



Transformers II

CS 4391 Introduction Computer Vision

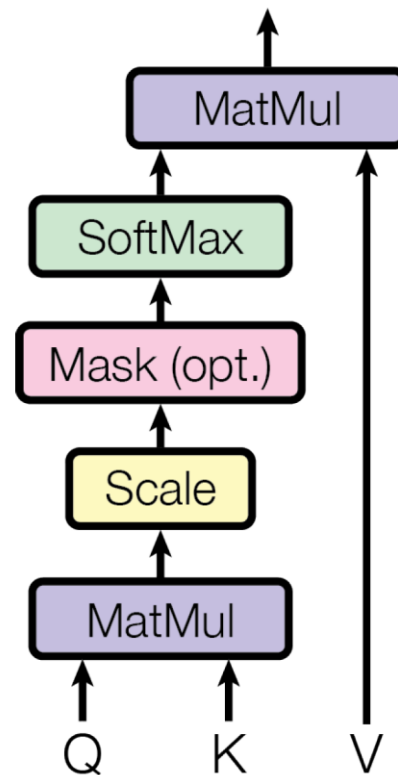
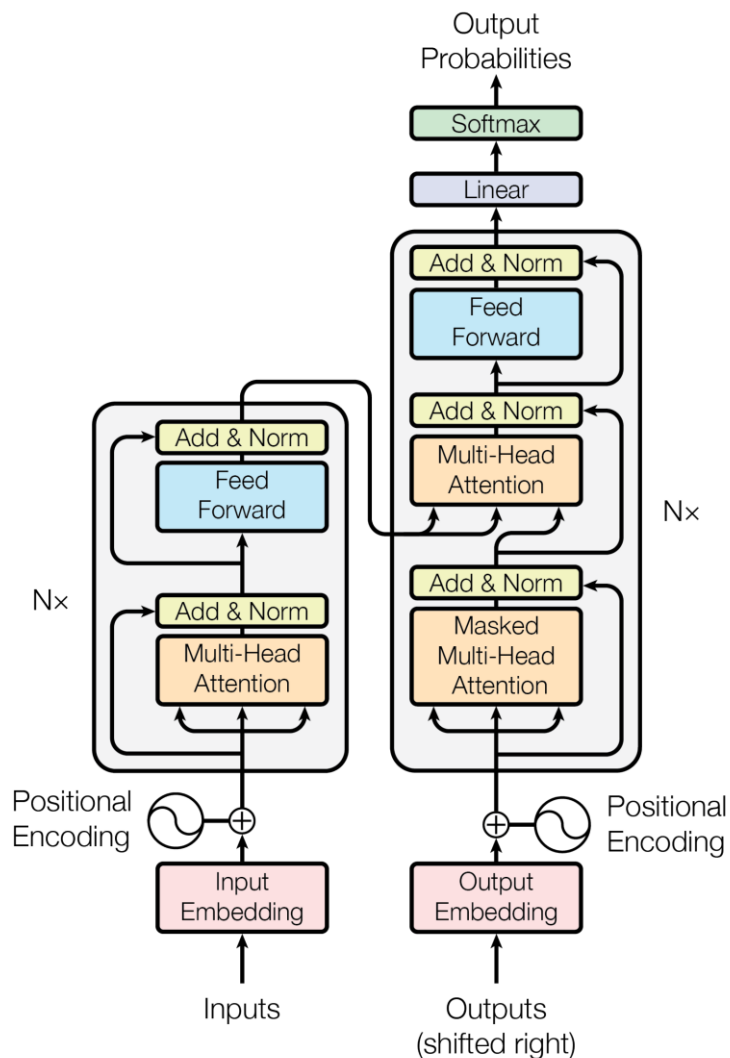
Professor Yu Xiang

The University of Texas at Dallas

Transformer

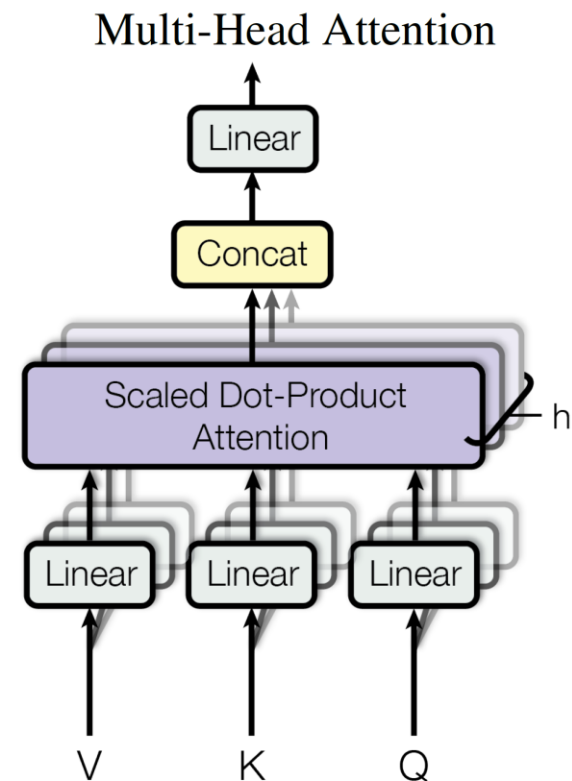
- No recurrence
- Attention only
 - Global dependencies between input and output
 - More parallelization compared to RNNs

