



# **AutoFocusFormer:**

## **Image Segmentation off the Grid**

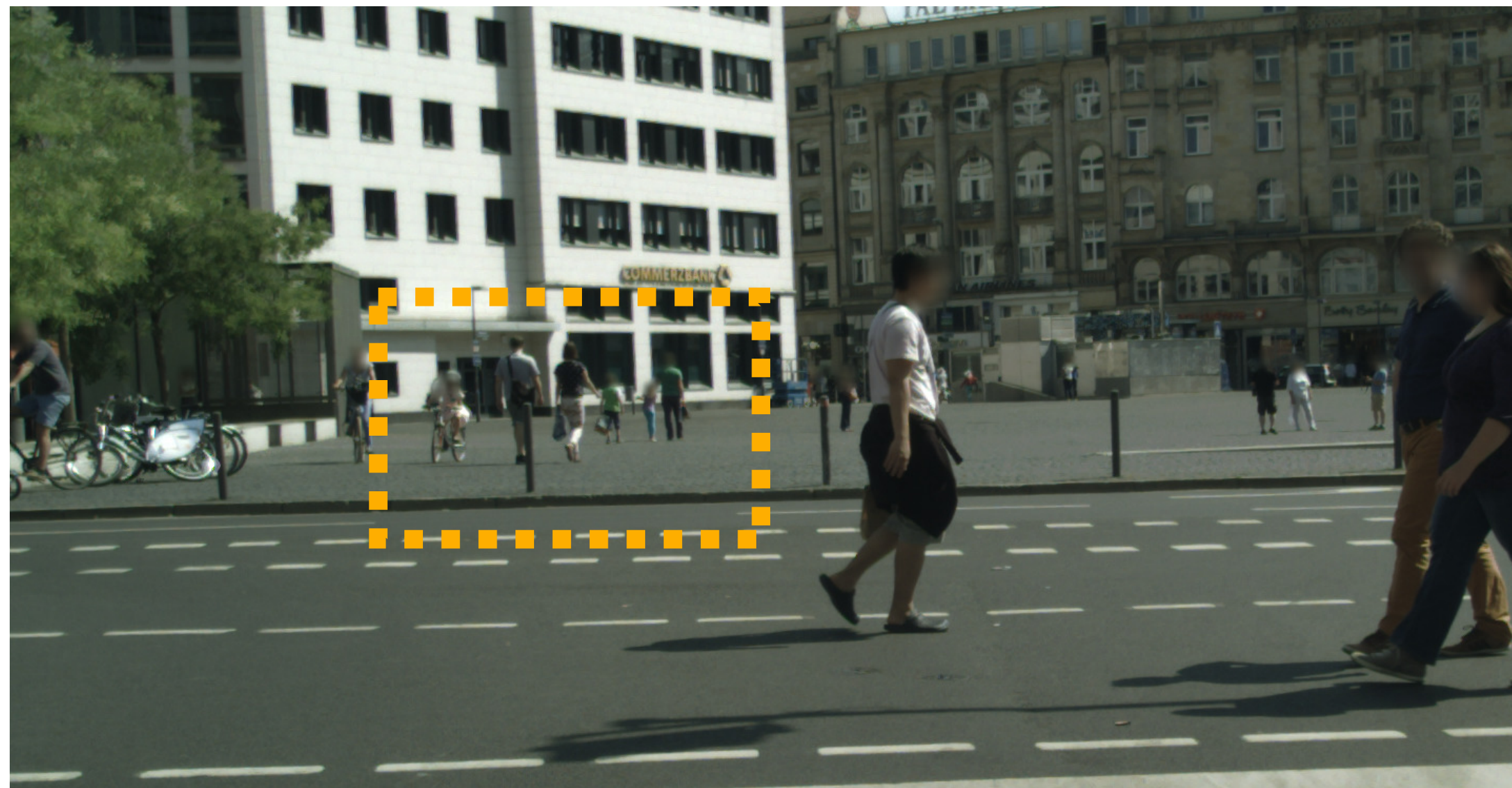
**THU-AM-167**

**Chen Ziwen, Kaushik Patnaik, Shuangfei Zhai, Alvin Wan, Zhile Ren, Alex Schwing,  
Alex Colburn, Li Fuxin**

2023.6 | Apple Inc. & Oregon State University

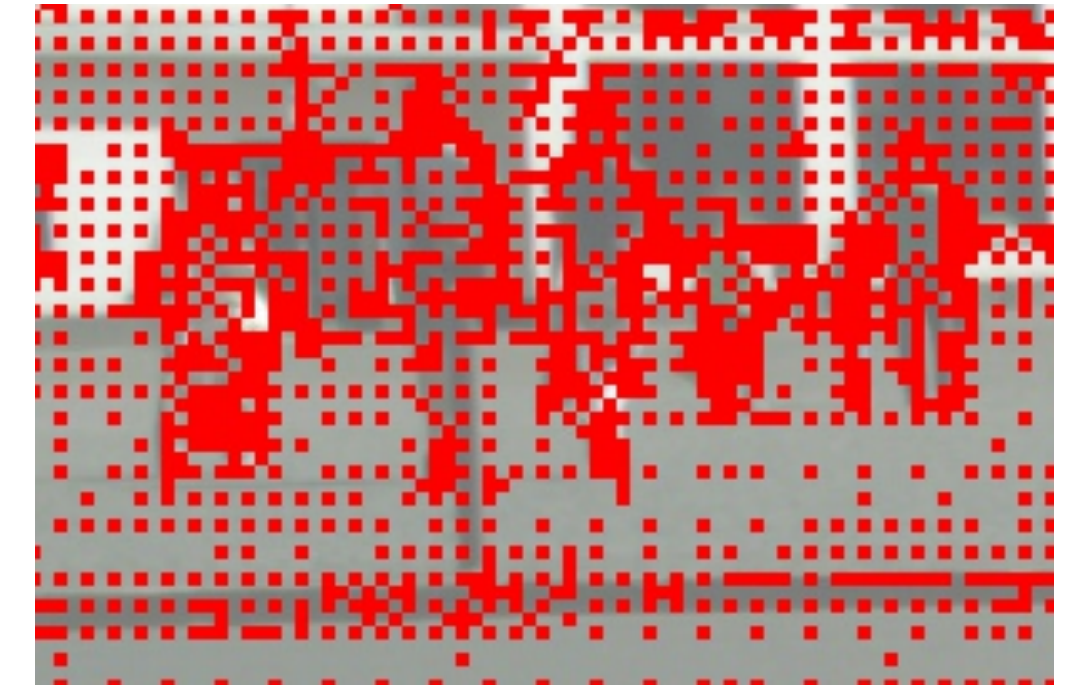
# AutoFocusFormer(AFF)

Multi-stage, local-attention transformer, equipped with successive adaptive downsampling layers

