

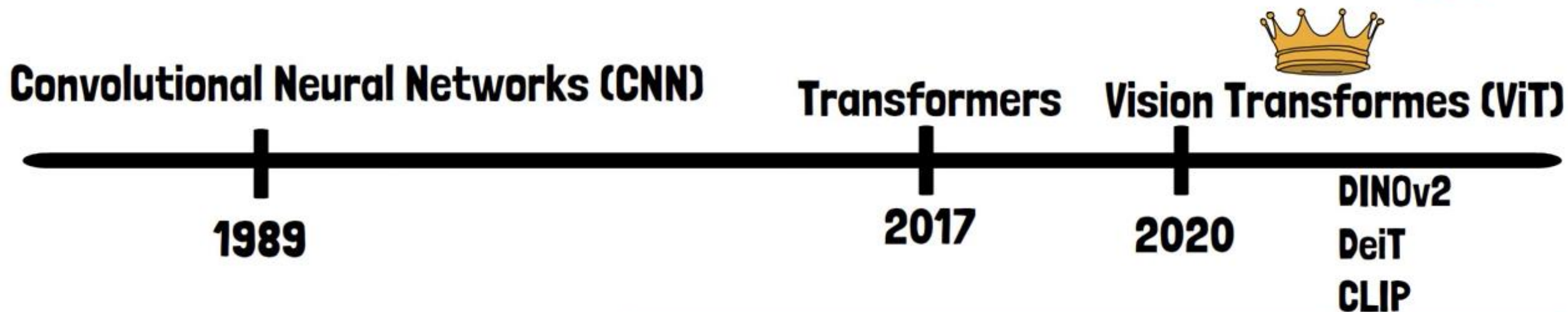
# Vision *Transformer*





# Vision Transformers

Computer vision dominating architecture



AN IMAGE IS WORTH 16x16 WORDS:  
TRANSFORMERS FOR IMAGE RECOGNITION AT SCALE

Alexey Dosovitskiy<sup>\*,†</sup>, Lucas Beyer<sup>\*</sup>, Alexander Kolesnikov<sup>\*</sup>, Dirk Weissenborn<sup>\*</sup>,  
Xiaohua Zhai<sup>\*</sup>, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer,  
Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, Neil Houlsby<sup>\*,†</sup>


<sup>\*</sup>equal technical contribution, <sup>†</sup>equal advising

Google Research, Brain Team

{adosovitskiy, neilhoulby}@google.com

# Using a Transformer as-is?

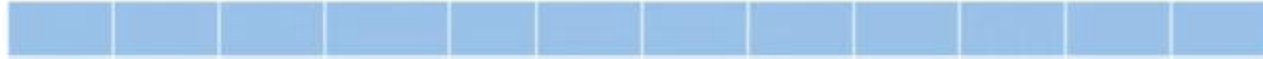
 **Transformer Self-Attention** – each token attends all other tokens

