

# A Theoretical Analysis of Mamba's Training Dynamics: Filtering Relevant Features for Generalization in State Space Models

Mugunthan Shandirasegaran<sup>1</sup>, Hongkang Li<sup>2</sup>, Songyang Zhang<sup>3</sup>,  
Meng Wang<sup>4</sup>, Shuai Zhang<sup>1</sup>

<sup>1</sup>New Jersey Institute of Technology

<sup>2</sup>University of Pennsylvania

<sup>3</sup>University of Louisiana at Lafayette

<sup>4</sup>Rensselaer Polytechnic Institute



International Conference on Learning Representations (ICLR 2026)  
April 2026

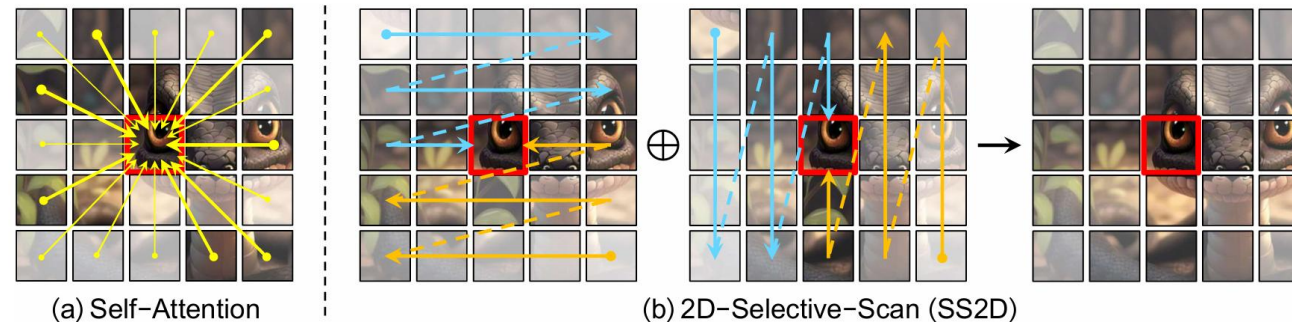




# Lack of Theoretical Understanding

- Visual State Space Model highlighted the **role of token order and scanning strategies**.

VMamba: Visual State Space Model (NeurIPS 2024)



Source: [Liu et al., 2024](#)

- Despite the empirical success, the theoretical understanding remains less investigated.
  - What types of data distributions can Mamba learn?
  - How does Mamba's gating affect representation learning compared to Transformers?

# Highlights of Our Theoretical Insights

- Characterization of the optimization trajectories of the gating weights
  - We prove that the gating mechanism selects class-relevant features while ignoring irrelevant ones.