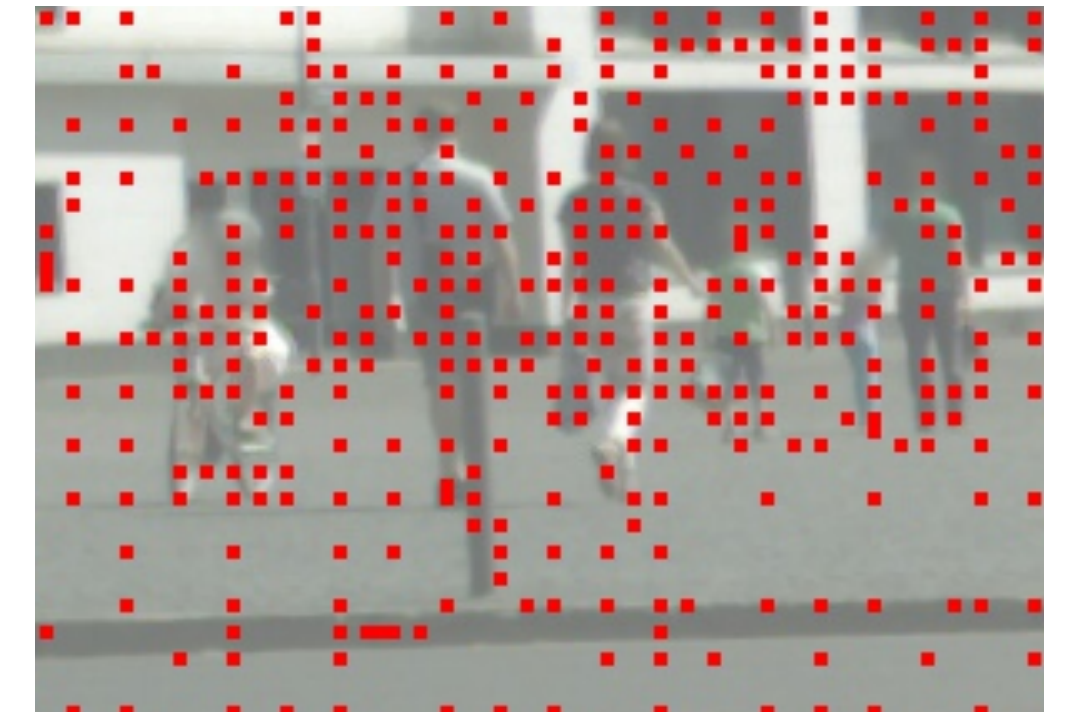


Image

Stage 2



Stage 3

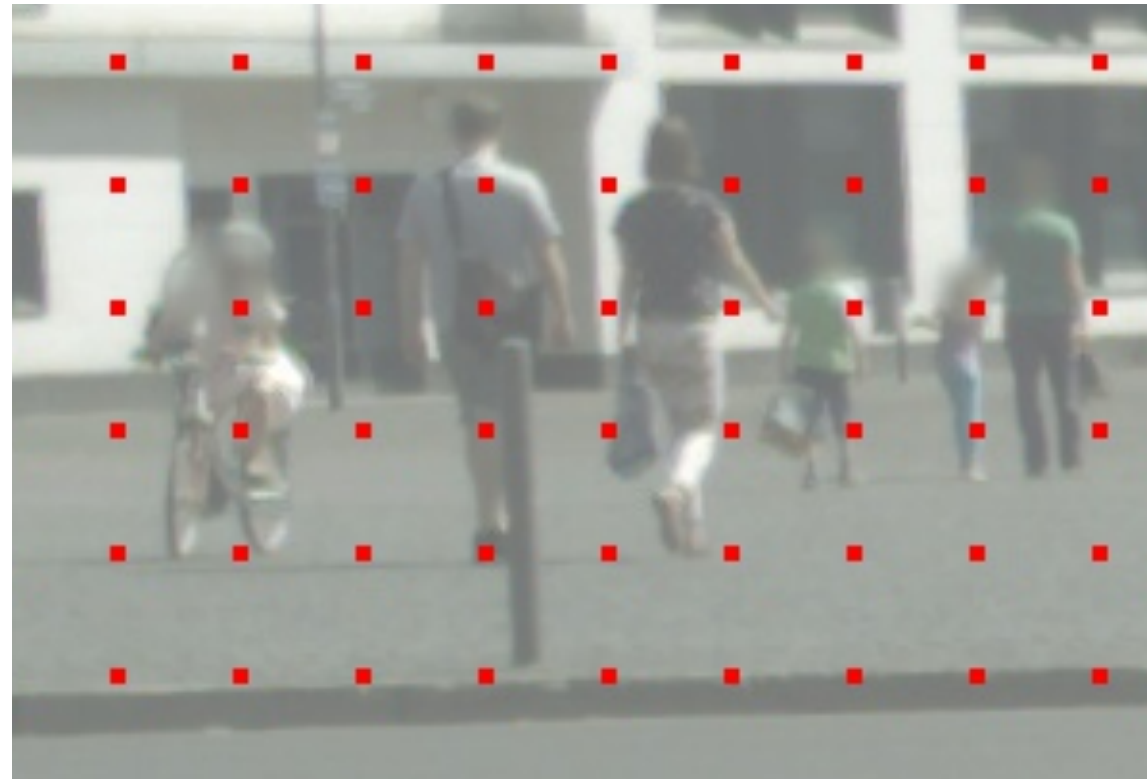


Stage 4

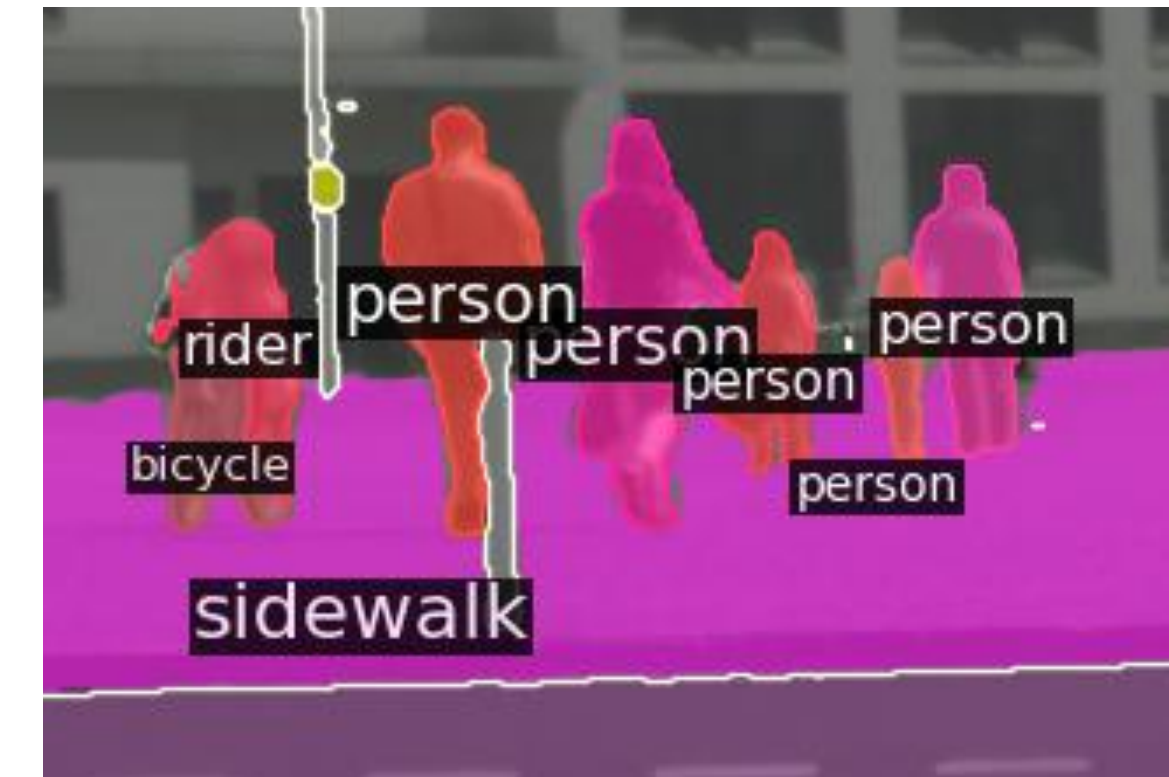


Red pixels indicate remaining tokens

# AutoFocusFormer(AFF)



Panoptic Segmentation  
With Swin backbone



Panoptic Segmentation  
With AFF backbone



# Motivation

Natural images often have highly imbalanced content density

Mainstream networks use grid downsampling