 **Feeding an image to transformer** – each pixel attends all other pixels

**256 X 256**  
image

**512 X 512**  
image

Input sequence

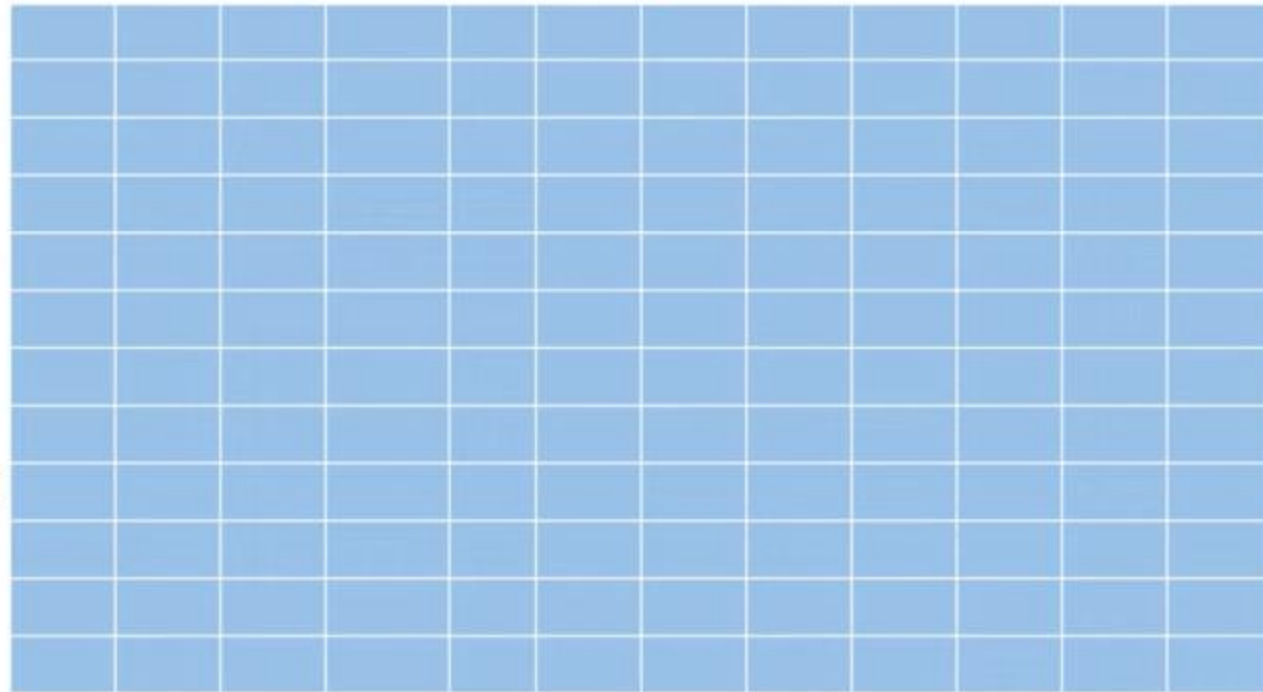


**65k pixels**

**262k pixels**



Attention matrix



**65k X 65k**  
attention  
matrix

**262k X 262k**  
attention  
matrix



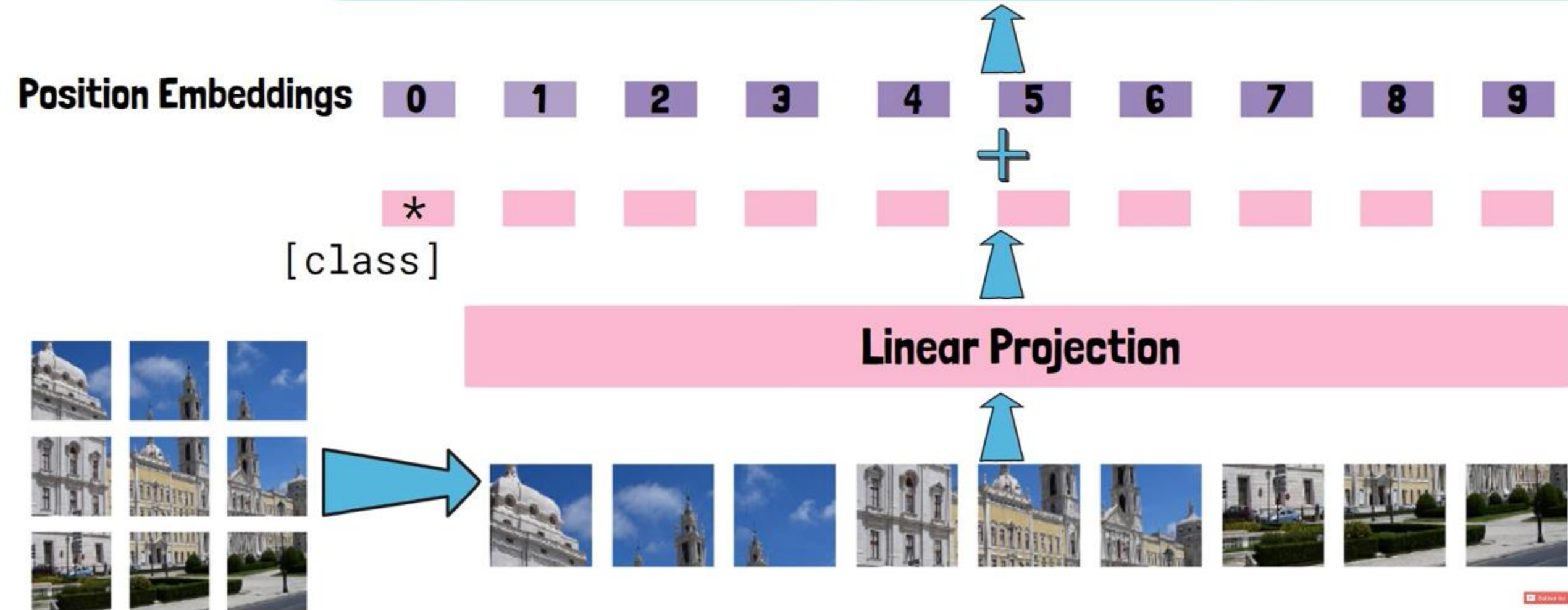
Quadratic dependency  
on the sequence length

