
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@cripac.ia.ac.cn

Abstract

Convolutional Neural Networks (CNNs) and Vision Transformers (ViTs) stand as the two most popular foundation models for visual representation learning. While CNNs exhibit remarkable scalability with linear complexity w.r.t. image resolution, ViTs surpass them in fitting capabilities despite contending with quadratic complexity. A closer inspection reveals that ViTs achieve superior visual modeling performance through the incorporation of global receptive fields and dynamic weights. This observation motivates us to propose a novel architecture that inherits these components while enhancing computational efficiency. To this end, we draw inspiration from the recently introduced state space model and propose the Visual State Space Model (VMamba), which achieves linear complexity without sacrificing global receptive fields. To address the encountered direction-sensitive issue, we introduce the Cross-Scan Module (CSM) to traverse the spatial domain and convert any non-causal visual image into order patch sequences. Extensive experimental results substantiate that VMamba not only demonstrates promising capabilities across various visual perception tasks, but also exhibits more pronounced advantages over established benchmarks as the image resolution increases. Source code has been available at <https://github.com/MzeroMiko/VMamba>.

1 Introduction

Visual representation learning is one of the most fundamental research topics in computer vision, which has experienced significant breakthroughs since the onset of the deep learning era. Two primary categories of deep foundation models, *i.e.*, Convolution Neural Networks (CNNs) [38, 19, 22, 29, 42] and Vision Transformers (ViTs) [10, 28, 45, 56], have been extensively employed in a variety of visual tasks. While both have achieved remarkable success in computing expressive visual representations, ViTs generally exhibit superior performance compared to CNNs, which could be attributed to global receptive fields and dynamic weights facilitated by the attention mechanism.

However, the attention mechanism requires quadratic complexity in terms of image sizes, resulting in expensive computational overhead when addressing downstream dense prediction tasks, such as object detection, semantic segmentation, *etc.* To tackle this issue, substantial effort has been dedicated to improving the efficiency of attention by constraining the size or stride of computing windows [43], albeit at the cost of imposing restrictions on the scale of receptive fields. This motivates us to design a novel visual foundation model with linear complexity, while still preserving the advantages associated with global receptive fields and dynamic weights.

Drawing inspiration from the recently proposed state space model [12, 34, 47], we introduce the Visual State Space Model (denoted as VMamba) for efficient visual representation learning. The pivotal concept behind VMamba’s success in effectively reducing attention complexity is inherited from the Selective Scan Space State Sequential Model (S6) [12], originally devised to address Natural Language Processing (NLP) tasks. In contrast to the conventional attention computation approach, S6 enables each element in a 1-D array (*e.g.*, text sequence) to interact with any of the previously scanned samples through a compressed hidden state, effectively reducing the quadratic complexity to linear.

However, due to the non-causal nature of visual data, directly applying such a strategy to a patchified and flattened image would inevitably result in restricted receptive fields, as the relationships against unscanned patches could not be estimated. We term this issue as the ‘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 four 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, which renders a global receptive field without increasing the linear computational complexity.

Extensive experiments on diverse visual tasks are conducted to verify the effectiveness of VMamba. As shown in Figure 1, VMamba models show superior or at least competitive performance on ImageNet-1K in comparison with benchmark vision models including Resnet [19], ViT [10], and Swin [28] ¹. We also report the results on downstream dense prediction tasks. For example, VMamba-Tiny/Small/Base (with 22/44/75 M parameters respectively) achieves 46.5%/48.2%/48.5% mAP on COCO using the MaskRCNN detector (1× training schedule) and 47.3%/49.5%/50.0% mIoU on ADE20K using UperNet with 512×512 inputs, demonstrating its potential to serve as a powerful foundation model. Furthermore, when larger images are used as input, the FLOPs of ViT increase significantly faster than those of CNN models, despite usually still exhibiting superior performance. However, it is intriguing that VMamba, being essentially a foundation model based on the Transformer architecture, is able to attain performance comparable to ViT with a steady increase in FLOPs.

We summarize the contributions below:

- We propose VMamba, a visual state space model with global receptive fields and dynamic weights for visual representation learning. VMamba presents a novel option for vision foundation models, extending beyond the existing choices of CNNs and ViTs.
- The Cross-Scan Module (CSM) is introduced to bridge the gap between 1-D array scanning and 2-D plain traversing, facilitating the extension of S6 to visual data without compromising the field of reception.
- Without bells and whistles, we show that VMamba achieves promising results across various visual tasks including image classification, object detection, and semantic segmentation. These findings underscore the potential of VMamba to serve as a robust vision foundation model.

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 [23, 38, 41, 19, 42] and ViTs [10, 28, 48, 9, 6, 56]. Recently, the success of State Space Models (SSMs) [12, 34, 47] has

¹We encounter a bug during the training of VMamba-B, and we will update the latest result as soon as possible

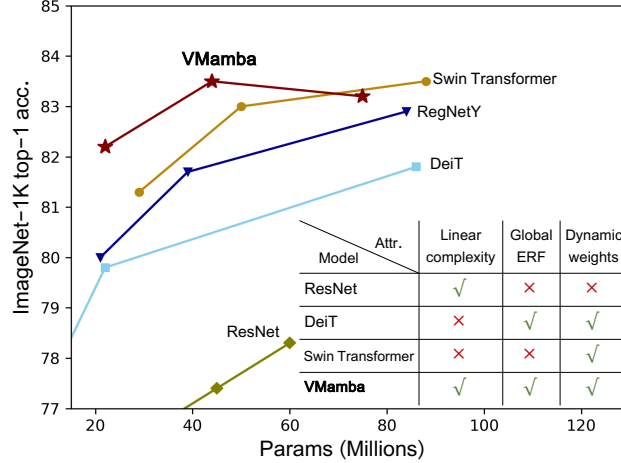


Figure 1: **Performance comparison on ImageNet-1K.** VMamba series achieves superior top-1 accuracy compared to popular counterparts. We note that the proposed VMamba has the capability of showing global effective reception field (ERF), dynamic weights with linear complexity.