- Mis-classify small objects
- Waste computation on large objects

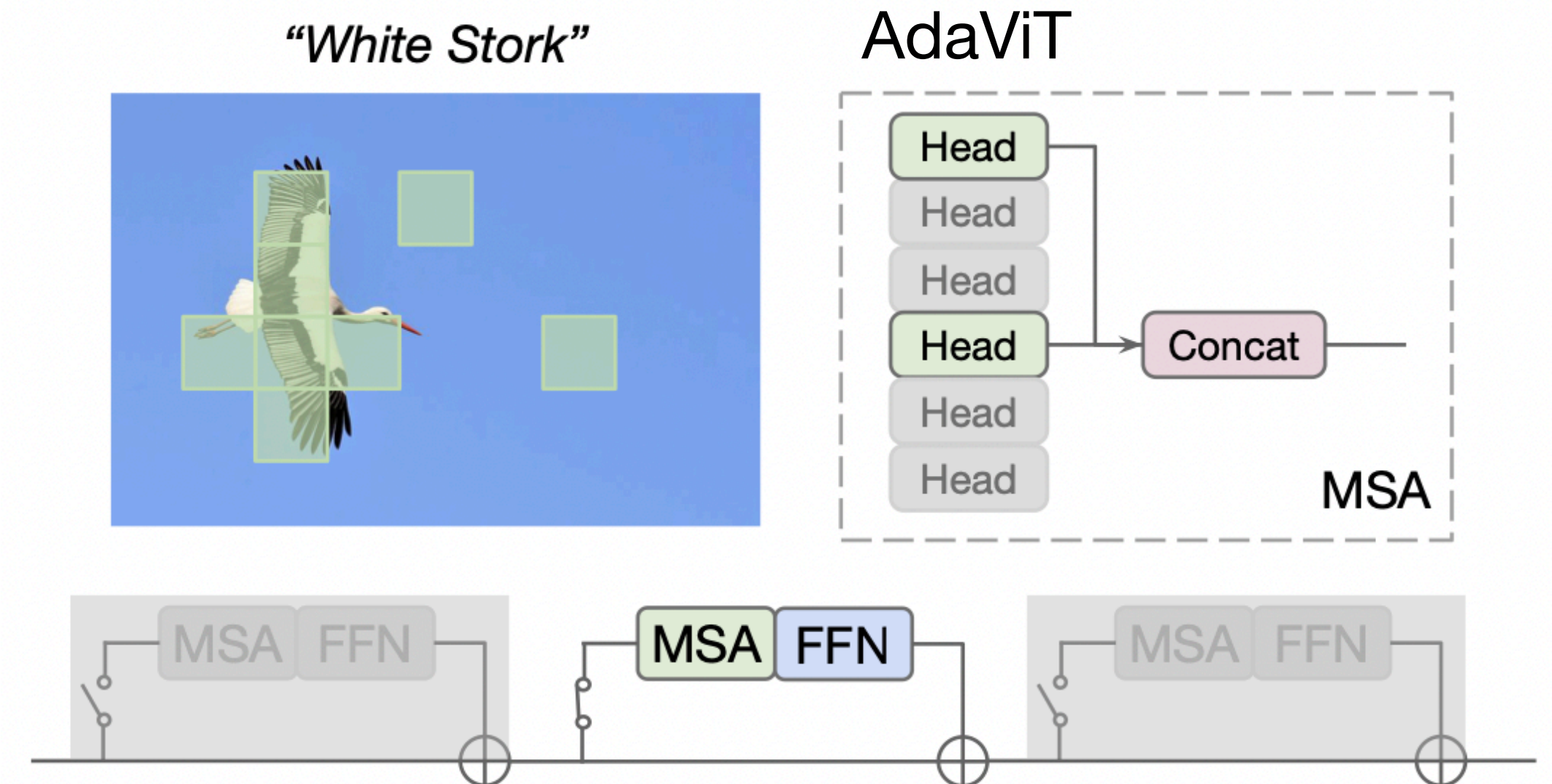


# Adaptive Models on Transformers

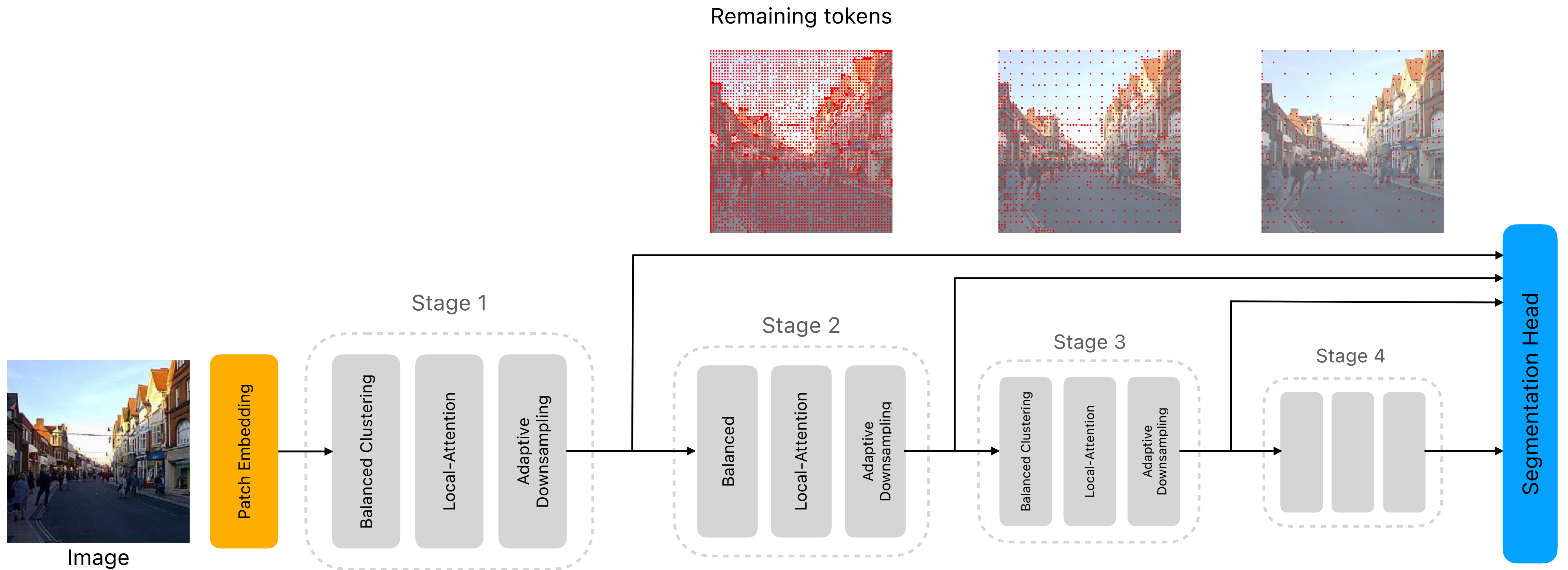
E.g., AdaViT, DynamicViT, A-ViT...

- Adopt global attention (quadratic complexity!)
- No actual downsampling in training (need gradient)

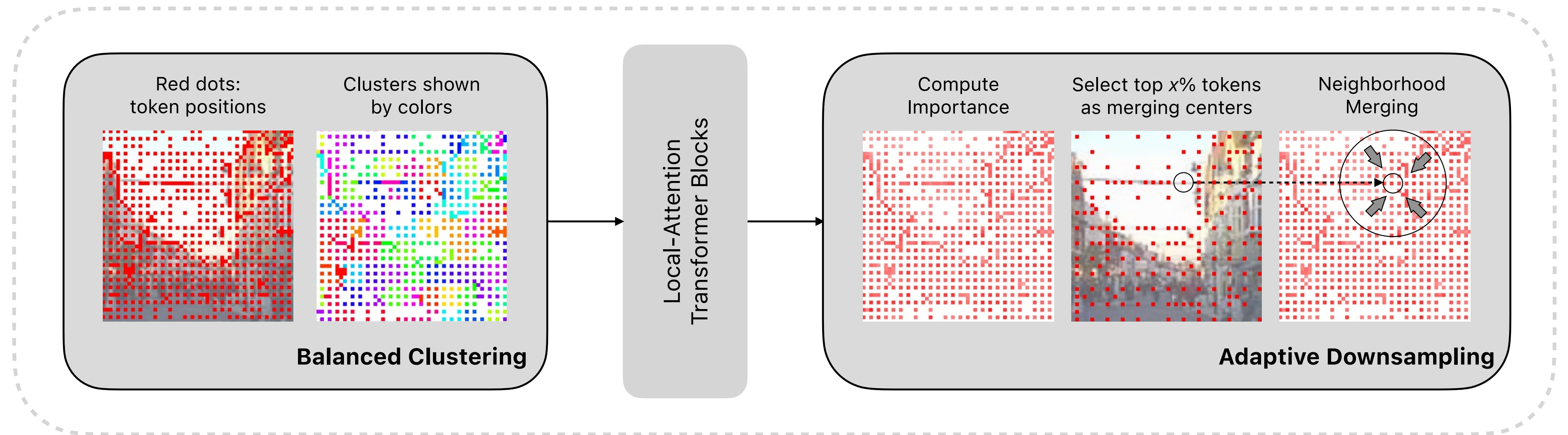
Thus, they cannot scale to high-resolution segmentation tasks!