 **Handling images as-is results in too long sequence lengths**

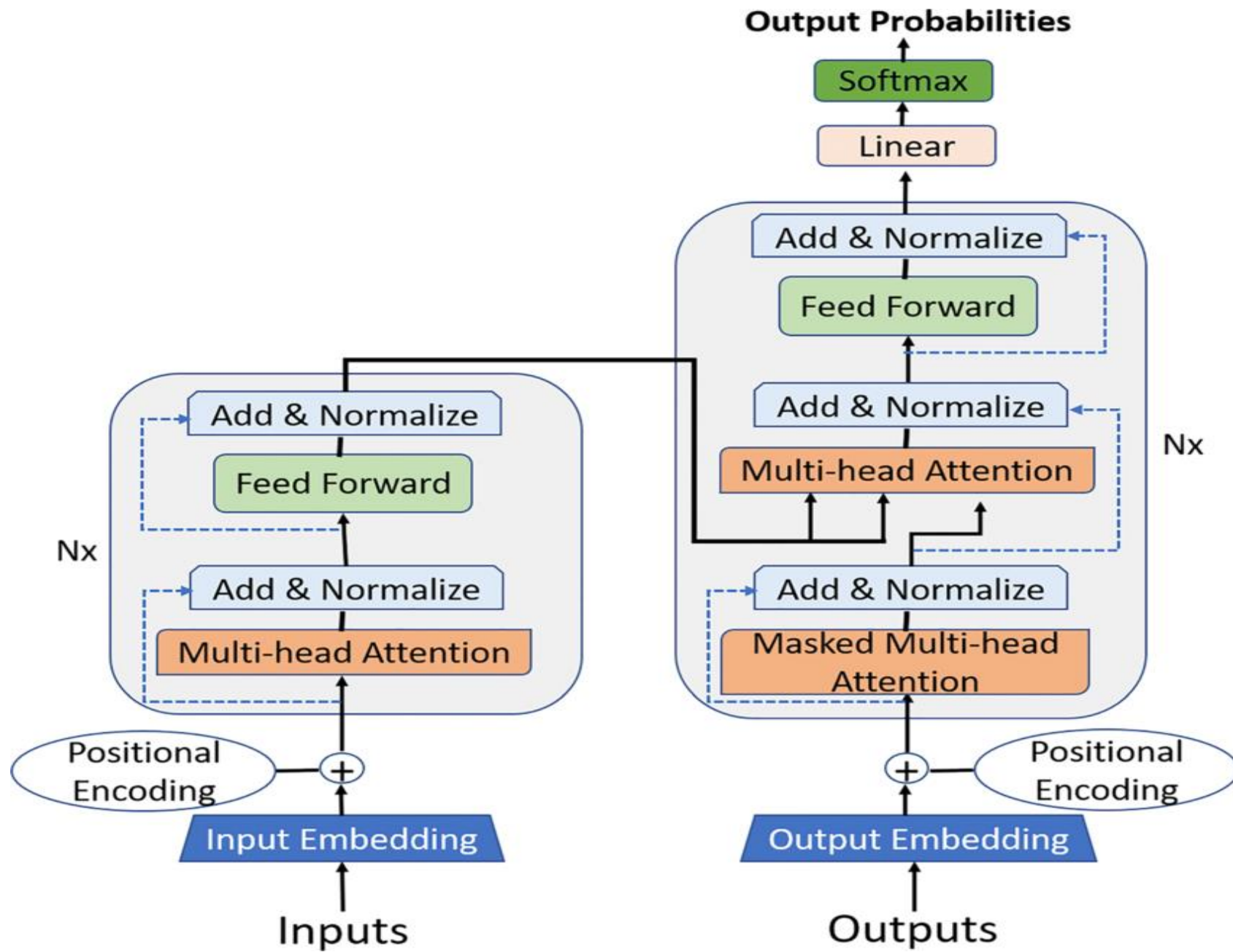


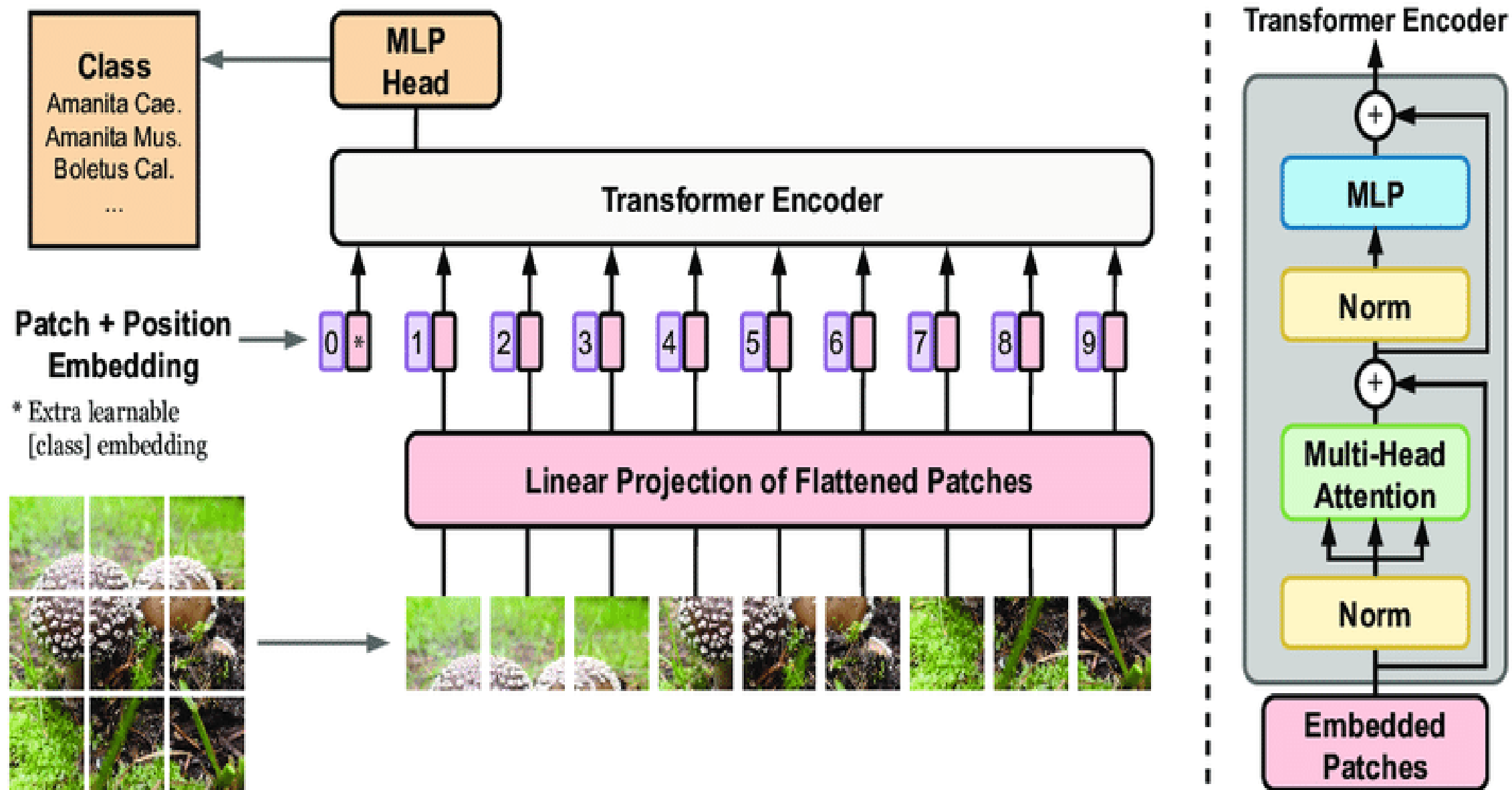
# How Vision Transformer Works?

## Transformer









# Vision Transformers (ViTs)

224\*224



patch size =  $16 \times 16$

Number of patches =  $(224 / 16) * (224 / 16) = 14 * 14 = 196 \text{ patches}$

Stride = 16

# Vision Transformers (ViTs)

196 patches



Total patches are 196 and each patch is having size of  $16 \times 16 \times 3$ , assuming an RGB image