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 [25, 23] are designed for basic tasks, such as recognizing handwritten digits [24] and classifying character categories [55]. 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 [7], increasingly deeper [38, 41, 19, 22] and efficient models [20, 42, 52, 36] 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 [4, 21, 53, 5] or more efficient network architectures [59, 3, 51, 20].

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 [10] and appear in a plain configuration [54, 58, 1, 31]. Later, DeiT [45] 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 [28, 9, 48, 31, 56, 44, 6, 8, 57] to gradually decrease the feature resolution throughout the backbone. Moreover, other studies propose to harness the advantages of CNNs, such as introducing convolution operations [49, 6, 46], designing hybrid architectures by combining CNN and ViT modules [6, 40, 31], *etc.*

State Space Models (SSMs) are recently proposed models that are introduced into deep learning as state space transforming [16, 15, 39]. Inspired by continuous state space models in control systems, combined with HiPPO [13] initialization, LSSL [16] 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 [15] proposes to normalize the parameter into diagonal structure. Since then, many flavors of structured state space models sprang up with different structures like complex-diagonal structure [17, 14], multiple-input multiple output supporting [39], decomposition of diagonal plus low-rank operations [18], selection mechanism [12]. These models are then integrated into large representation models [34, 33, 11].

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 [33, 34], content-based reasoning [12],

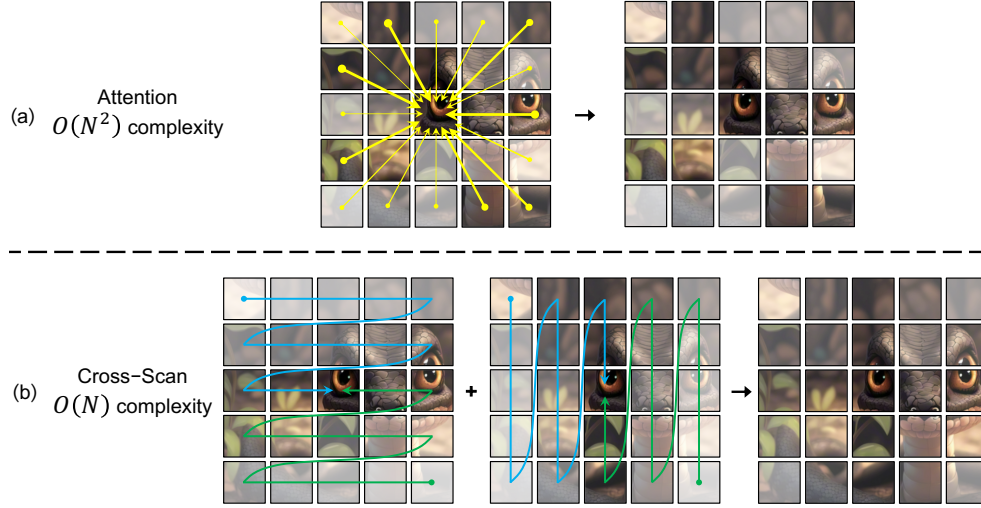


Figure 2: **Comparison of information flow: Attention vs. Cross-Scan Module (CSM).** (a) The attention mechanism uniformly integrates all pixels for the center pixel, resulting in $\mathcal{O}(N^2)$ complexity. (b) CSM integrates pixels from top-left, bottom-right, top-right, and bottom-left with $\mathcal{O}(N)$ complexity.

pixel-level 1-D image classification [15], few of them pay attention in visual recognition. The most similar work to ours is S4ND [35]. S4ND is the first work applying state space mechanism into visual tasks and showing the potential that its performance may compete with ViT [10]. However, S4ND expands the S4 model in a simple manner, fails on efficiently capturing image information in an input-dependent manner. We demonstrate that with selective scan mechanism introduced by mamba [12], the proposed VMamba is able to match existing popular vision foundation models like ResNet [19], ViT [10], swin [27], and convnext [29], showcasing the potential of VMamba to be the powerful foundation model.

3 Method

In this section, we start by introducing the preliminary concepts related to VMamba, including the state space models, the discretization process, and the selective scan mechanism. We then provide detailed specifications of the 2D state space model which serves as the core element of VMamba. Finally, we present a comprehensive discussion of the overall VMamba architecture.

3.1 Preliminaries

State Space Models. State Space Models (SSMs) are commonly considered as linear time-invariant systems that map stimulation $x(t) \in \mathbb{R}^L$ to response $y(t) \in \mathbb{R}^L$. Mathematically, these models are typically formulated as linear ordinary differential equations (ODEs) (Eq. 1), where the parameters include $A \in \mathbb{C}^{N \times N}$, $B, C \in \mathbb{C}^N$ for a state size N , and the skip connection $D \in \mathbb{C}^1$.

$$\begin{aligned} h'(t) &= Ah(t) + Bx(t) \\ y(t) &= Ch(t) + Dx(t) \end{aligned} \tag{1}$$

Discretization. State Space Models (SSMs), as continuous-time models, face great challenges when integrated into deep learning algorithms. To overcome this obstacle, the discretization process becomes imperative.