Transformer



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

Attention is all you need. Vaswani et al., NeurIPS'17



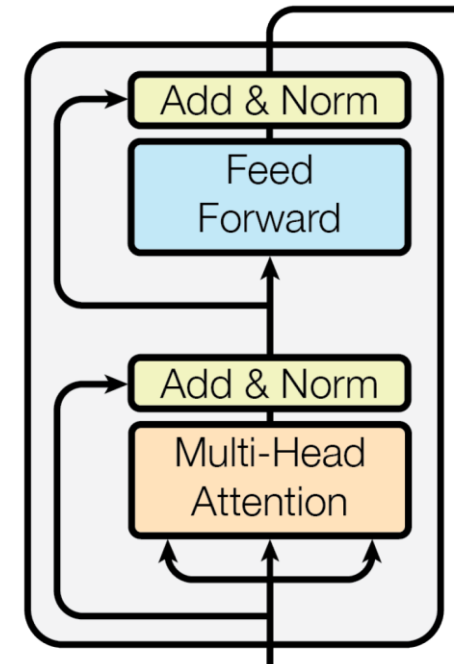
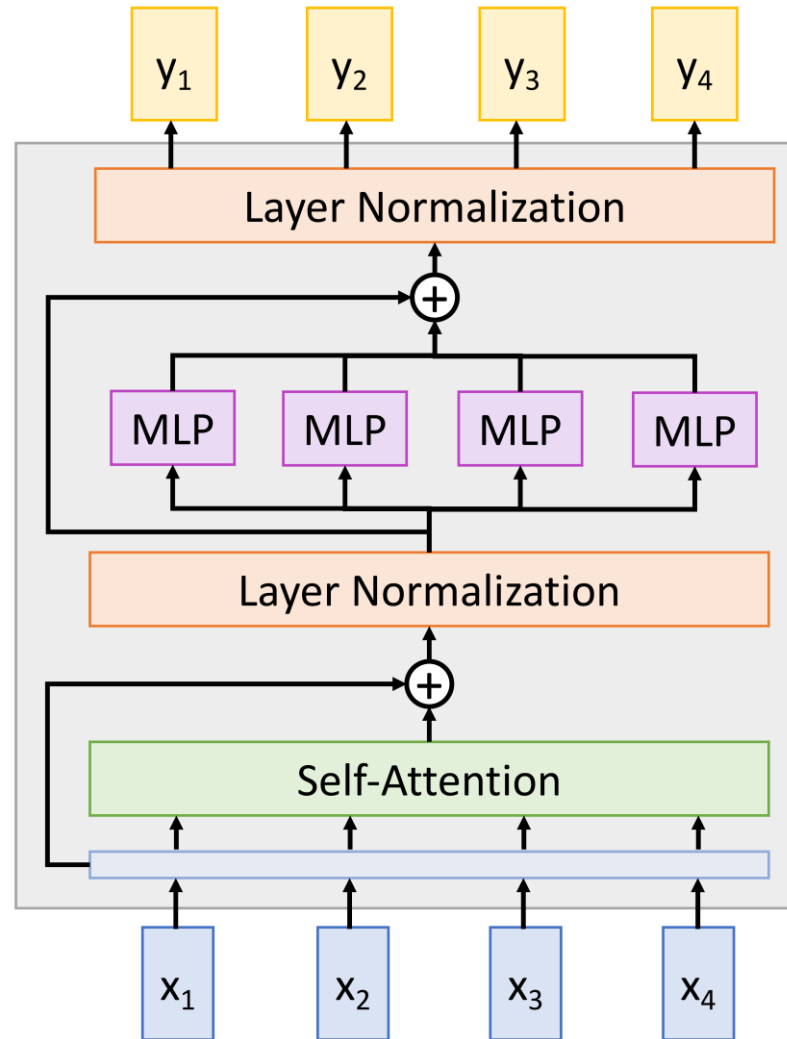
$$\text{head}_i = \text{Attention}(QW_i^Q, KW_i^K, VW_i^V)$$

$$\text{MultiHead}(Q, K, V)$$

$$= \text{Concat}(\text{head}_1, \dots, \text{head}_h)W^O$$

Transformer

- Transformer block
 - Input: a set of vectors
 $n \times d_{\text{model}}$
 - Output: a set of vectors
 $n \times d_{\text{model}}$

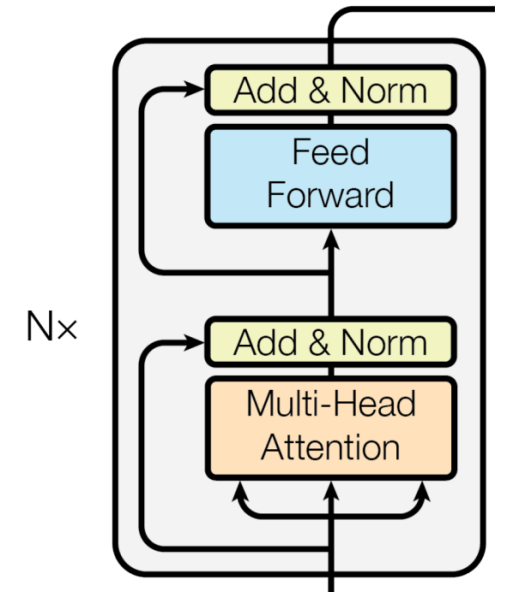
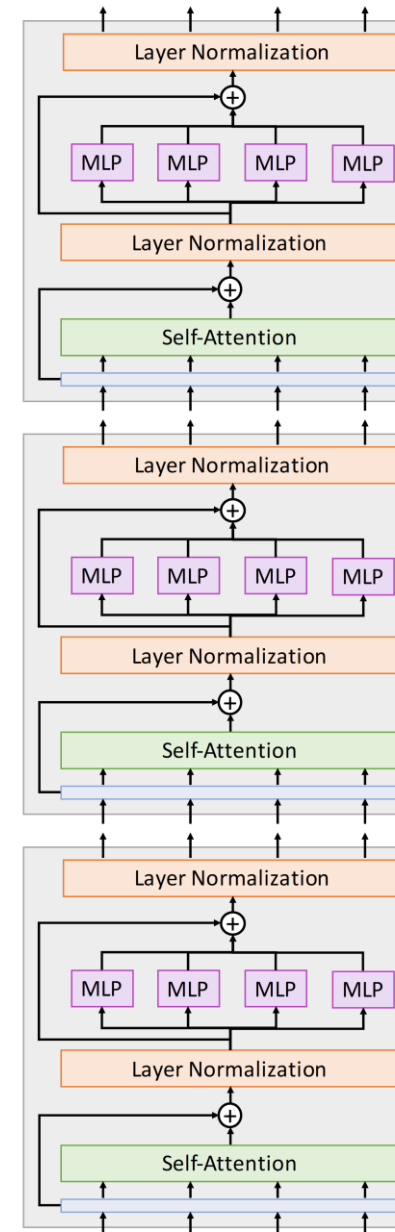


Transformer

- Hyper-parameters
 - Number of blocks
 - Number of heads per block

$$\text{MultiHead}(Q, K, V) = \text{Concat}(\text{head}_1, \dots, \text{head}_h)W^O$$

- Width (channels per head, FFN width)



Vision Transformer (ViT)

- Convert an image into a sequence of “token”



- Input embedding by linear projection

$$\mathbf{x}_p^1 \mathbf{E}; \mathbf{x}_p^2 \mathbf{E}; \dots; \mathbf{x}_p^N \mathbf{E}$$