Flatten the patches from 2D to 1D



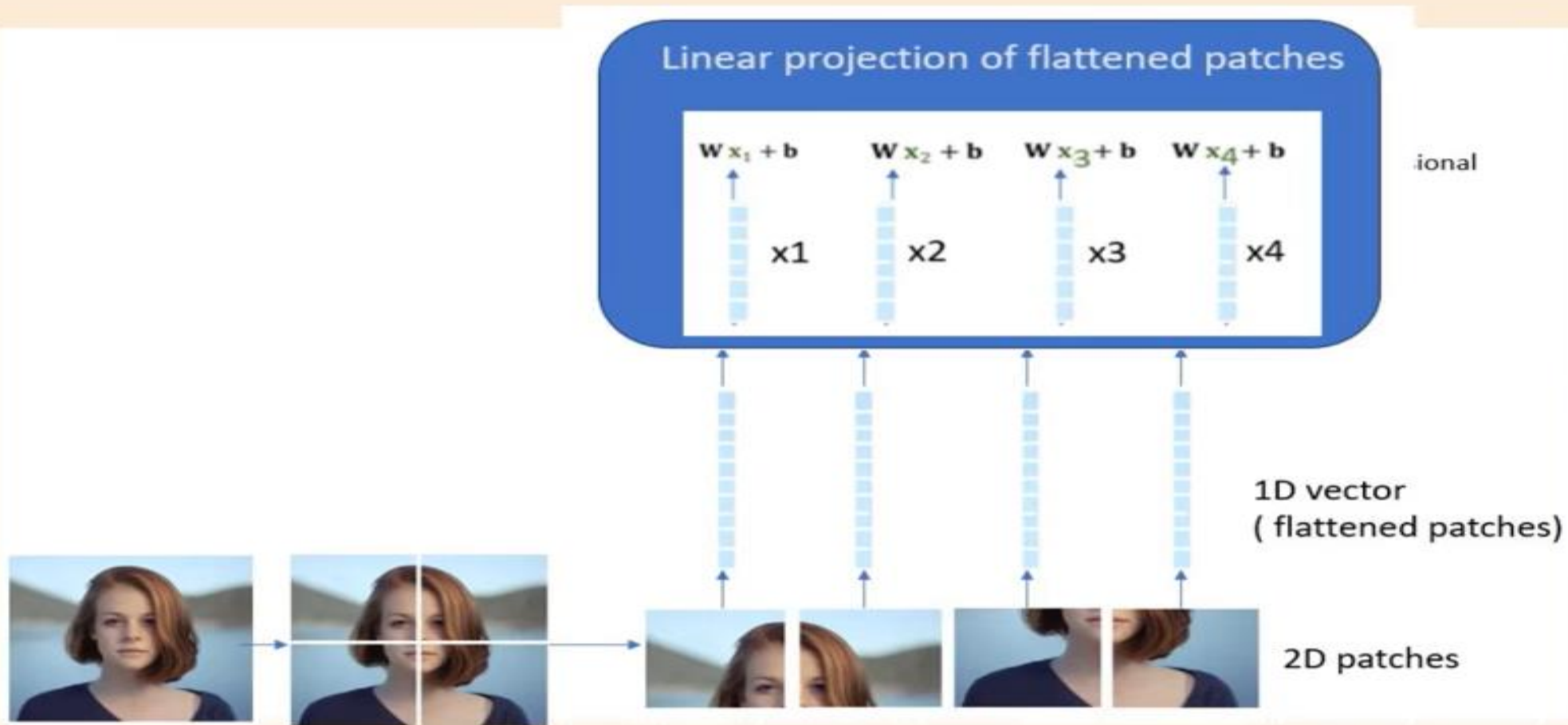
Each patch is flattened into a 1D vector



