# Vision Transformers are Circulant Attention Learners

**Dongchen Han**[1*]**, Tianyu Li**[2*]**, Ziyi Wang**[1]**, Gao Huang**[1†]

[1]Department of Automation, BNRist, Tsinghua University
[2]Institute for Interdisciplinary Information Sciences, Tsinghua University

## Abstract

The self-attention mechanism has been a key factor in the advancement of vision Transformers. However, its quadratic complexity imposes a heavy computational burden in high-resolution scenarios, restricting the practical application. Previous methods attempt to mitigate this issue by introducing handcrafted patterns such as locality or sparsity, which inevitably compromise model capacity. In this paper, we present a novel attention paradigm termed **Circulant Attention** by exploiting the inherent efficient pattern of self-attention. Specifically, we first identify that the self-attention matrix in vision Transformers often approximates the Block Circulant matrix with Circulant Blocks (BCCB), a kind of structured matrix whose multiplication with other matrices can be performed in $\mathcal{O}(N \log N)$ time. Leveraging this interesting pattern, we explicitly model the attention map as its nearest BCCB matrix and propose an efficient computation algorithm for fast calculation. The resulting approach closely mirrors vanilla self-attention, differing only in its use of BCCB matrices. Since our design is inspired by the inherent efficient paradigm, it not only delivers $\mathcal{O}(N \log N)$ computation complexity, but also largely maintains the capacity of standard self-attention. Extensive experiments on diverse visual tasks demonstrate the effectiveness of our approach, establishing circulant attention as a promising alternative to self-attention for vision Transformer architectures.

**Appendix**: github.com/LeapLabTHU/Circulant-Attention

## 1 Introduction

Transformer models have rapidly gained prominence in the field of computer vision in recent years. The superior capacity of self-attention enables vision Transformers to effectively learn from large-scale data, achieving significant success in image classification (Dosovitskiy et al. 2021), object detection (Carion et al. 2020), semantic segmentation (Xie et al. 2021), and multimodal tasks (Xia et al. 2024).

However, integrating self-attention into vision architectures also poses a challenge. The quadratic complexity $\mathcal{O}(N^2)$ of self-attention leads to prohibitively high computational cost when applied over a global receptive field. Previous works (Wang et al. 2021; Liu et al. 2021; Zhu et al.
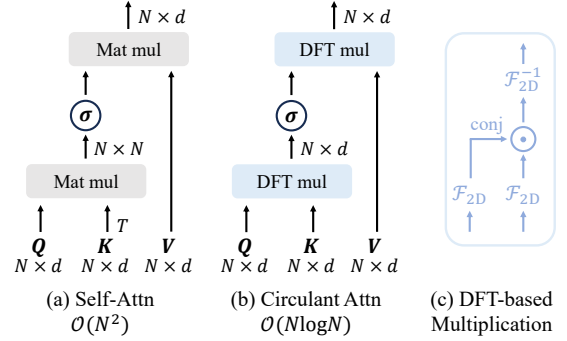
---

Figure 1: An illustration of vanilla self-attention and the proposed circulant attention. The $\sigma$ represents Softmax function, $\odot$ is the Hadamard product, and $\mathcal{F}_{2D}, \mathcal{F}_{2D}^{-1}$ denote the 2D discrete Fourier transform (DFT) and its inversion, respectively. Our circulant attention largely inherits the paradigm of self-attention, except for employing BCCB attention matrices. This simple modification enables our design to be efficiently calculated through DFT-based multiplication, thereby achieving $\mathcal{O}(N \log N)$ complexity. The scaling factor and head-wise summation are omitted for simplicity. Please refer to Section 4 for details.

2023; Han et al. 2023; Wang et al. 2025) address this challenge by introducing *handcrafted* patterns, such as restricting receptive fields or introducing sparsity. While effectively reducing computational demands, these handcrafted designs practically function as *external constraints* imposed on self-attention mechanism, which inevitably compromise long-range modeling capability and limit scalability.

In this paper, we identify an interesting phenomenon that the attention maps in vision Transformers frequently approximate a kind of special matrix: the Block Circulant matrix with Circulant Blocks (BCCB). This kind of matrix is an extension of the circulant matrix in 2D scenarios, whose multiplication with other matrices can be efficiently implemented by 2D discrete Fourier transform (DFT). This indicates that while standard self-attention operates at $\mathcal{O}(N^2)$ computational cost, it *inherently* learns *efficient* patterns that can be calculated with $\mathcal{O}(N \log N)$ complexity. This motivates us to rethink the design of self-attention and come up with a compelling research question:

*Can we explicitly set the attention map as a BCCB matrix to facilitate efficient computation, while preserving the high expressiveness of vanilla self-attention?*

To answer this question, we delve into the essence of self-attention operation, presenting a novel paradigm named **Circulant Attention** to fully exploit the observed pattern. Specifically, we explicitly transform attention maps into BCCB matrices by vertically projecting the original self-attention matrices onto the BCCB matrix subspace. We demonstrate that this mathematical reformulation enables efficient computation of attention scores and output features via 2D discrete Fourier transform, thus reducing the computation complexity from $\mathcal{O}(N^2)$ to $\mathcal{O}(N \log N)$ with the fast Fourier transform algorithm (FFT). As illustrated in Fig. 1, the resulting method highly resembles vanilla self-attention, except that it directly produces BCCB attention matrices and replaces dense matrix multiplication with DFT-based operations. Building on this design, we incorporate a token reweighting module to overcome the inherent limitations of the BCCB structure, which further increases model capacity.

Extensive experiments on diverse visual recognition tasks are conducted to validate the effectiveness of our design, including image classification, object detection, and semantic segmentation. The results confirm that the proposed Circulant Attention offers high efficiency and expressiveness, suggesting it as a promising alternative to self-attention for vision Transformer designs.

Our main contributions and takeaways are as follows:

- We reveal that the attention maps in vision Transformers frequently approximate the Block Circulant matrix with Circulant Blocks (BCCB), an extension of the circulant matrix in 2D scenarios. We provide detailed analyses and visualizations of this phenomenon.
- We present a novel mechanism dubbed Circulant Attention, which utilizes the observed BCCB pattern to achieve efficient computation with $\mathcal{O}(N \log N)$ complexity. Our method serves as a plug-in attention module and can be applied to various vision Transformer models.
- Extensive experiments across image classification, object detection, and semantic segmentation confirm that circulant attention delivers a favorable balance between efficiency and expressiveness, establishing it as a promising alternative to the widely employed self-attention.

## 2   Related Work

**Vision Transformer.** Transformer architectures and attention mechanisms have seen significant advances in computer vision in recent years (Vaswani et al. 2017). However, the quadratic computation complexity of self-attention leads to unmanageable cost when processing global feature maps. To address this problem, various approaches have been proposed to reduce the computational overhead by introducing handcrafted patterns, such as locality or sparsity. Local methods such as Swin Transformer (Liu et al. 2021) restrict attention to non-overlapping windows, reducing computational cost. CSwin (Dong et al. 2022) extends this idea with cross-shaped windows for richer context. Neighborhood Attention Transformer (NAT) (Hassani et al. 2023)

| Notations | Descriptions |
|---|---|
| $\mathcal{F}_{1D}, \mathcal{F}_{1D}^{-1}$ | 1D discrete Fourier transform (DFT) and inversion. |
| $\mathcal{F}_{2D}, \mathcal{F}_{2D}^{-1}$ | 2D discrete Fourier transform (DFT) and inversion. |
| $\overline{(\cdot)}$ | The complex conjugate. |
| $\sigma$ | The Softmax operation on each row. |
| $\odot$ | Hadamard product, element-wise product. |
| $\circledast$ | The DFT-based matrix multiplication we defined. |
| $\| \cdot \|$ | The Frobenius norm. |
| $\langle \cdot, \cdot \rangle$ | The Frobenius inner product. |

Table 1: Important notations used in this paper.

further mimics convolution by limiting attention to local neighborhoods of each query. Apart from these designs, another line of research employs sparse attention paradigms. PVT (Wang et al. 2021) employs downsampling of keys and values to reduce computational complexity, and DAT (Xia et al. 2022) presents an input-dependent sparse attention pattern. BiFormer (Zhu et al. 2023) uses bi-level routing attention to dynamically determine areas of interest for each query. Despite their efficiency, these handcrafted attention patterns inevitably compromise the expressiveness of global self-attention. In this paper, we leverage an intrinsic efficient structure of self-attention to strike a favorable balance between efficiency and expressiveness.

