

Reducing Token Redundancy in LVLMs: A Systematic Review of Token Pruning Methods

Anonymous ACL submission

Abstract

Large Vision-Language Models (LVLMs) excel at visual understanding but face severe computational bottlenecks when processing high-resolution images and long videos due to massive visual token counts. Token pruning mitigates this by selectively removing less informative tokens while maintaining performance. However, existing methods vary widely in pruning location (vision encoder vs. LLM decoder), importance criteria (attention vs. similarity vs. learned scores), and application strategy, lacking systematic comparison. This survey presents the first comprehensive review of token pruning for LVLMs. We propose a taxonomy categorizing methods into vision-side, LLM-side, and hybrid paradigms, systematically analyze token selection mechanisms and pruning strategy. We further discuss evaluation protocols and identify key challenges including prompt-adaptive pruning and hardware-aware design. Our survey provides a structured foundation for this rapidly growing research area.

1 Introduction

Large Vision–Language Models (LVLMs), such as GPT-4V (Achiam et al., 2023), the LLaVA family (Liu et al., 2023a, 2024a,b), Qwen-VL (Wang et al., 2024), and the BLIP family (Li et al., 2022), have demonstrated remarkable capabilities in multimodal understanding, reasoning, and generation. By integrating powerful visual encoders with large language models, LVLMs enable a wide range of applications, including visual question answering (Singh et al., 2019), video understanding (Zhou et al., 2025), multimodal retrieval (Abootorabi et al., 2025), and agentic reasoning (Yao et al., 2025). These advances make LVLMs a promising foundation for real-world systems that require both perception and reasoning.

However, this performance comes at a significant computational cost. Modern LVLMs typically rely on Vision Transformer (ViT)-based encoders

(e.g., CLIP (Radford et al., 2021), SigLIP (Zhai et al., 2023)) that partition an image into a large number of patch tokens. The number of visual tokens scales rapidly with image resolution, video length, and multi-image inputs, leading to long multimodal input sequences. During inference, these tokens must be processed by the language model in the *prefill* stage, where full self-attention is computed over all visual and textual tokens. As a result, inference cost scales quadratically with the total token length, making visual tokens a dominant bottleneck and limiting the scalability of LVLMs in latency-sensitive and resource-constrained settings.

A key challenge underlying this inefficiency is **visual token redundancy**: many visual tokens contribute little to the final prediction for a given prompt or task, yet still incur substantial computational and memory cost. Visual token pruning has emerged as a distinct and promising direction that directly targets redundancy at the token level. It selects a compact subset of visual tokens from the original visual sequence in order to reduce computational and memory costs while preserving task-relevant information. This process can be implemented through token dropping, masking, routing, replacement, or reweighting mechanisms, and can be applied at different stages of the LVLm pipeline.

Despite the rapid growth of research in this area, there is currently no systematic and comprehensive survey that focuses specifically on visual token pruning in LVLMs. Existing surveys primarily address efficiency in large language models through model compression (Zhu et al., 2024; Cheng et al., 2025), or discuss token compression in multimodal models at a high level without detailed analysis of pruning mechanisms and design choices (Shao et al., 2025). This leaves an important gap in understanding the design space, trade-offs, and practical implications of visual token pruning. In this survey, we aim to fill this gap by providing the first structured and in-depth review of visual token pruning

methods for LVLMs. We organize existing approaches into a principled taxonomy consisting of *vision-side token pruning*, *LLM-side visual token pruning*, and *hybrid vision–LLM pruning*, reflecting where pruning is applied and what signals it relies on. This novel classification highlights fundamental differences in pruning granularity, adaptivity, and computational impact, and provides a unified framework for comparing methods across architectures and tasks. We further analyze token importance estimation strategies, compare pruning with related efficiency techniques, review evaluation protocols, and discuss open challenges and future directions.

The remainder of this article is organized as follows. Section 2 introduces background on LVLMs and visual token pruning. Section 3 reviews pruning methods according to our taxonomy. Section 4 compares pruning with related efficiency techniques and discusses evaluation practices. Section 5 presents widely used benchmarks and experiment protocols. Finally, Sections 6 and 7 discuss future directions and limitations, respectively.

2 Background

In this section, we introduce background knowledge on Large Vision–Language Models and the corresponding problem formulation for Visual Token Pruning in LVLMs.

2.1 Large Vision-Language Models

Given an input image I and a text prompt T , LVLM inference proceeds in three stages. First, a vision encoder E_v maps the image into a sequence of visual tokens, where N equals the patch number:

$$V = E_v(I) = \{v_1, v_2, \dots, v_N\}, \quad v_i \in \mathbb{R}^d.$$

Second, a multimodal projector f , which project the vision tokens to align with the language space: $H_v = f(V) \in \mathbb{R}^{N \times d}$. Finally, a language model g_θ accepts the concatenation of visual tokens H_v and text tokens $H_t = \{t_1, \dots, t_M\}$ and autoregressively generates an output sequence $Y = \{y_1, \dots, y_L\}$:

$$p(Y | I, T) = \prod_{\ell=1}^L p(y_\ell | y_{<\ell}, H_v, H_t).$$

From a computational perspective, inference in LVLMs can be divided into two phases. During the *prefill* stage, the language model processes the entire input sequence, comprising both visual tokens

from the vision encoder and textual prompt tokens, in a single forward pass to compute their hidden representations, after which the first output token is generated. This stage incurs high computational of $\mathcal{O}(N_{N+M}^2)$ and activation memory cost because self-attention scales quadratically with the total number of input tokens. In the subsequent *autoregressive decoding* stage, new tokens are generated one at a time using cached key–value (KV) states; at decoding step ℓ , the forward pass only computes attention for the newly generated token against all previously cached tokens, resulting in linear per-step complexity $\mathcal{O}(N_{N+M} + \ell)$. Consequently, the prefill stage is substantially more expensive than decoding when the input contains many visual tokens. As a result, reducing the number of visual tokens *before or during prefill* yields significantly larger efficiency gains than pruning strategies that operate only during the decoding stage.

2.2 Visual Token Pruning

Visual token pruning refers to the process of selecting a compact subset of visual tokens from the original visual sequence in order to reduce computational and memory costs while preserving minimal performance loss through reserving task-relevant information.

Formally, given an image encoded as a sequence of N visual tokens $H_v = \{h_1, h_2, \dots, h_N\}$, $h_i \in \mathbb{R}^d$, pruning constructs a reduced sequence $\tilde{H}_v = \{h_i | i \in \mathcal{I}\}$, $\mathcal{I} \subseteq \{1, \dots, N\}$ with a *token budget* enforced either by a fixed number $|\mathcal{I}| = K$, $K \ll N$, or by a keep ratio $r \in (0, 1]$, $|\mathcal{I}| = \lceil rN \rceil$.

