

Iwin Transformer: Hierarchical Vision Transformer using Interleaved Windows

Simin Huo, Ning Li

Abstract—We introduce Iwin Transformer, a novel position-embedding-free hierarchical vision transformer, which can be fine-tuned directly from low to high resolution, through the collaboration of innovative interleaved window attention and depthwise separable convolution. This approach uses attention to connect distant tokens and applies convolution to link neighboring tokens, enabling global information exchange within a single module, overcoming Swin Transformer’s limitation of requiring two consecutive blocks to approximate global attention. Extensive experiments on visual benchmarks demonstrate that Iwin Transformer exhibits strong competitiveness in tasks such as image classification (87.4 top-1 accuracy on ImageNet-1K), semantic segmentation and video action recognition. We also validate the effectiveness of the core component in Iwin as a standalone module that can seamlessly replace the self-attention module in class-conditional image generation. The concepts and methods introduced by the Iwin Transformer have the potential to inspire future research, like Iwin 3D Attention in video generation. The code and models are available at <https://github.com/cominder/Iwin-Transformer>.

Index Terms—Iwin Transformer, interleaved window attention, position-embedding-free.

I. INTRODUCTION

VISION Transformers (ViTs) [1] have transformed computer vision by borrowing the transformer architecture from natural language models [2]. Unlike Convolutional Neural Networks (CNNs) [3], which rely on local receptive fields to capture image features, ViTs leverage self-attention mechanisms to get global dependencies, demonstrating remarkable performance on vision tasks. However, its quadratic computational complexity $\mathcal{O}(N^2)$ with respect to input sequence length N presents significant scalability challenges, particularly for high-resolution image processing applications that are increasingly common in computer vision.

To tackle the challenge of quadratic complexity in Vision Transformers (ViTs) and enhance their efficiency while maintaining performance, various approaches have been proposed. Hierarchical Designs such as PVT [4] and Twins [5] utilize multi-scale feature pyramids to progressively reduce spatial dimensions. Hybrid CNN-Transformer Architectures like ConViT [6] and CoAtNet [7] combine convolutional operations with self-attention to leverage the strengths of both paradigms. Efficient Token Fusion strategies such as TokenLearner [8] dynamically aggregate tokens to reduce sequence length, while Sparse Attention Patterns exemplified by Reformer [9] utilize locality-sensitive hashing to attend only to relevant tokens. Additionally, efficient implementations like Performer [10] approximate attention through kernel methods to achieve linear complexity. Diverse strategies are employed to mitigate the computational demands of vision transformers.

One of the most promising approaches to tackling these challenges is the Swin Transformer [11], which introduces a hierarchical architecture with shifted window-based self-attention. By constraining attention computation within local windows and enabling cross-window connection through a window-shifting mechanism, Swin Transformer successfully reduces the quadratic complexity to linear complexity with respect to image size. This elegant design maintains the model’s capability to capture both local and global dependencies while significantly improving computational efficiency. Moreover, Swin Transformer adopts a hierarchical structure that progressively merges image patches in deeper layers, generating multi-scale feature maps similar to conventional CNN backbones, which facilitates its application across various vision tasks including object detection and semantic segmentation. The impressive performance of Swin Transformer across benchmarks has established it as a milestone in efficient vision transformer design and demonstrated the viability of window-based attention mechanisms for large-scale vision applications.

Despite its pioneering design and impressive performance, the Swin Transformer exhibits several noteworthy limitations. First, the shifted window mechanism introduces extra computational overhead due to the complex masking operations required during attention computation, complicating implementation and reducing hardware efficiency. Second, Swin’s architecture requires two consecutive transformer blocks: one with regular windows and another with shifted windows to achieve global information exchange, resulting in computational redundancy as certain features are processed multiple times. This two-block requirement poses particular challenges in the era of AI generated content (AIGC), where conditioning information such as text prompts must be injected into the model; there is no obvious optimal placement for cross-attention between text and images within this rigid two-block structure, explaining Swin’s limited adoption in modern text-to-image diffusion models. Furthermore, as acknowledged in Swin Transformer v2 [12], the model faces scalability issues when fine-tuned for higher-resolution inputs. The bi-cubic interpolation of relative position encodings for larger windows leads to significant performance degradation, necessitating the introduction of complex alternatives such as log-spaced continuous position bias (Log-CPB). This reliance on sophisticated position encoding schemes ultimately hinders the model’s scaling capabilities and broader applicability.

To address these limitations while preserving the computational efficiency of window-based attention, we introduce the **Interleaved Window Transformer** (Iwin Transformer). The key innovation of Iwin lies in its incorporation of depth-

wise separable convolutions alongside the interleaved window mechanism, which rearranges features before applying window attention such that each window contains pixels from different regions of the image. This elegant approach enables global information interaction in a single transformer block without the complex masking operations required by Swin. Additionally, convolution introduces inductive biases that are beneficial for vision tasks and provides implicit positional information. This hybrid approach not only enhances feature representation but also reduces the reliance on explicit position encodings, addressing a key limitation of Swin. The combined design enables Iwin to achieve an equivalent global receptive field of two consecutive Swin blocks with approximately half the computational cost, making it particularly advantageous for high-resolution vision applications and more amenable to integration with text-conditioning mechanisms in generative models.

The primary contributions of this work are summarized as follows:

- 1) **Interleaved window attention:** We propose a novel Reshape-Transpose-Reshape (RTR) operation that systematically reorders feature sequences into an interleaved pattern for applying window self-attention, and then restores the original spatial arrangement. This mechanism realized linear complexity.
- 2) **Hybrid attention-convolution module:** We elegantly integrate depthwise separable convolutions with interleaved window attention to create a computationally efficient module that leverages the complementary strengths of both mechanisms.
- 3) **Theoretical analysis:** We provide mathematical proof that Iwin achieves global information exchange through hybrid attention-convolution module.
- 4) **Position-embedding-free:** Iwin Transformer does not require explicit position coding, ensuring its strong scalability across varying input resolutions without performance degradation, overcoming a key limitation in previous transformer architectures.
- 5) **Comprehensive empirical validation:** We provide extensive experimental evidence showing that Iwin maintains or improves upon the performance of Swin across various vision tasks including image classification, semantic segmentation and video recognition, proving its effectiveness as a vision backbone.
- 6) **Extensibility to other domains:** In our discussion, we present clues to extend Iwin's interleaved window attention to 1D for large language models and to 3D for video generation, offering a third alternative to conventional 3D full attention and spatial-temporal attention mechanisms, with potential benefits for computational efficiency.

II. RELATED WORKS

Inspired by the success of Vision Transformers (ViTs) [1], Transformer architectures have drawn significant attention in computer vision research [13]–[18]. While they all encounter one common problem, which is the heavy computational overhead of quadratic complexity. Various approaches have been

proposed to enhance the efficiency of transformer structure while maintaining good performance. Some works [6], [7], [19] integrate CNNs and Transformers to leverage the advantages of both structures. Other approaches [4], [11], [20] focus on modifying the structure of ViTs to better suit vision tasks. The following subsections briefly review these related works categorized by their methodologies.

A. Linear and Sparse Attention

Self-attention operations introduce quadratic computational complexity with respect to sequence length, presenting significant challenges for high-resolution visual inputs. Linformer [21] achieved linear complexity through low-rank factorization of attention matrices, decomposing the $N \times N$ attention matrix into two smaller matrices. Performer [10] introduced Fast Attention Via positive Orthogonal Random features (FAVOR+), using random feature maps to approximate the attention kernel. Luna [22] proposed linear unified nested attention by introducing a set of fixed-length projected embeddings that serve as an intermediate representation for attention computation. BigBird [23] combined random, window, and global attention patterns to maintain linear complexity. Longformer [24] employed dilated sliding window attention with select global attention tokens. Sparse Transformer [25] introduced factorized attention patterns that reduce complexity through structured sparsity. [26], [27] selectively compute attention for the most relevant token pairs. These approaches collectively demonstrate that full global attention is often unnecessary for effective visual representation learning, enabling more efficient transformer designs without substantial performance degradation.

B. Hierarchical Vision Transformers

Hierarchical vision transformers adopt multi-scale feature representations to enhance computational efficiency while maintaining modeling capabilities. Pyramid Vision Transformer (PVT) [4] introduces a progressive shrinking pyramid structure that reduces sequence length at deeper layers through spatial-reduction attention. Swin Transformer [11] presents a hierarchical architecture with shifted windows, restricting self-attention computation to local windows and establishing cross-window connections through window shifting between layers. This design reduces computational complexity from quadratic to linear with respect to image size. MViT [20] employs pooling attention that progressively expands channel capacity while reducing spatial resolution. Twins combines locally-grouped self-attention with global sub-sampled attention to balance local and global interactions efficiently. CSWin [28] utilizes cross-shaped window self-attention to capture horizontal and vertical dependencies separately. These hierarchical designs have proven particularly effective for dense prediction tasks like object detection and semantic segmentation, where multi-scale feature representations are essential.

C. Hybrid CNN-Transformer Architectures

Hybrid architectures combine convolutional operations with self-attention to leverage the strengths of both paradigms. Con-

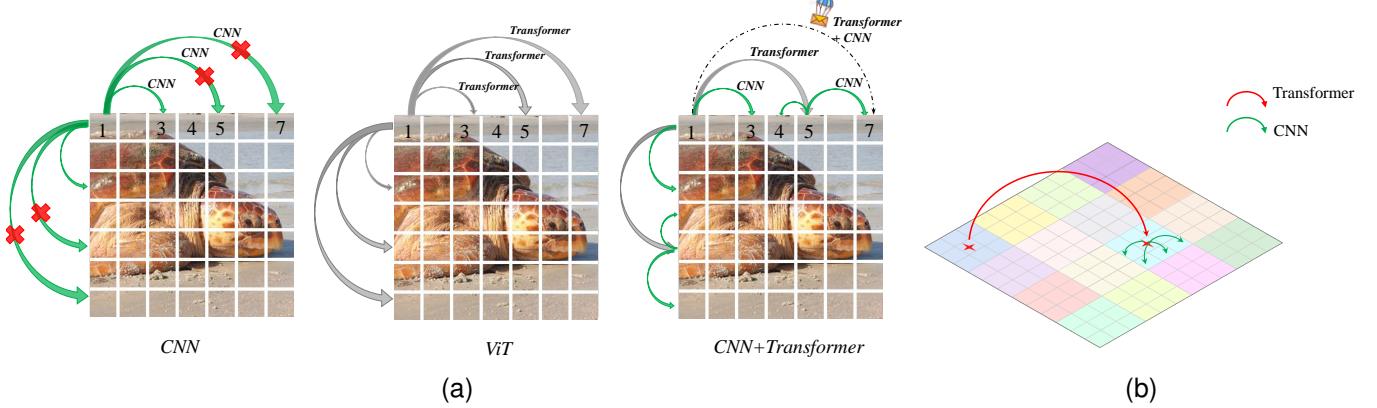


Fig. 1. Diagram of the proposed pattern. In (a), token 1 within the CNN can only interact with token 3 nearby and cannot reach token 7 over a long distance. Therefore, CNN is restricted to capturing local features. In contrast, token 1 in the ViT can be associated with any token, enabling the capture of global features but with a quadratic complexity. In the third proposed CNN+Transformer pattern, token 1 first connects with token 5 at a short distance through attention, and token 5 is related to token 7 via convolution. In this way, tokens 1 and 7, despite being far away, communicate indirectly. In (b) shows an intuitive top view of the proposed CNN+Transformer pattern.

ViT [6] incorporates a gated positional self-attention mechanism that transitions smoothly from convolution to transformer behavior. CoAtNet [7] unifies depthwise convolutions with self-attention in a relative attention framework, employing convolutions in earlier stages and self-attention in later stages. MobileViT [19] integrates the local processing of convolutions with the global processing of transformers for lightweight models. LocalViT [29] enhances vision transformers with depth-wise convolutions, introducing locality to self-attention layers. These hybrid approaches typically achieve better parameter efficiency and performance on smaller datasets compared to pure transformer architectures, while maintaining the global modeling capacity essential for complex vision tasks.

D. Dynamic Computation Strategies

Dynamic computation strategies adjust computational resources based on input complexity. DynamicViT [30] introduces a token sparsification framework that progressively prunes redundant tokens based on their importance scores, reducing sequence length as the network deepens. A-ViT [31] employs adaptive computation depth, allowing different tokens to exit the network at different layers based on their complexity. Token Merging (ToMe) [32] dynamically merges similar tokens throughout the network, preserving information while reducing sequence length. Adaptive Token Sampling [33] introduces learnable modules that sample important tokens based on the input, preserving critical information while reducing computation. These approaches enable more efficient processing by allocating computational resources where they are most needed, making them particularly suitable for real-time applications with varying input complexities.

