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.

GroupMamba: Efficient Group-Based Visual State Space Model

Abdelrahman Shaker¹ Syed Talal Wasim^{2,3} Salman Khan¹ Juergen Gall^{2,3} Fahad Shahbaz Khan^{1,4}

¹Mohamed Bin Zayed University of Artificial Intelligence ²University of Bonn ³Lamarr Institute for Machine Learning and Artificial Intelligence ⁴Linköping University

Abstract

State-space models (SSMs) have recently shown promise in capturing long-range dependencies with subquadratic computational complexity, making them attractive for various applications. However, purely SSM-based models face critical challenges related to stability and achieving state-ofthe-art performance in computer vision tasks. Our paper addresses the challenges of scaling SSM-based models for computer vision, particularly the instability and inefficiency of large model sizes. We introduce a parameter-efficient modulated group mamba layer that divides the input channels into four groups and applies our proposed SSM-based efficient Visual Single Selective Scanning (VSSS) block independently to each group, with each VSSS block scanning in one of the four spatial directions. The Modulated Group Mamba layer also wraps the four VSSS blocks into a channel modulation operator to improve cross-channel communication. Furthermore, we introduce a distillationbased training objective to stabilize the training of large models, leading to consistent performance gains. Our comprehensive experiments demonstrate the merits of the proposed contributions, leading to superior performance over existing methods for image classification on ImageNet-1K, object detection, instance segmentation on MS-COCO, and semantic segmentation on ADE20K. Our tiny variant with 23M parameters achieves state-of-the-art performance with a classification top-1 accuracy of 83.3% on ImageNet-1K, while being 26% efficient in terms of parameters, compared to the best existing Mamba design of same model size. Code and models are available at: https://github.com/Amshaker/GroupMamba

1. Introduction

Various context modeling methods have emerged in the domains of language and vision understanding. These include Convolution [21, 66], Attention [60], and, more recently, State Space Models [16, 17]. Transformers with their multi-



Figure 1. Comparison in terms of Parameters vs. Top-1 Accuracy on ImageNet-1k [9]. Our GroupMamba-B achieves superior top-1 classification accuracy while reducing parameters by 36% compared to VMamba [35].

headed self-attention mechanism [60] have been central to both language models such as GPT-3 [2] and vision models such as Vision Transformers [10, 36]. However, challenges arose due to the quadratic computational complexity of attention mechanisms particularly for longer sequences, leading to the recent emergence of State Space models such as S4 [17].

While being effective in handling extended input sequences due to their linear complexity in terms of sequence lengths, S4 [17] encountered limitations in global context processing in information-dense data, especially in domains like computer vision due to the data-independent nature of the model. Alternatively, approaches such as global convolutions-based state space models [14] and Liquid S4 [20] have been proposed to mitigate the aforementioned limitations. The recent Mamba [16] introduces the S6 architecture which aims to enhance the ability of statespace models to handle long-range dependencies efficiently. The selective-scan algorithm introduced by Mamba uses input-dependent state-space parameters, which allow for better in-context learning while still being computationally efficient compared to self-attention.

However, Mamba, specifically the S6 algorithm, is known to be unstable for e.g., image classification, especially when scaled to large sizes [46]. Additionally, the Mamba model variant used in image classification, generally called the VSS (Visual State Space) block, can be more efficient in terms of parameters and compute requirements based on the number of channels. The VSS block includes extensive input and output projections along with depth-wise convolutions, whose parameters and compute complexities are directly proportional to the number of channels in the input. To address this issue, we propose a hierarchical-based *Modulated Group Mamba* layer that mitigates the aforementioned issues in a computation and parameter-efficient manner. The main contributions of our paper are:

- 1. We introduce a *Modulated Group Mamba* layer, inspired by Group Convolutions, which enhances computational efficiency and interaction in state-space models by using a multi-direction scanning method for comprehensive spatial coverage and effective modeling of local and global information.
- 2. We introduce a *Channel Affinity Modulation (CAM)* operator, which enhances communication across channels to improve feature aggregation, addressing the limited interaction inherent in the grouping operation.
- To address the instability issue in the SSM-based architecture, we introduce a distillation-based training objective designed to stabilize models with a large number of parameters, leading to better performance and a smooth loss convergence trend.
- 4. We build a series of parameter-efficient generic classification models called "GroupMamba", based on the proposed *Modulated Group Mamba* layer. Our *tiny* variant achieves 83.3% top-1 accuracy on ImageNet-1k [9] with 23M parameters and 4.5G FLOPs. Additionally, our *base* variant achieves top-1 accuracy of 84.5% with 57M parameters and 14G FLOPs, outperforming all recent SSM methods (see Fig. 1).

2. Related Work

Convolutional Neural Networks (ConvNets) have been the popular choice for computer vision tasks since the introduction of AlexNet [30]. The field has rapidly evolved with several landmark ConvNet architectures [21, 25, 52, 56, 57]. Alongside these architectural advances, significant efforts have been made to refine individual convolution layers, including depthwise convolution [65], group convolution [7], and deformable convolution [8]. Recently, ConvNeXt variants [37, 63] have taken concrete steps towards modernizing traditional 2D ConvNets by incorporating macro designs with advanced settings and training recipes to achieve onpar performance with the state-of-the-art models.

