
VMamba: Visual State Space Model

Yue Liu
UCAS

liuyue171@mailsucas.ac.cn

Yunjie Tian
UCAS

tianyunjie19@mailsucas.ac.cn

Yuzhong Zhao
UCAS

zhaoyuzhong20@mailsucas.ac.cn

Hongtian Yu
UCAS

yuhongtian17@mailsucas.ac.cn

Lingxi Xie
Huawei Inc.
198808xc@gmail.com

Yaowei Wang
Pengcheng Lab.
wangyw@pcl.ac.cn

Qixiang Ye
UCAS
qxye@ucas.ac.cn

Yunfan Liu
UCAS
yunfan.liu@ucas.ac.cn

Abstract

Convolutional Neural Networks (CNNs) and Vision Transformers (ViTs) have long been the predominant backbone networks for visual representation learning. While ViTs have recently gained prominence over CNNs due to their superior fitting capabilities, their scalability is largely constrained by the quadratic complexity of attention computation. Inspired by the capability of Mamba in efficiently modeling long sequences, we propose VMamba, a generic vision backbone model aiming to reduce the computational complexity to linear while retaining ViTs' advantageous features. To enhance VMamba's adaptability in processing vision data, we introduce the Cross-Scan Module (CSM) to enable 1D selective scanning in 2D image space with global receptive fields. Additionally, we make further improvements in implementation details and architectural designs to enhance VMamba's performance and boost its inference speed. Extensive experimental results demonstrate VMamba's promising performance across various visual perception tasks, highlighting its pronounced advantages in input scaling efficiency compared to existing benchmark models. Source code is available at <https://github.com/MzeroMiko/VMamba>.

1 Introduction

Visual representation learning stands as one of the fundamental research topics in computer vision, which has witnessed significant breakthroughs in the era of deep learning. To represent the complex patterns in vision data, two primary categories of backbone networks, *i.e.*, Convolution Neural Networks (CNNs) [42, 22, 25, 33, 47] and Vision Transformers (ViTs) [12, 32, 50, 64], have been proposed and extensively employed in a variety of visual tasks. Compared to CNNs, ViTs generally exhibit superior modeling capabilities owing to the incorporation of the self-attention mechanism [52, 12] which realizes global receptive fields and dynamically predicted weighting parameters [20].

However, self-attention also brings a side effect, *i.e.*, the complexity of self-attention grows quadratically as the size of input increases, resulting in significant computational overhead in downstream tasks with large spatial resolutions. To address this issue, considerable efforts have been dedicated to enhancing the efficiency of attention, primarily by imposing constraints on the size or stride of computing windows [48, 32, 11]. Although effective, these approaches inevitably make trade-offs between effective receptive field and computational efficiency, thereby limiting the capability to establish long-distance dependencies in vision data. This motivates us to develop a novel architecture for efficient modeling of vision data, while preserving the advantages associated with the vanilla self-attention mechanism, namely, global receptive fields and dynamic weights.

Recently, State Space Models (SSMs) [14, 39, 53] have demonstrated great potential for long sequence modeling with linear complexity in natural language processing (NLP) tasks. Inspired by this, we introduce VMamba, a generic vision backbone featuring SSM-based blocks for efficient visual representation learning. VMamba’s effectiveness in reducing the complexity of attention computation is largely attributed to the selective scan mechanism present in S6 models, also referred to as selective SSMs [14]. Unlike conventional attention computation methods that permit dense information routing within the context, S6 mandates each element in a 1D array (*e.g.*, a text sequence) to acquire contextual knowledge only through a compressed hidden state, thereby reducing the quadratic complexity to linear.

However, due to the two-dimensional nature of vision data, it is challenging for a single scanning process to simultaneously capture the dependency information across different directions, resulting in restricted receptive fields. We refer to this issue as a ‘direction-sensitive’ problem and propose to address it through the newly introduced Cross-Scan Module (CSM). Instead of traversing the spatial domain of image feature maps in a unidirectional pattern (either column-wise or row-wise), CSM adopts a four-way scanning strategy, *i.e.*, from the upper-left and lower-right corners all across the feature map to the opposite location (see Figure 2 (b)). This strategy ensures that each element in a feature map integrates information from all other locations in different directions, thereby achieving a global receptive field without increasing the computational complexity. Aside from introducing CSM, we also conduct a systematic and comprehensive analysis of both the architectural design and implementation details of VMamba, leading to substantial enhancements in computational efficiency while retaining the modeling capacity.

We validate the effectiveness of VMamba through extensive experiments across various visual recognition tasks. As illustrated in Figure 1, the family of VMamba models exhibits superior or at least competitive performance on ImageNet-1K when compared with benchmark vision models, including ResNet [22], ViT [12], and Swin [32]. In downstream tasks with dense prediction, VMamba-Tiny/Small/Base achieves 47.4%/48.7%/49.2% mAP on COCO [30] with the MaskRCNN detector ($1\times$ training schedule) and 48.3%/50.6%/51.0% mIoU (Single-Scale) on ADE20K [66] using UperNet with 512×512 input, demonstrating its great potential to serve as a powerful backbone model. Moreover, unlike ViT-based models whose computational complexity grows quadratically with input size, VMamba exhibits linear growth in FLOPs while still maintaining performance comparable to ViT, showcasing its state-of-the-art input scalability.

The main contributions of this study are summarized as follows:

- We propose VMamba, a generic vision backbone network built on State Space Models (SSMs), featuring global receptive fields and dynamic weights for efficient visual representation learning. VMamba introduces a novel architecture for vision foundation models, expanding upon existing choices of CNNs and ViTs.
- The Cross-Scan Module (CSM) is involved to bridge the gap between 1D array scanning and 2D plain traversing, facilitating the extension of S6 to vision data without compromising the size of receptive fields. In-depth analysis and improvement are conducted on both the implementation details and architectural design of VMamba to enhance its performance and efficiency.
- Without bells and whistles, VMamba achieves promising performance across various visual tasks including image classification, object detection, and semantic segmentation. It also demonstrates remarkable scalability in terms of input size, exhibiting linear growth in computational complexity.

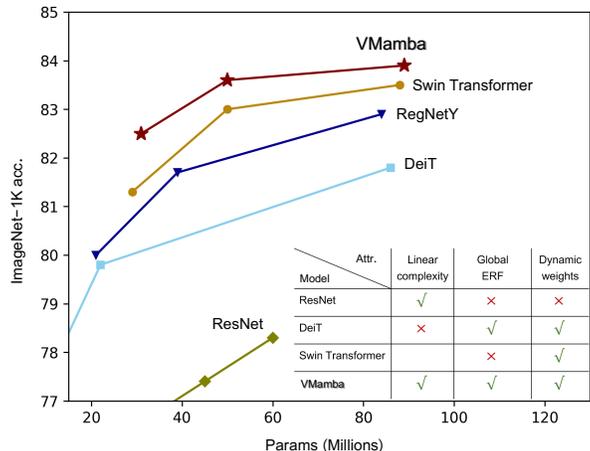


Figure 1: Performance comparison on ImageNet-1K. VMamba achieves superior top-1 classification accuracy compared to popular benchmark models. Remarkably, VMamba leverages the benefits of a global effective receptive field (ERF) and dynamic weights, all while maintaining linear computational complexity.

2 Related Work

Deep neural networks have substantially advanced the research in machine visual perception. There are primarily two prevalent types of visual foundation models, *i.e.*, CNNs [26, 42, 46, 22, 47] and ViTs [12, 32, 54, 11, 7, 64]. Recently, the success of State Space Models (SSMs) [14, 39, 53] has illustrated their efficacy in efficient long sequence modeling, which has attracted extensive attention in both the NLP and CV communities. Our study sticks to this line of work and proposes VMamba, a SSM-based architecture for data modeling in the vision domain. VMamba contributes as an alternative foundation model to the community, alongside CNNs and ViTs.

Convolution Neural Networks (CNNs) serve as the landmark models in the history of visual perception. Early CNN-based models [29, 26] are designed for basic tasks, such as recognizing handwritten digits [28] and classifying character categories [63]. The distinctive characteristics of CNNs are encapsulated in the convolution kernels, which employ receptive fields to capture visual information of interest from images. With the aid of powerful computing devices (*GPU*) and large-scale datasets [8], increasingly deeper [42, 46, 22, 25] and efficient models [23, 47, 59, 41] have been proposed to enhance performance across a spectrum of visual tasks. In addition to these efforts, progress has been made to propose more advanced convolution operators [5, 24, 61, 6] or more efficient network architectures [69, 4, 58, 23].

Vision Transformers (ViTs) are adapted from the NLP community, showcasing a potent perception model for visual tasks and swiftly evolving into one of the most promising visual foundation models. Early ViT-based models usually require large-scale datasets [12] and appear in a plain configuration [62, 67, 2, 35]. Later, DeiT [50] employs training techniques to address challenges encountered in the optimization process, and subsequent studies tend to incorporate inductive bias of visual perception into network design. For example, the community proposes hierarchical ViTs [32, 11, 54, 35, 64, 49, 7, 9, 65] to gradually decrease the feature resolution throughout the backbone. Moreover, other studies propose to harness the advantages of CNNs, such as introducing convolution operations [55, 7, 51], designing hybrid architectures by combining CNN and ViT modules [7, 45, 35], *etc.*

State Space Models (SSMs) are recently proposed models that are introduced into deep learning as state space transforming [18, 17, 43]. Inspired by continuous state space models in control systems, combined with HiPPO [15] initialization, LSSL [18] showcases the potential in handling long range dependency problems. However, due to the prohibitive computation and memory requirements induced by the state representation, LSSL is infeasible to use in practice. To solve this problem, S4 [17] proposes to normalize the parameter into diagonal structure. Since then, many flavors of

