Why 1 + 1 < 1 in Visual Token Pruning: Beyond Naïve Integration via Multi-Objective Balanced Covering

Yangfu Li , Hongjian Zhan , Xing Chen, Qi Liu, Yu-jie Xiong, Yue Lu
East China Normal University, Shanghai Jiao Tong University
Chongqing Institute of East China Normal University
Shanghai University of Engineering Science
{yfli_cee, qiliu}@stu.ecnu.edu.cn, hjzhan@cee.ecnu.edu.cn
guleimurray@sjtu.edu.cn, xiong@sues.edu.cn, ylu@cs.ecnu.edu.cn

Abstract

Existing visual token pruning methods target prompt alignment and visual preservation with static strategies, overlooking the varying relative importance of these objectives across tasks, which leads to inconsistent performance. To address this, we derive the first closed-form error bound for visual token pruning based on the Hausdorff distance, uniformly characterizing the contributions of both objectives. Moreover, leveraging ϵ -covering theory, we reveal an intrinsic trade-off between these objectives and quantify their optimal attainment levels under a fixed budget. To practically handle this trade-off, we propose Multi-Objective Balanced Covering (MoB), which reformulates visual token pruning as a bi-objective covering problem. In this framework, the attainment trade-off reduces to budget allocation via greedy radius trading. MoB offers a provable performance bound and linear scalability with respect to the number of input visual tokens, enabling adaptation to challenging pruning scenarios. Extensive experiments show that MoB preserves 96.4% of performance for LLaVA-1.5-7B using only 11.1% of the original visual tokens and accelerates LLaVA-Next-7B by 1.3-1.5× with negligible performance loss. Additionally, evaluations on Qwen2-VL and Video-LLaVA confirm that MoB integrates seamlessly into advanced MLLMs and diverse vision-language tasks. The code is available at https://github.com/YChenL/MoB.

1 Introduction

Multimodal large language models (MLLMs) have shown impressive performance across a variety of vision-language tasks, including visual understanding [30, 27, 20], visual question answering [40, 16, 37], and visual-language reasoning [9, 47, 45]. Since visual data exhibits much higher spatial redundancy than language, MLLMs are typically required to encode visual inputs as numerous tokens, resulting in substantial computational overhead.

To address this issue, visual token pruning methods are proposed to accelerate MLLMs by selecting representative subsets of visual tokens. Most pruning methods focus on two distinct objectives: Visual Preservation (VP) [6, 8, 59, 46], which retains tokens by minimizing redundancy or maximizing visual salience, and Prompt Alignment (PA) [58, 51, 48], which selects tokens most relevant to the prompt. Recently, several multi-objective approaches [31, 51, 42] have been proposed to integrate VP and PA through various complex strategies. Counterintuitively, these methods do not exhibit dominant superiority compared to single-objective approaches, as shown in Figure 1(a). This observation naturally raises a question: *Does integrating different objectives offer fundamental advantages?*

^{*}Corresponding Author.

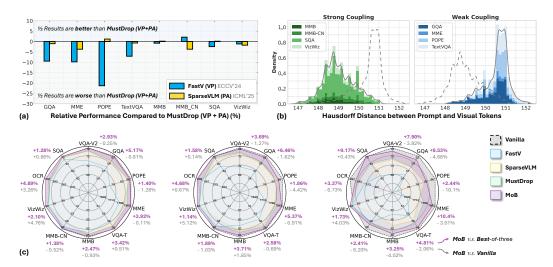


Figure 1: (a) Comparison of single- vs. bi-objective pruning methods on LLaVA-1.5-7B at a 66.7% pruning rate; (b) distribution of the prompt-visual coupling, revealing two distinct patterns across various tasks: weak coupling (large distance) and strong coupling (small distance); (c) radar charts of LLaVA-1.5-7B with visual tokens reduced from 576 to 192, 128, and 64 (*left-to-right*), demonstrating the consistent improvements of MoB across 10 well-recognized benchmarks.

Inspired by this question, we formulate preservation using the *Hausdorff distance* between the original and pruned token sets and derive the first closed-form error bound for visual token pruning (Lemma 1). This bound depends on VP and PA, while it is also affected by a prompt-visual coupling, measured by the Hausdorff distance between prompt and visual tokens. Notably, we identify two patterns of this coupling across popular benchmarks, as presented in Figure 1(b): weak coupling with large distance (*e.g.*, TextVQA, POPE) and strong coupling with small distance (*e.g.*, MMB, VizWiz). Our further analysis reveals that the effectiveness of the pruning objectives varies under distinct coupling patterns (Lemma 2). However, existing multi-objective methods overlook this variation and integrate VP and PA via constant strategies, yielding inconsistent improvements over single-objective baselines.

To quantify the effect of prompt-visual coupling, we reexamine visual token pruning from a geometric covering perspective. In this view, the retained tokens can be thought of as the union of two disjoint covers for prompt and visual tokens, where each objective corresponds to a Hausdorff covering radius, and the prompt-visual coupling is represented by the inter-cover diameter. By analyzing the geometric relationship between the radii and the diameter, we reveal an intrinsic trade-off between the two objectives (Theorem 1), which identifies the optimal attainment level of each objective to achieve the performance ceiling under a fixed pruning budget and prompt-visual coupling.

For a practical solution to this trade-off, we propose Multi-objective Balanced Covering (MoB), a training-free visual token pruning method with provable performance guarantees and multilinear complexity (Theorem 2). MoB partitions the retained tokens into two disjoint subsets for PA and VP, employing greedy radius-trading strategies to reduce the trade-off in objective attainment to a budget allocation problem. This allows MoB to achieve the optimal balance under each coupling pattern by selecting appropriate subset sizes. As shown in Figure 1(c), MoB consistently outperforms both single-objective and multi-objective baselines by a clear margin at identical pruning rates. Besides, MoB accelerates LLaVA-Next-7B by $1.3-1.5\times$ with negligible performance loss. Ablation studies further validate our theoretical analysis. Our key contributions are summarized as follows:

- To our knowledge, we present the first closed-form error bound for visual token pruning and its practical relaxation, characterizing the contributions of the two objectives to preservation quality.
- **②** We quantify the trade-off between the objectives and identify their optimal attainment level under a fixed budget and prompt-visual coupling, offering valuable insights into visual token pruning.
- **3** We propose Multi-objective Balanced Covering (MoB) for training-free visual token pruning, which reduces the trade-off of objective attainment to a budget allocation problem via two greedy radiustrading strategies, yielding both a provable performance guarantee and multilinear scalability.

• Extensive experiments across 14 public benchmarks demonstrate the superiority of MoB. For instance, it retains 96.4% and 97.9% performance for LLaVA-1.5-7B and Video-LLaVA-7B with an 88.9% reduction ratio, outperforming the second-best method by 2.7% and 1.6%, respectively. MoB can also be readily incorporated into advanced MLLMs, such as LLaVA-Next and Qwen2-VL.

2 Background

2.1 Related Work

Multimodal Large Language Model (MLLM). MLLMs [30, 21, 60, 28] have achieved remarkable progress in vision-language reasoning, owing to their robust cross-modality modeling via attention mechanisms [43, 34]. However, the spatial redundancy inherent in visual signals typically leads to a large number of input tokens [25, 22, 29, 44], particularly in high-resolution images and multi-frame videos (*e.g.*, 2048 tokens in Video-LLaVA [27]). This issue exacerbates the quadratic scaling problem of attention mechanisms, posing significant computational challenges. Moreover, to further enhance the visual capability by incorporating high-quality details, advanced MLLMs are now designed to support higher resolution images [24, 11, 10, 4], thereby necessitating the processing of even more visual tokens (*e.g.*, 2880 tokens in LLaVA-NEXT [29]). In these scenarios, effectively selecting representative visual tokens becomes a critical requirement for the real-world application of MLLMs.

Visual Token Pruning. Due to the spatial redundancy, inputs to MLLMs contain numerous less informative visual tokens. Visual token pruning accelerates MLLMs by selectively retaining only the most critical tokens during inference. Existing methods typically focus on either visual preservation (VP) [6, 38, 8, 52, 57, 32, 46] or prompt alignment (PA) [58, 51, 48]. VP-driven methods, such as ToMe [6] and LLaVA-PruMerge [38], reduce redundancy by merging similar tokens, while FastV [8] and FasterVLM [57] select tokens based on visual salience. PA-driven approaches like SparseVLM [58] rely on cross-modal attention to identify prompt-relevant tokens. More recently, MustDrop [31] integrates VP and PA through a multi-stage pruning pipeline, reporting notable improvements. Despite these advances, existing methods largely overlook the varying relative importance of VP and PA across different scenarios. In this paper, we formally characterize the contribution of each objective under a fixed pruning budget, and propose an algorithm that balances these objectives per scenario, yielding consistent improvements across diverse pruning conditions.

2.2 Preliminaries

Pipeline of MLLM. MLLMs perform vision-language reasoning by jointly processing multimodal inputs in a shared representation space. Formally, given visual tokens $\mathcal{V}^{(1)}$ extracted from the visual inputs and prompt tokens $\mathcal{P}^{(1)}$ encoded from user prompts, the multimodal input is defined as

$$\mathcal{X}^{(1)} \ = \ \mathcal{V}^{(1)} \ \sqcup \ \mathcal{P}^{(1)}, \quad \ \mathcal{V}^{(1)} = \{v_1^{(1)}, \dots, v_N^{(1)}\}, \ \mathcal{P}^{(1)} = \{p_1^{(1)}, \dots, p_L^{(1)}\} \subseteq \mathbb{R}^d,$$

where N and L denote the numbers of visual and prompt tokens, respectively. We regard both $\mathcal{V}^{(1)}$ and $\mathcal{P}^{(1)}$ as compact sets on d-dimensional Euclidean space $(\mathbb{R}^d, \|\cdot\|)$. The input $\mathcal{X}^{(1)}$ is then fed into a language model $\mathcal{F}_{[1,I]}$ with I transformer block, and the final output is given by

$$y = \mathcal{F}_{[1,I]}(\mathcal{X}^{(1)})$$
 where $\mathcal{F}_{[1,I]} = f_I \circ f_{I-1} \circ \ldots \circ f_1$,

In particular, each f_ℓ follows the standard Transformer (e.g., multi-head self-attention [43], layer normalization [3, 50]). The intermediate feature for any layer $\ell \in \{2, \dots, I\}$ is defined as

$$\mathcal{X}^{(\ell)} \coloneqq \mathcal{F}_{[1,\ell-1]}(\mathcal{X}^{(1)}) = \mathcal{V}^{(\ell)} \sqcup \mathcal{P}^{(\ell)}, \qquad \mathcal{F}_{[1,\ell-1]} \coloneqq f_{\ell-1} \circ \ldots \circ f_1,$$

with $\mathcal{V}^{(\ell)}$ and $\mathcal{P}^{(\ell)}$ representing the visual and prompt tokens after $\ell-1$ layers, respectively.

Visual Token Pruning. To accelerate MLLMs with minimal performance loss, visual token pruning selectively removes less-informative visual tokens at chosen intermediate layers of the language model $\mathcal{F}_{[1,I]}$. Specifically, for any chosen layer f_ℓ , $\ell \in \{2,\ldots,I\}$, pruning algorithms first select a subset $\mathcal{S}^{(\ell)} \subseteq \mathcal{V}^{(\ell)}$ of size K (i.e., pruning budget) and form the pruned input $\mathcal{X}_{\mathrm{s}}^{(\ell)} = \mathcal{S}^{(\ell)} \sqcup \mathcal{P}^{(\ell)}$. The corresponding output before and after pruning are then defined as

$$y = \mathcal{F}_{[\ell,I]} ig(\mathcal{X}^{(\ell)} ig), \quad y_{\mathrm{s}} = \mathcal{F}_{[\ell,I]} ig(\mathcal{X}^{(\ell)}_{\mathrm{s}} ig) \quad ext{where} \quad \mathcal{F}_{[\ell,I]} \coloneqq f_{I} \circ \cdots \circ f_{\ell}.$$

Notation	Description	Notation	Description
ℓ , I	Pruning layer index; Final layer index.	f_{ℓ}	Transformer block at layer ℓ .
d	Embedding dimension.	$\mathcal{V}^{(\ell)}$	Visual tokens at layer ℓ , <i>i.e.</i> , $\mathcal{V} \subseteq \mathbb{R}^d$.
$\mathcal{P}^{(\ell)}$	Prompt tokens at layer ℓ , <i>i.e.</i> , $\mathcal{P} \subseteq \mathbb{R}^d$.	$\mathcal{S}^{(\ell)}$	Retained visual tokens at layer ℓ , <i>i.e.</i> , $S \subseteq V$.
$\mathcal{X}^{(\ell)}$	All tokens at layer ℓ , <i>i.e.</i> , $\mathcal{X} = \mathcal{V} \sqcup \mathcal{P}$.	$\mathcal{X}_{\mathrm{s}}^{(\ell)}$	Retained tokens at layer ℓ , <i>i.e.</i> , $\mathcal{X}_s = \mathcal{S} \sqcup \mathcal{P}$.
N, L	#visual tokens, $ \mathcal{V} = N$; #prompt tokens, $ \mathcal{P} = L$.	$K^{^{\!$	Pruning budget, $ S = K$.
$\mathcal{F}[1,I]$	Full model (layers $1 \dots I$).	$\mathcal{F}[\ell,I]$	Submodel from layer ℓ to I .
$y,\ y_{ m s}$	Outputs with full tokens \mathcal{X} / pruned tokens \mathcal{X}_{s} .	$d_H(A, B)$	Hausdorff distance between sets A and B .
C_ℓ	Lipschitz constant of $\mathcal{F}[\ell, I]$ w.r.t. d_H .	η	Visual-prompt coupling bound: $d_H(\mathcal{V}, \mathcal{P}) \leq \eta$.
$\mathcal{S}_{\mathrm{p}},~\mathcal{S}_{\mathrm{v}}$	Prompt center / Visual center set, $S=S_p \sqcup S_v$.	$K_{\rm p},~K_{\rm v}$	Budgets for S_p and S_v , $K=K_p+K_v$.
$\epsilon_{ m p}$	Covering radius for \mathcal{P} , $d_H(\mathcal{S}_p, \mathcal{P})$.	$\epsilon_{ m v}$	Covering radius for V , $d_H(S_v, V)$.
$\mathcal{N}(\mathcal{X},\epsilon)$	Covering number of \mathcal{X} at radius ϵ .	$d_{ m eff}$	Effective dimension of V , P .
a, b	Covering-number lower/upper constants for \mathcal{P} .	a',b'	Covering-number lower/upper constants for V .
ϵ_0	Validity radius for covering bounds.	δ	Small dilation radius ($\delta \ll \eta$).
$\mathcal{V}_{\delta},~\mathcal{P}_{\delta}$	δ -dilation of $\mathcal V$ and $\mathcal P$.	$B(c,\epsilon)$	Ball $\{x: x-c \le \epsilon\}$.
z	Radius scaling factor (> 1) .	D_1	Trade-off constant $(4aa')^{1/d_{\text{eff}}}$.
D_2	Trade-off constant $1/z^2$.	ϵ^*	Optimal radius $\max\{\eta/z, \sqrt{D_1}K^{-1/d_{\text{eff}}}\}.$
k	Fold for the proposed nearest-neighbor covering.	$\mathcal{S}_{\mathrm{p}}'$	Candidate set before final truncation by $K_{\rm p}$.
$\alpha(\eta, k, L)$	Alignment constant $\eta(bkL/a)^{1/d_{\text{eff}}}$.	β^{r}	Preservation constant $2 (b')^{1/d_{\text{eff}}}$.
$\Theta(\cdot)$	Asymptotically equal (same order); i.e., $\exists c_1, c_2 >$	$\Omega(\cdot)$	Asymptotic lower bound (at least on the order of g),
	$0, n_0: c_1g(n) \le f(n) \le c_2g(n) \text{ for } n \ge n_0.$		i.e., $\exists c > 0, n_0 : f(n) \ge c g(n) \text{ for } n \ge n_0.$

Table 1: Summary of notation used in the theoretical framework.

Finally, the objective of visual token pruning is formulated as

$$\mathcal{S}^{(\ell)*} = \operatorname{argmin}_{\mathcal{S}^{(\ell)} \subset \mathcal{V}^{(\ell)}, \ |\mathcal{S}^{(\ell)}| = K} \|y - y_s\|_2.$$

Notation. For brevity we omit the layer index (ℓ) and simply write $\mathcal{X} = \mathcal{V} \sqcup \mathcal{P}$ and $\mathcal{X}_s = \mathcal{S} \sqcup \mathcal{P}$ to denote the input and its pruned counterpart at an arbitrary layer f_{ℓ} . We use \mathcal{F} to denote any composition mapping of the full model $\mathcal{F}_{[1,I]}$. Finally, we let $\|\cdot\|$ denote the Euclidean norm.

3 Methodology

3.1 Revisiting Visual Token Pruning: Insights into Prompt-Visual Coupling

As shown in Fig. 1(a), multi-objective pruning methods fail to achieve the expected improvements, and objective-specific methods exhibit inconsistent performance across benchmarks. These observations motivate us to reexamine the problem of visual token pruning. We begin by introducing Assumption 1, which quantifies pruning performance in terms of the preservation of the original token set.

Assumption 1 (Lipschitz Continuity w.r.t. the Hausdorff Distance). *Assume every partial composition* \mathcal{F} (from layer ℓ to I) of the language model is Lipschitz continuous w.r.t. the Hausdorff distance with constant $C_{\ell} \geq 1$. Formally, for any intermediate token sets $\mathcal{X}, \mathcal{X}_s \subset \mathbb{R}^d$,

$$\|\mathcal{F}(\mathcal{X}) - \mathcal{F}(\mathcal{X}_s)\| \le C_\ell d_H(\mathcal{X}, \mathcal{X}_s),$$

where d_H is the Hausdorff distance induced by the Euclidean norm:

$$d_{H}(\mathcal{X}, \mathcal{X}_{s}) := \max \left\{ \sup_{x \in \mathcal{X}} \inf_{x_{s} \in \mathcal{X}_{s}} \|x - x_{s}\|, \sup_{x_{s} \in \mathcal{X}_{s}} \inf_{x \in \mathcal{X}} \|x - x_{s}\| \right\}.$$
 (1)

Subsequently, we measure the preservation of the original token set \mathcal{X} using three pairwise distances among visual tokens \mathcal{V} , retained tokens \mathcal{S} , and prompt tokens \mathcal{P} , thereby establishing a unified performance bound for various visual token pruning algorithms, as presented in Lemma 1.

Lemma 1 (An Error Bound for Visual Token Pruning). *Under Assump 1, given a token set with its pruned counterpart* $\mathcal{X} = \mathcal{V} \sqcup \mathcal{P}$, $\mathcal{X}_s = \mathcal{S} \sqcup \mathcal{P} \subseteq \mathbb{R}^d$, the pruning error bound is given by:

$$\|\mathcal{F}(\mathcal{X}) - \mathcal{F}(\mathcal{X}_{\mathrm{s}})\| \le C_{\ell} \max \Big\{ \min \big\{ d_H(\mathcal{S}, \mathcal{V}), d_H(\mathcal{V}, \mathcal{P}) \big\}, \min \big\{ d_H(\mathcal{S}, \mathcal{V}), d_H(\mathcal{S}, \mathcal{P}) \big\} \Big\}.$$

Remark. Here $d_H(S, P)$ and $d_H(S, V)$ describe the prompt alignment and visual preservation, while $d_H(V, P)$ is an inherent term that describes the prompt-visual coupling of input data.

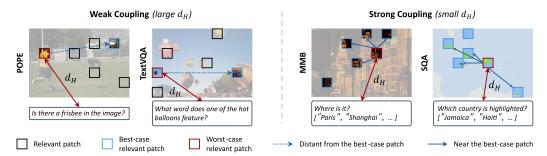


Figure 2: Illustration of prompt-visual coupling with two distinct patterns: In fine-grained tasks (*e.g.* POPE), only a few patches are critical, so the worst-case patch lies far from best-case ones, resulting in a large Hausdorff distance and making prompt alignment valuable. In coarse-grained tasks (*e.g.* MMB), many relevant patches contain the answer cues; thus, the worst-case patch remains close to best-case ones, yielding a small Hausdorff distance and making visual preservation more efficient.

