



# Foundations of Edge AI

# Lecture 18

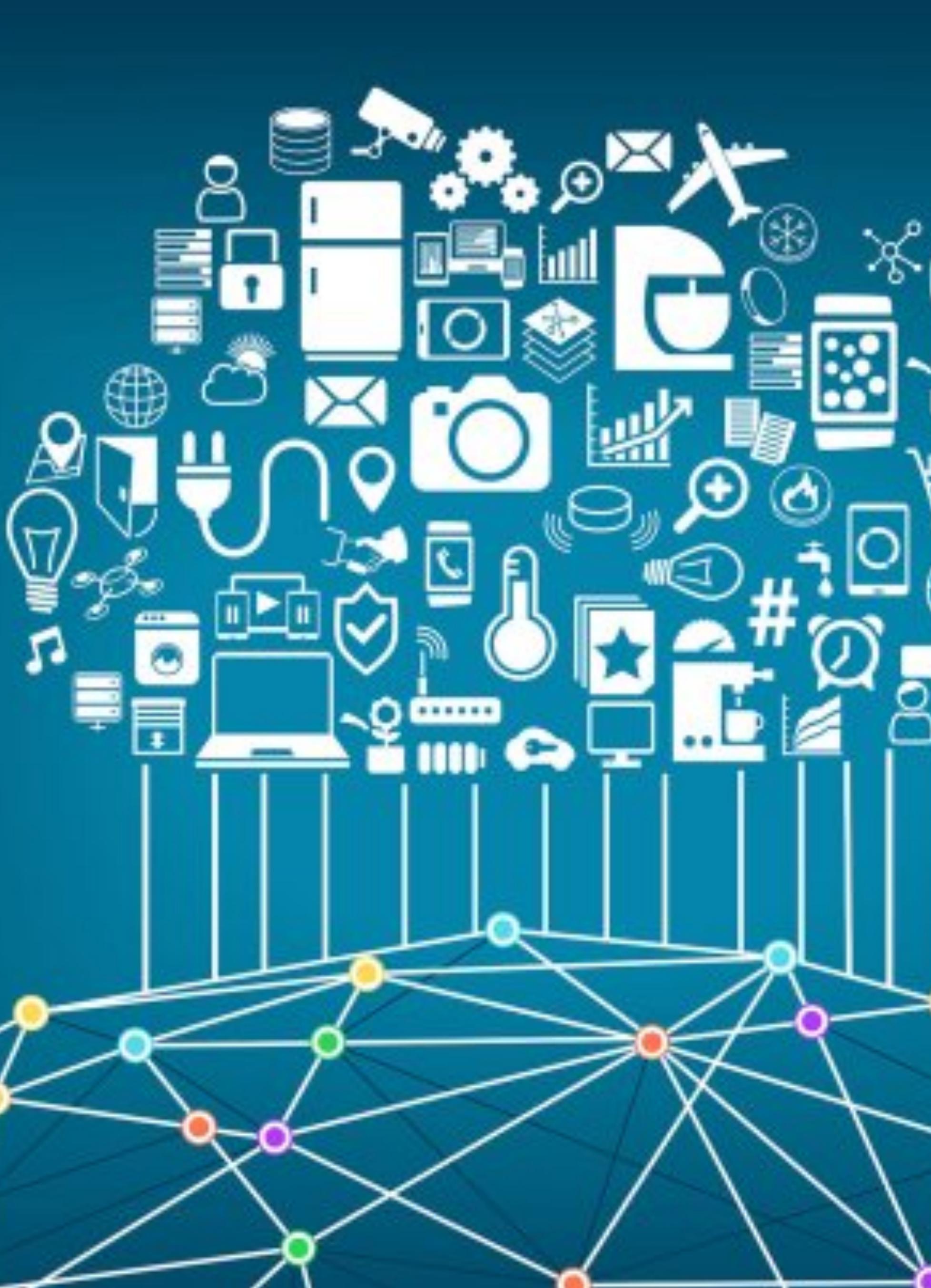
# Vision Transformer

# Lanyu (Lori) Xu

Email: lxu@oakland.edu

Homepage: <https://lori930.github.io/>

Office: EC 524

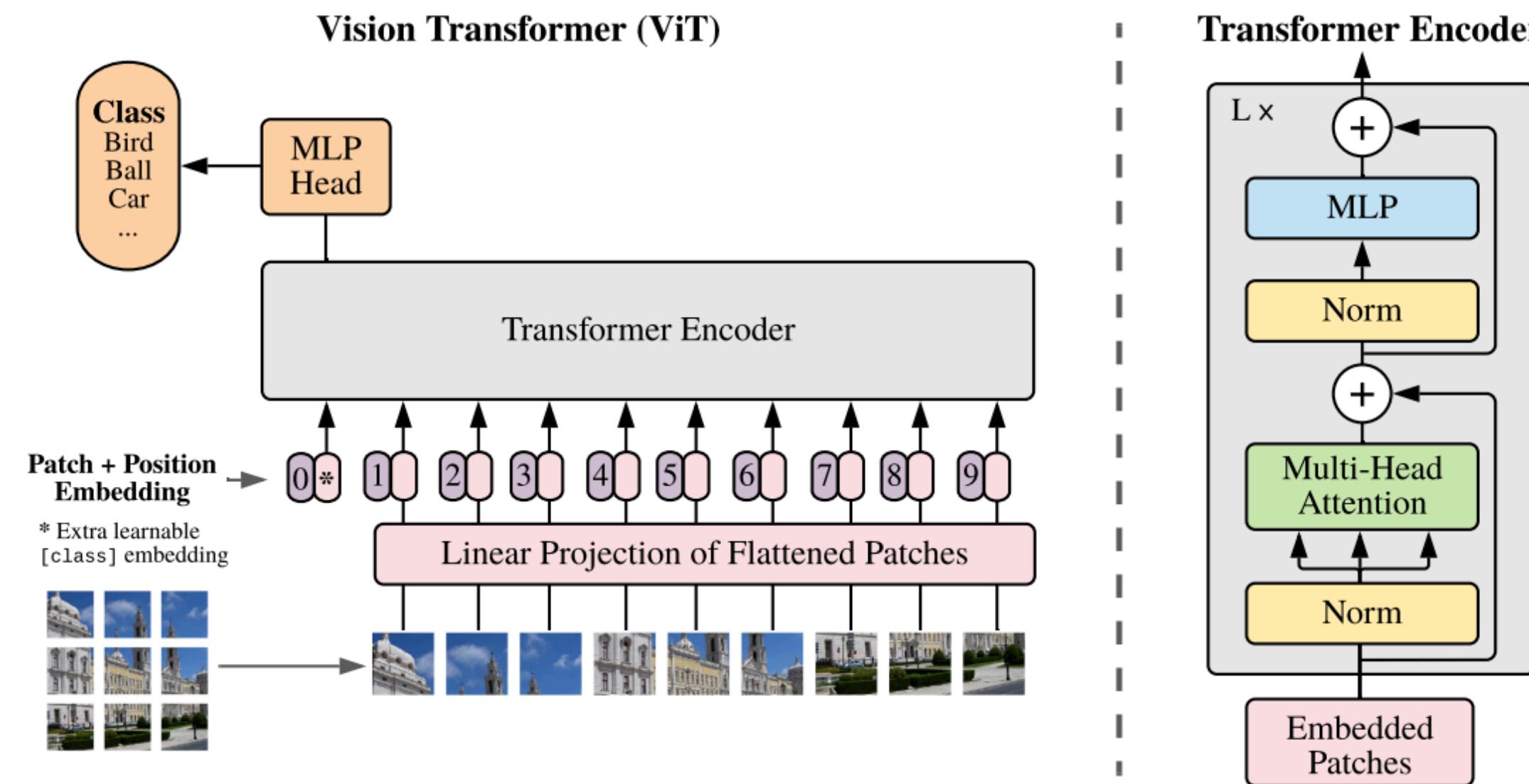


# Lecture Plan

Today we will...

- Introduce the basics of Vision Transformer (ViT)
- Introduce efficient ViT and acceleration techniques
  - Window attention: SwinTransformer
  - Linear attention: EfficientViT
  - Sparse attention: SparseViT
- Self-supervised learning for ViT
  - Contrastive learning
  - Masked image modeling

# Basics of Vision Transformer (ViT)



# Visio Transformer

Convert 2D Images to a Sequence of Patches



- Convert the 2D image to a sequence of patches

# Visio Transformer

Convert 2D Images to a Sequence of Patches

Each patch is a token



Image size: 96x96

Patch size: 32x32

Number of tokens:  $3 \times 3 = 9$

Dimension of each token:  $3 \times 32 \times 32 = 3072$

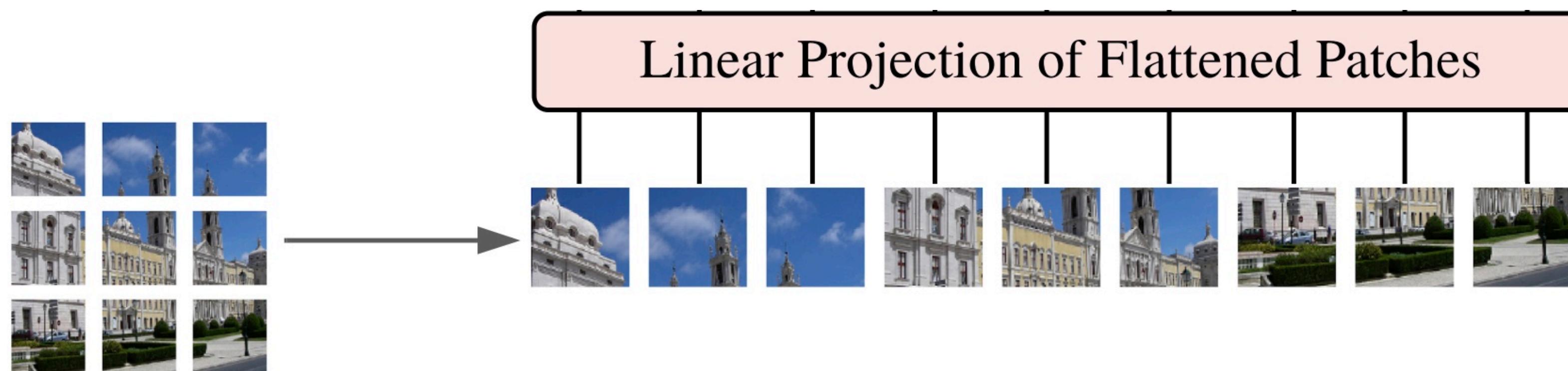
- Convert the 2D image to a sequence of patches

# Visio Transformer

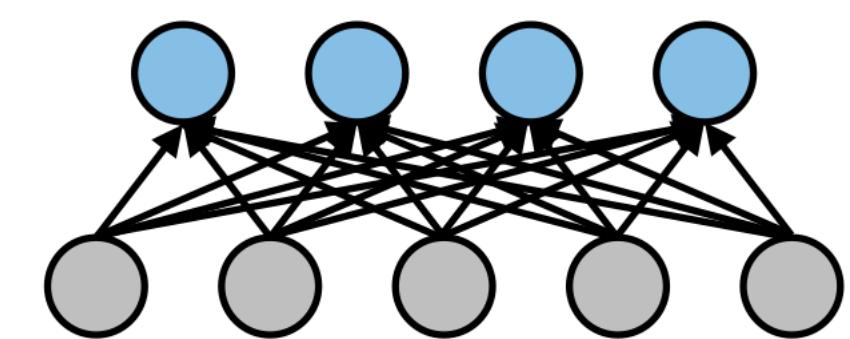
Convert 2D Images to a Sequence of Patches

Each patch is a token

Output dim = hidden size of ViT = 768



- Convert the 2D image to a sequence of patches



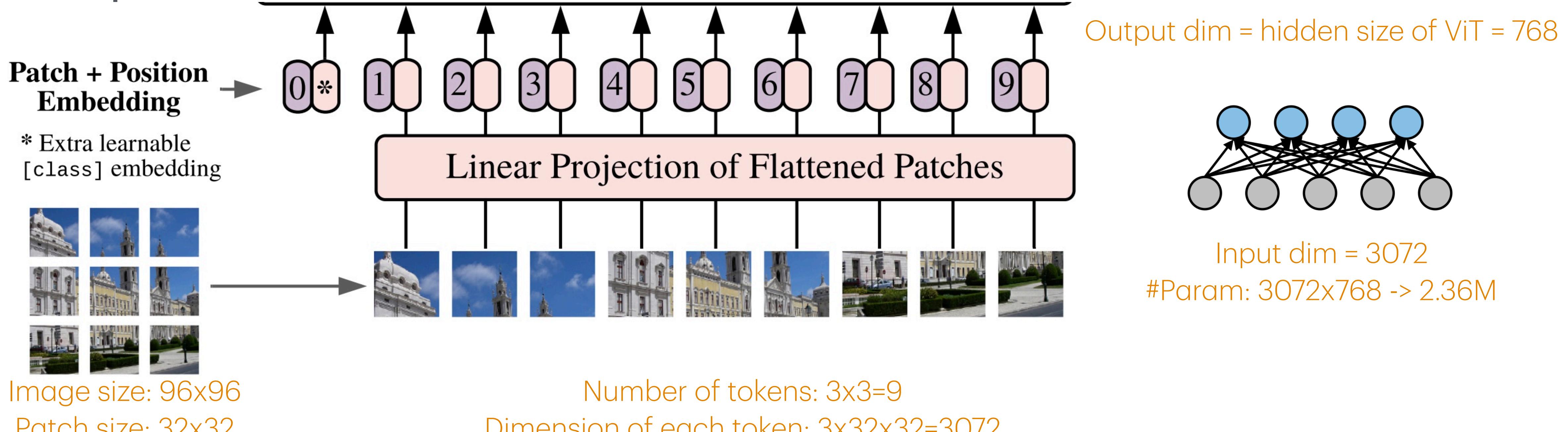
Input dim = 3072

#Param:  $3072 \times 768 \rightarrow 2.36M$

# Visio Transformer

Convert 2D Images to a Sequence of Patches

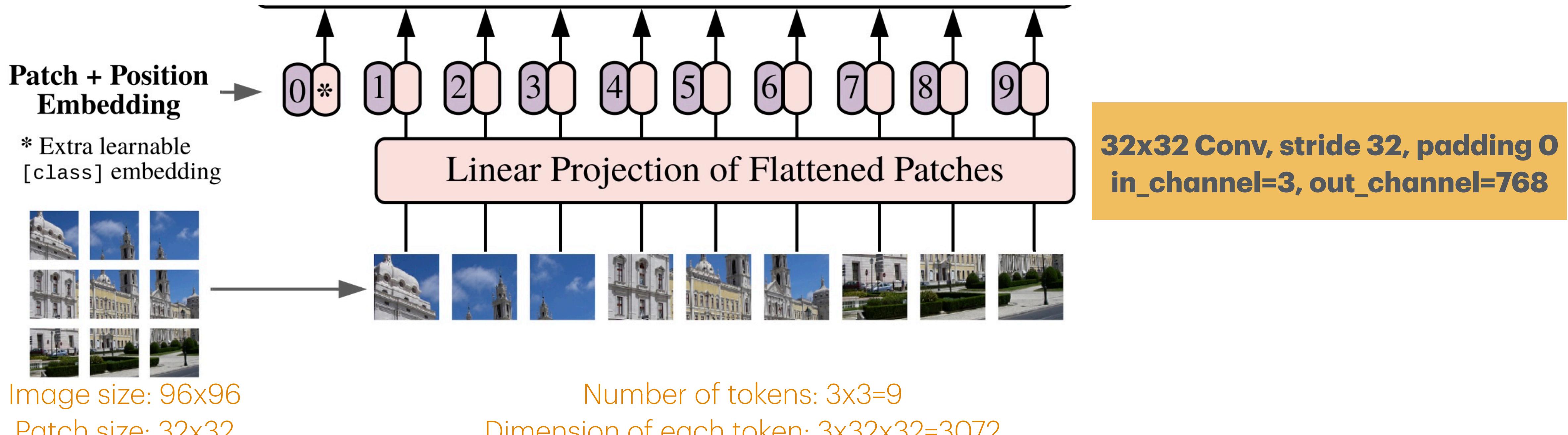
Each patch is a token



- Convert the 2D image to a sequence of patches

# Visio Transformer

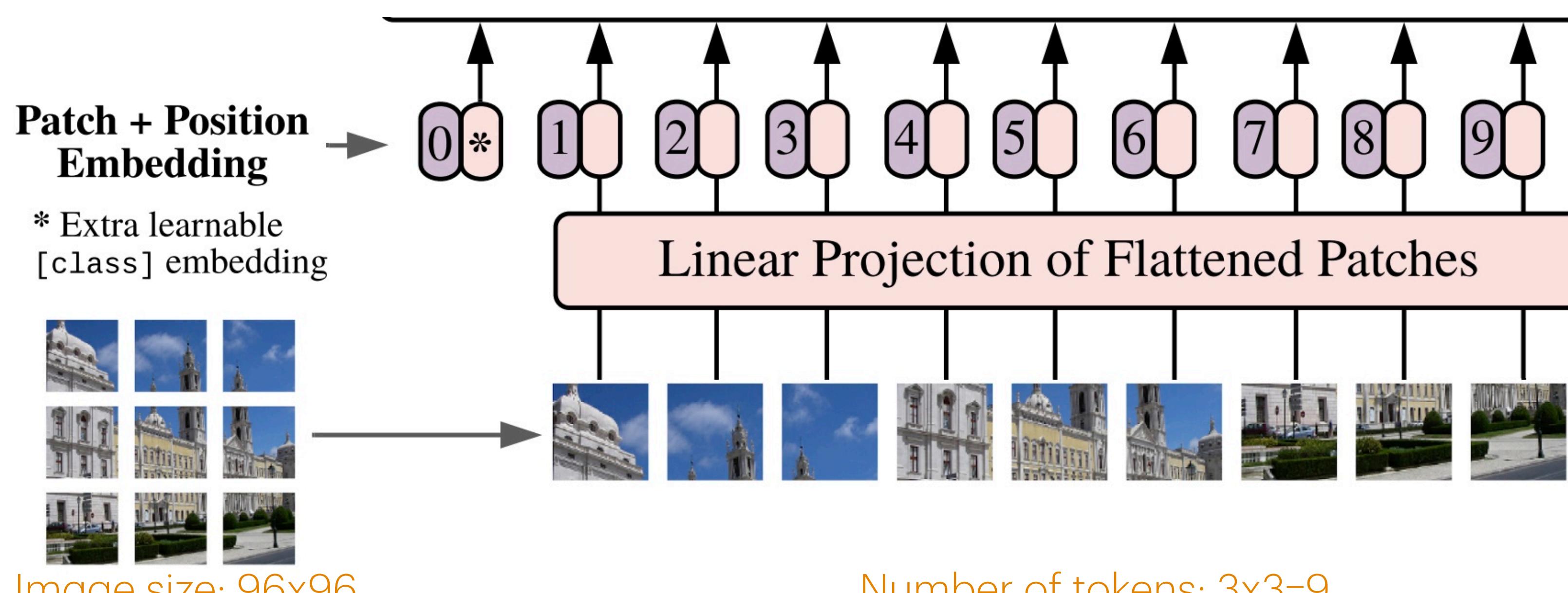
## Practical Implementation



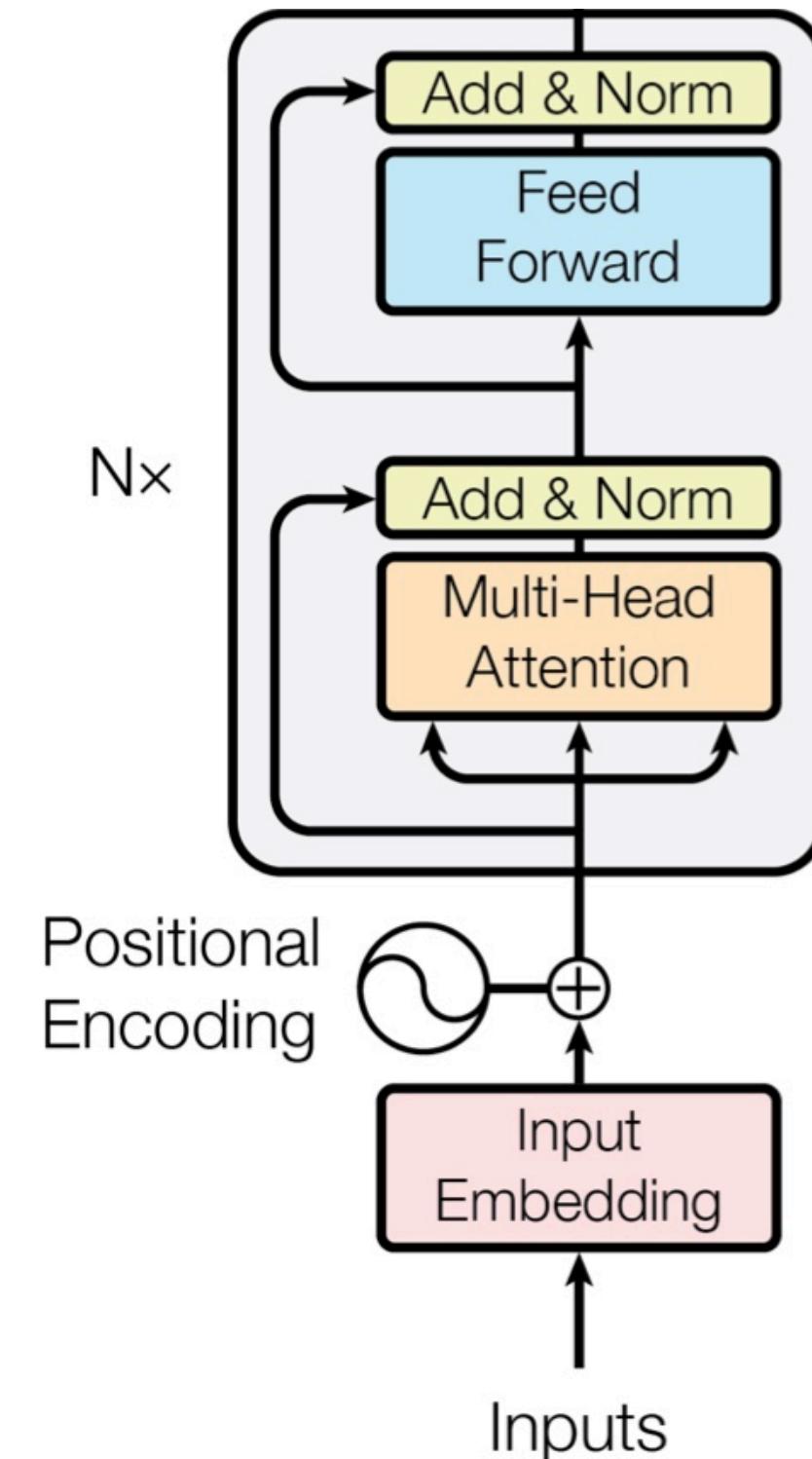
Dosovitskiy, A. (2020). An image is worth 16x16 words: Transformers for image recognition at scale. arXiv preprint arXiv:2010.11929.