$$\mathbf{E} \in \mathbb{R}^{(P^2 \cdot C) \times D}$$

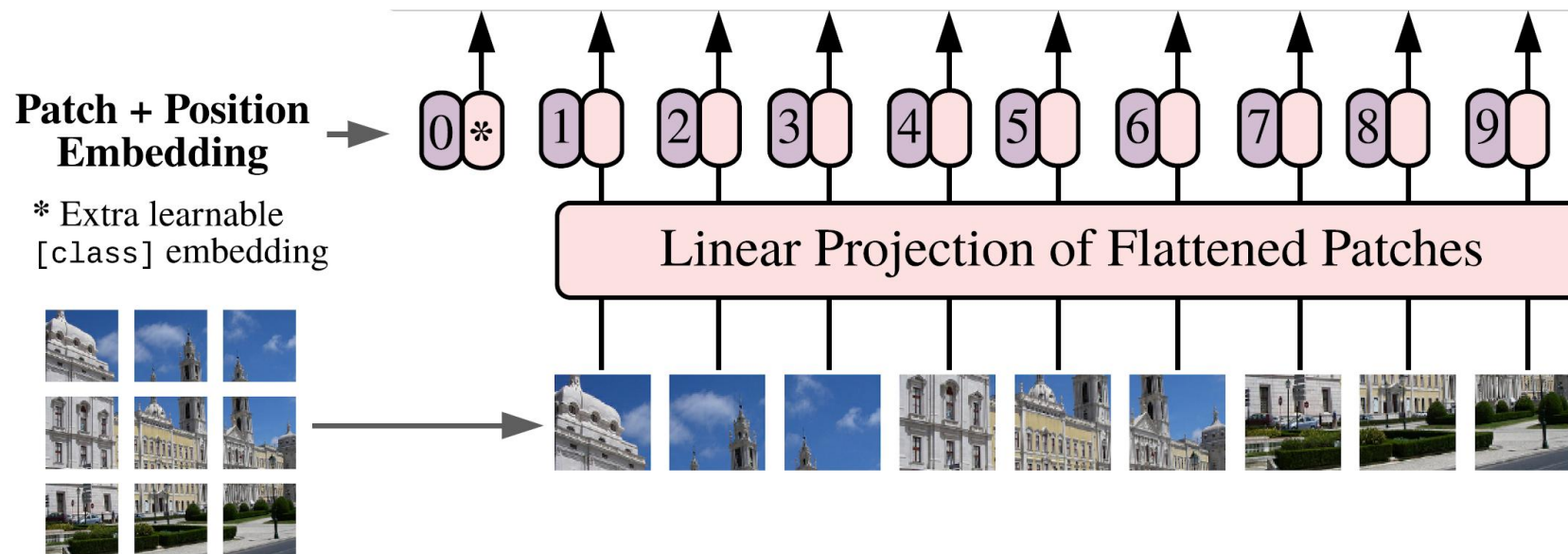
d_{model}

AN IMAGE IS WORTH 16x16 WORDS: TRANSFORMERS FOR IMAGE RECOGNITION AT SCALE. Dosovitskiy et al., ICLR'21

Vision Transformer (ViT)

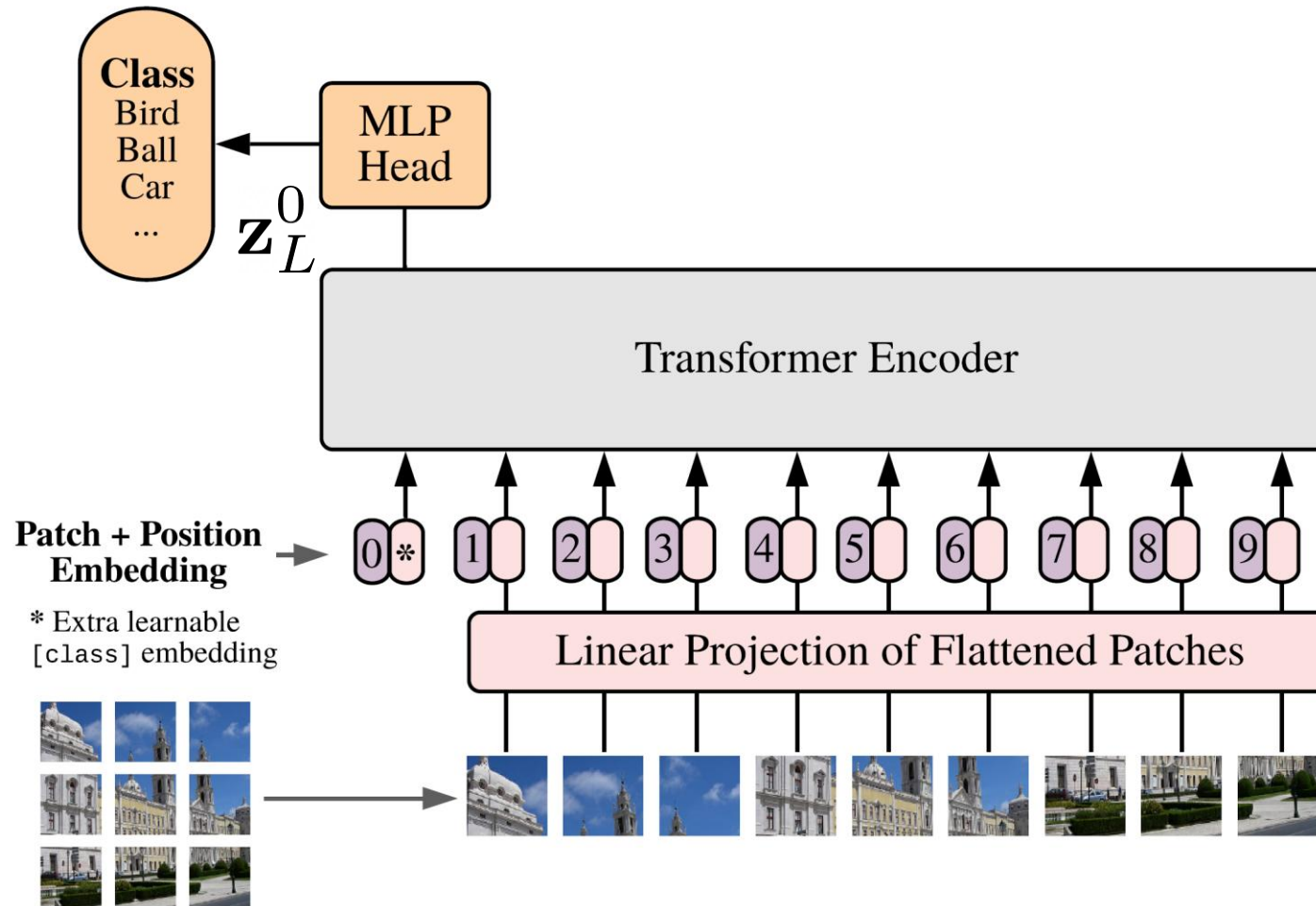
- Adding positional embedding
- Prepend a learnable embedding \mathbf{z}_0^0