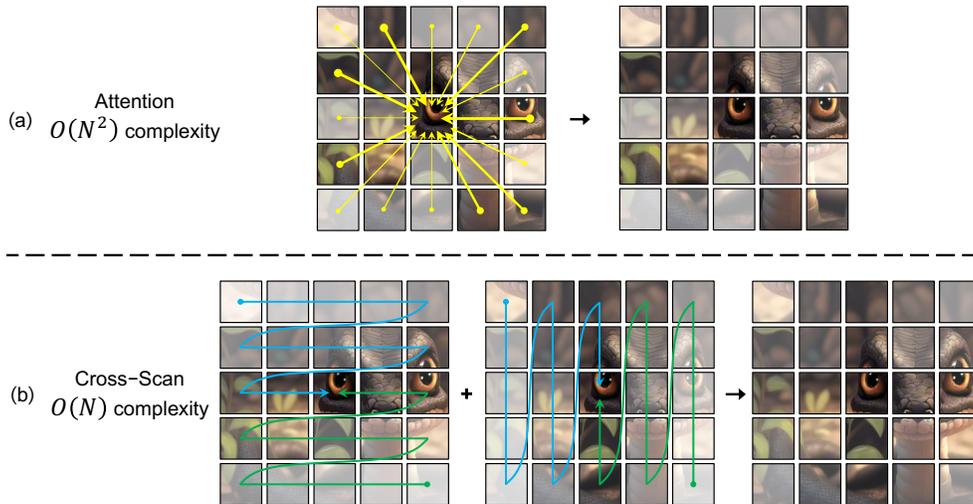


Figure 2: Comparison of information flow: Attention vs. Cross-Scan Module (CSM). (a) The attention mechanism calculates the correlation between the target image patch (the central one in the figure) and all other patches in the image., resulting in $\mathcal{O}(N^2)$ complexity; (b) The proposed CSM scans the image along four different paths with a computational complexity of $\mathcal{O}(N)$.

structured state space models sprang up with different structures like complex-diagonal structure [19, 16], multiple-input multiple output supporting [43], decomposition of diagonal plus low-rank operations [21], selection mechanism [14]. These models are then integrated into large representation models [39, 37, 13].

Those models are mainly focuses on the how state space models are applied on long-range and casual data like language and speech, such as language understanding [37, 39], content-based reasoning [14], pixel-level 1-D image classification [17], few of them pay attention in visual recognition. The most similar work to ours is S4ND [40]. S4ND is the first work applying state space mechanism into visual tasks and showing the potential that its performance may compete with ViT [12]. However, S4ND expands the S4 model in a simple manner, fails on efficiently capturing image information in an input-dependent manner. We demonstrates that with selective scan mechanism introduced by mamba [14], the proposed VMamba is able to match existing popular vision foundation models like ResNet [22], ViT [12], swin [31], and convnext [33], showcasing the potential of VMamba to be the powerful foundation model.

3 Method

In this section, we first introduce the preliminary concepts associated with State Space Models (SSMs), including their continuous and discretized formulations, along with an overview of the efficient computational approach of selective SSMs. We then provide detailed specifications of the proposed Cross-Scan Module (CSM), which is essential for implementing the selective scan operation within the two-dimensional image domain. Subsequently, we conduct theoretical analysis to elucidate the internal operational mechanism of CSM and its correlation with self-attention and gate linear attention, which are supported by the accompanying visualization results.

3.1 Preliminaries

State Space Model (SSM). SSMs can be regarded as linear time-invariant (LTI) systems that map the input stimulation $u(t) \in \mathbb{R}^L$ to the output response $y(t) \in \mathbb{R}^L$ through a hidden state $h(t) \in \mathbb{C}^N$. They are commonly formulated as linear ordinary differential equations (ODEs) as follows,

$$\begin{aligned} h'(t) &= \mathbf{A}h(t) + \mathbf{B}u(t) \\ y(t) &= \mathbf{C}h(t) + \mathbf{D}u(t) \end{aligned} \quad (1)$$

where $\mathbf{A} \in \mathbb{C}^{N \times N}$, $\mathbf{B}, \mathbf{C} \in \mathbb{C}^N$, and $D \in \mathbb{C}^1$ are the weighting parameters.

Discretization of SSM. To be integrated into deep models, continuous-time SSMs need to be discretized beforehand, which can be achieved by solving the ODE followed by a straightforward discretization process. Concretely, the analytic solution to Eq. 1 can be expressed as

$$h(t_b) = e^{\mathbf{A}(t_b - t_a)} h(t_a) + e^{\mathbf{A}(t_b - t_a)} \int_{t_a}^{t_b} \mathbf{B}(\tau) u(\tau) e^{-\mathbf{A}(\tau - t_a)} d\tau \quad (2)$$

Then, by sampling with step size Δ (i.e., $d\tau|_{t_i}^{t_{i+1}} = \Delta_i$), $h(t_b)$ can be discretized by

$$h_b = e^{\mathbf{A}(\Delta_a + \dots + \Delta_{b-1})} \left(h_a + \sum_{i=a}^{b-1} \mathbf{B}_i u_i e^{-\mathbf{A}(\Delta_a + \dots + \Delta_i)} \Delta_i \right) \quad (3)$$

Notably, this discretization formulation is approximately equivalent to the result obtained by the zero-order hold (ZOH) method, which is frequently utilized in the literature related to SSM. To be concrete, let $b = a + 1$, then Eq. 3 can be written as

$$h_{a+1} = e^{\mathbf{A}\Delta_a} (h_a + \mathbf{B}_a u_a e^{-\mathbf{A}\Delta_a} \Delta_a) \quad (4)$$

$$= e^{\mathbf{A}\Delta_a} h_a + \mathbf{B}_a \Delta_a u_a \quad (5)$$

$$= \overline{\mathbf{A}}_a h_a + \overline{\mathbf{B}}_a u_a \quad (6)$$

where $\overline{\mathbf{A}}_a = e^{\mathbf{A}\Delta_a}$ aligns with the discretization result of ZOH [14], while $\overline{\mathbf{B}}_a = \mathbf{B}_a \Delta_a$ is basically the first-order Taylor expansion of the counterpart acquired through ZOH.

Selective Scan. Notably, the weight matrix \mathbf{B} in Eq. 2 and 3 is configured to be input-dependent to address the inadequacy of LTI SSMs (Eq. 1) in capturing contextual information [14] (same for \mathbf{C} , \mathbf{D} , and Δ). However, the resulting time-varying SSMs pose a challenge to efficient computation, as convolutions do not support dynamic weights and are therefore no longer applicable. Nevertheless, if we can derive the recurrence relation of h_b in Eq. 3, it can still be efficiently computed. Concretely, let us denote $e^{\mathbf{A}(\Delta_a + \dots + \Delta_{i-1})}$ as $\mathbf{p}_{\mathbf{A},a}^i$, then its recurrence relation can be directly written as

$$\mathbf{p}_{\mathbf{A},a}^i = e^{\mathbf{A}\Delta_{i-1}} \mathbf{p}_{\mathbf{A},a}^{i-1} \quad (7)$$

For the second term of Eq. 3, we have

$$\mathbf{p}_{\mathbf{B},a}^b = e^{\mathbf{A}(\Delta_a + \dots + \Delta_{b-1})} \sum_{i=a}^{b-1} \mathbf{B}_i u_i e^{-\mathbf{A}(\Delta_a + \dots + \Delta_i)} \Delta_i \quad (8)$$

$$= e^{\mathbf{A}\Delta_{b-1}} \mathbf{p}_{\mathbf{B},a}^{b-1} + \mathbf{B}_{b-1} u_{b-1} \Delta_{b-1} \quad (9)$$

Therefore, with the associations derived in Eq. 7 and Eq. 8, $h_b = \mathbf{p}_{\mathbf{A},a}^b h_a + \mathbf{p}_{\mathbf{B},a}^b$ can be efficiently computed in parallel using associative scan algorithms [1, 38, 44], which are supported by numerous modern programming libraries. This approach effectively reduces the overall computational complexity to linear, and VMamba further accelerates the computation by adopting a hardware-aware implementation [14].

3.2 2D-Selective-Scan for Vision Data

The engaging characteristic of the S6 model, namely, simultaneously attaining global receptive fields, dynamic weights (i.e., selectivity), and linear computational complexity, inspires us to incorporate it for learning visual representations. While the sequential nature of S6 aligns well with NLP tasks involving temporal data, it poses significant challenges when applied to vision data, which is inherently non-sequential and contains spatial information (e.g. local texture and global structure). To address this issue, S4ND [40] reformulates SSM with convolutional operations and straightforwardly extending the kernel from 1D to 2D through outer-product. However, such modification prevents the weights from being dynamic (i.e., input independent), resulting in a limited capacity for context-aware data modeling. Therefore, we opt to adhere to the selective scan approach [14] for data processing, and propose the 2D-Selective-Scan (SS2D) module for adapting S6 to vision data without compromising its advantages.

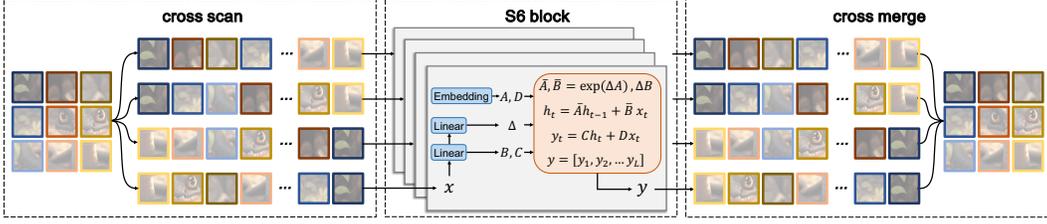


Figure 3: **Illustration of the 2D-Selective-Scan (SS2D) operation.** The input image is first divided into patches and then flattened along four different scanning paths (*Cross-Scan*). The resulting sequences of image patches are then separately processed by distinct S6 blocks. Subsequently, the outputs are merged to construct the 2D feature map as the final output (*Cross-Merge*).

As depicted in Figure 3, passing data through the SS2D module involves three steps: cross-scan, selective scanning with S6 blocks, and cross-merge. We collectively refer to cross-scan and cross-merge as the Cross Scan Module (CSM). Given the input data, SS2D first unfolds image patches into sequences along four distinct traversal paths (*i.e.*, cross-scan), processes each patch sequence using a separate S6 block in parallel, and subsequently reshapes and merges the resultant sequences to form the output map (*i.e.*, cross-merge). Through the adoption of complementary traversal paths, SS2D enables each pixel in the image to efficiently integrate information from all other pixels in different directions, thereby facilitating the establishment of global receptive fields.