# Visio Transformer

Apply the Standard Transformer Encoder



- Convert the 2D image to a sequence of patches
- Feed patch embeddings to the standard transformer encoder



# Visio Transformer

## Model Variants

Size of ResNet50: 25M

Model	Layers	Hidden size $D$	MLP size	Heads	Params
ViT-Base	12	768	3072	12	86M
ViT-Large	24	1024	4096	16	307M
ViT-Huge	32	1280	5120	16	632M

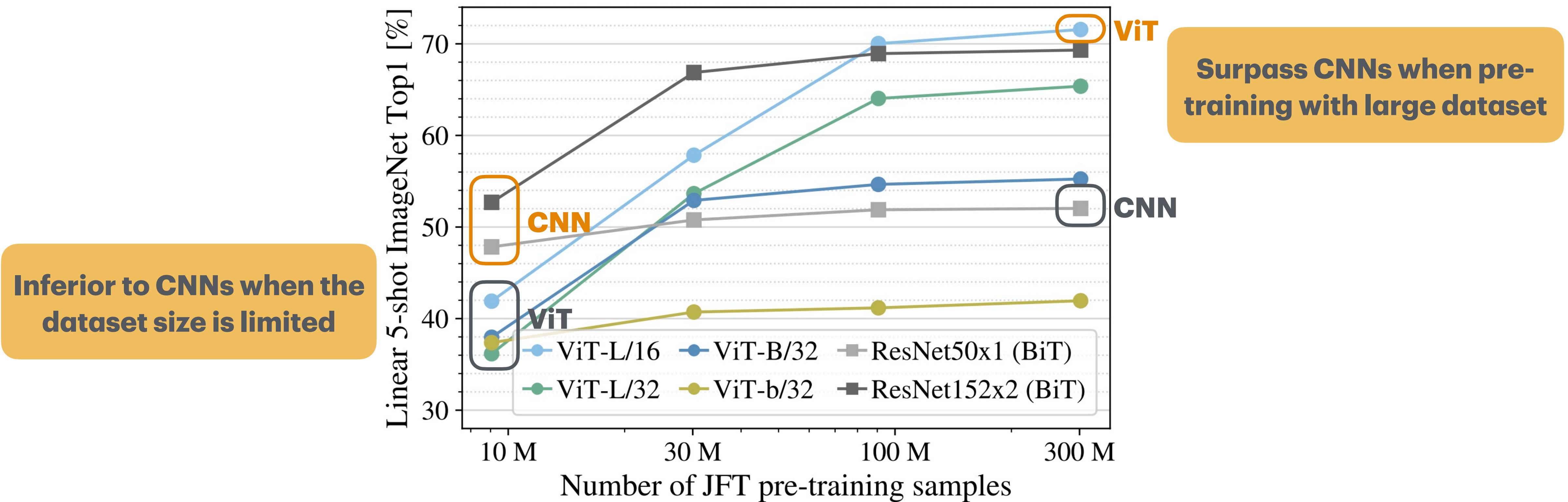


**Patch size:**

**2, 4, 8, 16, 32, ...**

- Common annotation: ViT-L/16: ViT-Large with patch size 16x16

# Image Classification Results



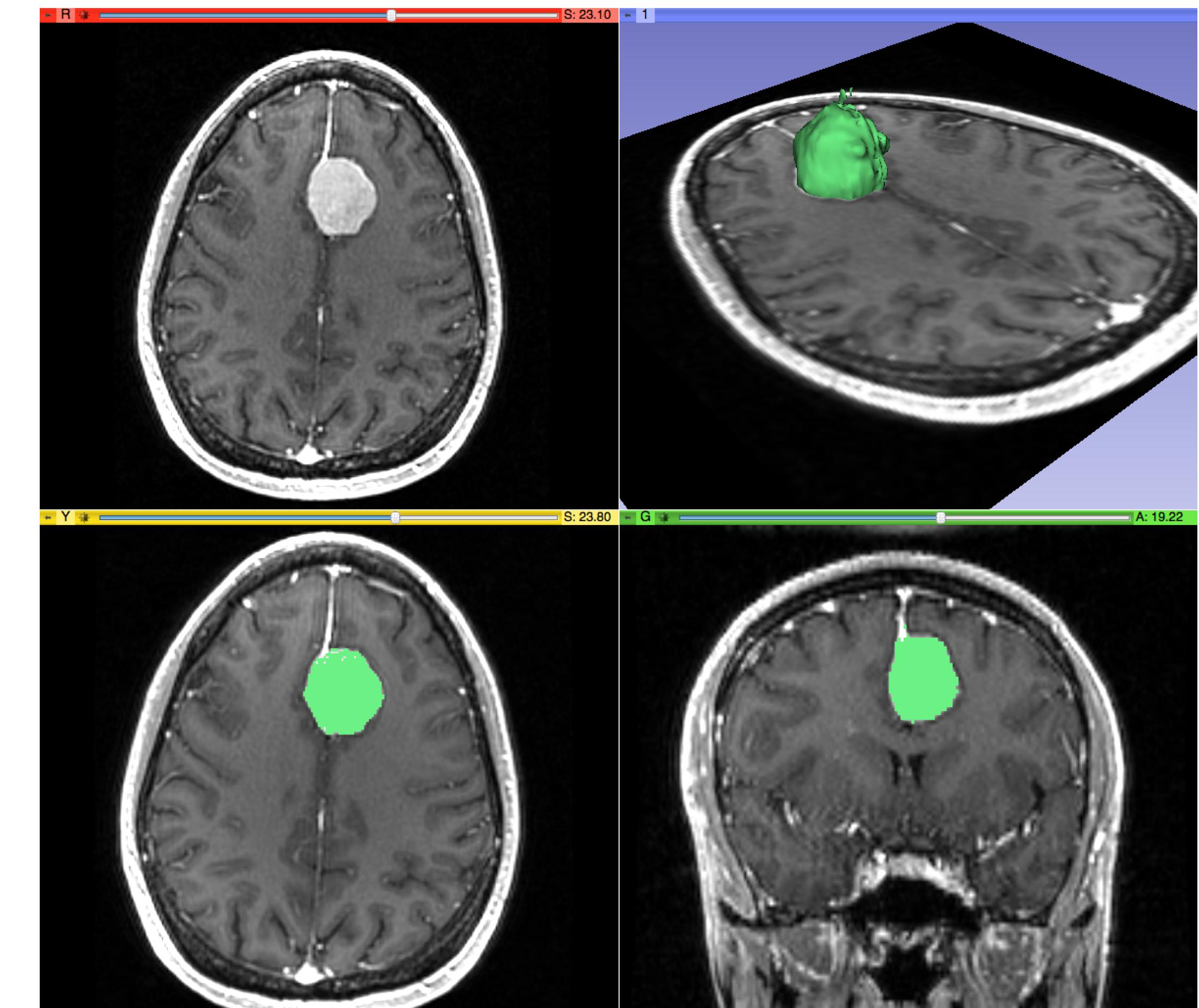
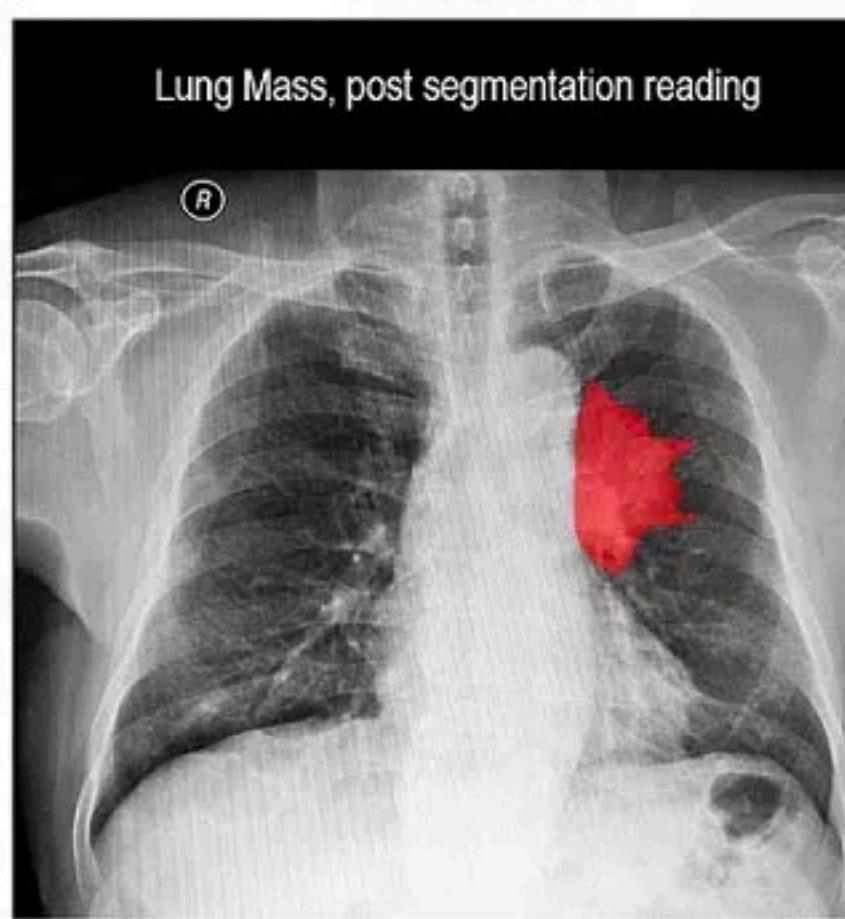
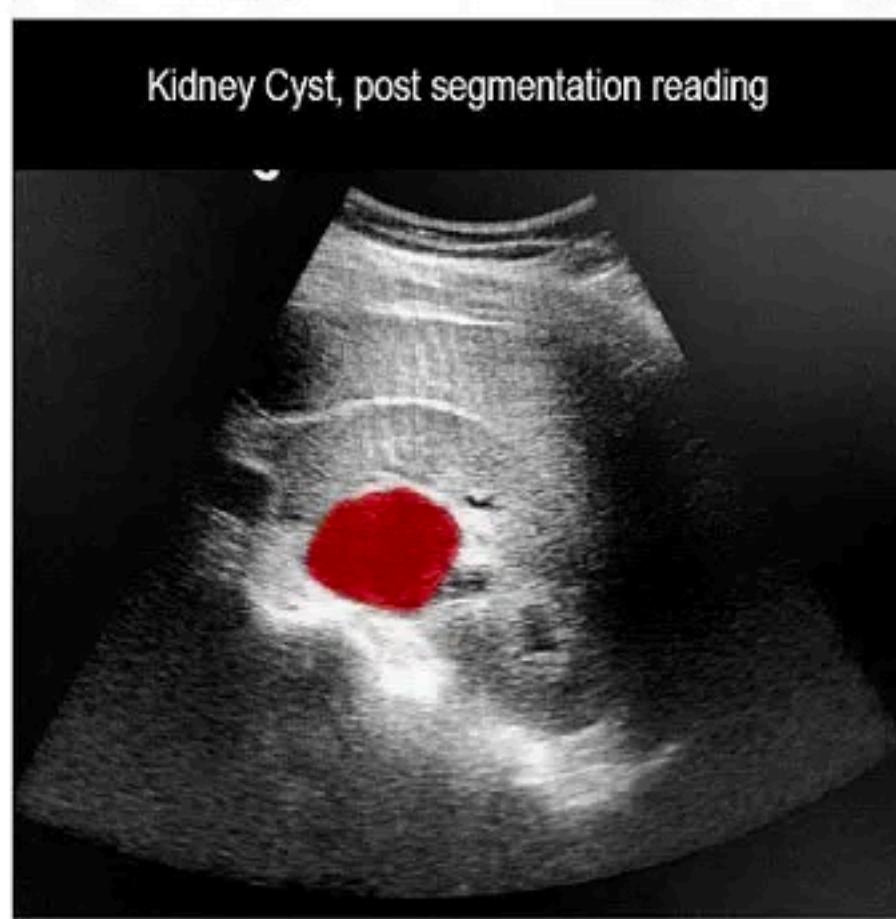
Inferior to CNNs when the dataset size is limited

Surpass CNNs when pre-training with large dataset

High-resolution images are a challenge  
for ViT

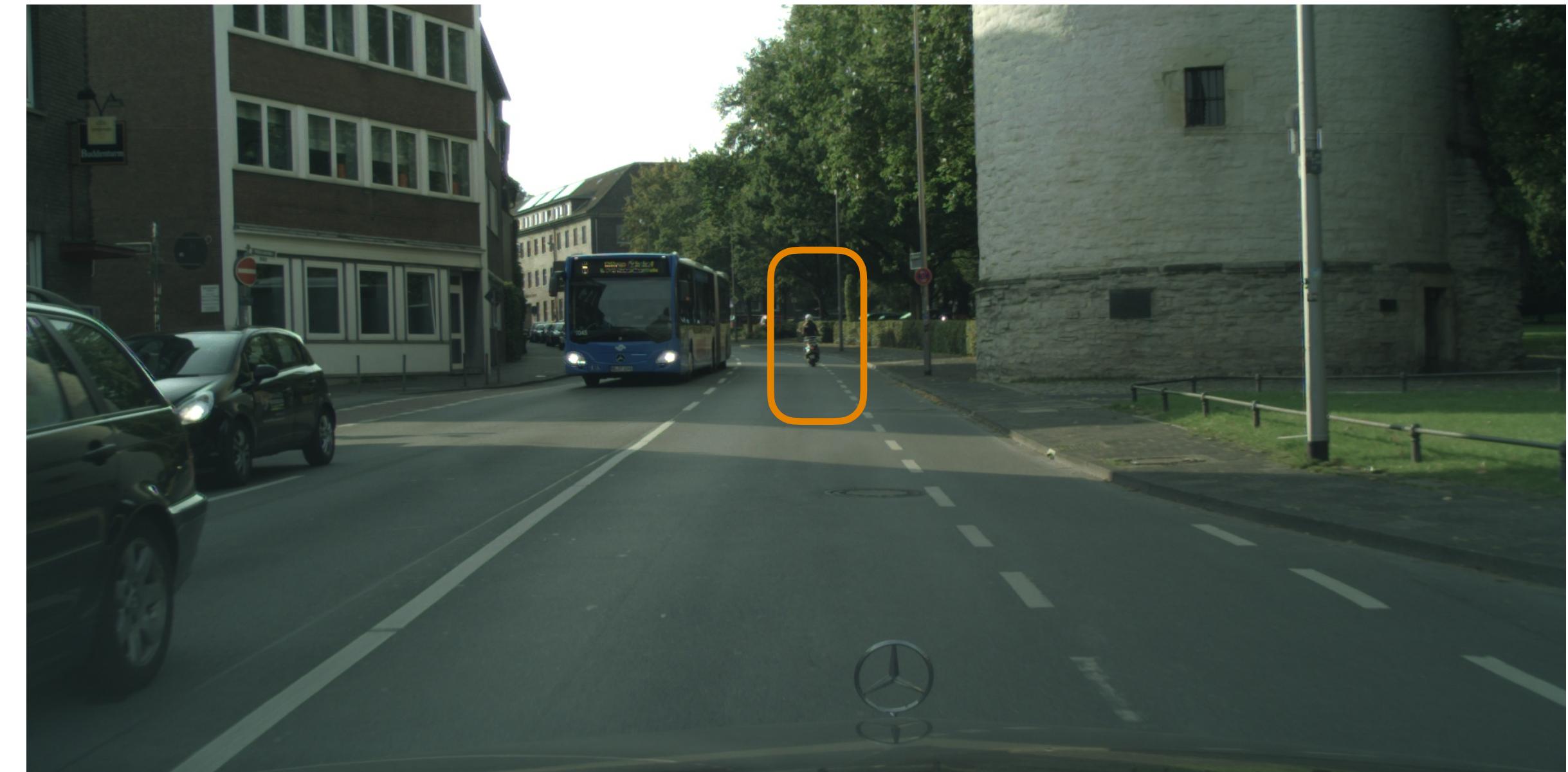
# High-Resolution Dense Prediction is Important

## Medical Imaging



# High-Resolution Dense Prediction is Important

## Autonomous Driving



# Challenges of High-Resolution Images in ViT

- **Increased sequence length:** quadratically with the resolution

- Sequence length:  $N = \frac{H}{P} \times \frac{W}{P}$ , where  $P$  is the patch size. Assume  $P = 16$

