This CVPR paper is the Open Access version, provided by the Computer Vision Foundation. Except for this watermark, it is identical to the accepted version; the final published version of the proceedings is available on IEEE Xplore.

# DefMamba: Deformable Visual State Space Model

Leiye Liu<sup>1</sup> Miao Zhang<sup>1,\*</sup> Jihao Yin<sup>1</sup> Tingwei Liu<sup>1</sup> Wei Ji<sup>2</sup> Yongri Piao<sup>1,\*</sup> Huchuan Lu<sup>1</sup> <sup>1</sup>Dalian University of Technology <sup>2</sup>Yale University

leiyeliu@mail.dlut.edu.cn, wei.ji@yale.edu, yrpiao@dlut.edu.cn

# Abstract

Recently, state space models (SSM), particularly Mamba, have attracted significant attention from scholars due to their ability to effectively balance computational efficiency and performance. However, most existing visual Mamba methods flatten images into 1D sequences using predefined scan orders, which results the model being less capable of utilizing the spatial structural information of the image during the feature extraction process. To address this issue, we proposed a novel visual foundation model called Def-Mamba. This model includes a multi-scale backbone structure and deformable mamba (DM) blocks, which dynamically adjust the scanning path to prioritize important information, thus enhancing the capture and processing of relevant input features. By combining a deformable scanning (DS) strategy, this model significantly improves its ability to learn image structures and detects changes in object details. Numerous experiments have shown that Def-Mamba achieves state-of-the-art performance in various visual tasks, including image classification, object detection, instance segmentation, and semantic segmentation. The code is open source on DefMamba.

# 1. Introduction

Most existing visual foundation models primarily rely on convolutional neural networks (CNNs) [25, 28, 32] and Transformer architectures [9, 24, 33]. However, CNNs are constrained by their sliding window structure, which limits the receptive field and significantly impedes global information aggregation across the input data. In contrast, Transformers excel in global information aggregation due to their attention mechanism, but their high computational complexity poses a challenge in achieving a balance between efficiency and performance. State space models (SSMs) [12] provide a potential solution to this trade-off. SSMs aggregate previous features through a hidden state matrix to update current features, thereby reducing the computational



Figure 1. Comparison of deformable scanning and other scanning methods. (a) Raster scanning [23, 47], (b) Local scanning [19], (c) Continuous scanning [42], (d) Our designed deformable scanning. The blue dots represent the reference points, and the red dots represent the deformable points. The yellow arrows represent the fixed scan order, and the red gradient arrows represent the deformable scan order. Our method exhibits an enhanced capacity to accurately capture the structural characteristics of objects, thereby enabling the development of a more refined scanning approach.

complexity to a linear relationship with the sequence length. Although SSMs process sequences in a recurrent manner, SSMs can perform calculations on sequences in parallel after simplification. Despite these advantages, SSMs struggle to capture long-range dependencies due to the lack of content-aware perception in the state matrix update process.

Recently, Mamba [11] proposes an improved selection mechanism designed to optimize the training process of SSMs. This innovative mechanism introduces content awareness into the feature extraction pipeline, expands the effective receptive field, and achieves remarkable performance enhancements in various NLP tasks. Consequently, numerous studies have attempted to extend this approach to

<sup>\*</sup>Corresponding authors.

the broader field of computer vision [1-4, 10, 22, 30, 40, 43, 46]. The principal challenge in this endeavor is how to effectively map 2D image feature maps into 1D sequences, without losing essential information. Most existing methods employ predefined strategies for mapping, such as raster scanning [23, 47], local scanning [19], and continuous scanning [42]. However, as illustrated in Figure 1, these methods all rely on a fixed scanning path. This results in adjacent tokens no longer being adjacent after flattening. Thereby, they neglect the inherent spatial structure of the image, leading to the loss of structural information. To address this issue, QuadMamba [39] determines the scanning window size based on the amount of information contained in different areas of the image. However, the scanning order in each window is fixed. This leads to an incomplete solution of the aforementioned issues. GrootV [38] adaptively constructs a tree topology based on input features and subsequently extracts features from this topology. Nevertheless, it employs only adjacent features in constructing the topology and distributes attention uniformly across the patch. The aforementioned methods are either based on a fixed scanning order, leading to the loss of structural information, or treat the information in the perception area equally, resulting in insensitivity to variations in object details.

To solve this problem, we proposed a novel framework called DefMamba inspired by deformable mechanisms [6, 36]. However, intuitively applying deformable mechanisms to SSMs still causes structural information loss and increases computational complexity. Therefore, we proposed a deformable state space model and a deformable scanning strategy (DS) to prioritize deformable tokens based on essential information and slide reference points towards the important area. This approach enables the SSMs to capture and process relevant features related to the input more effectively. Specifically, we shifted the reference point by the generated offset from a fixed position to an adjustable one that provided more useful information, thereby facilitating the awareness of changes in object details. On the other hand, we also dynamically adjusted the scanning order by the offset vector for obtaining a structureaware sequence. In this way, our framework adaptively perceives the variations in object details to find the most suitable feature points, and determines the optimal scanning order consistent with the object structure based on the input image features.

We conducted extensive experiments to validate the effectiveness of DefMamba across multiple visual benchmarks, including image classification on ImageNet [8], object detection and instance segmentation on COCO [21], and semantic segmentation on ADE20K [45]. These results demonstrate that our method outperforms existing SSMs based approaches on all benchmarks and remains competitive with CNN and transformer-based methods.