The primary objective of discretization is to transform the ODE into a discrete function. This transformation is crucial to align the model with the sample rate of the underlying signal embodied in the input data, enabling computationally efficient operations [16]. Considering the input $x_k \in \mathbb{R}^{L \times D}$,

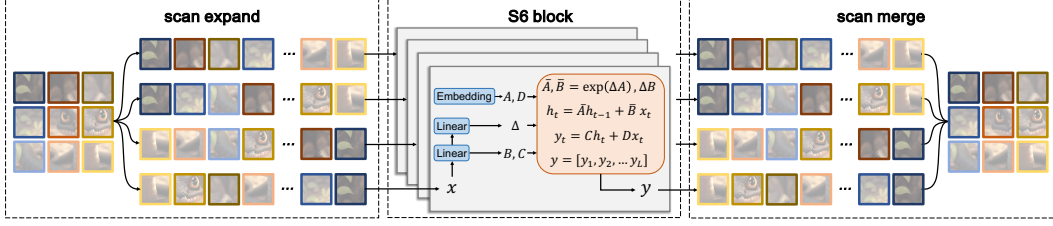


Figure 3: **Illustration of the 2D-Selective-Scan on an image.** We commence by scanning an image using CSM (*scan expand*). The four resulting features are then individually processed through the S6 block, and the four output features are merged (*scan merge*) to construct the final 2D feature map.

a sampled vector within the signal flow of length L following [17], the ODE (Eq. 1) could be discretized as follows using the zeroth-order hold rule:

$$\begin{aligned}
 h_k &= \bar{A}h_{k-1} + \bar{B}x_k, \\
 y_k &= \bar{C}h_k + \bar{D}x_k, \\
 \bar{A} &= e^{\Delta A}, \\
 \bar{B} &= (e^{\Delta A} - I)A^{-1}B, \\
 \bar{C} &= C
 \end{aligned} \tag{2}$$

where $B, C \in \mathbb{R}^{D \times N}$ and $\Delta \in \mathbb{R}^D$. In practice, following [12], we refine the approximation of \bar{B} using the first-order Taylor series:

$$\bar{B} = (e^{\Delta A} - I)A^{-1}B \approx (\Delta A)(\Delta A)^{-1}\Delta B = \Delta B \tag{3}$$

Selective Scan Mechanism. Diverging from the prevalent approach that predominantly focuses on linear time-invariant (LTI) SSMs, the proposed VMamba sets itself apart by incorporating the selective scan mechanism (S6) [12] as the core SSM operator. In S6, the matrices $B \in \mathbb{R}^{B \times L \times N}$, $C \in \mathbb{R}^{B \times L \times N}$, and $\Delta \in \mathbb{R}^{B \times L \times D}$ are derived from the input data $x \in \mathbb{R}^{B \times L \times D}$. This implies that S6 is aware of the contextual information embedded in the input, ensuring the dynamism of weights within this mechanism.

3.2 2D Selective Scan

Despite its distinctive characteristics, S6 causally processes the input data, and thus can only capture information within the scanned part of the data. This naturally aligns S6 with NLP tasks that involve temporal data but poses significant challenges when adapting to non-causal data such as image, graph, set, *etc.* A straightforward solution to this problem would be to scan data along two different directions (*i.e.*, forward and backward), allowing them to compensate for the receptive field of each other without increasing the computational complexity.

Despite the non-causal nature, images differ from texts in that they contain 2D spatial information (*e.g.* local texture and global structure). To tackle this problem, S4ND [35] suggests reformulating SSM with convolution and straightforwardly expanding the kernel from 1-D to 2-D via outer-product. However, such modification prevents the weights from being dynamic (*i.e.*, input independent), resulting in a loss of the context-based data modeling capability. Therefore, we choose to preserve dynamic weights by sticking to the selective scan approach [12], which unfortunately disallows us to follow [35] and integrate convolution operations.

To address this problem, we propose the Cross-Scan Module (CSM) as shown in Figure 2. We choose to unfold image patches along rows and columns into sequences (*scan expand*), and then proceed with scanning along four different directions: top-left to bottom-right, bottom-right to top-left, top-right to bottom-left, and bottom-left to top-right. In this way, any pixel (such as the center pixel in Figure 2) integrates information from all other pixels in different directions. We then reshape each sequence into a single image, and all sequences are merged to form a new one as illustrated in Figure 3 (*scan merge*).

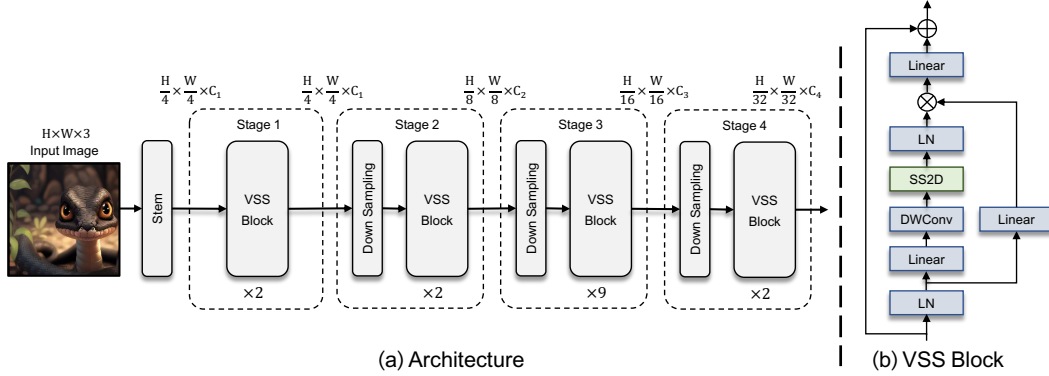


Figure 4: (a) The overall architecture of a VMamba model (VMamba-T); (b) The fundamental building block of VMamba, namely the VSS block.

The integration of S6 with CSM, referred to as the S6 block, serves as the core element to construct the Visual State Space (VSS) block, which constitutes the fundamental building block of VMamba (further detailed in the next subsection). We emphasize that the S6 block inherits the linear complexity of the selective scan mechanism while retaining a global receptive field, which aligns with our motivation to construct such a vision model.

3.3 VMamba Model

3.3.1 Overall Architecture

An overview of the architecture of VMamba-Tiny is illustrated in Figure 4 (a). VMamba begins the process by partitioning the input image into patches using a stem module, similar to ViTs, but without further flattening the patches into a 1-D sequence. This modification preserves the 2D structure of images, resulting in a feature map with dimensions of $\frac{H}{4} \times \frac{W}{4} \times C_1$.

