

Harnessing Transformers for various Computer Vision tasks

Rajesh Shreedhar Bhat
Sr. Data Scientist, Walmart & GDE ML

Section 01

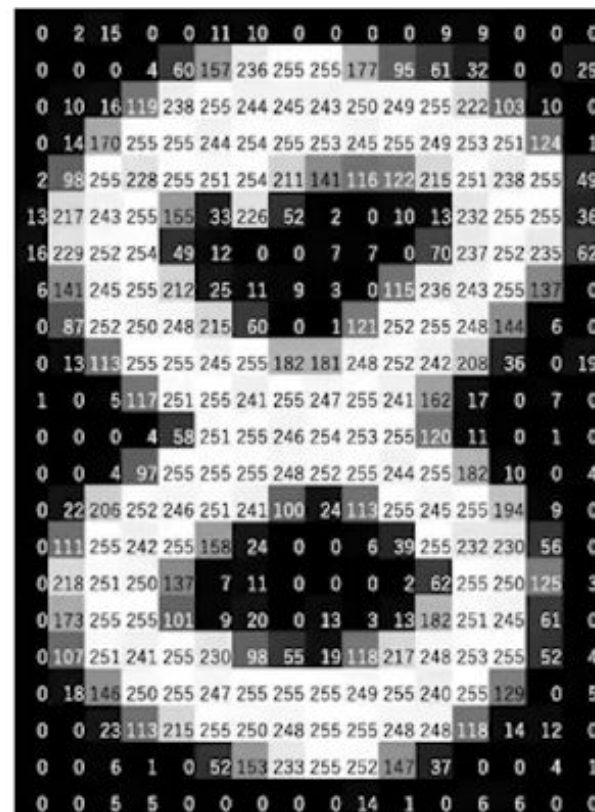
Computer Vision and various vision tasks

Computer Vision

Computer vision is a field of artificial intelligence (AI) that enables computers and systems to derive meaningful information from digital images and videos.

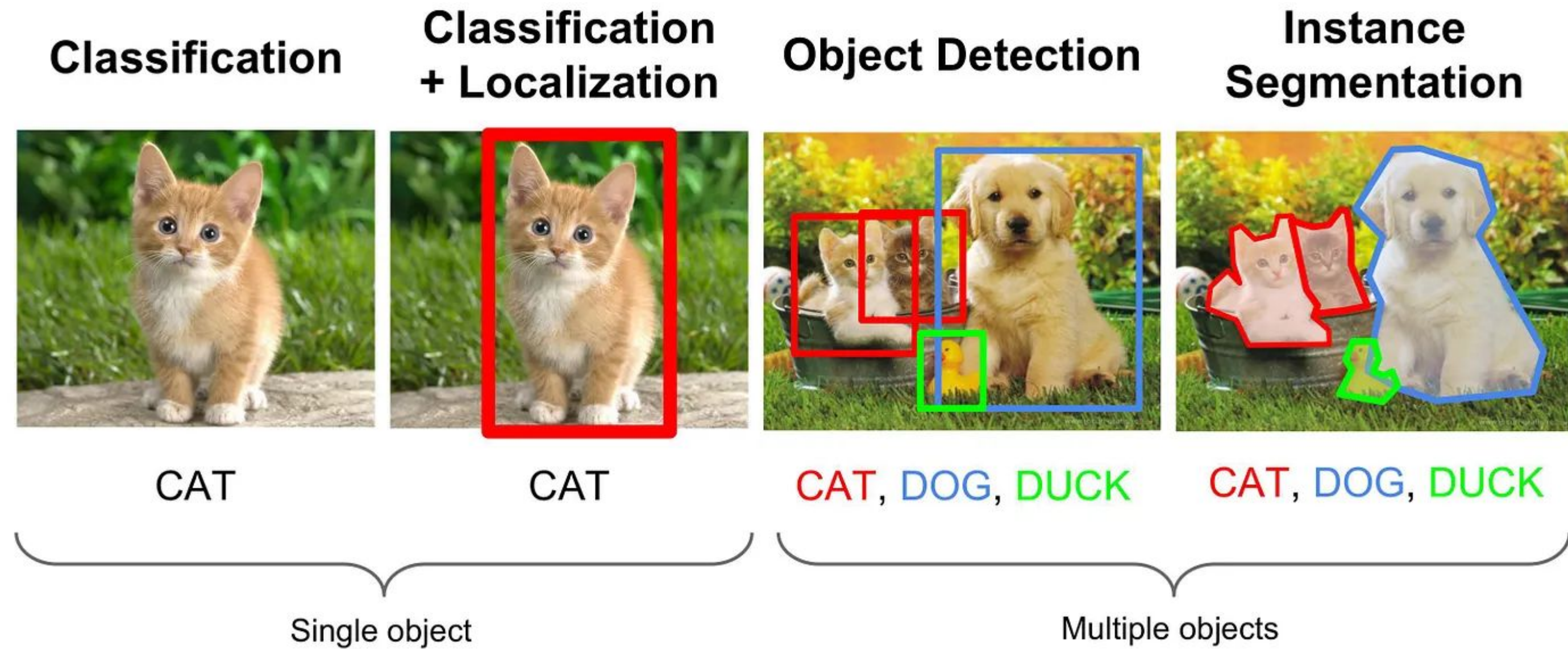


Grayscale Image

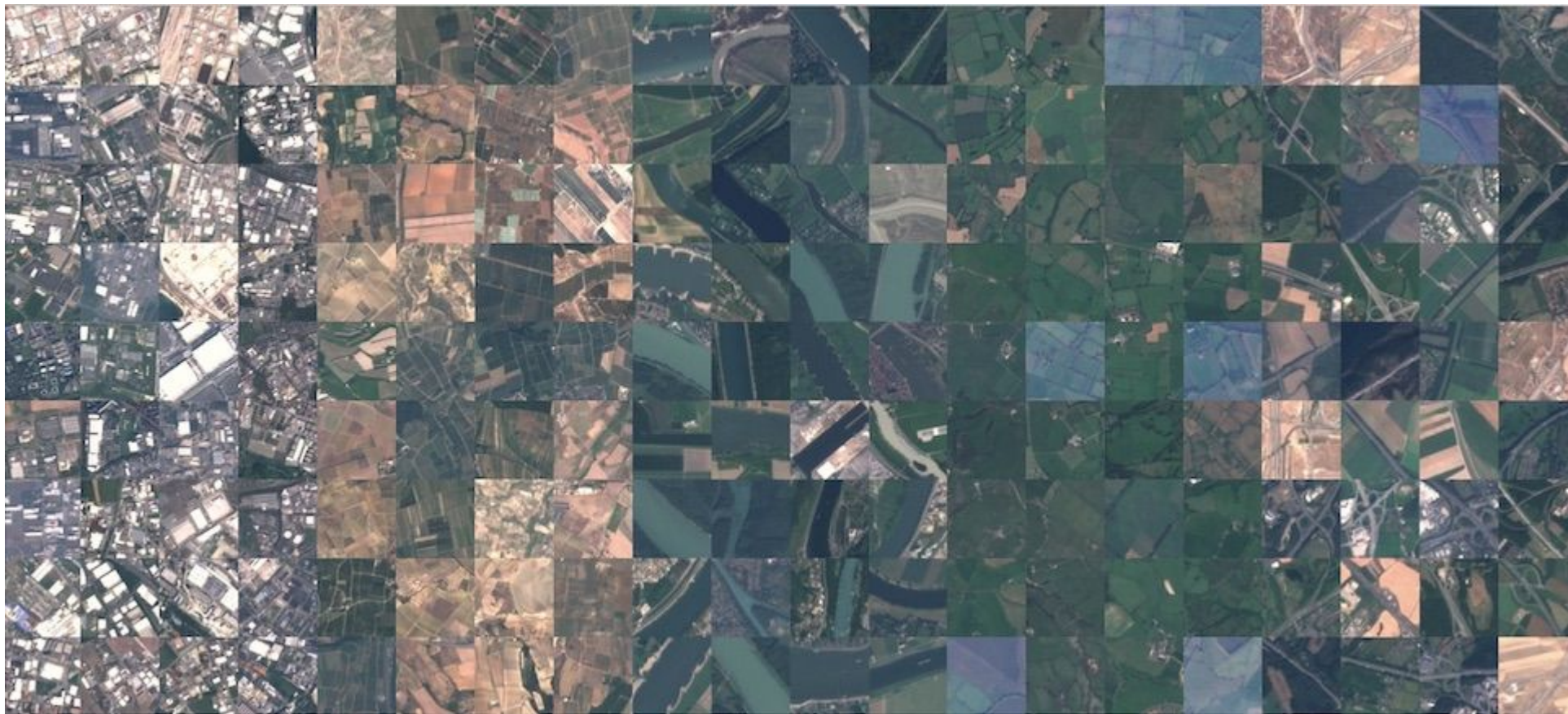


Color Image

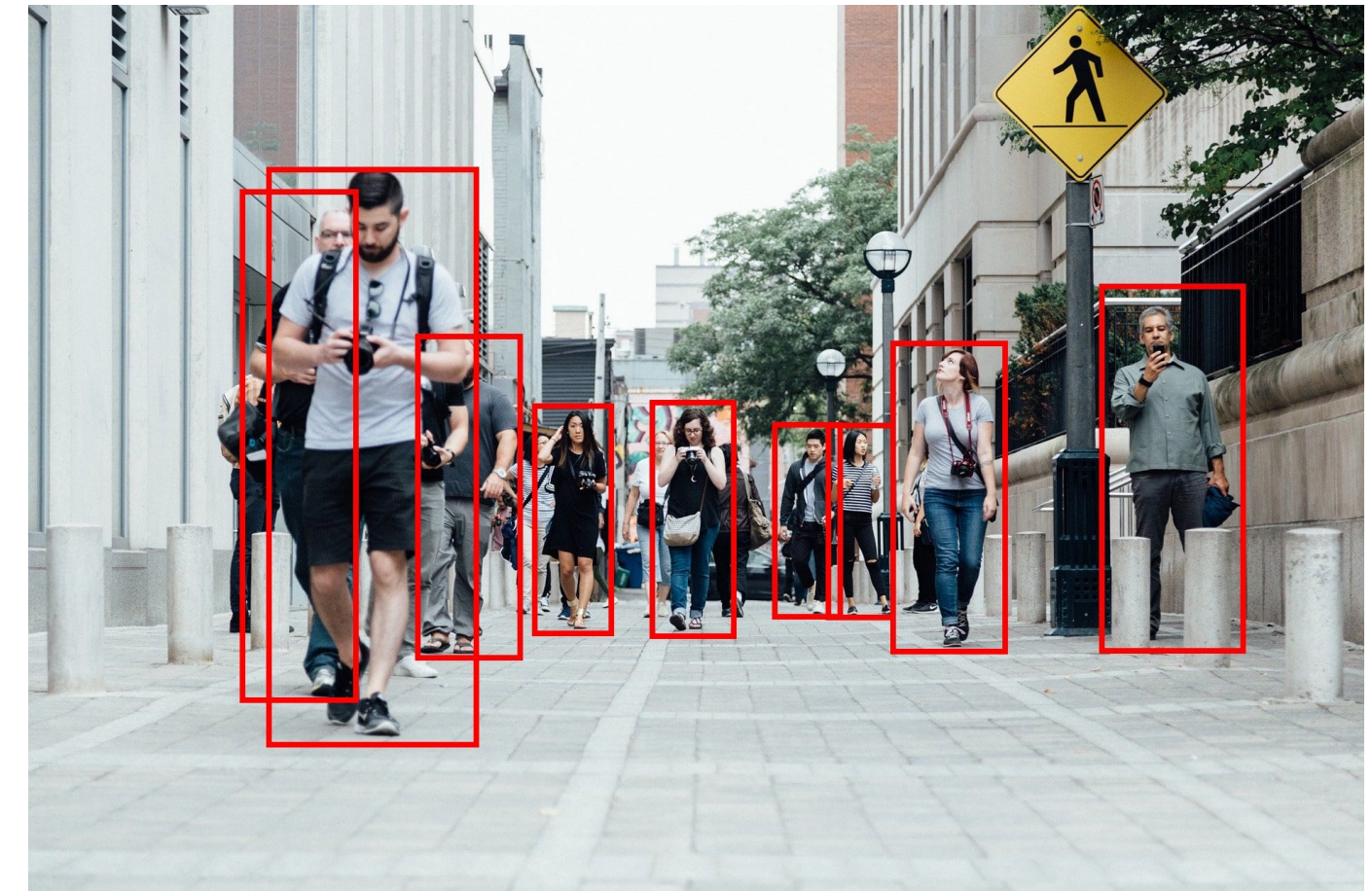
Various Computer Vision tasks



Other examples



Classification: Land Cover Classification Dataset
with classes like Forest area, River, Highway etc ..



Object Detection: People Detection

Other examples ..