Proof in Appendix E.1. By Lemma 1, in practical settings where $|\mathcal{S}| \ll |\mathcal{V}|$, pruning performance is governed by a non-trivial interaction among visual preservation, prompt alignment, and prompt-visual coupling. However, existing multi-objective methods typically overlook the coupling term $d_H(\mathcal{V}, \mathcal{P})$ and statically combine the two objectives across tasks, limiting their effectiveness. Our empirical evidence across popular benchmarks validates two distinct patterns of $d_H(\mathcal{V}, \mathcal{P})$, each favoring different pruning objectives, as shown in Figure 2. To further explicate the effect of prompt-visual coupling, we introduce Assumption 2 and propose a practical relaxed error bound in Lemma 3.

Assumption 2 (Prompt-Visual Coupling Bound). We assume the input visual data and prompts are not entirely unrelated; hence, there exists a constant $\eta > 0$ for any intermediate token set $\mathcal{X} = \mathcal{V} \sqcup \mathcal{P} \subseteq \mathbb{R}^d$ such that $d_H(\mathcal{V}, \mathcal{P}) \leq \eta$, ensuring the reasonability of vision-language reasoning.

Lemma 2 (A Relaxed Error Bound under Practical Budgets). *Under Assumptions 1 and 2,* let $\mathcal{X} = \mathcal{V} \sqcup \mathcal{P}$, $\mathcal{X}_s = \mathcal{S} \sqcup \mathcal{P} \subseteq \mathbb{R}^d$ with $|\mathcal{S}| = K \ll N$. Partition the retained token set \mathcal{S} into two disjoint subsets: $\mathcal{S} = \mathcal{S}_p \sqcup \mathcal{S}_v$, devoted to prompt alignment $d_H(\mathcal{S}_p, \mathcal{P})$ and visual preservation $d_H(\mathcal{S}_v, \mathcal{V})$, respectively. Then, the pruning error bound reduces to

$$\|\mathcal{F}(\mathcal{X}) - \mathcal{F}(\mathcal{X}_{s})\| \le C_{\ell} \max \{d_{H}(\mathcal{S}_{p}, \mathcal{P}), d_{H}(\mathcal{S}_{v}, \mathcal{V})\} + C_{\ell} \eta.$$

Proof in Appendix E.2. As Lemma 2 indicates, under weak coupling (large η), most visual regions are distant from prompt tokens in the semantic space. Consequently, if \mathcal{S}_p misses the critical patches, $d_H(\mathcal{S}_p, \mathcal{P})$ dominates the pruning error, making the selection of \mathcal{S}_p *i.e.*, prompt alignment, more significant. Conversely, under strong coupling (small η), $d_H(\mathcal{S}_p, \mathcal{P})$ tends to decrease in tandem with $d_H(\mathcal{S}_v, \mathcal{V})$, reducing the marginal benefit of prompt alignment. To further guide pruning methods design, we next quantify this trade-off governed by η through an ϵ -covering argument.

3.2 Quantifying Prompt-Visual Trade-Off: A Geometric Covering Perspective

We first introduce some geometric metrics in Definition 1, recasting each objective term $d_H(\mathcal{S}_p, \mathcal{P})$ and $d_H(\mathcal{S}_v, \mathcal{V})$ as covering radii and the coupling term $d_H(\mathcal{V}, \mathcal{P})$ as an inter-cover diameter. Next, we relate each recasted objective to its token budget $|\mathcal{S}_p|, |\mathcal{S}_v|$ via covering regularity in Lemma 3. Finally, by loading the budget constraint and applying the triangle inequality between radii and diameter, we derive a quantitative trade-off jointly governed by K and η in Theorem 1.

Definition 1 (ϵ -cover, Covering Number, and Covering Regularity). Let $(\mathbb{R}^d, \|\cdot\|)$ be the d-dimensional Euclidean space and let $\mathcal{X} \subseteq \mathbb{R}^d$ be a compact set.

(a) ϵ -cover. if there exists a finite set $\mathcal{C} = \{c_1, \ldots, c_M\} \subset \mathbb{R}^d$, an ϵ -cover of \mathcal{X} is given by

$$\mathcal{X} \subseteq \bigcup_{c \in \mathcal{C}} B(c, \epsilon), \qquad B(c, \epsilon) \coloneqq \{x \in \mathbb{R}^d : ||x - c|| \le \epsilon\},$$

where C is the collection of covering centers, and ϵ is the covering radius.

(b) Covering number. The minimum cardinality of C is the covering number of X at radius ϵ :

$$\mathcal{N}(\mathcal{X}, \epsilon) := \min \Big\{ M \in \mathbb{N} : \exists \mathcal{C} \subset \mathbb{R}^d, |\mathcal{C}| = M, \ \mathcal{X} \subseteq \bigcup_{c \in \mathcal{C}} B(c, \epsilon) \Big\}.$$

(c) Covering regularity. We say that \mathcal{X} satisfies d-dimensional covering regularity if there exist constants $0 < A \leq B$ and $\epsilon_0 > 0$ such that

$$A \epsilon^{-d} \leq \mathcal{N}(\mathcal{X}, \epsilon) \leq B \epsilon^{-d}, \quad \forall \epsilon \in (0, \epsilon_0].$$

Based on Definition 1(a) (b), S_p , $S_v \subseteq V$ can be thought of as two collections of centers such that

$$\mathcal{P} \subseteq \bigcup_{i=1}^{K_{p}} B(s_{p}^{(i)}, \epsilon_{p}), \ \mathcal{V} \subseteq \bigcup_{i=1}^{K_{v}} B(s_{v}^{(j)}, \epsilon_{v}),$$

where the radii are given by $\epsilon_{\rm p}\coloneqq d_H(\mathcal{S}_{\rm p},\mathcal{P}),\ \ \epsilon_{\rm v}\coloneqq d_H(\mathcal{S}_{\rm v},\mathcal{V}),$ and the covering numbers satisfy $\mathcal{N}(\mathcal{P},\epsilon_{\rm p})\le |\mathcal{S}_{\rm p}|,\ \mathcal{N}(\mathcal{V},\epsilon_{\rm v})\le |\mathcal{S}_{\rm v}|.$ Thereby, we derive a lower bound of the required budget, $i.e., |\mathcal{S}_{\rm p}|, |\mathcal{S}_{\rm v}|,$ to improve each objective, $i.e.,\epsilon_{\rm p},\epsilon_{\rm v},$ based on $d_{\rm eff}$ -dimensional covering regularity.

Lemma 3 (Covering Number Bounds). Given $\mathcal{P}, \mathcal{V} \subset \mathbb{R}^d$ with an effective dimension d_{eff} . Suppose their δ -dilations $\mathcal{V}_{\delta} := \bigcup_{v \in \mathcal{V}} B(v, \delta)$, $\mathcal{P}_{\delta} := \bigcup_{p \in \mathcal{P}} B(p, \delta)$ ($\delta \ll \eta$) satisfy d_{eff} -dimensional covering regularity; thus, there exist constants b > a > 0, b' > a' > 0 and $\epsilon_0 > \delta$ such that

$$a \, \epsilon_{\mathrm{p}}^{-d_{\mathrm{eff}}} \leq \mathcal{N}(\mathcal{P}, \epsilon_{\mathrm{p}}) \leq b \, \epsilon_{\mathrm{p}}^{-d_{\mathrm{eff}}}, \qquad a' \, \epsilon_{\mathrm{v}}^{-d_{\mathrm{eff}}} \leq \mathcal{N}(\mathcal{V}, \epsilon_{\mathrm{v}}) \leq b' \, \epsilon_{\mathrm{v}}^{-d_{\mathrm{eff}}}, \qquad \forall \, \epsilon_{\mathrm{p}}, \epsilon_{\mathrm{v}} \in (\delta, \epsilon_{0}],$$

Remark. Previous work suggests that both visual and language embeddings concentrate on a low-dimensional manifold, so the effective covering dimension satisfies the typical relation $d_{\text{eff}} \ll d$.

Proof in Appendix E.3. Lemma 3 demonstrates that once the radius (*i.e.*, the objective) falls below ϵ_0 , any further improvement of it demands a $\Theta(\epsilon^{-d_{\rm eff}})$ increase in the number of selected token.

By loading Lemma 3 into the budget constraint: $|S_p| + |S_v| = K$, and applying a two-step triangle inequality between the covering radii ϵ_p , ϵ_v and the inter-cover diameter η , we establish a K- η -bound in Theorem 1(b), which quantifies the trade-off governed by the budget and prompt-visual coupling.

Theorem 1 (Trade-off between Prompt Alignment and Visual Preservation). *Under Assumption 2 and the covering-regularity hypothesis of Lemma 3 with constants* $a, a', d_{\text{eff}} > 0$, there exist a radius-scaling factor z > 1 such that $\eta/z > \delta$ and $K < \mathcal{N}(\mathcal{P}, \eta/z) + \mathcal{N}(\mathcal{V}, \eta/z)$, for every pruning results $S = (S_p \sqcup S_v) \subseteq \mathcal{V}$ with budget K satisfying

$$\max \{ D_1 K^{-2/d_{\text{eff}}}, \ D_2 \eta^2 \} \le d_H(\mathcal{S}_{\mathbf{p}}, \mathcal{P}) \ d_H(\mathcal{S}_{\mathbf{v}}, \mathcal{V}),$$

where
$$D_1 := (4 \, a \, a')^{1/d_{\text{eff}}} > 0$$
, $D_2 := 1/z^2 > 0$.

Remark (Optimal Attainment Level). $D_1 K^{-2/d_{\rm eff}}$ is completely determined by the pruning budget, while $D_2 \eta^2$ quantifies the effect of prompt-visual coupling. The optimal attainment level per objective is given by $\epsilon^* = \max\{\eta/z, \sqrt{D_1} K^{-1/d_{\rm eff}}\}$. Any attempt to reduce one objective below ϵ^* forces the other above ϵ^* , thereby increasing the overall pruning error.

Remark (Effect of Budget and Coupling Strength). As K decreases, z correspondingly shrinks, ultimately making $D_2 \eta^2$ dominate the bound; while as K increases, both of the terms reduce, thereby diminishing the trade-off and tightening the overall error bound.

Proof in Appendix E.4. Theorem 1 characterizes the optimal attainment level for each objective under a fixed pruning budget and prompt-visual coupling. However, it is actually very challenging to dynamically determine the attainment level per objective during the pruning process. To address this, we propose Multi-objective Balanced Covering, which leverages the monotonic relationship between covering radii and numbers to reduce the trade-off of attainment to a budget-allocation problem.

3.3 Multi-Objective Balanced Covering: From Trade-Off to Budget Allocation

Motivated by the insights in §3.2, Multi-objective Balanced Covering (MoB) recasts visual token pruning as bi-objective covering. Specifically, given a token set $\mathcal{X} = \mathcal{V} \sqcup \mathcal{P} \subseteq \mathbb{R}^d$ with a budget K, the retained token set \mathcal{S} is defined as the union of a prompt center set \mathcal{S}_p and a visual center set \mathcal{S}_v :

$$\textstyle \mathcal{S} = \mathcal{S}_{\mathrm{p}} \, \sqcup \, \mathcal{S}_{\mathrm{v}} \subseteq \mathcal{V} \subseteq \mathbb{R}^d \quad \text{where} \quad \mathcal{P} \subset \bigcup_{i=1}^{K_{\mathrm{p}}} B(s_{\mathrm{p}}^{(i)}, \epsilon_{\mathrm{p}}), \ \ \mathcal{V} \subset \bigcup_{j=1}^{K-K_{\mathrm{p}}} B(s_{\mathrm{v}}^{(j)}, \epsilon_{\mathrm{v}}).$$

MoB then selects the cover centers (i.e., retained tokens) by minimizing the overall maximum radius:

$$(\mathcal{S}_{\mathrm{p}}^*,\,\mathcal{S}_{\mathrm{v}}^*) = \mathop{\arg\min}_{\mathcal{S}_{\mathrm{p}} \sqcup \mathcal{S}_{\mathrm{v}} \subseteq \mathcal{V},\, |\mathcal{S}_{\mathrm{p}}| = K_{\mathrm{p}},\, |\mathcal{S}_{\mathrm{v}}| = K - K_{\mathrm{p}}} \max\{\epsilon_{\mathrm{p}}(\mathcal{S}_{\mathrm{p}}), \epsilon_{\mathrm{v}}(\mathcal{S}_{\mathrm{v}})\}.$$

In practice, MoB solves this problem approximately by two sequential greedy covering procedures: selection of prompt center set S_p with budget K_p , and selection of visual center set S_v with the remaining budget $K - K_p$. By the covering number bounds given in Lemma 3, we have

$$K_{\rm p} = \Theta(\epsilon_{\rm p}^{-d_{\rm eff}}), \quad K - K_{\rm p} = \Theta(\epsilon_{\rm v}^{-d_{\rm eff}}),$$

where $d_{\rm eff}$ is the effective dimension of \mathcal{V} , \mathcal{P} . Accordingly, by selecting the unique budget $K_{\rm p}$ (i.e., fixing the remaining budget $K-K_{\rm p}$) under each coupling pattern, MoB ensures $\epsilon_{\rm p}$, $\epsilon_{\rm v}=\Omega(\max\{\eta/z,\,\sqrt{D_1}\,K^{-1/d_{\rm eff}}\})$, thus yielding provable performance guarantees across scenarios.

Normalization. For efficiency, MoB applies L2 normalization to each $x \in \mathcal{X}$ so that ||x|| = 1. Hence, for any token pair $x_1, x_2 \in \mathcal{X}$, the Euclidean distance can be induced by their cosine similarity:

$$||x_1 - x_2|| = \sqrt{2 - 2\cos(x_1, x_2)}.$$

Selection of Prompt Center Set \mathcal{S}_p . Since all $s_p \in \mathcal{V}$ lie outside \mathcal{P} , a typical solution for minimizing the radius ϵ_p is Nearest-Neighbor covering (NN covering) [15], which uniformly allocates the nearest $s_p \in \mathcal{V}$ for each prompt token. However, the contribution of each prompt token is inequivalent, especially under weak prompt-visual coupling; thus, equal allocation risks missing the "best-case tokens." To remedy this, we introduce a k-fold NN covering procedure. Formally, let $L = |\mathcal{P}|$ and k > 1 be a hyperparameter; we first utilize a temporary budget of kL to form a candidate set.

$$S'_{p} = \bigcup_{p \in \mathcal{P}} \arg \operatorname{topk}_{s \in \mathcal{V}} (\cos(s, p), k), \quad |S'_{p}| \ge K_{p},$$

thereby over-sampling the k nearest visual tokens for each prompt token. Subsequently, we refine the candidate set by selecting the final $K_{\rm D}$ centers that maximize their worst-case alignment with \mathcal{P} :

$$S_{p} = \operatorname{arg} \operatorname{topk}_{s \in S'_{p}} (\max_{p \in \mathcal{P}} \cos(s, p), K_{p}).$$

By concentrating the limited budget on those visual tokens most strongly aligned with the key prompt tokens, this strategy ensures a better preservation of the critical regions in the visual input. We determine the appropriate k by ablation to avoid the oversampling of a few salient prompt tokens.

Selection of Visual Center Set S_v . Unlike the prompt center selection, each visual center s_v lies in V. Thereby, we employ *Farthest Point Sampling* (FPS) [36] on the remaining tokens, *i.e.*, $V \setminus S$, to select the visual centers, which makes the visual centers S_v well-spread over V, minimizing the covering radius ϵ_v . Concretely, FPS operates by iteratively selecting the token farthest (*i.e.*, the most different) from the current centers S_v , where the distance is given by

$$\operatorname{dist}_{\mathrm{FPS}}(s_{\mathrm{v}}, \mathcal{S}) = \min_{s \in \mathcal{S}} (1 - \cos(s_{\mathrm{v}}, s)), \quad \forall s_{\mathrm{v}} \in \mathcal{V} \setminus \mathcal{S}.$$

Subsequently, we initialize the visual centers with the empty set, i.e., $\mathcal{S}_{\mathbf{v}}^{(1)} \coloneqq \emptyset$. We then successively add the farthest visual token to the current centers $\mathcal{S}_{\mathbf{v}}^{(i)} \sqcup \mathcal{S}_{\mathbf{p}}$ until it contains a total of K elements. Hence, the visual centers at the subsequent iteration, $\mathcal{S}_{\mathbf{v}}^{(i+1)}$, is given by:

$$\mathcal{S}_{\mathbf{v}}^{(i+1)} = \mathcal{S}_{\mathbf{v}}^{(i)} \sqcup \arg\max_{s_{\mathbf{v}} \in \mathcal{V} \setminus \left(\mathcal{S}_{\mathbf{v}}^{(i)} \sqcup \mathcal{S}_{\mathbf{p}}\right)} \operatorname{dist}_{\mathrm{FPS}}(s_{\mathbf{v}}, \ \mathcal{S}_{\mathbf{v}}^{(i)} \sqcup \mathcal{S}_{\mathbf{p}}), \quad \text{for } i \in [1, \dots, K - K_{\mathbf{p}}].$$

More details of the proposed MoB algorithm are provided in Appendix B.

Theorem 2 (Performance Guarantee). Under Assump 1 and the covering-regularity of Lem 3 with consts $a, a', d_{\text{eff}} > 0$ and b > a, b' > a', for any budget split $(K_p, K - K_p)$, covering fold k, and token set $\mathcal{X} = \mathcal{V} \sqcup \mathcal{P} \subseteq \mathbb{R}^d$ with $|\mathcal{V}| = N$, $|\mathcal{P}| = L$, $d_H(\mathcal{V}, \mathcal{P}) \leq \eta$, the following hold: (a) **Performance bound:** The Performance degradation caused by MoB is upper bounded by

$$\|\mathcal{F}(\mathcal{X}) - \mathcal{F}(\text{MoB}(\mathcal{X}))\| \le C_{\ell} \max \left\{ \alpha(\eta, k, L) \left(K_{\text{p}} \right)^{-1/d_{\text{eff}}}, \ \beta \left(K - K_{\text{p}} \right)^{-1/d_{\text{eff}}} \right\} + C_{\ell} \eta,$$

where
$$\alpha(\eta, k, L) \; = \; \eta \left(b \, k \, L/a \right)^{1/d_{\rm eff}}, \quad \beta \; = \; 2 \, (b')^{1/d_{\rm eff}}.$$

(b) Multilinear complexity: The complexity of MoB is given by $T_{\text{MoB}} = \mathcal{O}(N(L+K)d)$.

Remark (Coupling Trade-off). Under weak coupling (large α), minimizing the bound requires a larger $K_{\rm p}$. Conversely, under strong coupling (small α), the alignment term decays rapidly, favoring visual preservation (increasing $K-K_{\rm p}$). Specially, under perfect coupling ($\eta=0$), the bound simplifies to $\|\Delta y\| \leq C_\ell \beta (K-K_{\rm p})^{-1/d_{\rm eff}}$, i.e., MoB reduces to pure visual preservation.

Remark (Budget Scaling). As the budget K increases, the preservation term $\beta (K - K_p)^{-1/d_{eff}}$ decays, requiring a corresponding increase in K_p (and thus a reduction in the alignment term) to re-balance the trade-off and further lower the overall error bound.

Remark (Scalability). *MoB exhibits a multilinear scalability* w.r.t #visual tokens N, #prompt tokens L, and #retained tokens K (K, $L \ll N$), making it easily adaptable to more challenging scenarios involving large token counts, e.g., higher-resolution inputs or multi-frame video.

Proof in Appendix E.5.

