

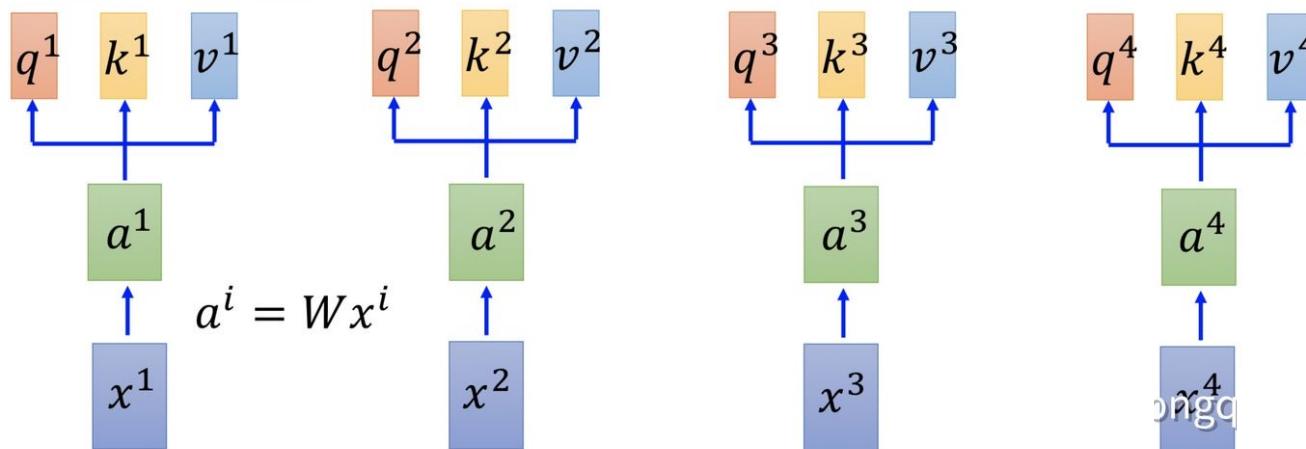


DIVE INTO  
VISION  
TRANSFORMERS

Tianyi

# Self-attention

<https://arxiv.org/abs/1706.03762>



$q$ : query (to match others)

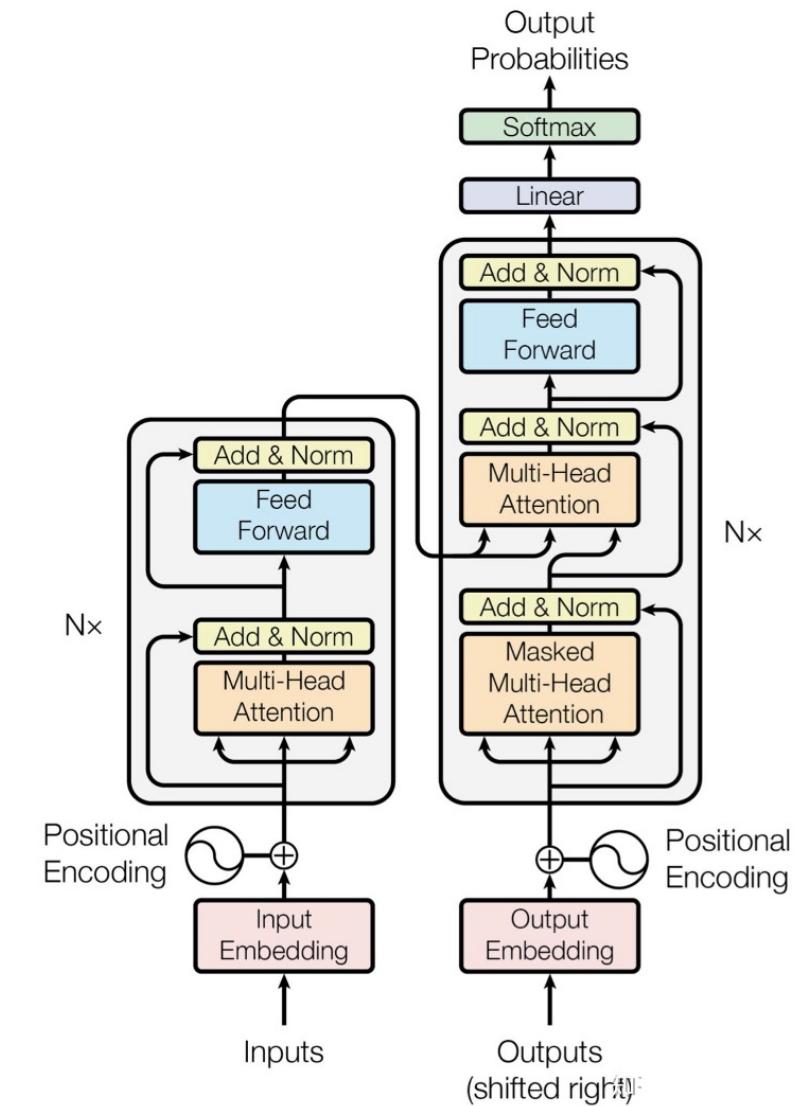
$$q^i = W^q a^i$$

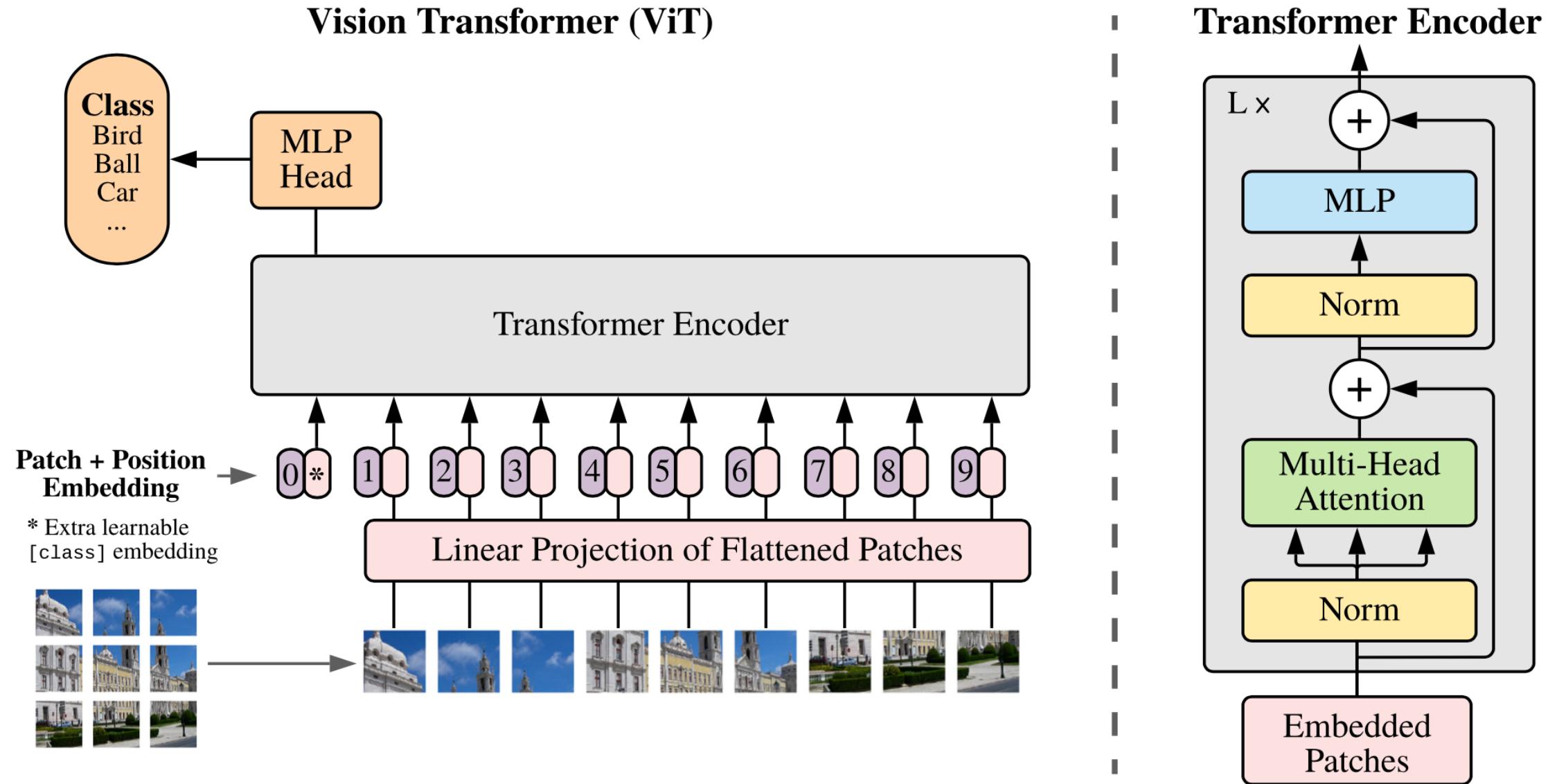
$k$ : key (to be matched)