Segmentation: Lungs Segmentation



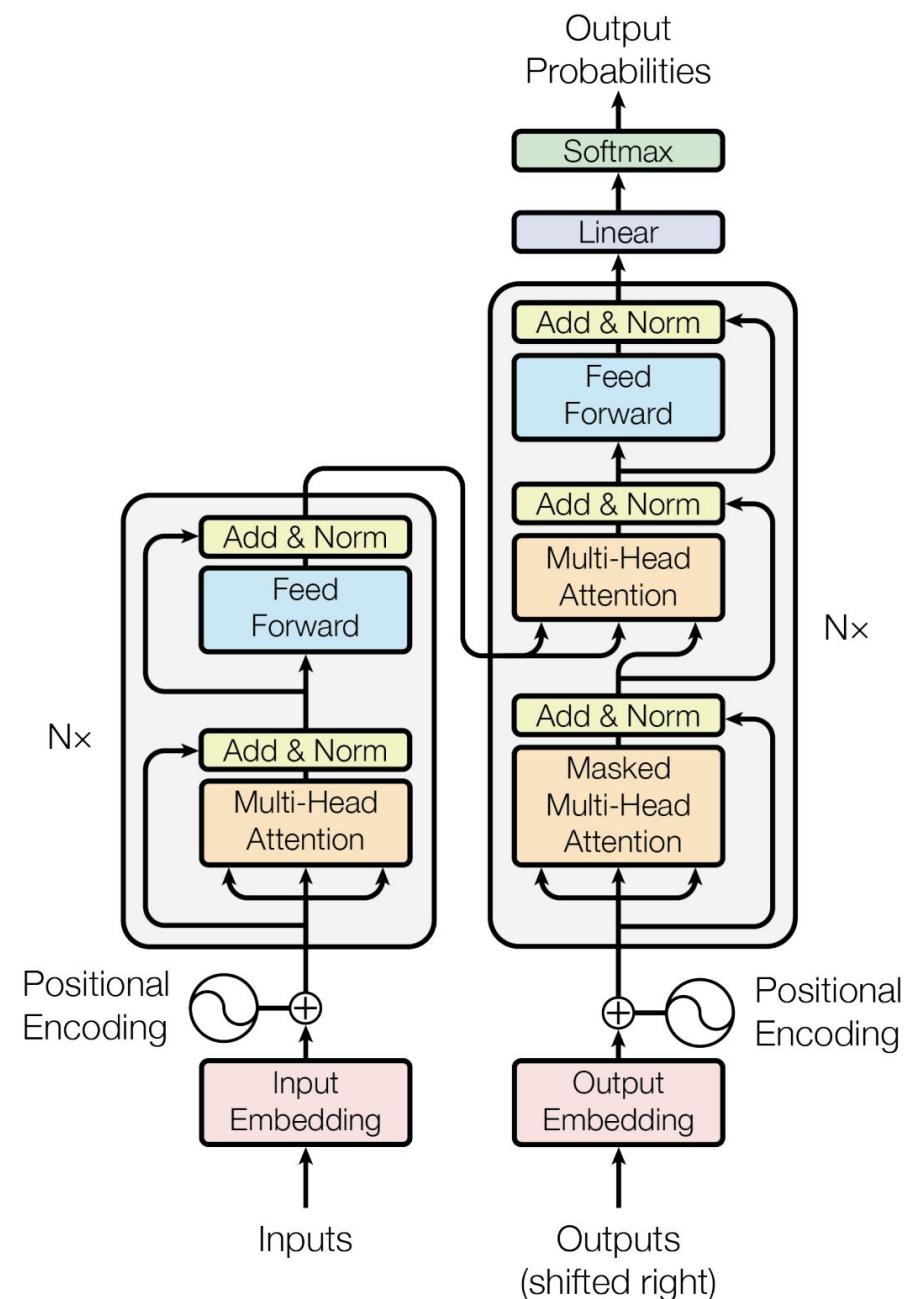
Instance Segmentation: cars, road, pedestrian, cyclists etc ..

Section 02

Evolution of Transformer models

Section 3

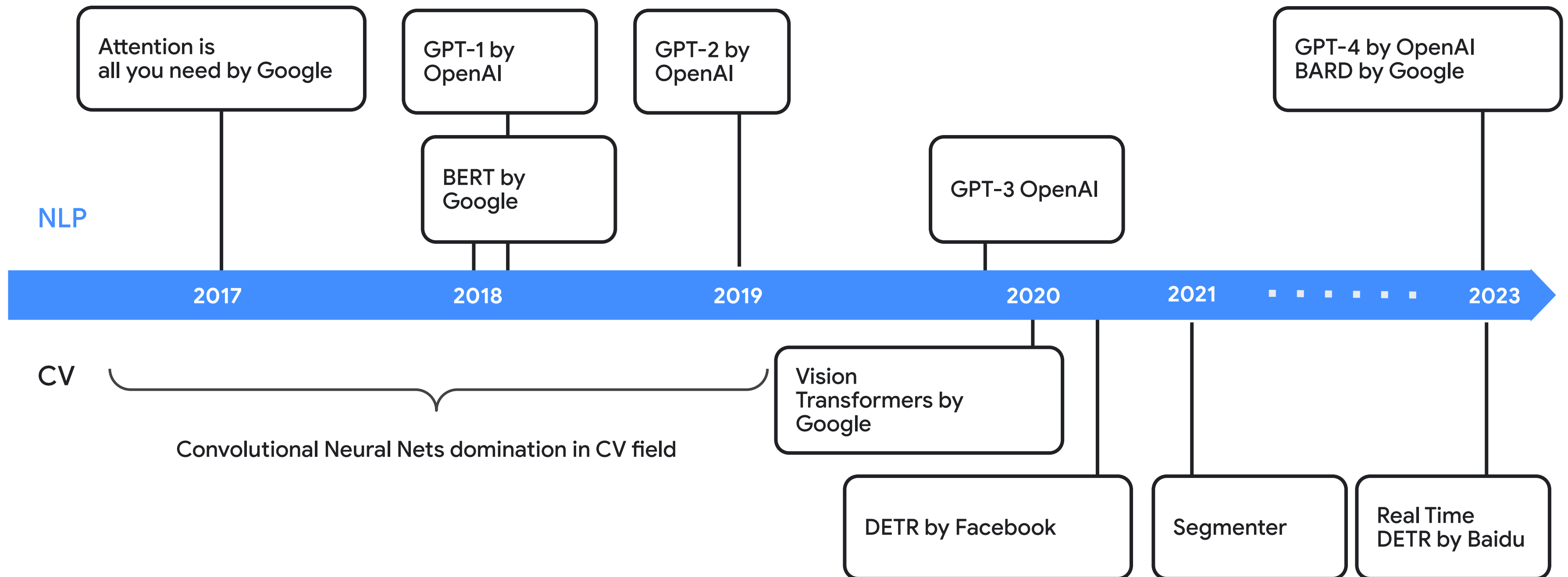
Attention is all you need !



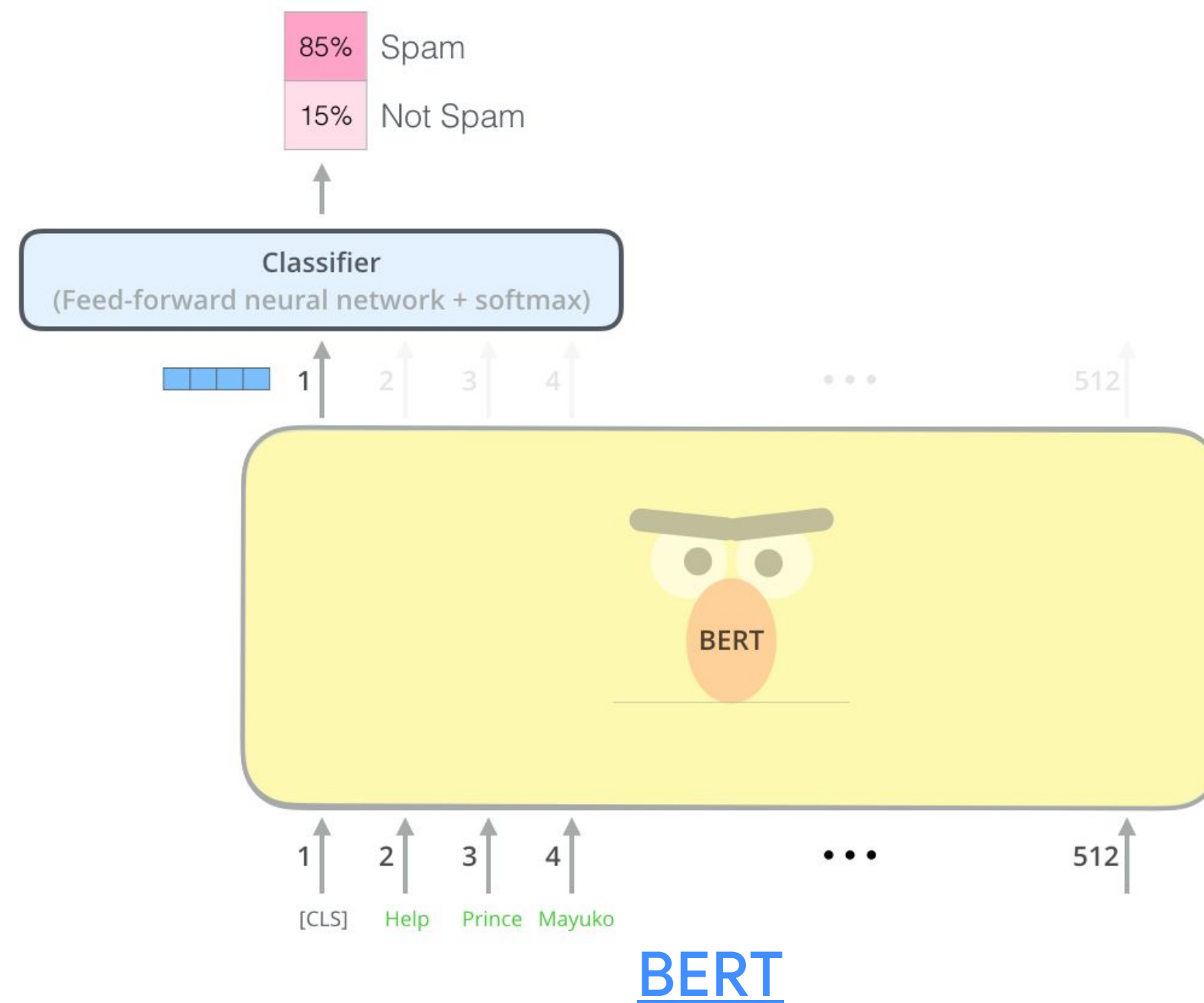
- Outperformed previous state-of-the-art models on a variety of natural language processing tasks
- Landmark paper in the field of natural language processing(NLP) by Vaswani et.al

Section 3

Timeline of different transformer based models



BERT: Bidirectional Encoder Representations from Transformer

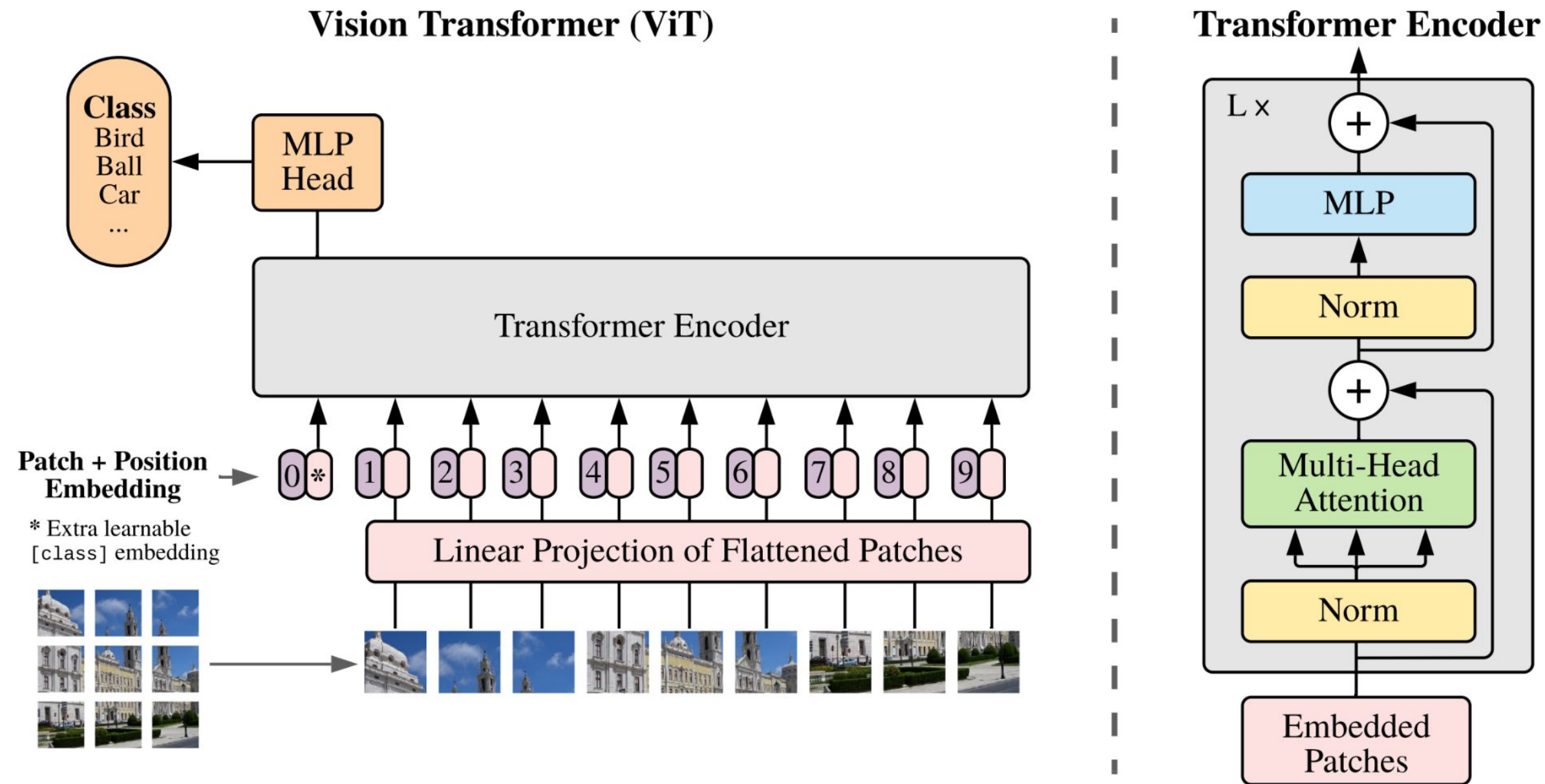


Section 03

Vision Transformer(ViT)