In recent years, the pioneering Vision Transformer (ViT) [10] has significantly impacted the computer vision field, including tasks such as image classification [12, 36, 38, 58], object detection [3, 44, 68, 71], and segmentation [5, 28, 51]. ViT [10] introduces a monolithic design that approaches an image as a series of flattened 2D patches without image-specific inductive bias. The remarkable performance of ViT for computer vision tasks, along with its scalability, has inspired numerous subsequent endeavors to design better architectures. The early ViT-based models usually require large-scale datasets (e.g., JFT-300M [55]) for pretraining. Later, DeiT [58] proposes advanced training techniques in addition to integrating a distillation token into the architecture, enabling effective training on smaller datasets (e.g., ImageNet-1K [9]). Since then, subsequent studies have designed hierarchical and hybrid architectures by combining CNN and ViT modules to improve performance on different vision tasks [11, 12, 41, 50, 54]. Another line of work is to mitigate the quadratic complexity inherent in self-attention, a primary bottleneck of ViTs. This effort has led to significant improvements and more efficient and approximated variants [6, 29, 43, 45, 50, 59, 62], offering reduced complexity while maintaining effectiveness.

Recently, State Space Models (SSMs) have emerged as an alternative to ViTs [60], capturing the intricate dynamics and inter-dependencies within language sequences [17]. One notable method in this area is the structured state-space sequence model (S4) [17], designed to tackle long-range dependencies while maintaining linear complexity. Following this direction, several models have been proposed, including S5 [53], H3 [13], and GSS [42]. More recently, Mamba [16] introduces an input-dependent SSM layer and leverages a parallel selective scan mechanism (S6).

In the visual domain, various works have applied SSMs to different tasks. In particular for image classification, VMamba [35] uses Mamba with bidirectional scans across both spatial dimensions in a hierarchical Swin-Transformer [36] style design to build a global receptive field efficiently. A concurrent work, Vision Mamba (Vim) [70], instead proposed a monolithic design with a single bidirectional scan for the entire image, outperforming traditional vision transformers like DeiT. LocalVMamba [27] addresses the challenge of capturing detailed local information by introducing a scanning methodology within distinct windows (inspired from Swin-Transformer [36]), coupled with dynamic scanning directions across network layers. EfficientVMamba [47] integrates atrous-based selective scanning and dual-pathway modules for efficient global and local feature extraction, achieving competitive results with reduced computational complexity. These models have been applied for image classification, as well as image segmentation [15, 34, 40, 49], video understanding [4, 31, 67], and various other tasks [18, 19, 23, 32, 61]. Their wide applicability shows the effectiveness of SSMs [13, 17, 42, 53], and in particular Mamba [16], in the visual domain. In this paper, we propose a *Modulated Group Mamba* layer that mitigates the drawbacks of the default vision Mamba block, such as lack of stability [46] and the increased number of parameters with respect to the number of channels.

3. Method

Motivation: Our method is motivated based on the observations with respect to the limitations of existing Visual State-Space models.

- Lack of Stability for Larger Models: We observe from [46] that Mamba [16] based image classification models with an MLP channel mixer are unstable when scaled to a large number of parameters. This instability can be seen in SiMBA-L (MLP) [46], which leads to suboptimal classification results of 49% accuracy. We mitigate this issue by introducing a *Modulated Group Mamba* design alongside a distillation objective (as presented in Sec. 3.4) that stabilizes the Mamba SSM training without modifying the channel mixer.
- *Efficient Improved Interaction*: Given the computational impact of Mamba-based design on the number of channels, the proposed *Modulated Group Mamba* layer is computationally inexpensive and parameter efficient than the default Mamba and able to model both local and global information from the input tokens through multi-direction scanning. An additional *Channel Affinity Modulation* operator is proposed in this work to compensate for the limited channel interaction due to the grouped operation and enhance their interactions.

3.1. Preliminaries

State-Space Models: State-space models (SSMs) like S4 [17] and Mamba [16] are structured sequence architectures inspired by a combination of recurrent neural networks (RNNs) and convolutional neural networks (CNNs), with linear or near-linear scaling in sequence length. Derived from continuous systems, SSMs define and 1D *function-to-function map* for an input $x(t) \in \mathbb{R}^L \rightarrow y(t) \in \mathbb{R}^L$ via a hidden state $h(t) \in \mathbb{R}^N$. More formally, SSMs are described by the continuous time Ordinary Differential Equation (ODE) in Eq. 1.

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

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

where h(t) is the current hidden state, h'(t) is the updated hidden state, x(t) is the current input, y(t) is the

output, $\mathbf{A} \in \mathbb{R}^{N \times N}$ is SSM's evolution matrix, and $\mathbf{B} \in \mathbb{R}^{1 \times N}$, $\mathbf{C} \in \mathbb{R}^{1 \times N}$ are the input and output projection matrices, respectively.

Discrete State-Space Models: To allow these models to be used in sequence modeling tasks in deep learning, they need to be discretized, converting the SSM from a continuous time *function-to-function map* into a discrete-time *sequence-to-sequence map*. S4 [17] and Mamba [16] are among the discrete adaptations of the continuous system, incorporating a timescale parameter Δ to convert the continuous parameters A, B into their discrete equivalents $\overline{A}, \overline{B}$. This discretization is typically done through the Zero-Order Hold (ZOH) method given in Eq. 2.

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

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

$$h_t = \overline{\mathbf{A}}h_{t-1} + \overline{\mathbf{B}}x_t,$$

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

While both S4 [17] and Mamba [16] utilize a similar discretization step as stated above in Eq. 2, Mamba differentiates 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 $x \in \mathbb{R}^{B \times L \times D}$, through the S6 Selective Scan Mechanism, where *B* is the batch size, *L* is the sequence length, and *D* is the feature dimension.

