

Abstract

*Multimodal large language models (MLLMs) have shown remarkable capabilities in a wide range of vision-language tasks. However, the large number of visual tokens introduces significant computational overhead. To address this issue, visual token pruning has emerged as a key technique for enhancing the efficiency of MLLMs. In cognitive science, humans tend to first determine which regions of a scene to attend to (“where to look”) before deciding which specific elements within those regions to process in detail (“what to select”). This two-stage strategy enables the visual system to efficiently allocate attention at a coarse spatial level before performing fine-grained selection. However, existing pruning methods primarily focus on directly optimizing “what to select”, typically using attention scores or similarity metrics. They rarely consider “where to look”, which has been shown to lead to inefficient spatial allocation, positional bias, and the retention of irrelevant or redundant tokens. In this paper, we propose **GridPrune**, a method that replaces the global Top-K mechanism with a “guide-globally, select-locally” zonal selection system. GridPrune splits the pruning process into two steps: first, it uses text-conditional guidance to dynamically allocate a token budget across spatial zones; and then, it performs local selection within each budgeted zone. Experimental results demonstrate that GridPrune achieves superior performance across various MLLM architectures. On LLaVA-NeXT-7B, GridPrune retains 96.98% of the full performance while using 11.1% of the tokens, outperforming the best-performing baseline by 2.34% at the same pruning rate.*

1. Introduction

Multimodal Large Language Models (MLLMs)[3, 10, 11, 25, 26, 47] have shown remarkable capabilities in a wide range of vision-language tasks, such as visual question answering (VQA) and complex reasoning[14, 17, 19]. In these models, an image is processed by a vision encoder[31, 44]

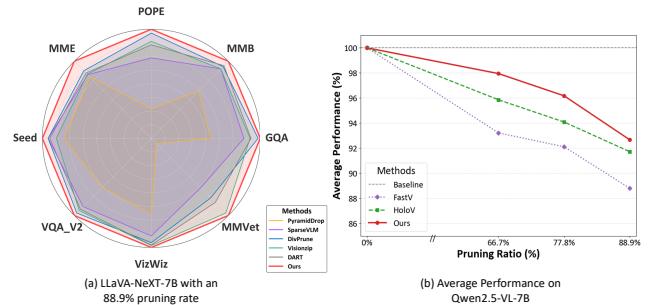


Figure 1. Performance comparison of GridPrune against state-of-the-art methods across various MLLM architectures. (a) presents results on the high-resolution LLaVA-NeXT-7B, with 11.1% of visual tokens retained. (b) shows the average performance trend on Qwen2.5-VL-7B as the token retention ratio varies.

and a projector[6, 22] to generate a sequence of visual tokens, which are then fed into the LLM[2, 5, 34]. However, this introduces significant computational overhead. In a typical input, the number of visual tokens can reach hundreds, often far exceeding the length of text tokens. Since the computational complexity of the self-attention mechanism scales quadratically with the sequence length, a large number of visual tokens makes using MLLMs costly. When processing high-resolution images[9, 21, 27] or video streams[24, 28, 37], the number of tokens increases even further. Therefore, developing effective visual token pruning strategies is crucial for enhancing the inference efficiency of MLLMs.

Extensive efforts have been made in visual token pruning to reduce the inference cost of MLLMs[8, 18, 46]. Early mainstream methods can be broadly categorized into two types: attention-based[8, 40, 46] and similarity-based[1, 4, 39]. Attention-based methods typically rely on attention scores as an importance signal for token selection. However, they suffer from two issues: first, they are susceptible to position bias in global ranking, leading to unreliable selection results[38, 45, 48], as shown in Fig. 2 (a); second, the selected tokens are often highly redundant[48], because nearby tokens with similar visual features tend to

*Corresponding author.

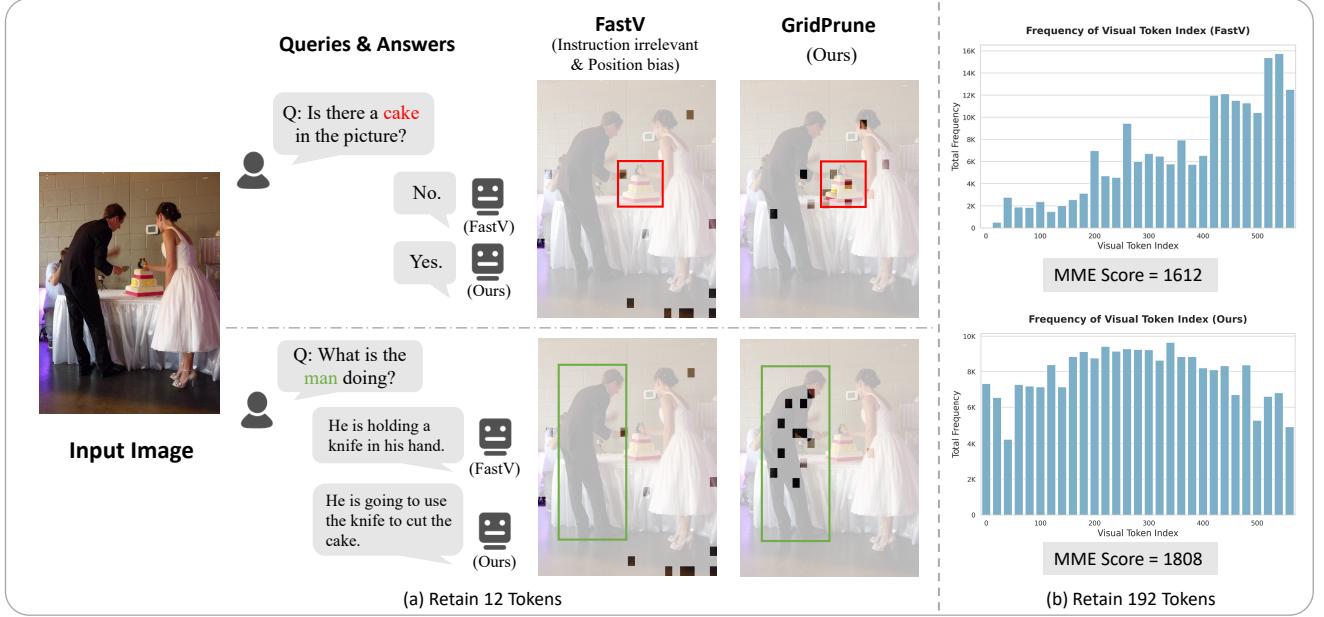


Figure 2. Comparison of GridPrune with FastV. (a) In a direct comparison, FastV’s selection is guided by a positional bias towards final tokens, while GridPrune’s is guided by the query’s semantic content. (b) Statistical analysis on the MME benchmark at scale shows that the histogram of selected indices exhibits a massive spike for FastV at the end of the sequence, revealing a strong positional bias. This is inefficient, as important content in images is typically centered or evenly distributed, rather than confined to one corner. In contrast, GridPrune’s distribution is more balanced.

receive same high attention scores. On the other hand, similarity-based methods, which mainly aim to construct a generic token subset, are often instruction-irrelevant, resulting in low informational efficiency when handling specific tasks[41, 45]. Subsequent work attempted to improve this from different angles. However, these methods either pursue diversity through task conditioning at the global level[45], which lacks effective spatial planning, or rely on crops that are task-unaware and cannot be dynamically adjusted to specific questions[48]. **Overall, existing research has focused mainly on how to optimize the “what to select” problem, but has overlooked another problem, “where to look”, which can lead to inefficient spatial allocation, positional bias, and the retention of irrelevant or redundant tokens.**

Studies in cognitive science provide inspiration for this challenge. Research shows that when observing a scene, humans tend to first determine “where to look” before deciding “what to select” [35]. This two-stage strategy allows the visual system to efficiently allocate attention at a coarse spatial level before performing fine-grained selection. We hypothesize that the problems in current methods, like positional bias and redundancy, are caused by neglecting the “where to look” step. This is because deciding “where to look” first allows the model to focus its budget on important regions, instead of wasting it on irrelevant backgrounds.