E. Difference With Previous Works

Unlike Swin Transformer [11], which requires two consecutive blocks with regular and shifted window patterns to establish cross-window connections, Iwin Transformer achieves global information exchange within a single block through a synergistic combination of interleaved window attention and

depthwise separable convolution. Any token can interact with others through an intermediary connection, similar to a flat organization where individuals can reach others without going through multiple people. Another benefit is the module composed of interleaved window attention and depthwise separable convolution can seamlessly replace the standard attention module in generative models without affecting subsequent cross-attention operations with text conditions. Swin can not do it due to its two-block dependency.

In contrast to previous hybrid CNN-Transformer architectures like LocalViT [29] and MobileViT [19], which typically use convolutions merely to capture local features as the motivation. In Iwin Transformer, depthwise separable convolution and interleaved window attention are interdependent components forming a cohesive information processing unit. Depthwise separable convolutions establish connections between certain tokens that are not established through interleaved window attention. Another advantage is that, because convolutions naturally carry position information, Iwin no longer requires explicit position encoding. This makes Iwin a position-embedding-free transformer, which allows models trained on low resolution to be easily fine-tuned to high resolution while maintaining performance.

In summary, Iwin is a position-embedding-free transformer realizing global information exchange through a synergistic combination of interleaved window attention and depthwise separable convolution.

III. METHODOLOGY

A. Overall Architecture

An overview of the Iwin Transformer architecture is presented in Figure 7. The detailed configuration is shown in Table III. Iwin Transformer follows a hierarchical architecture similar to Swin Transformer [11], progressively reduces spatial resolution while expanding channel dimensions across four stages. Given an input image, Iwin first splits it into non-overlapping patches by a patch splitting module. Each patch is treated as a token or feature, then the architecture processes

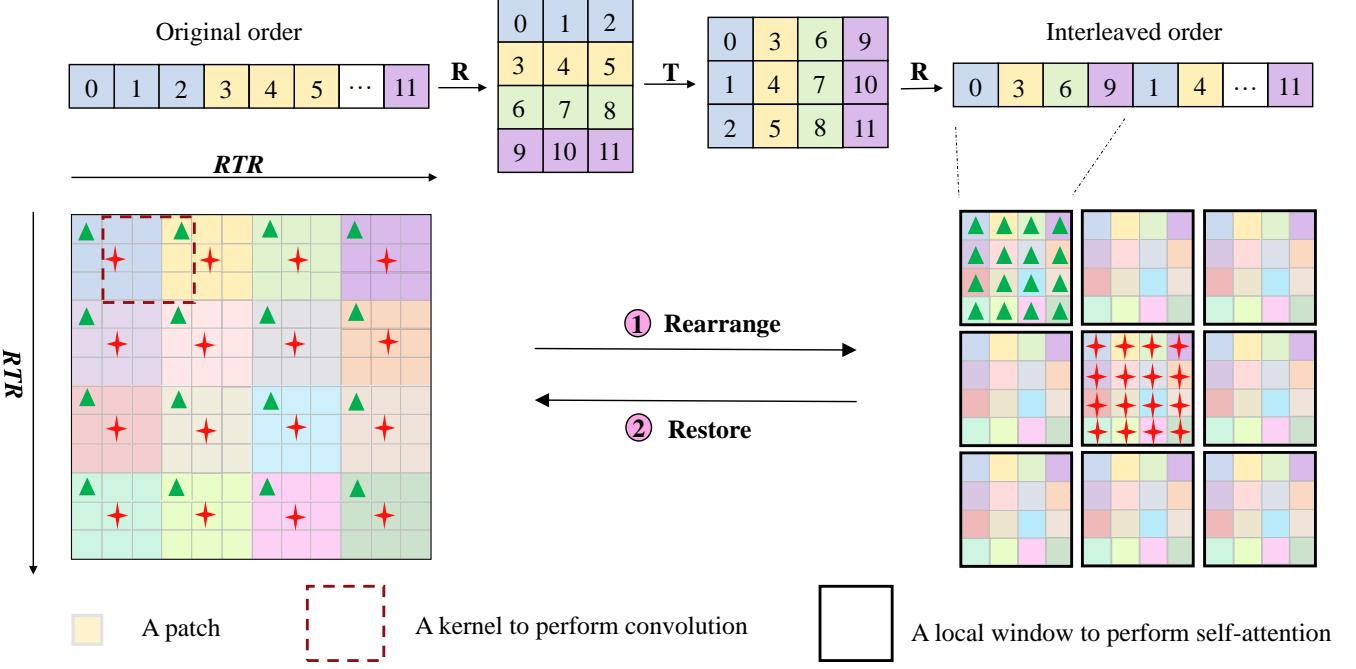


Fig. 2. Illustration of Iwin attention. In the left image, the green triangles and red stars representing tokens are connected through convolutions in the original image. In the right image, all green triangles representing tokens are assigned to the same window through the RTR (Reshape-Transpose-Reshape) operation and window segmentation, executing window attention to establish connections among them. All red stars representing tokens do the same thing. The result is that global convolution and window attention on the interleaved sequence work together to effectively approximate standard global attention, which means that connections are established between any tokens in the original image.

features through stages $\{S_1, S_2, S_3, S_4\}$ with resolutions $\{H/4, H/8, H/16, H/32\}$ and channels $\{C, 2C, 4C, 8C\}$ respectively. Although pyramid structures were used in this study, the core modules of Iwin can be used in flat structures.

B. Interleaved Window Attention

Interleaved Window Attention (IWA), as shown in Figure 2, is the core innovation of the Iwin Transformer. Unlike standard window attention, which evenly divides the feature map into non-overlapping windows, IWA rearranges the feature map before window partition such that tokens from different regions are grouped into the same window for attention computation.

For a feature map $\mathbf{X} \in \mathbb{R}^{H \times W \times C}$, where H and W are the spatial dimensions and C is the channel dimension, IWA operates as follows:

- 1) **Rearrange:** The input feature map is rearranged such that tokens from different regions are grouped into the same window
- 2) **Attention:** Standard multi-head self-attention is applied to the rearranged tokens
- 3) **Restore:** The tokens are restored to their original spatial arrangement

Supposing window size $M \times M$, we first divide the feature map into non-overlapping $M \times M$ windows containing $H_g \times W_g$ tokens, where $H_g = H/M$ and $W_g = W/M$ are the number of windows along the height and width dimensions, respectively.

1) **Rearrange:** The rearrangement can be expressed as: for a token at position (i, j) , its new position in the rearranged feature map is:

$$\begin{aligned} i' &= (i \bmod H_g) \times M + \lfloor i/H_g \rfloor \\ j' &= (j \bmod W_g) \times M + \lfloor j/W_g \rfloor \end{aligned} \quad (1)$$

Luckily, we can elegantly achieve this process through the RTR (Reshape-Transpose-Reshape) operation as clearly shown in Figure 2 and Algorithm 1.

Afterwards, the new feature map is evenly divided into $H_g \times W_g$ non-overlapping windows containing $M \times M$ tokens.

For a token originally at position (i, j) , it will be assigned to the window represented by $(\text{window_row}, \text{window_col})$:

$$\begin{aligned} \text{window_row} &= \lfloor i'/M \rfloor = \lfloor ((i \bmod H_g) \times M + \lfloor i/H_g \rfloor)/M \rfloor \\ &= \lfloor (i \bmod H_g) + \lfloor i/H_g \rfloor / M \rfloor \\ &= i \bmod H_g \quad (\text{since } \lfloor i/H_g \rfloor / M < 1) \end{aligned} \quad (2)$$

Note that $i \bmod H_g$ ranges from 0 to $H_g - 1$, which exactly corresponds to the row index of the window in the grid of windows. Similarly, $\text{window_col} = j \bmod W_g$ ranges from 0 to $W_g - 1$, corresponding to the column index of the window.

This ensures that tokens with the same $(i \bmod H_g, j \bmod W_g)$ are grouped into the same window for attention computation. That means, if two tokens at (i_1, j_1) and (i_2, j_2) are in the same window, they must satisfy:

$$i_1 \bmod H_g = i_2 \bmod H_g \text{ and } j_1 \bmod W_g = j_2 \bmod W_g \quad (3)$$

2) *Self-Attention*: Within each window, we apply the standard self-attention mechanism [2]:

$$\begin{aligned} \mathbf{Q} &= \mathbf{XW}_Q \quad \mathbf{K} = \mathbf{XW}_K \quad \mathbf{V} = \mathbf{XW}_V \\ \text{Attention}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) &= \text{Softmax}\left(\frac{\mathbf{QK}^T}{\sqrt{d_k}}\right) \mathbf{V} \end{aligned} \quad (4)$$

where \mathbf{W}_Q , \mathbf{W}_K , and \mathbf{W}_V are learnable projection matrices, and d_k is the dimension of the key vectors.

3) *Restore*: Finally, the tokens are restored to their original spatial arrangement using the inverse RTR, which is also a Reshape-Transpose-Reshape operation as shown in Figure 2 and Algorithm 1 to realize $(i', j') \mapsto (i, j)$:

$$\begin{aligned} i &= (i' \bmod M) \times H_g + \lfloor i'/M \rfloor \\ j &= (j' \bmod M) \times W_g + \lfloor j'/M \rfloor \end{aligned} \quad (5)$$

Algorithm 1 Pseudocode of **Rearrange** and **Restore** Operations in a PyTorch-like style.

```
# H_num_win: H // window_size, which means the number of
# windows along H
# W_num_win: w // window_size, which means the number of
# windows along W

def rearrange(x, H_num_win, W_num_win):
    B, H, W, C = x.shape

    x = x.reshape(B, -1, H_num_win, W, C).transpose(1, 2)
    x = x.reshape(B, -1, W, C)

    x = x.reshape(B, H, -1, W_num_win, C).transpose(2, 3)
    x = x.reshape(B, H, -1, C)
    return x

def restore(x, H_num_win, W_num_win):
    B, H, W, C = x.shape

    x = x.reshape(B, H, W_num_win, -1, C).transpose(2, 3)
    x = x.reshape(B, H, -1, C)

    x = x.reshape(B, H_num_win, -1, W, C).transpose(1, 2)
    x = x.reshape(B, -1, W, C)
    return x
```

C. Depthwise Separable Convolution

Depthwise Separable Convolution (DWConv) [34] is used to help build missing relationships of certain tokens which are not in the same attention window and provide implicit positional information by the way.

D. Downsampling Layer

Iwin Transformer employs standard convolution to progressively reduce spatial resolution and increase channel dimensions, following the hierarchical design principle [35] common in vision backbones.

$$\mathcal{D}(\mathbf{X}) = \text{Conv}_{3 \times 3, \text{stride}=2}(\mathbf{X}) \quad (6)$$

We tested four downsampling methods as shown in Table VII: average pooling, patch merging, standard convolution, and depthwise separable convolution. They all worked very well, with only a 0.2% performance difference between them. We chose standard convolution, which had the highest accuracy.

E. Iwin Transformer Block

As shown in Figure 3a, the Iwin Transformer block consists of a unified module integrating Interleaved Window Multi-Head Self-Attention (IW-MSA) and Depthwise Separable Convolution (DWConv) modules in parallel, followed by a two-layer MLP with GELU [36] activation between layers. A LayerNorm layer precedes each unified module and MLP, with residual connections after each. The forward pass of the Iwin Transformer block can be formulated as:

$$\begin{aligned} \mathbf{X}' &= \text{LayerNorm}(\mathbf{X}) \\ \mathbf{X}'' &= \mathbf{X} + \text{IW-MSA}(\mathbf{X}') + \text{DWConv}(\mathbf{X}') \\ \mathbf{X}''' &= \mathbf{X}'' + \text{MLP}(\text{LayerNorm}(\mathbf{X}'')) \end{aligned} \quad (7)$$

The computation cost of the unified module combining IW-MSA and DWConv as follows:

$$\begin{aligned} \mathcal{O}_{Iwin} &= \underbrace{\frac{HW}{M^2} \times 3M^2C^2}_{\text{QKV projection}} + \underbrace{2\frac{HW}{M^2} \times M^4C}_{\text{Attention computation}} \\ &\quad + \underbrace{\frac{HW}{M^2} \times M^2C^2}_{\text{Output projection}} + \underbrace{HW \times C \times k^2}_{\text{Convolution}} \\ &= 4HWC^2 + (2M^2 + k^2)HWC \end{aligned} \quad (8)$$

Compared to Swin

$$\mathcal{O}_{Swin} = 4HWC^2 + 2M^2HWC \quad (9)$$

Although Iwin introduces additional k^2HWC computations compared to Swin, this is well worth it. This allows Iwin's layer configuration to be $\{2, 2, 7, 2\}$, whereas Swin relies on two consecutive blocks to approach global attention, and can only increase its depth from $\{2, 2, 6, 2\}$ to $\{2, 2, 8, 2\}$. Therefore, Iwin has great flexibility. Furthermore, when $M \gg k$, M dominates the formula 8, so the computational complexity of Iwin and Swin is nearly the same. In addition, the unified module as a standalone module can seamlessly replace the self-attention module in some generation models.