$$k^i = W^k a^i$$

$v$ : information to be extracted

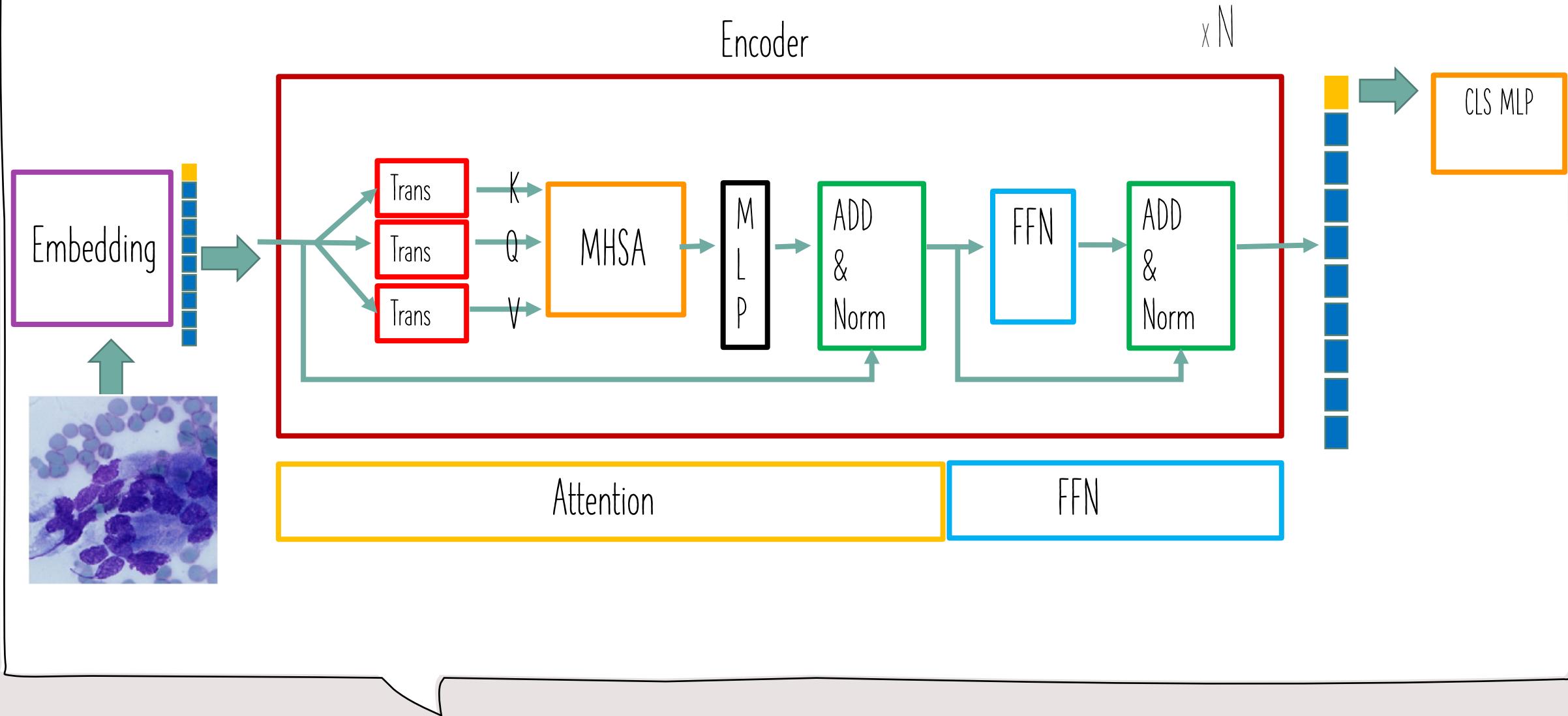
$$v^i = W^v a^i$$





# ViT Vision Transformer

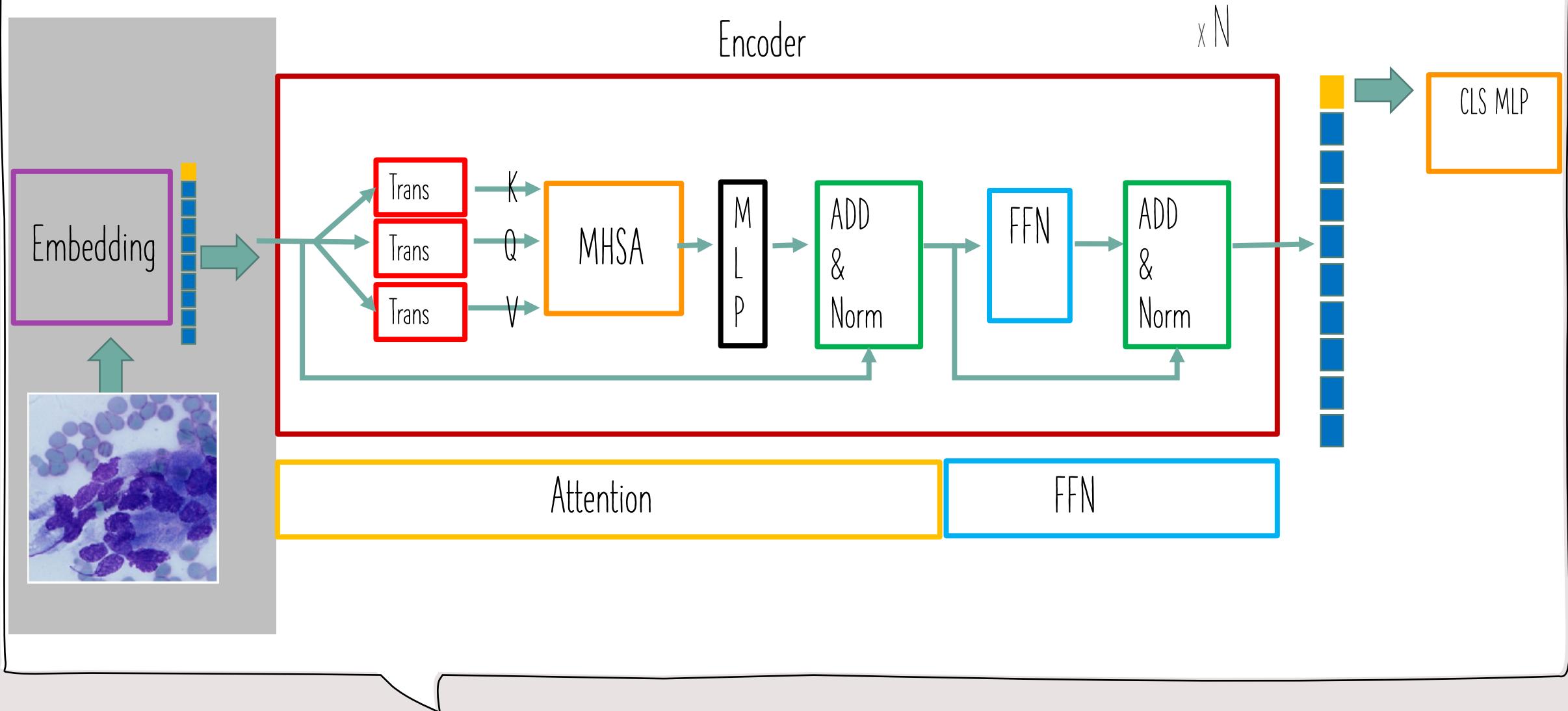
AN IMAGE IS WORTH 16X16 WORDS: TRANSFORMERS FOR IMAGE RECOGNITION AT SCALE  
Arxiv 2010.11929



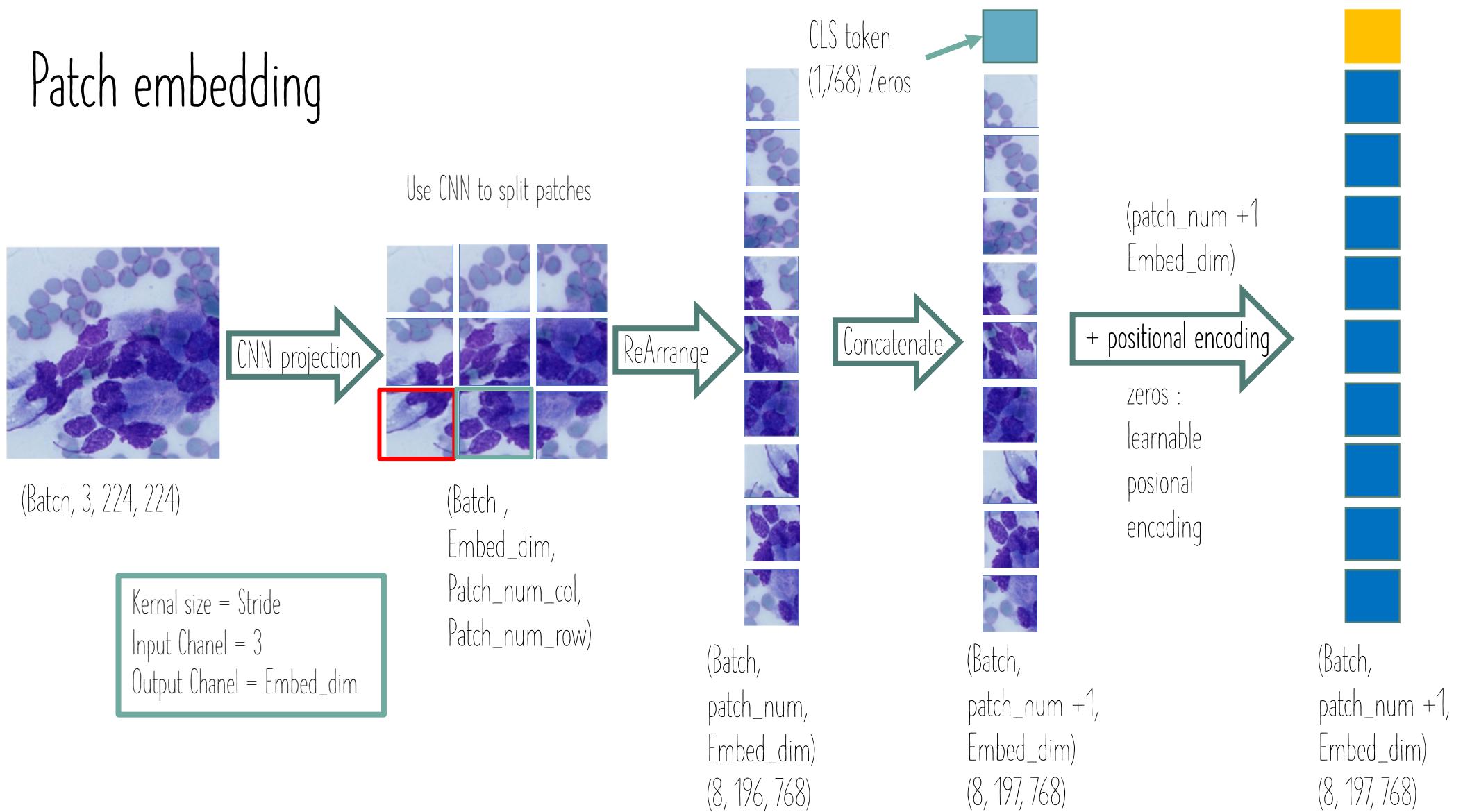


# ViT Vision Transformer

AN IMAGE IS WORTH 16X16 WORDS: TRANSFORMERS FOR IMAGE RECOGNITION AT SCALE  
Arxiv 2010.11929

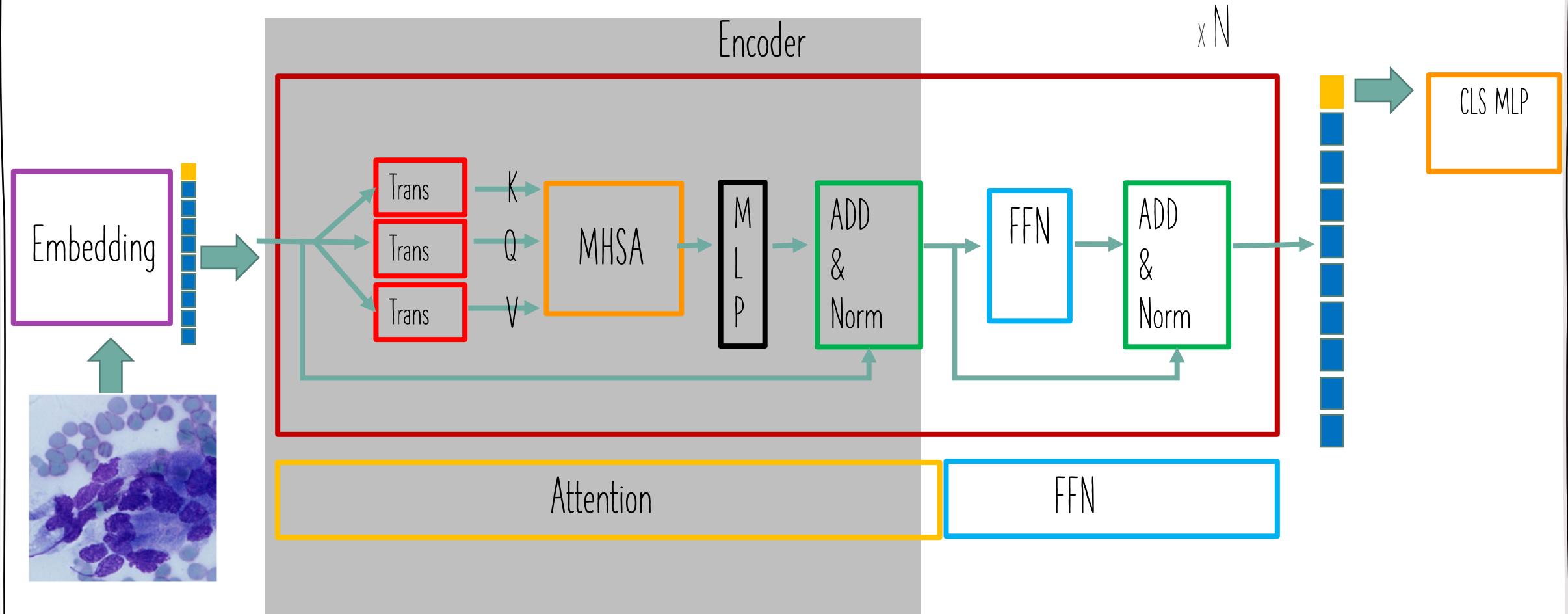


# Patch embedding



# ViT Vision Transformer

AN IMAGE IS WORTH 16X16 WORDS: TRANSFORMERS FOR IMAGE RECOGNITION AT SCALE  
Arxiv 2010.11929



# Self duplicate (MLP)

Each token vector of  $D$  enter FC **seperately**

(D) (D) (D) (D)

... (D) (D) (D)

(N,D)

Why transformer can fit  
different length ?

