# Swin Transformer

Mengxue

# 1 Transformer

## Attention Is All You Need

**Ashish Vaswani***
Google Brain
avaswani@google.com

**Noam Shazeer***
Google Brain
noam@google.com

**Niki Parmar***
Google Research
nikip@google.com

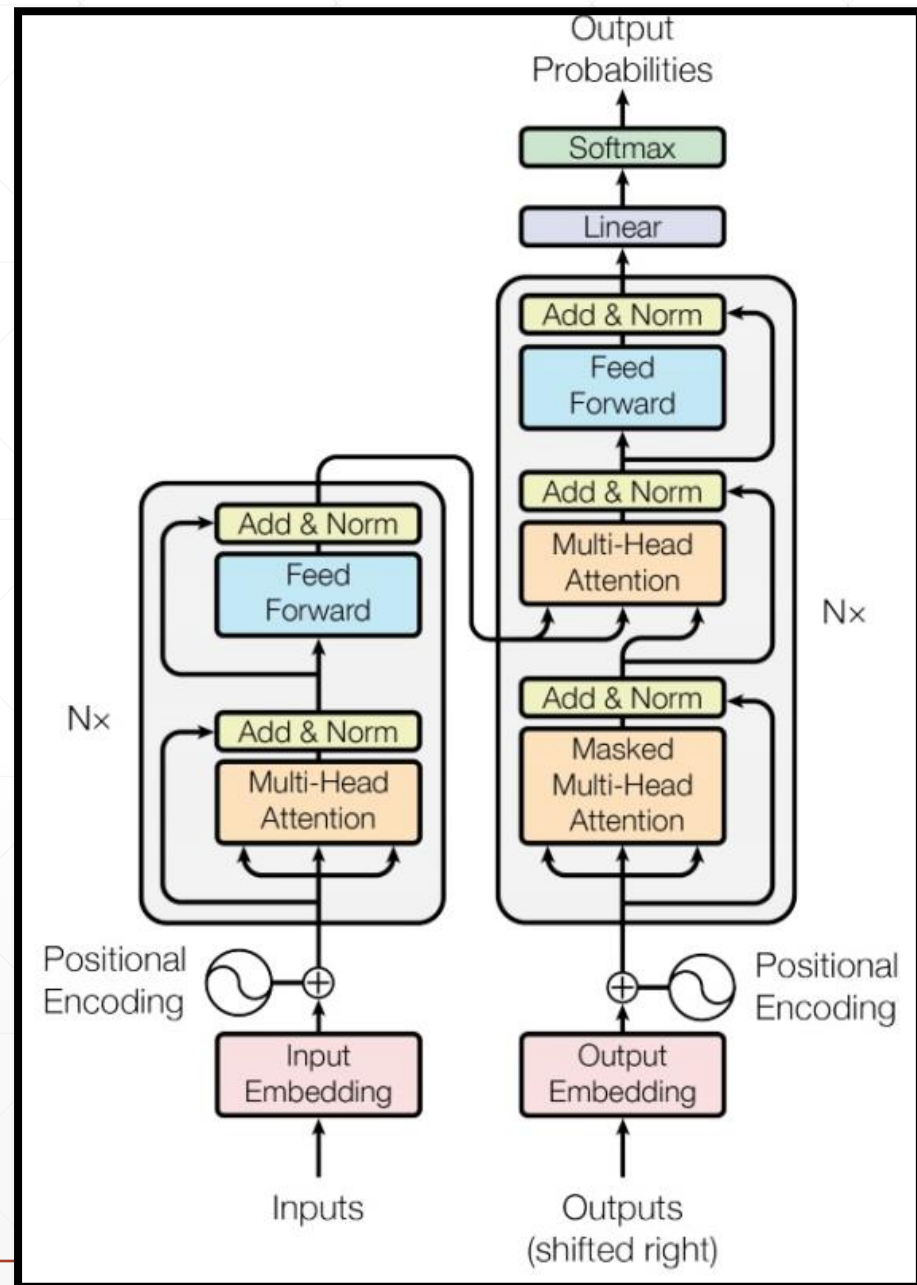**Jakob Uszkoreit***
Google Research
usz@google.com

