



# TransMIL: Transformer based Correlated Multiple Instance Learning for Whole Slide Image Classification

Zhuchen Shao<sup>\*,1</sup>, Hao Bian<sup>\*,1</sup>, Yang Chen<sup>\*,1</sup>, Yifeng Wang<sup>2</sup>, Jian Zhang<sup>3</sup>, Xiangyang Ji<sup>4</sup>  
Yongbing Zhang<sup>†,2</sup>

<sup>1</sup>Tsinghua Shenzhen International Graduate School, Tsinghua University

<sup>2</sup>Harbin Institute of Technology (Shenzhen)

<sup>3</sup>School of Electronic and Computer Engineering, Peking University

<sup>4</sup>Department of Automation, Tsinghua University



# CONTENTS



**Background**



**Method**



**Experiments and Results**

# Whole Slide Image (WSI)

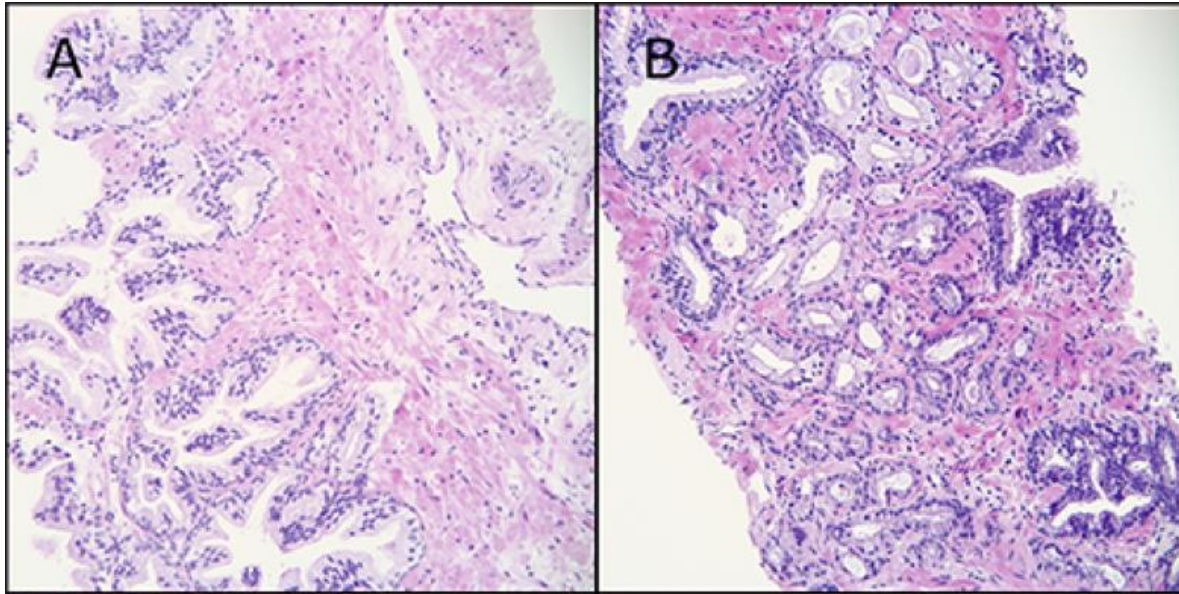


Fig. 1 Whole slide image

- Huge size ( $\sim 50000 \times 50000$  pixels at  $20\times$ )
- Lack of pixel-level annotations