This makes the “what to select” step more effective.

In this paper, we propose GridPrune, a pruning method that incorporates this two-stage strategy into MLLMs. GridPrune splits the pruning process into two steps. First, it uses the text query as a high-level command to make a task-driven decision, dynamically giving the limited token budget to image zones that are relevant to the task, which solves the “where to look” problem. Second, inside each zone that gets a budget, a selection process is carried out to pick the most informative tokens based on a fused score of text relevance and visual saliency, which solves the “what to select” problem. GridPrune changes a global optimization problem into a set of local decisions, which not only reduces problems like position bias but also improves information efficiency by giving the token budget to instruction-relevant locations, as shown in Fig. 2. Additionally, our method adapts to high-resolution application scenarios. We test the effectiveness of GridPrune on ten benchmarks, demonstrating that GridPrune outperforms existing state-of-the-art approaches, especially under aggressive pruning ratios, as shown in Fig. 1. For instance, on LLaVA-NeXT-7B (Table 4), GridPrune retains 96.98% of the original performance while using only 11.1% of the tokens, which is 2.34% higher than the best-performing baseline. At a lower retention rate of 5.6%, GridPrune outperforms the best-performing baseline by 3.1% in average performance.

In summary, our contributions are summarized as follows:

- We identify the limitation of existing visual token pruning methods: they primarily focus on the “what to select” problem, while rarely considering “where to look”, which leads to issues such as inefficient spatial allocation and positional bias.
- We propose GridPrune, a training-free method based on above two-stage method. It divides the pruning process into a task-driven budget allocation and a series of intra-zone selections.
- Extensive experiments show that GridPrune performs better than state-of-the-art methods across diverse MLLM architectures. It can keep most of the model’s original performance even with very few tokens.

2. Related Work

2.1. Multimodal large language models.

The remarkable achievements of Large Language Models (LLMs)[2, 5, 34] have encouraged the development of Multimodal Large Language Models (MLLMs)[3, 10, 11, 25, 26, 47]. The dominant paradigm for these models is to connect a pretrained vision encoder[31, 44] to an LLM via a projector[6, 22]. This architecture converts visual inputs into a token sequence to take advantage of the powerful capabilities of the LLM. However, this introduces significant computational challenges. The token sequence encoded from the images is typically much longer than the text input[8, 40]. For example, LLaVA-1.5[26] converts a standard image into 576 visual tokens. The number of tokens increases further in scenarios that involve higher-resolution images[9, 21, 27] or video streams[24, 28, 37]. These lengthy visual tokens lead to high inference costs, because the self-attention mechanism in LLMs has a computational complexity that scales quadratically with sequence length. Therefore, developing efficient visual token pruning strategies is crucial for improving the inference efficiency of MLLMs.

2.2. Token Pruning.

To reduce the computational burden of MLLMs, researchers have made many efforts in visual token pruning [8, 18, 46]. Early mainstream methods can be roughly divided into two types: attention-based [8, 40, 46] and similarity-based [1, 4, 39].

Attention-based methods rely on attention scores as the main importance signal. FastV [8] computes the average attention score that a token receives from all other tokens to judge its importance; PyramidDrop [40] uses a multi-stage pruning strategy; and SparseVLM [46] uses the attention between instruction tokens and visual tokens. However, this type of method has two main problems: first, it is easily

affected by position bias in global ranking, which leads to unreliable selection results [38, 45, 48]; second, the selected tokens are often highly redundant, because nearby tokens tend to receive similarly high scores [48].

Similarity-based methods focus on building a general token subset. ToMe [4] progressively combines similar tokens through token merging; DART [39] aims to select a representative subset from the feature space; and DivPrune [1] frames the problem as a Max-Min Diversity Problem to maximize the diversity among selected tokens. VisionZip [41] also merges similar tokens to reduce redundancy. The main limitation of these methods is that they are often instruction-irrelevant, which leads to low efficiency when processing task-specific queries [41, 45].

Later studies attempted to improve this from different perspectives. CDPruner [45] introduced text conditioning to guide diversity selection, but the text only serves as a way to adjust for diversity. This limits the potential for the task instruction to perform spatial planning. HoloV [48] uses spatial crops to preserve context, but its strategy is task-unaware. Overall, existing research has mainly focused on how to score or diversify individual tokens (“what to select”). In contrast, how to allocate the limited token budget spatially (“where to look”) has received less attention, which can lead to inefficient spatial allocation, positional bias, and the retention of irrelevant or redundant tokens.

3. Method

This section introduces the proposed GridPrune method. We first redefine the visual token pruning problem, and then introduce the two core components of GridPrune: a Dual-Source Importance Scoring function and the “guide-globally, select-locally” zonal selection system, as shown in Fig. 3.

3.1. Preliminaries and Problem Formulation

In MLLMs, a vision encoder transforms an image into a large set of N patch tokens, $V = \{v_1, \dots, v_N\}$. The goal of visual token pruning is to select a small subset $V' \subset V$ of size $k \ll N$ to reduce the computational load, which is dominated by the $O(N^2)$ self-attention complexity in the subsequent LLM.

The current pruning paradigm operates on the principle of global importance. It first assigns a scalar importance score s_i to each token v_i and then performs a global Top-K selection:

$$V' = \operatorname{argmax}_{V_{sub} \subset V, |V_{sub}|=k} \sum_{v_i \in V_{sub}} s_i \quad (1)$$

This global sorting method, although intuitive, has its limitations. For example, attention-based methods often derive the scores s_i from ViT[12], which are affected by