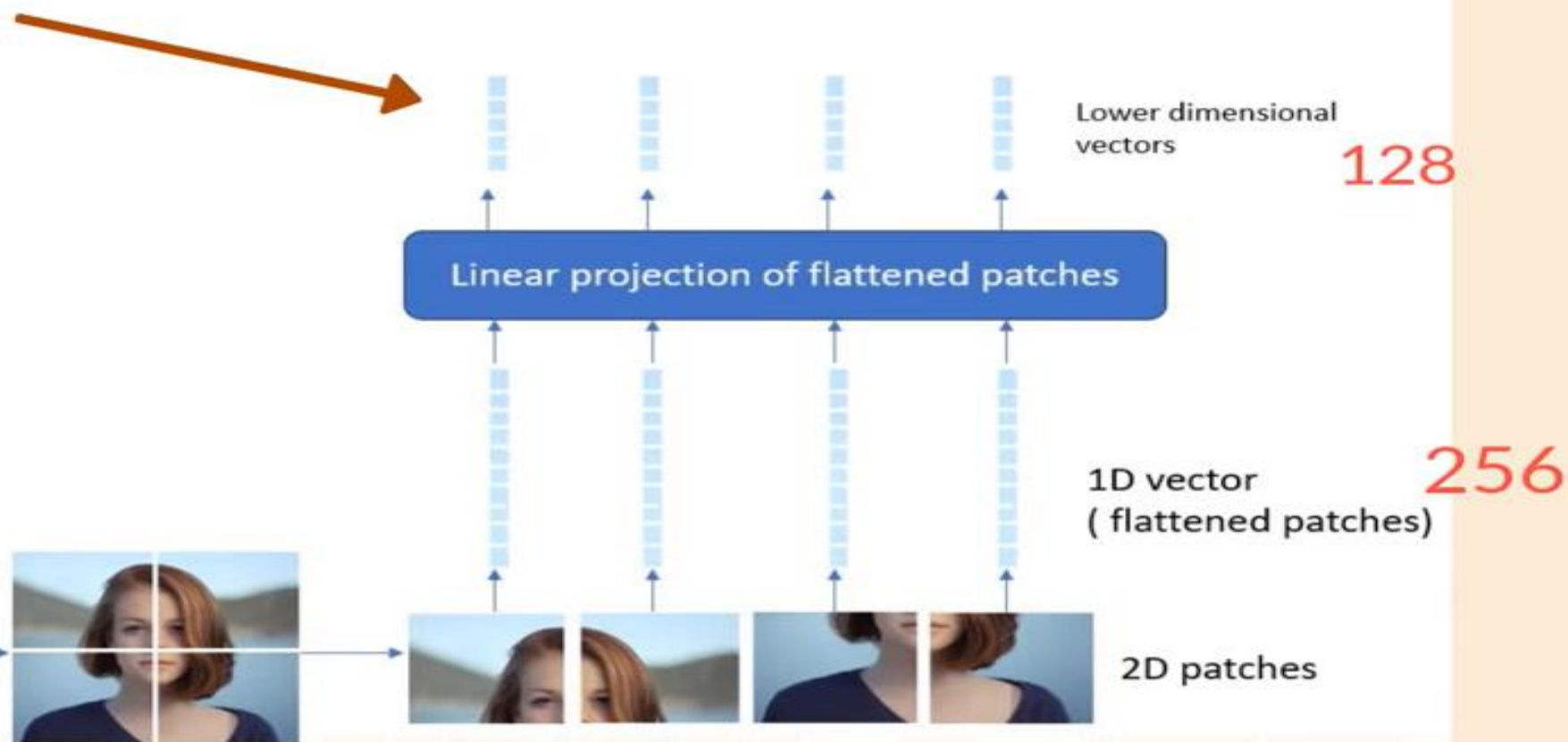
1D vector



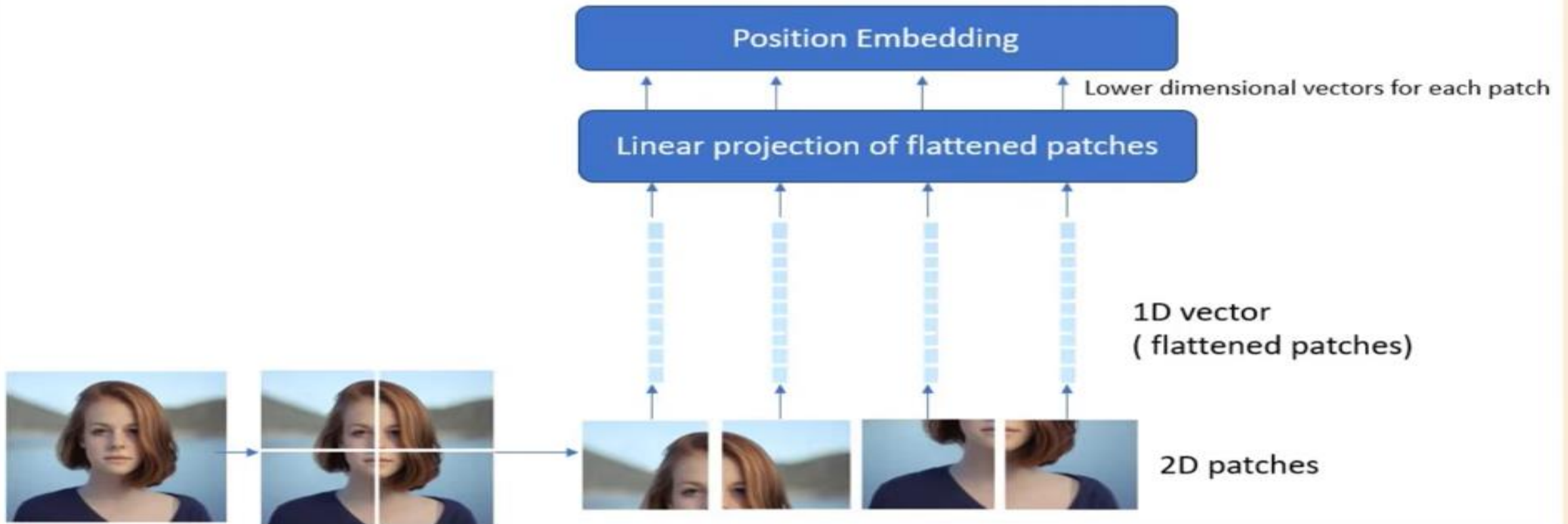
# Vision Transformers (ViTs)