One FC to project all tokens

Self-duplicate projection obtain each q,k,v token **togather**

(D)



q (N,D)



k (N,D)

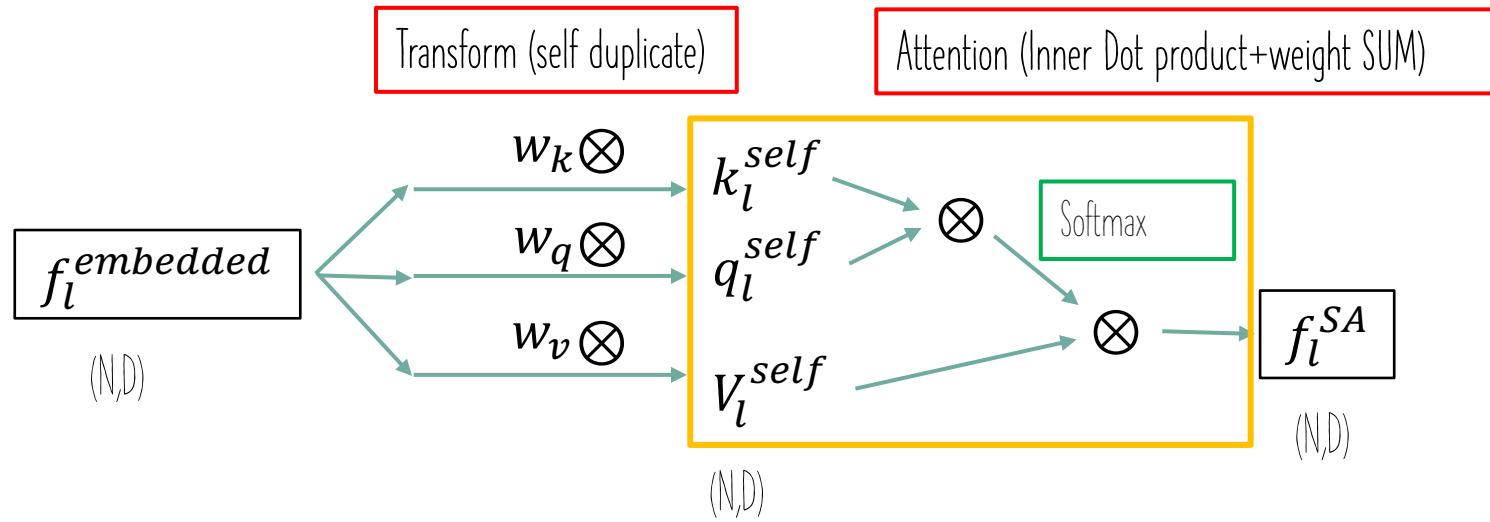


v (N,D)

What is q,k,v ?

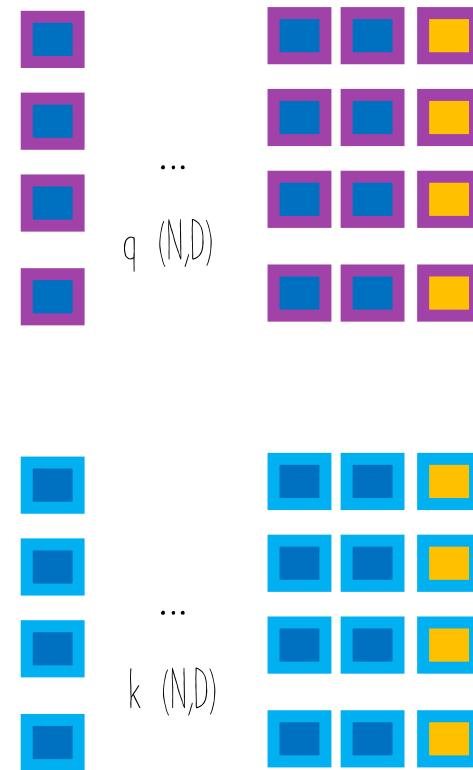
Tokens projected to different representation latent space

## SA self attention

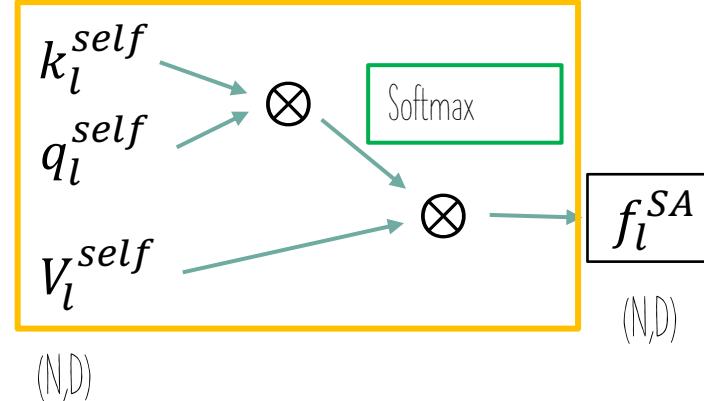


For a given embedded patch  $f_l^{\text{embedded}}$ ,  $f_l^{SA} = \text{SoftMax}(q_l^{\text{self}}{}^T \cdot k_l^{\text{self}})^T \cdot V_l^{\text{self}}$ , where  $q_l^{\text{self}} = w_q \cdot f_l^{\text{embedded}}$ ,  $k_l^{\text{self}} = w_k \cdot f_l^{\text{embedded}}$ ,  $V_l^{\text{self}} = w_v \cdot f_l^{\text{embedded}}$ .

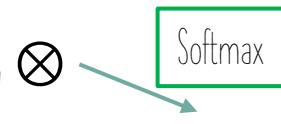
# SA self attention



Attention (Inner Dot product+weight SUM)



In place dot product

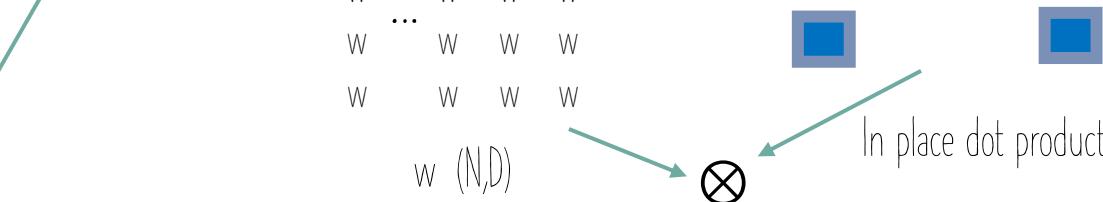


w      w      w      w  
w      w      w      w  
w      ...      w      w  
w      w      w      w

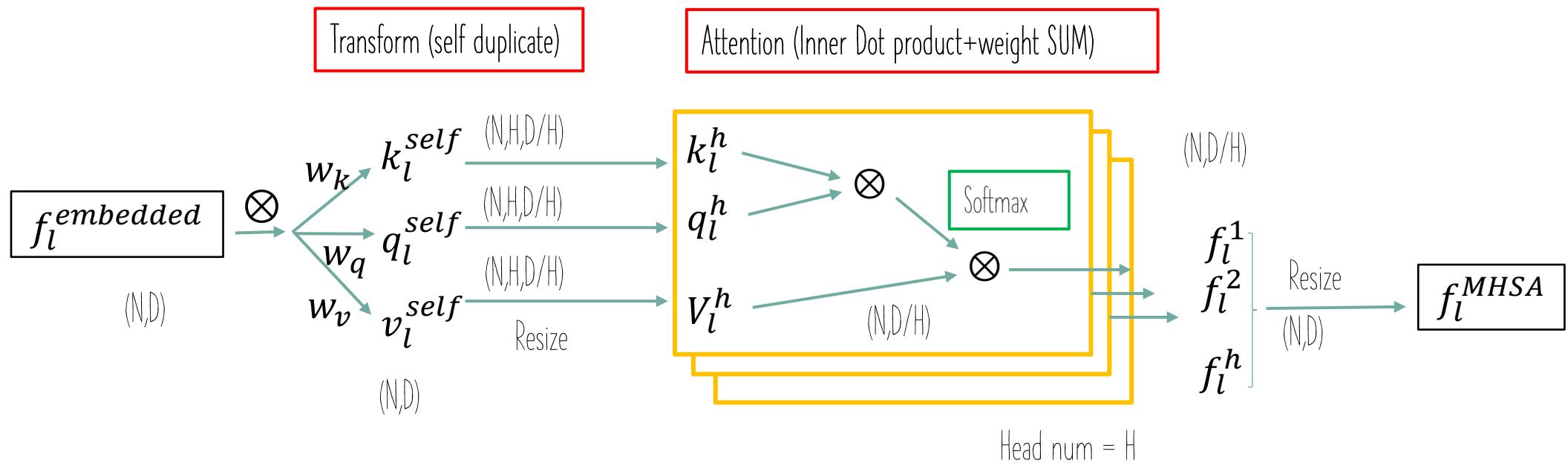
w (N,D)

$f_l^{SA}$  ( $N,D$ )

Weight sum (attention)

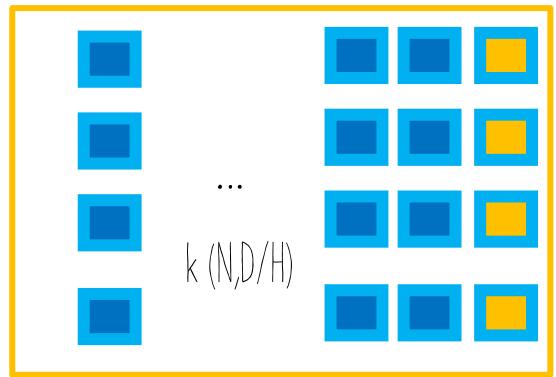
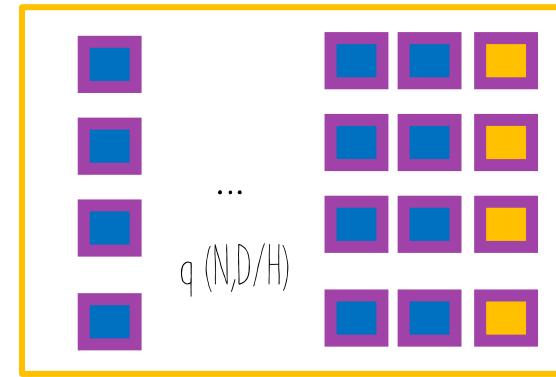


## MHSA multi-head self attention



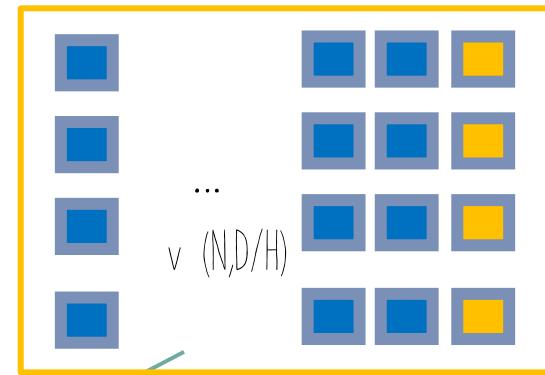
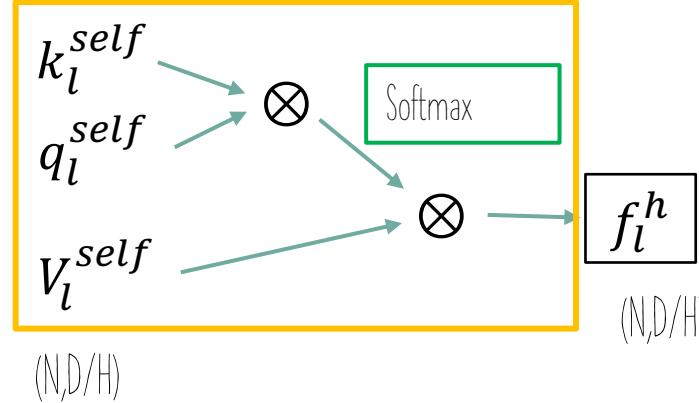
For a given embedded patch in each head  $f_l^{embedded}$ ,  $f_l^{MHSAs} = \text{SoftMax}(q_l^{selfT} \cdot k_l^{self})^T \cdot V_l^{self}$ , where  $q_l^{self} = w_q \cdot f_l^{embedded}$ ,  $k_l^{self} = w_k \cdot f_l^{embedded}$ ,  $V_l^{self} = w_v \cdot f_l^{embedded}$ .