\mathbf{z}_L^0 Will be used as the
image representation
After L attention layers

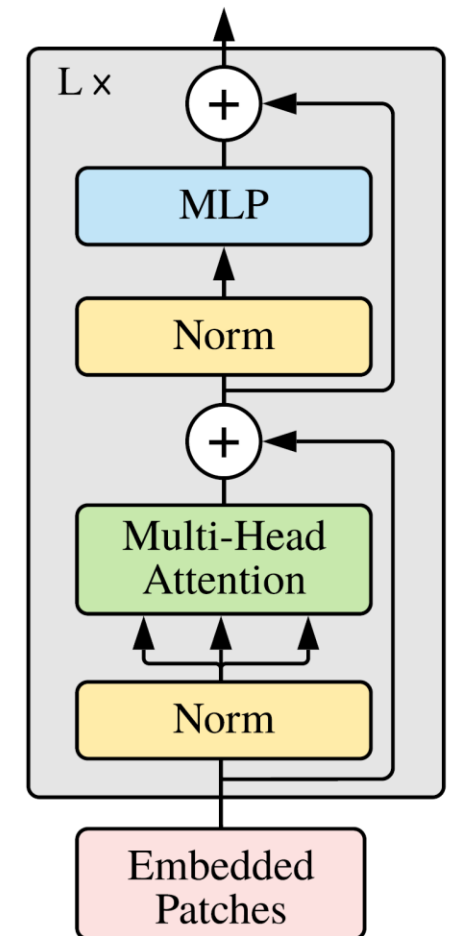


AN IMAGE IS WORTH 16x16 WORDS: TRANSFORMERS FOR IMAGE RECOGNITION AT SCALE. Dosovitskiy et al., ICLR'21

Vision Transformer (ViT)



Transformer Encoder



AN IMAGE IS WORTH 16x16 WORDS: TRANSFORMERS FOR IMAGE RECOGNITION AT SCALE. Dosovitskiy et al., ICLR'21

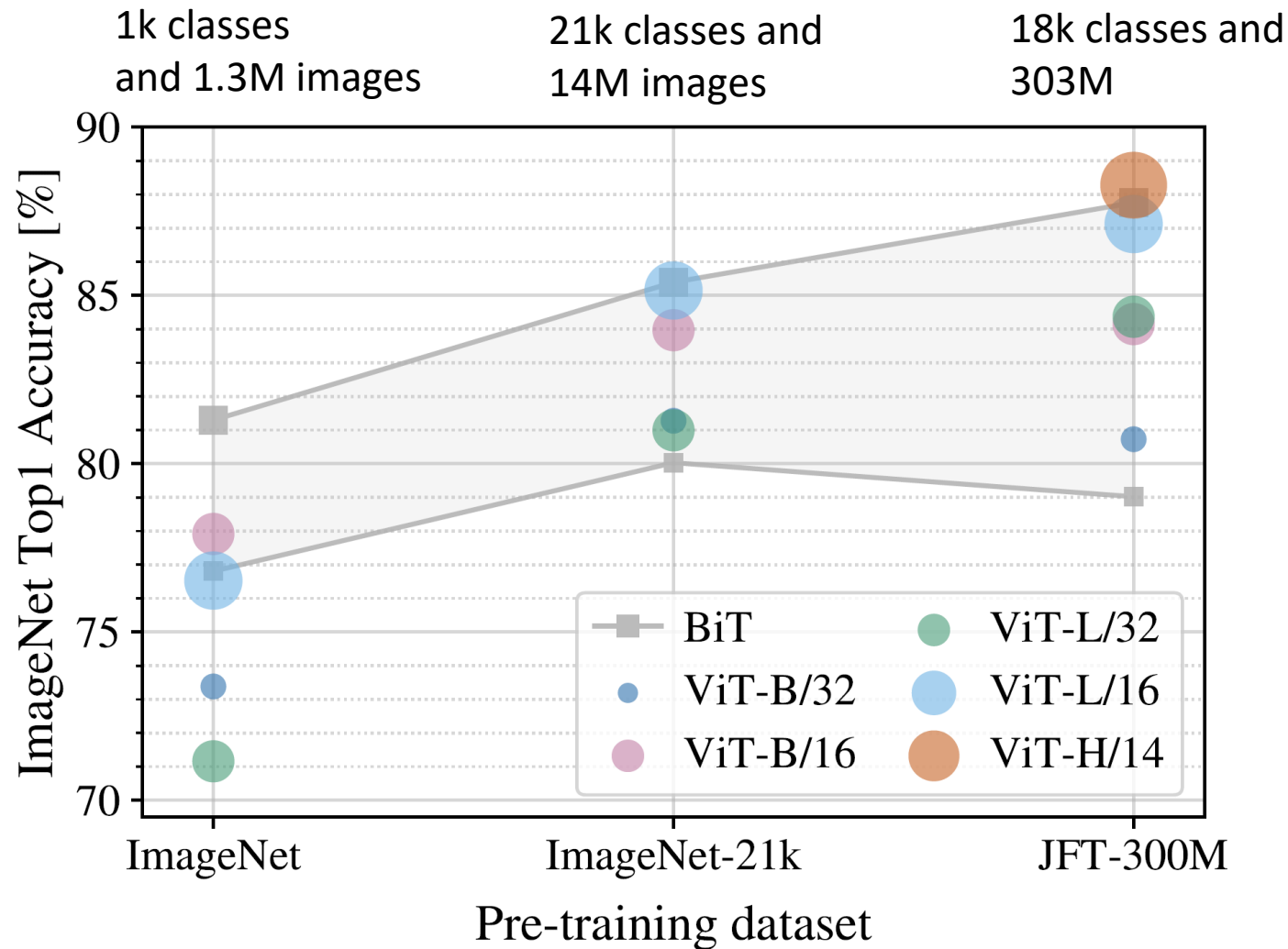
Vision Transformer (ViT)

- Pretrain on a large-scale dataset
- Fine-tune on different tasks

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

AN IMAGE IS WORTH 16x16 WORDS: TRANSFORMERS FOR IMAGE RECOGNITION AT SCALE. Dosovitskiy et al., ICLR'21

Vision Transformer (ViT)



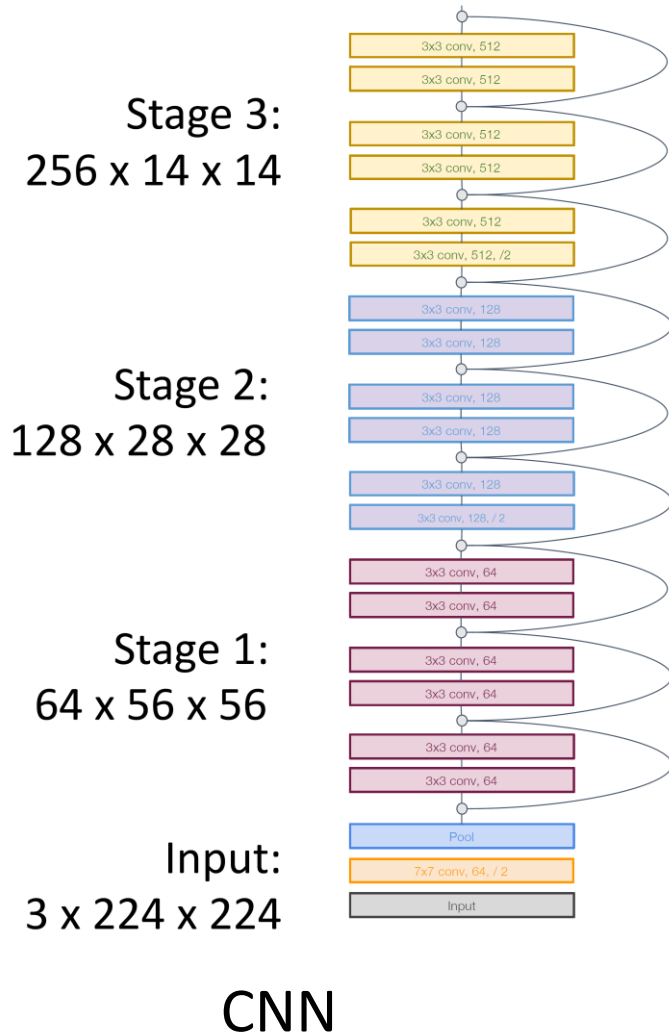
Big Transfer (BiT)

