# Swin Transformer: Hierarchical Vision Transformer using Shifted Windows

Ze Liu, Yutong Lin, Yue Cao, Han Hu et al.

Microsoft Research Asia

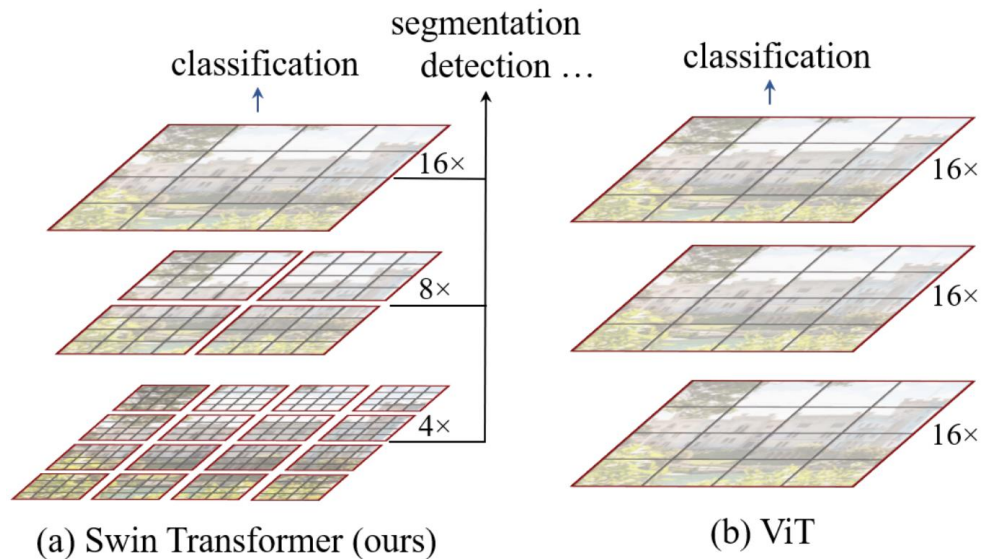ICCV 2021 Best Paper Award

Presented by Minho Park

# Contribution

- **Propose Swin Transformer:** A hierarchical Transformer whose representation is computed with *S*hifted *win*dows.

- **Use stage structure (for FPN).**

- Ablation studies include **relative positional encoding.**

# Shifted Window Approach

- CNN (VGG, ResNet, etc.)

- Global attention (ViTs)

- Local attention (Sliding Window, Axial attention, Swin)

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

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



(a) Swin Transformer (ours)

(b) ViT

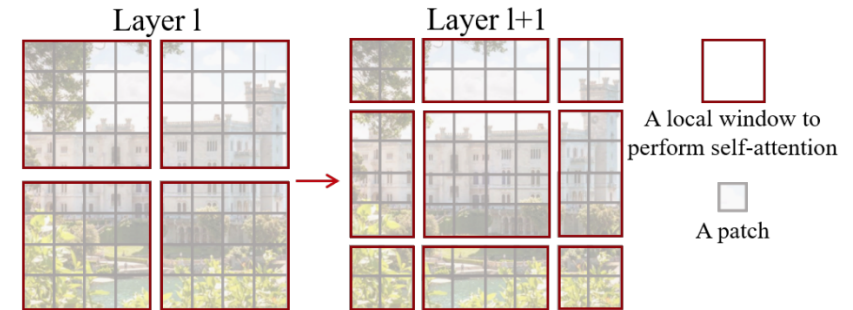Architectures of Swin Transformer and Vision Transformer.



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

- Two successive Swin Transformer Block.