## MHSA multi-head self attention

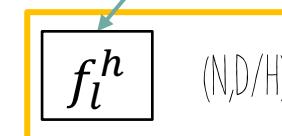


In each head

In place dot product

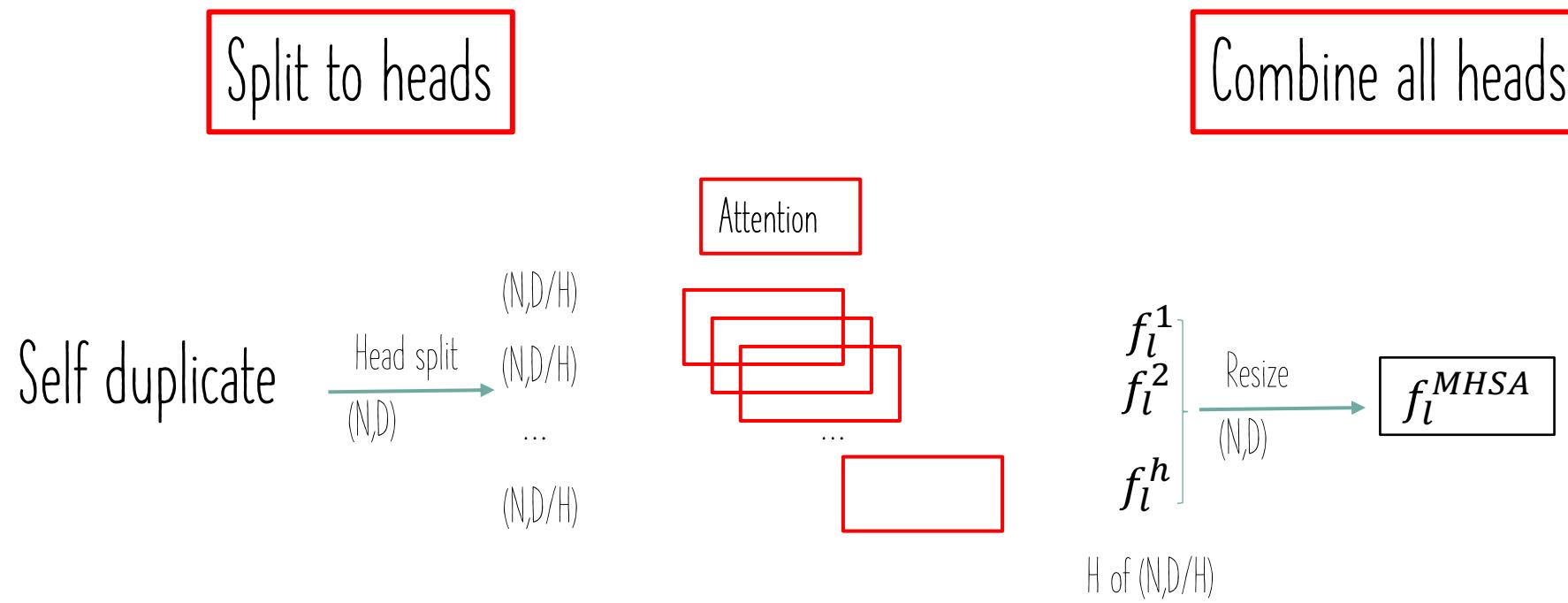


In place dot product



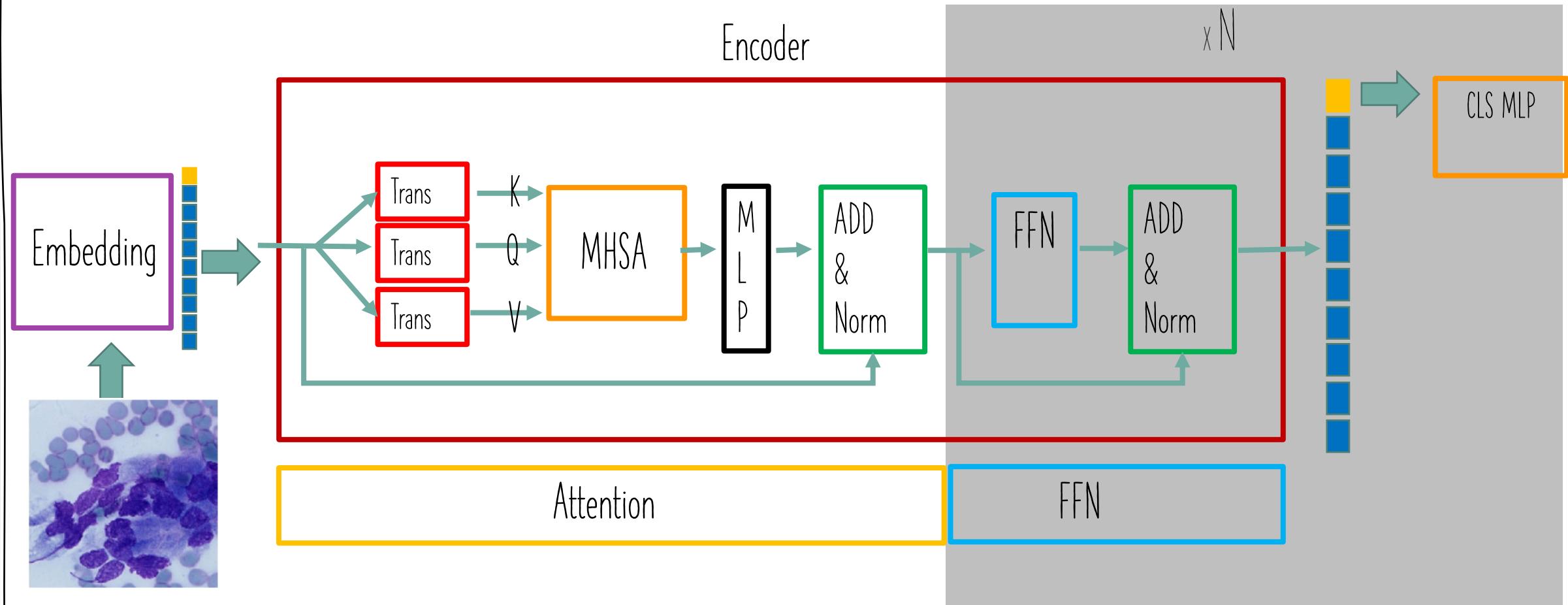
Weight sum (attention)

## MHSA multi-head self attention



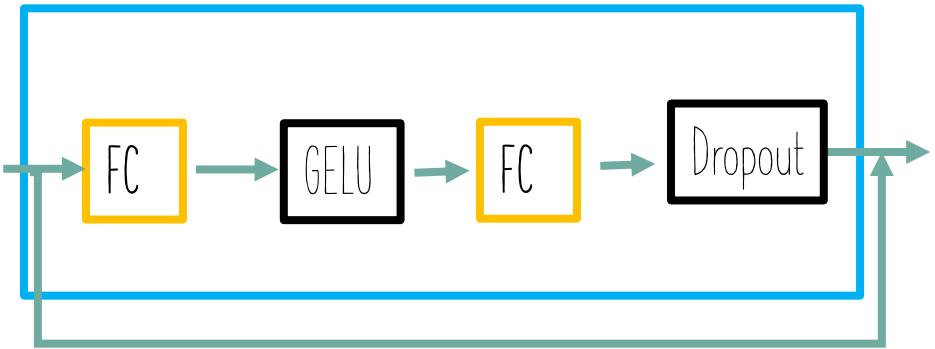
# ViT Vision Transformer

AN IMAGE IS WORTH 16X16 WORDS: TRANSFORMERS FOR IMAGE RECOGNITION AT SCALE  
Arxiv 2010.11929



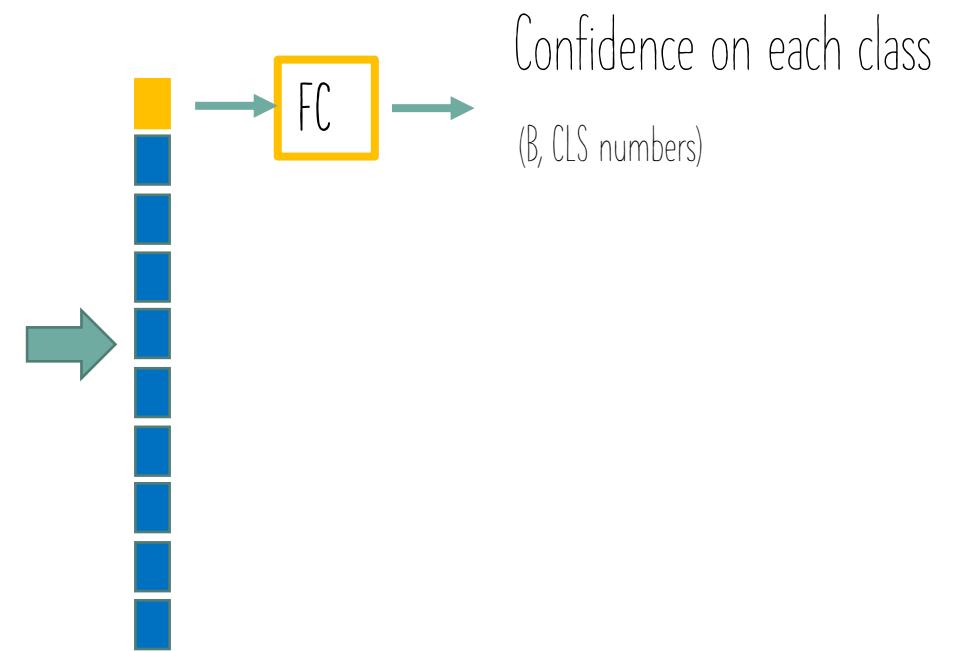
# FFN

(Batch,  
patch\_num,  
Embed\_dim)  
(8, 196, 768)



# CLS MLP

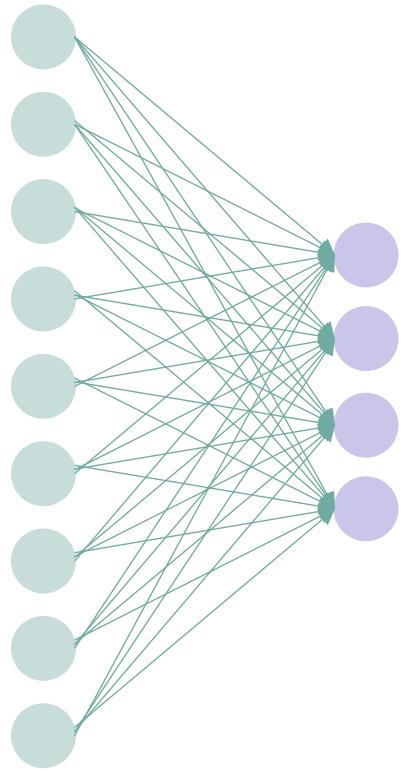
(Batch,  
patch\_num,  
Embed\_dim)  
(8, 196, 768)





