

Enabling True Global Perception in State Space Models for Visual Tasks

Jie Hui, Zhenxiang Zhang, Wenyu Mi, Jianji Wang



西安交通大学
XI'AN JIAOTONG UNIVERSITY

IAIR Est.
1986

Institute of
Artificial Intelligence
and Robotics



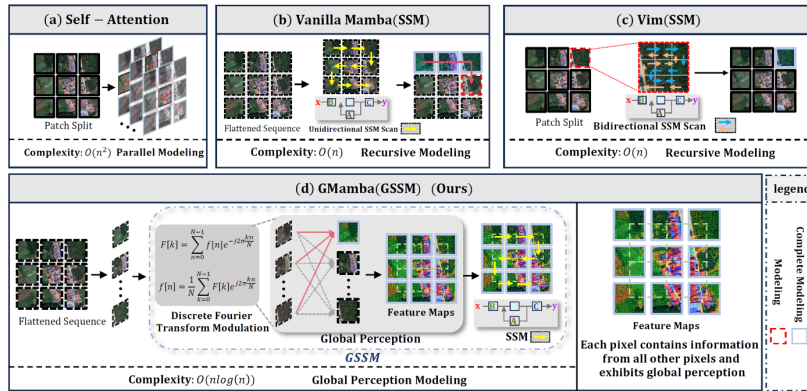
ICLR2026



Motivation

Global modeling is crucial for capturing long-range dependencies in visual recognition and semantic understanding.

- Existing global modeling approaches (Transformer-based and Mamba-based) are mainly evaluated through empirical evidence (e.g., ablations or visualizations), while a **rigorous definition of true global perception** is still missing.
- Transformer-based methods enable explicit global interactions, but their **quadratic computational complexity** limits efficiency in high-resolution vision tasks.
- Mamba/SSM-based methods achieve linear complexity, yet rely on sequential state propagation and causal assumptions, **leading to localized perception and structural mismatch with non-causal image data**. Although sophisticated scanning strategies can partially alleviate this issue, they do not fundamentally overcome **the intrinsic limitations of SSM-based modeling** and often **introduce additional complexity**.



Main ideas

- First formal definition of global image modeling**, providing a rigorous mathematical foundation for analyzing and guaranteeing true global perception in visual tasks.
- Theoretical framework for frequency-modulated SSM**, connecting frequency domain analysis with state-space modeling to guide global feature propagation.
- Design of a plug-and-play GMamba (GSSM) module** that injects frequency-guided global semantics into SSM, adaptable to diverse CNN backbones and multiple vision tasks while maintaining computational efficiency.

Mathematical Definition of Global Image Modeling

Definition: For an image modeling function $f: \mathbb{R}^{H \times W \times C} \rightarrow \mathbb{R}^{H \times W \times C}$ which is differentiable almost everywhere in \mathcal{D}_f , and a global influence function $\mathcal{I}: \{1, \dots, H\} \times \{1, \dots, W\} \times \{1, \dots, C\} \rightarrow \mathbb{R}^+$ whose derivative exists at all points in \mathcal{D}_f , if the Frobenius norm of the derivative of f with respect to any input pixel is greater than \mathcal{I} , then f exhibits global gradient dependency.

$$\left\| \frac{\partial f(\mathbf{X})}{\partial X_{i,j,c}} \right\|_F \geq \mathcal{I}(i,j,c) > 0, \quad \forall (i,j,c) \in \mathcal{D}_f; \quad \inf_{(i,j,c)} \mathcal{I}(i,j,c) \geq \tau > 0, \quad (1)$$

where $X_{i,j,c}$ represents the pixel value at position (i,j) and channel c in the input image $\mathbf{X} \in \mathbb{R}^{H \times W \times C}$, $\|\cdot\|_F$ denotes the Frobenius norm of the gradient tensor which measures the degree of influence of input pixels at different positions on the output image, \mathcal{D}_f denotes the set of points where both f and \mathcal{I} are differentiable, \inf represents the infimum, and τ ensures that each input pixel exerts a stable and non-negligible influence on the reconstructed output.