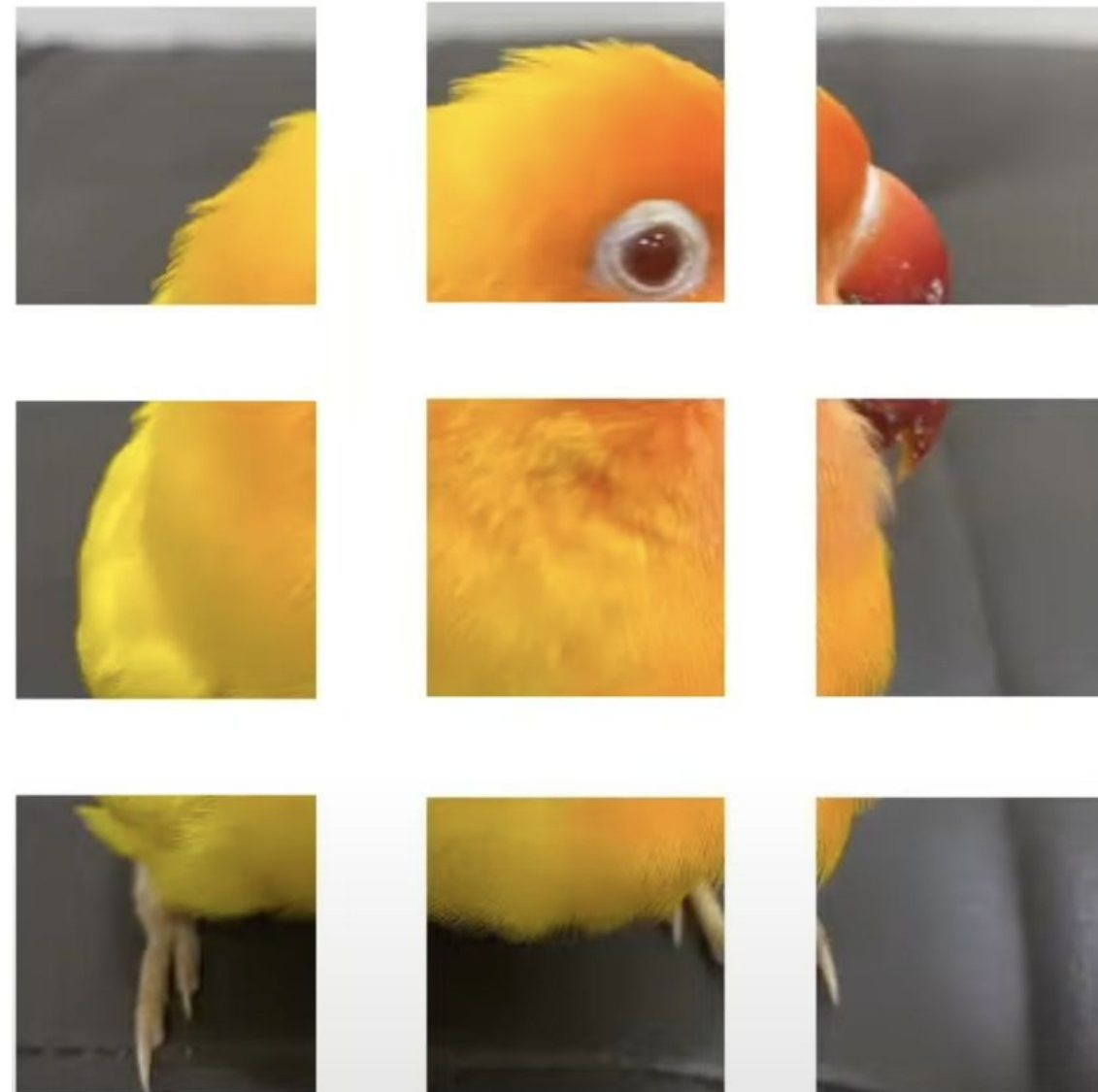
Details of ViT model architecture

ViT Architecture



Ref: [An image is worth 16 x 16 words: Transformers for Image Recognition at Scale](#)

Image to Patches



Number of patches
 $= W * H / (P * P)$

where

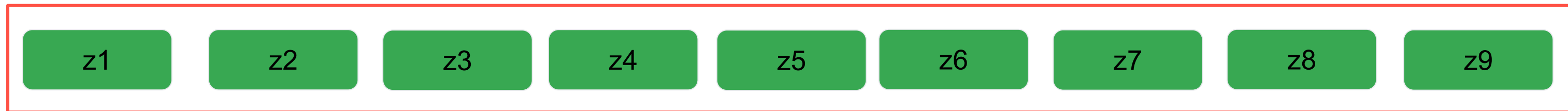
W : width of the image

H : height of the image

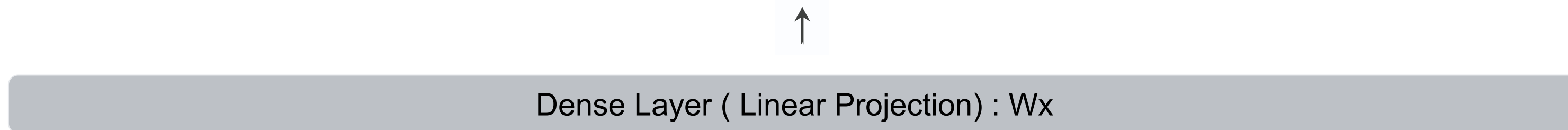
P : Patch size

Section 3

Patch embeddings



$z(i) \in \mathbb{R}^{1 \times D}$
where i from 1 to 9



$W \in \mathbb{R}^{(P \times P \times C) \times D}$



$x(i) \in \mathbb{R}^{1 \times (P \times P \times C)}$
where i from 1 to 9

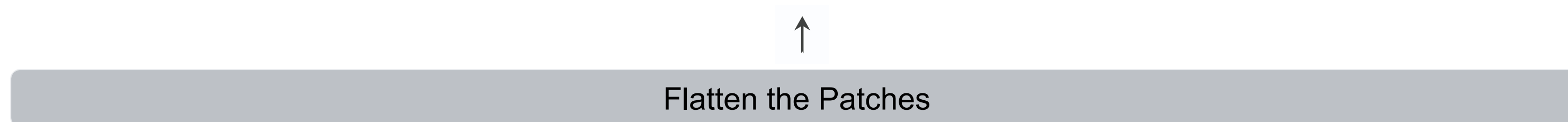
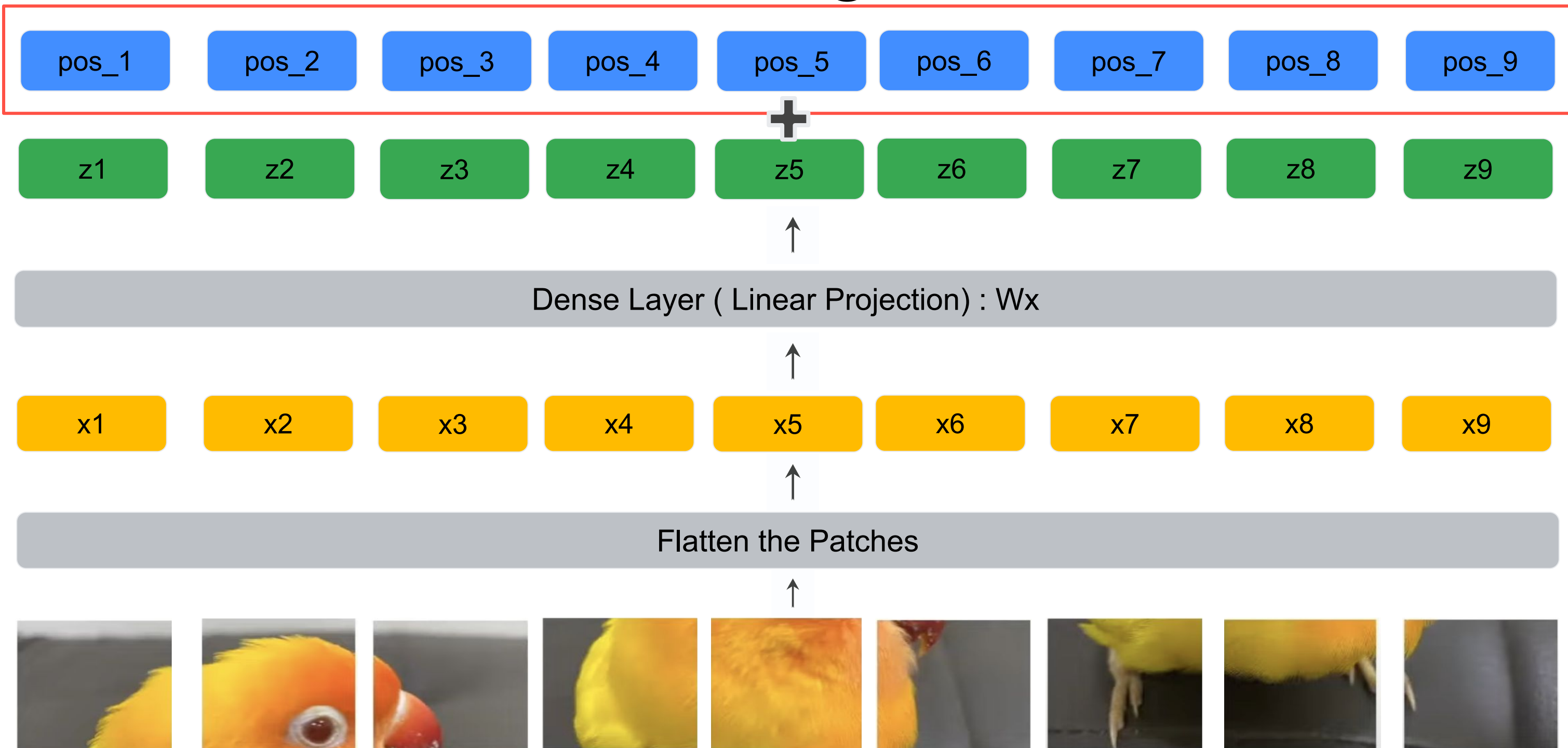


Image Patches
of size $P \times P \times C$

Section 3

Positional embeddings



$pos_i \in \mathbb{R}^{1 \times D}$
where i from 1 to 9

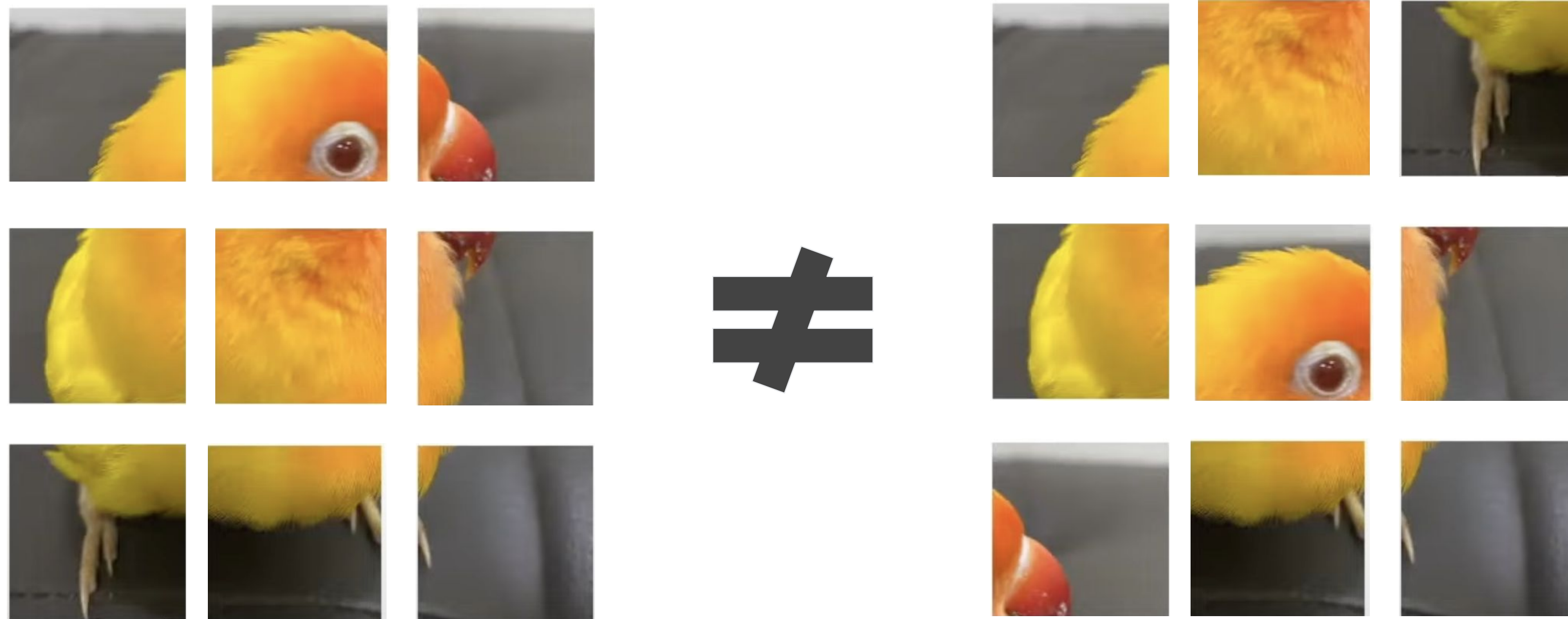
$z(i) \in \mathbb{R}^{1 \times D}$
where i from 1 to 9