3.2. Overall Architecture

As shown in Fig. 2 (a), our model uses a hierarchical architecture, similar to Swin-Transformer [36], with four stages to efficiently process images at varying resolutions. Assuming an input image, $\mathbf{I} \in \mathbb{R}^{H \times W \times 3}$, we first apply a Patch Embedding layer to divide the image into nonoverlapping patches of size 4×4 and embed each patch into a C_1 -dimensional feature vector. The patch embedding layer is implemented using two 3×3 convolutions with a stride of 2. This produces features maps of size $\frac{H}{4} \times \frac{W}{4} \times C_1$ at the first stage. These feature maps are passed through N_1 blocks of our Modulated Grouped Mamba (as detailed in Sec. 3.3). In each subsequent stage, a down-sampling layer merges patches in a 2×2 region, followed by another N blocks of our Modulated Grouped Mamba layer. Hence, feature size at stages two, three and four are $\frac{H}{8} \times \frac{W}{8} \times C_2$, $\frac{H}{16} \times \frac{W}{16} \times C_3$, and $\frac{H}{32} \times \frac{W}{32} \times C_4$, respectively.

3.3. Modulated Group Mamba Layer

We present the overall operations of the proposed *Modulated Group Mamba* layer (Fig. 2 (b)) for an input sequence \mathbf{X}_{in} , with dimensions (B, H, W, C), where B is the batch size, C is the number of input channels and H/W are the



Figure 2. Overview of the proposed method. **Top Row:** The overall architecture of our framework with a consistent hierarchical design comprising four stages. **Bottom Row:** We present (**b**) The design of the modulated group mamba layer. The input channels are divided into four groups with a single scanning direction for each VSSS block. This significantly reduces the computational complexity compared to the standard mamba layer, with similar performance. Channel Affinity Modulation mechanism is introduced to address the limited interactions within the VSSS blocks. (**c**) The design of VSSS block. It consists of Mamba block with 1D Selective Scanning block followed by FFN. (**d**) The four scanning directions used for the four VSSS blocks are illustrated.

width and height of the feature map, in Eq. 3.

$$\begin{split} \mathbf{X}_{GM} &= \text{GroupedMamba}(\mathbf{X}_{\text{in}}, \Theta) \\ \mathbf{X}_{CAM} &= \text{CAM}(\mathbf{X}_{GM}, \text{Affinity}(\mathbf{X}_{\text{in}})) \\ \mathbf{X}_{\text{out}} &= \mathbf{X}_{\text{in}} + \text{FFN}(\text{LN}(\mathbf{X}_{\text{CAM}})) \end{split} \tag{3}$$

Here, X_{GM} is the output of Eq. 6, X_{CAM} is the output of Eq. 9, LN is the Layer Normalization [1] operation, FFN is the Feed-Forward Network as described by Eq. 5, and X_{out} is the final output of the Modulated Group Mamba block. The individual operations, namely the GroupedMamba operator, the VSSS block used inside the GroupedMamba operator, and the CAM operator, are presented in Sec. 3.3.1, Sec. 3.3.2 and Sec. 3.3.3, respectively.

3.3.1. Visual Single Selective Scan (VSSS) Block

The VSSS block (Fig. 2 (c)) is a token and channel mixer based on the Mamba operator, comprising of a Mamba block followed by a Feed-Forward Network, each with a LayerNorm before it. Mathematically, for an input token sequence \mathbf{Z}_{in} , the VSSS block performs the operations as described in Eq. 4.

$$\begin{split} \mathbf{Z}_{out}' &= \mathbf{Z}_{in} + \text{Mamba}(\text{LN}(\mathbf{Z}_{in})) \\ \mathbf{Z}_{out} &= \mathbf{Z}_{out}' + \text{FFN}(\text{LN}(\mathbf{Z}_{out}')) \end{split} \tag{4}$$

Where \mathbf{Z}_{out} is the output sequence, Mamba is the discretized Mamba SSM operator as described in Eq. 2.

$$FFN(LN(\mathbf{Z}'_{out})) = GELU(LN(\mathbf{Z}'_{out})\mathbf{W}_1 + \mathbf{b}_1)\mathbf{W}_2 + \mathbf{b}_2$$
(5)

Where GELU [24] is the activation function and W_1 , W_2 , b_1 , and b_2 are weights and biases for the linear projections.

3.3.2. Grouped Mamba Operator

Considering the motivation presented earlier in Sec. 3, we aim to design a variant of the Mamba [16] that is both computationally efficient and can effectively model the spatial dependencies of the input sequence. Given that Mamba is computationally inefficient on large number of channels Cin the input sequence, we propose a grouped variant of the operator, inspired by Grouped Convolutions. The Grouped Mamba operation is a variant of the VSSS block presented in Sec. 3.3.1, where the input channels are divided into groups, and the VSSS operator is applied separately to each group. Specifically, we divide the input channels into four groups, each of size $\frac{C}{4}$, and an independent VSSS block is applied to each group. Hence, the proposed grouped mamba operator enhances the model efficiency by splitting channels into smaller groups. To better model spatial dependencies in the input, each of the four groups scans in one of four directions across the input: left-to-right, right-to-left, bottom-to-top, and top-to-bottom as outlined in Fig. 2 (d).

Let G = 4 be the number of groups representing four scanning directions: left-to-right, right-to-left, top-tobottom, and bottom-to-top. We form four sequences from the input sequence \mathbf{X}_{in} , namely \mathbf{X}_{LR} , \mathbf{X}_{RL} , \mathbf{X}_{TB} , and \mathbf{X}_{BT} , each of shape $(B, H, W, \frac{C}{4})$, representing one of the four directions specified earlier. These are then flattened to form a single token sequence of shape $(B, N, \frac{C}{4})$, where $N = W \times H$ is the number of tokens in the sequence. The parameters for each of the four groups can be specified by θ_{LR} , θ_{RL} , θ_{TB} , and θ_{BT} , respectively, for each of the four groups, representing the parameters for the VSSS blocks.