Method	Ohioativaa	Strong Coupling				Weak Coupling						
Method	Objectives	MMB	$\mathrm{MMB}_{\mathrm{CN}}$	SQA	VizWiz	GQA	MME	POPE	VQA^{T}	$\mathrm{VQA}^{\!\mathrm{V2}}$	OCR	Avg.
LLaVA-1.5-7B Vanilla [28]	-	64.7	w/o P 58.1	runing, 69.5	N = 57	6; Toke	n Reduc 1862	tion Rat 85.9	e = 0.0% 58.2	78.5	297	100%
LLaVA-1.5-7B			Pruning	g budge	t K = 19			ction Ra				
FastV (ECCV'24) [8]	VP	61.2	57.0	67.3	50.8	52.7	1612	64.8	52.5	67.1	291	91.2%
SparseVLM (ICML'25) [58]	PA	62.5	53.7	69.1	50.5	57.6	1721	83.6	56.1	75.6	292	96.3%
MustDrop (24.11) [31]	PA VP	62.3	55.8	69.2	51.4	58.2	1787	82.6	56.5	76.0	289	97.2%
DART (EMNLP'25) [46]	VP	63.6	57.0	69.8	51.2	60.0	1856	82.8	57.4	76.7	296	98.8%
MoB (w/o η -prior)	PA VP	63.8 64.1	57.5 57.8	70.0 70.1	52.4 52.5	61.2	1858 1860	84.5 84.8	58.2 58.5	77.9 78.3	304 307	100.2% 100.6%
+ η-prior		04.1	37.6	70.1	32.3	01.4	1800	04.0	36.3	16.5	307	100.0%
LLaVA-1.5-7B					t K = 12							
FastV (ECCV'24)	VP	56.1	56.4	60.2	51.3	49.6	1490	59.6	50.6	61.8	285	86.4%
SparseVLM (ICML'25)	PA	60.0	51.1	67.1	51.4	56.0	1696	80.5	54.9	73.8	280	93.8%
MustDrop (24.11)	PA VP	61.1	55.2	68.5	52.1	56.9	1745	78.7	56.3	74.6	281	95.6%
DART (EMNLP'25)	VP	63.2	57.5 57.3	69.1 69.3	51.7 52.8	58.7 60.7	1840	80.1 81.7	56.4 57.5	75.9 77.2	296	98.0% 99.2%
MoB (w/o η-prior) + η-prior	PA VP	63.2 63.5	57.5	69.6	52.8	60.7	1842 1845	81.7	57.8	77.5	299 299	99.2%
		05.5									277	22.4 /b
LLaVA-1.5-7B					et $K=6$							
FastV (ECCV'24)	VP	48.0	52.7	51.1	50.8	46.1	1256	48.0	47.8	55.0	245	77.3%
SparseVLM (ICML'25)	PA	56.2	46.1	62.2	50.1	52.7	1505	75.1	51.8	68.2	180	84.6%
MustDrop (24.11) DART (EMNLP'25)	PA VP VP	60.0 60.6	53.1 53.2	63.4	51.2 51.6	53.1 55.9	1612 1765	68.0 73.9	54.2 54.4	69.3 72.4	267 270	90.1% 93.7%
MoB (w/o η -prior)	PA VP	61.7	54.2	69.7	52.0	59.0	1806	77.2	57.0	75.5	277	96.3%
$+ \eta$ -prior	ra vr	62.1	54.5	69.8	52.1	59.0	1806	77.2	57.0	75.5	277	96.4%
- 17 piloi		02.1	5 1.5	07.0	32.1	37.0	1000	77.2	37.0	15.5	2//	70.170
LLaVA-Next-7B			w/o Pi	uning,	N = 288	80; <i>Tok</i>	en Redu	ction Ra	te = 0.0%	,		
Vanilla [29]	-	67.4	60.6	70.1	57.6	64.2	1851	86.5	64.9	81.8	517	100%
LLaVA-Next-7B			Pruning	hudoe	t K = 3	20: <i>Toki</i>	en Redu	ction Ra	$t_e = 88.9$	%		
FastV (ECCV'24)	VP	61.6	51.9	62.8	53.1	55.9	1661	71.7	55.7	71.9	374	86.4%
SparseVLM (ICML'25)	PA	60.6	54.5	66.1	52.0	56.1	1533	82.4	58.4	71.5	270	85.9%
MustDrop (24.11)	PA VP	62.8	55.1	68.0	54.0	57.3	1641	82.1	59.9	73.7	382	90.4%
FasterVLM (24.12) [57]	VP	61.6	53.5	66.5	52.6	56.9	1701	83.6	56.5	74.0	401	89.8%
DART (EMNLP'25)	VP	65.3	58.2	68.4	56.1	61.7	1710	84.1	58.7	79.1	406	93.9%
MoB (with η -prior)	PA VP	65.8	58.9	68.7	57.0	62.6	1760	84.4	60.2	80.1	418	95.4%

Table 2: Partial comparison of image understanding on the LLaVA-7B series. For MoB, we set $K_{\rm p} \in \{64,48,32\}$ and $k \in \{4,6,8\}$, corresponding to token-reduction rates of $\{88.9\%,77.8\%,66.7\%\}$. For MoB with the η prior, we use $K_{\rm p} \in \{\frac{3K}{8},\frac{K}{4},\frac{K}{4}\}$ with $k=\frac{3K_{\rm p}}{40}$ for strong-coupling benchmarks and $K_{\rm p} \in \{\frac{K}{2},\frac{7K}{16},\frac{5K}{12}\}$ with $k=\frac{K_{\rm p}}{8}$ for weak-coupling benchmarks, corresponding to the same token-reduction rates; the pruning layer is fixed at $\ell=2$. Blue and Orange denote the best and the second. See Appendix C.4 for the detailed setting, and see Appendix D.1 for the full results.

4 Experimental Results

Experiment Setting. We perform a comprehensive evaluation of the proposed MoB and several representative methods on two visual tasks: image understanding and visual understanding, together with an efficiency analysis. Our experiments employ four popular MLLMs and include a total of 14 widely recognized benchmarks. For further details regarding the benchmarks, models, baselines, and implement details please refer to Appendix C.

Image Understanding. Table 2 and Table 3 report the evaluation results across a variety of image-understanding tasks on LLaVA series and Qwen2-VL, respectively. We highlight five key observations: (a) MoB consistently outperforms all base-

Method	GQA	MME	POPE	$\mathbf{V}\mathbf{Q}\mathbf{A}^{\mathrm{T}}$	MMB	SQA	Avg.
Qwen2-VL-7B Vanilla [44]			n g; Toke 86.1	n Reduc 82.1			.0% 100%
Owen2-VL-7B		Toke	en Redu	ction Ra	te = 66 .	7%	_
FastV	58.0	2130	82.1	77.3	76.1	80.0	94.0%
DART	60.2	2245	83.9	80.5	78.9	81.4	97.0%
MoB (with η)	61.8	2268	84.7	81.1	79.5	82.3	98.4%
Owen2-VL-7B		Toke	n Redu	ction Ra	te = 77.	8%	
FastV	56.7	2031	79.2	72.0	74.1	78.3	91.0%
DART	58.5	2175	82.1	75.3	77.3	79.6	94.3%
MoB (with η)	59.4	2203	82.8	75.8	78.1	80.4	95.2%
Qwen2-VL-7B		Toke	n Redu	ction Ra	te = 88 .	9%	
FastV	51.9	1962	76.1	60.3	70.1	75.8	84.4%
DART	55.5	2052	77.9	61.8	72.0	77.6	87.4%
MoB (with η)	56.5	2094	78.5	62.7	72.8	78.4	88.6%

Table 3: Comparative experiments on image understanding with Owen2-VL-7B.

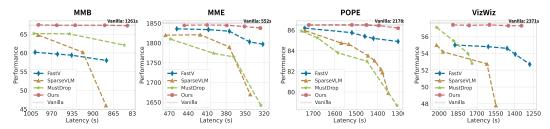


Figure 3: Performance-Latency trade-off comparisons across four benchmarks on LLaVA-Next-7B.

lines on LLaVA-1.5-7B in most cases. This will be more pronounced when incorporating the η -prior, which highlights the inherent advantage of our approach; (b) single-objective baselines exhibit complementary strengths under different coupling patterns, whereas MoB consistently outperforms all baselines, demonstrating the benefit of balanced objectives; (c) the superiority of MoB becomes even more significant under aggressive token reduction. Specifically, the improvement of MoB over the best baseline in average scores increases from 1.8% at a 66.7% token reduction to 2.7% at an 88.8% reduction on LLaVA-1.5-7B; (d) MoB matches the performance of the vanilla LLaVA-1.5-7B with only 33.3% of visual tokens, which may be attributed to the mitigation of hallucinations caused by redundant tokens; and (e) MoB scales seamlessly to advanced models, preserving 95.2% performance on Qwen2-VL-7B using only 22.2% of visual tokens. These observations demonstrate the superiority of MoB in leveraging limited visual tokens while minimizing performance degradation.

Video Understanding. As presented in Table 4, MoB is general and can be readily extended to more challenging video scenarios without incurring additional cost. Specifically, MoB preserves 97.9% of average performance for Video-LLaVA-7B using only 6.6% of the visual tokens, which sets new records in most VideoQA benchmarks, achieving 1.6% and 4.7% improvements over TwigVLM and VisionZip, respectively. These results validate the generalization ability of MoB.

Method	TGIF	MSVD	MSRV	ActNet	Avg.
Video-LLaVA-7B	To	ken Red	uction R	ate = 0. 0)%
Vanilla [27]	47.1	69.8	56.7	43.1	100%
Video-LLaVA-7B	To	ken Redu	ction Re	te = 93.	4%
FastV (ECCV'24)	23.1	38.0	19.3	30.6	52.1%
SparseVLM (ICML'25)	44.7	68.2	31.0	42.6	86.5%
VisionZip (24.12) [51]	42.4	63.5	52.1	43.0	93.2%
TwigVLM (ICCV'25) [39]	44.7	68.3	54.6	41.5	96.3%
MoB (with η -prior)	45.3	68.8	55.2	42.8	97.9%

Table 4: Comparative experiments on video understanding with Video-LLaVA-7B.

Efficiency Analysis. We present a performance-latency trade-off measured on an NVIDIA A800-80GB GPU in Figure 3. The results show that (a) MoB achieves a strong performance-latency trade-off, delivering a $1.3-1.5\times$ speed-up for LLaVA-NEXT-7B with negligible performance loss; (b) due to ignoring the $K-\eta$ trade-off, the multi-stage method MustDrop is outperformed by single-objective methods FastV and SparseVLM on MME and POPE, and suffers significant performance drops as token budgets shrink (*i.e.*, latency decreases). In contrast, MoB consistently maintains a robust trade-off across all benchmarks, surpassing all the baselines by a clear margin; (c) MoB does not rely on attention scores to identify important tokens, making it compatible with flash attention and more efficient than attention-based methods such as SparseVLM and FastV.

5 Ablation and Discussion

Impact of $\langle K, \eta, K_{\rm p}, \rangle$. We study the impact of K, η , and $K_{\rm p}$ on pruning performance across four benchmarks: GQA and TextVQA (weak coupling); VizWiz and MMB (strong coupling). As shown in Figure 4, the results can be interpreted by Theorem 1 and Theorem 2(a), respectively.

A. Theorem 1 Perspective: When K is large, e.g., K=192, the trade-off is governed by $D_1K^{-2/d_{\rm eff}}$, hence the trade-off intensity remains nearly identical across benchmarks. Conversely, When K is small, especially K=64, in weak-coupling benchmarks, the trade-off turns to be governed by $D_2\eta^2$; thus, the trade-off intensity is obviously more pronounced in GQA and TextVQA than that in VizWiz and MMB. These observations exactly confirm the validity of Theorem 1.

B. Theorem 2(a) Perspective. (a) Under weak coupling, the alignment term $\alpha(\eta,k,L)(K_{\rm p})^{-1/d_{\rm eff}}$ is amplified, which requires a larger $K_{\rm p}$ to suppress the overall error. However, across benchmarks sharing the same coupling pattern, the optimal $K_{\rm p}$ values exhibit only minor variation. (b) Increasing the total budget K pushes the optimal $K_{\rm p}$ upward to rebalance the two bound terms. Since the prompt length L is fixed, adding more tokens yields diminishing returns for prompt alignment, which is reflected in the declining ratio $K_{\rm p}/K$. These validate the performance bound in Theorem 2(a).

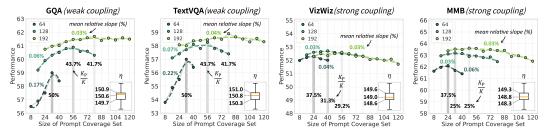


Figure 4: Comprehensive ablation on the budget configuration $\langle K_{\rm p},K\rangle$ across four benchmarks with distinct prompt-visual coupling η on LLaVA-1.5-7B, where $K=\{64,128,192\}$; the *mean relative slope* (%) is given by $\frac{100}{x_n-x_1}\sum_{i=1}^{n-1}\frac{y_{i+1}-y_i}{y_i}$, quantifying the trade-off intensity; the ratio $\frac{K_{\rm p}}{K}$ reflects the cost-effectiveness of prompt alignment, and the box plot presents the distribution of η .

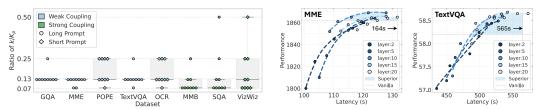


Figure 5: Ablation on the ratio of k/K_p .

Figure 6: Ablation on the pruning layer.

Remarkably, the experimental results suggest that simply determining the optimal K_p for each of the two coupling patterns suffices to guarantee effective generalization across all scenarios.

Impact of Covering Fold k. We chose the covering fold k by examining the normalized ratio $k/K_{\rm p}$ across eight benchmarks and nine budget configurations. As shown in Figure 5, (a) weak-coupling benchmarks generally require a larger k to ensure critical region coverage, whereas strong-coupling settings suffice with a smaller k; (b) benchmarks with longer prompts impose a lower cap on k to preserve sampling diversity and avoid redundant selection of salient tokens. Notably, weak-coupling benchmarks with long prompts (e.g., GQA, TextVQA) exhibit a narrowly clustered optimal $k/K_{\rm p}$ range, reflecting their strict requirement to cover key tokens without excessive redundancy.

Impact of Pruning Layer. As shown in Figure 6, (a) models with visual token pruning consistently achieve a more favorable performance-efficiency trade-off than the vanilla model on both benchmarks. (b) Pruning in deeper layers provides more significant benefits for the weak-coupling TextVQA than strong-coupling MME. We attribute this to stronger cross-modal interactions in deeper MLLM layers, which facilitate identification of answer-relevant tokens under weak coupling, whereas pruning in shallow layers disrupts these interactions and incurs greater performance degradation.

6 Conclusion

In this paper, we present a comprehensive analysis of visual token pruning, deriving the first closed-form error bound with a practical relaxation. Leveraging ϵ -covering theory, we quantify the intrinsic trade-off between the fundamental pruning objectives, *i.e.*, visual preservation and prompt alignment, and identify their optimal attainment levels under a fixed pruning budget. Building on these insights, we introduce MoB, a training-free algorithm for visual token pruning. Based on greedy radius trading, MoB ensures the near-optimal attainment per objective via budget allocation, offering a provable performance bound and multilinear scalability. Experimental results indicate that MoB matches the performance (100.6%) of LLaVA-1.5-7B with only 33.3% of visual tokens and can be seamlessly integrated into advanced MLLMs, such as LLaVA-Next-7B and Qwen2-VL-7B. Our work advances the understanding of visual token pruning and offers valuable insights for future MLLM compression.

Limitations. Our theoretical guarantees rely on assumption 1, which is generally satisfied in practice but may not hold for all MLLMs. Besides, MoB applies a preliminary search to select the proper $K_{\rm p}$, which potentially introduces extra tuning overhead in practical applications. Future work will focus on developing an adaptive $K_{\rm p}$ selection mechanism driven by online estimation of the coupling η .

Acknowledgments

The work was performed at the Shanghai Key Laboratory of Multidimensional Information Processing, East China Normal University; the Institute of Natural Sciences and School of Mathematical Sciences, Shanghai Jiao Tong University; and the Chongqing Key Laboratory of Precision Optics, Chongqing Institute of East China Normal University, with joint support from the National Natural Science Foundation of China (62176091), the Natural Science Foundation of Chongqing (CSTB2024NSCQ-MSX0877), the Science and Technology Commission of Shanghai Municipality (21DZ2203100) and the Fundamental Research Funds for the Central Universities.

References

- [1] Yash Akhauri, Ahmed F AbouElhamayed, Yifei Gao, Chi-Chih Chang, Nilesh Jain, and Mohamed S Abdelfattah. Tokenbutler: Token importance is predictable. *arXiv preprint arXiv:2503.07518*, 2025.
- [2] Kazi Hasan Ibn Arif, JinYi Yoon, Dimitrios S Nikolopoulos, Hans Vandierendonck, Deepu John, and Bo Ji. Hired: Attention-guided token dropping for efficient inference of high-resolution vision-language models in resource-constrained environments. arXiv preprint arXiv:2408.10945, 2024.
- [3] Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E Hinton. Layer normalization. *arXiv preprint* arXiv:1607.06450, 2016.
- [4] Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Sibo Song, Kai Dang, Peng Wang, Shijie Wang, Jun Tang, et al. Qwen2. 5-vl technical report. *arXiv preprint arXiv:2502.13923*, 2025.
- [5] Kaichen Zhang* Fanyi Pu* Xinrun Du Yuhao Dong Haotian Liu Yuanhan Zhang Ge Zhang Chunyuan Li Bo Li*, Peiyuan Zhang* and Ziwei Liu. Lmms-eval: Accelerating the development of large multimoal models, March 2024.
- [6] Daniel Bolya, Cheng-Yang Fu, Xiaoliang Dai, Peizhao Zhang, Christoph Feichtenhofer, and Judy Hoffman. Token merging: Your vit but faster. In *The Eleventh International Conference on Learning Representations*, 2022.
- [7] David Chen and William B Dolan. Collecting highly parallel data for paraphrase evaluation. In *Proceedings of the 49th annual meeting of the association for computational linguistics: human language technologies*, pages 190–200, 2011.
- [8] 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*, pages 19–35. Springer, 2024.
- [9] Liangyu Chen, Bo Li, Sheng Shen, Jingkang Yang, Chunyuan Li, Kurt Keutzer, Trevor Darrell, and Ziwei Liu. Large language models are visual reasoning coordinators. *Advances in Neural Information Processing Systems*, 36:70115–70140, 2023.
- [10] Zhe Chen, Weiyun Wang, Yue Cao, Yangzhou Liu, Zhangwei Gao, Erfei Cui, Jinguo Zhu, Shenglong Ye, Hao Tian, Zhaoyang Liu, et al. Expanding performance boundaries of open-source multimodal models with model, data, and test-time scaling. *arXiv preprint arXiv:2412.05271*, 2024.
- [11] Zhe Chen, Jiannan Wu, Wenhai Wang, Weijie Su, Guo Chen, Sen Xing, Muyan Zhong, Qinglong Zhang, Xizhou Zhu, Lewei Lu, et al. Internvl: Scaling up vision foundation models and aligning for generic visual-linguistic tasks. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 24185–24198, 2024.
- [12] 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*, pages 6904–6913, 2017.