Constraint: For function f , due to the absence of inherent sequential regularity in image features, f should not impose strict order-dependent constraints on the input.

Frequency-Domain Modulated SSM Framework

Discrete-time SSM formulation: Continuous-time SSMs are discretized using a step size Δ and the zero-order hold method:

$$x_t = \bar{A}(\Delta)x_{t-1} + \bar{B}(\Delta)u_t, \quad y_t = Cx_t,$$

enabling efficient sequence modeling in neural networks.

SSM as dynamic convolution: The output can be equivalently expressed as

$$y_t = \sum_{k=0}^t K_k u_{t-k}, \quad K_k = C\bar{A}^k\bar{B},$$

showing that SSM performs step-by-step filtering on input sequences.

Frequency-domain equivalence: The discrete-time transfer function of SSM is

$$H(\omega) = C(e^{j\omega}I - \bar{A})^{-1}\bar{B},$$

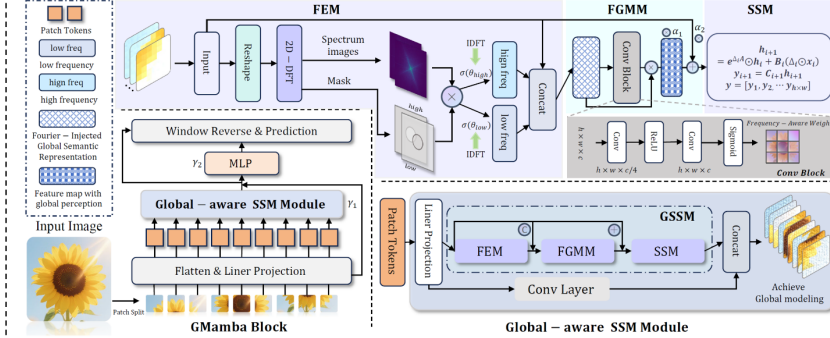
which is mathematically equivalent to the discrete-time Fourier transform of the convolution kernel K_k , linking time-domain convolution with frequency-domain representation.

$$\hat{K}(\omega) = \sum_{k=0}^{\infty} K_k e^{-j\omega k} = C(I - \bar{A}e^{-j\omega})^{-1}\bar{B} = e^{j\omega}C(e^{j\omega}I - \bar{A})^{-1}\bar{B} = e^{j\omega}H(\omega).$$

Theoretical Foundation of Frequency-Domain Modulated SSM.

- Information Preservation:** The encoded hidden state information x_t in SSMs can be transformed into a convolutional form K_k , where frequency-transformed information of K_k directly maps with original information, ensuring frequency domain operations preserve the representational capacity of the original model.
- Bidirectional Convertibility:** The dynamic convolution kernel K_k of SSM and its frequency domain transfer function $H(\omega)$ form a Fourier transform pair, enabling lossless bidirectional conversion between time-domain and frequency-domain representations.
- Linear Information Propagation:** Both the discrete convolution $y_t = \sum_{k=0}^t K_k u_{t-k}$ and the frequency-domain multiplication $\hat{Y}(\omega) = H(\omega) \cdot \hat{U}(\omega)$ are linear relations, ensuring consistent and stable linear information propagation through the network in either domain.

GMamba (GSSM) Block



Experiments

Experiments are conducted on five datasets across three backbones and vision tasks to validate GMamba's effectiveness. For semantic segmentation, four remote sensing datasets are used: Vaihingen, Potsdam, LoveDA, and UAVid. For instance segmentation and object detection, the MS-COCO dataset is employed. Results demonstrate GMamba's strong generalization and consistent improvements across different tasks and architectures.