- $H = W = 224 \rightarrow N = 14 \times 14 = 196$

- $H = W = 448 \rightarrow N = 14 \times 14 = 784$

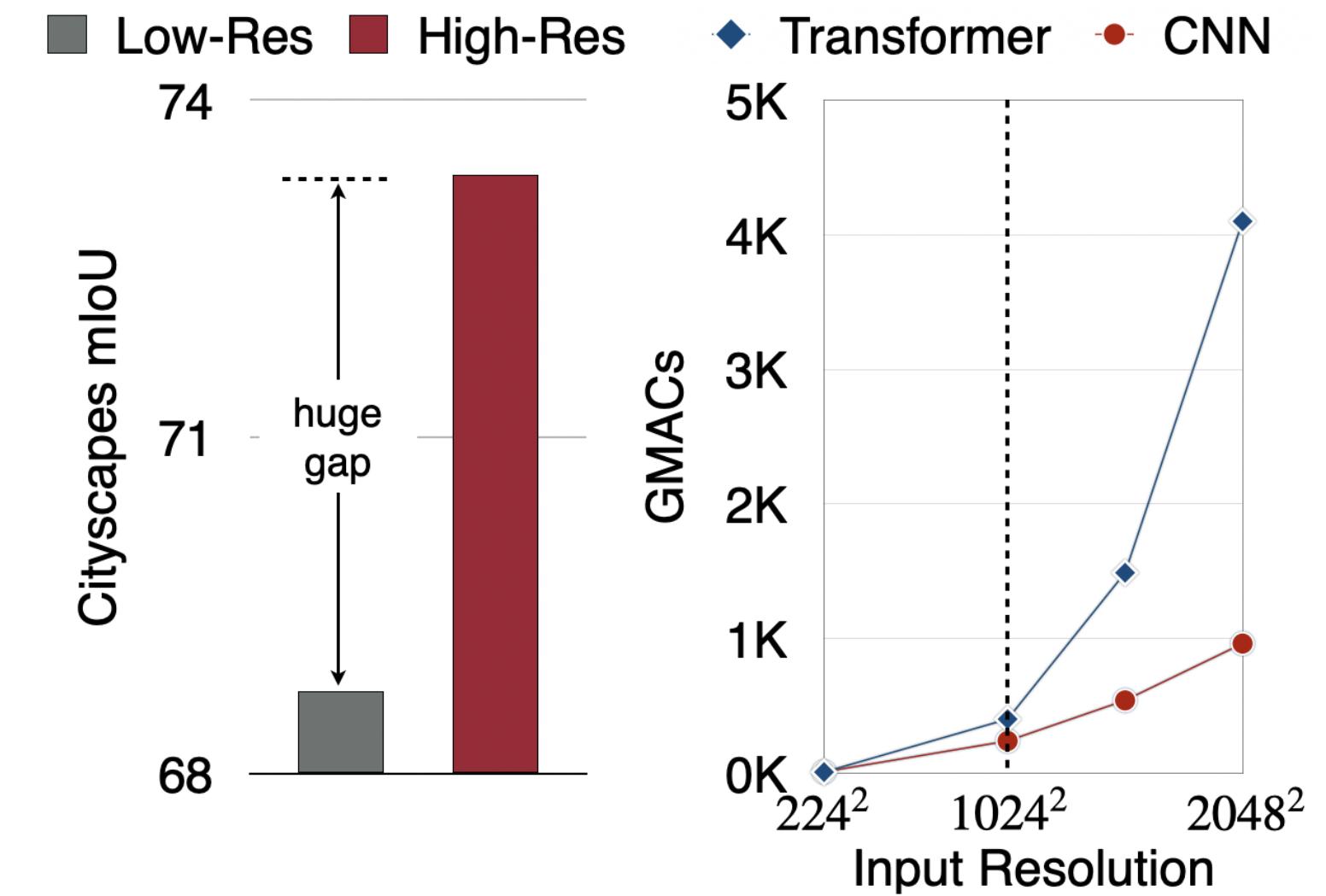
- **Growing memory requirements**

- The attention matrix:  $\text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)$ , with shape  $N \times N$

- Larger intermediate matrix during the self-attention computation

- **Quadratic complexity in self-attention:** 16 times more expensive

- Computational complexity of self-attention mechanism:  $O(N^2 \cdot D)$ , where  $D$  is the embedding size (#channels per token)
- Doubling the resolution → Quadrupling the  $N \rightarrow 16$  times more expensive in self-attention computation

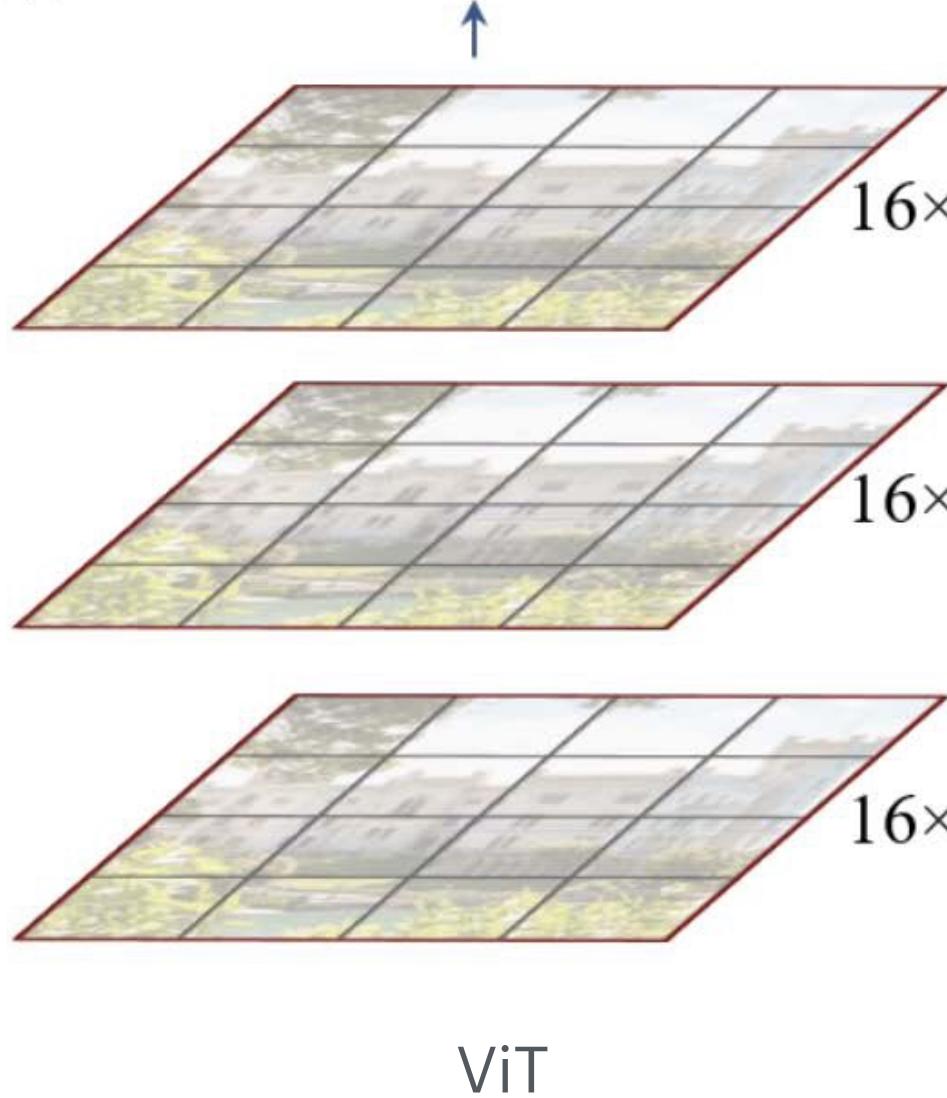


Layer Type	Complexity per Layer
Self-Attention	$O(n^2 \cdot d)$
Recurrent	$O(n \cdot d^2)$
Convolutional	$O(k \cdot n \cdot d^2)$
Self-Attention (restricted)	$O(r \cdot n \cdot d)$

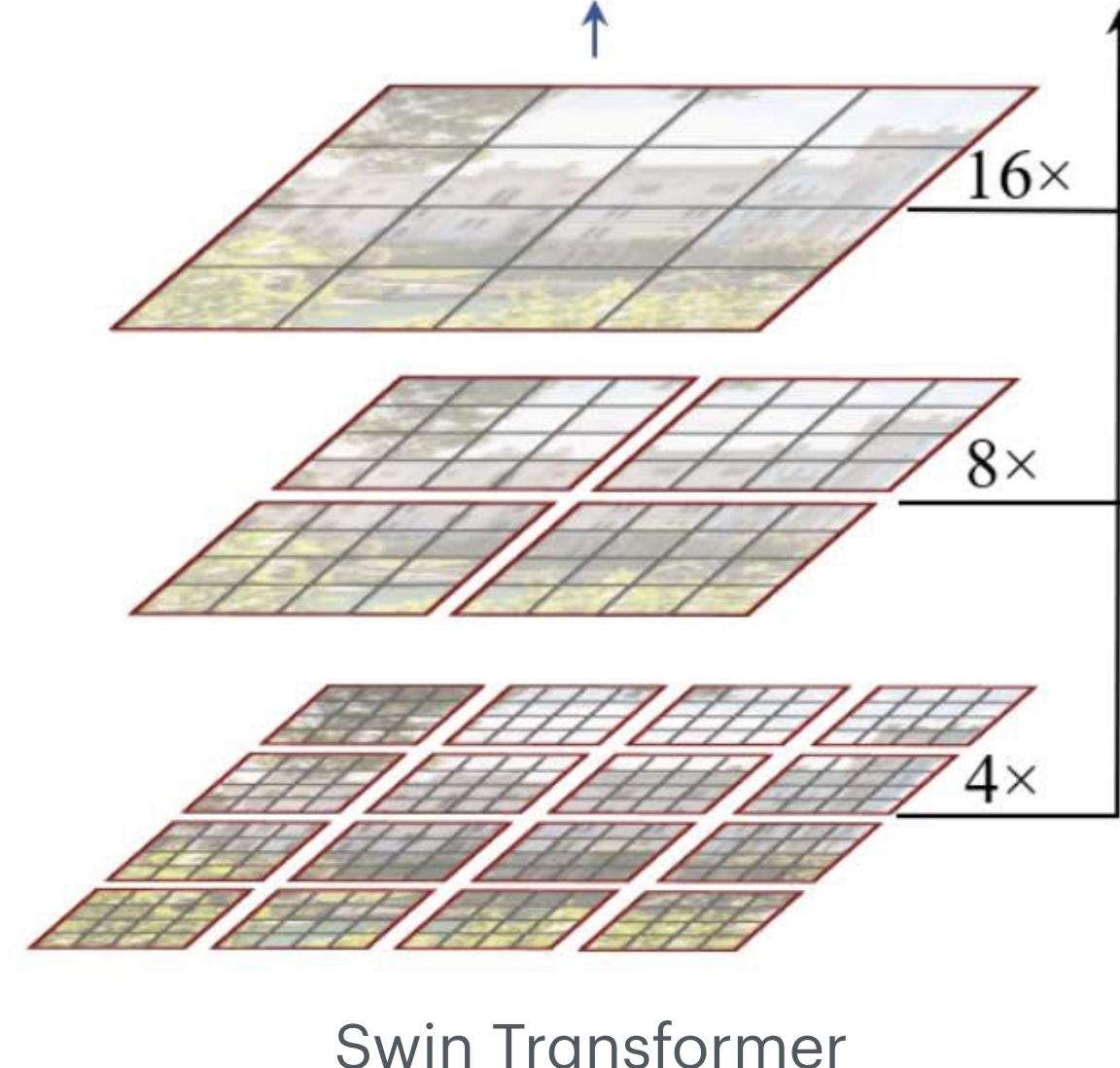
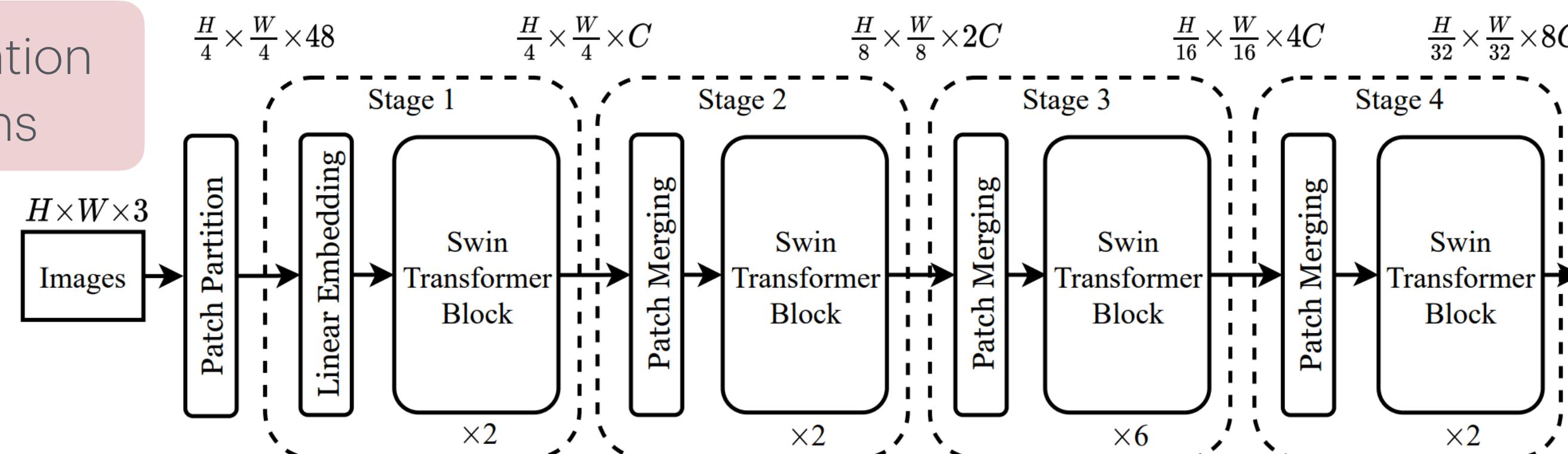
# Efficient ViT and Acceleration Techniques

# Window Attention

Restrict attention computation within local windows



Compute attention over all tokens

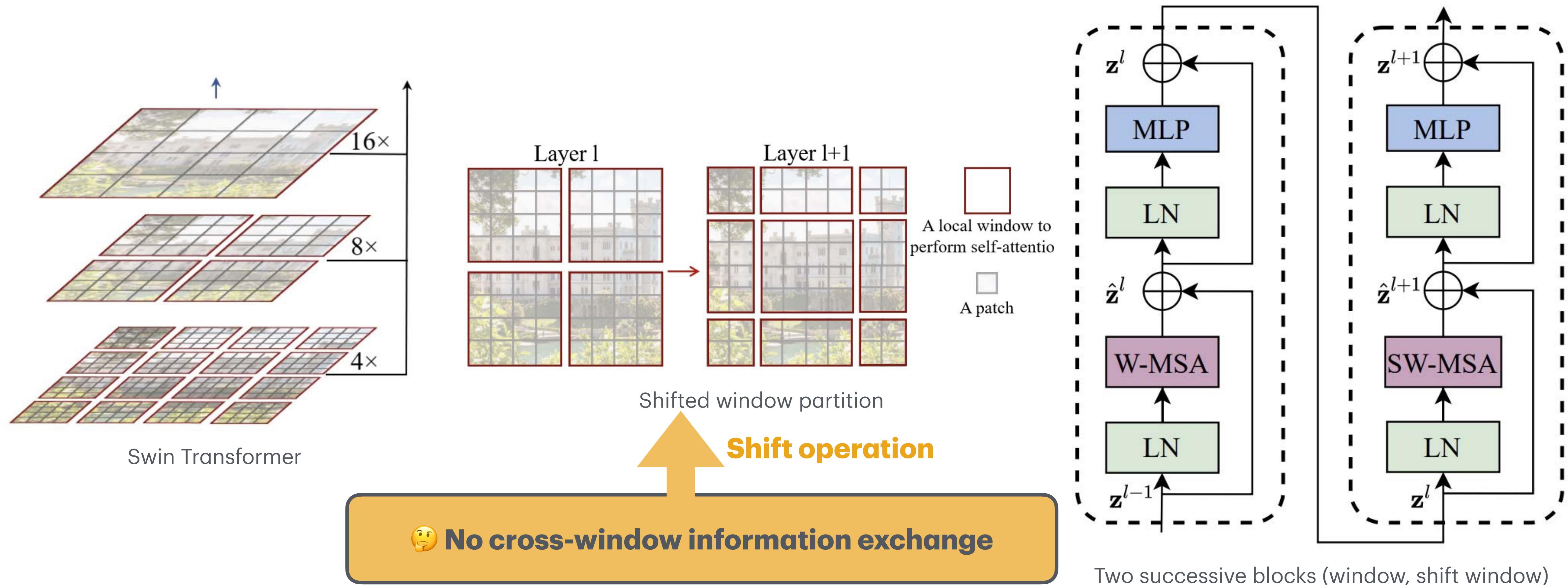


Compute attention within fixed-size local windows

🤔 **But there is no cross-window information exchange**

# Window Attention

Restrict attention computation within local windows



# Window Attention

## SwinTransformer results

(a) Regular ImageNet-1K trained models					
method	image size	#param.	FLOPs	throughput (image / s)	ImageNet top-1 acc.
ViT-B/16 [20]	384 <sup>2</sup>	86M	55.4G	85.9	77.9
ViT-L/16 [20]	384 <sup>2</sup>	307M	190.7G	27.3	76.5
DeiT-S [63]	224 <sup>2</sup>	22M	4.6G	940.4	79.8
DeiT-B [63]	224 <sup>2</sup>	86M	17.5G	292.3	81.8
DeiT-B [63]	384 <sup>2</sup>	86M	55.4G	85.9	83.1
Swin-T	224 <sup>2</sup>	29M	4.5G	755.2	81.3
Swin-S	224 <sup>2</sup>	50M	8.7G	436.9	83.0
Swin-B	224 <sup>2</sup>	88M	15.4G	278.1	83.5
Swin-B	384 <sup>2</sup>	88M	47.0G	84.7	84.5

Better ImageNet accuracy than ViT

Method	mini-val		test-dev		#param. FLOPs
	AP <sup>box</sup>	AP <sup>mask</sup>	AP <sup>box</sup>	AP <sup>mask</sup>	
RepPointsV2* [12]	-	-	52.1	-	- -
GCNet* [7]	51.8	44.7	52.3	45.4	- 1041G
RelationNet++* [13]	-	-	52.7	-	- -
SpineNet-190 [21]	52.6	-	52.8	-	164M 1885G
ResNeSt-200* [78]	52.5	-	53.3	47.1	- -
EfficientDet-D7 [59]	54.4	-	55.1	-	77M 410G
DetectoRS* [46]	-	-	55.7	48.5	- -
YOLOv4 P7* [4]	-	-	55.8	-	- -
Copy-paste [26]	55.9	47.2	56.0	47.4	185M 1440G
X101-64 (HTC++)	52.3	46.0	-	-	155M 1033G
Swin-B (HTC++)	56.4	49.1	-	-	160M 1043G
Swin-L (HTC++)	57.1	49.5	57.7	50.2	284M 1470G
Swin-L (HTC++)*	<b>58.0</b>	<b>50.4</b>	<b>58.7</b>	<b>51.1</b>	284M -

Applicable to downstream tasks (e.g., detection)

Method	Backbone	val	test	#param. FLOPs	FPS
		mIoU	score		
DANet [23]	ResNet-101	45.2	-	69M	1119G 15.2
DLab.v3+ [11]	ResNet-101	44.1	-	63M	1021G 16.0
ACNet [24]	ResNet-101	45.9	38.5	-	-
DNL [71]	ResNet-101	46.0	56.2	69M	1249G 14.8
OCRNet [73]	ResNet-101	45.3	56.0	56M	923G 19.3
UperNet [69]	ResNet-101	44.9	-	86M	1029G 20.1
OCRNet [73]	HRNet-w48	45.7	-	71M	664G 12.5
DLab.v3+ [11]	ResNeSt-101	46.9	55.1	66M	1051G 11.9
DLab.v3+ [11]	ResNeSt-200	48.4	-	88M	1381G 8.1
SETR [81]	T-Large <sup>†</sup>	50.3	61.7	308M	- -
UperNet	DeiT-S <sup>†</sup>	44.0	-	52M	1099G 16.2
UperNet	Swin-T	46.1	-	60M	945G 18.5
UperNet	Swin-S	49.3	-	81M	1038G 15.2
UperNet	Swin-B <sup>‡</sup>	51.6	-	121M	1841G 8.7
UperNet	Swin-L <sup>‡</sup>	<b>53.5</b>	<b>62.8</b>	234M	3230G 6.2