**Llion Jones***
Google Research
llion@google.com

**Aidan N. Gomez*** [†]
University of Toronto
aidan@cs.toronto.edu

**Łukasz Kaiser***
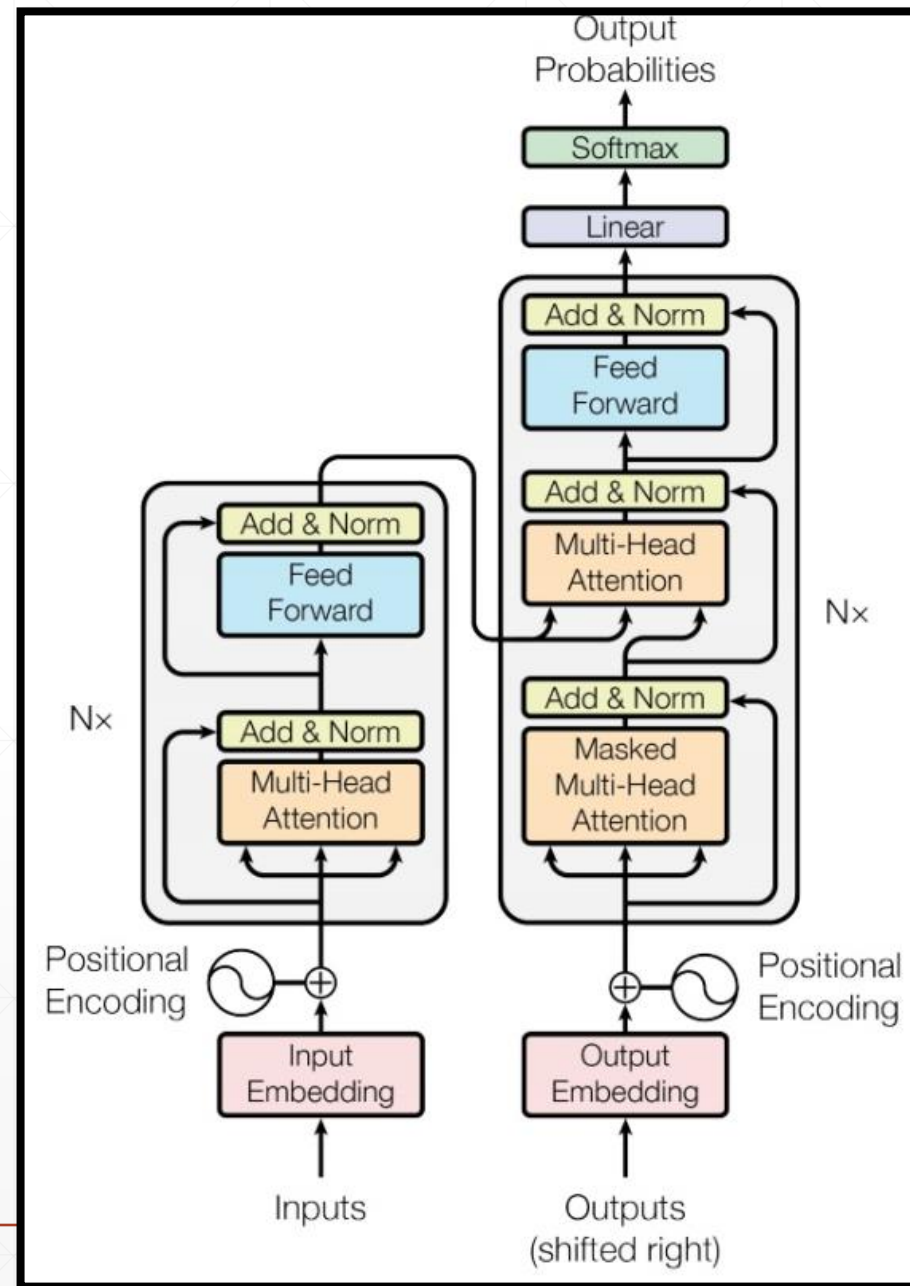Google Brain
lukaszkaiser@google.com

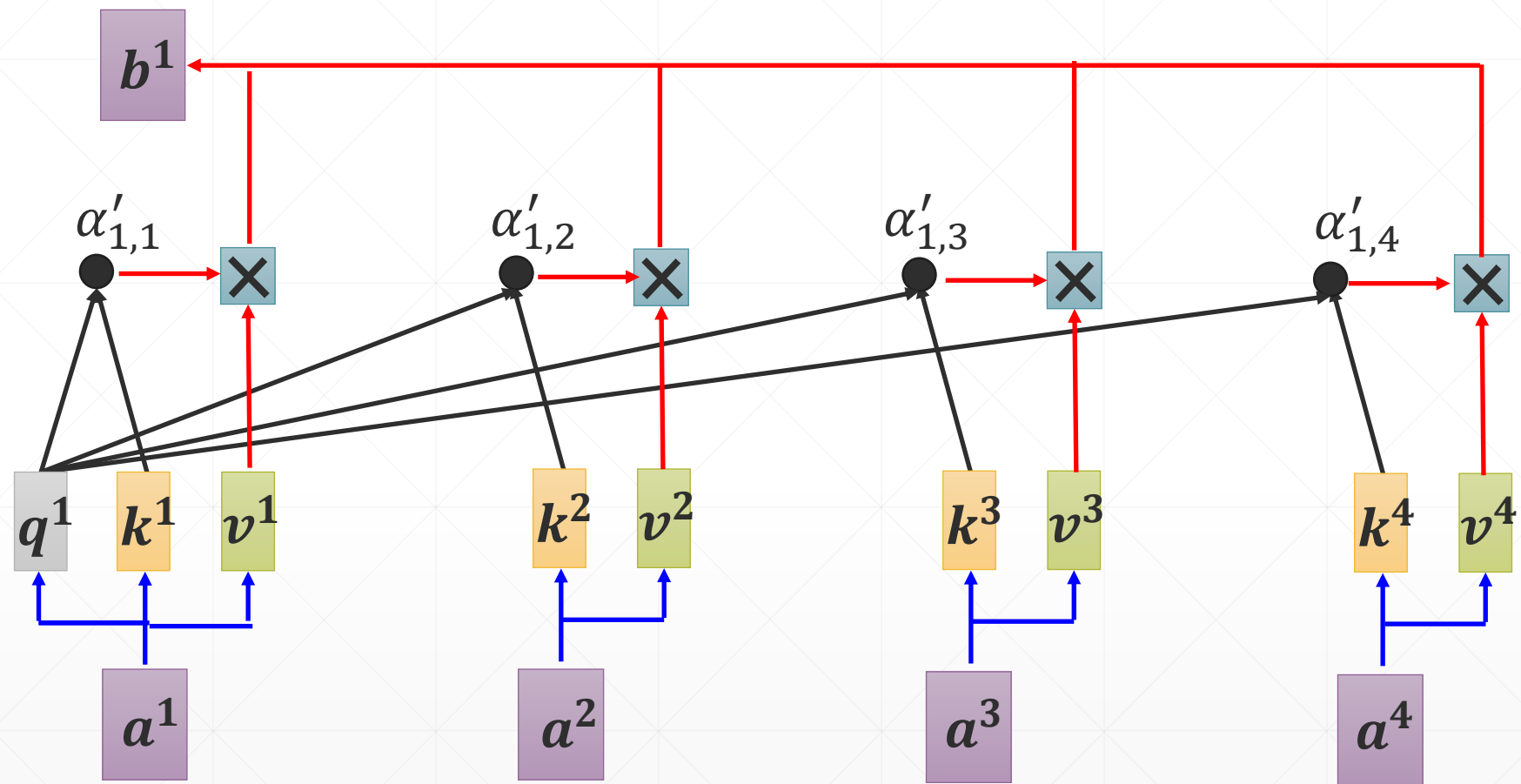**Illia Polosukhin*** [‡]
illia.polosukhin@gmail.com

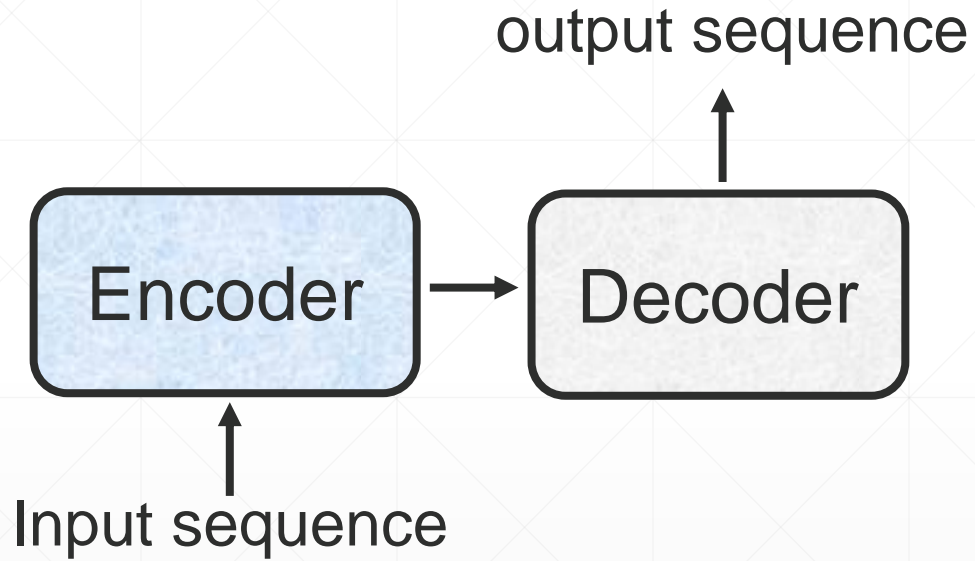**Neural Information Processing Systems (NIPS 2017)**