- ResNets-based transfer

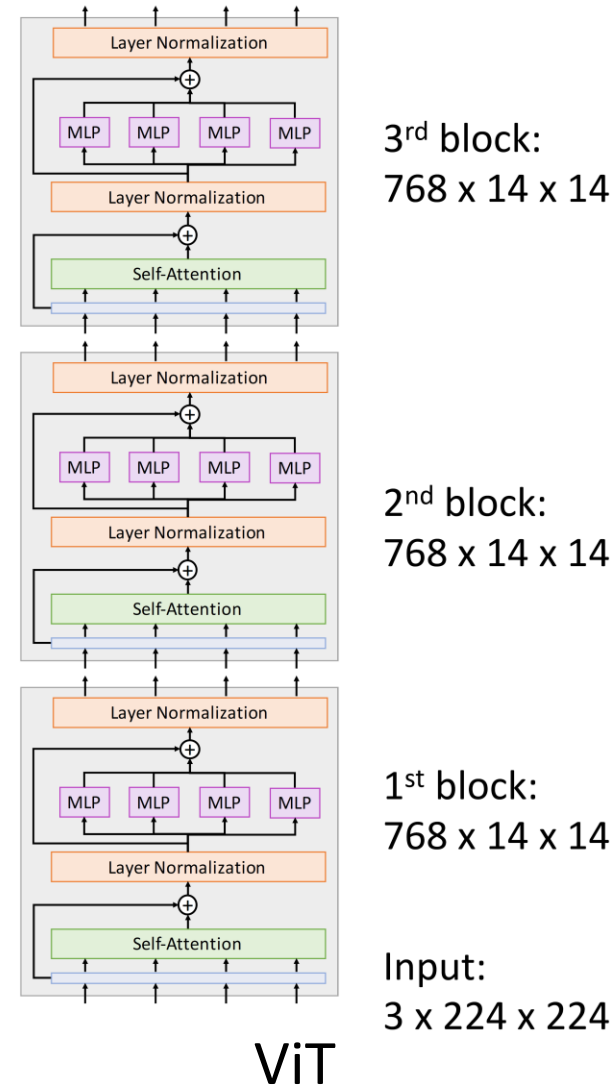
Vision transformer works better when pre-trained on large-scale dataset

AN IMAGE IS WORTH 16x16 WORDS: TRANSFORMERS FOR IMAGE RECOGNITION AT SCALE. Dosovitskiy et al., ICLR'21

ViT vs CNN

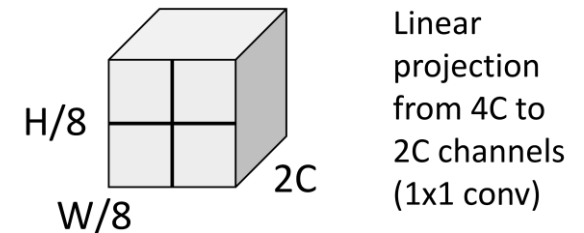
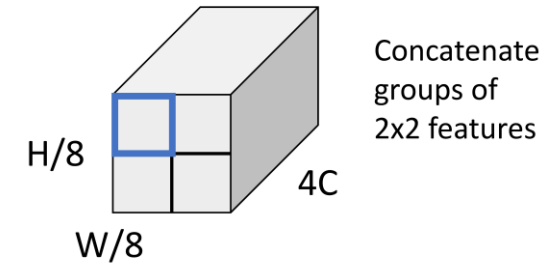
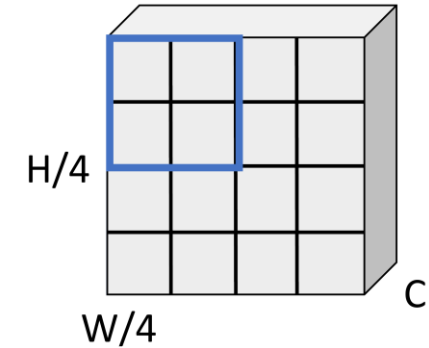
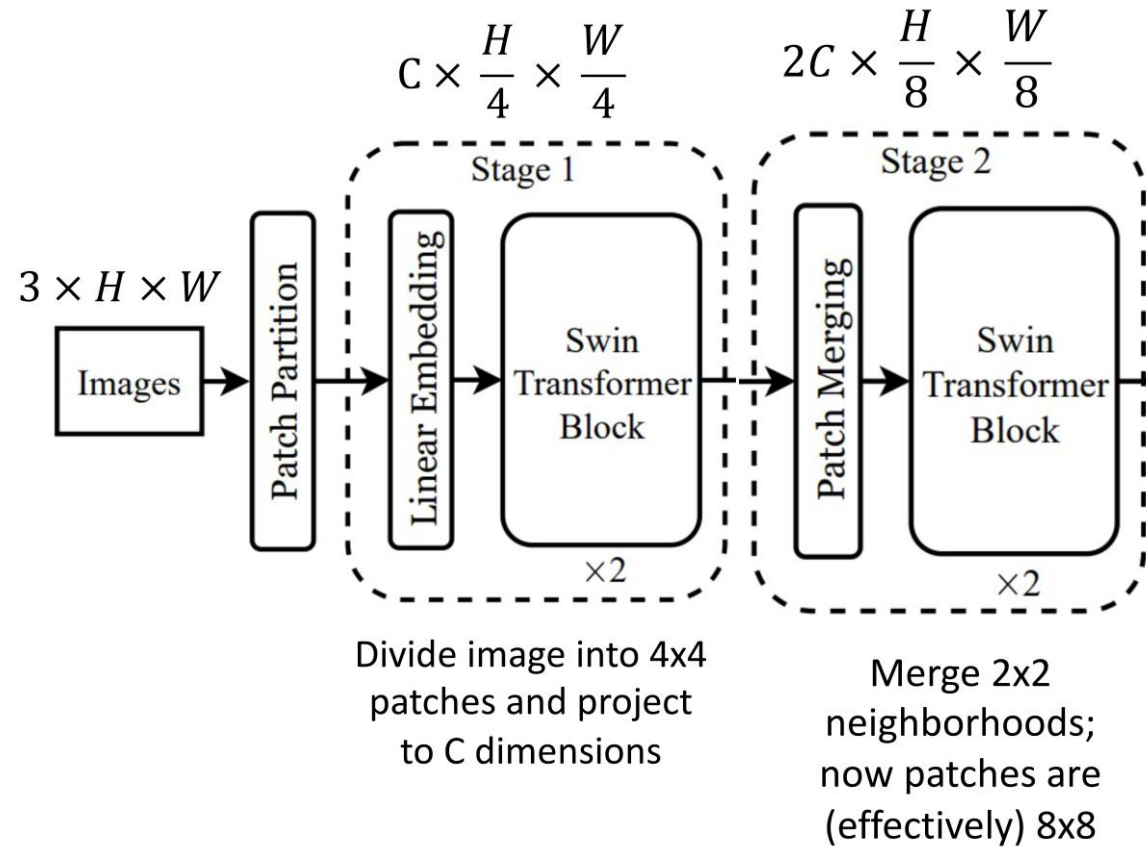


