DTP: Delta-Guided Two Stage Pruning for Mamba-based Multimodal Large Language Models

Anonymous authorsPaper under double-blind review

ABSTRACT

Multimodal large language models built on the Mamba architecture offer efficiency advantages, yet remain hampered by redundant visual tokens that inflate inference cost, with the prefill stage accounting for the majority of total inference time. We introduce Delta-guided Two stage Pruning (DTP), a method that progressively reduces token redundancy through selective pruning at early layer and complete pruning at late layer. Unlike Transformer-oriented pruning methods, our approach derives token importance directly from Mamba's internal parameters. The statistical distribution of these importance scores, combined with implicit attention patterns, then provides the basis for determining both the pruning layers and the tokens to be removed. Extensive evaluation across diverse benchmarks demonstrates that DTP reduces computation by nearly 50% while preserving task performance more effectively than existing pruning methods under the same reduction setting. Beyond efficiency, our analysis reveals previously underexplored behaviors of visual tokens within Mamba layers, suggesting a principled perspective for designing future pruning techniques in Mamba-based Multimodal Large Language Models.

1 Introduction

Multimodal Large Language Models (MLLMs) are capable of jointly understanding and generating across different modalities such as images and text, thereby handling complex tasks that are difficult for single-modality models (Liu et al., 2023; Dai et al., 2023; Peng et al., 2023; Wu et al., 2024). These capabilities have demonstrated outstanding performance in diverse tasks such as visual question answering (Guo et al., 2023; Hu et al., 2024; Fang et al., 2025; Wang et al., 2025; Dong et al., 2025) and reasoning segmentation (Lai et al., 2024; Ren et al., 2024; Xia et al., 2024).

Most existing MLLMs are built upon the Transformer architecture (Vaswani et al., 2017), which has shown strong performance in a wide range of multimodal tasks through the self-attention mechanism. However, self-attention requires computing interactions between all token pairs, leading to $O(n^2)$ time complexity that must be repeatedly incurred at every step when generating new tokens. To mitigate this redundant computation, the KV-Cache technique has been introduced, which effectively accelerates autoregressive generation during inference. Nevertheless, KV-Cache requires storing key-value pairs at every generation step, which results in substantial memory consumption. Furthermore, it does not reduce the computational cost of the prefill stage.

To overcome these structural limitations of Transformers, the recently proposed Mamba architecture (Gu & Dao, 2023) leverages State Space Models (SSMs) to recurrently update hidden states, thereby achieving linear-time complexity. Figure 1 contrasts the inference process between Transformer-based and Mamba-based MLLMs. As shown in Figure 1 (a), Transformers must recompute self-attention with all previous tokens whenever a new token is generated, causing decoding costs to increase linearly with sequence length. In contrast, as illustrated in Figure 1 (b), Mamba generates the next token through a single-step hidden state update without revisiting the entire input sequence. This structural difference allows Mamba-based MLLMs to achieve much lower memory usage and faster decoding compared to Transformer-based models, a benefit that has been empirically validated in recent studies (Liu et al., 2024; Qiao et al., 2024; Zhao et al., 2025).

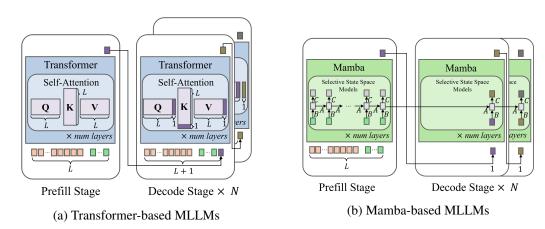


Figure 1: Comparison of inference structures between Transformer-based and Mamba-based MLLMs.

Thus, while Mamba provides clear efficiency advantages in the decoding stage, the majority of inference time is still spent in the prefill stage, where all input tokens must be processed initially. The inefficiency of the prefill stage is particularly pronounced in multimodal settings, since the number of visual tokens far exceeds that of text tokens, greatly increasing the overall input length. However, many of these visual tokens are redundant or uninformative, and often do not contribute to the final output (Chen et al., 2024). Based on this observation, various pruning methods have been proposed in Transformer-based MLLMs to reduce computational cost by removing unnecessary visual tokens (Alvar et al., 2025; Yang et al., 2025; Ye et al., 2025a;b; Lin et al., 2025). While effective in improving efficiency, these approaches rely on attention scores to estimate token importance, making them inherently specific to the Transformer architecture. As a result, they cannot be directly applied to Mamba.

To address this limitation, we propose a novel visual token pruning framework tailored for Mambabased MLLMs. The proposed DTP (Delta-guided Two stage Pruning) leverages the input-dependent parameter Δ_t to estimate token importance and performs pruning during inference without any retraining. Pruning is applied at two specific layers, with their positions determined from the observed distribution of token importance and the analysis of Mamba's implicit attention matrices. These findings indicate that redundant visual tokens can be reliably identified in the early layer, while in the late layer they no longer contribute meaningfully. Based on this rationale, our design effectively removes unnecessary tokens while incurring smaller performance degradation compared to other approaches.

Through extensive experiments on representative Mamba-based MLLMs, Cobra (Zhao et al., 2025) and RoboMamba (Liu et al., 2024), we demonstrate that DTP reduces computational cost by nearly half while maintaining comparatively high performance relative to existing methods. Furthermore, we identify the optimal internal parameters in Mamba for evaluating token importance, thereby maximizing the effectiveness of pruning. In addition to empirical results, our study provides new insights into the role and significance of visual tokens across layers in Mamba-based MLLMs, offering an effective methodology for visual token pruning.