F. Architecture Variants

Referring to the Swin Transformer [11], we build Iwin-T, Iwin-S, Iwin-B, and Iwin-L with identical network depth and width as Swin for fair comparison. And it is also supported by our ablation studies on Iwin-T, which revealed that the layer configuration $\{2, 2, 6, 2\}$ achieves highest accuracy (see Table VII). Therefore, the subsequent Iwin-S, Iwin-B, and Iwin-L all follow the same settings as Swin. For input resolutions of 224, 384, 512, and 1024, window sizes are 7, 12, 16, and 16 respectively. Table I shows the model size, computational complexity (FLOPs), and performance of different variants on ImageNet.

G. Global Information Exchange

We believe that global information exchange is achieved when there is a path in the feature map along which information can flow from one position (i_1, j_1) to another (i_2, j_2) .

A key theoretical property of the Iwin Transformer is its ability to achieve global information exchange with linear computational complexity. We analyze this property by examining the information flow between any two positions in the feature map.

Lemma 1 (Modular Property of Interleaved Window Attention): In interleaved window attention, tokens at positions (i_1, j_1) and (i_2, j_2) are in the same attention window if and only if:

$$i_1 \bmod H_g = i_2 \bmod H_g \text{ and } j_1 \bmod W_g = j_2 \bmod W_g$$

Lemma 2 (Locality of Depthwise Separable Convolution): For depthwise separable convolution with kernel size $K \times K$, tokens at positions (i_1, j_1) and (i_2, j_2) can directly exchange information if and only if:

$$|i_1 - i_2| \leq K \text{ and } |j_1 - j_2| \leq K$$

Based on these lemmas, we prove the following theorem:

Theorem 3 (Global Information Exchange Condition): If $KM \geq \max(H, W)$, where K is kernel size and M is window size, then the combination of interleaved window attention and depthwise separable convolution in the Iwin Transformer block enables information exchange between any two positions (i_1, j_1) and (i_2, j_2) in the feature map.

Consider arbitrary two positions (i_1, j_1) and (i_2, j_2) in the feature map. We need to prove that there exists a path for information to flow from (i_1, j_1) to (i_2, j_2) .

We discuss three cases:

Case 1: $(i_1 \bmod H_g = i_2 \bmod H_g) \text{ and } (j_1 \bmod W_g = j_2 \bmod W_g)$

In this case, according to Lemma 1, positions (i_1, j_1) and (i_2, j_2) are in the same attention window, so they can directly exchange information through the attention mechanism.

Case 2: $(|i_1 - i_2| \leq K \text{ and } |j_1 - j_2| \leq K)$

In this case, according to Lemma 2, positions (i_1, j_1) and (i_2, j_2) are in the same convolution kernel, so they can directly exchange information through the convolution mechanism.

Case 3: Otherwise (i.e., when (i_1, j_1) and (i_2, j_2) are not in the same attention window and convolution kernel)

In this case, we need to find an intermediate position (i_3, j_3) to bridge (i_1, j_1) and (i_2, j_2) .

We construct such a position (i_3, j_3) as follows:

$$\begin{aligned} i_3 &= (i_1 \bmod H_g) + H_g \cdot \lfloor i_2 / H_g \rfloor \\ j_3 &= (j_1 \bmod W_g) + W_g \cdot \lfloor j_2 / W_g \rfloor \end{aligned} \quad (10)$$

Now we have

$$\begin{aligned} i_3 \bmod H_g &= ((i_1 \bmod H_g) + H_g \cdot \lfloor i_2 / H_g \rfloor) \bmod H_g \\ &= (i_1 \bmod H_g) \bmod H_g + (H_g \cdot \lfloor i_2 / H_g \rfloor) \bmod H_g \\ &= (i_1 \bmod H_g) + 0 \\ &= i_1 \bmod H_g \end{aligned} \quad (11)$$

Similarly, $j_3 \bmod W_g = j_1 \bmod W_g$. This means that positions (i_1, j_1) and (i_3, j_3) are in the same attention window according to Lemma 1.

Now we check:

$$\begin{aligned} |i_2 - i_3| &= |i_2 - ((i_1 \bmod H_g) + H_g \cdot \lfloor i_2 / H_g \rfloor)| \\ &= |i_2 - (i_1 \bmod H_g) - H_g \cdot \lfloor i_2 / H_g \rfloor| \\ &= |i_2 \bmod H_g - i_1 \bmod H_g| \\ &\leq H_g - 1 \\ &= H/M - 1 \\ &< H/M \end{aligned} \quad (12)$$

Similarly, $|j_2 - j_3| < W/M$.

When $KM \geq \max(H, W)$, we have $|i_2 - i_3| \leq K$ and $|j_2 - j_3| \leq K$, so positions (i_3, j_3) and (i_2, j_2) can directly exchange information through depthwise separable convolution.

Therefore, the Iwin Transformer block enables information exchange between any two positions in the feature map when $KM \geq \max(H, W)$. There always exists (i_3, j_3) such that

$$\begin{aligned} i_1 \bmod H_g &= i_3 \bmod H_g \\ j_1 \bmod W_g &= j_3 \bmod W_g \\ |i_2 - i_3| &\leq K \\ |j_2 - j_3| &\leq K \end{aligned} \quad (13)$$

This means that positions (i_1, j_1) and (i_3, j_3) are connected through interleaved window attention, while positions (i_3, j_3) and (i_2, j_2) are connected through depthwise separable convolution, and (i_1, j_1) and (i_2, j_2) established a connection through (i_3, j_3) as an intermediary bridge.

At first, we followed the rule $KM \geq \max(H, W)$, assigning different convolution kernel sizes to each stage. However, ablation experiments (see in Table VII) on Iwin-T showed that kernel sizes of 7, 5, and 3 for stages 1, 2, and 3 yielded the worst performance, with lowest accuracy and slowest training speed, compared to using a consistent kernel size of 7, 5, or 3 across stages. This indicates consistent kernel sizes lead to faster training and better optimization. This aligns with observations in [37] that balanced network outperform theoretically optimal but imbalanced configurations. We think that as the network deepens and downsampling increases the effective receptive field (ERF) [38], resulting in $K_{ERF} \cdot M \geq \max(H, W)$.

Therefore, we believe that if the network is deep enough, after sufficient consecutive Iwin Transformer blocks, the initially small kernel size can expand to large enough so that $K_{ERF} \cdot M \geq \max(H, W)$ at a certain depth and beyond, and here, the model see the whole world.

IV. EXPERIMENTS

We conduct experiments on ImageNet-1K image classification [39], COCO object detection [40], ADE20K semantic segmentation [41], Kinetics-400 [42] video recognition and class-conditional image generation. In the following, we first compare the proposed Iwin Transformer with the previous state-of-the-arts. Then, we ablate the important design elements of Iwin Transformer.

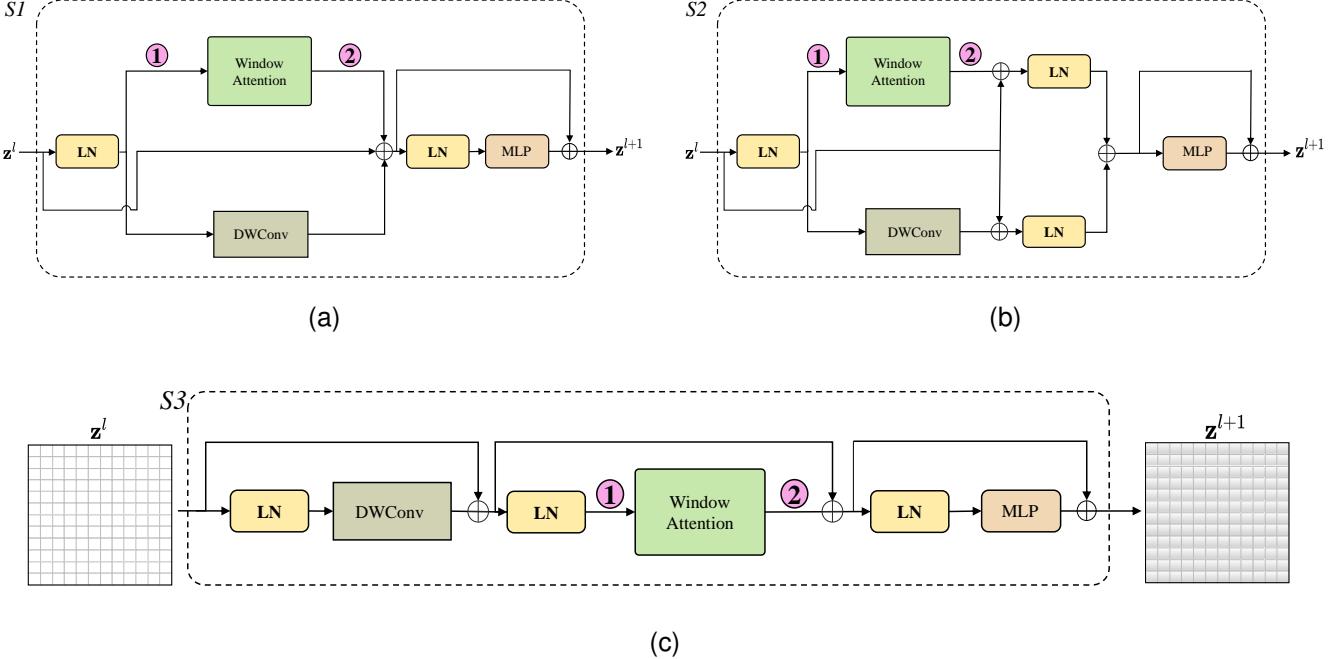


Fig. 3. Diagram of Iwin Block. (a) S1 shows a parallel structure where convolution and attention results are directly combined, as implemented in this study (most fast). (b) S2 is a parallel scheme with independent convolution and attention connections to input, exhibiting the poorest performance. (c) S3 is a serial configuration, where attention input receives convolution output, performing slightly better than S1 but requires one more layer normalization, increasing computation.

A. Image Classification on ImageNet-1K

a) **Settings:** For image classification, we benchmark the proposed Iwin Transformer on ImageNet-1K [39], which contains 1.28M training images and 50K validation images from 1,000 classes. Experimental settings closely follow those of Swin Transformer [11]. We report top-1 accuracy on a single crop under two scenarios: training from scratch on ImageNet-1K and pre-training on ImageNet-22K, which contains 14.2M images and 22K classes, followed by fine-tuning on ImageNet-1K.

From scratch on ImageNet-1K. We use an AdamW optimizer [43] for 300 epochs with a 20-epoch linear warmup and a cosine learning rate decay. We set the batch size to 512, weight decay to 0.05, and use most of the augmentation strategies from DeiT [44]. The initial learning rate is set to 0.0005, and the drop path rates are 0.2, 0.3, and 0.5 for Iwin-T, Iwin-S, and Iwin-B, respectively.

Pre-training on ImageNet-22K. Training lasts for 90 epochs with a 5-epoch warmup. We use a batch size of 4096, an initial learning rate of $1.25e - 4$, and a weight decay of 0.05. The pre-trained models are then fine-tuned on ImageNet-1K for 10 30 epochs at a 224×224 resolution, using a batch size of 1024, a constant learning rate of $2e - 05$, and a weight decay of $1e - 8$.

Cross-Resolution Fine-tuning. A key advantage of Iwin Transformer is its simple to transfer to higher resolutions. Unlike conventional practices that require staged fine-tuning or interpolate absolute/relative position biases, our models pre-trained at 224^2 resolution can be directly fine-tuned on higher

resolutions such as 384^2 , 512^2 , and 1024^2 . The only thing need to do is to change the hyperparameter window size.