**Efficient architecture with DFT.** The discrete Fourier transform (DFT) has long been an important tool in digital image processing (Pitas 2000). Recent work applies DFT and the Fast Fourier Transform (FFT) to design efficient network components (Li et al. 2021; Rao et al. 2021; Guibas et al. 2021; Huang et al. 2023; Kong et al. 2023). GFNet (Rao et al. 2021) utilizes the convolution theorem of DFT to build global depth-wise convolution module, achieving $\mathcal{O}(N \log N)$ complexity via FFT. FNO (Li et al. 2021) further generalizes this approach to dense global convolution. AFNO (Guibas et al. 2021) and AFFNet (Huang et al. 2023) achieve equivalent dynamic depth-wise convolution through input-dependent frequency filters. In this paper, we employ 2D DFT and the FFT algorithm to efficiently implement our circulant attention mechanism.

## 3   Preliminaries

This section revisits the formulation of self-attention, circulant matrix, and BCCB matrix. To facilitate reading, we summarize important notations in Table 1.

### 3.1   Attention Formulation

We first briefly review the calculation of vanilla self-attention (Vaswani et al. 2017) in vision Transformer models. Consider an image token sequence $x \in \mathbb{R}^{N \times C}$, where $N = H \times W$ and $H, W, C$ are the height, width, and dimension of the feature map, respectively. In each attention head, $x$ is transformed into $Q, K, V \in \mathbb{R}^{N \times d}$ through projection matrices $W_{Q/K/V} \in \mathbb{R}^{C \times d}$, where $d$ is the head dimension. Based on this, self-attention computes the attention weights

and calculates the output as a weighted sum of values using normalized attention score:

$$A = QK^\top/\sqrt{d}, \ O = \sigma(A)V, \qquad (1)$$

where $A \in \mathbb{R}^{N \times N}$ is the raw attention matrix and $\sigma$ represents the Softmax function.

## 3.2 Circulant Matrix

An $N \times N$ matrix $C$ is a circulant matrix if and only if each row is a cyclic shift of the previous one. It has the form:

$$C = \begin{pmatrix} c_0 & c_1 & \cdots & c_{N-1} \\ c_{N-1} & c_0 & \cdots & c_{N-2} \\ \vdots & \vdots & \ddots & \vdots \\ c_1 & c_2 & \cdots & c_0 \end{pmatrix}. \qquad (2)$$

This type of matrix can be fully determined by its first row $c = [c_0, c_1, \ldots, c_{N-1}]$, where $C_{i,j} = c_{j-i \pmod N}$.

The multiplication of a circulant matrix $C \in \mathbb{R}^{N \times N}$ and a vector $x \in \mathbb{R}^N$ can be expressed as:

$$(Cx)_i = \sum_{j=0}^{N-1} C_{i,j}x_j = \sum_{j=0}^{N-1} c_{(j-i) \bmod N} \cdot x_j. \qquad (3)$$

Let $k = (j - i) \bmod N$, which implies $j = (i + k) \bmod N$. We can rewrite the expression as:

$$(Cx)_i = \sum_{k=0}^{N-1} c_k \cdot x_{(i+k) \bmod N}. \qquad (4)$$

This is precisely the definition of the 1D circular cross-correlation between $c$ and $x$, which is the 1D depth-wise convolution with circular padding in deep learning (where the concept of convolution does not involve flipping the kernel). Therefore, the multiplication $y = Cx$ can be achieved with the *Cross-Correlation Theorem* (Wang 2019), which states that the Fourier transform of a cross-correlation result is equivalent to the element-wise product of the first signal's conjugated Fourier transform and the second signal's Fourier transform. Mathematically:

$$\mathcal{F}_{1D}(Cx) = \overline{\mathcal{F}_{1D}(c)} \odot \mathcal{F}_{1D}(x). \qquad (5)$$

Thus, we have:

$$Cx = \mathcal{F}_{1D}^{-1}\left(\overline{\mathcal{F}_{1D}(c)} \odot \mathcal{F}_{1D}(x)\right), \qquad (6)$$

where $\mathcal{F}_{1D}, \mathcal{F}_{1D}^{-1}$ denotes the 1D discrete Fourier transform (DFT) and its inversion, $\odot$ is the element-wise (Hadamard) product, and $\overline{(\cdot)}$ represents the complex conjugate. Leveraging the fast Fourier transform algorithm, we can compute $y = Cx$ with a time complexity of $O(N \log N)$.

## 3.3 BCCB Matrix

Block Circulant matrix with Circulant Blocks (BCCB) is the 2D generalization of circulant matrix (Davis 1979). A BCCB matrix $B \in \mathbb{R}^{N \times N}, N = H \times W$ has a block circulant structure with $H \times H$ blocks:

$$B = \begin{pmatrix} C_0 & C_1 & \cdots & C_{H-1} \\ C_{H-1} & C_0 & \cdots & C_{H-2} \\ \vdots & \vdots & \ddots & \vdots \\ C_1 & C_2 & \cdots & C_0 \end{pmatrix}, \qquad (7)$$
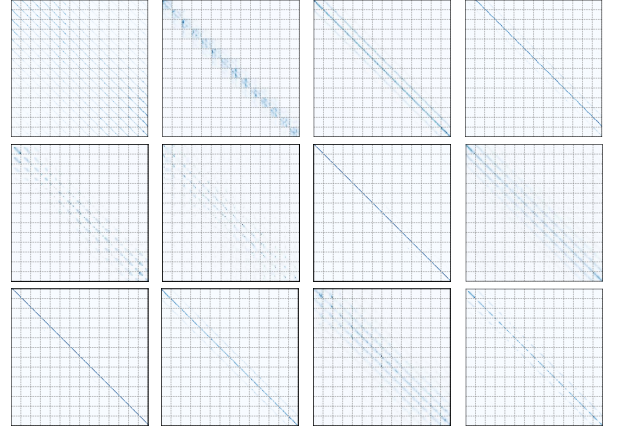


Figure 2: $N \times N$ attention maps from DeiT (Touvron et al. 2021), where $N = H \times W, H = W = 14$. Without handcrafted constraints, the self-attention module in vision Transformer learns near BCCB patterns. For better observation, we divide the matrix into $H \times H$ blocks of shape $W \times W$ using gray dashed lines. Zoom in for best view.

where each block $C_i$ is a circulant matrix of size $W \times W$. Similar to the circulant matrix, the BCCB matrix $B$ is fully determined by its first row $b = [c_0, c_1, \ldots, c_{HW-1}]$. Let $\hat{b}, \hat{x} \in \mathbb{R}^{H \times W}$ be the 2D reshaped versions of $b, x$, respectively. As proved in the Appendix, the multiplication $y = Bx$ is equivalent to the 2D circular cross-correlation between $\hat{b}$ and $\hat{x}$, i.e., the 2D depth-wise convolution with circular padding in deep learning. Similar to the 1D scenario, this operation can be implemented by 2D DFT (Davis 1979):

$$Bx = \mathcal{F}_{2D}^{-1}\left(\overline{\mathcal{F}_{2D}(b)} \odot \mathcal{F}_{2D}(x)\right) \triangleq b \circledast x, \qquad (8)$$

where $\mathcal{F}_{2D}, \mathcal{F}_{2D}^{-1}$ denotes the 2D DFT and its inversion. The $\circledast$ is our defined DFT-based multiplication. For simplicity, the reshaping operations between 1D $\mathbb{R}^N$ sequences and 2D $\mathbb{R}^{H \times W}$ feature maps are not demonstrated, and we define the inputs/outputs of $\mathcal{F}_{2D}, \mathcal{F}_{2D}^{-1}$ to be 1D sequences.

# 4 Method

## 4.1 Efficient Pattern in Vision Transformer

The self-attention mechanism calculates pairwise similarities between each query and key at quadratic complexity $\mathcal{O}(N^2)$. However, we find an interesting phenomenon that, in practice, vision Transformer tends to learn attention patterns that are much more structured and efficient. Specifically, in Fig. 2, we visualize several attention matrices extracted from the DeiT (Touvron et al. 2021) model. These matrices closely resemble block circulant matrices with circulant blocks (BCCB). As discussed in Section 3, multiplication by a BCCB matrix is equivalent to a 2D global convolution operation, which can be carried out in $\mathcal{O}(N \log N)$ time via fast Fourier transform algorithm. This suggests that although self-attention formally incurs quadratic complexity, it implicitly inherits an efficient structure.
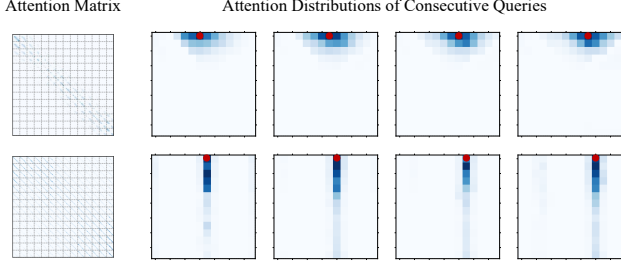
Figure 3: The $N \times N$ attention matrices and $H \times W$ attention distributions of consecutive queries, where $N = H \times W, H = W = 14$. The red points correspond to the query tokens. When exhibiting BCCB pattern, the attention distributions show convolution-like translation invariance.