3.3 The Relationship between SS2D and Self-attention

In this section, we clarify the relationship between SS2D and the self-attention mechanism commonly employed in existing backbone models. Subsequently, visualization results are provided to substantiate our explanation.

Notation Definition. Let T denote the length of the sequence with indices from a to b , we define the following variables

$$\mathbf{V} := [\mathbf{V}_1; \dots; \mathbf{V}_T] \in \mathbb{R}^{T \times D_v}, \text{ where } \mathbf{V}_i := \mathbf{u}_{a+i-1} \Delta_{a+i-1} \in \mathbb{R}^{1 \times D_v} \quad (10)$$

$$\mathbf{K} := [\mathbf{K}_1; \dots; \mathbf{K}_T] \in \mathbb{R}^{T \times D_k}, \text{ where } \mathbf{K}_i := \mathbf{B}_{a+i-1} \in \mathbb{R}^{1 \times D_k} \quad (11)$$

$$\mathbf{Q} := [\mathbf{Q}_1; \dots; \mathbf{Q}_T] \in \mathbb{R}^{T \times D_k}, \text{ where } \mathbf{Q}_i := \mathbf{C}_{a+i-1} \in \mathbb{R}^{1 \times D_k} \quad (12)$$

$$\mathbf{w} := [\mathbf{w}_1; \dots; \mathbf{w}_T] \in \mathbb{R}^{T \times D_k \times D_v}, \text{ where } \mathbf{w}_i := \prod_{j=1}^i e^{A \Delta_{a-1+j}} \in \mathbb{R}^{D_k \times D_v} \quad (13)$$

$$\mathbf{H} := [\mathbf{h}_a; \dots; \mathbf{h}_b] \in \mathbb{R}^{T \times D_k \times D_v}, \text{ where } \mathbf{h}_i \in \mathbb{R}^{D_k \times D_v} \quad (14)$$

$$\mathbf{Y} := [\mathbf{y}_a; \dots; \mathbf{y}_b] \in \mathbb{R}^{T \times D_v}, \text{ where } \mathbf{y}_i \in \mathbb{R}^{D_v} \quad (15)$$

Mathematical Derivation. Based on these notations, the discretized solution to time-varying SSMs (Eq. 3) can be written as

$$\mathbf{h}_b = \mathbf{w}_T \odot \mathbf{h}_a + \sum_{i=1}^T \frac{\mathbf{w}_T}{\mathbf{w}_i} \odot \left(\mathbf{K}_i^\top \mathbf{V}_i \right) \quad (16)$$

where \odot denotes the element-wise product between matrices, and the division is also elements-wise. Therefore, the first term of the output of SSM can be computed by

$$\mathbf{y}_b = \mathbf{Q}_T \mathbf{h}_b \quad (17)$$

$$= \mathbf{Q}_T (\mathbf{w}_T \odot \mathbf{h}_a) + \mathbf{Q}_T \sum_{i=1}^T \frac{\mathbf{w}_T}{\mathbf{w}_i} \odot \left(\mathbf{K}_i^\top \mathbf{V}_i \right) \quad (18)$$

Therefore, the j -th slice along dimension D_v of \mathbf{Y} , denoted as $\mathbf{X} := \mathbf{Y}^{(j)} \in \mathbb{R}^{T \times 1}$, can be expressed as

$$\mathbf{X} = (\mathbf{Q}_T \odot \mathbf{w}^{(j)}) \mathbf{h}_a^{(j)} + \left[(\mathbf{Q}_T \odot \mathbf{w}^{(j)}) \left(\frac{\mathbf{K}}{\mathbf{w}^{(j)}} \right)^\top \odot \mathbf{M} \right] \mathbf{V}^{(j)} \quad (19)$$

where $M := \text{tril}(T, T) \in \{0, 1\}^{T \times T}$ denotes the temporal mask matrix with the lower triangular portion of a $T \times T$ matrix set to 1 and elsewhere 0. It is evident that the manner in which matrices Q , K , and V are multiplied in Eq. 19 closely resembles the process in the self-attention module of Vision Transformers. Moreover, if w is in shape (T, D_k) rather than (T, D_k, D_v) , then Eq. 19 reduces to the form of Gated Linear Attention (GLA) [60], which indicates that GLA is also a special case of Mamba.

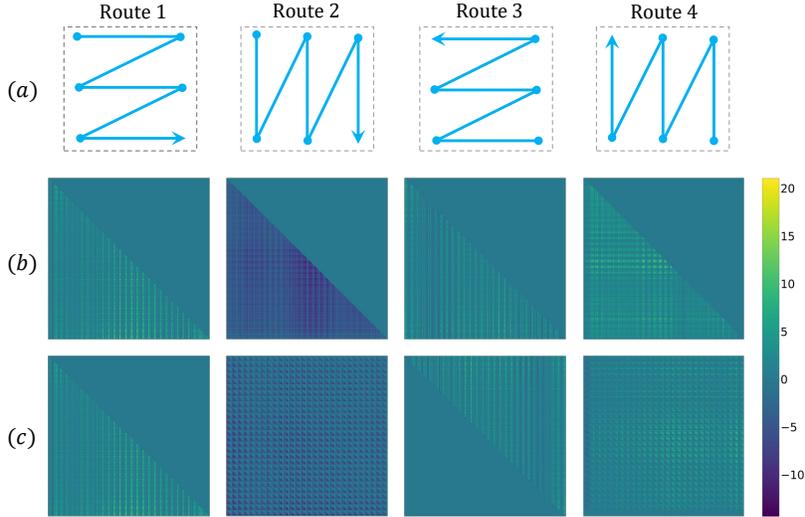


Figure 4: The visualization of activation maps associated with the four scanning routes in CSM. (a) illustration of the traversing path for each scanning route; (b) the activation map of each route; (c) the transformed activation map of each route.

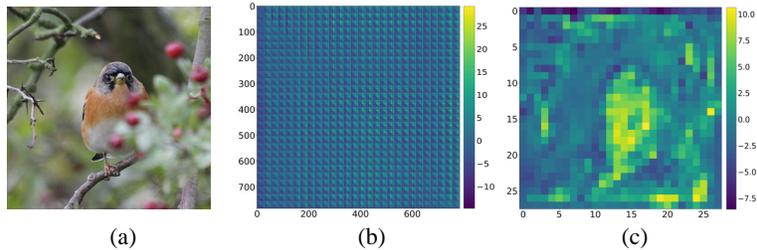


Figure 5: The visualization of the (a) input image, (b) overall activation map, and (c) heat map of CSM.

Activation Map Visualization. We have demonstrated that the computational process of selective SSMS encompasses mechanisms resembling self-attention, and this enables us to probe the internal operations of SS2D by visualizing its weight matrices.

Given the input image as depicted in Figure 5 (a), the illustration of the four scanning routes in CSM and the visualization of corresponding activation maps (that are obtained by calculating QK^T like the attention mechanism) are shown in Figure 4 (a) and (b), respectively. These results underscore the effectiveness of CSM in capturing and retaining the traversed information, as each row in a single attention map corresponds to all previously scanned foreground tokens impartially. Additionally, in Figure 4 (c), we showcase the transformed activation maps, where the pixel order corresponds to that of the first route, traversing the image row-wise from the upper-left to the bottom-right. Combining these individual maps produces the complete activation map of CSM (see Figure 5 (b)), which collectively enables VMamba to acquire the global receptive field. Moreover, by rearranging the diagonal elements of the obtained activation map in the image space, we derive the heat map illustrated in Figure 5 (c). This map illustrates VMamba’s effectiveness in accurately distinguishing between foreground and background pixels within an image.

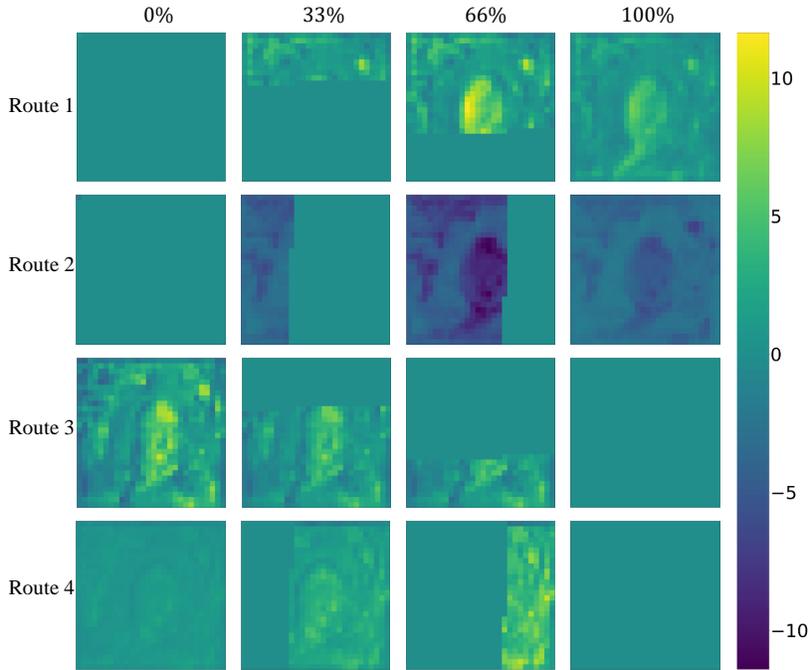


Figure 6: Illustration of the heat map for different scanning routes, sampled at four equally spaced progression stages.

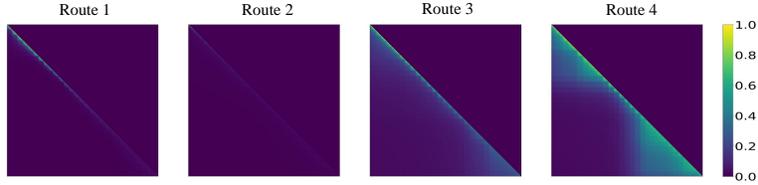


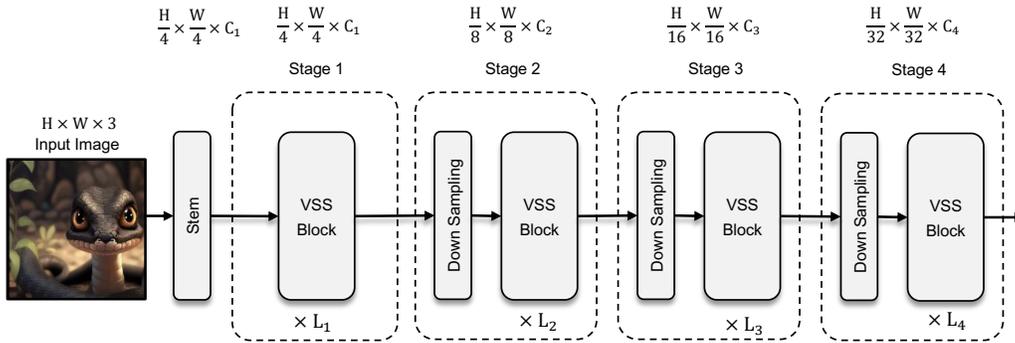
Figure 7: Illustration of the w matrix of different scanning routes.