During this cross-resolution fine-tuning, we adjust the window size to match the input resolution—7 for 224, 12 for 384, and 16 for 512 and 1024—while keeping the architectures unchanged. The fine-tuning process runs for 30 epochs with a constant learning rate of $2e - 05$, a weight decay of $1e - 8$, and a warm-up of 5 to 10 epochs.

b) **Results:** Table I presents comprehensive results of the Iwin Transformer’s competitive performance and resolution adaptability.

ImageNet-1K. For models trained from scratch on ImageNet-1K at a resolution of 224×224 , Iwin Transformer variants consistently demonstrate competitive Top-1 accuracy. Iwin-T achieves **82.0%**, matching ConvNeXt-T (82.1%) and surpassing Swin-T (81.3%) by **0.7%**. The deeper variant, Iwin-S, reaches highest **83.4%**, outperforming Swin-S (83.0%) by **0.4%** and ConvNeXt-S (83.1%) by **0.3%**, while maintaining similar parameters (51.6M vs. 50.0M) and FLOPs (9.0G vs. 8.7G) compared to both Swin-T and ConvNeXt-T. By increasing the embedding dimension, Iwin-B achieves 83.5%, which is close to the 83.6% of PVT-v2-B4 and falls between the 83.3% of Swin-B and the 83.8% of ConvNeXt-B. The throughput for Iwin models (e.g., Iwin-T at 729 img/s) is slightly lower than that of Swin and ConvNeXt (e.g., Swin-T at 758 img/s and ConvNeXt-T at 775 img/s).

ImageNet-22K. By leveraging large-scale ImageNet-22K pre-training, Iwin significantly enhances its performance. The pre-trained Iwin-B achieves an impressive **85.5%** Top-1 accuracy at a resolution of 224×224 . When fine-tuned to a

resolution of 384×384 , it reaches a robust **86.6%**, surpassing Swin-B (86.4%) by **0.2%** and closely aligning with the carefully configured ConvNeXt-B (86.8%). The largest variant, Iwin-L, demonstrates cutting-edge performance with **86.4%** at 224×224 and an impressive **87.4%** at 384×384 , competing favorably with top-tier models like Swin-L (87.3%) and ConvNeXt-L (87.5%).

Cross-Resolution Fine-tuning. Iwin models have the capability to fine-tune directly from a 224^2 pre-trained state to higher resolutions, a feature particularly notable in models trained from scratch on ImageNet-1K. For example, Iwin-S, after its initial training at 224^2 , can be fine-tuned to achieve impressive Top-1 accuracies of **84.3%** at 384^2 (a **0.9%** increase from its 224^2 baseline of 83.4%) and **84.4%** at 512^2 (a **1.0%** gain). Similarly, Iwin-B reaches **84.9%** at 384^2 (a **1.4%** improvement from its 224^2 baseline of 83.5%), surpassing Swin-B (84.5%) by **0.4%** at 384^2 and closely aligning with ConvNeXt-B (85.1%). Furthermore, Iwin-B demonstrates robust performance at even higher resolutions, achieving **85.1%** at 512^2 (a **1.6%** gain) and a remarkable **85.0%** even at 1024^2 (a **1.5%** gain).

For models pre-trained on ImageNet-22K, this advantage is equally pronounced, showcasing Iwin's seamless adaptability. Iwin-B fine-tunes from 224^2 to 384^2 (86.6%, a **1.1%** gain from its 224^2 baseline of 85.5%), directly outperforming Swin-B (86.4%) by **0.2%** at 384^2 . It continues to perform strongly at 512^2 (86.1%) and 1024^2 (85.6%). Meanwhile, Iwin-L seamlessly transitions from 224^2 to 384^2 (**87.4%**, a 1.0% gain), standing with Swin-L (87.3%) and ConvNeXt-L (87.5%).

This high-resolution fine-tuning capability is attributed to the collaboration of Iwin's interleaved window attention and depthwise separable convolution, which allows the model to enjoy a global view without relying on positional encoding. This robust performance at high resolutions and the potential of collaboration to replace standard attention in models such as high-resolution image and video generation are far more important than the small gap in throughput compared to Swin.

B. Object Detection on COCO

a) **Settings:** We follow Swin [11] settings and evaluate the Iwin backbone on the COCO [40] for both object detection and instance segmentation tasks using Mask R-CNN [52] and Cascaded Mask R-CNN [53] with the MMDetection [54] toolbox. We use multi-scale training, AdamW optimizer and pre-trained models on ImageNet-1K.

b) **Results:** As shown in the table II, the Iwin Transformer consistently underperforms the Swin Transformer across various frameworks and training schedules. For the Mask-RCNN $3 \times$ schedule, Swin-T achieves a bounding box AP of 46.0, surpassing Iwin-T's 44.7. When using the Cascade Mask-RCNN with a $3 \times$ schedule, Swin-T reaches 50.4 AP^{box}, which is 1.0 AP higher than the 49.4 AP^{box} achieved by Iwin-S. Notably, under this setup, the performance of Iwin-S shows no improvement over Iwin-T.

To find the reason for the gap, we compared Iwin-T and Swin-T performance under Cascade Mask-RCNN $3 \times$ schedule

TABLE I
CLASSIFICATION ACCURACY ON IMAGENET-1K. THROUGHPUT IS TESTED
BASED ON THE PYTORCH FRAMEWORK WITH A V100 GPU

Method	Image Size (px)	Param (M)	FLOPs (G)	Throughput (img/s)	Top-1 Acc (%)
(a) ImageNet-1K trained models					
DeiT-Small/16 [44]	224^2	22.0	4.6	406	79.9
T2T-ViT-14 [45]	224^2	22.0	5.2	-	81.5
PVT-Small [46]	224^2	24.5	3.8	794	79.8
Twins-S [5]	224^2	24.0	2.9	979	81.7
PVT-v2-B2 [47]	224^2	25.4	4.0	664	82.0
PoolFormer-S36 [48]	224^2	31.0	5.1	764	81.4
Swin-T [11]	224^2	29.0	4.5	758	81.3
ConvNeXt-T [49]	224^2	29	4.5	775	82.1
Iwin-T(ours)	224^2	30.2	4.7	729	82.0
T2T-ViT-19 [45]	224^2	39.2	8.9	-	81.9
PVT-Medium [46]	224^2	44.2	6.7	511	81.2
Twins-B [5]	224^2	56.0	8.6	433	83.2
PVT-v2-B3 [47]	224^2	45.2	6.9	443	83.2
PoolFormer-M36 [48]	224^2	56.0	9.0	494	82.1
Swin-S [11]	224^2	50.0	8.7	437	83.0
ConvNeXt-S [49]	224^2	50	8.7	447	83.1
Iwin-S(ours)	224^2	51.6	9.0	410	83.4
Iwin-S(ours)	384^2	51.6	27.7	142	84.3
Iwin-S(ours)	512^2	51.6	52.0	78	84.4
Iwin-S(ours)	1024^2	51.6	207.9	20	83.8
DeiT-Base/16 [44]	224^2	86.6	17.6	273	81.8
T2T-ViT-24 [45]	224^2	64.1	14.1	-	82.3
PVT-Large [46]	224^2	61.4	9.8	357	81.7
Twins-L [5]	224^2	99.2	15.1	271	83.7
PVT-v2-B4 [47]	224^2	62.6	10.1	298	83.6
PoolFormer-M48 [48]	224^2	73.0	11.8	337	82.5
Swin-B [11]	224^2	88.0	15.4	287	83.3
Swin-B [11]	384^2	88.0	47.0	85	84.5
ConvNeXt-B [49]	224^2	89	15.4	292	83.8
ConvNeXt-B [49]	384^2	89	45.0	96	85.1
Iwin-B(ours)	224^2	91.2	15.9	271	83.5
Iwin-B(ours)	384^2	91.2	48.3	78	84.9
Iwin-B(ours)	512^2	91.3	89.5	51	85.1
Iwin-B(ours)	1024^2	91.3	358.2	12	85.0
(b) ImageNet-22K pre-trained models					
R-101x3 [50]	384^2	388	204.6	-	84.4
R-152x4 [50]	480^2	937	840.5	-	85.4
ViT-B/16 [51]	384^2	86	55.4	86	84.0
ViT-L/16 [51]	384^2	307	190.7	27	85.2
Swin-B [11]	224^2	88	15.4	278	85.2
Swin-B [11]	384^2	88	47.0	85	86.4
ConvNeXt-B [49]	224^2	89	15.4	292	85.8
ConvNeXt-B [49]	384^2	89	45.1	96	86.8
Iwin-B(ours)	224^2	91.2	15.9	271	85.5
Iwin-B(ours)	384^2	91.2	48.3	79	86.6
Iwin-B(ours)	512^2	91.2	89.5	51	86.1
Iwin-B(ours)	1024^2	91.2	358.2	12	85.6
Swin-L [11]	224^2	197	34.5	145	86.3
Swin-L [11]	384^2	197	103.9	46	87.3
ConvNeXt-L [49]	224^2	198	34.4	147	86.6
ConvNeXt-L [49]	384^2	198	101.0	50	87.5
Iwin-L(ours)	224^2	204.3	35.4	138	86.4
Iwin-L(ours)	384^2	204.3	106.6	43	87.4

by plotting (AP^{box}) at each validation epoch. Both models followed identical learning rate schedules with step-wise decays.

Initially, Iwin-T showed competitive performance, tracking closely with Swin-T and occasionally surpassing it during the first 27 epochs. However, a significant divergence occurred at epoch 28, coinciding with the first major learning rate decay. Swin-T's AP^{box} jumped from 45.9 to 49.2, while Iwin-T showed modest improvement (from 46.2 to 48.7) and lagged behind thereafter. The learning rate drop appears more beneficial for Swin.

To address this performance gap, we explored several alternative training configurations for the Iwin model, including different learning rate strategies (e.g., Cosine Annealing) and architectural enhancements with relative position encoding

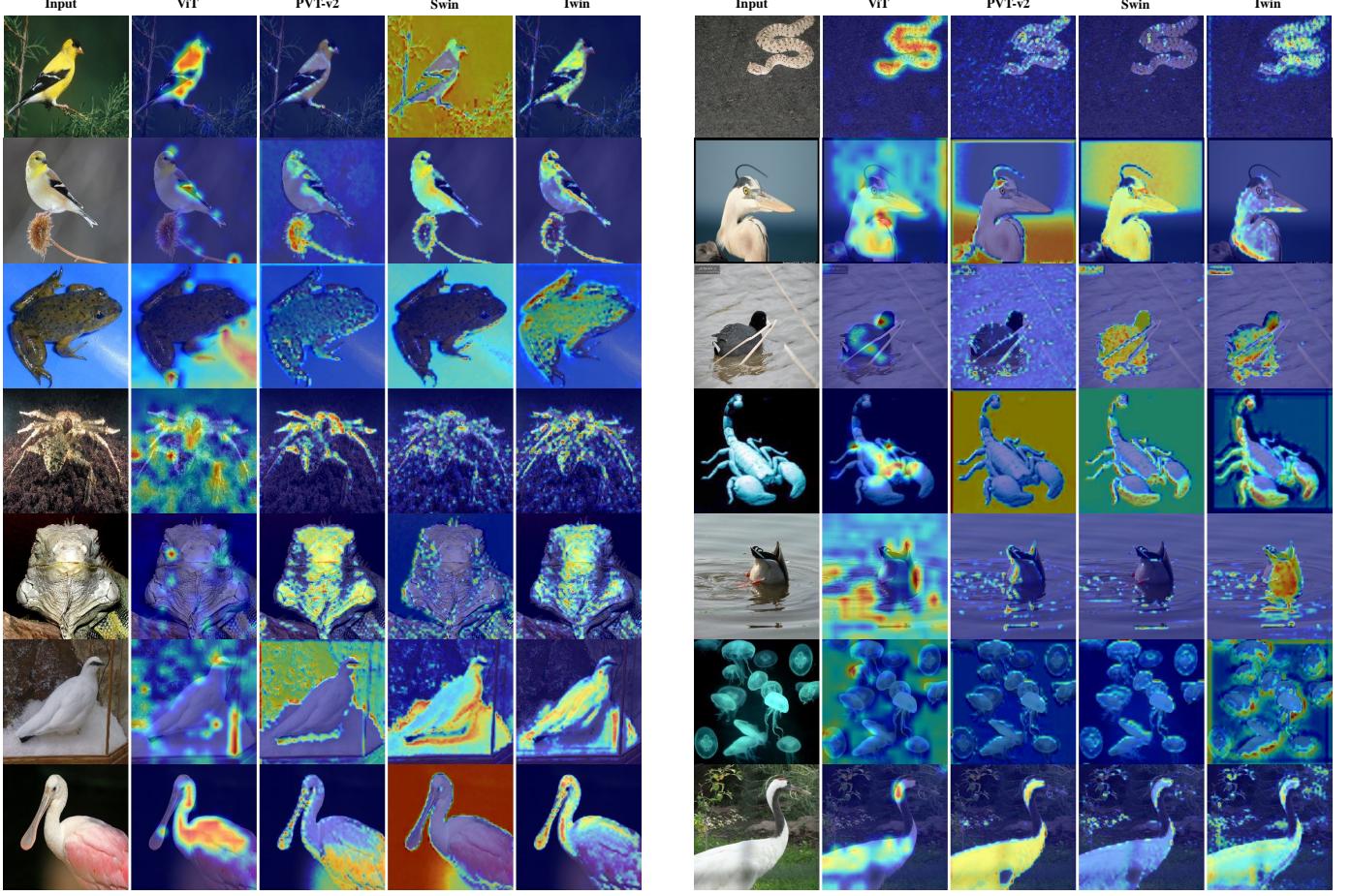


Fig. 4. The visualization of heatmap. The left column shows input images, while subsequent columns show results from native VIT, PVTv2, Swin, and Iwin (our method). Results demonstrate that Iwin effectively concentrates activation on target objects.

