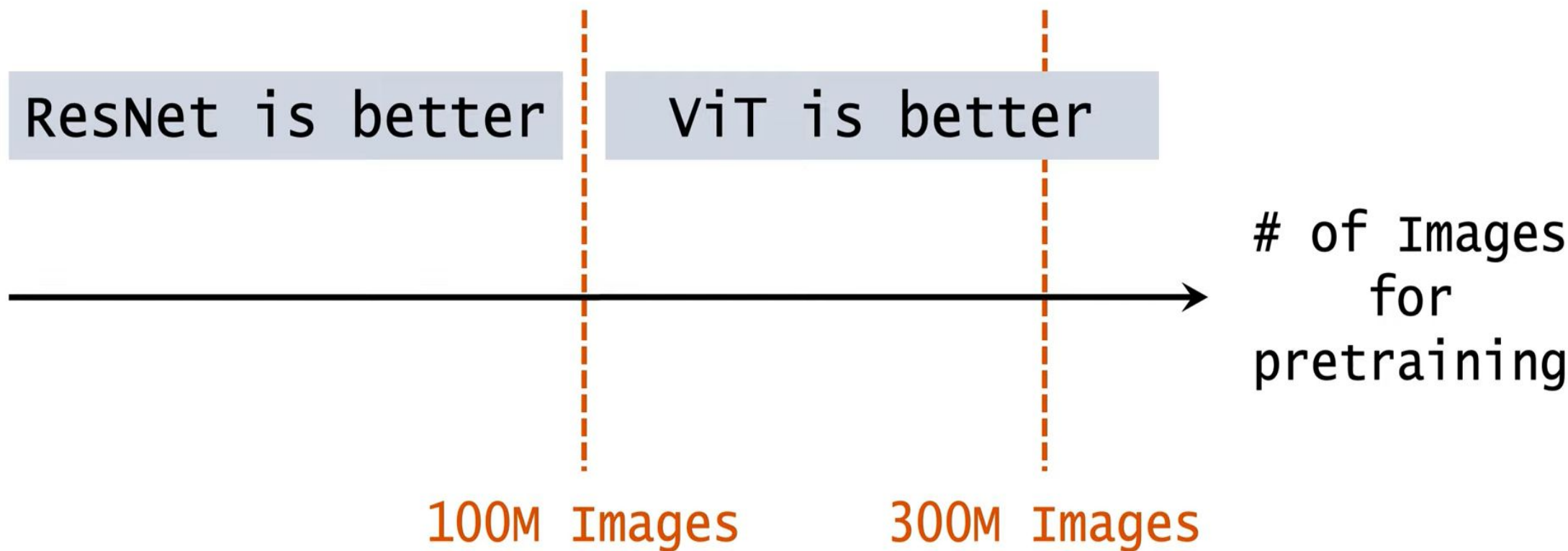
ResNet is better

ViT is better

# of Images  
for  
pretraining

100M Images

300M Images





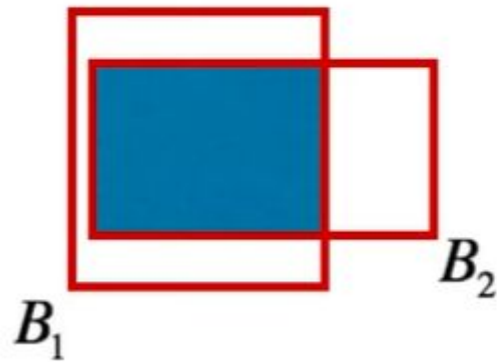
# Architectures used for object - localization

- R-CNN
- Fast R-CNN
- Faster R-CNN
- YOLO
- SSD

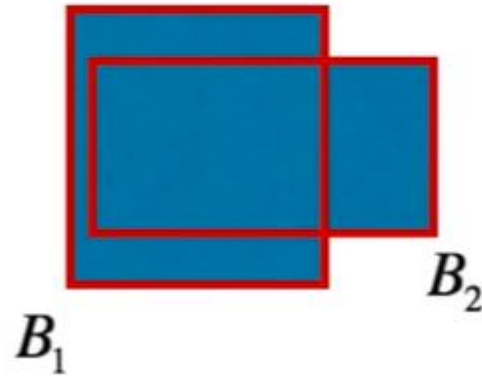
# IOU

- iou = 1 (best fit), 0 (worst fit)