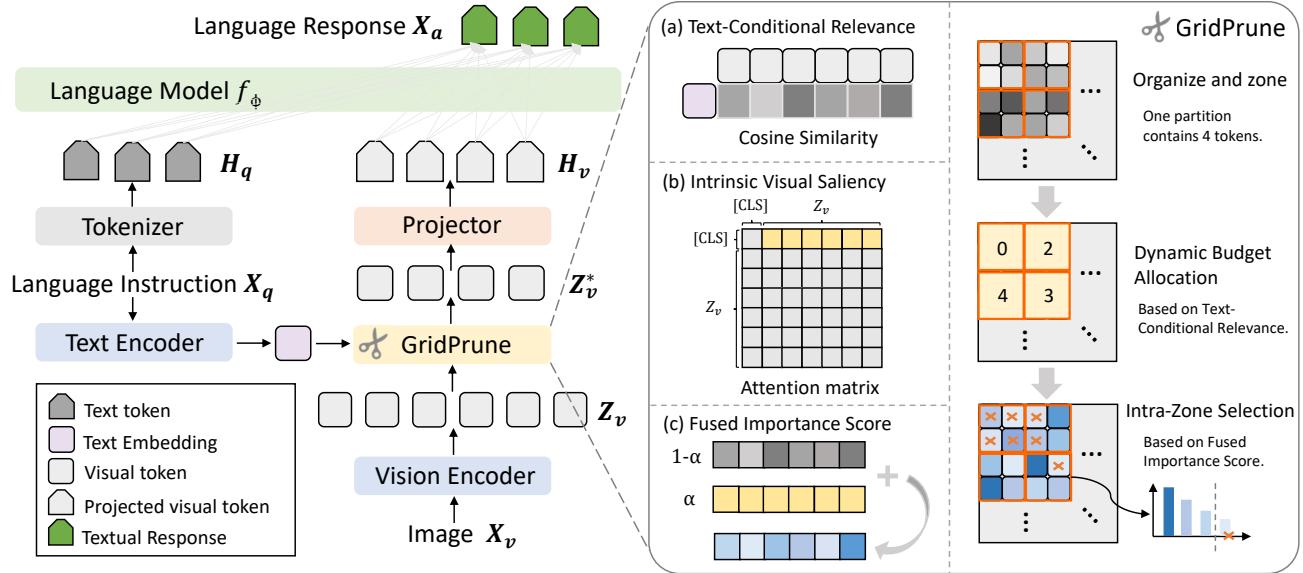


Figure 3. An overview of the GridPrune method. We first calculate two scores for each visual token: (a) Text-Conditional Relevance, derived from the cosine similarity between token embeddings and the text embedding (obtained from the CLIP text encoder using the user’s prompt as input), and (b) Intrinsic Visual Saliency, extracted from the vision encoder’s attention matrix. These are combined into (c) Fused Importance Score via α . GridPrune follows a “guide-globally, select-locally” process: (1) the tokens are partitioned into zones; (2) a token budget is dynamically allocated to these zones based on their aggregate text-conditional relevance; and (3) a local Top-K selection is performed within each zone using the fused importance score to select the final token set. This mechanism ensures a selection that is both query-aware and spatially balanced.

Algorithm 1 Budget Allocation

Require: Vector of float budgets $B = \{b_1, \dots, b_M\}$; retained tokens k ; zone capacity $c = \text{block_size}^2$.
Ensure: Vector of integer budgets $K = \{k_1, \dots, k_M\}$.

```

1:  $K_j \leftarrow \min(\lfloor B_j \rfloor, c)$ ,  $\forall j \in \{1, \dots, M\}$   $\triangleright$  Initialize integer budgets with capped floor values
2:  $k' \leftarrow k - \sum_{j=1}^M K_j$   $\triangleright$  Calculate the number of unallocated tokens
3:  $F_j \leftarrow B_j - \lfloor B_j \rfloor$ ,  $\forall j \in \{1, \dots, M\}$ 

4: for  $i = 1 \rightarrow k'$  do
5:    $j^* \leftarrow \operatorname{argmax}_{j \text{ s.t. } K_j < c} F_j$   $\triangleright$  Find the zone with the highest priority
6:    $K_{j^*} \leftarrow K_{j^*} + 1$   $\triangleright$  Allocate one token to the selected zone
7:    $F_{j^*} \leftarrow -\infty$   $\triangleright$  Prevent re-selection of the same zone
8: end for

9: return  $K$ 

```

positional bias, leading the selection process to systematically over-sample from certain spatial regions while ne-

glecting others. Instead of seeking the globally “most important” tokens, we propose a two-stage process: first, use text-conditional guidance to dynamically allocate a token budget across spatial zones, and then perform local selection within each budgeted zone.

3.2. Dual-Source Importance Score

In order to achieve effective selection within our two-stage system, we first require a robust scoring function. Therefore, we propose a Dual-Source Importance Scoring function that combines two complementary sources of information: the conditional relevance to the user’s query and the model’s intrinsic visual saliency.

Text-Conditional Relevance. To make the pruning process respond to the user’s prompt, we compute a relevance score r_i for each patch token. We leverage the hidden states $\{h_i\}$ from CLIP’s penultimate layer, as they are already enriched with contextual information. These hidden states are projected using the vision tower’s own projection layer, g_v , to obtain multimodal visual embeddings. The relevance score r_i is then the cosine similarity between this projected embedding and the text embedding (the [EOS] token [36, 43]) from the CLIP text encoder f_t :

$$r_i = \frac{g_v(h_i) \cdot f_t(Q)}{\|g_v(h_i)\| \cdot \|f_t(Q)\|} \quad (2)$$

Algorithm 2 GridPrune for a High-Resolution Image

Require: A high-resolution image I , processed into a set of S sub-images (e.g., 4 high-res tiles + 1 low-res global view), $I = \{p_1, \dots, p_S\}$. A single text prompt T .

Ensure: The final visual feature sequence F' .

```

1: for all  $i \in \{1, \dots, S\}$  do
2:    $p'_i \leftarrow \text{GridPrune}(p_i, T)$      $\triangleright$  Apply pruning to each
      sub-image
3: end for

4:  $P'' \leftarrow \{p'_1, \dots, p'_S\}$          $\triangleright$  Collect pruned token sets
5:  $\tilde{P} \leftarrow g_{\text{mm\_proj}}(P'')$      $\triangleright$  Project all features in one batch
6:  $F' \leftarrow \bigoplus_{i=1}^S \tilde{p}_i$          $\triangleright$  Concatenate features into final
      sequence

7: return  $F'$ 

```

This operation is carried out completely independently of the visual tower module. By using a highly text-aware representation method, we ensure that the relevance score depends on the user’s intention.

Intrinsic Visual Saliency. While text relevance is important, the model’s intrinsic understanding of salient regions provides another valuable signal. We derive this signal from the multi-head self-attention mechanism in the penultimate layer of the CLIP vision encoder. Specifically, we obtain the attention weights from the [CLS] token to all patch tokens, and average these weights across all attention heads. The scores are then min–max normalized to the $[0, 1]$ range. The final score a_i reflects the model’s assessment of which vision token holds the most information for the overall image representation.