# 2. Related Work

# 2.1. Mamba for Visual Applications

Numerous studies have successfully integrated Mamba [11] into visual tasks [19, 23, 27, 29, 31, 38, 39, 42, 47], demonstrating preliminary achievements. ViM [47] introduces a bidirectional scanning approach to transform 2D image into 1D sequences, which are then fed into SSM for global context modeling, marking the first integration of Mamba into visual tasks. VMamba [23] employs a fourway scanning algorithm to convert 2D image into 1D sequences. PlainMamba [42] modifies the scanning method from raster to continuous, preserving the spatial dependencies of the image. MSVMamba [31] downsamples the sequence based on four scans to reduce computational redundancy and mitigates the issue of information loss. GrootV [38] constructs a minimum spanning tree on a four-way plane graph using the differences between adjacent features, dynamically adjusting the scanning order according to different inputs. QuadMamba [39] adaptively adjusts the window granularity during the scanning process based on the information content of the image to better aggregate local information. While GrootV and QuadMamba can adapt their scanning methods based on input data, GrootV only considers the relationships between adjacent elements when generating the minimum spanning tree, neglecting global information. On the other hand, QuadMamba still relies on predefined scanning methods and does not achieve true dynamic scanning. In contrast, our DefMamba introduces a content-aware deformable scanning strategy that allows the network to dynamically learn the scanning order and the reference points position.

### 2.2. Deformable CNNs and Attention Mechanism

Deformable convolution [6, 20, 41, 48] employs a convolution kernel that adapts to geometric variations in the input feature map, thereby overcoming the limitations of traditional convolution, which performs poorly when dealing with complex targets. Recently, the deformable mechanism has been extended to visual transformers [9] to enhance their ability to capture local features and adapt to geometric variations. DPT [5] proposes an adaptive patch embedding method that dynamically adjusts the position and size of patches while preserving their semantic information. PS-ViT [44] introduces a progressive sampling module prior to ViT [9], which iteratively identifies the most suitable deformable point positions for the current image. DAT [36] integrates the deformable mechanism with the self-attention mechanism in Transformers for the first time, incorporating deformable attention into the visual backbone. This approach learns a set of features corresponding to global keypoints and adapts to spatial variations. Previous approaches have explored various ways to effectively incorporate de-



Figure 2. Overview of DefMamba. (a) depicts the overall architecture of our network. (b) illustrates the structure of the deformable Mamba block. LN means LayerNorm and FFN is a feed-forward network.

formable mechanism into transformer architectures. With the recent popularity of Mamba [11], we attempted to introduce deformable mechanism into Mamba. However, directly applying these mechanism has resulted in issues like the loss of structural information and the necessity for additional modules. In this context, our designed DS strategy stands out by effectively prioritizing deformable tokens while directing reference points toward key areas.

# 3. Method

In this section, we first summarized the SSMs in Section 3.1. Then, in Section 3.2, we described the overall structure of the proposed network. Section 3.3 introduces a Deformable State Space Model (DSSM). Finally, Section 3.4 presents the designed model configurations across multiple scales.

#### 3.1. Preliminaries

SSMs, including notable implementations like S4 [12] and Mamba [11], are structured sequence architectures that combine elements of recurrent neural networks (RNNs) and CNNs, enabling linear or near-linear scaling with respect to sequence length. These models, derived from continuous systems, define a 1D function-to-function map for an input  $u(t) \in \mathbb{R}^L$  to an output  $y(t) \in \mathbb{R}^L$  through a hidden state  $h(t) \in \mathbb{R}^N$ . Where t represents time. More formally, SSMs are characterized by the continuous-time Ordinary Differential Equation (ODE) [11] presented:

$$h'(t) = \mathbf{A}h(t) + \mathbf{B}u(t),$$
  

$$y(t) = \hat{C}h(t),$$
(1)

where h(t) is the current hidden state. h'(t) is the updated hidden state. u(t) is the current input. y(t) is the output.  $\mathbf{A} \in \mathbb{R}^{N \times N}$  is SSM's evolution matrix, and  $\mathbf{B} \in \mathbb{R}^{N \times 1}$ ,  $\hat{C} \in \mathbb{R}^{1 \times N}$  are the input and output projection matrices, respectively.

To enable the application of SSMs in sequence modeling tasks within deep learning, they must be discretized, transforming the SSM from a continuous-time functionto-function map into a discrete-time sequence-to-sequence map. S4 [12] and Mamba [11] are examples of discrete adaptations of the continuous system, incorporating a timescale parameter  $\Delta$  to convert the continuous parameters **A**, **B** into their discrete counterparts  $\bar{\mathbf{A}}$ ,  $\bar{\mathbf{B}}$ . This discretization is typically achieved using the Zero-Order Hold (ZOH) [11] method:

$$\bar{\mathbf{A}} = \exp(\Delta \mathbf{A}),$$
  

$$\bar{\mathbf{B}} = (\Delta \mathbf{A})^{-1} (\exp(\Delta \mathbf{A}) - \mathbf{I}) \cdot \Delta \mathbf{B},$$
  

$$h_t = \bar{\mathbf{A}} h_{t-1} + \bar{\mathbf{B}} u_t,$$
  

$$y_t = \mathbf{C} h_t.$$
(2)

While both S4 [12] and Mamba [11] employ a similar discretization process as outlined in Equation 2, Mamba distinguishes itself from S4 by conditioning the parameters  $\Delta \in \mathbb{R}^{b \times L \times D}$ ,  $\mathbf{B} \in \mathbb{R}^{b \times L \times N}$ , and  $\mathbf{C} \in \mathbb{R}^{b \times L \times N}$  on the input  $u \in \mathbb{R}^{b \times L \times D}$  through the S6 Selective Scan Mechanism. Here, *b* represents the batch size, *L* denotes the sequence length, and *D* signifies the feature dimension.

