



# Open-Vocabulary Segmentation

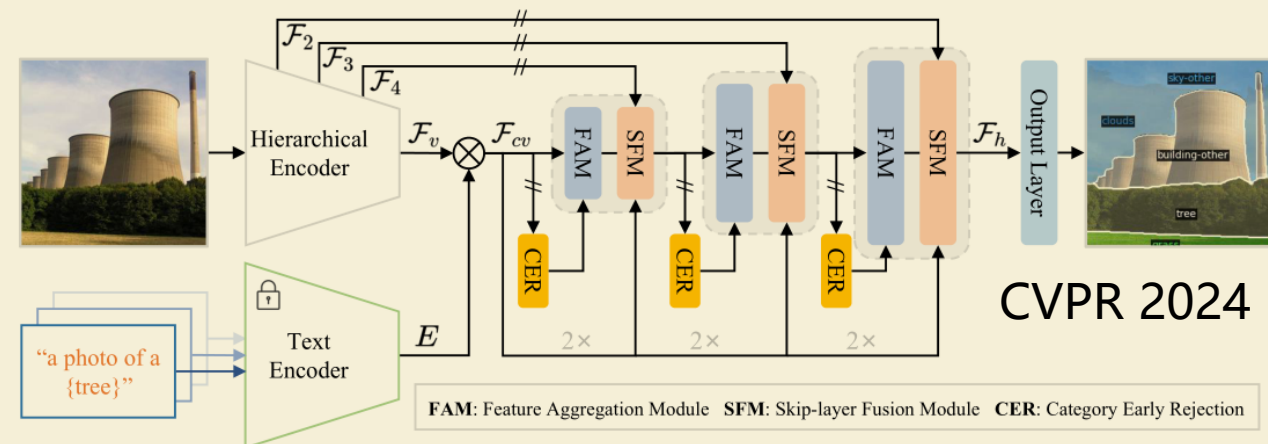
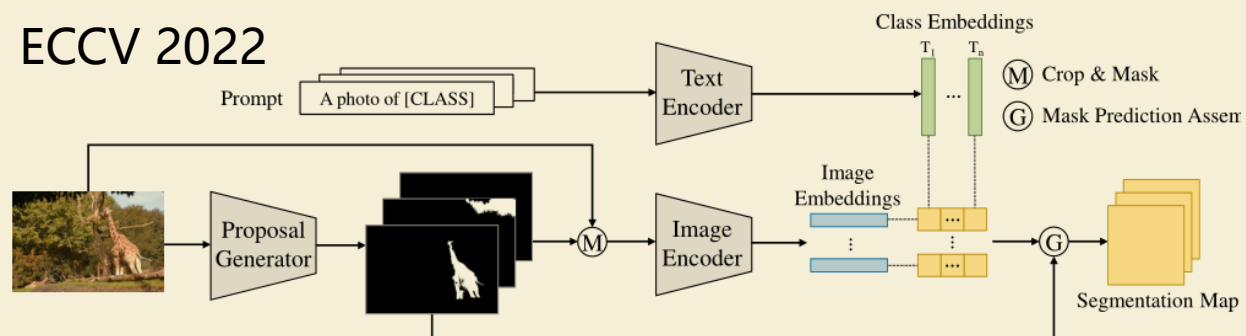
常世杰

2023.3.31

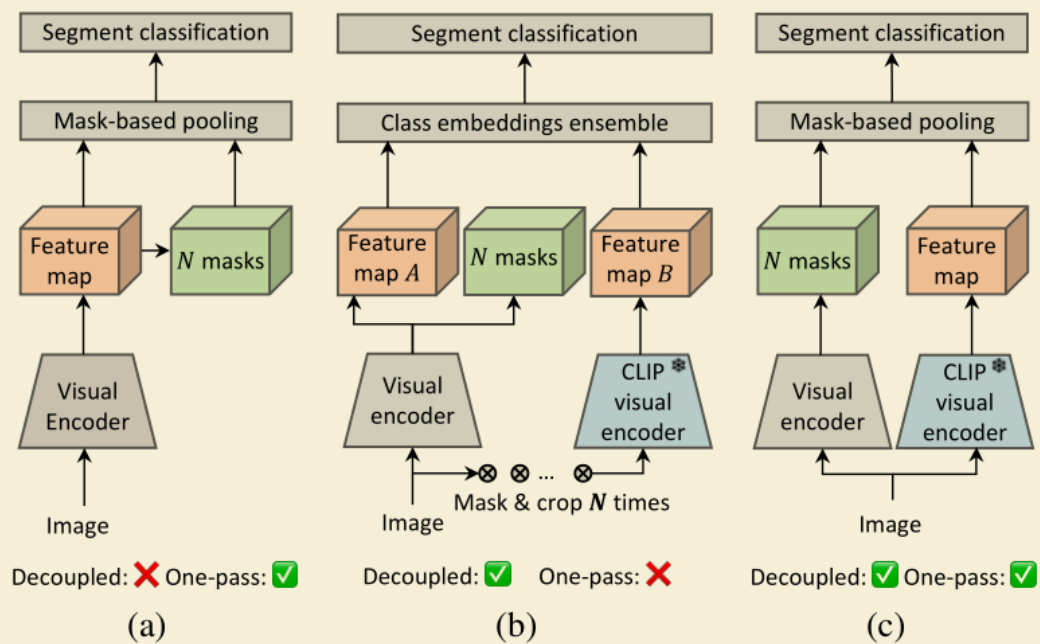


# Fully-Supervised OV SS

ECCV 2022

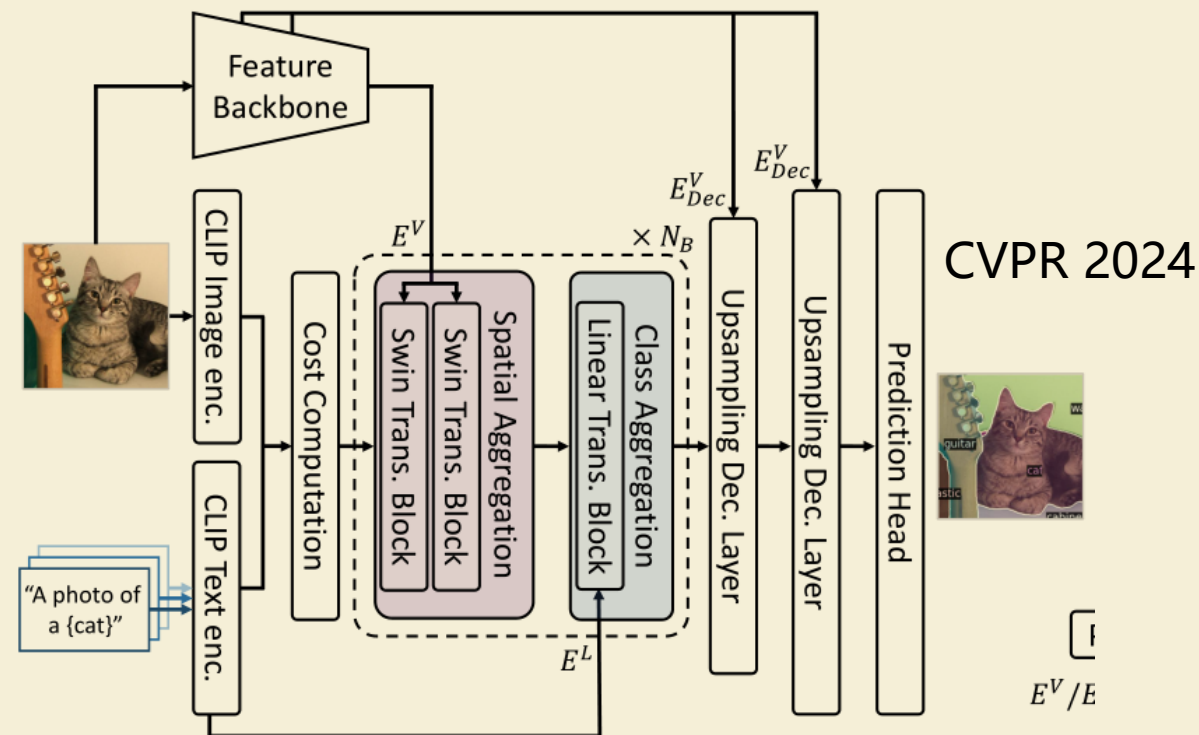


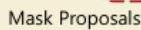
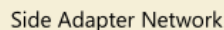
## Two-stage OVSeg



ICCV 2023

## One-stage OVSeg



A photograph showing two people riding a tandem bicycle on a paved street. The person in the front is wearing a red jacket and a grey cap, and the person in the back is also wearing a red jacket. They are both looking forward. In the background, there is a brick building and a parked car.

