

Vision Transformer

Weihua Zhang

Exported on 02/13/2022

Table of Contents

1 ViT(link https://arxiv.org/abs/2010.11929), AN IMAGE IS WORTH 16X16 WORDS: TRANSFORMERS FOR IMAGE RECOGNITION AT SCALE, google.....	3
2 Swin Transformer: Hierarchical Vision Transformer using ShiftedWindows, Microsoft.....	7

1 ViT(link¹), AN IMAGE IS WORTH 16X16 WORDS: TRANSFORMERS FOR IMAGE RECOGNITION AT SCALE, google.

1. Motivation

1) In vision, attention is either applied in conjunction with convolutional networks, or used to replace certain components of convolutional networks while keeping their overall structure in place. A pure transformer applied directly to sequences of image patches can perform very well on image classification tasks.

2) Explore image recognition at larger scales than the standard ImageNet dataset

2. Model overview

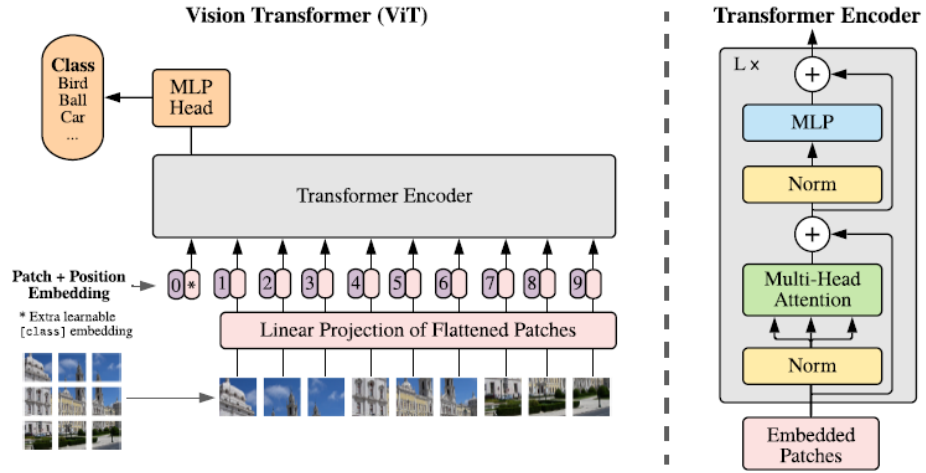


Figure 1: Model overview. We split an image into fixed-size patches, linearly embed each of them, add position embeddings, and feed the resulting sequence of vectors to a standard Transformer encoder. In order to perform classification, we use the standard approach of adding an extra learnable “classification token” to the sequence. The illustration of the Transformer encoder was inspired by Vaswani et al. (2017).

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

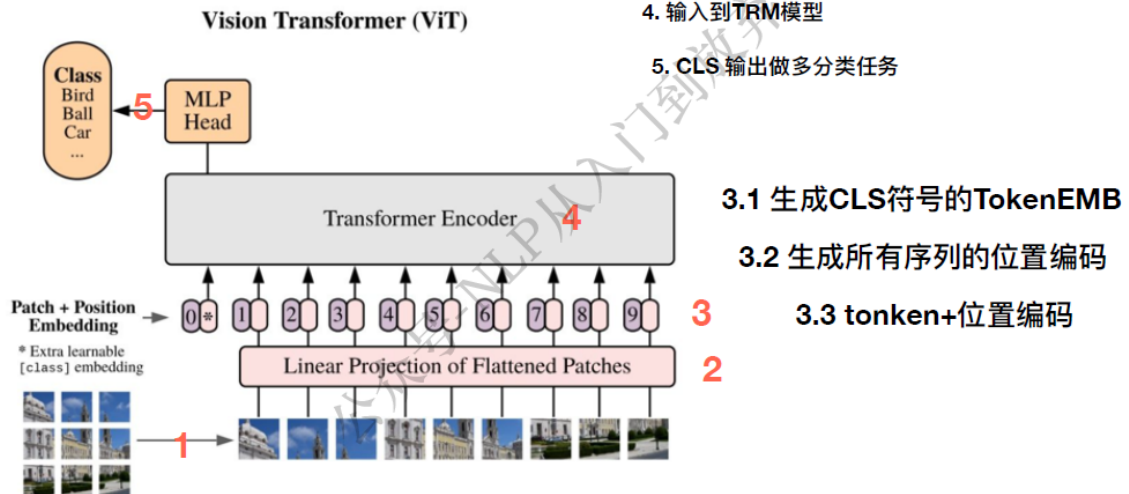
$$\mathbf{z}'_\ell = \text{MSA}(\text{LN}(\mathbf{z}_{\ell-1})) + \mathbf{z}_{\ell-1}, \quad \ell = 1 \dots L \quad (2)$$

