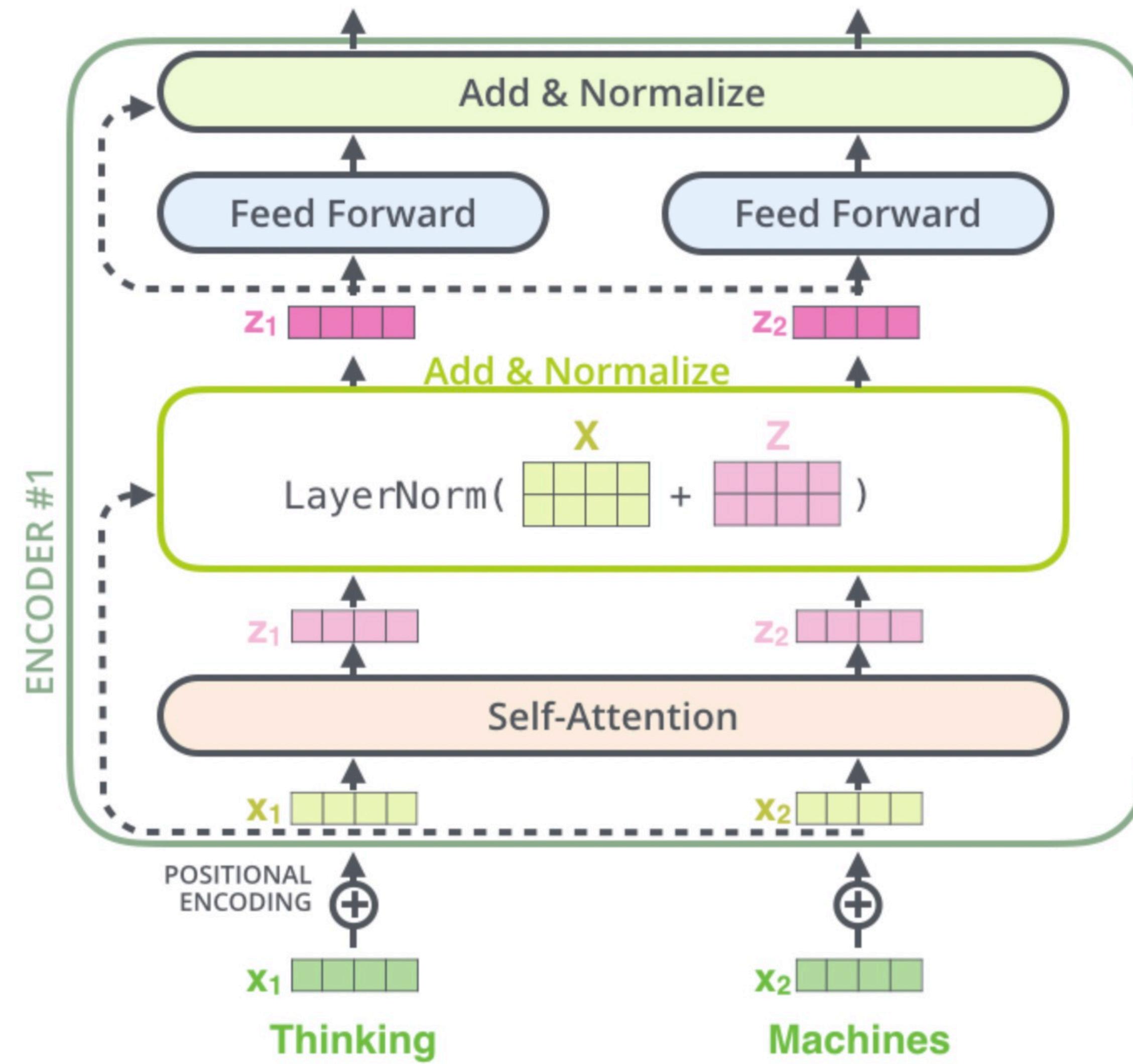


COMP5212: Machine Learning

Lecture 17

Minhao Cheng

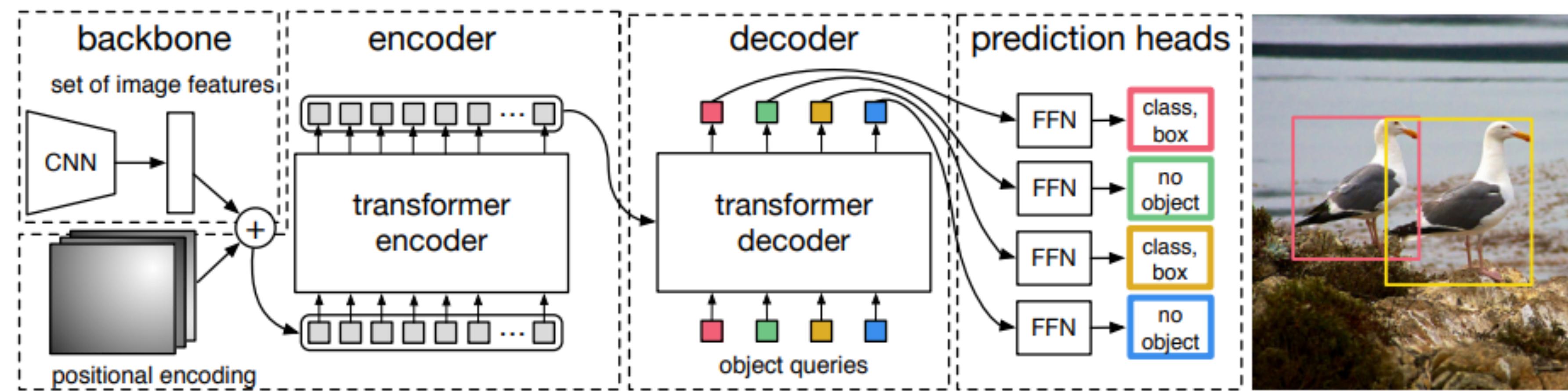
Transformer Overview



Vision Transformer (ViT)