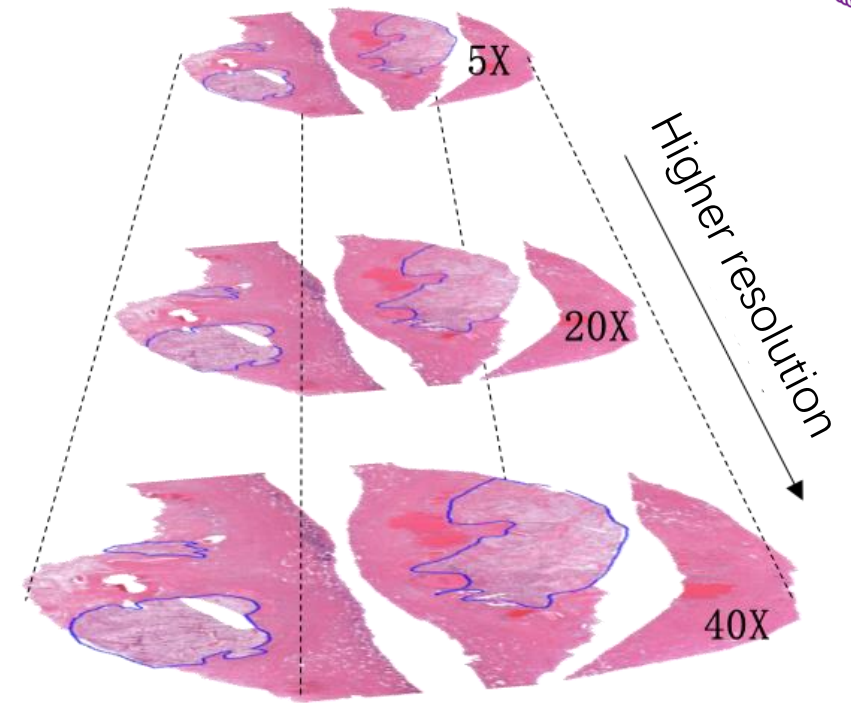


Fig. 2 WSI storage

- Multiple resolution images
- Tissue to cell level information



# Multiple Instance Learning (MIL)

## Description of MIL problem

- MIL is a weak supervision problem
- Each bag contains an unequal number of instances
- Bag level label is known, instance label is unknown

## Traditional assumption in MIL problem

- All the instances in each bag are independent and identically distributed (i.i.d.)

# Correlated Multiple Instance Learning

- **Difference:** Consider the correlation and spatial information between different instances in a bag.