(see ablation Table X). Despite these efforts, we were unable to match or exceed the performance of the Swin on the COCO object detection task, which is the only benchmark where the Iwin Transformer did not surpass Swin, indicating a task-specific optimization challenge rather than a general architectural deficiency. This task is left to future research.

C. Semantic Segmentation on ADE20K

a) Settings: We evaluate the Iwin backbone on the ADE20K [41] semantic segmentation benchmark using UperNet [55] with the MMSegmentation [56] toolbox. The experimental settings follow Swin [11]. Training is performed for 160K iterations with a total batch size of 16 (2 images per GPU across 8 GPUs).

b) Results: As shown in Table IV, Iwin-B achieves a 48.9% mIoU, surpassing Swin-B’s 48.1% mIoU by 0.8%, while maintaining nearly identical FLOPs (1189G vs 1188G) and parameters (124.8M vs 121.0M). This performance is close to the leading ConvNeXt-B, which records a 49.1% mIoU with a similar computational cost. For smaller models, Iwin-T achieves a 44.7% mIoU, slightly exceeding Swin-T’s 44.5% mIoU with comparable FLOPs. These results demonstrate Iwin’s effectiveness and competitiveness in semantic segmentation tasks.

D. Video Recognition on Kinetics-400

a) Settings: In line with the Video Swin Transformer [58] settings, we propose the Video Iwin Transformer for action recognition with the MMaction2 [59] toolbox on the Kinetics-400 [42] dataset. The Video Iwin Transformer is initialized with its pre-trained Iwin model from ImageNet. Unlike the Video Swin Transformer’s intricate adaptations from 2D to 3D, which involve window shifting, masking, and relative positioning, the Iwin model only necessitates the addition of a learnable absolute position encoding in the temporal dimension, leaving other components unchanged. As illustrated in Figure 9, in Iwin 3D Attention, tokens from interleaved windows across all frames are collected for attention calculation, while depthwise separable convolution maintains its 2D operation on individual frames.

b) Results: According to Table V, the Iwin model outperforms the Swin in performance and efficiency. At a comparable model scale, Iwin-T achieves a Top-1 accuracy of 79.1% and a Top-5 accuracy of 93.8%, slightly surpassing Swin-T’s 78.8% and 93.6%. More importantly, Iwin-T’s computational cost is significantly lower, requiring only 74 GFLOPs compared to Swin-T’s 88 GFLOPs, marking a **15.9%** reduction. This demonstrates that Iwin-T not only exceeds Swin-T’s performance while offering greater computational

TABLE II

RESULTS FOR COCO OBJECT DETECTION AND SEGMENTATION RESULTS USING MASK-RCNN AND CASCADE MASK-RCNN. FLOPS ARE CALCULATED WITH IMAGE SIZE (1280, 800).

Backbone	FLOPs	AP ^{box}	AP ^{box} ₅₀	AP ^{box} ₇₅	AP ^{mask}	AP ^{mask} ₅₀	AP ^{mask} ₇₅
Mask-RCNN 1× schedule							
PVTv2-B1	-	41.8	64.3	45.9	38.8	61.2	41.6
Swin-T	267G	43.7	66.6	47.7	39.8	63.3	42.7
Iwin-T	268G	42.2	65.3	45.8	38.9	62.1	41.6
Iwin-S	358G	43.7	67.0	47.4	40.0	63.9	42.5
Mask-RCNN 3× schedule							
PVTv2-B2	-	47.8	-	-	43.1	-	-
Swin-T	267G	46.0	68.1	50.3	41.6	65.1	44.9
ConvNeXt-T	262G	46.2	67.9	50.8	41.7	65.0	44.9
Iwin-T	268G	44.7	67.2	48.8	40.9	64.1	43.6
Iwin-S	358G	45.5	67.5	49.6	41.0	64.3	44.0
Cascade Mask-RCNN 1× schedule							
Swin-T	745G	48.1	67.1	52.2	41.7	64.4	45.0
Iwin-T	747G	47.2	66.1	51.3	40.9	63.5	44.1
Cascade Mask-RCNN 3× schedule							
ResNet-50	739G	46.3	64.3	50.5	40.1	61.7	43.4
X101-32	819G	48.1	66.5	52.4	41.6	63.9	45.2
X101-64	972G	48.3	66.4	52.3	41.7	64.0	45.1
PVTv2-B2	788G	51.1	69.8	55.3	44.4	-	-
Swin-T	745G	50.4	69.2	54.7	43.7	66.6	47.3
ConvNeXt-T	741G	50.4	69.1	54.8	43.7	66.5	47.3
Iwin-T	747G	49.4	68.4	53.5	42.9	65.8	46.4
Swin-S	838G	51.9	70.7	56.3	45.0	68.2	48.8
ConvNeXt-S	827G	51.9	70.8	56.5	45.0	68.4	49.1
Iwin-S	837G	49.4	68.1	53.3	43.0	65.6	46.4

efficiency. For the larger-scale Iwin-S, although its Top-1 and Top-5 accuracies (80.0% and 94.1%, respectively) are lower than those of Swin-S (80.6% and 94.5%), its computational cost (140 GFLOPs) remains considerably lower than Swin-S's (166 GFLOPs), reflecting a reduction of approximately **15.7%**. This indicates that Iwin significantly reduces computational cost while maintaining competitive performance.

E. Image Generation

a) **Settings:** We follow LightningDiT [65] settings and build a FlashDiT model to validate the effectiveness of the key component in Iwin for class-conditional image generation task on ImageNet. We replace standard self-attention with a proposed combination of interleaved window attention and depthwise separable convolution. We remove position encodings and set the convolution kernel size to 3×3 and the window size to 4×4 , based on the latent feature map dimensions of 16×16 (derived from 256×256 input images through $16 \times$ downsampling), where $3 \times 4 = 12$ closely approximates 16.

b) **Results:** As shown in Table VI, our proposed FlashDiT demonstrates efficiency in image generation. It achieves competitive performance in only 56 training epochs, a fraction of the epochs required by previous state-of-the-art models such as DiT (1400) and MAR (800). Unlike these models have complexity of 32^2 or 16^2 , FlashDiT has a computational complexity of only $3^2 + 4^2 = 25$ while maintaining high generation quality, achieving a gFID of 3.08 and an IS of 223.2 without classifier guidance. Therefore, FlashDiT verifies the effectiveness of the key component in Iwin as a standalone module that can seamlessly replace the self-attention module in generation models.



Fig. 5. The visualization of object detection on the COCO2017. The leftmost column shows the input images. From left to right, the results generated by PVTv2-based, Swin-based, and Iwin-based Mask R-CNN are shown.

F. Ablation Study

Our extensive ablation studies, summarized in Table VII, are all conducted on Iwin-T and ImageNet by default.

a) **Attention and Convolution Combination:** We investigated the effects of integrating depthwise separable convolutions (DWConv) and different attention mechanisms. As demonstrated, the collaboration of DWConv and IW-MSA (Interleaved Window Multi-head Self-Attention) yields the best performance, achieving a Top-1 accuracy of 82.0%. This verified the superiority of our proposed method.

TABLE III

DETAILED ARCHITECTURE SPECIFICATIONS FOR RESOLUTION 224^2 . WIN. SZ. ARE 12, 16, 16 FOR RESOLUTIONS 384^2 , 512^2 , 1024^2 .

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

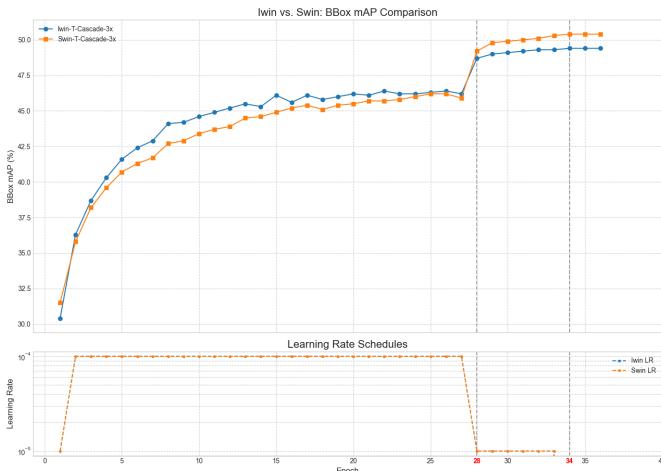


Fig. 6. BBox mAP and Learning Rate Progression for Iwin-T vs. Swin-T on the COCO (Cascade Mask-RCNN 3x schedule).

TABLE IV
RESULTS FOR ADE20K SEMANTIC SEGMENTATION TASK. FLOPS ARE MEASURED WITH THE INPUT SIZE OF 512×2048 .

Backbone	Semantic FPN 80k			UperNet 160k		
	Param(M)	FLOPs(G)	mIoU(%)	Param(M)	FLOPs(G)	mIoU(%)
ResNet50 [35]	28.5	183	36.7	-	-	-
PVTv2 B2 [46]	29.1	165	45.2	-	-	-
ConvNeXt-T [49]	-	-	-	60.0	939	46.0
Swin-T [11]	-	-	-	59.9	945	44.5
Iwin-T(ours)	-	-	-	61.9	946	44.7
ResNet101 [35]	47.5	260	38.8	-	-	-
PVTv2 B3 [46]	49.0	224	47.3	-	-	-
ConvNeXt-S [49]	-	-	-	82.0	1027	48.7
Swin-S [11]	-	-	-	81.3	1038	47.6
Iwin-S(ours)	-	-	-	83.2	1038	47.5
ResNeXt101-64×4d [57]	86.4	-	40.2	-	-	-
PVTv2 B4 [46]	66.3	285	47.9	-	-	-
ConvNeXt-B [49]	-	-	-	122.0	1170	49.1
Swin-B [11]	-	-	-	121.0	1188	48.1
Iwin-B(ours)	-	-	-	124.8	1189	48.9

b) **Downsampling Methods:** We evaluated various downsampling methods between network stages. The results show that using Standard Convolution (Std Conv), adopted in our final Iwin design, achieves the highest accuracy of 82.0%.

TABLE V
COMPARISON ON KINETICS-400. “VIEWS” INDICATES # TEMPORAL CLIP × # SPATIAL CROP. THE MAGNITUDES ARE GIGA (10^9) AND MEGA (10^6) FOR FLOPS AND PARAM RESPECTIVELY.