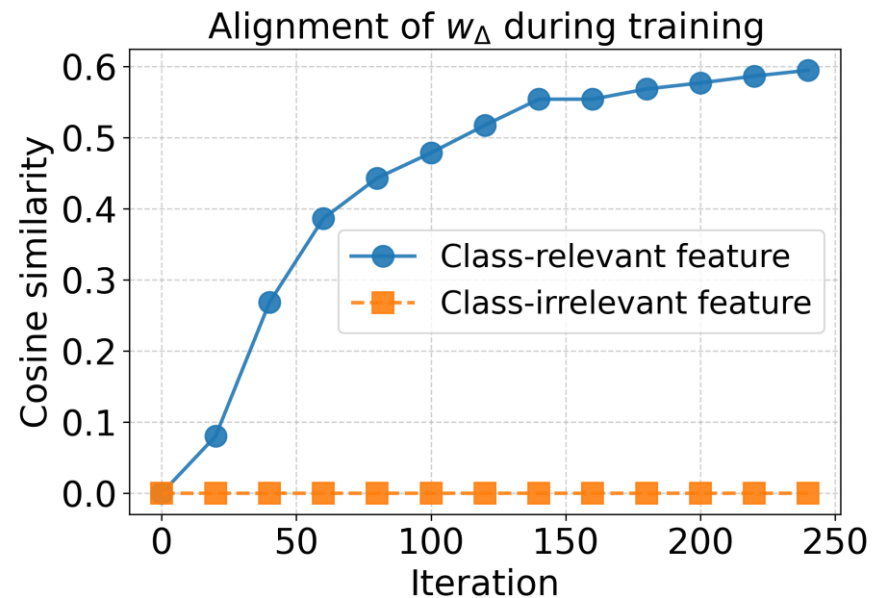
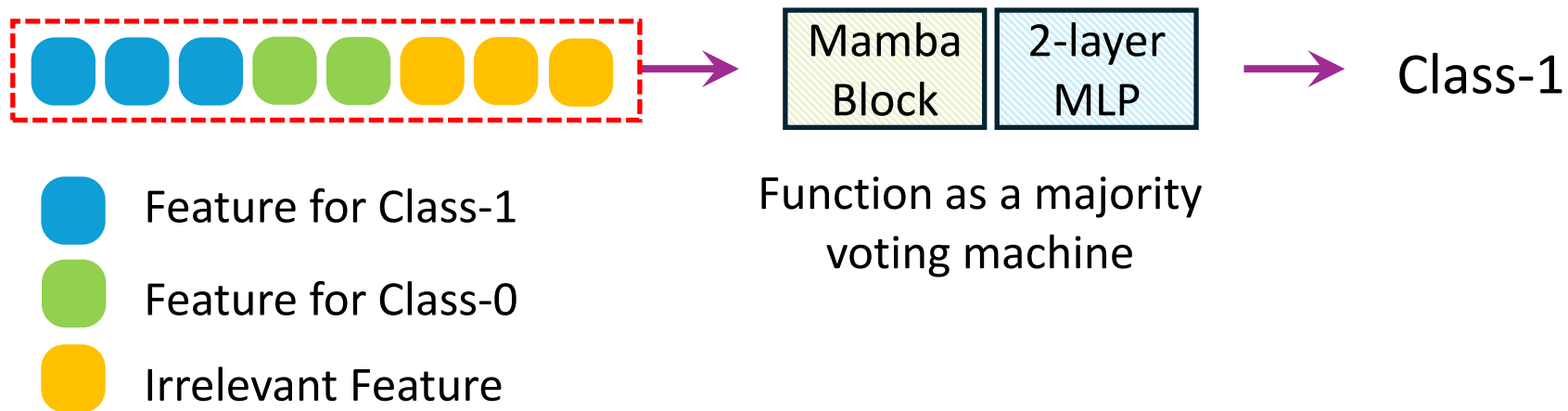


Fig 1: Alignment of  $w_\Delta$  during training

# Highlights of Our Theoretical Insights

- Characterization of the optimization trajectories of the gating weights
  - We prove that the gating mechanism selects class-relevant features while ignoring irrelevant ones.
- Mamba can learn from data where the label is determined by a **majority vote over class-relevant features**.



# Highlights of Our Theoretical Insights

- Mamba can learn locality-structured data, where class-relevant features appear in spatially or temporally concentrated regions.
  - Higher local concentration of class-relevant features accelerates learning.