# AFF Architecture



# AFF Architecture



Three modules that form AFF's stage

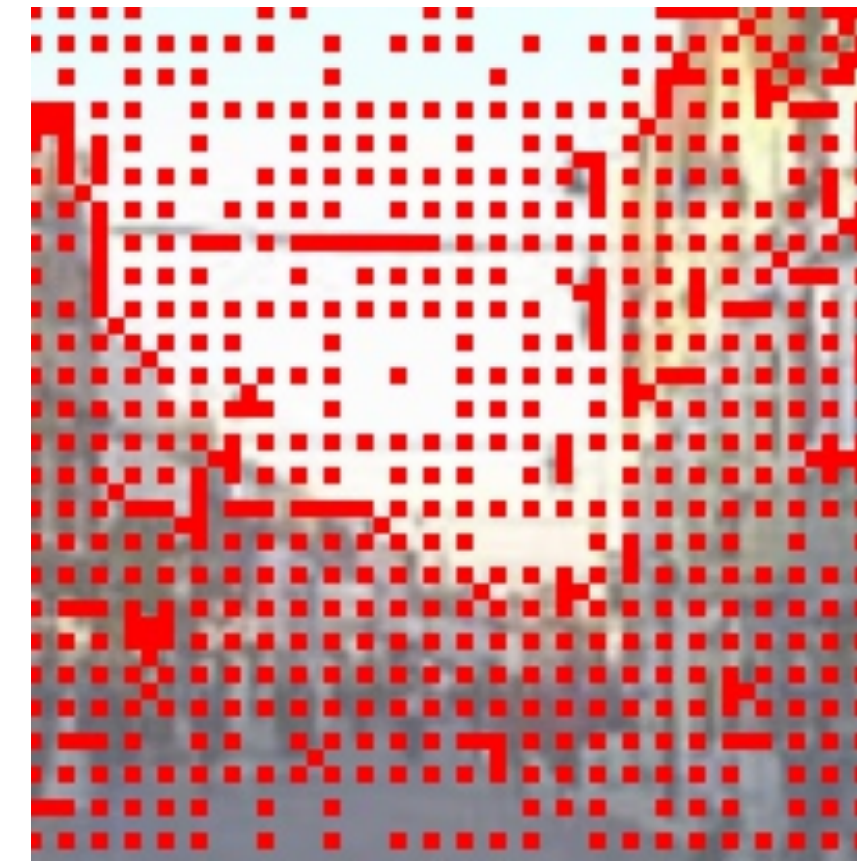


# AFF Clustering

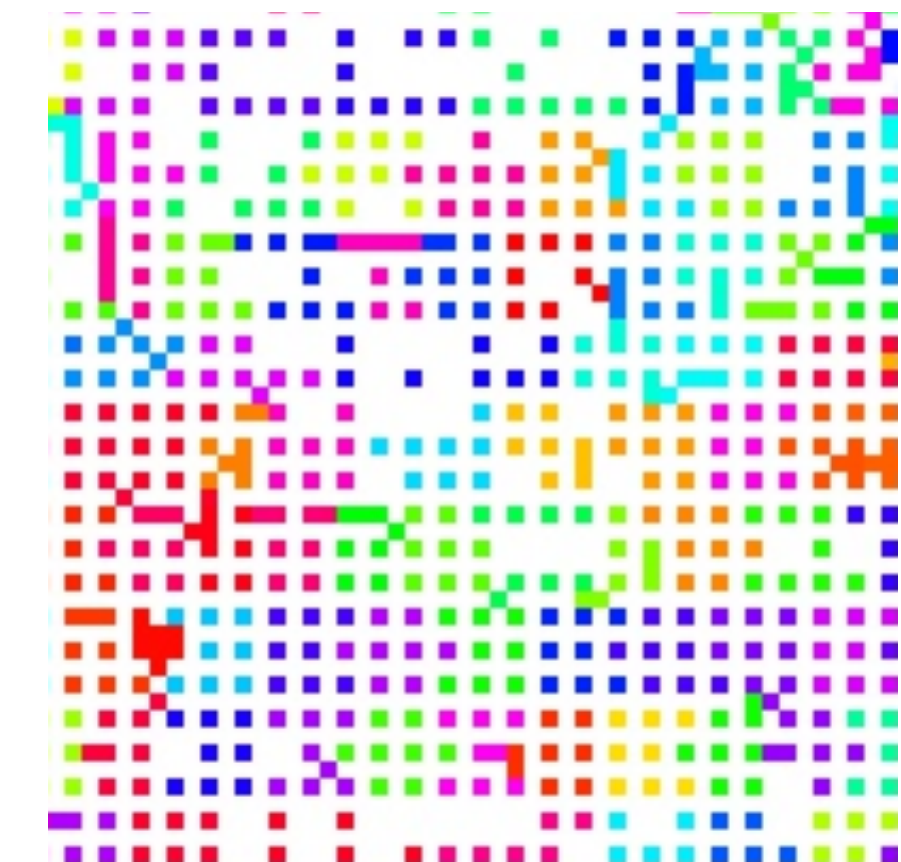
We propose a fast, non-iterative, equal-sized clustering method for 2D tokens

The method is based on space-filling curves

Please refer to our paper for algorithm details!

