

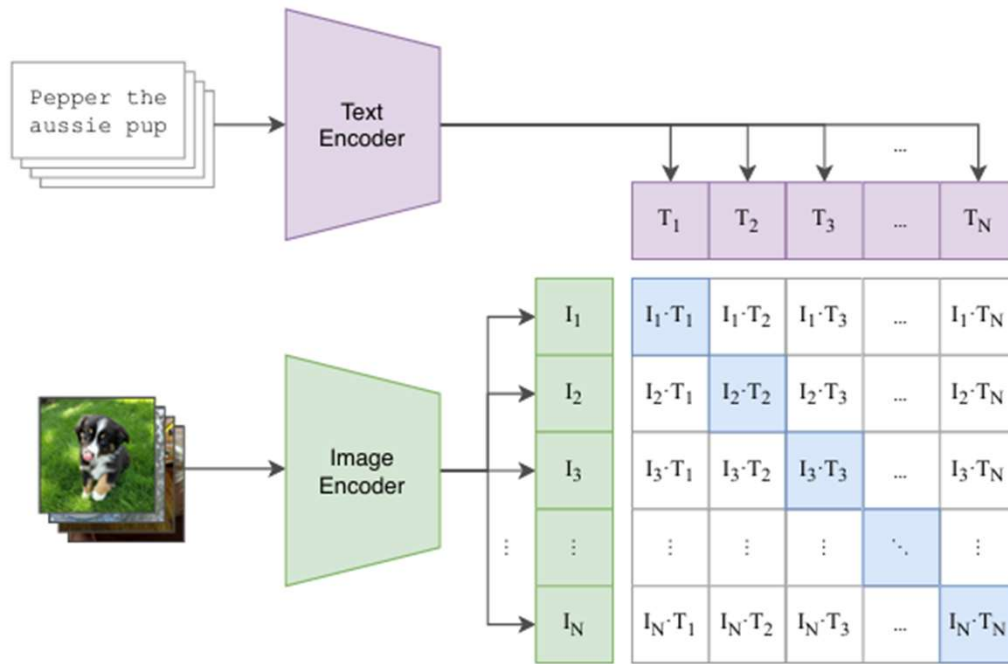
# TagCLIP: A Local-to-Global Framework to Enhance Open-Vocabulary Multi-Label Classification of CLIP without Training

---

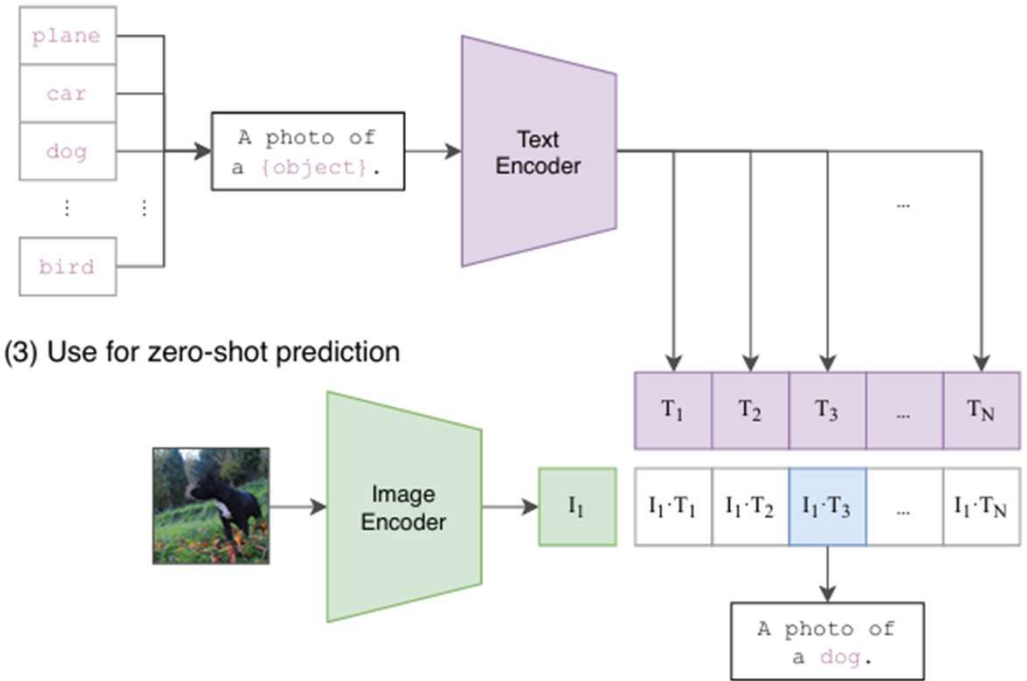
Yuqi Lin<sup>1,3</sup>, Minghao Chen<sup>2\*</sup>, Kaipeng Zhang<sup>3\*</sup>, Hengjia Li<sup>1</sup>, Mingming Li<sup>1</sup>,  
Zheng Yang<sup>4</sup>, Dongqin Lv<sup>6</sup>, Binbin Lin<sup>5</sup>, Haifeng Liu<sup>1</sup>, Deng Cai<sup>1,4</sup>