For a better understanding of this phenomenon, we provide the attention distributions of adjacent query tokens (marked in red) in the attention matrices of Fig. 3. A key observation is that the attention distributions exhibit an approximate shift invariance corresponding to the movement of the query positions, which is exactly the behavior of a 2D global convolution implemented by a BCCB matrix. This further confirms that the learned attention in vision Transformer exhibits strong similarity to BCCB matrix.

Motivated by the above discoveries, a critical research question emerges: *Can we explicitly enforce a BCCB structure on the attention matrix?*

In the next section, we answer this question with **Circulant Attention**, a novel attention mechanism that employs BCCB attention matrices. By imposing this structural prior, we aim to explicitly encode the beneficial properties of BCCB matrices, i.e., $\mathcal{O}(N \log N)$ complexity and inherent shift invariance, into the attention mechanism, thereby benefiting from both efficiency and expressiveness.

## 4.2 Circulant Attention

Our circulant attention approximates the original attention map $A \in \mathbb{R}^{N \times N}$ using its nearest Block Circulant with Circulant Blocks (BCCB) matrix $\tilde{A}$. Formally, we have

$$\tilde{A} = \arg\min_{B \in \mathcal{B}} \|A - B\|, \quad (9)$$

where $\mathcal{B}$ denotes the $N \times N$ BCCB matrix subspace and $\|\cdot\|$ represents the Frobenius norm. Therefore, $\tilde{A}$ is the orthogonal projection of $A$ in the BCCB matrix subspace. As discussed in Section 3, a $N \times N$ BCCB matrix could be fully determined by its first row. Let $B_k$ denotes the BCCB matrix whose first row is a one-hot vector with 1 at the $k$-th position. It is easy to see that $\{B_0, \cdots, B_{N-1}\}$ forms the basis for BCCB matrix subspace. Furthermore, we have

$$\langle B_k, B_k \rangle = N, \langle B_k, B_j \rangle = 0, k \neq j, \quad (10)$$

where $\langle \cdot, \cdot \rangle$ denotes the Frobenius inner product, i.e., the sum of element-wise products of two matrices. Eq. (10) indicates that $\{B_0, \cdots, B_{N-1}\}$ is an orthogonal basis. Therefore, the orthogonal projection $\tilde{A}$ in the BCCB matrix subspace can

be expressed as:

$$\begin{aligned} \tilde{A} &= \sum_{k=0}^{N-1} \frac{\langle A, B_k \rangle}{\langle B_k, B_k \rangle} B_k \\ &= \frac{1}{N} \sum_{k=0}^{N-1} \langle A, B_k \rangle B_k \quad (11) \\ &= \frac{1}{N} \sum_{k=0}^{N-1} \langle \frac{QK^\top}{\sqrt{d}}, B_k \rangle B_k. \end{aligned}$$

With the BCCB attention matrix $\tilde{A}$, our circulant attention computes the output as $O = \sigma(\tilde{A})V$, akin to the vanilla self-attention mechanism (see Eq. (1)).

**Efficient computation algorithm.** The proposed circulant attention still incurs $\mathcal{O}(N^2)$ complexity when computed using the projection formula above. Here, we demonstrate that our method can be equivalently calculated in $\mathcal{O}(N \log N)$ time employing the 2D DFT.

Specifically, since $\tilde{A}$ is a BCCB matrix, its entire structure is determined by the first row. We denote this row as $a$, which can be expressed as:

$$a = \frac{1}{N\sqrt{d}}[\langle QK^\top, B_0 \rangle, \cdots, \langle QK^\top, B_{N-1} \rangle]. \quad (12)$$

According to the properties of the BCCB matrix, $B_k$ actually corresponds to a spatial shift in the 2D space. Define

$$\Delta h = \lfloor k/W \rfloor, \ \Delta w = k \bmod W. \quad (13)$$

Then we have

$$\begin{aligned} a_k &= \frac{1}{N\sqrt{d}} \langle QK^\top, B_k \rangle \\ &= \frac{1}{N\sqrt{d}} \sum_{h=0}^{H-1} \sum_{w=0}^{W-1} \langle \hat{Q}_{h,w}, \hat{K}_{h+\Delta h, w+\Delta w} \rangle, \end{aligned} \quad (14)$$

where $\hat{Q}, \hat{K} \in \mathbb{R}^{H \times W \times d}$ are the 2D reshaped versions of $Q, K \in \mathbb{R}^{N \times d}$ with circular padding (to handle boundary conditions). It could be observed that $a_k$ is the circular cross-correlation result between $\hat{Q}$ and $\hat{K}$, which is exactly the convolution concept in deep learning (see Section 3). Therefore, $a$ can be computed efficiently using 2D DFT:

$$\begin{aligned} a &= \frac{1}{N\sqrt{d}} \left[ \mathcal{F}_{2D}^{-1} \left( \overline{\mathcal{F}_{2D}(Q)} \odot \mathcal{F}_{2D}(K) \right) \right] \cdot \mathbf{1}_{d \times 1} \\ &= \frac{1}{N\sqrt{d}} (Q \circledast K) \cdot \mathbf{1}_{d \times 1}, \end{aligned} \quad (15)$$

where $\circledast$ is the DFT-based matrix multiplication we defined in Section 3.3 and $Q \circledast K \in \mathbb{R}^{N \times d}$. The $\mathbf{1}_{d \times 1} \in \mathbb{R}^{d \times 1}$ is an all-one vector, which corresponds to the implicit summation over dimension $d$ of the inner product between $\hat{Q}_{h,w}$ and $\hat{K}_{h+\Delta h, w+\Delta w}$ in Eq. (14).

Since $\tilde{A}$ is a BCCB matrix and $\sigma$ is the Softmax function on each row, $\sigma(\tilde{A})$ is also a BCCB matrix and $\sigma(a)$ represents its first row. As discussed in Section 3, the output $O = \sigma(\tilde{A})V$ can thus be efficiently computed with:

$$\begin{aligned} O &= \mathcal{F}_{2D}^{-1} \left( \mathcal{F}_{2D}(\sigma(a)) \odot \mathcal{F}_{2D}(V) \right) \\ &= \sigma(a) \circledast V. \end{aligned} \quad (16)$$

| Method | Reso | #Params | FLOPs | Top-1 |
|---|---|---|---|---|
| DeiT-T (Touvron et al. 2021) | $224^2$ | 5.7M | 1.2G | 72.2 |
| **CA-DeiT-T** | $224^2$ | 6.1M | 1.2G | **75.0** (+2.8) |
| DeiT-S (Touvron et al. 2021) | $224^2$ | 22.1M | 4.6G | 79.8 |
| **CA-DeiT-S** | $224^2$ | 23.8M | 4.8G | **81.0** (+1.2) |
| DeiT-B (Touvron et al. 2021) | $224^2$ | 86.6M | 17.6G | 81.8 |
| **CA-DeiT-B** | $224^2$ | 93.6M | 18.9G | **82.3** (+0.5) |
| PVT-T (Wang et al. 2021) | $224^2$ | 13.2M | 1.9G | 75.1 |
| **CA-PVT-T** | $224^2$ | 12.2M | 2.0G | **78.1** (+3.0) |
| PVT-S (Wang et al. 2021) | $224^2$ | 24.5M | 3.8G | 79.8 |
| **CA-PVT-S** | $224^2$ | 22.8M | 4.0G | **81.7** (+1.9) |
| PVT-M (Wang et al. 2021) | $224^2$ | 44.2M | 6.7G | 81.2 |
| **CA-PVT-M** | $224^2$ | 42.5M | 6.8G | **82.6** (+1.4) |
| PVT-L (Wang et al. 2021) | $224^2$ | 61.4M | 9.8G | 81.7 |
| **CA-PVT-L** | $224^2$ | 58.6M | 10.1G | **82.9** (+1.2) |
| Swin-T (Liu et al. 2021) | $224^2$ | 29M | 4.5G | 81.3 |
| **CA-Swin-T** | $224^2$ | 28M | 4.6G | **82.2** (+0.9) |
| Swin-S (Liu et al. 2021) | $224^2$ | 50M | 8.7G | 83.0 |
| **CA-Swin-S** | $224^2$ | 50M | 8.8G | **83.6** (+0.6) |
| Swin-B (Liu et al. 2021) | $224^2$ | 88M | 15.4G | 83.5 |
| **CA-Swin-B** | $224^2$ | 88M | 15.7G | **83.9** (+0.4) |
| Swin-B (Liu et al. 2021) | $384^2$ | 88M | 47.0G | 84.5 |
| **CA-Swin-B** | $384^2$ | 88M | 47.1G | **85.1** (+0.6) |