Fused Importance Score. Finally, we fuse these two scores into a single importance score s_i . To ensure both scores are on a comparable scale, we first normalize the text-conditional relevance score r_i from its native range of $[-1, 1]$ to $[0, 1]$:

$$\hat{r}_i = (r_i + 1)/2 \quad (3)$$

The final fused score s_i is then a weighted linear combination of the normalized relevance score \hat{r}_i and the saliency score a_i :

$$s_i = (1 - \alpha) \cdot \hat{r}_i + \alpha \cdot a_i \quad (4)$$

where the hyperparameter $\alpha \in [0, 1]$ balances the influence of text-conditional guidance and intrinsic visual saliency. As suggested by prior work [38], the relevance and saliency score scales are not consistent across different architectures, and the ideal trade-off shifts with the pruning rate. Therefore, we tune this parameter for each model

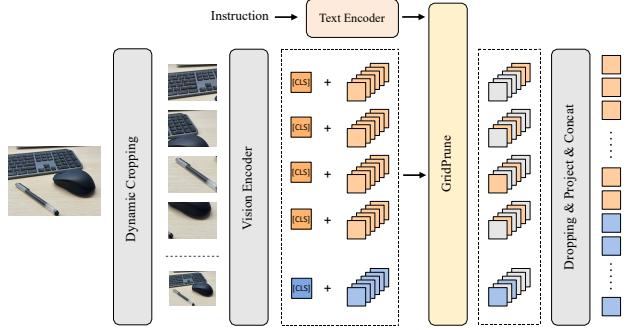


Figure 4. The processing flow of GridPrune applied to LLaVA-NeXT. The image is first dynamically cropped into multiple sub-images, and each sub-image independently passes through the vision encoder. Then, guided by the instruction, GridPrune prunes the tokens of each sub-image separately. Finally, all the retained tokens are projected and concatenated to form the final sequence.

and rate, which leads to a more effective balance and better overall performance.

These scores lay the foundation for our two-stage selection strategy. The text-conditional relevance (r_i) is used in the “where to look” stage, while the fused importance score (s_i) is used in the “what to select” stage. This design mimics how the brain operates: r_i first provides a task-driven answer for “where to look”, preventing irrelevant but visually prominent areas from competing for the budget. Then, s_i solves the problem of “what to select” by integrating all the information in these areas.

3.3. Zonal Selection System

The proposed “guide-globally, select-locally” zonal selection system replaces the conventional global Top-K mechanism, enabling a more spatially aware and efficient allocation of the token budget. First, it uses text-conditional guidance to dynamically allocate a token budget across spatial zones, and then performs local selection within each budgeted zone.

Zone Partitioning. The image’s N patch tokens are organized into a 2D grid and partitioned into M non-overlapping square zones of a predefined block size. To determine the budget for each zone, we compute an aggregate relevance score, \bar{r}_j , by averaging the text-conditional relevance scores r_i of all tokens within zone j .

Dynamic Budget Allocation. We dynamically allocate the total token budget k across the M zones. This is accomplished by converting the zone scores $\{\bar{r}_j\}$ into a probability distribution using a standard softmax function:

$$P_j = \frac{\exp(\bar{r}_j)}{\sum_{m=1}^M \exp(\bar{r}_m)} \quad (5)$$

The budget k_j for each zone is then determined by a

rounding procedure applied to $P_j \cdot k$ that ensures $\sum k_j = k$, detailed in Algorithm 1. This step translates text relevance into a spatial allocation plan.

Intra-Zone Selection with Fused Score. With the budget k_j for each zone finalized, the final selection is made locally. Within each zone j , we select the k_j tokens that have the highest fused importance scores s_i . By separating the query-aware ‘where to look’ step from the local ‘what to select’ step, this two-stage process addresses the limitations of global Top-K selection, such as inefficient spatial allocation, positional bias, and redundancy.

Table 1. Optimal settings for the fusion hyperparameter α used across all main experiments.

Base Model	Retention ratio	Value of α
<i>Standard Resolution</i>		
LLaVA-1.5-7B	33.3%	0.8
	22.2%	0.7
	11.1%	0.7
<i>High Resolution</i>		
LLaVA-NeXT-7B	22.2%	0.8
	11.1%	0.7
	5.6%	0.7
<i>Qwen Architecture</i>		
Qwen2.5-VL-7B	33.3%	0.7
	22.2%	0.7
	11.1%	0.7

4. Experiments

4.1. Experimental Setup

Models and Baselines. To validate the effectiveness and generalizability of GridPrune, we conduct experiments on a diverse set of three multimodal large language models. LLaVA-1.5-7B[26] serves as a widely adopted benchmark model for standard-resolution inputs. LLaVA-NeXT-7B[27] represents models designed for high-resolution imagery, processing a larger sequence of tokens. Qwen2.5-VL-7B[3], featuring a distinct architecture from the LLaVA series, is included to assess the architectural generalization of our method. We compare GridPrune against a comprehensive set of state-of-the-art (SOTA) training-free pruning methods, including ToMe[4], FastV[8], PyramidDrop[40], SparseVLM[46], DivPrune[1], VisionZip[41], LLaVA-PruMerge[32], DART[39], and HoloV[48].

Benchmarks. We conduct a comprehensive evaluation across a suite of ten challenging benchmarks to assess a wide range of multimodal capabilities. The suite includes general visual question answering (VQAv2[13], GQA[16], VizWiz[15]); text-intensive VQA (TextVQA[33]); scientific reasoning (ScienceQA[30]); compositional understanding (MME[7], SEED-Bench[20]); object hallucination detection (POPE[23]); and holistic model assessment

Table 2. Performance comparison on Qwen2.5-VL-7B. Avg.(%) represents the average percentage of performance maintained.

Methods	Avg.(%)	SQA ^{IMG}	POPE	MMB ^{EN}	MME
<i>Upper Bound (100%)</i>					
Qwen2.5-VL-7B	100%	84.7	86.1	82.8	2304
<i>Retain 33.3% Tokens</i>					
FastV (ECCV24)	92.4%	78.5	82.2	75.7	2072
HoloV (NeurIPS25)	94.6%	79.8	85.0	78.3	2093
GridPrune (Ours)	97.6%	84.4	84.2	79.8	2221
<i>Retain 22.2% Tokens</i>					
FastV (ECCV24)	91.2%	78.0	80.7	74.9	2036
HoloV (NeurIPS25)	92.7%	79.8	82.3	76.5	2043
GridPrune (Ours)	95.6%	83.1	82.8	78.0	2161
<i>Retain 11.1% Tokens</i>					
FastV (ECCV24)	87.6%	77.4	78.6	69.2	1940
HoloV (NeurIPS25)	90.5%	79.5	80.7	72.4	2006
GridPrune (Ours)	91.3%	80.7	78.1	76.2	2009

(MMBench[29], MM-Vet[42]).

Implementation Details. Across all experiments, GridPrune is configured with a block size of 2, partitioning the standard 24×24 patch grid into $M = 144$ zones. The fusion hyperparameter α , which balances text-conditional guidance with intrinsic visual saliency, is tuned for different models and retention ratios, with specific values reported in Table 1. The performance of the unpruned, full-token model is reported as the Upper Bound. All methods are evaluated at multiple token retention ratios, and all experiments are conducted on RTX 3090 GPUs (24 GB).

4.2. Main Results

We present the main results of our comparative experiments in Tables 2, 3, and 4. The findings across all models and benchmarks demonstrate the superior performance of GridPrune compared to existing SOTA methods.

Performance on LLaVA-1.5-7B. As shown in Table 3, the average performance of GridPrune surpasses all other methods on the LLaVA-1.5-7B model. At a retention ratio of 33.3%, GridPrune achieves an average score of 99.81% relative to the unpruned model, matching its performance while using only one-third of the visual tokens. The advantage of GridPrune becomes even more pronounced under higher pruning settings. When retaining 11.1% of the tokens, GridPrune maintains a remarkable 96.76% of the original performance. In contrast, other methods experience a more evident performance drop, with GridPrune outperforming the best-performing baseline by 1.76 percentage points. This demonstrates the robustness of our method in preserving critical visual information, even with a very small token budget.