CVPR 2023 SAN

of Tech. (DUT), P. R. China.

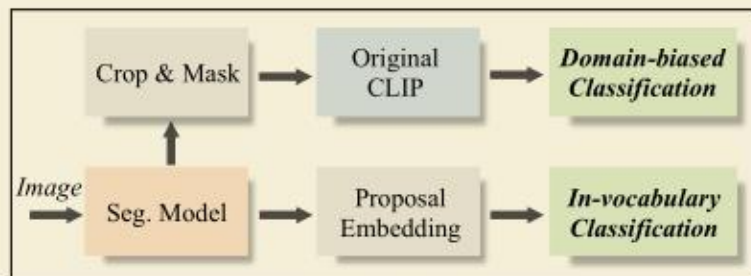


A segmented image of a street scene. The image is divided into colored regions representing different objects. Labels are placed over these regions: 'sky' (blue), 'house' (tan), 'tree' (green), 'car' (red), 'wall' (blue), 'road, route' (grey), 'bus' (black), 'sidewalk, pavement' (red), 'bridge, span' (green), and 'signboard, sign' (black).

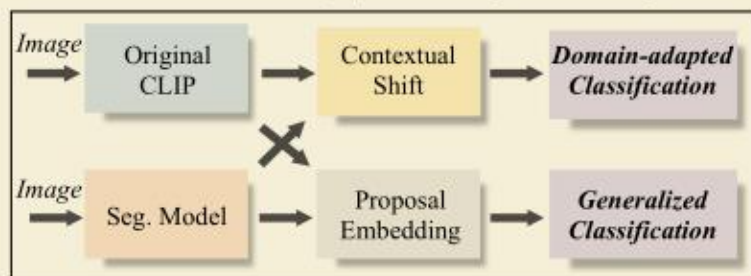
CVPR 2024 SCAN

## Mask proposal (mask2former) + CLIP Classification

## 23 SAN



(a) existing two-stage methods



natural image



masked image

Ori. CLIP	0.41	0.30	0.17	..
Pred:	van	truck	car	
CS. CLIP	0.71	0.07	0.03	..
Pred:	road	dirt	sidewalk	

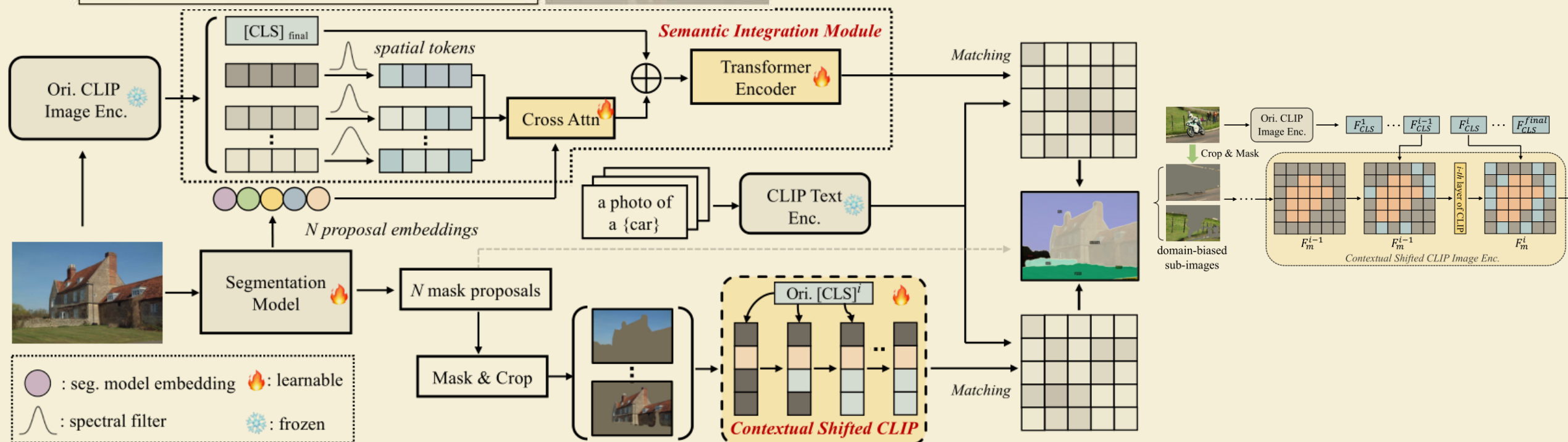


natural image



masked image

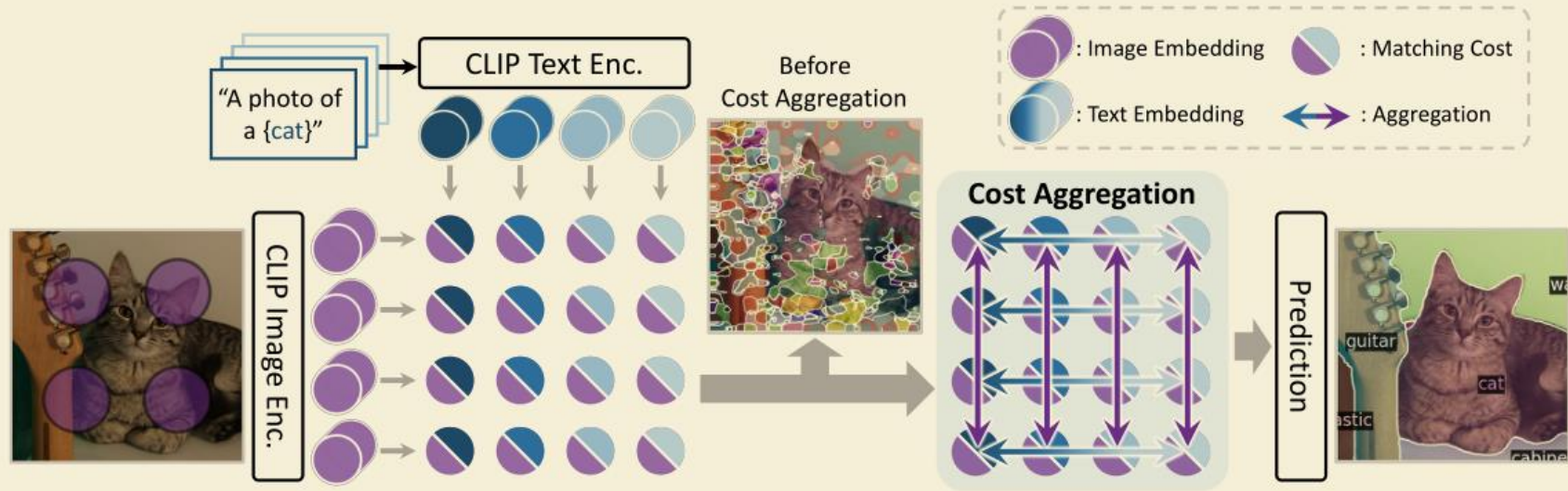
Ori. CLIP	0.999	7e-4	6e-5	..
Pred:	plane	sky	runway	
CS. CLIP	0.989	5e-3	1e-3	..
Pred:	sky	plane	canopy	