To delve deeper into the internal mechanism of the scanning process of CSM, we visualize the heat map of a single scanning trajectory (row-wise from the upper-left to the lower-right) at four equally spaced progression stages from the beginning to the end. As shown in Figure 6, VMamba incrementally retains the cumulative scanning history, thus effectively establishing long-term dependencies across image patches. This behavior aligns with the characteristics of HiPPO [15], a prior research serving as the theoretical basis for SSM-based models.

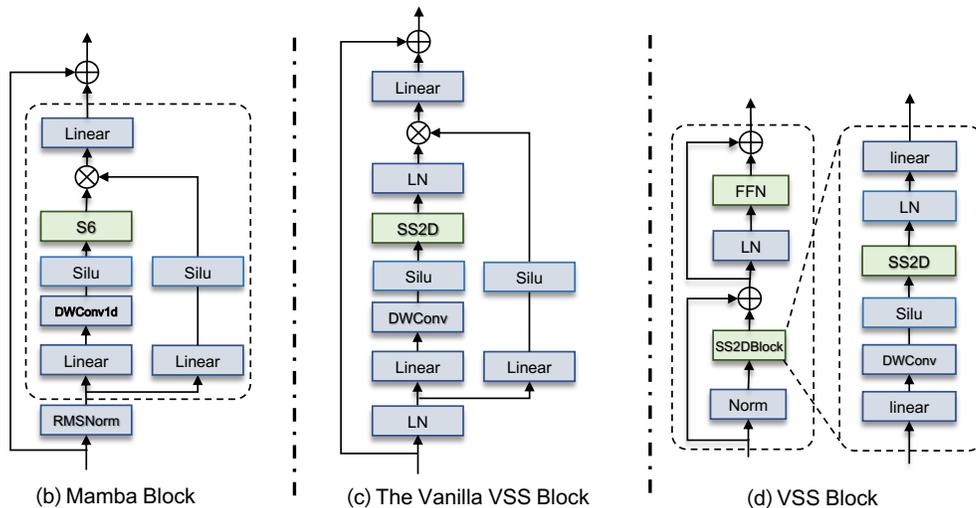
We further visualize the w matrix (Eq. 13) of four scanning paths in Figure 7. It is evident that each token selectively responds to a small subset of previously scanned tokens rather than uniformly responding to all others, thus limiting the memorization capacity of the hidden state in SSMs. Notably, this is distinct from the activation maps derived from \mathbf{BC}^\top (or \mathbf{QK}^\top , as shown in Figure 4 (b) and (c)), which retain all corresponding details prior to the scanning process. The property of w directs the hidden state to focus more on states from recent time steps, a behavior consistent with the criss-cross pattern of ERF illustrated in Figure 10.

4 The VMamba Model Family

In this section, we first introduce the overall architecture of VMamba, followed by presenting the VSS block, which serves as the core module in building VMamba. Subsequently, we enhance the speed of VMamba by refining the implementation details and architectural configurations.



(a) Architecture of VMamba



(b) Mamba Block

(c) The Vanilla VSS Block

(d) VSS Block

Figure 8: Illustration of network architectures.

4.1 Overall Architecture

An overview of the architecture of VMamba-Tiny is illustrated in Figure 8 (a). VMamba starts with partitioning the input image into patches using a stem module, resulting in a 2D feature map with the spatial dimension of $\frac{H}{4} \times \frac{W}{4}$.

Subsequently, multiple network stages, each consisting of VSS blocks preceded by down-sampling layers (except for the first stage), are utilized to create hierarchical representations with resolutions of $\frac{H}{8} \times \frac{W}{8}$, $\frac{H}{16} \times \frac{W}{16}$, and $\frac{H}{32} \times \frac{W}{32}$. The down-sampling operation is performed via patch merging [31], and the detailed structure of the VSS block is shown in Figure 8 (b). Detail discuss on variations of the VSS block and the relationship with Mamba will be presented in the following subsection.

We develop VMamba in three scales, *i.e.*, VMamba-Tiny, VMamba-Small, and VMamba-Base (referred to as VMamba-T, VMamba-S, and VMamba-B, respectively). FLOPs for all models are estimated at a resolution of 224×224 . Further architectures, including a large-scale model, will be introduced in future updates.

4.2 The Vanilla VSS Block

The structure of the vanilla VSS block is presented in Figure 8 (c), resembling the overall architecture of the Mamba block (shown in Figure 8 (b)) with some variations. Specifically, both of these two blocks can be seen as a residual network with a skip connection. The residual network comprises

two branches: one for feature extraction using a 3×3 depth-wise convolution layer, and the other consisting of a linear mapping followed by an activation layer, which computes the multiplicative gating signal. The primary difference between the Mamba and vanilla VSS block is the replacement of the S6 module with the SS2D module, which enables the adaptation of selective scanning to 2D vision data.

Notably, although position embedding bias is a common practice in vision transformers-based models, we refrain from utilizing it due to the sequential nature of selective scanning in SS2D. Moreover, the architectural design of the vanilla VSS block also diverges from that of vision transformer blocks, which typically follow the Norm \rightarrow attention \rightarrow Norm \rightarrow MLP routine (no MLP layer in the vanilla VSS block). Consequently, this version of the VSS block is shallower than the ViT block, enabling us to stack more blocks within a similar budget of total model depth.

4.3 Accelerating VMamba

Despite their efficiency in long sequence modeling, SSM-based architectures [14] often encounter reduced computational speed when dealing with inputs of smaller scales, which could limit the practical utility of VMamba. In this subsection, we present our efforts to enhance the throughput of VMamba, boosting it from 400 images/s to 1,336 images/s at an input resolution of 224×224 , without significant sacrifice in performance.

As reported in Figure 9, the vanilla VMamba-Tiny model proposed in Section 4 (**V0**¹) achieves a throughput of 426 images/s and comprises 22.9M parameters with 5.6G FLOPs (the FLOPs will drop to 4.5G if the selective scan operation could be implemented by a single for-loop). The low throughput and high memory overhead pose challenges for the practical deployment of VMamba. Therefore, to boost its inference speed, significant efforts have been made, primarily focusing on advancements in both implementation details and architectural design.

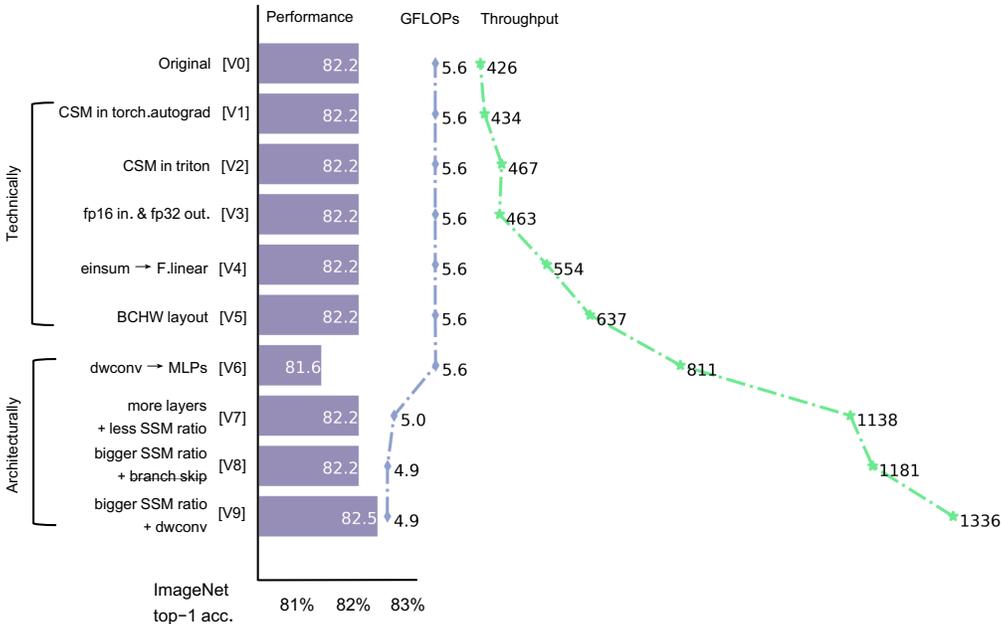


Figure 9: Illustration of the process of VMamba updates, measured in FLOPs, throughput, and classification accuracy.

Improvements in Implementation Details. From VMamba **V0** to **V2**, we successively implemented CSM in `torch.autograd.Function`, then re-implemented it in `Triton`. These modifications help

¹We note that the V0 version of VMamba maintains a similar design to Mamba [14], where the input data is first processed using a convolution layer followed by the selective scan operation.

increase the throughput from 426 to 467. Then, in **V3**, we adjusted the CUDA implementation associated with the selective scan operation to accommodate float16 input tensors and generate output tensors with float32 datatype. This adjustment enhances performance, particularly during training, compared to the implementation that processes tensors of float32 datatype, while also achieving improved numerical stability compared to using float16 for both input and output tensors. Additionally, in **V4** and **V5**, we substituted the relatively slow `einsum` operation in the selective scan with a linear transformation (*i.e.*, `torch.nn.functional.linear`). We also adopted the tensor layout of (B, C, H, W) to eliminate unnecessary data permutations. These changes resulted in a 49.5% increase in throughput (from 426 to 637) without affecting other metrics such as the number of parameters, FLOPs, and the classification performance on ImageNet-1K.