Our main contributions are summarized as follows:

- We propose DTP (Delta-guided Two stage Pruning), the novel visual token pruning framework designed for Mamba-based MLLMs. DTP leverages the input-dependent parameter Δ_t to estimate token importance and performs visual token pruning during inference without requiring additional training.
- Pruning is applied at two specific layers, with their positions determined based on token importance distribution and the analysis of Mamba's implicit attention matrices. This design enables the effective removal of redundant visual tokens while maintaining stable performance.

• Extensive experiments on Cobra and RoboMamba demonstrate that DTP reduces FLOPs by nearly half with less performance degradation compared to existing methods. In addition, we identify the optimal internal parameters for token importance estimation, providing further insights into the role of visual tokens across layers in Mamba-based MLLMs.

2 RELATED WORK

Token Reduction in Mamba-based Model. Token reduction has been studied in Transformerbased models as a way to alleviate high computational cost and accelerate inference by removing unnecessary input tokens without causing significant performance degradation (Bolya et al., 2022; Kong et al., 2022; Haurum et al., 2023; Wei et al., 2023; Kim et al., 2024; 2025). However, since Mamba is based on State Space Models (SSMs) rather than the self-attention mechanism central to Transformers, token reduction techniques designed for Transformer architectures cannot be directly applied. Accordingly, recent studies have proposed token reduction methods specifically tailored to the structural characteristics of Mamba. Zhan et al. (2024) proposed a token pruning method for Vision Mamba (ViM) (Zhu et al., 2024) and PlainMamba (Yang et al., 2024), introducing a pruningaware hidden state alignment approach to stabilize the neighborhoods of the remaining tokens and a token importance estimation mechanism specific to Mamba, thereby improving inference speed while minimizing performance degradation. Shen et al. (2024) proposed Famba-V, a cross-layer token fusion method for ViM that identifies and merges similar tokens across layers, improving training efficiency while maintaining a balance with accuracy. Although these works explore diverse token reduction strategies such as pruning and fusion in Mamba-based models, they remain limited to unimodal vision tasks, and token reduction in Mamba-based MLLMs has not yet been sufficiently explored.

Mamba-based MLLMs. Recent studies have sought to leverage the structural efficiency of Mamba by extending Transformer-based MLLMs into Mamba-based MLLMs, aiming to achieve faster inference speed and improved efficiency in long-sequence processing. Qiao et al. (2024) proposed VL-Mamba, the first multimodal architecture that replaces the Transformer-based language model with a Mamba language model. It adopts SigLIP (Zhai et al., 2023) as the vision encoder and introduces the Vision Selective Scan (VSS) module within a multimodal connector to enhance representational capacity, demonstrating competitive results across a variety of multimodal benchmarks. Liu et al. (2024) introduced RoboMamba, which combines a CLIP vision encoder (Radford et al., 2021) with Mamba and adds it with a lightweight policy head to enable SE(3) pose prediction and vision-language-action modeling. In addition to its strong performance in robotic manipulation tasks, it also exhibits remarkable multimodal reasoning capability. Furthermore, Zhao et al. (2025) presented Cobra, which integrates a pre-trained Mamba language model with visual encoders such as DINOv2 (Oquab et al., 2023) and SigLIP (Zhai et al., 2023). Compared to Transformer-based MLLMs, it achieves both faster inference speed and superior performance.

3 Preliminaries

3.1 STATE SPACE MODELS AND MAMBA

State Space Models (SSMs) are the core structure of Mamba (Gu & Dao, 2023), which transform an input sequence $x(t) \in \mathbb{R}$ into an output sequence $y(t) \in \mathbb{R}$ through a hidden state $h(t) \in \mathbb{R}^N$, and are defined as:

$$h'(t) = Ah(t) + Bx(t), \quad y(t) = Ch(t),$$
 (1)

where A governs the state transitions, B maps the input into the hidden state, and C projects the hidden state into the output sequence.

In this basic form, SSMs have the limitation of linear time invariance (LTI), in which the same fixed parameters are applied to every time step of the input sequence, making it impossible to selectively process important or redundant information. In addition, while this formulation is designed for continuous systems, deep learning models generally operate on discrete systems, and therefore the continuous parameters must be discretized.

To address this issue, Mamba introduces an input-dependent parameter Δ_t . This parameter is computed from the input at each time step through a linear transformation followed by the softplus

function, which enables different state updates at each step. By leveraging this input-dependent parameter, Mamba discretizes SSMs and proposes Selective SSMs (S6), which allow inputs to be processed selectively at each time step, as formulated below:

$$\bar{A}_t = \exp(\Delta_t A), \quad \bar{B}_t = (\Delta_t A)^{-1} \left(\exp(\Delta_t A) - I \right) \Delta_t B$$

$$h_t = \bar{A}_t h_{t-1} + \bar{B}_t x_t, \quad y_t = C_t h_t$$
(2)

3.2 ATTENTION MATRICES IN MAMBA

Ali et al. (2024) demonstrated that the selective SSM layer can be unfolded into a causal kernel that closely resembles the attention matrix in Transformers. This finding provides an interpretation of Mamba's selective SSMs as implicitly incorporating attention-like behavior, even without an explicit attention mechanism. Specifically, the selective SSM formulation can be expanded into a convolutional form, yielding a kernel that functions as implicit attention weights. From Equation 2, assuming the initial state $h_0=0$, the output y_t can be written as:

$$y_t = \sum_{j=1}^t C_t \left(\prod_{k=j+1}^t \bar{A}_k \right) \bar{B}_j x_j \tag{3}$$

By defining

$$K_{t,j} = C_t \left(\prod_{k=j+1}^t \bar{A}_k \right) \bar{B}_j, \quad y_t = \sum_{j=1}^t K_{t,j} x_j,$$
 (4)

we see that $K_{t,j}$ represents the coefficient through which the input x_j is linearly transformed and incorporated into the output at time step t.