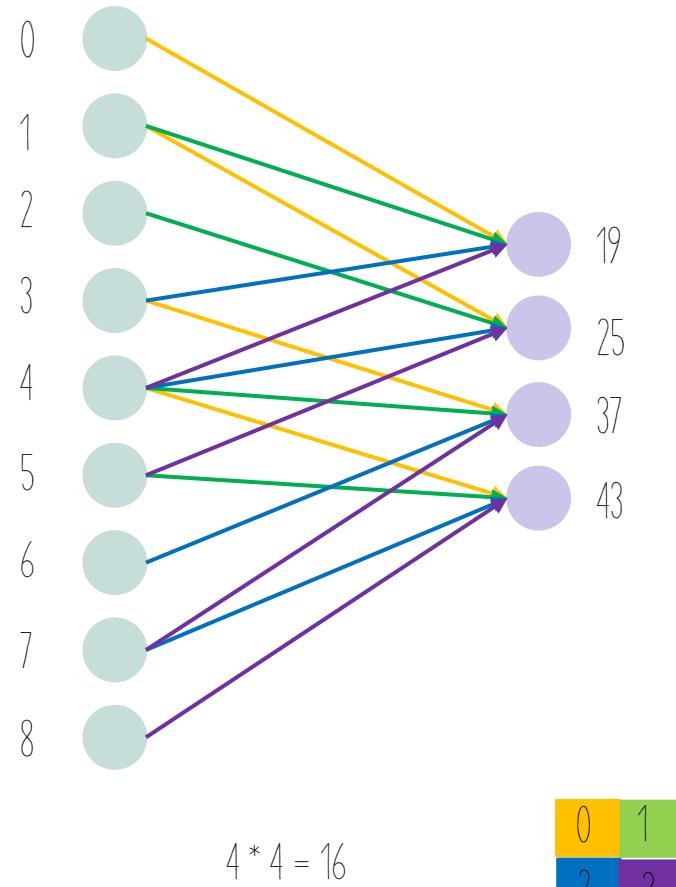
The projection issue

# What is the neurons of CNN looks like?

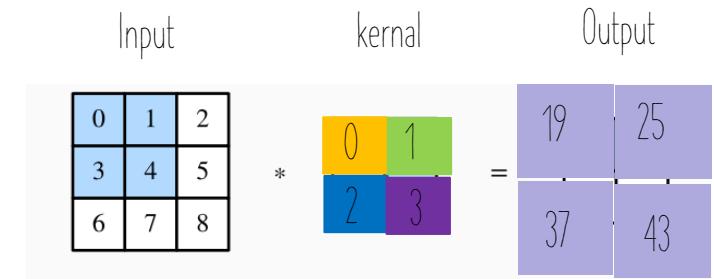


$$9 * 4 = 36$$

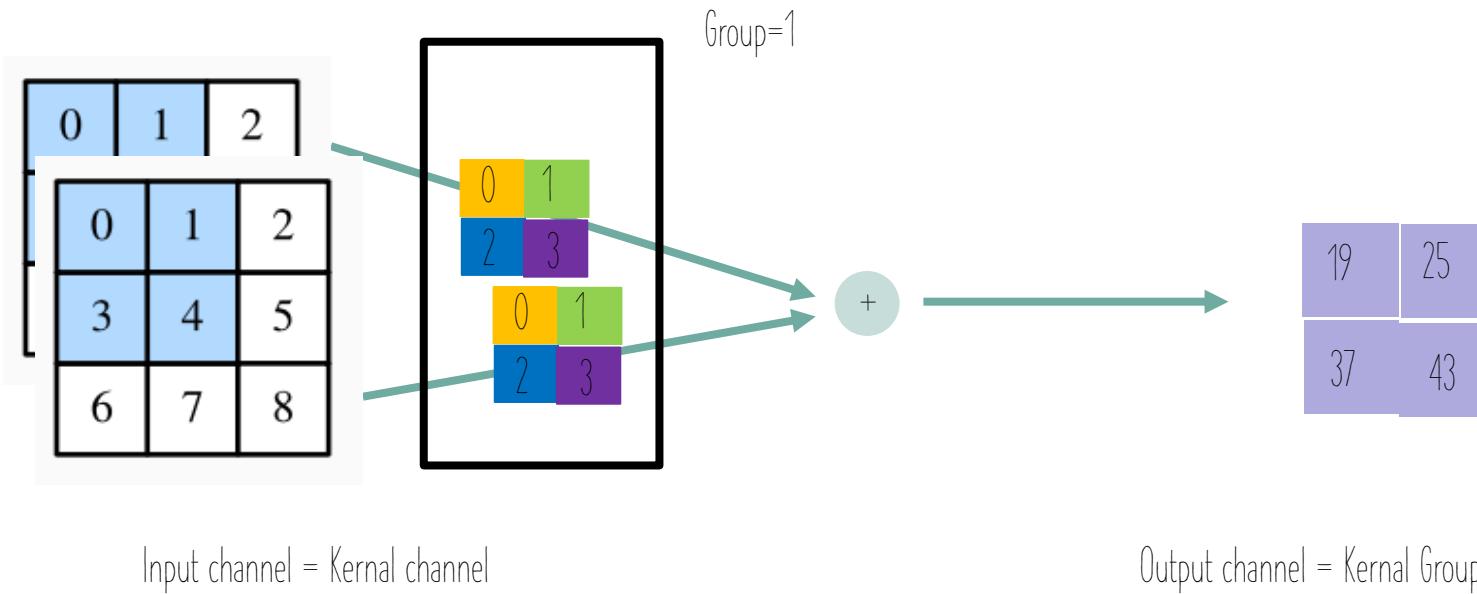
$$F(X) = \sigma(A X + B)$$



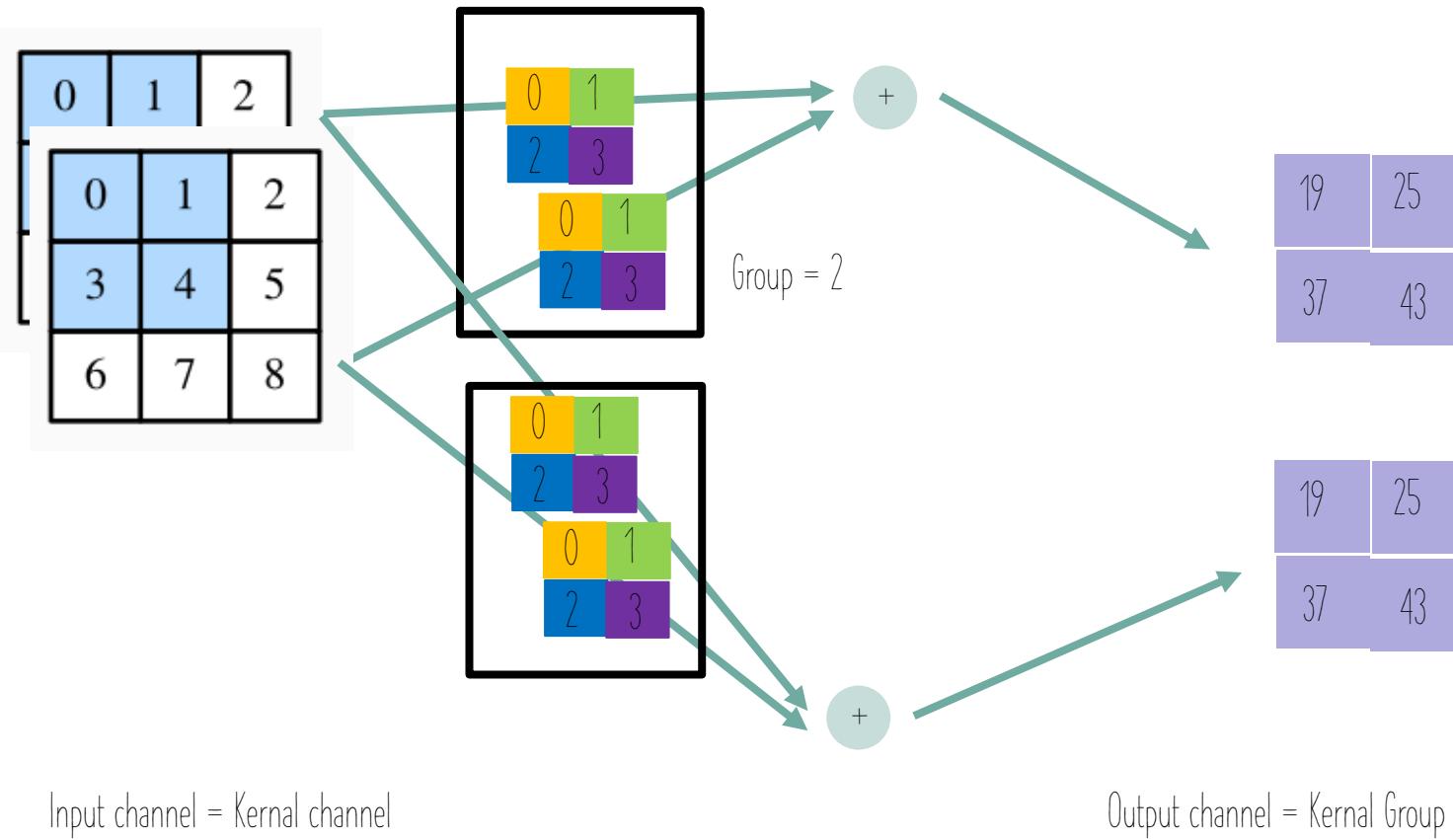
The number of real parameters is only 4



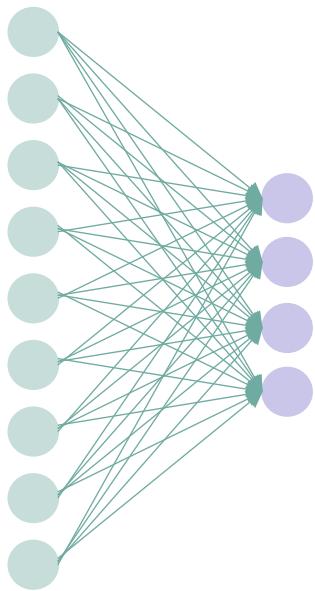
# What is CNN's channel?



