

Swin Transformer

(Shifted Window Transformer)

Limitation of ViT:

1) Quadratic computational complexity,

* Attention complexity scales $O(N^2)$

* N = number of patches.

* For images, $N \propto H \times W$

} \rightarrow Imagine image with 9 patches
} \rightarrow Total attention scores
will be 9×9

Goal of Swin Transformer:

* Make transformers a general purpose backbone for vision task while achieving,

\rightarrow linear computational complexity w.r.t img size

\rightarrow strong performance on classification, detection and segmentation.

* It restricts self attention to local windows.

Architecture:

21p Image ($H \times W \times 3$)



Patch Partition (4×4) → Result: $H/4 \times W/4$ resolution



stage-1

After flattening = $\frac{H}{4} \times \frac{W}{4} \times 48$

Linear Embedding

→ Project to C, Result: $H/4 \times W/4 \times C$



Swan Transformer Block $\times 2$ → Result: $H/4 \times W/4 \times C$



stage-2

Patch Merging → Result: $H/8 \times W/8 \times 2C$



Swan Transformer Block $\times 2$ → Result: $H/8 \times W/8 \times 2C$



stage-3

Patch Merging → Result: $H/16 \times W/16 \times 4C$



Swan Transformer Block $\times 6$ → Result: $H/16 \times W/16 \times 4C$



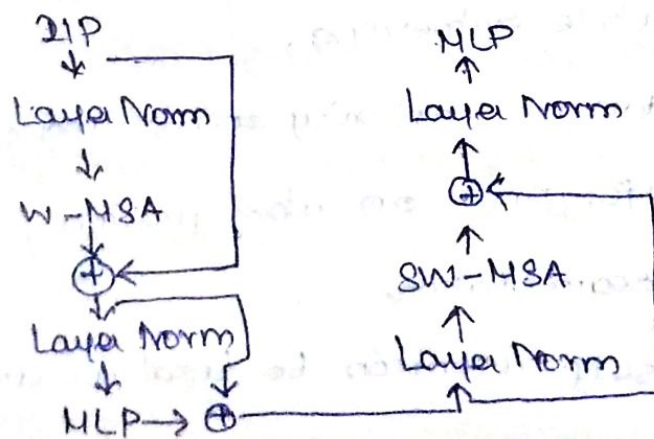
stage-4

Patch Merging → Result: $H/32 \times W/32 \times 8C$



Swan Transformer Block $\times 2$ → Result: $H/32 \times W/32 \times 8C$

Swan Transformer Block:



Patch Partition

Step 1: Image \rightarrow Partition

* DIP image: $H \times W \times 3$

* Patch size: 4×4

* Each patch contains: $4 \times 4 \times 3 = 48$ values.

Step 2: Flattening:

* Each patch \rightarrow 48 dim vector

* Resulting dim $\rightarrow H/4 \times W/4 \times 48$

Linear Embedding:

* Its a linear projection converts $48 \rightarrow C$

* Resulting dim: $H/4 \times W/4 \times C$

Patch Merging:

* It reduces spatial size while increasing channels.

Working:

\rightarrow Take 2×2 groups of neighboring patches

\rightarrow Concatenate to channels: $C + C + C + C = 4C$

\rightarrow Apply linear projection: $4C \rightarrow 2C$

* Result = Height $\downarrow 2$, Width $\downarrow 2$, Channels $\uparrow 2$

Window-Based Self Attention (W-MSA):

* Instead of global attention, divide images into non-overlapping windows and compute self attention within each window only.

\rightarrow Let window size = $N \times N$

\rightarrow No of windows $\approx \frac{H \times W}{N^2}$

Total attention complexity $O\left(\frac{HW}{N^2} \cdot N^2\right)$

Shifted Window Attention (SW-USA):

* Shifting the window in a cyclic manner let patches that move out of one side and re-enter from the opposite side.

* This enables cross window communication end-to-end.

Attention Formula in GPT-4:

Regular Window Attention,

$$\text{Attention} = \text{Softmax}\left(\frac{QK^T}{\sqrt{d}} + B\right)V$$

where,

$B \rightarrow$ relative position bias.

Shifted Window Attention,

$$\text{Attention} = \text{Softmax}\left(\frac{QK^T}{\sqrt{d}} + B + M\right)V$$

where,

- $M = 0$ if the token are same window
- $M = -\infty$ otherwise.

Relative Position Bias,

* Unlike ViT where absolute position embedding was implemented, in such transformer a learnable bias is added

* It is based upon relative displacement $(\Delta x, \Delta y)$

For a 7×7 window,

$$\rightarrow \Delta x, \Delta y \in [-6, +6]$$

\rightarrow Total unique relative position: $13 \times 13 = 169$

* There is only 169 learnable parameters, instead of $49 \times 49 = 2401$

Classification in Swan Transformer: [Although Swan Transformer is General Purpose]

* There is no CLS Token, instead the final stage features are global average pooled and passed to classification head.

— x —