VMamba then stacks several VSS blocks on the feature map, maintaining the same dimension, constituting “Stage 1”. Hierarchical representations in VMamba are built by down-sampling the feature map in “Stage 1” through a patch merge operation [27]. Subsequently, more VSS blocks are involved, resulting in an output resolution of $\frac{H}{8} \times \frac{W}{8}$ and forming “Stage 2”. This procedure is repeated to create “Stage 3” and “Stage 4” with resolutions of $\frac{H}{16} \times \frac{W}{16}$ and $\frac{H}{32} \times \frac{W}{32}$, respectively. All these stages collectively construct hierarchical representations akin to popular CNN models [19, 22, 41, 29, 42], and some ViTs [27, 48, 6, 56]. The resulting architecture can serve as a versatile replacement for other vision models in practical applications with similar requirements.

We develop VMamba in three distinct scales, *i.e.*, VMamba-Tiny, VMamba-Small, and VMamba-Base (referred to as VMamba-T, VMamba-S, and VMamba-B, respectively). Detailed architectural specifications are outlined in Table 1. The FLOPs for all models are assessed using a 224×224 input size. Additional architectures, such as a large-scale model, will be introduced in future updates.

3.3.2 VSS Block

The structure of VSS block is illustrated in Figure 4 (b). The input undergoes an initial linear embedding layer, and the output splits into two information flows. One flow passes through a 3×3 depth-wise convolution layer, followed by a Silu activation function [37] before entering the core SS2D module. The output of SS2D goes through a layer normalization layer and is then added to the output of the other information flow, which has undergone a Silu activation. This combination produces the final output of the VSS block.

Unlike vision transformers, we refrain from utilizing position embedding bias in VMamba due to its causal nature. Our design diverges from the typical vision transformer structure, which employs the following sequence of operations: Norm \rightarrow attention \rightarrow Norm \rightarrow MLP in a block, and discards

layer name	output size	Tiny	Small	Base
stem	112×112	conv 4×4, 96, stride 4	conv 4×4, 96, stride 4	conv 4×4, 128, stride 4
stage 1	56×56	$\begin{bmatrix} \text{linear } 96 \rightarrow 2 \times 96 \\ \text{DWConv } 3 \times 3, 2 \times 96 \\ \text{SS2D, dim } 2 \times 96 \\ \text{linear } 2 \times 96 \rightarrow 96 \end{bmatrix} \times 2$	$\begin{bmatrix} \text{linear } 96 \rightarrow 2 \times 96 \\ \text{DWConv } 3 \times 3, 2 \times 96 \\ \text{SS2D, dim } 2 \times 96 \\ \text{linear } 2 \times 96 \rightarrow 96 \end{bmatrix} \times 2$	$\begin{bmatrix} \text{linear } 128 \rightarrow 2 \times 128 \\ \text{DWConv } 3 \times 3, 2 \times 128 \\ \text{SS2D, dim } 2 \times 128 \\ \text{linear } 2 \times 128 \rightarrow 128 \end{bmatrix} \times 2$
		patch merging → 192	patch merging → 192	patch merging → 256
stage 2	28×28	$\begin{bmatrix} \text{linear } 192 \rightarrow 2 \times 192 \\ \text{DWConv } 3 \times 3, 2 \times 192 \\ \text{SS2D, dim } 2 \times 192 \\ \text{linear } 2 \times 192 \rightarrow 192 \end{bmatrix} \times 2$	$\begin{bmatrix} \text{linear } 192 \rightarrow 2 \times 192 \\ \text{DWConv } 3 \times 3, 2 \times 192 \\ \text{SS2D, dim } 2 \times 192 \\ \text{linear } 2 \times 192 \rightarrow 192 \end{bmatrix} \times 2$	$\begin{bmatrix} \text{linear } 256 \rightarrow 2 \times 256 \\ \text{DWConv } 3 \times 3, 2 \times 256 \\ \text{SS2D, dim } 2 \times 256 \\ \text{linear } 2 \times 256 \rightarrow 256 \end{bmatrix} \times 2$
		patch merging → 384	patch merging → 384	patch merging → 512
stage 3	14×14	$\begin{bmatrix} \text{linear } 384 \rightarrow 2 \times 384 \\ \text{DWConv } 3 \times 3, 2 \times 384 \\ \text{SS2D, dim } 2 \times 384 \\ \text{linear } 2 \times 384 \rightarrow 384 \end{bmatrix} \times 9$	$\begin{bmatrix} \text{linear } 384 \rightarrow 2 \times 384 \\ \text{DWConv } 3 \times 3, 2 \times 384 \\ \text{SS2D, dim } 2 \times 384 \\ \text{linear } 2 \times 384 \rightarrow 384 \end{bmatrix} \times 27$	$\begin{bmatrix} \text{linear } 512 \rightarrow 2 \times 512 \\ \text{DWConv } 3 \times 3, 2 \times 512 \\ \text{SS2D, dim } 2 \times 512 \\ \text{linear } 2 \times 512 \rightarrow 512 \end{bmatrix} \times 27$
		patch merging → 768	patch merging → 768	patch merging → 1024
stage 4	7×7	$\begin{bmatrix} \text{linear } 768 \rightarrow 2 \times 768 \\ \text{DWConv } 3 \times 3, 2 \times 768 \\ \text{SS2D, dim } 2 \times 768 \\ \text{linear } 2 \times 768 \rightarrow 768 \end{bmatrix} \times 2$	$\begin{bmatrix} \text{linear } 768 \rightarrow 2 \times 768 \\ \text{DWConv } 3 \times 3, 2 \times 768 \\ \text{SS2D, dim } 2 \times 768 \\ \text{linear } 2 \times 768 \rightarrow 768 \end{bmatrix} \times 2$	$\begin{bmatrix} \text{linear } 1024 \rightarrow 2 \times 1024 \\ \text{DWConv } 3 \times 3, 2 \times 1024 \\ \text{SS2D, dim } 2 \times 1024 \\ \text{linear } 2 \times 1024 \rightarrow 1024 \end{bmatrix} \times 2$
		1×1	average pool, 1000-d fc, softmax	
Param. (M)		22	44	75
FLOPs		4.5×10 ⁹	9.1×10 ⁹	15.2×10 ⁹