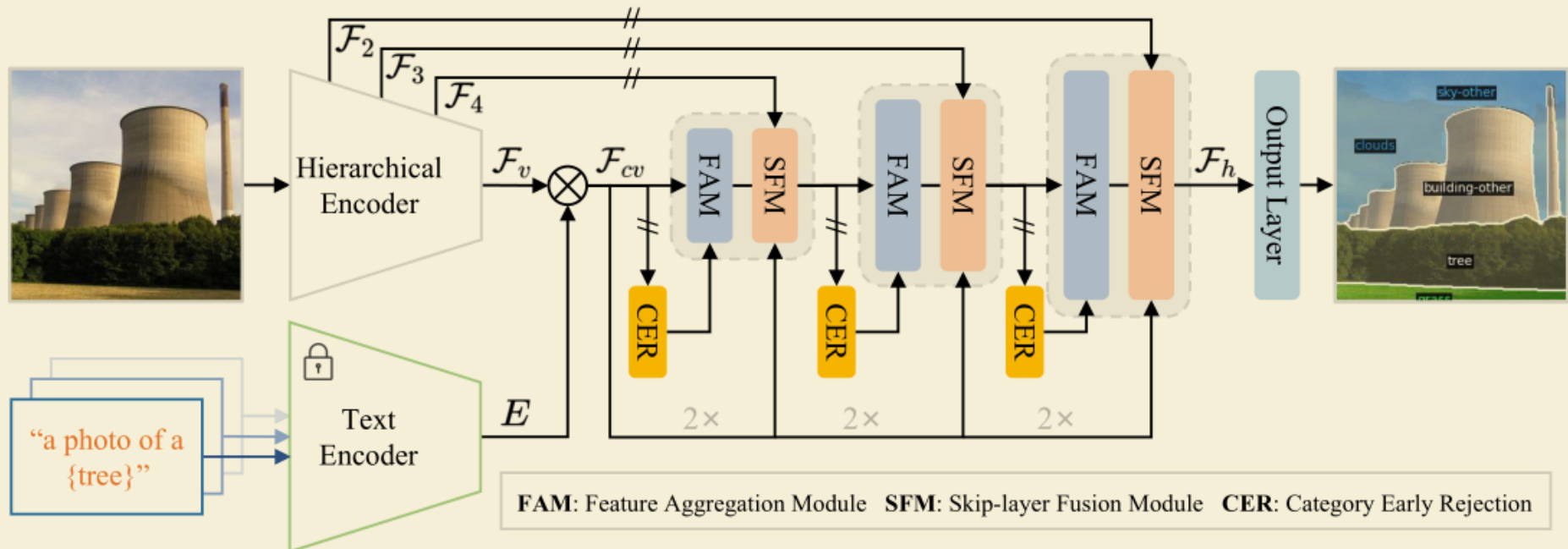


# One-stage OVSeg

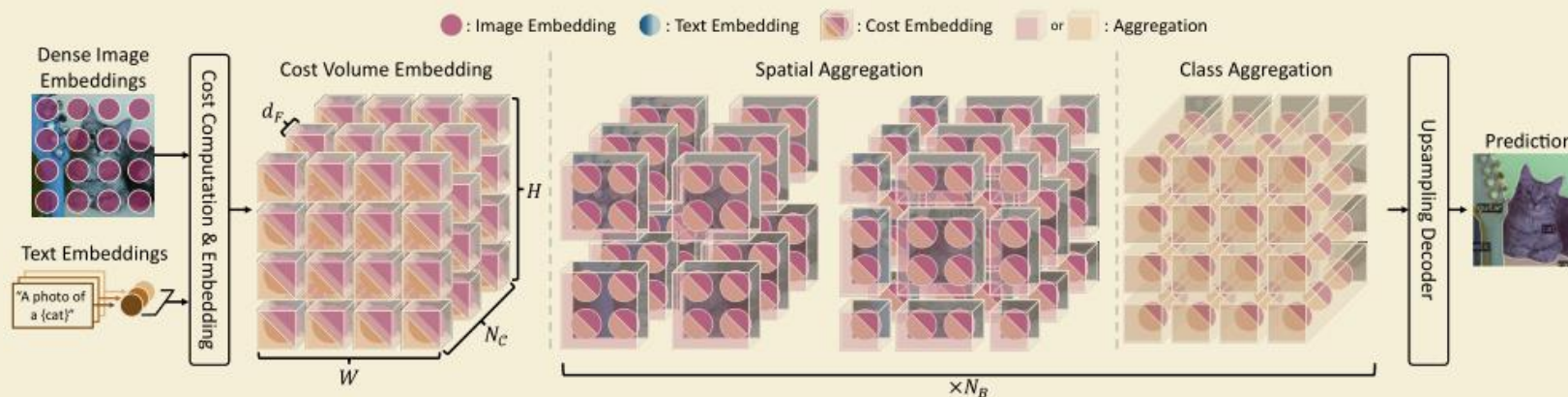
CVPR 2024  
CAT-Seg



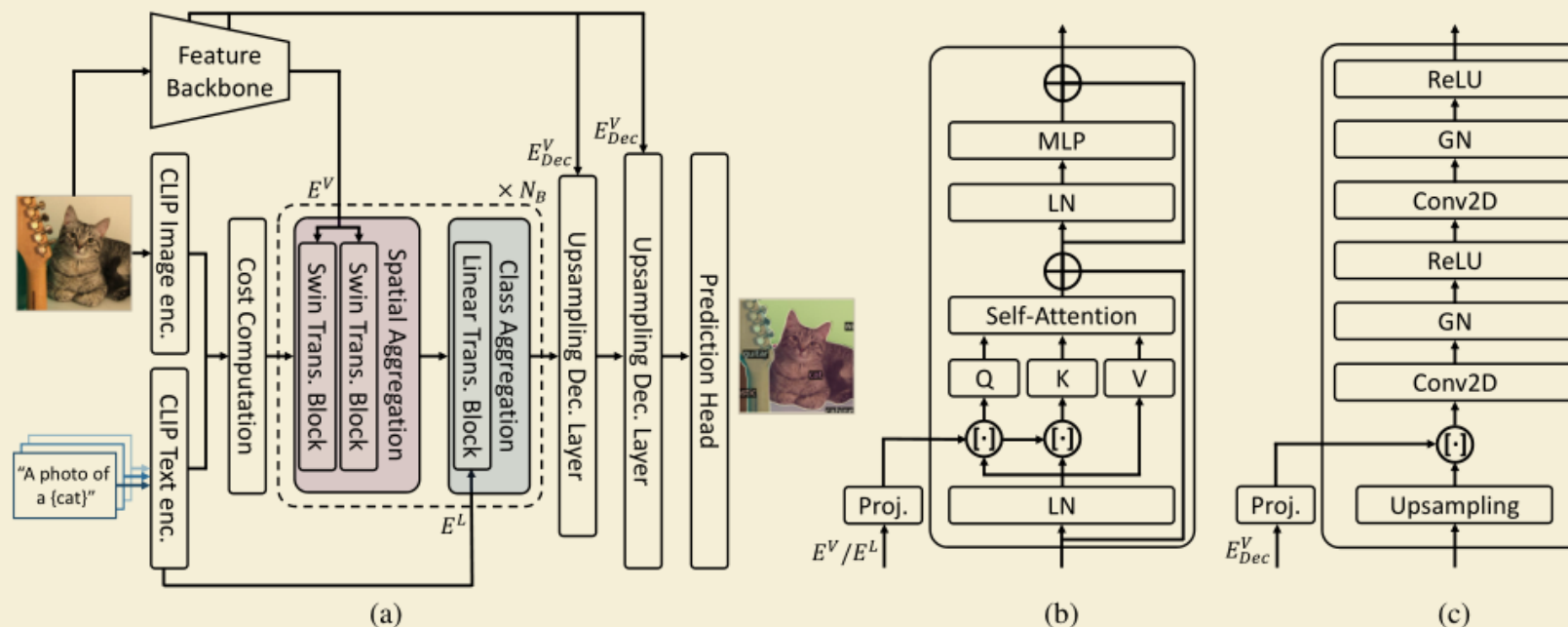
CVPR 2024 SED



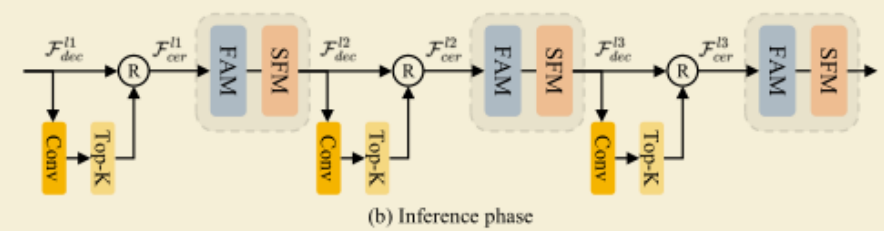
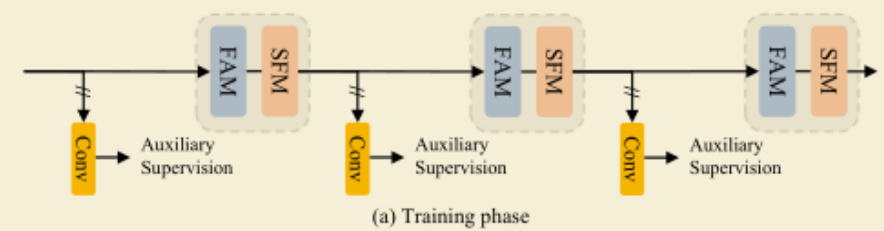
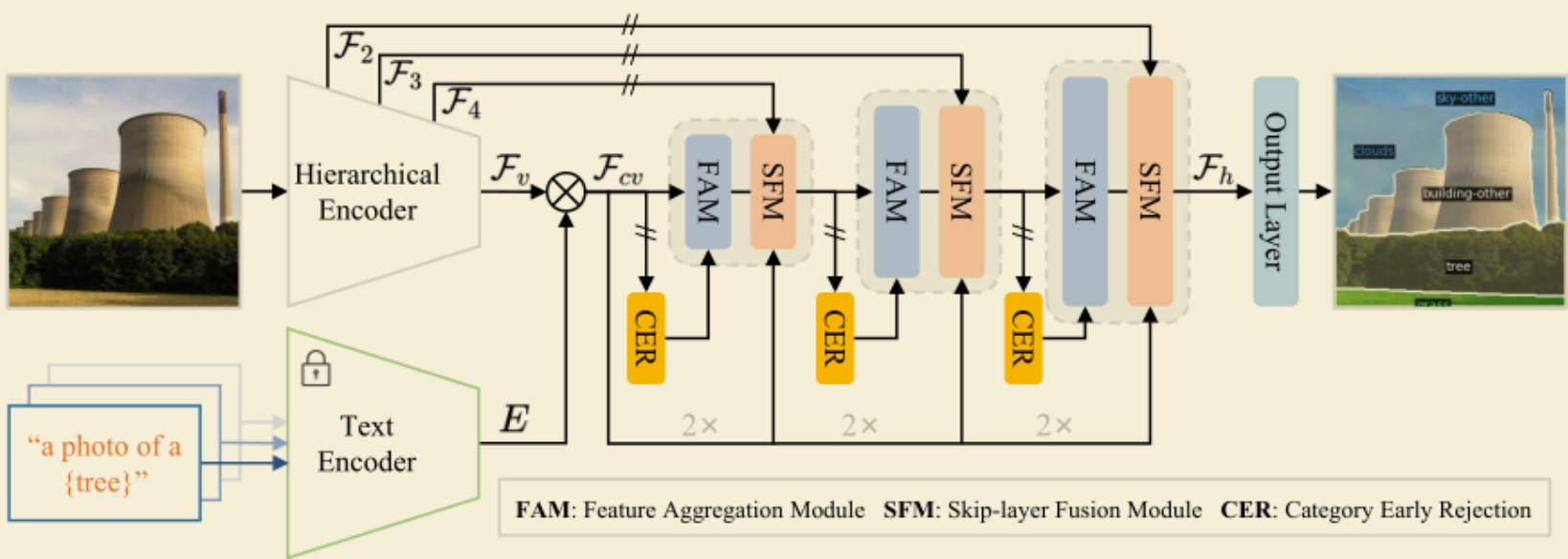
Seokju Cho<sup>\*,1</sup> Heeseong Shin<sup>\*,1</sup> Sunghwan Hong<sup>1</sup> Seungjun An<sup>1</sup> Seungjun Lee<sup>1</sup>,  
 Anurag Arnab<sup>2</sup> Paul Hongsuck Seo<sup>2</sup> Seungryong Kim<sup>1</sup>  
<sup>1</sup>Korea University <sup>2</sup>Google Research



$N, H, W$



Bin Xie<sup>1</sup>, Jiale Cao<sup>1</sup>, Jin Xie<sup>2</sup>, Fahad Shahbaz Khan<sup>3</sup>, and Yanwei Pang<sup>1,4</sup>  
<sup>1</sup>Tianjin University <sup>2</sup>Chongqing University