By arranging all coefficients, we obtain the following lower-triangular implicit attention matrix:

$$K = \begin{bmatrix} C_1 \bar{B}_1 & 0 & 0 & \cdots & 0 \\ C_2 \bar{A}_2 \bar{B}_1 & C_2 \bar{B}_2 & 0 & \cdots & 0 \\ C_3 \bar{A}_2 \bar{A}_3 \bar{B}_1 & C_3 \bar{A}_3 \bar{B}_2 & C_3 \bar{B}_3 & \cdots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ C_L \left(\prod_{k=2}^L \bar{A}_k \right) \bar{B}_1 & C_L \left(\prod_{k=3}^L \bar{A}_k \right) \bar{B}_2 & C_L \left(\prod_{k=4}^L \bar{A}_k \right) \bar{B}_3 & \cdots & C_L \bar{B}_L \end{bmatrix}$$
 (5)

This matrix $K \in \mathbb{R}^{L \times L}$, where L denotes the sequence length, can be interpreted in a way similar to the attention matrix of Transformers, where each row corresponds to the output at time t and each column shows how an input x_j propagates its influence to subsequent outputs.

4 Method

In this section, we propose DTP (Delta-guided Two stage Pruning), a pruning strategy for Mambabased MLLMs. As shown in Figure 2, DTP follows a two stage pruning strategy, performing selective pruning in the early layer and complete pruning in the late layer. To determine the specific layers where pruning is applied and to conduct selective pruning in the early layer, we leverage Δ_t , the key parameter that enables the selectivity of Mamba, to evaluate the importance of visual tokens in each Mamba block.

4.1 TOKEN IMPORTANCE FROM Δ_t

As described in Section 3.1, LTI SSMs apply the same parameters to all input sequences and therefore lack the selectivity to distinguish the relative importance of tokens. Accordingly, Mamba (Gu & Dao, 2023) introduces an input-dependent parameter Δ_t derived from the input sequence, and Δ_t serves as the key mechanism that enables selectivity by discretizing continuous SSMs and controlling the state transition matrix \bar{A}_t and the input mapping matrix \bar{B}_t . Building on this role, we directly leverage Δ_t to obtain token importance scores for visual token pruning in Mamba-based

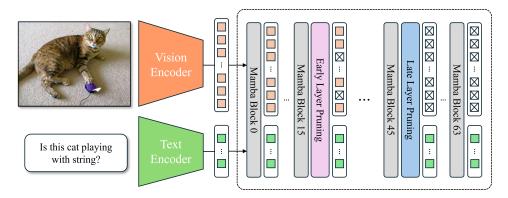


Figure 2: Overview of the proposed DTP(Delta-guided Two stage Pruning) method.

MLLMs. The importance score of the j-th token, denoted as s_i , is defined as follows:

$$s_j = \frac{1}{D} \sum_{d=1}^{D} \Delta_{j,d},\tag{6}$$

where D is the dimension size. This allows token importance to be directly estimated from Δ_t during inference, without requiring any additional training or modification. Based on the token importance scores s derived from Δ_t , we adopt a top-k selection strategy to retain the most informative visual tokens. At the pruning layers, visual tokens are ranked according to their importance scores, and only the top k tokens are preserved while the others are discarded. This reduces redundant visual information and enables computation to focus on the more critical tokens during reasoning. For a more comprehensive analysis, we further compared the proposed token importance parameter Δ_t with other internal parameters such as y_t , B_t , and C_t , and the results confirmed that Δ_t is the most suitable criterion for evaluating the importance of tokens in pruning.

4.2 PRUNING STRATEGY

Drawing on prior studies of token reduction in Transformer-based MLLMs (Chen et al., 2024; Lin et al., 2025; Ye et al., 2025a), which observed that visual tokens contain redundant information and exhibit minimal attention in deeper layers, we apply pruning specifically to visual tokens. Pruning is performed at two specific layers, one in the early layer and one in the late layer, and the positions are determined based on the token importance measure proposed in Section 4.1 and the analysis of Mamba's implicit attention matrix discussed in Section 3.2.

Selective pruning at the early layer. Pruning at very early layers has the advantage of achieving high computational efficiency, but it carries a significant risk of discarding tokens that could later serve as meaningful information in deeper layers. In addition, at such early layers, the distinction of token importance is not yet clear, making reliable selection difficult. To address this, we determine the appropriate layer for applying selective pruning by using the standard deviation of delta-guided token importance at each layer, which are defined as follows:

$$\operatorname{Std}_{\ell} = \sqrt{\frac{1}{N} \sum_{j=1}^{N} (s_{j,\ell} - \bar{s}_{\ell})^2}$$
 (7)

where, $s_{j,\ell}$ denotes the importance score of the j-th token at layer ℓ , \bar{s}_ℓ is the average importance at that layer, and N is the number of tokens. The standard deviation of token importance scores across layers is shown in Figure 3, where we observe that the 15th layer exhibits the first global peak. To examine this in more detail, Figure 4 visualizes the token importance distributions and the top 50% tokens separately at the 5th, 15th, 35th, and 45th layers. The results show that at the 15th and 35th layers, the standard deviation is low, with most tokens clustered around similar importance values while a small number of tokens retain relatively higher scores. In contrast, at the 5th and 45th layers, such separation is less evident, and tokens are more uniformly included within the top 50% across all positions. Taken together, these analyses indicate that the 15th layer, where the standard deviation

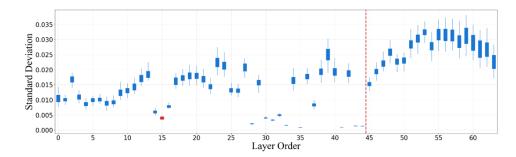


Figure 3: Standard deviation analysis of token importance across layers in a Mamba-based MLLM, Cobra (Zhao et al., 2025) The statistics were computed using a subset of the VQAv2 dataset, showing that the standard deviation reaches its first local peak at the 15th layer and increases sharply after the 45th layer.