It is important to distinguish visual token pruning from related but different techniques: **token merging** combines similar tokens instead of removing them; **KV-cache pruning** reduces memory and attention cost during decoding by discarding cached key–value states without affecting prefill (Zhang et al., 2023, 2025b); **pooling or condensation** compresses groups of semantically related tokens into compact representations (Han et al., 2025); and **structural weight pruning** removes model parameters (e.g., attention heads or feed-forward blocks) to reduce FLOPs (Liang et al., 2025a). While complementary, these methods operate at different levels of the model and target different efficiency bottlenecks.

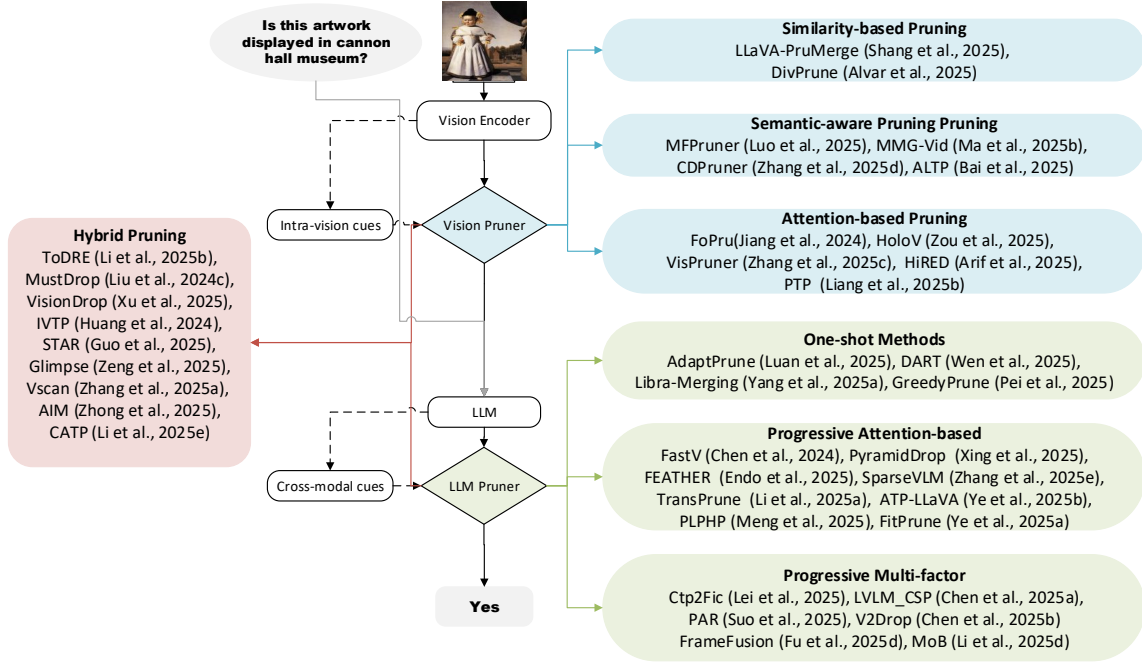


Figure 1: Comparison of three pruning paradigms. Pipelines flow top-to-bottom. Dashed arrows indicate what signals the pruner uses. Green blocks contain the LLM-side pruning methods; blue blocks contain the vision-side pruning methods; red blocks summarize the hybrid pruning methods.

3 LVLMS Token Pruning

In this section, we discuss token pruning strategy based on *where* and *how* the token number is reduced, divided the methods into **vision-side**, **LLM-side** and **hybrid** token pruning.

3.1 Vision-side Token Pruning

Vision-side token pruning operate entirely before cross-modal interaction, typically within or right after the ViT-based visual encoders, and aim to reduce the number of visual tokens injected into the LLM without relying on language guidance. In this work, we define the vision-side pruning as any pruning operation that reduces or transforms prior to the prefill stage.

3.1.1 Similarity-based Pruning

Similarity-based pruning mainly focuses on reducing token redundancy by identifying a minimal subset of tokens with high diversity, such that the selected tokens collectively preserve information close to that of the original token set. LLaVA-PruMerge combines [cls] attention sparsity with key-similarity clustering to prune and merge redundant visual tokens (Shang et al., 2025). DivPrune can also be seen as a similarity-based method with a diversity-driven criterion that maximizes the min-

imum pairwise distance among retained visual tokens (Alvar et al., 2025).

3.1.2 Semantic-aware Pruning

Beyond similarity-based redundancy reduction, another line of work focuses on visual-encoder-side pruning guided by alternative importance criteria. These methods explicitly incorporate task semantics or cross-modal signals to inform token selection, rather than relying solely on visual similarity.

Some approaches leverage text-image relevance as a complementary pruning signal. CDPPruner formulates visual token pruning as a diversity maximization problem under textual guidance, using Determinantal Point Processes (DPP) to jointly balance visual token similarity and instruction relevance—measured by cosine similarity between text embeddings and image tokens—within a unified probabilistic framework (Zhang et al., 2025d). MFPruner fuses [CLS] attention, token similarity, and instruction relevance and make pruning decision via voting mechanism (Luo et al., 2025). ALTP targets grounded conversational generation by introducing locality-aware pruning, preserving object-centric tokens via region partitioning and density-adaptive token allocation (Bai et al., 2025). In video task, temporal becomes another semantic dimension. MMG-Vid exploits temporal coher-

ence via three-stage pruning: it segments videos by frame similarity, dynamically allocates budgets to maximize marginal gain (reducing allocation for redundant segments), and selects tokens with a strategy which prioritizes tokens novel to the selection history while salient in the current frame (Ma et al., 2025b).