#### 3.2. Overall Model Architecture

DefMamba uses a common multi-scale backbone structure similar to many CNNs [25, 28] and Transformers [24]. As illustrated in Figure 2(a), the image  $I \in \mathbb{R}^{H \times W \times 3}$  is first divided into patches through a patch embedding layer, which produces a 2D feature map with spatial dimensions  $H/4 \times W/4$  and channel dimensions C. Subsequently, multiple network stages are used to create hierarchical representations of dimensions  $H/8 \times W/8 \times 2C$ ,  $H/16 \times W/16 \times 4C$ , and  $H/32 \times W/32 \times 8C$ . Each stage consists of a stack of Deformable Mamba (DM) blocks followed by a downsampling layer (except the fourth stage). Finally, the features are average-pooled and sent to the classification head to obtain the prediction results. In particular, we followed [38] and used overlapping forms in the patch embedding layer and the downsampling layers. For specific structural details, please refer to the Appendix.



Figure 3. Illustration of Deformable State Space Model. (a) illustrates the processing flow of the deformable state space model for feature extraction. (b) depicts the processing flow of the deformable scan. The upper part primarily shifts the feature points to enable the model to focus on more salient regions, while the lower part shifts the token positions to facilitate the discovery of a scanning order that is better suited to the current input. To clearly illustrate the process, only nine points are depicted in the figure, however, the actual processing involves a greater number of points. (c) presents the detailed structure of the offset network.

Different from the Mamba structure used in language models [11], the DM block follows the popular structure of the Transformer block [9, 24], which consists of two Layer Norm (LN) layers, an FFN, a DSSM (to be introduced in the following subsection) and residual connections, as shown in the Figure 2(b).

#### **3.3. Deformable State Space Model**

The overall architecture of the deformable state space model is presented in Figure 3(a). Inspired by [23, 31, 38], we employed depthwise convolution to replace the original 1D convolution in the vision mamba block [47] and incorporated a deformable branch comprising a deformable scanning and a deformable SSM (DSSM). We maintained the standard forward and backward branches to ensure stable model convergence, as our method introduces more spatial token jumps compared to previous scanning methods, potentially complicating model training and learning. Subsequent experiments validate this (Table 4).

**Deformable Scanning.** Given the issue of mutual interference between multiple deformable points, we constrained the offset of the deformable points in a certain range. This ensures the relationship between deformable points and reference points remains invariant, allowing us to approximate the relationship after deformation using the relationship before deformation. Furthermore, considering the simplicity of the computation, we employed a parallel approach to simultaneously adjust the reference points and scanning order, thereby reducing the computational burden. The specific structure is illustrated in the Figure 3(b).

Given an input feature  $x \in \mathbb{R}^{H \times W \times C}$ , where C denotes

the channel dimension and H, W represent the spatial resolution. We first generated the offset  $o \in \mathbb{R}^{H \times W \times 3}$  using a subnetwork that utilizes x to output the offset values o for the reference point and the reference token.

We initially implemented the subnetwork as depicted in Figure 3(c). The input features are first processed through a  $K \times K$  depthwise convolution to capture local features. Subsequently, GELU, Layer Normalization (LN), and a  $1 \times 1$  convolution are employed to derive the offset values, which encompass three dimensions in total. The first two dimensions represent the offset of the reference point in two-dimensional space, while the third dimension signifies the offset of the reference token index within the entire patch. During our experiments, we observed that the offset of the token necessitates global perception of the features within the patch, which cannot be achieved solely through convolution. In light of this, and considering the previous method's findings [29, 31] that mamba has redundancy in the channel dimension, we incorporated a Channel Attention (CA) mechanism [18] following the depthwise convolution layer. This mechanism mitigates channel redundancy and facilitates the integration of global contextual information. Notably, following the configuration of DAT [36], we set K to [9, 7, 5, 3] across four stages and omitted the bias in the  $1 \times 1$  convolution.

To stabilize the training process, we employed the tanhfunction to mitigate the impact of extreme values in o, where  $\hat{o} = tanh(o)$ . Subsequently, we split  $\hat{o}$  to 2 parts along the channel dimension, one with 2 channels and one with 1 channel, to obtain the point offset  $\Delta p$  and token index offset  $\Delta t$ . As previously mentioned, we needed to fur-

Method	Туре	#Param.	#FLOPs	Top-1 Acc.	Method	Туре	#Param.	#FLOPs	Top-1 Acc.
RegNetY-800M [28]	С	6M	0.8G	76.3	LocalVim-S [19]	S	28M	4.8G	81.2
GhostNet $1.3 \times [13]$	С	7M	0.2G	75.7	LocalVMamba-T [19]	S	26M	5.7G	82.7
DeiT-Ti [33]	Т	6M	1.3G	72.2	EffVMamba-B [27]	S	33M	4.6G	83.0
ViM-T [47]	S	7M	1.5G	76.1	MSVMamba-T [31]	S	33M	4.6G	82.8
EffVMamba-T [27]	S	6M	0.8G	76.5	PlainMamba-L2 [42]	S	25M	8.1G	81.6
LocalVim-T [19]	S	8M	1.5G	76.2	GrootV-T [38]	S	30M	4.8G	83.4
PlainMamba-L1 [42]	S	7M	3.0G	77.9	QuadMamba-S [39]	S	31M	5.5G	82.4
MSVMamba-N [31]	S	7M	0.9G	77.3	DefMamba-S	S	32M	4.8G	83.5
QuadMamba-Li [39]	S	5M	0.8G	74.2	ConvNeXt-S [25]	С	50M	8.7G	83.1
DefMamba-T	S	8M	1.2G	78.6	Conv2Former [17]	C	50M	8.7G	84.1
RegNetY-4G [28]	С	21M	4.0G	80.0	Swin-S [24]	T	50M	8.7G	83.0
ConvNeXt-T [25]	С	29M	4.5G	82.1	CoAtNet-1 [7]	Т	42M	8.0G	83.3
Conv2Former-T [17]	С	27M	4.4G	83.2	CrossFormer-B [35]	Т	52M	9.2G	83.4
DeiT-S [33]	Т	22M	4.6G	79.8	VMamba-S [23]	S	50M	8.7G	83.6
Swin-T [24]	Т	29M	4.5G	81.3	LocalVMamba-B [19]	S	50M	11.4G	83.7
CoAtNet-0 [7]	Т	25M	4.0G	81.6	PlainMamba-L3 [42]	S	50M	14.4G	82.3
CrossFormer-S [35]	Т	31M	4.9G	82.5	GrootV-S [38]	S	51M	8.5G	84.2
ViM-S [47]	S	26M	5.1G	81.0	QuadMamba-B [39]	S	50M	9.3G	83.8
VMamba-T [23]	S	22M	5.6G	82.2	DefMamba-B	S	51M	8.5G	84.2