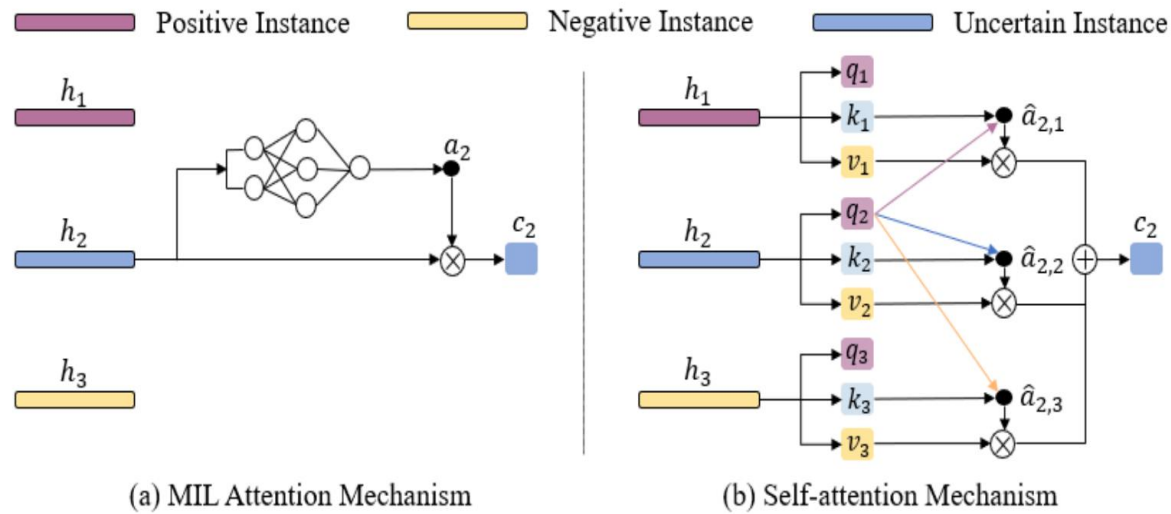


Fig. 3 Difference between MIL attention mechanism and self-attention mechanism

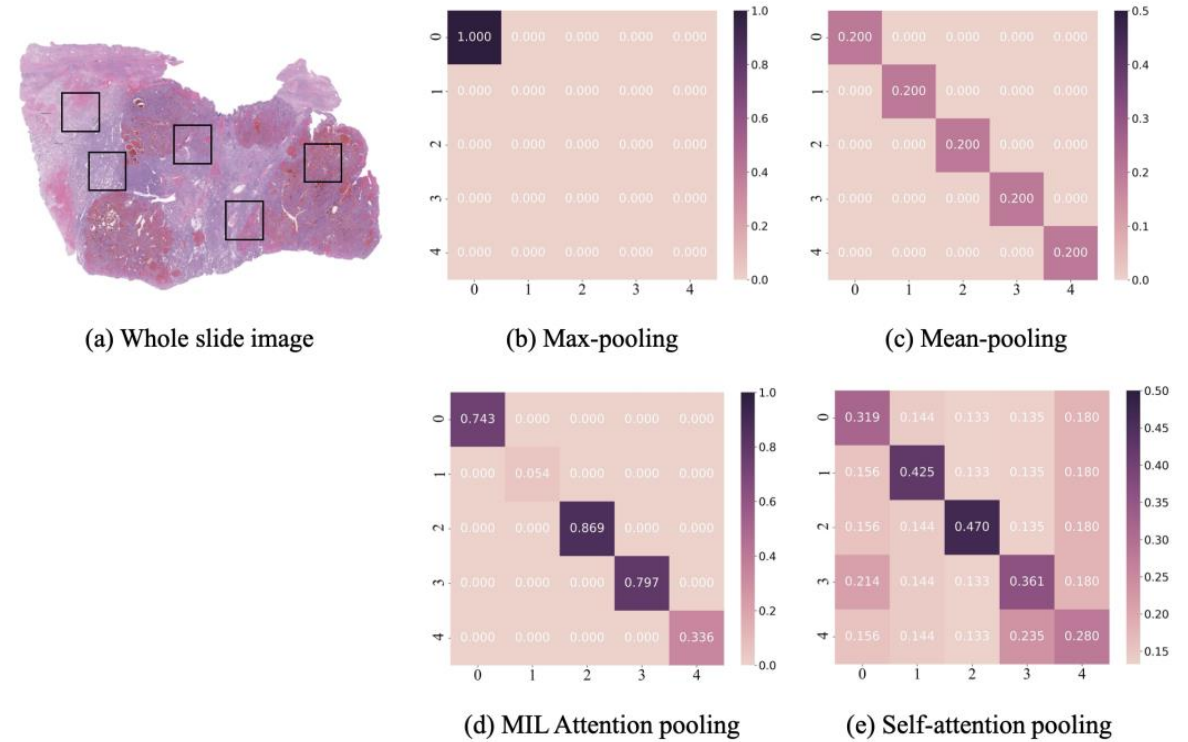


Fig.4 The difference between different Pooling Matrix

# Method

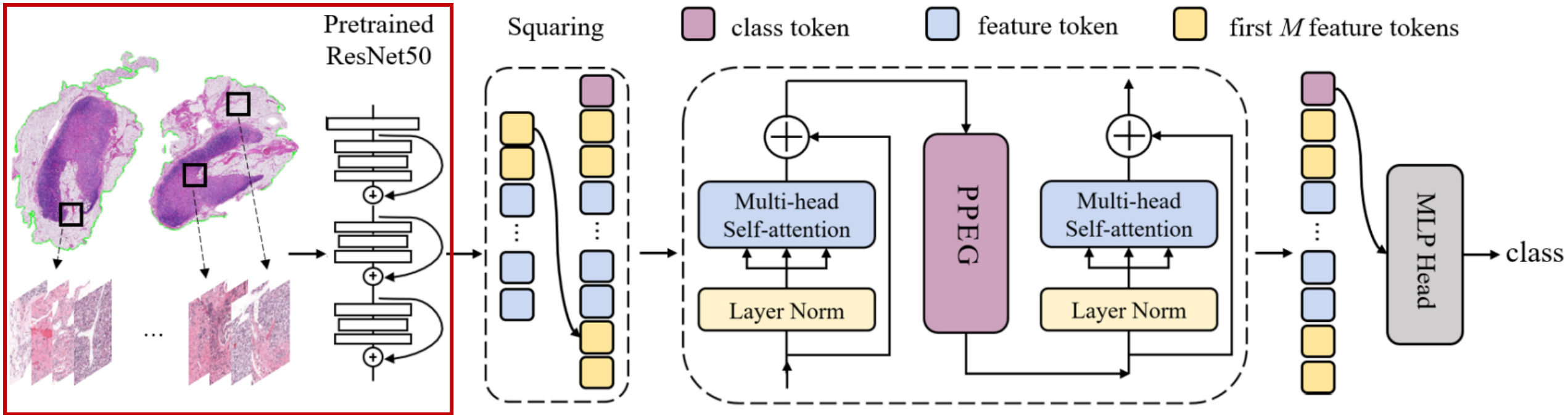


Fig. 5 Overview of our TransMIL

## - Preprocessing :

- 1) Each WSI is cropped into a series of  $256 \times 256$  non-overlapping patches, where background region (saturation  $< 15$ ) is discarded.
- 2) The feature of each patch is embedded in a 1024-dimensional vector by a ResNet50 model pre-trained on ImageNet.