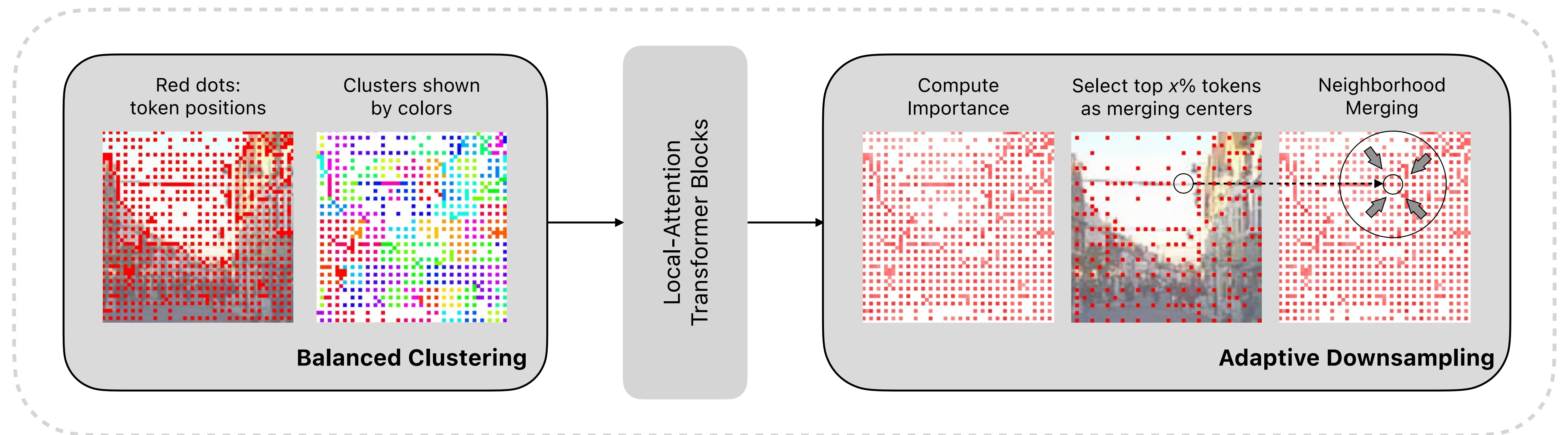


Red dots: tokens



Clusters

# AFF Architecture



Three modules that form AFF's stage

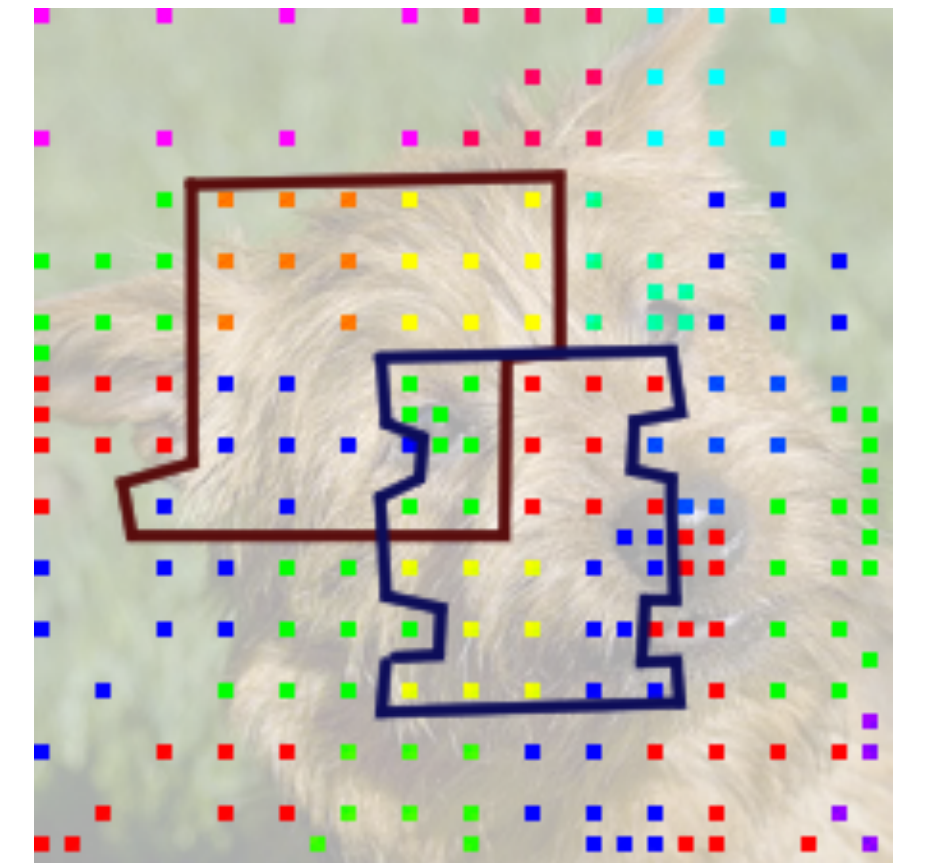
# AFF Local Attention

Using the equal-sized clusters, we define the neighborhood of a token by its nearest  $R$  clusters

Overlapping neighborhoods enable the smooth propagation of information among tokens



On-grid



Off-grid

# AFF Local Attention

On each neighborhood, we apply the standard transformer self-attention:

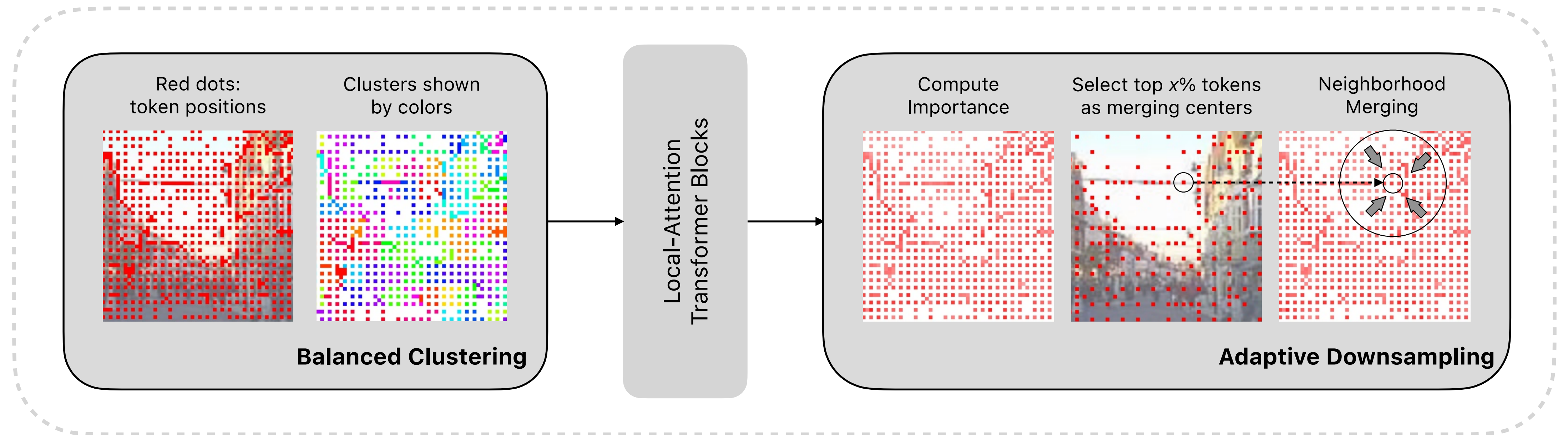
$$A = \text{softmax}(QK^T + P)$$

where  $P_{i,j} = w(p_i - p_j)$  is the **positional embedding**

We further make this embedding aware of potential **rotation/scale invariances** by expanding the relative position vector:

$$\left( \Delta x, \Delta y, \sqrt{\Delta x^2 + \Delta y^2}, \frac{\Delta x}{\sqrt{\Delta x^2 + \Delta y^2}}, \frac{\Delta y}{\sqrt{\Delta x^2 + \Delta y^2}} \right)$$