$$\mathbf{z}_\ell = \text{MLP}(\text{LN}(\mathbf{z}'_\ell)) + \mathbf{z}'_\ell, \quad \ell = 1 \dots L \quad (3)$$

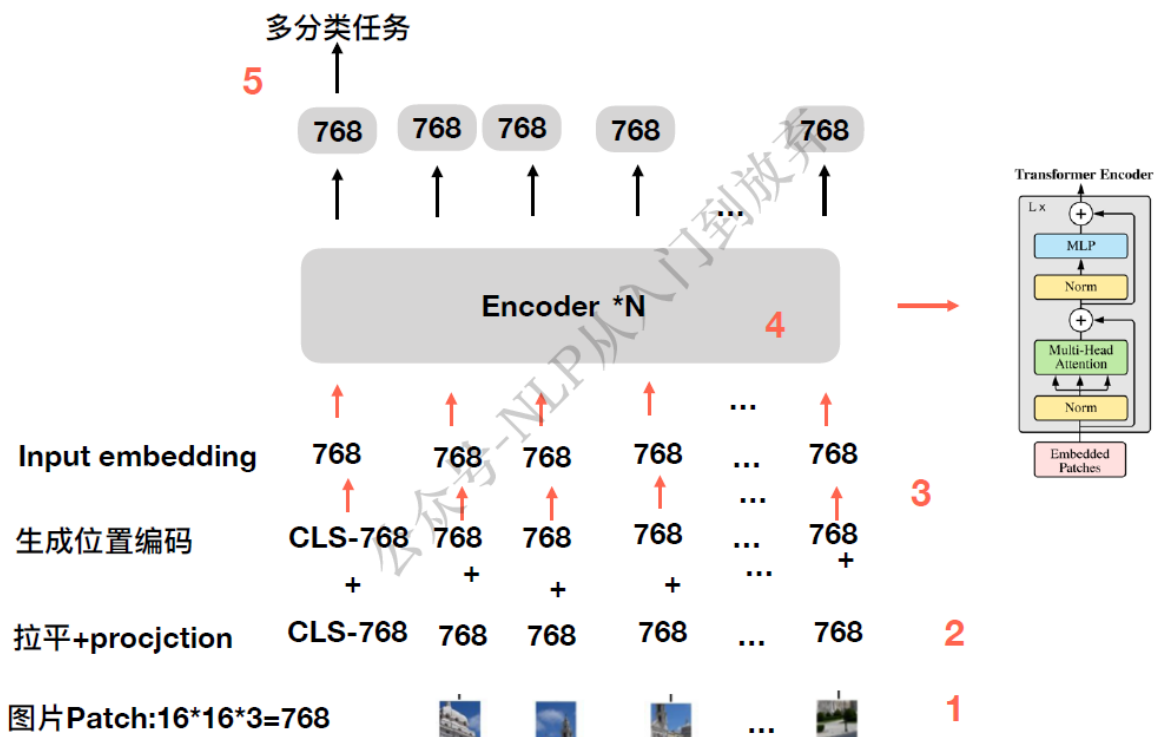
$$\mathbf{y} = \text{LN}(\mathbf{z}_L^0) \quad (4)$$

¹ <https://arxiv.org/abs/2010.11929>

ViT模型架构图：



知乎 @DASOU



3. Novelty

1) How to deal with the computational complexity when the input image is large ?