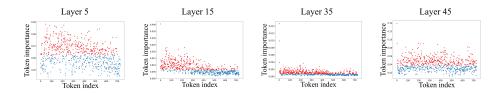


Figure 4: Visualization of token importance distributions at different layers of a Mamba-based MLLM, Cobra (Zhao et al., 2025). Red and blue points denote tokens within and outside the top 50% by importance, respectively.

reaches its first global peak, represents an optimal point for pruning redundant visual tokens while reliably preserving salient ones. Furthermore, as shown in Figure 5, the 15th layer is also the first depth at which attention patterns become significantly stronger compared to earlier layers, providing additional evidence for applying selective pruning at this layer. Therefore, we perform selective pruning at the 15th layer to preserve meaningful tokens while reducing unnecessary computation at an early layer, thereby improving efficiency.

Complete pruning at the late layer. As shown in Figure 3, the standard deviation of token importance scores sharply decreases at several points in the middle layers, while Figure 5 confirms that strong token interactions are still present at these depths. However, beyond the 45th layer, the standard deviation increases steeply and remains consistently high. Furthermore, an analysis of implicit attention patterns across all layers, presented in the Appendix B, reveals that after the 45th layer interactions with neighboring tokens almost completely vanish, indicating that visual tokens no longer make meaningful contributions. As also illustrated in Figure 4, in this regime it becomes difficult to reliably distinguish redundant tokens. Based on these observations, we adopt complete pruning at the 45th layer, where all visual tokens are removed and only text tokens are retained.

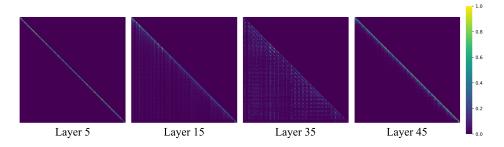


Figure 5: Visualization of implicit attention patterns across different layers in Mamba-based MLLM, Cobra (Zhao et al., 2025).

Table 1: Comparison of pruning methods on the Cobra model. The table presents FLOPs, FLOPs ratio, and evaluation scores across six benchmarks including GQA (Hudson & Manning, 2019), VQAv2 (Goyal et al., 2017), TextVQA (Singh et al., 2019), POPE (Li et al., 2023), VSR (Liu et al., 2023), and VizWiz (Gurari et al., 2018), along with the averaged performance, for different pruning methods and settings.

Method	FLOPs	FLOPs ratio	GQA	VQAv2	TextVQA	POPE	VSR	VizWiz	Avg
Baseline (Cobra)	2.01	100%	62.3	77.8	58.2	88.4	58.4	49.7	65.8
FastV (k=2, r=0.7)	1.45	72%	62.1 (-0.2)	77.4 (-0.4)	56.9 (-1.3)	87.7 (-0.7)	58.0 (-0.4)	49.8 (+0.1)	65.3 (-0.5)
VTW (k=45)	1.43	71%	62.1 (-0.2)	77.7 (-0.1)	58.2 (+0.0)	88.3 (-0.1)	58.5 (+0.1)	49.5 (-0.2)	65.7 (-0.1)
Ours (r=0.9)	1.35	67%	62.0 (-0.3)	77.7 (-0.1)	57.9 (-0.3)	88.3 (-0.1)	58.9 (+0.5)	49.7 (+0.0)	65.8 (+0.0)
FastV (k=2, r=0.5)	1.06	53%	61.7 (-0.6)	76.8 (-1.0)	55.0 (-3.2)	87.4 (-1.0)	57.3 (-1.1)	50.1	64.7
VTW (k=32)	1.04	52%	47.1 (-15.2)	54.1 (-23.7)	42.6 (-15.6)	74.1 (-14.3)	57.9 (-0.5)	48.5 (-1.2)	54.0 (-11.8)
Ours (r=0.5)	0.97	48%	61.4 (-0.9)	77.1 (-0.7)	56.1 (-2.1)	87.3 (-1.1)	57.9 (-0.5)	49.6 (-0.1)	64.9 (-0.9)

5 EXPERIMENT

5.1 EXPERIMENTAL SETUP

We evaluate our pruning strategy on two representative Mamba-based MLLMs: Cobra (Zhao et al., 2025) and RoboMamba (Liu et al., 2024). These models are representative Mamba-based MLLMs and are used as baselines since their pretrained weights are publicly available and they provide strong performance across multimodal tasks. For Cobra, we apply our method to six benchmarks that cover diverse aspects of multimodal reasoning, including GQA (Hudson & Manning, 2019), VQAv2 (Goyal et al., 2017), TextVQA (Singh et al., 2019), POPE (Li et al., 2023), VSR (Liu et al., 2023), and VizWiz (Gurari et al., 2018). For RoboMamba, we evaluate on five benchmarks, including OKVQA (Marino et al., 2019), GQA (Hudson & Manning, 2019), VQAv2 (Goyal et al., 2017), POPE (Li et al., 2023), and VSR (Liu et al., 2023).

The baselines considered in this study are two representative token pruning methods originally proposed for Transformer-based MLLMs: (1) FastV (Chen et al., 2024) prunes visual tokens at a designated layer k based on attention scores, with the pruning ratio controlled by r. Since Mamba-based MLLMs do not expose explicit attention scores, we replace them with our proposed Δ_t -based token importance measure for benchmarking. (2) VTW (Lin et al., 2025) determines the optimal layer k for withdrawing visual tokens by sampling a small subset of the dataset, comparing the original output with the output after token withdrawal, and selecting the earliest layer where the KL divergence between the two logits falls below a predefined threshold. As this method provides a general criterion that is independent of model architecture, it also serves as an appropriate baseline for our study.