AAAI 2024

(1) Contrastive pre-training

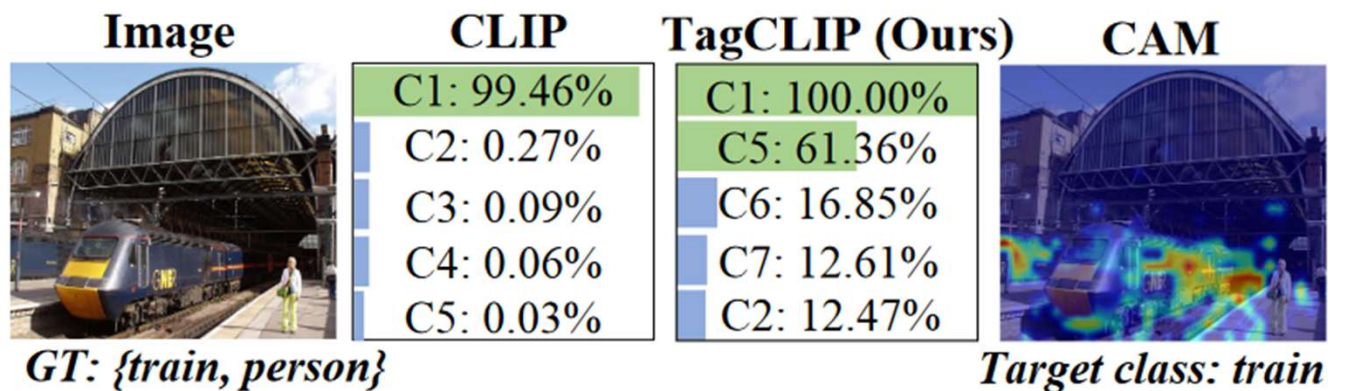


(2) Create dataset classifier from label text

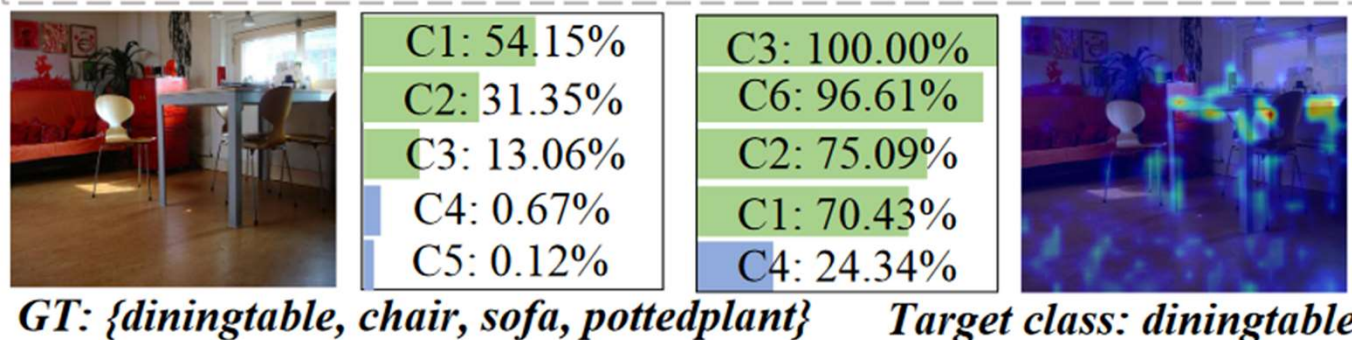


(3) Use for zero-shot prediction

# Motivation



C1: train C2:aerolane C3:bus C4:car C5:person C6:bird C7:bicycle



C1: diningtable C2:sofa C3:chair C4:person C5:train C6:pottedplant

# Motivation

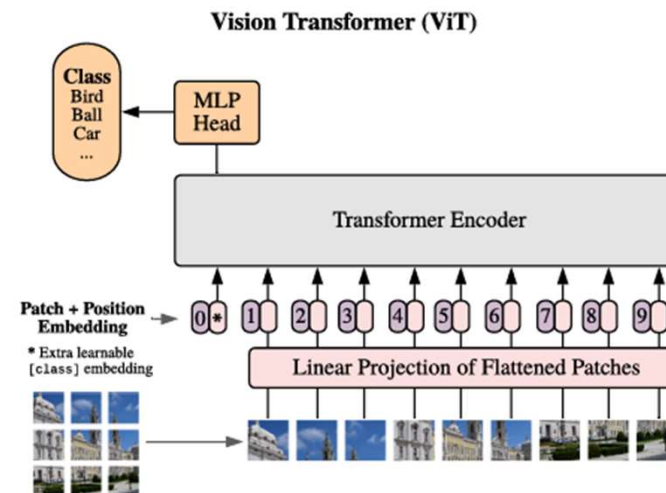
The forward propagation of the last transformer layer