Attempts on applying self-attention to vision

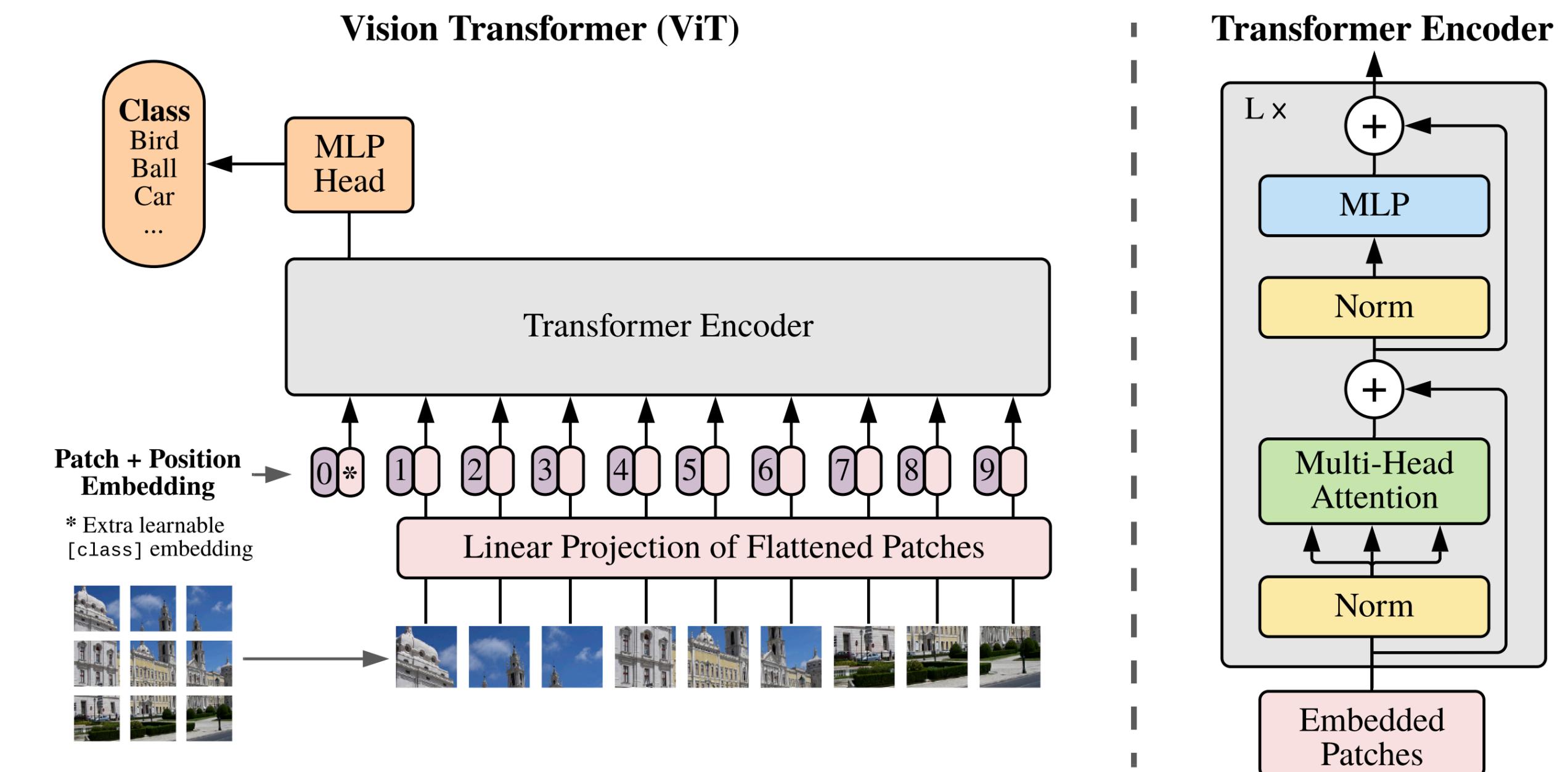
- DETR (Carion et al., 2020): CNN + Self-attention for object detection
- Stand-alone self-attention (Ramachandran et al., 2020)
- ...



Vision Transformer (ViT)

Vision Transformer (ViT)

- Partition input image into $K \times K$ patches
- A linear projection to transform each patch to feature (no convolution)
- Pass tokens into Transformer



Vision Transformer (ViT)

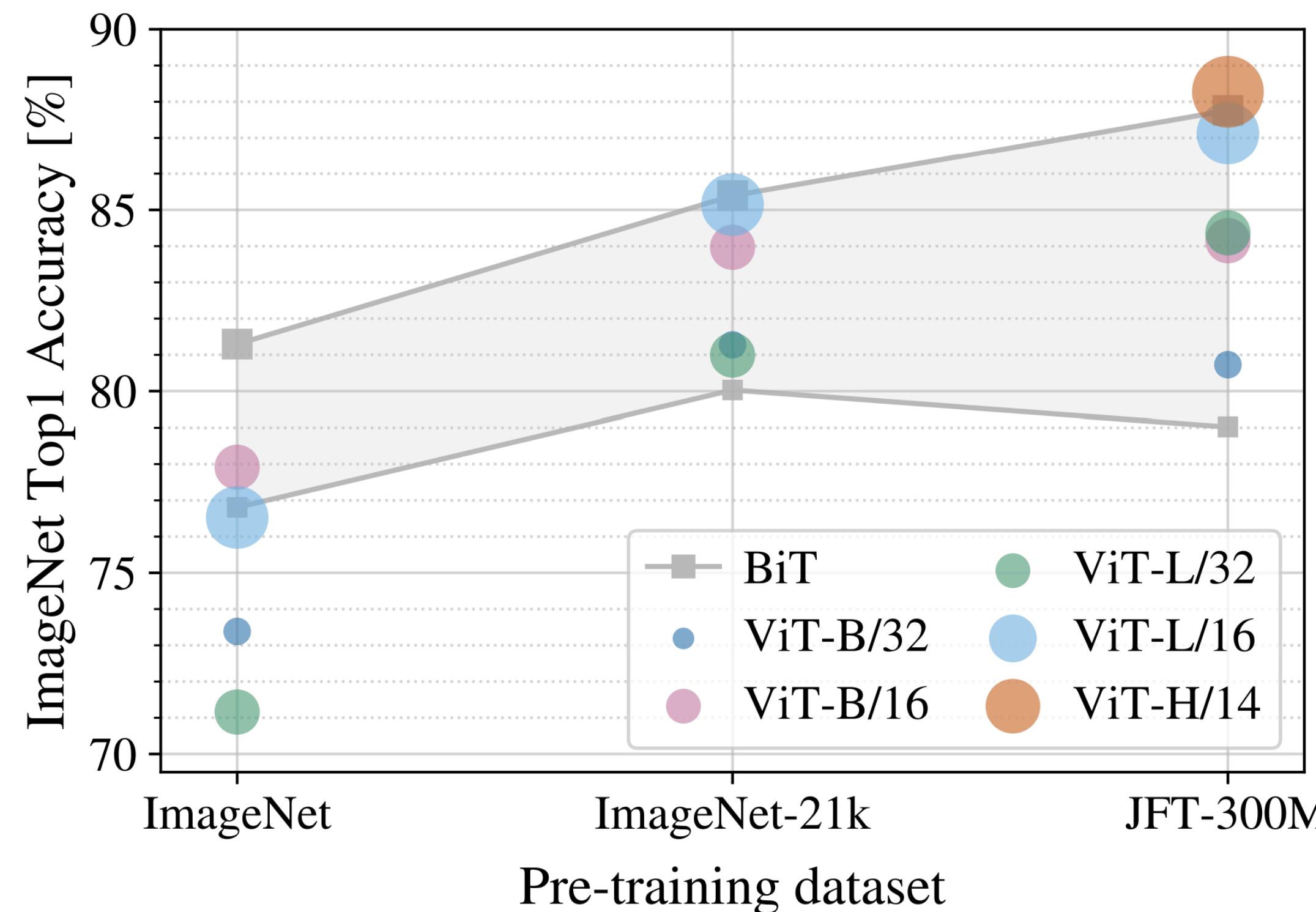
Vision Transformer (ViT)