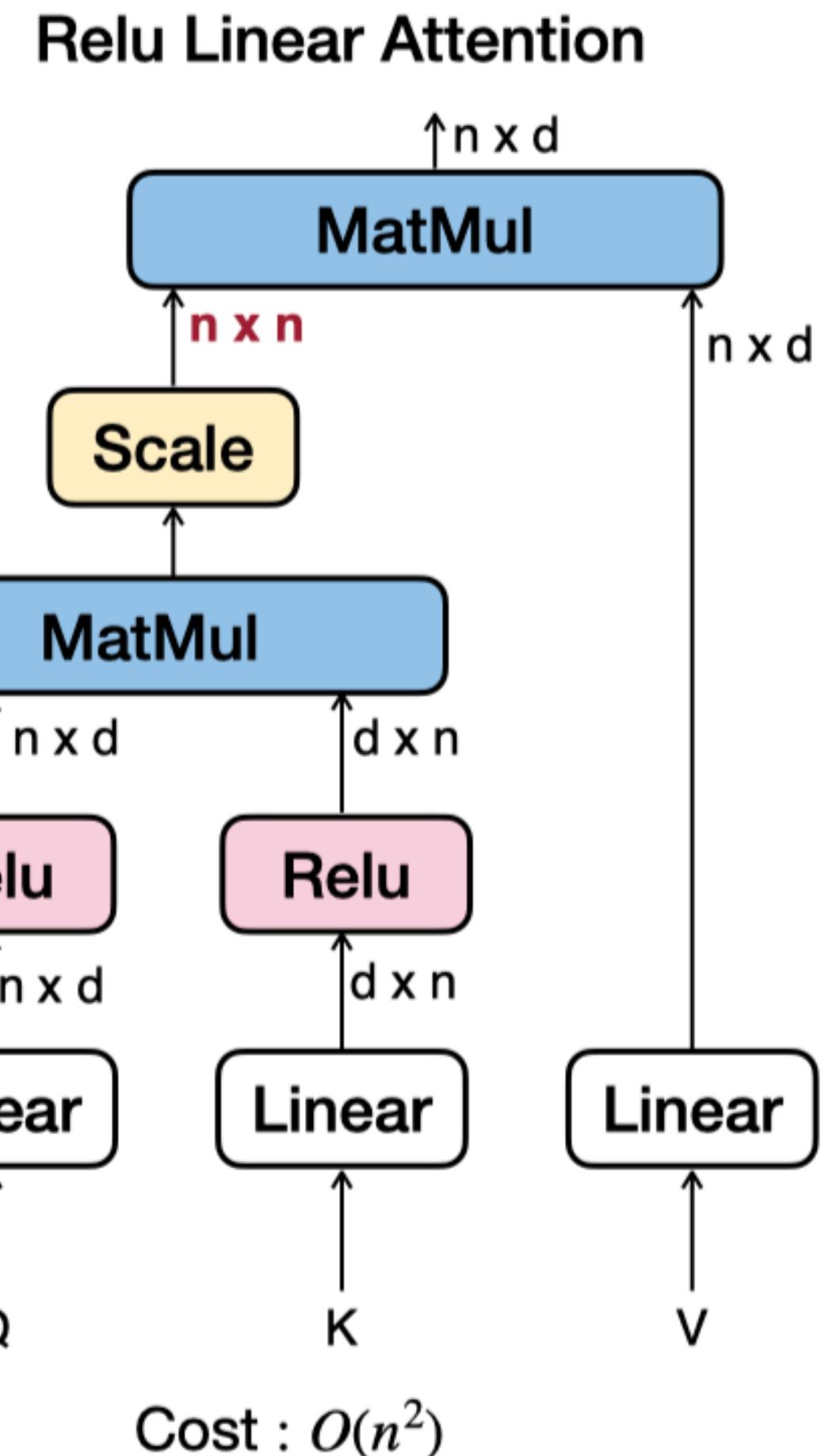
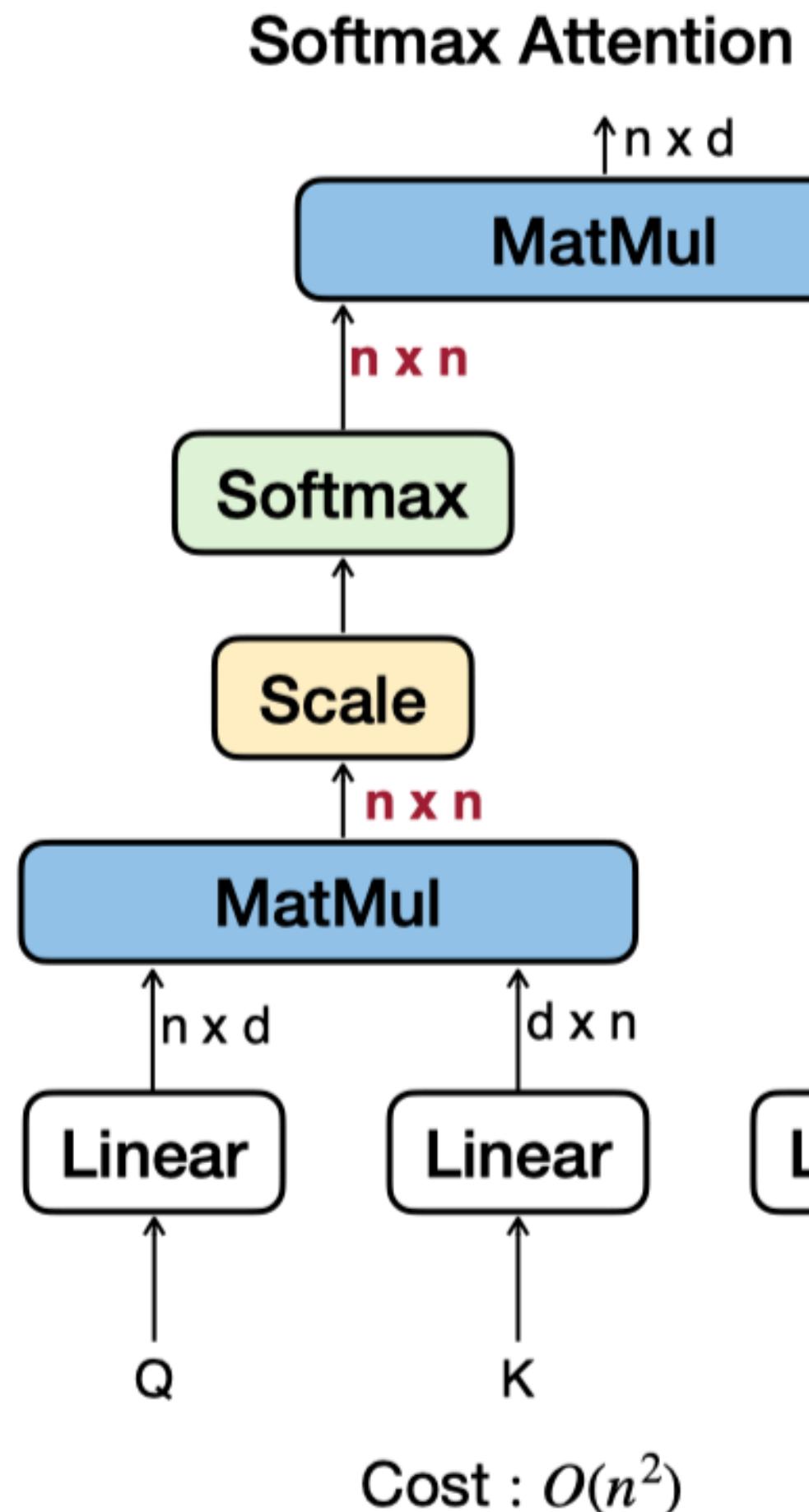
Applicable to downstream tasks (e.g., segmentation)

# Linear Attention

Replace softmax attention with linear attention

**Softmax Attention**

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$



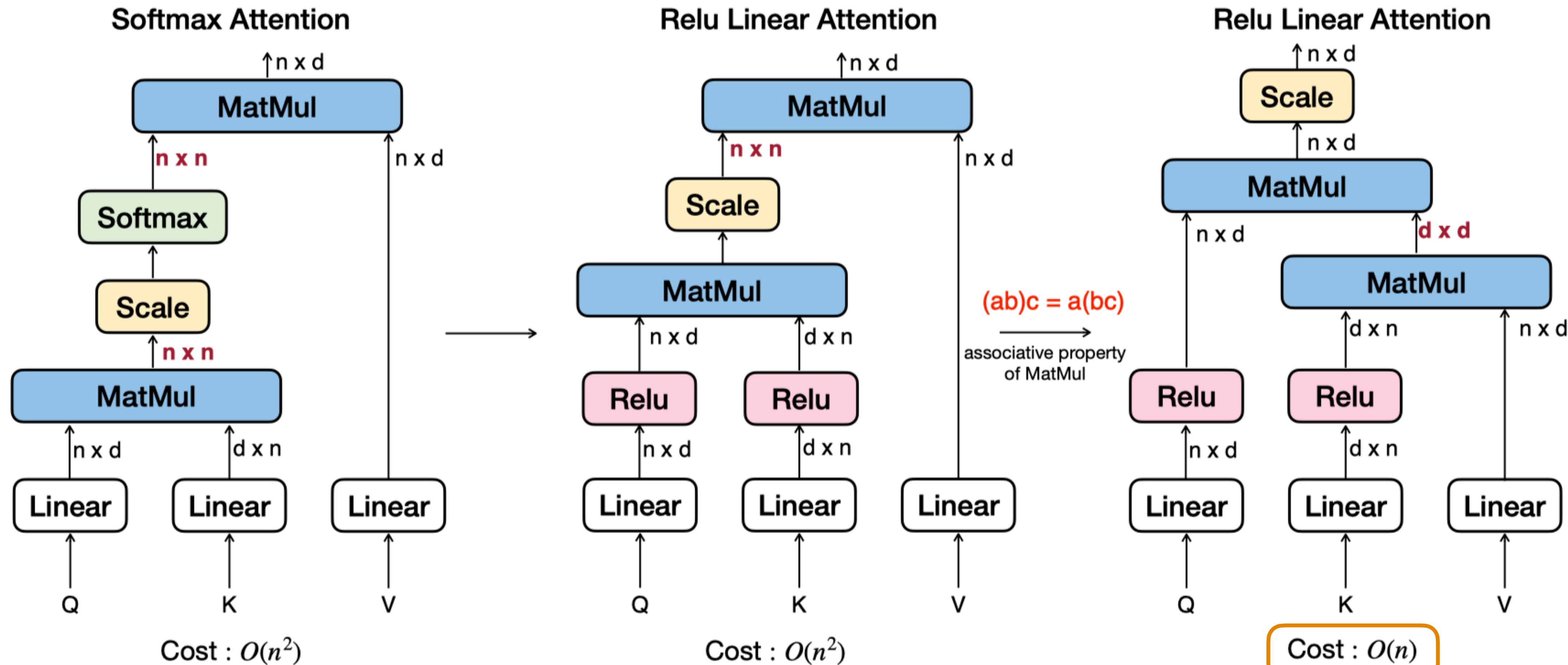
- Replace non-linear similarity function with linear similarity function

$$Sim(Q, K) = \exp\left(\frac{QK^T}{\sqrt{d}}\right)$$

$$Sim(Q, K) = \text{ReLU}(Q)\text{ReLU}(K)^T$$

# Linear Attention

Replace softmax attention with linear attention

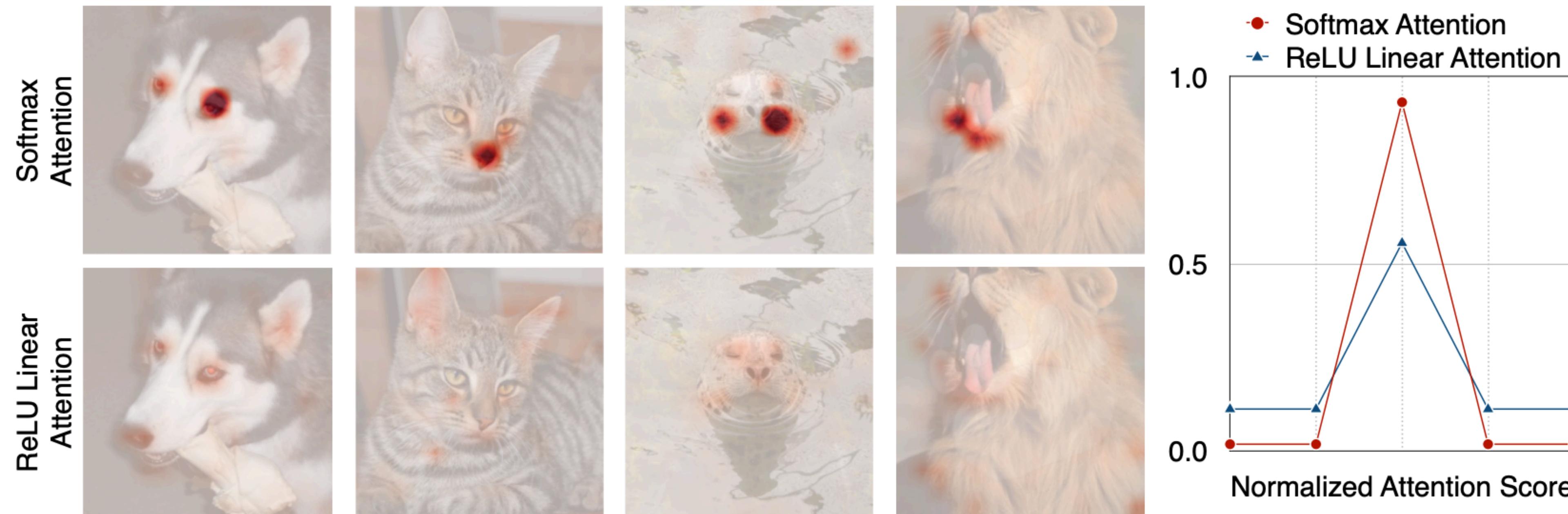


# Linear Attention

## Comparison with softmax attention

- BUT, ReLU linear attention cannot produce sharp distributions. It is good at capturing global context information, but bad at capturing local information.
- It also lacks multi-scale learning ability.

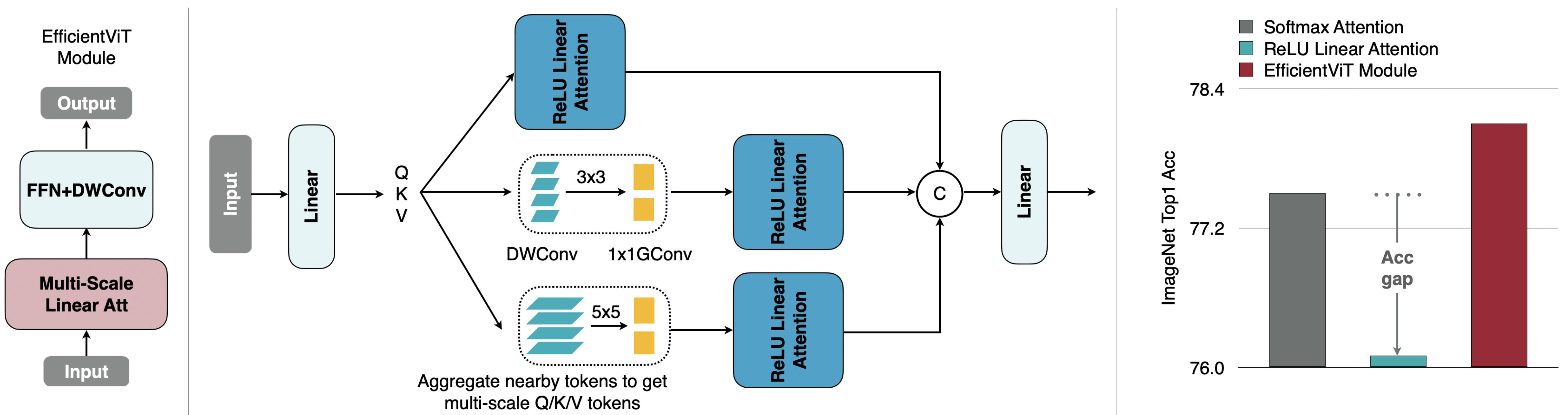
💡 How to address linear attention's limitation w/o losing efficiency?



# Linear Attention

## EfficientViT

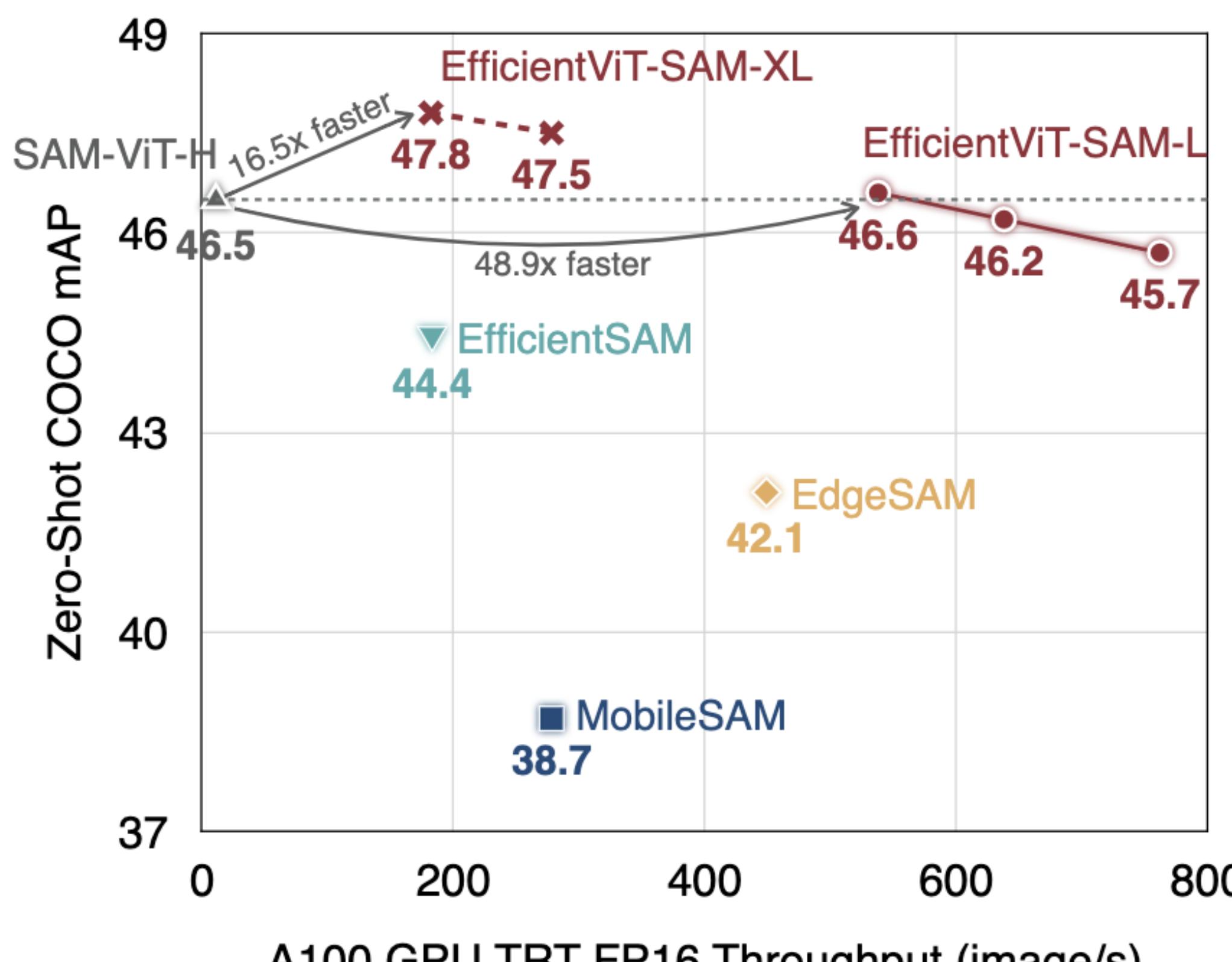
- Enhance linear attention by multi-scale aggregation
- Add depthwise convolution in FFN to further enhance linear attention's local information extraction ability



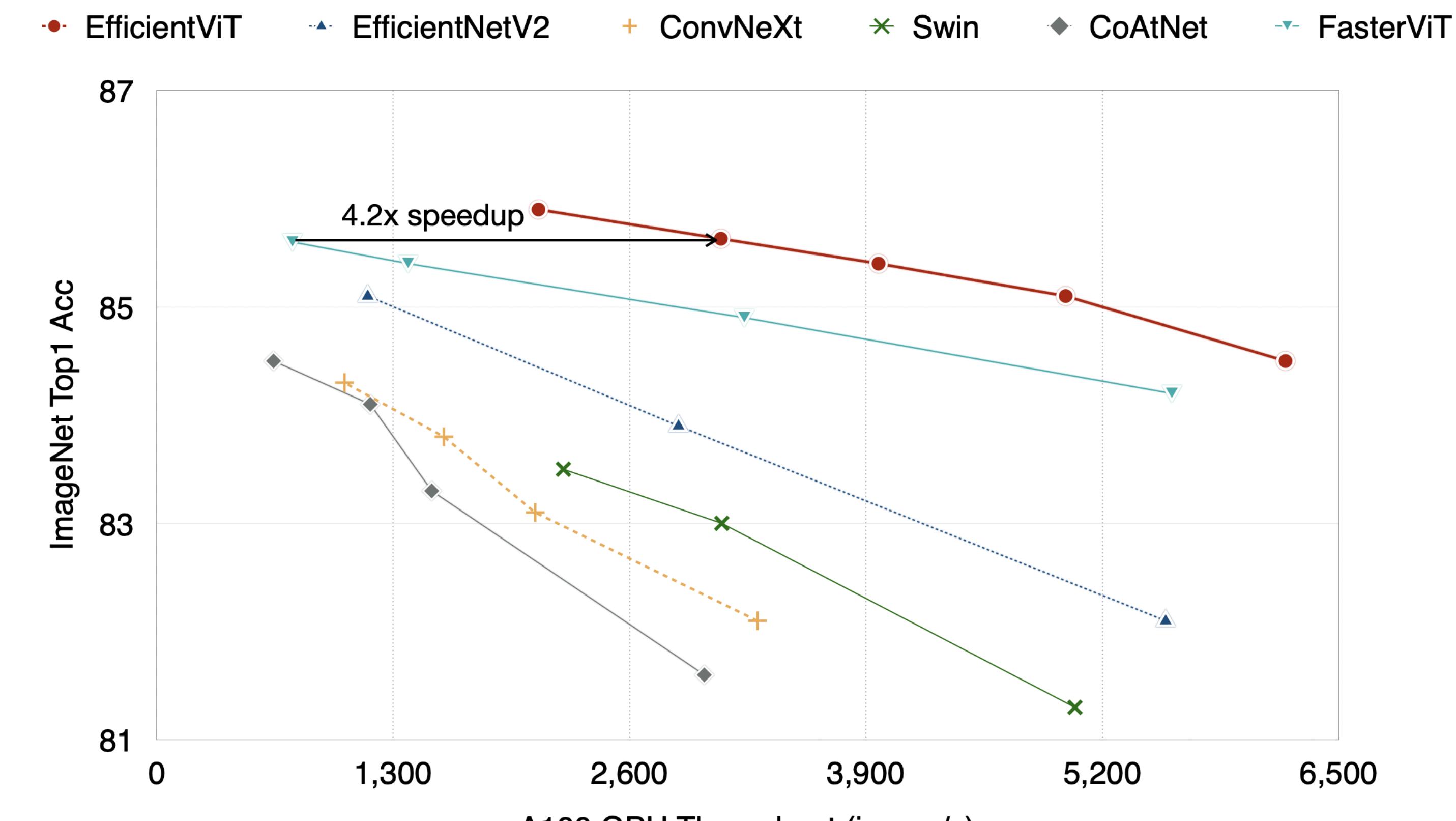
Cai, H., Li, J., Hu, M., Gan, C., & Han, S. (2023). EfficientViT: Lightweight multi-scale attention for on-device semantic segmentation. arXiv preprint arXiv:2205.14756, 2.