Intersection



Union

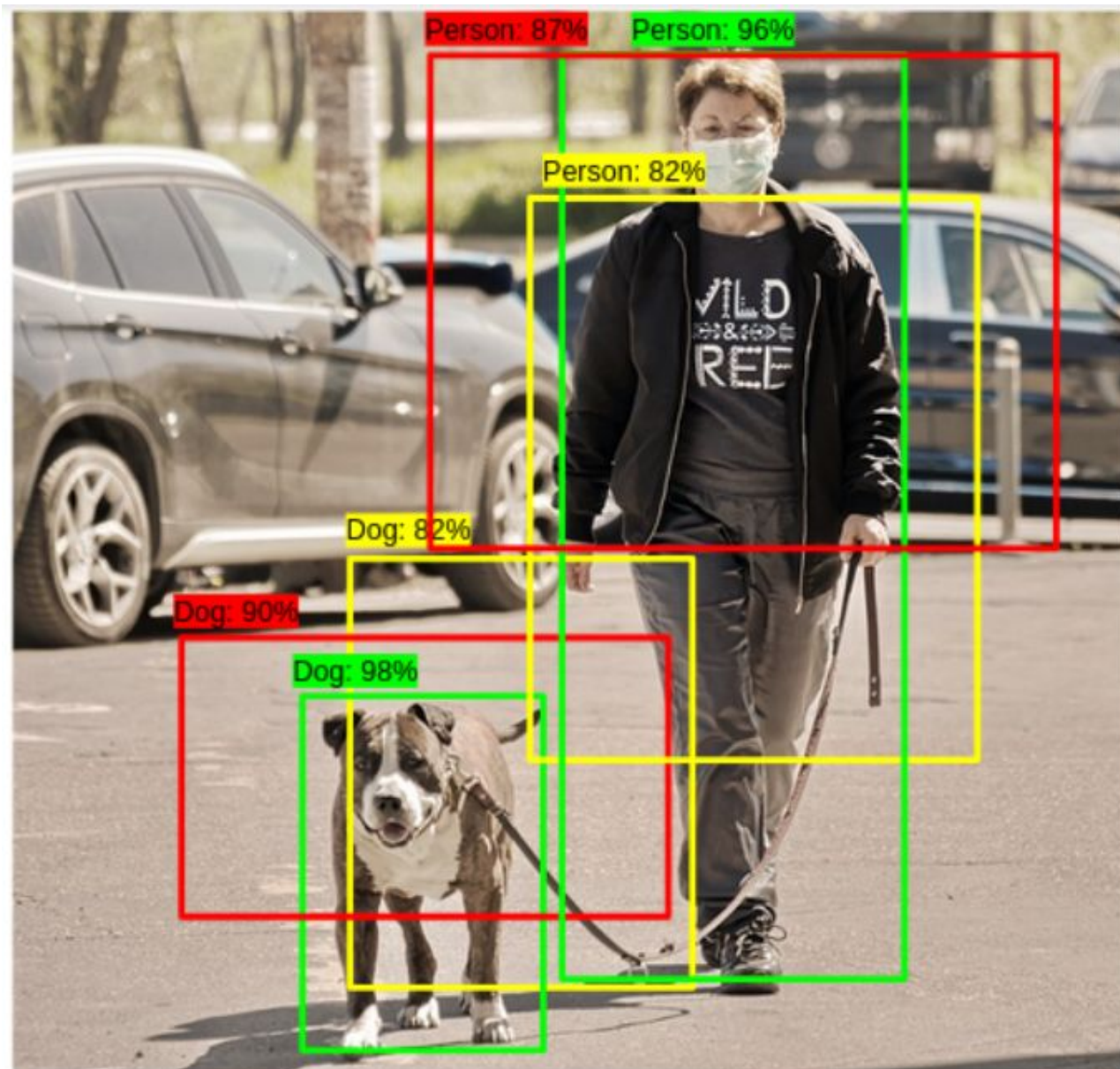


Intersection over Union

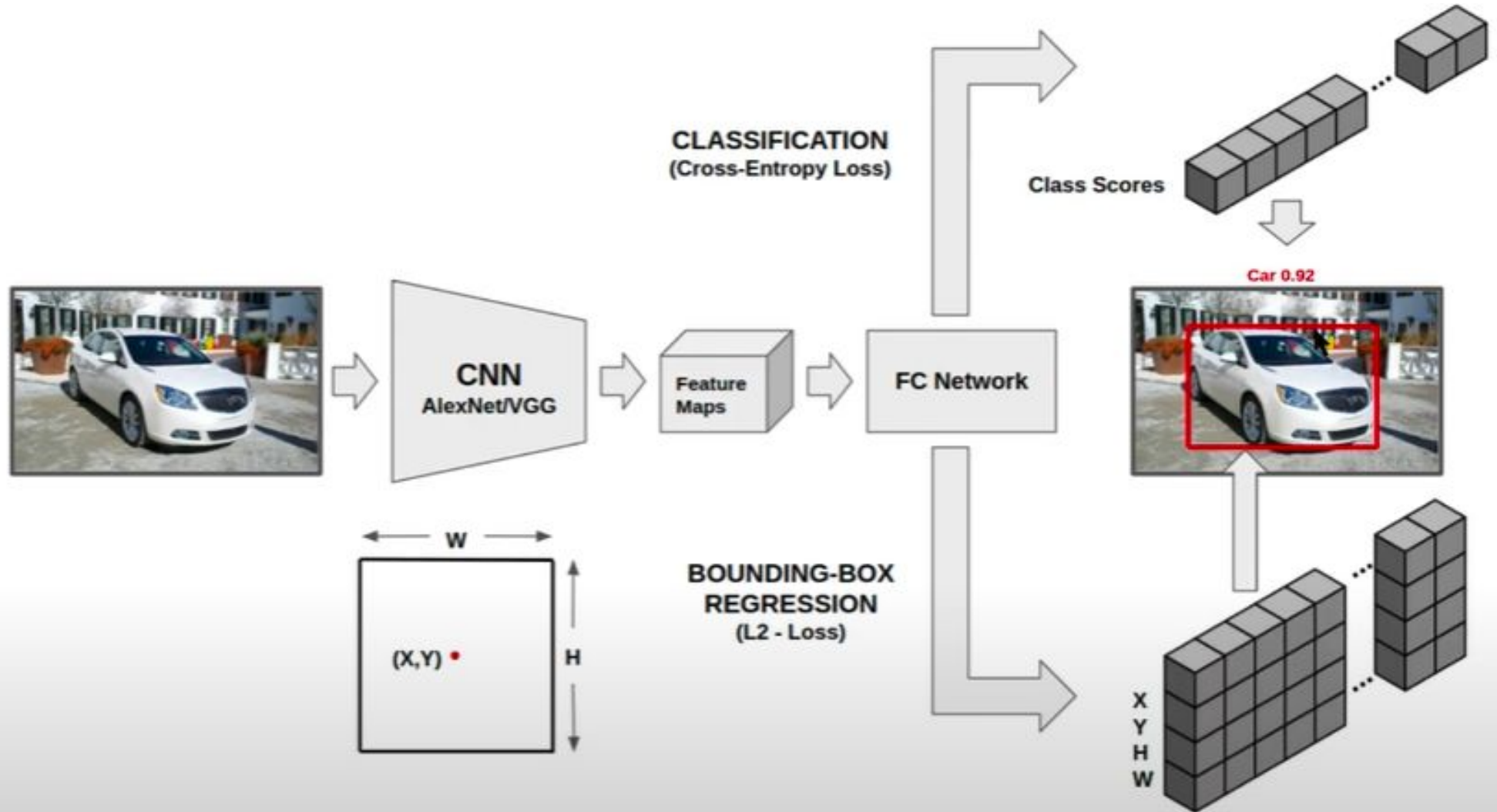
$$IoU = \frac{B_1 \cap B_2}{B_1 \cup B_2} = \frac{\text{Intersection}}{\text{Union}}$$

The diagram shows the formula  $IoU = \frac{B_1 \cap B_2}{B_1 \cup B_2} = \frac{\text{Intersection}}{\text{Union}}$ . Below the fraction, there are two small diagrams. The top diagram shows the intersection of two overlapping rectangles, with the overlapping area filled blue. The bottom diagram shows the union of two overlapping rectangles, with the entire area covered by both rectangles filled blue.

