756 A SUPPLEMENTARY MATERIAL

This section contains supplemental material, offering further results and analysis to complement the main paper. We provide additional details on the following topics:

- Architectural Details (Section A.1)
- Ablations (Section A.2)
 - Qualitative Results (Section A.3)
- 765
- 766

773

764

761

- Discussion (Section A.4)
- 767 A.1 ARCHITECTURAL DETAILS 768

We develop three variants of our GroupMamba backbones, each tailored to different performance
and efficiency requirements: GroupMamba-T (Tiny), GroupMamba-S (Small), and GroupMamba-B
(Base). These variants differ in their channel dimensions and the number of layers per stage, as
detailed in Table 4.

Table 4: GroupMamba Architectures. Description of the configurations of the model variants for
the embedding size, the number of layers, and the model's GFLOPs and Parameters. Between two
consecutive stages, we incorporate a downsampling layer to increase the number of channels and
reduce the resolution by two.

778	Stage	Output Resolution	Type	Config	GroupMamba		
79 780	Stuge		-540		T	S B	
781		$\underline{H} \times \underline{W}$	Patch Embedding	Patch Size	k=3x3, s=2		
782	Stem	$2 \land 2$		Embed. Dim.	32	32 48	
783 794		$\frac{H}{4} \times \frac{W}{4}$	Patch Embedding	Patch Size	k=3x3, s=2		
785				Embed. Dim.	64	64 96	
786	1	$\frac{H}{4} \times \frac{W}{4}$	Modulated Group Mamba	Layers	3	3 3	
787		$\underline{H} \times \underline{W}$	Down-Sampling	Patch Size	k=3x3, s=2		
′88 789		8 ^ 8	Down Sampling	Embed. Dim.	128	128 192	
'90	2	$\frac{H}{8} \times \frac{W}{8}$	Modulated Group Mamba	Layers	4	4 6	
'91 '22	$H \times W$		Down-Sampling	Patch Size	k=3x3, s=2		
92 93		16 ^ 16	Down bumphing	Embed. Dim.	320	320 384	
'94	3	$\frac{H}{16} \times \frac{W}{16}$	Modulated Group Mamba	Layers	6	12 18	
95	$H \searrow W$		Down-Sampling	Patch Size	k=3x3, s=2		
96 97		$\overline{32} \wedge \overline{32}$	Down Sampling	Embed. Dim.	448	512 512	
98	4	$\frac{H}{32} \times \frac{W}{32}$	Modulated Group Mamba	Layers	3	3 3	
99 00	Parameters FLOPs				23M 4.5G	34M 57M 7.0G 14.0G	

801

802 803 804

A.2 ABLATIONS

In Table 5, we provide additional ablation results regarding the distillation training objective. For
the GroupMamba-T and GroupMamba-S variants, the distilled loss improves performance by an
absolute gain of 0.8% and 0.9%, respectively. For the largest variant, GroupMamba-B, the distilled
loss improves performance by 1.3%. This demonstrates that larger Mamba-based models with MLP
tend to saturate and struggle to converge effectively without distillation. Incorporating distillation for
the large model boosts its performance from 83.2% to 84.5%.

811	Table 5: Ablation study on GroupMamba variants with and without the Distilled Lo						
812		Method	#Param.	FLOPs	Top-1 acc.		
813 814	Gr	oupMamba-T w/o Distilled Loss oupMamba-T with Distilled Loss	23M 23M	4.6G 4.6G	82.5 83.3 (+0.8)		
815 816 817	Gi	roupMamba-S w/o Distilled Loss oupMamba-S with Distilled Loss	34M 34M	7.0G 7.0G	83.0 83.9 (+0.9)		
818 819	Gi Gr	oupMamba-B w/o Distilled Loss oupMamba-B with Distilled Loss	57M 57M	14G 14G	83.2 84.5 (+1.3)		
820							

. **D** : SS.

QUALITATIVE RESULTS A.3

823 In Figure 5, we present the qualitative results of GroupMamba-T on samples from the COCO 824 validation set Lin et al. (2014), demonstrating its performance in instance segmentation and object 825 detection. Our model accurately localizes objects and correctly segments them across diverse scenes 826 and varying scales. In Figure 6, we show additional qualitative results of GroupMamba-T on samples 827 from the ADE20K Zhou et al. (2017) validation set for semantic segmentation. The first row shows 828 the ground truth masks, while the second row displays the predicted masks. It is notable that our 829 model delineates the masks accurately, highlighting the effectiveness for semantic segmentation. The 830 quantitative and qualitative results of GroupMamba demonstrate the robust generalization capability 831 of our GroupMamba backbones across diverse downstream tasks, including semantic segmentation, object detection, and instance segmentation. 832

A.4 DISCUSSION

835 Our main contributions include introducing the Modulated Group Mamba layer, which enhances 836 computational efficiency and interaction in state-space models through a multi-direction scanning 837 method. We also introduce the Channel Affinity Modulation (CAM) operator to improve feature 838 aggregation across channels, addressing limitations in grouping operations. Additionally, we employ 839 a distillation-based training objective to stabilize the training of models with a large number of 840 parameters. These contributions enable us to achieve competitive performance with recent state-space 841 models in image classification, object detection, instance segmentation, and semantic segmentation 842 with fewer number of parameters.

843 This can further facilitate the development of vision foundation models based on Mamba that can be 844 scaled to a large number of parameters efficiently and stably. The Modulated Group Mamba layer and 845 CAM operator enhance computational efficiency and feature interaction, allowing models to manage 846 more extensive and complex datasets without excessive resource demands. The distillation-based 847 training objective ensures stability during training, which is crucial for maintaining performance as 848 model sizes increase. Together, these advancements enable the creation of scalable, reliable vision 849 models that can be deployed effectively in various real-world applications.

850

810

821

822

833

834

- 851
- 852
- 853
- 854 855
- 856
- 857
- 858
- 859
- 861

862

863



Figure 5: Qualitative results of GroupMamba-T for object detection and instance segmentation on the COCO validation set.



Figure 6: Qualitative results of GroupMamba-T for semantic segmentation on ADE20K validation set. The first row shows the ground truth for the masks, while the second and second show the corresponding predictions of our model.