# Tips

- Self-attention

- Multihead Self-attention

- Position encoding
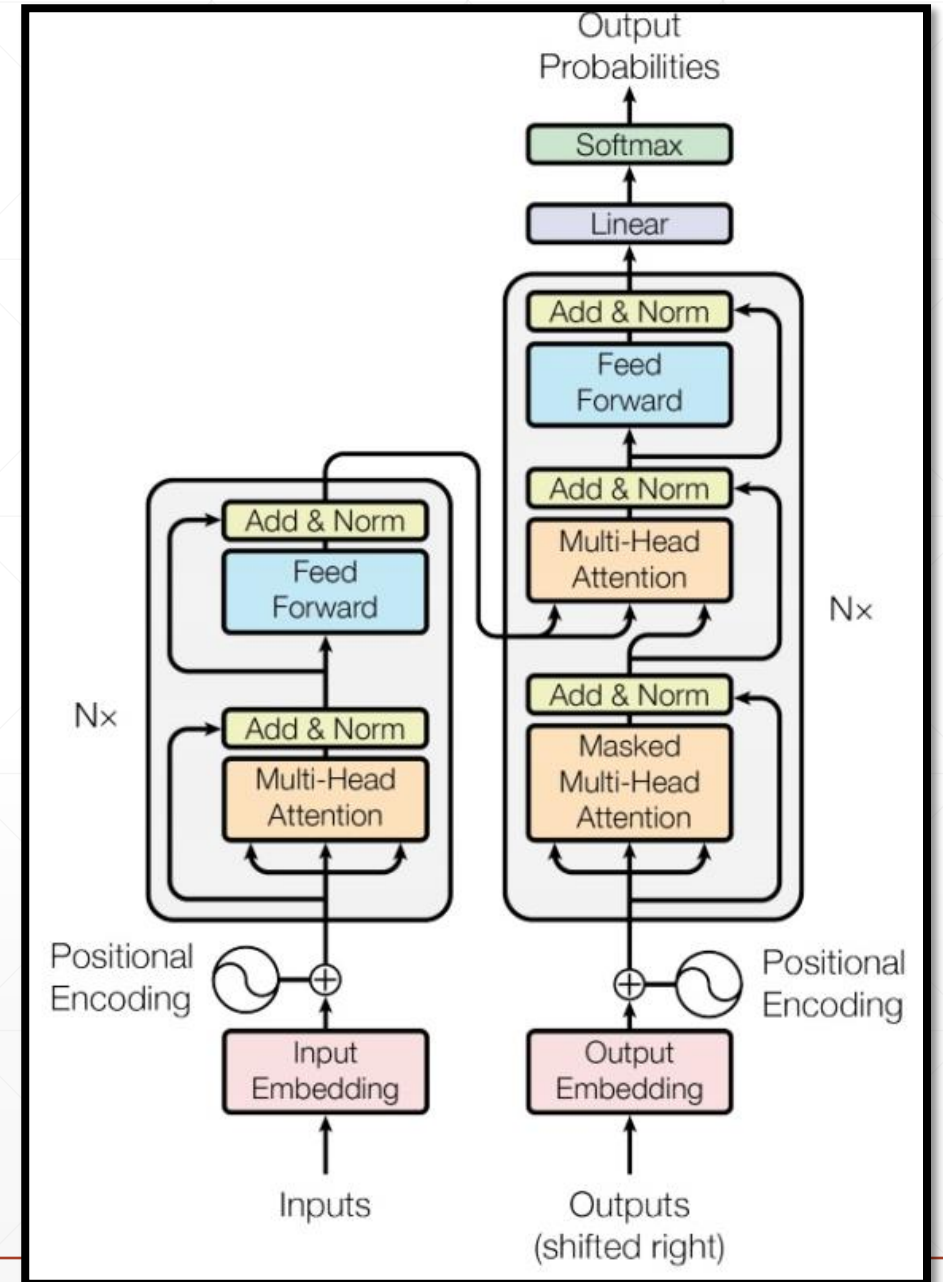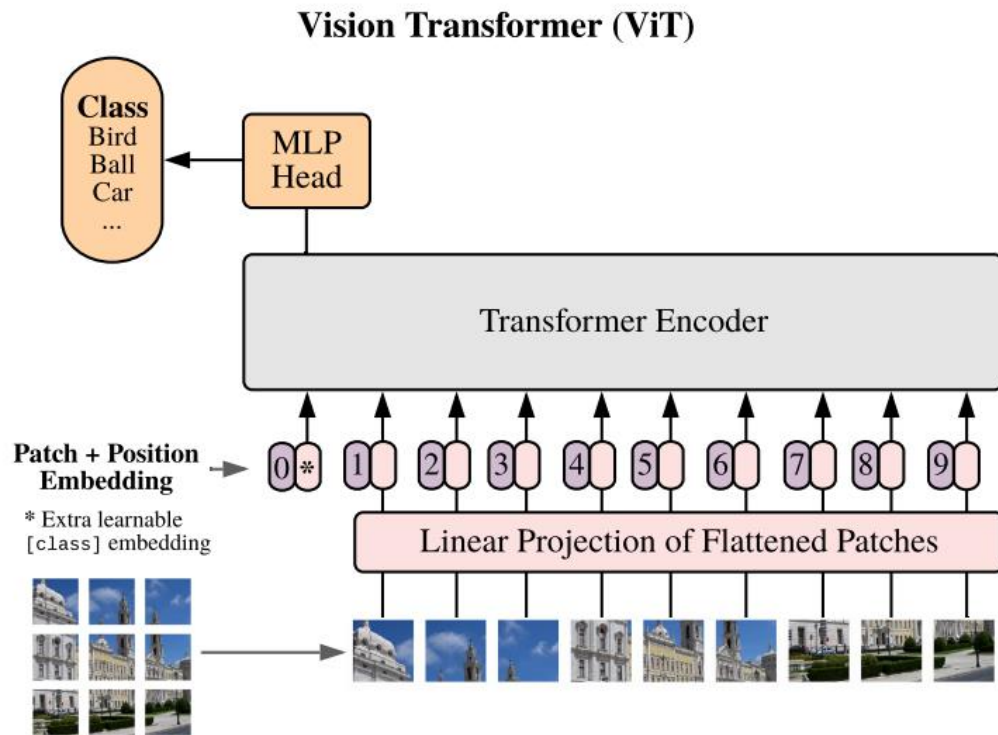
- Masked Multi-Head Attention

- Cross Attention

# Seq2seq Model



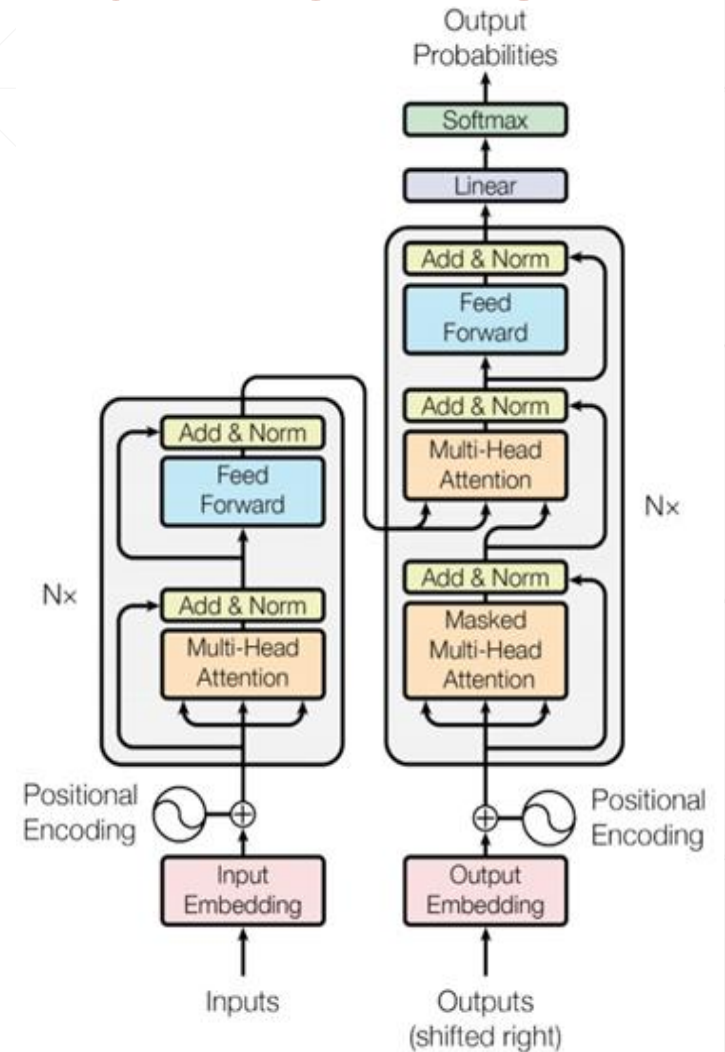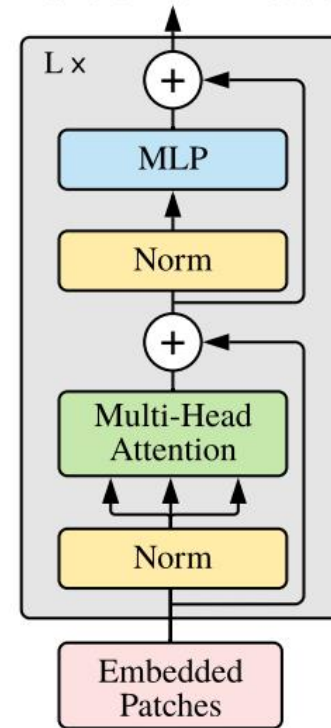output sequence