- Patches are non-overlapping in the original ViT
- $N \times N$ image $\Rightarrow (N/K)^2$ tokens
- Smaller patch size \Rightarrow more input tokens
 - Higher computation (memory) cost, (usually) higher accuracy
- Use 1D (learnable) positional embedding
- Inference with higher resolution:
 - Keep the same patch size, which leads to longer sequence
 - Interpolation for positional embedding

Vision Transformer (ViT)

ViT Performance

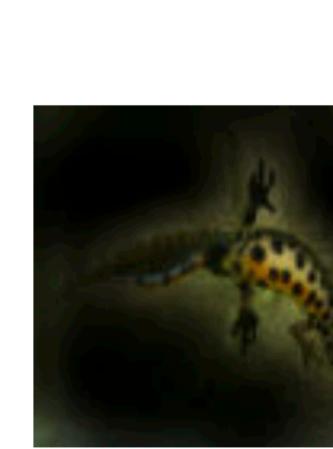
- ViT outperforms CNN with large pretraining



Vision Transformer (ViT)

ViT Performance

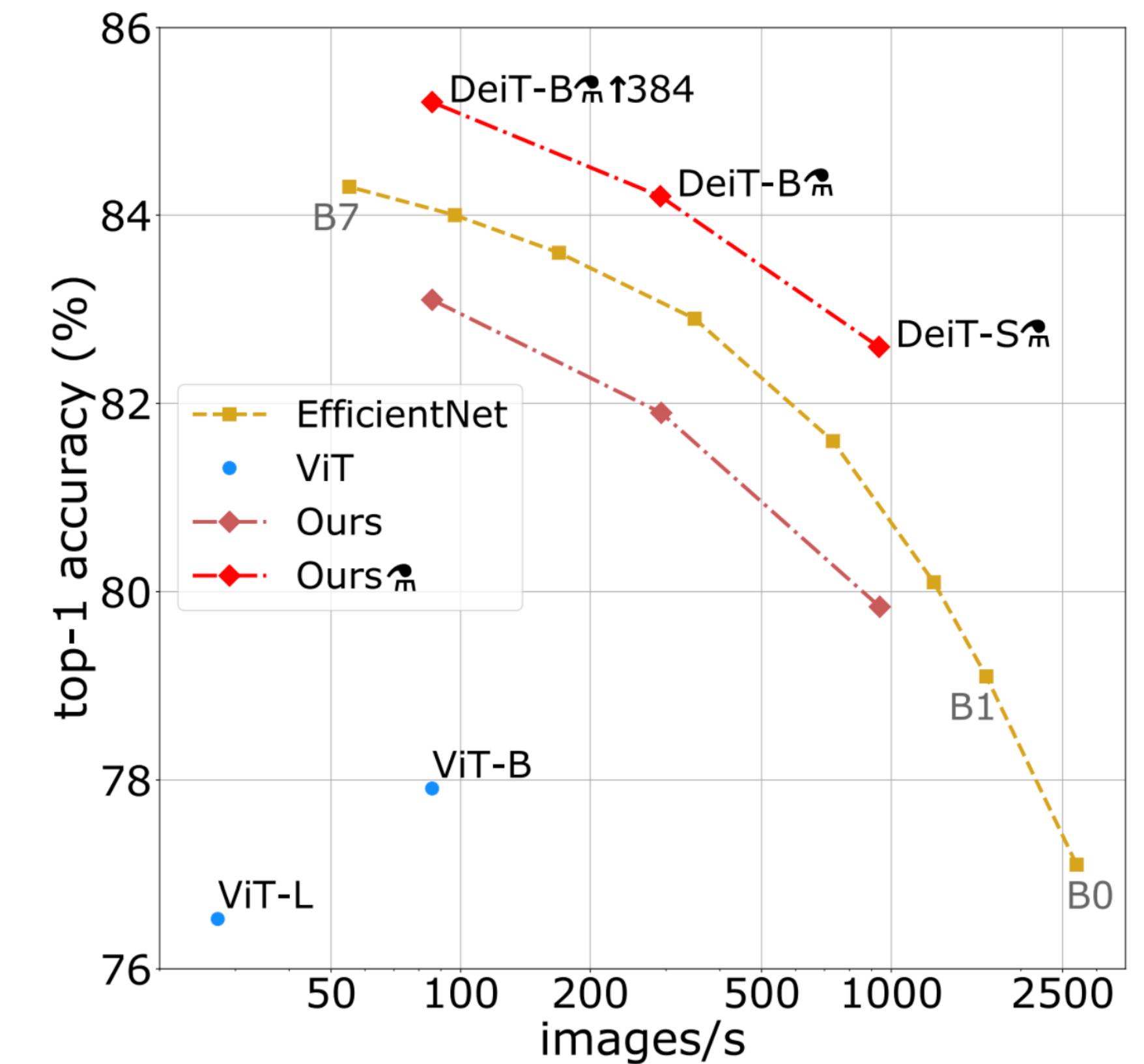
- Attention maps of ViT (to input)



Vision Transformer (ViT)