$$\begin{aligned}\hat{X}^L &= X^{L-1} + \mathbf{a}^L, \\ &= X^{L-1} + A^L (X^{L-1} W_V^L), \\ A^L &= \sigma\left(\frac{(X^{L-1} W_Q^L)(X^{L-1} W_V^L)^T}{\sqrt{d}} + M^L\right),\end{aligned}$$

$$X^L = \hat{X}^L + \text{MLP}(\hat{X}^L),$$

$X^L$  consists of the [cls] token and remaining tokens (denoted as dense tokens)

$$X^L = [x_{cls}^L, x^L].$$



Let the dense token outputted by the penultimate layer pass the rest layer without self-attention

$$\begin{aligned}\hat{x}_{dense} &= x^{L-1} + \mathbf{c}^L, \\ &= x^{L-1} + x^{L-1} W_V^L, \\ x_{dense} &= \hat{x}_{dense} + \text{MLP}(\hat{x}_{dense}).\end{aligned}$$

# Motivation

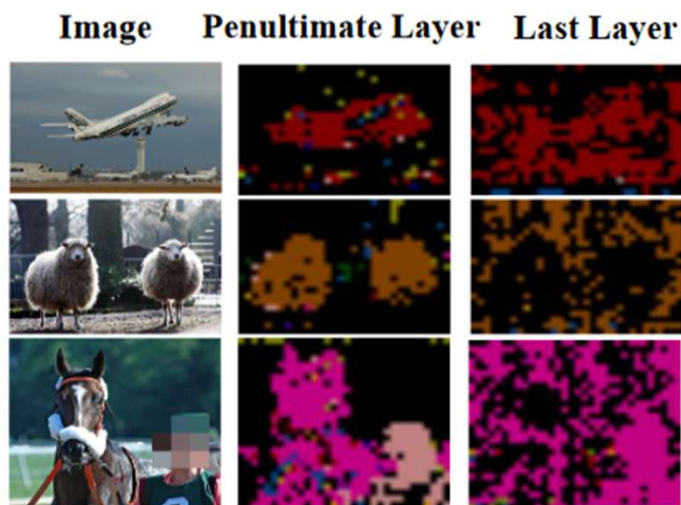
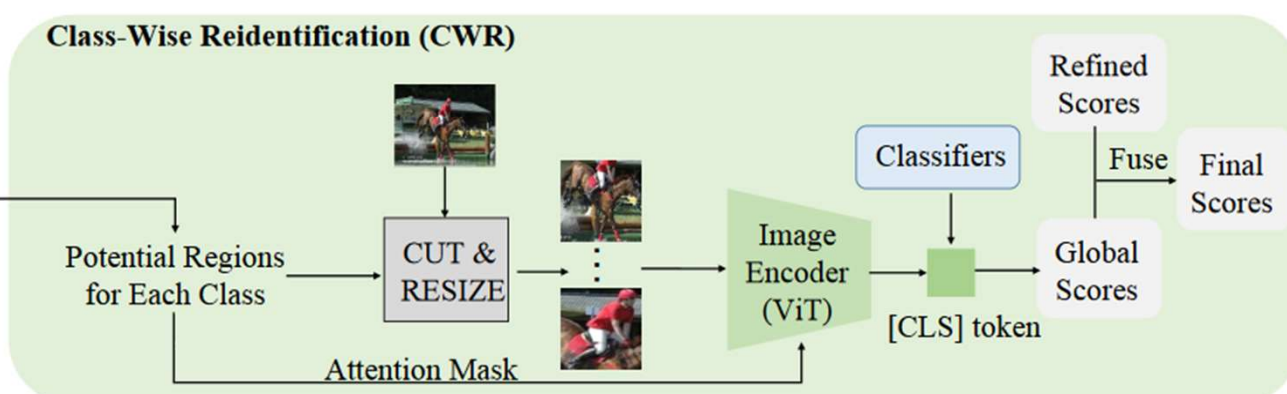
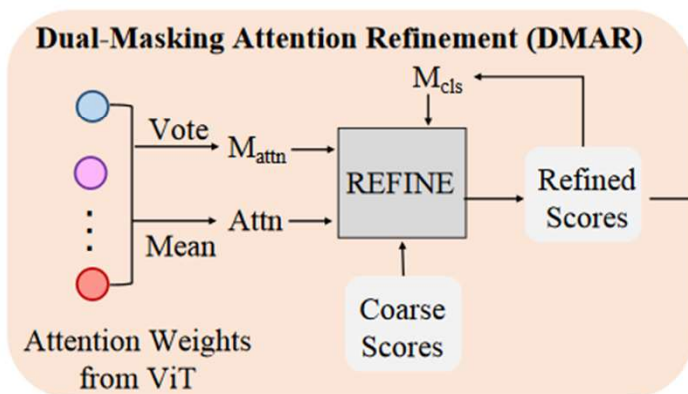
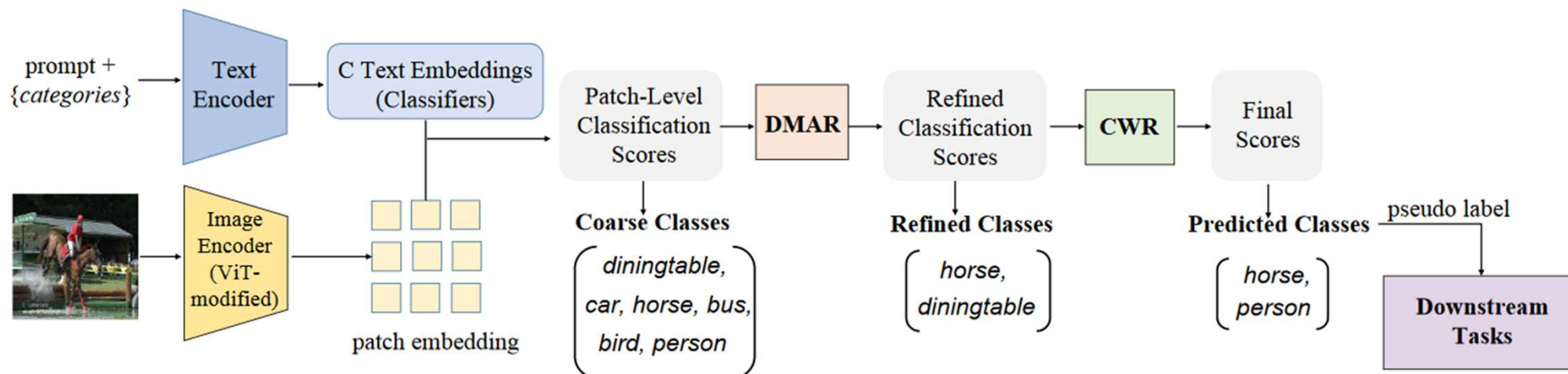