- [13] Zhiyu Guo, Hidetaka Kamigaito, and Taro Watanabe. Attention score is not all you need for token importance indicator in kv cache reduction: Value also matters. *arXiv preprint arXiv:2406.12335*, 2024.
- [14] 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*, pages 3608–3617, 2018.
- [15] Dorit S Hochbaum and David B Shmoys. A best possible heuristic for the k-center problem. *Mathematics of operations research*, 10(2):180–184, 1985.
- [16] 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, pages 2256–2264, 2024.
- [17] 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*, pages 6700–6709, 2019.
- [18] Yunseok Jang, Yale Song, Youngjae Yu, Youngjin Kim, and Gunhee Kim. Tgif-qa: Toward spatio-temporal reasoning in visual question answering. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2758–2766, 2017.
- [19] Kasper Green Larsen and Jelani Nelson. The johnson-lindenstrauss lemma is optimal for linear dimensionality reduction. *arXiv* preprint arXiv:1411.2404, 2014.
- [20] Bo Li, Yuanhan Zhang, Dong Guo, Renrui Zhang, Feng Li, Hao Zhang, Kaichen Zhang, Peiyuan Zhang, Yanwei Li, Ziwei Liu, et al. Llava-onevision: Easy visual task transfer. arXiv preprint arXiv:2408.03326, 2024.
- [21] Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models. In *International conference on machine learning*, pages 19730–19742. PMLR, 2023.
- [22] Yanwei Li, Yuechen Zhang, Chengyao Wang, Zhisheng Zhong, Yixin Chen, Ruihang Chu, Shaoteng Liu, and Jiaya Jia. Mini-gemini: Mining the potential of multi-modality vision language models. *arXiv preprint arXiv:2403.18814*, 2024.
- [23] 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.
- [24] Zhang Li, Biao Yang, Qiang Liu, Zhiyin Ma, Shuo Zhang, Jingxu Yang, Yabo Sun, Yuliang Liu, and Xiang Bai. Monkey: Image resolution and text label are important things for large multi-modal models. In *proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 26763–26773, 2024.
- [25] Youwei Liang, GE Chongjian, Zhan Tong, Yibing Song, Jue Wang, and Pengtao Xie. Evit: Expediting vision transformers via token reorganizations. In *International Conference on Learning Representations*, 2022.
- [26] Zijing Liang, Yanjie Xu, Yifan Hong, Penghui Shang, Qi Wang, Qiang Fu, and Ke Liu. A survey of multimodel large language models. In *Proceedings of the 3rd International Conference on Computer, Artificial Intelligence and Control Engineering*, pages 405–409, 2024.
- [27] Bin Lin, Yang Ye, Bin Zhu, Jiaxi Cui, Munan Ning, Peng Jin, and Li Yuan. Video-llava: Learning united visual representation by alignment before projection, 2024.
- [28] Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. Improved baselines with visual instruction tuning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 26296–26306, 2024.
- [29] Haotian Liu, Chunyuan Li, Yuheng Li, Bo Li, Yuanhan Zhang, Sheng Shen, and Yong Jae Lee. Llava-next: Improved reasoning, ocr, and world knowledge, January 2024.

- [30] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning, 2023.
- [31] Ting Liu, Liangtao Shi, Richang Hong, Yue Hu, Quanjun Yin, and Linfeng Zhang. Multi-stage vision token dropping: Towards efficient multimodal large language model. arXiv preprint arXiv:2411.10803, 2024.
- [32] Xuyang Liu, Ziming Wang, Yuhang Han, Yingyao Wang, Jiale Yuan, Jun Song, Bo Zheng, Linfeng Zhang, Siteng Huang, and Honggang Chen. Compression with global guidance: Towards training-free high-resolution mllms acceleration. arXiv preprint arXiv:2501.05179, 2025.
- [33] Yuliang Liu, Zhang Li, Mingxin Huang, Biao Yang, Wenwen Yu, Chunyuan Li, Xu-Cheng Yin, Cheng-Lin Liu, Lianwen Jin, and Xiang Bai. Ocrbench: on the hidden mystery of ocr in large multimodal models. *Science China Information Sciences*, 67(12):220102, 2024.
- [34] Jiasen Lu, Dhruv Batra, Devi Parikh, and Stefan Lee. Vilbert: Pretraining task-agnostic visiolinguistic representations for vision-and-language tasks. *Advances in neural information processing systems*, 32, 2019.
- [35] Pan Lu, Swaroop Mishra, Tony Xia, Liang Qiu, Kai-Wei Chang, Song-Chun Zhu, Oyvind Tafjord, Peter Clark, and Ashwin Kalyan. Learn to explain: Multimodal reasoning via thought chains for science question answering. In *The 36th Conference on Neural Information Process*ing Systems (NeurIPS), 2022.
- [36] Carsten Moenning and Neil A Dodgson. Fast marching farthest point sampling. Technical report, University of Cambridge, Computer Laboratory, 2003.
- [37] Yingzhe Peng, Xinting Hu, Jiawei Peng, Xin Geng, Xu Yang, et al. Live: Learnable in-context vector for visual question answering. Advances in Neural Information Processing Systems, 37:9773–9800, 2024.
- [38] Yuzhang Shang, Mu Cai, Bingxin Xu, Yong Jae Lee, and Yan Yan. Llava-prumerge: Adaptive token reduction for efficient large multimodal models. *arXiv preprint arXiv:2403.15388*, 2024.
- [39] Zhenwei Shao, Mingyang Wang, Zhou Yu, Wenwen Pan, Yan Yang, Tao Wei, Hongyuan Zhang, Ning Mao, Wei Chen, and Jun Yu. Growing a twig to accelerate large vision-language models. *arXiv preprint arXiv:2503.14075*, 2025.
- [40] Zhenwei Shao, Zhou Yu, Meng Wang, and Jun Yu. Prompting large language models with answer heuristics for knowledge-based visual question answering. In *Proceedings of the IEEE/CVF Conference on computer vision and pattern recognition*, pages 14974–14983, 2023.
- [41] Amanpreet Singh, Vivek Natarajan, Meet Shah, Yu Jiang, Xinlei Chen, Dhruv Batra, Devi Parikh, and Marcus Rohrbach. Towards vqa models that can read. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 8317–8326, 2019.
- [42] Xudong Tan, Peng Ye, Chongjun Tu, Jianjian Cao, Yaoxin Yang, Lin Zhang, Dongzhan Zhou, and Tao Chen. Tokencarve: Information-preserving visual token compression in multimodal large language models. *arXiv* preprint arXiv:2503.10501, 2025.
- [43] 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.
- [44] Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhihao Fan, Jinze Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, et al. Qwen2-vl: Enhancing vision-language model's perception of the world at any resolution. *arXiv preprint arXiv:2409.12191*, 2024.
- [45] Yiqi Wang, Wentao Chen, Xiaotian Han, Xudong Lin, Haiteng Zhao, Yongfei Liu, Bohan Zhai, Jianbo Yuan, Quanzeng You, and Hongxia Yang. Exploring the reasoning abilities of multimodal large language models (mllms): A comprehensive survey on emerging trends in multimodal reasoning. *arXiv preprint arXiv:2401.06805*, 2024.

- [46] Zichen Wen, Yifeng Gao, Shaobo Wang, Junyuan Zhang, Qintong Zhang, Weijia Li, Conghui He, and Linfeng Zhang. Stop looking for important tokens in multimodal language models: Duplication matters more. *arXiv preprint arXiv:2502.11494*, 2025.
- [47] 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.
- [48] Long Xing, Qidong Huang, Xiaoyi Dong, Jiajie Lu, Pan Zhang, Yuhang Zang, Yuhang Cao, Conghui He, Jiaqi Wang, Feng Wu, et al. Pyramiddrop: Accelerating your large vision-language models via pyramid visual redundancy reduction. *arXiv preprint arXiv:2410.17247*, 2024.
- [49] Dejing Xu, Zhou Zhao, Jun Xiao, Fei Wu, Hanwang Zhang, Xiangnan He, and Yueting Zhuang. Video question answering via gradually refined attention over appearance and motion. In *Proceedings of the 25th ACM international conference on Multimedia*, pages 1645–1653, 2017.
- [50] Jingjing Xu, Xu Sun, Zhiyuan Zhang, Guangxiang Zhao, and Junyang Lin. Understanding and improving layer normalization. *Advances in neural information processing systems*, 32, 2019.
- [51] Senqiao Yang, Yukang Chen, Zhuotao Tian, Chengyao Wang, Jingyao Li, Bei Yu, and Jiaya Jia. Visionzip: Longer is better but not necessary in vision language models. *arXiv preprint arXiv:2412.04467*, 2024.
- [52] 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. *arXiv preprint arXiv:2409.10197*, 2024.
- [53] Zhou Yu, Dejing Xu, Jun Yu, Ting Yu, Zhou Zhao, Yueting Zhuang, and Dacheng Tao. Activitynet-qa: A dataset for understanding complex web videos via question answering. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pages 9127–9134, 2019.
- [54] Yuanhan Zhang Bo Li Songyang Zhang Wangbo Zhao Yike Yuan Jiaqi Wang Conghui He Ziwei Liu Kai Chen Dahua Lin Yuan Liu, Haodong Duan. Mmbench: Is your multi-modal model an all-around player? *arXiv:2307.06281*, 2023.
- [55] Zheng Zhan, Yushu Wu, Zhenglun Kong, Changdi Yang, Yifan Gong, Xuan Shen, Xue Lin, Pu Zhao, and Yanzhi Wang. Rethinking token reduction for state space models. In *Proceedings* of the 2024 Conference on Empirical Methods in Natural Language Processing, pages 1686– 1697, 2024.
- [56] Kaichen Zhang, Bo Li, Peiyuan Zhang, Fanyi Pu, Joshua Adrian Cahyono, Kairui Hu, Shuai Liu, Yuanhan Zhang, Jingkang Yang, Chunyuan Li, and Ziwei Liu. Lmms-eval: Reality check on the evaluation of large multimodal models, 2024.
- [57] Qizhe Zhang, Aosong Cheng, Ming Lu, Zhiyong Zhuo, Minqi Wang, Jiajun Cao, Shaobo Guo, Qi She, and Shanghang Zhang. [cls] attention is all you need for training-free visual token pruning: Make vlm inference faster. *arXiv preprint arXiv:2412.01818*, 2024.
- [58] Yuan Zhang, Chun-Kai Fan, Junpeng Ma, Wenzhao Zheng, Tao Huang, Kuan Cheng, Denis Gudovskiy, Tomoyuki Okuno, Yohei Nakata, Kurt Keutzer, et al. Sparsevlm: Visual token sparsification for efficient vision-language model inference. *arXiv preprint arXiv:2410.04417*, 2024.
- [59] Yiwu Zhong, Zhuoming Liu, Yin Li, and Liwei Wang. Aim: Adaptive inference of multi-modal llms via token merging and pruning. *arXiv preprint arXiv:2412.03248*, 2024.
- [60] Deyao Zhu, Jun Chen, Xiaoqian Shen, Xiang Li, and Mohamed Elhoseiny. Minigpt-4: Enhancing vision-language understanding with advanced large language models. *arXiv* preprint *arXiv*:2304.10592, 2023.

NeurIPS Paper Checklist

1. Claims

Question: Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope?

Answer: [Yes]

Justification: See Section 1 for the main claims; see Sections 3 to 5 and Appendices D and E for the detailed contributions and scope.

Guidelines:

- The answer NA means that the abstract and introduction do not include the claims made in the paper.
- The abstract and/or introduction should clearly state the claims made, including the contributions made in the paper and important assumptions and limitations. A No or NA answer to this question will not be perceived well by the reviewers.
- The claims made should match theoretical and experimental results, and reflect how much the results can be expected to generalize to other settings.
- It is fine to include aspirational goals as motivation as long as it is clear that these goals are not attained by the paper.

2. Limitations

Question: Does the paper discuss the limitations of the work performed by the authors?

Answer: [Yes]

Justification: See Section 6 for the discussion on the limitations.

Guidelines:

- The answer NA means that the paper has no limitation while the answer No means that the paper has limitations, but those are not discussed in the paper.
- The authors are encouraged to create a separate "Limitations" section in their paper.
- The paper should point out any strong assumptions and how robust the results are to violations of these assumptions (e.g., independence assumptions, noiseless settings, model well-specification, asymptotic approximations only holding locally). The authors should reflect on how these assumptions might be violated in practice and what the implications would be.
- The authors should reflect on the scope of the claims made, e.g., if the approach was only tested on a few datasets or with a few runs. In general, empirical results often depend on implicit assumptions, which should be articulated.
- The authors should reflect on the factors that influence the performance of the approach. For example, a facial recognition algorithm may perform poorly when image resolution is low or images are taken in low lighting. Or a speech-to-text system might not be used reliably to provide closed captions for online lectures because it fails to handle technical jargon.
- The authors should discuss the computational efficiency of the proposed algorithms and how they scale with dataset size.
- If applicable, the authors should discuss possible limitations of their approach to address problems of privacy and fairness.
- While the authors might fear that complete honesty about limitations might be used by reviewers as grounds for rejection, a worse outcome might be that reviewers discover limitations that aren't acknowledged in the paper. The authors should use their best judgment and recognize that individual actions in favor of transparency play an important role in developing norms that preserve the integrity of the community. Reviewers will be specifically instructed to not penalize honesty concerning limitations.

3. Theory assumptions and proofs

Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof?

Answer: [Yes]

Justification: The assumptions are clearly stated in the statement of the theorems: any partial composition \mathcal{F} of the language model satisfies a Lipschitz Continuity w.r.t. the Hausdorff Distance with constant $C_\ell \geq 1$, the δ -dilations $\mathcal{V}_\delta \coloneqq \bigcup_{v \in \mathcal{V}} B(v, \delta)$, $\mathcal{P}_\delta \coloneqq \bigcup_{p \in \mathcal{P}} B(p, \delta)$ ($\delta \ll \eta$) satisfy d_{eff} -dimensional covering regularity with constants a, a', b, b'. See Theorems 1 and 2. The proof of Lemma 1 can be found in Appendix E.1, the proof of Lemma 2 can be found in Appendix E.2, the proof of Lemma 3 can be found in Appendix E.3, the proof of Theorems 1 and 2 are provided in Appendices E.4 and E.5, respectively.

Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and crossreferenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

4. Experimental result reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [Yes]

Justification: See Appendices B and C for the information needed to reproduce the main experimental results.

Guidelines:

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
- (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
- (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
- (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
- (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility.

In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

5. Open access to data and code

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

Answer: [Yes]

Justification: The code is provided in the supplemental material.

Guidelines:

- The answer NA means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (https://nips.cc/public/guides/CodeSubmissionPolicy) for more details.
- While we encourage the release of code and data, we understand that this might not be possible, so "No" is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (https://nips.cc/public/guides/CodeSubmissionPolicy) for more details.
- The authors should provide instructions on data access and preparation, including how
 to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

6. Experimental setting/details

Question: Does the paper specify all the training and test details (e.g., data splits, hyperparameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

Answer: [Yes]

Justification: The experiment does not involve any training process. See Section 5 and appendix C for all the test details.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- The full details can be provided either with the code, in appendix, or as supplemental material.

7. Experiment statistical significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [Yes]

Justification: The experimental results are accompanied by significance tests, and cross-validation conducted using a publicly available third-party framework.

Guidelines:

• The answer NA means that the paper does not include experiments.

- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).
- It should be clear whether the error bar is the standard deviation or the standard error of the mean.
- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
- For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
- If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

8. Experiments compute resources

Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

Answer: [Yes]

Justification: See Appendix C for the computer resources needed to reproduce the experiments.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
- The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
- The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).

9. Code of ethics

Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics https://neurips.cc/public/EthicsGuidelines?

Answer: [Yes]

Justification: The research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics.

Guidelines:

- The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.
- The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).

10. **Broader impacts**

Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?

Answer: [Yes]

Justification: See Section 6 and appendix A for the discussion on the broader impacts.

Guidelines:

- The answer NA means that there is no societal impact of the work performed.
- If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.
- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

11. Safeguards

Question: Does the paper describe safeguards that have been put in place for responsible release of data or models that have a high risk for misuse (e.g., pretrained language models, image generators, or scraped datasets)?

Answer: [NA]

Justification: The paper poses no such risks.

Guidelines:

- The answer NA means that the paper poses no such risks.
- Released models that have a high risk for misuse or dual-use should be released with
 necessary safeguards to allow for controlled use of the model, for example by requiring
 that users adhere to usage guidelines or restrictions to access the model or implementing
 safety filters.
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

12. Licenses for existing assets

Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?

Answer: [Yes]

Justification: The assets used in this paper are properly credited, and the license and terms are respected.

Guidelines:

- The answer NA means that the paper does not use existing assets.
- The authors should cite the original paper that produced the code package or dataset.
- The authors should state which version of the asset is used and, if possible, include a URL.
- The name of the license (e.g., CC-BY 4.0) should be included for each asset.

- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
- If assets are released, the license, copyright information, and terms of use in the
 package should be provided. For popular datasets, paperswithcode.com/datasets
 has curated licenses for some datasets. Their licensing guide can help determine the
 license of a dataset.
- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
- If this information is not available online, the authors are encouraged to reach out to the asset's creators.

13. New assets

Question: Are new assets introduced in the paper well documented and is the documentation provided alongside the assets?

Answer: [NA]

Justification: The paper does not release new assets.

Guidelines:

- The answer NA means that the paper does not release new assets.
- Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
- The paper should discuss whether and how consent was obtained from people whose asset is used.
- At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

14. Crowdsourcing and research with human subjects

Question: For crowdsourcing experiments and research with human subjects, does the paper include the full text of instructions given to participants and screenshots, if applicable, as well as details about compensation (if any)?

Answer: [NA]

Justification: The paper does not involve crowdsourcing nor research with human subjects Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
- According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

15. Institutional review board (IRB) approvals or equivalent for research with human subjects

Question: Does the paper describe potential risks incurred by study participants, whether such risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or an equivalent approval/review based on the requirements of your country or institution) were obtained?

Answer: [NA]

Justification: The paper does not involve crowdsourcing nor research with human subjects Guidelines:

 The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.

- Depending on the country in which research is conducted, IRB approval (or equivalent)
 may be required for any human subjects research. If you obtained IRB approval, you
 should clearly state this in the paper.
- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.

16. Declaration of LLM usage

Question: Does the paper describe the usage of LLMs if it is an important, original, or non-standard component of the core methods in this research? Note that if the LLM is used only for writing, editing, or formatting purposes and does not impact the core methodology, scientific rigorousness, or originality of the research, declaration is not required.

Answer: [Yes]

Justification: MoB is an *inference-time* algorithm that prunes visual tokens in pre-trained multimodal LLMs—LLaVA-1.5-7B, LLaVA-Next-7B, Qwen2-VL-7B and Video-LLaVA-7B. At Transformer layer $\ell=2$, it removes redundant vision tokens using a bi-objective covering rule (see Section 3 and appendices B and C). The LLM **weights remain frozen**; no additional data, gradient updates, or prompt engineering are used. Thus the LLMs serve as essential yet unmodified back-bones whose intermediate embeddings are the input to MoB.

Guidelines:

- The answer NA means that the core method development in this research does not involve LLMs as any important, original, or non-standard components.
- Please refer to our LLM policy (https://neurips.cc/Conferences/2025/LLM) for what should or should not be described.

Appendix

In the appendix, we provide additional information as listed below:

- §A provides the broader impacts of MoB
- §B provides the algorithm details and pseudocode of MoB
- §C provides the overview of the data, models, baselines and implementation details.
- §D provides the additional experimental results.
- §E provides the omitted technical details.

A Broader Impacts and Limitations

Theory Impacts. Beyond the *visual* setting, MoB's theoretical lens—balancing *Visual Preservation* (retaining sufficient context) and *Prompt Alignment* (isolating "golden evidence")—naturally transfers to *language* domain. It makes the key challenging (context vs. evidence) in long-context LLM explicit and offers actionable guidance for token-level compression and scheduling under fixed context budgets. In practice, this perspective informs RAG (calibrating recall vs. precision) and summarization/LLM memory (trading coherence vs. conciseness).

Application Impacts. The proposed MoB yields substantial acceleration of MLLMs with negligible performance loss, thereby enabling high-resolution vision-language models to operate on resource-constrained platforms such as edge devices and mobile systems while supporting low-latency applications—including assistive technologies for the visually impaired, autonomous navigation, and AR/VR. Besides, MoB potentially benefits other redundancy-heavy domains (*e.g.*, point clouds and multi-sensor fusion), guiding efficient token-level compression beyond vision.

Theory Limitations. The theoretical analysis (lemma 1 and theorem 1) and the performance guarantees (theorem 2) rely on assumption 1 (Lipschitz Continuity) and lemma 3 (Covering Regularity). In embedding spaces that violate metric properties or exhibit highly irregular token distributions, these conditions may fail to hold, and the provable performance bounds may no longer apply.

Application Limitations. Our deployment presently requires an *a priori* estimate of η to set the pruning hyperparameters K_p and k. When η is misestimated for a new model, domain, or input distribution, the selected K_p and k can deviate from their optimum, leading to suboptimal speed–accuracy trade-offs.