ViT v.s. ResNet

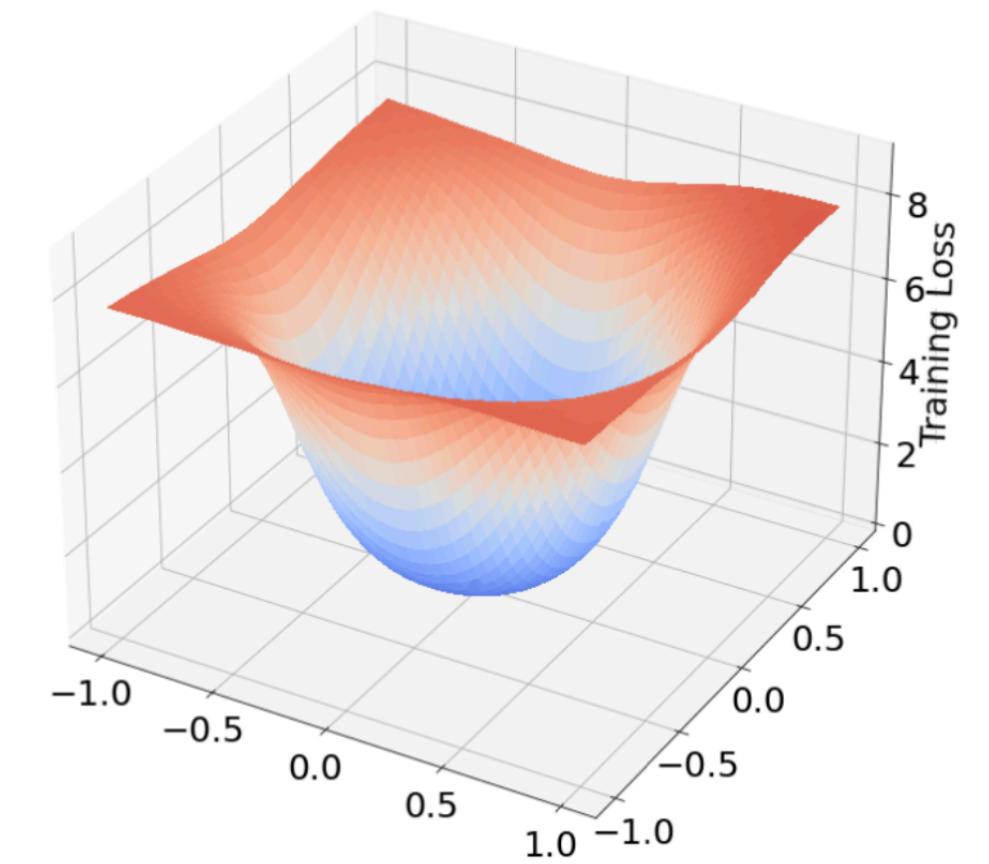
- Can ViT outperform ResNet on ImageNet without pretraining?
- DeiT (Touvron et al., 2021):
 - Use very strong data augmentation
 - Use a ResNet teacher and distill to ViT



Vision Transformer (ViT)

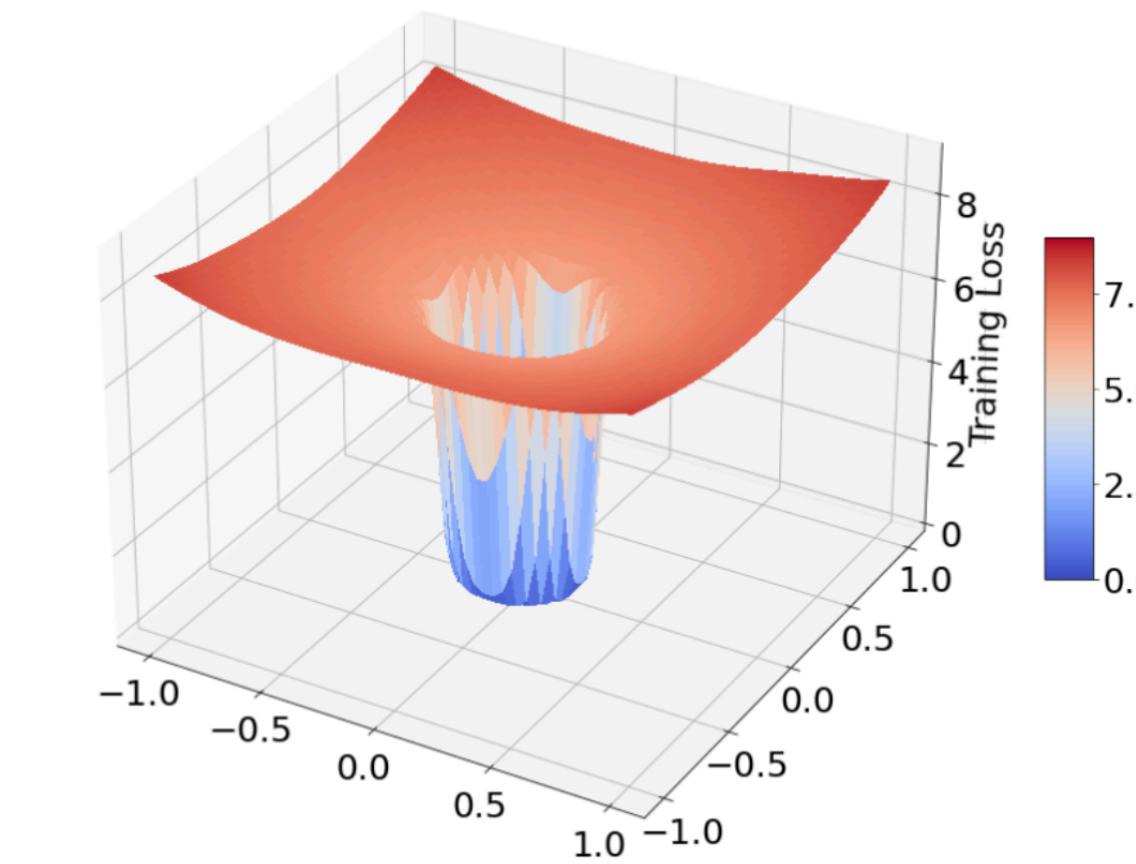
ViT v.s. ResNet

- ViT tends to converge to sharper regions than ResNet



(a) ResNet

Leading eigenvalue of
Hessian: **179.8**



(b) ViT

Leading eigenvalue of
Hessian: **738.8**

Vision Transformer (ViT)