# Linear Attention

## EfficientViT results



On Segment Anything



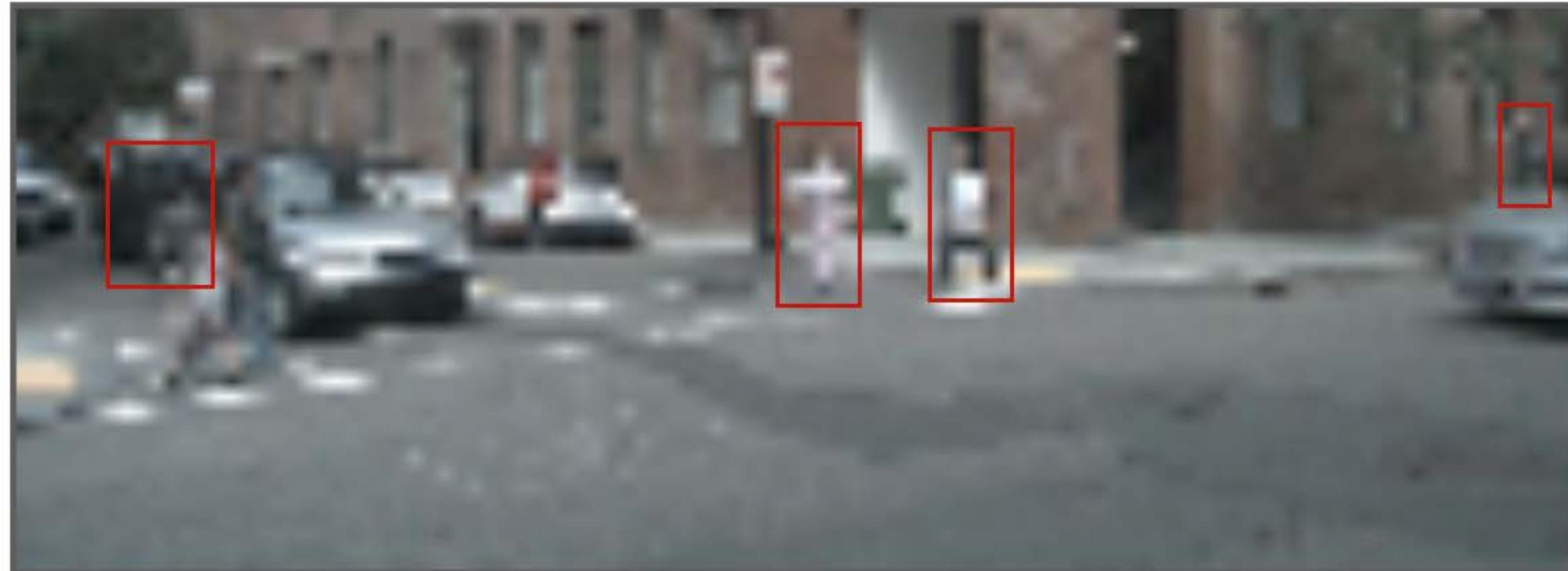
On Image Classification

# Sparse Attention

Motivation: sparse, high-resolution, or dense, low-resolution?

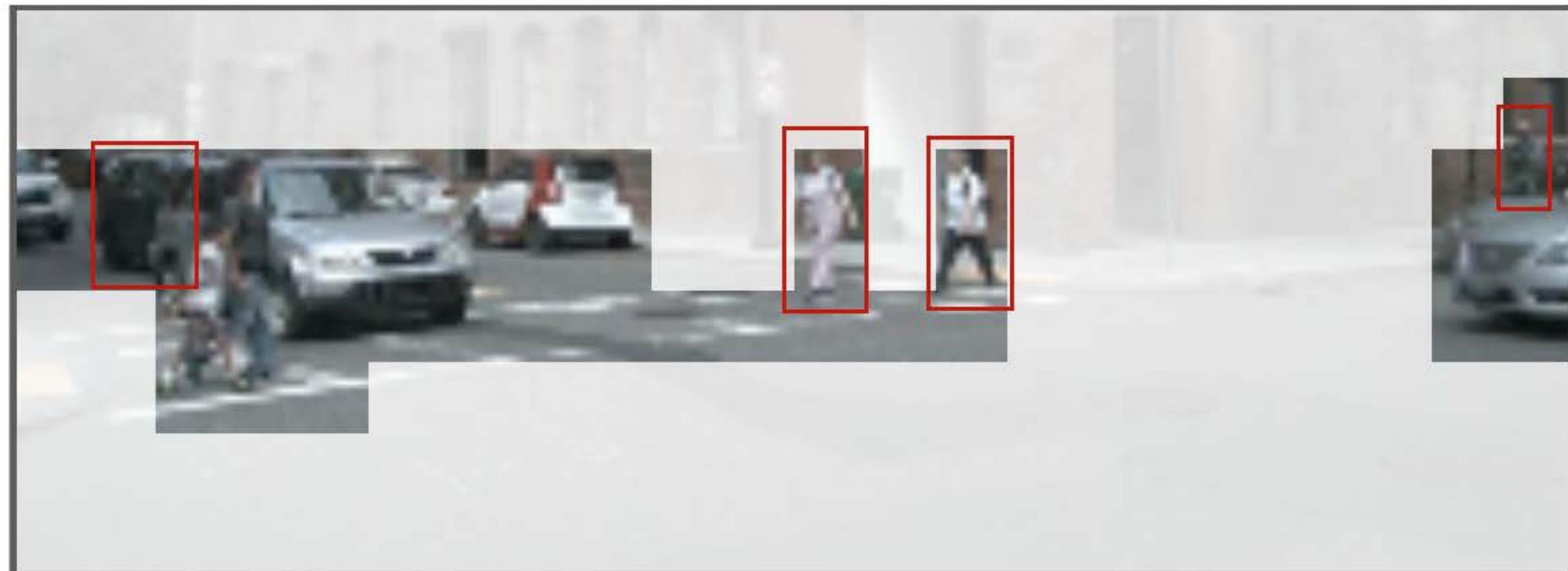
## Uniform Resizing

**Low Resolution (0.5X)**  
**Dense Pixels (100%)**



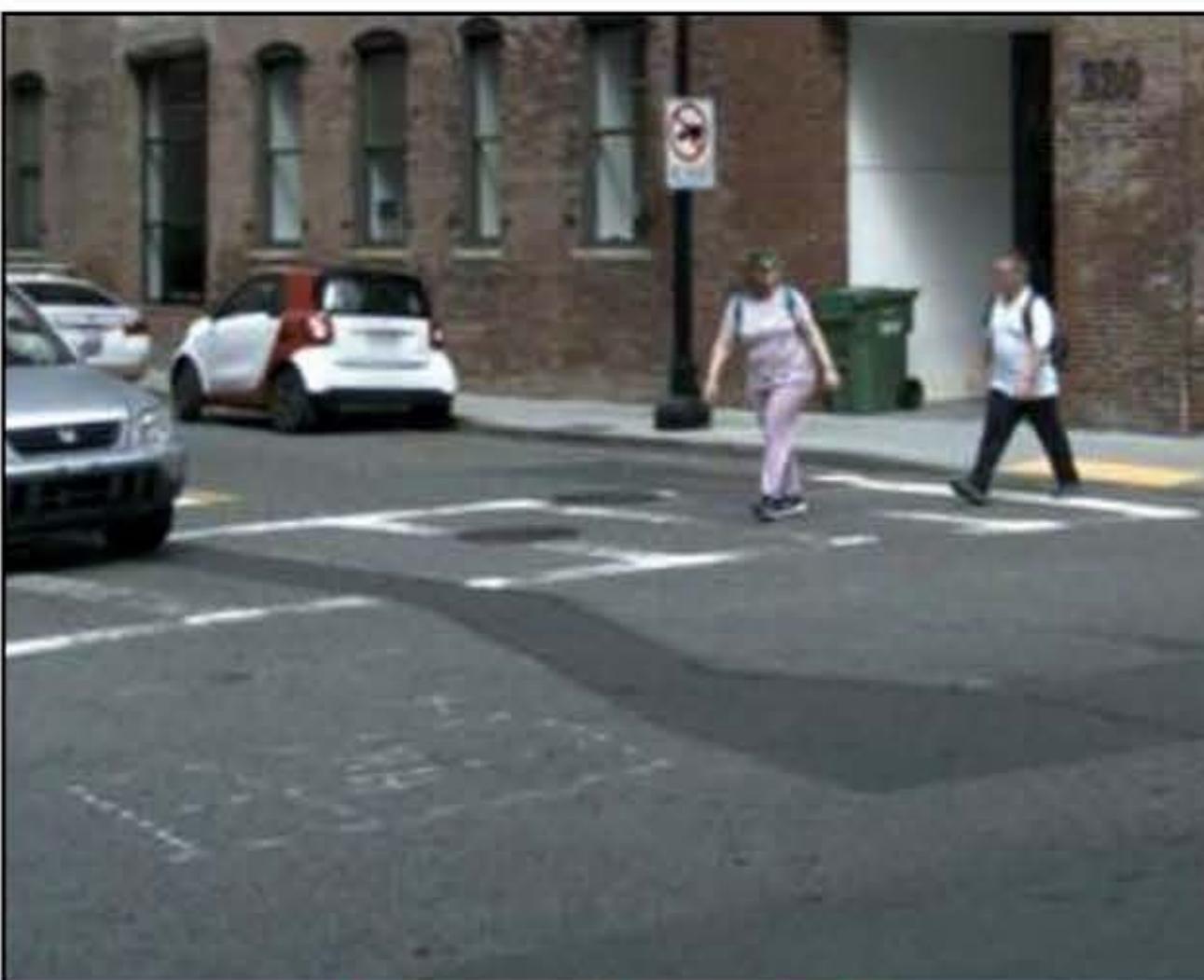
## Activation Pruning

**High Resolution (1X)**  
**Sparse Pixels (25%)**

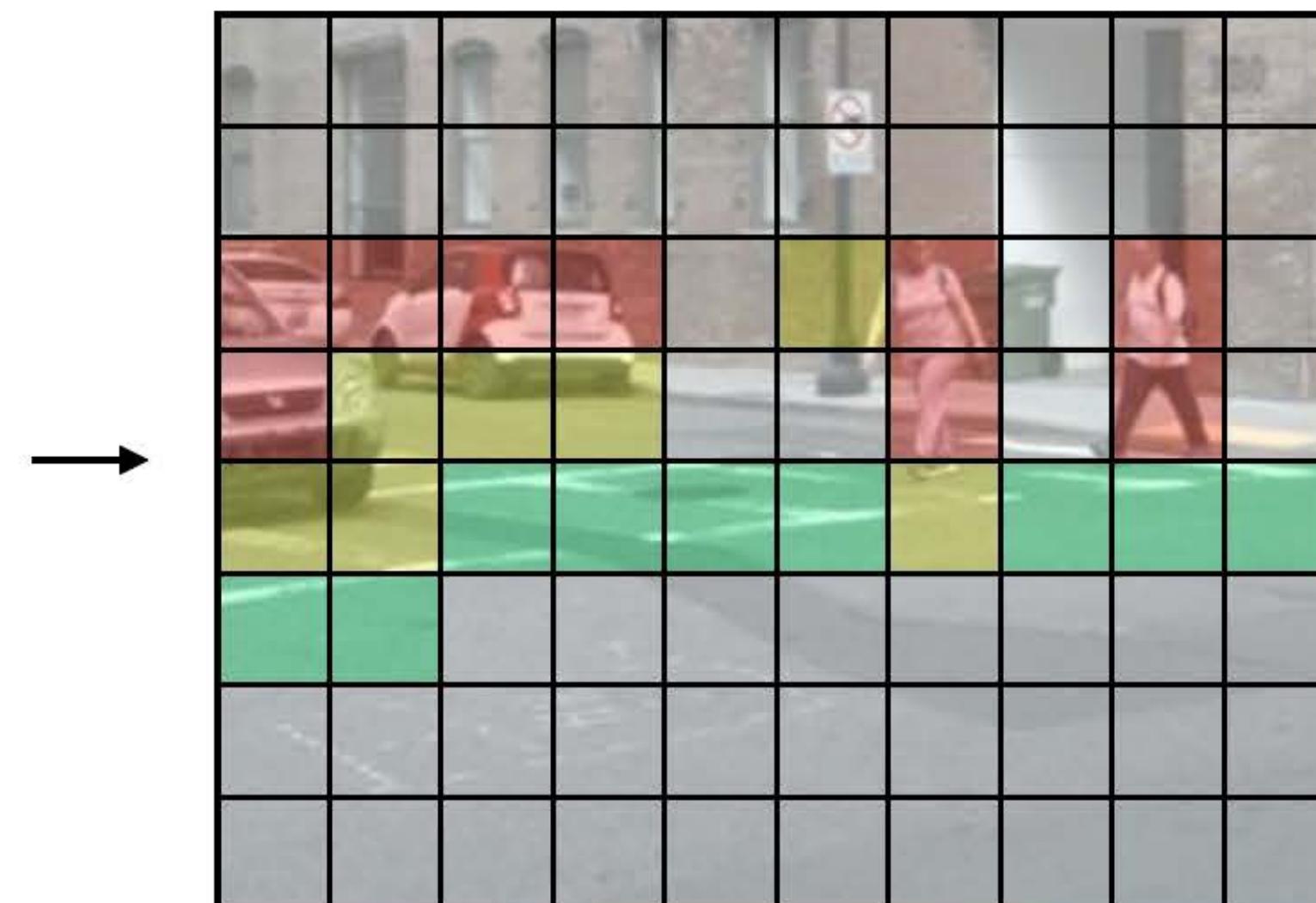


# Sparse Attention

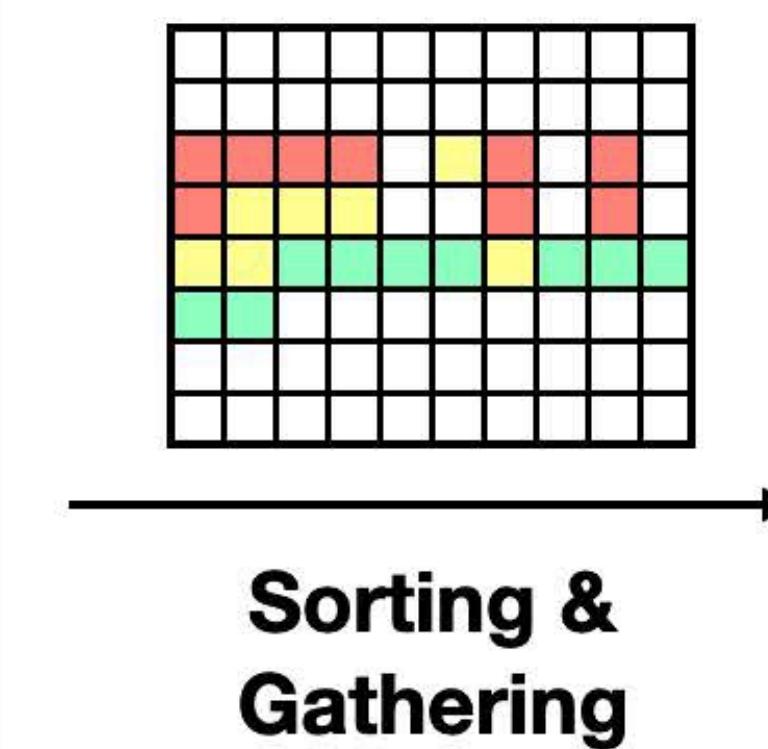
Step 1: window activation pruning (with non-uniform sparsity)



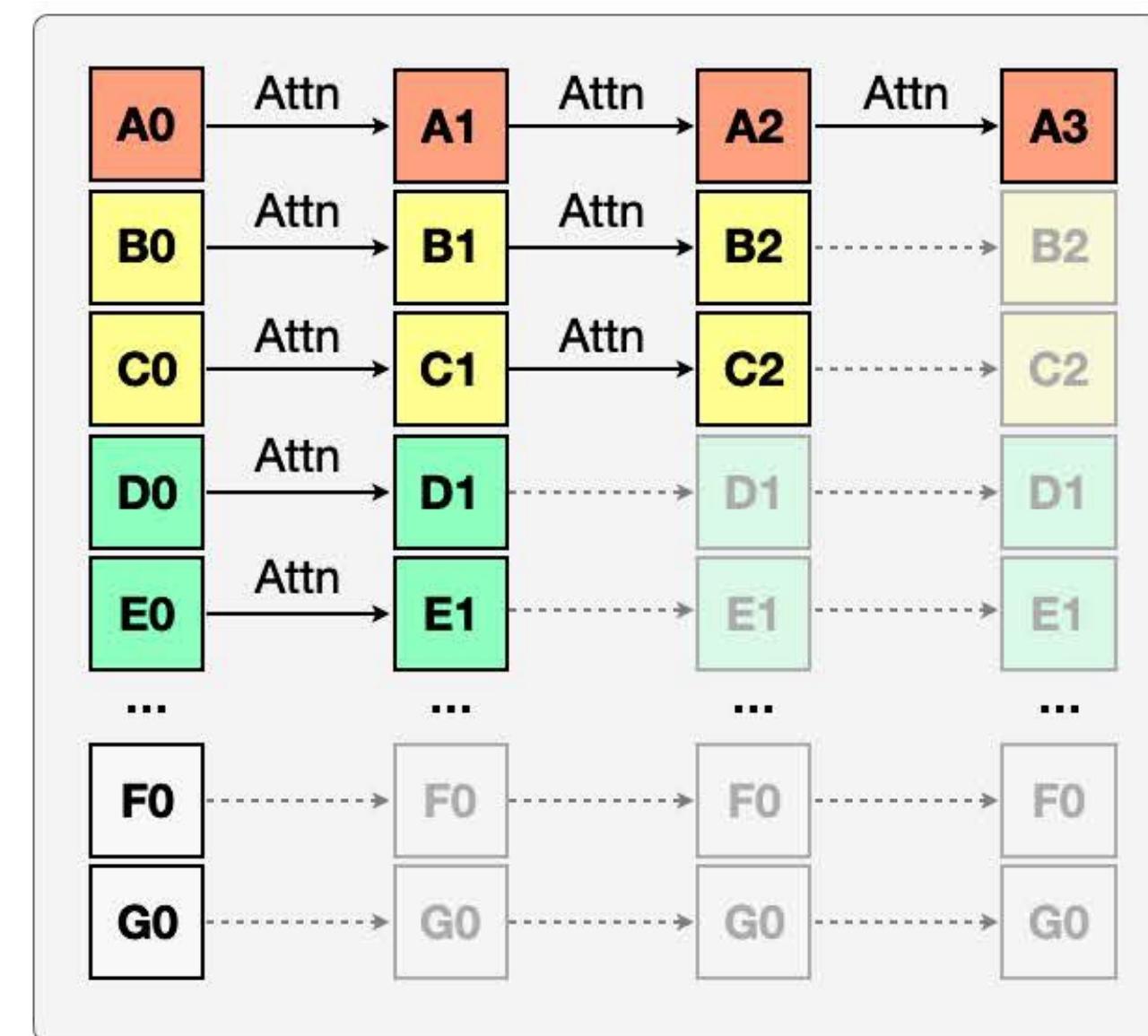
**Input Image**  
(or Input Feature Map)



**Window Importance**  
(L2 Activation Magnitude)



**Sorting &  
Gathering**

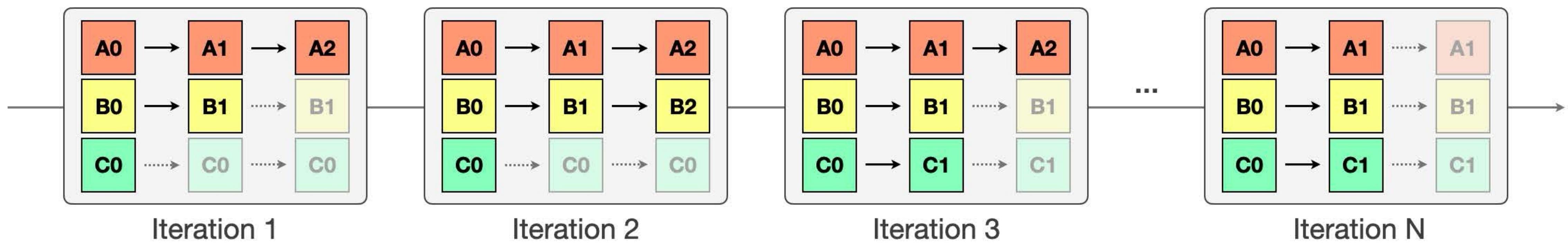


**Sparse Window Attention**

# Sparse Attention

## Step 2: sparsity-aware adaption

**Goal:** Assess the model's accuracy under different activation sparsity settings both **efficiently** and **accurately**.