Encoder → Decoder

Input sequence

**Transformer→**

# 2 ViT—AN IMAGE IS WORTH 16×16 WORDS

# Patch Embedding+Positional Encoding

标准的接受token的一维嵌入向量作为输入。为了处理二维数据，要进行reshape。

原始图像输入： （H,W）是图片分辨率，C是通道数

$$\mathbf{x} \in \mathbb{R}^{H \times W \times C}$$

reshape（分割patch）：P是patch的大小，N是patch的个数

$$\mathbf{x} \in \mathbb{R}^{N \times \left(P^2 \cdot C\right)}, \quad N = HW / P^2$$

→ 分块

flatten(拍平，映射成Transformer接受的固定大小D，映射E是可学习的):

$$\mathbf{z}_0 = \left[\mathbf{x}_{class}; \mathbf{x}_p^1 \mathbf{E}; \mathbf{x}_p^2 \mathbf{E}; \cdots ; \mathbf{x}_p^N \mathbf{E}\right] + \mathbf{E}_{pos}, \quad \mathbf{E} \in \mathbb{R}^{\left(P^2 \cdot C\right) \times D}, \quad \mathbf{E}_{pos} \in \mathbb{R}^{(N+1) \times D}$$

→ 维度转换

映射后的结果称为 patch embeddings。

在patch前面添加一个可学习的xclass，代表着图片的标签信息（全局信息）

# 3 Swin Transformer: Hierarchical Vision Transformer using Shifted Windows

Ze Liu[†*]    Yutong Lin[†*]    Yue Cao[*]    Han Hu[*‡]    Yixuan Wei[†]

Zheng Zhang    Stephen Lin    Baining Guo
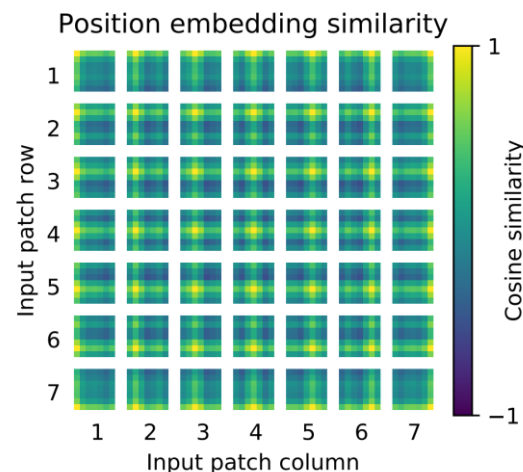
Microsoft Research Asia

{v-zeliu1,v-yutlin,yuecao,hanhu,v-yixwe,zhez,stevelin,bainguo}@microsoft.com

**CVPR 2021**

**Motivation:**

Tokens are all of a fixed scale ,which is unsuitable for vision applications

Higher resolution of pixels in images compared to words in passages of text.

# Overall Architecture



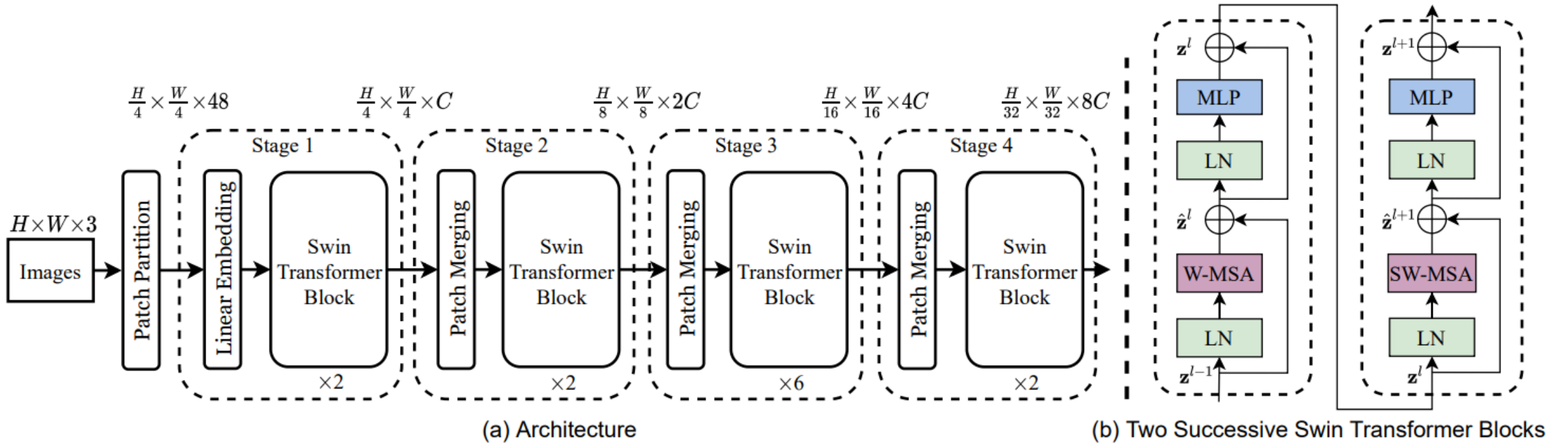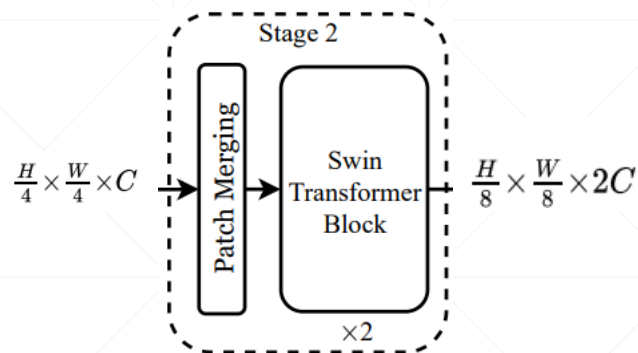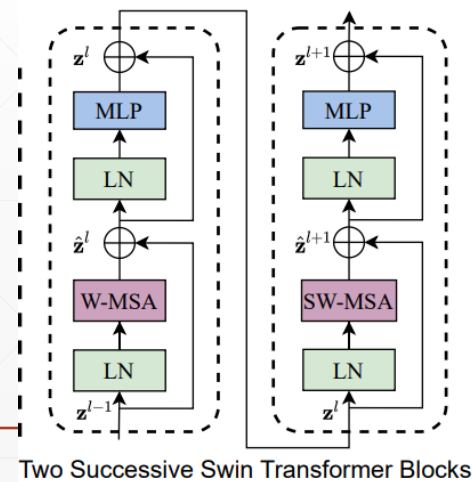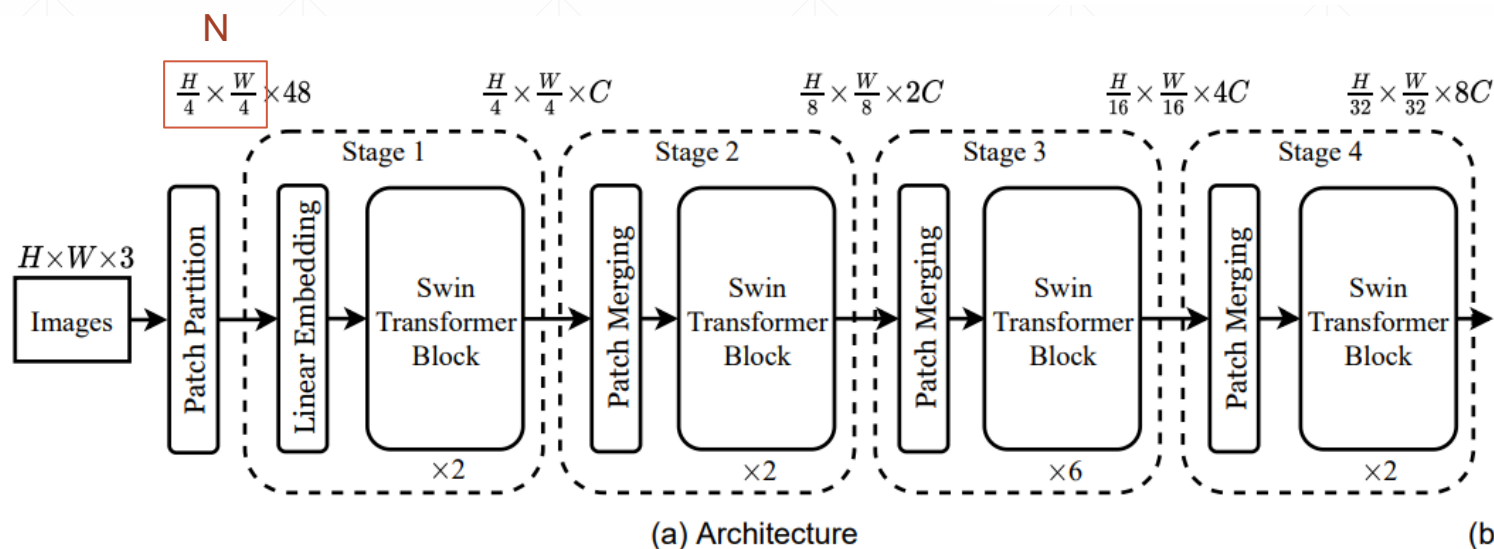Figure 3. (a) The architecture of a Swin Transformer (Swin-T); (b) two successive Swin Transformer Blocks (notation presented with Eq. (3)). W-MSA and SW-MSA are multi-head self attention modules with regular and shifted windowing configurations, respectively.