$W \in \mathbb{R}^{(P \times P \times C) \times D}$

$x(i) \in \mathbb{R}^{1 \times (P \times P \times C)}$
where i from 1 to 9

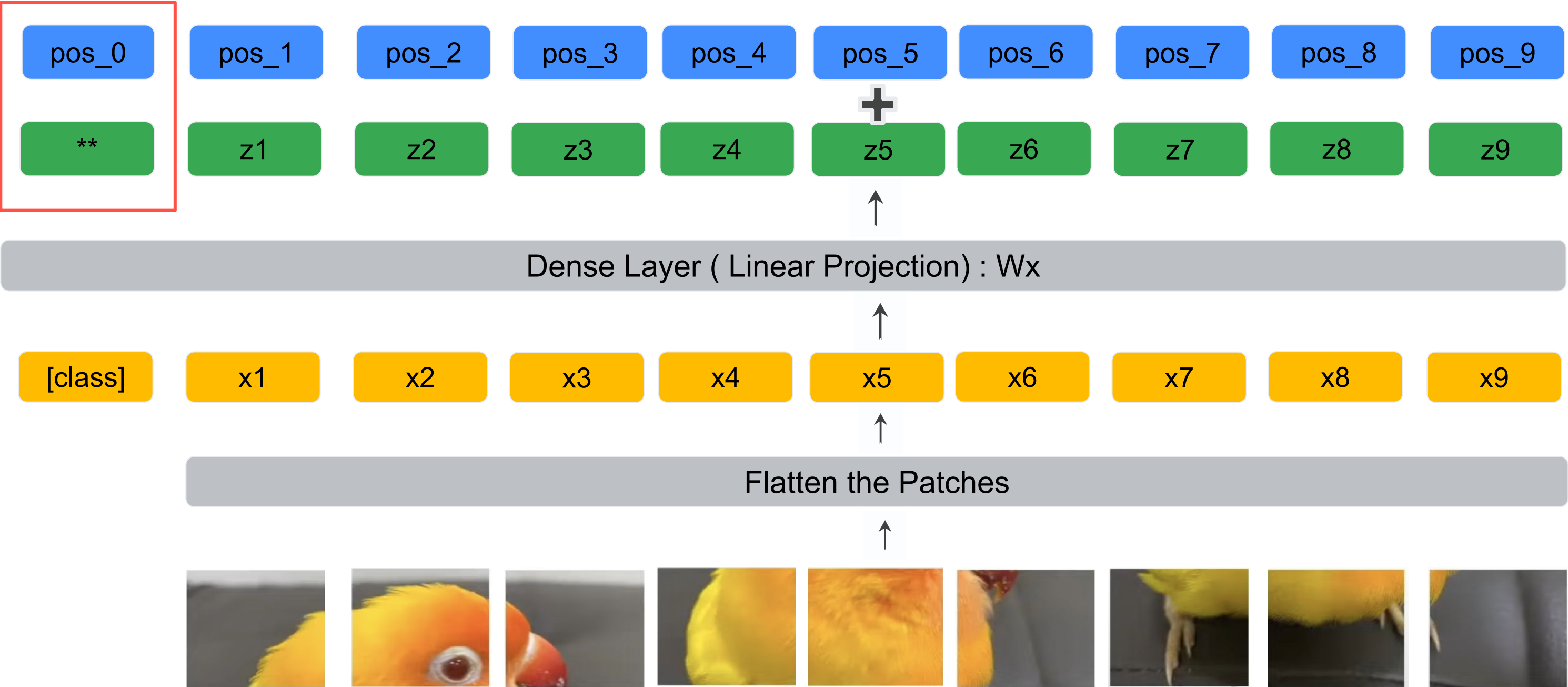
Image Patches
of size $P \times P \times C$

Why Positional Embeddings ?



Allows the model to learn the relative position of each patch in the image
3% drop in accuracy observed without positional embeddings