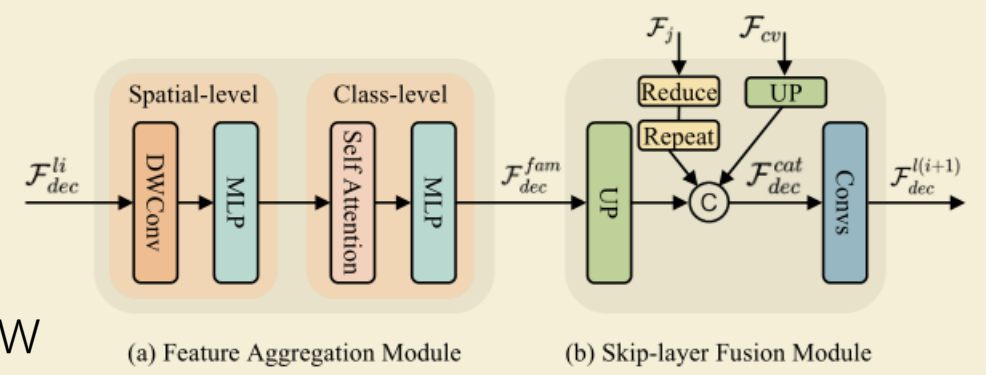


双阶段法用两个独立的网络预测mask和类别——速度慢  
 现有单阶段法用ViT作为视觉backbone——局部表示差

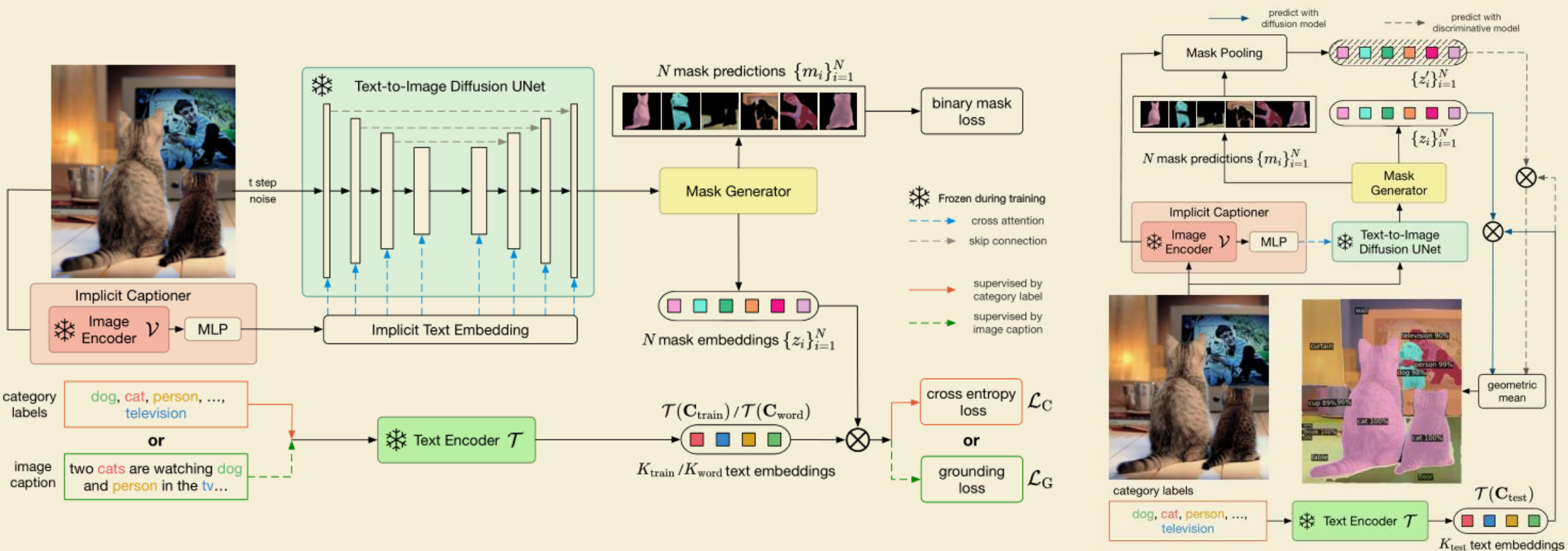
利用层次化结构带来更多局部信息

提出CER，提前过滤图片中不存在的类别，提高推理速度

NPC  
 CHW  
 NDHW   N1HW

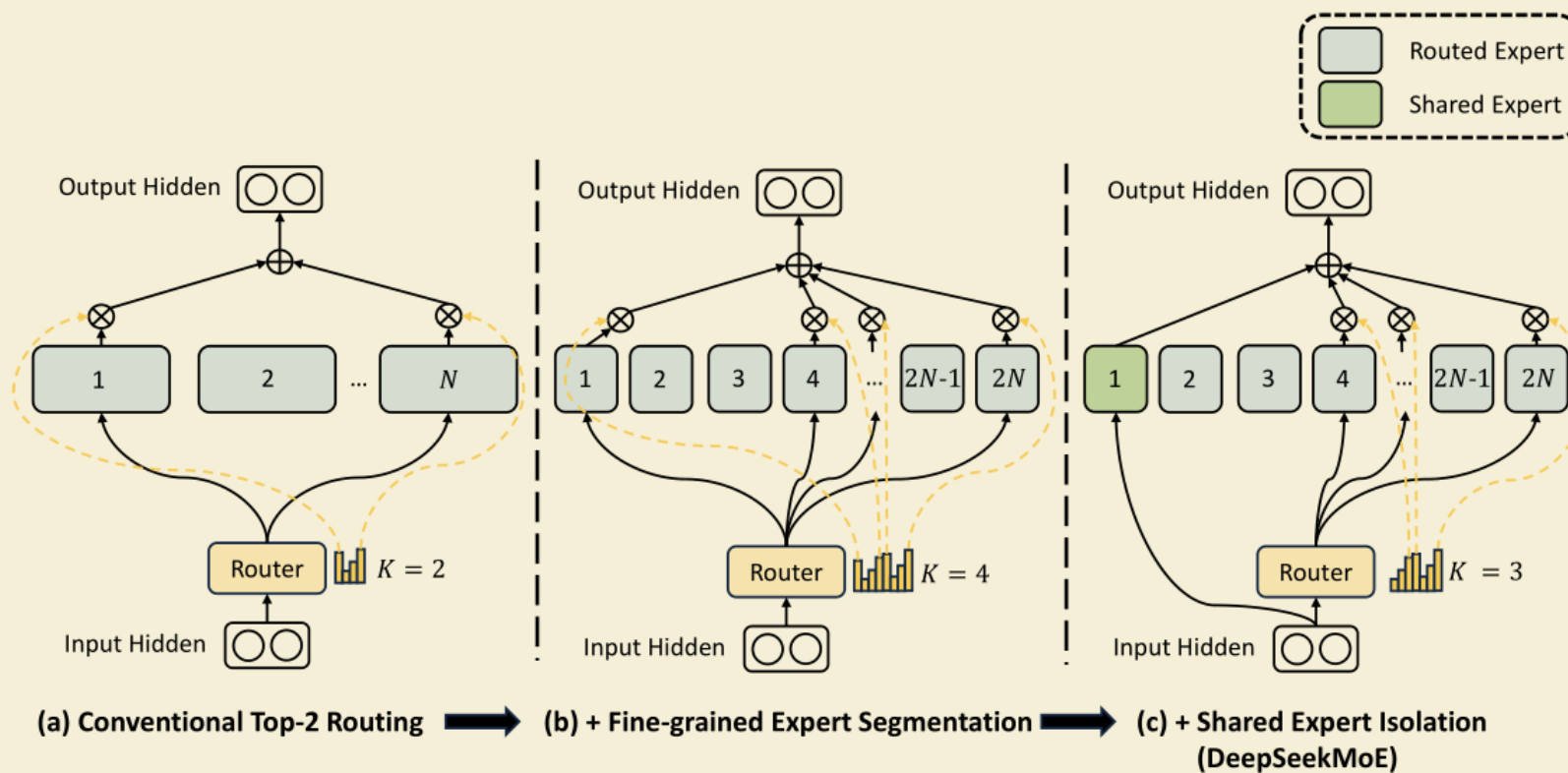




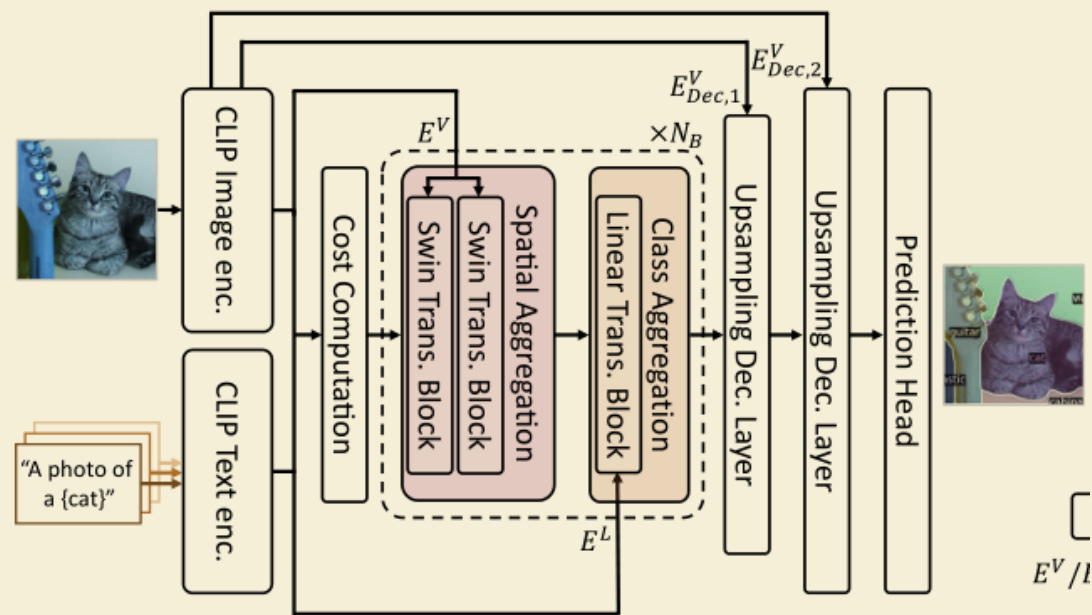
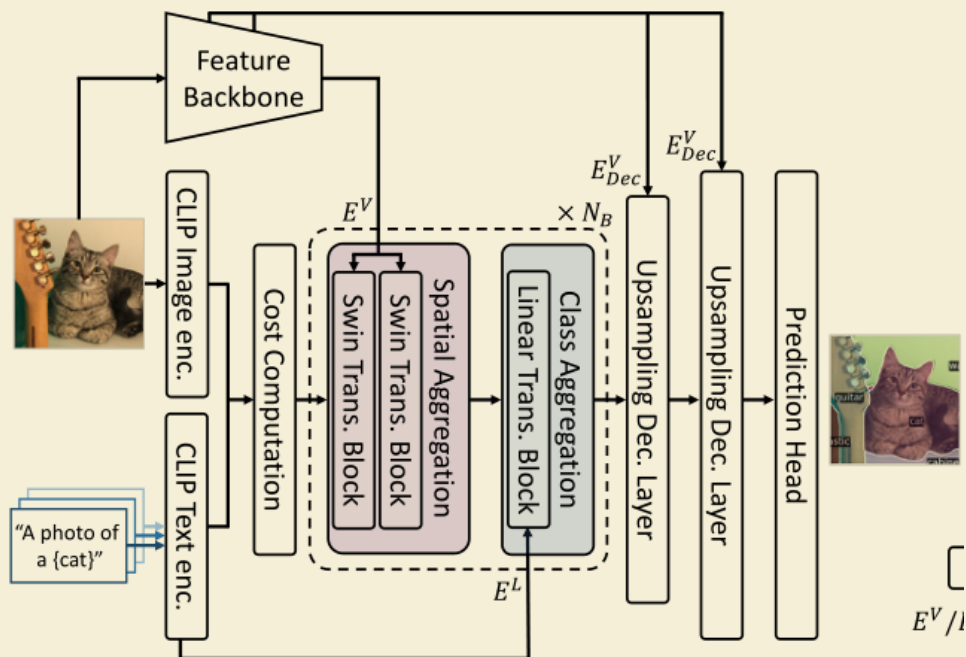
Jiarui Xu<sup>1\*</sup> Sifei Liu<sup>2†</sup> Arash Vahdat<sup>2†</sup> Wonmin Byeon<sup>2</sup>Xiaolong Wang<sup>1</sup> Shalini De Mello<sup>2</sup><sup>1</sup>UC San Diego <sup>2</sup>NVIDIA