Hierarchical features are useful since objects in images can occur at various scales



In a ViT, all blocks have same resolution and number of channels (Isotropic architecture)

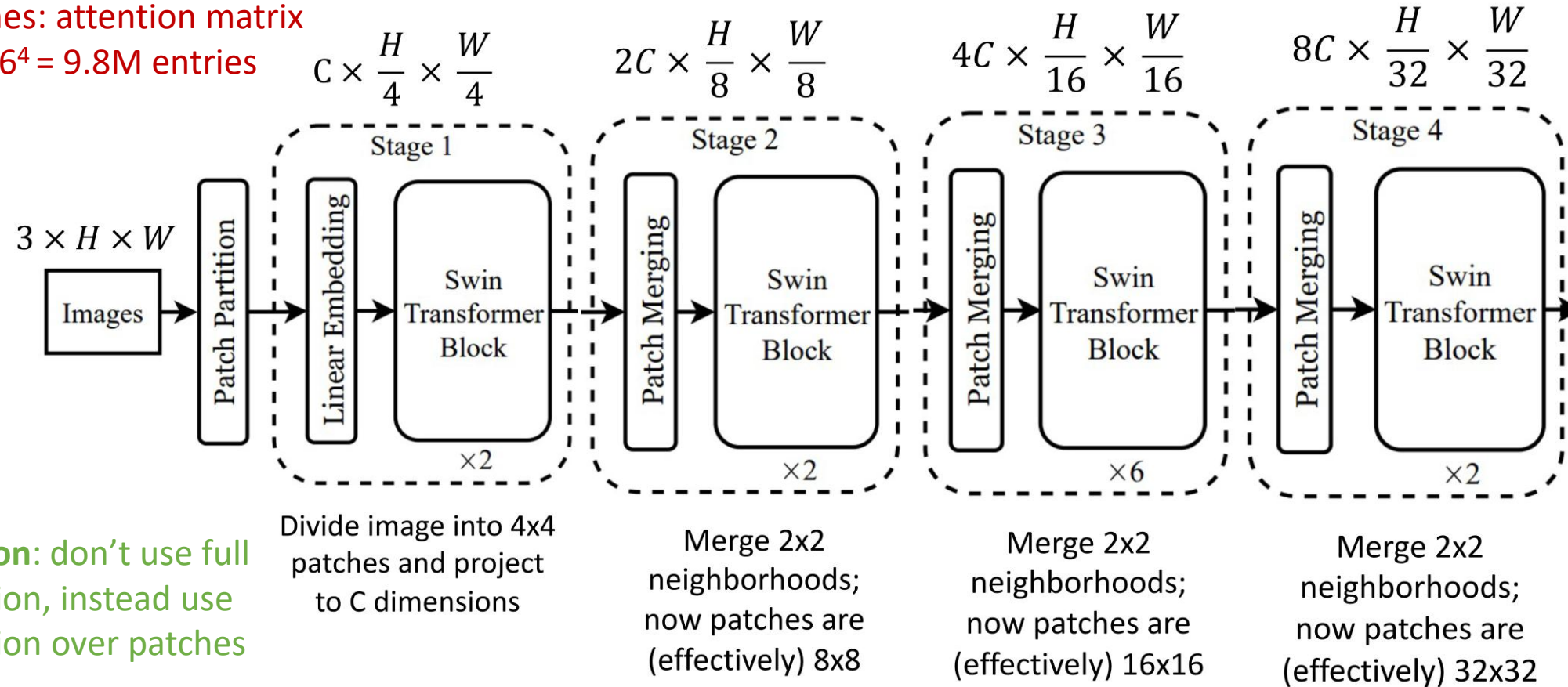
Hierarchical ViT: Swin Transformer



Liu et al, "Swin Transformer: Hierarchical Vision Transformer using Shifted Windows", CVPR 2021

Hierarchical ViT: Swin Transformer

Problem: 224x224 image
with 56x56 grid of 4x4
patches: attention matrix
has $56^4 = 9.8\text{M}$ entries



Liu et al, "Swin Transformer: Hierarchical Vision Transformer using Shifted Windows", CVPR 2021