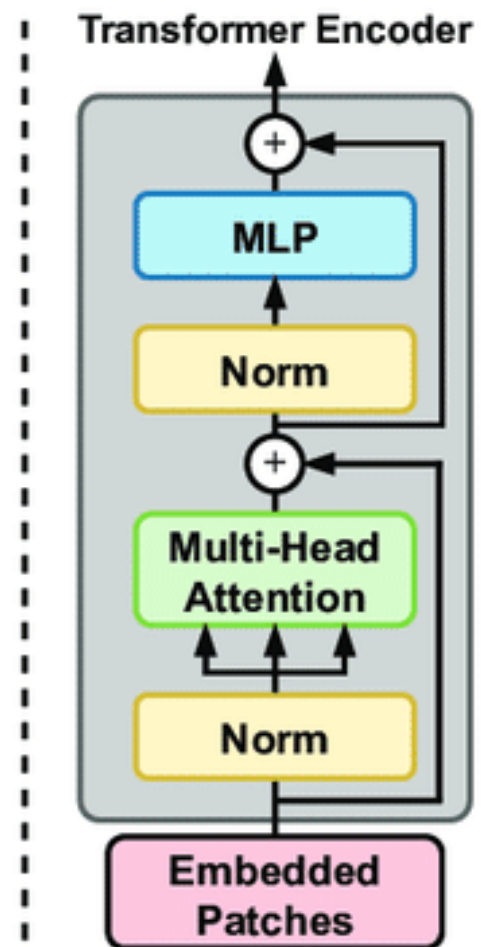
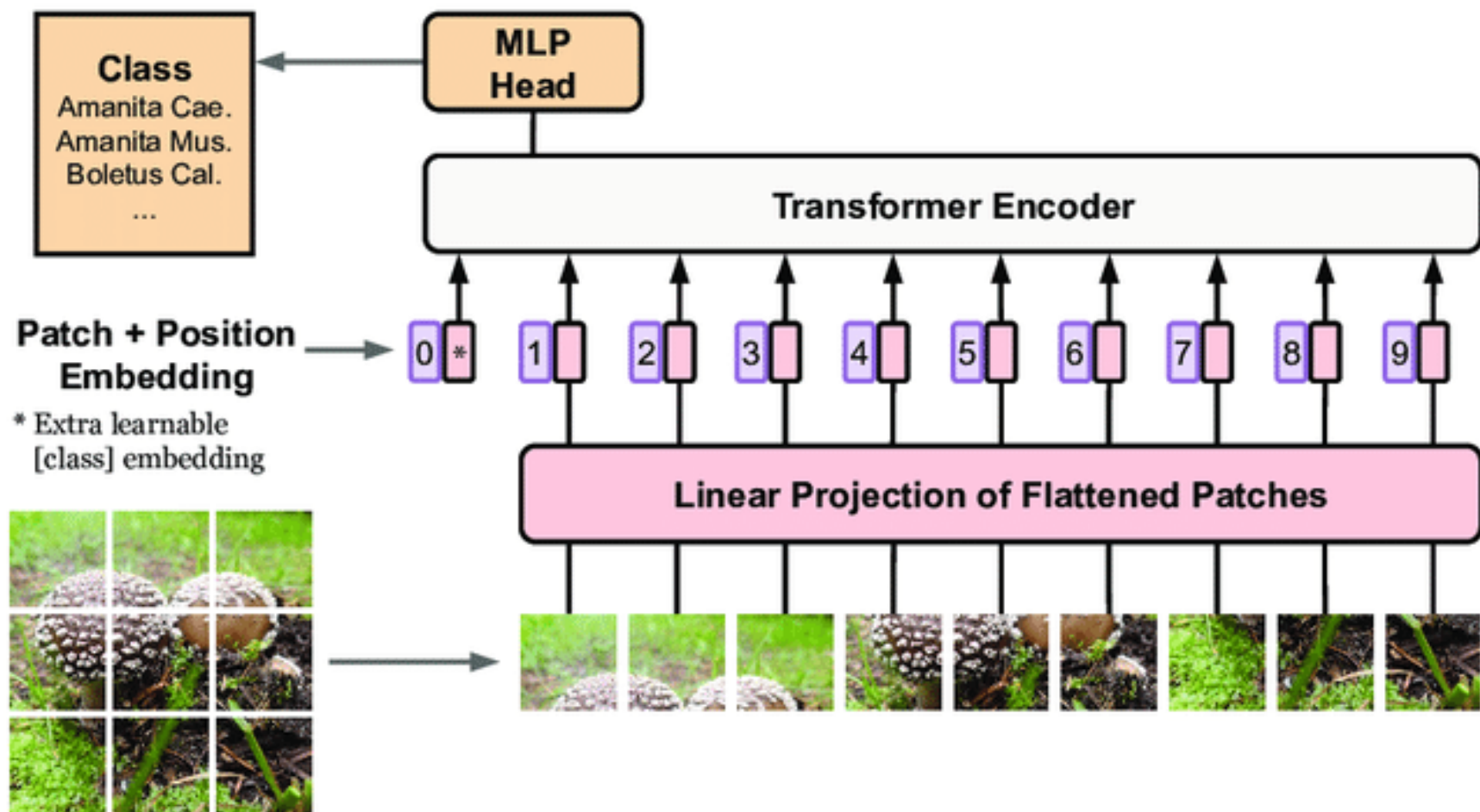
# Vision Transformers (ViTs)