3.1.3 Attention-based Pruning

Attention statistics from the visual encoder have also been widely adopted as importance indicators. FoPru estimates token importance by averaging attention maps across heads and selecting the row or column with higher variance to identify salient tokens, as these tokens may have more influence towards others (Jiang et al., 2024). PTP extends attention-based pruning with a pyramid-style importance modeling strategy by first computing region-level saliency using [CLS] attention to allocate token budgets across spatial regions, then performs token-level selection within each region based on [CLS]-to-patch attention, and finally refines the selected tokens using instruction-aware relevance (Liang et al., 2025b). In contrast, HoloV addresses the over-localization issue inherent in attention-first pruning by combining [CLS] attention for saliency estimation with intra-crop token variance to encourage semantic diversity. It further allocates token budgets adaptively across spatial crops to preserve holistic visual context (Zou et al., 2025). VisPruner performs two-stage pruning after the projector by first keeping important tokens based on [CLS] attention, and then iteratively removes redundant tokens by token similarity (cosine similarity) to retain a diverse complement under a fixed budget (Zhang et al., 2025d). HiRED similarly adopts a two-stage strategy, before the projector, first using early-layer [CLS] attention to estimate content distribution across image partitions, and then selecting top-K informative tokens based on [CLS] attention of layer 22 (Arif et al., 2025).

By aggregating complementary signals, such approaches mitigate common failure modes of single-criterion pruning, such as attention collapse or semantic bias toward dominant visual regions.

3.2 LLM-side Token Pruning

Since vision-side token pruning reduces the number of visual tokens *before the LLM generation process*, it directly lowers the computational cost at the prefill stage. In contrast, LLM-side pruning methods do not modify the visual encoder; instead,

they leverage language-model behavior during the prefill or decoding stages, typically by exploiting prompt-visual interactions to identify and drop less informative visual tokens.

3.2.1 One-shot Pruning

This line of work in LLM-side token pruning make the pruning decision for one time.

Several methods exploit token similarity to identify redundancy. Specifically, they identify and prune similar tokens while retaining a diverse subset to preserve information for generation. For example, AdaptPrune reframes token pruning as an adaptive non-maximum suppression process that jointly considers attention, spatial distance, and token similarity, pruning visual tokens for one time in the early layer of LLM (Luan et al., 2025). In contrast, DART reframes the problem from token importance to token duplication. Rather than relying on attention scores, DART removes redundant visual tokens in one early layer by measuring embedding similarity to a small set of pivot tokens, ensuring diverse token retention (Wen et al., 2025). Libra-Merging is a hybrid pruning and merge method, also resolving the importance-redundancy dilemma by selecting representative tokens from spatial intervals and performing similarity-aware grouped merging with compensation tokens (Yang et al., 2025a). GreedyPrune formulates token selection as a combinatorial optimization problem that jointly optimizes semantic saliency and visual diversity, solving it via an efficient greedy strategy and finish the token pruning after the first layer for only one time (Pei et al., 2025).

3.2.2 Progressive Pruning

Progressive pruning reduces computation by exploiting the fact that token utility evolves during inference, allowing different pruning operations to be applied progressively rather than in a single step. **Attention-based.** In shallow layers (layers 1-2), attention in relatively balances across all token types, with image tokens actively aggregating visual information through self-attention, while in deeper layers, attention becomes extremely imbalances – system prompts receive $472\times$ higher attention efficiency than image tokens, capturing 85% of total attention (Chen et al., 2024).

Based on this observation, FastV removes visual tokens with persistently low cross-modal attention after layer 2, as it noticed in deeper LVLM layers the visual tokens receives much less attention

than system prompts in deep layers, significantly reducing computation with minimal accuracy loss (Chen et al., 2024). PyramidDrop extended the idea by exploiting the layer-wise growth of visual redundancy by retaining all tokens in shallow layers and progressively dropping less informative tokens in deeper layers (Xing et al., 2025). To eliminate the positional bias within LLM layers, FEATHER uses RoPE-free attention from the last text token as the primary criteria, ensembles it with uniform sampling criteria in early layers for coverage, and applies aggressive pruning with the refined attention criteria in later layers (Endo et al., 2025). While FastV successfully identifies inefficient visual attention in deep layers, its pruning strategy remains fundamentally text-agnostic, SparseVLM addresses this limitation through text-aware guidance, arguing that pruning should be question-adaptive (Zhang et al., 2025e).

Both FastV and SparseVLM posit that visual tokens become progressively less important in deeper layers, justifying increasingly aggressive pruning. However, this assumption of monotonic attention decline lacks empirical validation across diverse LLM architectures. PLPHP discovers the Vision Token Re-attention phenomenon where visual attention resurges in deep layers and further introduces per-layer, per-head retention rates, enabling more adaptive pruning strategy (Meng et al., 2025). FitPrune is also a progressive pruning method, which take the self-attention score and cross-attention score as criterion and prunes the visual tokens during every layer during decoding time (Ye et al., 2025a).

Besides using attention score itself as an importance metric, there are also some variant based on attention. TransPrune evaluates token importance via Token Transition Variation through layers of token’s self attention, capturing both magnitude and directional changes across layers, and enables training-free, multi-stage pruning guided by representation dynamics (Li et al., 2025a). Related method such as ATP-LLaVA Leverages dual criteria of redundancy scores (averaged self-modal and cross-modal attention) and spatial scores (2D RoPE-enhanced uniform sampling) with learnable thresholds to achieve instance-wise and layer-wise adaptive pruning (Ye et al., 2025b).

Multi-factor. Progressive pruning can also be lead by different pruning criteria through the stages with different retention goals. Multi-factor methods jointly consider multiple pruning goals such as

diversity and importance, and use different computation algorithms to balance between criteria. V2Drop adopts a variation-aware criterion, measuring token importance by representation changes between consecutive transformer layers. Tokens exhibiting low variation are considered redundant and progressively removed by three times during inference, enabling pruning without relying on attention statistics (Chen et al., 2025b). Ctp2Fic combines shallow-layer text-guided pruning with deep-layer semantic clustering in a coarse-to-fine manner (Lei et al., 2025). MoB formulates visual token pruning as a bi-objective covering problem and theoretically characterizes the trade-off between visual preservation and prompt alignment under fixed budgets, which first choose visual tokens nearest to prompts tokens and then choose the farthest tokens to the chosen tokens to increase the diversity (Li et al., 2025d). LVM_CSP adopts a three-stage progressive pruning framework across LLM decoder layers with first the clustering stage using Seg-First criteria or [cls] attention scores, scattering stage re-activates all tokens for fine-grained reasoning, and finally retains top tokens ranked by [SEG] token’s attention scores (Chen et al., 2025a).

In video tasks, FrameFusion performs similarity-based token merging at shallow layers, where visual redundancy across adjacent frames is most pronounced, and permanently removes merged tokens. At deeper layers, where semantic importance becomes more discriminative, FrameFusion applies importance-based pruning using cumulative self-attention scores to further satisfy computational budgets (Fu et al., 2025d).