| Method | Reso | #Params | FLOPs | Top-1 |
|---|---|---|---|---|
| SLAB-T (Guo et al. 2024) | $224^2$ | 29M | 4.5G | 81.8 |
| VMamba-T (Liu et al. 2024) | $224^2$ | 31M | 4.9G | 82.5 |
| SOFT-S++ (Lu et al. 2024) | $224^2$ | 27M | 4.5G | 82.6 |
| PolaFormer-T (Meng et al. 2025) | $224^2$ | 29M | 4.5G | 82.6 |
| LocalVMamba-T (Huang et al. 2024) | $224^2$ | 26M | 5.7G | 82.7 |
| VVT-S (Sun et al. 2023) | $224^2$ | 26M | 5.6G | 82.7 |
| Agent-CSwin-T (Han et al. 2024c) | $224^2$ | 21M | 4.3G | 83.1 |
| FasterViT-1 (Hatamizadeh et al. 2024) | $224^2$ | 53M | 5.3G | 83.2 |
| EfficientViT-B3 (Cai et al. 2023) | $224^2$ | 49M | 4.0G | 83.5 |
| **CAT-T** | $224^2$ | 27M | 4.3G | **83.6** |
| LocalVMamba-S (Huang et al. 2024) | $224^2$ | 50M | 11.4G | 83.7 |
| $QFormer_h$-S (Zhang et al. 2024) | $224^2$ | 51M | 8.9G | 84.0 |
| BiFormer-B (Zhu et al. 2023) | $224^2$ | 57M | 9.8G | 84.3 |
| MILA-S (Han et al. 2024b) | $224^2$ | 43M | 7.3G | 84.4 |
| TransXNet-B (Lou et al. 2025) | $224^2$ | 48M | 8.3G | 84.6 |
| **CAT-S** | $224^2$ | 51M | 7.9G | **84.5** |
| VMamba-B (Liu et al. 2024) | $224^2$ | 89M | 15.4G | 83.9 |
| InLine-Swin-B (Han et al. 2024a) | $224^2$ | 88M | 15.4G | 84.1 |
| SOFT-L++ (Lu et al. 2024) | $224^2$ | 85M | 15.4G | 84.1 |
| FasterViT-3 (Hatamizadeh et al. 2024) | $224^2$ | 160M | 18.2G | 84.9 |
| OverLock-B (Lou and Yu 2025) | $224^2$ | 95M | 16.7G | 85.1 |
| **CAT-B** | $224^2$ | 90M | 15.2G | **85.0** |

Table 2: Comparison with baseline models (left) and advanced methods (right) on ImageNet-1K.

In Fig. 1, we provide an illustration of the proposed circulant attention paradigm. It can be observed that our design is structurally similar to vanilla self-attention, with the key differences being the use of BCCB attention matrices and DFT-based efficient multiplication.

**Complexity analysis.** Leveraging the efficient computation algorithm, the overall complexity of our circulant attention is expressed as:

$$\Omega(\text{CA}) = \underbrace{2N(\log_2 N)d + 2Nd + N\log_2 N}_{\text{DFT, Product, IDFT in Eq. (7)}} +$$

$$\underbrace{N(\log_2 N)(d+1) + 2Nd + N(\log_2 N)d}_{\text{DFT, Product, IDFT in Eq. (8)}} \quad (17)$$

$$= N(\log_2 N)(4d+2) + 4Nd,$$

which is much more efficient than the computation complexity of vanilla self-attention:

$$\Omega(\text{SA}) = N^2 d + N^2 d = 2N^2 d. \quad (18)$$

**Token reweighting module.** In standard self-attention, a row-wise Softmax is applied to the raw attention map $A \in \mathbb{R}^{N \times N}$, ensuring that each row sums to one. Notably, this operation does not constrain the column sums of $\sigma(A)$, allowing certain keys to accumulate higher total attention scores across queries. Consequently, different queries can emphasize similar keys and values, helping the model distinguish more informative tokens from others. However, the block circulant structure of our circulant attention map $\sigma(\tilde{A})$ enforces both row and column sums to equal one, thereby limiting its ability to highlight salient tokens. As a remedy, we introduce a simple yet effective token reweighting method to further improve our design. Specifically, there are two ways to implement this design: the pre-reweighting

$$O = \text{CirAttn}(Q, K, V \odot T), \quad (19)$$

and the post-reweighting

$$O = \text{CirAttn}(Q, K, V) \odot T, \quad (20)$$

where $T = \text{SiLU}(xW_T) \in \mathbb{R}^{N \times d}$ is an input-dependent token reweighting factor, and CirAttn represents our circulant attention operator.

### 4.3 Implementation

Our circulant attention serves as a plug-in module and can be applied to various vision Transformer models. As a showcase, we employ three representative Transformer architectures: DeiT (Touvron et al. 2021), PVT (Wang et al. 2021), and Swin Transformer (Liu et al. 2021) to implement our approach. These three models represent global self-attention, sparse attention and local attention, respectively. We replace the original attention module in these models with circulant attention to establish our models. Since our method benefits from $N \log N$ complexity, it is possible to directly process a high-resolution feature map in the early stages without incurring high computational cost. Therefore, for hierarchical models like Swin and PVT, we mainly restrict the attention replacement to the first two stages. Beyond baseline improvements, we have also developed a family of specialized models, termed Circulant Attention Transformer (CAT), to compare with various advanced vision Transformers. Detailed model architectures are shown in the Appendix.

**Mask R-CNN Object Detection on COCO**

| Backbone | FLOPs | Sch. | $AP^b$ | $AP^b_{50}$ | $AP^b_{75}$ | $AP^m$ | $AP^m_{50}$ | $AP^m_{75}$ |
|---|---|---|---|---|---|---|---|---|
| PVT-T | 240G | 1x | 36.7 | 59.2 | 39.3 | 35.1 | 56.7 | 37.3 |
| **CA-PVT-T** | 218G | 1x | 40.5 | 63.4 | 43.9 | 37.9 | 60.4 | 40.7 |
| PVT-S | 305G | 1x | 40.4 | 62.9 | 43.8 | 37.8 | 60.1 | 40.3 |
| **CA-PVT-S** | 269G | 1x | 44.2 | 66.9 | 48.3 | 40.7 | 63.9 | 43.7 |
| PVT-M | 392G | 1x | 42.0 | 64.4 | 45.6 | 39.0 | 61.6 | 42.1 |
| **CA-PVT-M** | 356G | 1x | 45.2 | 67.7 | 49.4 | 41.2 | 64.7 | 44.1 |
| PVT-L | 494G | 1x | 42.9 | 65.0 | 46.6 | 39.5 | 61.9 | 42.5 |
| **CA-PVT-L** | 444G | 1x | 46.3 | 68.8 | 50.8 | 41.9 | 65.4 | 45.0 |
| Swin-T | 267G | 1x | 43.7 | 66.6 | 47.7 | 39.8 | 63.3 | 42.7 |
| **CA-Swin-T** | 269G | 1x | 44.5 | 67.3 | 48.9 | 40.5 | 64.2 | 43.5 |
| Swin-S | 358G | 1x | 45.7 | 67.9 | 50.4 | 41.1 | 64.9 | 44.2 |
| **CA-Swin-S** | 361G | 1x | 46.8 | 69.3 | 51.7 | 42.2 | 66.2 | 45.3 |
| Swin-B | 503G | 1x | 46.9 | - | - | 42.3 | - | - |
| **CA-Swin-B** | 507G | 1x | 47.5 | 70.1 | 52.2 | 42.8 | 66.8 | 45.9 |
| Swin-T | 267G | 3x | 46.0 | 68.1 | 50.3 | 41.6 | 65.1 | 44.9 |
| **CA-Swin-T** | 269G | 3x | 46.7 | 68.8 | 51.2 | 42.5 | 65.9 | 45.9 |

Table 3: Results on COCO dataset. The FLOPs are computed with an input resolution of 1280×800.