“Sharpness” is related to generalization

- Testing can be viewed as a slightly perturbed training distribution
- Sharp minimum \Rightarrow performance degrades significantly from training to testing

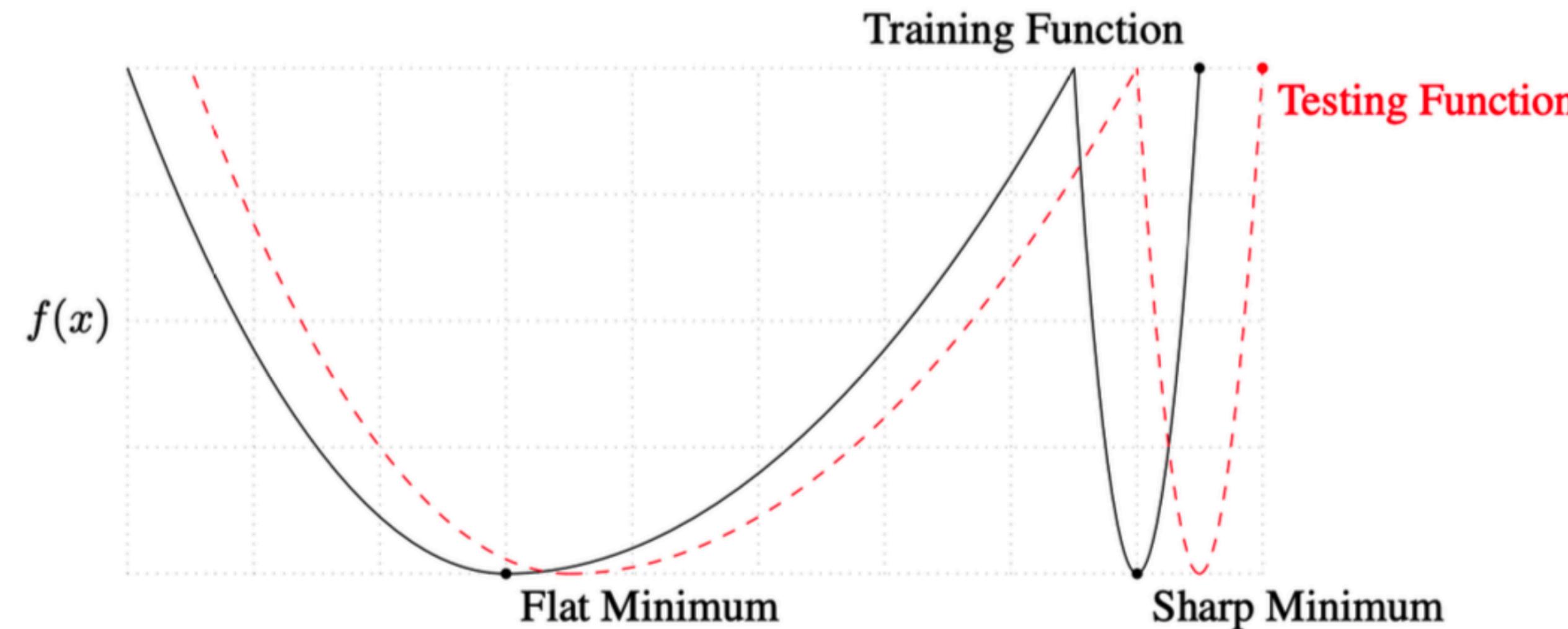


Figure from (Keskar et al., 2017)

Vision Transformer (ViT)

Sharpness Aware Minimization (SAM)

- Optimize the worst-case loss within a small neighborhood

- $\min_w \max_{\|\delta\|_2 \leq \epsilon} L(w + \delta)$

- ϵ is a small constant (hyper-parameter)

- Use 1-step gradient ascent to approximate inner max:

- $\hat{\delta} = \arg \max_{\|\delta\|_2 \leq \epsilon} L(w) + \nabla L(w)^T \delta = \epsilon \frac{\nabla L(w)}{\|\nabla L(w)\|}$

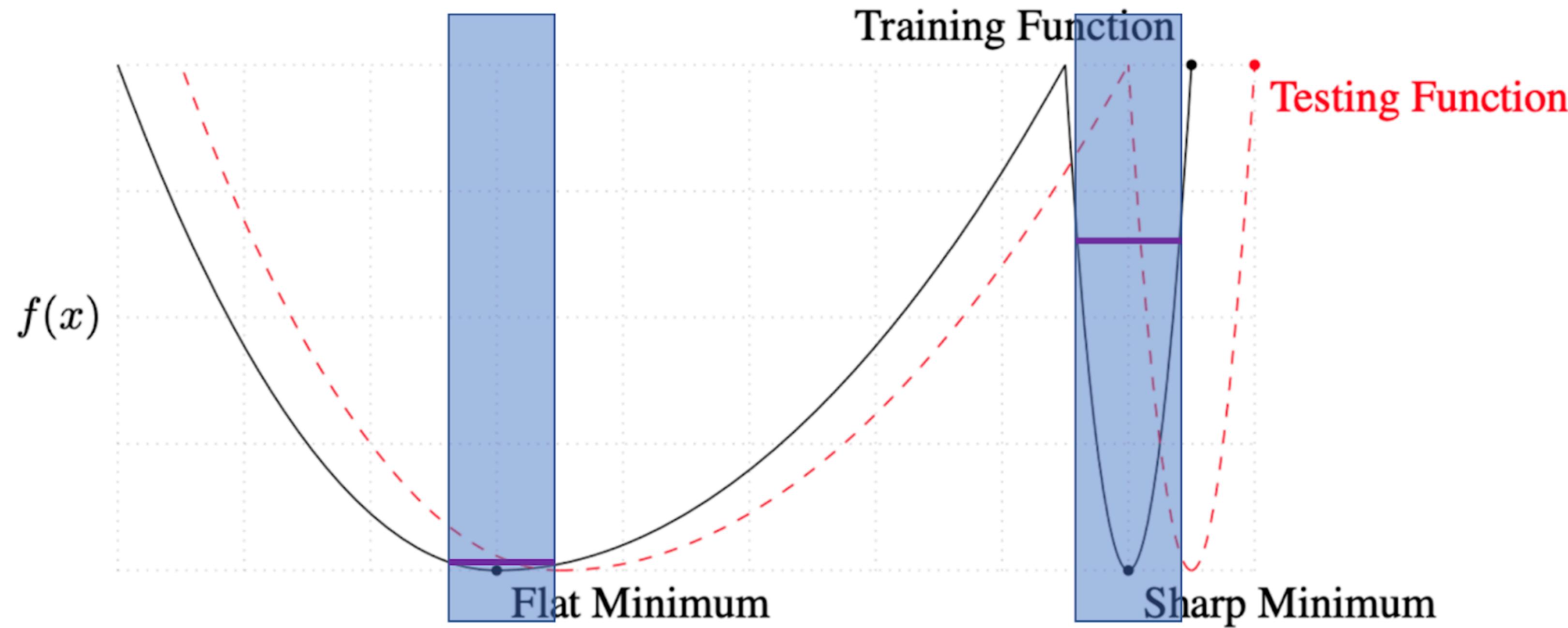
- Conduct the following update for each iteration:

- $w \leftarrow w - \alpha \nabla L(w + \hat{\delta})$

Vision Transformer (ViT)

Sharpness Aware Minimization (SAM)

- SAM is a natural way to penalize sharpness region (but requires some computational overhead)



Vision Transformer (ViT)

ViT v.s. ResNet

Model	#params	Throughput (img/sec/core)	ImageNet	Real	V2	ImageNet-R	ImageNet-C
ResNet							
ResNet-50-SAM	25M	2161	76.7 (+0.7)	83.1 (+0.7)	64.6 (+1.0)	23.3 (+1.1)	46.5 (+1.9)
ResNet-101-SAM	44M	1334	78.6 (+0.8)	84.8 (+0.9)	66.7 (+1.4)	25.9 (+1.5)	51.3 (+2.8)
ResNet-152-SAM	60M	935	79.3 (+0.8)	84.9 (+0.7)	67.3 (+1.0)	25.7 (+0.4)	52.2 (+2.2)
ResNet-50x2-SAM	98M	891	79.6 (+1.5)	85.3 (+1.6)	67.5 (+1.7)	26.0 (+2.9)	50.7 (+3.9)
ResNet-101x2-SAM	173M	519	80.9 (+2.4)	86.4 (+2.4)	69.1 (+2.8)	27.8 (+3.2)	54.0 (+4.7)
ResNet-152x2-SAM	236M	356	81.1 (+1.8)	86.4 (+1.9)	69.6 (+2.3)	28.1 (+2.8)	55.0 (+4.2)
Vision Transformer							
ViT-S/32-SAM	23M	6888	70.5 (+2.1)	77.5 (+2.3)	56.9 (+2.6)	21.4 (+2.4)	46.2 (+2.9)
ViT-S/16-SAM	22M	2043	78.1 (+3.7)	84.1 (+3.7)	65.6 (+3.9)	24.7 (+4.7)	53.0 (+6.5)
ViT-S/14-SAM	22M	1234	78.8 (+4.0)	84.8 (+4.5)	67.2 (+5.2)	24.4 (+4.7)	54.2 (+7.0)
ViT-S/8-SAM	22M	333	81.3 (+5.3)	86.7 (+5.5)	70.4 (+6.2)	25.3 (+6.1)	55.6 (+8.5)
ViT-B/32-SAM	88M	2805	73.6 (+4.1)	80.3 (+5.1)	60.0 (+4.7)	24.0 (+4.1)	50.7 (+6.7)
ViT-B/16-SAM	87M	863	79.9 (+5.3)	85.2 (+5.4)	67.5 (+6.2)	26.4 (+6.3)	56.5 (+9.9)
MLP-Mixer							
Mixer-S/32-SAM	19M	11401	66.7 (+2.8)	73.8 (+3.5)	52.4 (+2.9)	18.6 (+2.7)	39.3 (+4.1)
Mixer-S/16-SAM	18M	4005	72.9 (+4.1)	79.8 (+4.7)	58.9 (+4.1)	20.1 (+4.2)	42.0 (+6.4)
Mixer-S/8-SAM	20M	1498	75.9 (+5.7)	82.5 (+6.3)	62.3 (+6.2)	20.5 (+5.1)	42.4 (+7.8)
Mixer-B/32-SAM	60M	4209	72.4 (+9.9)	79.0 (+10.9)	58.0 (+10.4)	22.8 (+8.2)	46.2 (12.4)
Mixer-B/16-SAM	59M	1390	77.4 (+11.0)	83.5 (+11.4)	63.9 (+13.1)	24.7 (+10.2)	48.8 (+15.0)
Mixer-B/8-SAM	64M	466	79.0 (+10.4)	84.4 (+10.1)	65.5 (+11.6)	23.5 (+9.2)	48.9 (+16.9)

Vision Transformer (ViT)

ViT v.s. ResNet

- Let's compare one ViT layer vs one convolution layer
- Reception field: (which input neurons can affect an output neuron)
 - CNN: some subarea of image (kernel size)
 - Self-attention: the whole image
 - \Rightarrow there exists self-attention function that cannot be captured by convolution
- Is the function set of self-attention strictly larger than convolution?
 - Yes, given enough attention heads

Vision Transformer (ViT)

How can self-attention do convolution?

- Consider self-attention with relative positional encoding

$$\text{Output} = \text{Softmax}\left(\underbrace{\frac{QK^T}{\sqrt{d}}}_{\begin{array}{c} \text{context aware} \\ \cdot \end{array}} + \underbrace{B}_{\begin{array}{c} \text{context agnostic} \\ \cdot \end{array}}\right)V$$

- Q, K, V : query, key, value matrices
- $B_{i,j} = b_{(x_i-x_j, y_i-y_j)}$: relative positional encoding (trainable scalars)
- To perform convolution: Set $Q, K = 0$ and purely rely on B
- Implication: the positional encoding can capture CNN; the query/key matrices can capture context-aware information beyond convolution

Vision Transformer (ViT)