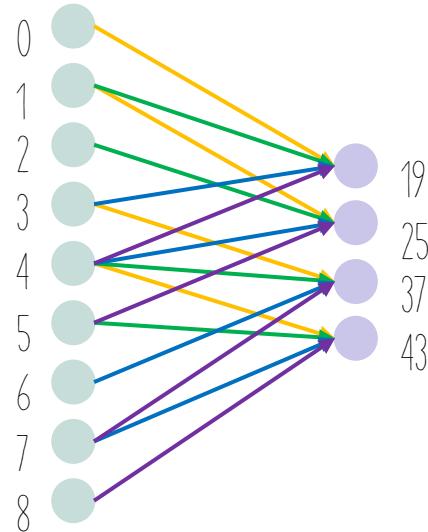
# What is CNN's channel?



```
import torch  
import torch.nn as nn  
  
# FC  
layer = nn.Linear(9, 4)  
img = torch.randn(1, 9)  
  
preds = layer(img) # (1, class_number)  
print('test model output: ', preds.shape)
```



```
img = torch.randn(1, 10)  
preds = layer(img) # (1, class_number)  
print('test model output: ', preds.shape)  
# RuntimeError: mat1 and mat2 shapes cannot be multiplied (1x10 and 9x4)
```



```
import torch  
import torch.nn as nn  
  
layer = nn.Conv2d(3, 768, kernel_size=(16,16),  
                 stride=(16,16))  
  
edge_size = 384  
img = torch.randn(1, 3, edge_size, edge_size)  
preds = layer(img) # (1, class_number)  
print('test model output: ', preds.shape)
```

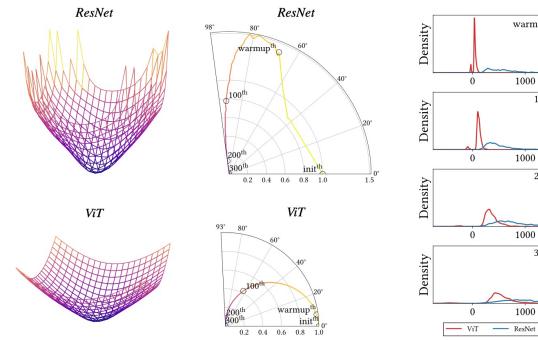
```
edge_size = 224  
img = torch.randn(1, 3, edge_size, edge_size)  
preds = layer(img) # (1, class_number)  
print('test model output: ', preds.shape)
```

# POINTS TO DECLEAR

1. Who's learning the most of the hard works ?
2. Who's the heavyest ?
3. Why its soooo stirring ?
4. A slight view of the future ? (structure, strategy, ...)

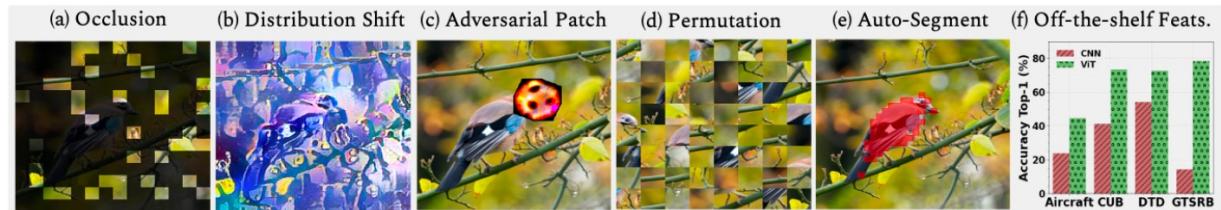
## Hard Math

HOW DO VISION TRANSFORMERS WORK?  
Arxiv 2202.06709



## Inspiring Experiments

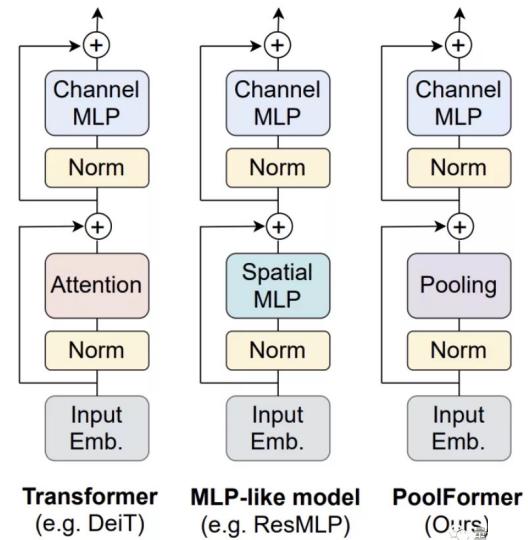
NeurPS2021 Intriguing Properties of Vision Transformers  
Arxiv 2105.10497



## Inspiring Model play

MetaFormer is Actually What You Need for Vision  
Poolformer Arxiv 2111.11418

Patches Are All You Need  
<https://openreview.net/pdf?id=TVHS5Y4dNvM>



Under review as a conference paper at ICLR 2022

CONVOLUTIONS ATTENTION MLPs  
PATCHES ARE ALL YOU NEED? 🤔

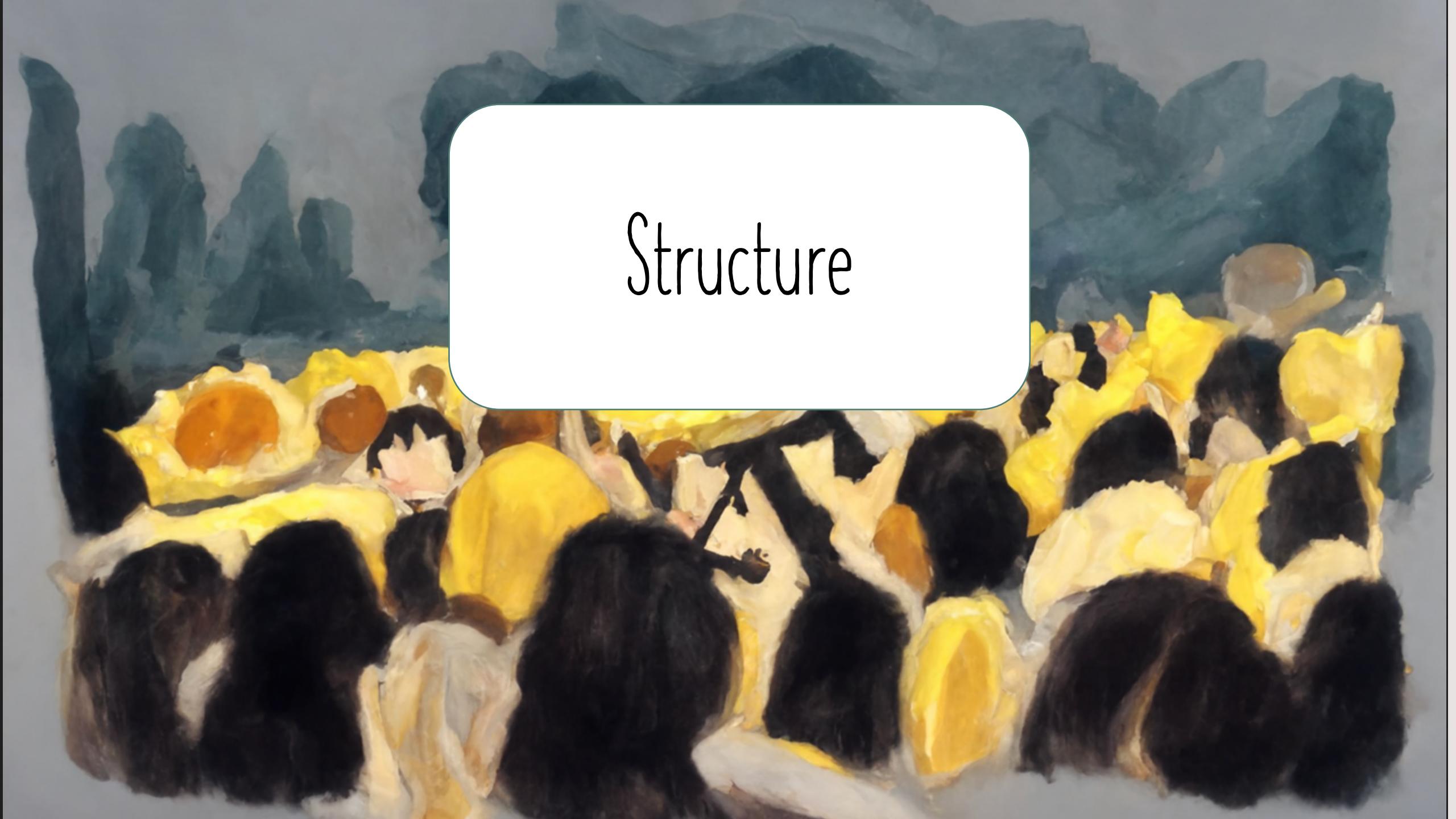
Anonymous authors  
Paper under double-blind review