Table 1. Image classification performance on ImageNet-1K validation set. C, T and S indicate the model type of CNNs, Transformer and SSM. The best results are shown in **bold** font.

ther constrain  $\Delta p$ , in order to stabilize training and simplify the structure. We divided the horizontal and vertical dimensions of  $\Delta p$  by W and H, respectively, thereby limiting the offset to the range of a single token. The detailed process is outlined as follows:

$$\Delta p, \Delta t = Split(tanh(Offset(x)), dims = [2, 1]),$$
  

$$\hat{\Delta p} = Norm(\Delta p).$$
(3)

Then, we sent the point offset  $\Delta p$ , token index offset  $\Delta t$ , and input feature x to the point offset branch and index branch, to obtain the final output, respectively.

**Point Offset**. To obtain feature representations that are more sensitive to changes in objects, we dynamically adapted the network's reference points to deformable points that contain more relevant information based on the input. Firstly, we generated reference points  $p \in \mathbb{R}^{H \times W \times 2}$ . The values in p correspond to the two-dimensional coordinates of points ranging from (0, 0) to (H - 1, W - 1). To simplify network calculations, we normalized p from its original range to [-1, 1], where [-1, 1] represents the point in the upper left corner and [1, 1] represents the point in the lower right corner. We then added the reference point p and the offset to obtain the deformable point  $\hat{p} = p + \Delta p$ . Since the offset  $\hat{p}$  contains a decimal part, it cannot be used directly. Therefore, we used bilinear interpolation to extract features at the spatial position corresponding to the offset point  $\hat{p}$ from the input x.

When our model performs point offset, the features will move in space, potentially causing the position encoding added in the initial stage to become ineffective and thereby decreasing model performance. To address this, we designed an offset bias based on the relative position encoding in the Swin Transformer. Specifically, given the feature map of size  $H \times W$ , the relative coordinate displacement of points lies within the range of [-H, H] and [-W, W] in two dimensions, respectively. Therefore, we set a learnable relative offset bias matrix  $R \in \mathbb{R}^{(2H-1)\times(2W-1)}$ . However, considering that such a matrix would result in a significant increase in parameter count, we performed a downsampling operation on this matrix to obtain  $R \in \mathbb{R}^{H \times W}$ . At the same time, we divided the point displacement by 2 to accommodate this change. Then we used the point displacement to calculate the corresponding compensation by interpolating on R. Finally, this compensation is added to the interpolated features. The specific process is as follows:

$$\hat{p} = p + \Delta p,$$
  

$$\hat{x} = \phi(x, \hat{p}) + \phi(R, \hat{p}),$$
(4)

where  $\phi(x, \hat{p})$  represents the use of bilinear interpolation function to extract the feature corresponding to position  $\hat{p}$  on x.

**Index Offset**. We modified the scanning order to enable the model to perceive the structure of the input object effectively by varying the reference token index and the deformable token index. We initially generated the reference token index  $t_r \in \mathbb{R}^{N \times 1}$ , where  $N = H \times W$ . The values in t indicate the token positions within the current patch, ranging from 0 to N - 1. To simplify network computations, we normalized  $t_r$  from its original range to [-1, 1]. This enables us to compute the deformable token index  $t_d = t_r + \Delta t$ .

Mask R-CNN 1× schedule									
Backbone	#FLOPs	APb	APb 50	AP <sup>b</sup> 75	APm	APm <sub>50</sub>	APm <sub>75</sub>		
ResNet-50 [15]	260G	38.2	58.8	41.4	34.7	55.7	37.2		
Swin-T [24]	267G	42.7	65.2	46.8	39.3	62.2	42.2		
ConvNeXt-T [25]	262G	44.2	66.6	48.3	40.1	63.3	42.8		
PVTv2-B2 [34]	309G	45.3	67.1	49.4	41.2	64.2	44.4		
VMamba-T [23]	286G	47.4	69.5	52.0	42.7	66.3	46.0		
LocalVMamba-T [19]	291G	46.7	68.7	50.8	42.2	65.7	45.5		
QuadMamba-S [39]	301G	46.7	69.0	51.3	42.4	65.9	45.6		
MSVMamba-T [31]	252G	46.9	68.8	51.4	42.2	65.6	45.4		
GrootV-T [38]	265G	47.0	69.4	51.5	42.7	66.4	46.0		
DefMamba-S	268G	47.5	69.6	51.7	42.8	66.3	46.2		

Table 2. Object detection and instance segmentation performance on MSCOCO 2017 val set. All using Mask R-CNN framework.  $AP^b$  and  $AP^m$  indicate the mean Average Precision (mAP) of detection and segmentation, respectively

Since the derived  $t_d$  contain decimal components, they cannot be utilized directly. Consequently, for  $t_d$ , we applied a sorting algorithm to determine the indices post-offset based on the magnitude of their values. Finally, we transformed the offset features  $\hat{x}$  into a 1D sequence according to these indices, thereby obtaining a content-adaptive image feature sequence. It is important to note that the sorting algorithm truncates gradients, rendering the network untrainable. To address this, we averaged the gradients of the final image sequence across the dimension and replicate them to  $\Delta t$  to approximate the gradient of the scanning order offset.