# Hierarchical feature maps → Patch Merging



$\frac{H}{4} \times \frac{W}{4} \times C$ → Patch Merging → Swin Transformer Block ×2 → $\frac{H}{8} \times \frac{W}{8} \times 2C$ (Stage 2)

patch size: 4 × 4 × 48(48=4 × 4 × 3)
patch num: H/4 × W/4

(a) Swin Transformer (ours)

(b) ViT

$H \times W \times 3$ Images → Patch Partition → Linear Embedding → Swin Transformer Block ×2 (Stage 1) → Patch Merging → Swin Transformer Block ×2 (Stage 2) → Patch Merging → Swin Transformer Block ×6 (Stage 3) → Patch Merging → Swin Transformer Block ×2 (Stage 4)

N $\frac{H}{4} \times \frac{W}{4} \times 48$   $\frac{H}{4} \times \frac{W}{4} \times C$   $\frac{H}{8} \times \frac{W}{8} \times 2C$   $\frac{H}{16} \times \frac{W}{16} \times 4C$   $\frac{H}{32} \times \frac{W}{32} \times 8C$

(a) Architecture

(b) Two Successive Swin Transformer Blocks

# W-MSA and SW-MSA



Figure 2. An illustration of the *shifted window* approach for computing self-attention in the proposed Swin Transformer architecture. In layer *l* (left), a regular window partitioning scheme is adopted, and self-attention is computed within each window. In the next layer *l* + 1 (right), the window partitioning is shifted, resulting in new windows. The self-attention computation in the new windows crosses the boundaries of the previous windows in layer *l*, providing connections among them.

# W-MSA and SW-MSA

Window contains M × M patches
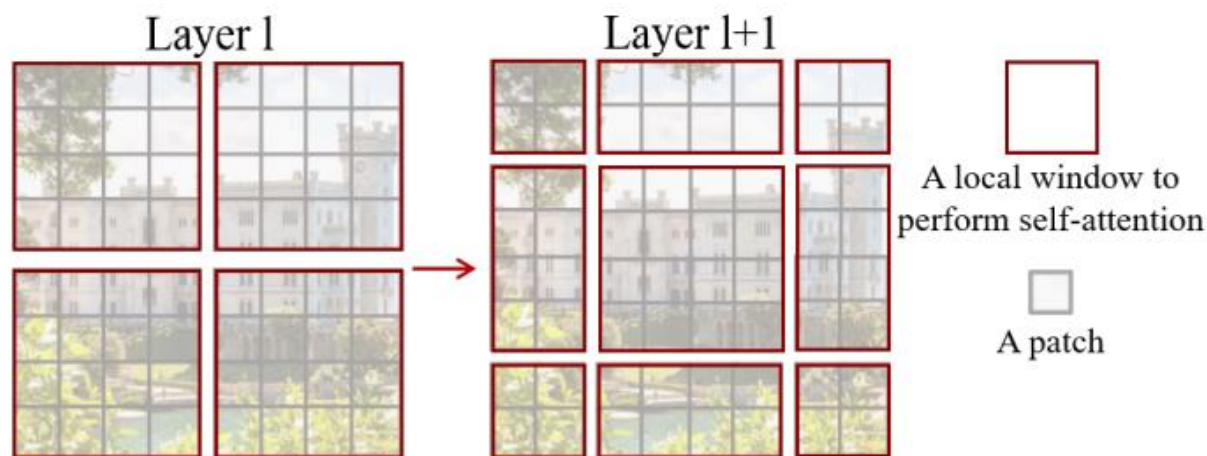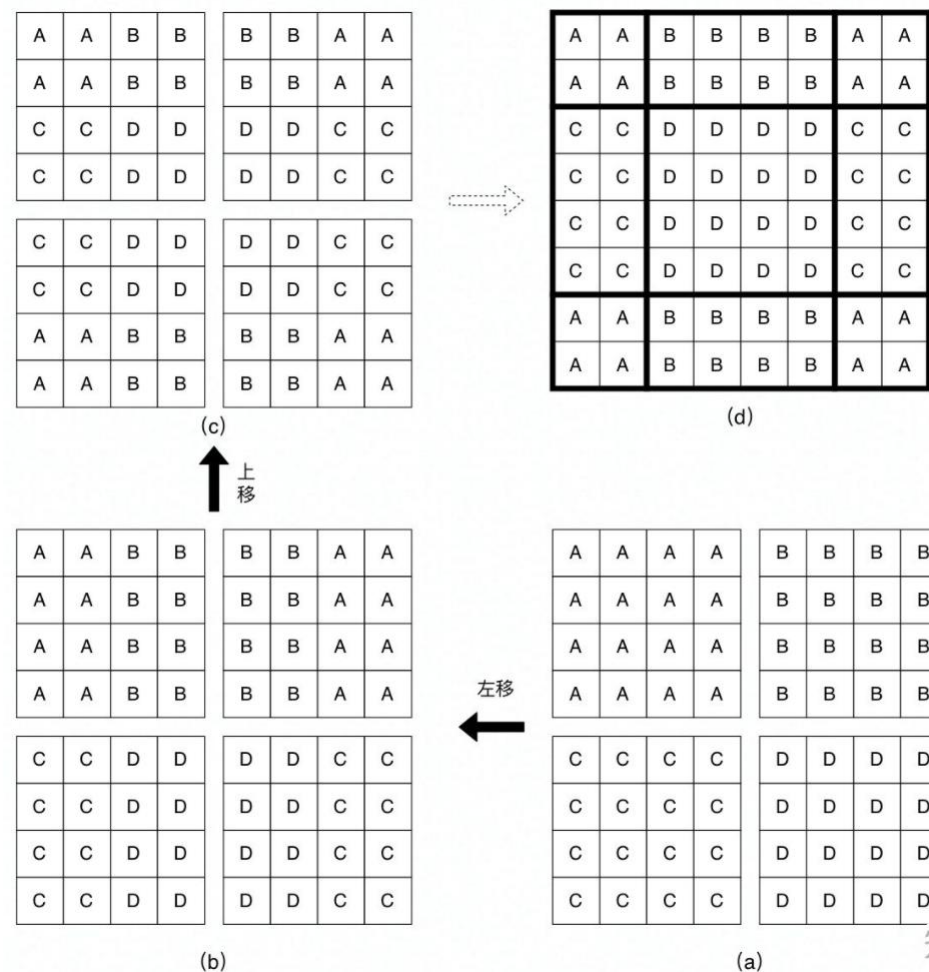If M=4, feature map=8 × 8,