B Algorithm

Algorithm 1 Multi-Objective Balanced Covering (MoB)

Require: Visual token $\mathbf{V} \in \mathbb{R}^{N \times d}$, Prompt token $\mathbf{P} \in \mathbb{R}^{L \times d}$, Budget K_p, K_v , Covering fold k **Ensure:** Index list for select tokens $\mathbf{S} \in \mathbb{N}^{K_p + K_v}$

1: Normalize all token embeddings to unit ℓ_2 norm: $\mathbf{V} \leftarrow \mathbf{V}/\|\mathbf{V}\|_{2,row}$, $\mathbf{P} \leftarrow \mathbf{P}/\|\mathbf{P}\|_{2,row}$

Step 1. Select Prompt Centers via Nearest-Neighbor Covering

2: Compute cosine-similarity matrix via $\mathbf{P} \mathbf{V}^{\top} : \mathbf{M} \leftarrow \mathbf{P} \mathbf{V}^{\top}$

 $\triangleright \mathbf{M} \in \mathbb{R}^{L \times N}$.

3: Retrieve k nearest token indices per prompt:

$$\mathbf{C}_{\mathtt{idx}} \leftarrow \mathtt{ArgTopK}(\mathbf{M}, \mathtt{k}, \mathtt{axis} = 1), \quad \mathbf{C}_{\mathtt{sim}} \leftarrow \mathtt{TopK}(\mathbf{M}, \mathtt{k}, \mathtt{axis} = 1)$$

 $ho \ C_{idx}, C_{sim} \in \mathbb{R}^{L \times k}$ collects index and similarity of k closest centers per prompt token. # Deduplicate candidate indices

- 4: Flatten index and similarity arrays: $\mathbf{C}_{\mathtt{idx}} \leftarrow \mathtt{Flatten}(\mathbf{C}_{\mathtt{idx}}), \quad \mathbf{C}_{\mathtt{sim}} \leftarrow \mathtt{Flatten}(\mathbf{C}_{\mathtt{sim}}) \quad \triangleright \quad \mathbf{C}_{\mathtt{idx}} \in \mathbb{N}^{\mathtt{Lk}}, \, \mathbf{C}_{\mathtt{sim}} \in \mathbb{R}^{\mathtt{Lk}}$
- 5: Remove duplicate indices, preserving associated similarities:

$$\langle \mathbf{C}^*_{\mathtt{idx}}, \mathbf{C}^*_{\mathtt{sim}} \rangle \leftarrow \mathtt{UniqueIndices}(\mathbf{C}_{\mathtt{idx}}, \, \mathbf{C}_{\mathtt{sim}})$$

- $\triangleright \, \mathtt{K_p} \leq |\mathbf{C}^*_{\mathtt{idx}}| \leq \mathtt{Lk}$
- 6: Identify top- $K_{\rm p}$ prompt centers by similarity: $\mathbf{i_p} \leftarrow \texttt{ArgTopK}(\mathbf{C_{sim}^*}, \mathtt{K_p})$
- 7: Form the prompt-center index list: $\mathbf{S}_{\mathrm{p}} \leftarrow \mathbf{C}^*_{\mathtt{idx}}[\mathbf{i}_{\mathrm{p}}]$

$hd \mathbf{S}_{\mathrm{p}} \in \mathbb{N}^{\mathtt{K}_{\mathrm{p}}}$

Step 2. Select Visual Centers via Farthest-Point Sampling

- 8: Initialize selected centers: S ← S_p
 # Initialize token-to-prompt minimum distances
- 9: Compute pairwise minimum distances between all tokens and selected prompt centers:

$$\mathbf{d} \leftarrow \mathbf{1}_{\texttt{N} \times \texttt{K}_{\texttt{n}}} - \mathbf{V} \, \mathbf{V} [\mathbf{S}_{\texttt{D}}]^\top, \quad \mathbf{d} \leftarrow \texttt{Min}(\mathbf{d}, \, \texttt{axis} = 1)$$

 \triangleright Selected centers have zero distance in $\mathbf{d} \in \mathbb{R}^N$.

Farthest-Point Sampling

- 10: **for** t = 1 to K_{v} **do**
- 11: Select the token farthest from current centers: $i^* \leftarrow \text{ArgMax}(d)$, $S \leftarrow \text{Concat}(S, i^*) \triangleright \text{Selected tokens are excluded (distance = 0) from further sampling.}$
- 12: Compute cosine distances to the newly selected token: $\mathbf{d}_{\Delta} \leftarrow \mathbf{1}_{\mathbb{N}} \mathbf{V} \mathbf{V} [\mathbf{i}^*]^{\top}$
- 13: Update each token's minimum distance: d ← ElementwiseMin(d, d_Δ) ▷ Distance of newly selected token i* set to zero in d.
- 14: **end for**
- 15: return S

Algorithm 2 Compute Prompt-Visual Coupling

Require: Visual embeddings $\mathbf{V} \in \mathbb{R}^{n_v \times d}$, Prompt embeddings $\mathbf{P} \in \mathbb{R}^{n_p \times d}$ **Ensure:** Hausdorff distance $h(\mathbf{V}, \mathbf{P})$

Step 1. Compute Pairwise Euclidean Distances

- 1: Compute distance matrix via cdist: $\mathbf{D} \leftarrow \operatorname{cdist}(\mathbf{V}, \mathbf{P}, p = 2)$
- $riangleright \mathbf{D} \in \mathbb{R}^{n_{ extsf{v}} imes n_{ extsf{p}}}$

Step 2. Directed Hausdorff Distances

2: Visual-to-prompt directed distance:

$$d_{v \to p}$$
, $\underline{} \leftarrow \min(\mathbf{D}, \text{ axis} = 2)$, $h_{v \to p} \leftarrow \max(d_{v \to p})$

3: Prompt-to-visual directed distance:

$$d_{p \to v}$$
, $\underline{} \leftarrow \min(\mathbf{D}, \text{ axis} = 1)$, $h_{p \to v} \leftarrow \max(d_{p \to v})$

Step 3. Final Hausdorff Distance

4: **return** $\max(h_{v\to p}, h_{p\to v})$

C Experiment Details

C.1 Benchmarks

Our experiments evaluate the vision-language reasoning abilities of multimodal large language models using a comprehensive suite of widely recognized benchmarks. For image understanding tasks, we assess performance on ten public benchmarks: GQA, MMBench (MMB) and MMBench-CN (MMB $_{\rm CN}$), MME, POPE, VizWiz, ScienceQA (SQA), VQA $^{\rm V2}$, TextVQA (VQA $^{\rm T}$), and OCRBench (OCR). For video understanding tasks, we conduct experiments on four popular benchmarks: TGIF-QA (TGIF), MSVD-QA (MSVD), MSRVTT-QA (MSRV), and ActivityNet-QA (ActNet). The following section provides a concise overview of these benchmarks:

GQA [17] leverages scene graphs, questions, and images to evaluate visual scene understanding and reasoning. By incorporating detailed spatial relationships and object-level attributes, it poses significant challenges for models to perform accurate visual reasoning in complex environments.

MMBench [54] introduces a hierarchical evaluation framework where model capabilities are dissected into three levels. Level-1 focuses on basic perception and reasoning; Level-2 subdivides these abilities into six distinct sub-skills; and Level-3 further refines the evaluation into 20 specific dimensions. Its Chinese counterpart, **MMBench-CN**, adopts a similar structure.

MME [26] rigorously tests perceptual and cognitive abilities across 14 sub-tasks. By employing carefully crafted instruction-answer pairs and succinct instructions, MME minimizes data leakage and provides a robust, fair assessment of a model's multifaceted performance.

POPE [23] targets the evaluation of object hallucination by posing binary questions about object presence in images. It quantifies hallucination levels using metrics, *e.g.*, accuracy, recall, precision, and F1 score, offering a precise and focused measure of model reliability.

VizWiz [14] is a visual question answering benchmark derived from interactions with blind users. Comprising over 31,000 image-question pairs with 10 human-annotated answers per query, it encapsulates the challenges of low-quality image capture and conversational spoken queries, thereby emphasizing real-world visual understanding.

ScienceQA [35] spans multiple scientific domains by organizing questions into 26 topics, 127 categories, and 379 skills. This hierarchical categorization provides a diverse and rigorous testbed for evaluating multimodal understanding, multi-step reasoning, and interpretability across natural, language, and social sciences.

 \mathbf{VQA}^{V2} [12] challenges models with open-ended questions based on 265,016 images that depict a variety of real-world scenes. Each question is paired with 10 human-annotated answers, facilitating a thorough evaluation of a model's capacity to interpret and respond to diverse visual queries.

TextVQA [41] focuses on the integration of text within visual content. It evaluates a model's proficiency in reading and reasoning about textual information embedded in images, thereby requiring a balanced understanding of both visual and linguistic cues.

OCRBench [33] is a comprehensive benchmark for evaluating the OCR capabilities of multi-modal language models across five key tasks: text recognition, scene text-centric and document-oriented VQA, key information extraction, and handwritten mathematical expression recognition.

TGIF-QA [18] adapts the visual question answering task to the video domain by focusing on GIFs. With 165K question-answer pairs, it incorporates tasks, *e.g.*, counting repetitions, identifying repeating actions, detecting state transitions, and frame-specific question answering, thereby demanding detailed spatio-temporal analysis.

MSVD-QA [49] builds upon the MSVD dataset by pairing 1,970 video clips with approximately 50.5K QA pairs. Questions are categorized into five distinct types, *e.g.*, what, who, how, when, and where, making it a versatile tool for evaluating video understanding.

MSRVTT-QA [7] features 10K video clips and 243K QA pairs designed to test the integration of visual and temporal information. Its structure, which parallels that of MSVD-QA through the inclusion of five question types, further enriches the evaluation landscape for video-based tasks.

ActivityNet-QA [53] provides 58K human-annotated question-answer pairs drawn from 5.8K videos. Its focus on questions related to motion, spatial relationships, and temporal dynamics necessitates long-term spatio-temporal reasoning, thus serving as a benchmark for advanced video understanding.

C.2 Multi-modal Large Language Models

We evaluate MoB using various open-source multimodal large language models (MLLMs). For image understanding tasks, experiments are conducted on the LLaVA series, including LLaVA-1.5-7B and LLaVA-Next-7B, as well as the Qwen-VL series, such as Qwen2-VL-7B. Specifically, LLaVA-Next and Qwen2-VL are utilized to validate performance on high-resolution images, *i.e.*, those with a larger number of visual tokens. For video understanding tasks, we employ Video-LLaVA-7B as the baseline model, following the settings reported in its original paper to ensure a fair comparison.

LLaVA-1.5-7B [28] is a robust vision-language model built on the LLaVA framework. It processes images resized to 224×224 and tokenizes them into roughly 572 visual tokens using a patch-based vision encoder. This design balances fine-grained visual representation with computational efficiency, making it effective for diverse multimodal tasks.

LLaVA-Next-7B [29] extends the LLaVA-1.5 by incorporating refined training strategies and data curation. It supports higher-resolution inputs (up to 448×448), yielding up to 2880 visual tokens. These enhancements improve its visual reasoning capabilities and enable more precise alignment between visual content and language but also incur significantly increased computational cost.

Qwen2-VL-7B [44] augments the Qwen2 language model with visual input capabilities. This model leverages cross-modal pretraining to seamlessly merge vision and language, demonstrating strong performance in complex visual question answering and comprehensive scene understanding.

Video-LLaVA-7B [27] extends the LLaVA framework into the temporal domain by processing video inputs. It is designed to capture both spatial and temporal dynamics, enabling effective video comprehension and video-based question answering with coherent and context-aware responses.

C.3 Baselines

To validate the superiority of the proposed MoB, we construct a robust baseline that integrates a comprehensive set of representative existing methods, which encompass single-stage methods with both two distinct objectives and several multi-stage methods.

ToMe [6] employs a lightweight token-matching scheme to merge visually similar tokens across transformer layers, thereby reducing computation without additional training. Its simple yet effective design makes it well suited for real-time applications.

FastV [8] leverages attention maps in the early layers to identify and prune non-critical tokens, significantly reducing initial computational overhead. This focus on early-stage reduction allows the model to operate more efficiently while maintaining performance.

SparseVLM [58] ranks tokens based on cross-modal attention to assess image-prompt relevance and adopts adaptive sparsity ratios to retain key information. It further incorporates a token recycling mechanism to balance the trade-off between efficiency and accuracy.

HiRED [2] allocates token budgets across image partitions by using CLS token attention and then selects the most informative tokens within each partition. This spatially aware approach ensures balanced reduction while preserving contextual details.

LLaVA-PruMerge [38] combines pruning and merging strategies by dynamically removing less important tokens using sparse CLS-visual attention. It then clusters the retained tokens based on key similarity, ensuring that crucial visual features remain intact.

PyramidDrop [48] adopts a progressive token-dropping strategy across different model stages, resulting in a pyramid-like token structure. This method carefully balances the reduction of tokens with the preservation of performance as the processing advances.

MustDrop [31] integrates several token-reduction strategies including spatial merging, text-guided pruning, and output-aware cache policies. Its multi-faceted approach efficiently reduces token counts across various stages of the model.

VisionZip [51] first selects dominant tokens that capture the majority of an image's information and then merges the remaining tokens based on semantic similarity. This approach dramatically reduce token redundancy while accelerating inference and maintaining robust performance.

FasterVLM [57] evaluates token importance using CLS attention in the encoder and prunes tokens before they interact with the language model. This preemptive reduction streamlines the overall process and enhances model efficiency.

GlobalCom² [32] employs a hierarchical strategy by coordinating thumbnail tokens to allocate adaptive retention ratios for high-resolution crops. This approach successfully preserves local details while providing effective global context reduction.

DART [46] leverages token duplication to guide its pruning process instead of relying solely on attention scores. By selecting a small set of pivot tokens and retaining only those with minimal redundancy, DART achieves significant acceleration in a training-free manner.

TokenCarve [42] implements a two-stage, training-free compression framework that preserves critical visual information during aggressive token reduction. It first prunes low-information tokens using an information-preservation guided selection and then merges the remaining tokens based on similarity to minimize accuracy loss.

TwigVLM [39] accelerates large vision-language models by appending a lightweight twig block to an early layer of a frozen base VLM. It utilizes twig-guided token pruning coupled with self-speculative decoding to boost generation speed while retaining high accuracy even under aggressive token reduction.

C.4 Implement Details

From Theorems 1 and 2, the balance between the visual preservation and prompt alignment, *i.e.*, the optimal budget K_p applied for covering prompt tokens \mathcal{P} , is jointly determined by the total budget K and the visual-prompt coupling η . To ensure fair comparison, we evaluate two settings.

- (i) Without η prior. This setting deliberately avoids any benchmark-specific prior (w/o η prior). MoB adjusts $K_{\rm p}$ solely as a function of K to balance the two objectives. Based on an ablation over $\langle K, K_{\rm p} \rangle$, we set $K_{\rm p} \in \{64, 48, 32\}$ and $k \in \{4, 6, 8\}$, corresponding to token-reduction rates of $\{88.9\%, 77.8\%, 66.7\%\}$.
- (ii) With η prior. To verify the $K-\eta-K_{\rm p}$ relationship formulated in Theorems 1 and 2, we introduce a coarse benchmark prior on η . Specifically, we **do not** meticulously search the optimal hyperparameters for MoB, *i.e.*, $K_{\rm p}$ and the covering fold k, per benchmark. Instead, we partition benchmarks by their empirical η distribution into two groups (strong v.s. weak coupling) and employ **the same configuration per group**. From a joint ablation over $\langle K, \eta, K_{\rm p} \rangle$, for image understanding we set

strong coupling:
$$K_{\mathrm{p}} \in \left\{\frac{3K}{8}, \frac{K}{4}, \frac{11K}{24}\right\}, \quad k = \frac{3K_{\mathrm{p}}}{40};$$

weak coupling:
$$K_{\mathrm{p}} \in \left\{ \frac{K}{2}, \frac{7K}{16}, \frac{5K}{12} \right\}, \quad k = \frac{K_{\mathrm{p}}}{8},$$

which again yield token-reduction rates of $\{88.9\%, 77.8\%, 66.7\%\}$.

As for video understanding, we set $K_{\rm p}=\frac{3K}{8},\ k=\frac{3K_{\rm p}}{40}$ for MSVD, MSRV, and ActNet, and $K_{\rm p}=\frac{K}{2},\ k=\frac{K_{\rm p}}{8}$ for TGIF. Unless otherwise stated, the pruning layer index is fixed to $\ell=2$ for both image and video tasks. The same configurations are applied across all MLLMs, and all baselines are run with their default settings.

To ensure reproducibility, we cross-validated our experimental results using the publicly available MLLMs evaluation tool lmms-eval (v0.3.0) [56, 5], with the random seed set to 1234. All experiments were conducted on $4 \times$ Nvidia A800-80GB GPUs paired with $2 \times$ Intel Xeon® Gold 6348 CPUs. The implementation was carried out in Python 3.10 using PyTorch 2.1.2 and CUDA 11.8.