Model Variant	Params (M)	FLOPs (G)	mIoU (%)	mF1 (%)	OA (%)
UNet ResNet54 (Baseline)	25.33	30.32	81.65	89.24	91.86
+ Swin (x7)	35.81	43.43	83.24 (16.96)	90.63 (11.99)	93.08 (11.23)
+ SwinV2 (x7)	35.86	43.46	83.10 (11.48)	90.54 (11.36)	93.04 (11.18)
+ ViM (x7)	25.20	35.81	83.02 (11.36)	90.51 (11.29)	92.93 (11.07)
+ VMamba (x7)	32.45	38.49	83.24 (11.96)	90.62 (11.96)	93.04 (11.80)
+ TinyViM (x7)	31.19	36.51	83.17 (11.52)	90.59 (11.59)	92.99 (11.23)
+ Mamba Version (x7)	34.87	39.73	83.20 (11.95)	90.62 (11.96)	93.00 (11.84)
+ Spatial Mamba (x7)	30.92	35.72	82.90 (11.25)	90.40 (11.36)	92.85 (10.99)
+ FreqMamba (x7)	31.98	36.21	83.00 (11.35)	90.50 (11.26)	92.90 (11.06)
+ Group Mamba (x7)	29.72	35.15	82.99 (11.36)	90.41 (11.19)	92.86 (11.08)
+ GMamba (x7) (Ours)	30.96	36.30	84.74 (13.08)	91.56 (12.32)	93.72 (11.86)
UNet SwinT (Baseline)	36.48	44.46	82.41	89.78	92.18
+ Swin (x7)	60.01	72.27	83.70 (11.26)	90.92 (11.47)	93.12 (10.66)
+ SwinV2 (x7)	60.07	73.28	83.73 (11.26)	90.94 (11.49)	93.19 (11.61)
+ ViM (x7)	42.83	57.99	83.62 (11.80)	90.88 (11.89)	93.05 (10.63)
+ VMamba (x7)	52.06	61.99	84.07 (11.61)	91.14 (11.36)	93.32 (11.14)
+ TinyViM (x7)	49.38	58.01	84.04 (11.60)	91.13 (11.36)	93.28 (11.35)
+ Mamba Version (x7)	52.47	62.95	83.68 (11.26)	90.90 (11.69)	93.10 (10.62)
+ Spatial Mamba (x7)	48.93	56.88	83.40 (10.96)	90.70 (10.95)	92.92 (10.74)
+ FreqMamba (x7)	49.95	57.98	83.55 (11.10)	90.82 (11.09)	93.00 (10.63)
+ Group Mamba (x7)	46.31	55.15	83.51 (11.05)	90.80 (11.05)	92.97 (10.70)
+ GMamba (x7) (Ours)	49.13	57.81	84.83 (11.26)	91.61 (11.86)	93.65 (11.47)
UNet ConvNeXt (Baseline)	42.83	68.88	83.11	90.19	92.30
+ Swin (x7)	81.95	97.69	84.82 (11.71)	91.59 (11.48)	93.57 (11.27)
+ SwinV2 (x7)	81.77	93.61	84.26 (11.25)	91.31 (11.62)	93.32 (11.02)
+ ViM (x7)	66.76	78.09	84.24 (11.16)	91.22 (11.08)	93.37 (11.05)
+ VMamba (x7)	73.99	86.41	84.56 (11.40)	91.45 (11.26)	93.56 (11.26)
+ TinyViM (x7)	71.31	82.43	84.38 (11.27)	91.33 (11.11)	93.41 (11.23)
+ Mamba Version (x7)	74.87	88.47	84.80 (11.46)	91.55 (11.36)	93.55 (11.25)
+ Spatial Mamba (x7)	70.93	84.92	84.50 (11.36)	91.32 (11.15)	93.35 (11.05)
+ FreqMamba (x7)	71.95	85.98	84.60 (11.16)	91.42 (11.29)	93.42 (11.23)
+ Group Mamba (x7)	68.20	79.60	84.56 (11.40)	91.44 (11.25)	93.38 (11.06)
+ GMamba (x7) (Ours)	71.08	85.86	86.00 (11.98)	92.31 (11.12)	93.99 (11.68)