Method	Pretrain	Top-1	Top-5	Views	FLOPs	Param
R(2+1)D [60]	-	72.0	90.0	10×1	75	61.8
I3D [61]	ImageNet-1K	72.1	90.3	-	108	25.0
NL I3D-101 [62]	ImageNet-1K	77.7	93.3	10×3	359	61.8
SlowFast R101+NL [63]	-	79.8	93.9	10×3	234	59.9
X3D-XXL [64]	-	80.4	94.6	10×3	144	20.3
MViT-B, 32×3 [20]	-	80.2	94.4	1×5	170	36.6
MViT-B, 64×3 [20]	-	81.2	95.1	3×3	455	36.6
ViViT-L/16x2 [16]	ImageNet-21K	80.6	94.7	4×3	1446	310.8
ViViT-L/16x2 320 [16]	ImageNet-21K	81.3	94.7	4×3	3992	310.8
Swin-T [11]	ImageNet-1K	78.8	93.6	4×3	88	28.2
Swin-S [11]	ImageNet-1K	80.6	94.5	4×3	166	49.8
Iwin-T(ours)	ImageNet-1K	79.1	93.8	4×3	74	29.8
Iwin-S(ours)	ImageNet-1K	80.0	94.1	4×3	140	51.1

We struggled between standard convolution and average pooling offering highest throughput, but for higher accuracy, we chose standard convolution.

c) **Kernel Size Choice:** We explored the impact of varying kernel sizes for DWConv across different stages. Our final configuration, which utilizes a fixed kernel size of {3, 3, 3, None}, yields a accuracy of 82.0% with a throughput of 736 img/s. While a {5, 5, 5, None} kernel size achieves a slightly higher accuracy of 82.2%, the smaller {3, 3, 3, None} kernels offer a more favorable trade-off between performance and computational efficiency (higher throughput). Additionally, using different kernel sizes at different stages to satisfy $KM \geq \max(H, W)$ did not yield the best results. This aligns with observations in [37] that balanced network outperform theoretically optimal but imbalanced configurations.

d) **Block Number Configuration:** We examined different distributions of block numbers across the four stages of the network. Initially, we explored the {4, 3, 2, 2} configuration, which aimed to approximate larger kernel sizes (e.g., a 7×7 kernel with four 3×3 convolutions) by stacking smaller convolutional blocks. However, this setup resulted in the lowest accuracy (80.5%) and throughput (473 img/s). In contrast, the {2, 2, 6, 2} configuration achieved the highest accuracy of 82.0% with a throughput of 736 img/s. This

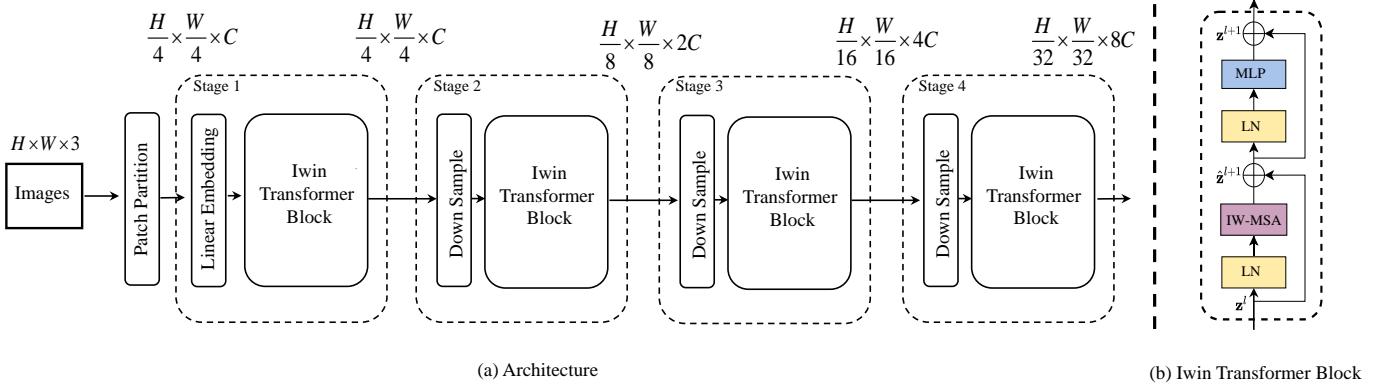


Fig. 7. (a) The overall architecture of the Iwin Transformer (Iwin-T). (b) Single Iwin Transformer Block. IW-MSA involves applying interleaved window multi-head self-attention and depthwise separable convolution in parallel.

TABLE VI
COMPARISON WITH STATE-OF-THE-ART MODELS ON IMAGENET FOR IMAGE GENERATION.

Method	Epoches	Params	256×256 w/o CFG					256×256 w/ CFG				
			gFID ↓	sFID ↓	IS ↑	Pre. ↑	Rec. ↑	gFID ↓	sFID ↓	IS ↑	Pre. ↑	Rec. ↑
<i>AutoRegressive (AR)</i>												
MaskGIT [?]	555	227M	6.18	-	182.1	0.80	0.51	-	-	-	-	-
LlamaGen [66]	300	3.1B	9.38	8.24	112.9	0.69	0.67	2.18	5.97	263.3	0.81	0.58
VAR [67]	350	2.0B	-	-	-	-	-	1.80	-	365.4	0.83	0.57
MagViT-v2 [68]	1080	307M	3.65	-	200.5	-	-	1.78	-	319.4	-	-
MAR [69]	800	945M	2.35	-	227.8	0.79	0.62	1.55	-	303.7	0.81	0.62
<i>Latent Diffusion Models</i>												
MaskDiT [70]	1600	675M	5.69	10.34	177.9	0.74	0.60	2.28	5.67	276.6	0.80	0.61
DiT [18]	1400	675M	9.62	6.85	121.5	0.67	0.67	2.27	4.60	278.2	0.83	0.57
SiT [71]	1400	675M	8.61	6.32	131.7	0.68	0.67	2.06	4.50	270.3	0.82	0.59
MDTv2 [72]	1080	675M	-	-	-	-	-	1.58	4.52	314.7	0.79	0.65
REPA [73]	800	675M	5.90	-	-	-	-	1.42	4.70	305.7	0.80	0.65
LightningDiT [65]	64	675M	5.14	4.22	130.2	0.76	0.62	2.11	4.16	252.3	0.81	0.58
FlashDiT(Ours)	56	675M	7.88	5.91	116.2	0.73	0.60	3.08	6.00	223.2	0.78	0.57

configuration, which allocates more blocks to the deeper stages, proves beneficial for maximizing performance while maintaining optimal speed.

e) **Position Embedding:** We explored the effect of positional encoding on the models. In the Iwin-T model, relative position embedding achieves the highest accuracy at 82.4%. However, in deeper models like Iwin-S, the no-position-embedding approach reaches the highest Top-1 accuracy of 83.4%, surpassing the relative position embedding (83.3%) and processing more images per second (410 img/s). This indicates that in very deep networks, position embeddings might be unnecessary or even harmful. Additionally, during training, models with absolute or relative position embeddings take more time to learn compared to those without position embedding.

As shown in the Table VIII, our baseline Iwin-T model achieves $42.2 AP^{box}$, which is 1.5 lower than the $43.7 AP^{box}$ from Swin-T under the same settings. Altering the learning rate strategy to smoother Cosine Annealing or increasing the initial learning rate did not yield improvements. However, incorporating relative position encoding provided a notable gain of $0.7 AP^{box}$ for the Iwin-T model. Encouraged by this

enhancement, we further investigated its impact on the scaled-up Iwin-S model. Iwin-S successfully bridged the performance gap, reaching $43.7 AP^{box}$, which is comparable to Swin-T. However, adding relative position encoding to Iwin-S caused a slight performance decline to $43.5 AP^{box}$. These results indicate that the Iwin architecture encounters a complex task-specific optimization challenge on the COCO benchmark.

V. DISCUSSION

We believe that some fields of work can benefit and be inspired by Iwin Transformer.

A. Migration to Large Language Models

The Iwin Transformer's position-embedding-free design principle offers promising opportunities for application in Large Language Models (LLMs). Currently, LLMs heavily depend on position embeddings to preserve sequence order information. By integrating interleaved window attention with depthwise separable convolution, it may be possible to achieve more natural length generalization. This approach relies on structural rather than parametric position information, facilitating easier length extrapolation. As shown in Figure 10,

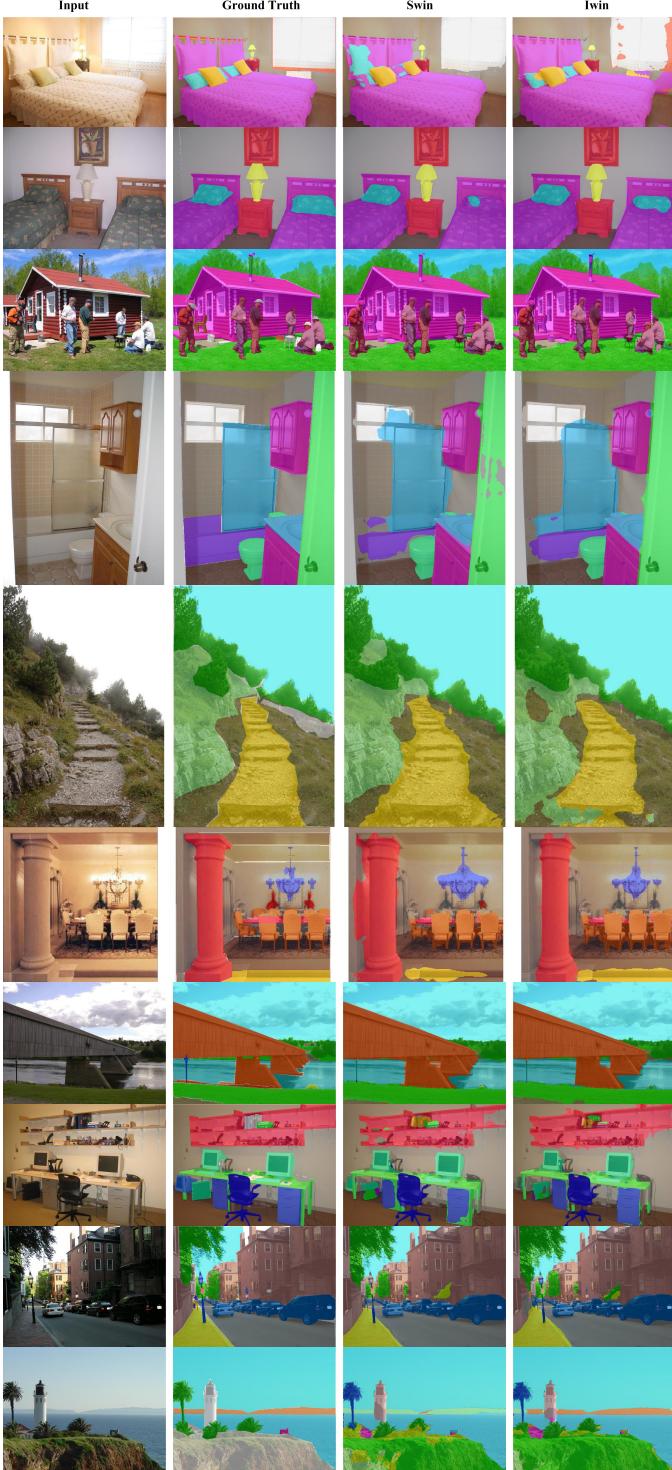


Fig. 8. Results for semantic segmentation on ADE20K. The first column shows the input images. From left to right: Ground Truth, Swin-based, and Iwin-based UperNet.

the computation is divided into two components: 1D causal depthwise separable convolution and 1D interleaved window causal attention. Both ensure tokens relate only to preceding tokens, giving Iwin 1D Attention causality. Furthermore, we can replace the depthwise separable convolution with normal window causal attention by setting two window sizes M_1 and

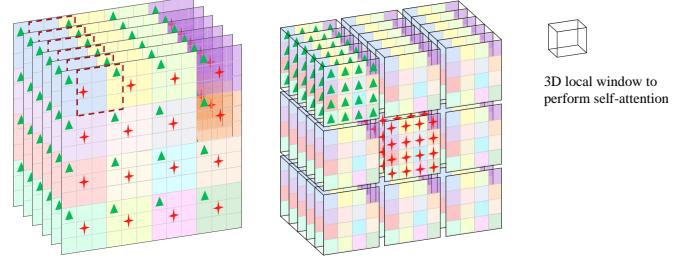


Fig. 9. Illustration of Iwin 3D Attention. In the left image, depth-wise separable convolution is performed on each frame, as in 2D. On the right, Interleaved Window Attention is performed in 3D. This effectively approximates 3D full attention, allowing connections between any tokens within a video. The only disadvantage is that position coding is required for the time dimension. To get rid of it completely, we propose two options: perform RTR operations in the time dimension and cooperate with one-dimensional convolution; the second is to increase convolution kernel size in the time dimension to span frames. Due to paper limitations, this part will be reserved for future work.

TABLE VII
ABLATION STUDIES ON VARIOUS ARCHITECTURAL COMPONENTS. FLOPS ARE GIGA (10^9) AND PARAM ARE MEGA (10^6).