D Additional Experimental Results

D.1 Quantitative Comparison

M-4-3	Obt. "		Strong Co	oupling	3	Weak Coupling						Avia
Method	Objectives	MMB	MMB_{CN}	SQA	VizWiz	GQA	MME	POPE	VQA^{T}	VQA^{V2}	OCR	Avg.
LLaVA-1.5-7B					N = 57							
Vanilla [28]	-	64.7	58.1	69.5	50.0	61.9	1862	85.9	58.2	78.5	297	100%
LLaVA-1.5-7B			Pruning	budge	tK = 19	92; <i>Tok</i>	en Redu	ction Ra	te = 66.7	1%		
ToMe (ICLR'23) [6]	VP	60.5		65.2		54.3	1563	72.4	52.1	68.0	. . .	88.5%
FastV (ECCV'24) [8]	VP	61.2	57.0	67.3	50.8	52.7	1612	64.8	52.5	67.1	291	91.2%
HiRED (AAAI'25) [2]	VP	62.8	54.7	68.4	50.1	58.7	1737	82.8	47.4	74.9	190	91.5%
LLaVA-PruMerge (24.05) [38]	VP	59.6	52.9	67.9	50.1	54.3	1632	71.3	54.3	70.6	253	90.8%
SparseVLM (ICML'25) [58]	PA	62.5	53.7	69.1 68.8	50.5	57.6	1721 1797	83.6	56.1	75.6	292 290	96.3% 96.7%
PyramidDrop (CVPR'25) [48] FiCoCo-V (EMNLP'24) [55]	PA VP	63.3 62.3	56.8 55.3	67.8	51.1 51.0	57.1 58.5	1732	82.3 82.5	56.1 55.7	75.1 74.4	290	96.1%
MustDrop (24.11) [31]	PA VP	62.3	55.8	69.2	51.4	58.2	1787	82.5 82.6	56.5	76.0	289	97.2%
VisionZip (24.12) [51]	VP	63.0	-	68.9	J1. 4	59.3	1783	85.3	57.3	76.8	209	97.7%
DART (EMNLP'25) [46]	VP	63.6	57.0	69.8	51.2	60.0	1856	82.8	57.4	76.7	296	98.8%
TokenCarve (25.03) [42]	PA VP	63.0	-	69.1	50.9	-	1830	84.9	58.4	78.0	-	99.3%
TwigVLM (ICCV'25) [39]	PA	64.0	-	68.8	-	61.2	1848	87.2	58.0	78.1		99.5%
MoB (w/o η-prior)	PA VP	63.8	57.5	70.0	52.4	61.2	1858	84.5	58.2	77.9	304	100.2%
+ η -prior	-	64.1	57.8	70.1	52.5	61.4	1860	84.8	58.5	78.3	307	100.6%
LLaVA-1.5-7B			Pruning	hudae	tK = 12	28. Tak	en Redu	ction Ro	te = 77 \$	8%		
ToMe (ICLR'23)	VP	53.3		59.6	- 12	52.4	1343	62.8	49.1	63.0	-	80.4%
FastV (ECCV'24)	VP	56.1	56.4	60.2	51.3	49.6	1490	59.6	50.6	61.8	285	86.4%
HiRED (AAAI'25)	VP	61.5	53.6	68.1	51.3	57.2	1710	79.8	46.1	73.4	191	90.2%
LLaVA-PruMerge (24.05)	VP	58.1	51.7	67.1	50.3	53.3	1554	67.2	54.3	68.8	248	88.8%
SparseVLM (ICML'25)	PA	60.0	51.1	67.1	51.4	56.0	1696	80.5	54.9	73.8	280	93.8%
PyramidDrop (CVPR'25)	PA	61.6	56.6	68.3	51.0	56.0	1761	82.3	55.1	72.9	287	95.1%
FiCoCo-V (EMNLP'24)	VP	61.1	54.3	68.3	49.4	57.6	1711	82.2	55.6	73.1	-	94.9%
MustDrop (24.11)	PA VP	61.1	55.2	68.5	52.1	56.9	1745	78.7	56.3	74.6	281	95.6%
VisionZip (24.12)	VP	62.0	-	68.9	-	57.6	1762	83.2	56.8	75.6	-	96.2%
DART (EMNLP'25)	VP	63.2	57.5	69.1	51.7	58.7	1840	80.1	56.4	75.9	296	98.0%
TokenCarve (25.03)	PA VP	62.7	-	68.9	51.0	-	1829	84.5	58.1	77.3	-	99.0%
TwigVLM (ICCV'25)	PA	63.5	-	69.5	-	60.6	1818	86.6	57.8	77.9	-	99.0%
MoB (w/o η -prior) + η -prior	PA VP	63.2 63.5	57.3 57.5	69.3 69.6	52.8 52.7	60.7	1842 1845	81.7 82.1	57.5 57.8	77.2 77.5	299 299	99.2% 99.4%
		05.5									233	99. 4 //
LLaVA-1.5-7B		40.7	Pruning		et K = 6							70.16
ToMe (ICLR'23)	VP	43.7 48.0	52.7	50.0 51.1	50.8	48.6 46.1	1138	52.5 48.0	45.3 47.8	57.1 55.0	245	70.1% 77.3%
FastV (ECCV'24) HiRED (AAAI'25)	VP VP	60.2	51.4	68.2	50.8	54.6	1256 1599	73.6	44.2	69.7	191	87.0%
LLaVA-PruMerge (24.05)	VP	55.3	49.1	68.1	50.2	51.9	1549	65.3	54.0	67.4	250	87.0% 87.4%
SparseVLM (ICML'25)	PA	56.2	46.1	62.2	50.1	52.7	1505	75.1	51.8	68.2	180	84.6%
PyramidDrop (CVPR'25)	PA	58.8	50.5	68.6	50.7	41.9	1561	55.9	45.9	69.2	250	78.1%
FiCoCo-V (EMNLP'24)	VP	60.3	53.0	68.1	49.8	52.4	1591	76.0	53.6	71.3	-	91.5%
MustDrop (24.11)	PA VP	60.0	53.1	63.4	51.2	53.1	1612	68.0	54.2	69.3	267	90.1%
VisionZip (24.12)	VP	60.1	-	69.0	-	55.1	1690	77.0	55.5	72.4	-	92.8%
DART (EMNLP'25)	VP	60.6	53.2	69.8	51.6	55.9	1765	73.9	54.4	72.4	270	93.7%
TokenCarve (25.03)	PA VP	62.0	-	69.7	51.4	-	1754	79.9	57.0	74.8	-	97.0%
TwigVLM (ICCV'25)	PA	60.4	-	70.0	-	58.8	1760	82.7	55.8	75.6	-	96.1%
MoB (w/o η-prior)	PA VP	61.7	54.2	69.7	52.0	59.0	1806	77.2	57.0	75.5	277	96.3%
+ η -prior	-	62.1	54.5	69.8	52.1	59.0	1806	77.2	57.0	75.5	277	96.4%
LLaVA-Next-7B			w/o Pr	unina	N = 288	RO: Tok	on Rodu	ction Ro	ıte – 0 09	70		
Vanilla [29]	-	67.4	60.6	70.1	57.6	64.2	1851	86.5	64.9	81.8	517	100%
LLaVA-Next-7B FastV (ECCV'24)	VP	61.6	Pruning 51.9	62.8	t K = 32 53.1	20; <i>10k</i> i 55.9	en Redu 1661	ction Ra 71.7	te = 88. 9 55.7	71.9	374	86.4%
HiRED (AAAI'25)	VP	64.2	55.9	66.7	54.2	59.3	1690	83.3	58.8	75.7	404	91.8%
LLaVA-PruMerge (24.05)	VP	61.3	55.3	66.4	54.0	53.6	1534	60.8	50.6	69.7	146	79.9%
SparseVLM (ICML'25)	PA	60.6	54.5	66.1	52.0	56.1	1533	82.4	58.4	71.5	270	85.9%
PyramidDrop (CVPR'25)	PA	63.4	56.2	67.5	54.1	56.4	1663	77.6	54.4	73.5	259	86.8%
MustDrop (24.11)	PA VP	62.8	55.1	68.0	54.0	57.3	1641	82.1	59.9	73.7	382	90.4%
VisionZip (24.12)	VP	63.1	-	67.3	-	59.3	1702	-	58.9	76.2	-	93.0%
FasterVLM (24.12) [57]	VP	61.6	53.5	66.5	52.6	56.9	1701	83.6	56.5	74.0	401	89.8%
GlobalCom ² (25.01) [32]	VP	61.8	53.4	67.4	54.6	57.1	1698	83.8	57.2	76.7	375	90.3%
DART (EMNLP'25)	VP	65.3	58.2	68.4	56.1	61.7	1710	84.1	58.7	79.1	406	93.9%
TwigVLM (ICCV'25)	PA	65.0	-	68.7	-	62.2	1758	-	57.4	79.7	-	95.4%
MoB (with η-prior)	PA VP	65.8	58.9	68.7	57.0	62.6	1760	84.4	60.2	80.1	418	95.4%
(00.0	23.7	00.7	27.0	02.0	1,00	U 1. 1	00.2	00.1	.10	75.170

Table 5: Full results on image understanding with the LLaVA-7B Series. For MoB, we set $K_{\rm p} \in \{64,48,32\}$ and $k \in \{4,6,8\}$, corresponding to token-reduction rates of $\{88.9\%,77.8\%,66.7\%\}$. For MoB with the η prior, we use $K_{\rm p} \in \{\frac{3K}{8},\frac{K}{4},\frac{K}{4}\}$ with $k=\frac{3K_{\rm p}}{40}$ for strong-coupling benchmarks and $K_{\rm p} \in \{\frac{K}{2},\frac{7K}{16},\frac{5K}{12}\}$ with $k=\frac{K_{\rm p}}{8}$ for weak-coupling benchmarks, corresponding to the same token-reduction rates; the pruning layer is fixed at $\ell=2$. B and O denote the best and the second.

D.2 Visualization

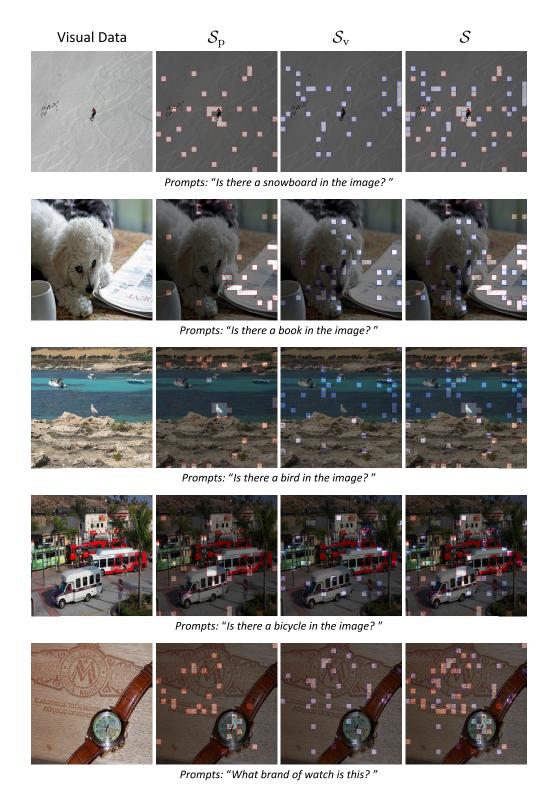
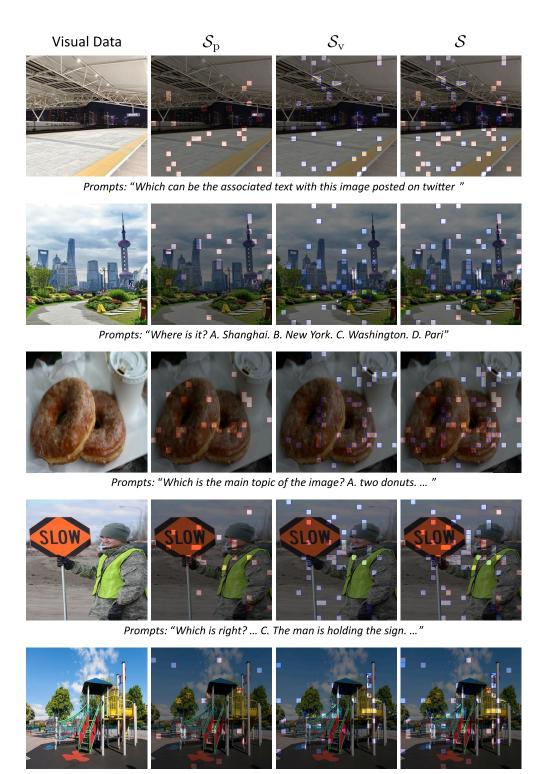


Figure 7: Visualization of the selected prompt and visual centers under weak coupling.



Prompts: "What type of environment is depicted in the picture? ... B. Children's playground. ..."

Figure 8: Visualization of the selected prompt and visual centers under strong coupling.

MoB formulates visual token pruning as a bi-objective covering problem over $(\mathcal{V}, \mathcal{P})$, which is expected to gather query-relevant, fine-grained evidence with \mathcal{S}_p while preserving global scene context with \mathcal{S}_v . The visualizations (Figures 7 and 8) qualitatively validate this design: tokens in \mathcal{S}_p concentrate in regions aligned with the text query and key visual evidence, whereas elements of \mathcal{S}_v spread more uniformly across the image to maintain the overall context. Together, this complementary allocation enables MoB to retain the most informative visual content for each image–query pair, accounting for its strong empirical performance.

D.3 Additional Ablation & Discussion

GQA												
$\langle K, K_{ m p} \rangle$	0	2	4	6	8	12	16	24	32	48	64	96
$\langle 64, 32 \rangle$	58.3	58.8	59.0	-	58.7	-	58.2	-	57.4	-	-	-
$\langle 128, 64 \rangle$	60.2	-	60.5	-	60.7	-	60.6	-	60.0	-	59.5	
$\langle 192, 96 \rangle$	60.6	-	-	61.1	-	61.2	-	60.9	-	60.7	-	60.5
	TextVQA											
$\langle K, K_{\rm p} \rangle$	0	2	4	6	8	12	16	24	32	48	64	96
$\langle 64, 32 \rangle$	56.5	56.9	57.0	-	56.8	-	56.5	-	56.2	-	-	-
$\langle 128, 64 \rangle$	57.1	-	57.5	-	57.7	-	57.7	-	57.2	-	56.8	-
$\langle 192, 96 \rangle$	57.8	-	-	58.2	-	58.2	-	58.1	-	57.7	-	57.5

Table 6: Detailed ablation on the covering fold k for GQA and TextVQA.

To assess MoB's sensitivity to the covering-fold parameter k—particularly under weak coupling with long prompts—we conduct a detailed ablation on k using GQA and TextVQA.

As Table 6 demonstrates, MoB is not overly sensitive to the choice of k, particularly within a clear optimal range. For instance, in both two benchmarks, performance only varies by approximately 0.3% for k values between [2,8] under $\langle K=64,K_{\rm p}=32\rangle$ setting.

There is also a principled, theoretical reason for this robustness, which stems from the relationship between the covering fold k, the budget $K_{\rm p}$, and the length L of prompt tokens \mathcal{P} . From covering theory, every prompt token $p \in \mathcal{P}$ is covered by at least one visual token $v \in \mathcal{V}$ under the condition $K_{\rm p} \geq kL$, thereby ensuring the performance guarantee of MoB. Therefore, as selected k satisfies $k \leq K_{\rm p}/L$, the performance will remain stable.

Heuristic for estimating k. In practice, a robust range for k can be inferred from the prompt length L. Given the analysis above, we expect an adaptive, per-sample search for a fine-grained k to yield only limited gains, so we rely on this length-based heuristic instead.

D.4 Real-life Application

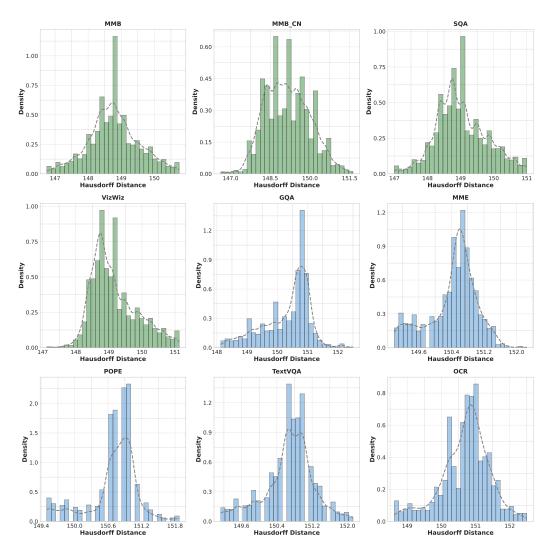


Figure 9: Observations of prompt-visual coupling η across 9 popular benchmarks.

Open-domain recipe. MoB is task-agnostic, which does not require pre-defined task labels and can operate online by classifying each sample's coupling pattern. For a given target model, we adopt a two-stage strategy:

- Offline calibration. Analyze the empirical η distributions on a set of representative benchmarks (as shown in Figure 9) and set a robust threshold τ that separates *weak* vs. *strong* coupling.
- Online classification and inference. For each incoming query, compute its Hausdorff distance using Algorithm 2 with tractable bilinear complexity $\mathcal{O}(NLd)$. Classify the sample by comparing this value to τ , then apply the corresponding budget configuration (e.g., $K_{\rm p}$, k) and run MoB + forward inference. In practice, this online cost is negligible relative to the pruned forward pass.

Computational Overhead. We provide a detailed cost breakdown for online computation of the Hausdorff distance using Algorithm 2 with complexity $\mathcal{O}(NLd)$ on LLaVA-1.5-7B and LLaVA-Next-7B, where N, L, and d denote the numbers of visual tokens, prompt tokens, and the feature dimension, respectively. As shown in Table 7, the measured cost (TFLOPs) of exact Hausdorff computation is orders of magnitude smaller than that of MoB itself and the model's forward pass, yielding a negligible overhead.

Model L	LaVA-1.5 (<i>N</i> Vanilla	= 576, L = 1 K = 64	0, d = 4096) K = 128	K = 192						
Forward Compute d_H MoB	8.2 $2.3e - 5$	1.0 $2.3e - 5$ $1.7e - 4$	1.9 $2.3e - 5$ $3.3e - 4$	2.8 $2.3e - 5$ $4.8e - 4$						
LLaVA-Next ($N=2880, L=10, d=4096$) Model Vanilla $K=320$ $K=640$ $K=960$										
Forward Compute d_H MoB	40.5 1.2e – 4 –	4.6 $1.2e - 4$ $3.9e - 3$	9.1 $1.2e - 4$ $7.6e - 3$	$ \begin{array}{c} 13.6 \\ 1.2e - 4 \\ 1.1e - 2 \end{array} $						

Table 7: Computation cost in LLaVA-7B series (TFLOPs)

Concretely, computing d_H (e.g., $\sim 1.2 \times 10^{-4}$ TFLOPs on LLaVA-Next) is insignificant relative to the pruned forward pass (e.g., ~ 4.6 TFLOPs at K=320) and, more importantly, to the savings from pruning (~ 35.9 TFLOPs). Thus, exact online estimation is not a practical bottleneck; its cost is dwarfed by the efficiency gains of our method. Further acceleration is possible with standard techniques (e.g., heuristic support sampling or low-dimensional random projections), although it is unnecessary in our settings.

- Heuristic Sampling: It computes the distance on smaller support sets of the tokens $(\mathcal{V}' \subset \mathcal{V}, \mathcal{P}' \subset \mathcal{P})$, which can be constructed via random sampling [46] or more advanced heuristics such as Key-Norm selection [1, 13]. This reduces complexity to $\mathcal{O}(N'L'd)$, where $|\mathcal{V}'| = N'$, $|\mathcal{P}'| = L'$.
- Random Projections: For a more theoretically grounded approach, the Johnson–Lindenstrauss (JL) lemma [19] allows us to project embeddings to a much lower dimension $(d' \ll d)$ while preserving geometric structure, reducing complexity to $\mathcal{O}(NL\,d')$.

Potential Extensions. A natural extension is to maintain an online estimate of the coupling statistic η during inference—e.g., a running summary of an approximate $\hat{\eta}$ computed from shallow-layer tokens. As more samples are processed, we *expect* the empirical distribution of $\hat{\eta}$ to become bimodal (consistent with the benchmark patterns in Figure 9), enabling a data-driven threshold to be derived on the fly that separates weak vs. strong coupling regimes. Using this live threshold, MoB could *adapt* K_p (and k) per sample or per mini-batch by selecting from a small budget pool or by scheduling K_p as a function of $\hat{\eta}$, with conservative warm-up and safeguards for distribution shift.

E Omitted Technical Details

E.1 Proof of Lemma 1

Restatement of Lemma 1 (An Error Bound for Visual Token Pruning). *Under Assumption 1, given any token set with its pruned counterpart* $\mathcal{X} = \mathcal{V} \sqcup \mathcal{P}$, $\mathcal{X}_s = \mathcal{S} \sqcup \mathcal{P} \subseteq \mathbb{R}^d$, the pruning error bound is given by:

$$\|\mathcal{F}(\mathcal{X}) - \mathcal{F}(\mathcal{X}_{s})\| \le C_{\ell} \max \Big\{ \min \big\{ d_{H}(\mathcal{S}, \mathcal{V}), d_{H}(\mathcal{V}, \mathcal{P}) \big\}, \min \big\{ d_{H}(\mathcal{S}, \mathcal{V}), d_{H}(\mathcal{S}, \mathcal{P}) \big\} \Big\}.$$

Remark. Here $d_H(S, P)$ and $d_H(S, V)$ describe the prompt alignment and visual preservation, while $d_H(V, P)$ is an inherent term that describes the prompt-visual coupling of input data.

Proof. The intermediate input for any layer and its pruned counterpart are given by

$$\mathcal{X} = \mathcal{V} \sqcup \mathcal{P}$$
 and $\mathcal{X}_s = \mathcal{S} \sqcup \mathcal{P}$.

By Equation (1), the Hausdorff distance is symmetric, i.e.,

$$d_H(\mathcal{S}, \mathcal{V}) = d_H(\mathcal{V}, \mathcal{S}), \tag{E1-1}$$

and induced by Euclidean distance.

Step 1. Bound the one-sided distances.

We analyze the distances by considering the membership of the points in the subsets.