Performance on LLaVA-NeXT-7B. The algorithm 2 and Fig. 4 show the logic of applying GridPrune to LLaVA-NeXT. The effectiveness of GridPrune is evident in the high-resolution setting of LLaVA-NeXT-7B, where the ini-

Table 3. Main results on the standard-resolution benchmark LLaVA-1.5-7B.

Methods	Avg.(%)	SQA ^{IMG}	GQA	MMB ^{EN}	POPE	TextVQA	MME	Seed ^I	VQA ^{V2}	VizWiz	MMVet
<i>Upper Bound, 576 Tokens (100%)</i>											
LLaVA-1.5-7B	100%	69.5	61.9	64.6	85.9	58.2	1862	66.2	78.5	50.1	31.1
<i>Retain 192 Tokens (33.3%)</i>											
ToMe (ICLR23)	-	65.2	54.3	60.5	72.4	52.1	1563	-	68.0	-	-
FastV (ECCV24)	89.13%	67.3	52.7	61.2	64.8	52.5	1612	57.2	67.1	50.8	27.7
PyramidDrop (CVPR25)	96.20%	69.2	57.1	63.2	82.3	56.1	1797	58.2	75.1	51.1	30.5
DivPrune (CVPR25)	98.95%	68.9	59.9	62.3	87.0	56.9	1762	64.2	76.8	54.9	30.8
Visionzip (CVPR25)	98.30%	68.8	59.3	62.9	85.5	57.2	1769	63.2	76.8	51.5	31.7
SparseVLM (ICML25)	97.45%	69.1	57.6	62.5	83.6	56.1	1787	64.2	75.6	50.6	31.5
LLaVA-PruMerge (ICCV25)	-	67.9	54.3	59.6	71.3	54.3	1632	-	70.6	50.1	-
DART (EMNLP25)	98.61%	69.8	58.9	63.6	82.8	57.4	1834	64.6	76.7	51.1	31.5
HoloV (NeurIPS25)	-	69.8	59.0	65.4	85.6	57.4	1820	-	76.7	50.9	-
GridPrune (Ours)	99.81%	68.6	60.7	63.7	86.3	57.0	1808	64.5	77.4	52.1	33.3
<i>Retain 128 Tokens (22.2%)</i>											
ToMe (ICLR23)	-	59.6	52.4	53.3	62.8	49.1	1343	-	63.0	-	-
FastV (ECCV24)	86.94%	68.5	54.0	56.1	68.2	56.4	1490	52.2	71.0	51.9	27.0
PyramidDrop (CVPR25)	93.30%	68.4	57.1	62.3	77.5	56.7	1761	54.1	74.3	49.4	27.6
DivPrune (CVPR25)	97.43%	68.6	59.4	61.5	87.0	55.9	1718	62.4	76.0	52.8	30.6
Visionzip (CVPR25)	97.47%	68.9	57.7	62.0	83.2	56.8	1757	61.3	75.6	52.0	32.6
SparseVLM (ICML25)	96.24%	69.0	57.3	62.6	83.1	56.3	1746	63.6	75.1	49.7	29.7
LLaVA-PruMerge (ICCV25)	-	67.1	53.3	58.1	67.2	54.3	1554	-	68.8	50.3	30.4
DART (EMNLP25)	97.50%	69.1	57.9	63.2	80.1	56.4	1845	63.4	75.9	51.7	30.9
HoloV (NeurIPS25)	-	69.8	57.7	63.9	84.0	56.8	1802	-	75.5	51.5	-
GridPrune (Ours)	98.12%	68.5	59.6	62.4	86.2	54.9	1744	62.9	76.2	52.7	32.4
<i>Retain 64 Tokens (11.1%)</i>											
ToMe (ICLR23)	-	50.0	48.6	43.7	52.5	45.3	1138	-	57.1	-	-
FastV (ECCV24)	74.02%	68.4	46.0	50.1	35.5	51.6	1255	41.4	55.9	49.1	18.9
PyramidDrop (CVPR25)	76.11%	69.0	46.1	48.0	40.8	50.6	1561	48.8	56.3	46.3	17.7
DivPrune (CVPR25)	95.00%	68.0	57.5	60.1	85.5	54.5	1674	60.5	74.1	53.6	28.1
Visionzip (CVPR25)	94.13%	69.0	55.1	60.1	77.0	55.5	1687	57.7	72.4	52.9	30.9
SparseVLM (ICML25)	87.78%	69.2	52.0	58.3	69.7	52.1	1589	56.7	66.9	49.4	24.4
LLaVA-PruMerge (ICCV25)	-	68.1	51.9	55.3	65.3	54.0	1549	-	67.4	50.1	28.0
DART (EMNLP25)	92.88%	69.8	55.9	60.6	73.9	54.4	1765	59.3	72.4	51.6	26.5
HoloV (NeurIPS25)	-	69.5	55.3	63.3	80.3	55.4	1715	-	72.8	52.8	-
GridPrune (Ours)	96.76%	68.2	58.7	62.3	85.8	54.3	1719	62.1	75.3	54.5	29.3

tial number of tokens is much larger. Table 4 and Fig. 1 (a) show that the average performance of GridPrune outperforms all baselines. For instance, when retaining 320 tokens, GridPrune achieves an average performance of 96.98%, which is 2.34% higher than that of DivPrune. When only 160 tokens are kept, GridPrune still maintains 94.65% of the upper-bound performance. This result shows that GridPrune enables high-resolution models to operate at a fraction of their original computational cost.

Performance on Qwen2.5-VL-7B. To verify the architectural generalizability of our approach, we also evaluate GridPrune on Qwen2.5-VL-7B. We calculate text relevance using the cosine similarity between the instruction vector (from the LLM’s word embedding layer) and each visual token. To obtain visual saliency, we extract the average self-attention score of each token from the last layer of the vision encoder. The results are presented in Table 2 and

Fig. 1 (b). The average performance of GridPrune surpasses that of other methods across all tested retention ratios. For example, at a 33.3% retention rate, GridPrune achieves a 97.6% average score, outperforming HoloV by 3 percentage points. These results show that the principles behind GridPrune are not tailored to a specific model architecture but offer a more effective approach to token pruning.

4.3. Efficiency Analysis.

Our efficiency evaluation demonstrates the practical benefits of GridPrune, as shown in Table 5. On LLaVA-1.5-7B, our method achieves a 2.14x speedup, reducing inference latency to 90.8 ms. On LLaVA-NeXT-7B, GridPrune delivers a 5.09x speedup, cutting latency to just 113.4 ms. These results show the substantial real-world acceleration provided by our method.

This empirical speedup is theoretically grounded in the

Table 4. Performance on the high-resolution LLaVA-NeXT-7B model.