Setting	Param(M)	FLOPs(G)	Throughput(img/s)	Top-1 Acc(%)
Ablation on Attention and Convolution Combination				
DwConv	21.60	3.20	861	79.4
W-MSA	30.20	4.71	758	80.2
IW-MSA	30.20	4.71	756	80.4
DwConv + W-MSA	30.23	4.72	737	81.8
DwConv + IW-MSA	30.23	4.72	736	82.0
Ablation on Downsampling Methods				
DWConv	27.14	4.57	746	81.9
Avg Pooling	27.13	4.51	762	81.8
Patch Merging	28.29	4.51	741	81.8
Std Conv	30.23	4.72	736	82.0
Ablation on Kernel Size for Depthwise Convolution				
{7, 5, 3, None}	30.24	4.75	714	81.0
{7, 7, 7, None}	30.34	4.78	729	82.1
{5, 5, 5, None}	30.27	4.75	731	82.2
{3, 3, 3, None}	30.23	4.72	736	82.0
Ablation on Block Number Configuration				
{4, 3, 2, 2}	23.79	4.43	473	80.5
{3, 3, 3, 3}	32.56	4.80	588	81.8
{2, 2, 6, 2}	30.23	4.72	736	82.0
Ablation on Position Embedding				
abs. pos.	30.53	4.72	735	82.1
rel. pos.	30.25	4.72	724	82.4
no pos.	30.23	4.72	736	82.0
rel. pos. (Iwin-S)	51.60	8.98	403	83.3
no pos. (Iwin-S)	51.60	8.98	410	83.4

M_2 equal, with $M_1M_2 = N$, resulting in window size \sqrt{N} . This reduces complexity for sequence length N from N^2 to N .

B. Application to Generation Models

Generation Models, such as denoising diffusion models, present another domain where Iwin's design could offer substantial benefits. Iwin's absence of position embedding facilitates seamless adaptation to various resolutions without requiring parameter adjustments or interpolation, which is crucial for progressive generation strategies. This feature can help image generation models create higher resolution or any size images. Moreover, the faster convergence properties observed in Iwin could reduce the training time for diffusion models, thanks to the inductive bias introduced by the depthwise separable convolution in Iwin.

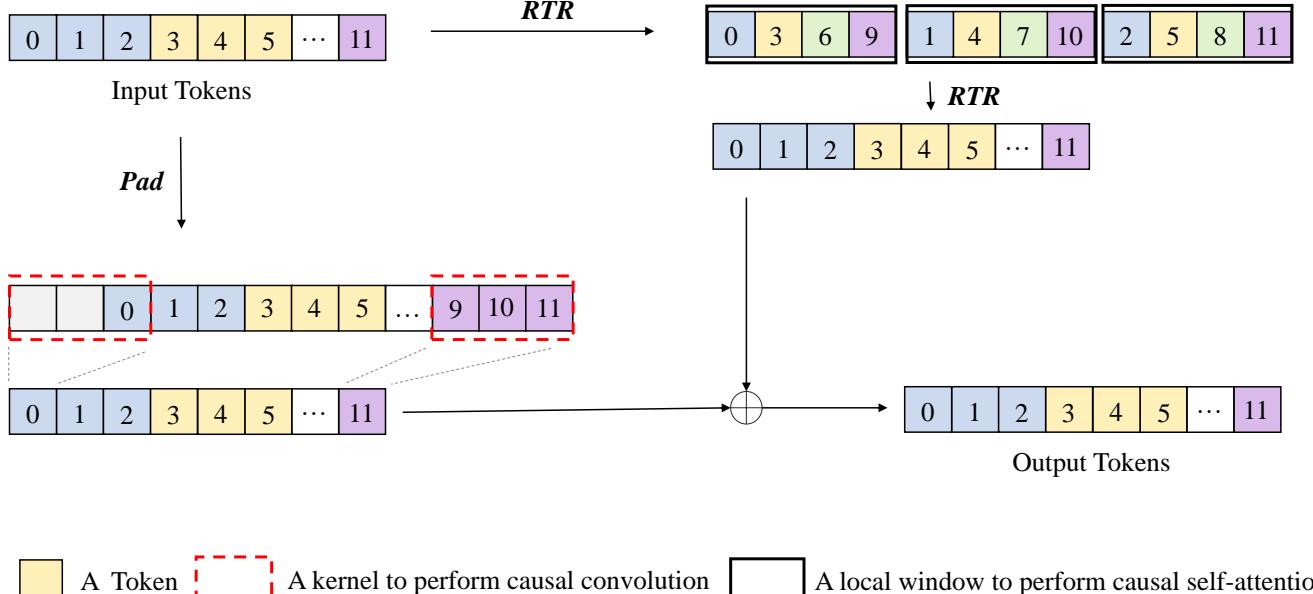


Fig. 10. The illustration of Iwin 1D Attention, shows how to apply Iwin Transformer concepts into large language models (LLMs). The red boxes represent causal depthwise separable convolution, while the black boxes denote causal attention within a window. Neither process leaks future information. The final output, derived from the combination of these two causal operations, adheres to the principle of causality. This approach might be also helpful in solving the problem of high complexity with long sequences in LLMs.

TABLE VIII
ABLATION STUDIES FOR IWIN-T ON THE COCO WITH THE MASK R-CNN 1 × SCHEDULE. "LR" DENOTES THE INITIAL LEARNING RATE. "REL. POS." INDICATES THE RELATIVE POSITION ENCODING.

Method	AP ^{box}	AP ^{box} ₅₀	AP ^{box} ₇₅	AP ^{mask}	AP ^{mask} ₅₀	AP ^{mask} ₇₅
Swin-T (lr=1e-4, step)	43.7	66.6	47.7	39.8	63.3	42.7
Iwin-T (lr=1e-4, step)	42.2	65.3	45.8	38.9	62.1	41.6
CosineAnnealing	42.0	65.2	45.7	38.7	61.9	41.2
Increase lr to 2e-4	42.2	64.9	46.2	39.1	62.0	41.9
rel. pos.	42.9	66.0	46.7	39.4	62.7	42.2
Iwin-S	43.7	67.0	47.4	40.0	63.9	42.5
Iwin-S + rel. pos.	43.5	66.7	47.4	40.1	63.4	42.7

Our proposed Iwin 3D Attention shown in Figure 9, which forms windows spanning both spatial and temporal domains, works in conjunction with 2D depthwise separable convolution in the spatial domain and has already demonstrated effectiveness in action recognition. We believe it can serve as a third option in video generation, alongside 3D Full Attention and Spatiotemporal Attention mechanisms. While 3D Full Attention is an expensive ideal solution, the serial structure of Spatiotemporal Attention mechanisms can cause temporal attention to disrupt the distribution formed by spatial attention, potentially leading to disharmony in frame images. In contrast, Iwin 3D Attention uses one attention operation and one depthwise separable convolution to establish relationships among all tokens within a video, so we can anticipate that the quality of generated videos will be higher.

C. Limitations and Future Work

Iwin's performance in object detection is not as good as Swin's, and the reasons for this remain unclear. Our limited

computational resources prevented us from conducting extensive experiments to identify an effective learning strategy or to enhance its object detection capabilities through careful optimization. This task is left to future research. Additionally, we did not verify whether Iwin adheres to the scaling law. Future work will involve extending the proposed Iwin Attention and its variants to applications in large language models, image generation, and 3D video generation.

VI. CONCLUSION

In this paper, we introduced Iwin Transformer, a novel position-embedding-free vision Transformer, leveraging the collaboration of innovative interleaved window attention and depthwise separable convolution. Extensive experimental evaluations across vision benchmarks demonstrate Iwin's competitive performance in tasks such as image classification, semantic segmentation, and video action recognition.

Most importantly, the idea of using attention to capture long-range dependencies and convolution to grasp local relationships for building global connections, along with its implementation method, can inspire future work. The core component of the Iwin Transformer can be directly applied to 2D generative models and have been proven its effectiveness in class-conditional image generation task, and shows potential for extension to 3D data (e.g., video generation) with Iwin 3D Attention. Iwin 1D Attention might also be effective for 1D data in large language models, left for future work.