Improvements in Architectural Design. By balancing efficiency and effectiveness, we further enhanced the speed of VMamba based on the aforementioned improvements by modifying the architecture. Specifically, in **V6**, we introduced MLP into VMamba due to its computational efficiency, discarded the depth-wise convolution operation, and changed the layer configuration of VMamba-Tiny from [2, 2, 9, 2] to [2, 2, 2, 2]. We additionally adjusted the configuration of convolutional layers by replacing the downsampling kernel from $\text{conv}2 \times 2$ to $\text{conv}3 \times 3$, and substituted $\text{conv}4 \times 4$ in the stem module with two consecutive $\text{conv}3 \times 3$ kernels. This results in a significant increase in inference throughput by 27.3% (637 to 811). However, there is a slight decrease in classification accuracy from 82.2% to 81.6% and an increase in the number of parameters from 23M to 29M.

From **V6** to **V8**, we adjusted the dimension expansion parameter (`SSM_ratio`) of SSM from 2.0 to 1.0, and discarded the entire multiplicative branch in the vanilla VSS Block as the gating mechanism has already been implemented by the selectivity of SS2D. These modifications significantly lowered the FLOPs of VMamba (from 5.6G to 4.9G), allowing us to increase the depth of the model by changing the layer configuration from [2, 2, 2, 2] to [2, 2, 5, 2]. This enlargement of the network scale has resulted in an enhanced data modeling capacity, as evidenced by the notable improvement in classification accuracy from 81.6% to 82.2%. Finally, in **V9**, we re-introduced the depth-wise convolutional layers, and reduced the `dstate` parameter from 16 to 1 to conserve parameters for larger `SSM_ratio`. Compared to **V5**, these enhancements in the architectural design of VMamba result in an inference throughput improvement of 109.7% (from 637 to 1336). Additionally, they improve the classification performance on ImageNet-1K from 82.2% to 82.5% and keeps FLOPs as 4.9G.

5 Experiments

In this section, we conduct a series of experiments to evaluate and compare VMamba with popular benchmark models, including architectures based on both CNNs and Vision Transformers. Our assessment spans diverse visual tasks, including image classification, object detection, and semantic segmentation. Following this evaluation, we delve into a comprehensive analysis of VMamba’s characteristics, with a particular focus on its primary advantage: the remarkable capability to efficiently scale to increasingly large input resolutions.

5.1 Image Classification on ImageNet-1K

Settings. We assess VMamba’s performance in image classification using the ImageNet-1K dataset [8]. Following the evaluation protocol outlined in [31], VMamba-T/S/B models are trained from scratch for 300 epochs, with the first 20 epochs dedicated to warm-up, using a batch size of 1024. The training process utilizes the AdamW optimizer [34] with betas set to (0.9, 0.999), a momentum of 0.9, a cosine decay learning rate scheduler, an initial learning rate of 1×10^{-3} , and a weight decay of 0.05. Additional techniques such as label smoothing (0.1) and exponential moving average (EMA) are also applied. No further training techniques are employed beyond these.

Results. The comparison results of VMamba against benchmark backbone models on ImageNet-1K are summarized in Table 1. It is evident that, with similar FLOPs, VMamba-T achieves a performance of 82.5%, surpassing RegNetY-4G by 2.5%, DeiT-S by 2.7%, and Swin-T by 1.2%. Remarkably, these performance advantages of VMamba persist consistently across both small and base-scale models. Specifically, VMamba-S attains a top-1 accuracy of 83.6%, outperforming RegNetY-8G by

Method	Image size	#Param.	FLOPs	Throughput	Train Throughput	ImageNet top-1 acc.
RegNetY-4G [41]	224 ²	21M	4.0G	–	–	80.0
RegNetY-8G [41]	224 ²	39M	8.0G	–	–	81.7
RegNetY-16G [41]	224 ²	84M	16.0G	–	–	82.9
EffNet-B3 [47]	300 ²	12M	1.8G	–	–	81.6
EffNet-B4 [47]	380 ²	19M	4.2G	–	–	82.9
EffNet-B5 [47]	456 ²	30M	9.9G	–	–	83.6
EffNet-B6 [47]	528 ²	43M	19.0G	–	–	84.0
ViT-B/16 [12]	384 ²	86M	55.4G	–	–	77.9
ViT-L/16 [12]	384 ²	307M	190.7G	–	–	76.5
DeiT-S [50]	224 ²	22M	4.6G	1759	2397	79.8
DeiT-B [50]	224 ²	86M	17.5G	500	1024	81.8
DeiT-B [50]	384 ²	86M	55.4G	498	344	83.1
ConvNeXt-T [33]	224 ²	29M	4.5G	1189	701	82.1
ConvNeXt-S [33]	224 ²	50M	8.7G	682	444	83.1
ConvNeXt-B [33]	224 ²	89M	15.4G	435	334	83.8
HiViT-T [64]	224 ²	19M	4.6G	1391	1300	82.1
HiViT-S [64]	224 ²	38M	9.1G	711	697	83.5
HiViT-B [64]	224 ²	66M	15.9G	456	541	83.8
Swin-T [32]	224 ²	28M	4.6G	1247	985	81.3
Swin-S [32]	224 ²	50M	8.7G	719	640	83.0
Swin-B [32]	224 ²	88M	15.4G	457	494	83.5
S4ND-ConvNeXt-T [40]	224 ²	30M	–	684	331	82.2
S4ND-ViT-B [40]	224 ²	89M	–	404	340	80.4
ViM-S [68]	224 ²	26M	–	811	232 [†]	80.5
VMamba-T	224 ²	31M	4.9G	1335	464	82.5
VMamba-S	224 ²	50M	8.7G	874	313	83.6
VMamba-B	224 ²	89M	15.4G	645	246	83.9

Table 1: **Performance comparison on ImageNet-1K.** Throughput values are measured with an A100 GPU and an AMD EPYC 7542 CPU, using the toolkit released by [56], following the protocol proposed in [32]. † denotes that the training process only supports float32 datatype.

1.9% and Swin-S by 0.6%. Meanwhile, VMamba-B achieves a top-1 accuracy of 83.9%, surpassing RegNetY-16G by 1.0% and DeiT-B by 0.6%.

In terms of computational efficiency, while existing SSM-based vision models typically exhibit significantly better throughput only with large-scale inputs [68] such as 1024×1024 , VMamba-T achieves a throughput of 1,335 images/s even at an input resolution of 224×224 . This performance is either better or at least comparable to state-of-the-art methods, and this advantage persists across VMamba-S and VMamba-B. Notably, as the input size scales from 224×224 to 1024×1024 , the advantage of VMamba over existing methods becomes even more pronounced, as illustrated in Table 4. Further discussions on this topic will be provided in the subsequent sections.

5.2 Object Detection on COCO

Settings. In this section, we assess the performance of VMamba on object detection using the MSCOCO 2017 dataset [30]. Our training framework is constructed using the MMDetection library [3], and we adhere to the hyper-parameters used in Swin [31] with the Mask-RCNN detector. Specifically, we employ the AdamW optimizer [34] and fine-tune the pre-trained classification models (on ImageNet-1K) for both 12 and 36 epochs. The drop path rates are set to 0.2%/0.3%/0.5% for VMamba-T/S/B, respectively. The learning rate is initialized with 1×10^{-4} and is reduced by a factor

Mask R-CNN 1× schedule								
Backbone	AP ^b	AP ₅₀ ^b	AP ₇₅ ^b	AP ^m	AP ₅₀ ^m	AP ₇₅ ^m	#param.	FLOPs
ResNet-50	38.2	58.8	41.4	34.7	55.7	37.2	44M	260G
Swin-T	42.7	65.2	46.8	39.3	62.2	42.2	48M	267G
ConvNeXt-T	44.2	66.6	48.3	40.1	63.3	42.8	48M	262G
PVTv2-B2	45.3	67.1	49.6	41.2	64.2	44.4	45M	309G
ViT-Adapter-S	44.7	65.8	48.3	39.9	62.5	42.8	48M	403G
VMamba-T	47.4	69.5	52.0	42.7	66.3	46.0	50M	270G
ResNet-101	38.2	58.8	41.4	34.7	55.7	37.2	63M	336G
Swin-S	44.8	66.6	48.9	40.9	63.2	44.2	69M	354G
ConvNeXt-S	45.4	67.9	50.0	41.8	65.2	45.1	70M	348G
PVTv2-B3	47.0	68.1	51.7	42.5	65.7	45.7	65M	397G
VMamba-S	48.7	70.0	53.4	43.7	67.3	47.0	64M	357G
Swin-B	46.9	–	–	42.3	–	–	107M	496G
ConvNeXt-B	47.0	69.4	51.7	42.7	66.3	46.0	108M	486G
PVTv2-B5	47.4	68.6	51.9	42.5	65.7	46.0	102M	557G
ViT-Adapter-B	47.0	68.2	51.4	41.8	65.1	44.9	102M	557G
VMamba-B	49.2	70.9	53.9	43.9	67.7	47.6	108M	485G
Mask R-CNN 3× MS schedule								
Swin-T	46.0	68.1	50.3	41.6	65.1	44.9	48M	267G
ConvNeXt-T	46.2	67.9	50.8	41.7	65.0	44.9	48M	262G
PVTv2-B2	47.8	69.7	52.6	43.1	66.8	46.7	45M	309G
ViT-Adapter-S	48.2	69.7	52.5	42.8	66.4	45.9	48M	403G
VMamba-T	48.9	70.6	53.6	43.7	67.7	46.8	50M	270G
Swin-S	48.2	69.8	52.8	43.2	67.0	46.1	69M	354G
ConvNeXt-S	47.9	70.0	52.7	42.9	66.9	46.2	70M	348G
PVTv2-B3	48.4	69.8	53.3	43.2	66.9	46.7	65M	397G
VMamba-S	49.9	70.9	54.7	44.2	68.2	47.7	70M	384G

Table 2: Results of object detection and instance segmentation on COCO dataset. FLOPs are calculated with input size 1280×800 . AP^b and AP^m denote box AP and mask AP, respectively. The notation ‘1×’ indicates models fine-tuned for 12 epochs, while ‘3×MS’ denotes the utilization of multi-scale training for 36 epochs.

of $10\times$ at the 9th and 11th epoch. We implement multi-scale training and random flip with a batch size of 16, which are consistent with established practices for object detection evaluations.