(a) Architecture
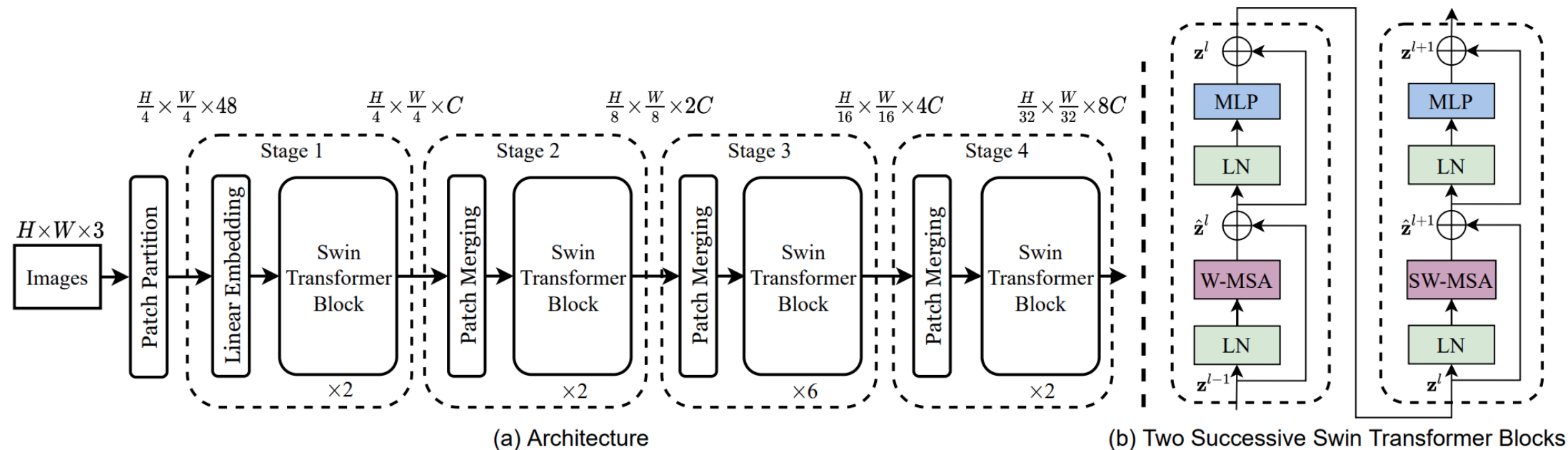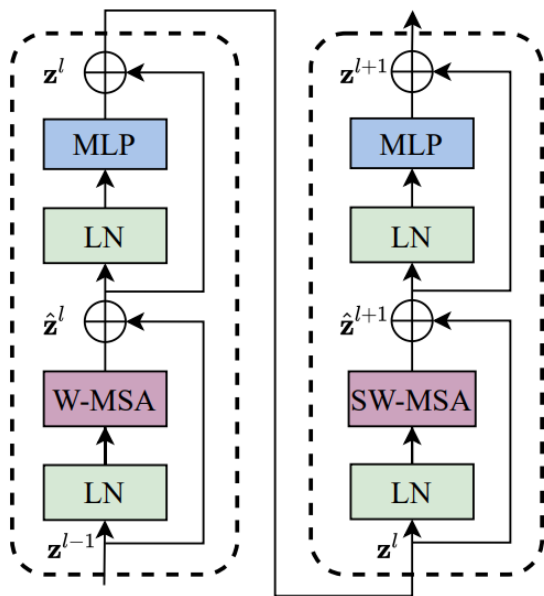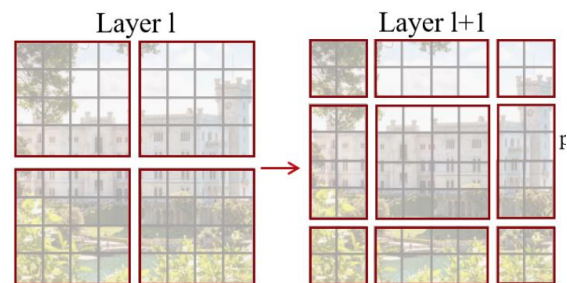
(b) Two Successive Swin Transformer Blocks

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.

# Efficient Batch Computation

- The number of batched windows remains the same as that of regular window partitioning.



Two successive Swin Transformer Block.



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.

# Relative Positional Bias

- Positional biases: absolute, relative

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

- A different window size through bi-cubic interpolation.

| | 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** |

Ablation studies of Swin and different positional biases

# Architecture Variants

- Swin-T, Swin-S, Swin-B, Swin-L: 0.25×, 0.5×, 1×, and 2× the model size and computational complexity, respectively.