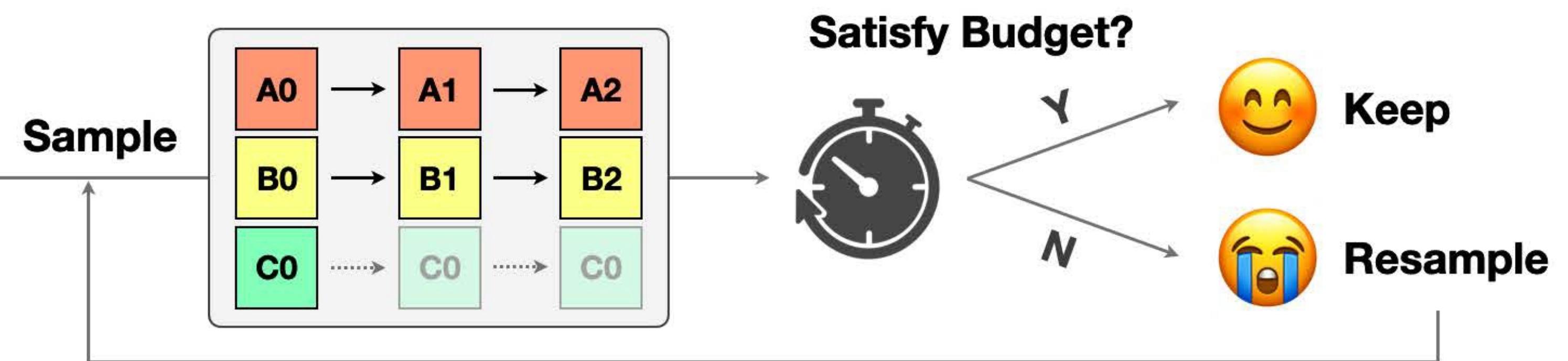


As the original model is trained with **only dense activations**, we **improve its sparsity awareness** by finetuning it with **randomly sampled** layerwise activation sparsity configurations at each training iteration.

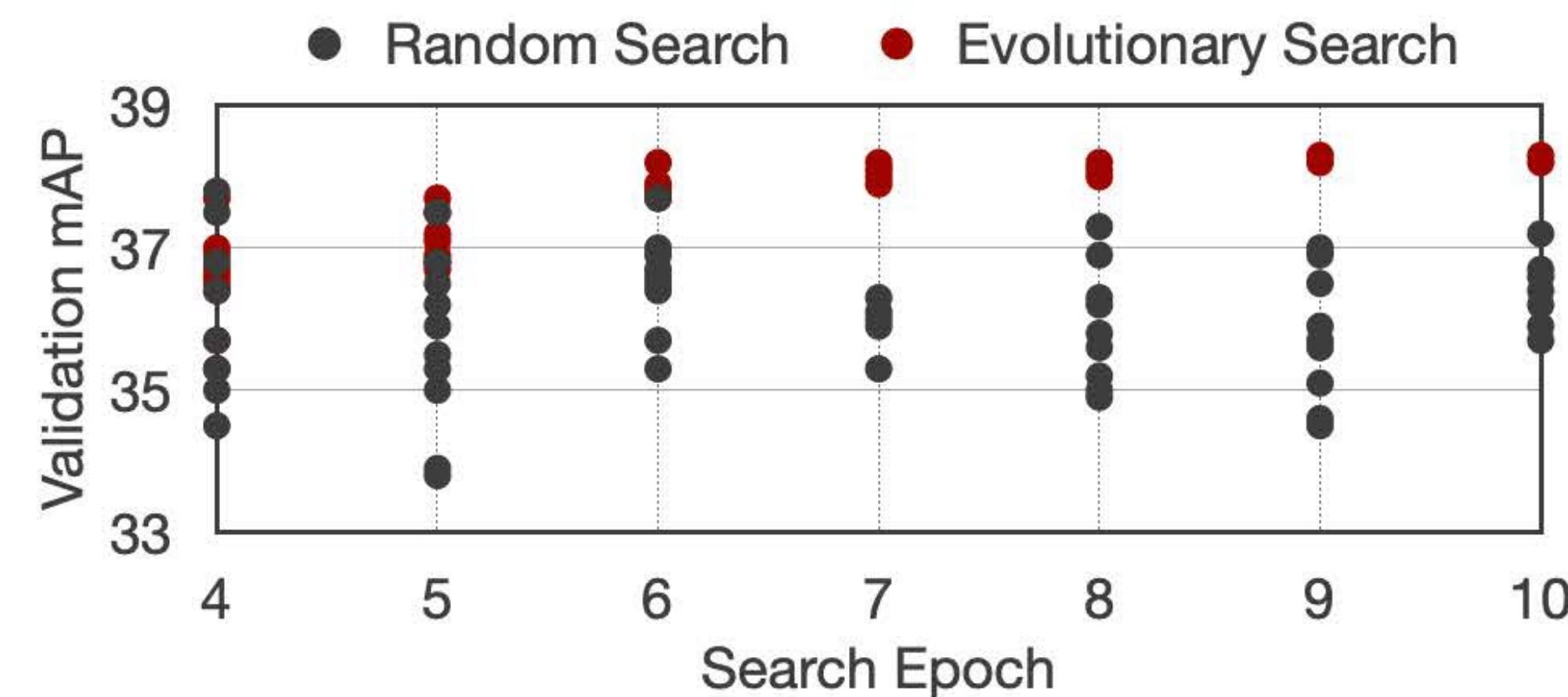
# Sparse Attention

## Step 3: resource-constrained search

**Goal:** Discover the **optimal** layerwise activation sparsity configuration **under a given latency budget**.



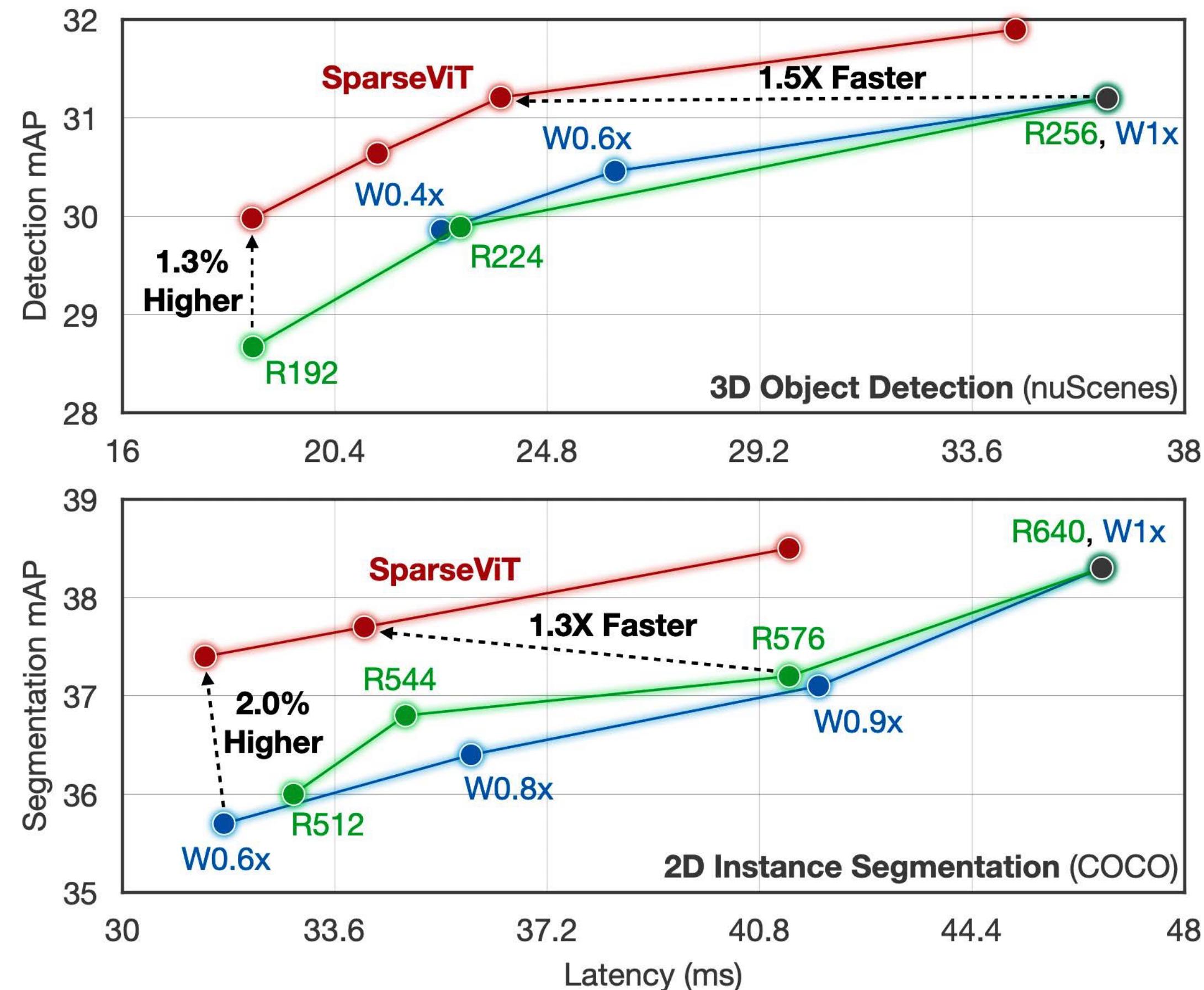
We enforce the latency constraint using **rejection sampling** (repeated resampling until satisfaction).



Evolutionary search is **sample-efficient**!

# Sparse Attention

## SparseViT results

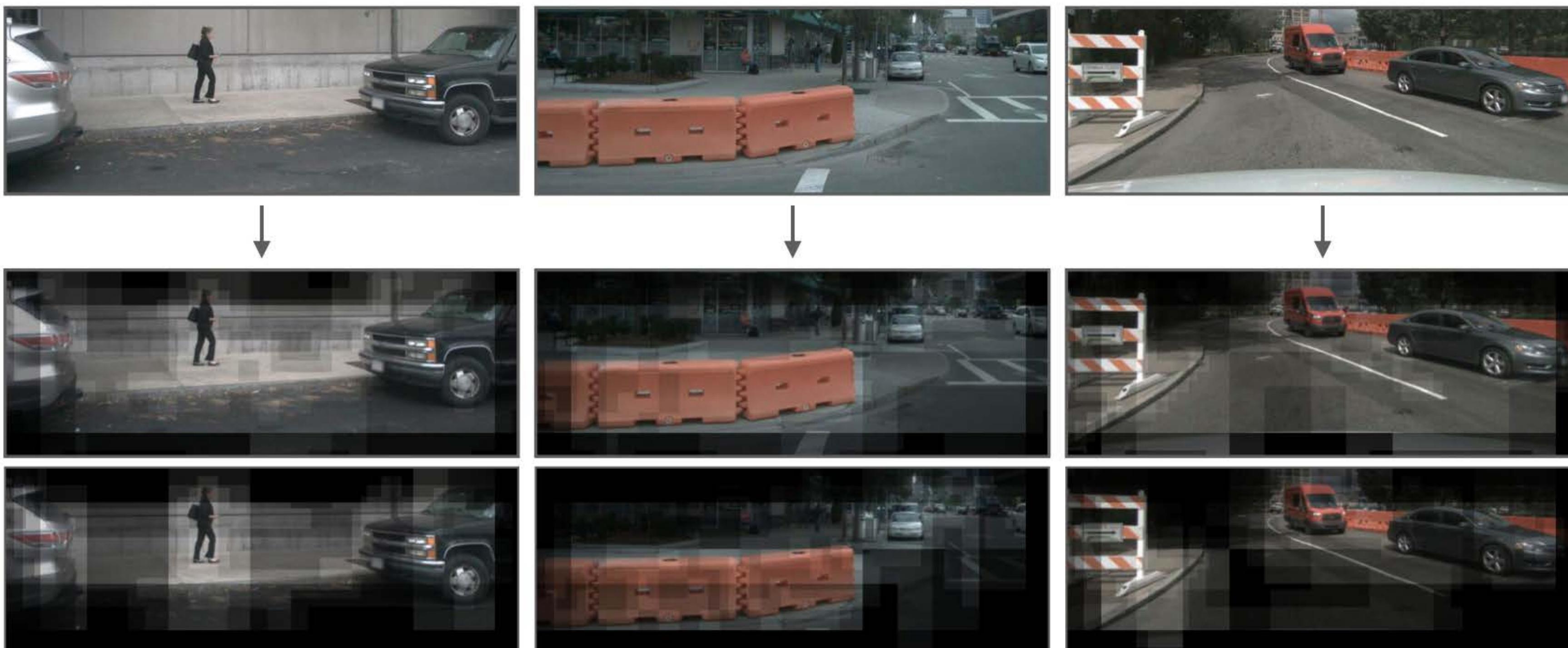


3D object detection and 2D instance segmentation

# Sparse Attention

SparseViT results: sparsity masks under different latency budgets

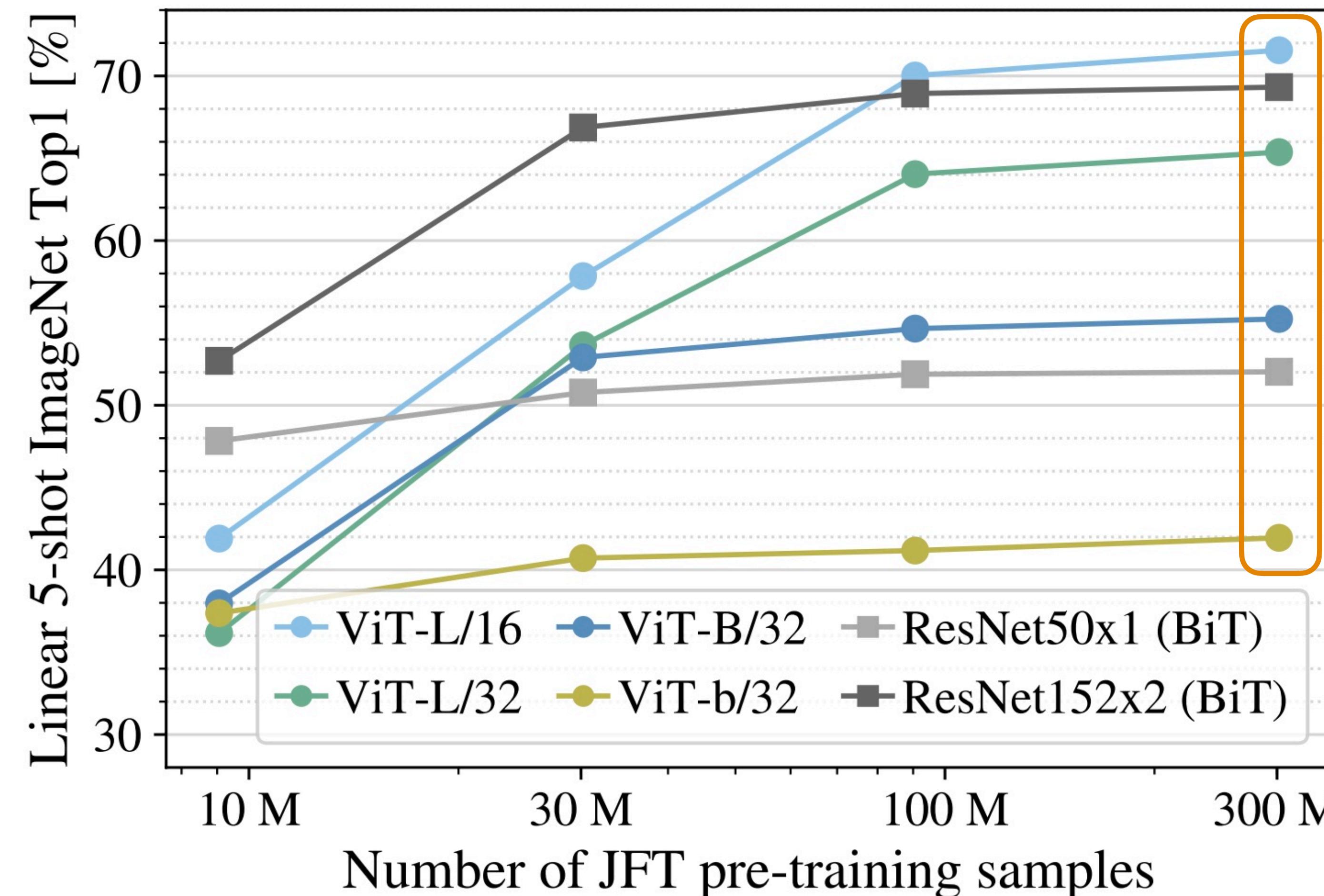
**SparseViT learns to prune **irrelevant background windows** while retaining **informative foreground ones!****



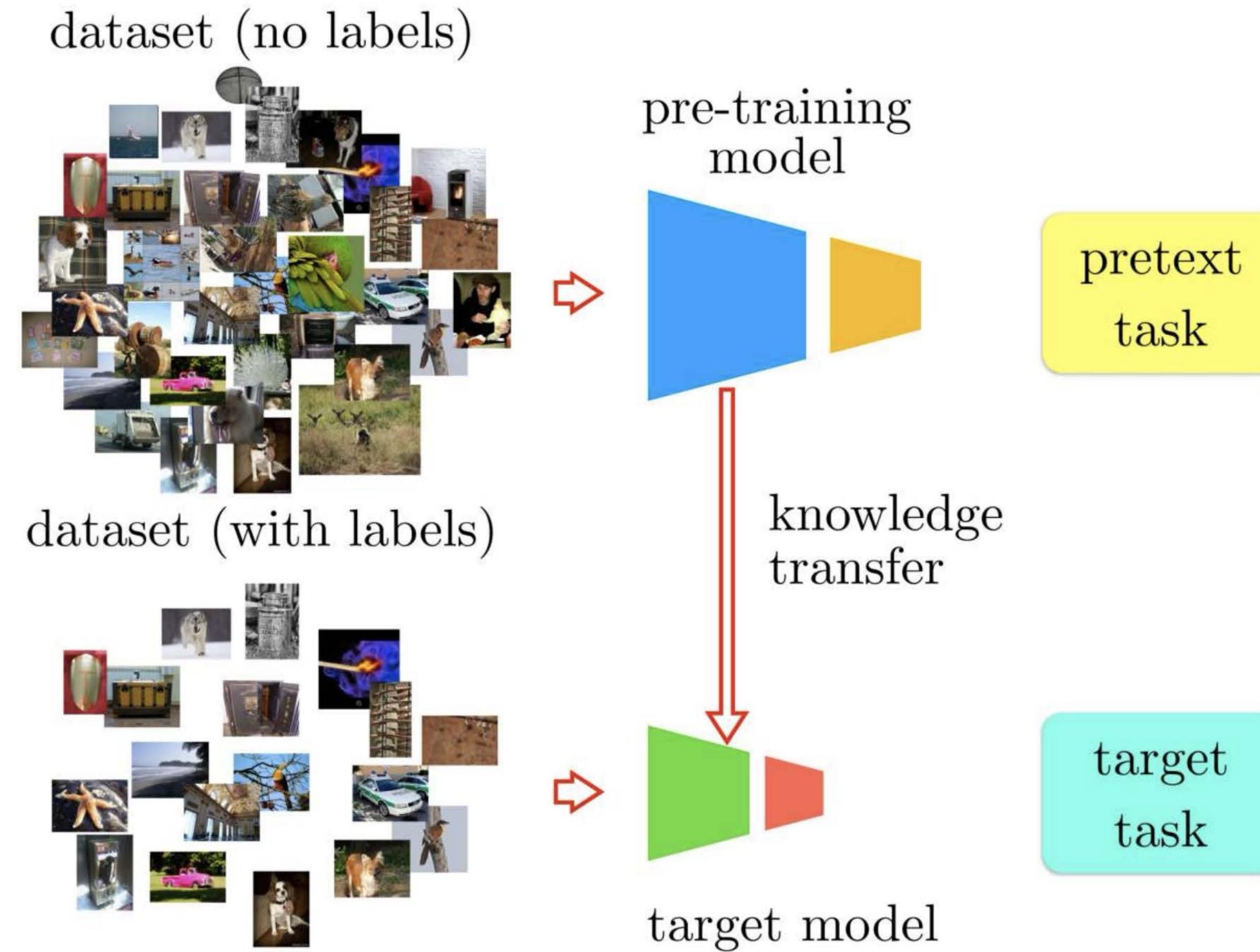
# Self-supervised Learning for ViT

# ViT Needs Large Datasets to Work Well

Labeling large datasets is costly!