# Method

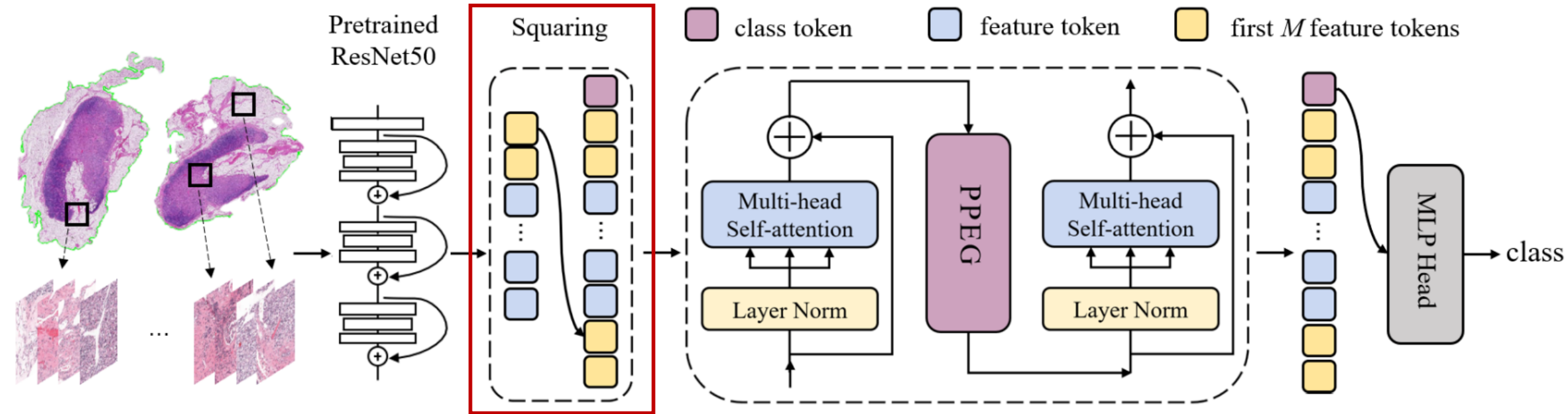


Fig. 5 Overview of our TransMIL

- **Squaring of sequence:**
  - Square the length of the sequence, and add the class token, then reduce the dimension of each feature embedding from 1024 to 512.