Methods	Avg.(%)	SQA ^{IMG}	GQA	MMB ^{EN}	POPE	TextVQA	MME	Seed ^I	VQA ^{V2}	VizWiz	MMVet
<i>Upper Bound, 2880 Tokens (100%)</i>											
LLaVA-NeXT-7B	100%	70.2	64.3	67.9	86.5	61.3	1842	70.2	80.1	55.2	40.0
<i>Retain 640 Tokens (22.2%)</i>											
FastV (ECCV24)	94.60%	67.4	58.9	63.1	79.5	58.1	1807	61.9	77.0	53.9	39.5
PyramidDrop (CVPR25)	95.21%	66.7	60.0	64.1	83.8	57.8	1782	65.6	79.1	53.8	36.7
SparseVLM (ICML25)	96.44%	67.6	61.2	65.9	85.3	59.7	1772	68.4	79.2	53.6	36.1
DivPrune (CVPR25)	97.07%	67.8	61.9	65.8	86.9	57.0	1773	67.6	79.3	55.7	38.0
Visionzip (CVPR25)	97.83%	68.1	61.3	66.3	86.2	59.9	1782	66.7	79.1	57.1	38.8
DART (EMNLP25)	97.10%	68.2	61.3	64.9	85.0	59.5	1793	68.1	78.3	57.0	36.9
GridPrune (Ours)	98.30%	68.2	62.8	67.1	86.9	57.4	1815	67.8	79.3	57.1	39.1
<i>Retain 320 Tokens (11.1%)</i>											
FastV (ECCV24)	77.72%	66.6	49.8	53.4	49.5	52.2	1539	56.6	61.5	51.3	20.0
PyramidDrop (CVPR25)	81.72%	66.7	50.4	55.5	60.8	49.0	1672	61.5	66.8	49.7	24.0
SparseVLM (ICML25)	92.11%	67.2	57.9	63.1	76.9	56.5	1747	65.4	74.6	54.2	32.8
DivPrune (CVPR25)	94.64%	67.7	61.1	63.9	84.7	56.2	1731	65.4	77.2	55.6	34.8
Visionzip (CVPR25)	94.60%	67.3	59.3	63.1	82.1	58.9	1698	63.4	76.2	56.2	37.8
DART (EMNLP25)	94.17%	67.5	59.5	64.2	81.0	57.6	1710	65.0	75.7	56.1	35.7
GridPrune (Ours)	96.98%	67.3	61.6	65.7	85.9	56.5	1821	66.8	78.2	56.6	38.3
<i>Retain 160 Tokens (5.6%)</i>											
DivPrune (CVPR25)	91.55%	67.1	59.3	62.9	80.0	54.1	1658	62.5	75.0	56.1	32.0
Visionzip (CVPR25)	89.22%	68.3	55.5	60.1	74.8	55.7	1630	58.3	71.4	55.5	32.6
DART (EMNLP25)	90.30%	67.8	56.8	62.0	75.3	54.9	1700	59.1	72.5	56.7	32.2
GridPrune (Ours)	94.65%	66.7	60.0	64.2	85.5	54.0	1761	65.2	76.1	55.6	37.0

reduction of total Floating-Point Operations (TFLOPs). To quantify the computational cost associated with visual inputs, we follow the convention of prior works [40, 46] and the total computational cost for the visual component across the L -layer decoder is approximated by:

$$\text{TFLOPs} \approx L \cdot (2N^2d + 4Nd^2 + 3Ndm) \quad (6)$$

where N is the number of visual tokens, d is the hidden dimension, and m is the FFN intermediate dimension. By pruning visual tokens, GridPrune directly reduces N , primarily reducing the slowdown from the quadratic $2N^2d$ term. This analysis of the visual processing cost formally justifies the empirical latency reduction observed in the full system, as the overhead of GridPrune is dwarfed by the computational savings in the decoder.

4.4. Ablation Studies

To validate GridPrune, we conduct a series of ablation studies focusing on its two core components: the Zoned Selection Mechanism and the Fused Importance Score.

Effectiveness of the Zoned Selection Mechanism. To verify this, we vary the block size to adjust the granularity of zone partitioning. A larger block size results in fewer, coarser zones, while a smaller size creates more, finer-grained zones for local selection. The results are presented in Table 6.

Table 5. Efficiency and performance trade-off analysis.

Methods	Token \downarrow	TFLOPs \downarrow	Latency (ms) \downarrow	MME Score
LLaVA-1.5-7B	576	3.82	194.2	1862
+ PyramidDrop	192	1.30	112.4	1797
+ SparseVLM	192	1.33	100.5	1787
+ Visionzip	192	1.25	91.2	1769
+ GridPrune (Ours)	192	1.25	90.8	1808
LLaVA-NeXT-7B	2880	20.83	576.9	1842
+ PyramidDrop	320	3.02	227.3	1672
+ SparseVLM	320	3.04	150.9	1747
+ Visionzip	320	2.10	115.8	1698
+ GridPrune (Ours)	320	2.10	113.4	1821

The configuration with a block size of 24, which is equivalent to a standard global Top-K selection, serves as our baseline. As the block size decreases, the overall performance improves. Decreasing the block size from 24 to 2 leads to a steady increase in the average performance from 97.6% to 98.6%. We argue that this is because finer partitioning provides a more granular budget allocation, which is crucial for preserving small but important details that might be overlooked in a coarser selection process. While the optimal block size can exhibit minor variations across different benchmarks (e.g., size 3 slightly outperforms on MME), a

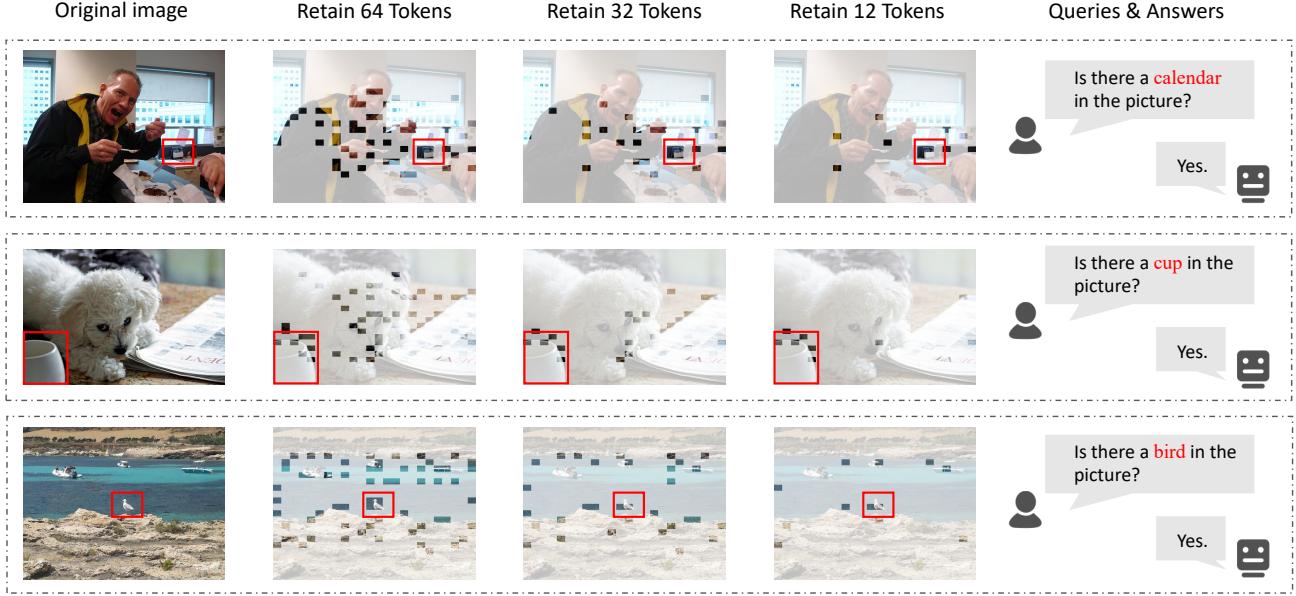


Figure 5. Visualization of GridPrune’s token selection. The token focus dynamically shifts to the calendar, cup, or bird based on the user’s query, even when the number of reserved tokens is limited to 12.