Table 1: **Architectural overview of the VMamba series.** Down-sampling is executed through patch merging operations in stages 1, 2, and 3. The term “linear” refers to a linear layer. The “DWConv” denotes a depth-wise convolution operation. The proposed 2D-selective-scan is labeled as “SS2D”.

the MLP operation. Consequently, the VSS block is shallower than the ViT block, which allows us to stack more blocks with a similar budget of total model depth.

4 Experiment

In this section, we perform a series of experiments to assess and compare VMamba against popular models, including CNNs and vision transformers. Our evaluation spans diverse tasks, including image classification on ImageNet-1K, object detection on COCO, and semantic segmentation on ADE20K. Subsequently, we delve into analysis experiments to gain deeper insights into the architecture of VMamba.

4.1 Image Classification on ImageNet-1K

Settings We evaluate VMamba’s classification performance on ImageNet-1K [7]. Following the configuration in [27], VMamba-T/S/B undergo training from scratch for 300 epochs (with the first 20 epochs to warmup), utilizing a batch size of 1024. The training process incorporates the AdamW optimizer 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 employed. Beyond these, no further training techniques are applied.

Results Table 2 summarizes results on ImageNet-1K, comparing VMamba with popular CNN models and vision transformers. The comparison reveals that, with similar FLOPs, VMamba-T achieves a performance of 82.2%, surpassing RegNetY-4G by 2.2%, DeiT-S by 2.4%, and Swin-T by 0.9%. Notably, the performance advantages of VMamba persist across small and base scale models. For instance, at the small scale, VMamba-S attains a top-1 accuracy of 83.5%, outperforming RegNetY-8G by 1.8% and Swin-S by 0.5%. Meanwhile, VMamba-B achieves a top-1 accuracy of

method	image size	#param.	FLOPs	ImageNet top-1 acc.
RegNetY-4G [36]	224 ²	21M	4.0G	80.0
RegNetY-8G [36]	224 ²	39M	8.0G	81.7
RegNetY-16G [36]	224 ²	84M	16.0G	82.9
EffNet-B3 [42]	300 ²	12M	1.8G	81.6
EffNet-B4 [42]	380 ²	19M	4.2G	82.9
EffNet-B5 [42]	456 ²	30M	9.9G	83.6
EffNet-B6 [42]	528 ²	43M	19.0G	84.0
ViT-B/16 [10]	384 ²	86M	55.4G	77.9
ViT-L/16 [10]	384 ²	307M	190.7G	76.5
DeiT-S [45]	224 ²	22M	4.6G	79.8
DeiT-B [45]	224 ²	86M	17.5G	81.8
DeiT-B [45]	384 ²	86M	55.4G	83.1
Swin-T [28]	224 ²	29M	4.5G	81.3
Swin-S [28]	224 ²	50M	8.7G	83.0
Swin-B [28]	224 ²	88M	15.4G	83.5
S4ND-ViT-B [35]	224 ²	89M	-	80.4
VMamba-T	224 ²	22M	4.5G	82.2
VMamba-S	224 ²	44M	9.1G	83.5
VMamba-B	224 ²	75M	15.2G	83.2 [†]

Table 2: **Accuracy comparison across various models on ImageNet-1K.** The symbol [†] indicates that a bug is encountered during the training of VMamba-B, and we will update the correct number in the near future.

83.2%, surpassing RegNetY-16G by 0.3% and DeiT-B by 0.1%. These promising results underscore VMamba’s potential as a robust foundational model, extending its superiority beyond traditional CNN models and vision transformers.

4.2 Object Detection on COCO

Settings In this section, we assess the performance of the proposed VMamba on object detection using the MSCOCO 2017 dataset [26]. Our training framework is built on the mmdetection library [2], and we adhere to the hyperparameters in Swin [27] with the Mask-RCNN detector. Specifically, we employ the AdamW optimizer 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%/0.2%/0.2%² for VMamba-T/S/B, respectively. The learning rate is initialized at 1×10^{-4} and is reduced by a factor of $10\times$ at the 9th and 11th epochs. We implement multi-scale training and random flip with a batch size of 16. These choices align with established practices for object detection evaluations.

Results The results for COCO are summarized in Table 3. 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 46.5%/48.2%/48.5%, surpassing Swin-T/S/B by 3.8%/3.6%/1.6% mAP, and ConvNeXt-T/S/B by 2.3%/2.8%/1.5% mAP. Using the same configuration, VMamba-T/S/B achieves instance segmentation mIoUs of 42.1%/43.0%/43.1%, outperforming Swin-T/S/B by 2.8%/2.1%/0.8% mIoU, and ConvNeXt-T/S/B by 2.0%/1.2%/0.7% mIoU, respectively.

Furthermore, the advantages of VMamba persist under the 36-epoch fine-tuning schedule with multi-scale training, as indicated in Table 3. When compared to counterparts, including Swin [28], ConvNeXt [29], PVTv2 [49], and ViT [10] (with Adapters), VMamba-T/S exhibit superior per-

²All being 0.2 is due to our oversight, and we will update the latest experiments.