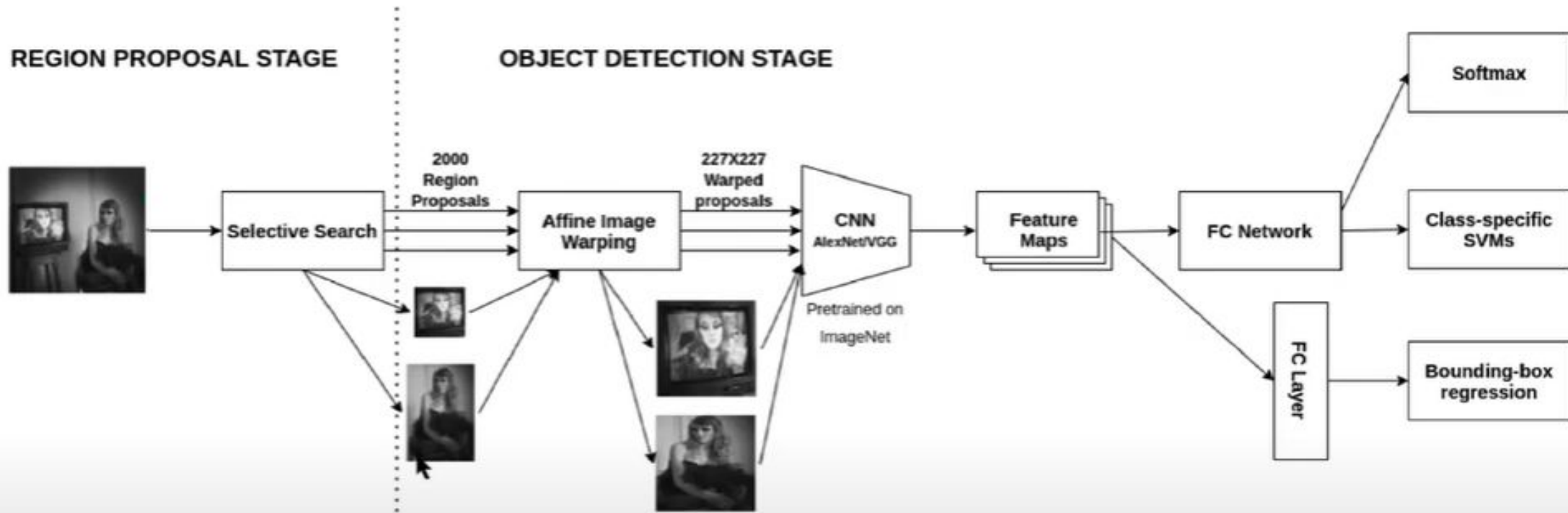
# Non Maximum Suppression



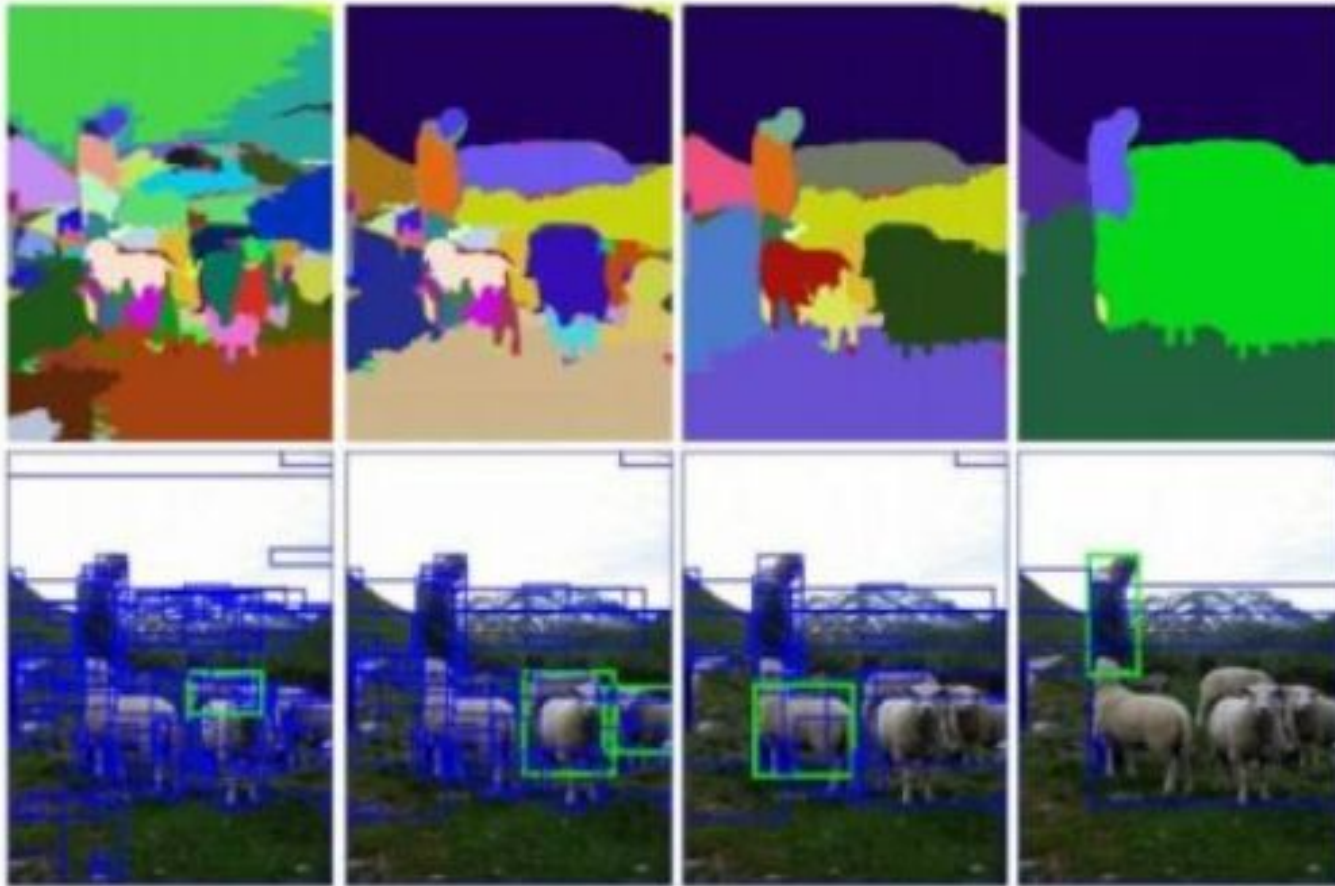
# Object detection using CNN's and Bounding box regression