Table 6. Ablation study on the effect of selection granularity, conducted on LLaVA-1.5-7B with a fixed retention of 192 tokens and $\alpha = 0.8$. We vary the block size, where a size of 24 corresponds to global selection, and a size of 2 represents the most localized selection.

Size	Number	Avg.(%)	MME	MMB ^{EN}	POPE	SQA ^{IMG}	GQA
24	1	97.6%	1767	62.5	86.6	68.7	59.8
12	4	97.6%	1767	63.3	86.5	68.1	59.8
8	9	97.9%	1777	63.2	86.6	68.5	60.1
4	36	98.3%	1788	63.6	86.1	69.0	60.3
3	64	98.4%	1809	63.9	86.0	68.2	60.6
2	144	98.6%	1808	63.7	86.3	68.6	60.7

block size of 2 achieves the highest average score and shows robust high performance across all tested benchmarks. Considering its superior overall performance and stability, we adopt a block size of 2 as the default configuration for GridPrune.

Analysis of the Fused Importance Score. GridPrune uses a fused score that combines text-conditional relevance and intrinsic visual saliency. To analyze the contribution of each component, we ablate the fusion hyperparameter α from 0.0 to 1.0. As shown in Table 7, relying solely on text relevance or visual saliency does not produce the best results. The optimal performance is achieved when α is set to 0.8 (for this specific model and retention rate), showing that both information sources are valuable. This suggests that the model’s own assessment of important visual features acts as a beneficial regularizer, preventing the model

Table 7. Ablation study on the fusion hyperparameter α , conducted on LLaVA-1.5-7B with a fixed block size of 2 and 192 retained tokens. This parameter balances text-conditional relevance ($\alpha = 0.0$) against intrinsic visual saliency ($\alpha = 1.0$). The optimal average performance is achieved at $\alpha=0.8$, showing the synergistic effect between the two components.

α	Value	Avg.(%)	MME	MMB ^{EN}	TextVQA	SQA ^{IMG}
	0.0	97.0%	1794	64.2	54.7	68.3
	0.1	97.5%	1791	64.3	55.5	68.8
	0.2	97.4%	1768	64.2	56.0	68.7
	0.3	97.3%	1773	63.4	56.4	68.7
	0.4	97.1%	1779	63.2	56.4	68.3
	0.5	97.0%	1778	62.7	56.5	68.3
	0.6	97.4%	1787	63.4	56.6	68.2
	0.7	97.8%	1801	63.4	56.8	68.5
0.8	98.1%	1808	63.7	57.0	68.6	
0.9	97.7%	1786	63.8	56.8	68.4	
1.0	97.8%	1796	63.7	56.7	68.6	

from overfitting to the text query and thus leading to a more robust and comprehensive final token set.

4.5. Visualization

To provide an intuitive understanding of the dynamics of our method, we visualize the token selection process of GridPrune in Fig. 5. The figure shows that GridPrune is highly adaptive to user queries. When asked to identify an object on the table (top row), the model focuses its limited

token budget precisely on the calendar. Similarly, when the query shifts to the object beside the dog (middle row) or a distant bird (bottom row), GridPrune dynamically reallocates its attention to the cup and the bird. With an extremely sparse budget of only 12 tokens, the selection remains accurate, enabling the model to answer correctly.

This shows the advantage of GridPrune over instruction-irrelevant methods. Rather than generating a generic summary of the image, our method performs task-driven active filtering of visual information. This ability to form a targeted and relevant representation by preserving only the most important details is the key reason for its superior performance, allowing it to maintain better accuracy under high pruning conditions.

5. Conclusion

In this paper, we identify the limitation in existing visual token pruning methods: they primarily focus on the “what to select” problem, while paying less attention to the “where to look”, which can lead to issues like inefficient spatial allocation, positional bias and so on. So we propose GridPrune, a training-free method that divides pruning into a two-stage, “guide-globally, select-locally” process. GridPrune addresses the “where to look” problem by using text-conditional guidance to dynamically allocate a token budget across spatial zones. It then solves the “what to select” problem through a localized selection within each budgeted zone. Extensive experiments demonstrate that GridPrune achieves superior performance over state-of-the-art methods, particularly under aggressive pruning ratios. Our work suggests that incorporating the “where to look” problem into the pruning method is a promising direction for building more efficient and robust MLLMs.