Other works: a) Self-attention only in local neighborhoods; b) Sparse Transformers;

ViT: Patch Embedding (Split an image into fixed-size patches, linearly embed each of them, add position embeddings, and feed the resulting sequence of vectors to a standard Transformer encoder)

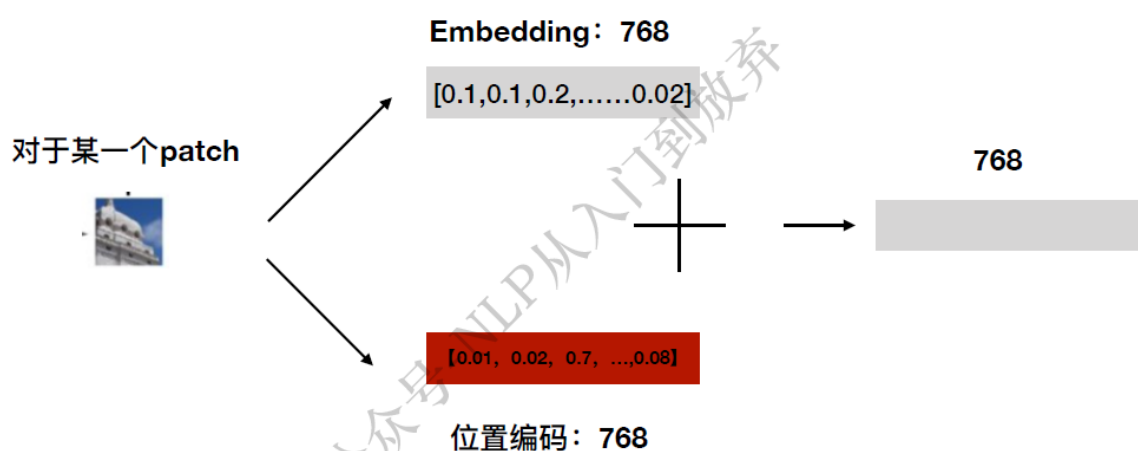
一个简单的改进方式：图像化整为零，切分patch

也就说原来是一个像素点代表一个token，
现在是一大块的token一个patch作为一个token



2) How to add the positional embedding?

Standard learnable 1D position embeddings.



4.Dataset

- 1) ILSVRC-2012 ImageNet dataset with 1k classes and 1.3M images;
- 2) Superset ImageNet-21k with 21k classes and 14M images;
- 3) JFT with 18k classes and 303M high-resolution images;

5.Experiment

	Ours-JFT (ViT-H/14)	Ours-JFT (ViT-L/16)	Ours-I21K (ViT-L/16)	BiT-L (ResNet152x4)	Noisy Student (EfficientNet-L2)
ImageNet	88.55 ± 0.04	87.76 ± 0.03	85.30 ± 0.02	87.54 ± 0.02	88.4/88.5*
ImageNet ReaL	90.72 ± 0.05	90.54 ± 0.03	88.62 ± 0.05	90.54	90.55
CIFAR-10	99.50 ± 0.06	99.42 ± 0.03	99.15 ± 0.03	99.37 ± 0.06	—
CIFAR-100	94.55 ± 0.04	93.90 ± 0.05	93.25 ± 0.05	93.51 ± 0.08	—
Oxford-IIIT Pets	97.56 ± 0.03	97.32 ± 0.11	94.67 ± 0.15	96.62 ± 0.23	—
Oxford Flowers-102	99.68 ± 0.02	99.74 ± 0.00	99.61 ± 0.02	99.63 ± 0.03	—
VTAB (19 tasks)	77.63 ± 0.23	76.28 ± 0.46	72.72 ± 0.21	76.29 ± 1.70	—
TPUv3-core-days	2.5k	0.68k	0.23k	9.9k	12.3k

Table 2: Comparison with state of the art on popular image classification benchmarks. We report mean and standard deviation of the accuracies, averaged over three fine-tuning runs. Vision Transformer models pre-trained on the JFT-300M dataset outperform ResNet-based baselines on all datasets, while taking substantially less computational resources to pre-train. ViT pre-trained on the smaller public ImageNet-21k dataset performs well too. *Slightly improved 88.5% result reported in [Touvron et al. \(2020\)](#).