[class] embedding



$pos_{(i)} \in R^{1 \times D}$
where i from 1 to 9

$z(i) \in R^{1 \times D}$
where i from 1 to 9

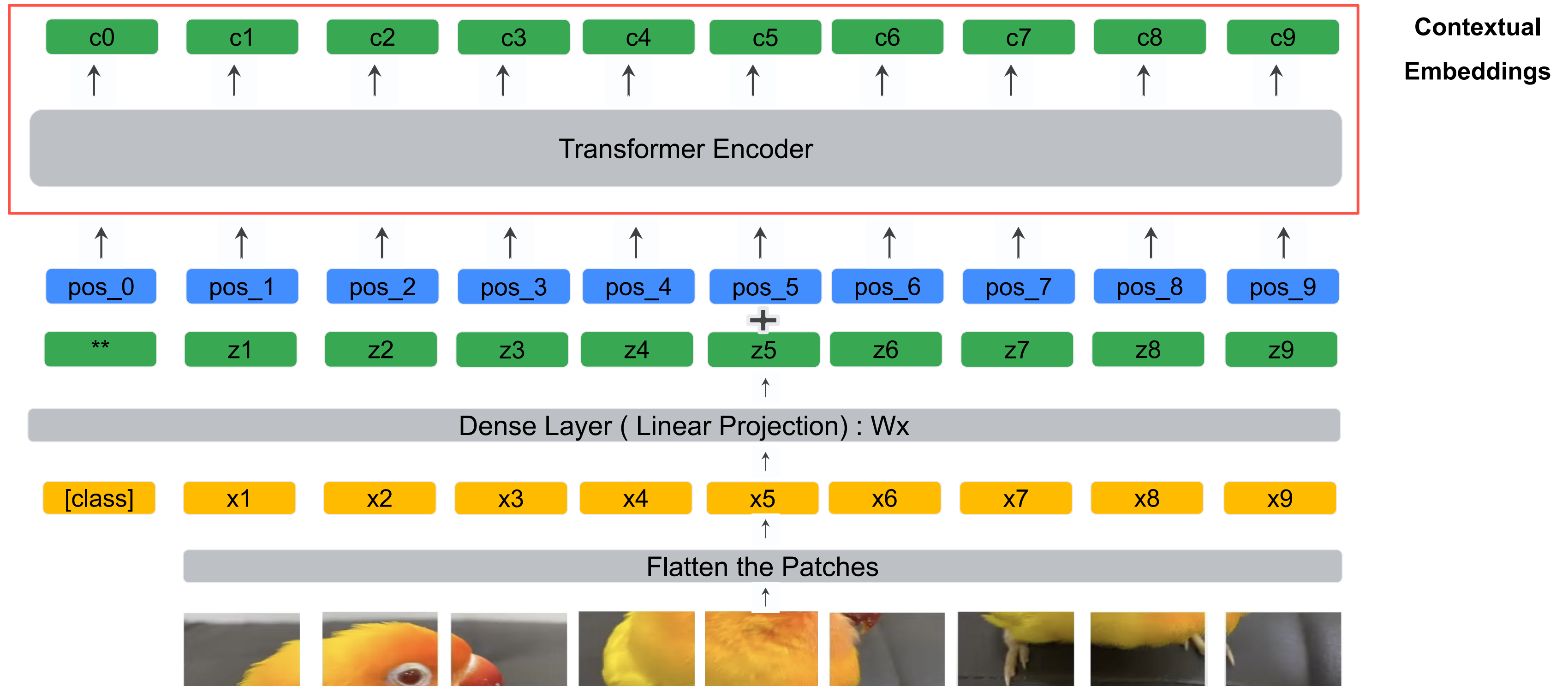
$W \in R^{(P*P*C) \times D}$

$x(i) \in R^{1 \times (P * P * C)}$
where i from 1 to 9

Image Patches
of size $P * P * C$

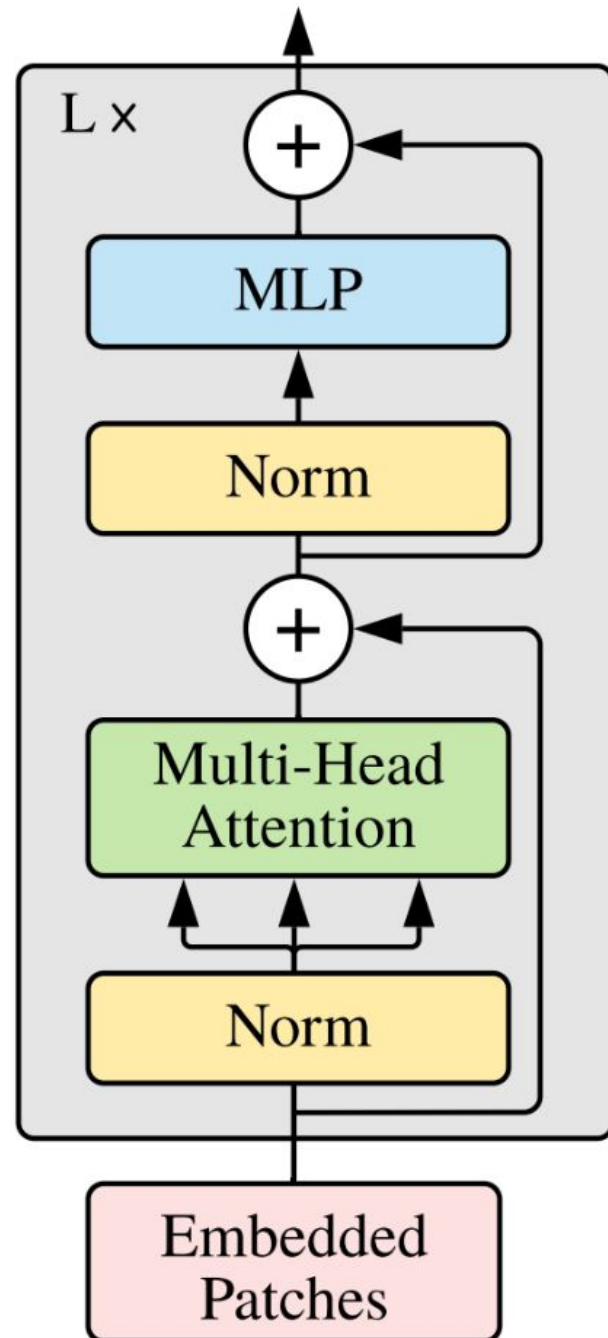
Section 3

Transformer Encoder

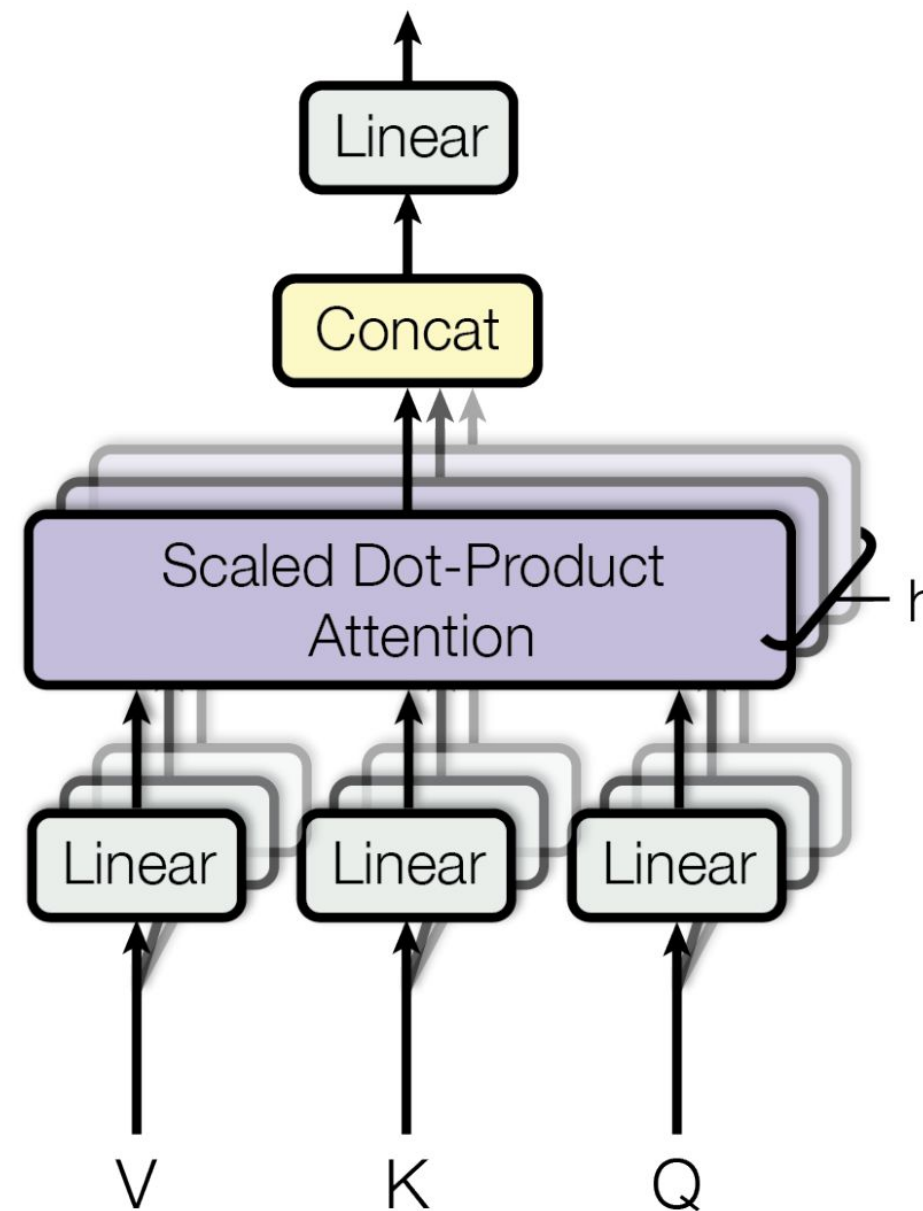


Section 3

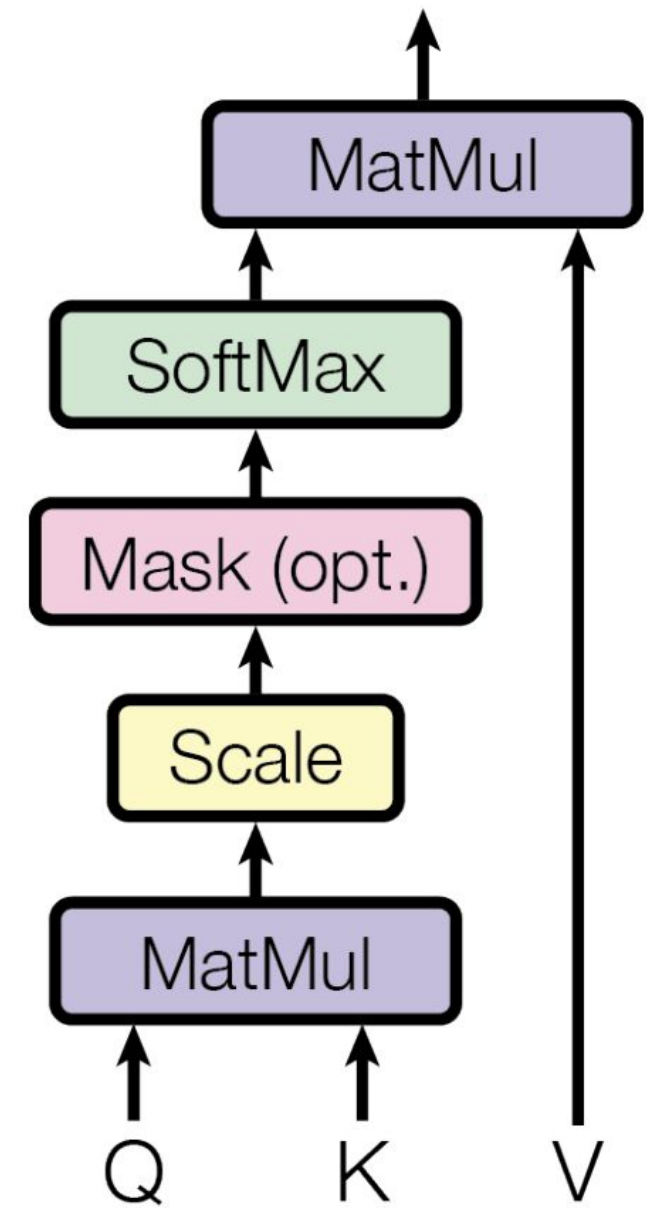
Transformer Encoder continued ..



Single Encoder Block



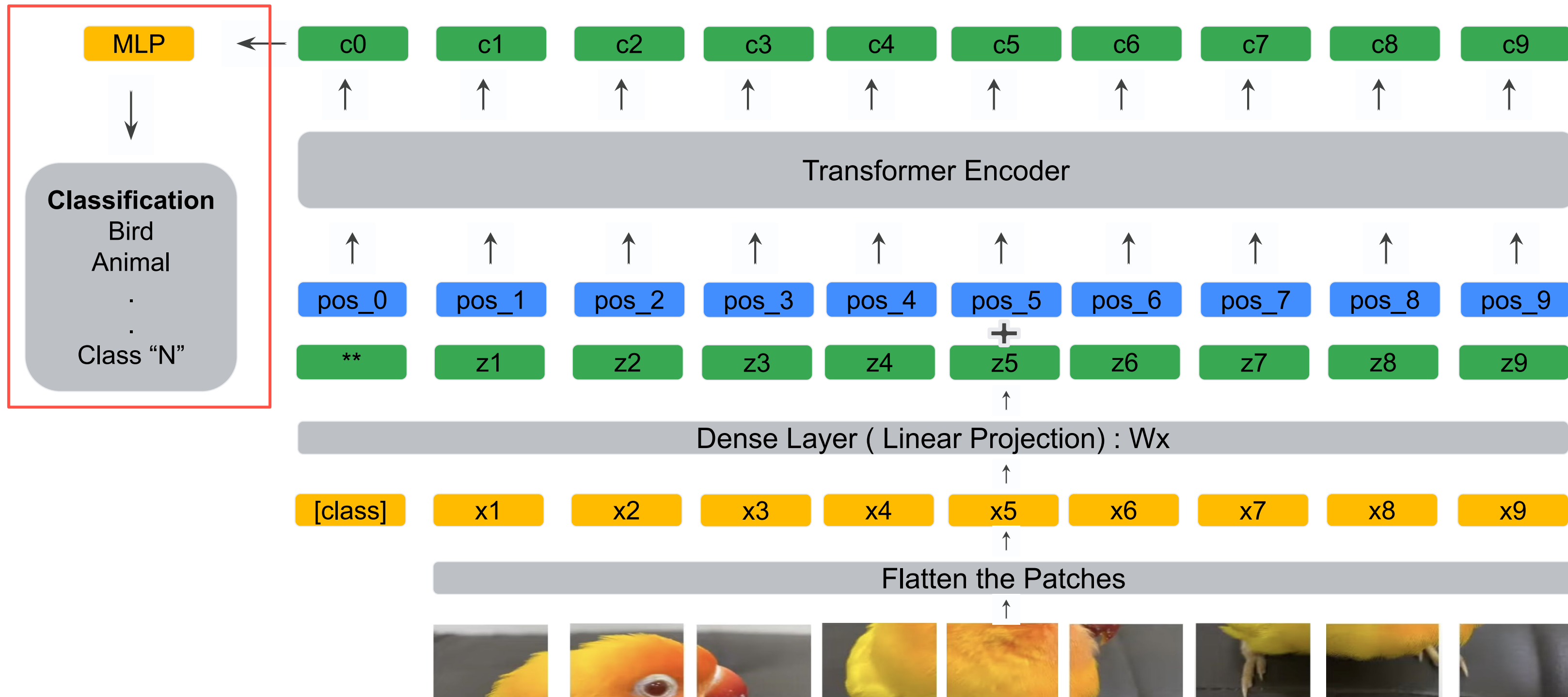
Multi head attention



Scaled Dot Product Attention

Section 3

ViT: MLP head and classification



Section 04

Keras Implementation for Vision Transformer(ViT)

https://github.com/rajesh-bhat/vision_transformers_googleio_extended



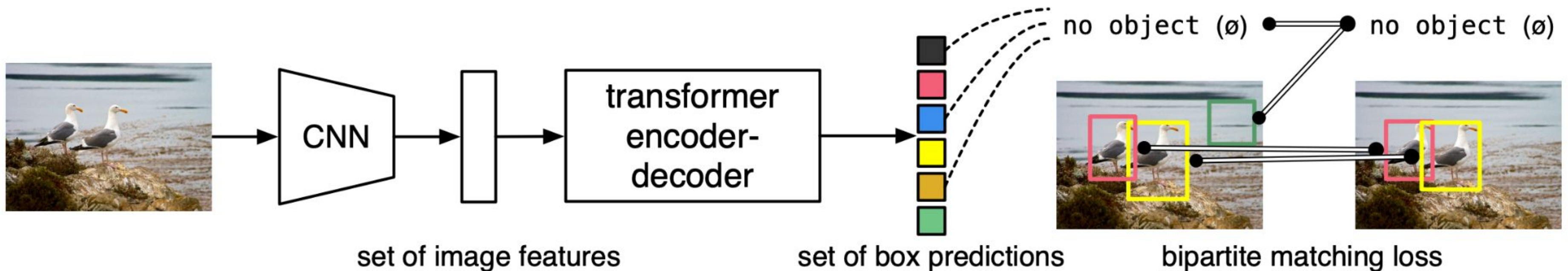
Transformer Encoder Keras NLP: code snippet

```
keras_nlp.layers.TransformerEncoder(  
    intermediate_dim,  
    num_heads,  
    dropout=0,  
    activation="relu",  
    layer_norm_epsilon=1e-05,  
    kernel_initializer="glorot_uniform",  
    bias_initializer="zeros",  
    normalize_first=False,  
    name=None,  
    **kwargs  
)
```

Transformers for other Vision tasks

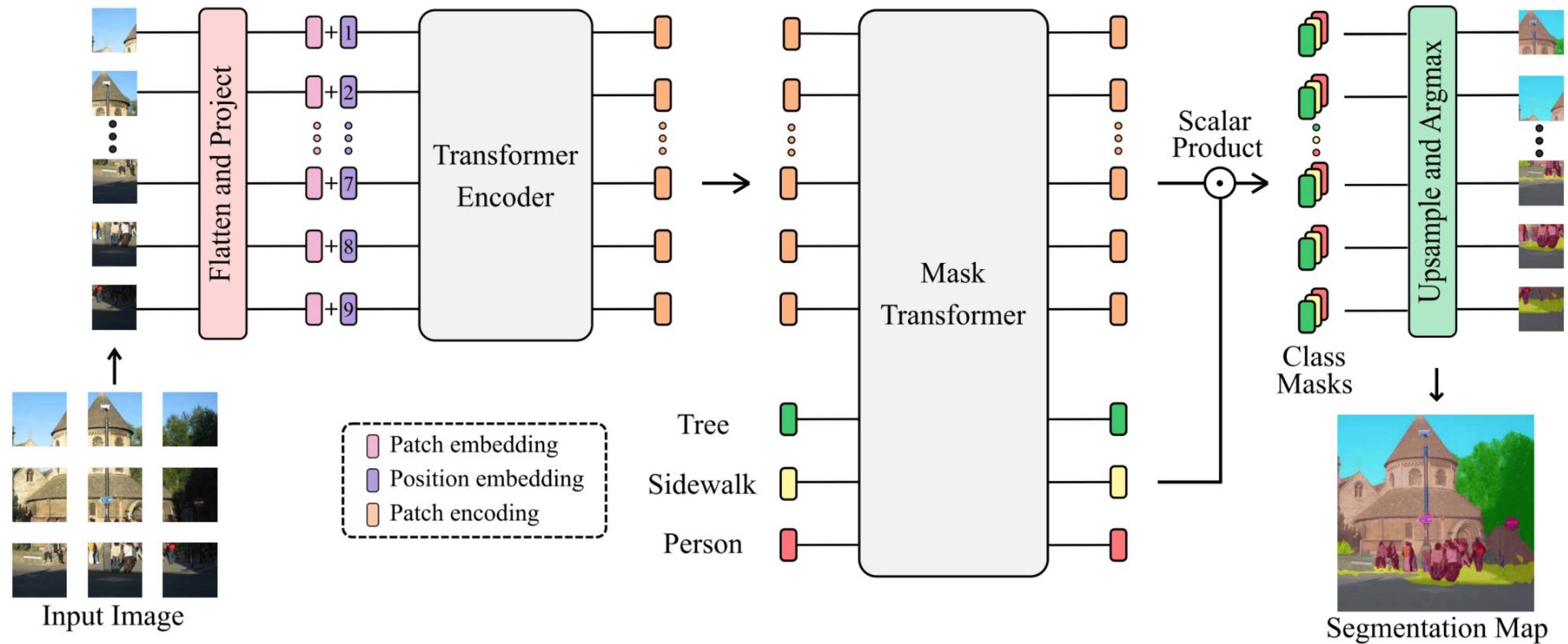
Object Detection, segmentation ..