# DeepSeekMoE: Towards Ultimate Expert Specialization in Mixture-of-Experts Language Models



	A-847	PC-459	A-150	PC-59	PAC-20
SED-B	11.4	18.6	31.6	57.3	94.4
moe	10.9	17.7	31.5	57.5	94.1
deepseek	10.8	18.1	31.5	57.1	93.6



	A-847	PC-459	A-150	PC-59	PAC-20
CAT-v1-B	8.4	16.6	27.2	57.5	93.7
SED-B	11.4	18.6	31.6	57.3	94.4
CAT-v2-B	12.0	19.0	31.8	57.5	94.6
SED-L	13.9	22.6	35.2	60.6	96.1
CAT-v2-L	16.0	23.8	37.9	63.3	97.0

特征图基于cost volume, 导致特征图大小受类别数影响( $B, C, H, W$ ) -> ( $B, T, C, H, W$ ) solov2  
 预测逻辑: 为每个类别T预测一个二值mask  
 训练集和测试集类别域偏移太大, 测试集很多类别  
 $IoU=0$

# CAT-Seg-v1 v.s. SED v.s. CAT-Seg-v2

CAT-Seg-v1

CLIP-ViT + feature backbone

Fine-tune visual encoder attention

Train 384

Test sliding inference 384/640

Imagenet templates

SED

CLIP-ConvNext

Fine-tune full visual encoder

Train 768

Test 640

Imagenet templates

CAT-Seg-v2

CLIP-ViT

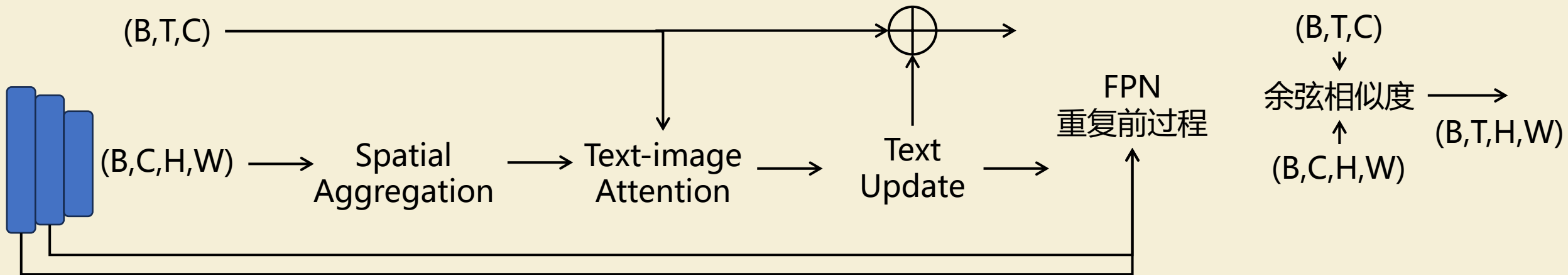
Fine-tune q/v visual&textual encoder

Train 384

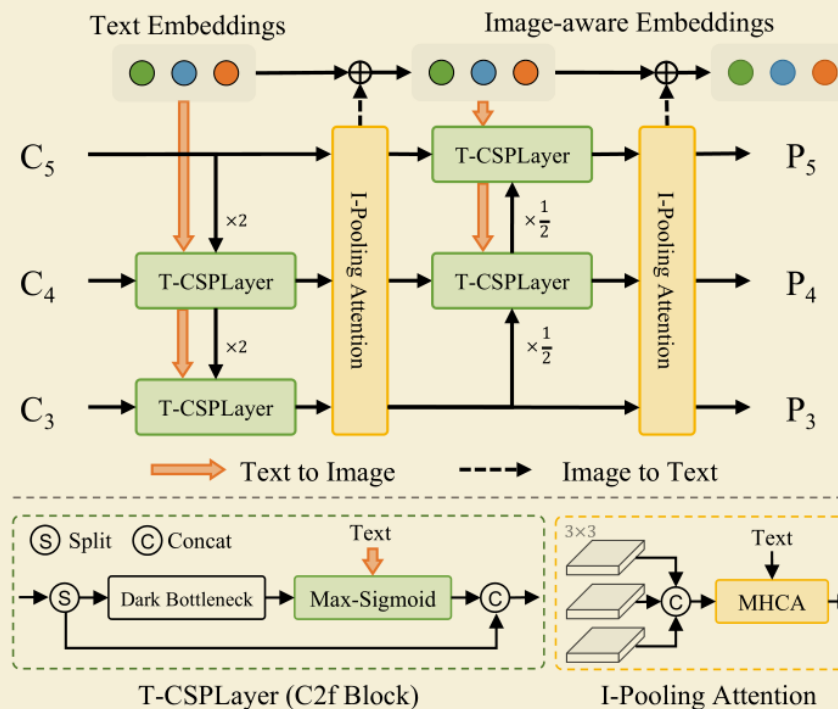
Test sliding inference 384/640

Single template

ViT-Base 2d8h



	A-847	PC-459	A-150	PC-59	PAC-20
SED-B	11.4	18.6	31.6	57.3	94.4
	3.1	9.0	17.3	49.3	89.1

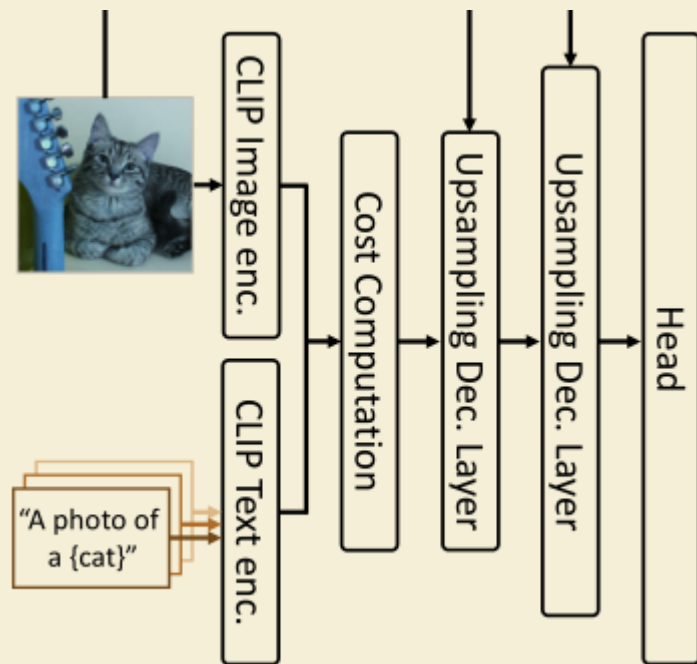


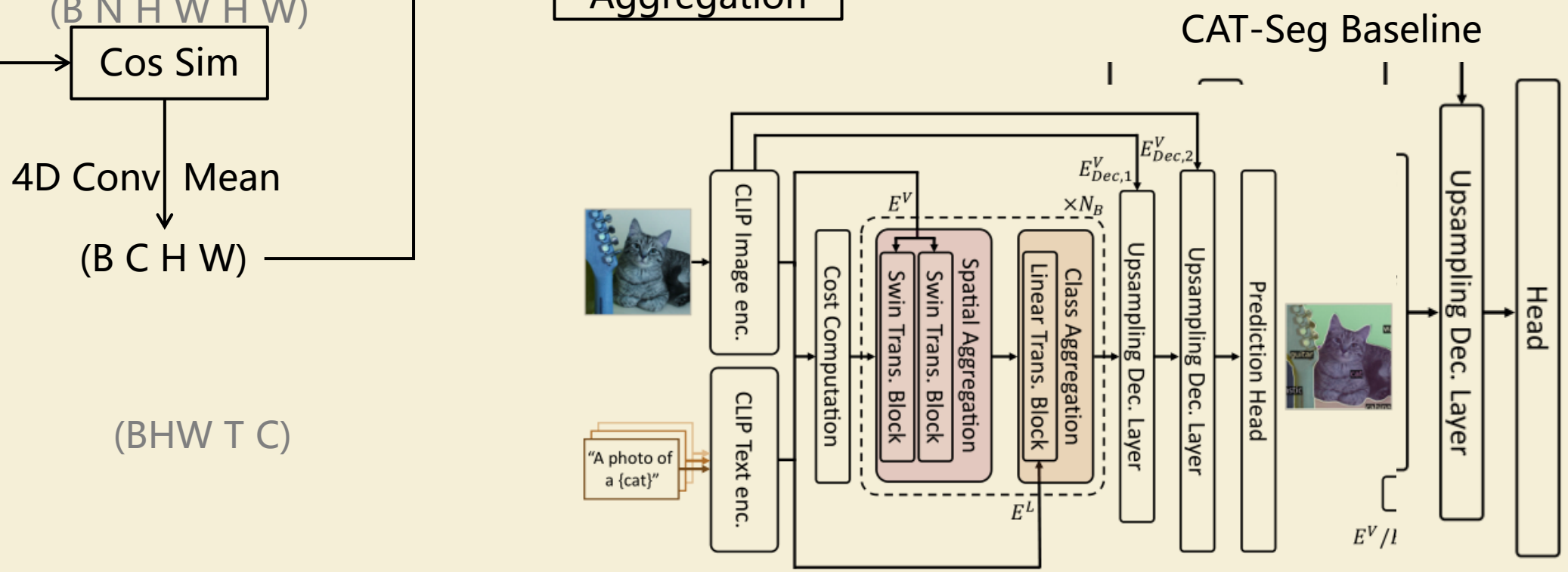
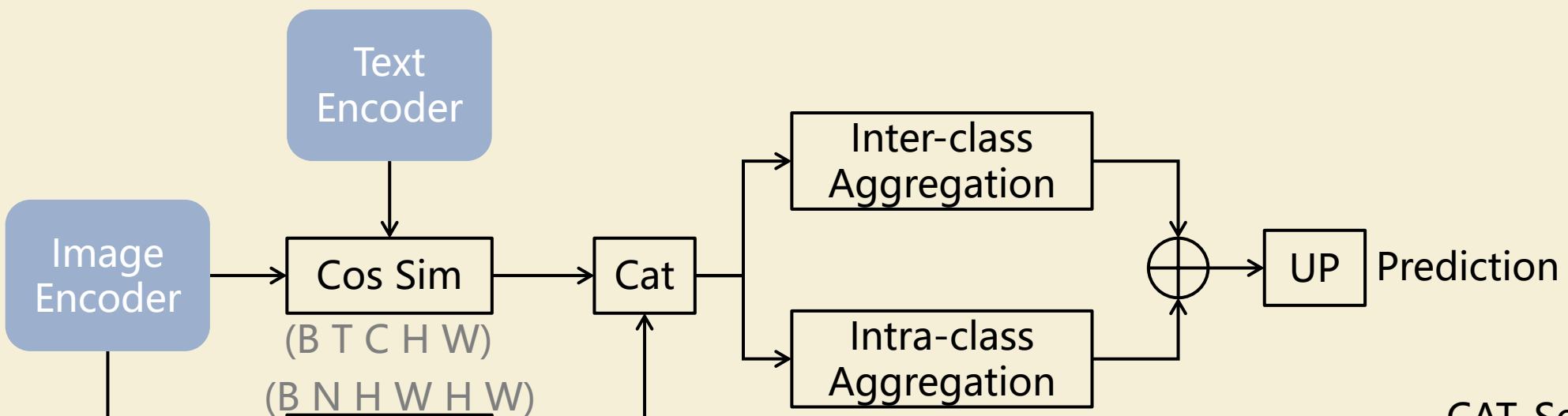


	A-847 8.00	PC-459 20	A-150 5.31	PC-59 9.01	PAC-20 2.0
CAT-B	12.1	19.1	31.6	57.4	94.7
baseline	10.0	16.6	27.0	52.6	93.8
baseline2	10.8	17.5	28.6	54.0	94.2
	11.1	17.8	30.7	56.0	94.4
	11.6	18.6	31.1	56.8	94.8
dino	10.1	16.5	27.1	52.5	93.9
deformcl	11.6	18.9	31.5	57.2	94.8

	A-847 8.00	PC-459 20	A-150 5.31	PC-59 9.01	PAC-20 2.0