6. Conclusion:

We have explored the direct application of Transformers to image recognition. Unlike prior works using self-attention in computer vision, we do not introduce any image-specific inductive biases into the architecture. Instead, we interpret an image as a sequence of patches and process it by a standard Transformer encoder as used in NLP. This simple, yet scalable, strategy works surprisingly well when coupled with pre-training on large datasets. Thus, Vision Transformer matches or exceeds the state of the art on many image classification datasets, whilst being relatively cheap to pre-train.

2 Swin Transformer: Hierarchical Vision Transformer using Shifted Windows, Microsoft.

1. Motivation

To address the challenges in adapting Transformer from language to vision:

- a) large variations in the scale of visual entities;
- b) the high resolution of pixels in images compared to words in text.

2. Model overview

Propose a **hierarchical Transformer** whose representation is computed with **Shifted windows**. The shifted windowing scheme brings greater efficiency by **limiting self-attention computation to non-overlapping local windows** while also **allowing for cross-window connection**.

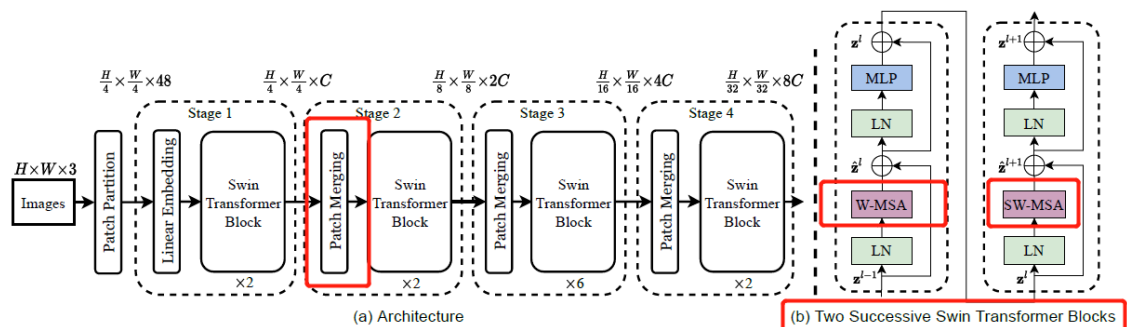
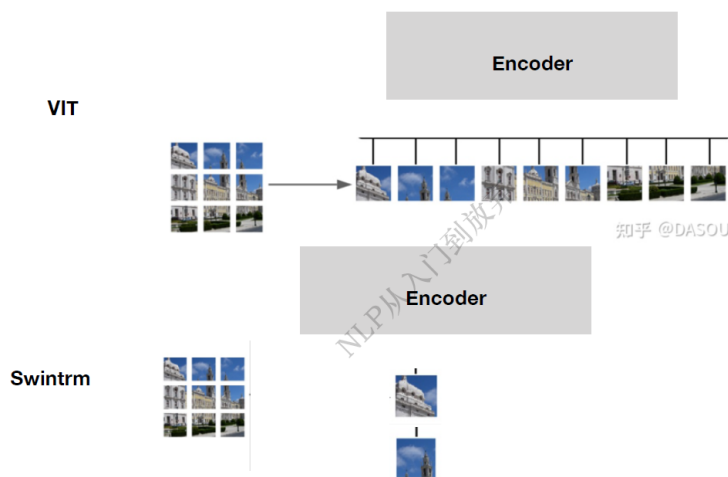


Figure 3. (a) The architecture of a Swin Transformer (Swin-T); (b) two successive Swin Transformer Blocks (notation presented with Eq. (3)). W-MSA and SW-MSA are multi-head self attention modules with regular and shifted windowing configurations, respectively.