Figure 2. An illustration of the *shifted window* approach for computing self-attention in the proposed Swin Transformer architecture. In layer $l$ (left), a regular window partitioning scheme is adopted, and self-attention is computed within each window. In the next layer $l + 1$ (right), the window partitioning is shifted, resulting in new windows. The self-attention computation in the new windows crosses the boundaries of the previous windows in layer $l$, providing connections among them.

# Swin Transformer Block



- W-MSA (Window-Multihead Self Attention)
- SW-MSA (Shifted Window-Multihead Self Attention)

# Transformer Decoder

# Efficient batch computation for shifted configuration



Figure 4. Illustration of an efficient batch computation approach for self-attention in shifted window partitioning.

# Efficient batch computation for shifted configuration



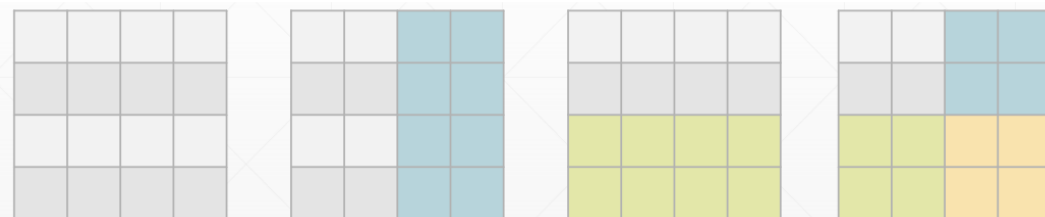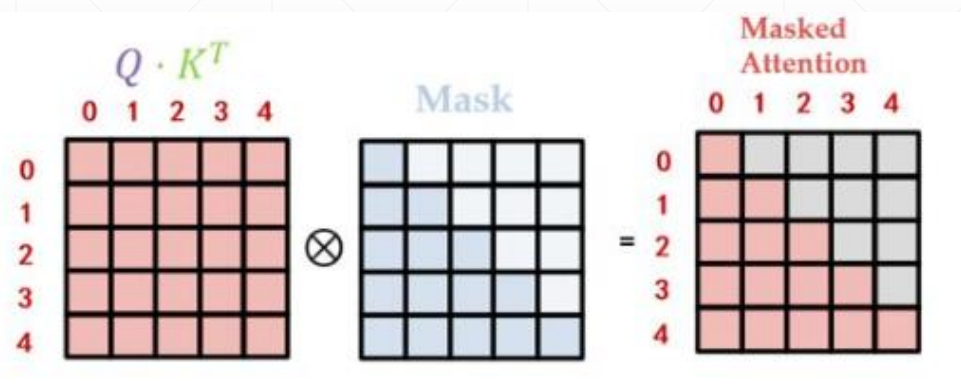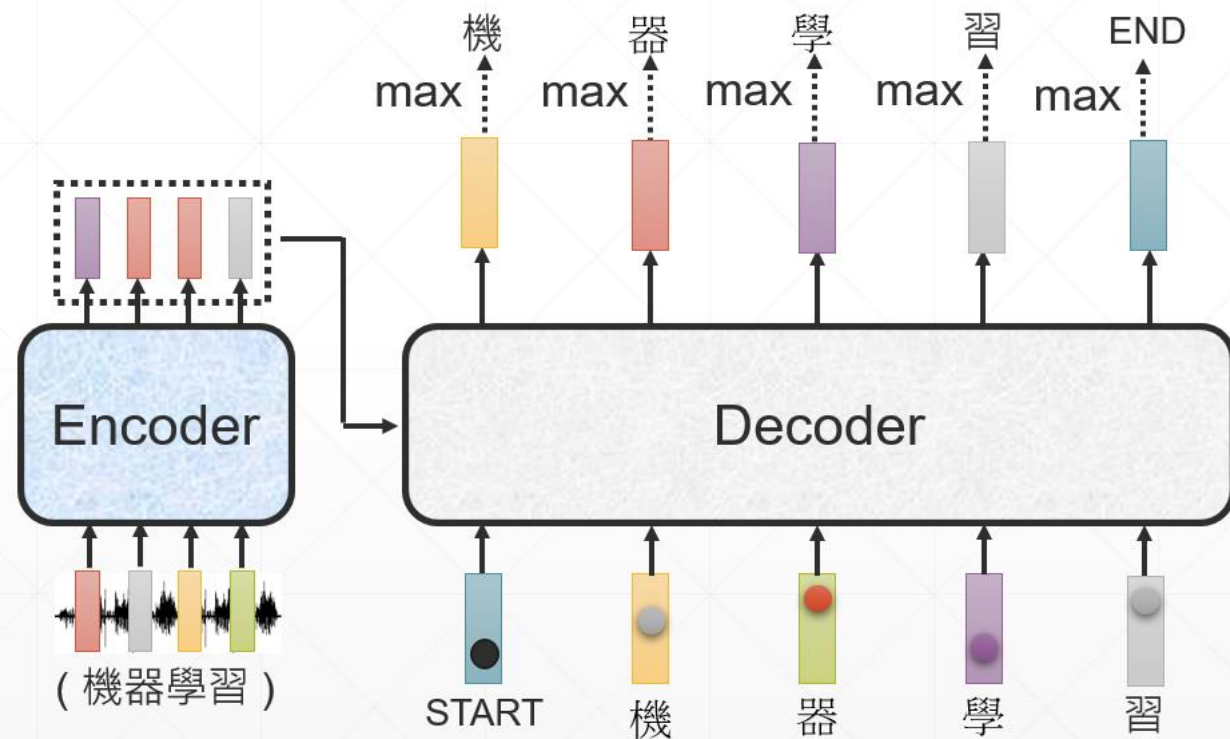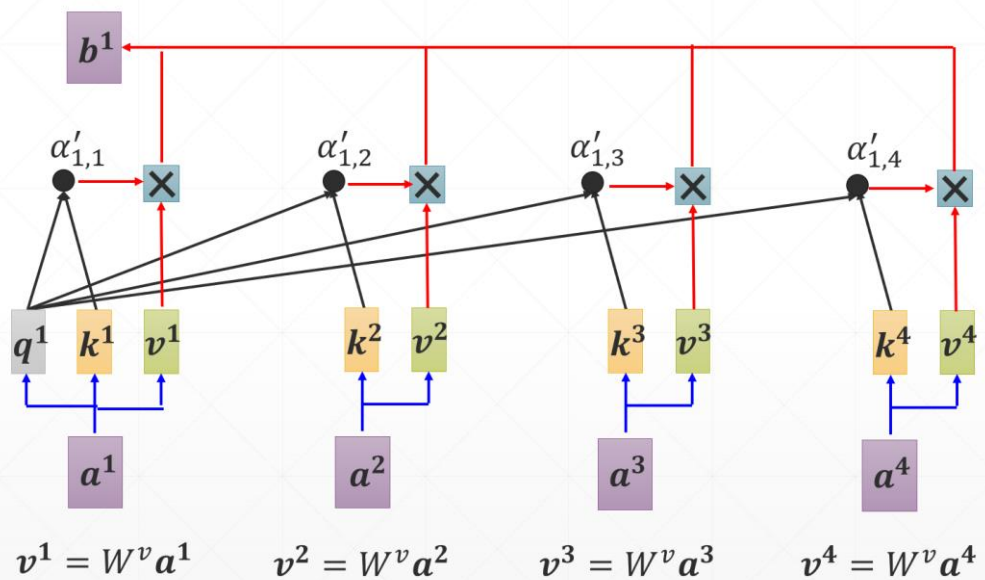| method | MSA in a stage (ms) | | | | Arch. (FPS) | | |
|---|---|---|---|---|---|---|---|
| | S1 | S2 | S3 | S4 | T | S | B |
| sliding window (naive) | 122.5 | 38.3 | 12.1 | 7.6 | 183 | 109 | 77 |
| sliding window (kernel) | 7.6 | 4.7 | 2.7 | 1.8 | 488 | 283 | 187 |
| Performer [13] | 4.8 | 2.8 | 1.8 | 1.5 | 638 | 370 | 241 |
| window (w/o shifting) | 2.8 | 1.7 | 1.2 | 0.9 | 770 | 444 | 280 |
| shifted window (padding) | 3.3 | 2.3 | 1.9 | 2.2 | 670 | 371 | 236 |
| shifted window (cyclic) | 3.0 | 1.9 | 1.3 | 1.0 | 755 | 437 | 278 |