Figure 3: Qualitative results of the patch-level classification upon  $x_{dense}$  and  $x^L$  outputted by the last two layers of CLIP-ViT respectively. The last self-attention operation breaks spatial information in ViT.

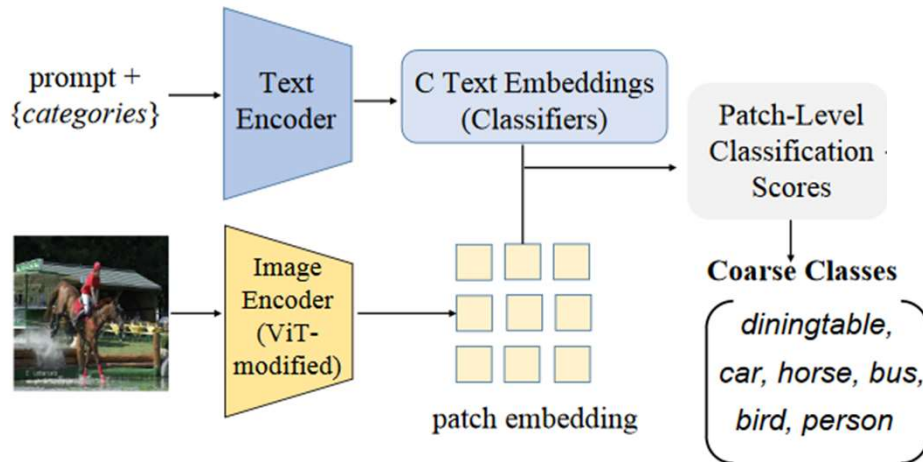
Last Self-Attention	mAP	mIoU
✓	82.7	16.2
✗	85.4	41.6

Table 1: Quantitative results for the effect of last self-attention operation in terms of classification (mAP) and segmentation (mIoU) on PASCAL VOC 2012 validation set.