Results. The results on COCO are summarized in Table 2. VMamba maintains superiority in box/mask Average Precision (AP) on COCO, regardless of the training schedule employed (12 or 36 epochs). Specifically, with a 12-epoch fine-tuning schedule, VMamba-T/S/B models achieve object detection mAPs of 47.4%/48.7%/49.2%, surpassing Swin-T/S/B by 4.7%/3.9%/2.3% mAP, and ConvNeXt-T/S/B by 3.2%/3.3%/2.2% mAP, respectively. With the same configuration, VMamba-T/S/B achieves instance segmentation mIoUs of 42.7%/43.7%/43.9%, outperforming Swin-T/S/B by 3.4%/2.8%/1.6% mIoU, and ConvNeXt-T/S/B by 2.6%/1.9%/1.3% mIoU, respectively.

Furthermore, the advantages of VMamba persist under the 36-epoch fine-tuning schedule with multi-scale training, as indicated in Table 2. When compared to counterparts including Swin [32], ConvNeXt [33], PVTv2 [55], and ViT [12] (with Adapters), VMamba-T/S exhibit superior performance, achieving 48.9%/49.9% mAP on object detection and 43.7%/44.2% mIoU on instance segmentation, respectively. These results underscore the potential of VMamba to achieve promising performance in downstream tasks with dense prediction.

5.3 Semantic Segmentation on ADE20K

Settings. Following Swin [32], we construct a UperHead [57] on top of the pre-trained model. We employ the AdamW optimizer [34] and set the learning rate to 6×10^{-5} . The fine-tuning process

method	crop size	mIoU (SS)	mIoU (MS)	#param.	FLOPs
ResNet-50	512 ²	42.1	42.8	67M	953G
DeiT-S + MLN	512 ²	43.8	45.1	58M	1217G
Swin-T	512 ²	44.4	45.8	60M	945G
ConvNeXt-T	512 ²	46.0	46.7	60M	939G
VMamba-T	512 ²	48.3	48.6	62M	948G
ResNet-101	512 ²	42.9	44.0	85M	1030G
DeiT-B + MLN	512 ²	45.5	47.2	144M	2007G
Swin-S	512 ²	47.6	49.5	81M	1039G
ConvNeXt-S	512 ²	48.7	49.6	82M	1027G
VMamba-S	512 ²	50.6	51.2	82M	1039G
Swin-B	512 ²	48.1	49.7	121M	1188G
ConvNeXt-B	512 ²	49.1	49.9	122M	1170G
VMamba-B	512 ²	51.0	51.6	122M	1170G

Table 3: Results of semantic segmentation on ADE20K using UperNet [57]. FLOPs are calculated with input size of 512×2048 . ‘SS’ and ‘MS’ denote single-scale and multi-scale testing, respectively.

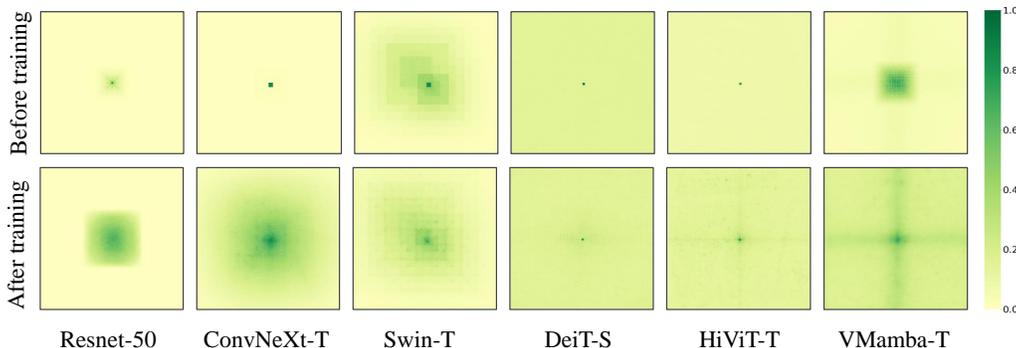


Figure 10: Visualization of the Effective Receptive Field (ERF) [36] for ResNet50 [22], ConvNeXt-T [33], Swin-T [32], DeiT-S [50] (ViT), and the proposed VMamba-T. Pixels of higher intensity indicate larger responses with the central pixel. Only DeiT [50], HiViT [64], and VMamba exhibit a global ERF.

spans a total of $160k$ iterations with a batch size of 16. The default input resolution is 512×512 , and we additionally present experimental results using 640×640 inputs and multi-scale (MS) testing.

Results. The results of semantic segmentation on ADE20K are summarized in Table 3. Consistent with the conclusions drawn from prior experiments, VMamba demonstrates superior accuracy, as VMamba-T achieves a 48.3% mIoU with a resolution of 512×512 and 48.6% mIoU with multi-scale (MS) input. These results outperform those of all benchmark methods, including ResNet [22], DeiT [50], Swin [32], and ConvNeXt [33]. Notably, the advantage persists in the experiments involving VMamba-S/B models.

5.4 Analytical Experiments

Visualization of the Effective Receptive Field. The Effective Receptive Field (ERF) [36] of an output unit denotes the region of input that contains elements with a non-negligible influence on that unit. We conduct a comparative analysis on the the central pixel’s ERFs [36, 10] across various prominent visual foundation models, including ResNet50 [22], ConvNeXt-T [33], Swin-T [32], and DeiT-S [50] (ViT), before and after training. According to the results shown in Figure 10, we make the following observations:

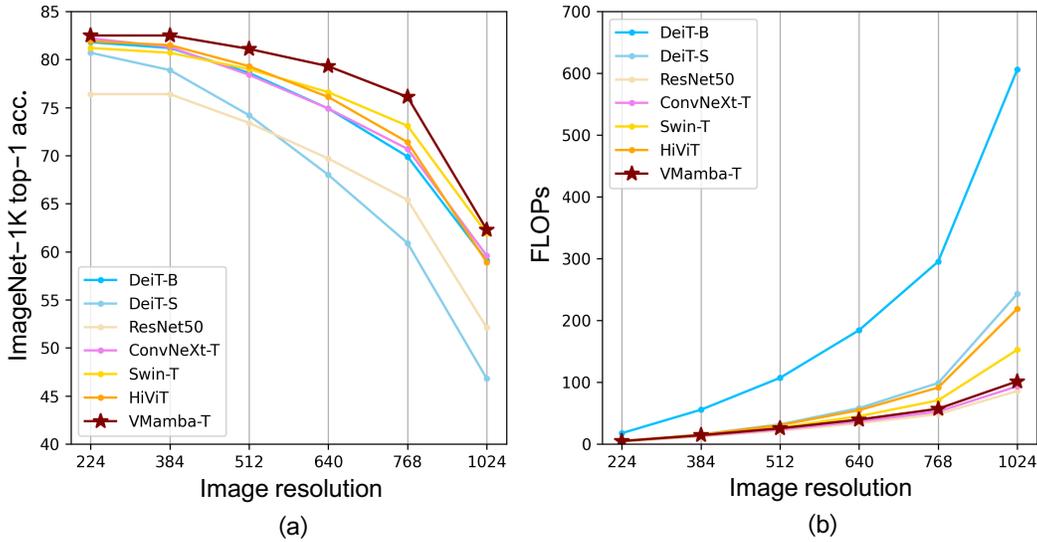


Figure 11: Illustration of the change in (a) classification accuracy and (b) FLOPs with progressively larger test image resolutions.

1. Among the models examined, only ViT-based methods, including DeiT, HiViT, and VMamba, showcase global ERFs, while the others exhibit local ERFs, despite their theoretical potential for global coverage. It is worth noting that DeiT (ViT) and HiViT incur quadratic complexity cost (please refer to Figure 11).
2. Unlike DeiT (ViT), which uniformly activates all pixels through the attention mechanism, VMamba activates all pixels while notably emphasizing cross-shaped activations. The scanning mechanism of the Cross-Scan Module ensures that the central pixel is primarily influenced by pixels along the cross, prioritizing long-range dependency context over local information for each pixel.
3. Notably, during the training process, VMamba’s ERF coverage expands from local to global, demonstrating significant adaptability and generalizability. We believe this adaptive process contributes to the model’s enhanced perception of images. This is in contrast to DeiT, which maintains nearly identical ERFs before and after training.

Computational Efficiency with Increasing Resolutions. Noting that Mamba is characterized by its remarkable capability in long sequence modeling, we conduct experiments to check if VMamba inherits this feature by assessing its computational efficiency with increasingly larger input images. Specifically, given models trained with input size 224×224 , we directly apply them to images with resolutions ranging from 384×384 to 1024×1024 , and measure their generalization performance in terms of the number of parameters, FLOPs, throughput in both training and inference time, as well as the top-1 classification accuracy on ImageNet-1K.

Based on the results listed in Table 4, VMamba demonstrates the most stable performance across different input image sizes, achieving a top-1 classification accuracy of 62.3% with the throughput of 66 when the input resolution reaches 1024×1024 . In comparison, with the same size of input, Swin [32] achieves the second-highest performance with a top-1 accuracy of 61.9%, but its throughput drastically drops to 22. Furthermore, while ResNet50 [22] maintains a relatively high inference speed (*i.e.*, a throughput of 181) with the largest input resolution, its classification accuracy drops to 52.1%, indicating its limited generalization capability. The performance change across various resolutions for different models is more intuitively illustrated in Figure 11. Notably, VMamba exhibits a linear growth in computational complexity, as measured by FLOPs, which is comparable to CNN-based architectures. This closely aligns with the theoretical conclusions of selective SSMs [14].