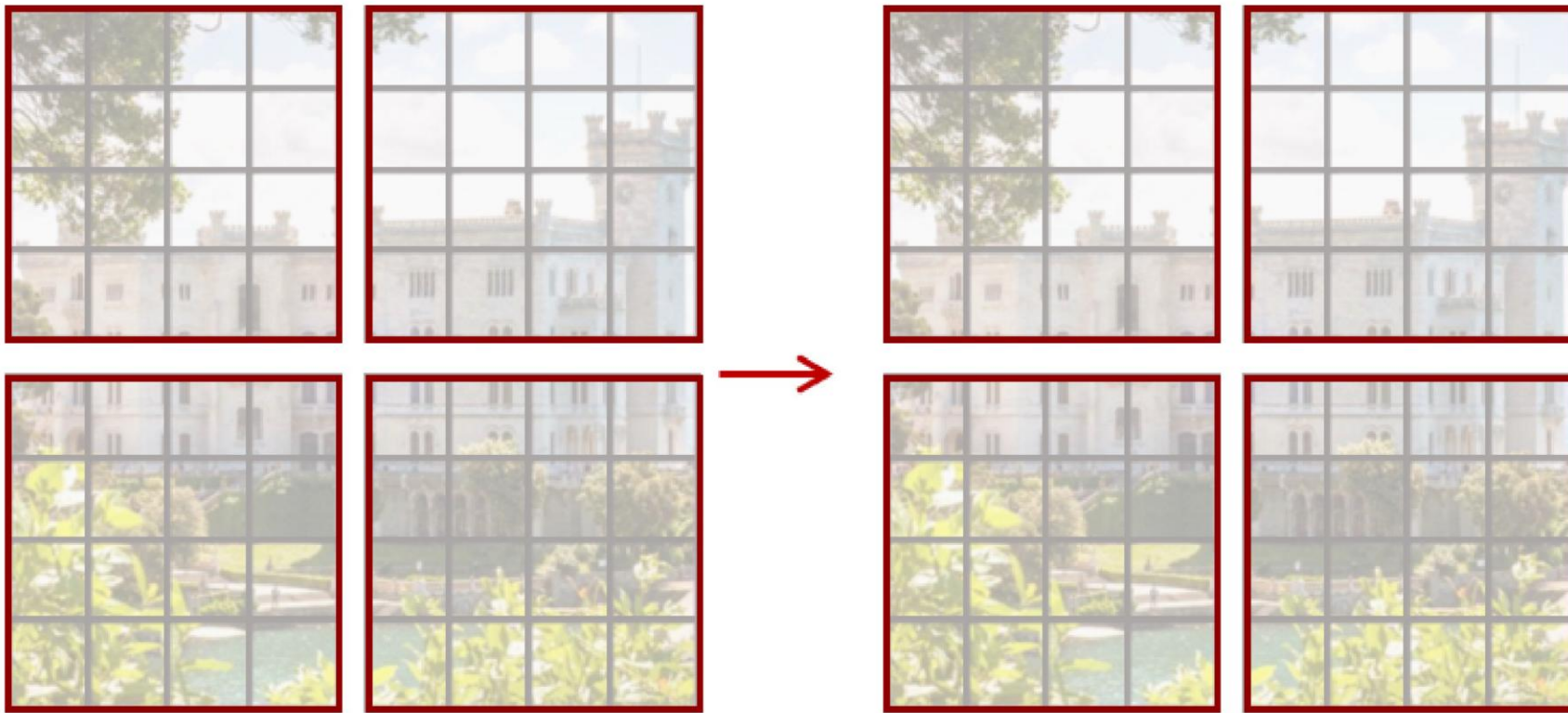
Hierarchical ViT: Swin Transformer

- With $H \times W$ grid of **tokens**, each attention matrix is $H \times H \times W \times W$ – **quadratic** in image size
- Window attention
 - Divide the image into windows of $M \times M$ tokens (here $M=4$)
 - Only compute attention within each window
 - Total size of attention matrices $M^4(H/M)(W/M) = M^2HW$
 - Linear in image size for fixed M ! Swin uses $M=7$ throughout the network



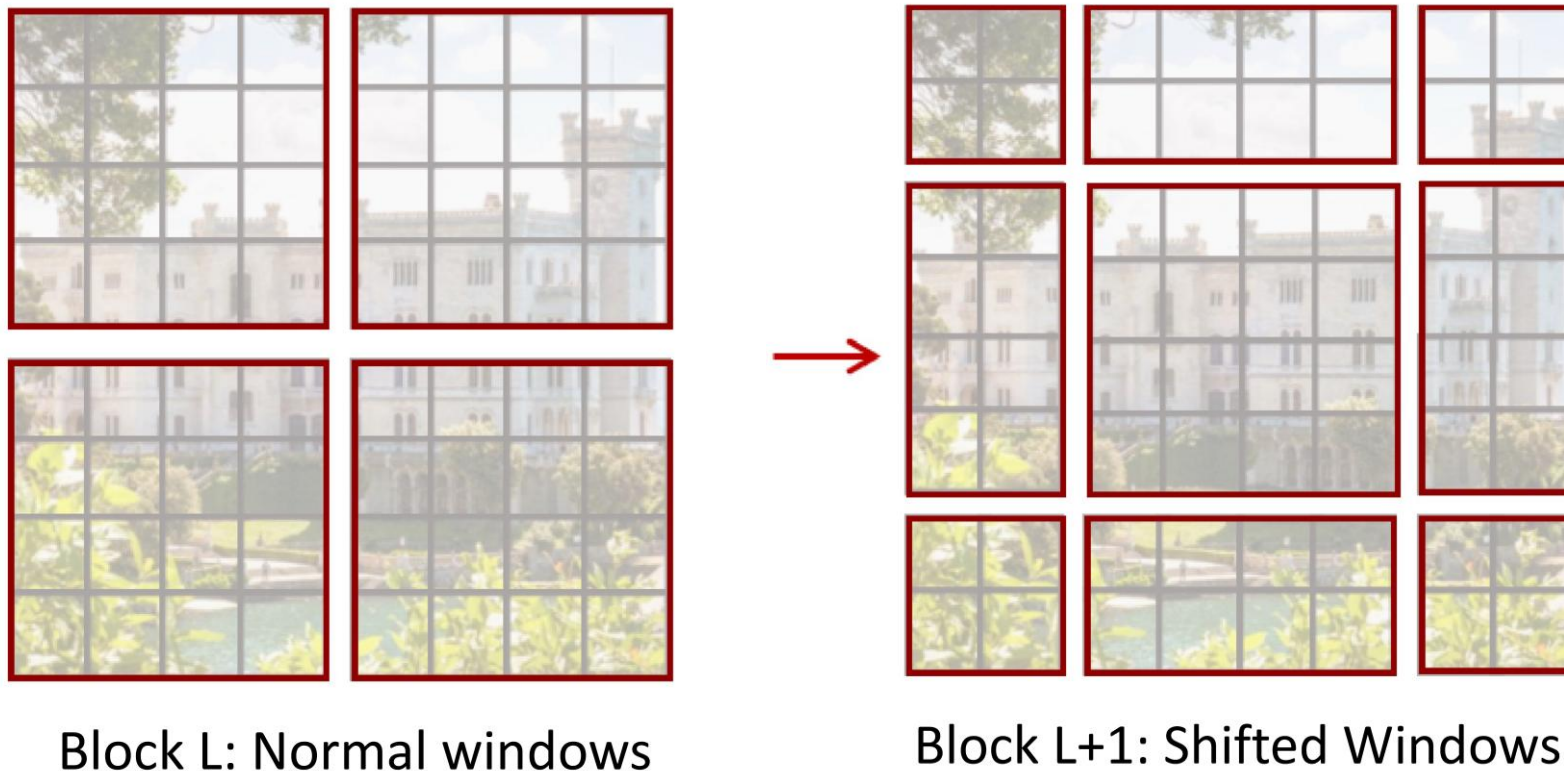
Hierarchical ViT: Swin Transformer

Problem: tokens only interact with other tokens within the same window; no communication across windows