Mask R-CNN 1× schedule								
Backbone	AP ^b	AP ^b ₅₀	AP ^b ₇₅	AP ^m	AP ^m ₅₀	AP ^m ₇₅	#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	46.5	68.5	50.7	42.1	65.5	45.3	42M	262G
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	48.2	69.7	52.5	43.0	66.6	46.4	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	48.5	69.6	53.0	43.1	67.0	46.4	96M	482G
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	48.5	69.9	52.9	43.2	66.8	46.3	42M	262G
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	49.7	70.4	54.2	44.0	67.6	47.3	64M	357G

Table 3: **Object detection and instance segmentation results on COCO dataset.** The FLOPs are calculated using inputs of size 1280×800 . Here, AP^b and AP^m denote box AP and mask AP, respectively. "1×" indicates models fine-tuned for 12 epochs, while "3×MS" signifies the utilization of multi-scale training for 36 epochs.

formance, achieving 48.5%/49.7% mAP on object detection and 43.2%/44.0% mIoU on instance segmentation. These results underscore the potential of VMamba in downstream dense prediction tasks.

4.3 Semantic Segmentation on ADE20K

Settings Following Swin [28], we construct a UperHead [50] on top of the pre-trained model. Employing the AdamW optimizer [30], we set the learning rate as 6×10^{-5} . The fine-tuning process 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 are presented in Table 4. Once again, VMamba exhibits superior accuracy, particularly with the VMamba-T model achieving 47.3% mIoU with a resolution of 512×512 and 48.3% mIoU using multi-scale (MS) input. These scores surpass all competitors, including ResNet [19], DeiT [45], Swin [28], and ConvNeXt [29]. Notably, the advantages extend to VMamba-S/B models, even when using 640×640 inputs.

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	47.3	48.3	55M	939G
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	49.5	50.5	76M	1037G
Swin-B	512^2	48.1	49.7	121M	1188G
ConvNeXt-B	512^2	49.1	49.9	122M	1170G
VMamba-B	512^2	50.0	51.3	110M	1167G
Swin-S	640^2	47.9	48.8	81M	1614G
ConvNeXt-S	640^2	48.8	48.9	82M	1607G
VMamba-S	640^2	50.8	50.8	76M	1620G

Table 4: **Semantic segmentation results on ADE20K using UperNet [50].** We evaluate the performance of semantic segmentation on the ADE20K dataset with UperNet [50]. The FLOPs are calculated with input sizes of 512×2048 or 640×2560 based on the crop size. "SS" and "MS" denote single-scale and multi-scale testing, respectively.

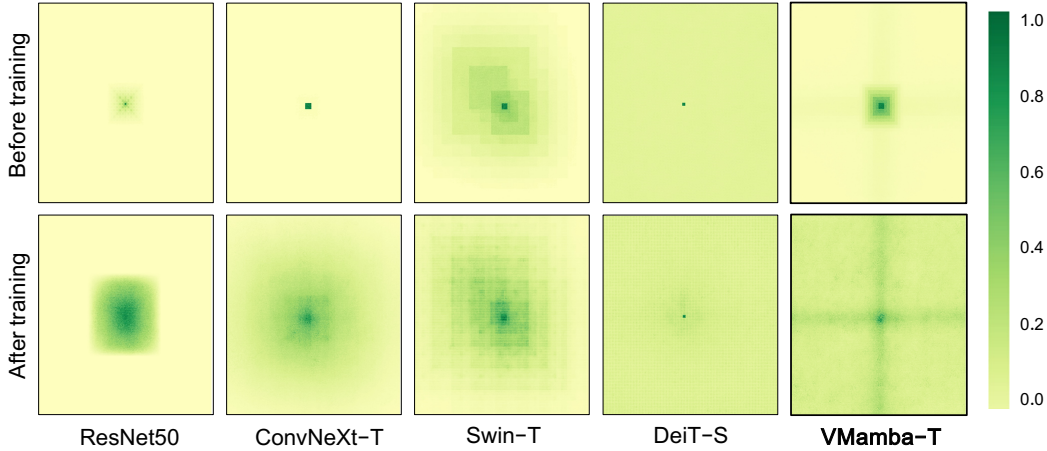


Figure 5: The **Effective Receptive Field (ERF)** is visualized for ResNet50 [19], ConvNeXt-T [29], Swin-T [28], DeiT-S [45] (ViT), and the proposed VMamba-T. A larger ERF is indicated by a more extensively distributed dark area. **Only DeiT [45] and the proposed VMamba exhibit a global ERF.** The inspiration for this visualization is drawn from [32].

4.4 Analysis Experiments

Effective Receptive Field To assess the effective receptive fields (ERFs) [32] across various models, we present a comparative analysis in Figure 5. The ERF measures the significance of model input concerning its output. Visualizing the ERF of the central pixel with an input size of 1024×1024 , we compare VMamba with four prominent visual foundation models: ResNet50 [19], ConvNeXt-T [29], Swin-T [28], and DeiT-S [45] (ViT) at both the `Before training` and `After training` stages. Key observations from Figure 5 include: 1) Only DeiT (ViT) and VMamba exhibit global ERFs, while other models demonstrate local ERFs, despite their theoretical global potential. It's important to note that the DeiT (ViT) model incurs quadratic complexity costs (refer to Figure 6). 2) In contrast to DeiT