# AFF Architecture



Three modules that form AFF's stage

# Adaptive Downsampling

First, for each token, we compute an importance score  $s_i$

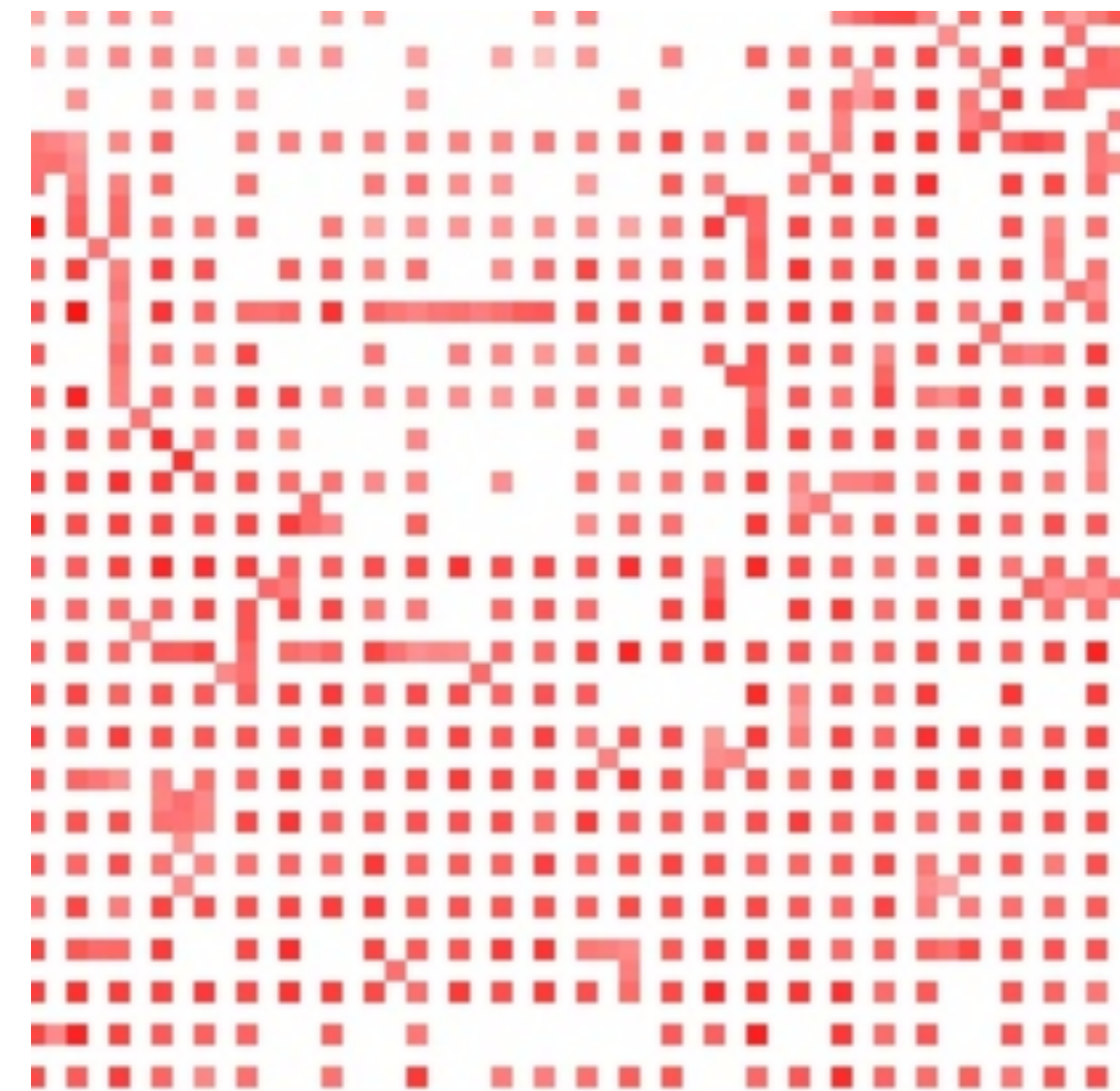
$f_i$ : token feature

$l$ : a fully-connected layer

$$\text{Importance score} = \overbrace{\sigma(l(f_i))}^{s_i} \cdot$$

↑  
Sigmoid

Compute Importance

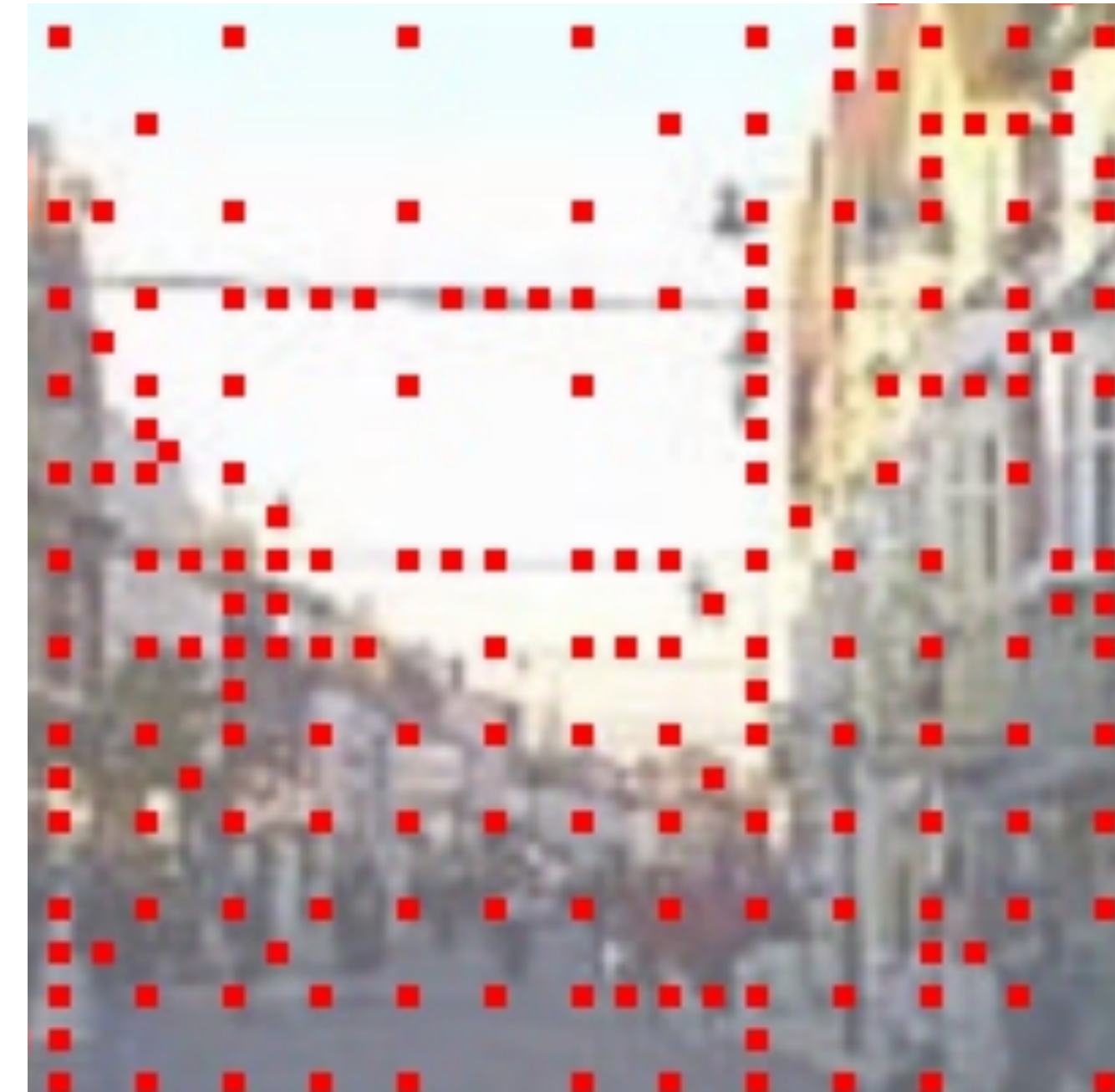


# Merging center selection

Second, we select top  $x\%$  tokens according to the importance scores

$x\%$  is the downsampling rate (we show experiment results with 1/4 and 1/5)