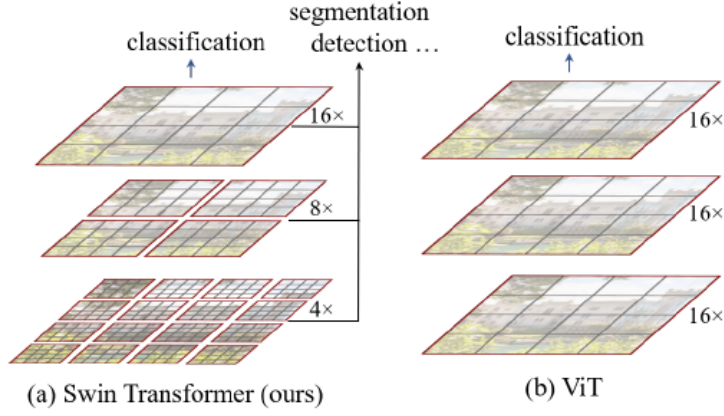


3. Novelty

1) How to deal with the computational complexity when the input image is large ?

A: a) Build hierarchical feature maps by merging image patches in deeper layers (through Patch Merging layer);

Patch Merging: The first patch merging layer concatenates the features of each group of 2×2 neighboring patches, and applies a linear layer on the 4C-dimensional concatenated features. This reduces the number of tokens by a multiple of $2 \times 2 = 4$ (2x downsampling of resolution), and the output dimension is set to $2C$.



	downsp. rate (output size)	Swin-T	Swin-S	Swin-B	Swin-L
stage 1	4× (56×56)	concat 4×4, 96-d, LN win. sz. 7×7, dim 96, head 3 × 2	concat 4×4, 96-d, LN win. sz. 7×7, dim 96, head 3 × 2	concat 4×4, 128-d, LN win. sz. 7×7, dim 128, head 4 × 2	concat 4×4, 192-d, LN win. sz. 7×7, dim 192, head 6 × 2
stage 2	8× (28×28)	concat 2×2, 192-d, LN win. sz. 7×7, dim 192, head 6 × 2	concat 2×2, 192-d, LN win. sz. 7×7, dim 192, head 6 × 2	concat 2×2, 256-d, LN win. sz. 7×7, dim 256, head 8 × 2	concat 2×2, 384-d, LN win. sz. 7×7, dim 384, head 12 × 2
stage 3	16× (14×14)	concat 2×2, 384-d, LN win. sz. 7×7, dim 384, head 12 × 6	concat 2×2, 384-d, LN win. sz. 7×7, dim 384, head 12 × 18	concat 2×2, 512-d, LN win. sz. 7×7, dim 512, head 16 × 18	concat 2×2, 768-d, LN win. sz. 7×7, dim 768, head 24 × 18
stage 4	32× (7×7)	concat 2×2, 768-d, LN win. sz. 7×7, dim 768, head 24 × 2	concat 2×2, 768-d, LN win. sz. 7×7, dim 768, head 24 × 2	concat 2×2, 1024-d, LN win. sz. 7×7, dim 1024, head 32 × 2	concat 2×2, 1536-d, LN win. sz. 7×7, dim 1536, head 48 × 2

Table 7. Detailed architecture specifications.

b) Compute the self-attention within each local window, and allowing for cross-window connection(as Figure 2). In order to compute efficiently, cyclic-shifting is proposed instead of padding.

After this shift, a batched window may be composed of several sub-windows that are not adjacent in the feature map, so a masking mechanism is employed to limit self-attention computation to within each sub-window.

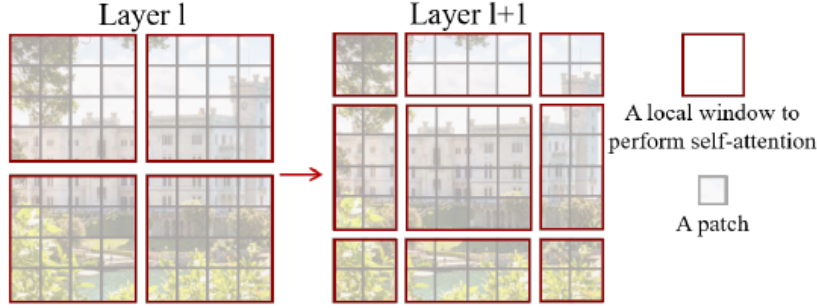


Figure 2. An illustration of the *shifted window* approach for computing self-attention in the proposed Swin Transformer architecture. In layer l (left), a regular window partitioning scheme is adopted, and self-attention is computed within each window. In the next layer $l + 1$ (right), the window partitioning 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.