REFERENCES

- [1] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly *et al.*,

- “An image is worth 16x16 words: Transformers for image recognition at scale,” *International Conference on Learning Representations (ICLR)*, 2021.
- [2] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin, “Attention is all you need,” *Advances in Neural Information Processing Systems (NeurIPS)*, vol. 30, 2017.
 - [3] A. Krizhevsky, I. Sutskever, and G. E. Hinton, “Imagenet classification with deep convolutional neural networks,” *Advances in neural information processing systems*, vol. 25, 2012.
 - [4] W. Wang, E. Xie, X. Li, D.-P. Fan, K. Song, D. Liang, T. Lu, P. Luo, and L. Shao, “Pyramid vision transformer: A versatile backbone for dense prediction without convolutions,” *IEEE International Conference on Computer Vision (ICCV)*, pp. 568–578, 2021.
 - [5] X. Chu, Z. Tian, Y. Wang, B. Zhang, H. Ren, X. Wei, H. Xia, and C. Shen, “Twins: Revisiting the design of spatial attention in vision transformers,” in *Advances in Neural Information Processing Systems (NeurIPS)*, vol. 34, 2021, pp. 9355–9366.
 - [6] S. d’Ascoli, H. Touvron, M. Leavitt, A. Morcos, G. Biroli, and L. Sagun, “Convit: Improving vision transformers with soft convolutional inductive biases,” *International Conference on Machine Learning (ICML)*, pp. 2286–2296, 2021.
 - [7] Z. Dai, H. Liu, Q. V. Le, and M. Tan, “Coatnet: Marrying convolution and attention for all data sizes,” *Advances in Neural Information Processing Systems (NeurIPS)*, vol. 34, pp. 3965–3977, 2021.
 - [8] M. S. Ryoo, A. Piergiovanni, A. Arnab, M. Dehghani, and A. Angelova, “Tokenlearner: What can 8 learned tokens do for images and videos?” *arXiv preprint arXiv:2106.11297*, 2021.
 - [9] N. Kitaev, Ł. Kaiser, and A. Levskaya, “Reformer: The efficient transformer,” *arXiv preprint arXiv:2001.04451*, 2020.
 - [10] K. Choromanski, V. Likhoshesterov, D. Dohan, X. Song, A. Gane, T. Sarlos, P. Hawkins, J. Davis, A. Mohiuddin, Ł. Kaiser *et al.*, “Rethinking attention with performers,” *International Conference on Learning Representations (ICLR)*, 2021.
 - [11] Z. Liu, Y. Lin, Y. Cao, H. Hu, Y. Wei, Z. Zhang, S. Lin, and B. Guo, “Swin transformer: Hierarchical vision transformer using shifted windows,” *IEEE International Conference on Computer Vision (ICCV)*, pp. 10 012–10 022, 2021.
 - [12] Z. Liu, H. Hu, Y. Lin, Z. Yao, Z. Xie, Y. Wei, J. Ning, Y. Cao, Z. Zhang, L. Dong, F. Wei, and B. Guo, “Swin transformer v2: Scaling up capacity and resolution,” *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 12 009–12 019, 2022.
 - [13] H. Touvron, M. Cord, M. Douze, F. Massa, A. Sablayrolles, and H. Jegou, “Training data-efficient image transformers & distillation through attention,” *International Conference on Machine Learning (ICML)*, 2021.
 - [14] N. Carion, F. Massa, G. Synnaeve, N. Usunier, A. Kirillov, and S. Zagoruyko, “End-to-end object detection with transformers,” *European Conference on Computer Vision (ECCV)*, pp. 213–229, 2020.
 - [15] S. Zheng, J. Lu, H. Zhao, X. Zhu, Z. Luo, Y. Wang, Y. Fu, J. Feng, T. Xiang, P. H. Torr *et al.*, “Rethinking semantic segmentation from a sequence-to-sequence perspective with transformers,” *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 6881–6890, 2021.
 - [16] A. Arnab, M. Dehghani, G. Heigold, C. Sun, M. Lučić, and C. Schmid, “Vivit: A video vision transformer,” *IEEE International Conference on Computer Vision (ICCV)*, pp. 6836–6846, 2021.
 - [17] X. Zhu, W. Su, L. Lu, B. Li, X. Wang, and J. Dai, “Deformable detr: Deformable transformers for end-to-end object detection,” *arXiv preprint arXiv:2010.04159*, 2020.
 - [18] W. Peebles and S. Xie, “Scalable diffusion models with transformers,” in *Proceedings of the IEEE/CVF international conference on computer vision*, 2023, pp. 4195–4205.
 - [19] S. Mehta and M. Rastegari, “Mobilevit: Light-weight, general-purpose, and mobile-friendly vision transformer,” *International Conference on Learning Representations (ICLR)*, 2022.
 - [20] H. Fan, B. Xiong, K. Mangalam, Y. Li, K. He, and J. Malik, “Multiscale vision transformers,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 6824–6835.
 - [21] S. Wang, B. Z. Li, M. Khabsa, H. Fang, and H. Ma, “Linformer: Self-attention with linear complexity,” *arXiv preprint arXiv:2006.04768*, 2020.
 - [22] X. Ma, X. Kong, S. Wang, C. Zhou, J. May, H. Ma, and L. Zettlemoyer, “Luna: Linear unified nested attention,” *Advances in Neural Information Processing Systems*, vol. 34, pp. 2441–2453, 2021.
 - [23] M. Zaheer, G. Guruganesh, K. A. Dubey, J. Ainslie, C. Alberti, S. Ontanon, P. Pham, A. Ravula, Q. Wang, L. Yang *et al.*, “Big bird: Transformers for longer sequences,” *Advances in Neural Information Processing Systems (NeurIPS)*, vol. 33, pp. 17 283–17 297, 2020.
 - [24] I. Beltagy, M. E. Peters, and A. Cohan, “Longformer: The long-document transformer,” *arXiv preprint arXiv:2004.05150*, 2020.
 - [25] R. Child, S. Gray, A. Radford, and I. Sutskever, “Generating long sequences with sparse transformers,” *arXiv preprint arXiv:1904.10509*, 2019.
 - [26] Y. Tay, D. Bahri, L. Yang, D. Metzler, and D.-C. Juan, “Sparse sinkhorn attention,” *International Conference on Machine Learning (ICML)*, pp. 9438–9447, 2020.
 - [27] A. Roy, M. Saffar, A. Vaswani, and D. Grangier, “Efficient content-based sparse attention with routing transformers,” *Transactions of the Association for Computational Linguistics*, vol. 9, pp. 53–68, 2021.
 - [28] X. Dong, J. Bao, D. Chen, W. Zhang, N. Yu, L. Yuan, D. Chen, and B. Guo, “Cswin transformer: A general vision transformer backbone with cross-shaped windows,” *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 12 124–12 134, 2022.
 - [29] Y. Li, K. Zhang, J. Cao, R. Timofte, and L. Van Gool, “Localvit: Bringing locality to vision transformers,” *arXiv preprint arXiv:2104.05707*, 2021.
 - [30] Y. Rao, W. Zhao, B. Liu, J. Lu, J. Zhou, and C.-J. Hsieh, “Dynamivit: Efficient vision transformers with dynamic token sparsification,” *Advances in Neural Information Processing Systems (NeurIPS)*, vol. 34, 2021.
 - [31] H. Yin, A. Vahdat, J. Alvarez, A. Mallya, J. Kautz, and P. Molchanov, “A-vit: Adaptive tokens for efficient vision transformer,” *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 10 809–10 818, 2022.
 - [32] D. Bolya, C.-Y. Fu, X. Dai, P. Zhang, C. Feichtenhofer, and J. Hoffman, “Token merging: Your vit but faster,” *arXiv preprint arXiv:2210.09461*, 2022.
 - [33] M. Fayyaz, S. A. Koohpayegani, F. Rezaei, S. Somayaji, H. Pirsavash, and J. Gall, “Adaptive token sampling for efficient vision transformers,” *European Conference on Computer Vision (ECCV)*, pp. 209–226, 2022.
 - [34] F. Chollet, “Xception: Deep learning with depthwise separable convolutions,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 1251–1258.
 - [35] K. He, X. Zhang, S. Ren, and J. Sun, “Deep residual learning for image recognition,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 770–778.
 - [36] D. Hendrycks and K. Gimpel, “Gaussian error linear units (gelus),” *arXiv preprint arXiv:1606.08415*, 2016.
 - [37] M. Tan and Q. Le, “Efficientnet: Rethinking model scaling for convolutional neural networks,” in *International conference on machine learning*. PMLR, 2019, pp. 6105–6114.
 - [38] W. Luo, Y. Li, R. Urtasun, and R. Zemel, “Understanding the effective receptive field in deep convolutional neural networks,” *Advances in neural information processing systems*, vol. 29, 2016.
 - [39] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, “Imagenet: A large-scale hierarchical image database,” in *2009 IEEE conference on computer vision and pattern recognition*. Ieee, 2009, pp. 248–255.
 - [40] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick, “Microsoft coco: Common objects in context,” in *Computer vision–ECCV 2014: 13th European conference, Zurich, Switzerland, September 6–12, 2014, proceedings, part v 13*. Springer, 2014, pp. 740–755.
 - [41] B. Zhou, H. Zhao, X. Puig, T. Xiao, S. Fidler, A. Barriuso, and A. Torralba, “Semantic understanding of scenes through the ade20k dataset,” *International Journal of Computer Vision*, vol. 127, pp. 302–321, 2019.
 - [42] W. Kay, J. Carreira, K. Simonyan, B. Zhang, C. Hillier, S. Vijayanarasimhan, F. Viola, T. Green, T. Back, P. Natsev, M. Suleyman, and A. Zisserman, “The kinetics human action video dataset,” 2017.
 - [43] D. P. Kingma and J. Ba, “Adam: A method for stochastic optimization,” 2017. [Online]. Available: <https://arxiv.org/abs/1412.6980>
 - [44] H. Touvron, M. Cord, M. Douze, F. Massa, A. Sablayrolles, and H. Jegou, “Training data-efficient image transformers & distillation through attention,” in *International Conference on Machine Learning (ICML)*, 2021, pp. 10 347–10 357.
 - [45] L. Yuan, Y. Chen, T. Wang, W. Yu, Y. Shi, Z. Jiang, F. E. Tay, J. Feng, and S. Yan, “Tokens-to-token vit: Training vision transformers from scratch on imagenet,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, 2021, pp. 558–567.
 - [46] W. Wang, E. Xie, X. Li, D. Fan, M.-M. Song, D. Liang, T. Lu, P. Luo, and L. Shao, “Pyramid vision transformer: A versatile backbone for dense prediction without convolutions,” *arXiv preprint arXiv:2102.12122*, 2021.

- [47] ——, “Pvtv2: Improved baselines with pyramid vision transformer,” in *Computational Visual Media (CVM)*, 2022, pp. 100–110.
- [48] W. Yu, C. Si, Z. Zhou, X. Tan, and J. Wang, “Metaformer is actually what you need for vision,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, 2021, pp. 6658–6668.
- [49] Z. Liu, H. Mao, C.-Y. Wu, C. Feichtenhofer, T. Darrell, and S. Xie, “A convnet for the 2020s,” in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2022, pp. 11 976–11 986.
- [50] A. Kolesnikov, L. Beyer, X. Zhai, J. Puigcerver, J. Yung, S. Gelly, and N. Houlsby, “Big transfer (bit): General visual representation learning,” 2020. [Online]. Available: <https://arxiv.org/abs/1912.11370>
- [51] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly, J. Uszkoreit, and N. Houlsby, “An image is worth 16x16 words: Transformers for image recognition at scale,” 2021. [Online]. Available: <https://arxiv.org/abs/2010.11929>
- [52] K. He, G. Gkioxari, P. Dollár, and R. Girshick, “Mask r-cnn,” in *Proceedings of the IEEE international conference on computer vision*, 2017, pp. 2961–2969.
- [53] Z. Cai and N. Vasconcelos, “Cascade r-cnn: Delving into high quality object detection,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 6154–6162.
- [54] K. Chen, J. Wang, J. Pang, Y. Cao, Y. Xiong, X. Li, S. Sun, W. Feng, Z. Liu, J. Xu, Z. Zhang, D. Cheng, C. Zhu, T. Cheng, Q. Zhao, B. Li, X. Lu, R. Zhu, Y. Wu, J. Dai, J. Wang, J. Shi, W. Ouyang, C. C. Loy, and D. Lin, “MMDetection: Open mmlab detection toolbox and benchmark,” *arXiv preprint arXiv:1906.07155*, 2019.
- [55] T. Xiao, Y. Liu, B. Zhou, Y. Jiang, and J. Sun, “Unified perceptual parsing for scene understanding,” in *Proceedings of the European conference on computer vision (ECCV)*, 2018, pp. 418–434.
- [56] M. Contributors, “MMSegmentation: Openmmlab semantic segmentation toolbox and benchmark,” <https://github.com/open-mmlab/mmsegmentation>, 2020.
- [57] S. Xie, R. Girshick, P. Dollár, Z. Tu, and K. He, “Aggregated residual transformations for deep neural networks,” 2017. [Online]. Available: <https://arxiv.org/abs/1611.05431>
- [58] Z. Liu, J. Ning, Y. Cao, Y. Wei, Z. Zhang, S. Lin, and H. Hu, “Video swin transformer,” *arXiv preprint arXiv:2106.13230*, 2021.
- [59] M. Contributors, “Openmmlab’s next generation video understanding toolbox and benchmark,” <https://github.com/open-mmlab/mmaction2>, 2020.
- [60] D. Tran, H. Wang, L. Torresani, J. Ray, Y. LeCun, and M. Paluri, “A closer look at spatiotemporal convolutions for action recognition,” in *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*, 2018, pp. 6450–6459.
- [61] J. Carreira and A. Zisserman, “Quo vadis, action recognition? a new model and the kinetics dataset,” in *proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 6299–6308.
- [62] X. Wang, R. Girshick, A. Gupta, and K. He, “Non-local neural networks,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 7794–7803.
- [63] C. Feichtenhofer, H. Fan, J. Malik, and K. He, “Slowfast networks for video recognition,” in *Proceedings of the IEEE/CVF international conference on computer vision*, 2019, pp. 6202–6211.
- [64] C. Feichtenhofer, “X3d: Expanding architectures for efficient video recognition,” in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2020, pp. 203–213.
- [65] J. Yao, B. Yang, and X. Wang, “Reconstruction vs. generation: Taming optimization dilemma in latent diffusion models,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2025.
- [66] P. Sun, Y. Jiang, S. Chen, S. Zhang, B. Peng, P. Luo, and Z. Yuan, “Autoregressive model beats diffusion: Llama for scalable image generation,” *arXiv preprint arXiv:2406.06525*, 2024.
- [67] K. Tian, Y. Jiang, Z. Yuan, B. Peng, and L. Wang, “Visual autoregressive modeling: Scalable image generation via next-scale prediction,” *arXiv preprint arXiv:2404.02905*, 2024.
- [68] L. Yu, J. Lezama, N. B. Gundavarapu, L. Versari, K. Sohn, D. Minnen, Y. Cheng, V. Birodkar, A. Gupta, X. Gu *et al.*, “Language model beats diffusion–tokenizer is key to visual generation,” *arXiv preprint arXiv:2310.05737*, 2023.
- [69] T. Li, Y. Tian, H. Li, M. Deng, and K. He, “Autoregressive image generation without vector quantization,” *arXiv preprint arXiv:2406.11838*, 2024.
- [70] H. Zheng, W. Nie, A. Vahdat, and A. Anandkumar, “Fast training of diffusion models with masked transformers,” *arXiv preprint arXiv:2306.09305*, 2023.
- [71] N. Ma, M. Goldstein, M. S. Albergo, N. M. Boffi, E. Vandenberghe, and S. Xie, “Sit: Exploring flow and diffusion-based generative models with scalable interpolant transformers,” *arXiv preprint arXiv:2401.08740*, 2024.
- [72] S. Gao, P. Zhou, M.-M. Cheng, and S. Yan, “Mdvtv2: Masked diffusion transformer is a strong image synthesizer,” *arXiv preprint arXiv:2303.14389*, 2023.
- [73] S. Yu, S. Kwak, H. Jang, J. Jeong, J. Huang, J. Shin, and S. Xie, “Representation alignment for generation: Training diffusion transformers is easier than you think,” *arXiv preprint arXiv:2410.06940*, 2024.



Simin Huo received his B.S. degree from Nanjing University of Science and Technology, China, Nanjing, in 2018, and his M.S. degree from Bau-man Moscow State Technical University, Russia, Moscow, in 2021. He is currently a Ph.D. student with the Department of Automatics, Shanghai Jiao Tong University. His main research interests include computer vision and deep learning , with a focus on efficient vision transformer and generation models.



Ning Li (Member, IEEE) received the B.S. and M.S. degrees from Qingdao University of Science and Technology, Qingdao, China, in 1996 and 1999, respectively, and the Ph.D. degree from Shanghai Jiao Tong University, Shanghai, China, in 2002. She is currently a Professor with the Department of Automation, Shanghai Jiao Tong University. Her research interests include modeling and control of complex systems, artificial intelligence, and big data analysis.