baseline2	10.8	17.5	28.6	54.0	94.2
+aug	10.9	17.7	28.6	54.4	94.6
+ms	10.8	17.5	28.4	53.9	94.5
dcn	11.4	18.3	31.1	55.7	94.8
+cls	12.0	19.5	31.9	58.2	94.4
Oracle	28.7	43.2	51.1	75.3	97.7

Components		A-847	PC-459	A-150	PC-59	PAS-20	PAS-20 <sup>b</sup>
(I)	Feature Agg.	5.6	12.8	23.6	58.1	96.3	77.7
(II)	Cost Agg.	14.7	<u>23.2</u>	35.3	60.3	<u>96.7</u>	78.9
(III)	(II) + Spatial agg.	14.9	23.1	35.9	60.3	<u>96.7</u>	79.5
(IV)	(II) + Class agg.	14.7	21.5	36.6	60.6	95.5	80.5
(V)	(II) + Spatial and Class agg.	<u>15.5</u>	<u>23.2</u>	<u>37.0</u>	<u>62.3</u>	<u>96.7</u>	<u>81.3</u>
(VI)	(V) + Embedding guidance	<b>16.0</b>	<b>23.8</b>	<b>37.9</b>	<b>63.3</b>	<b>97.0</b>	<b>82.5</b>
		-	-	-	-	52.3	-
		8.1	11.5	26.4	44.8	-	<u>70.2</u>
		9.0	12.4	29.6	55.7	94.5	-
		<u>12.4</u>	<u>15.7</u>	<u>32.1</u>	<u>57.7</u>	<u>94.6</u>	-
		11.1	14.5	29.9	57.3	-	-
		<b>16.0</b>	<b>23.8</b>	<b>37.9</b>	<b>63.3</b>	<b>97.0</b>	<b>82.5</b>
		(+3.6)	(+8.1)	(+5.8)	(+5.6)	(+2.4)	(+12.3)





	A-847 8.00	PC-459 20	A-150 5.31	PC-59 9.01	PAC-20 2.0
CAT-B	12.1	19.1	31.6	57.4	94.7
baseline	10.0	16.6	27.0	52.6	93.8
baseline2	10.8	17.5	28.6	54.0	94.2
i-i corr	11.1	17.8	30.7	56.0	94.4
i-i corr2	11.6	18.6	31.1	56.8	94.8
dino	10.1	16.5	27.1	52.5	93.9
maskMH A	11.6	18.9	31.5	57.2	94.8
dcn	11.4	18.3	31.1	55.7	94.8

去掉spatial和class aggregation和decoder中的浅层特征

去掉spatial和class aggregation

加入i-i correlation

在baseline2的基础上将浅层特征替换为dino对应层特征

Class aggregation时使用mask MHA，仅和相近的K个类别做MHA  
K初步设置为16，感觉应该增大K试一下

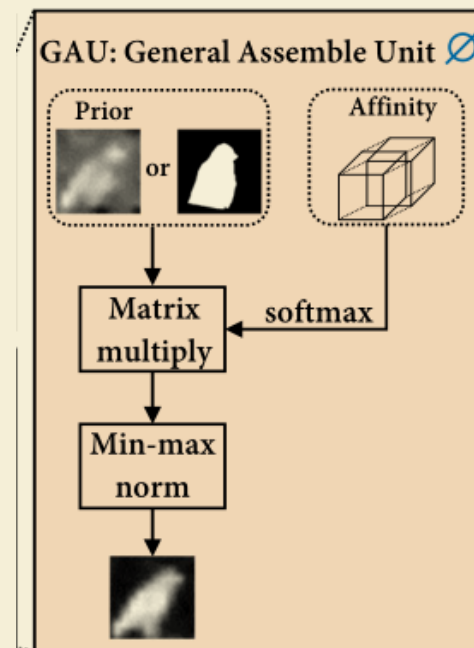
降点原因可能是原spatial agg都使用了大核 CAT-Seg 12, SED 9  
初始DCNv3核设置为3