Select top  $x\%$  tokens

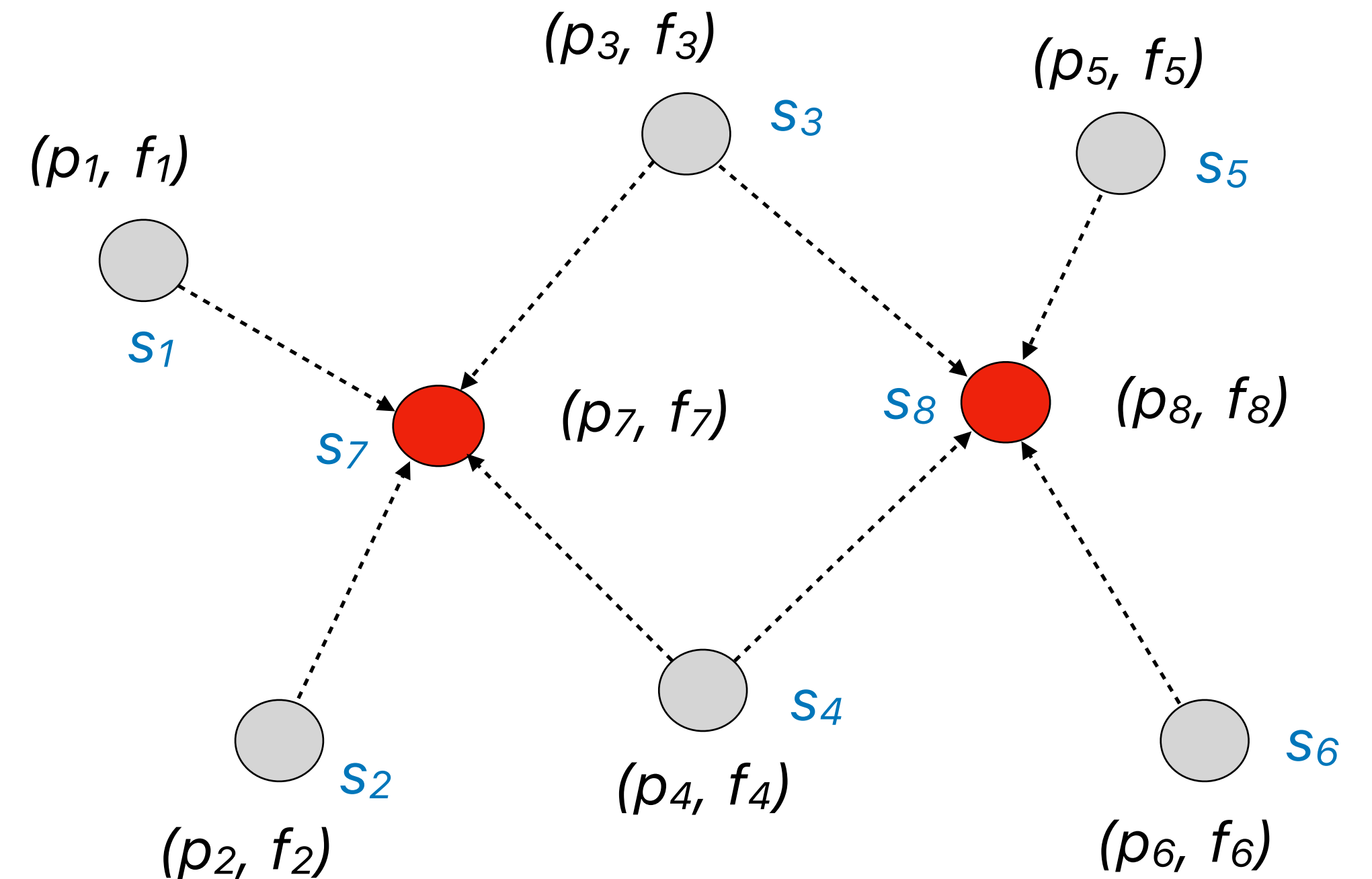


# Neighborhood Merging

Lastly, we merge the neighborhoods of the selected tokens

We use a PointConv layer, modulated by the learnable importance score  $s_i$ 's

The output is  $x\%$  merged tokens





# Experiments

## ImageNet classification

*1/5: downsampling rate*

Model	Top-1 Acc	# Params	FLOPs
Swin-Mini	76.9%	6.76M	1.07G
AFF-Mini	<b>78.2%</b>	6.75M	1.08G
AFF-Mini-1/5	77.5%	6.75M	<b>0.72G</b>
Swin-Tiny	81.3%	28M	4.5G
AFF-Tiny	<b>83%</b>	27M	4G
AFF-Tiny-1/5	82.4%	27M	<b>2.74G</b>
Swin-Small	83%	50M	8.7G
AFF-Small	<b>83.5%</b>	42.6M	8.16G
AFF-Small-1/5	<b>83.4%</b>	42.6M	<b>5.69G</b>

# Experiments

## ADE20K semantic segmentation

Model	mIoU	FLOPs
Swin-Mini	44.1	48.9G
AFF-Mini	<b>46.5</b>	48.3G
AFF-Mini-1/5	46.0	<b>39.9G</b>
Swin-Tiny	47.7	74G
AFF-Tiny	<b>50.2</b>	64.6G
AFF-Tiny-1/5	<b>50.0</b>	<b>51.1G</b>
Swin-Small	51.3	98G
AFF-Small	51.2	87G
AFF-Small-1/5	<b>51.9</b>	<b>67.2G</b>

*Segmentation head: Mask2Former*

# Experiments

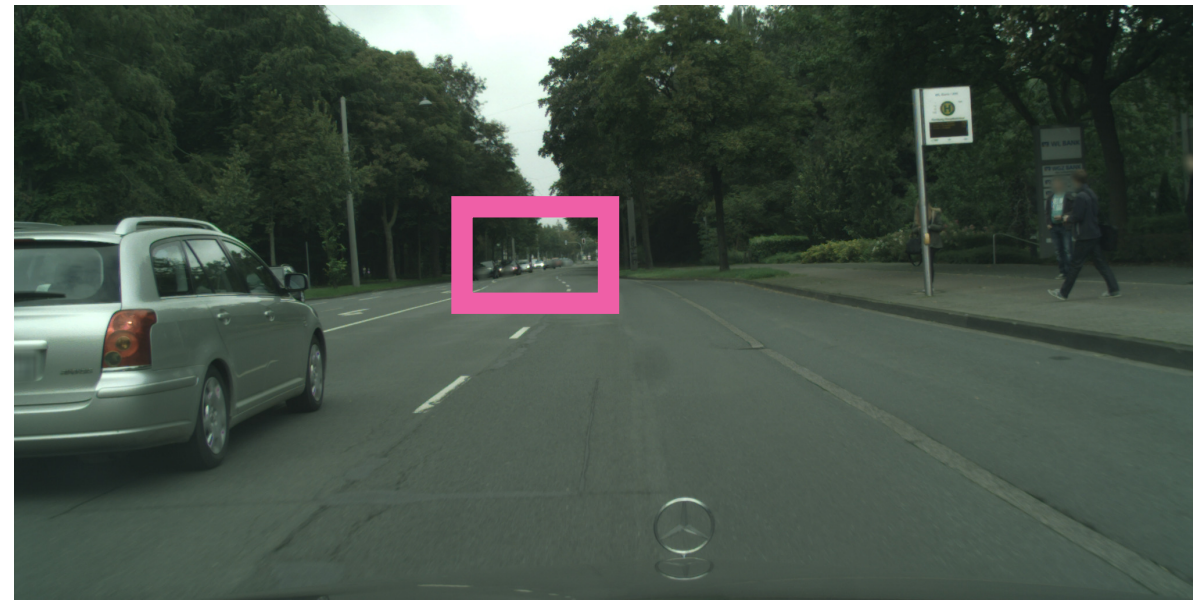
## Cityscapes instance & panoptic segmentation

Model	Instance AP	Panoptic PQ (s.s.)	Backbone # Params	
AFF-Mini	40.0	62.7	6.75M	
Swin-Tiny	39.7	63.9	28M	
AFF-Tiny	<b>42.7</b>	<b>65.7</b>	27M	
Swin-Small	41.8	64.8	50M	} <b>3.3x</b>
AFF-Small	<b>44.0</b>	<b>66.9</b>	42.6M	
Swin-Base	42	66.1	88M	} <b>4.6x</b>
Swin-Large	43.7	66.6	197M	