All experiments are implemented in PyTorch and executed on a single NVIDIA RTX 5090 GPU to ensure fair comparisons.

5.2 MAIN RESULTS

Table 1 presents the results on the Cobra model. For a fair comparison, we adjusted the pruning layer k and keep ratio r so that each method has a similar FLOPs ratio. Our proposed DTP sets the keep ratio of all visual tokens to 0.9 in the early layer, reducing FLOPs to 67% of the baseline while maintaining the same average score of 65.8 as the unpruned baseline. For FastV, pruning was applied at the 2nd layer with a keep ratio of 0.7, achieving 72% FLOPs efficiency. However, this setting resulted in performance drops of 1.3 points on TextVQA and 0.7 points on POPE, showing a larger degradation compared to the baseline. In contrast, VTW identified the 45th layer as the optimal depth where the influence of visual tokens vanishes based on KL divergence, while showing only a 0.1 point decrease on average. This indicates that the optimal k determined by the withdrawal

Table 2: Comparison of pruning methods on the RoboMamba model. The table presents FLOPs, FLOPs ratio, and evaluation scores on five benchmarks including OKVQA (Marino et al., 2019), GQA (Hudson & Manning, 2019), VQAv2 (Goyal et al., 2017), POPE (Li et al., 2023), and VSR (Liu et al., 2023), together with the averaged performance across different pruning methods and settings

Method	FLOPs	FLOPs ratio	OKVQA	GQA	VQAv2	POPE	VSR	Avg
Baseline (RoboMamba)	0.70	100%	64.4	56.4	74.9	85.2	54.2	67.6
FastV ($k=2, r=0.7$)	0.50	71%	64.1	56.1	74.4	85.3	53.0	67.2
$1 \text{ ast } \mathbf{v} \ (h-2, 1-0.1)$			(-0.3)	(-0.3)	(-0.5)	(+0.1)	(-1.2)	(-0.4)
VTW $(k=45)$	0.50	71%	64.0	55.9	74.7	85.3	53.3	67.3
V I W (n=40)			(-0.4)	(-0.5)	(-0.2)	(+0.1)	(-0.9)	(-0.3)
Ours (r=0.9)	0.46	66%	63.8	56.1	74.7	85.2	53.2	67.3
Ours (7=0.9)			(-0.6)	(-0.3)	(-0.2)	(0.0)	(-1.0)	(-0.3)
FastV ($k=2, r=0.5$)	0.37	53%	63.4	55.1	73.4	84.2	52.1	65.6
rast v $(k=2, T=0.5)$			(-1.0)	(-1.3)	(-1.5)	(-1.0)	(-2.1)	(-2.0)
VTW (k=32)	0.36	51%	40.0	44.7	55.2	82.1	45.2	53.4
			(-24.4)	(-11.7)	(-19.7)	(-3.1)	(-9.0)	(-14.2)
0 (0.5)	0.24	49%	63.8	54.9	73.6	84.4	52.8	65.9
Ours (r=0.5)	0.34		(-0.6)	(-1.5)	(-1.3)	(-0.8)	(-1.4)	(-1.7)

criterion of VTW indirectly supports the validity of our complete pruning at the late layer strategy described in Section 4.2. When FLOPs were reduced to about half of the baseline, FastV was configured by setting the keep ratio to 0.5 at the 2nd layer for comparison with existing methods. Although VTW is a method for finding the optimal k, we fixed k=32 to enable comparison under this specific FLOPs condition. In this setting, FastV showed significant performance degradation, including a 3.2 point drop on TextVQA and a 1.1 point drop on VSR, while VTW exhibited large performance losses across all datasets. In contrast, our proposed method showed only a 0.9 point decrease on average compared to the baseline, demonstrating smaller performance degradation than the existing methods.

Table 2 presents the results on RoboMamba, where the pruning parameters k and r were set to the same values as used for Cobra. Our proposed method shows only a 0.3 point performance drop even when FLOPs are reduced to 66%. When FLOPs are further reduced to around 50%, VTW exhibits a severe performance degradation, similar to the case of Cobra. Both FastV and our method experience some decrease in performance, but our method shows relatively smaller drops on OKVQA and VSR, achieving an average score of 65.9 with the least overall degradation. This effect arises because RoboMamba employs the CLIP vision encoder (Radford et al., 2021) with only 256 image tokens. Compared to Cobra, which processes 729 tokens, the smaller number of tokens makes RoboMamba more susceptible to information loss under aggressive pruning.

5.3 ABLATION STUDY

All ablation studies in this section are conducted under the setting where FLOPs are reduced to approximately 50% of the baseline.

Identifying effective internal parameters for token importance. Table 3 presents the ablation study results when different internal parameters of the selective SSM are used to compute token importance. Comparing the output term y_t , the input coefficient B_t , the state coefficient C_t , and the temporal delta term Δ_t , it was found that Δ_t provides the most stable and effective signal in both Cobra and RoboMamba. The output term y_t shows generally competitive performance but falls short of Δ_t on TextVQA, while B_t and C_t exhibit relatively lower performance, particularly on TextVQA for Cobra and on POPE for RoboMamba. In contrast, Δ_t achieves the highest or comparable scores across all datasets in the ablation study, demonstrating that it serves as the most reliable criterion for distinguishing salient from redundant tokens. These results justify our choice of adopting Δ_t as the default measure for token importance.