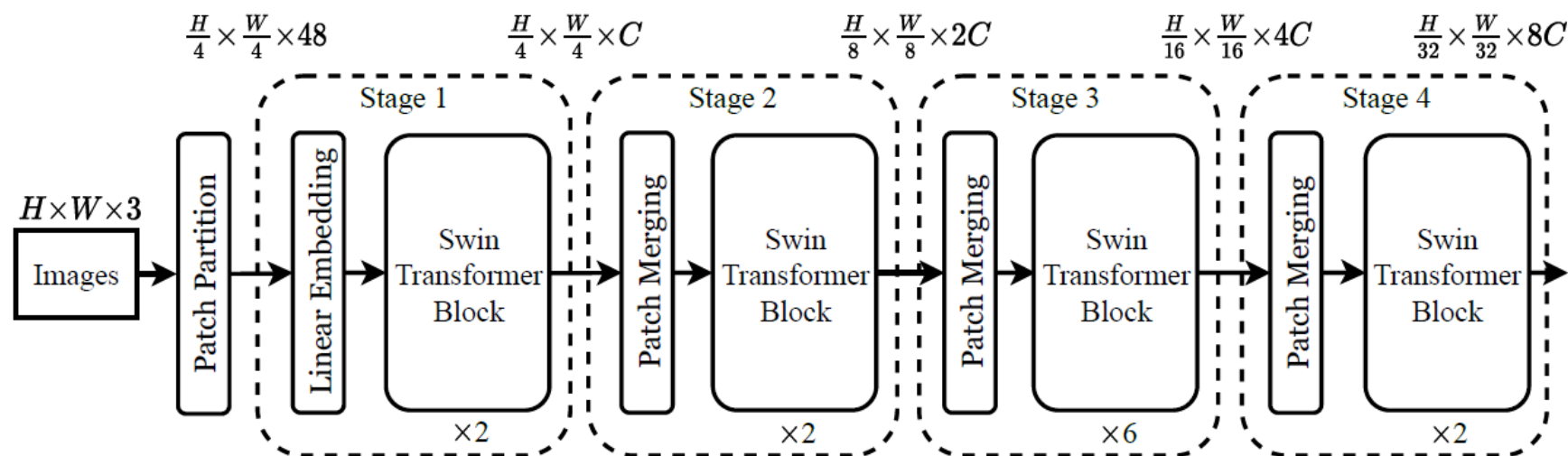
Hierarchical ViT: Swin Transformer

- Shifted Window Attention
- **Solution:** Alternate between normal windows and shifted windows in successive Transformer blocks



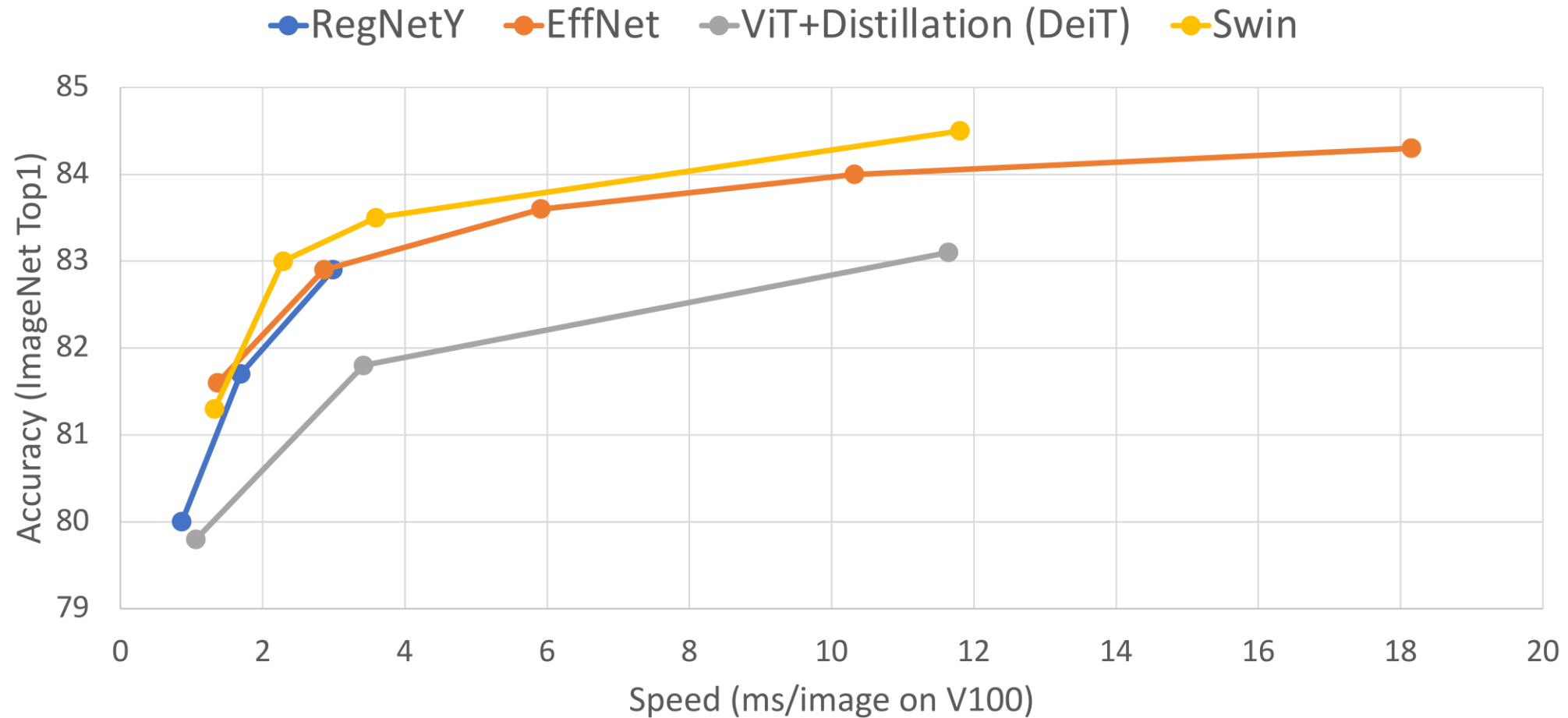
Hierarchical ViT: Swin Transformer

- Architecture variants



- Swin-T: $C = 96$, layer numbers = $\{2, 2, 6, 2\}$
- Swin-S: $C = 96$, layer numbers = $\{2, 2, 18, 2\}$
- Swin-B: $C = 128$, layer numbers = $\{2, 2, 18, 2\}$
- Swin-L: $C = 192$, layer numbers = $\{2, 2, 18, 2\}$

Hierarchical ViT: Swin Transformer



Summary

- Transformers
 - Can capture long-distance dependencies (global attention)
 - Computationally efficient, more parallelizable
- Vision transformers
 - Works better when pre-trained on large scale datasets (e.g., 300M images)
 - Swin transformer

Further Reading

- Transformer: Attention is all you need
<https://arxiv.org/abs/1706.03762>
- Vision transformer: An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale <https://arxiv.org/abs/2010.11929>
- Swin Transformer: Hierarchical Vision Transformer using Shifted Windows <https://arxiv.org/abs/2103.14030>