Model structure

Learning strategy

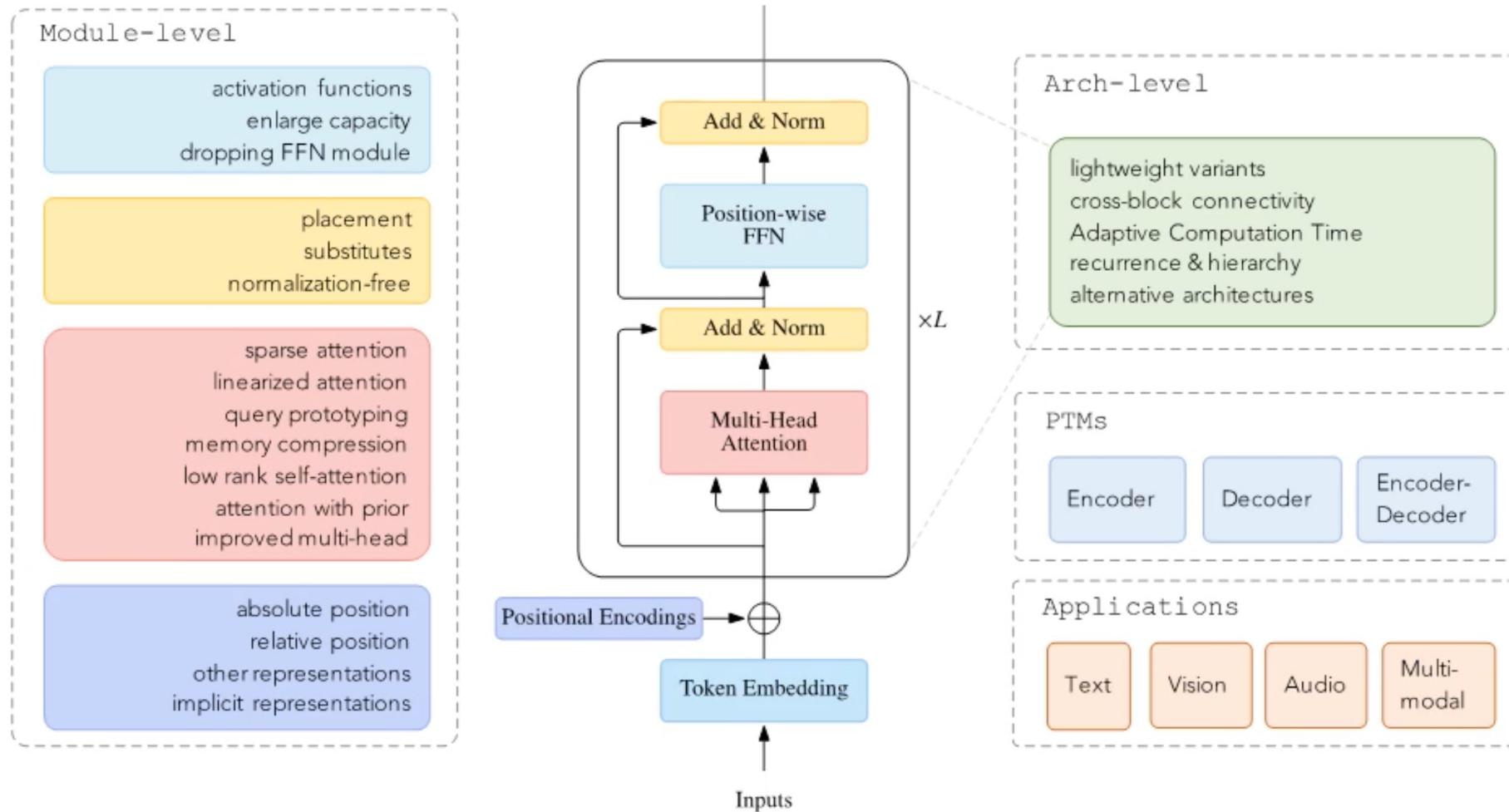
Initialization



Structure

# A taxonomy of X-formers (Transformer variants)

arXiv 2106.04554



MSHT

arXiv:2112.13513

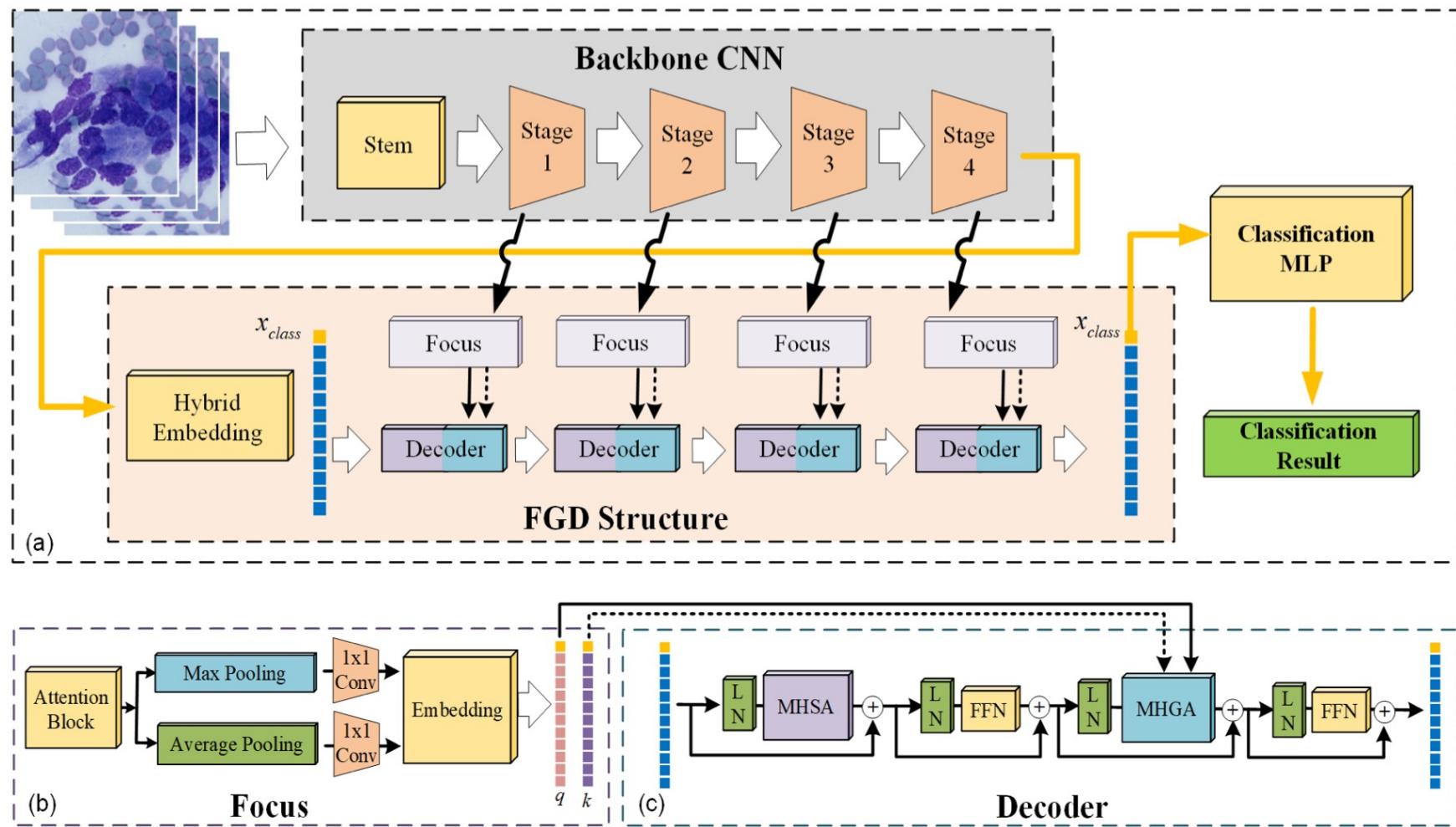
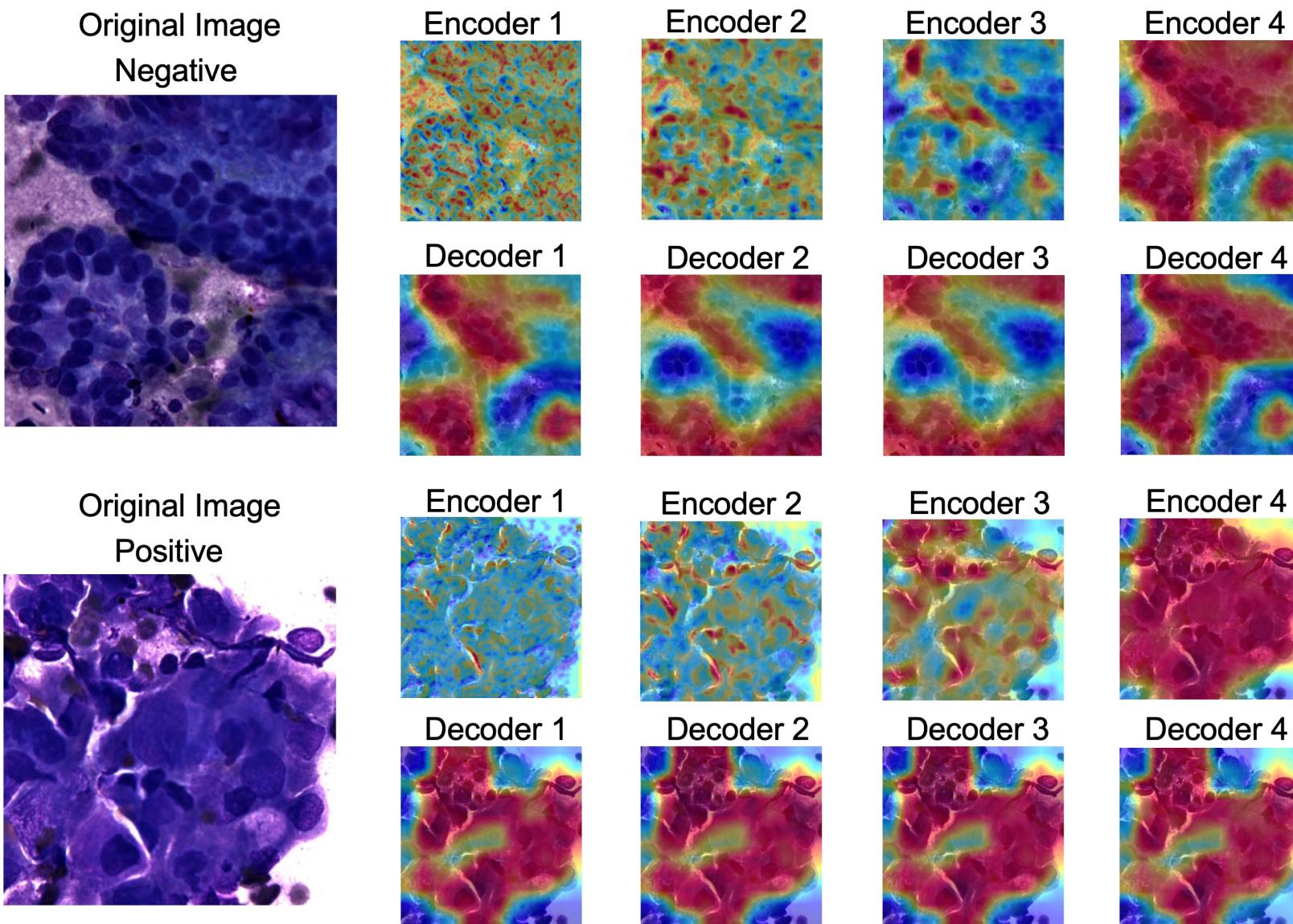


Fig. 1. The architecture of the proposed Multi-stage Hybrid Transformer (MSHT) model for the classification of ROSE images of pancreatic cancer. (a) The architecture of MSHT. (b) The focus block of the FGD structure (c) The decoder of the FGD structure. MHSA denotes multi-head self-attention, MHGA denotes multi-head guided-attention, LN denotes layer norm block, FFN denotes the feed-forward network, and MLP denotes multi-layer perceptron.

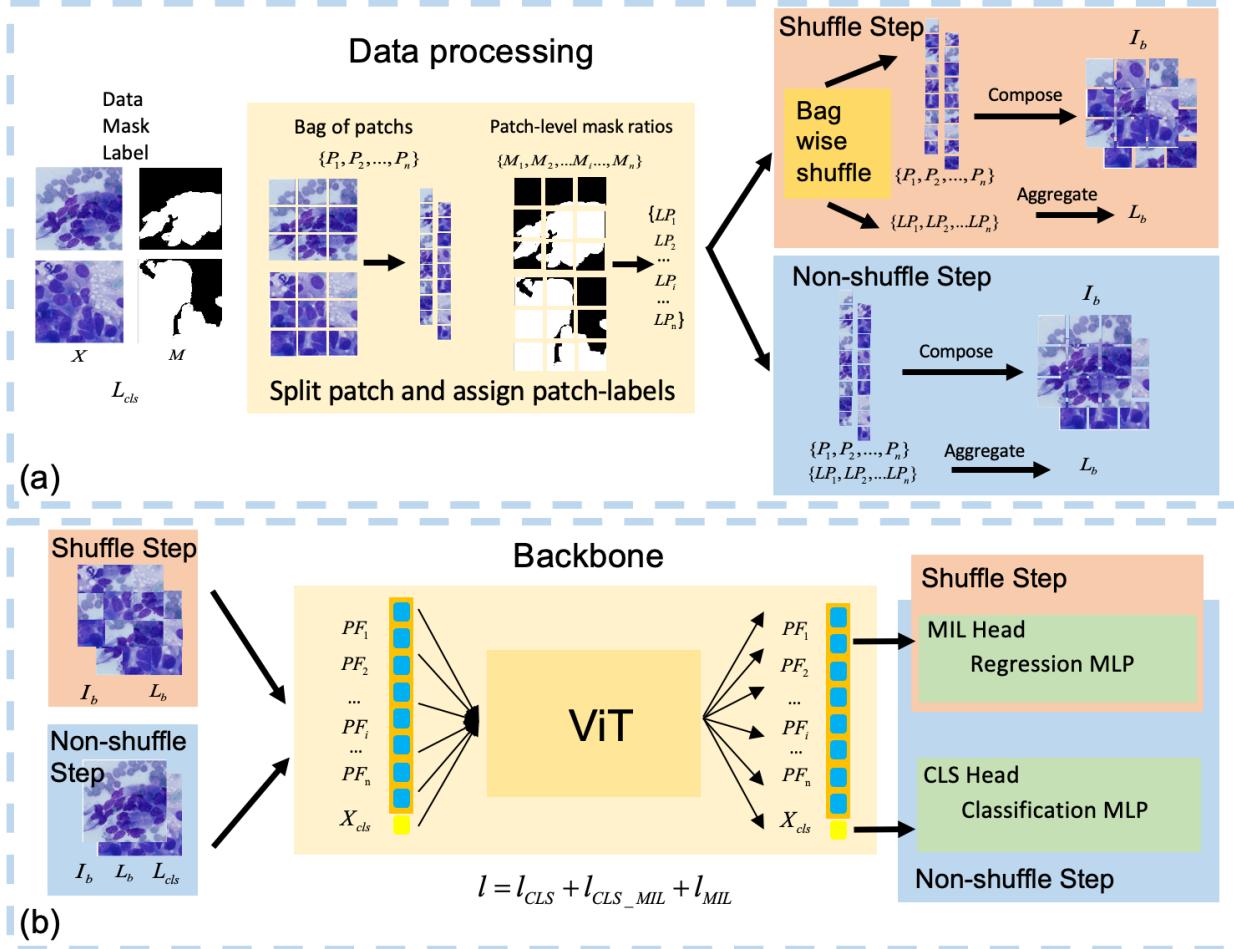
Tianyi Zhang, Yunlu Feng, Yu Zhao, Guangda Fan, Aiming Yang, Shangqin Lyu, Peng Zhang, Fan Song, Chenbin Ma, Yangyang Sun, Youdan Feng, and Guanglei Zhang\*, "MSHT: Multi-stage hybrid transformer for the ROSE image analysis of pancreatic cancer arXiv:2112.13513 . (JBHI under review)