**Semantic Segmentation on ADE20K**

| Backbone | Method | FLOPs | #Params | mIoU |
|---|---|---|---|---|
| PVT-T | S-FPN | 158G | 17M | 35.7 |
| **CA-PVT-T** | S-FPN | 135G | 16M | 39.4 |
| PVT-S | S-FPN | 225G | 28M | 39.8 |
| **CA-PVT-S** | S-FPN | 187G | 26M | 42.3 |
| PVT-M | S-FPN | 315G | 48M | 41.6 |
| **CA-PVT-M** | S-FPN | 278G | 46M | 43.7 |
| PVT-L | S-FPN | 420G | 65M | 42.1 |
| **CA-PVT-L** | S-FPN | 369G | 62M | 44.2 |
| Swin-T | UperNet | 945G | 60M | 44.5 |
| **CA-Swin-T** | UperNet | 947G | 59M | 45.2 |
| Swin-S | UperNet | 1038G | 81M | 47.6 |
| **CA-Swin-S** | UperNet | 1040G | 80M | 48.6 |

Table 4: Results of semantic segmentation. The FLOPs are computed over encoders and decoders with an input image at the resolution of 512×2048. S-FPN is short for SemanticFPN (Kirillov et al. 2019) model.
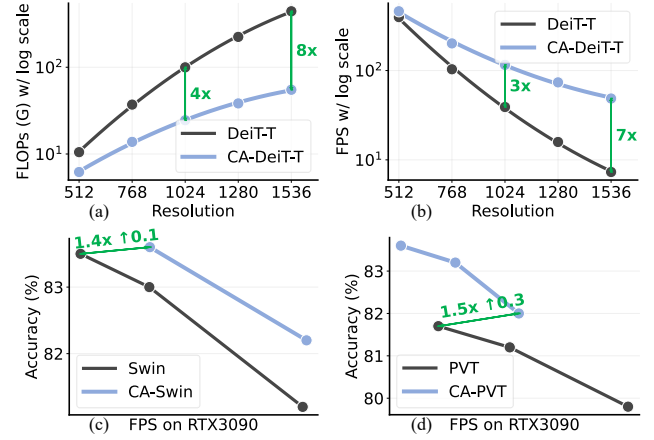


Figure 4: Comparisons between self-attention and the proposed circulant attention in (a) FLOPs, (b) inference FPS, and (c, d) throughput-accuracy trade-off. Throughput is measured on a RTX3090 GPU.

# 5 Experiments

In this section, we conduct extensive experiments to evaluate the effectiveness of our method, including ImageNet classification (Deng et al. 2009), COCO object detection (Lin et al. 2014), and ADE20K semantic segmentation (Zhou et al. 2019). Visualization and analysis are also provided.

## 5.1 ImageNet Classification

The ImageNet-1K (Deng et al. 2009) dataset comprises 1.28 million training and 50 thousand validation images spanning 1,000 classes. To ensure a fair comparison, we adopt identical training settings as the baseline models. Specifically, models are trained from scratch for 300 epochs using the AdamW (Loshchilov and Hutter 2018) optimizer, with cosine learning rate decay and a linear warm-up over the first 20 epochs. The initial learning rate is set to $1 \times 10^{-3}$, and the weight decay is 0.05. Augmentation and regularization strategies include RandAugment (Cubuk et al. 2020), Mixup (Zhang et al. 2018), CutMix (Yun et al. 2019) and random erasing (Zhong et al. 2020).

From Table 2 left, we observe that replacing self-attention in the three representative models with circulant attention leads to consistent improvements. For example, the CA-DeiT outperforms the global self-attention baseline DeiT (Touvron et al. 2021). This indicates that our method not only enjoys high efficiency, but also facilitates the learning of vision Transformers. Furthermore, compared to the sparse attention PVT (Wang et al. 2021) and local attention Swin Transformer (Liu et al. 2021), circulant attention demonstrates significant advantages. Our CA-PVT-S matches PVT-L's accuracy using 30% of the parameters and 40% of the FLOPs. These results establish circulant attention as a promising alternative to self-attention.

Table 2 right compares our CAT model against various state-of-the-art vision Transformers and Mamba variants.

We see that circulant attention delivers better results than highly optimized methods, validating the superior capacity of our approach.

## 5.2 Object Detection

The COCO (Lin et al. 2014) object detection and instance segmentation dataset contains 118K training images and 5K validation images. We follow the training and testing strategies of the baseline models. The results are shown in Table 3. Circulant attention offers effective global modeling with $\mathcal{O}(N \log N)$ complexity, making it ideally suitable for high-resolution image modeling scenarios. Notably, CA-PVT-S outperforms the larger PVT-L model by 1.3 box AP with substantially fewer FLOPs. This gain is more pronounced than in the classification task, where CA-PVT-S and PVT-L yield the same accuracy. These results highlight the superiority of our method in high-resolution scenarios.
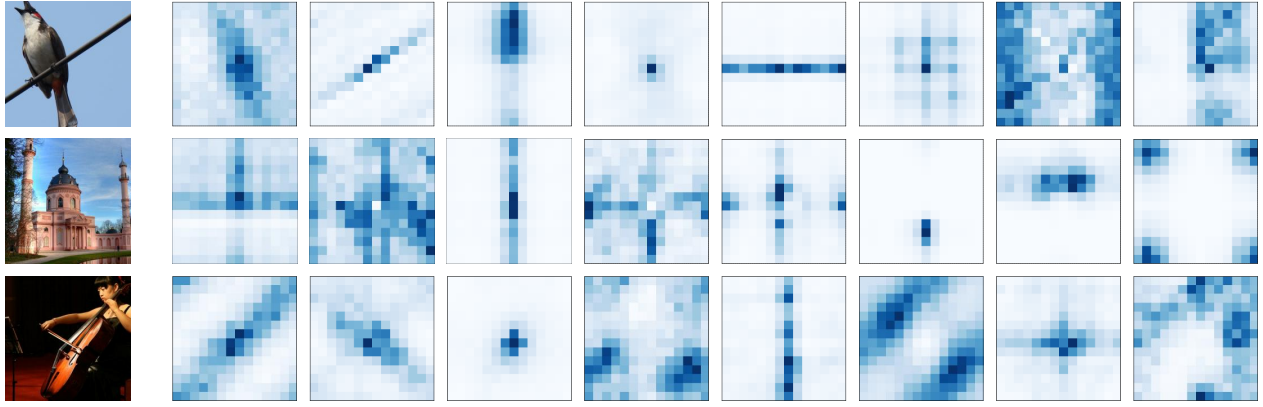
Figure 5: An illustration of the equivalent global convolution kernels from CA-DeiT model.

## 5.3 Semantic Segmentation

ADE20K (Zhou et al. 2019) is a well-established benchmark for semantic segmentation, consisting of 20K training images, 2K validation images, and 150 semantic categories. We use the same settings as baseline models. Similar to the object detection task, the results in Table 4 demonstrate that circulant attention consistently enhances the performance of PVT and Swin models. We observe up to 3.7% mIoU gain with comparable or less computation cost and parameters.

## 5.4 Analysis and Visualization

**Efficiency analysis.** Benefiting from $\mathcal{O}(N \log N)$ complexity, our circulant attention shows remarkable efficiency advantages over self-attention operation. As demonstrated in Fig. 4(a), CA-DeiT-T requires $8\times$ fewer FLOPs at a $1536^2$ image resolution, highlighting its potential for high-resolution modeling tasks. Moreover, these theoretical savings directly translate into practical gains, with our method delivering a $7\times$ speedup at $1536^2$ resolution in Fig. 4(b). Additionally, Fig. 4(c, d) shows that our model achieves a better trade-off between throughput and accuracy at the default ImageNet resolution $224^2$, enjoying up to $1.5\times$ faster inference speed with improved performance.

**Visualization.** As discussed in Section 3.3, the BCCB matrices correspond to 2D global convolution operations. Therefore, it is easy to visualize and interpret the equivalent global convolution kernels. As shown in Fig. 5, our model can generate diverse equivalent kernels based on the input image. For example, given the first image, our model generates bird- and wire-shaped kernels (the first two). There are also kernels focusing more on spatial patterns, such as local, global, half-plane, strip, and cross-shaped ones.