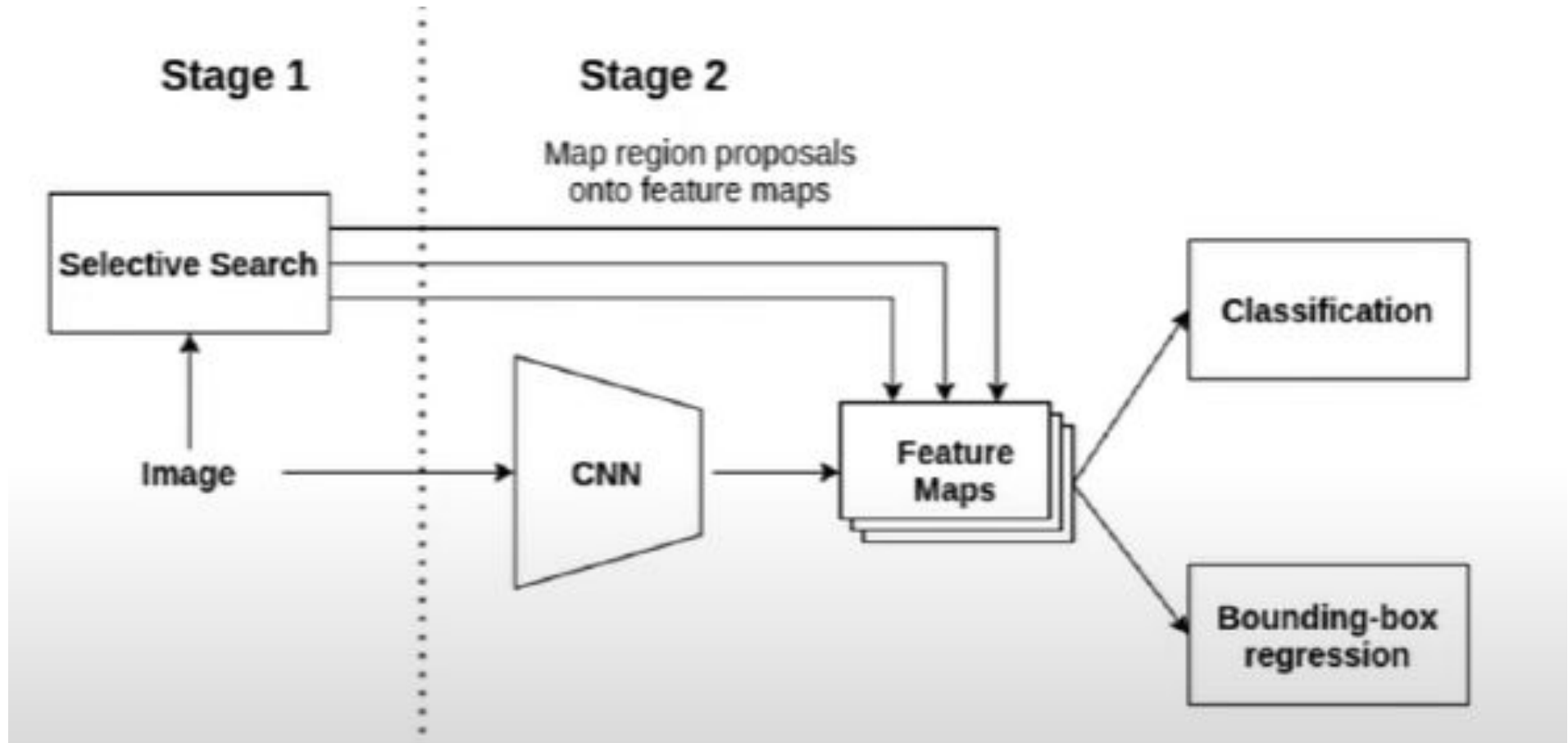
# Object detection methods



# Selective search

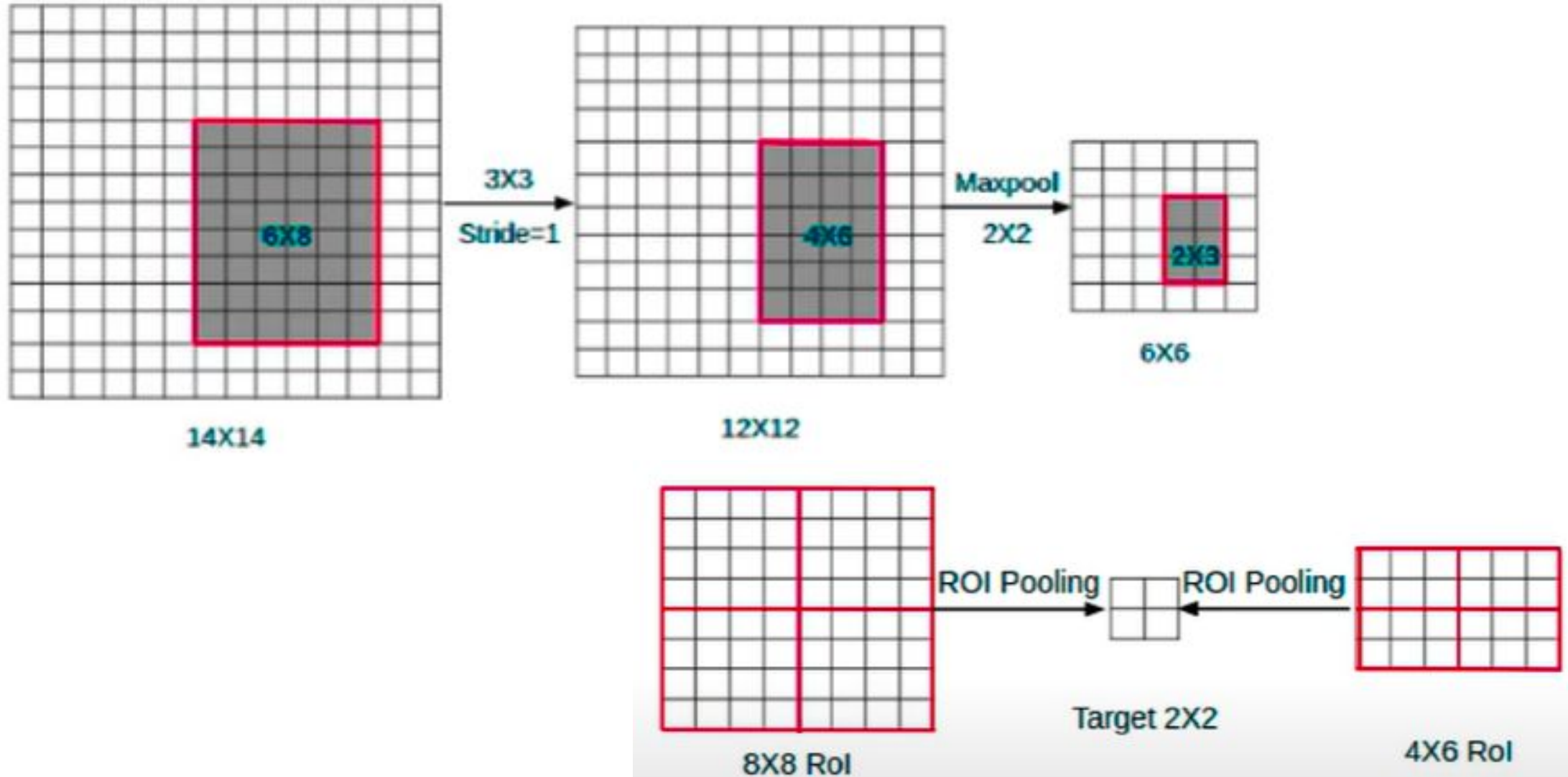


# Fast - RCNN