# 4. Experiments

#### 4.1. Image Classification

Settings. The image classification experiments are conducted using the ImageNet-1K [8] dataset, which comprises over 1.28 million training images and 50,000 validation images across 1,000 categories. Our training setup closely follows the methodology of previous practices [23, 31, 38], incorporating various data augmentations such as random cropping, random horizontal flipping, label-smoothing regularization, mixup, autoaugment, and random erasing. The models are trained for 300 epochs using the AdamW [26] optimizer with a cosine decay learning rate scheduler, including a 20-epoch warm-up period. The total batch size is set to 1,024, with the models trained on  $8 \times$  A800 GPUs. The optimizer parameters are configured with betas set to (0.9, 0.999), momentum set to 0.9, an initial learning rate of  $1 \times 10^{-3}$ , a weight decay of 0.05, and Exponential Moving Average (EMA).

**Results.** Table 1 presents a comparison of our proposed DefMamba models (T, S, B) with various state-of-theart (SOTA) methods. Specifically, DefMamba-T achieves 78.6% Top-1 Acc., outperforming CNNs based RegNetY-800M [28] and transformer based DeiT-Ti [33] by 2.3% and 6.4%, respectively. Moreover, DefMamba-T outperformes the recently introduced SSMs models, achieving 2.5%, 2.4%, and 1.3% higher performance than ViM-T

ADE20K with crop size 512								
Backbone	mIOU (SS)	mIOU (MS)	#Param.	#FLOPs				
ResNet-50 [15]	42.1	42.8	67M	953G				
Swin-T [24]	44.5	45.8	60M	945G				
ConvNeXt-T [25]	46.0	46.7	60M	939G				
NAT-T [14]	47.1	48.4	58M	934G				
Vim-S [47]	44.9	-	46M	-				
LocalVim-S [19]	46.4	47.5	58M	297G				
VMamba-T [23]	47.9	48.8	62M	949G				
PlainMamba-L2 [42]	46.8	-	55M	285G				
LocalVMamba-T [19]	47.9	49.1	57M	970G				
MSVMamba-T [31]	47.6	48.5	65M	942G				
QuadMamba-S [39]	47.2	48.1	62M	961G				
GrootV-T [38]	48.5	49.4	60M	941G				
DefMamba-S	48.8	49.6	65M	946G				

Table 3. Semantic segmentation performance on ADE20K val set. The crop size is all set to  $512^2$ . SS and MS denote single-scale and multi-scale testing, respectively.

[47], LocalViM-T [19], and MSVMamba-N [31] in terms of parameter and computational complexity. Moreover, it reduces the computational burden by 60% while achieving a performance improvement of 0.7% over PlainMamba-L1 [42]. DefMamba-S achieves 83.5% Top-1 Acc., surpassing GrootV-T [38] and EfficientVMamba-B [27]. Furthermore, DefMamba-B achieves an accuracy of 84.2%, exceeding VMamba-S [23] by 0.6%, demonstrating the effectiveness of our methods.

# 4.2. Object Detection

Settings. We evaluated DefMamba on the MSCOCO 2017 dataset [21] using the Mask R-CNN framework [16] for object detection and instance segmentation tasks. Following prior works [24, 31, 38], we utilized backbones pretrained on ImageNet-1K for initialization. We employed standard training strategies, including  $1 \times (12 \text{ epochs})$  with Multi-Scale (MS) training, to ensure a fair comparison.

**Results.** As depicted in Table 2, our method outperforms existing methods on most evaluation metrics. Specifically, under  $1 \times$  schedule, DefMamba-S achieves 47.5 in box mAP (AP<sup>b</sup>) and 42.8 in mask mAP (AP<sup>m</sup>). Our model surpassed ResNet-50 [15], Swin-T [24] and ConvNeXt-T [25]. At the same time, our method improved AP<sup>m</sup> by 0.6 points compared to LocalVMamba-T [19] and QuadMamba-S [39]. Furthermore, our method exhibited comparable performance to the previous SOTA method, VMamba-T [23], while reducing the computational load by 4%.

# 4.3. Semantic Segmentation

**Settings.** To evaluate the semantic segmentation performance of DefMamba, we trained our models using Uper-Net [37] initialized with pre-trained classification weights on ADE20K [45] for 160,000 iterations. We employed the AdamW optimizer [26] with a learning rate set at  $6 \times 10^{-5}$ .

Index	FB-BB	CB [42]	LB [19]	DB	#Param.	#FLOPs	Top-1
(1)	$\checkmark$				7.3M	1.1G	76.9
(2)					8.1M	1.1G	76.5
(3)	$\checkmark$				7.5M	1.1G	77.3
(4)	$\checkmark$				7.5M	1.1G	77.1
(5)	$\checkmark$				8.3M	1.2G	78.6

Table 4. Comparisons with the different scanning branches in the setting of tiny model size. DB means our deformable branch. The best results are shown in **bold** font.

Index	DP	DT	OB	CA	#Param.	#FLOPs	Top-1
(1)					7.5M	1.1G	77.0
(2)					7.5M	1.2G	77.4
(3)					7.5M	1.2G	77.2
(4)					7.5M	1.2G	77.9
(5)					8.2M	1.2G	78.2
(6)					8.3M	1.2G	78.6

Table 5. Comparisons with the different components of proposed deformable scanning in tiny model size settings on ImageNet-1k. The best results are shown in **bold** font.