# Method



# Coarse Classification

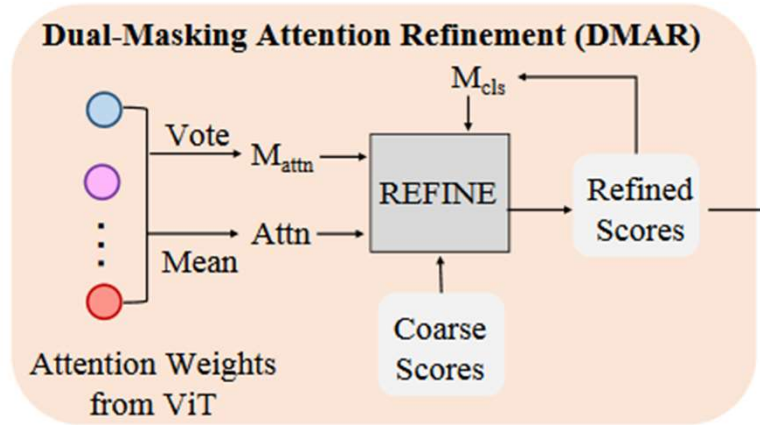


The output feature map based on the penultimate layer  $x_{dense} \in R^{N \times D}$  is leveraged. The output of the text encoder is denoted as  $T \in R^{D \times C}$ , which acts as the classifier based on the text inputs.  $N$ ,  $D$ ,  $C$  represent token length, token dimension and class number, respectively.

$$s_i = \text{Linear}(x_{dense,i}) * T,$$

$$P_{coarse}(i, c) = \frac{\exp(s_i^c)}{\sum_{k=1}^C \exp(s_i^k)}.$$

# Dual-Masking Attention Refinement (DMAR)



$$P_{refined} = \frac{1}{|\psi|} \sum_{l \in \psi} A_l * P_{coarse},$$

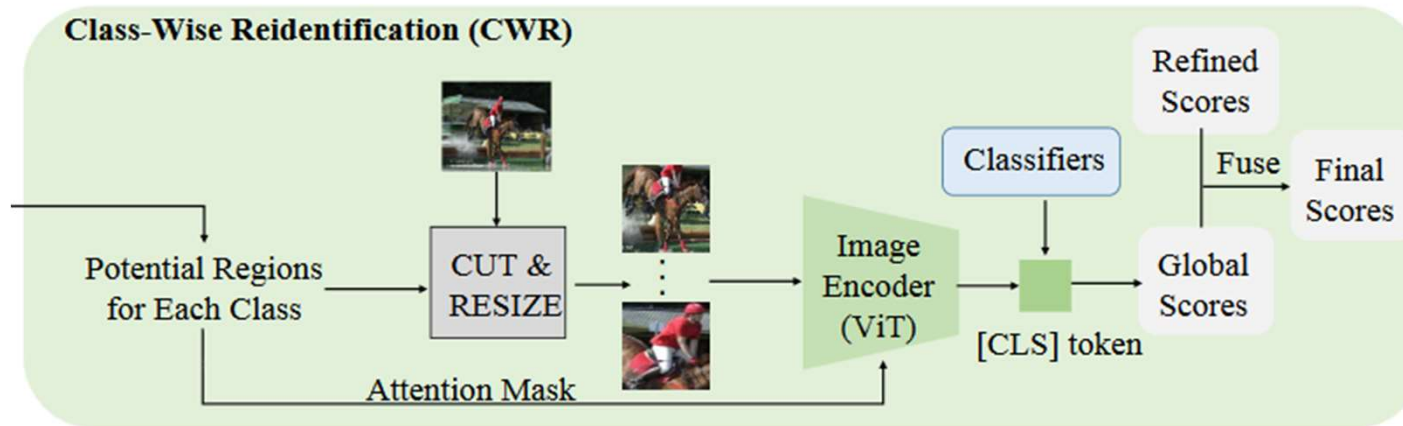
$P_{coarse} \in R^{N \times C}$  denotes the coarse score map,  $A_l \in R^{N \times N}$  represents the attention weight in the  $l$ -th layer of ViT,  $\psi$  represents the index set of used attention layer and  $|\psi|$  is its number of elements.

$$M_{attn}(i, j) = 1, \text{ if } \sum_{l=1}^L \mathbb{I}_{(A(i, j, l) > \bar{A}_l)}(A) > K,$$

$$\hat{P}_{refined} = \frac{1}{|\psi|} \sum_{l \in \psi} M_{attn} \odot A_l * P_{coarse},$$

$$P_{refined}(c) = \frac{1}{|\psi|} \sum_{l \in \psi} M_{attn} \odot A_l \odot M_{cls}(c) * P_{coarse}(c).$$