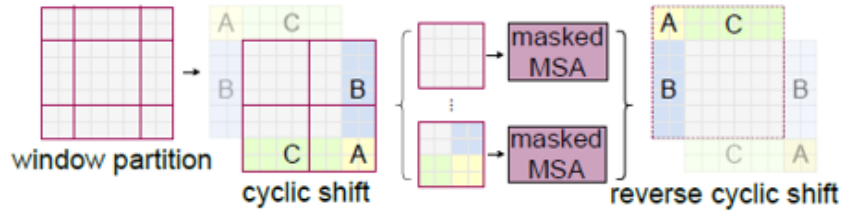


Figure 4. Illustration of an efficient batch computation approach for self-attention in shifted window partitioning.

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 [14]	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.

The window based self-attention is scalable, **linear computational complexity with respect to input image size**.

The windows are arranged to evenly partition the image in a **non-overlapping manner**. Supposing each window contains **$M \times M$ patches**, the computational complexity of a global MSA module and a window based one on an image of $H \times W$ patches are:

$$\Omega(\text{MSA}) = 4hwC^2 + 2(hw)^2C, \quad (1)$$

$$\Omega(\text{W-MSA}) = 4hwC^2 + 2M^2hwC, \quad (2)$$

2) How to deal with the pixel position ?

a) Use the **relative position bias**. b) Put the position bias into the attention multiply.

Relative position bias In computing self-attention, we follow [49, 1, 32, 33] by including a relative position bias $B \in \mathbb{R}^{M^2 \times M^2}$ to each head in computing similarity:

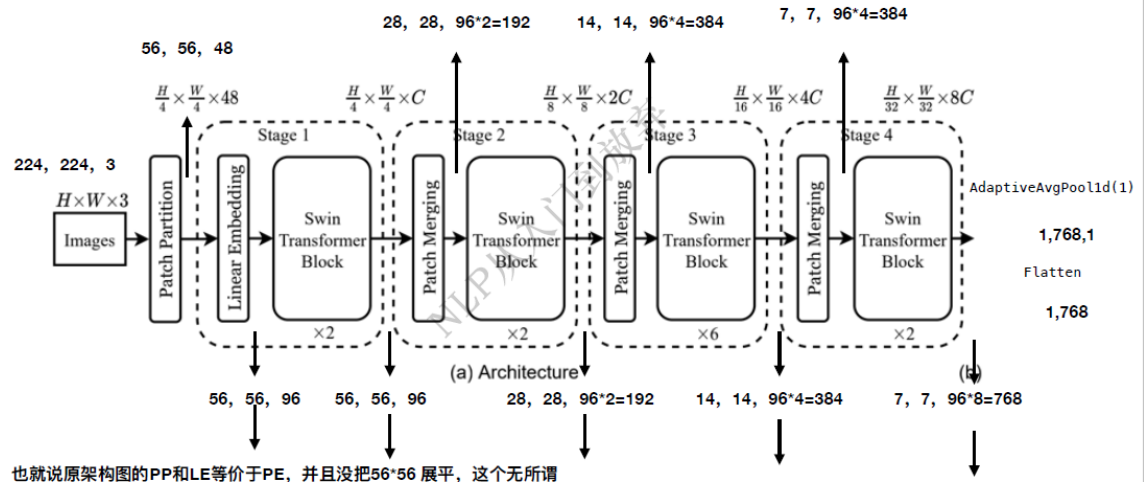
$$\text{Attention}(Q, K, V) = \text{SoftMax}(QK^T / \sqrt{d} + B)V, \quad (4)$$

where $Q, K, V \in \mathbb{R}^{M^2 \times d}$ are the *query*, *key* and *value* matrices; d is the *query/key* dimension, and M^2 is the number of patches in a window. Since the relative position along each axis lies in the range $[-M + 1, M - 1]$, we parameterize a smaller-sized bias matrix $\hat{B} \in \mathbb{R}^{(2M-1) \times (2M-1)}$, and values in B are taken from \hat{B} .