Diagnostic Study on Selective Scan Patterns. We perform diagnostic experiments to verify the effectiveness of the proposed scanning pattern (*i.e.*, CSM) for 2D-Selective-Scan. Specifically, three alternative image traversal methods are designed and used as benchmarks, with feature dimensions

Model	Image size	#Param.	FLOPs	Throughput	Train throughput	ImageNet top-1 acc.
ResNet50 [22]	224 ²	26M	4.1G	2988	1005	76.4
	384 ²	26M	12.1G	1192	518	76.5
	512 ²	26M	21.5G	705	305	73.4
	640 ²	26M	33.5G	435	200	69.7
	768 ²	26M	48.3G	309	142	65.3
	1024 ²	26M	85.9G	181	81	52.1
ConvNeXt [33]	224 ²	29M	4.5G	1107	613	82.0
	384 ²	29M	13.1G	405	240	81.0
	512 ²	29M	23.3G	225	140	78.0
	640 ²	29M	36.5G	147	90	74.3
	768 ²	29M	52.5G	103	63	69.5
	1024 ²	29M	93.3G	60	36	55.4
Deit [50]	224 ²	22M	4.6G	1559	1130	80.7
	384 ²	22M	15.5G	494	695	78.9
	512 ²	22M	31.8G	256	386	74.2
	640 ²	22M	58.2G	144	243	68.0
	768 ²	22M	98.7G	87	156	70.0
	1024 ²	22M	243.1G	36	73	46.9
Swin [32]	224 ²	28M	4.5G	1061	740	81.2
	384 ²	28M	14.5G	305	246	80.7
	512 ²	28M	26.6G	168	122	79.0
	640 ²	28M	45.0G	85	64	76.6
	768 ²	28M	70.7G	52	36	73.1
	1024 ²	28M	152.5G	22	13	61.9
HiViT [64]	224 ²	19M	4.6G	1259	948	81.9
	384 ²	19M	15.2G	392	332	81.5
	512 ²	19M	30.6G	186	150	79.3
	640 ²	19M	54.8G	93	71	76.0
	768 ²	19M	91.4G	54	37 [†]	71.4
	1024 ²	19M	218.9G	21	11 [†]	58.9
VMamba	224 ²	31M	4.9G	1226	398	82.5
	384 ²	31M	14.3G	450	154	82.5
	512 ²	31M	25.4G	272	100	81.1
	640 ²	31M	39.6G	170	60	79.3
	768 ²	31M	57.1G	117	42 [†]	76.1
	1024 ²	31M	101.5G	66	25 [†]	62.3

Table 4: Comparison of generalizability to inputs with increased spatial resolutions. The throughput is measured with batch size of 32 using Pytorch 2.0. † denotes that the batch size ≤ 16 due to out-of-memory (OOM). we re-implemented the HiViT-T, as the checkpoint of HiViT-T has not been released.

adjusted to maintain similar architectural parameters and FLOPs for fair comparison. As illustrated in Figure 12, these approaches include unidirectional scanning (*Unidi-Scan*), bidirectional scanning (*Bidi-Scan*), and cascade scanning (*Cascade-Scan*, scans the data row-wise and column-wise successively). Both Unidi-Scan and Bidi-Scan move across the image in between the upper-left and lower-right corners.

Based on the results summarized in Table 5.4, Notably, Unidi-Scan, Bidi-Scan, and CSM are all implemented in Triton, and they exhibit little difference in training or inference speed. However, CSM demonstrates superior data modeling capacity, as evidenced by its higher classification accuracy, likely benefiting from the two-dimensional prior introduced by the Cross-Scan design. Regarding

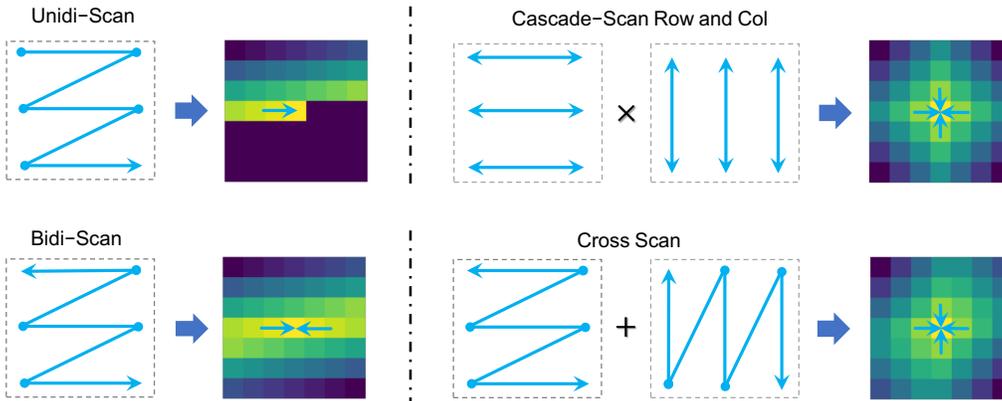


Figure 12: Illustration of different scanning methods for selective scan.

Method	#Param.	FLOPs	Throughput	Train Throughput.	ImageNet top-1 acc.
Unidi-Scan	30.70	4.86	1342	464	82.2
Bidi-Scan	30.70	4.86	1344	465	74.8 [†]
Cascade-Scan	30.70	4.86	817	253	–
CSM	30.70	4.86	1343	464	82.5

Table 5: Performance comparison of different scanning approaches. The proposed CSM achieves superior performance in both speed, while maintaining the same number of parameters and FLOPs. The relatively poor performance obtained by models configured with Bidi-Scan is caused by the numerical instability, and thus we report the result before encountering NaN outputs ([†]). We also omitted the results of models with Cascade-Scan due to the low training throughput.

Cascade-Scan, its practical application is severely hindered by its extremely slow computational speed, mainly due to the poor compatibility between selective scanning and large-dimensional data, as well as the multi-step scanning mechanism. We also note that models configured with Bidi-Scan consistently experience unstable training during experiments, often encountering numerical errors (NaN) after approximately 50 epochs, even with float32 datatype for the s6 process.

6 Conclusion

We introduce VMamba, a versatile backbone network designed for efficient visual representation learning using State Space Models (SSMs). The primary goal of VMamba is to incorporate the advantages of selective SSMS, including global receptive field, input-dependent weighting parameters, and linear computational complexity, into vision data processing. Specifically, we propose the Cross-Scan Module (CSM) to bridge the gap between 1D selective scanning and 2D vision data, and we illustrate its relationship with the attention mechanism and its effectiveness in achieving a global receptive field through both mathematical derivation and qualitative visualization. Furthermore, we have significantly enhanced the inference speed of VMamba by improving both its technical implementation and architectural design. The effectiveness of the VMamba family, including the VMamba-T/S/B models, has been demonstrated through extensive experiments and ablation studies, surpassing the performance of popular CNNs and vision transformers. Moreover, VMamba demonstrates remarkable scalability with increasing input resolution, exhibiting minimal performance degradation while maintaining linear computational complexity.

References

- [1] Guy E Blelloch. Prefix sums and their applications. 1990.

- [2] Chun-Fu Chen, Quanfu Fan, and Rameswar Panda. Crossvit: Cross-attention multi-scale vision transformer for image classification. In *ICCV*, 2021.
- [3] Kai Chen, Jiaqi Wang, Jiangmiao Pang, Yuhang Cao, Yu Xiong, Xiaoxiao Li, Shuyang Sun, Wansen Feng, Ziwei Liu, Jiarui Xu, Zheng Zhang, Dazhi Cheng, Chenchen Zhu, Tianheng Cheng, Qijie Zhao, Buyu Li, Xin Lu, Rui Zhu, Yue Wu, Jifeng Dai, Jingdong Wang, Jianping Shi, Wanli Ouyang, Chen Change Loy, and Dahua Lin. Mmdetection: Open mmlab detection toolbox and benchmark. *arXiv preprint arXiv:1906.07155*, 2019.
- [4] Yunpeng Chen, Jianan Li, Huaxin Xiao, Xiaojie Jin, Shuicheng Yan, and Jiashi Feng. Dual path networks. *NeurIPS*, 30, 2017.
- [5] François Chollet. Xception: Deep learning with depthwise separable convolutions. In *CVPR*, pages 1251–1258, 2017.
- [6] Jifeng Dai, Haozhi Qi, Yuwen Xiong, Yi Li, Guodong Zhang, Han Hu, and Yichen Wei. Deformable convolutional networks. In *ICCV*, pages 764–773, 2017.
- [7] Zihang Dai, Hanxiao Liu, Quoc V Le, and Mingxing Tan. Coatnet: Marrying convolution and attention for all data sizes. *NeurIPS*, 34:3965–3977, 2021.
- [8] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Fei-Fei Li. Imagenet: A large-scale hierarchical image database. In *CVPR*, pages 248–255, 2009.
- [9] Mingyu Ding, Bin Xiao, Noel Codella, Ping Luo, Jingdong Wang, and Lu Yuan. Davit: Dual attention vision transformers. In *ECCV*, pages 74–92, 2022.
- [10] Xiaohan Ding, Xiangyu Zhang, Jungong Han, and Guiguang Ding. Scaling up your kernels to 31x31: Revisiting large kernel design in cnns. In *CVPR*, pages 11963–11975, 2022.
- [11] Xiaoyi Dong, Jianmin Bao, Dongdong Chen, Weiming Zhang, Nenghai Yu, Lu Yuan, Dong Chen, and Baining Guo. Cswin transformer: A general vision transformer backbone with cross-shaped windows. In *CVPR*, pages 12124–12134, 2022.
- [12] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. In *ICLR*, 2021.
- [13] Daniel Y Fu, Tri Dao, Khaled Kamal Saab, Armin W Thomas, Atri Rudra, and Christopher Re. Hungry hungry hippos: Towards language modeling with state space models. In *ICLR*, 2022.
- [14] Albert Gu and Tri Dao. Mamba: Linear-time sequence modeling with selective state spaces. *arXiv preprint arXiv:2312.00752*, 2023.
- [15] Albert Gu, Tri Dao, Stefano Ermon, Atri Rudra, and Christopher Ré. Hippo: Recurrent memory with optimal polynomial projections. *NeurIPS*, 33:1474–1487, 2020.
- [16] Albert Gu, Karan Goel, Ankit Gupta, and Christopher Ré. On the parameterization and initialization of diagonal state space models. *NeurIPS*, 35:35971–35983, 2022.
- [17] Albert Gu, Karan Goel, and Christopher Re. Efficiently modeling long sequences with structured state spaces. In *ICLR*, 2021.
- [18] Albert Gu, Isys Johnson, Karan Goel, Khaled Saab, Tri Dao, Atri Rudra, and Christopher Ré. Combining recurrent, convolutional, and continuous-time models with linear state space layers. *NeurIPS*, 34:572–585, 2021.
- [19] Ankit Gupta, Albert Gu, and Jonathan Berant. Diagonal state spaces are as effective as structured state spaces. *NeurIPS*, 35:22982–22994, 2022.
- [20] Qi Han, Zejia Fan, Qi Dai, Lei Sun, Ming-Ming Cheng, Jiaying Liu, and Jingdong Wang. Demystifying local vision transformer: Sparse connectivity, weight sharing, and dynamic weight. *arXiv preprint arXiv:2106.04263*, 2021.
- [21] Ramin Hasani, Mathias Lechner, Tsun-Hsuan Wang, Makram Chahine, Alexander Amini, and Daniela Rus. Liquid structural state-space models. In *ICLR*, 2022.
- [22] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, pages 770–778, 2016.