Exploring strategies for token selection in pruning. Table 4 compares three strategies for selecting pruning candidates. The first baseline simply selects tokens at random. The second strategy applies a Top-k policy but also allows text tokens to be pruned. This strategy leads to catastrophic degradation

432 433 434

436

437

438

439 440 441

442

450

451

452

453

454

455

456

457

458

459

460

461

462 463

464 465

466

467

468

469

470

471

472

473

474

475 476 477

478

479

480 481

482

483

484

485

Table 3: Ablation study on internal parameters of Mamba for token importance estimation.

Cobra RoboMamba GQA TextVQA POPE OKVQA Parameter Vizwiz GOA 61.1 48.0 49.0 54.5 83.8 63.8 B_t 60.2 47.4 49.2 54.2 82.9 62.6 58.9 44.9 49.6 53.0 62.0 C_t 81.6

Table 4: Ablation study on token selection methods for Cobra and RoboMamba

		Cobra		RoboMamba			
Method	GQA	TextVQA	Vizwiz	GQA	POPE	OKVQA	
Random	61.4	53.7	49.4	54.1	83.6	63.2	
Top- k (all tokens)	7.01	10.3	30.7	21.2	56.8	46.5	
Top-k (visual only)	61.4	56.1	49.6	54.9	84.4	63.8	

Table 5: Ablation study on the effect of complete pruning at the late layer for Cobra and Robo-Mamba.

		Cobra					RoboMamba				
Complete pruning	FLOPs	FLOPs Ratio	GQA	TextVQA	Vizwiz	FLOPs	FLOPs Ratio	GQA	POPE	OKVQA	
Disable	1.26	100%	61.5	56.1	49.6	0.44	100%	55.1	84.4	63.9	
Enable	1.04	83%	61.4	56.1	49.6	0.34	77%	54.9	84.4	63.8	

across all benchmarks, with Cobra and RoboMamba both showing dramatic performance drops that make the results far below usable levels. These results highlight that text tokens are essential for reasoning and must not be removed. The third approach applies Top-k only to visual tokens. This strategy reduces redundancy while maintaining overall performance stably, showing that visual tokens can be pruned safely without affecting model performance. These results support the validity of a ranked Top-k pruning policy that targets only visual tokens.

Effect of complete pruning at the late layer. Table 5 reports the ablation results of complete pruning at the late layer. For Cobra, enabling complete pruning reduces FLOPs from 1.26 to 1.04, which corresponds to a 17% reduction, while accuracy remains nearly identical across GQA, TextVQA, and Vizwiz. Likewise, RoboMamba achieves a reduction from 0.44 to 0.34, amounting to a 23% decrease, with no meaningful accuracy drop on GQA, POPE, and OKVQA. These results emphasize that complete pruning at the late layer offers significant computational savings without sacrificing performance, validating it as an effective strategy for improving efficiency.

6 Conclusion

In this paper, we introduced Delta-guided Two stage Pruning (DTP), a novel framework for token pruning in Mamba-based MLLMs. Token importance is derived from the internal parameter Δ_t in Mamba's selective SSM, and the statistical distribution of these scores is analyzed together with implicit attention patterns to determine where pruning should occur and which tokens should be removed. Extensive experiments revealed that, under the same FLOPs budget, DTP preserves task performance more effectively than alternative pruning approaches. Ablation studies further demonstrated that Δ_t provides the most reliable criterion for estimating token importance, and that selecting the top-k visual tokens based on these scores is a reasonable pruning strategy. Moreover, we confirmed that applying complete pruning at late layers maintains performance equivalent to retaining tokens at those layers. Taken together, these findings establish DTP as an effective method for pruning visual tokens in Mamba-based MLLMs.

REFERENCES

Ameen Ali, Itamar Zimerman, and Lior Wolf. The hidden attention of mamba models. *arXiv* preprint arXiv:2403.01590, 2024.

Saeed Ranjbar Alvar, Gursimran Singh, Mohammad Akbari, and Yong Zhang. Divprune: Diversity-based visual token pruning for large multimodal models. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pp. 9392–9401, 2025.

Daniel Bolya, Cheng-Yang Fu, Xiaoliang Dai, Peizhao Zhang, Christoph Feichtenhofer, and Judy Hoffman. Token merging: Your vit but faster. *arXiv preprint arXiv:2210.09461*, 2022.