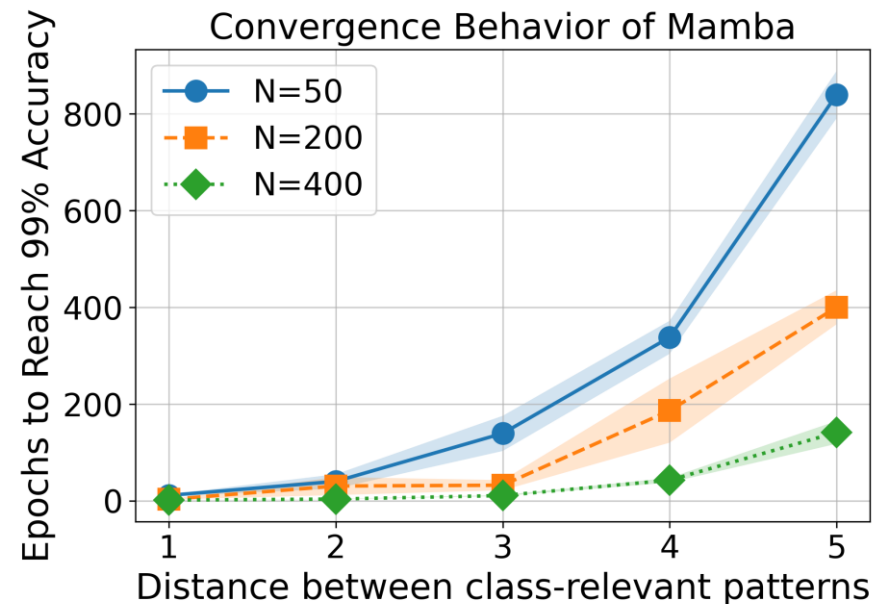


Fig 2: Convergence for **Locality-structured data**

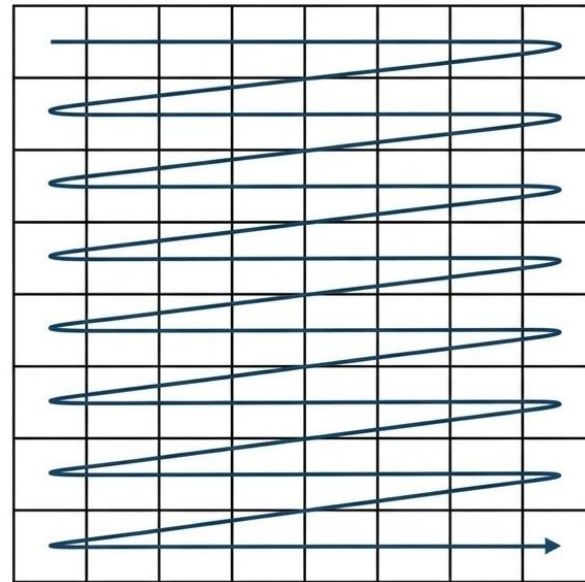
# Which Scan is better?

- Our results help explain why scan orders that preserve local structure, such as Hilbert curves, lead to improved performance.

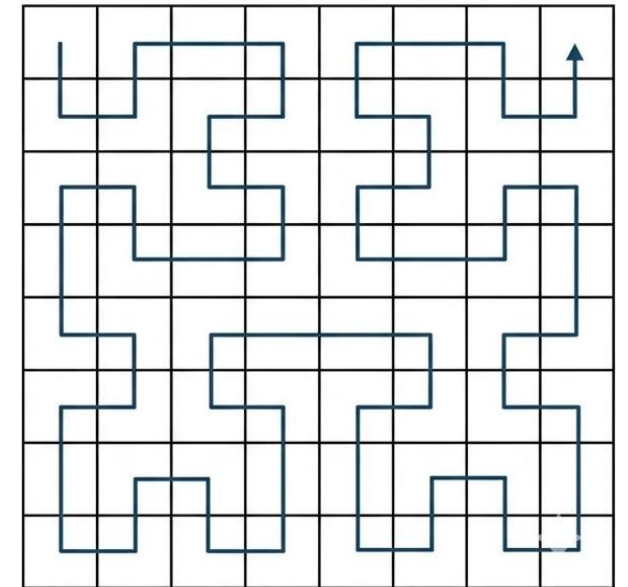
Order	Accuracy
Horizontal	$0.870 \pm 0.024$
Hilbert	$0.893 \pm 0.019$

Results on whole slide image classification from (Jiang et al. 2026).

Horizontal (row-wise)  
Scan



Hilbert Curve  
Scan



# Performance on Locality-Structured Data

- Mamba outperforms Transformer and local attention.
  - Global attention performs near random.
  - Local attention learns meaningful patterns but remains less effective than Mamba.

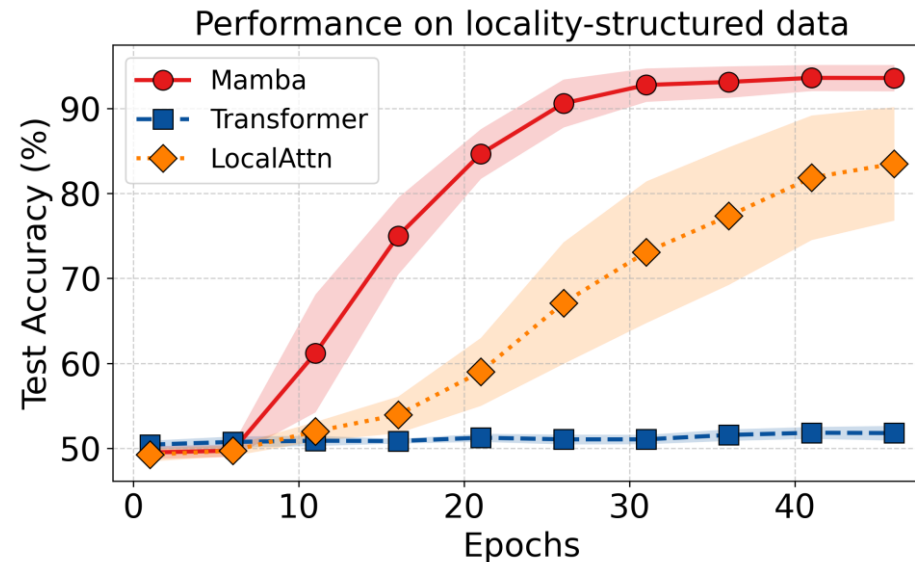


Fig 3: Mamba outperforms on **Locality-structured data**