Besides the training free methods, PAR prunes visual tokens across all 32 LLM layers using a meta-router and simultaneously skips redundant layers in the last 16 layers based on learnable layer controller embeddings’ importance scores, optimized via self-supervised DPO by minimizing KL divergence between pruned and original outputs. (Suo et al., 2025).

3.3 Hybrid Vision-LLM Pruning

Hybrid vision-LLM pruning methods determine visual token importance by *jointly leveraging visual structure and language semantics*, explicitly coupling visual token reduction with linguistic relevance signals. The pruning strategy may happen on both vision-side and LLM-side with different guidance. By integrating cues from both modalities, these methods aim to improve task alignment

and robustness under aggressive pruning budgets, at the cost of increased complexity or additional supervision. These methods adopt a multi-stage pruning strategy across both vision-side and LLM-side, leveraging information from both modalities.

MustDrop is a typical hybrid pruning method with three pruning stages with different goals: it removes spatially redundant tokens and retain key tokens across vision encoder, utilizes text-to-image attention score to guide pruning of text-irrelevant tokens during prefilling stage, and removes output-irrelevant tokens during decoding stage (Liu et al., 2024c). STAR performs two-stage pruning—early conservative visual self-attention pruning and later aggressive cross-modal attention pruning—balancing feature preservation with task relevance (Guo et al., 2025). ToDRE also adopts a two-stage strategy that first retains a diverse subset of visual tokens via greedy k -center selection before LLM decoder and then removes all remaining visual tokens once cross-modal attention becomes negligible in deeper layers inside LLM layers (Li et al., 2025b). IVTP is another two-stage pruning methods which stabilizes token importance using group-wise attention rollout on vision-side, and then pruning the text-irrelevant visual tokens again inside the LLM (Huang et al., 2024). CATP targets multimodal in-context learning by pruning image tokens based on cross-example and query-conditioned relevance. It adopts a two-stage pruning strategy: the first stage happen after projector and before decoder, maximizing text alignment and diversity; and the second stage progressively prunes tokens based on variation of token attention and query relevance (Li et al., 2025e).

Some methods jointly perform token pruning and token merging. We summarize them as a complementary category of hybrid approaches, since they involve pruning and operate across both the vision- and LLM-side. AIM combines token merging with token pruning, employing cosine similarity between embeddings as merging criteria before the LLM, then uses PageRank scores computed from self-attention weights as pruning criteria for progressive layer-wise token reduction within the LLM (Zhong et al., 2025). VisionDrop addresses cross-modal misalignment in LLMs by performing training-free, visual-only token pruning across multiple stages (in both visual encoder and LLM decoder), combining dominant token selection with contextual merging to preserve fine-grained visual information (Xu et al., 2025). VScan demonstrates

that the effectiveness of visual token pruning critically depends on when pruning is performed along the LVLM pipeline, rather than solely on how tokens are scored. By revealing distinct roles of shallow vision layers and middle LLM layers, VScan reframes token reduction as a stage-aware optimization problem (Zhang et al., 2025a).

4 Comparison with Related Paradigms

Visual token pruning in LVLMs is closely related to and sometimes cooperates with several efficiency-oriented paradigms that also aim to reduce computation or memory cost. However, these paradigms differ fundamentally from token pruning in terms of *what* is reduced, *where* the reduction is applied, and *when* the reduction takes effect. We summarize the most relevant paradigms below and clarify their conceptual differences from visual token pruning.

Token Merging and Token Compression. Token merging and compression methods reduce the effective token count by *aggregating or replacing* tokens rather than explicitly discarding them. A representative example is ToMe, which introduces a training-free token merging mechanism that progressively fuses similar tokens via fast bipartite matching and proportional attention (Bolya et al., 2023). iLLaVA extends token merging to large vision–language models by performing attention-guided token merging in both the visual encoder and the language model (Hu et al., 2024). Related approaches such as VisionZip (Yang et al., 2025b), FiCoCo (Han et al., 2025), LaCo (Liu et al., 2025), LightVLM (Hu et al., 2025), and Fwd2Bot (Bulat et al., 2025) further compress dense visual tokens into compact semantic representations through clustering, pooling, or learned summarization. Unlike token pruning, these methods preserve information through aggregation rather than removal, typically offering higher stability but more limited aggressiveness under extreme compression budgets.

KV-cache Management and Adaptive Attention. KV-cache based methods reduce *decoding-time* memory footprint and attention computation by selectively retaining or computing key–value states, without modifying the input token sequence. Twilight generalizes sparse attention by replacing fixed-budget top- k selection with adaptive top- p retention (Lin et al., 2025). A-VL proposes a plug-and-play adaptive attention mechanism that separately manages visual and textual KV caches, dynamically retaining only critical visual states and a small local

536 text window (Zhang et al., 2025b). These meth- 586
537 ods primarily accelerate the decoding stage and 587
538 are orthogonal to token pruning, which focuses on 588
539 reducing prefilling cost via token selection. 589

540 **Structural Pruning.** Structural pruning reduces 590
541 computation by removing or compressing *model* 591
542 *structures* such as weights, modules, attention 592
543 heads, or transformer layers. Representative exam- 593
544 ples include EfficientLLaVA (Liang et al., 2025a) 594
545 and UKMP (Wu et al., 2025), which perform 595
546 parameter-level pruning with learned importance 596
547 metrics, and Short-LVLM (Ma et al., 2025a), which 597
548 removes redundant transformer layers in a training- 598
549 free manner. While structural pruning provides 599
550 consistent speedups across both prefilling and de- 600
551 coding stages, it lacks the input adaptivity and 601
552 instance-level flexibility offered by token pruning. 602

553 Overall, these paradigms are complementary 603
554 rather than competing. Token pruning focuses on 604
555 dynamic, input-adaptive token selection before or 605
556 during language interaction, while token merging, 606
557 KV-cache management, patch merging, structural 607
558 pruning, and semantic compensation address effi- 608
559 ciency from orthogonal dimensions. In practice, 609
560 these techniques can be combined to further im- 610
561 prove the efficiency of large vision–language mod- 611
562 els. 612

563 5 Benchmarks and Experimental 613 564 Protocols 614

565 **Benchmarks.** We introduce several benchmarks 615
566 to which most of the selected papers adapt. The 616
567 detailed information of each benchmarks are pre- 617
568 sented in Table 1. Notably, these tasks differ in to- 618
569 ken redundancy, reasoning depth, and dependency 619
570 on fine-grained features. For instance, text-oriented 620
571 VQA and detailed visual reasoning typically de- 621
572 mand high token fidelity, whereas global caption- 622
573 ing may tolerate higher sparsity. Consequently, this 623
574 diverse benchmark suite is essential for evaluating 624
575 the generalization of pruning methods across vary- 625
576 ing sensitivities to information loss, ensuring that 626
577 efficiency gains do not come at the cost of failing 627
578 specific task distributions. 628