	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).

c) How does Swin Transformer block work?

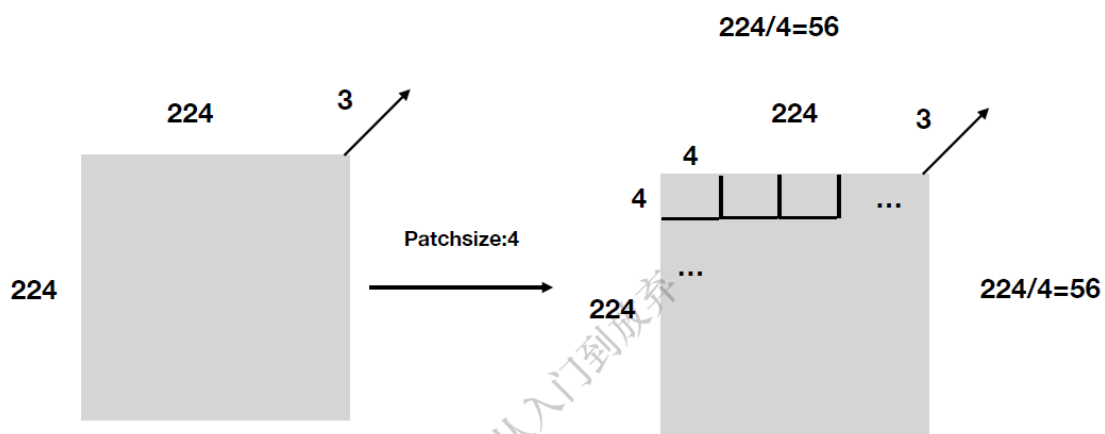


Take 'swin_tiny_patch4_window7_224.yaml' as an example.

```
MODEL:
  TYPE: swin
  NAME: swin_tiny_patch4_window7_224
  DROP_PATH_RATE: 0.2
  SWIN:
    EMBED_DIM: 96
    DEPTHS: [ 2, 2, 6, 2 ]
    NUM_HEADS: [ 3, 6, 12, 24 ]
    WINDOW_SIZE: 7
```

Patch_size: How many pixels does one patch have? (width pixels=height pixels)

Window_size: Fixed. How many patches does one window have? (width patches = height patches)