Given the above definitions, the overall relation for the Grouped Mamba operator can be written as shown in Eq. 6.

$$\begin{split} \mathbf{X}_{GM} = & \text{GroupedMamba}(\mathbf{X}_{\text{in}}, \Theta) = \text{Concat}(\\ & \text{VSSS}(\mathbf{X}_{\text{LR}}, \Theta_{\text{LR}}), \text{VSSS}(\mathbf{X}_{\text{RL}}, \Theta_{\text{RL}}), \quad (6) \\ & \text{VSSS}(\mathbf{X}_{\text{TB}}, \Theta_{\text{TB}}), \text{VSSS}(\mathbf{X}_{\text{BT}}, \Theta_{\text{BT}})) \end{split}$$

Where:

- X_{LR}, X_{RL}, X_{TB}, and X_{BT} represent the input tensors scanned in the respective directions.
- Θ_{LR} , Θ_{RL} , Θ_{TB} , and Θ_{BT} represents the parameters of the VSSS block for each direction.
- The output of each Mamba operator is reshaped again to $(B, H, W, \frac{C}{4})$, and concatenated back to form the token sequence \mathbf{X}_{GM} , again of the size (B, H, W, C).

3.3.3. Channel Affinity Modulation (CAM)

On its own, the Grouped Mamba operator may have a disadvantage in the form of limited information exchange across channels, given the fact that each operator in the group only operates over $\frac{C}{4}$ channels. To encourage the exchange of information across channels, we propose a Channel Affinity Modulation operator, which recalibrates channel-wise feature responses to enhance the representation power of the network. In this block, we first average pool the input to calculate the channel statistics as shown in Eq. 7.

$$ChannelStat(\mathbf{X}_{in}) = AvgPool(\mathbf{X}_{in})$$
(7)

where \mathbf{X}_{in} is the input tensor, and AvgPool represents the global average pooling operation. Next comes the affinity calculation operation as shown in Eq. 8.

$$\mathsf{Affinity}(\mathbf{X}_{\mathsf{in}}) = \sigma \left(W_2 \delta \left(W_1 \mathsf{ChannelStat}(\mathbf{X}_{\mathsf{in}}) \right) \right) \quad (8)$$

where δ and σ represent non-linearity functions, and W_1 and W_2 are learnable weights. The role of σ is to assign an importance weight to each channel to compute the affinity. The result of the affinity calculation is used to recalibrate the output of the Grouped Mamba operator, as shown in Eq. 9.

$$\mathbf{X}_{\mathsf{CAM}} = \mathsf{CAM}(\mathbf{X}_{\mathsf{GM}}, \mathsf{Affinity}(\mathbf{X}_{\mathsf{in}})) = \mathbf{X}_{\mathsf{GM}} \cdot \mathsf{Affinity}(\mathbf{X}_{\mathsf{in}})$$
(9)

where \mathbf{X}_{CAM} is the recalibrated output, \mathbf{X}_{GM} is the concatenated output of the four VSSS groups from Eq. 6, \mathbf{X}_{in} is the input tensor, and Affinity(\mathbf{X}_{in}) are the channel-wise attention scores obtained from the channel affinity calculation operation in Eq. 8.

While the average pooling and affinity procedure employed by the CAM module resembles the Squeeze-and-Excitation (SE) block [26], it introduces a distinct mechanism tailored explicitly for cross-channel attention within multi-group transformations. Specifically, CAM allows inter-group information exchange to overcome the inherent limitations of the "Grouped Mamba Operator," which inherently restricts interactions within individual groups. In contrast, SE blocks typically focus on recalibrating a single feature group and have not yet been investigated within the context of Mamba-based architectures.

3.4. Distilled Loss Function

As mentioned earlier in the motivation in Sec. 3, the Mamba training is unstable when scaled to large models [46]. To mitigate this issue, we propose to utilize a distillation objective alongside the standard cross-entropy objective. Knowledge distillation involves training a student model to learn from a teacher model's behavior by minimizing a combination of the classification loss and distillation loss. The distillation loss is computed using the cross-entropy objective between the logits of the teacher and student models. Given the logits (Z_s) from the student model, logits (Z_t) from a teacher model (RegNetY-16G [48] in our case), the ground truth label y, and the hard decision of the teacher $y_t = \operatorname{argmax}_c Z_t(c)$, the joint loss function is defined as shown in Eq. 10.

$$\mathcal{L}_{\text{total}} = \alpha \mathcal{L}_{\text{CE}}(Z_s, y) + (1 - \alpha) \mathcal{L}_{\text{CE}}(Z_s, y_t).$$
(10)

where \mathcal{L}_{CE} is the cross-entropy objective and α is the weighting parameter. We demonstrate in the supplementary material that incorporating the distilled loss enhances training stability, resulting in consistent performance improvements for larger model variants.

4. Experiments

4.1. Image Classification

Settings: The image classification experiments are based on ImageNet-1K [9], which comprising of over 1.28 million training images and 50K validation images, spanning 1,000 categories. Following [38], we train our models for using the AdamW [39] optimizer and a cosine decay learning rate scheduler for 300 epochs, including a 20 epoch warm-up. The total batch size is set to 1024, with models trained on 8x A100 GPUs, each with 80GB of CUDA memory. Optimizer betas are set to (0.9, 0.999); momentum is set to 0.9, and