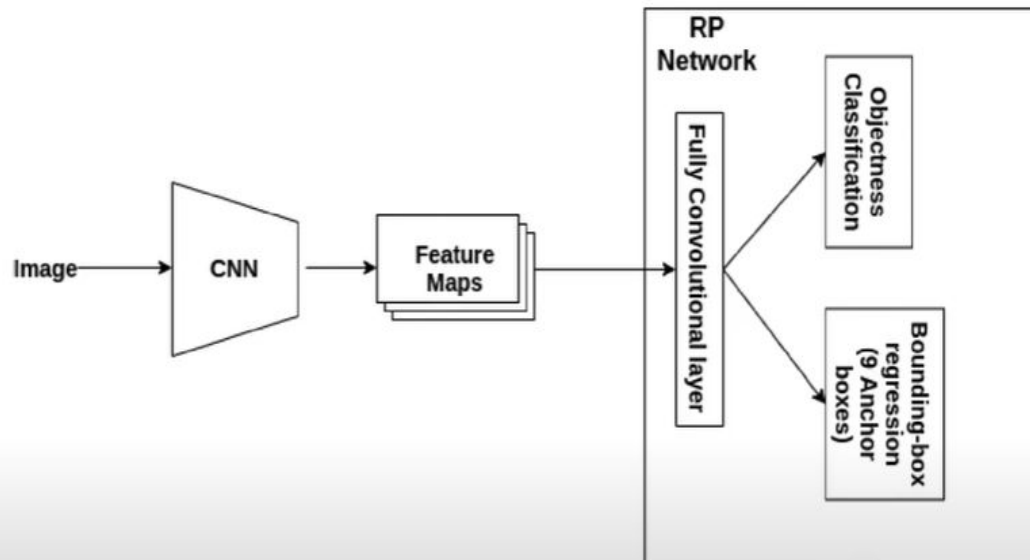
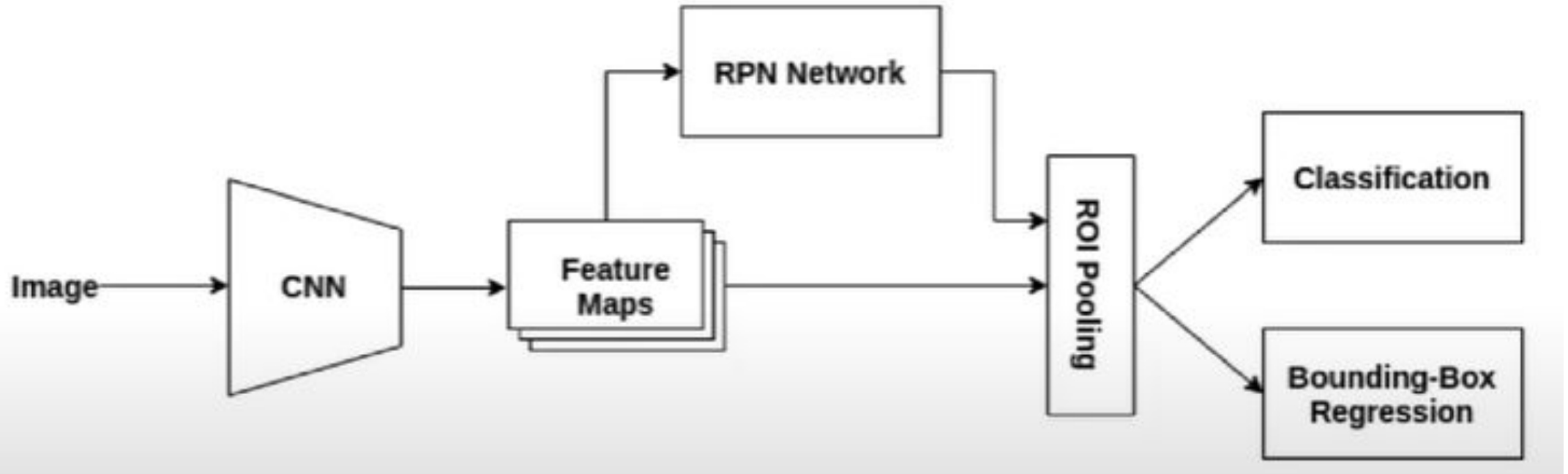


# ROI projection and pooling





# Faster - RCNN



Input Image  
(e.g. 3 x 640 x 480)