# Class Scores

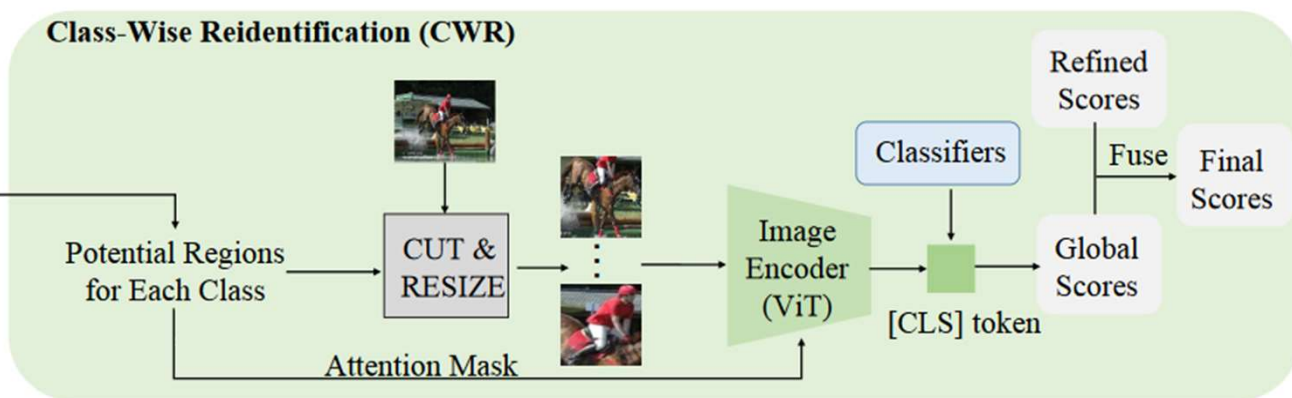
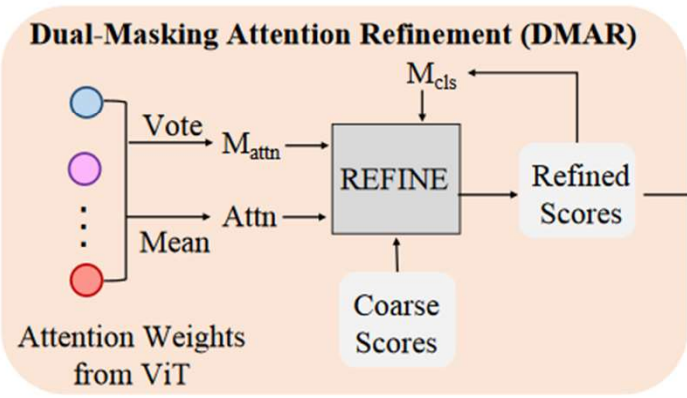
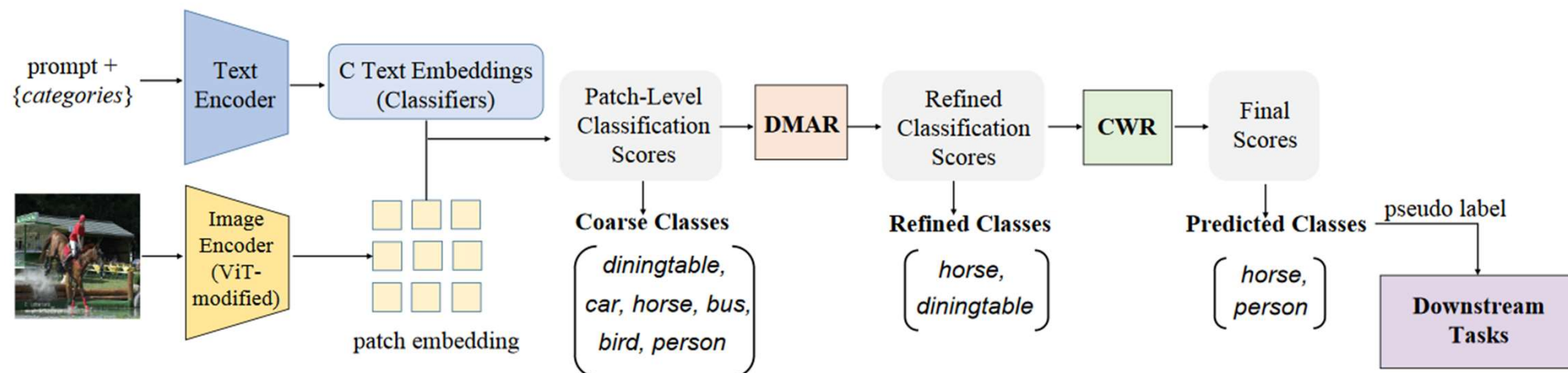