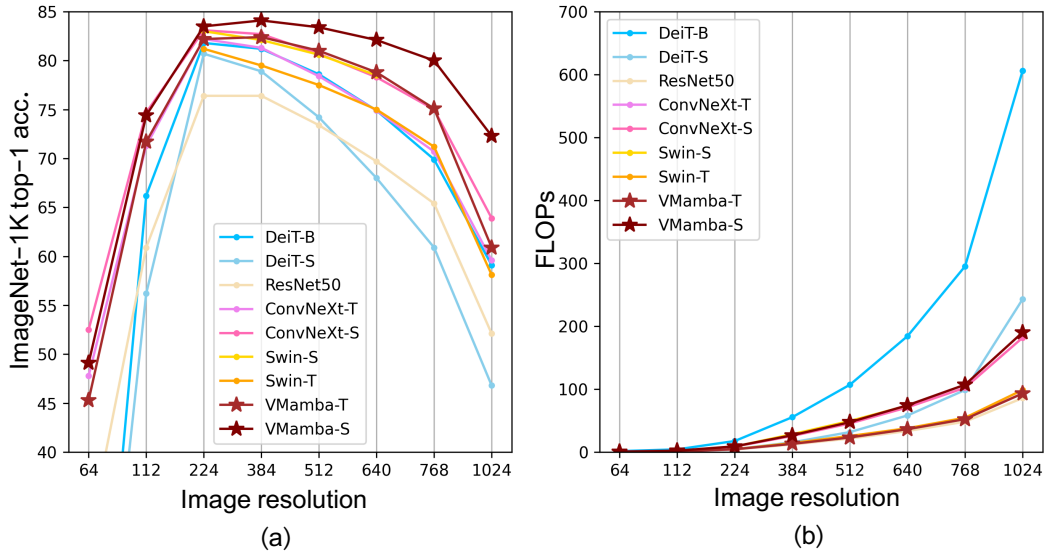


Figure 6: **A comparison of input scaling evaluation for popular models trained with 224×224 inputs.** We assess performance (a) and FLOPs (b) of various popular models trained with 224×224 inputs across different input sizes, ranging from 64×64 to 1024×1024 .

(ViT), which evenly activates all pixels using the attention mechanism, VMamba activates all pixels and notably emphasizes cross-shaped activations. The Cross-Scan Module’s scanning mechanism ensures the central pixel is most influenced by pixels along the cross, prioritizing long-dependency context over local information for each pixel. 3) Intriguingly, VMamba initially exhibits only a local ERF at `Before training`. However, After training transforms the ERF to global, signifying an adaptive process in the model’s global capability. We believe this adaptive process contributes to the model’s enhanced perception of images. This stands in contrast to DeiT, which maintains nearly identical ERFs at both `Before training` and `After training`.

Input Scaling We proceed to perform experiments on input scaling, measuring top-1 accuracy on ImageNet-1K and FLOPs, as illustrated in Figure 6. In Figure 6 (a), we assess the inference performance of popular models (trained with a 224×224 input size) across various image resolutions (ranging from 64×64 to 1024×1024). In comparison to counterparts, VMamba demonstrates the most stable performance across different input image sizes. Notably, as the input size increases from 224×224 to 384×384 , only VMamba exhibits an upward trend in performance (VMamba-S achieving 84%), highlighting its robustness to changes in input image size. In Figure 6 (b), we evaluate FLOPs using different image resolutions (also ranging from 64×64 to 1024×1024). As anticipated, the VMamba series report a linear growth in complexity, aligning with CNN models. VMamba’s complexity is consistent with carefully designed vision transformers like Swin [28]. However, it’s crucial to note that only VMamba achieves a global effective receptive field (ERF). DeiT, which also exhibits global ERF capability, experiences a quadratic growth in complexity.

5 Conclusion

Convolutional Neural Networks (CNNs) and Vision Transformers (ViTs) represent the predominant foundation models for visual representation learning. While CNNs exhibit linear complexity with respect to image resolution, ViTs excel in fitting capabilities despite quadratic complexity. Our investigation reveals that ViTs achieve superior visual modeling through global receptive fields and dynamic weights. Motivated by this, we propose the Visual State Space Model (VMamba), drawing inspiration from the state space model to achieve linear complexity without sacrificing global receptive fields. To address direction sensitivity, we introduce the Cross-Scan Module (CSM) for spatial traversal, converting non-causal visual images into ordered patch sequences. Extensive experiments demonstrate VMamba’s promising performance across visual tasks, with pronounced advantages as image resolution increases, surpassing established benchmarks.

References

- [1] Chun-Fu Chen, Quanfu Fan, and Rameswar Panda. Crossvit: Cross-attention multi-scale vision transformer for image classification. In *IEEE ICCV*, 2021.
- [2] 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.
- [3] Yunpeng Chen, Jianan Li, Huaxin Xiao, Xiaojie Jin, Shuicheng Yan, and Jiashi Feng. Dual path networks. *Advances in neural information processing systems*, 30, 2017.
- [4] François Chollet. Xception: Deep learning with depthwise separable convolutions. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1251–1258, 2017.
- [5] Jifeng Dai, Haozhi Qi, Yuwen Xiong, Yi Li, Guodong Zhang, Han Hu, and Yichen Wei. Deformable convolutional networks. In *Proceedings of the IEEE international conference on computer vision*, pages 764–773, 2017.
- [6] Zihang Dai, Hanxiao Liu, Quoc V Le, and Mingxing Tan. Coatnet: Marrying convolution and attention for all data sizes. *NeurIPS*, 34:3965–3977, 2021.
- [7] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Fei-Fei Li. Imagenet: A large-scale hierarchical image database. In *IEEE CVPR*, pages 248–255, 2009.
- [8] Mingyu Ding, Bin Xiao, Noel Codella, Ping Luo, Jingdong Wang, and Lu Yuan. Davit: Dual attention vision transformers. In *European Conference on Computer Vision*, pages 74–92. Springer, 2022.
- [9] 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 *IEEE CVPR*, pages 12124–12134, 2022.
- [10] 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.
- [11] 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 *The Eleventh International Conference on Learning Representations*, 2022.
- [12] Albert Gu and Tri Dao. Mamba: Linear-time sequence modeling with selective state spaces. *arXiv preprint arXiv:2312.00752*, 2023.
- [13] Albert Gu, Tri Dao, Stefano Ermon, Atri Rudra, and Christopher Ré. Hippo: Recurrent memory with optimal polynomial projections. *Advances in neural information processing systems*, 33:1474–1487, 2020.
- [14] Albert Gu, Karan Goel, Ankit Gupta, and Christopher Ré. On the parameterization and initialization of diagonal state space models. *Advances in Neural Information Processing Systems*, 35:35971–35983, 2022.
- [15] Albert Gu, Karan Goel, and Christopher Re. Efficiently modeling long sequences with structured state spaces. In *International Conference on Learning Representations*, 2021.
- [16] 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. *Advances in neural information processing systems*, 34:572–585, 2021.
- [17] Ankit Gupta, Albert Gu, and Jonathan Berant. Diagonal state spaces are as effective as structured state spaces. *Advances in Neural Information Processing Systems*, 35:22982–22994, 2022.
- [18] Ramin Hasani, Mathias Lechner, Tsun-Hsuan Wang, Makram Chahine, Alexander Amini, and Daniela Rus. Liquid structural state-space models. In *The Eleventh International Conference on Learning Representations*, 2022.
- [19] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *IEEE CVPR*, pages 770–778, 2016.