Semantic Segmentation and Object Detection Results

Model Variant	mIoU (%)					mAP (%)					Params (M)	FLOPs (G)
	AP ₅₀ (%)	AP ₇₅ (%)	AP _s (%)	AP _m (%)	AP _l (%)	AP ₅₀ (%)	AP ₇₅ (%)	AP _s (%)	AP _m (%)	AP _l (%)		
Baseline (ResNet50)	37.2	57.8	40.4	21.5	40.6	40.0	43.88 (3.51)	20.07 (7.50)				
+ Swin (x7)	39.0	59.1	41.0	23.0	41.5	48.4	76.70 (6.44)	218.57 (6.00)				
+ SwinV2 (x7)	37.7	58.7	40.9	21.5	41.3	48.8	76.90 (6.61)	218.37 (6.00)				
+ TinyViM (x7)	37.6	58.6	40.6	21.6	41.3	48.4	61.40 (4.11)	208.87 (6.30)				
+ VMamba (x7)	37.6	58.8	40.9	21.5	41.4	48.8	60.94 (4.09)	209.07 (6.10)				
+ Mamba Version (x7)	37.7	58.9	41.0	21.6	41.4	48.5	60.80 (4.11)	209.97 (6.40)				
+ Spatial Mamba (x7)	37.6	58.7	40.9	21.5	41.3	48.4	60.90 (4.09)	209.57 (6.10)				
+ FreqMamba (x7)	37.5	58.5	40.5	21.4	41.0	48.5	60.80 (4.09)	209.57 (6.10)				
+ GMamba (x7) (Ours)	38.5	59.6	42.2	22.1	42.6	48.9	61.40 (4.11)	210.22 (7.65)				
Baseline (SwinT)	45.2	62.9	41.8	23.9	41.8	45.6	61.82 (7.57)	214.00 (6.10)				
+ Swin (x7)	42.9	60.6	40.9	22.1	40.4	46.4	55.25 (7.76)	214.50 (6.25)				
Baseline (ConvNeXt)	37.5	58.1	41.3	21.2	40.9	48.3	64.55 (5.21)	201.16 (7.50)				
+ Swin (x7)	38.4	59.2	41.7	21.6	42.2	49.3	79.55 (6.44)	203.64 (6.00)				
+ SwinV2 (x7)	38.3	59.1	41.7	22.0	41.9	49.0	79.55 (6.61)	203.44 (6.00)				
+ TinyViM (x7)	38.2	58.8	41.6	21.5	41.6	48.9	64.05 (4.11)	203.94 (7.30)				
+ VMamba (x7)	38.4	59.1	41.9	22.3	41.8	49.2	68.05 (4.11)	203.06 (7.40)				
+ Mamba Version (x7)	38.2	58.8	41.2	22.3	42.0	49.2	68.05 (4.11)	203.06 (7.40)				
+ Spatial Mamba (x7)	38.3	59.0	42.0	22.6	41.4	49.2	67.54 (4.11)	202.44 (7.80)				
+ FreqMamba (x7)	38.1	58.6	41.1	22.1	41.9	49.1	67.54 (4.11)	202.44 (7.80)				
+ GMamba (x7) (Ours)	39.1	60.1	42.8	23.3	42.8	50.3	64.05 (4.11)	203.26 (7.65)				
Baseline (SwinT)	42.2	64.7	48.0	26.7	43.5	55.6	47.78 (7.57)	206.00 (6.10)				
+ GMamba (x7) (Ours)	43.7	66.4	47.6	27.6	47.5	57.0	57.08 (7.65)	209.15 (6.25)				

Receptive field