Method	Token mixing	Image size	#Param.	FLOPs	Top-1 acc.
RegNetY-8G [48]	Conv	224^{2}	39M	8.0G	81.7
RegNetY-16G [48]	Conv	224^{2}	84M	16.0G	82.9
EffNet-B4 [57]	Conv	380 ²	19M	4.2G	82.9
EffNet-B5 [57]	Conv	456^{2}	30M	9.9G	83.6
DeiT-S [58]	Attention	224^{2}	22M	4.6G	79.8
DeiT-B [58]	Attention	224^{2}	86M	17.5G	81.8
DeiT-B [58]	Attention	384^{2}	86M	55.4G	83.1
ConvNeXt-T [37]	Conv	224^{2}	29M	4.5G	82.1
ConvNeXt-S [37]	Conv	224^{2}	50M	8.7G	83.1
ConvNeXt-B [37]	Conv	224^{2}	89M	15.4G	83.8
Swin-T [36]	Attention	224^{2}	28M	4.6G	81.3
Swin-S [36]	Attention	224^{2}	50M	8.7G	83.0
Swin-B [36]	Attention	224^{2}	88M	15.4G	83.5
ViM-S [70]	SSM	224^{2}	26M	-	80.5
VMamba-T [35]	SSM	224^{2}	31M	4.9G	82.5
VMamba-S [35]	SSM	224^{2}	50M	8.7G	83.6
VMamba-B [35]	SSM	224^{2}	89M	15.4G	83.9
LocalVMamba-T [27]	SSM	224^{2}	26M	5.7G	82.7
LocalVMamba-S [27]	SSM	224^{2}	50M	11.4G	83.7
EfficientVMamba-B [47]	SSM	224^{2}	33M	4.0G	81.8
GroupMamba-T	SSM	224^{2}	23M	4.5G	83.3
GroupMamba-S	SSM	224^{2}	34M	7.0G	83.9
GroupMamba-B	SSM	224^{2}	57M	14G	84.5

Table 1. Performance comparison of GroupMamba models with state-of-the-art convolution-based, attention-based, and SSM-based models on ImageNet-1K [9]. Our models demonstrate better trade-off between accuracy and parameters.

an initial learning rate of 1×10^{-3} is used with a weight decay of 0.05. Label smoothing of 0.1 is used alongside the distillation objective (see Sec. 3.4).

Tab. 1 presents a comparison of our pro-**Results:** posed GroupMamba models (T, S, B) with various stateof-the-art methods. The GroupMamba models exhibit a notable balance of accuracy and computational efficiency. GroupMamba-T achieves a top-1 accuracy of 83.3% with 23 million parameters and 4.5 GFLOPs, outperforming ConvNeXt-T [37] and Swin-T [36] by 1.2% and 2.0%, respectively, with fewer parameters. Additionally, GroupMamba-T surpasses the recently introduced SSM models, outperforming VMamba-T [35] and LocalVMamba-T [27] by 0.8% and 0.6%, respectively, while using 26% fewer parameters than VMamba-T. GroupMamba-S, with 34 million parameters and 7.0 GFLOPs, achieves an accuracy of 83.9%, surpassing VMamba-S [35], Swin-S [36], and EfficientVMamba-B [47]. The performance is better than LocalVMamba-S [27] by 0.2% with 32% fewer parameters. Furthermore, GroupMamba-B achieves an accuracy of 84.5% with only 57 million parameters and 14 GFLOPs, exceeding VMamba-B [35] by 0.6% while using 36% fewer parameters.

4.2. Object Detection and Instance Segmentation

Settings: We evaluate the performance of GroupMamba-T for object detection on the MS-COCO 2017 dataset [33]. Our method is based on the Mask R-CNN 1× schedule [22] detector with the hyperparameters as used for Swin [36]. We use the AdamW [39] optimizer and train Mask-RCNN with GroupMamba-T backbone for 12 epochs. The backbone is initialized and fine-tuned from the ImageNet-1K [9]. We use an initial learning rate of 1×10^{-4} and decay by a factor of 10 at epochs 9 and 11. FLOPs are computed for an input dimension of 1280×800 .

Results: Tab. 2 shows the results of GroupMamba-T, comparing it against various state-of-the-art models for object detection and instance segmentation using the Mask R-CNN framework on the MS-COCO dataset. Our model

Detection & Instance Segmentation								Semantic Segmentation		
Backbone	AP ^b	AP_{50}^{b}	AP^b_{75}	AP ^m	AP_{50}^m	AP^m_{75}	#param.	FLOPs	mIoU (SS)	mIoU (MS)
ResNet-50 [21]	38.2	58.8	41.4	34.7	55.7	37.2	44M	260G	42.1	42.8
Swin-T [36]	42.7	65.2	46.8	39.3	62.2	42.2	48M	267G	44.4	45.8
ConvNeXt-T [37]	44.2	66.6	48.3	40.1	63.3	42.8	48M	262G	46.0	46.7
VMamba-T [35]	47.4	69.5	52.0	42.7	66.3	46.0	50M	270G	48.3	48.6
LocalVMamba-T [27]	46.7	68.7	50.8	42.2	65.7	45.5	45M	291G	47.9	49.1
GroupMamba-T	47.6	69.8	52.1	42.9	66.5	46.3	40M	279G	48.6	49.2

Table 2. Comparison of model performance on dense prediction tasks: Object detection and instance segmentation results on MS-COCO [33] using Mask R-CNN 1× schedule [22], and semantic segmentation results on ADE20K [69] using UperNet [64]. 'SS' and 'MS' denote single-scale and multi-scale evaluations, respectively. AP^b and AP^m represent box and mask AP.

achieves box AP (AP^b) of 47.6 and mask AP (AP^m) of 42.9. It surpasses ResNet-50 [21], Swin-T [38], ConvNeXt-T [37]. In addition, GroupMamba-T has competitive performance compared to VMamba-T [35] and LocalVMamba-T [27], with less 20% parameters compared to VMamba-T. Fig. 3 (first row) displays qualitative examples of object detection and instance segmentation. GroupMamba-T accurately detects and segments the targets in various scenes. More qualitative examples are presented in the supplementary material.