Our experiments are primarily conducted using a default input resolution of 512 × 512. Additionally, we incorporated Multi-Scale (MS) testing to assess performance variations. **Results.** The DefMamba-S model demonstrates favorable performance in semantic segmentation compared to various SOTA methods, as presented in Table 3. DefMamba-S achieves a mIOU of 48.8 in single-scale and 49.6 in multi-scale evaluation. This outperforms ResNet-50 [15], Swin-T [24], and ConvNeXt-T [25]. Additionally, DefMamba-S exceeds the performance of the recent SSM methods, including GrootV-T [38], QuadMamba-S [39] and MSVMamba-T [31]. Our method achieves improvements of 0.3 points, 1.6 points, and 1.2 points, respectively, on the single-scale mIoU metric.

# 4.4. Ablation Study

Effect of the Proposed DSSM Structure. To give the evidence for the effectiveness of the proposed deformable branch, we conducted a series of experiments over different branch settings in Table 4. As shown in Table 4, FB-BB refers to the forward and the backward branches for feature extraction. CB represents continuous scanning [42] branch. LB is local scanning [19] branch and DB denotes our proposed deformable branch in Figure 2 (a). By comparing the results in Table 4(1) and (5), we observed that the proposed deformable branch significantly increased the accuracy of the ImageNet dataset by 1.7% within a reasonable computational budget, to demonstrate the effectiveness of the proposed deformable branch. Moreover, compared with Table 4 (3), (4) and (5), our method only increased the computational cost by 0.1G, while significantly improving the accuracy on ImageNet dataset by 1.4%. This further demon-



Figure 4. Visualization of activation maps in the specific position. The position is marker by red and orange point. RS stands for raster scanning, DS stands for our deformable scanning.

strates that our deformable scanning approach, as opposed to other fixed scanning methods, is more capable of capturing structural information of objects and enhancing model performance. From Table 4 (1) and (2), we noted that incorporating only the proposed deformable branch would lead to a decrease in performance due to the increased spatial token jumps. To achieve more stable training and higher model performance, we adhered to the previous paradigm by combining the FB-BB and DB in our proposed model.

Effect of the Deformable Scanning Components. To thoroughly evaluate the contributions of each component in our proposed deformable scanning method, we conducted an ablation study in Table 5. As shown in Table 5, DP, DT, OB and CA represent the components of the deformable scanning in Figure 2 (b), respectively. DP (Deformable Points) involves the process of generating deformable points, which includes initializing reference points, an offsets network and a bilinear interpolation. DT (Deformable Tokens) refers to the process of dynamically changing the token order based on predicted offsets, which includes initializing token index and deformable token index, as well as an offsets network. OB represents the operation of generating offset bias. CA denotes channel attention operation in the offset network. By comparing the results in Table 5 (1), (2) and (3), we observed that adding either the DP operation or the DT operation boosting the baseline (1) performance by 0.2-0.4%, with only a minor increases in computational costs (0.1G), demonstrated that the effectiveness of both DP and DT. Moreover, when combining both DP and DT operations (4), the performance further improves by 1% compared to the baseline, as shown in Table 5 (1) and (4). These results strongly prove that whether adding DT or DP individually, or combining both, the model can better learn the structure of images and facilitate the awareness of changes in object details, thus improving performance. Furthermore, we run ablation experiments to confirm the effectiveness of OB and CA by comparing Table 5 (4) and (5), and Table 5 (5) and (6), respectively. The experiments validate the effectiveness of our methods.



Figure 5. Visualization of deformable points and deformable token index. In (a), the orange dots represent deformable points, the green dots represent reference points, and the red arrows represent the offset path of the points. In (b) and (c), the gradient from yellow to green represents the scanning path, with the yellow dots being scanned first and the green dots being scanned later.

**Visualization Results.** To better demonstrate the superiority of our deformable scanning strategy, we present the activation maps of images at different positions for various methods in Figure 4, marked clearly by red and orange dots in the image. Specifically, we visualized the activation map corresponding to the final layer of the second stage using the method outlined in [23]. As illustrated in the activation map of the pen in the last line of Figure 4, our method demonstrates an enhanced ability to focus on the structural and shape information of objects, even when dealing with complex scenes that contain multiple overlapping objects. This capability allows for more precise recognition and segmentation, further highlighting the effectiveness of our approach in capturing essential details.

We also visualized the deformable points and deformable token index to intuitively demonstrate the performance of our method, as shown in the Figure 5. In the red box of Figure 5 (a), we can observe that some focus points outside the object shift towards the object, Such movements allow the network to attend to more object information. Compared with Figure 5 (b), our method (c) adjusts the scanning order to emphasize important tokens. For example, as shown in the first row of the Figure 5, the token corresponding to the snake's head is shifted from the middle position in a raster scan to being the first position in our method. Such offsets are beneficial for the network to learn relevant features.

# 5. Limitation

Despite the strong results of our approach, there is still a limitation. In cases where the image contains incomplete object structures or multiple objects arranged in a regular



Regularly spaced objects

Figure 6. Failure cases. The blue dots represent deformable points, while the black arrows indicate the scanning order after deformation.

pattern, the deformable scanning strategy may be less effective. As illustrated in Figure 6, when the image shows only a portion of a baseball, the deformable mechanism does not capture complete structural information, resulting in offsets that are too small and converge towards the predefined scanning method. Meanwhile, when multiple objects are arranged according to a certain rule, the information variation between adjacent tokens is minimal. Results in the model remaining in an indolent learning process.