DETR: End to End Object detection with Transformers



[DETR](#) directly predicts (in parallel) the final set of detections by combining a common CNN with a transformer architecture

Segmenter



Overview of [Segementer](#) model

Summary

- Intro to Computer Vision and various tasks
- Evolution of Transformer models
- Vision Transformer(ViT) architecture in detail
 - Images to Patches
 - Patch Embeddings
 - Positional Embeddings
 - Transformer Encoder
 - MLP head and classification
- Vision Transformers in Keras
- Transformers for Object Detection(DETR) and Segmentation tasks(Segmenter)

Thank You



Rajesh Shreedhar Bhat

Sr. Data Scientist, Walmart & GDE ML

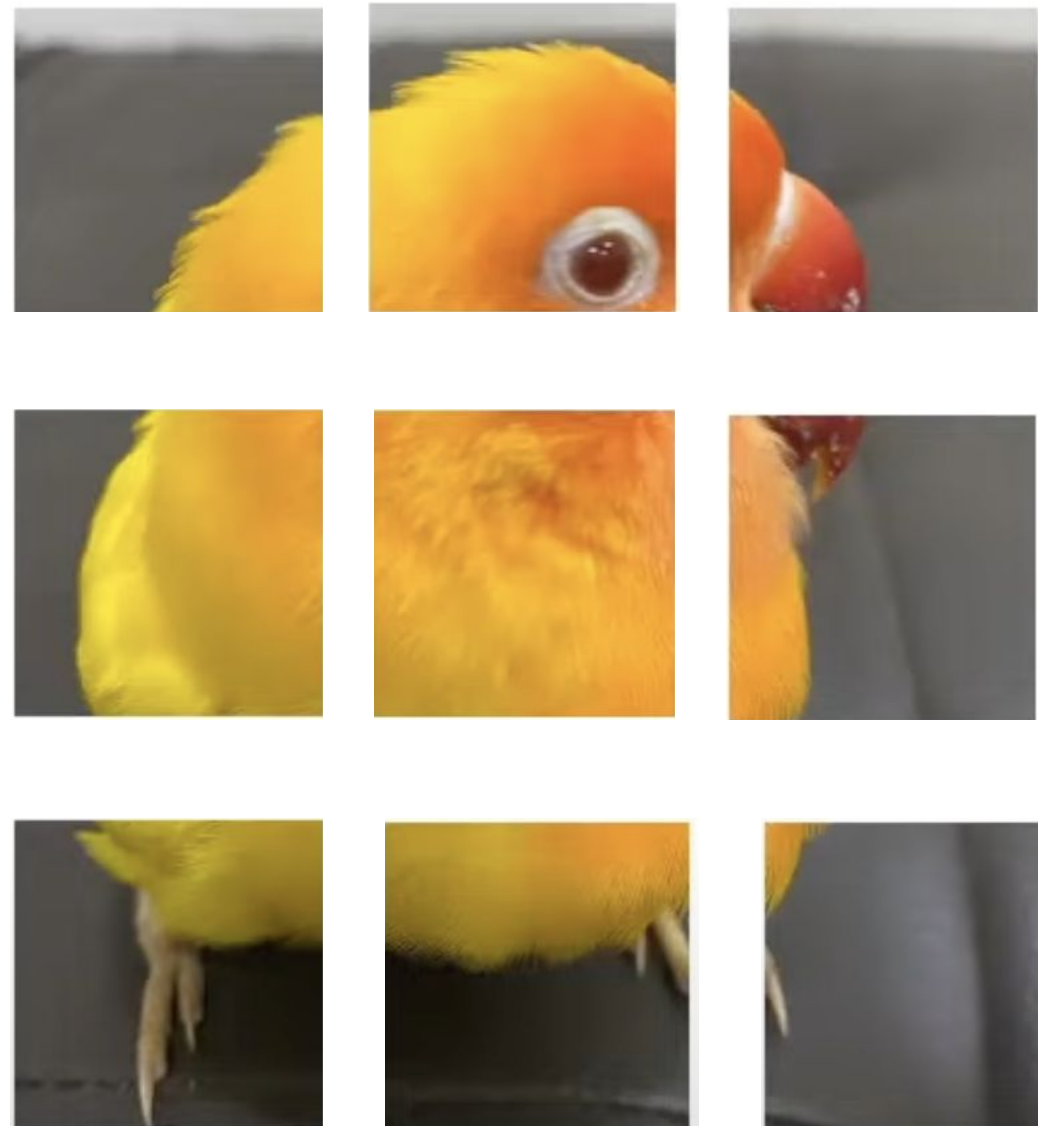


Section 06

Appendix

Section #

Patch slide

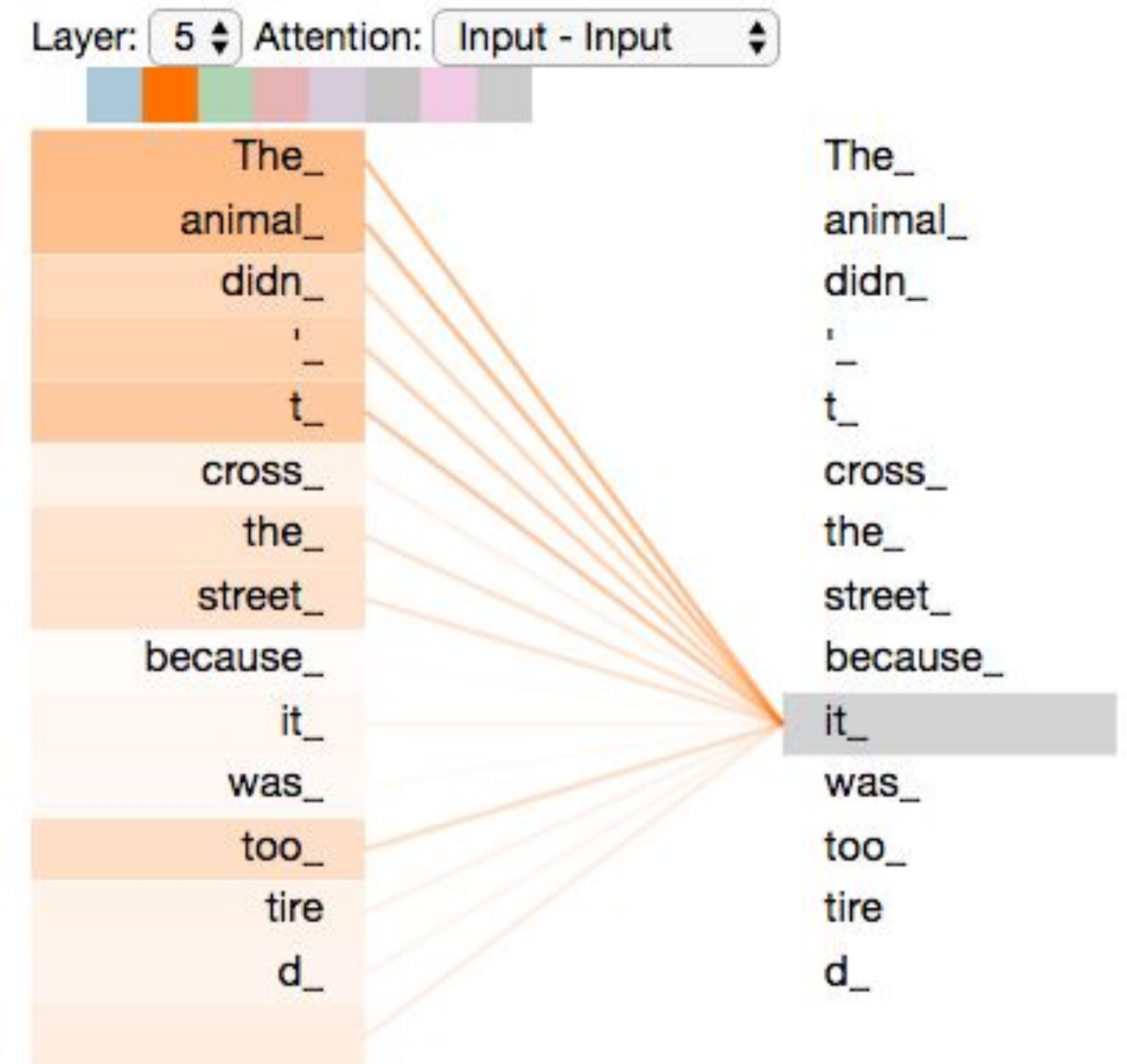


Ref: [An image is worth 16 x 16 words: Transformers for Image Recognition at Scale](#)

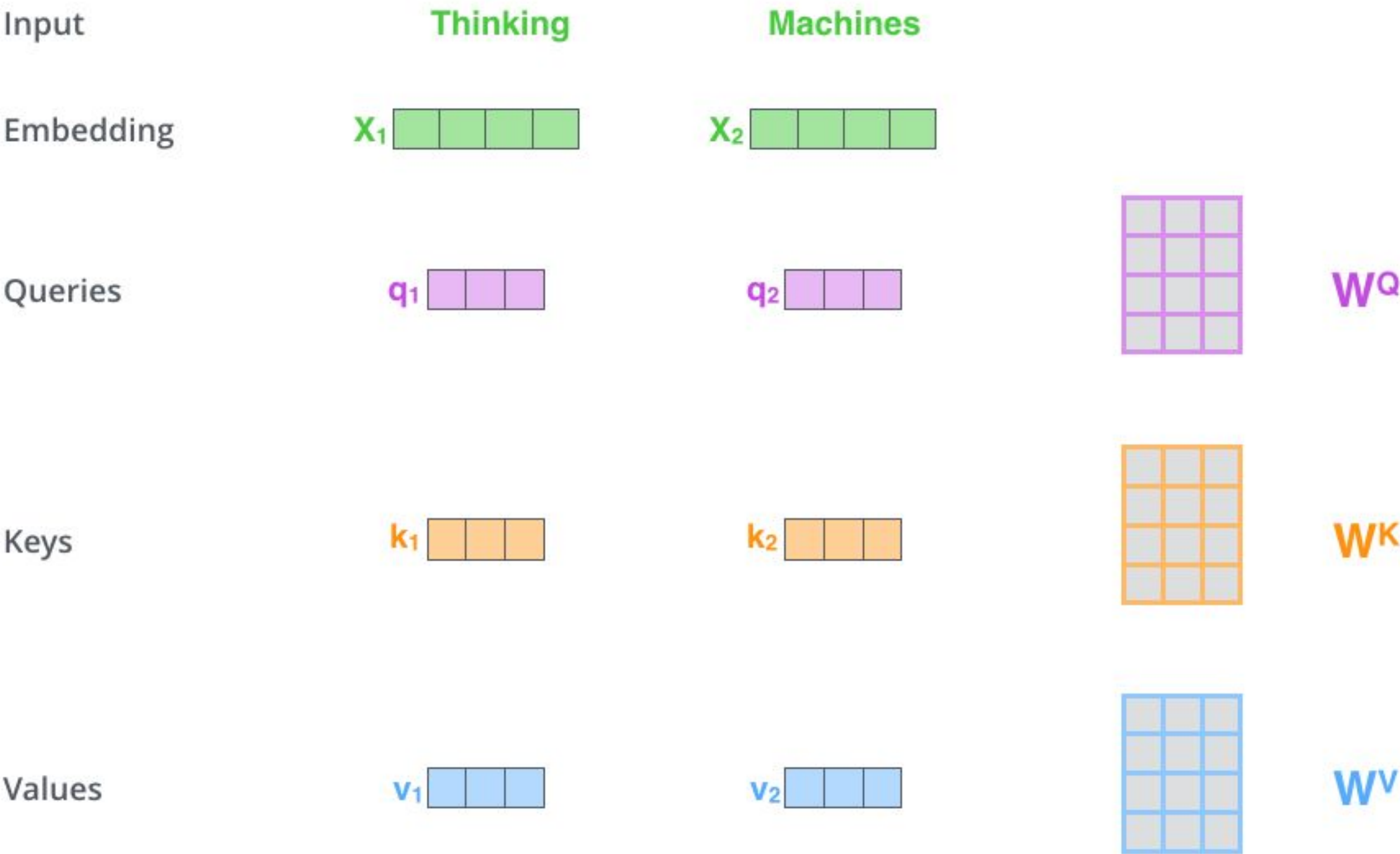
Self Attention

Translate : “The animal didn't cross the street because it was too tired”.