579 **Experimental Protocols.** To ensure fair compari- 629
580 son and reproducibility, rigorous protocols are re- 630
581 quired beyond simple metric reporting. Standard 631
582 evaluations in the papers typically adhere to three 632
583 key aspects: 633

584 • **Backbone Consistency:** Comparisons are 634
585 strictly conducted on identical LVLM architec- 635

586 tures (e.g., LLaVA (Liu et al., 2023b), BLIP-2 587
587 (Li et al., 2023a)) to isolate the efficacy of the 588
588 pruning algorithm from the underlying model 589
589 capability. 590

591 • **Pruning Paradigms:** Distinctions are explicitly 591
592 drawn between training-free (zero-shot) methods 592
593 and fine-tuning-based approaches, as they oper- 593
594 ate under fundamentally different computational 594
594 budgets. 595

595 • **Inference Constraints:** Critical hyperparame- 595
596 ters, including input resolution, batch size, and 596
597 maximum generation length, are fixed to stan- 597
598 dardize the evaluation. This control is vital when 598
599 comparing static pruning ratios against dynamic 599
600 token budgets. 600

601 To assess the efficiency of token pruning meth- 601
602 ods, evaluations typically measure the trade-off 602
603 between model sparsity and downstream task per- 603
604 formance. Researchers primarily examine the ca- 604
605 pabilities of LVLMs across varying pruning ratios 605
606 (i.e., the number or proportion of preserved tokens). 606
607 The fundamental performance metrics include ab- 607
608 solute Accuracy and the Performance Retention 608
609 Rate, which quantifies the percentage of perfor- 609
610 mance maintained relative to the original, unpruned 610
611 baseline. 611

612 In terms of computational efficiency, quantitative 612
613 comparisons rely on a multi-dimensional suite of 613
614 metrics covering theoretical complexity, temporal 614
615 latency, and spatial memory footprint: 615

616 • **FLOPs** (Floating Point Operations) serve as a 616
617 standard proxy for evaluating the theoretical com- 617
618 putational complexity of the model. 618

619 • **CUDA Latency** measures the actual wall-clock 619
620 time required for kernel launches and data trans- 620
621 fers. This metric is frequently decomposed into 621
622 two phases to isolate stage-specific benefits: (1) 622
623 Prefill Time, the duration required to process 623
624 all input tokens and compute the embedding for 624
625 the first generated token; and (2) Decode Time, 625
626 the latency incurred during the subsequent auto- 626
627 regressive generation process. 627

628 • **KV Cache and GPU Memory** are utilized to 628
629 evaluate the spatial efficiency gains. These met- 629
630 rics assess the reduction in video random access 630
631 memory (VRAM) usage and Key-Value cache 631
632 overhead, highlighting the method’s ability to 632
633 alleviate hardware bottlenecks during inference. 633

634 These metrics provide a comprehensive view of the 634
635 efficiency gains achieved by token pruning from 635

Task Category	Description	Benchmarks
VQA	Answer questions based on images	VQAv2(Goyal et al., 2016), GQA(Hudson and Manning, 2019), VizWiz(Gurari et al., 2018), ScienceQA-IMG(Lu et al., 2022), HallBench(Guan et al., 2024), POPE(Li et al., 2023b), MME(Fu et al., 2025a), MMBench(Liu et al., 2024d), MMBench-CN(Liu et al., 2024d), MM-Vet(Yu et al., 2023)
Text-oriented VQA	VQA requiring text recognition in images	TextVQA(Singh et al., 2019), ChartQA(Masry et al., 2022), AI2D(Kembhavi et al., 2016), OCRBench(Liu et al., 2023c)
Video Understanding	Answer questions based on videos	MLVU(Zhou et al., 2025), MVBench(Li et al., 2024), LongVideoBench(Wu et al., 2024), Video-MME(Fu et al., 2025b)
Visual Captioning	Generate descriptive captions for images	COCO Caption(Lin et al., 2015), Flickr30k(Young et al., 2014)
Document & OCR	Extract and understand textual information	DocVQA(Mathew et al., 2020), IIIT5K(Mishra et al., 2012), ICDAR(Pfitzmann et al., 2022)

Table 1: Evaluation Benchmarks for Different Vision-Language Tasks

both temporal (operation time) and spatial (memory footprint) perspectives.

6 Future Directions and Conclusion

Previous token pruning methods has achieved good balancing efficiency and performance. Here we outline some practical limitations and open challenges in token pruning methods.

Vision-side pruning reduces visual redundancy using intra-visual signals such as similarity and attention offering efficient acceleration since the token number is reduced before the LLM decoding. However, its task-agnostic nature and lack of language awareness might fundamentally limit its effectiveness in more complex multimodal reasoning tasks. LLM-side pruning introduces query-relevance token selection by leveraging cross-modal signals during decoding. While effective for suppressing visually redundant but task-irrelevant information, its reliance on noisy attention cues and late-stage pruning fundamentally limits efficiency gains and exposes it to cross-modal misalignment risks (Zhang et al., 2025c). Hybrid pruning methods combines both vision- and LLM-side methods by multi-stage or multi-factor pruning strategy, offering a principled compromise between efficiency and task relevance.

Task-agnostic Pruning. Recent studies suggest that aggressive visual token pruning may introduce failure modes beyond accuracy degradation, including spatial misalignment and hallucination amplification. GAP demonstrates that preserving positional consistency during pruning is crucial for maintaining grounding performance, indicating that future pruning methods should explicitly account for geometric and spatial constraints (Chien et al., 2025). Similarly, VASparse reveals that naive

sparsification strategies can exacerbate visual hallucinations, motivating hallucination-aware pruning objectives that go beyond attention-based importance estimation (Zhuang et al., 2025). These findings highlight the need for behavior-aware pruning frameworks that optimize not only efficiency but also grounding fidelity and hallucination robustness. Also, while LVLMs are widely used in a wide range of tasks, today’s LVLMs pruning methods mainly limit on VQA tasks. How to extend the pruning strategy and cooperate with a wider range of practical phenomena is a trending research area. Researchers may explore the different ways for token importance evaluation and criteria for token retaining for different practical tasks.