## 5.5 Ablation Study

We conduct ablation studies to assess the contribution of each component in our circulant attention design. Experiments are performed on ImageNet-1K under the CA-DeiT framework. As shown in Table 5, we begin with the DeiT-S baseline and introduce our designs in turn. Simply introducing circulant attention leads to a negligible 0.1% performance drop, indicating that the BCCB pattern is suitable for

| | FLOPs | #Param | Acc. | Diff. |
|---|---|---|---|---|
| DeiT-S (Touvron et al. 2021) | 4.6G | 22M | 79.8 | -1.2 |
| + Circulant Attention | 4.4G | 22M | 79.7 | -1.3 |
| + Head Dim $d = 1$ | 4.4G | 22M | 80.2 | -0.8 |
| + Token Pre-Reweighting | 4.8G | 24M | 80.9 | -0.1 |
| Token Post-Reweighting | 4.8G | 24M | **81.0** | **Ours** |
| DeiT-S + Token Reweighting | 5.0G | 24M | 80.0 | -1.0 |

Table 5: Ablation on the key designs based on DeiT-S.

vision Transformers and does not sacrifice expressiveness. On this basis, we find that our design benefits from a small head dimension and more heads. Setting the head dimension $d = 1$ improves the accuracy to 80.2%, outperforming the DeiT-S baseline. This can be attributed to the fact that circulant attention score is a summation of $N$ query-key pairs (see Eq. (14)), thus having an equivalent head dimension of $Nd$, which is still large enough when $d = 1$. Additionally, we study the two token reweighting methods defined in Eq. (19) and Eq. (20). Both designs lead to further improvement, while the post-reweighting delivers a slightly better result. Notably, token reweighting offers limited gains on the DeiT-S baseline. In this paper, we employ post-reweighting as default, achieving 81.0% accuracy on DeiT-S. These results confirm the effectiveness of our design.

## 6 Conclusion

This paper reveals an interesting pattern in vision Transformers, where the self-attention maps frequently exhibit nearly block circulant with circulant blocks (BCCB) structures. This observation directly inspires our design of Circulant Attention, a novel attention paradigm that explicitly leverages this inherent pattern to optimize computational efficiency while preserving expressive power. Extensive experiments across classification, object detection, and semantic segmentation fully demonstrate the effectiveness of the proposed circulant attention, establishing it as an efficient and competitive alternative to widely adopted self-attention for vision Transformer architectures.

## Acknowledgements

## References

Cai, H.; Li, J.; Hu, M.; Gan, C.; and Han, S. 2023. Efficientvit: Lightweight multi-scale attention for high-resolution dense prediction. In *ICCV*.

Carion, N.; Massa, F.; Synnaeve, G.; Usunier, N.; Kirillov, A.; and Zagoruyko, S. 2020. End-to-end object detection with transformers. In *ECCV*.

Chu, X.; Tian, Z.; Zhang, B.; Wang, X.; and Shen, C. 2023. Conditional Positional Encodings for Vision Transformers. In *ICLR*.

Cubuk, E. D.; Zoph, B.; Shlens, J.; and Le, Q. V. 2020. Randaugment: Practical automated data augmentation with a reduced search space. In *CVPRW*.

Davis, P. J. 1979. *Circulant Matrices*. John Wiley & Sons. ISBN 9780471057710.

Deng, J.; Dong, W.; Socher, R.; Li, L.-J.; Li, K.; and Fei-Fei, L. 2009. Imagenet: A large-scale hierarchical image database. In *CVPR*.

Dong, X.; Bao, J.; Chen, D.; Zhang, W.; Yu, N.; Yuan, L.; Chen, D.; and Guo, B. 2022. Cswin transformer: A general vision transformer backbone with cross-shaped windows. In *CVPR*.

Dosovitskiy, A.; Beyer, L.; Kolesnikov, A.; Weissenborn, D.; Zhai, X.; Unterthiner, T.; Dehghani, M.; Minderer, M.; Heigold, G.; Gelly, S.; et al. 2021. An image is worth 16x16 words: Transformers for image recognition at scale. In *ICLR*.

Guibas, J.; Mardani, M.; Li, Z.; Tao, A.; Anandkumar, A.; and Catanzaro, B. 2021. Adaptive fourier neural operators: Efficient token mixers for transformers. *arXiv preprint arXiv:2111.13587*.

Guo, J.; Chen, X.; Tang, Y.; and Wang, Y. 2024. SLAB: Efficient Transformers with Simplified Linear Attention and Progressive Re-parameterized Batch Normalization. In *ICML*.

Han, D.; Pan, X.; Han, Y.; Song, S.; and Huang, G. 2023. FLatten Transformer: Vision Transformer using Focused Linear Attention. In *ICCV*.

Han, D.; Pu, Y.; Xia, Z.; Han, Y.; Pan, X.; Li, X.; Lu, J.; Song, S.; and Huang, G. 2024a. Bridging the divide: Reconsidering softmax and linear attention. In *NeurIPS*.

Han, D.; Wang, Z.; Xia, Z.; Han, Y.; Pu, Y.; Ge, C.; Song, J.; Song, S.; Zheng, B.; and Huang, G. 2024b. Demystify Mamba in Vision: A Linear Attention Perspective. In *NeurIPS*.

Han, D.; Ye, T.; Han, Y.; Xia, Z.; Song, S.; and Huang, G. 2024c. Agent attention: On the integration of softmax and linear attention. In *ECCV*.

Hassani, A.; Walton, S.; Li, J.; Li, S.; and Shi, H. 2023. Neighborhood attention transformer. In *CVPR*.

Hatamizadeh, A.; Heinrich, G.; Yin, H.; Tao, A.; Alvarez, J. M.; Kautz, J.; and Molchanov, P. 2024. FasterViT: Fast Vision Transformers with Hierarchical Attention. In *ICLR*.

Huang, T.; Pei, X.; You, S.; Wang, F.; Qian, C.; and Xu, C. 2024. Localmamba: Visual state space model with windowed selective scan. In *ECCVW*.

Huang, Z.; Zhang, Z.; Lan, C.; Zha, Z.-J.; Lu, Y.; and Guo, B. 2023. Adaptive frequency filters as efficient global token mixers. In *ICCV*.

Kirillov, A.; Girshick, R.; He, K.; and Dollár, P. 2019. Panoptic feature pyramid networks. In *CVPR*.

Kong, L.; Dong, J.; Ge, J.; Li, M.; and Pan, J. 2023. Efficient frequency domain-based transformers for high-quality image deblurring. In *CVPR*.

Li, Z.; Kovachki, N. B.; Azizzadenesheli, K.; Bhattacharya, K.; Stuart, A.; Anandkumar, A.; et al. 2021. Fourier Neural Operator for Parametric Partial Differential Equations. In *ICLR*.

Lin, T.-Y.; Maire, M.; Belongie, S.; Hays, J.; Perona, P.; Ramanan, D.; Dollár, P.; and Zitnick, C. L. 2014. Microsoft coco: Common objects in context. In *ECCV*.

Liu, Y.; Tian, Y.; Zhao, Y.; Yu, H.; Xie, L.; Wang, Y.; Ye, Q.; Jiao, J.; and Liu, Y. 2024. Vmamba: Visual state space model. In *NeurIPS*.

Liu, Z.; Lin, Y.; Cao, Y.; Hu, H.; Wei, Y.; Zhang, Z.; Lin, S.; and Guo, B. 2021. Swin transformer: Hierarchical vision transformer using shifted windows. In *ICCV*.

Loshchilov, I.; and Hutter, F. 2018. Decoupled weight decay regularization. In *ICLR*.

Lou, M.; and Yu, Y. 2025. OverLoCK: An Overview-first-Look-Closely-next ConvNet with Context-Mixing Dynamic Kernels. In *CVPR*.

Lou, M.; Zhang, S.; Zhou, H.-Y.; Yang, S.; Wu, C.; and Yu, Y. 2025. TransXNet: learning both global and local dynamics with a dual dynamic token mixer for visual recognition. *TNNLS*.

Lu, J.; Zhang, J.; Zhu, X.; Feng, J.; Xiang, T.; and Zhang, L. 2024. Softmax-free linear transformers. *IJCV*.

Meng, W.; Luo, Y.; Li, X.; Jiang, D.; and Zhang, Z. 2025. PolaFormer: Polarity-aware Linear Attention for Vision Transformers. In *ICLR*.

Pitas, I. 2000. *Digital image processing algorithms and applications*. John Wiley & Sons.

Rao, Y.; Zhao, W.; Zhu, Z.; Lu, J.; and Zhou, J. 2021. Global filter networks for image classification. In *NeurIPS*.