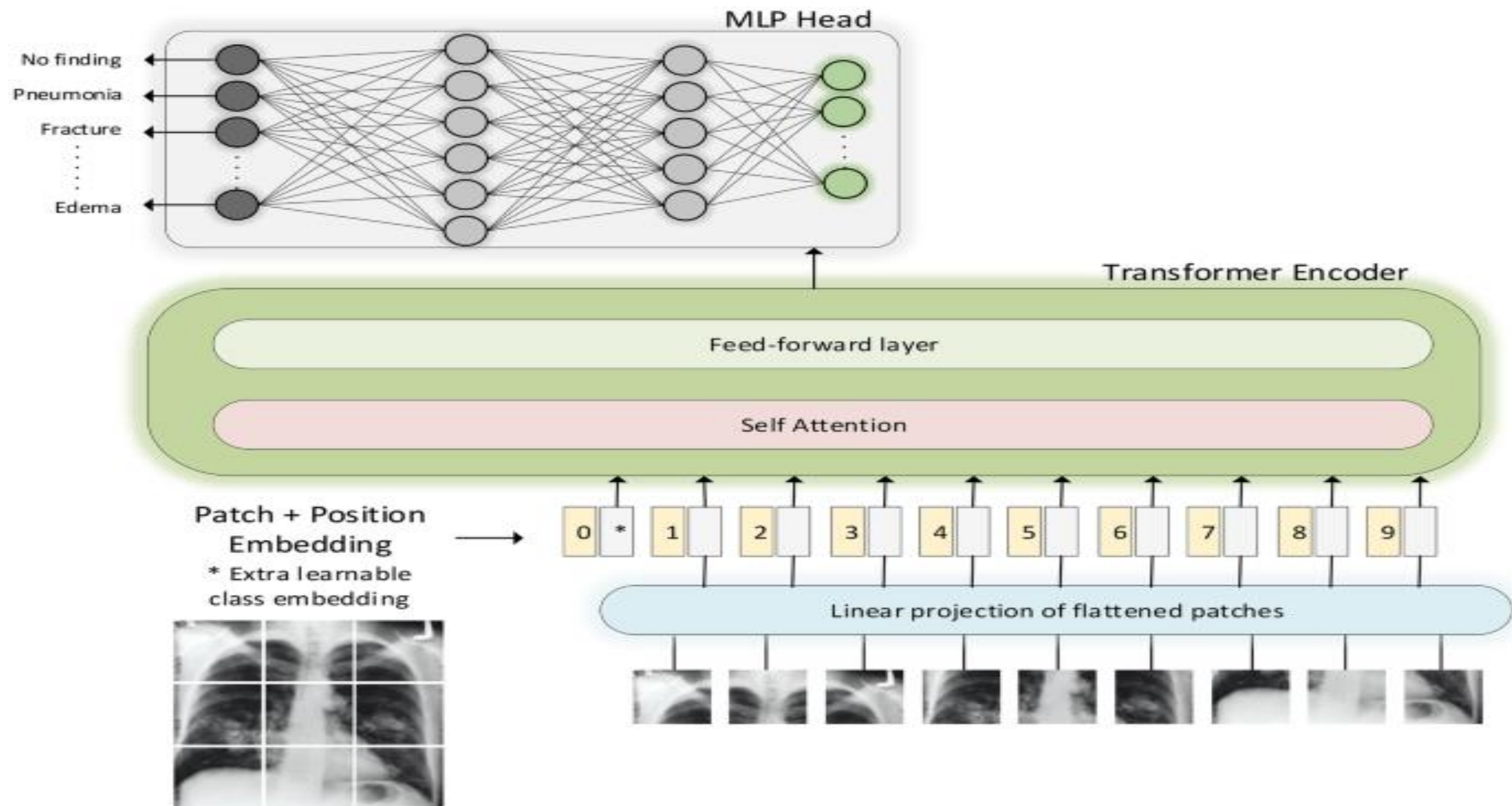


# Vision Transformers (ViTs)









# CNN – high inductive bias

A lot of guidance is given to the model

Input

## Inductive Bias



Transformer – each token can attend to any token in every layer

Much less guidance