Direction 1 ($\mathcal{X} \to \mathcal{X}_s$) For any $x \in \mathcal{X}$:

Case (i): If $x \in \mathcal{P}$, then since $\mathcal{P} \subset \mathcal{X}_s$,

$$\inf_{y \in \mathcal{X}_{\mathbf{s}}} \|x - y\| = 0.$$

Case (ii): If $x \in \mathcal{V}$, then the candidate points in $\mathcal{X}_s = \mathcal{S} \sqcup \mathcal{P}$ can be chosen either from \mathcal{S} or \mathcal{P} . Thus,

$$\inf_{y \in \mathcal{X}_s} \|x - y\| \le \min \Big\{ \inf_{s \in \mathcal{S}} \|x - s\|, \inf_{p \in \mathcal{P}} \|x - p\| \Big\}.$$

Taking the supremum over $x \in \mathcal{V}$ yields

$$\sup_{x \in \mathcal{V}} \inf_{y \in \mathcal{X}_s} \|x - y\| \le \min \Big\{ \sup_{x \in \mathcal{V}} \inf_{s \in \mathcal{S}} \|x - s\|, \sup_{x \in \mathcal{V}} \inf_{p \in \mathcal{P}} \|x - p\| \Big\}.$$

$$\sup_{x \in \mathcal{V}} \inf_{p \in \mathcal{P}} \|x - p\| \le \max \left\{ \sup_{x \in \mathcal{V}} \inf_{p \in \mathcal{P}} \|x - p\|, \sup_{p \in \mathcal{P}} \inf_{x \in \mathcal{V}} \|p - x\| \right\} = d_H(\mathcal{V}, \mathcal{P}),$$

By Equation (1), we derive the distance in direction 1:

$$\sup_{x \in \mathcal{X}} \inf_{y \in \mathcal{X}_{s}} \|x - y\| \le \min \Big\{ d_{H}(\mathcal{V}, \mathcal{S}), d_{H}(\mathcal{V}, \mathcal{P}) \Big\}.$$
 (E1-2)

Direction 2 $(\mathcal{X}_s \to \mathcal{X})$ For any $y \in \mathcal{X}_s$:

Case (i): If $y \in \mathcal{P}$, then as $\mathcal{P} \subset \mathcal{X}$,

$$\inf_{x \in \mathcal{X}} \|y - x\| = 0.$$

Case (ii): If $y \in \mathcal{S}$, the candidate points in $\mathcal{X} = \mathcal{V} \sqcup \mathcal{P}$ can be chosen from either \mathcal{V} or \mathcal{P} ; hence

$$\inf_{x \in \mathcal{X}} \|y - x\| \le \min \Bigl\{ \inf_{v \in \mathcal{V}} \|y - v\|, \inf_{p \in \mathcal{P}} \|y - p\| \Bigr\}.$$

Taking the supremum over $y \in \mathcal{S}$ yields

$$\sup_{y \in \mathcal{S}} \inf_{x \in \mathcal{X}} \|y - x\| \leq \min \Big\{ \sup_{y \in \mathcal{S}} \inf_{v \in \mathcal{V}} \|y - v\|, \sup_{y \in \mathcal{S}} \inf_{p \in \mathcal{P}} \|y - p\| \Big\}.$$

$$\sup_{y \in \mathcal{S}} \inf_{p \in \mathcal{P}} \|y - p\| \leq \max \left\{ \sup_{y \in \mathcal{S}} \inf_{p \in \mathcal{P}} \|y - p\|, \sup_{p \in \mathcal{P}} \inf_{y \in \mathcal{S}} \|p - y\| \right\} = d_H(\mathcal{S}, \mathcal{P}),$$

By Equation (1), we derive the distance in direction 2:

$$\sup_{y \in \mathcal{X}_s} \inf_{x \in \mathcal{X}} \|y - x\| \le \min \Big\{ d_H(\mathcal{S}, \mathcal{V}), \ d_H(\mathcal{S}, \mathcal{P}) \Big\}.$$
 (E1-3)

Step 2. Combine the bounds.

By Equation (1), combining the bounds in (E1-2) and (E1-3), we obtain

$$d_H(\mathcal{X}, \mathcal{X}_s) \le \max \Big\{ \min \{ d_H(\mathcal{V}, \mathcal{S}), d_H(\mathcal{V}, \mathcal{P}) \}, \min \{ d_H(\mathcal{S}, \mathcal{V}), d_H(\mathcal{S}, \mathcal{P}) \} \Big\}.$$

Based on (E1-1), we have

$$d_H(\mathcal{X}, \mathcal{X}_s) \leq \max \Big\{ \min \big\{ d_H(\mathcal{S}, \mathcal{V}), \ d_H(\mathcal{V}, \mathcal{P}) \big\}, \min \big\{ d_H(\mathcal{S}, \mathcal{V}), \ d_H(\mathcal{S}, \mathcal{P}) \big\} \Big\}.$$

Loading the Assumption 1, we have the output discrepancy is bounded by

$$\|\mathcal{F}(\mathcal{X}) - \mathcal{F}(\mathcal{X}_{s})\| \leq C_{\ell} d_{H}(\mathcal{X}, \mathcal{X}_{s}).$$

$$= C_{\ell} \max \Big\{ \min \big\{ d_{H}(\mathcal{S}, \mathcal{V}), d_{H}(\mathcal{V}, \mathcal{P}) \big\}, \min \big\{ d_{H}(\mathcal{S}, \mathcal{V}), d_{H}(\mathcal{S}, \mathcal{P}) \big\} \Big\}.$$

This completes the proof.

E.2 Proof of Lemma 2

Restatement of Lemma 2 (A Relaxed Error Bound under Practical Budgets). *Under Assumptions 1* and 2, let $\mathcal{X} = \mathcal{V} \sqcup \mathcal{P}$, $\mathcal{X}_s = \mathcal{S} \sqcup \mathcal{P} \subseteq \mathbb{R}^d$ with $|\mathcal{S}| = K \ll N$. Partition the retained token set \mathcal{S} into two disjoint subsets: $\mathcal{S} = \mathcal{S}_p \sqcup \mathcal{S}_v$, devoted to prompt alignment $d_H(\mathcal{S}_p, \mathcal{P})$ and visual preservation $d_H(\mathcal{S}_v, \mathcal{V})$, respectively. Then, the pruning error bound reduces to

$$\|\mathcal{F}(\mathcal{X}) - \mathcal{F}(\mathcal{X}_{s})\| \le C_{\ell} \max \{d_{H}(\mathcal{S}_{p}, \mathcal{P}), d_{H}(\mathcal{S}_{v}, \mathcal{V})\} + C_{\ell} \eta.$$

Proof. By Lemma 1, we obtain

$$\|\mathcal{F}(\mathcal{X}) - \mathcal{F}(\mathcal{X}_{s})\| \le C_{\ell} \max \Big\{ \min \big\{ d_{H}(\mathcal{S}, \mathcal{V}), d_{H}(\mathcal{V}, \mathcal{P}) \big\}, \min \big\{ d_{H}(\mathcal{S}, \mathcal{V}), d_{H}(\mathcal{S}, \mathcal{P}) \big\} \Big\}.$$

Since $\min\{a, b\} \le \max\{a, b\}$, we have

$$\|\mathcal{F}(\mathcal{X}) - \mathcal{F}(\mathcal{X}_{s})\| \le C_{\ell} \max \left\{ d_{H}(\mathcal{S}, \mathcal{P}), d_{H}(\mathcal{S}, \mathcal{V}), d_{H}(\mathcal{V}, \mathcal{P}) \right\}.$$
 (E2-1)

For any $p \in \mathcal{P}$, we have

$$\inf_{s \in \mathcal{S}} \|p - s\| = \min \left\{ \inf_{s \in \mathcal{S}_{\mathbf{p}}} \|p - s\|, \inf_{s \in \mathcal{S}_{\mathbf{v}}} \|p - s\| \right\} \le \inf_{s \in \mathcal{S}_{\mathbf{p}}} \|p - s\|.$$

Taking the supremum over $p \in \mathcal{P}$ yields

$$\sup_{p \in \mathcal{P}} \inf_{s \in \mathcal{S}} \|p - s\| \le \sup_{p \in \mathcal{P}} \inf_{s \in \mathcal{S}_{p}} \|p - s\|.$$

Similarly, since $S_v \subset S$,

$$\sup_{s \in \mathcal{S}_{v}} \inf_{p \in \mathcal{P}} \|s - p\| \le \sup_{s \in \mathcal{S}} \inf_{p \in \mathcal{P}} \|s - p\|.$$

Thus, by Equation (1),

$$d_H(\mathcal{S}, \mathcal{P}) \le \max \Big\{ d_H(\mathcal{S}_{p}, \mathcal{P}), \ d_H(\mathcal{S}_{v}, \mathcal{P}) \Big\}.$$

Using Assumption 2 $(d_H(\mathcal{V}, \mathcal{P}) \leq \eta)$ and the triangle inequality for Hausdorff distance, we have

$$d_H(\mathcal{S}_{\mathbf{v}}, \mathcal{P}) \le d_H(\mathcal{S}_{\mathbf{v}}, \mathcal{V}) + d_H(\mathcal{V}, \mathcal{P}) \le d_H(\mathcal{S}_{\mathbf{v}}, \mathcal{V}) + \eta_{\mathcal{P}}$$

$$d_H(\mathcal{S}_{\mathbf{p}}, \mathcal{V}) \le d_H(\mathcal{S}_{\mathbf{p}}, \mathcal{P}) + d_H(\mathcal{P}, \mathcal{V}) \le d_H(\mathcal{S}_{\mathbf{p}}, \mathcal{P}) + \eta_{\mathcal{P}}$$

Hence,

$$d_H(\mathcal{S}, \mathcal{P}) \le \max \left\{ d_H(\mathcal{S}_{\mathbf{p}}, \mathcal{P}), \ d_H(\mathcal{S}_{\mathbf{v}}, \mathcal{V}) + \eta \right\}.$$
 (E2-2)

Similarly, one can show that

$$d_H(\mathcal{S}, \mathcal{V}) \le \max \left\{ d_H(\mathcal{S}_{\mathbf{v}}, \mathcal{V}), \ d_H(\mathcal{S}_{\mathbf{p}}, \mathcal{P}) + \eta \right\}.$$
 (E2-3)

Loading the maximum of (E2-2), (E2-3) and $d_H(\mathcal{V}, \mathcal{P})$ into (E2-1), we obtain

$$\|\mathcal{F}(\mathcal{X}) - \mathcal{F}(\mathcal{X}_{s})\| \leq C_{\ell} \max \left\{ d_{H}(\mathcal{S}, \mathcal{P}), d_{H}(\mathcal{S}, \mathcal{V}), d_{H}(\mathcal{V}, \mathcal{P}) \right\}$$

$$\leq C_{\ell} \max \left\{ d_{H}(\mathcal{S}_{p}, \mathcal{P}), d_{H}(\mathcal{S}_{v}, \mathcal{V}) + \eta, d_{H}(\mathcal{S}_{v}, \mathcal{V}), d_{H}(\mathcal{S}_{p}, \mathcal{P}) + \eta, \eta \right\}$$

Since $d_H(\mathcal{S}_p, \mathcal{P}) \geq 0$, $d_H(\mathcal{S}_v, \mathcal{V}) \geq 0$, $\eta \geq 0$, we have

$$\max\{d_H(\mathcal{S}_p, \mathcal{P}), d_H(\mathcal{S}_p, \mathcal{P}) + \eta, \eta\} = d_H(\mathcal{S}_p, \mathcal{P}) + \eta,$$

$$\max\{d_H(\mathcal{S}_v, \mathcal{V}), d_H(\mathcal{S}_v, \mathcal{V}) + \eta, \eta\} = d_H(\mathcal{S}_v, \mathcal{V}) + \eta.$$

Hence

$$\|\mathcal{F}(\mathcal{X}) - \mathcal{F}(\mathcal{X}_{\mathrm{s}})\| \le C_{\ell} \max \left\{ d_H(\mathcal{S}_{\mathrm{p}}, \mathcal{P}), d_H(\mathcal{S}_{\mathrm{v}}, \mathcal{V}) \right\} + C_{\ell} \eta.$$

This completes the proof.

E.3 Proof of Lemma 3

Restatement of Lemma 3 (d_{eff} -regular lower bound on covering numbers). Given $\mathcal{P}, \mathcal{V} \subset \mathbb{R}^d$ with an effective dimension d_{eff} . Suppose their δ -dilations $\mathcal{V}_{\delta} := \bigcup_{v \in \mathcal{V}} B(v, \delta)$, $\mathcal{P}_{\delta} := \bigcup_{p \in \mathcal{P}} B(p, \delta)$ ($\delta \ll \eta$) satisfy d_{eff} -dimensional covering regularity; thus, there exist constants b > a > 0, b' > a' > 0 and $\epsilon_0 > \delta$ such that

$$a \epsilon_{\mathrm{p}}^{-d_{\mathrm{eff}}} \leq \mathcal{N}(\mathcal{P}, \epsilon_{\mathrm{p}}) \leq b \epsilon_{\mathrm{p}}^{-d_{\mathrm{eff}}}, \qquad a' \epsilon_{\mathrm{v}}^{-d_{\mathrm{eff}}} \leq \mathcal{N}(\mathcal{V}, \epsilon_{\mathrm{v}}) \leq b' \epsilon_{\mathrm{v}}^{-d_{\mathrm{eff}}}, \qquad \forall \epsilon_{\mathrm{p}}, \epsilon_{\mathrm{v}} \in (\delta, \epsilon_{0}],$$

Remark Previous work suggests that both visual and language embeddings concentrate on a low-dimensional manifold, so the effective covering dimension satisfies the typical relation $d_{\rm eff} \ll d$.

Proof. We prove the two-sided bound for \mathcal{P} ; the argument for \mathcal{V} is identical.

Notation.

- $\mathcal{N}(X,r)$: minimal number of closed balls of radius r covering X.
- $X_{\delta} = \bigcup_{x \in X} B(x, \delta)$, with $B(x, \delta) = \{y : ||y x|| \le \delta\}$.

Step 1. Transfer trick for small ϵ .

Fix $\epsilon \in (\delta, \epsilon_0]$ and define $\epsilon' = \min\{\epsilon + \delta, \epsilon_0\}$.

If $\epsilon \leq \epsilon_0 - \delta$ (so $\epsilon' = \epsilon + \delta$), then any ϵ -cover $\{z_i\}_{i=1}^m$ of \mathcal{P} satisfies for each $y \in \mathcal{P}_{\delta}$:

$$\exists x \in \mathcal{P} : ||y - x|| \le \delta, \quad \exists i : ||x - z_i|| \le \epsilon \implies ||y - z_i|| \le \epsilon + \delta = \epsilon'.$$

Hence

$$\mathcal{P}_{\delta} \subseteq \bigcup_{i=1}^{m} B(z_{i}, \epsilon') \implies \mathcal{N}(\mathcal{P}_{\delta}, \epsilon') \leq \mathcal{N}(\mathcal{P}, \epsilon).$$
 (E3-1)

Note: For $\epsilon > \epsilon_0 - \delta$, the above transfer argument is not applied.

Step 2. Lower bound on $\mathcal{N}(\mathcal{P}, \epsilon)$.

Split into two cases:

Case I: $\epsilon \le \epsilon_0 - \delta$. Since \mathcal{P}_{δ} satisfies d_{eff} -dimensional covering regularity; loading the lower-bound for \mathcal{P}_{δ} at radius $\epsilon' = \epsilon + \delta$, there exists a constant $a_{\delta} \ge 0$ such that

$$\mathcal{N}(\mathcal{P}_{\delta}, \epsilon') = \mathcal{N}(\mathcal{P}_{\delta}, \epsilon + \delta) \geq a_{\delta} (\epsilon + \delta)^{-d_{\text{eff}}}.$$

Based on (E3-1), we obtain

$$a_{\delta}(\epsilon + \delta)^{-d_{\text{eff}}} < \mathcal{N}(\mathcal{P}_{\delta}, \epsilon') < \mathcal{N}(\mathcal{P}, \epsilon)$$

Since $\delta \leq \epsilon$, it follows that $\epsilon + \delta \leq 2\epsilon$; thus, we have

$$\mathcal{N}(\mathcal{P}, \epsilon) \ge a_{\delta} 2^{-d_{\text{eff}}} \epsilon^{-d_{\text{eff}}}.$$
 (E3-2)

Case II: $\epsilon > \epsilon_0 - \delta$. Define $\tilde{a} := (\epsilon_0 - \delta)^{d_{\rm eff}}$, such that

$$(\epsilon_0 - \delta)^{-d_{\text{eff}}} = \widetilde{a}^{-1}.$$

Since $\epsilon > \epsilon_0 - \delta$, we have

$$\epsilon^{-d_{\text{eff}}} \leq (\epsilon_0 - \delta)^{-d_{\text{eff}}}.$$

Hence

$$\epsilon^{-d_{\mathrm{eff}}} \leq \widetilde{a}^{-1} \iff \widetilde{a} \epsilon^{-d_{\mathrm{eff}}} \leq 1.$$

Since any nonempty set \mathcal{P} has covering number at least one, the following holds

$$\widetilde{a} \, \epsilon^{-d_{\text{eff}}} \leq 1 \leq \mathcal{N}(\mathcal{P}, \epsilon).$$
 (E3-3)

Therefore, set $a := \min\{a_{\delta}2^{-d_{\text{eff}}}, \widetilde{a}\} > 0$, combining (E3-2) and (E3-3) yields

$$\mathcal{N}(\mathcal{P}, \epsilon) \ge a \epsilon^{-d_{\text{eff}}}, \quad \forall \epsilon \in (\delta, \epsilon_0].$$
 (E3-4)

Similarly, \mathcal{V} holds $\mathcal{N}(\mathcal{V}, \epsilon) \geq a' \epsilon^{-d_{\mathrm{eff}}}, \quad \forall \epsilon \in (\delta, \epsilon_0].$

Step 3. Upper bound on $\mathcal{N}(\mathcal{P}, \epsilon)$.

Since \mathcal{P}_{δ} satisfies d_{eff} -dimensional covering regularity, there exists a constant $b_{\delta} \geq a_{\delta} \geq 0$ such that

$$\mathcal{N}(\mathcal{P}_{\delta}, \epsilon) \leq b_{\delta} \epsilon^{-d_{\text{eff}}}.$$

Since $\mathcal{P} \subseteq \mathcal{P}_{\delta}$, we have $\mathcal{N}(\mathcal{P}, \epsilon) \leq \mathcal{N}(\mathcal{P}_{\delta}, \epsilon)$; thus, the following holds

$$\mathcal{N}(\mathcal{P}, \epsilon) \le \mathcal{N}(\mathcal{P}_{\delta}, \epsilon) \le b_{\delta} \epsilon^{-d_{\text{eff}}}.$$

Based on the *monotonicity of covering numbers*, for every radius $\epsilon \geq \delta$, we have

$$\mathcal{N}(\mathcal{P}, \epsilon) \leq \mathcal{N}(\mathcal{P}, \delta).$$

Therefore, set $b := \max\{b_{\delta}, \mathcal{N}(\mathcal{P}, \delta)\}$, for all $\epsilon \in (\delta, \epsilon_0]$ we have

$$\mathcal{N}(\mathcal{P}, \epsilon) \leq b \, \epsilon^{-d_{\text{eff}}}.$$
 (E3-5)

Likewise for \mathcal{V} , the following holds $\mathcal{N}(\mathcal{V}, \epsilon) \leq b' \, \epsilon^{-d_{\mathrm{eff}}}, \quad \forall \, \epsilon \in (\delta, \epsilon_0].$

Step 4. Combine the bounds.

Based on (E3-4) and (E3-5), for all $\epsilon \in (\delta, \epsilon_0]$ the following holds

$$a \, \epsilon^{-d_{\text{eff}}} \leq \mathcal{N}(\mathcal{P}, \epsilon) \leq b \, \epsilon^{-d_{\text{eff}}}, \quad a' \, \epsilon^{-d_{\text{eff}}} \leq \mathcal{N}(\mathcal{V}, \epsilon) \leq b' \, \epsilon^{-d_{\text{eff}}}.$$

This completes the proof.

E.4 Proof of Theorem 1

Restatement of Theorem 1 (Trade-off between Prompt Alignment and Visual Preservation). *Under Assumption 2 and the covering-regularity hypothesis of Lemma 3 with constants* $a, a', d_{\text{eff}} > 0$, there exist a radius-scaling factor z > 1 such that $\eta/z > \delta$ and $K < \mathcal{N}(\mathcal{P}, \eta/z) + \mathcal{N}(\mathcal{V}, \eta/z)$, for every pruning results $\mathcal{S} = (\mathcal{S}_p \sqcup \mathcal{S}_v) \subseteq \mathcal{V}$ with budget K satisfying

$$\max \{D_1 K^{-2/d_{\text{eff}}}, D_2 \eta^2\} \le d_H(\mathcal{S}_p, \mathcal{P}) d_H(\mathcal{S}_v, \mathcal{V}),$$

where $D_1 := (4 \, a \, a')^{1/d_{\text{eff}}} > 0$, $D_2 := 1/z^2 > 0$.