Sun, W.; Qin, Z.; Deng, H.; Wang, J.; Zhang, Y.; Zhang, K.; Barnes, N.; Birchfield, S.; Kong, L.; and Zhong, Y. 2023. Vicinity vision transformer. *TPAMI*.

Touvron, H.; Cord, M.; Douze, M.; Massa, F.; Sablayrolles, A.; and Jégou, H. 2021. Training data-efficient image transformers & distillation through attention. In *ICML*.

Vaswani, A.; Shazeer, N.; Parmar, N.; Uszkoreit, J.; Jones, L.; Gomez, A. N.; Kaiser, Ł.; and Polosukhin, I. 2017. Attention is all you need. In *NeurIPS*.

Wang, C. 2019. *Kernel learning for visual perception*. Ph.D. thesis, Columbia University.

Wang, W.; Xie, E.; Li, X.; Fan, D.-P.; Song, K.; Liang, D.; Lu, T.; Luo, P.; and Shao, L. 2021. Pyramid vision transformer: A versatile backbone for dense prediction without convolutions. In *ICCV*.

Wang, Y.; Yue, Y.; Yue, Y.; Wang, H.; Jiang, H.; Han, Y.; Ni, Z.; Pu, Y.; Shi, M.; Lu, R.; et al. 2025. Emulating human-like adaptive vision for efficient and flexible machine visual perception. *Nature Machine Intelligence*.

Xia, Z.; Han, D.; Han, Y.; Pan, X.; Song, S.; and Huang, G. 2024. Gsva: Generalized segmentation via multimodal large language models. In *CVPR*.

Xia, Z.; Pan, X.; Song, S.; Li, L. E.; and Huang, G. 2022. Vision transformer with deformable attention. In *CVPR*.

Xie, E.; Wang, W.; Yu, Z.; Anandkumar, A.; Alvarez, J. M.; and Luo, P. 2021. SegFormer: Simple and efficient design for semantic segmentation with transformers. In *NeurIPS*.

Yun, S.; Han, D.; Oh, S. J.; Chun, S.; Choe, J.; and Yoo, Y. 2019. Cutmix: Regularization strategy to train strong classifiers with localizable features. In *ICCV*.

Zhang, H.; Cisse, M.; Dauphin, Y. N.; and Lopez-Paz, D. 2018. mixup: Beyond Empirical Risk Minimization. In *ICLR*.

Zhang, Q.; Zhang, J.; Xu, Y.; and Tao, D. 2024. Vision transformer with quadrangle attention. *TPAMI*.

Zhong, Z.; Zheng, L.; Kang, G.; Li, S.; and Yang, Y. 2020. Random erasing data augmentation. In *AAAI*.

Zhou, B.; Zhao, H.; Puig, X.; Xiao, T.; Fidler, S.; Barriuso, A.; and Torralba, A. 2019. Semantic understanding of scenes through the ade20k dataset. *IJCV*.

Zhu, L.; Wang, X.; Ke, Z.; Zhang, W.; and Lau, R. W. 2023. BiFormer: Vision Transformer with Bi-Level Routing Attention. In *CVPR*.

## A  BCCB Matrix

Block Circulant matrix with Circulant Blocks (BCCB) is the 2D generalization of circulant matrix (Davis 1979). A BCCB matrix $B \in \mathbb{R}^{N \times N}, N = H \times W$ has a block circulant structure with $H \times H$ blocks:

$$B = \begin{pmatrix} C_0 & C_1 & \cdots & C_{H-1} \\ C_{H-1} & C_0 & \cdots & C_{H-2} \\ \vdots & \vdots & \ddots & \vdots \\ C_1 & C_2 & \cdots & C_0 \end{pmatrix}, \quad (21)$$

where each block $C_i$ is a circulant matrix of size $W \times W$. Similar to the circulant matrix, the BCCB matrix $B$ is fully determined by its first row $b = [c_0, c_1, \ldots, c_{HW-1}]$.

Let $\hat{b}, \hat{x} \in \mathbb{R}^{H \times W}$ be the 2D reshaped versions of $b, x$, respectively. We now prove that the multiplication $y = Bx$ is equivalent to the 2D circular cross-correlation between $\hat{b}$ and $\hat{x}$, i.e., the 2D global convolution with circular padding in deep learning. Notably, the concept of convolution in deep learning does not involve flipping the kernel, thus corresponding to the cross-correlation operation.

We interpret the index $0 \leq k < N$ as a pair of two-dimensional coordinates $(k_x, k_y)$, where $k_x = \lfloor k/W \rfloor$ and $k_y = k \bmod W$. Using this mapping, the matrix entry $B_{i,j}$ can be written as $B_{(i_x,i_y),(j_x,j_y)}$ and $b_i$ can be written as $b_{(i_x,i_y)}$.

Then, for the matrix entry $B_{i,j} = B_{(i_x,i_y),(j_x,j_y)}$, we can see that:

- The indices $i_x$ and $j_x$ (from 0 to $H-1$) specify the **block-row** and **block-column**.
- The indices $i_y$ and $j_y$ (from 0 to $W-1$) specify the **intra-block row** and **intra-block column**.

Thus, the element $B_{(i_x,i_y),(j_x,j_y)}$ is located in the row $i_y$ and column $j_y$ of the matrix block in the block row $i_x$ and block column $j_x$.

Since the matrix $B$ has a block circulant structure, the block in block row $i_x$ and block column $j_x$ is given by $C_{(j_x-i_x) \bmod H}$.

We can see that

$$\begin{aligned} B_{(i_x,i_y),(j_x,j_y)} &= (C_{(j_x-i_x) \bmod H})_{i_y,j_y} \\ &= (C_{(j_x-i_x) \bmod H})_{(j_y-i_y) \bmod W} \quad (22) \\ &= b_{((j_x-i_x) \bmod H,(j_y-i_y) \bmod W)} \end{aligned}$$

We can express the matrix multiplication $Bx$, where $B \in \mathbb{R}^{N \times N}$ and $x \in \mathbb{R}^N$, as follows:

$$\begin{aligned} (Bx)_{(i_x,i_y)} &= \sum_{j=0}^{N-1} B_{(i_x,i_y),(j_x,j_y)} x_{(j_x,j_y)} \\ &= \sum_{j_x=0}^{H-1} \sum_{j_y=0}^{W-1} B_{(i_x,i_y),(j_x,j_y)} x_{(j_x,j_y)} \\ &= \sum_{j_x=0}^{H-1} \sum_{j_y=0}^{W-1} b_{((j_x-i_x) \bmod H,(j_y-i_y) \bmod W)} x_{(j_x,j_y)} \end{aligned}$$

$$(23)$$

Let $k_x = (j_x - i_x) \bmod H$ and $k_y = (j_y - i_y) \bmod W$. This implies $j_x = (i_x + k_x) \bmod H$ and $j_y = (i_y + k_y) \bmod W$. Substituting these into the summation yields:

$$y_{(i_x,i_y)} = \sum_{k_x=0}^{H-1} \sum_{k_y=0}^{W-1} b_{(k_x,k_y)} x_{((i_x+k_x) \bmod H,(i_y+k_y) \bmod W)}$$

$$(24)$$

This expression is the definition of the **2D circular cross-correlation** between the 2D kernel $b$ and the 2D signal $x$.

We leverage the **2D Cross-Correlation Theorem**, which is a direct extension of the 1D theorem. It states that the 2D Fourier transform of a cross-correlation is the element-wise product of the first signal's conjugated 2D Fourier transform and the second signal's 2D Fourier transform:

$$\mathcal{F}_{2D}(Bx) = \overline{\mathcal{F}_{2D}(b)} \odot \mathcal{F}_{2D}(x) \quad (25)$$

Thus, we have:

$$Bx = \mathcal{F}_{2D}^{-1}\left(\overline{\mathcal{F}_{2D}(b)} \odot \mathcal{F}_{2D}(x)\right) \quad (26)$$

where $\mathcal{F}_{2D}, \mathcal{F}_{2D}^{-1}$ denotes the 2D DFT and its inversion. For simplicity, the reshaping operations between 1D $\mathbb{R}^N$ sequences and 2D $\mathbb{R}^{H \times W}$ feature maps are not demonstrated, and we define the inputs/outputs of $\mathcal{F}_{2D}, \mathcal{F}_{2D}^{-1}$ to be 1D sequences. Therefore, we can also compute it using the fast Fourier transform with a time complexity of $O(N \log N)$.