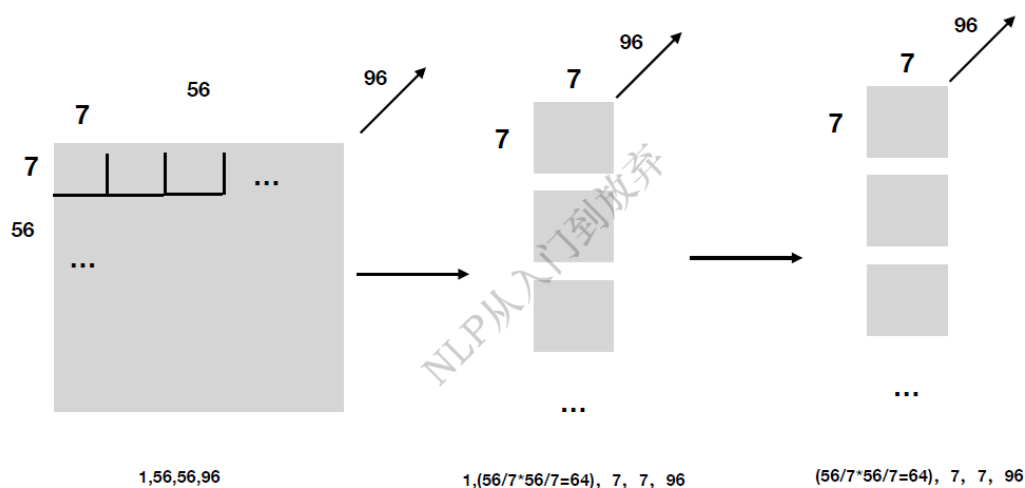


Patch数目就是： $56 \times 56 = 3136$

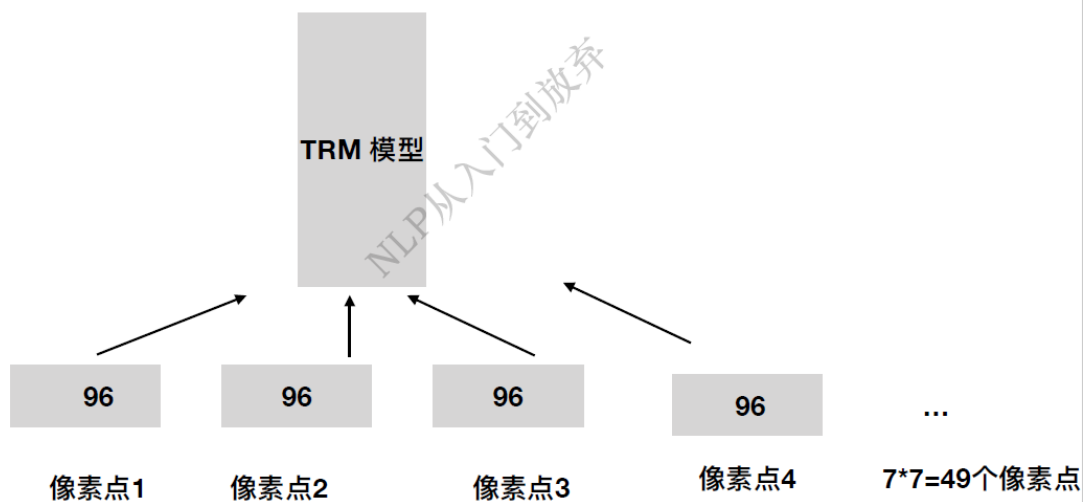
1: $224 \times 224 \times 3$ 到 $(224/4) \times (224/4) \times (4 \times 4 \times 3)$

2. $(224/4) \times (224/4) \times (4 \times 4 \times 3)$ 到 $(224/4) \times (224/4) \times (96)$

切分窗口：win_size=7 一直保持不变



很简单：把每个元素点作为token，96个通道数作为token维度



4.Dataset

Unlike ViT, Imagenet 1k can train a good model (Regular ImageNet-1K training is supported). No need for a huge dataset. But it could work better with a huge dataset.

5.Experiment

SOTA on classification, object detection and segmentation.

(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

Table 1. Comparison of different backbones on ImageNet-1K classification. Throughput is measured using the GitHub repository of [68] and a V100 GPU, following [63].

(a) Various frameworks							
Method	Backbone	AP ^{box}	AP ^{box} ₅₀	AP ^{box} ₇₅	#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 ^{box} ₅₀	AP ^{box} ₇₅	AP ^{mask}	AP ^{mask} ₅₀	AP ^{mask} ₇₅	paramFLOPsFPS
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		test-dev		#param.		FLOPs
	AP ^{box}	AP ^{mask}	AP ^{box}	AP ^{mask}			
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++)*	58.0	50.4	58.7	51.1	284M	-	-

Table 2. Results on COCO object detection and instance segmentation. [†]denotes that additional deconvolution layers are used to upsample feature maps.

ADE20K		val	test	#param.	FLOPs	FPS
Method	Backbone	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 [†]	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.

6. Conclusion

A new vision Transformer which produces a **hierarchical feature** representation and has **linear computational complexity with respect to input image size**. Swin Transformer achieves the state-of-the-art performance on COCO object detection and ADE20K semantic segmentation. **The shifted window based self-attention** is shown to be effective and efficient on vision problems.

7. Code Implementation

Now two mainly implementations, Swin-Transformer and Swin-MLP(use MLP instead of SwinTransformerBlock) structure.