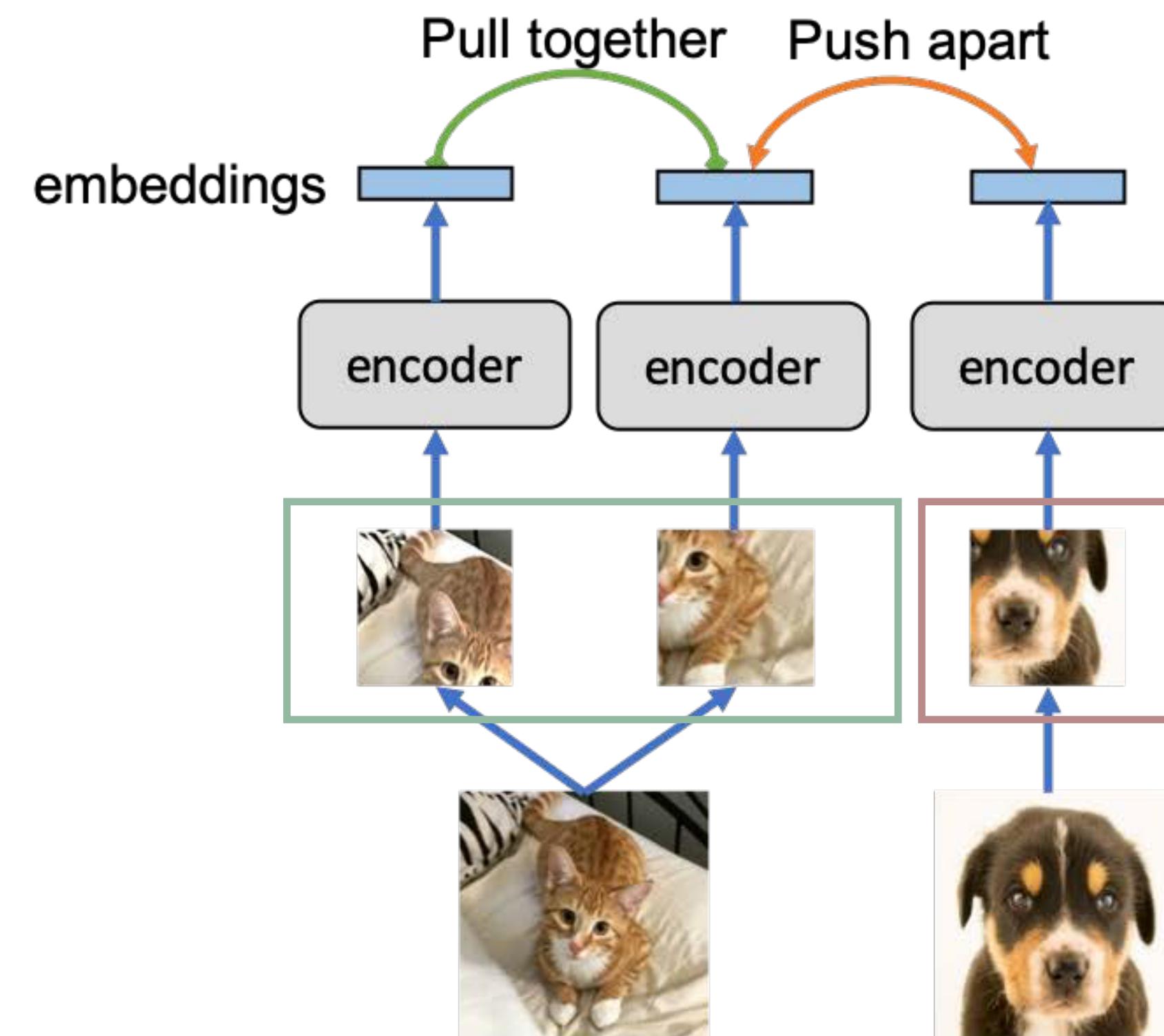
# Self-Supervised Learning for ViT



A knowledge transfer method to decouple pre-training and fune-tuning

# Contrastive Learning

- Positive samples: random views of the **same** images
- Negative samples: random views of **other** images



Prototypical contrastive learning: pushing the frontiers of  
unsupervised learning

Chen, X., Xie, S., & He, K. (2021). An empirical study of training self-supervised vision transformers. In Proceedings of the IEEE/CVF international conference on computer vision (pp. 9640-9649).

# Contrastive Learning

## Results

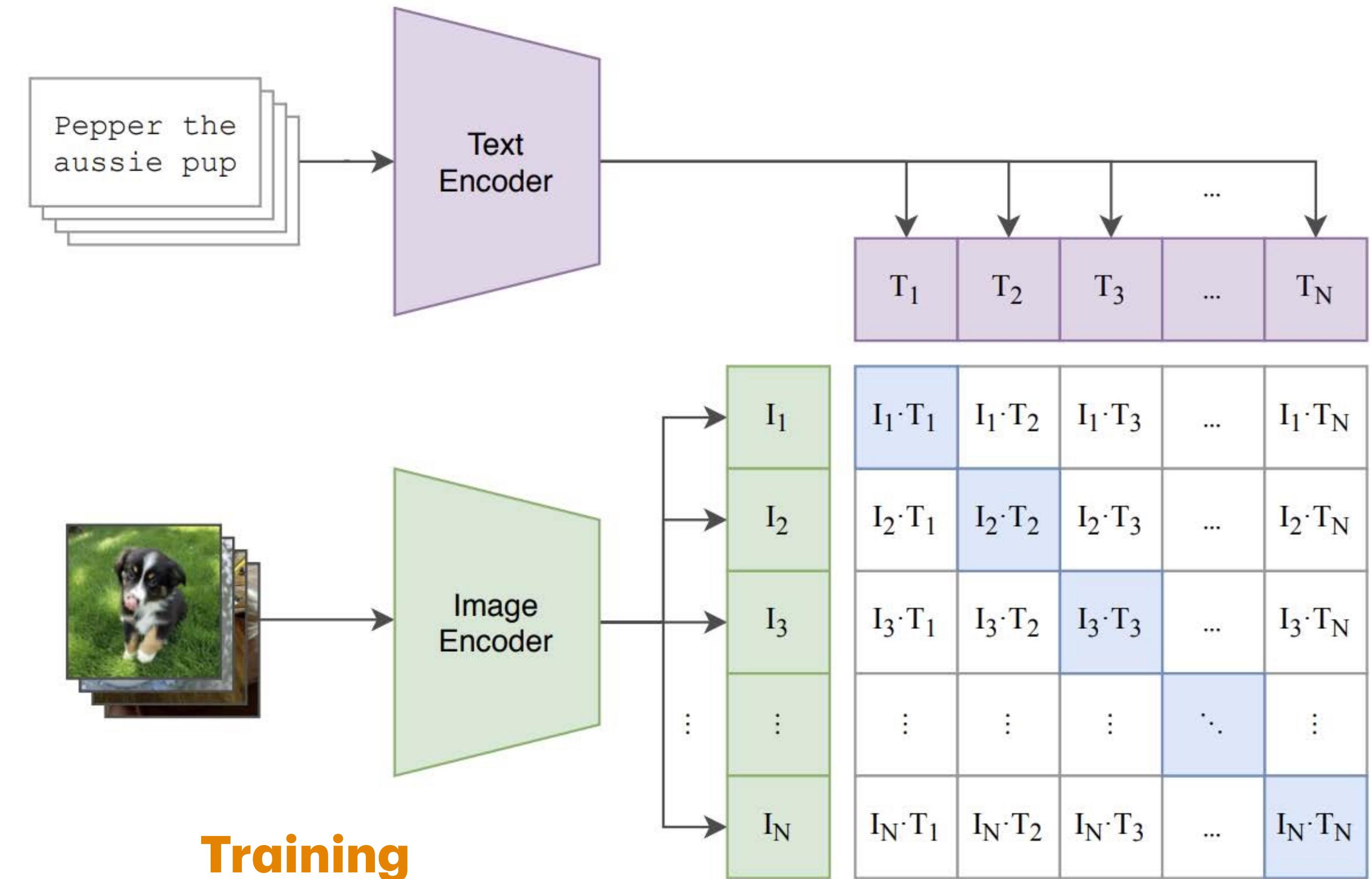
- Directly training large ViT models on small datasets **cannot achieve good results**. The accuracy decreases when increasing the model size for most cases.
- Self-supervised ViT models deliver **stronger** performances than their supervised counterparts. The accuracy increases when using large ViT models.

pre-train	CIFAR-10 [26]			CIFAR-100 [26]			Oxford Flowers-102 [33]			Oxford-IIIT-Pets [35]		
	ViT-B	ViT-L	ViT-H	ViT-B	ViT-L	ViT-H	ViT-B	ViT-L	ViT-H	ViT-B	ViT-L	ViT-H
random init.	77.8	77.1	75.9	48.5	48.3	48.0	54.4	54.3	52.8	40.1	42.8	40.4
ImNet supervised [16]	98.1	97.9	n/a	87.1	86.4	n/a	89.5	89.7	n/a	93.8	93.6	n/a
ImNet self-sup., MoCo v3	98.9 <span style="color: green;">↑0.8</span>	99.1 <span style="color: green;">↑1.2</span>	99.1	90.5 <span style="color: green;">↑3.4</span>	91.1 <span style="color: green;">↑4.7</span>	91.2	97.7 <span style="color: green;">↑8.2</span>	98.6 <span style="color: green;">↑8.9</span>	98.8	93.2 <span style="color: red;">↓0.6</span>	93.7 <span style="color: green;">↑0.1</span>	94.2

# Multi-Modal Contrastive Learning

## CLIP: Contrastive Language-Image Pre-Training

**Large-scale contrastive learning  
using paired image-text samples**

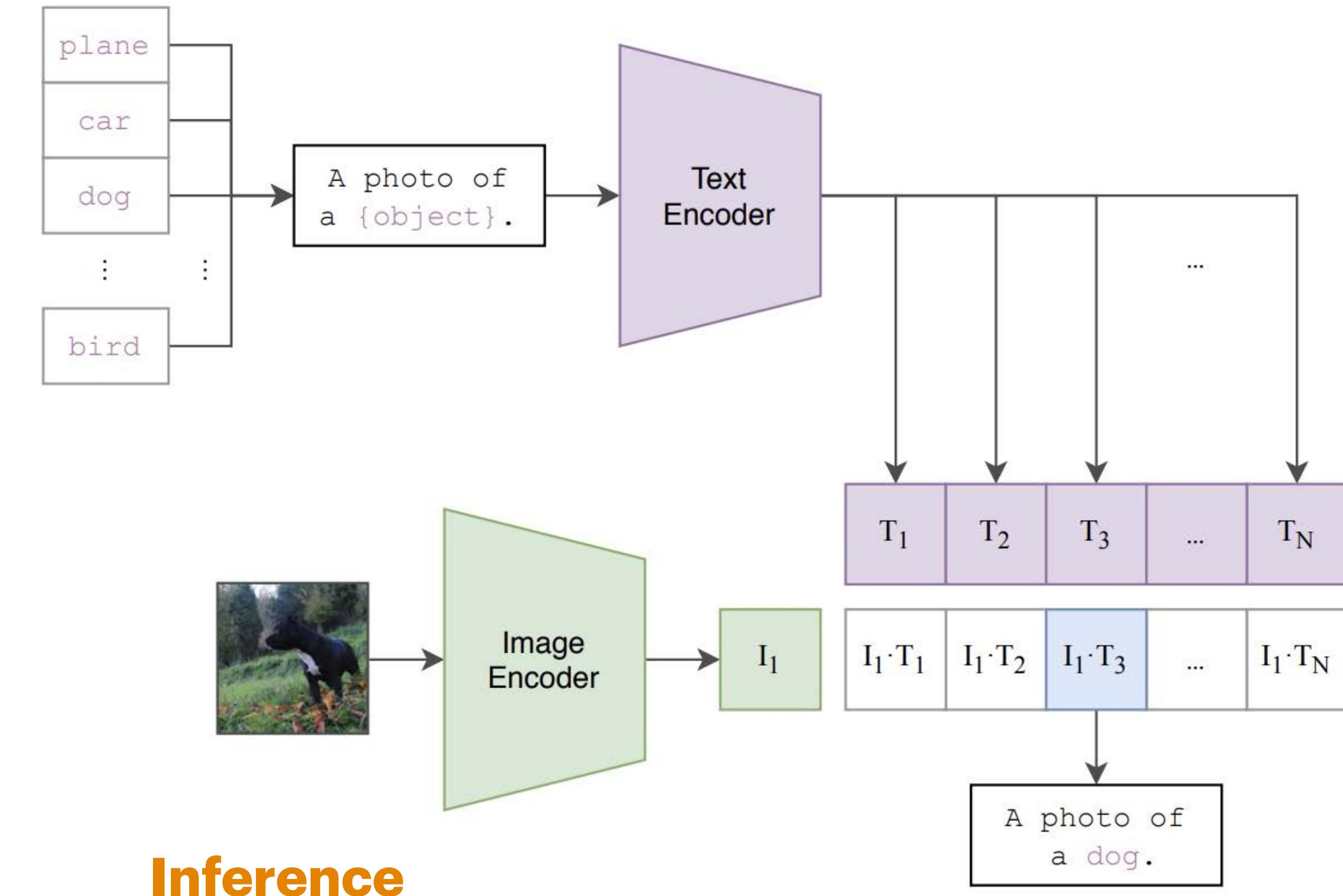


# Multi-Modal Contrastive Learning

## CLIP: Contrastive Language-Image Pre-Training

**Large-scale contrastive learning  
using paired image-text samples**

**Enable zero-shot and open-  
vocabulary classification (no fine-  
tuning, any #classes)**

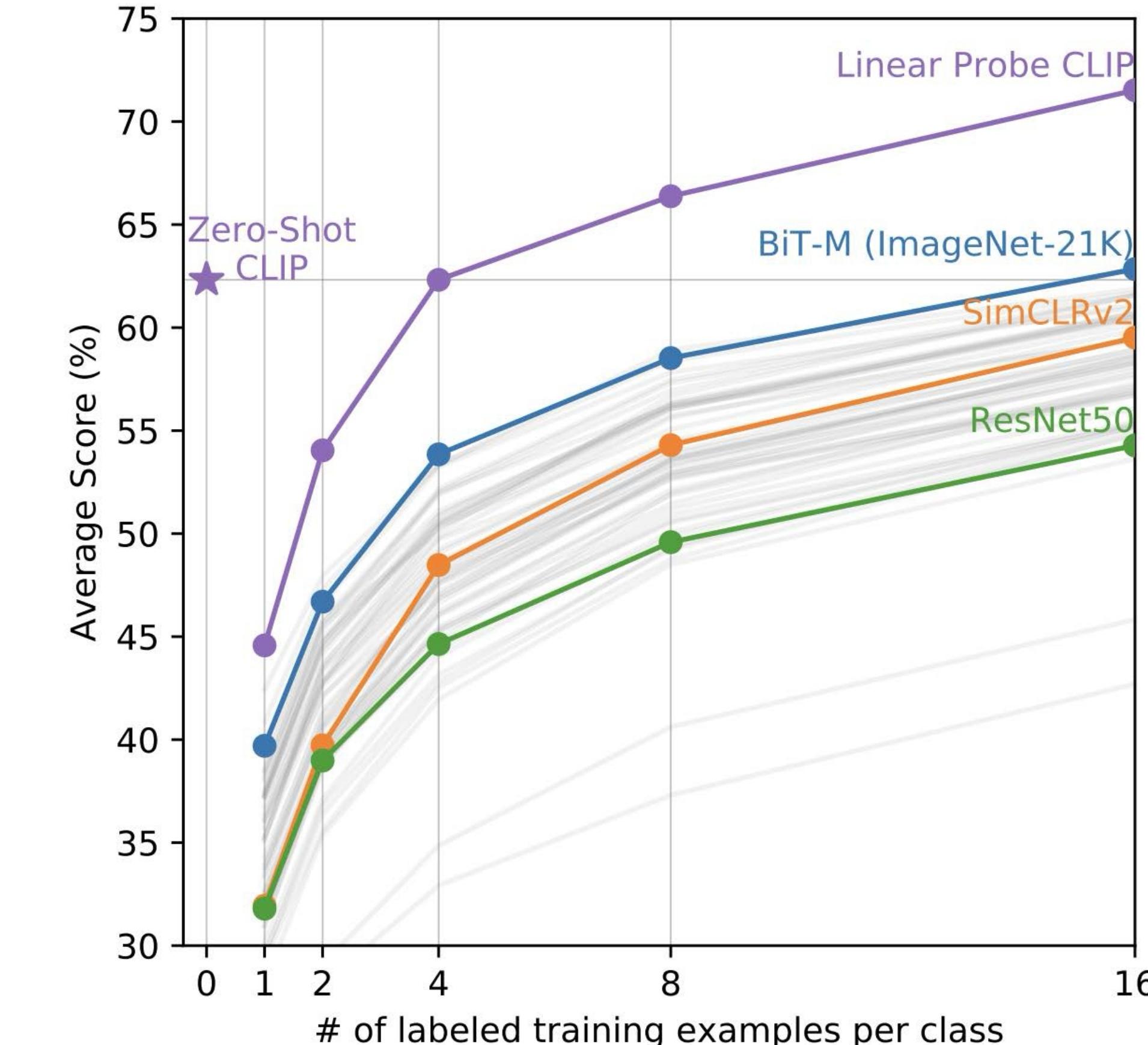


# Multi-Modal Contrastive Learning

## CLIP: Contrastive Language-Image Pre-Training

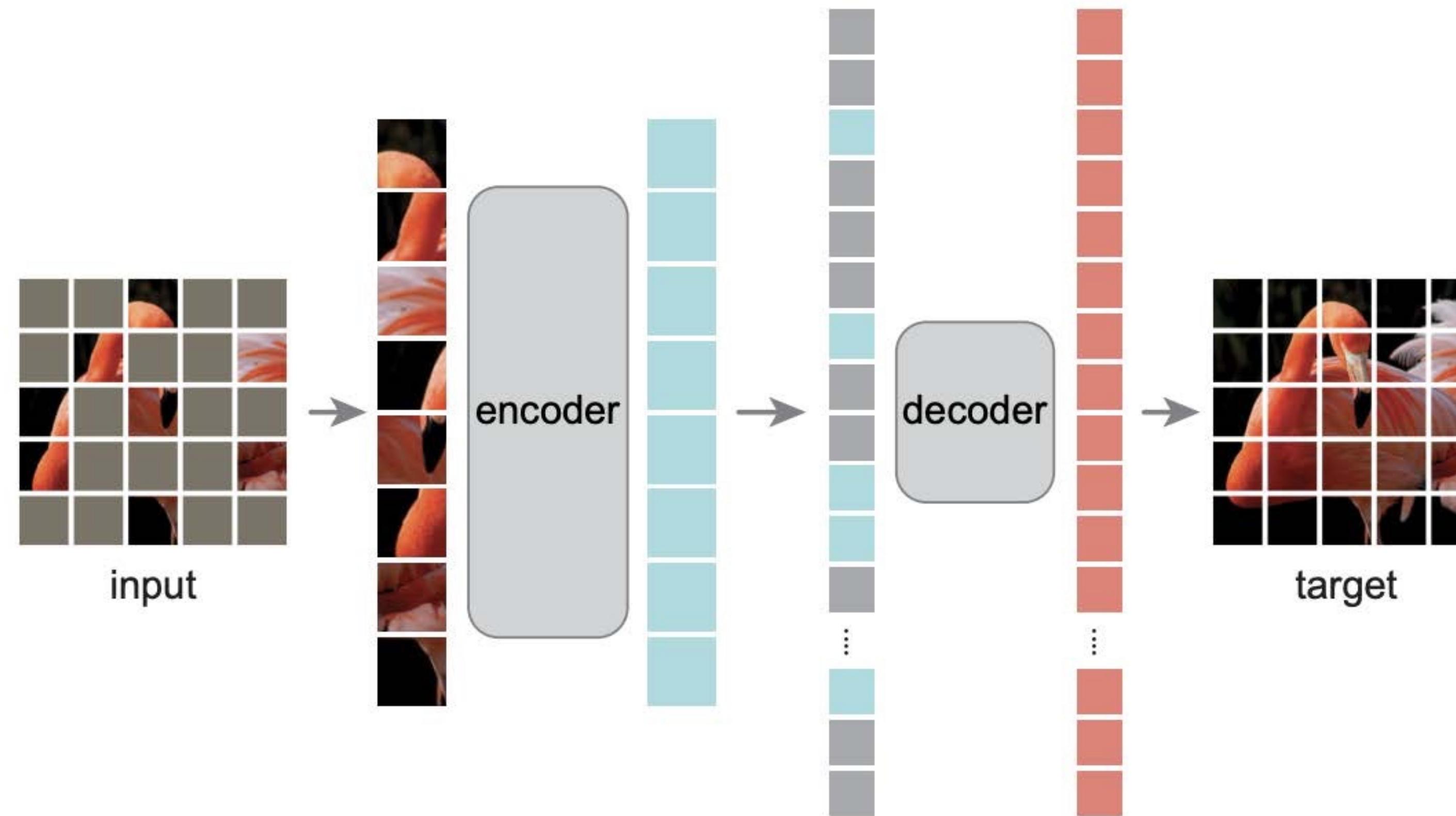
**Large-scale contrastive learning  
using paired image-text samples**