Strategy

# SI-ViT/MIL-SI



## Key Points

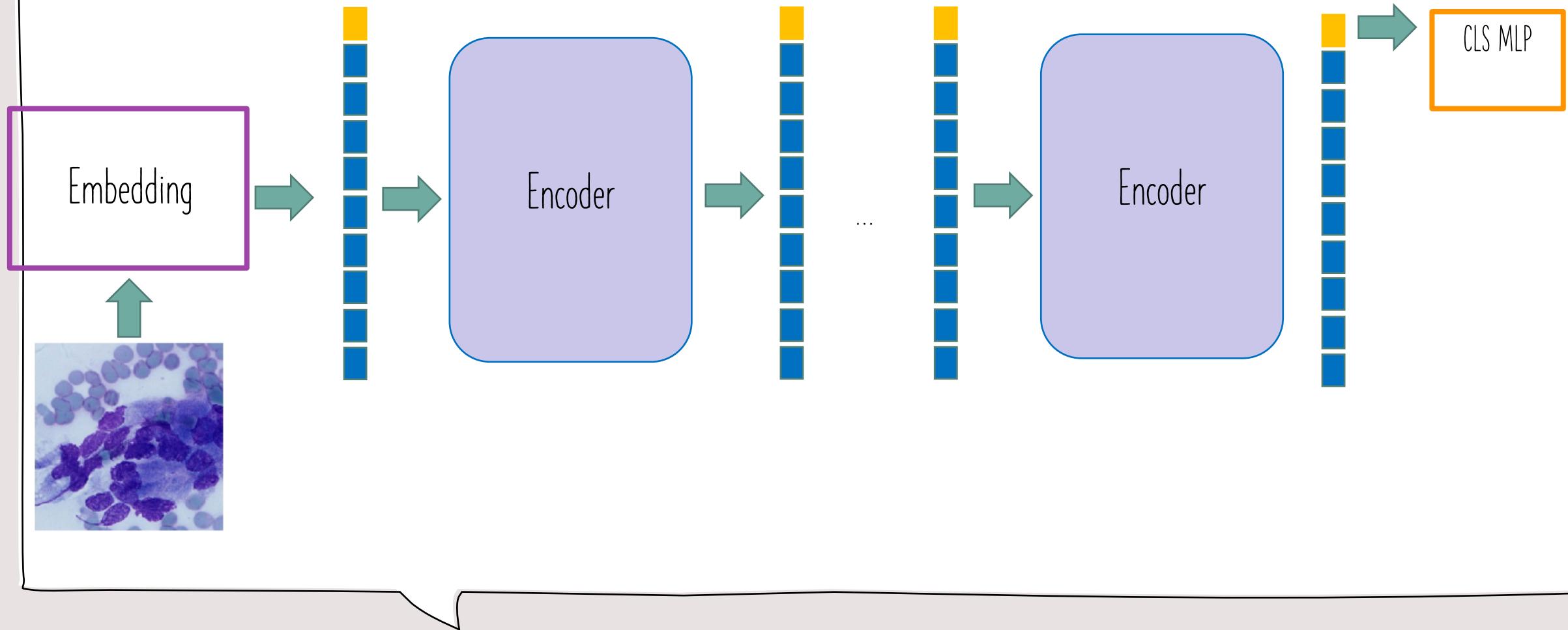
1. Fast and with only limited alterations
2. Instance re-grouping strategy
3. A soft-label with background
4. Alternative training strategy to balance CLS and MIL tasks
5. Two head design
6. Simple Loss design

Shuffle step to understand the differences of the patches and understand what each patch are, across the samples.

Non shuffle step to understand the relations of the patches in the original image.

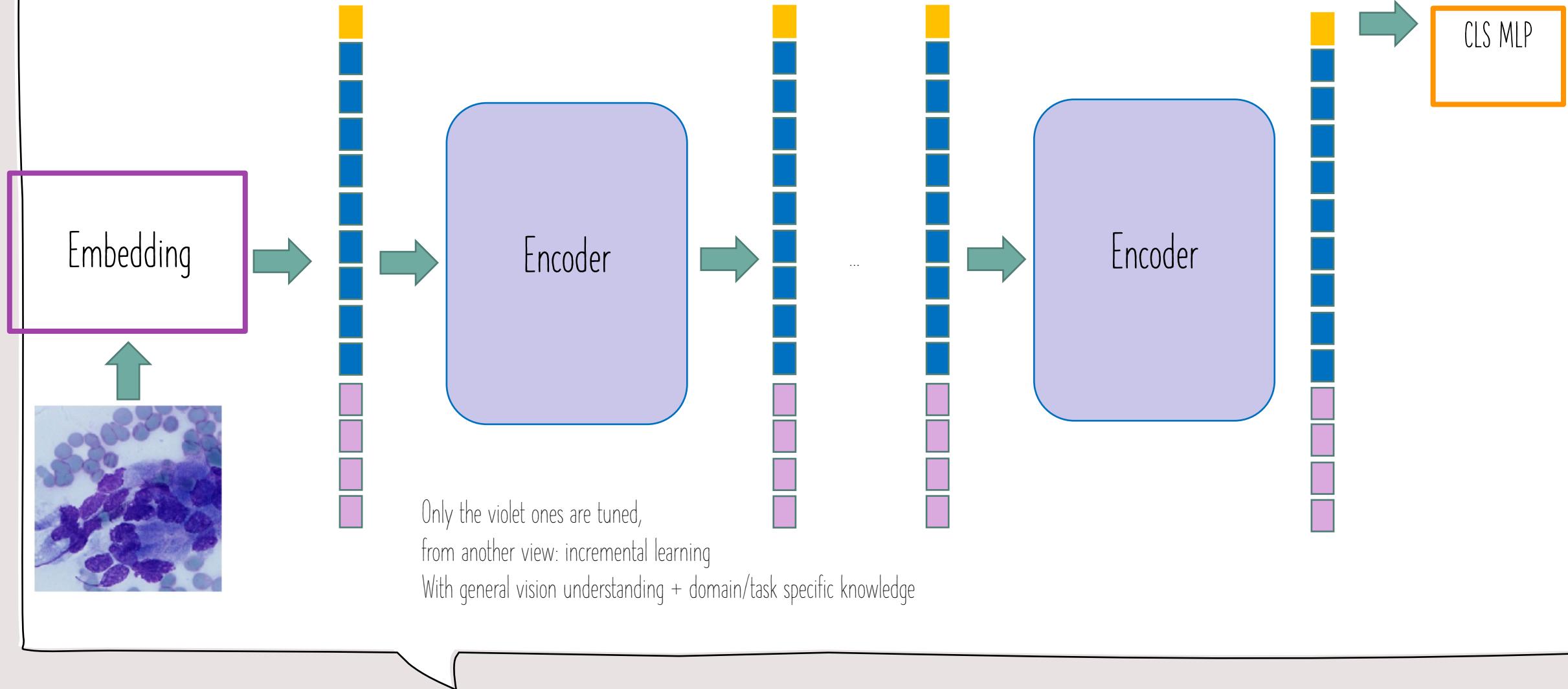
ViT

AN IMAGE IS WORTH 16X16 WORDS: TRANSFORMERS FOR IMAGE RECOGNITION AT SCALE  
arxiv 2010.11929

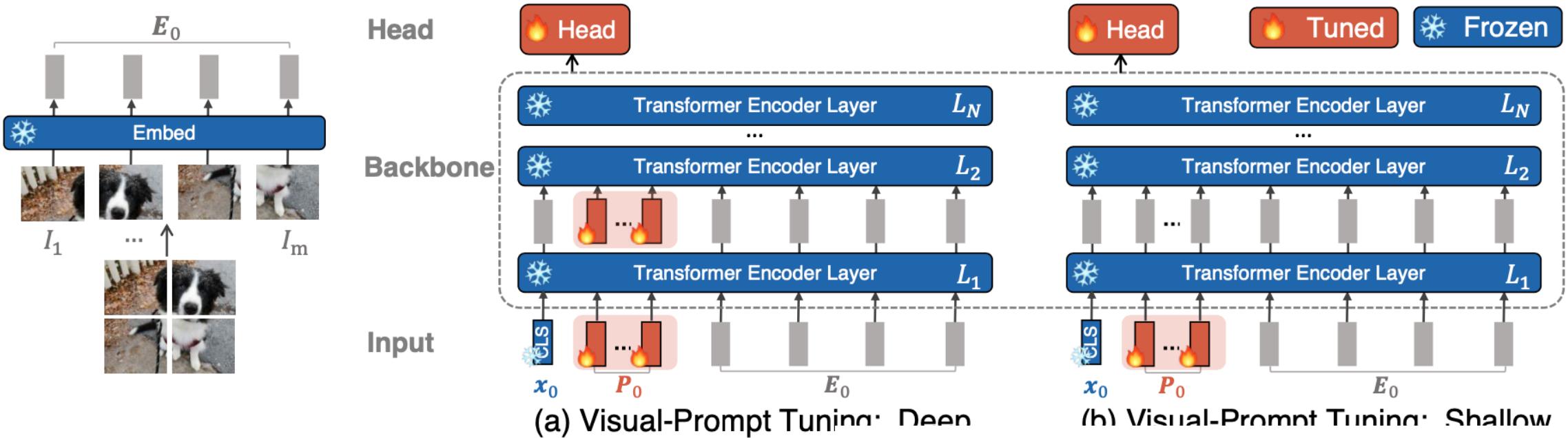


VPT

Jia, M., Tang, L., Chen, B.C., Cardie, C., Belongie, S., Hariharan, B. and Lim, S.N., 2022. Visual prompt tuning. *arXiv preprint arXiv:2203.12119*.



# VPT: use the Prompt tokens in the tuning process



$$[\mathbf{x}_i, \dots, \mathbf{E}_i] = L_i([\mathbf{x}_{i-1}, \mathbf{P}_{i-1}, \mathbf{E}_{i-1}])$$

$$\mathbf{y} = \text{Head}(\mathbf{x}_N) .$$

$$[\mathbf{x}_1, \mathbf{Z}_1, \mathbf{E}_1] = L_1([\mathbf{x}_0, \mathbf{P}, \mathbf{E}_0])$$

$$[\mathbf{x}_i, \mathbf{Z}_i, \mathbf{E}_i] = L_i([\mathbf{x}_{i-1}, \mathbf{Z}_{i-1}, \mathbf{E}_{i-1}])$$

$$\mathbf{y} = \text{Head}(\mathbf{x}_N) ,$$

My version of code: <https://github.com/sagizty/VPT>

# Thanks

PPT: <https://github.com/sagizty/Insight/blob/main/Dive%20into%20Transformer.pdf>  
Code: [https://github.com/sagizty/Insight/blob/main/ZTY\\_Torch/ViT.py](https://github.com/sagizty/Insight/blob/main/ZTY_Torch/ViT.py)

Tianyi