- Liang Chen, Haozhe Zhao, Tianyu Liu, Shuai Bai, Junyang Lin, Chang Zhou, and Baobao Chang. An image is worth 1/2 tokens after layer 2: Plug-and-play inference acceleration for large vision-language models. In *European Conference on Computer Vision*, pp. 19–35. Springer, 2024.
 - Wenliang Dai, Junnan Li, Dongxu Li, Anthony Tiong, Junqi Zhao, Weisheng Wang, Boyang Li, Pascale N Fung, and Steven Hoi. Instructblip: Towards general-purpose vision-language models with instruction tuning. *Advances in neural information processing systems*, 36:49250–49267, 2023.
 - Yuhao Dong, Zuyan Liu, Hai-Long Sun, Jingkang Yang, Winston Hu, Yongming Rao, and Ziwei Liu. Insight-v: Exploring long-chain visual reasoning with multimodal large language models. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pp. 9062–9072, 2025.
 - Wenlong Fang, Qiaofeng Wu, Jing Chen, and Yun Xue. guided mllm reasoning: Enhancing mllm with knowledge and visual notes for visual question answering. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pp. 19597–19607, 2025.
 - Yash Goyal, Tejas Khot, Douglas Summers-Stay, Dhruv Batra, and Devi Parikh. Making the v in vqa matter: Elevating the role of image understanding in visual question answering. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 6904–6913, 2017.
 - Albert Gu and Tri Dao. Mamba: Linear-time sequence modeling with selective state spaces. *arXiv* preprint arXiv:2312.00752, 2023.
 - Jiaxian Guo, Junnan Li, Dongxu Li, Anthony Meng Huat Tiong, Boyang Li, Dacheng Tao, and Steven Hoi. From images to textual prompts: Zero-shot visual question answering with frozen large language models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 10867–10877, 2023.
 - Danna Gurari, Qing Li, Abigale J Stangl, Anhong Guo, Chi Lin, Kristen Grauman, Jiebo Luo, and Jeffrey P Bigham. Vizwiz grand challenge: Answering visual questions from blind people. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 3608–3617, 2018.
 - Joakim Bruslund Haurum, Sergio Escalera, Graham W Taylor, and Thomas B Moeslund. Which tokens to use? investigating token reduction in vision transformers. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 773–783, 2023.
 - Wenbo Hu, Yifan Xu, Yi Li, Weiyue Li, Zeyuan Chen, and Zhuowen Tu. Bliva: A simple multimodal llm for better handling of text-rich visual questions. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pp. 2256–2264, 2024.
 - Drew A Hudson and Christopher D Manning. Gqa: A new dataset for real-world visual reasoning and compositional question answering. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 6700–6709, 2019.
 - Kwonyoung Kim, Jungin Park, Jin Kim, Hyeongjun Kwon, and Kwanghoon Sohn. Faster parameter-efficient tuning with token redundancy reduction. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pp. 30189–30198, 2025.
 - Minchul Kim, Shangqian Gao, Yen-Chang Hsu, Yilin Shen, and Hongxia Jin. Token fusion: Bridging the gap between token pruning and token merging. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pp. 1383–1392, 2024.
 - Zhenglun Kong, Peiyan Dong, Xiaolong Ma, Xin Meng, Wei Niu, Mengshu Sun, Xuan Shen, Geng Yuan, Bin Ren, Hao Tang, et al. Spvit: Enabling faster vision transformers via latency-aware soft token pruning. In *European conference on computer vision*, pp. 620–640. Springer, 2022.
 - Xin Lai, Zhuotao Tian, Yukang Chen, Yanwei Li, Yuhui Yuan, Shu Liu, and Jiaya Jia. Lisa: Reasoning segmentation via large language model. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 9579–9589, 2024.
 - Yifan Li, Yifan Du, Kun Zhou, Jinpeng Wang, Wayne Xin Zhao, and Ji-Rong Wen. Evaluating object hallucination in large vision-language models. *arXiv preprint arXiv:2305.10355*, 2023.

543

544

546

547

548

549

550 551

552

553 554

555

556

558

559

560 561

562

563

564 565

566

567

568

569

570

571

572 573

574

575

576

577

578 579

580

581

582

583

584

585

586

588

589

590 591

- 540 Zhihang Lin, Mingbao Lin, Luxi Lin, and Rongrong Ji. Boosting multimodal large language models with visual tokens withdrawal for rapid inference. In Proceedings of the AAAI Conference on 542 Artificial Intelligence, volume 39, pp. 5334–5342, 2025.
 - Fangyu Liu, Guy Emerson, and Nigel Collier. Visual spatial reasoning. Transactions of the Association for Computational Linguistics, 11:635–651, 2023.
 - Jiaming Liu, Mengzhen Liu, Zhenyu Wang, Pengju An, Xiaoqi Li, Kaichen Zhou, Senqiao Yang, Renrui Zhang, Yandong Guo, and Shanghang Zhang. Robomamba: Efficient vision-languageaction model for robotic reasoning and manipulation. Advances in Neural Information Processing Systems, 37:40085–40110, 2024.
 - Kenneth Marino, Mohammad Rastegari, Ali Farhadi, and Roozbeh Mottaghi. Ok-vqa: A visual question answering benchmark requiring external knowledge. In Proceedings of the IEEE/cvf conference on computer vision and pattern recognition, pp. 3195–3204, 2019.
 - Maxime Oquab, Timothée Darcet, Théo Moutakanni, Huy Vo, Marc Szafraniec, Vasil Khalidov, Pierre Fernandez, Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, et al. Dinov2: Learning robust visual features without supervision. arXiv preprint arXiv:2304.07193, 2023.
 - Zhiliang Peng, Wenhui Wang, Li Dong, Yaru Hao, Shaohan Huang, Shuming Ma, and Furu Wei. Kosmos-2: Grounding multimodal large language models to the world. arXiv preprint arXiv:2306.14824, 2023.
 - Yanyuan Qiao, Zheng Yu, Longteng Guo, Sihan Chen, Zijia Zhao, Mingzhen Sun, Qi Wu, and Jing Liu. Vl-mamba: Exploring state space models for multimodal learning. arXiv preprint arXiv:2403.13600, 2024.
 - Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pp. 8748-8763. PmLR, 2021.
 - Zhongwei Ren, Zhicheng Huang, Yunchao Wei, Yao Zhao, Dongmei Fu, Jiashi Feng, and Xiaojie Jin. Pixellm: Pixel reasoning with large multimodal model. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 26374–26383, 2024.
 - Hui Shen, Zhongwei Wan, Xin Wang, and Mi Zhang. Famba-v: Fast vision mamba with cross-layer token fusion. In European Conference on Computer Vision, pp. 268–278. Springer, 2024.
 - Amanpreet Singh, Vivek Natarajan, Meet Shah, Yu Jiang, Xinlei Chen, Dhruv Batra, Devi Parikh, and Marcus Rohrbach. Towards vga models that can read. In *Proceedings of the IEEE/CVF* conference on computer vision and pattern recognition, pp. 8317–8326, 2019.
 - Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. Advances in neural information processing systems, 30, 2017.
 - Zining Wang, Tongkun Guan, Pei Fu, Chen Duan, Qianyi Jiang, Zhentao Guo, Shan Guo, Junfeng Luo, Wei Shen, and Xiaokang Yang. Marten: Visual question answering with mask generation for multi-modal document understanding. In Proceedings of the Computer Vision and Pattern Recognition Conference, pp. 14460–14471, 2025.
 - Siyuan Wei, Tianzhu Ye, Shen Zhang, Yao Tang, and Jiajun Liang. Joint token pruning and squeezing towards more aggressive compression of vision transformers. In *Proceedings of the IEEE/CVF* conference on computer vision and pattern recognition, pp. 2092–2101, 2023.
 - Jiannan Wu, Muyan Zhong, Sen Xing, Zeqiang Lai, Zhaoyang Liu, Zhe Chen, Wenhai Wang, Xizhou Zhu, Lewei Lu, Tong Lu, et al. Visionllm v2: An end-to-end generalist multimodal large language model for hundreds of vision-language tasks. Advances in Neural Information Processing Systems, 37:69925-69975, 2024.