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

| layer name | output size | 18-layer | 34-layer | 50-layer | 101-layer | 152-layer |
|---|---|---|---|---|---|---|
| conv1 | 112×112 | 7×7, 64, stride 2 | | | | |
| | | 3×3 max pool, stride 2 | | | | |
| conv2_x | 56×56 | $\begin{bmatrix} 3\times3, 64 \\ 3\times3, 64 \end{bmatrix}\times2$ | $\begin{bmatrix} 3\times3, 64 \\ 3\times3, 64 \end{bmatrix}\times3$ | $\begin{bmatrix} 1\times1, 64 \\ 3\times3, 64 \\ 1\times1, 256 \end{bmatrix}\times3$ | $\begin{bmatrix} 1\times1, 64 \\ 3\times3, 64 \\ 1\times1, 256 \end{bmatrix}\times3$ | $\begin{bmatrix} 1\times1, 64 \\ 3\times3, 64 \\ 1\times1, 256 \end{bmatrix}\times3$ |
| conv3_x | 28×28 | $\begin{bmatrix} 3\times3, 128 \\ 3\times3, 128 \end{bmatrix}\times2$ | $\begin{bmatrix} 3\times3, 128 \\ 3\times3, 128 \end{bmatrix}\times4$ | $\begin{bmatrix} 1\times1, 128 \\ 3\times3, 128 \\ 1\times1, 512 \end{bmatrix}\times4$ | $\begin{bmatrix} 1\times1, 128 \\ 3\times3, 128 \\ 1\times1, 512 \end{bmatrix}\times4$ | $\begin{bmatrix} 1\times1, 128 \\ 3\times3, 128 \\ 1\times1, 512 \end{bmatrix}\times8$ |
| conv4_x | 14×14 | $\begin{bmatrix} 3\times3, 256 \\ 3\times3, 256 \end{bmatrix}\times2$ | $\begin{bmatrix} 3\times3, 256 \\ 3\times3, 256 \end{bmatrix}\times6$ | $\begin{bmatrix} 1\times1, 256 \\ 3\times3, 256 \\ 1\times1, 1024 \end{bmatrix}\times6$ | $\begin{bmatrix} 1\times1, 256 \\ 3\times3, 256 \\ 1\times1, 1024 \end{bmatrix}\times23$ | $\begin{bmatrix} 1\times1, 256 \\ 3\times3, 256 \\ 1\times1, 1024 \end{bmatrix}\times36$ |
| conv5_x | 7×7 | $\begin{bmatrix} 3\times3, 512 \\ 3\times3, 512 \end{bmatrix}\times2$ | $\begin{bmatrix} 3\times3, 512 \\ 3\times3, 512 \end{bmatrix}\times3$ | $\begin{bmatrix} 1\times1, 512 \\ 3\times3, 512 \\ 1\times1, 2048 \end{bmatrix}\times3$ | $\begin{bmatrix} 1\times1, 512 \\ 3\times3, 512 \\ 1\times1, 2048 \end{bmatrix}\times3$ | $\begin{bmatrix} 1\times1, 512 \\ 3\times3, 512 \\ 1\times1, 2048 \end{bmatrix}\times3$ |
| | 1×1 | average pool, 1000-d fc, softmax | | | | |
| FLOPs | | $1.8\times10^9$ | $3.6\times10^9$ | $3.8\times10^9$ | $7.6\times10^9$ | $11.3\times10^9$ |

ResNet Architectures for ImageNet.

He, Kaiming, et al. "Deep residual learning for image recognition." *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2016.

# Quantitative Results

**(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].

# Quantitative Results

**(a) Various frameworks**

| Method | Backbone | $AP^{box}$ | $AP^{box}_{50}$ | $AP^{box}_{75}$ | #param. | FLOPs | FPS |
|---|---|---|---|---|---|---|---|
| Cascade Mask R-CNN | R-50 | 46.3 | 64.3 | 50.5 | 82M | 739G | 18.0 |
| Cascade 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 |
| ATSS | 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 |
| RepPointsV2 | Swin-T | **50.0** | **68.5** | **54.2** | 45M | 283G | 12.0 |
| Sparse R-CNN | R-50 | 44.5 | 63.4 | 48.2 | 106M | 166G | 21.0 |
| Sparse R-CNN | Swin-T | **47.9** | **67.3** | **52.3** | 110M | 172G | 18.4 |

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

| | $AP^{box}$ | $AP^{box}_{50}$ | $AP^{box}_{75}$ | $AP^{mask}$ | $AP^{mask}_{50}$ | $AP^{mask}_{75}$ | 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* [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 decovolution layers are used to produce hierarchical feature maps. * indicates multi-scale testing.

# Quantitative Results

| ADE20K Method | Backbone | val mIoU | test score | #param. | FLOPs | FPS |
|---|---|---|---|---|---|---|
| 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$^{\ddagger}$ | 50.3 | 61.7 | 308M | - | - |
| UperNet | DeiT-S$^{\dagger}$ | 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$^{\ddagger}$ | 51.6 | - | 121M | 1841G | 8.7 |
| UperNet | Swin-L$^{\ddagger}$ | **53.5** | **62.8** | 234M | 3230G | 6.2 |

Table 3. Results of semantic segmentation on the ADE20K val and test set. $^{\dagger}$ indicates additional deconvolution layers are used to produce hierarchical feature maps. $^{\ddagger}$ indicates that the model is pre-trained on ImageNet-22K.

# Paperswithcode

- https://paperswithcode.com/task/object-detection

- https://paperswithcode.com/task/semantic-segmentation