- 获取更好的初始化cost volume
- 更好的cost aggregation方式

对spatial和class都使用deformable操作:

Spatial, 改一下DCNv3中生成offset和mask的过程

Class, 使用mask MHA仅和相近的K个类别做MHA



	A-847 8.00	PC-459 20	A-150 5.31	PC-59 9.01	PAC-20 2.0
CAT-B	12.1	19.1	31.6	57.4	94.7
oracle	28.7	43.2	51.1	75.3	97.9
Cls_head	12.2	19.2	31.9	58.1	94.7
APL	12.1	19.1	32.0	58.2	94.9
MLDecod er	12.1	19.5	31.5	58.0	95.0
+RAM	11.3	19.1	31.6	57.4	94.8

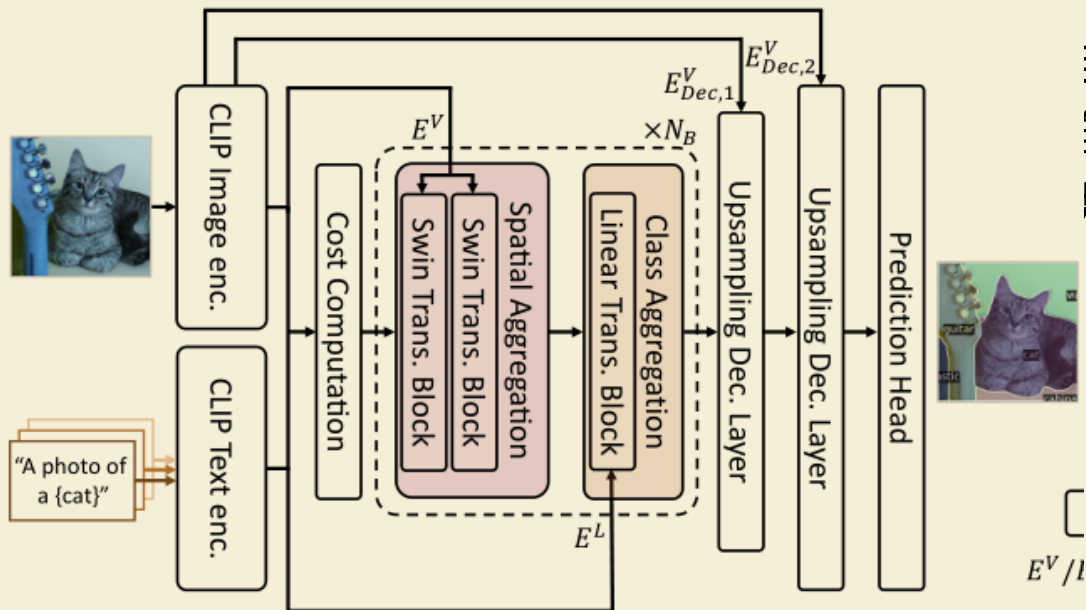
给定图片中存在的类别

Transformer cross-attention + BCE Loss

Transformer cross-attention + AP Loss

MLDecoder+ AP Loss

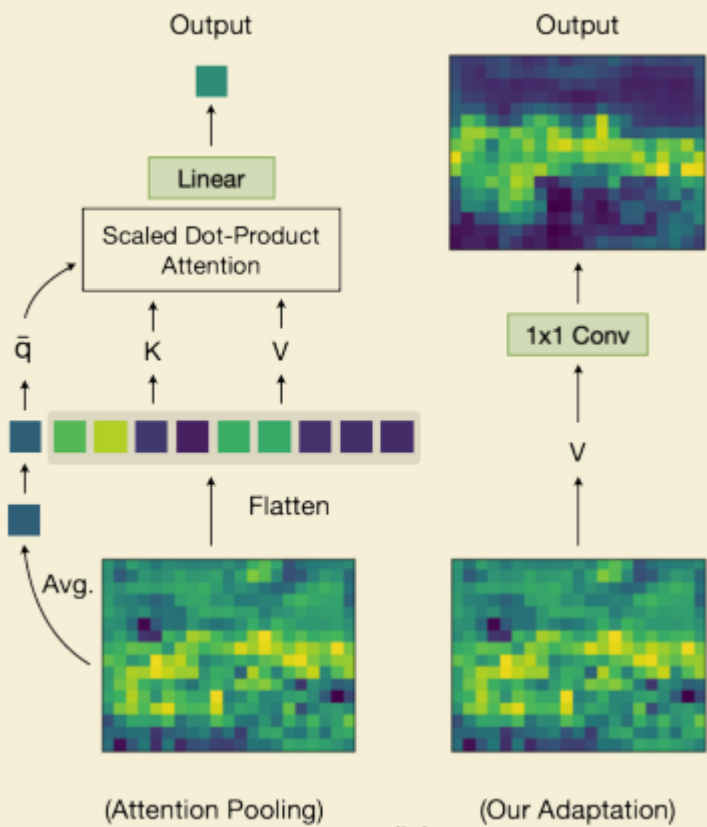
使用RAM进行分类直接测试性能



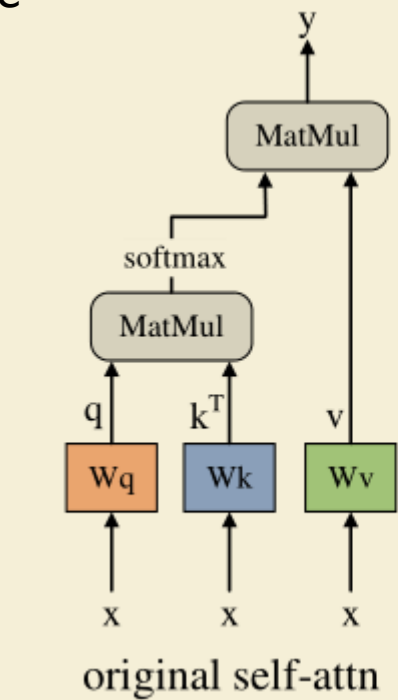
5分类联合训练后，性能会有一点提升，但过用分类头处理类别则效果不好。分类头不能做到囊括所有存在的类别



调整CLIP以生成更好的cost volume



ECCV 2022 MaskCLIP



$$Attn = \text{Softmax} \left( X W_q W_q^T X^T / \tau \right) + \text{Softmax} \left( X W_k W_k^T X^T / \tau \right),$$

SCLIP

Method		Fair	V21	PC60	C-Obj	V20	City	PC59	ADE	C-Stf	Avg
CLIP [27]	ICML'21	✓	18.6	7.8	6.5	49.1	6.7	11.2	3.2	7.2	13.8
MaskCLIP [46]	ECCV'22	✓	43.4	23.2	20.6	74.9	24.9	26.4	11.9	16.7	30.3
ReCo [28]	NeurIPS'22	✗	25.1	19.9	15.7	57.7	21.6	22.3	11.2	14.8	23.5
GroupViT [36]	CVPR'22	✗	52.3	18.7	27.5	79.7	18.5	23.4	10.4	15.3	30.7
TCL [7]	CVPR'23	✗	55.0	30.4	31.6	83.2	24.3	33.9	17.1	22.4	37.2
CLIP Surgery [18]	Arxiv'23	✓	41.2	30.5	-	-	31.4	-	12.9	21.9	-
SCLIP [32]	Arxiv'23	✓	61.7	31.5	32.1	<b>83.5</b>	34.1	36.1	17.8	23.9	40.1
GEM [4]	CVPR'24	✓	46.2	-	-	-	-	32.6	15.7	-	-
CLIP-DIY [35]	WACV'24	✗	59.0	-	30.4	-	-	-	-	-	-
FOSSIL [3]	WACV'24	✗	-	-	-	-	23.2	35.8	18.8	24.8	-
<b>NACLIP</b> [PAMR]	Ours	✓	58.8	32.2	33.2	77.1	35.5	35.2	17.4	22.9	39.0
<b>NACLIP</b>	Ours	✓	<b>62.4</b>	<b>35.0</b>	<b>36.2</b>	80.6	<b>38.3</b>	<b>38.4</b>	<b>19.1</b>	<b>25.2</b>	<b>41.9</b>

NACLIP