*Segmentation head: Mask2Former*

# Qualitative results

Image



Zoomed-in image



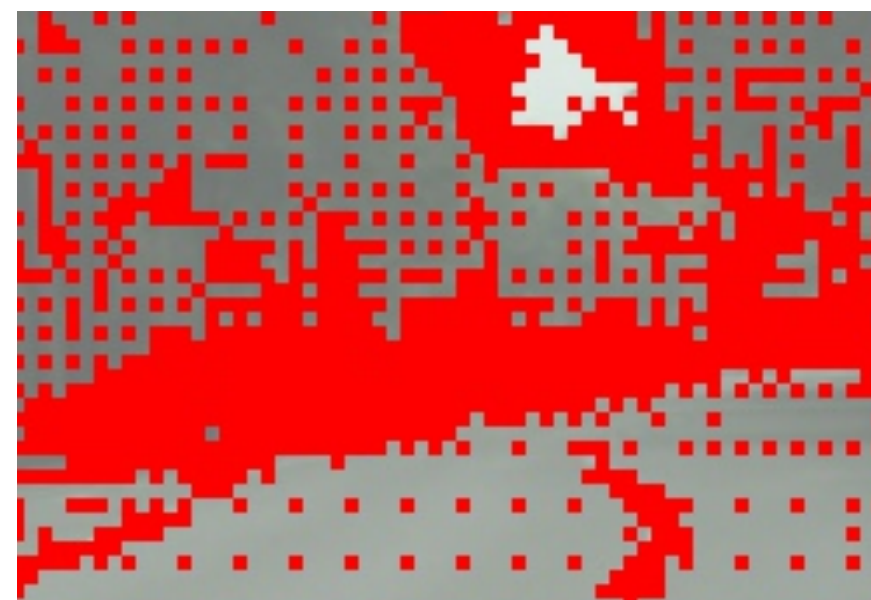
Swin's prediction



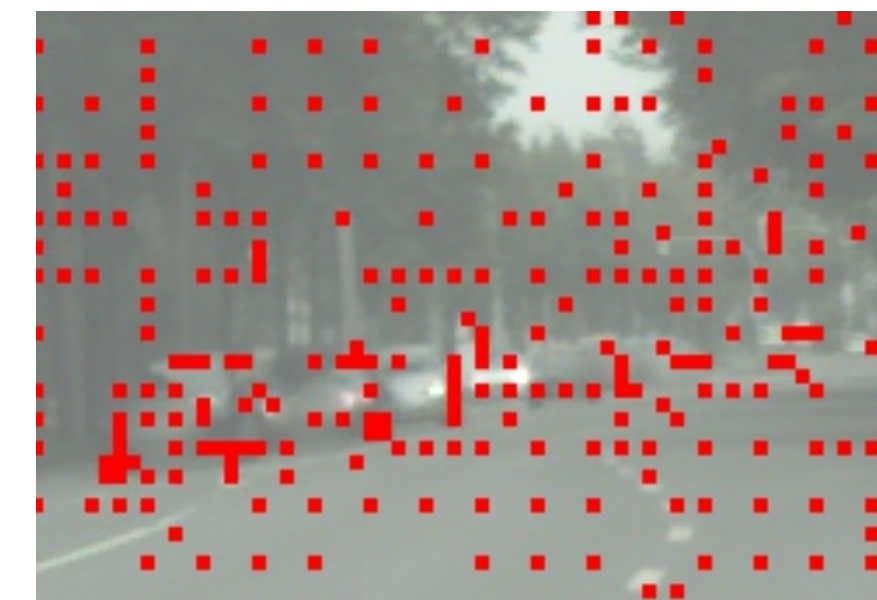
AFF's prediction



Stage 2



Stage 3



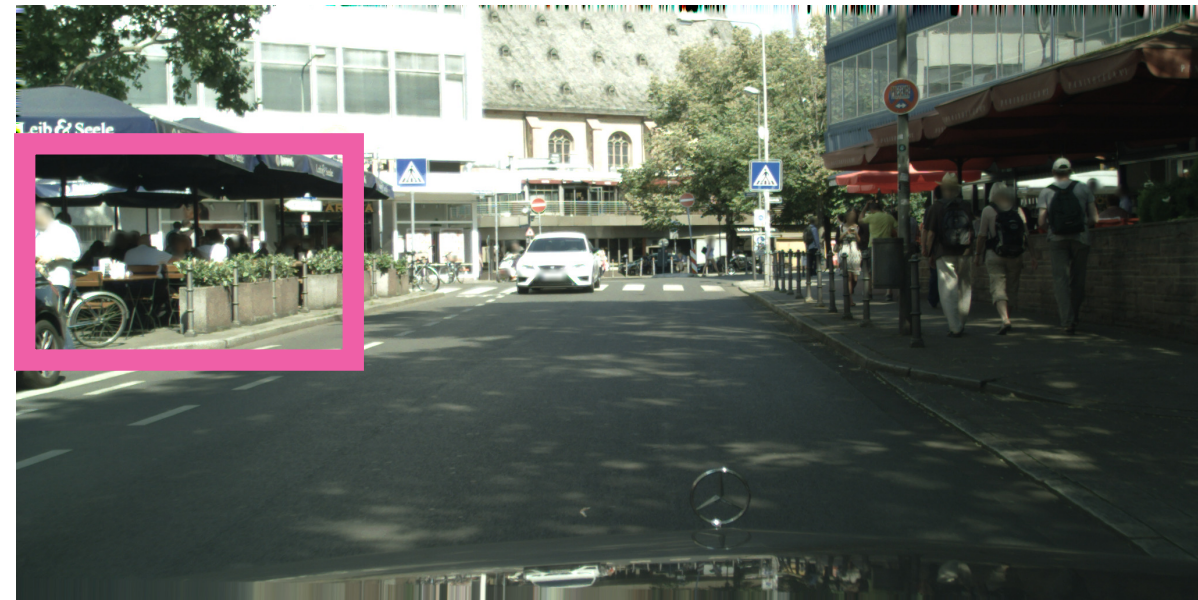
Stage 4



AFF's remaining tokens

# Qualitative results

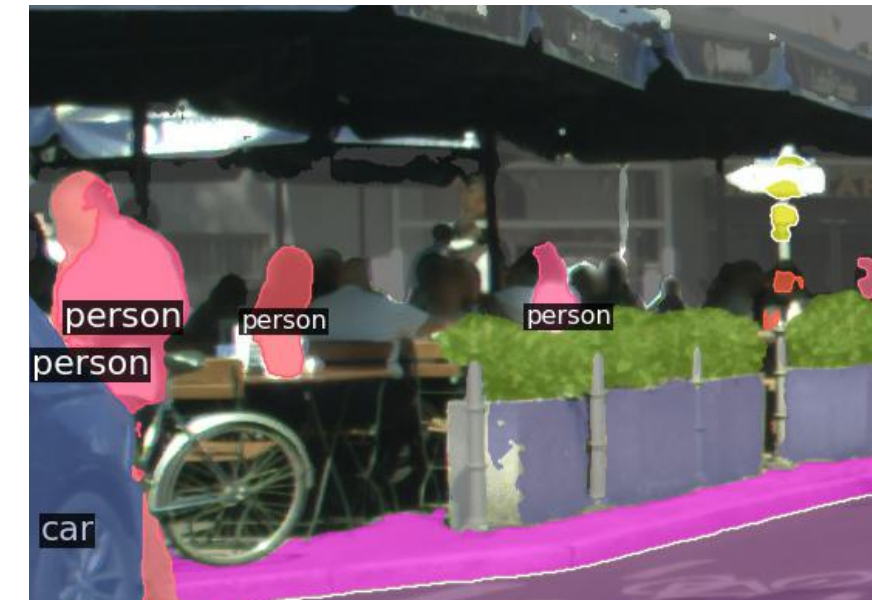
Image



Zoomed-in image



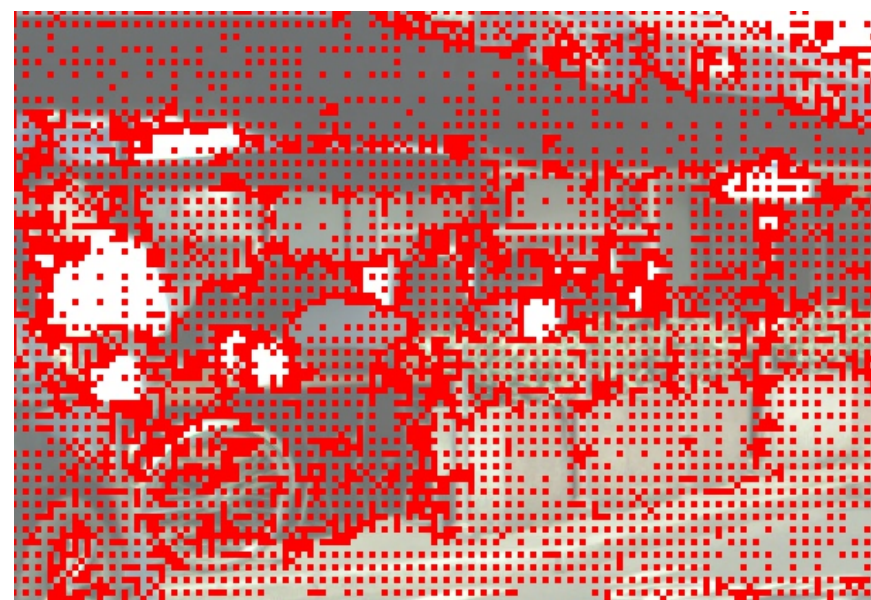
Swin's prediction



AFF's prediction



Stage 2



Stage 3



Stage 4



AFF's remaining tokens

# Conclusions

- We introduce the first adaptive-downsampling network capable of dense prediction tasks such as semantic/instance segmentation
- Flexible downsampling rate (e.g., 1/5 vs. traditional 1/4)
- Significant savings on FLOPs and significant improvement on recognition of small objects



Image



SWIN prediction



AFF prediction