References

- [1] Saeed Ranjbar Alvar, Gursimran Singh, Mohammad Akbari, and Yong Zhang. Divprune: Diversity-based visual token pruning for large multimodal models. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 9392–9401, 2025. [1](#), [3](#), [6](#)
- [2] Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenbin Ge, Yu Han, Fei Huang, et al. Qwen technical report. *arXiv preprint arXiv:2309.16609*, 2023. [1](#), [3](#)
- [3] Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Sibo Song, Kai Dang, Peng Wang, Shijie Wang, Jun Tang, et al. Owen2. 5-vl technical report. *arXiv preprint arXiv:2502.13923*, 2025. [1](#), [3](#), [6](#)
- [4] Daniel Bolya, Cheng-Yang Fu, Xiaoliang Dai, Peizhao Zhang, Christoph Feichtenhofer, and Judy Hoffman. Token merging: Your vit but faster. *arXiv preprint arXiv:2210.09461*, 2022. [1](#), [3](#), [6](#)
- [5] Zheng Cai, Maosong Cao, Haojong Chen, Kai Chen, Keyu Chen, Xin Chen, Xun Chen, Zehui Chen, Zhi Chen, Pei Chu, et al. Internlm2 technical report. *arXiv preprint arXiv:2403.17297*, 2024. [1](#), [3](#)
- [6] Junbum Cha, Wooyoung Kang, Jonghwan Mun, and Byungseok Roh. Honeybee: Locality-enhanced projector for multimodal llm. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13817–13827, 2024. [1](#), [3](#)
- [7] Fu Chaoyou, Chen Peixian, Shen Yunhang, Qin Yulei, Zhang Mengdan, Lin Xu, Yang Jinrui, Zheng Xiawu, Li Ke, Sun Xing, et al. Mme: A comprehensive evaluation benchmark for multimodal large language models. *arXiv preprint arXiv:2306.13394*, 3, 2023. [6](#)
- [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. [1](#), [3](#), [6](#)
- [9] Zhe Chen, Weiyun Wang, Hao Tian, Shenglong Ye, Zhangwei Gao, Erfei Cui, Wenwen Tong, Kongzhi Hu, Jiapeng Luo, Zheng Ma, et al. How far are we to gpt-4v? closing the gap to commercial multimodal models with open-source suites. *Science China Information Sciences*, 67(12):220101, 2024. [1](#), [3](#)
- [10] Zhe Chen, Jannan Wu, Wenhui 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. [1](#), [3](#)
- [11] Wenliang Dai, Junnan Li, Dongxu Li, Anthony Tiong, Junqi Zhao, Weisheng Wang, Boyang Li, Pascale N Fung, and Steven Hoi. Instructclip: Towards general-purpose vision-language models with instruction tuning. *Advances in neural information processing systems*, 36:49250–49267, 2023. [1](#), [3](#)
- [12] Alexey Dosovitskiy. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020. [3](#)
- [13] 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. [6](#)
- [14] Jiaxian Guo, Junnan Li, Dongxu Li, Anthony Meng Huat Tiong, Boyang Li, Dacheng Tao, and Steven Hoi. From images to textual prompts: Zero-shot visual question answering with frozen large language models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10867–10877, 2023. [1](#)
- [15] 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. [6](#)
- [16] 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. 6
- [17] Jiayi Kuang, Ying Shen, Jingyou Xie, Haohao Luo, Zhe Xu, Ronghao Li, Yinghui Li, Xianfeng Cheng, Xika Lin, and Yu Han. Natural language understanding and inference with mllm in visual question answering: A survey. *ACM Computing Surveys*, 57(8):1–36, 2025. 1
- [18] Ao Li, Yuxiang Duan, Jinghui Zhang, Congbo Ma, Yutong Xie, Gustavo Carneiro, Mohammad Yaqub, and Hu Wang. Transprune: Token transition pruning for efficient large vision-language model. *arXiv preprint arXiv:2507.20630*, 2025. 1, 3
- [19] Ao Li, Longwei Xu, Chen Ling, Jinghui Zhang, and Pengwei Wang. Emaverse: Enhancing multimodal large language models for affective computing via multitask learning. *Neurocomputing*, 650:130810, 2025. 1
- [20] Bohao Li, Rui Wang, Guangzhi Wang, Yuying Ge, Yixiao Ge, and Ying Shan. Seed-bench: Benchmarking multimodal llms with generative comprehension. *arXiv preprint arXiv:2307.16125*, 2023. 6
- [21] 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. 1, 3
- [22] 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. 1, 3
- [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. 6
- [24] 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. *arXiv preprint arXiv:2311.10122*, 2023. 1, 3
- [25] Haotian Liu, Chunyuan Li, Qingshang Wu, and Yong Jae Lee. Visual instruction tuning. *Advances in neural information processing systems*, 36:34892–34916, 2023. 1, 3
- [26] 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. 1, 3, 6
- [27] Haotian Liu, Chunyuan Li, Yuheng Li, Bo Li, Yuanhan Zhang, Sheng Shen, and Yong Jae Lee. Llavanext: Improved reasoning, ocr, and world knowledge, 2024. 1, 3, 6
- [28] Xiangrui Liu, Yan Shu, Zheng Liu, Ao Li, Yang Tian, and Bo Zhao. Video-xl-pro: Reconstructive token compression for extremely long video understanding. *arXiv preprint arXiv:2503.18478*, 2025. 1, 3
- [29] Yuan Liu, Haodong Duan, Yuanhan Zhang, Bo Li, Songyang Zhang, Wangbo Zhao, Yike Yuan, Jiaqi Wang, Conghui He, Ziwei Liu, et al. Mmbench: Is your multi-modal model an all-around player? In *European conference on computer vision*, pages 216–233. Springer, 2024. 6
- [30] Pan Lu, Swaroop Mishra, Tanglin 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. *Advances in Neural Information Processing Systems*, 35:2507–2521, 2022. 6
- [31] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PMLR, 2021. 1, 3
- [32] Yuzhang Shang, Mu Cai, Bingxin Xu, Yong Jae Lee, and Yan Yan. Llava-prumerge: Adaptive token reduction for efficient large multimodal models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 22857–22867, 2025. 6
- [33] 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. 6
- [34] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambo, Faisal Azhar, et al. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023. 1, 3
- [35] Anne M Treisman and Garry Gelade. A feature-integration theory of attention. *Cognitive psychology*, 12(1):97–136, 1980. 2
- [36] Xinlong Wang, Xiaosong Zhang, Zhengxiong Luo, Quan Sun, Yufeng Cui, Jinsheng Wang, Fan Zhang, Yueze Wang, Zhen Li, Qiying Yu, et al. Emu3: Next-token prediction is all you need. *arXiv preprint arXiv:2409.18869*, 2024. 4
- [37] Yi Wang, Xinhao Li, Ziang Yan, Yinan He, Jiashuo Yu, Xiangyu Zeng, Chenting Wang, Changlian Ma, Haian Huang, Jianfei Gao, et al. Internvideo2.5: Empowering video mllms with long and rich context modeling. *arXiv preprint arXiv:2501.12386*, 2025. 1, 3
- [38] Zichen Wen, Yifeng Gao, Weijia Li, Conghui He, and Linfeng Zhang. Token pruning in multimodal large language models: Are we solving the right problem? *arXiv preprint arXiv:2502.11501*, 2025. 1, 3, 5
- [39] 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. 1, 3, 6
- [40] 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. 1, 3, 6, 8
- [41] 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. In

- Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 19792–19802, 2025. [2](#), [3](#), [6](#)
- [42] Weihao Yu, Zhengyuan Yang, Linjie Li, Jianfeng Wang, Kevin Lin, Zicheng Liu, Xinchao Wang, and Lijuan Wang. Mm-vet: Evaluating large multimodal models for integrated capabilities. *arXiv preprint arXiv:2308.02490*, 2023. [6](#)
 - [43] Zihao Yue, Liang Zhang, and Qin Jin. Less is more: Mitigating multimodal hallucination from an eos decision perspective. *arXiv preprint arXiv:2402.14545*, 2024. [4](#)
 - [44] Xiaohua Zhai, Basil Mustafa, Alexander Kolesnikov, and Lucas Beyer. Sigmoid loss for language image pre-training. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 11975–11986, 2023. [1](#), [3](#)
 - [45] Qizhe Zhang, Mengzhen Liu, Lichen Li, Ming Lu, Yuan Zhang, Junwen Pan, Qi She, and Shanghang Zhang. Beyond attention or similarity: Maximizing conditional diversity for token pruning in mllms. *arXiv preprint arXiv:2506.10967*, 2025. [1](#), [2](#), [3](#)
 - [46] 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. [1](#), [3](#), [6](#), [8](#)
 - [47] 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. [1](#), [3](#)
 - [48] Xin Zou, Di Lu, Yizhou Wang, Yibo Yan, Yuanhuiyi Lyu, Xu Zheng, Linfeng Zhang, and Xuming Hu. Don't just chase "highlighted tokens" in mllms: Revisiting visual holistic context retention. *arXiv preprint arXiv:2510.02912*, 2025. [1](#), [2](#), [3](#), [6](#)