# Method

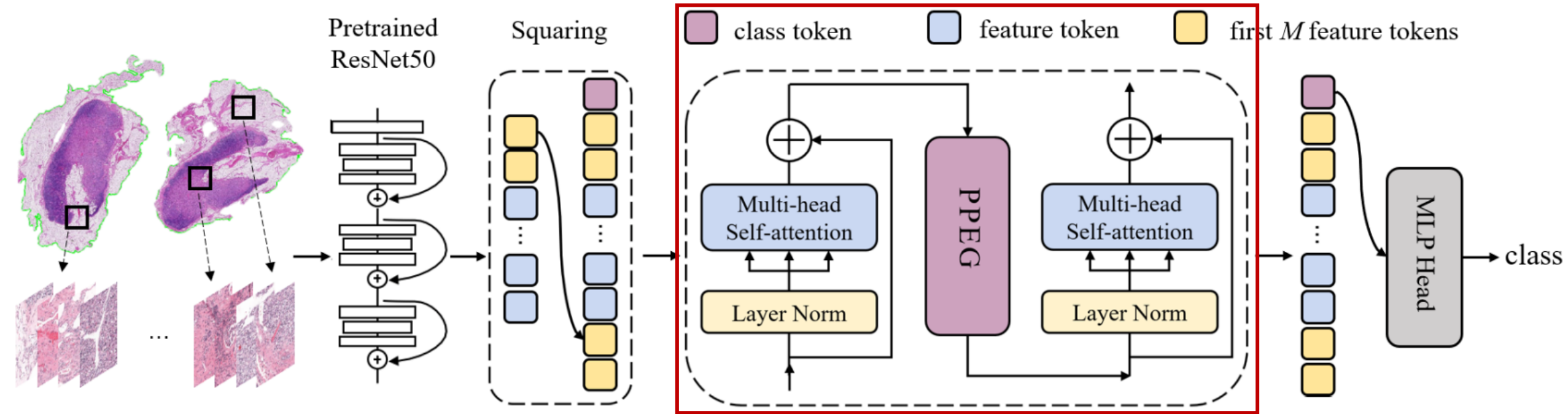


Fig. 5 Overview of our TransMIL

- 1) Correlation modelling of the sequence  $\mathbf{H}_S^\ell \leftarrow \text{MSA}(\mathbf{H}_S)$
- 2) Conditional position encoding and local information fusion  $\mathbf{H}_S^P \leftarrow \text{PPEG}(\mathbf{H}_S^\ell)$
- 3) Deep feature aggregation  $\mathbf{H}_S^{\ell+1} \leftarrow \text{MSA}(\mathbf{H}_S^P)$



# Position encoding with PPEG

- **Background** : Zero padding can provide an absolute position information to convolution<sup>[1]</sup>

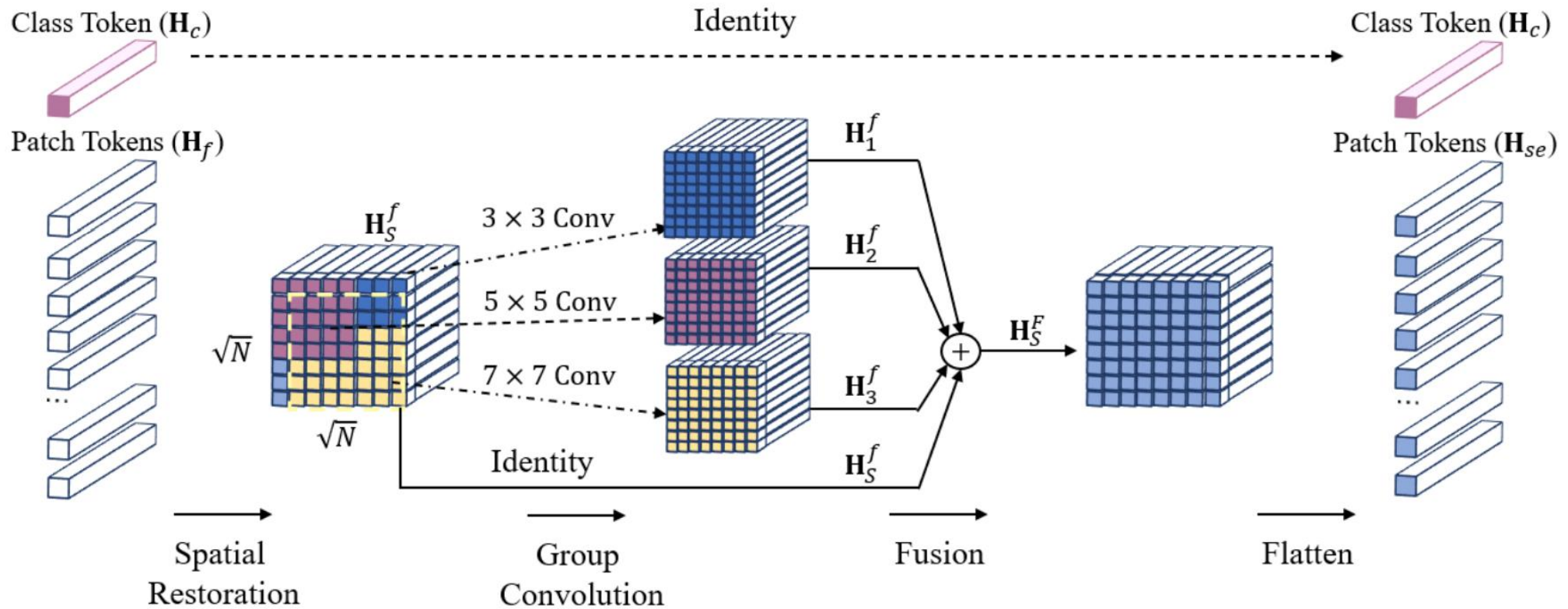


Fig.6 Pyramid Position Encoding Generator (PPEG)

[1] How much position information do convolutional neural networks encode? In International Conference on Learning Representations, 2020.