- [20] 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.
- [21] Binh-Son Hua, Minh-Khoi Tran, and Sai-Kit Yeung. Pointwise convolutional neural networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 984–993, 2018.
- [22] Gao Huang, Zhuang Liu, Laurens Van Der Maaten, and Kilian Q Weinberger. Densely connected convolutional networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4700–4708, 2017.
- [23] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E. Hinton. Imagenet classification with deep convolutional neural networks. In *NeurIPS*, pages 1106–1114, 2012.
- [24] 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.
- [25] 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.
- [26] 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.
- [27] 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 *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 12009–12019, 2022.
- [28] 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 *IEEE ICCV*, pages 10012–10022, 2021.
- [29] Zhuang Liu, Hanzi Mao, Chao-Yuan Wu, Christoph Feichtenhofer, Trevor Darrell, and Saining Xie. A convnet for the 2020s. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 11976–11986, 2022.
- [30] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*, 2017.
- [31] Jiasen Lu, Roozbeh Mottaghi, Aniruddha Kembhavi, et al. Container: Context aggregation networks. *NeurIPS*, 34:19160–19171, 2021.
- [32] Wenjie Luo, Yujia Li, Raquel Urtasun, and Richard Zemel. Understanding the effective receptive field in deep convolutional neural networks. *Advances in neural information processing systems*, 29, 2016.
- [33] Xuezhe Ma, Chunting Zhou, Xiang Kong, Junxian He, Liangke Gui, Graham Neubig, Jonathan May, and Luke Zettlemoyer. Mega: Moving average equipped gated attention. In *The Eleventh International Conference on Learning Representations*, 2022.
- [34] Harsh Mehta, Ankit Gupta, Ashok Cutkosky, and Behnam Neyshabur. Long range language modeling via gated state spaces. In *International Conference on Learning Representations*, 2023.
- [35] Eric Nguyen, Karan Goel, Albert Gu, Gordon Downs, Preety Shah, Tri Dao, Stephen Baccus, and Christopher Ré. S4nd: Modeling images and videos as multidimensional signals with state spaces. *Advances in neural information processing systems*, 35:2846–2861, 2022.
- [36] Ilija Radosavovic, Raj Prateek Kosaraju, Ross Girshick, Kaiming He, and Piotr Dollár. Designing network design spaces. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10428–10436, 2020.
- [37] Noam Shazeer. Glu variants improve transformer. *arXiv preprint arXiv:2002.05202*, 2020.
- [38] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014.
- [39] Jimmy TH Smith, Andrew Warrington, and Scott Linderman. Simplified state space layers for sequence modeling. In *The Eleventh International Conference on Learning Representations*, 2022.

- [40] Aravind Srinivas, Tsung-Yi Lin, Niki Parmar, Jonathon Shlens, Pieter Abbeel, and Ashish Vaswani. Bottleneck transformers for visual recognition. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 16519–16529, 2021.
- [41] Christian Szegedy, Wei Liu, Yangqing Jia, Pierre Sermanet, Scott Reed, Dragomir Anguelov, Dumitru Erhan, Vincent Vanhoucke, and Andrew Rabinovich. Going deeper with convolutions. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1–9, 2015.
- [42] Mingxing Tan and Quoc V. Le. Efficientnet: Rethinking model scaling for convolutional neural networks. In *ICML*, pages 6105–6114, 2019.
- [43] Yi Tay, Mostafa Dehghani, Dara Bahri, and Donald Metzler. Efficient transformers: A survey. volume 55, 2022.
- [44] 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 *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18610–18620, 2023.
- [45] 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.
- [46] Ashish Vaswani, Prajit Ramachandran, Aravind Srinivas, Niki Parmar, Blake Hechtman, and Jonathon Shlens. Scaling local self-attention for parameter efficient visual backbones. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12894–12904, 2021.
- [47] 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 *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6387–6397, 2023.
- [48] 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.
- [49] 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.
- [50] Tete Xiao, Yingcheng Liu, Bolei Zhou, Yuning Jiang, and Jian Sun. Unified perceptual parsing for scene understanding. In *ECCV*, pages 418–434, 2018.
- [51] Saining Xie, Ross Girshick, Piotr Dollár, Zhuowen Tu, and Kaiming He. Aggregated residual transformations for deep neural networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1492–1500, 2017.
- [52] 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.
- [53] Fisher Yu and Vladlen Koltun. Multi-scale context aggregation by dilated convolutions. *arXiv preprint arXiv:1511.07122*, 2015.
- [54] 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.
- [55] Xiang Zhang, Junbo Zhao, and Yann LeCun. Character-level convolutional networks for text classification. *Advances in neural information processing systems*, 28, 2015.
- [56] 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 *International Conference on Learning Representations*, 2023.
- [57] Weixi Zhao, Weiqiang Wang, and Yunjie Tian. Graformer: Graph-oriented transformer for 3d pose estimation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 20438–20447, 2022.
- [58] 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.
- [59] Barret Zoph and Quoc V Le. Neural architecture search with reinforcement learning. *arXiv preprint arXiv:1611.01578*, 2016.