**Enable zero-shot and open-  
vocabulary classification (no fine-  
tuning, any #classes)**



# Masked Autoencoder

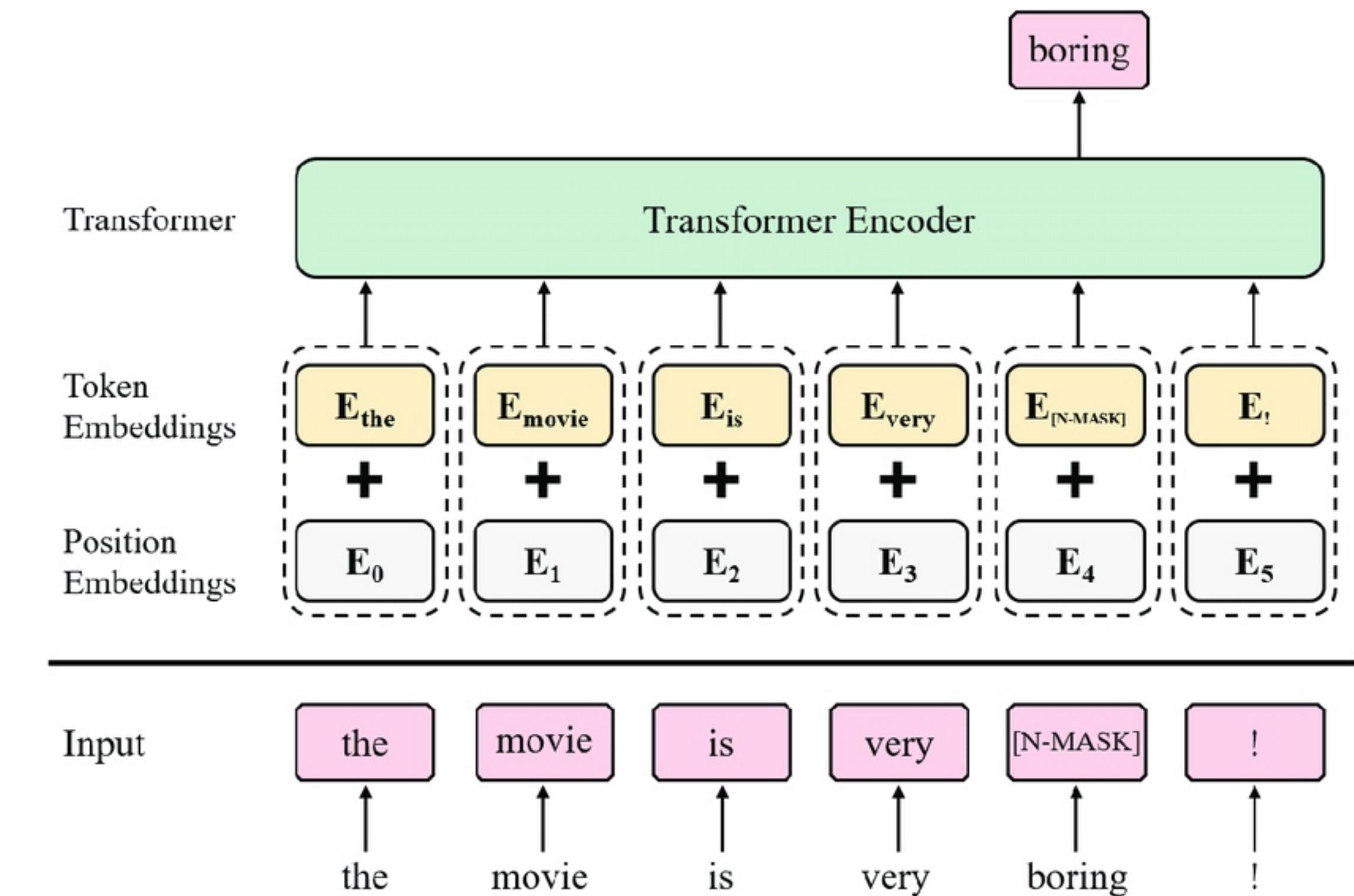
Randomly mask image patches, and predict them



# Review: Masked Language Models (MLM) — BERT

- Mask some percentage (15%) of the input tokens at random
- Predict those masked tokens

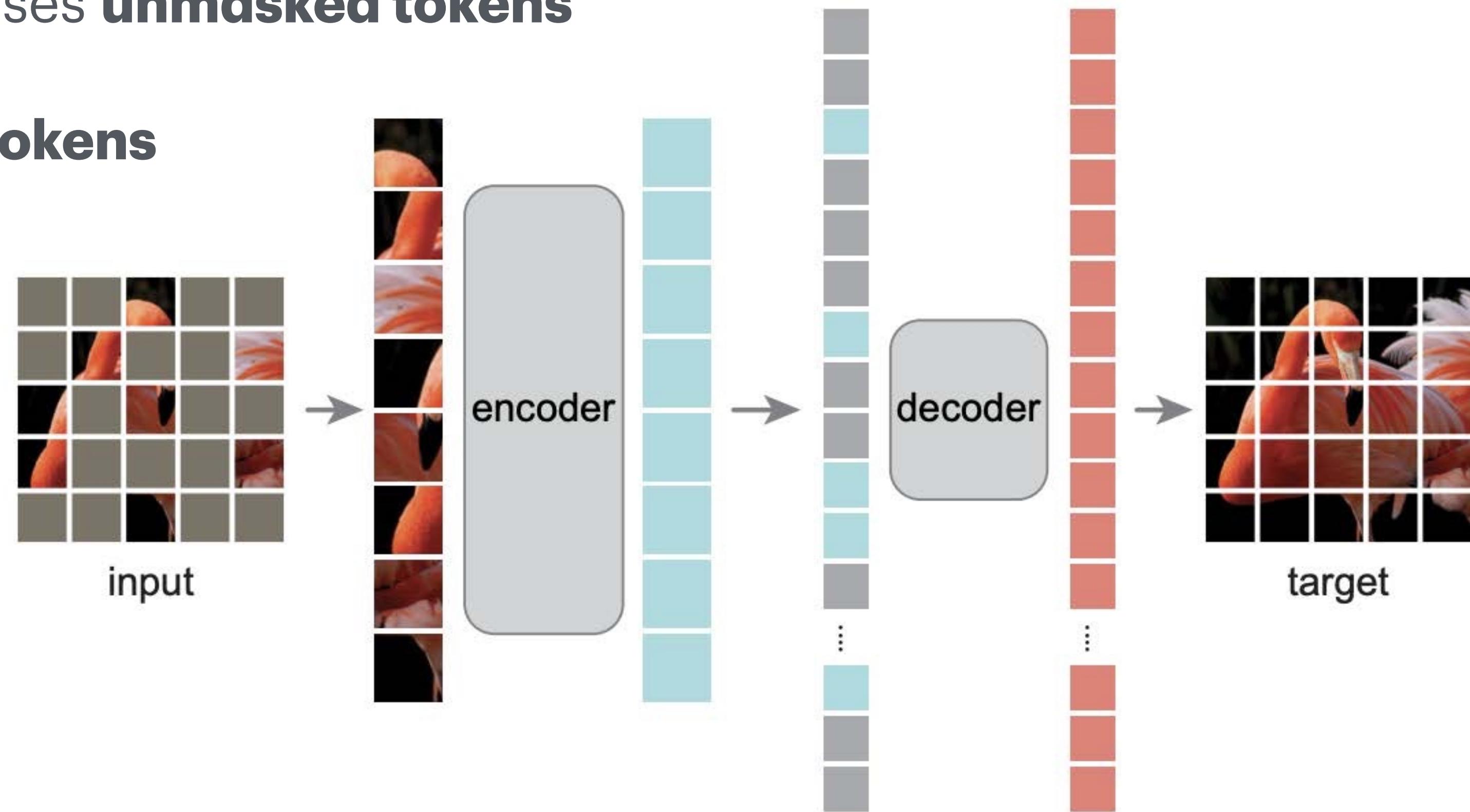
🤔 **Can we use a similar pre-training strategy in the vision domain?**



# Masked Autoencoder

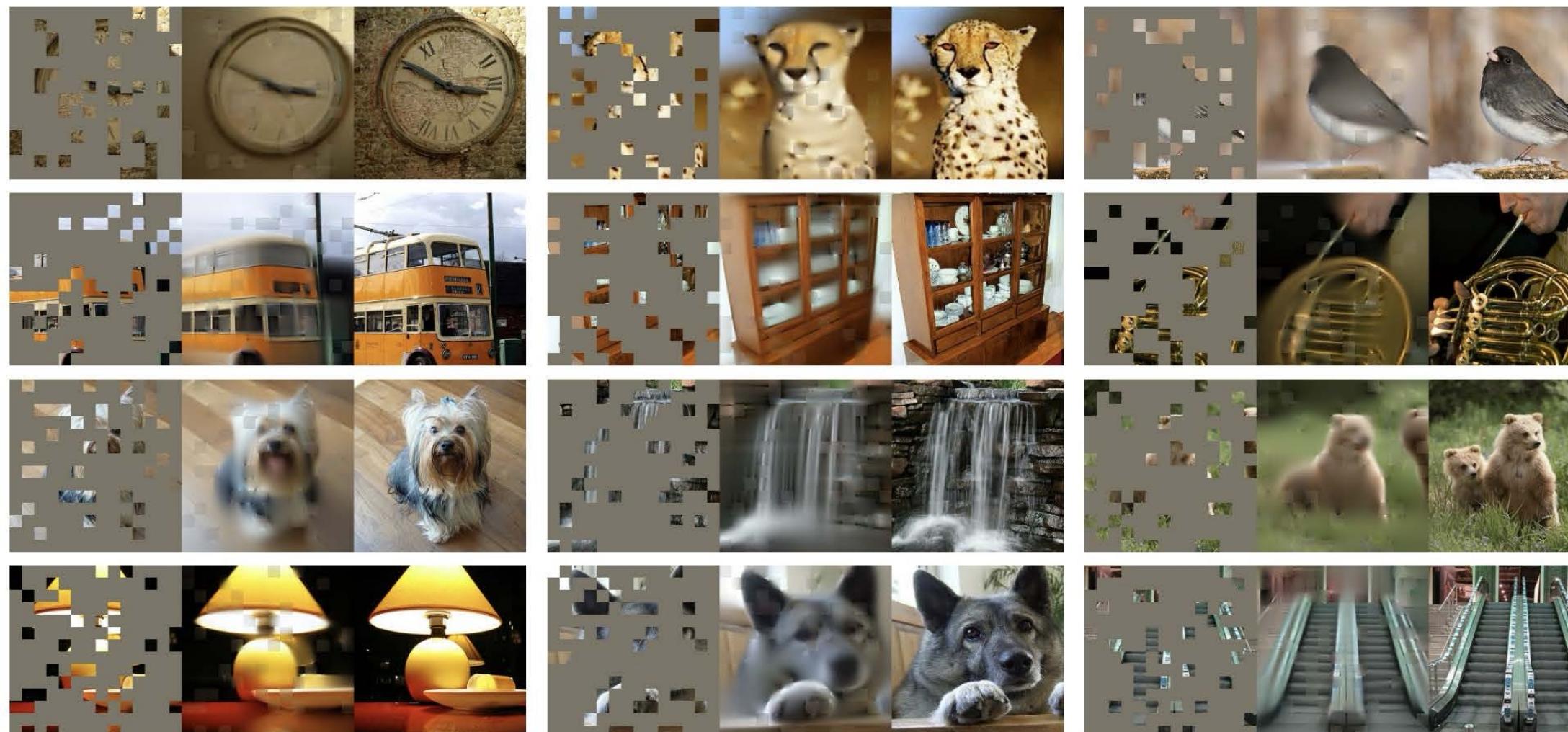
## Asymmetric Encoder-Decoder Design

- The **heavy** encoder only processes **unmasked tokens**
- The **lite** decoder processes **all tokens**

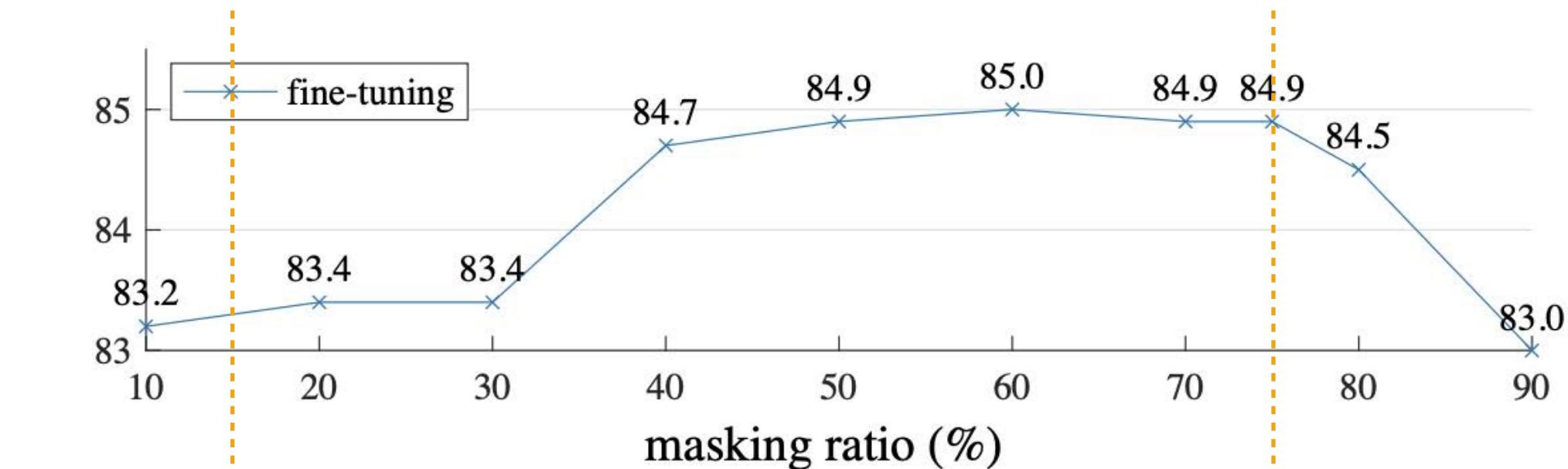


# Masked Autoencoder

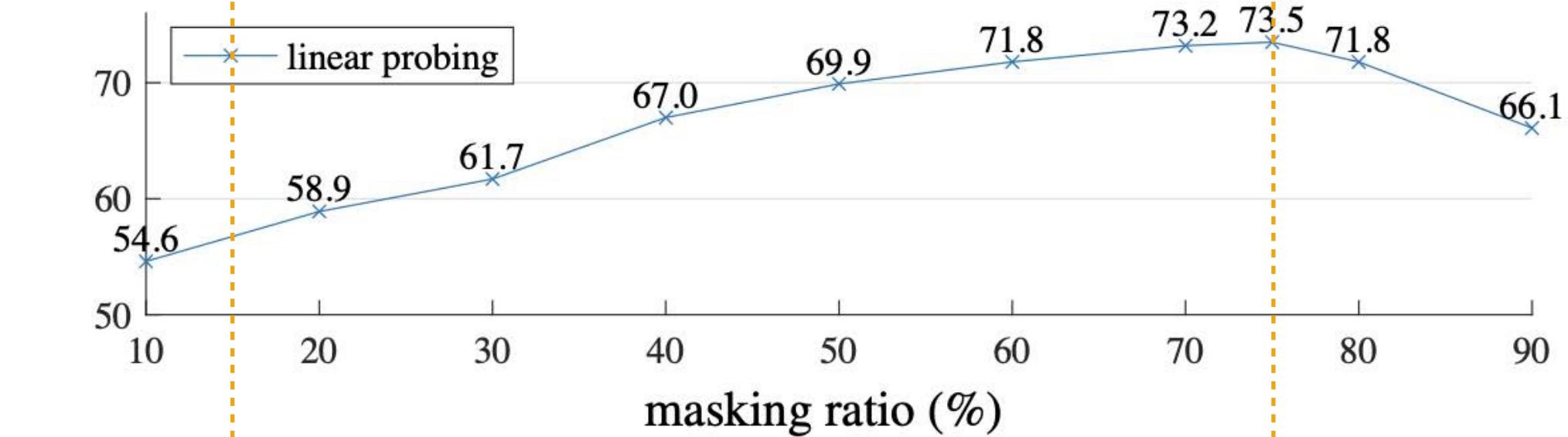
Higher mask ratio



BERT's mask ratio: 15%



MAE's mask ratio: 75%



- Use a higher mask ratio in vision than language: image data has lower information density than language data

# Reference

- Dosovitskiy, A. (2020). An image is worth 16x16 words: Transformers for image recognition at scale. arXiv preprint arXiv:2010.11929.
- Liu, Z., Lin, Y., Cao, Y., Hu, H., Wei, Y., Zhang, Z., ... & Guo, B. (2021). Swin transformer: Hierarchical vision transformer using shifted windows. In Proceedings of the IEEE/CVF international conference on computer vision (pp. 10012-10022).
- Cai, H., Li, J., Hu, M., Gan, C., & Han, S. (2023). EfficientViT: Lightweight multi-scale attention for on-device semantic segmentation. arXiv preprint arXiv:2205.14756, 2.
- Chen, X., Liu, Z., Tang, H., Yi, L., Zhao, H., & Han, S. (2023). Sparsevit: Revisiting activation sparsity for efficient high-resolution vision transformer. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (pp. 2061-2070).
- Noroozi, M., Vinjimoor, A., Favaro, P., & Pirsiavash, H. (2018). Boosting self-supervised learning via knowledge transfer. In Proceedings of the IEEE conference on computer vision and pattern recognition (pp. 9359-9367).
- Radford, A., Kim, J. W., Hallacy, C., Ramesh, A., Goh, G., Agarwal, S., ... & Sutskever, I. (2021, July). Learning transferable visual models from natural language supervision. In International conference on machine learning (pp. 8748-8763). PMLR.
- Devlin, J. (2018). Bert: Pre-training of deep bidirectional transformers for language understanding. arXiv preprint arXiv:1810.04805.
- He, K., Chen, X., Xie, S., Li, Y., Dollár, P., & Girshick, R. (2022). Masked autoencoders are scalable vision learners. In Proceedings of the IEEE/CVF conference on computer vision and pattern recognition (pp. 16000-16009).
- Vision Transformer [[MIT 6.5940](#)]