4.3. Semantic Segmentation

Settings: We also evaluate the performance of GroupMamba-T for semantic segmentation on the ADE20K [69] dataset. The framework is based on the UperNet [64] architecture, and we follow the same hyperparameters as used for the Swin [36] backbone. More specifically, we use the AdamW [39] optimizer for a total of 160k iterations with an initial learning rate of 6×10^{-5} . The default resolution in our experiments is 512×512 .

Results: The GroupMamba-T model demonstrates favorable performance in semantic segmentation compared to various state-of-the-art methods, as presented in Tab. 2. GroupMamba-T achieves a mIoU of 48.6 in single-scale and 49.2 in multi-scale evaluation. This outperforms ResNet-50 [21], Swin-T [36], and ConvNeXt-T [37]. Additionally, GroupMamba-T exceeds the performance of the recent SSM methods, including ViM-S [70], VMamba-T [35], and LocalVMamba [27]. Fig. 3 (second row) shows qualitative examples of GroupMamba-T. These examples demonstrate our model's ability to accurately segment various classes for indoor and outdoor scenes. More qualitative examples are presented in the supplementary material.

4.4. Ablation Study

Fig. 4 shows the impact of each proposed contribution in terms of top-1 accuracy, number of parameters, and throughput, compared to other SSM-based methods. GroupMamba-T with 4-D scanning, comprising 22M parameters, achieves a top-1 accuracy of 82.30% and a throughput of 803. By applying a unidirectional 1D scan across N/4 channels in four directions—left-to-right, right-to-left, top-to-bottom, and bottom-to-top instead of the full 4-D scanning across all N channels, the throughput significantly increased from 803 to 1125, with only 0.1% drop in accuracy and the same number of parameters.

The integration of the CAM module further elevates the top-1 accuracy from 82.20% to 82.50%, with a minor reduction in throughput (from 1125 to 1069). Finally, incorporating the proposed distillation-based loss pushes the top-1 accuracy to 83.30%, while preserving the throughput at 1069. Compared to Vim-S [70], GroupMamba-T demonstrates a more efficient design, achieving a 2.8% improvement in top-1 accuracy with $1.5 \times$ higher throughput, all while utilizing fewer parameters. Compared to LocalVMamba-T [27], GroupMamba-T has a 0.6% higher accuracy in top-1 accuracy, with $3 \times$ faster and smaller number of parameters. Regarding VMamba-T V1 [35], our model achieves a 1.1% gain in top-1 accuracy with a comparable number of parameters while being faster by $2.5 \times$. Likewise, when compared to VMamba-T V2 [35], GroupMamba-T shows marginally faster throughput, an increase of 0.8% in top-1 accuracy, and a 26% improvement in parameter efficiency.

5. Conclusion and Future Work

In this paper, we tackle the computational inefficiencies and stability challenges associated with visual SSMs for computer vision tasks by introducing a novel layer called *Modulated Group Mamba*. We also propose a multi-directional scanning method that improves parameter efficiency by scanning in four spatial directions and leveraging *Channel Affinity Modulation* (CAM) operator to enhance feature aggregation across channels. To stabilize training, especially for larger models, we employ a distillation-based training objective. Our experiments demonstrate that the proposed GroupMamba models outperform recent SSMs while being more efficient in terms of parameters and throughput.



Figure 3. Qualitative results of GroupMamba-T for object detection and instance segmentation (first row) on the MS-COCO val. set and semantic segmentation (second row) on ADE20k val. set.



Figure 4. Comparison of GroupMamba variants and SSM-based methods in top-1 accuracy on ImageNet-1k [9] and computational efficiency in terms of throughput and number of parameters. The throughput (number of predicted samples per second) is measured using a single NVIDIA A100 GPU with a batch size of 128 for all methods.

Our research has focused on image classification, object detection, and segmentation. To further validate and extend the generalization ability of our method, we aim to explore additional downstream tasks, such as video recognition and time-series data applications. Evaluating the Modulated Group Mamba layer in these contexts will help to uncover its potential benefits and limitations, providing deeper insights and guiding further improvements.

6. Acknowledgments

Syed Talal Wasim and Juergen Gall have been supported by the Federal Ministry of Education and Research (BMBF) under grant no. 01IS22094A WEST-AI and the ERC Consolidator Grant FORHUE (101044724).

References

- [1] Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E. Hinton. Layer normalization. *arxiv preprint*, *arXiv:1607.06450*, 2016. 4
- [2] Tom Brown, Benjamin Mann, Nick Ryder, et al. Language models are few-shot learners. In *NeurIPS*, 2020.
 1
- [3] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. End-to-end object detection with transformers. In *ECCV*, 2020. 2
- [4] Guo Chen, Yifei Huang, Jilan Xu, Baoqi Pei, Zhe Chen, Zhiqi Li, Jiahao Wang, Kunchang Li, Tong Lu, and Limin Wang. Video mamba suite: State space model as a versatile alternative for video understanding. arxiv preprint, arXiv:2403.09626, 2024. 3
- [5] Bowen Cheng, Ishan Misra, Alexander G. Schwing, Alexander Kirillov, and Rohit Girdhar. Maskedattention mask transformer for universal image segmentation. In CVPR, 2022. 2
- [6] Xiangxiang Chu et al. Twins: Revisiting the design of spatial attention in vision transformers. In *NIPS*, 2021.
- [7] Taco Cohen and Max Welling. Group equivariant convolutional networks. In *ICML*, 2016. 2
- [8] Jifeng Dai, Haozhi Qi, Yuwen Xiong, Yi Li, Guodong Zhang, Han Hu, and Yichen Wei. Deformable convolutional networks. In *ICCV*, 2017. 2
- [9] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. ImageNet: A large-scale hierarchical image database. In *CVPR*, 2009. 1, 2, 5, 6, 8
- [10] Alexey Dosovitskiy et al. An image is worth 16x16 words: Transformers for image recognition at scale. In *ICLR*, 2021. 1, 2
- [11] Stéphane d'Ascoli, Hugo Touvron, Matthew L Leavitt, Ari S Morcos, Giulio Biroli, and Levent Sagun. Convit: Improving vision transformers with soft convolutional inductive biases. In *ICML*, 2021. 2
- [12] Haoqi Fan et al. Multiscale vision transformers. In ICCV, 2021. 2
- [13] Daniel Y Fu, Tri Dao, Khaled K Saab, Armin W Thomas, Atri Rudra, and Christopher Ré. Hungry hungry hippos: Towards language modeling with state space models. In *ICLR*, 2023. 2, 3