Hallucination or Pruning. Along with the hallucination perspective, it’s also crucial to create pruning task specific benchmarks. For now the evaluation for pruning topic still builds on the existed tasks, researchers performs original and pruned models on the same benchmarks and compare the accuracy drop as well as efficiency gains. An ideal pruning strategy should maintain the informational and sufficient tokens to answer the questions. However, there still lacks of discussion on the hidden mechanisms of whether the model can answer the questions with remaining tokens or the insignificant accuracy drop comes from model hallucination, especially when the pruning rates are high. Li et.al (Li et al., 2025d) investigated the prompt-visual coupling effect across different benchmarks and found that budget allocation should differ by the coupling rate of tasks. This also highlighted the necessity of investigating whether the pruning strategy is aligned with the tasks. In the future, introducing pruning specific tasks and benchmarks is important.

7 Limitations

In this survey we reviewed visual token pruning method for LVLMs since 2024. Despite our effort to provide a comprehensive overview of token pruning in LVLMs, this survey has some limitations. First of all, as efficiency LVLMs nowadays are having great progress and there are more methods other than token pruning as well as methods combining token pruning with other compression technics. Secondly, most methods reviewed in this survey are developed for ViT-based visual encoders and late-fusion architectures. Emerging unified or early-fusion LVLMs may exhibit different redundancy patterns, which could limit the direct applicability of existing pruning principles and require rethinking pruning criteria and evaluation methodologies.

References

Mohammad Mahdi Abootorabi, Amirhosein Zobeiri, Mahdi Deghani, Mohammadali Mohammadkhani, Bardia Mohammadi, Omid Ghahroodi, Mahdiah Soleymani Baghshah, and Ehsaneddin Asgari. 2025. [Ask in any modality: A comprehensive survey on multimodal retrieval-augmented generation](#). *Preprint*, arXiv:2502.08826.

Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altschmidt, Sam Altman, Shyamal Anadkat, and 1 others. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*.

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

Kazi Hasan Ibn Arif, JinYi Yoon, Dimitrios S Nikolopoulos, Hans Vandierendonck, Deepu John, and Bo Ji. 2025. Hired: Attention-guided token dropping for efficient inference of high-resolution vision-language models. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 39, pages 1773–1781.

Bizhe Bai, Jianjian Cao, Yadan Luo, and Tao Chen. 2025. [Local information matters: Inference acceleration for grounded conversation generation models through adaptive local-aware token pruning](#). *Preprint*, arXiv:2503.23959.

Daniel Bolya, Cheng-Yang Fu, Xiaoliang Dai, Peizhao Zhang, Christoph Feichtenhofer, and Judy Hoffman. 2023. [Token merging: Your vit but faster](#). *Preprint*, arXiv:2210.09461.

Adrian Bulat, Yassine Ouali, and Georgios Tzimiropoulos. 2025. [Fwd2bot: Lvlm visual token compression with double forward bottleneck](#). *Preprint*, arXiv:2503.21757.

Hanning Chen, Yang Ni, Wenjun Huang, Hyunwoo Oh, Yezi Liu, Tamoghno Das, and Mohsen Imani. 2025a. [Lvlm_csp: Accelerating large vision language models via clustering, scattering, and pruning for reasoning segmentation](#). *Preprint*, arXiv:2504.10854.

Junjie Chen, Xuyang Liu, Zichen Wen, Yiyu Wang, Siteng Huang, and Honggang Chen. 2025b. [Variation-aware vision token dropping for faster large vision-language models](#). *Preprint*, arXiv:2509.01552.

Liang Chen, Haozhe Zhao, Tianyu Liu, Shuai Bai, Junyang Lin, Chang Zhou, and Baobao Chang. 2024. 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.

Jian Cheng, Haidong Kang, Yuxin Shao, Nan Li, Pengjun Chen, Rui Wang, Saiqin Long, Xiaochun Yang, and Lianbo Ma. 2025. [Survey on efficient large language models: Principles, algorithms, applications, and open issues](#). *IEEE Transactions on Neural Networks and Learning Systems*, pages 1–21.

Tzu-Chun Chien, Chieh-Kai Lin, Shiang-Feng Tsai, Ruei-Chi Lai, Hung-Jen Chen, and Min Sun. 2025. [Grounding-aware token pruning: Recovering from drastic performance drops in visual grounding caused by pruning](#). *Preprint*, arXiv:2506.21873.

Mark Endo, Xiaohan Wang, and Serena Yeung-Levy. 2025. Feather the throttle: Revisiting visual token pruning for vision-language model acceleration. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 22826–22835.

Chaoyou Fu, Peixian Chen, Yunhang Shen, Yulei Qin, Mengdan Zhang, Xu Lin, Jinrui Yang, Xiawu Zheng, Ke Li, Xing Sun, and 1 others. 2025a. Mme: A comprehensive evaluation benchmark for multimodal large language models. In *The Thirty-ninth Annual Conference on Neural Information Processing Systems Datasets and Benchmarks Track*.

Chaoyou Fu, Yuhan Dai, Yongdong Luo, Lei Li, Shuhuai Ren, Renrui Zhang, Zihan Wang, Chenyu Zhou, Yunhang Shen, Mengdan Zhang, and 1 others. 2025b. Video-mme: The first-ever comprehensive evaluation benchmark of multi-modal llms in video analysis. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 24108–24118.

Mingyu Fu, Wei Suo, Ji Ma, Lin Yuanbo Wu, Peng Wang, and Yanning Zhang. 2025c. Mitigating information loss under high pruning rates for efficient large vision language models. In *Proceedings of the 33rd ACM International Conference on Multimedia*, pages 4156–4165.