$$P_{local}(c) = \max_i(P_{coarse}(i, c)),$$

For each class, we pick out the highly responsive patches from  $P_{refined}$  and form the class-related region (class-wise mask). We crop the image by the bounding box of the region and resize it to a specific size, e.g.,  $224 \times 224$ . The class-wise mask serves as the attention mask in ViT to exclude patches that do not belong to the class. We input the class-wise image into original CLIP and use [cls] token for classification.

$$P_{final} = \lambda P_{local} + (1 - \lambda) P_{global},$$

# Method



# Experiment

Method	Extra Training Data	VOC	COCO
<i>Supervised specialist:</i>			
SARB	10% Data	83.5	75.5
DualCoOp	10% Data	90.3	78.7
TAI-DPT	10% Data	93.3	81.5
<i>Open-vocabulary generalist:</i>			
TAI-DPT	COCO captions	88.3	65.1
CLIP <sup>†</sup>	None	79.5	54.2
CLIP	None	85.8	63.3
DPT <sup>†</sup>	None	83.4	59.6
DPT	None	86.2	64.3
CLIPsurgery	None	85.4	61.2
<b>TagCLIP(Ours)</b>	None	<b>92.8</b>	<b>68.8</b>

Table 2: Experimental results of multi-label classification. <sup>†</sup> represents not using softmax on classification scores.

Method	VOC	COCO	COCO-Stuff
<i>Vanilla USS methods</i>			
IIC	9.8	-	6.7
MaskContrast	35.0	3.73	-
TransFGU	37.2	12.7	17.5
MaskDistill	45.8	-	-
PiCIE	-	-	13.8
PiCIE+H	-	-	14.4
<i>CLIP-based methods</i>			
MaskCLIP <sup>‡</sup>	42.1	20.2	23.9
CLIPsurgery <sup>‡</sup>	41.5	25.2	29.7
GroupViT	52.3	24.3	-
SegCLIP	52.6	26.5	-
ReCo	34.2	17.1	26.3
NamedMask	59.2	27.7	-
<b>CLS-SEG (Ours)</b>	<b>64.8</b>	<b>34.0</b>	<b>30.1</b>
<b>CLS-SEG*(Ours)</b>	<b>68.7</b>	<b>35.3</b>	<b>31.0</b>

Table 3: Results of annotation-free semantic segmentation. The vanilla USS results are based on K-means clustering. <sup>‡</sup> represents we re-implement it with the same experimental setting as ours. \* means using denseCRF to postprocess.

