

# Autoregressive Pretraining with Mamba in Vision



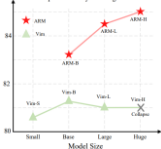
Sucheng Ren<sup>1</sup> Xianhang Li<sup>2</sup> Haoqin Tu<sup>2</sup> Feng Wang<sup>1</sup> Fangxun Shu<sup>3</sup> Lei Zhang<sup>3</sup>

Jieru Mei<sup>1</sup> Linjie Yang<sup>4</sup> Peng Wang<sup>4</sup> Heng Wang<sup>4</sup> Alan Yuille<sup>1</sup> Cihang Xie<sup>1</sup>

<sup>1</sup> Johns Hopkins University <sup>2</sup>UC Santa Cruz <sup>3</sup>Alibaba Group <sup>4</sup>ByteDance

## Scaling Mamba in Vision

Top-1 Accuracy on ImageNet



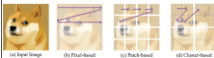
(1) attempts to scale the Vision Mamba (Vim) under supervised conditions often lead to either performance plateauing or even training collapse when pushed to very large sizes

(2) we primarily focus on the autoregressive pretraining paradigm for self-supervised visual representation learning

## Autoregressive Pretraining with Mamba in Vision

We propose ARM and study various designs in autoregressive modeling:

### Prediction Unit



we propose grouping spatially adjacent patches into larger clusters to serve as the prediction unit

### Order



Unlike the 1D sentences in NLP, which inherently have a clear sequence order. We explore four primary prediction orders:

- 1) Row-first and forward orders: the clusters row by row processing from the first to the last cluster within each row sequentially
- 2) Row-first and backward orders: the clusters row by row but inverts the processing direction.
- 3) Column-first and forward organizes the clusters column by column, processing sequentially within each column from top to bottom
- 4) Column-first and backward similarly sequences the clusters column by column but inverts the processing direction

## Experiments

We conduct experiments on ImageNet-1K classification.

Model	Token Mamba	Image Size	Params (M)	Throughput (img/s)	Top-1 (%)
<b>Base-size models</b>					
RegNet-Y160	28x28	224 <sup>2</sup>	34	870	82.0
ARM-B	Autoregressive	224 <sup>2</sup>	31	973	82.1
Vim-B	Mamba	224 <sup>2</sup>	36	950	81.2
MambaViT-B	Mamba	224 <sup>2</sup>	31	1360	82.2
VisionGPT-B	Mamba	224 <sup>2</sup>	31	131	82.4
<b>Large-size models</b>					
ARM-L	Mamba	256 <sup>2</sup>	101	1380	82.2
ARM-B	Mamba	256 <sup>2</sup>	91	1440	82.1
ARM-L	Mamba	448 <sup>2</sup>	351	56	82.3
<b>Large-size models</b>					
Vim-L	Mamba	224 <sup>2</sup>	148	340	81.0
MambaViT-L	Mamba	224 <sup>2</sup>	207	440	81.4
ARM-L	Mamba	224 <sup>2</sup>	207	840	82.5
ARM-L	Mamba	384 <sup>2</sup>	207	754	82.1
<b>Large-size models</b>					
Vim-B	Mamba	224 <sup>2</sup>	700	210	collapsed
ARM-B	Mamba	224 <sup>2</sup>	662	270	82.0
ARM-B	Mamba	384 <sup>2</sup>	662	64	82.1

(1) We are the first scaling Mamba to huge size.

## Ablation Study

Num of Prediction unit	Cluster size	Top-1 (%)
0 (Supervised)	N/A	81.2
144 (GPT)	1 × 1 (Pixel)	79.8
4	96 × 96	82.0
9	64 × 64	82.5
16	48 × 48	82.2
36	32 × 32	81.9
144	16 × 16	81.7

(1) Number of prediction units: the number of the prediction units set to 9.

Order	Direction	Top-1 (%)
Row-first	Forward	82.5
Row-first	Backward	82.3
Column-first	Forward	82.5
Column-first	Backward	82.4
Random	Random	81.5

(2) Prediction order: no significant difference with predefined order.