Table 5. Real speed of different self-attention computation methods and implementations on a V100 GPU.

| | Backbone | ImageNet | | COCO | | ADE20k |
|---|---|---|---|---|---|---|
| | | top-1 | top-5 | AP$^{box}$ | AP$^{mask}$ | mIoU |
| sliding window | Swin-T | 81.4 | 95.6 | 50.2 | 43.5 | 45.8 |
| Performer [13] | Swin-T | 79.0 | 94.2 | - | - | - |
| shifted window | Swin-T | 81.3 | 95.6 | 50.5 | 43.7 | 46.1 |

Table 6. Accuracy of Swin Transformer using different methods for self-attention computation on three benchmarks.

# Relative position bias

$$Attention(Q, K, V) = SoftMax(QK^T/\sqrt{d} + B)V$$

| | ImageNet | | COCO | | ADE20k |
|---|---|---|---|---|---|
| | top-1 | top-5 | $AP^{box}$ | $AP^{mask}$ | mIoU |
| w/o shifting | 80.2 | 95.1 | 47.7 | 41.5 | 43.3 |
| shifted windows | **81.3** | **95.6** | **50.5** | **43.7** | **46.1** |
| no pos. | 80.1 | 94.9 | 49.2 | 42.6 | 43.8 |
| abs. pos. | 80.5 | 95.2 | 49.0 | 42.4 | 43.2 |
| abs.+rel. pos. | 81.3 | 95.6 | 50.2 | 43.4 | 44.0 |
| rel. pos. w/o app. | 79.3 | 94.7 | 48.2 | 41.9 | 44.1 |
| rel. pos. | **81.3** | **95.6** | **50.5** | **43.7** | **46.1** |

Table 4. Ablation study on the *shifted windows* approach and different position embedding methods on three benchmarks, using the Swin-T architecture. w/o shifting: all self-attention modules adopt regular window partitioning, without *shifting*; abs. pos.: absolute position embedding term of ViT; rel. pos.: the default settings with an additional relative position bias term (see Eq. (4)); app.: the first scaled dot-product term in Eq. (4).

具体思想参考UniLMV2

# 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 = $\{2, 2, 18, 2\}$

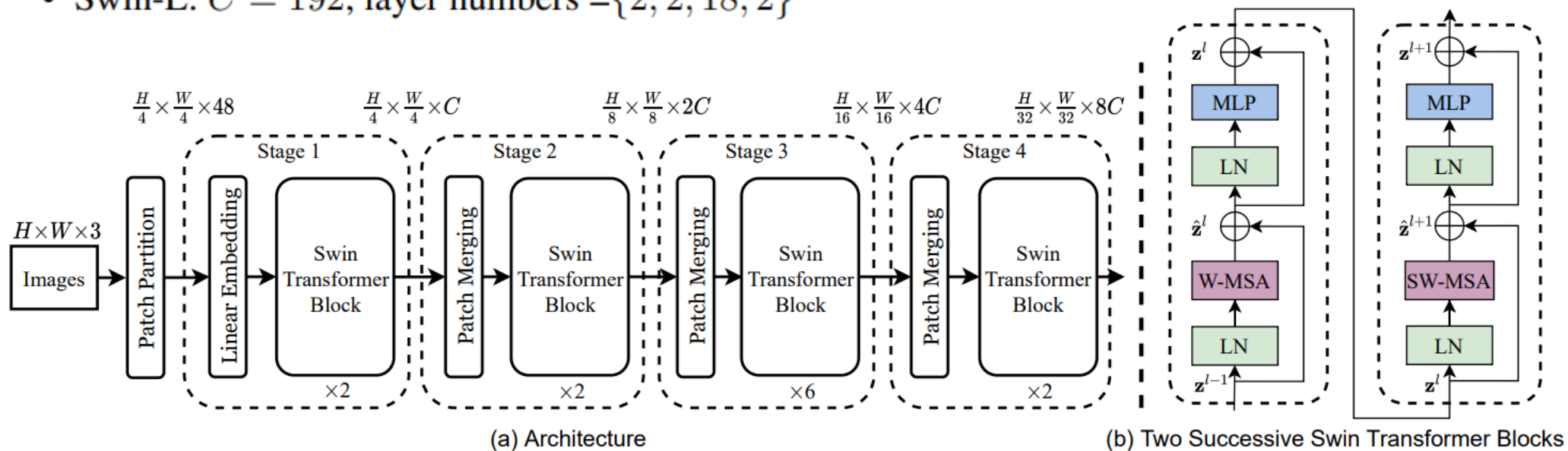- Swin-L: $C = 192$, layer numbers = $\{2, 2, 18, 2\}$



Figure 3. (a) The architecture of a Swin Transformer (Swin-T); (b) two successive Swin Transformer Blocks (notation presented with Eq. (3)). W-MSA and SW-MSA are multi-head self attention modules with regular and shifted windowing configurations, respectively.