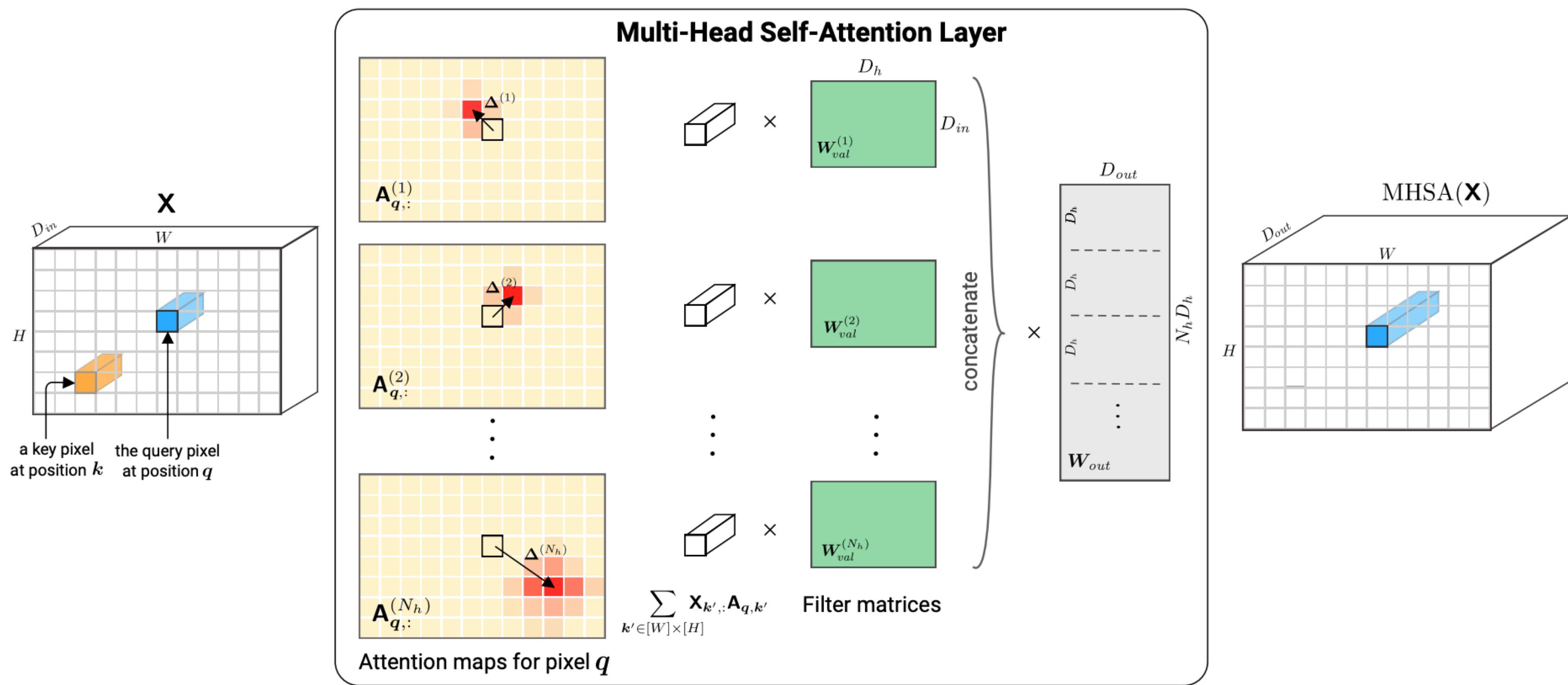
How can self-attention do convolution?

- (Cordonnier et al., 2020): If each input token is a **pixel**, a self-attention layer with K^2 heads can express any convolution with kernel $K \times K$
- If each input token is a **$P \times P$ patch**, then self-attention with $(2\lceil\frac{K-1}{2P}\rceil + 1)^2$ heads can express any convolution with kernel $K \times K$
- In patch setting, when $P > 2K$ (e.g., $P = 16$ in ViT), then 9 heads are enough for self-attention to express convolution
- In both settings, the number of heads above are also **necessary** for self-attention to perform convolution

Vision Transformer (ViT)

How can self-attention do convolution?

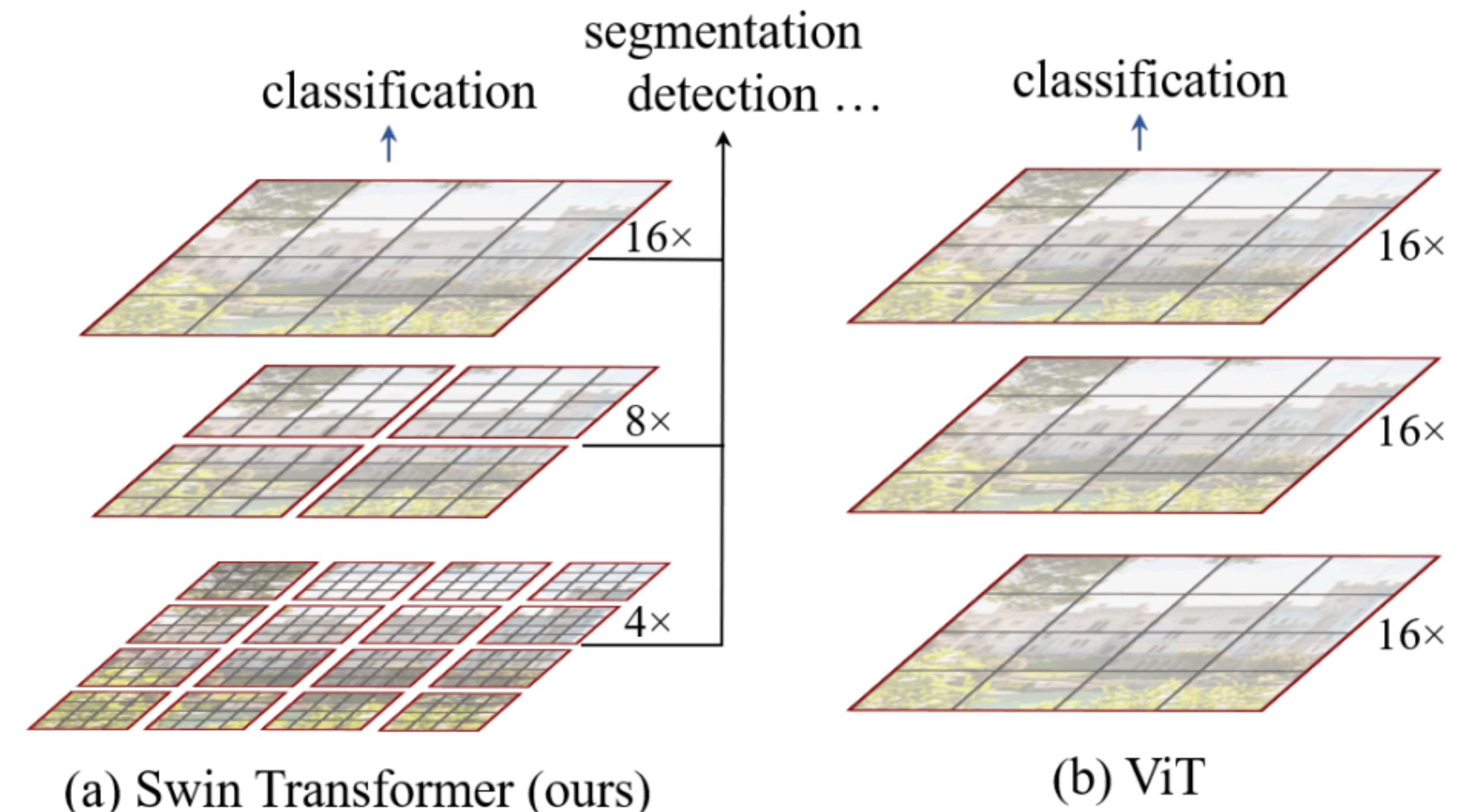
- Construct convolution using self-attention heads