- Zhuofan Xia, Dongchen Han, Yizeng Han, Xuran Pan, Shiji Song, and Gao Huang. Gsva: Generalized segmentation via multimodal large language models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 3858–3869, 2024.
- Cheng Yang, Yang Sui, Jinqi Xiao, Lingyi Huang, Yu Gong, Chendi Li, Jinghua Yan, Yu Bai, Ponnuswamy Sadayappan, Xia Hu, et al. Topv: Compatible token pruning with inference time optimization for fast and low-memory multimodal vision language model. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pp. 19803–19813, 2025.
- Chenhongyi Yang, Zehui Chen, Miguel Espinosa, Linus Ericsson, Zhenyu Wang, Jiaming Liu, and Elliot J Crowley. Plainmamba: Improving non-hierarchical mamba in visual recognition. *arXiv* preprint arXiv:2403.17695, 2024.
- Weihao Ye, Qiong Wu, Wenhao Lin, and Yiyi Zhou. Fit and prune: Fast and training-free visual token pruning for multi-modal large language models. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 39, pp. 22128–22136, 2025a.
- Xubing Ye, Yukang Gan, Yixiao Ge, Xiao-Ping Zhang, and Yansong Tang. Atp-llava: Adaptive token pruning for large vision language models. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pp. 24972–24982, 2025b.
- Xiaohua Zhai, Basil Mustafa, Alexander Kolesnikov, and Lucas Beyer. Sigmoid loss for language image pre-training. In *Proceedings of the IEEE/CVF international conference on computer vision*, pp. 11975–11986, 2023.
- Zheng Zhan, Zhenglun Kong, Yifan Gong, Yushu Wu, Zichong Meng, Hangyu Zheng, Xuan Shen, Stratis Ioannidis, Wei Niu, Pu Zhao, et al. Exploring token pruning in vision state space models. *Advances in Neural Information Processing Systems*, 37:50952–50971, 2024.
- Han Zhao, Min Zhang, Wei Zhao, Pengxiang Ding, Siteng Huang, and Donglin Wang. Cobra: Extending mamba to multi-modal large language model for efficient inference. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 39, pp. 10421–10429, 2025.
- 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.

A TOKEN IMPORTANCE IN ROBOMAMBA

We provide further analyses on token importance using another Mamba-based MLLM, Robomamba (Liu et al., 2024). The figures show a similar trend to those observed in Cobra, supporting the generality of our pruning strategy.

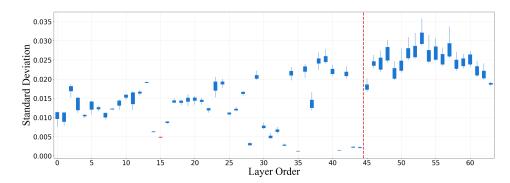


Figure 6: Standard deviation analysis of token importance across layers in a Mamba-based MLLM, Robomamba (Liu et al., 2024) The statistics were computed using a subset of the VQAv2 dataset, showing that the standard deviation reaches its first local peak at the 15th layer and increases sharply after the 45th layer.

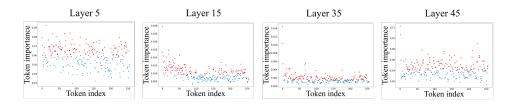


Figure 7: Visualization of token importance distributions at different layers of a Mamba-based MLLM, Robomamba (Liu et al., 2024). Red and blue points denote tokens within and outside the top 50% by importance, respectively.

B FULL IMPLICIT ATTENTION PATTERNS

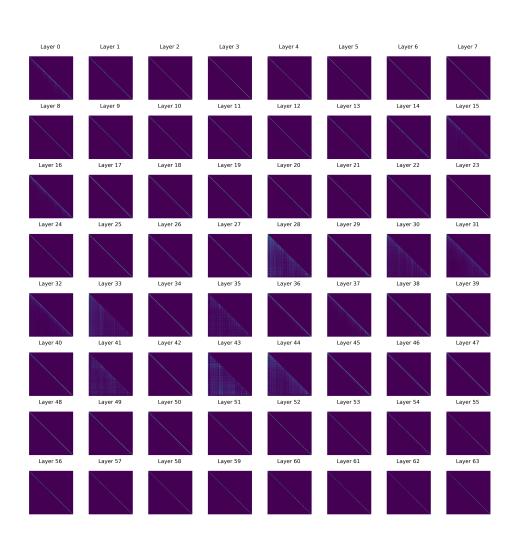


Figure 8: Full implicit attention patterns of each layer of Cobra.