# Experiments

## Table 1

### (a) Regular ImageNet-1K trained models

| method | image size | #param. | FLOPs | throughput (image / s) | ImageNet top-1 acc. |
|---|---|---|---|---|---|
| RegNetY-4G [47] | $224^2$ | 21M | 4.0G | 1156.7 | 80.0 |
| RegNetY-8G [47] | $224^2$ | 39M | 8.0G | 591.6 | 81.7 |
| RegNetY-16G [47] | $224^2$ | 84M | 16.0G | 334.7 | 82.9 |
| EffNet-B3 [57] | $300^2$ | 12M | 1.8G | 732.1 | 81.6 |
| EffNet-B4 [57] | $380^2$ | 19M | 4.2G | 349.4 | 82.9 |
| EffNet-B5 [57] | $456^2$ | 30M | 9.9G | 169.1 | 83.6 |
| EffNet-B6 [57] | $528^2$ | 43M | 19.0G | 96.9 | 84.0 |
| EffNet-B7 [57] | $600^2$ | 66M | 37.0G | 55.1 | 84.3 |
| ViT-B/16 [19] | $384^2$ | 86M | 55.4G | 85.9 | 77.9 |
| ViT-L/16 [19] | $384^2$ | 307M | 190.7G | 27.3 | 76.5 |
| DeiT-S [60] | $224^2$ | 22M | 4.6G | 940.4 | 79.8 |
| DeiT-B [60] | $224^2$ | 86M | 17.5G | 292.3 | 81.8 |
| DeiT-B [60] | $384^2$ | 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 | 83.3 |
| Swin-B | $384^2$ | 88M | 47.0G | 84.7 | 84.2 |

### (b) ImageNet-22K pre-trained models

| method | image size | #param. | FLOPs | throughput (image / s) | ImageNet top-1 acc. |
|---|---|---|---|---|---|
| R-101x3 [37] | $384^2$ | 388M | 204.6G | - | 84.4 |
| R-152x4 [37] | $480^2$ | 937M | 840.5G | - | 85.4 |
| ViT-B/16 [19] | $384^2$ | 86M | 55.4G | 85.9 | 84.0 |
| ViT-L/16 [19] | $384^2$ | 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 | 86.0 |
| Swin-L | $384^2$ | 197M | 103.9G | 42.1 | 86.4 |

Table 1. Comparison of different backbones on ImageNet-1K classification. Throughput is measured using the GitHub repository of [65] and a V100 GPU, following [60].

## Table 2

### (a) Various frameworks

| Method | Backbone | $AP^{box}$ | $AP_{50}^{box}$ | $AP_{75}^{box}$ | #param. | FLOPs | FPS |
|---|---|---|---|---|---|---|---|
| Cascade | R-50 | 46.3 | 64.3 | 50.5 | 82M | 739G | 18.0 |
| Mask R-CNN | Swin-T | 50.5 | 69.3 | 54.9 | 86M | 745G | 15.3 |
| ATSS | R-50 | 43.5 | 61.9 | 47.0 | 32M | 205G | 28.3 |
| | Swin-T | 47.2 | 66.5 | 51.3 | 36M | 215G | 22.3 |
| RepPointsV2 | R-50 | 46.5 | 64.6 | 50.3 | 42M | 274G | 13.6 |
| | Swin-T | 50.0 | 68.5 | 54.2 | 45M | 283G | 12.0 |
| Sparse | R-50 | 44.5 | 63.4 | 48.2 | 106M | 166G | 21.0 |
| R-CNN | Swin-T | 47.9 | 67.3 | 52.3 | 110M | 172G | 18.4 |

### (b) Various backbones w. Cascade Mask R-CNN

| | $AP^{box}$ | $AP_{50}^{box}$ | $AP_{75}^{box}$ | $AP^{mask}$ | $AP_{50}^{mask}$ | $AP_{75}^{mask}$ | param | FLOPs | FPS |
|---|---|---|---|---|---|---|---|---|---|
| DeiT-S† | 48.0 | 67.2 | 51.7 | 41.4 | 64.2 | 44.3 | 80M | 889G | 10.4 |
| R50 | 46.3 | 64.3 | 50.5 | 40.1 | 61.7 | 43.4 | 82M | 739G | 18.0 |
| Swin-T | 50.5 | 69.3 | 54.9 | 43.7 | 66.6 | 47.1 | 86M | 745G | 15.3 |
| X101-32 | 48.1 | 66.5 | 52.4 | 41.6 | 63.9 | 45.2 | 101M | 819G | 12.8 |
| Swin-S | 51.8 | 70.4 | 56.3 | 44.7 | 67.9 | 48.5 | 107M | 838G | 12.0 |
| X101-64 | 48.3 | 66.4 | 52.3 | 41.7 | 64.0 | 45.1 | 140M | 972G | 10.4 |
| Swin-B | 51.9 | 70.9 | 56.5 | 45.0 | 68.4 | 48.7 | 145M | 982G | 11.6 |

### (c) System-level Comparison

| Method | mini-val $AP^{box}$ | mini-val $AP^{mask}$ | test-dev $AP^{box}$ | test-dev $AP^{mask}$ | #param. | FLOPs |
|---|---|---|---|---|---|---|
| RepPointsV2* [11] | - | - | 52.1 | - | - | - |
| GCNet* [6] | 51.8 | 44.7 | 52.3 | 45.4 | - | 1041G |
| RelationNet++* [12] | - | - | 52.7 | - | - | - |
| SpineNet-190 [20] | 52.6 | - | 52.8 | - | 164M | 1885G |
| ResNeSt-200* [75] | 52.5 | - | 53.3 | 47.1 | - | - |
| EfficientDet-D7 [58] | 54.4 | - | 55.1 | - | 77M | 410G |
| DetectoRS* [45] | - | - | 55.7 | 48.5 | - | - |
| YOLOv4 P7* [3] | - | - | 55.8 | - | - | - |
| Copy-paste [25] | 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++)* | 58.0 | 50.4 | 58.7 | 51.1 | 284M | - |

Table 2. Results on COCO object detection and instance segmentation. †denotes that additional decovolution layers are used to produce hierarchical feature maps. * indicates multi-scale testing.

## Table 3

| ADE20K Method | Backbone | val mIoU | test score | #param. | FLOPs | FPS |
|---|---|---|---|---|---|---|
| DANet [22] | ResNet-101 | 45.2 | - | 69M | 1119G | 15.2 |
| DLab.v3+ [10] | ResNet-101 | 44.1 | - | 63M | 1021G | 16.0 |
| ACNet [23] | ResNet-101 | 45.9 | 38.5 | - | | |
| DNL [68] | ResNet-101 | 46.0 | 56.2 | 69M | 1249G | 14.8 |
| OCRNet [70] | ResNet-101 | 45.3 | 56.0 | 56M | 923G | 19.3 |
| UperNet [66] | ResNet-101 | 44.9 | - | 86M | 1029G | 20.1 |
| OCRNet [70] | HRNet-w48 | 45.7 | - | 71M | 664G | 12.5 |
| DLab.v3+ [10] | ResNeSt-101 | 46.9 | 55.1 | 66M | 1051G | 11.9 |
| DLab.v3+ [10] | ResNeSt-200 | 48.4 | - | 88M | 1381G | 8.1 |
| SETR [78] | T-Large‡ | 50.3 | 61.7 | 308M | - | - |
| UperNet | DeiT-S† | 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‡ | 51.6 | - | 121M | 1841G | 8.7 |
| UperNet | Swin-L‡ | 53.5 | 62.8 | 234M | 3230G | 6.2 |

Table 3. Results of semantic segmentation on the ADE20K val and test set. † indicates additional deconvolution layers are used to produce hierarchical feature maps. ‡ indicates that the model is pre-trained on ImageNet-22K.

# Thank you