Vision Transformer (ViT)

Swin Transformer (Liu et al., 2021)

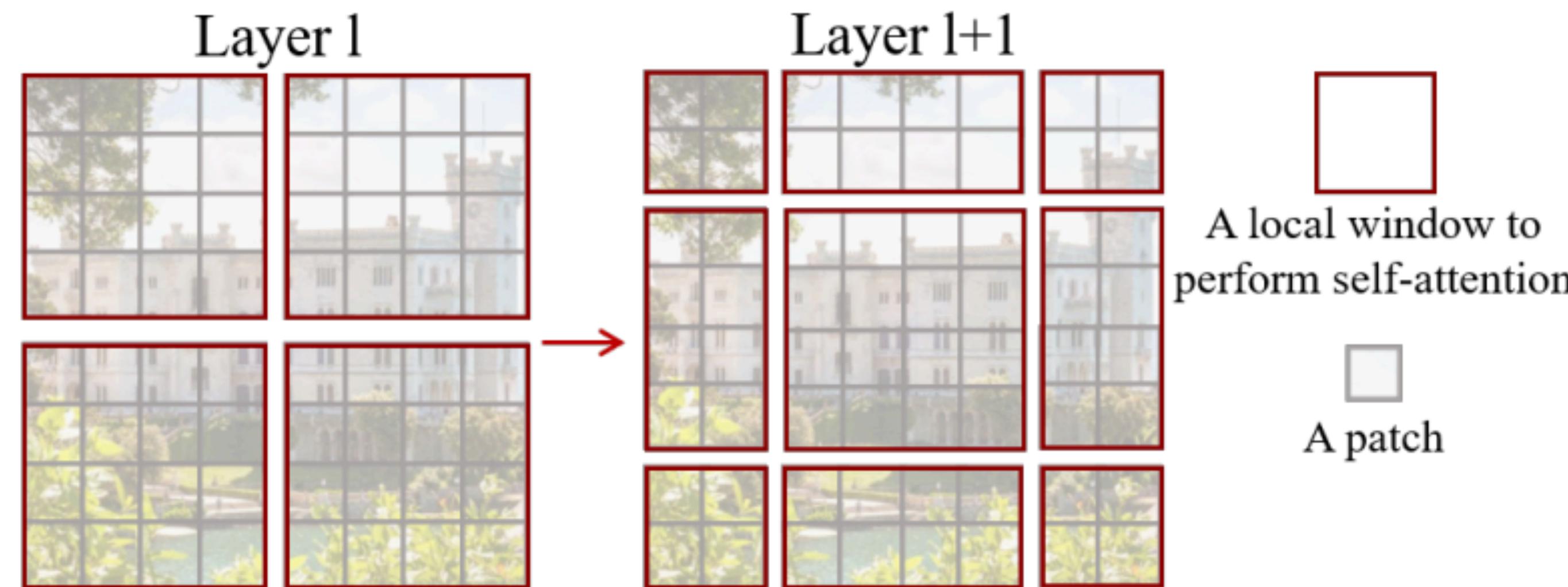
- Problems of the original ViT:
 - Non-overlapping partition
 - Only a single resolution
 - Quadratic complexity for attention computation
- Swin Transformer: hierarchical and sliding window partitions



Vision Transformer (ViT)

Swin Transformer

- Attention within sub-blocks with shifts to avoid huge attention matrix



Vision Transformer (ViT)

Swin Transformer

(a) Regular ImageNet-1K trained models					
method	image size	#param.	FLOPs	throughput (image / s)	ImageNet top-1 acc.
RegNetY-4G [48]	224 ²	21M	4.0G	1156.7	80.0
RegNetY-8G [48]	224 ²	39M	8.0G	591.6	81.7
RegNetY-16G [48]	224 ²	84M	16.0G	334.7	82.9
EffNet-B3 [58]	300 ²	12M	1.8G	732.1	81.6
EffNet-B4 [58]	380 ²	19M	4.2G	349.4	82.9
EffNet-B5 [58]	456 ²	30M	9.9G	169.1	83.6
EffNet-B6 [58]	528 ²	43M	19.0G	96.9	84.0
EffNet-B7 [58]	600 ²	66M	37.0G	55.1	84.3
ViT-B/16 [20]	384 ²	86M	55.4G	85.9	77.9
ViT-L/16 [20]	384 ²	307M	190.7G	27.3	76.5
DeiT-S [63]	224 ²	22M	4.6G	940.4	79.8
DeiT-B [63]	224 ²	86M	17.5G	292.3	81.8
DeiT-B [63]	384 ²	86M	55.4G	85.9	83.1
Swin-T	224 ²	29M	4.5G	755.2	81.3
Swin-S	224 ²	50M	8.7G	436.9	83.0
Swin-B	224 ²	88M	15.4G	278.1	83.5
Swin-B	384 ²	88M	47.0G	84.7	84.5

(b) ImageNet-22K pre-trained models					
method	image size	#param.	FLOPs	throughput (image / s)	ImageNet top-1 acc.
R-101x3 [38]	384 ²	388M	204.6G	-	84.4
R-152x4 [38]	480 ²	937M	840.5G	-	85.4
ViT-B/16 [20]	384 ²	86M	55.4G	85.9	84.0
ViT-L/16 [20]	384 ²	307M	190.7G	27.3	85.2
Swin-B	224 ²	88M	15.4G	278.1	85.2
Swin-B	384 ²	88M	47.0G	84.7	86.4
Swin-L	384 ²	197M	103.9G	42.1	87.3