818	Tianyu Fu, Tengxuan Liu, Qinghao Han, Guohao Dai,	Lei Jiang, Zixun Zhang, Yuting Zeng, Chunzhao Xie,	875
819	Shengen Yan, Huazhong Yang, Xuefei Ning, and	Tongxuan Liu, Zhen Li, Lechao Cheng, and Xiao-	876
820	Yu Wang. 2025d. Framefusion: Combining simi-	hua Xu. 2025. DCP: Dual-cue pruning for efficient	877
821	ilarity and importance for video token reduction on	large vision-language models . In <i>Proceedings of the</i>	878
822	large vision language models. In <i>Proceedings of the</i>	<i>2025 Conference on Empirical Methods in Natural</i>	879
823	<i>IEEE/CVF International Conference on Computer</i>	<i>Language Processing</i> , pages 21202–21215, Suzhou,	880
824	<i>Vision</i> , pages 22654–22663.	China. Association for Computational Linguistics.	881
825	Yash Goyal, Tejas Khot, Douglas Summers-Stay, Dhruv	Aniruddha Kembhavi, Mike Salvato, Eric Kolve, Min-	882
826	Batra, and Devi Parikh. 2016. Making the V in VQA	joon Seo, Hannaneh Hajishirzi, and Ali Farhadi.	883
827	matter: Elevating the role of image understanding in	2016. A diagram is worth a dozen images. In <i>Com-</i>	884
828	visual question answering . <i>CoRR</i> , abs/1612.00837.	<i>puter Vision – ECCV 2016</i> , pages 235–251, Cham.	885
829	Tianrui Guan, Fuxiao Liu, Xiyang Wu, Ruiqi Xian,	Springer International Publishing.	886
830	Zongxia Li, Xiaoyu Liu, Xijun Wang, Lichang Chen,	Yulong Lei, Zishuo Wang, Jinglin Xu, and Yuxin Peng.	887
831	Furong Huang, Yaser Yacoob, Dinesh Manocha, and	2025. Ctp2fic: From coarse-grained token pruning	888
832	Tianyi Zhou. 2024. Hallusionbench: An advanced	to fine-grained token clustering for lvm inference	889
833	diagnostic suite for entangled language hallucination	acceleration (chinamm 2025). <i>Available at SSRN</i>	890
834	and visual illusion in large vision-language models.	5545751.	891
835	In <i>CVPR</i> .	Ao Li, Yuxiang Duan, Jinghui Zhang, Congbo Ma,	892
836	Yichen Guo, Hanze Li, Zonghao Zhang, Jinhao You,	Yutong Xie, Gustavo Carneiro, Mohammad Yaqub,	893
837	Kai Tang, and Xiande Huang. 2025. Star: Stage-	and Hu Wang. 2025a. Transprune: Token transition	894
838	wise attention-guided token reduction for efficient	pruning for efficient large vision-language model .	895
839	large vision-language models inference . <i>Preprint</i> ,	<i>Preprint</i> , arXiv:2507.20630.	896
840	arXiv:2505.12359.	Duo Li, Zuhao Yang, Xiaoqin Zhang, Ling Shao, and	897
841	Danna Gurari, Qing Li, Abigale J Stangl, Anhong Guo,	Shijian Lu. 2025b. Todre: Effective visual to-	898
842	Chi Lin, Kristen Grauman, Jiebo Luo, and Jeffrey P	ken pruning via token diversity and task relevance .	899
843	Bigham. 2018. Vizwiz grand challenge: Answering	<i>Preprint</i> , arXiv:2505.18757.	900
844	visual questions from blind people. In <i>Proceedings of</i>	Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi.	901
845	<i>the IEEE conference on computer vision and pattern</i>	2023a. Blip-2: Bootstrapping language-image pre-	902
846	<i>recognition</i> , pages 3608–3617.	training with frozen image encoders and large lan-	903
847	Yuhang Han, Xuyang Liu, Zihan Zhang, Pengxiang	guage models . <i>Preprint</i> , arXiv:2301.12597.	904
848	Ding, Junjie Chen, Donglin Wang, Honggang Chen,	Junnan Li, Dongxu Li, Caiming Xiong, and Steven	905
849	Qingsen Yan, and Siteng Huang. 2025. Filter, cor-	Hoi. 2022. Blip: Bootstrapping language-image pre-	906
850	relate, compress: Training-free token reduction for	training for unified vision-language understanding	907
851	mllm acceleration . <i>Preprint</i> , arXiv:2411.17686.	and generation. In <i>International conference on ma-</i>	908
852	Lianyu Hu, Fanhua Shang, Wei Feng, and Liang Wan.	<i>chine learning</i> , pages 12888–12900. PMLR.	909
853	2025. Lightvlm: Accelerating large multimodal mod-	Kaiyuan Li, Xiaoyue Chen, Chen Gao, Yong Li, and	910
854	els with pyramid token merging and kv cache com-	Xinlei Chen. 2025c. Balanced token pruning: Ac-	911
855	pression . <i>Preprint</i> , arXiv:2509.00419.	celerating vision language models beyond local opti-	912
856	Lianyu Hu, Fanhua Shang, Liang Wan, and Wei Feng.	mization . <i>Preprint</i> , arXiv:2505.22038.	913
857	2024. illava: An image is worth fewer than 1/3 in-	Kunchang Li, Yali Wang, Yinan He, Yizhuo Li,	914
858	put tokens in large multimodal models . <i>Preprint</i> ,	Yi Wang, Yi Liu, Zun Wang, Jilan Xu, Guo Chen,	915
859	arXiv:2412.06263.	Ping Luo, and 1 others. 2024. Mvbench: A com-	916
860	Kai Huang, Hao Zou, Ye Xi, BoChen Wang, Zhen Xie,	prehensive multi-modal video understanding bench-	917
861	and Liang Yu. 2024. Ivtp: Instruction-guided vi-	mark . In <i>Proceedings of the IEEE/CVF Conference</i>	918
862	sual token pruning for large vision-language models .	<i>on Computer Vision and Pattern Recognition</i> , pages	919
863	In <i>Computer Vision – ECCV 2024: 18th European</i>	22195–22206.	920
864	<i>Conference, Milan, Italy, September 29–October 4,</i>	Yangfu Li, Hongjian Zhan, Tianyi Chen, Qi Liu, and	921
865	<i>2024, Proceedings, Part XVII</i> , page 214–230, Berlin,	Yue Lu. 2025d. Why 1 + 1 < 1 in visual token prun-	922
866	Heidelberg. Springer-Verlag.	ing: Beyond naive integration via multi-objective	923
867	Drew A. Hudson and Christopher D. Manning. 2019.	balanced covering . <i>Preprint</i> , arXiv:2505.10118.	924
868	GQA: a new dataset for compositional ques-	Yanshu Li, Jianjiang Yang, Zhennan Shen, Ligong Han,	925
869	tion answering over real-world images . <i>CoRR</i> ,	Haoyan Xu, and Ruixiang Tang. 2025e. Catp: Con-	926
870	abs/1902.09506.	textually adaptive token pruning for efficient and en-	927
871	Lei Jiang, Weizhe Huang, Tongxuan Liu, Yuting Zeng,	hanced multimodal in-context learning . <i>Preprint</i> ,	928
872	Jing Li, Lechao Cheng, and Xiaohua Xu. 2024. Fo-	arXiv:2508.07871.	929
873	pru: Focal pruning for efficient large vision-language		
874	models . <i>CoRR</i> .		