- [23] Andrew G Howard, Menglong Zhu, Bo Chen, Dmitry Kalenichenko, Weijun Wang, Tobias Weyand, Marco Andreetto, and Hartwig Adam. Mobilenets: Efficient convolutional neural networks for mobile vision applications. *arXiv preprint arXiv:1704.04861*, 2017.
- [24] Binh-Son Hua, Minh-Khoi Tran, and Sai-Kit Yeung. Pointwise convolutional neural networks. In *CVPR*, pages 984–993, 2018.
- [25] Gao Huang, Zhuang Liu, Laurens Van Der Maaten, and Kilian Q Weinberger. Densely connected convolutional networks. In *CVPR*, pages 4700–4708, 2017.
- [26] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E. Hinton. Imagenet classification with deep convolutional neural networks. In *NeurIPS*, pages 1106–1114, 2012.
- [27] P. Langley. Crafting papers on machine learning. In *ICML*, pages 1207–1216, 2000.
- [28] Yann LeCun, Bernhard Boser, John S Denker, Donnie Henderson, Richard E Howard, Wayne Hubbard, and Lawrence D Jackel. Backpropagation applied to handwritten zip code recognition. *Neural computation*, 1(4):541–551, 1989.
- [29] Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.
- [30] Tsung-Yi Lin, Michael Maire, Serge J. Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C. Lawrence Zitnick. Microsoft COCO: common objects in context. In *ECCV*, pages 740–755, 2014.
- [31] Ze Liu, Han Hu, Yutong Lin, Zhuliang Yao, Zhenda Xie, Yixuan Wei, Jia Ning, Yue Cao, Zheng Zhang, Li Dong, et al. Swin transformer v2: Scaling up capacity and resolution. In *CVPR*, pages 12009–12019, 2022.
- [32] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. Swin transformer: Hierarchical vision transformer using shifted windows. In *ICCV*, pages 10012–10022, 2021.
- [33] Zhuang Liu, Hanzi Mao, Chao-Yuan Wu, Christoph Feichtenhofer, Trevor Darrell, and Saining Xie. A convnet for the 2020s. In *CVPR*, pages 11976–11986, 2022.
- [34] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*, 2017.
- [35] Jiasen Lu, Roozbeh Mottaghi, Aniruddha Kembhavi, et al. Container: Context aggregation networks. *NeurIPS*, 34:19160–19171, 2021.
- [36] Wenjie Luo, Yujia Li, Raquel Urtasun, and Richard Zemel. Understanding the effective receptive field in deep convolutional neural networks. *NeurIPS*, 29, 2016.
- [37] Xuezhe Ma, Chunting Zhou, Xiang Kong, Junxian He, Liangke Gui, Graham Neubig, Jonathan May, and Luke Zettlemoyer. Mega: Moving average equipped gated attention. In *ICLR*, 2022.
- [38] Eric Martin and Chris Cundy. Parallelizing linear recurrent neural nets over sequence length. In *ICLR*, 2018.
- [39] Harsh Mehta, Ankit Gupta, Ashok Cutkosky, and Behnam Neyshabur. Long range language modeling via gated state spaces. In *ICLR*, 2023.
- [40] Eric Nguyen, Karan Goel, Albert Gu, Gordon Downs, Preey Shah, Tri Dao, Stephen Baccus, and Christopher Ré. S4nd: Modeling images and videos as multidimensional signals with state spaces. *NeurIPS*, 35:2846–2861, 2022.
- [41] Ilija Radosavovic, Raj Prateek Kosaraju, Ross Girshick, Kaiming He, and Piotr Dollár. Designing network design spaces. In *CVPR*, pages 10428–10436, 2020.
- [42] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014.
- [43] Jimmy TH Smith, Andrew Warrington, and Scott Linderman. Simplified state space layers for sequence modeling. In *ICLR*, 2022.
- [44] Jimmy TH Smith, Andrew Warrington, and Scott Linderman. Simplified state space layers for sequence modeling. In *ICLR*, 2022.

- [45] Aravind Srinivas, Tsung-Yi Lin, Niki Parmar, Jonathon Shlens, Pieter Abbeel, and Ashish Vaswani. Bottleneck transformers for visual recognition. In *CVPR*, pages 16519–16529, 2021.
- [46] Christian Szegedy, Wei Liu, Yangqing Jia, Pierre Sermanet, Scott Reed, Dragomir Anguelov, Dumitru Erhan, Vincent Vanhoucke, and Andrew Rabinovich. Going deeper with convolutions. In *CVPR*, pages 1–9, 2015.
- [47] Mingxing Tan and Quoc V. Le. Efficientnet: Rethinking model scaling for convolutional neural networks. In *ICML*, pages 6105–6114, 2019.
- [48] Yi Tay, Mostafa Dehghani, Dara Bahri, and Donald Metzler. Efficient transformers: A survey. *ACM Computing Surveys*, 55(6), 2022.
- [49] Yunjie Tian, Lingxi Xie, Zhaozhi Wang, Longhui Wei, Xiaopeng Zhang, Jianbin Jiao, Yaowei Wang, Qi Tian, and Qixiang Ye. Integrally pre-trained transformer pyramid networks. In *CVPR*, pages 18610–18620, 2023.
- [50] Hugo Touvron, Matthieu Cord, Matthijs Douze, Francisco Massa, Alexandre Sablayrolles, and Hervé Jégou. Training data-efficient image transformers & distillation through attention. In *ICML*, pages 10347–10357, 2021.
- [51] Ashish Vaswani, Prajit Ramachandran, Aravind Srinivas, Niki Parmar, Blake Hechtman, and Jonathon Shlens. Scaling local self-attention for parameter efficient visual backbones. In *CVPR*, pages 12894–12904, 2021.
- [52] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *NeurIPS*, 2017.
- [53] Jue Wang, Wentao Zhu, Pichao Wang, Xiang Yu, Linda Liu, Mohamed Omar, and Raffay Hamid. Selective structured state-spaces for long-form video understanding. In *CVPR*, pages 6387–6397, 2023.
- [54] Wenhai Wang, Enze Xie, Xiang Li, Deng-Ping Fan, Kaitao Song, Ding Liang, Tong Lu, Ping Luo, and Ling Shao. Pyramid vision transformer: A versatile backbone for dense prediction without convolutions. In *ICCV*, pages 568–578, 2021.
- [55] Wenhai Wang, Enze Xie, Xiang Li, Deng-Ping Fan, Kaitao Song, Ding Liang, Tong Lu, Ping Luo, and Ling Shao. Pvt v2: Improved baselines with pyramid vision transformer. *Computational Visual Media*, 8(3):415–424, 2022.
- [56] Ross Wightman. Pytorch image models. <https://github.com/rwightman/pytorch-image-models>, 2019.
- [57] Tete Xiao, Yingcheng Liu, Bolei Zhou, Yuning Jiang, and Jian Sun. Unified perceptual parsing for scene understanding. In *ECCV*, pages 418–434, 2018.
- [58] Saining Xie, Ross Girshick, Piotr Dollár, Zhuowen Tu, and Kaiming He. Aggregated residual transformations for deep neural networks. In *CVPR*, pages 1492–1500, 2017.
- [59] Jianwei Yang, Chunyuan Li, Pengchuan Zhang, Xiyang Dai, Bin Xiao, Lu Yuan, and Jianfeng Gao. Focal self-attention for local-global interactions in vision transformers. *arXiv preprint arXiv:2107.00641*, 2021.
- [60] Songlin Yang, Bailin Wang, Yikang Shen, Rameswar Panda, and Yoon Kim. Gated linear attention transformers with hardware-efficient training. *arXiv preprint arXiv:2312.06635*, 2023.
- [61] Fisher Yu and Vladlen Koltun. Multi-scale context aggregation by dilated convolutions. *arXiv preprint arXiv:1511.07122*, 2015.
- [62] Li Yuan, Yunpeng Chen, Tao Wang, Weihao Yu, Yujun Shi, Zi-Hang Jiang, Francis EH Tay, Jiashi Feng, and Shuicheng Yan. Tokens-to-token vit: Training vision transformers from scratch on imagenet. In *ICCV*, pages 558–567, 2021.
- [63] Xiang Zhang, Junbo Zhao, and Yann LeCun. Character-level convolutional networks for text classification. *NeurIPS*, 28, 2015.
- [64] Xiaosong Zhang, Yunjie Tian, Lingxi Xie, Wei Huang, Qi Dai, Qixiang Ye, and Qi Tian. Hivit: A simpler and more efficient design of hierarchical vision transformer. In *ICLR*, 2023.
- [65] Weixi Zhao, Weiqiang Wang, and Yunjie Tian. Graformer: Graph-oriented transformer for 3d pose estimation. In *CVPR*, pages 20438–20447, 2022.

- [66] Bolei Zhou, Hang Zhao, Xavier Puig, Sanja Fidler, Adela Barriuso, and Antonio Torralba. Scene parsing through ade20k dataset. In *CVPR*, 2017.
- [67] Daquan Zhou, Bingyi Kang, Xiaojie Jin, Linjie Yang, Xiaochen Lian, Zihang Jiang, Qibin Hou, and Jiashi Feng. Deepvit: Towards deeper vision transformer. *arXiv preprint arXiv:2103.11886*, 2021.
- [68] Lianghui Zhu, Bencheng Liao, Qian Zhang, Xinlong Wang, Wenyu Liu, and Xinggang Wang. Vision mamba: Efficient visual representation learning with bidirectional state space model. *arXiv preprint arXiv:2401.09417*, 2024.
- [69] Barret Zoph and Quoc V Le. Neural architecture search with reinforcement learning. *arXiv preprint arXiv:1611.01578*, 2016.