- [14] Daniel Y. Fu, Hermann Kumbong, Eric Nguyen, and Christopher Ré. FlashFFTConv: Efficient convolutions for long sequences with tensor cores. arXiv preprint, arXiv:2311.05908, 2023. 1
- [15] Haifan Gong, Luoyao Kang, Yitao Wang, Xiang Wan, and Haofeng Li. nnmamba: 3d biomedical image segmentation, classification and landmark detection with state space model. *arxiv preprint, arXiv:2402.03526*, 2024. 3
- [16] Albert Gu and Tri Dao. Mamba: Linear-time sequence modeling with selective state spaces. *arxiv preprint*, *arXiv:2312.00752*, 2023. 1, 2, 3, 4
- [17] Albert Gu, Karan Goel, and Christopher Ré. Efficiently modeling long sequences with structured state spaces. In *ICLR*, 2022. 1, 2, 3
- [18] Hang Guo, Jinmin Li, Tao Dai, Zhihao Ouyang, Xudong Ren, and Shu-Tao Xia. Mambair: A simple baseline for image restoration with state-space model. *arxiv preprint, arXiv:2402.15648*, 2024. 3
- [19] Tao Guo, Yinuo Wang, and Cai Meng. Mambamorph: a mamba-based backbone with contrastive feature learning for deformable mr-ct registration. *arxiv* preprint, arXiv:2401.13934, 2024. 3
- [20] Ramin Hasani, Mathias Lechner, Tsun-Huang Wang, Makram Chahine, Alexander Amini, and Daniela Rus. Liquid structural state-space models. *arXiv preprint*, *arXiv*:2209.12951, 2022. 1
- [21] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, 2016. 1, 2, 7
- [22] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. Mask r-cnn. In *ICCV*, 2017. 6, 7
- [23] Xuanhua He, Ke Cao, Keyu Yan, Rui Li, Chengjun Xie, Jie Zhang, and Man Zhou. Pan-mamba: Effective pan-sharpening with state space model. *arxiv preprint*, *arXiv*:2402.12192, 2024. 3
- [24] Dan Hendrycks and Kevin Gimpel. Gaussian error linear units (gelus). *arxiv preprint, arXiv:1606.08415*, 2016. 4
- [25] Andrew G. Howard et al. MobileNets: Efficient convolutional neural networks for mobile vision applications. *arxiv preprint, arXiv:1704.04861*, 2017. 2
- [26] Jie Hu, Li Shen, and Gang Sun. Squeeze-andexcitation networks. In *CVPR*, 2018. 5
- [27] 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. 2, 6, 7
- [28] Alexander Kirillov et al. Segment anything. In *ICCV*, 2023. 2
- [29] Nikita Kitaev, Lukasz Kaiser, and Anselm Levskaya. Reformer: The efficient transformer. In *ICML*, 2020.2

- [30] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E. Hinton. Imagenet classification with deep convolutional neural networks. In *NeurIPS*, 2012. 2
- [31] Kunchang Li, Xinhao Li, Yi Wang, Yinan He, Yali Wang, Limin Wang, and Yu Qiao. Videomamba: State space model for efficient video understanding. *arxiv* preprint, arXiv:2403.06977, 2024. 3
- [32] Dingkang Liang, Xin Zhou, Xinyu Wang, Xingkui Zhu, Wei Xu, Zhikang Zou, Xiaoqing Ye, and Xiang Bai. Pointmamba: A simple state space model for point cloud analysis. arxiv preprint, arXiv:2402.10739, 2024. 3
- [33] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, C. Lawrence Zitnick, and Piotr Dollár. Microsoft coco: Common objects in context. In *ECCV*, 2014. 6, 7, 13
- [34] Jiarun Liu, , et al. Swin-umamba: Mamba-based unet with imagenet-based pretraining. *arxiv preprint*, *arXiv:2402.03302*, 2024. 3
- [35] Yue Liu, Yunjie Tian, Yuzhong Zhao, Hongtian Yu, Lingxi Xie, Yaowei Wang, Qixiang Ye, and Yunfan Liu. Vmamba: Visual state space model. arxiv preprint, arXiv:2401.10166, 2024. 1, 2, 6, 7
- [36] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. Swin Transformer: Hierarchical vision transformer using shifted windows. In *ICCV*, 2021. 1, 2, 3, 6, 7
- [37] Zhuang Liu, Hanzi Mao, Chao-Yuan Wu, Christoph Feichtenhofer, Trevor Darrell, and Saining Xie. A convnet for the 2020s. In *CVPR*, 2022. 2, 6, 7
- [38] Ze Liu et al. Swin Transformer V2: Scaling up capacity and resolution. In *CVPR*, 2022. 2, 5, 7
- [39] Ilya Loshchilov and Frank Hutter. Fixing weight decay regularization in adam. *arxiv preprint, arXiv:1711.05101*, 2017. 5, 6, 7
- [40] Jun Ma, Feifei Li, and Bo Wang. U-mamba: Enhancing long-range dependency for biomedical image segmentation. *arxiv preprint, arXiv:2401.04722*, 2024.
 3
- [41] Muhammad Maaz et al. Edgenext: Efficiently amalgamated cnn-transformer architecture for mobile vision applications. In *International Workshop on Computational Aspects of Deep Learning at 17th European Conference on Computer Vision (CADL2022)*, 2022.
- [42] Harsh Mehta, Ankit Gupta, Ashok Cutkosky, and Behnam Neyshabur. Long range language modeling via gated state spaces. *arxiv preprint*, *arXiv*:2206.13947, 2022. 2, 3
- [43] Sachin Mehta and Mohammad Rastegari. Separable self-attention for mobile vision transformers. *Transactions on Machine Learning Research*, 2023. 2