# 6. Conclusion

In this work, we made efforts to tackle the challenge arising from existing Mamba-based methods, which rely on fixed scanning techniques to extract features. These methods do not fully leverage the spatial structure information inherent in images. To overcome this limitation, we proposed a novel foundation model named DefMamba, which is grounded in DS. This innovative model aims to enhance the capability of the network to learn and represent complex image structures, as well as to detect subtle changes in object details. The DS principally comprises two key operations: the shifting of focus points and the alteration of the scanning order. The first operation effectively repositions reference points toward significant regions of interest, thereby enhancing the model's sensitivity to variations in object details. The second operation modifies the scanning order to create a structure-aware sequence that is better aligned with the underlying object structure based on the input features. Extensive experimental evaluations on benchmark datasets, including ImageNet, COCO, and ADE20K, robustly demonstrate that our proposed method outperforms existing SSMs and remains competitive when compared to both CNNs and Transformer-based approaches.

### Acknowledgement

This work was supported by the National Natural Science Foundation of China (62172070, 62372080, 62376050, and U22B2052), and the Natural Science Foundation of Liaoning Province (2024-MSBA-24).

## References

- Qi Bi, Jingjun Yi, Hao Zheng, Wei Ji, Haolan Zhan, Yawen Huang, Yuexiang Li, and Yefeng Zheng. Samba: Severityaware recurrent modeling for cross-domain medical image grading. In Advances in Neural Information Processing Systems (NeurIPS), pages 75829–75852, 2024. 2
- [2] Qi Bi, Jingjun Yi, Haolan Zhan, Wei Ji, and Gui-Song Xia. Learning fine-grained domain generalization via hyperbolic state space hallucination. In *Proceedings of the AAAI Conference on Artificial Intelligence*, 2025.
- [3] Qi Bi, Jingjun Yi, Hao Zheng, Haolan Zhan, Wei Ji, Yawen Huang, and Yuexiang Li. Dgfamba: Learning flow factorized state space for visual domain generalization. In *Proceedings* of the AAAI Conference on Artificial Intelligence, 2025.
- [4] Keyan Chen, Bowen Chen, Chenyang Liu, Wenyuan Li, Zhengxia Zou, and Zhenwei Shi. Rsmamba: Remote sensing image classification with state space model. *IEEE Geo*science and Remote Sensing Letters, 2024. 2
- [5] Zhiyang Chen, Yousong Zhu, Chaoyang Zhao, Guosheng Hu, Wei Zeng, Jinqiao Wang, and Ming Tang. Dpt: Deformable patch-based transformer for visual recognition. In *Proceedings of the 29th ACM international conference on multimedia*, pages 2899–2907, 2021. 2
- [6] 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. 2
- [7] Zihang Dai, Hanxiao Liu, Quoc V Le, and Mingxing Tan. Coatnet: Marrying convolution and attention for all data sizes. Advances in neural information processing systems, 34:3965–3977, 2021. 5
- [8] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In 2009 IEEE conference on computer vision and pattern recognition, pages 248–255. Ieee, 2009. 2, 6
- [9] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. arXiv preprint arXiv:2010.11929, 2020. 1, 2, 4
- [10] Jiansong Fan, Tianxu Lv, Yicheng Di, Lihua Li, and Xiang Pan. Pathmamba: Weakly supervised state space model for multi-class segmentation of pathology images. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 500–509. Springer, 2024. 2
- [11] Albert Gu and Tri Dao. Mamba: Linear-time sequence modeling with selective state spaces. arXiv preprint arXiv:2312.00752, 2023. 1, 2, 3, 4

- [12] Albert Gu, Karan Goel, and Christopher Ré. Efficiently modeling long sequences with structured state spaces. arXiv preprint arXiv:2111.00396, 2021. 1, 3
- [13] Kai Han, Yunhe Wang, Qi Tian, Jianyuan Guo, Chunjing Xu, and Chang Xu. Ghostnet: More features from cheap operations. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 1580–1589, 2020. 5
- [14] Ali Hassani, Steven Walton, Jiachen Li, Shen Li, and Humphrey Shi. Neighborhood attention transformer. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 6185–6194, 2023. 6
- [15] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016. 6, 7, 1
- [16] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. Mask r-cnn. In Proceedings of the IEEE international conference on computer vision, pages 2961–2969, 2017. 6
- [17] Qibin Hou, Cheng-Ze Lu, Ming-Ming Cheng, and Jiashi Feng. Conv2former: A simple transformer-style convnet for visual recognition. *IEEE transactions on pattern analysis* and machine intelligence, 2024. 5
- [18] Jie Hu, Li Shen, and Gang Sun. Squeeze-and-excitation networks. In *Proceedings of the IEEE conference on computer* vision and pattern recognition, pages 7132–7141, 2018. 4
- [19] Tao Huang, Xiaohuan Pei, Shan You, Fei Wang, Chen Qian, and Chang Xu. Localmamba: Visual state space model with windowed selective scan. *arXiv preprint arXiv:2403.09338*, 2024. 1, 2, 5, 6, 7
- [20] Honghao Li, Yiwen Zhang, Yi Zhang, Hanwei Li, and Lei Sang. Dcnv3: Towards next generation deep cross network for ctr prediction. arXiv preprint arXiv:2407.13349, 2024. 2
- [21] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In Computer Vision–ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6-12, 2014, Proceedings, Part V 13, pages 740–755. Springer, 2014. 2, 6
- [22] Jiarun Liu, Hao Yang, Hong-Yu Zhou, Yan Xi, Lequan Yu, Cheng Li, Yong Liang, Guangming Shi, Yizhou Yu, Shaoting Zhang, et al. Swin-umamba: Mamba-based unet with imagenet-based pretraining. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 615–625. Springer, 2024. 2
- [23] Yue Liu, Yunjie Tian, Yuzhong Zhao, Hongtian Yu, Lingxi Xie, Yaowei Wang, Qixiang Ye, and Yunfan Liu. Vmamba: Visual state space model 2024. arXiv preprint arXiv:2401.10166, 2024. 1, 2, 4, 5, 6, 8
- [24] 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 Proceedings of the IEEE/CVF international conference on computer vision, pages 10012–10022, 2021. 1, 3, 4, 5, 6, 7
- [25] 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 com-