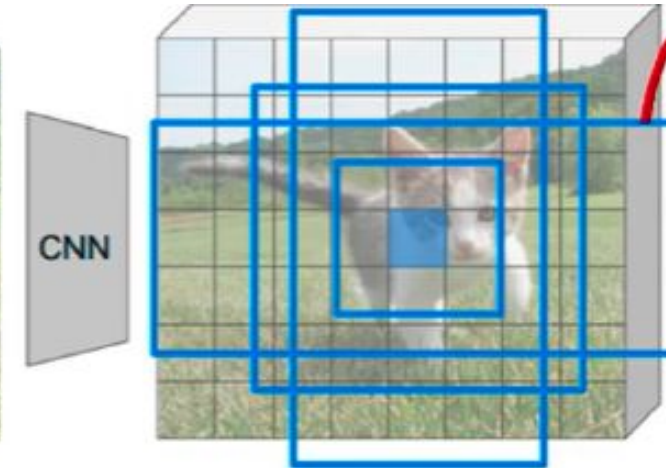
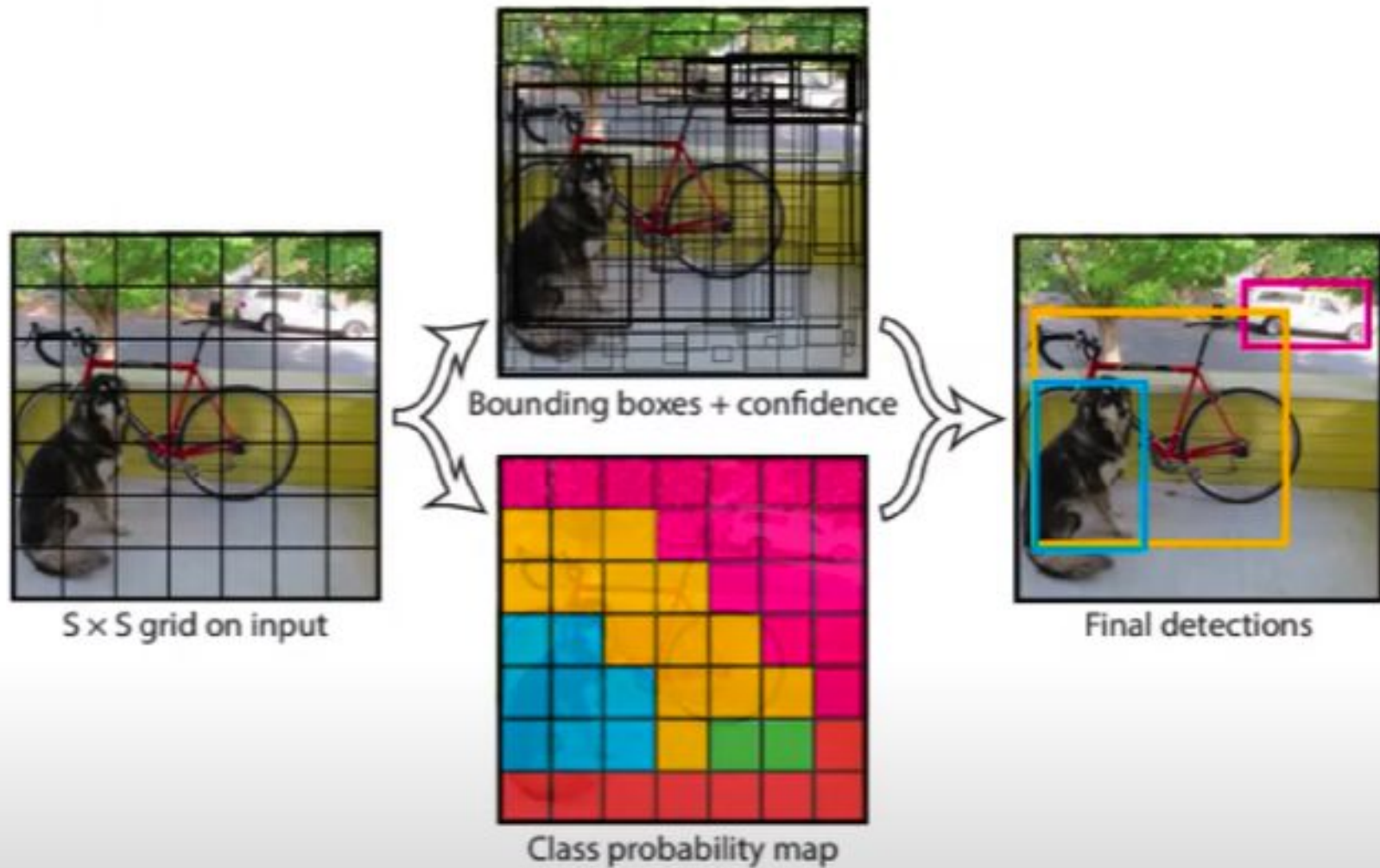
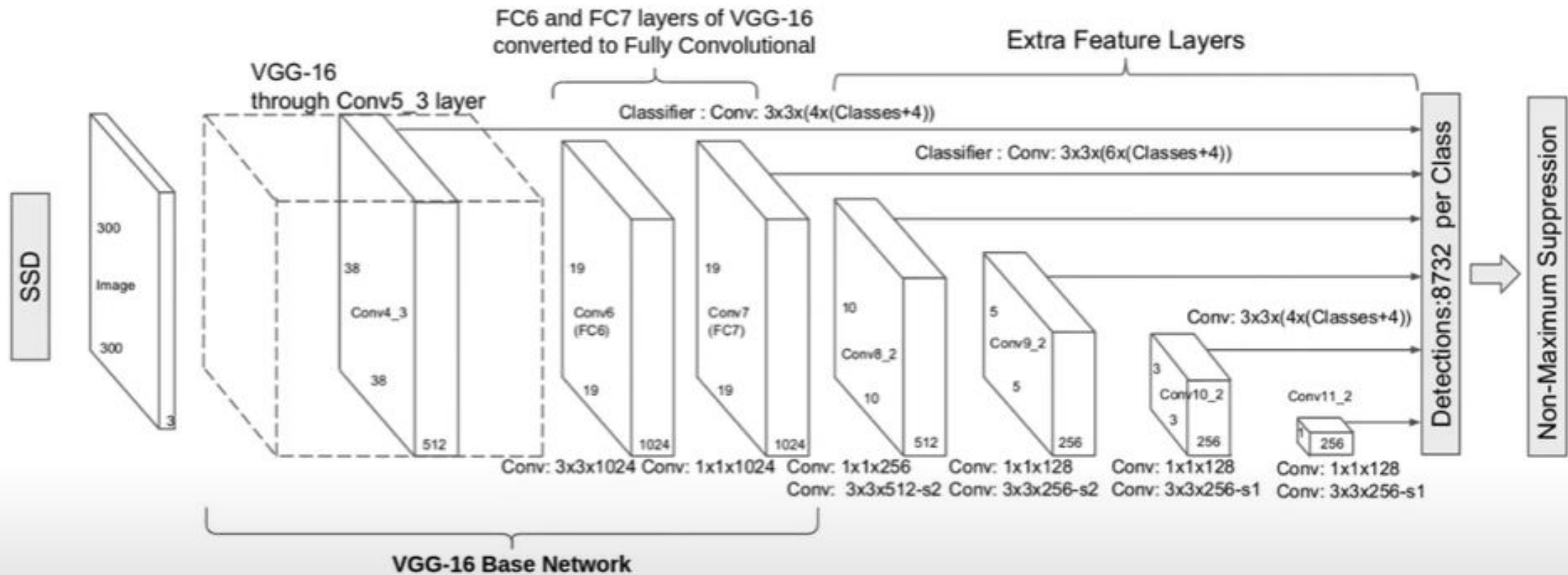