- [44] Depu Meng, Xiaokang Chen, Zejia Fan, Gang Zeng, Houqiang Li, Yuhui Yuan, Lei Sun, and Jingdong Wang. Conditional detr for fast training convergence. In *ICCV*, 2021. 2
- [45] Junting Pan et al. Edgevits: Competing light-weight cnns on mobile devices with vision transformers. In ECCV, 2022. 2
- [46] Badri N. Patro and Vijay S. Agneeswaran. Simba: Simplified mamba-based architecture for vision and multivariate time series. *arxiv preprint*, *arXiv*:2403.15360, 2024. 2, 3, 5
- [47] Xiaohuan Pei, Tao Huang, and Chang Xu. Efficientvmamba: Atrous selective scan for light weight visual mamba. arxiv preprint, arXiv:2403.09977, 2024. 2, 6
- [48] I. Radosavovic, R. Kosaraju, R. Girshick, K. He, and P. Dollar. Designing network design spaces. In *CVPR*, 2020. 5, 6
- [49] Jiacheng Ruan and Suncheng Xiang. Vm-unet: Vision mamba unit for medical image segmentation. *arxiv preprint, arXiv:*, 2024. 3
- [50] Abdelrahman Shaker et al. Swiftformer: Efficient additive attention for transformer-based real-time mobile vision applications. In *ICCV*, 2023. 2
- [51] Abdelrahman Shaker et al. Efficient video object segmentation via modulated cross-attention memory. *arXiv:2403.17937*, 2024. 2
- [52] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *ICLR*, 2015. 2
- [53] Jimmy TH Smith, Andrew Warrington, and Scott W Linderman. Simplified state space layers for sequence modeling. In *ICLR*, 2023. 2, 3
- [54] Aravind Srinivas, Tsung-Yi Lin, Niki Parmar, Jonathon Shlens, Pieter Abbeel, and Ashish Vaswani. Bottleneck transformers for visual recognition. In *CVPR*, 2021. 2
- [55] Chen Sun, Abhinav Shrivastava, Saurabh Singh, and Abhinav Gupta. Revisiting unreasonable effectiveness of data in deep learning era. In *ICCV*, 2017. 2
- [56] Christian Szegedy et al. Going deeper with convolutions. In *CVPR*, 2015. 2
- [57] Mingxing Tan and Quoc V. Le. EfficientNet: Rethinking model scaling for convolutional neural networks. In *ICML*, 2019. 2, 6
- [58] Hugo Touvron, Matthieu Cord, Matthijs Douze, Francisco Massa, Alexandre Sablayrolles, and Herve Jegou. Training data-efficient image transformers & distillation through attention. In *ICML*, 2021. 2, 6
- [59] Zhengzhong Tu, Hossein Talebi, Han Zhang, Feng Yang, Peyman Milanfar, Alan Bovik, and Yinxiao Li. Maxvit: Multi-axis vision transformer. In ECCV, 2022. 2

- [60] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *NeurIPS*, 2017. 1, 2
- [61] Chloe Wang, Oleksii Tsepa, Jun Ma, and Bo Wang. Graph-mamba: Towards long-range graph sequence modeling with selective state spaces. *arxiv preprint*, *arXiv:2402.00789*, 2024. 3
- [62] Sinong Wang, Belinda Z. Li, Madian Khabsa, Han Fang, and Hao Ma. Linformer: Self-attention with linear complexity. arXiv preprint arXiv:2006.04768, 2020. 2
- [63] Sanghyun Woo et al. Convnext v2: Co-designing and scaling convnets with masked autoencoders. In CVPR, 2023. 2
- [64] Tete Xiao, Yingcheng Liu, Bolei Zhou, Yuning Jiang, and Jian Sun. Unified perceptual parsing for scene understanding. In ECCV, 2018. 7
- [65] Saining Xie, Ross Girshick, Piotr Dollár, Zhuowen Tu, and Kaiming He. Aggregated residual transformations for deep neural networks. In *CVPR*, 2017. 2
- [66] Jianwei Yang, Chunyuan Li, Xiyang Dai, Lu Yuan, and Jianfeng Gao. Focal modulation networks. In *NeurIPS*, 2022. 1
- [67] Yijun Yang, Zhaohu Xing, and Lei Zhu. Vivim: a video vision mamba for medical video object segmentation. *arxiv preprint, arXiv:2401.14168*, 2024. 3
- [68] Hao Zhang et al. Dino: Detr with improved denoising anchor boxes for end-to-end object detection. In *ICLR*, 2022. 2
- [69] Bolei Zhou, Hang Zhao, Xavier Puig, Sanja Fidler, Adela Barriuso, and Antonio Torralba. Scene parsing through ade20k dataset. In CVPR, 2017. 7, 13
- [70] 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. 2, 6, 7
- [71] Xizhou Zhu, Weijie Su, Lewei Lu, Bin Li, Xiaogang Wang, and Jifeng Dai. Deformable detr: Deformable transformers for end-to-end object detection. In *ICLR*, 2021. 2