## B  Model Architectures

We provide the architectures of CA-DeiT, CA-PVT, CA-Swin and CAT in Table 6, Table 7, Table 8, Table 9 and Table 10. Our circulant attention serves as a plug-in module. We simply replace the original attention module with the proposed circulant attention, while keeping the other network components unchanged. To introduce positional information, we employ conditional positional encodings (Chu et al. 2023). Since our method benefits from $N \log N$ complexity, it is possible to directly process a high-resolution feature map in the early stages without incurring high computational cost. Therefore, for hierarchical models, we mainly restrict the attention replacement to the first two stages.

| | CA-DeiT-T | | CA-DeiT-S | | CA-DeiT-B | |
|---|---|---|---|---|---|---|
| | CA Block | DeiT Block | CA Block | DeiT Block | CA Block | DeiT Block |
| | $\begin{bmatrix} \text{res } 14\times14 \\ \text{dim } 192 \\ \text{head } 192 \end{bmatrix} \times 12$ | None | $\begin{bmatrix} \text{res } 14\times14 \\ \text{dim } 384 \\ \text{head } 384 \end{bmatrix} \times 12$ | None | $\begin{bmatrix} \text{res } 14\times14 \\ \text{dim } 768 \\ \text{head } 768 \end{bmatrix} \times 12$ | None |

Table 6: Architectures of CA-DeiT models.

| Stage | Output | CA-PVT-T | | CA-PVT-S | |
|---|---|---|---|---|---|
| | | CA Block | PVT Block | CA Block | PVT Block |
| res1 | 56 × 56 | Conv4×4, stride=4, 64, LN | | | |
| | | $\begin{bmatrix} \text{win } 56\times56 \\ \text{dim } 64 \\ \text{head } 64 \end{bmatrix} \times 2$ | None | $\begin{bmatrix} \text{win } 56\times56 \\ \text{dim } 64 \\ \text{head } 64 \end{bmatrix} \times 3$ | None |
| res2 | 28 × 28 | Conv2×2, stride=2, 128, LN | | | |
| | | $\begin{bmatrix} \text{win } 28\times28 \\ \text{dim } 128 \\ \text{head } 128 \end{bmatrix} \times 2$ | None | $\begin{bmatrix} \text{win } 28\times28 \\ \text{dim } 128 \\ \text{head } 128 \end{bmatrix} \times 4$ | None |
| res3 | 14 × 14 | Conv2×2, stride=2, 320, LN | | | |
| | | None | $\begin{bmatrix} \text{win } 14\times14 \\ \text{dim } 320 \\ \text{head } 5 \end{bmatrix} \times 2$ | None | $\begin{bmatrix} \text{win } 14\times14 \\ \text{dim } 320 \\ \text{head } 5 \end{bmatrix} \times 6$ |
| res4 | 7 × 7 | Conv2×2, stride=2, 512, LN | | | |
| | | None | $\begin{bmatrix} \text{win } 7\times7 \\ \text{dim } 512 \\ \text{head } 8 \end{bmatrix} \times 2$ | None | $\begin{bmatrix} \text{win } 7\times7 \\ \text{dim } 512 \\ \text{head } 8 \end{bmatrix} \times 3$ |

Table 7: Architectures of CA-PVT models (Part1).

| Stage | Output | CA-PVT-M | | CA-PVT-L | |
|---|---|---|---|---|---|
| | | CA Block | PVT Block | CA Block | PVT Block |
| res1 | 56 × 56 | Conv4×4, stride=4, 64, LN | | | |
| | | [win 56×56, dim 64, head 64] ×3 | None | [win 56×56, dim 64, head 64] ×3 | None |
| res2 | 28 × 28 | Conv2×2, stride=2, 128, LN | | | |
| | | [win 28×28, dim 128, head 128] ×4 | None | [win 28×28, dim 128, head 128] ×8 | None |
| res3 | 14 × 14 | Conv2×2, stride=2, 320, LN | | | |
| | | None | [win 14×14, dim 320, head 5] ×18 | None | [win 14×14, dim 320, head 5] ×27 |
| res4 | 7 × 7 | Conv2×2, stride=2, 512, LN | | | |
| | | None | [win 7×7, dim 512, head 8] ×3 | None | [win 7×7, dim 512, head 8] ×3 |

Table 8: Architectures of CA-PVT models (Part2).

| Stage | Output | CA-Swin-T | | CA-Swin-S | | CA-Swin-B | |
|---|---|---|---|---|---|---|---|
| | | CA Block | Swin Block | CA Block | Swin Block | CA Block | Swin Block |
| res1 | 56 × 56 | concat 4 × 4, 96, LN | | concat 4 × 4, 96, LN | | concat 4 × 4, 128, LN | |
| | | [win 56×56, dim 96, head 96] ×2 | None | [win 56×56, dim 96, head 96] ×2 | None | [win 56×56, dim 128, head 128] ×2 | None |
| res2 | 28 × 28 | concat 2 × 2, 192, LN | | concat 2 × 2, 192, LN | | concat 2 × 2, 256, LN | |
| | | [win 28×28, dim 192, head 192] ×2 | None | [win 28×28, dim 192, head 192] ×2 | None | [win 28×28, dim 256, head 256] ×2 | None |
| res3 | 14 × 14 | concat 2 × 2, 384, LN | | concat 2 × 2, 384, LN | | concat 2 × 2, 512, LN | |
| | | None | [win 7×7, dim 384, head 12] ×6 | None | [win 7×7, dim 384, head 12] ×18 | None | [win 7×7, dim 512, head 16] ×18 |
| res4 | 7 × 7 | concat 2 × 2, 768, LN | | concat 2 × 2, 768, LN | | concat 2 × 2, 1024, LN | |
| | | None | [win 7×7, dim 768, head 24] ×2 | None | [win 7×7, dim 768, head 24] ×2 | None | [win 7×7, dim 1024, head 32] ×2 |

Table 9: Architectures of CA-Swin models.

| Stage | Output | CAT-T | | CAT-S | | CAT-B | |
|---|---|---|---|---|---|---|---|
| | | CA Block | Attn Block | CA Block | Attn Block | CA Block | Attn Block |
| res1 | $56 \times 56$ | Stem, stride 4, 64, BN | | Stem, stride 4, 64, BN | | Stem, stride 4, 96, BN | |
| | | $\begin{bmatrix} \text{win } 56\times56 \\ \text{dim } 64 \\ \text{head } 64 \end{bmatrix} \times 1$ | None | $\begin{bmatrix} \text{win } 56\times56 \\ \text{dim } 64 \\ \text{head } 64 \end{bmatrix} \times 2$ | None | $\begin{bmatrix} \text{win } 56\times56 \\ \text{dim } 96 \\ \text{head } 96 \end{bmatrix} \times 2$ | None |
| res2 | $28 \times 28$ | Conv 3×3, stride 2, 128, BN | | Conv 3×3, stride 2, 128, BN | | Conv 3×3, stride 2, 192, BN | |
| | | $\begin{bmatrix} \text{win } 28\times28 \\ \text{dim } 128 \\ \text{head } 128 \end{bmatrix} \times 2$ | None | $\begin{bmatrix} \text{win } 28\times28 \\ \text{dim } 128 \\ \text{head } 128 \end{bmatrix} \times 4$ | None | $\begin{bmatrix} \text{win } 28\times28 \\ \text{dim } 192 \\ \text{head } 192 \end{bmatrix} \times 4$ | None |
| res3 | $14 \times 14$ | Conv 3×3, stride 2, 320, BN | | Conv 3×3, stride 2, 320, BN | | Conv 3×3, stride 2, 448, BN | |
| | | None | $\begin{bmatrix} \text{win } 14\times14 \\ \text{dim } 320 \\ \text{head } 10 \end{bmatrix} \times 9$ | None | $\begin{bmatrix} \text{win } 14\times14 \\ \text{dim } 320 \\ \text{head } 10 \end{bmatrix} \times 18$ | None | $\begin{bmatrix} \text{win } 14\times14 \\ \text{dim } 448 \\ \text{head } 14 \end{bmatrix} \times 18$ |
| res4 | $7 \times 7$ | Conv 3×3, stride 2, 512, BN | | Conv 3×3, stride 2, 512, BN | | Conv 3×3, stride 2, 640, BN | |
| | | None | $\begin{bmatrix} \text{win } 7\times7 \\ \text{dim } 512 \\ \text{head } 16 \end{bmatrix} \times 4$ | None | $\begin{bmatrix} \text{win } 7\times7 \\ \text{dim } 512 \\ \text{head } 16 \end{bmatrix} \times 8$ | None | $\begin{bmatrix} \text{win } 7\times7 \\ \text{dim } 640 \\ \text{head } 20 \end{bmatrix} \times 8$ |

Table 10: Architectures of CAT models.