Image features  
(e.g. 512 x 20 x 15)

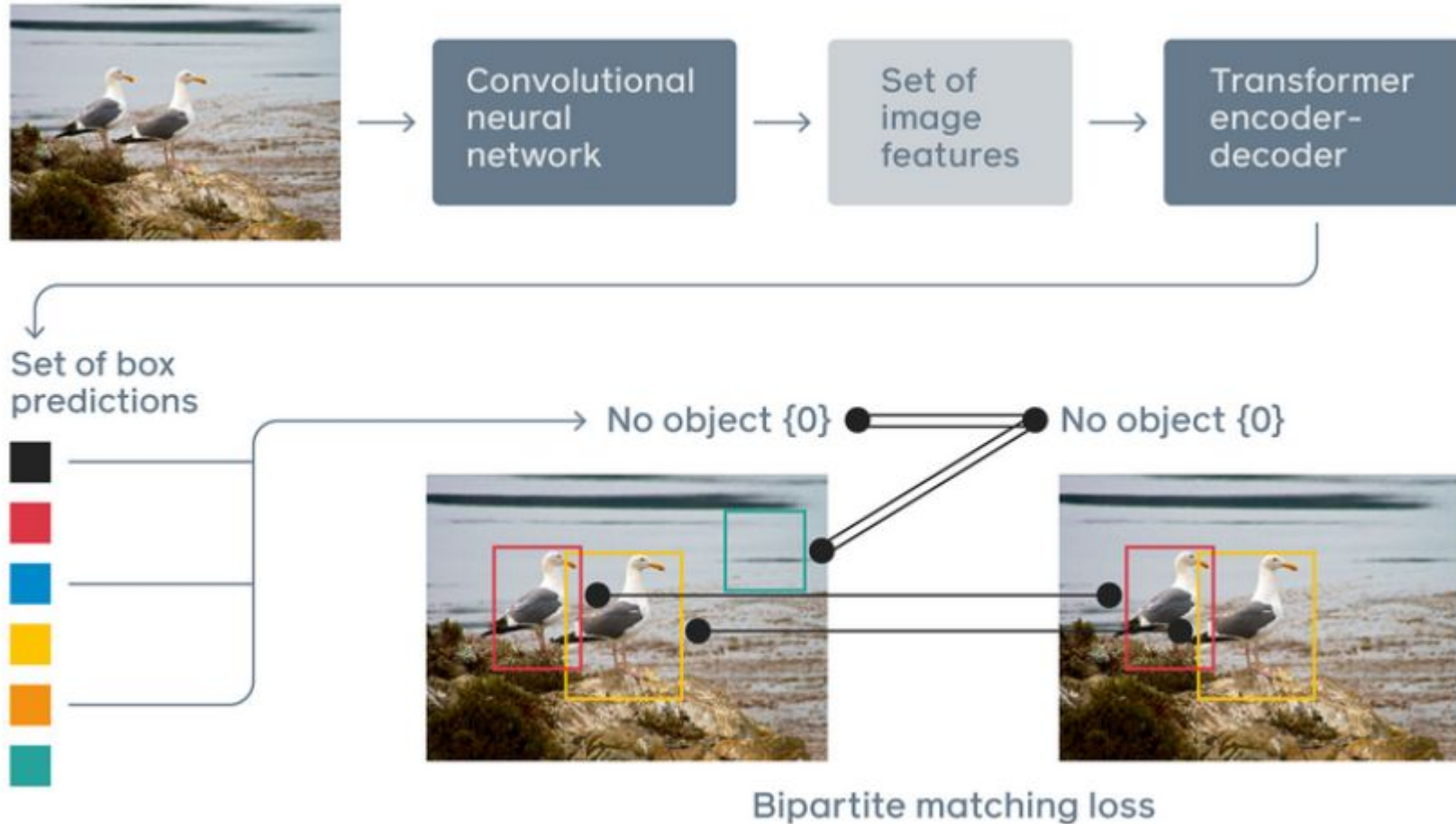
# YOLO

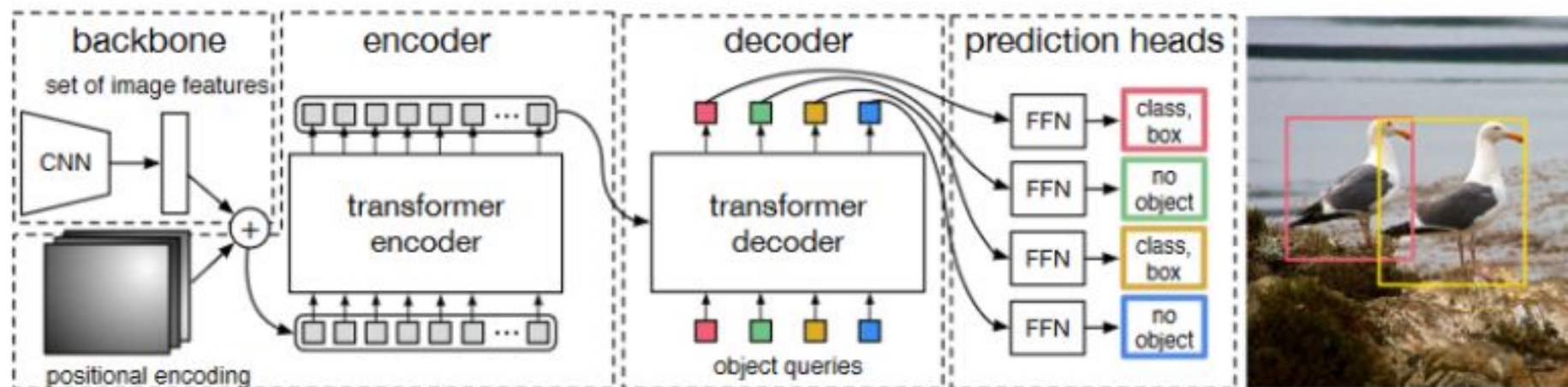


# SSD

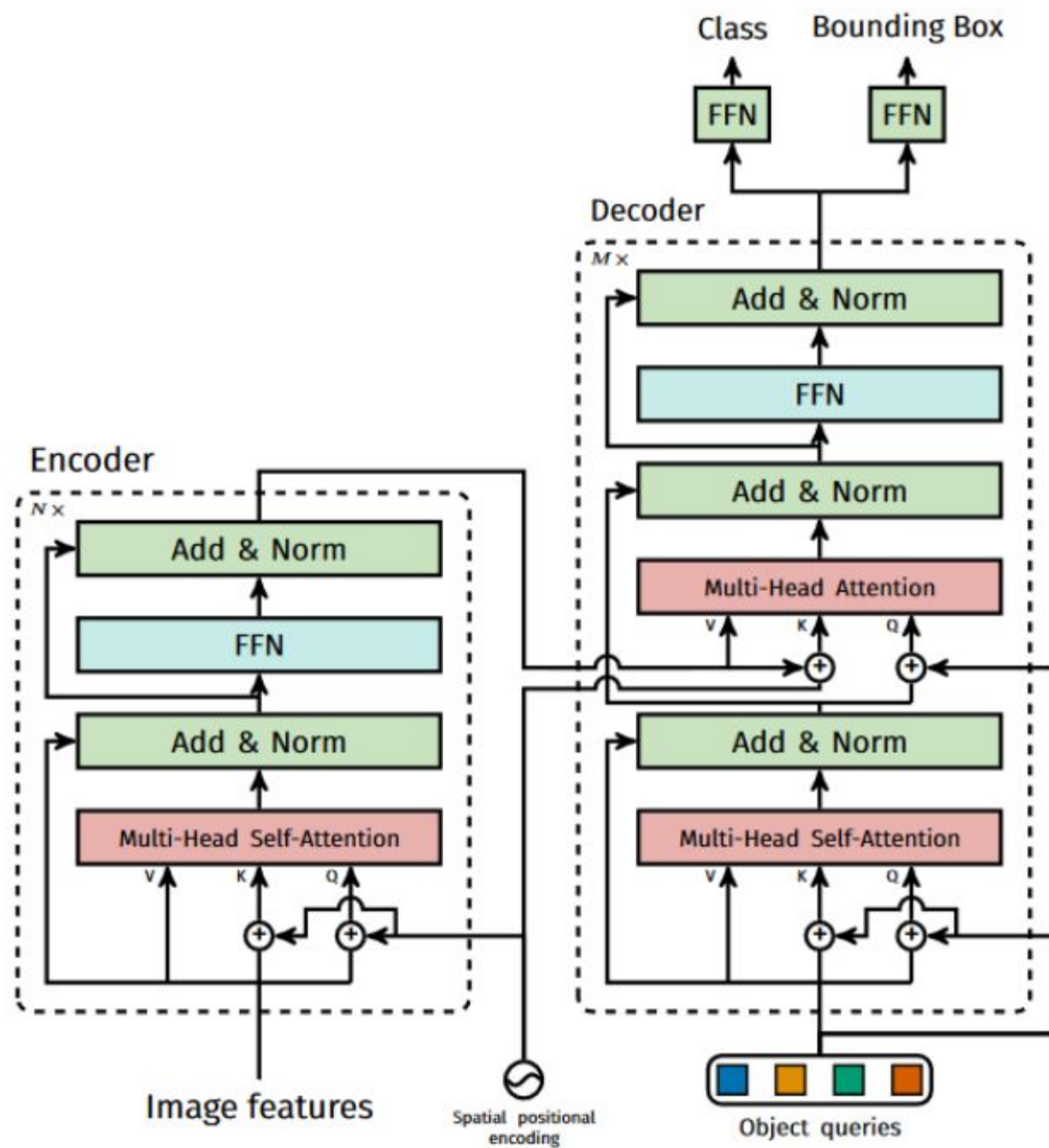


# DETR (DEtection TRansformer)







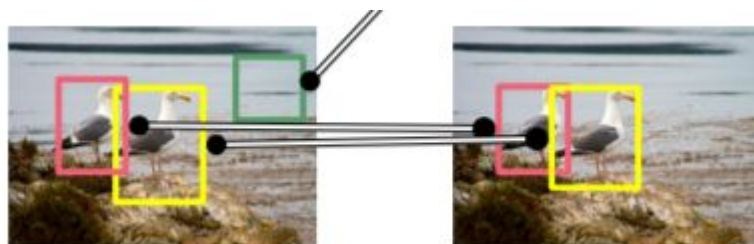


# Loss function

- Calculate the best match of predictions with respect to given ground truths using a graph technique with a cost function

$$\hat{\sigma} = \arg \min_{\sigma \in \mathfrak{S}_N} \sum_i^N \mathcal{L}_{\text{match}}(y_i, \hat{y}_{\sigma(i)}),$$

$$\mathcal{L}_{\text{Hungarian}}(y, \hat{y}) = \sum_{i=1}^N \left[ -\log \hat{p}_{\hat{\sigma}(i)}(c_i) + \mathbb{1}_{\{c_i \neq \emptyset\}} \mathcal{L}_{\text{box}}(b_i, \hat{b}_{\hat{\sigma}(i)}) \right]$$





THANK YOU