*puter vision and pattern recognition*, pages 11976–11986, 2022. 1, 3, 5, 6, 7

- [26] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *International Conference on Learning Representations*, 2017. 6
- [27] Xiaohuan Pei, Tao Huang, and Chang Xu. Efficientvmamba: Atrous selective scan for light weight visual mamba. arXiv preprint arXiv:2403.09977, 2024. 2, 5, 6
- [28] 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. 1, 3, 5, 6
- [29] Abdelrahman Shaker, Syed Talal Wasim, Salman Khan, Juergen Gall, and Fahad Shahbaz Khan. Groupmamba: Parameter-efficient and accurate group visual state space model. arXiv preprint arXiv:2407.13772, 2024. 2, 4, 1
- [30] Ruohua Shi, Qiufan Pang, Lei Ma, Lingyu Duan, Tiejun Huang, and Tingting Jiang. Shapemamba-em: Fine-tuning foundation model with local shape descriptors and mamba blocks for 3d em image segmentation. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 731–741. Springer, 2024. 2
- [31] Yuheng Shi, Minjing Dong, and Chang Xu. Multi-scale vmamba: Hierarchy in hierarchy visual state space model. *arXiv preprint arXiv:2405.14174*, 2024. 2, 4, 5, 6, 7, 1
- [32] Mingxing Tan and Quoc Le. Efficientnet: Rethinking model scaling for convolutional neural networks. In *International conference on machine learning*, pages 6105–6114. PMLR, 2019. 1
- [33] Hugo Touvron, Matthieu Cord, Matthijs Douze, Francisco Massa, Alexandre Sablayrolles, and Hervé Jégou. Training data-efficient image transformers & distillation through attention. In *International conference on machine learning*, pages 10347–10357. PMLR, 2021. 1, 5, 6
- [34] 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. 6
- [35] Wenxiao Wang, Wei Chen, Qibo Qiu, Long Chen, Boxi Wu, Binbin Lin, Xiaofei He, and Wei Liu. Crossformer++: A versatile vision transformer hinging on cross-scale attention. *IEEE Transactions on Pattern Analysis and Machine Intelli*gence, 46(5):3123–3136, 2023. 5
- [36] Zhuofan Xia, Xuran Pan, Shiji Song, Li Erran Li, and Gao Huang. Vision transformer with deformable attention. In Proceedings of the IEEE/CVF conference on computer vision and pattern recognition, pages 4794–4803, 2022. 2, 4
- [37] Tete Xiao, Yingcheng Liu, Bolei Zhou, Yuning Jiang, and Jian Sun. Unified perceptual parsing for scene understanding. In *Proceedings of the European conference on computer* vision (ECCV), pages 418–434, 2018. 6
- [38] Yicheng Xiao, Lin Song, Shaoli Huang, Jiangshan Wang, Siyu Song, Yixiao Ge, Xiu Li, and Ying Shan. Grootvl: Tree topology is all you need in state space model. *arXiv preprint arXiv:2406.02395*, 2024. 2, 3, 4, 5, 6, 7, 1

- [39] Fei Xie, Weijia Zhang, Zhongdao Wang, and Chao Ma. Quadmamba: Learning quadtree-based selective scan for visual state space model. arXiv preprint arXiv:2410.06806, 2024. 2, 5, 6, 7, 1
- [40] Zhaohu Xing, Tian Ye, Yijun Yang, Guang Liu, and Lei Zhu. Segmamba: Long-range sequential modeling mamba for 3d medical image segmentation. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 578–588. Springer, 2024. 2
- [41] Yuwen Xiong, Zhiqi Li, Yuntao Chen, Feng Wang, Xizhou Zhu, Jiapeng Luo, Wenhai Wang, Tong Lu, Hongsheng Li, Yu Qiao, et al. Efficient deformable convnets: Rethinking dynamic and sparse operator for vision applications. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5652–5661, 2024. 2
- [42] Chenhongyi Yang, Zehui Chen, Miguel Espinosa, Linus Ericsson, Zhenyu Wang, Jiaming Liu, and Elliot J Crowley. Plainmamba: Improving non-hierarchical mamba in visual recognition. arXiv preprint arXiv:2403.17695, 2024. 1, 2, 5, 6, 7
- [43] Jing Yao, Danfeng Hong, Chenyu Li, and Jocelyn Chanussot. Spectralmamba: Efficient mamba for hyperspectral image classification. *arXiv preprint arXiv:2404.08489*, 2024.
   2
- [44] Xiaoyu Yue, Shuyang Sun, Zhanghui Kuang, Meng Wei, Philip HS Torr, Wayne Zhang, and Dahua Lin. Vision transformer with progressive sampling. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 387–396, 2021. 2
- [45] Bolei Zhou, Hang Zhao, Xavier Puig, Sanja Fidler, Adela Barriuso, and Antonio Torralba. Scene parsing through ade20k dataset. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 633–641, 2017. 2, 6
- [46] Weilian Zhou, Sei-Ichiro Kamata, Haipeng Wang, Man Sing Wong, and Huiying Cynthia Hou. Mamba-in-mamba: Centralized mamba-cross-scan in tokenized mamba model for hyperspectral image classification. *Neurocomputing*, page 128751, 2024. 2
- [47] 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. 1, 2, 4, 5, 6
- [48] Xizhou Zhu, Han Hu, Stephen Lin, and Jifeng Dai. Deformable convnets v2: More deformable, better results. In Proceedings of the IEEE/CVF conference on computer vision and pattern recognition, pages 9308–9316, 2019. 2