Remark (Optimal Attainment Level). The term $D_1 K^{-2/d_{\rm eff}}$ is completely determined by the pruning budget, while $D_2 \eta^2$ quantifies the effect of prompt-visual coupling. Hence, the optimal attainment level per objective is given by $\epsilon^* = \max\{\eta/z, \sqrt{D_1} K^{-1/d_{\rm eff}}\}$. Any attempt to reduce one objective below ϵ^* forces the other above ϵ^* , thereby increasing the overall pruning error.

Remark (Effect of Budget and Coupling Strength). As K decreases, z correspondingly shrinks (D_2 growing as a power function), ultimately making $D_2 \eta^2$ dominate the bound; while as K increases, both of the terms reduce, thereby diminishing the trade-off and tightening the overall error bound.

Proof. We begin the proof by noting

$$\epsilon_{\mathrm{p}} = d_H(\mathcal{S}_{\mathrm{p}}, \mathcal{P}), \quad \epsilon_{\mathrm{v}} = d_H(\mathcal{S}_{\mathrm{v}}, \mathcal{V}), \quad K_{\mathrm{p}} = |\mathcal{S}_{\mathrm{p}}|, \quad K_{\mathrm{v}} = |\mathcal{S}_{\mathrm{v}}|, \quad K_{\mathrm{p}} + K_{\mathrm{v}} = K.$$

Step 1. Quantify the impact of budget K.

By Lemma 3, for all ϵ_p , $\epsilon_v \in (\delta, \epsilon_0]$, we have

$$a \epsilon_{\mathbf{p}}^{-d_{\text{eff}}} \leq \mathcal{N}(\mathcal{P}, \epsilon_{\mathbf{p}}) \leq K_{\mathbf{p}}, \quad a' \epsilon_{\mathbf{v}}^{-d_{\text{eff}}} \leq \mathcal{N}(\mathcal{V}, \epsilon_{\mathbf{v}}) \leq K_{\mathbf{v}}.$$
 (E4-1)

By AM-GM inequality, we have $K_pK_v \leq \left(\frac{K}{2}\right)^2$; thus, loading (E4-1) we have

$$(a\,a')\,(\epsilon_{\mathrm{p}}\,\epsilon_{\mathrm{v}})^{-d_{\mathrm{eff}}} \leq \left(\frac{K}{2}\right)^2 \implies \epsilon_{\mathrm{p}}\,\epsilon_{\mathrm{v}} \geq (4\,a\,a')^{1/d_{\mathrm{eff}}}\,K^{-2/d_{\mathrm{eff}}}.$$

Define $D_1 := (4 a a')^{1/d_{\text{eff}}} > 0$, the K-bound is established by

$$\epsilon_{\rm p} \, \epsilon_{\rm v} \geq D_1 \, K^{-2/d_{\rm eff}}.$$
 (E4-2)

Step 2. Quantify the impact of prompt-visual coupling η .

Based on the budget condition, the radius-scaling factor z holds

$$K < \mathcal{N}(\mathcal{P}, \frac{\eta}{z}) + \mathcal{N}(\mathcal{V}, \frac{\eta}{z}).$$
 (E4-3)

For contradiction, we suppose two covering radii is simultaneously small, such that $\epsilon_{\rm p} < \eta/z$ and $\epsilon_{\rm v} < \eta/z$. Then, the monotonicity of covering numbers gives

$$\mathcal{N}(\mathcal{P}, \epsilon_{\mathrm{p}}) \geq \mathcal{N}\big(\mathcal{P}, \tfrac{\eta}{z}\big), \quad \mathcal{N}(\mathcal{V}, \epsilon_{\mathrm{v}}) \geq \mathcal{N}\big(\mathcal{V}, \tfrac{\eta}{z}\big).$$

Hence

$$K \geq \mathcal{N}(\mathcal{P}, \epsilon_{\mathrm{p}}) + \mathcal{N}(\mathcal{V}, \epsilon_{\mathrm{v}}) \geq \mathcal{N}(\mathcal{P}, \frac{\eta}{z}) + \mathcal{N}(\mathcal{V}, \frac{\eta}{z}),$$

contradicting (E4-3). Therefore at least one of $\epsilon_{\rm p}$, $\epsilon_{\rm v}$ is $\geq \eta/z$. Consequently

$$\epsilon_{\mathrm{p}} \, \epsilon_{\mathrm{v}} \, \geq \, \left(\frac{\eta}{z} \right)^2,$$

Define $D_2 := \frac{1}{z^2} > 0$, the η -bound is given by

$$\epsilon_{\rm p} \, \epsilon_{\rm v} \geq D_2 \, \eta^2.$$
(E4-4)

Step 3. Combine the impacts.

By (E4-2) and (E4-4), we have

$$\epsilon_{\mathrm{p}}\epsilon_{\mathrm{v}} \geq D_1 K^{-2/d_{\mathrm{eff}}}$$
 and $\epsilon_{\mathrm{p}}\epsilon_{\mathrm{v}} \geq D_2 \eta^2 \implies \epsilon_{\mathrm{p}}\epsilon_{\mathrm{v}} \geq \max\{D_1 K^{-2/d_{\mathrm{eff}}}, D_2 \eta^2\}.$

This completes the proof.

E.5 Proof of Theorem 2

Restatement of Theorem 2 (Performance Guarantee). *Under Assumption 1 and the covering-regularity of Lemma 3 with constants* $a, a', d_{\text{eff}} > 0$ *and* b > a, b' > a', *for any budget split* $(K_p, K - K_p)$, *covering fold k, and token set* $\mathcal{X} = \mathcal{V} \sqcup \mathcal{P} \subseteq \mathbb{R}^d$ *with* $|\mathcal{V}| = N$, $|\mathcal{P}| = L$, and $d_H(\mathcal{V}, \mathcal{P}) \leq \eta$, *the following hold:*

(a) **Performance bound:** The Performance degradation caused by MoB is upper bounded by

$$\|\mathcal{F}(\mathcal{X}) - \mathcal{F}(\text{MoB}(\mathcal{X}))\| \leq C_{\ell} \max \left\{ \alpha(\eta, k, L) (K_{p})^{-1/d_{\text{eff}}}, \beta (K - K_{p})^{-1/d_{\text{eff}}} \right\} + C_{\ell} \eta,$$

where
$$\alpha(\eta, k, L) = \eta \left(b k L/a\right)^{1/d_{\text{eff}}}, \quad \beta = 2(b')^{1/d_{\text{eff}}}.$$

(b) Multilinear complexity: The complexity of MoB is given by $T_{\text{MoB}} = \mathcal{O}(N(L+K)d)$.

Remark (Coupling Trade-off). Under weak coupling (large $\alpha(\eta, k, L)$), minimizing the bound requires a larger K_p . Conversely, under strong coupling (small $\alpha(\eta, k, L)$), the alignment term decays rapidly, favoring visual preservation (increasing $K - K_p$). Specially, under perfect coupling ($\eta = 0$), the bound simplifies to $\|\Delta y\| \leq C_\ell \beta (K - K_p)^{-1/d_{\rm eff}}$, i.e., MoB reduces to pure visual preservation.

Remark (Budget Scaling). As the total budget K increases, the preservation term $\beta (K - K_p)^{-1/d_{eff}}$ decays, requiring a corresponding increase in K_p (and thus a reduction in the alignment term) to rebalance the trade-off and further lower the overall error bound.

Remark (Scalability). MoB exhibits a multilinear scalability with respect to visual tokens N, prompt tokens L, and retained tokens K (especially $K, L \ll N$), making it readily adaptable to more challenging scenarios, such as advanced MLLMs with higher-resolution inputs or multi-frame video.

Notation.

ullet The intermediate input ${\mathcal X}$ is formulated as

$$\mathcal{X} = \mathcal{V} \, \sqcup \, \mathcal{P} \, \subseteq \mathbb{R}^d \quad \text{where} \quad |\mathcal{V}| = N, \ |\mathcal{P}| = L, \ \text{and} \ N \gg L.$$

Particularly, V, P are compact sets with d_{eff} effective dimensions.

• We define the pruned intermediate input as

$$\mathrm{MoB}(\mathcal{X}) \coloneqq \mathcal{X}_{\mathrm{s}}, \quad \text{where} \quad \mathcal{X}_{\mathrm{s}} = \mathcal{S} \ \sqcup \ \mathcal{P} \quad \text{where} \quad |\mathcal{S}| = K.$$

• The budget configuration is given by $\langle K_p, K_v \rangle$, where $K_p + K_v = K$.

Proof. We separately proof the Performance Guarantee & Complexity in Part A & Part B

Part A: Performance Guarantee

Part A-1: Performance Guarantee of prompt alignment

Step A-1.1: Bound of the radius derived by k-fold NN-covering

Given any union set before K_p -truncation

$$\mathcal{S}_{\mathbf{p}}' \coloneqq \bigcup_{p \in \mathcal{P}} \underset{s_{\mathbf{p}} \in \mathcal{V}}{\arg \operatorname{top-k}}(\cos(s_{\mathbf{p}}, p), k) \quad \text{where} \quad |\mathcal{S}_{\mathbf{p}}'| = K_{\mathbf{p}}' \quad \text{and} \quad K_{\mathbf{p}} \leq K_{\mathbf{p}}' \leq kL,$$

we define

$$\epsilon_{\rm p}' = d_H(\mathcal{S}_{\rm p}', \mathcal{P}).$$

By previous work [15], NN-covering achieves a 1-approximation for the k-center problem with sufficient budget; i.e., specifically for any $p \in \mathcal{P}$ we have

$$\inf_{s_{p}' \in \mathcal{S}_{p}'} \|p - s_{p}'\| = \inf_{v \in \mathcal{V}} \|p - v\|.$$

Thus.

$$\sup_{p \in \mathcal{P}} \inf_{s_p' \in \mathcal{S}_p'} \|p - s_p'\| = \sup_{p \in \mathcal{P}} \inf_{v \in \mathcal{V}} \|p - v\|.$$

Based on Assumption 2, since $s \in \mathcal{S}'_p \subseteq \mathcal{V}$, the upper bound of the radius ϵ'_p is given by

$$\epsilon_{\mathbf{p}}' = d_{H}(\mathcal{S}_{\mathbf{p}}', \mathcal{P}) := \max\{ \sup_{s_{\mathbf{p}}' \in \mathcal{S}_{\mathbf{p}}'} \inf_{p \in \mathcal{P}} \|p - s_{\mathbf{p}}'\|, \sup_{p \in \mathcal{P}} \inf_{s_{\mathbf{p}}' \in \mathcal{S}_{\mathbf{p}}'} \|p - s_{\mathbf{p}}'\|\}$$

$$\leq \max\{ \sup_{v \in \mathcal{V}} \inf_{p \in \mathcal{P}} \|p - v\|, \sup_{p \in \mathcal{P}} \inf_{v \in \mathcal{V}} \|p - v\|\}$$

$$:= d_{H}(\mathcal{V}, \mathcal{P}) \leq \eta.$$
(E5-1)

Step A-1.2: Impact of K_p -truncation on the radius

Based on Lemma 3, we have

$$ar^{-d_{\text{eff}}} < \mathcal{N}(\mathcal{P}, r) < brack r^{-d_{\text{eff}}}$$

In particular:

$$b(\epsilon_{\rm p})^{-d_{\rm eff}} \geq K_{\rm p} \implies \epsilon_{\rm p} \leq \left(\frac{b}{K_{\rm p}}\right)^{1/d_{\rm eff}}$$

and also

$$a (\epsilon_{\mathbf{p}}')^{-d_{\mathrm{eff}}} \leq K_{\mathbf{p}}' \implies \epsilon_{\mathbf{p}}' \geq \left(\frac{a}{K_{\mathbf{p}}'}\right)^{1/d_{\mathrm{eff}}}.$$

Combining the upper and lower bound for ϵ_p and ϵ'_p , respectively in terms of b, K_p, K'_p , we obtain

$$\epsilon_{\rm p} \; \leq \; \left(\frac{b}{K_{\rm p}}\right)^{1/d_{\rm eff}} \; = \; \left(\frac{bK_{\rm p}'}{aK_{\rm p}}\right)^{1/d_{\rm eff}} \; \cdot \; \left(\frac{a}{K_{\rm p}'}\right)^{1/d_{\rm eff}} \; \leq \; \left(\frac{bK_{\rm p}'}{aK_{\rm p}}\right)^{1/d_{\rm eff}} \; \epsilon_{\rm p}'.$$

That is, truncating from $K'_{\rm D}$ to $K_{\rm D}$ centers increases the radius by at most the factor

$$\epsilon_{\rm p} \leq \left(bK_{\rm p}'/aK_{\rm p}\right)^{1/d_{\rm eff}} \epsilon_{\rm p}'.$$

Since $kL \geq K_{\rm p}'$, loading into above, we have

$$\epsilon_{\rm p} \leq \left(bkL/aK_{\rm p}\right)^{1/d_{\rm eff}} \epsilon_{\rm p}'.$$

By loading (E5-1) into the above, the performance guarantee of prompt alignment is given by

$$\epsilon_{\mathbf{p}} \coloneqq d_H \big(\mathcal{S}_{\mathbf{p}}, \, \mathcal{P} \big) \, \, \leq \, \, \alpha(\eta, k, L) \, (K_{\mathbf{p}})^{-1/d_{\mathrm{eff}}} \quad \text{where} \quad \alpha(\eta, k, L) \coloneqq \eta \, \big(bkL/a \big)^{1/d_{\mathrm{eff}}}. \quad \text{(E5-2)}$$

Part A-2: Performance Guarantee of Visual Preservation

By previous work [36], FPS achieves a 2-approximation for the k-center problem:

$$\epsilon_{\rm v} < 2 \, \epsilon^{\star}(K_{\rm v}),$$
 (E5-3)

where $\epsilon^*(K_v)$ is the optimal radius with K_v centers. Based on Lemma 3, we have

$$\mathcal{N}(\mathcal{V}, r) \leq b' r^{-d_{\text{eff}}},$$

thereby, the upper bound of optimal radius is given by

$$\epsilon^{\star}(K_{\rm v}) \leq (b'/K_{\rm v})^{1/d_{\rm eff}}$$
.

By loading the above into (E5-3), the performance guarantee of visual preservation is given by

$$\epsilon_{\mathbf{v}} \coloneqq d_H(\mathcal{S}_{\mathbf{v}}, \mathcal{V}) \le \beta (K_{\mathbf{v}})^{-1/d_{\mathrm{eff}}}, \quad \text{where} \quad \beta \coloneqq 2 b'^{1/d_{\mathrm{eff}}}.$$
 (E5-4)

Part A-3: Performance Guarantee of MoB

By substituting (E5-2) and (E5-4) into Lemma 2, the performance guarantee of the MoB is given by:

$$\|\mathcal{F}(\mathcal{X}) - \mathcal{F}(\mathrm{MoB}(\mathcal{X}))\| \leq C_{\ell} \, \max \left\{ \alpha(\eta, k, L) \, (K_{\mathrm{p}})^{-1/d_{\mathrm{eff}}}, \, \beta \, (K_{\mathrm{v}})^{-1/d_{\mathrm{eff}}} \right\} \, + \, C_{\ell} \, \eta,$$

where
$$\alpha(\eta, k, L) = \eta \left(b \, k \, L/a \right)^{1/d_{\rm eff}}, \quad \beta = 2 \, b'^{1/d_{\rm eff}}.$$

This completes the proof of Part A.

Part B: Complexity

Since $k \ll K_{\rm p} \leq K \sim L \ll N$, we restrict our complexity analysis to the leading-order terms.

Part B-1: Normalization

MoB do a L_2 normalization for each token $x \in \mathcal{X} \subseteq \mathbb{R}^d$; thus, the complexity is given by

$$T_{\text{norm}} = \mathcal{O}((N+L) d). \tag{E5-5}$$

Part B-2: Selection of Prompt Center

Firstly, MoB calculates the cosine similarity with each $p \in \mathcal{P}$ and $v \in \mathcal{V}$ via a matrix multiplication:

$$\mathbf{M}_{\mathrm{sim}} = \mathbf{P} \mathbf{V}^{\top} \in \mathbb{R}^{L \times N}$$
 where $\mathbf{V} \in \mathbb{R}^{N \times d}$ and $\mathbf{P} \in \mathbb{R}^{L \times d}$,

which leads a complexity of $T_{\text{step 1-1}} = \mathcal{O}(N \, L \, d)$. Subsequent, MoB do a top-k retrieval in the first dimension of \mathbf{M}_{sim} the select k most closed centers for each prompt token $p \in \mathcal{P}$, which can be reduced to a partial sorting, thereby leading to a complexity of $T_{\text{step 1-2}} = \mathcal{O}(N \, L \, \log k)$. Finally, MoB merge the selected result of each $p \in \mathcal{P}$, and truncated the top- K_p ones with largest similarity, leading to a $T_{\text{step 1-3}} = \mathcal{O}(L \, k \, \log K_p)$. Consequently, the total complexity $T_{p\text{-select}}$ of prompt center selection is given by:

$$T_{\text{p-select}} = T_{\text{step 1-1}} + T_{\text{step 1-2}} + T_{\text{step 1-3}},$$

$$= \mathcal{O}(N L d) + \mathcal{O}(N L \log k) + \mathcal{O}(L k \log K_{\text{p}}),$$

$$= \mathcal{O}(N L d).$$
(E5-6)

Part B-3: Selection of Visual Center

Initially, MoB calculates the minimum distance (used in FPS) with each visual token $v \in \mathcal{V} \setminus \mathcal{S}_p := \mathcal{V}'$ and the selected prompt centers via a matrix multiplication together with an argmin operator:

 $\mathbf{d}_{\mathrm{FPS}} = \arg\min \mathbf{V}'^{\top} \mathbf{S}_{\mathrm{p}} \in \mathbb{R}^{N-K_{\mathrm{p}}}$ where $\mathbf{V}' \in \mathbb{R}^{(N-K_{\mathrm{p}}) \times d}$ and $\mathbf{S}_{\mathrm{p}} \in \mathbb{R}^{K_{\mathrm{p}} \times d}$, thus, the complexity is given by

$$\begin{split} T_{\text{step 2--1}} &= \underbrace{\mathcal{O}((N-K_{\text{p}})\,K_{\text{p}}\,d)}_{\text{matrix multiplication}} + \underbrace{\mathcal{O}((N-K_{\text{p}})\,K_{\text{p}})}_{\text{argmin}}, \\ &= \mathcal{O}((N-K_{\text{p}})\,K_{\text{p}}\,d). \end{split}$$

Subsequently, in $K-K_{\rm p}$ iterations, MoB add the tokens with largest minimum distance with an argmax operator in $\mathbf{d}_{\rm FPS}$, and update the $\mathbf{d}_{\rm FPS}$ with an inner production together with an $N-K_{\rm p}$ -dimensional element-wise comparison; thus the complexity is given by

$$T_{\text{step 2-2}} = \underbrace{\mathcal{O}((N-K_{\text{p}})(K-K_{\text{p}}))}_{\text{argmax}} + \underbrace{\mathcal{O}((K-K_{\text{p}})\,N\,d)}_{\text{inner production}} + \underbrace{\mathcal{O}((K-K_{\text{p}})\,d)}_{\text{ele-wise comparision}} \,,$$

$$= \mathcal{O}((K-K_{\text{p}})\,N\,d).$$

Consequently, the total complexity $T_{v-select}$ of visual center selection is given by:

$$T_{\text{v-select}} = T_{\text{step 2-1}} + T_{\text{step 2-2}},$$

$$= \mathcal{O}((N - K_{\text{p}}) K_{\text{p}} d) + \mathcal{O}((K - K_{\text{p}}) N d),$$

$$= \mathcal{O}(N K d).$$
(E5-7)

Part B-4: Totally complexity

By (E5-5), (E5-6) and (E5-7), the totally complexity of MoB is given by

$$\begin{split} T_{\text{MoB}} &= T_{\text{norm}} + T_{\text{p-select}} + T_{\text{v-select}}, \\ &= \mathcal{O}((N+L)\,d) + \mathcal{O}(N\,L\,d) + \mathcal{O}(N\,K\,d), \\ &= \mathcal{O}(N\,L\,d) + \mathcal{O}(N\,K\,d), \\ &= \mathcal{O}(N\,(L+K)\,d). \end{split}$$

This completes the proof of Part B.

Combining the Part A & B, we complete the proof.