# Experiment

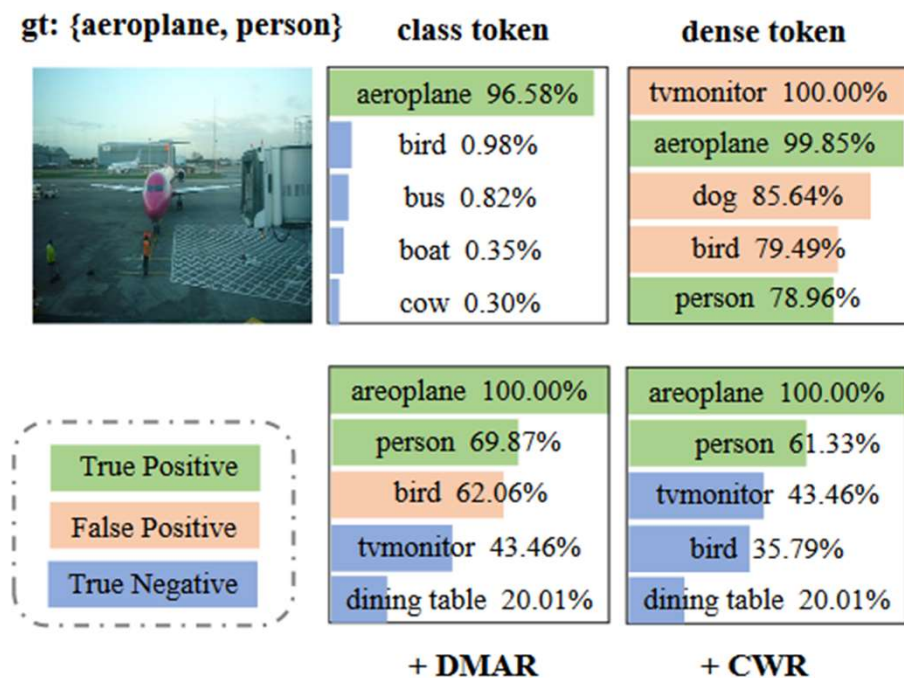
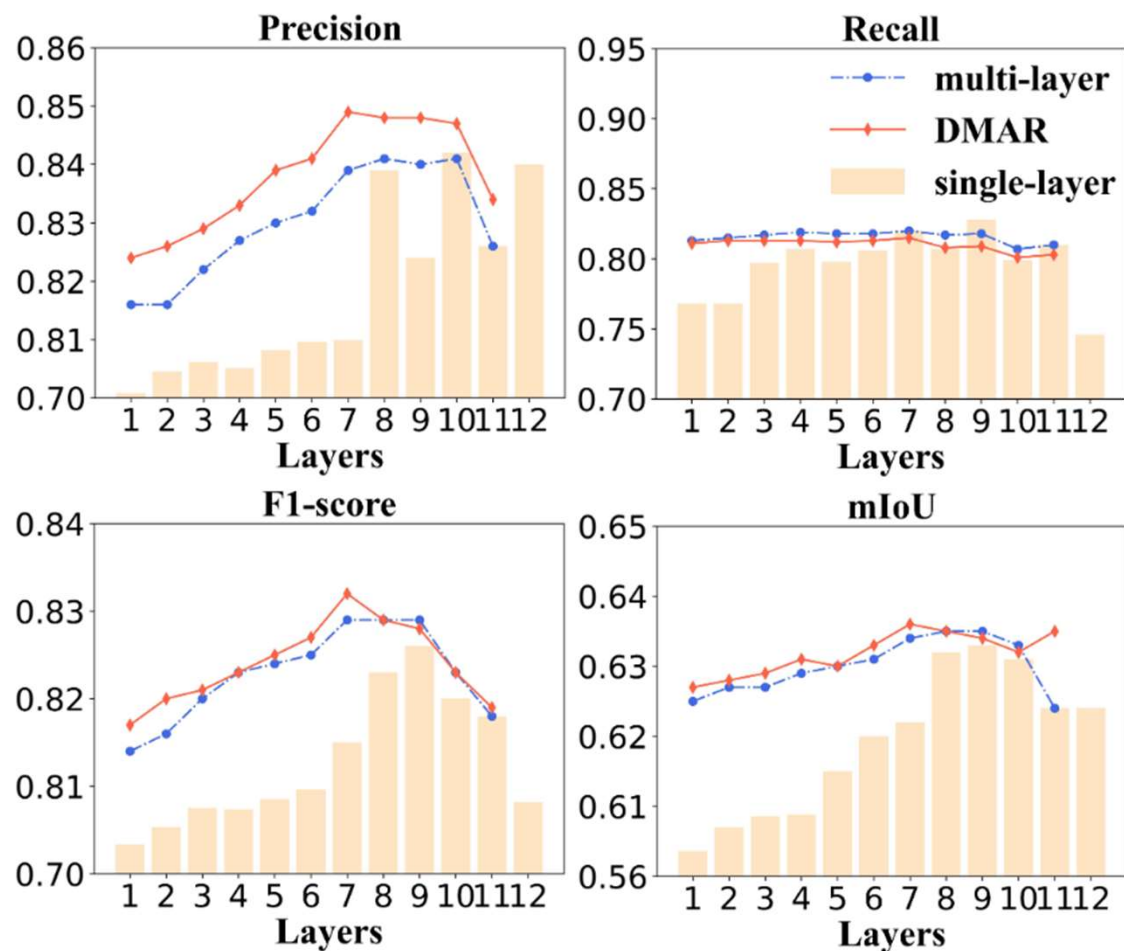


Figure 5: The classification results of different strategies. We use 0.5 as the threshold by default.

Coarse Score	DMAR	CWR	mAP	mIoU
✓			85.4	30.9
✓		✓	88.0	55.2
✓	✓		93.9	63.7
✓	✓	✓	94.1	64.8

Table 4: Results for the effectiveness of DMAR and CWR module in terms of classification and semantic segmentation. The results are evaluated on the VOC 2012 val set.

# Experiment



Comparison of single-layer and multi-layer attention refinement in terms of classification and segmentation tasks. For the single-layer setting, each tick  $i$  on the x-axis represents merely adopting attention weight in  $i$ -th layer. For the multi-layer setting, the  $i$ -th x-tick means fusing  $i$ -th to 11-th layers attention weights to refine coarse classification scores. We rule out the last attention layer during fusing.

**Thanks**