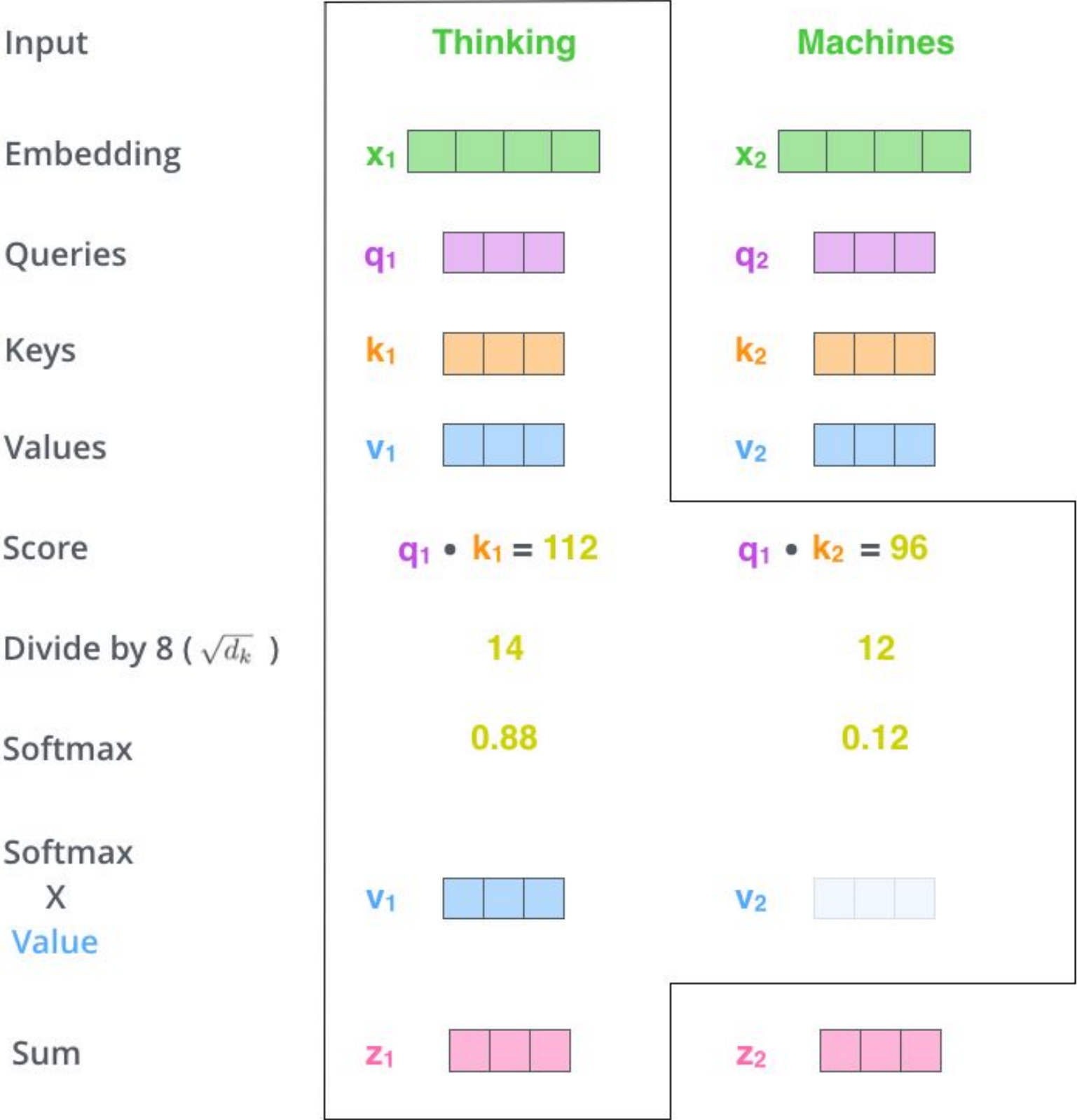
What does “it” in the above sentence refer to? Animal or Street?



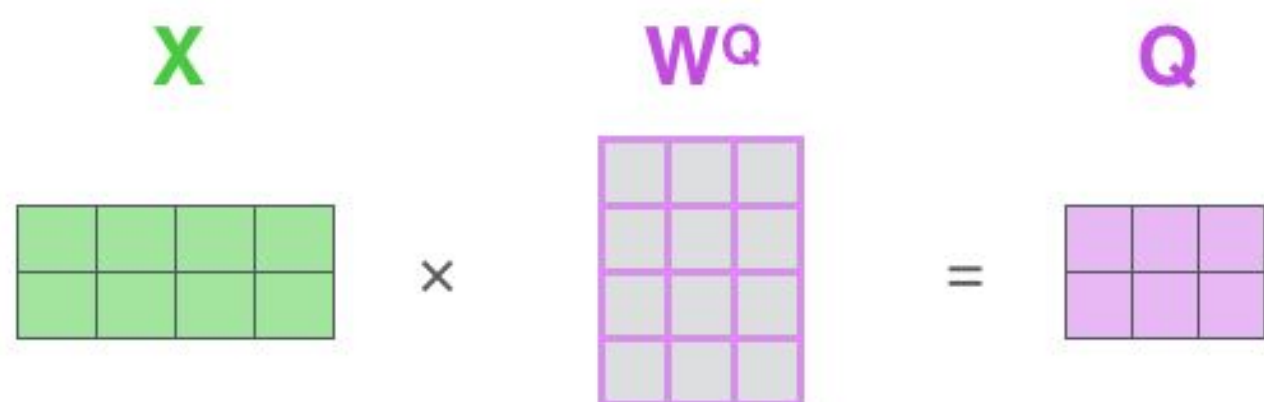
Self Attention in detail



Self-Attention: output

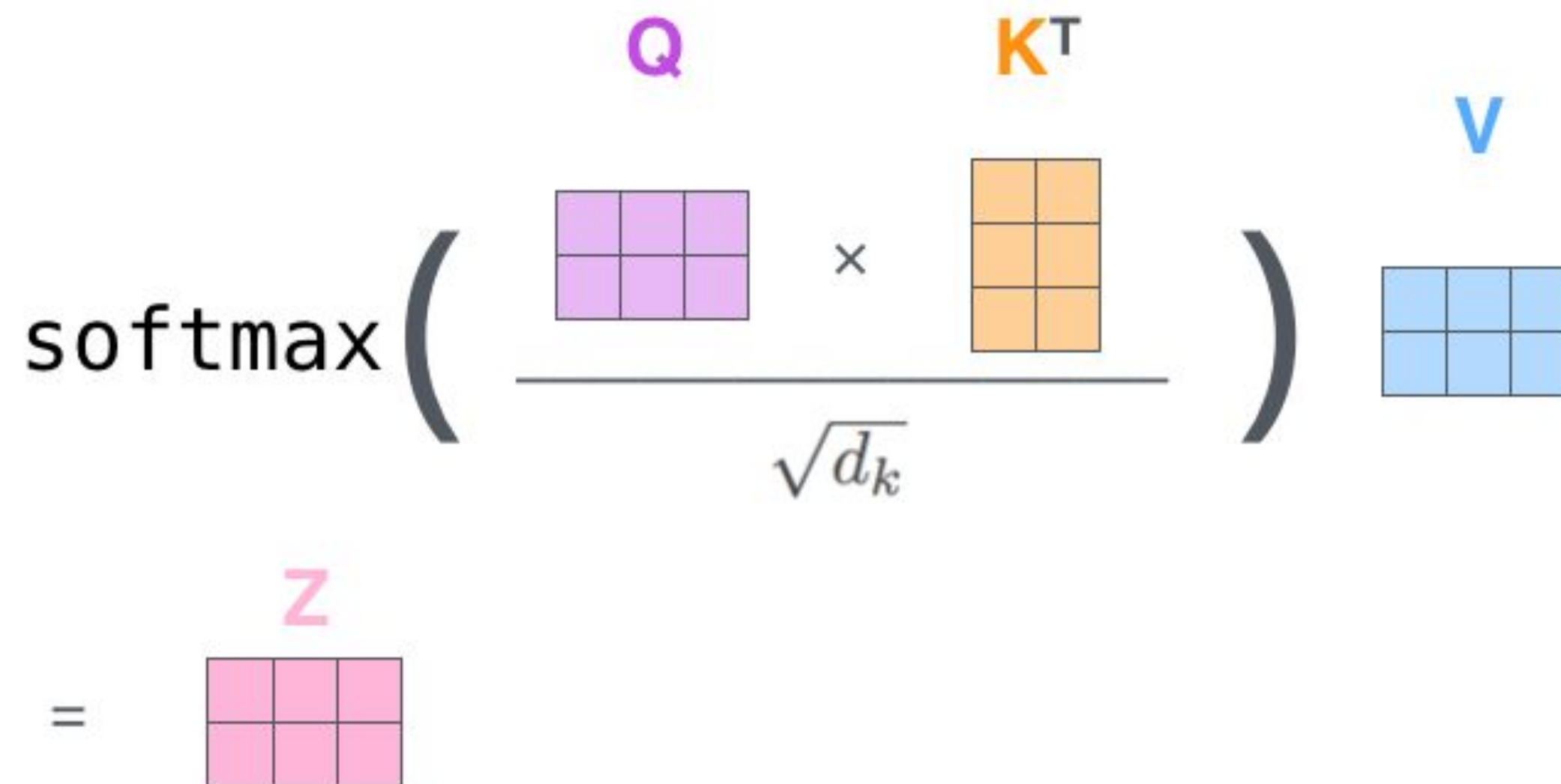


Matrix Calculation of Self-Attention

$$\mathbf{X} \times \mathbf{W}^Q = \mathbf{Q}$$


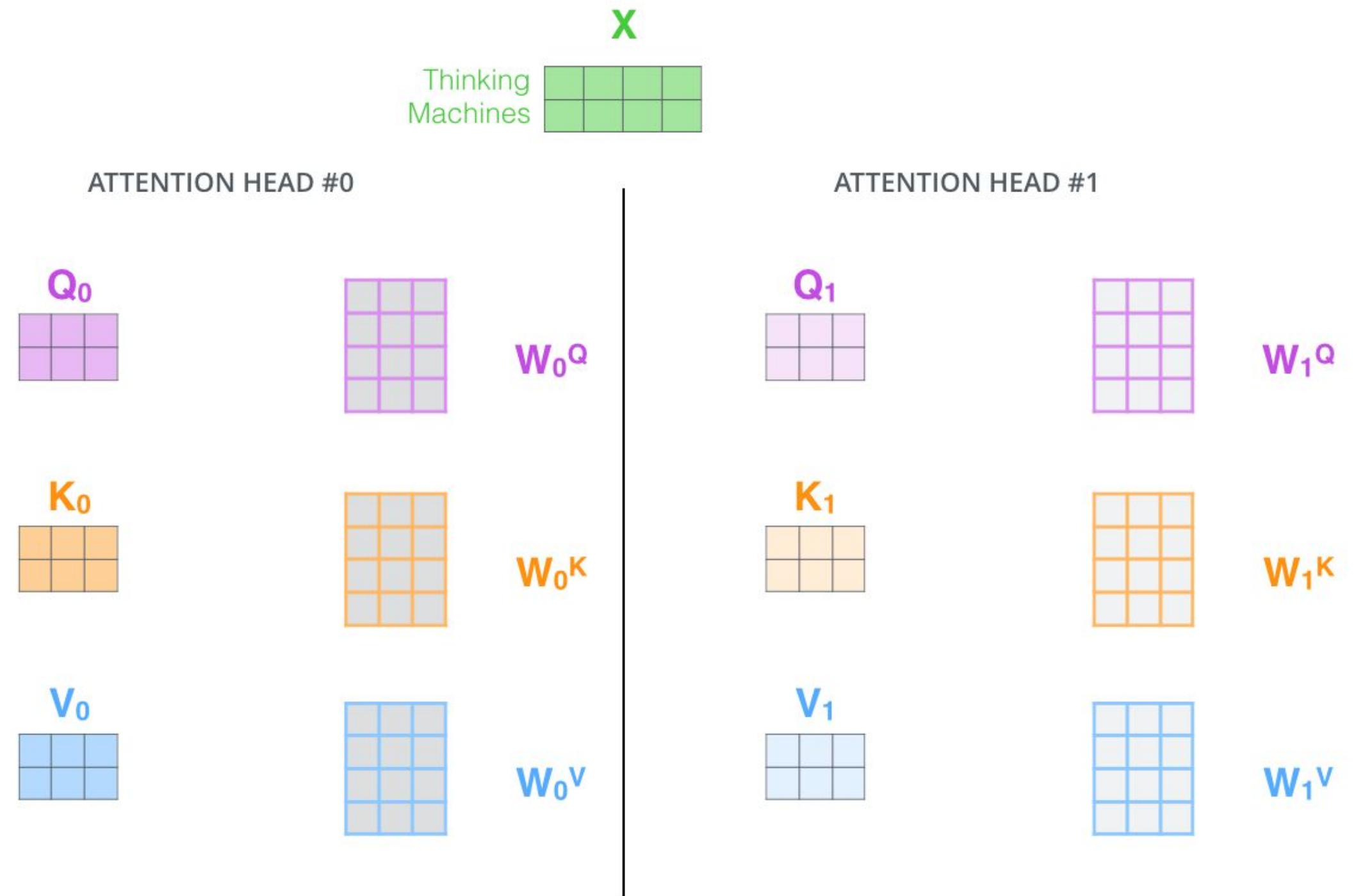
$$\mathbf{X} \times \mathbf{W}^K = \mathbf{K}$$


$$\mathbf{X} \times \mathbf{W}^V = \mathbf{V}$$


$$\text{softmax}\left(\frac{\mathbf{Q} \times \mathbf{K}^T}{\sqrt{d_k}}\right) \times \mathbf{V} = \mathbf{Z}$$


Multi-head attention

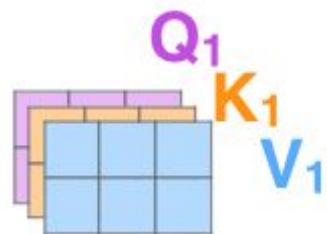
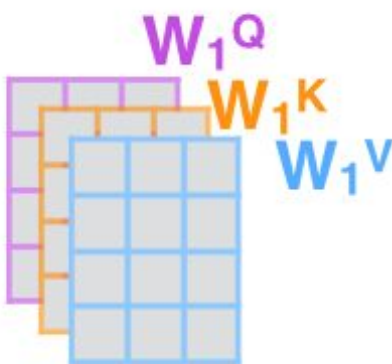
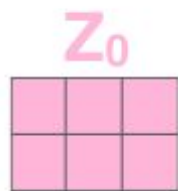
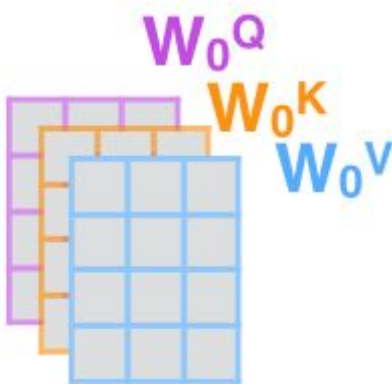
- Model ability to focus on different positions.
- Input embeddings/vectors from lower encoders/decoders into a different representation subspace.