# Method

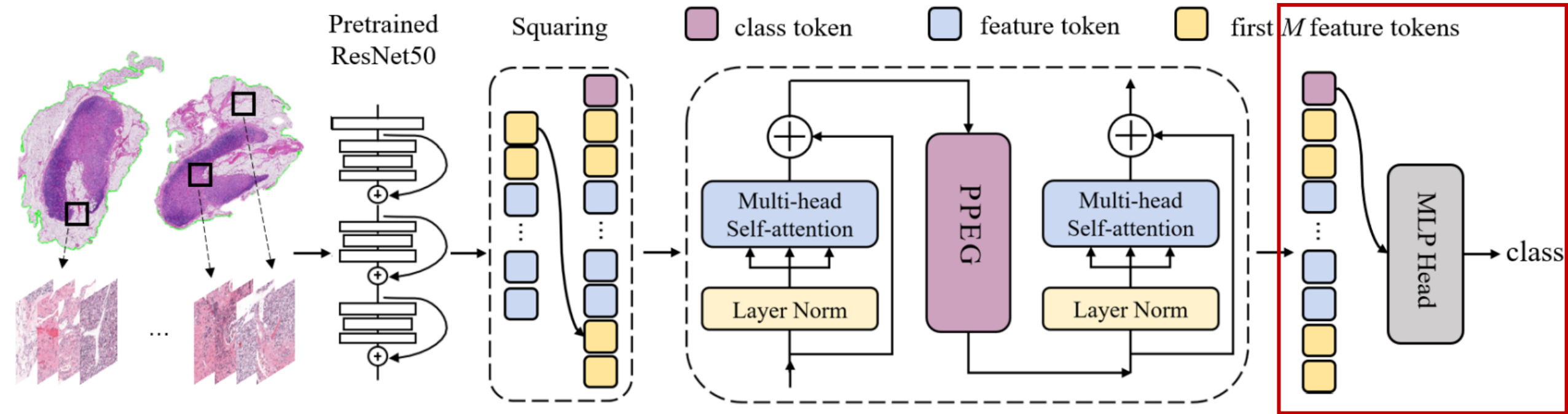


Fig. 5 Overview of our TransMIL

## - Mapping of $\mathbb{T} \rightarrow \mathcal{Y}$

- Use the class token to get the slide-level label of WSI

$$- \hat{Y} \leftarrow \text{MLP}(\text{LN}((\mathbf{H}_S^{\ell+1})^{(0)}))$$



# Experiments and Results

Table 1: Results on CAMELYON16, TCGA-NSCLC and TCGA-RCC.

|                              | CAMELYON16    |               | TCGA-NSCLC    |               | TCGA-RCC      |               |
|------------------------------|---------------|---------------|---------------|---------------|---------------|---------------|
|                              | Accuracy      | AUC           | Accuracy      | AUC           | Accuracy      | AUC           |
| <b>Mean-pooling</b>          | 0.6389        | 0.4647        | 0.7282        | 0.8401        | 0.9054        | 0.9786        |
| <b>Max-pooling</b>           | 0.8062        | 0.8569        | 0.8593        | 0.9463        | 0.9378        | 0.9879        |
| <b>ABMIL<sup>[1]</sup></b>   | 0.8682        | 0.8760        | 0.7719        | 0.8656        | 0.8934        | 0.9702        |
| <b>MIL-RNN<sup>[2]</sup></b> | 0.8450        | 0.8880        | 0.8619        | 0.9107        | \             | \             |
| <b>DSMIL<sup>[3]</sup></b>   | 0.7985        | 0.8179        | 0.8058        | 0.8925        | 0.9294        | 0.9841        |
| <b>CLAM-SB<sup>[4]</sup></b> | 0.8760        | 0.8809        | 0.8180        | 0.8818        | 0.8816        | 0.9723        |
| <b>CLAM-MB<sup>[4]</sup></b> | 0.8372        | 0.8679        | 0.8422        | 0.9377        | 0.8966        | 0.9799        |
| <b>TransMIL</b>              | <b>0.8837</b> | <b>0.9309</b> | <b>0.8835</b> | <b>0.9603</b> | <b>0.9466</b> | <b>0.9882</b> |

[1] Attention-based deep multiple instance learning. In International Conference on Machine Learning, 2018.

[2] Clinical-grade computational pathology using weakly supervised deep learning on whole slide images. Nature medicine, 2019.

[3] Dual-stream multiple instance learning network for whole slide image classification with self-supervised contrastive learning. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2021.

[4] Data-efficient and weakly supervised computational pathology on whole-slide images. Nature Biomedical Engineering, 2021.



THANK YOU