Multi-head attention ..

- 1) This is our input sentence*
- 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

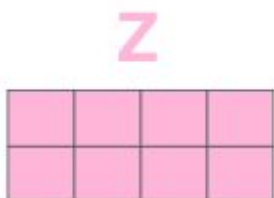
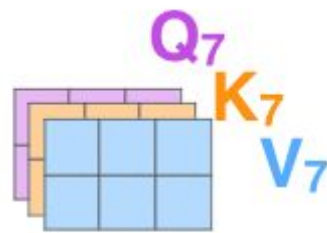
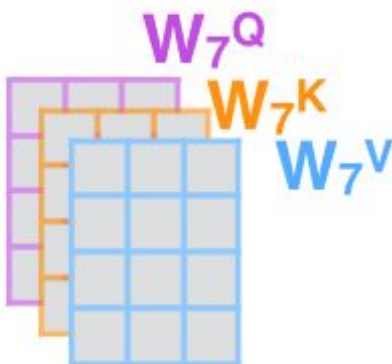
Thinking
Machines



...

...

...

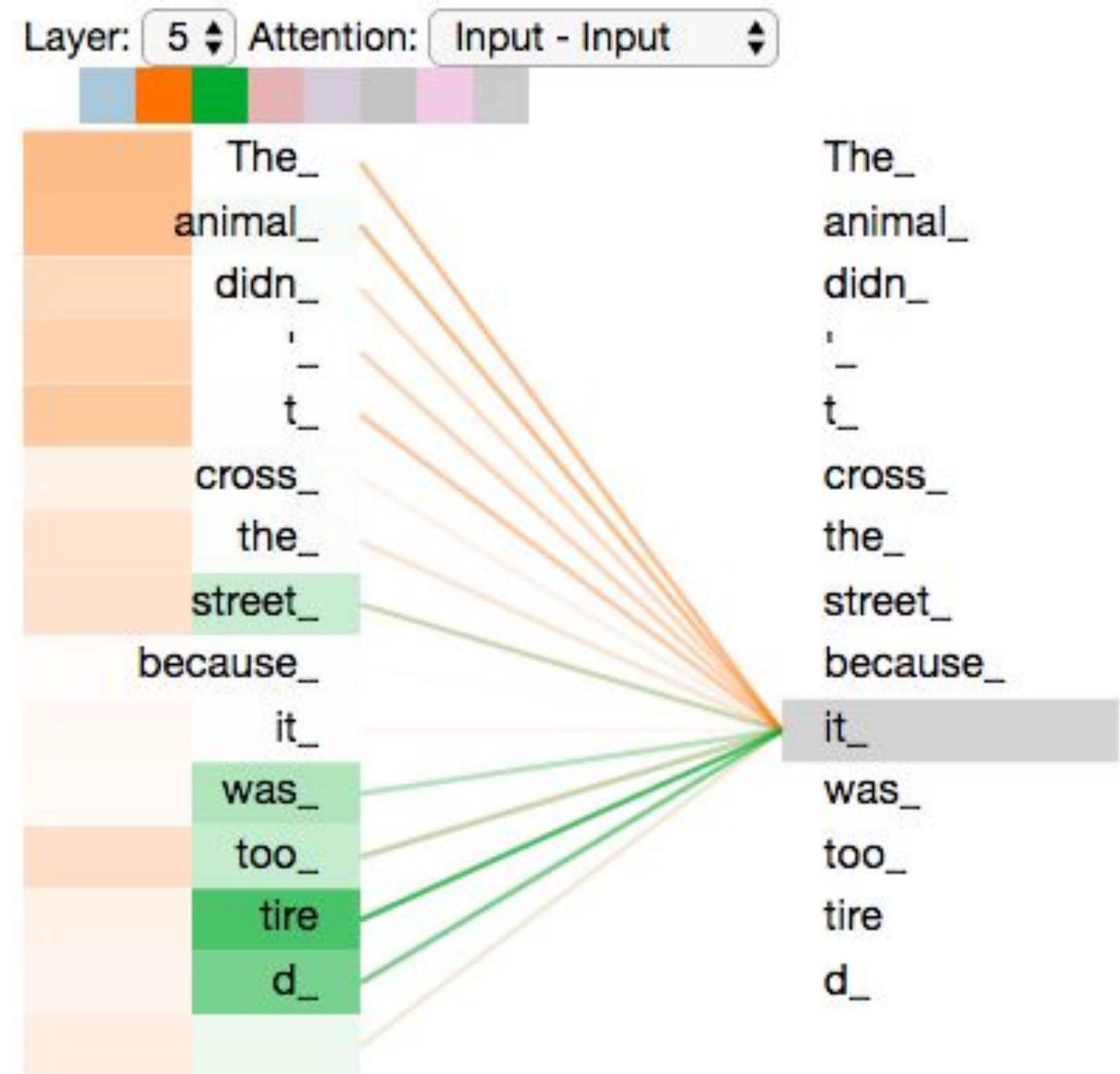


* In all encoders other than #0, we don't need embedding. We start directly with the output of the encoder right below this one



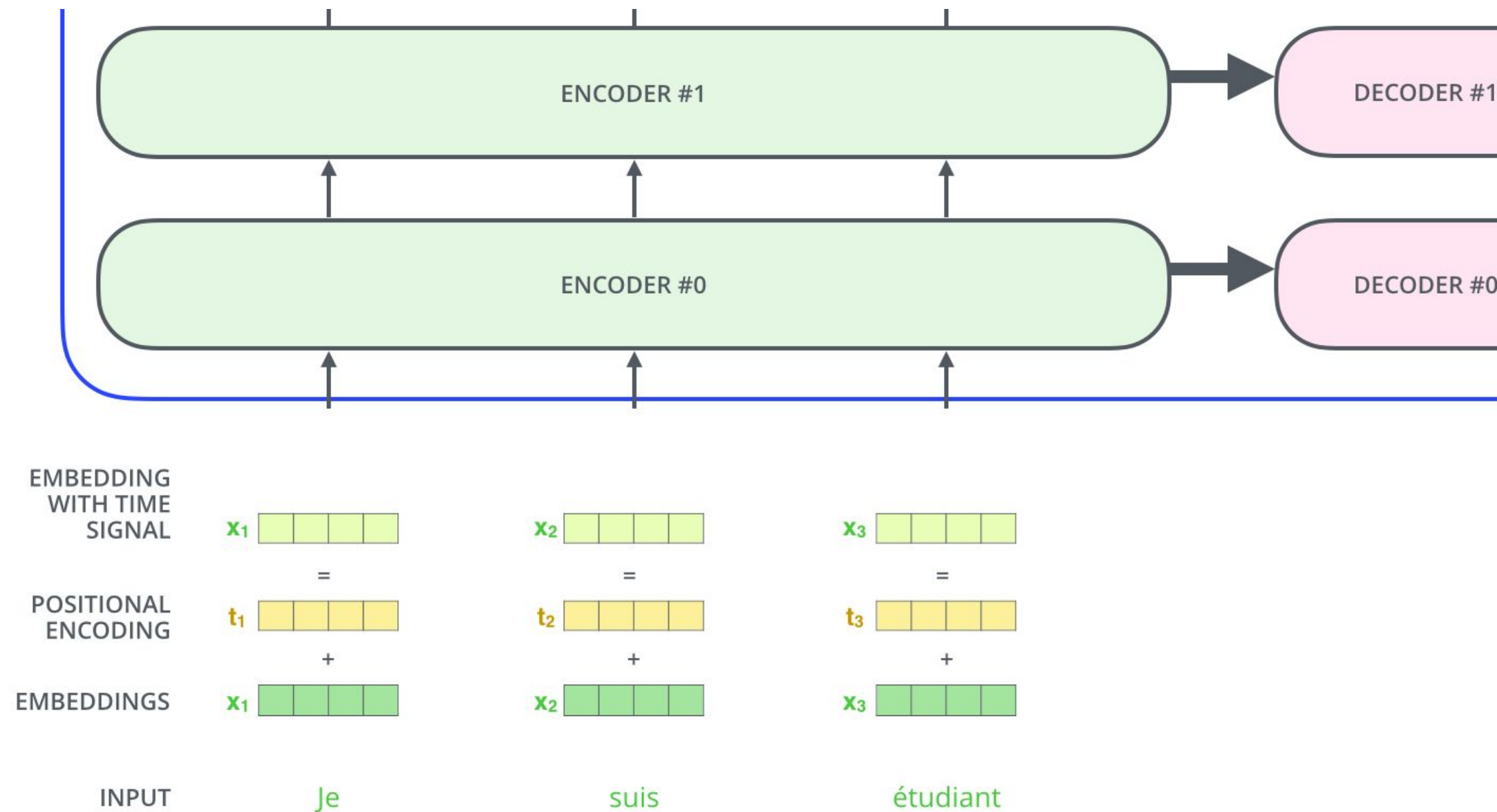
Interpreting - Multi Head Attention

- As we encode the word "it", one attention head is focusing most on "the animal", while another is focusing on "tired".
- In a sense, the model's representation of the word "it" bakes in some of the representation of both "animal" and "tired".



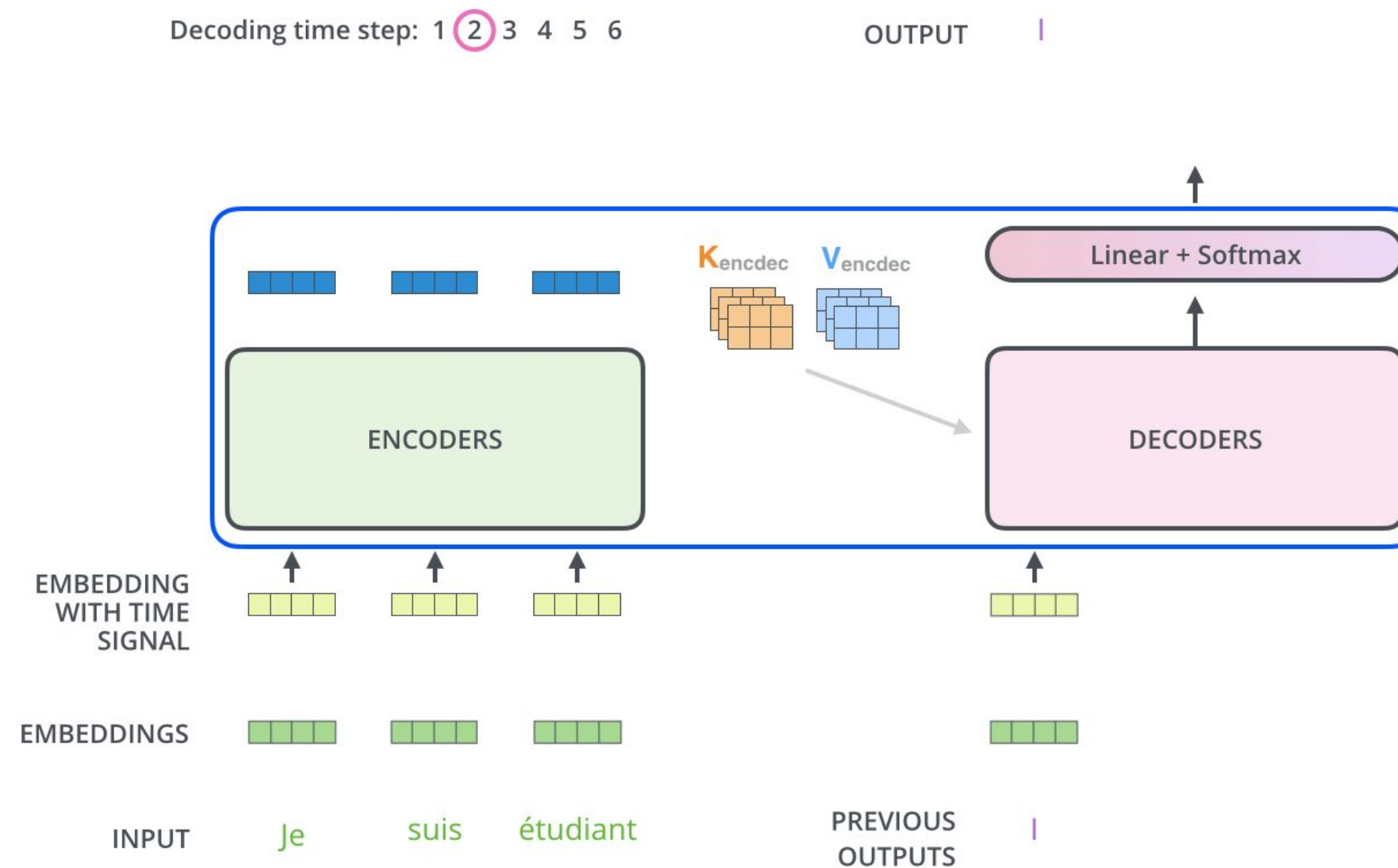
Positional Embeddings

- So far, order of the words in the input sequence is missing.
- **Positional embeddings:** vectors which follow a specific pattern that the model learns, which helps it determine the position of each word, or the distance between different words in the sequence.



Transformer Architecture

The decoder attends on the encoder's output and its own input (self-attention) to predict the next word.



BERT: Bidirectional Encoder Representations from Transformers