930	Yifan Li, Yifan Du, Kun Zhou, Jinpeng Wang, Xin Zhao, and Ji-Rong Wen. 2023b. Evaluating object hallucination in large vision-language models . In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pages 292–305, Singapore. Association for Computational Linguistics.	983
931		984
932		985
933		986
934		987
935		988
936	Yinan Liang, Ziwei Wang, Xiuwei Xu, Jie Zhou, and Jiwen Lu. 2025a. Efficientllava: Generalizable auto-pruning for large vision-language models . In <i>Proceedings of the Computer Vision and Pattern Recognition Conference</i> , pages 9445–9454.	989
937		990
938		991
939		992
940		993
941	Yuxuan Liang, Xu Li, Xiaolei Chen, Yi Zheng, Haotian Chen, Bin Li, and Xiangyang Xue. 2025b. Pyramid token pruning for high-resolution large vision-language models via region, token, and instruction-guided importance . <i>Preprint</i> , arXiv:2509.15704.	994
942		995
943		996
944		997
945		998
946	Chaofan Lin, Jiaming Tang, Shuo Yang, Hanshuo Wang, Tian Tang, Boyu Tian, Ion Stoica, Song Han, and Mingyu Gao. 2025. Twilight: Adaptive attention sparsity with hierarchical top-p pruning . <i>Preprint</i> , arXiv:2502.02770.	999
947		1000
948		1001
949		1002
950		1003
951	Tsung-Yi Lin, Michael Maire, Serge Belongie, Lubomir Bourdev, Ross Girshick, James Hays, Pietro Perona, Deva Ramanan, C. Lawrence Zitnick, and Piotr Dollár. 2015. Microsoft coco: Common objects in context . <i>Preprint</i> , arXiv:1405.0312.	1004
952		1005
953		1006
954		1007
955		1008
956	Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. 2024a. Improved baselines with visual instruction tuning . In <i>Proceedings of the IEEE/CVF conference on computer vision and pattern recognition</i> , pages 26296–26306.	1009
957		1010
958		1011
959		1012
960		1013
961	Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. 2023a. Visual instruction tuning . <i>Advances in neural information processing systems</i> , 36:34892–34916.	1014
962		1015
963		1016
964		1017
965	Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. 2023b. Visual instruction tuning . <i>Advances in neural information processing systems</i> , 36:34892–34916.	1018
966		1019
967		1020
968		1021
969	Juntao Liu, Liqiang Niu, Wenchao Chen, Jie Zhou, and Fandong Meng. 2025. Laco: Efficient layer-wise compression of visual tokens for multimodal large language models . <i>Preprint</i> , arXiv:2507.02279.	1022
970		1023
971		1024
972		1025
973	Shilong Liu, Hao Cheng, Haotian Liu, Hao Zhang, Feng Li, Tianhe Ren, Xueyan Zou, Jianwei Yang, Hang Su, Jun Zhu, and 1 others. 2024b. Llava-plus: Learning to use tools for creating multimodal agents . In <i>European conference on computer vision</i> , pages 126–142. Springer.	1026
974		1027
975		1028
976		1029
977		1030
978		1031
979	Ting Liu, Liangtao Shi, Richang Hong, Yue Hu, Quan-jun Yin, and Linfeng Zhang. 2024c. Multi-stage vision token dropping: Towards efficient multimodal large language model . <i>Preprint</i> , arXiv:2411.10803.	1032
980		1033
981		1034
982		1035
	Yuan Liu, Haodong Duan, Yuanhan Zhang, Bo Li, Songyang Zhang, Wangbo Zhao, Yike Yuan, Jiaqi Wang, Conghui He, Ziwei Liu, Kai Chen, and Dahua Lin. 2024d. Mmbench: Is your multi-modal model an all-around player? In <i>Computer Vision – ECCV 2024: 18th European Conference, Milan, Italy, September 29–October 4, 2024, Proceedings, Part VI</i> , page 216–233, Berlin, Heidelberg. Springer-Verlag.	1036
		1037
		1038
	Yuliang Liu, Zhang Li, Mingxin Huang, Biao Yang, Wenwen Yu, Chunyuan Li, Xucheng Yin, Cheng lin Liu, Lianwen Jin, and Xiang Bai. 2023c. Ocrbench: on the hidden mystery of ocr in large multimodal models . <i>Science China Information Sciences</i> , 67.	1039
		1040
	Pan Lu, Swaroop Mishra, Tanglin Xia, Liang Qiu, Kai-Wei Chang, Song-Chun Zhu, Oyvind Tafjord, Peter Clark, and Ashwin Kalyan. 2022. Learn to explain: Multimodal reasoning via thought chains for science question answering . <i>Advances in Neural Information Processing Systems</i> , 35:2507–2521.	1041
		1042
	Bozhi Luan, Wengang Zhou, Hao Feng, Zhe Wang, Xiaosong Li, and Houqiang Li. 2025. Multi-cue adaptive visual token pruning for large vision-language models . <i>Preprint</i> , arXiv:2503.08019.	1043
		1044
	Deng Luo, Dongyang Zhang, Qiuhaio Xie, Cencen Liu, Qiang Dong, and Xiurui Xie. 2025. Rethinking attention cues: Multi-factor guided token pruning for efficient vision-language understanding . <i>Available at SSRN 5615684</i> .	1045
		1046
	Ji Ma, Wei Suo, Peng Wang, and Yanning Zhang. 2025a. Short-llvm: Compressing and accelerating large vision-language models by pruning redundant layers . In <i>Proceedings of the 33rd ACM International Conference on Multimedia</i> , pages 3575–3584.	1047
		1048
	Junpeng Ma, Qizhe Zhang, Ming Lu, Zhibin Wang, Qiang Zhou, Jun Song, and Shanghang Zhang. 2025b. Mmg-vid: Maximizing marginal gains at segment-level and token-level for efficient video llms . <i>Preprint</i> , arXiv:2508.21044.	1049
		1050
	Ahmed Masry, Xuan Long Do, Jia Qing Tan, Shafiq Joty, and Enamul Hoque. 2022. Chartqa: A benchmark for question answering about charts with visual and logical reasoning . In <i>Findings of the association for computational linguistics: ACL 2022</i> , pages 2263–2279.	1051
		1052
	Minesh Mathew, Dimosthenis Karatzas, R Manmatha, and CV Jawahar. 2020. Docvqa: A dataset for vqa on document images . corr abs/2007.00398 (2020). <i>arXiv preprint arXiv:2007.00398</i> .	1053
		1054
	Yu Meng, Kaiyuan Li, Chenran Huang, Chen Gao, Xinlei Chen, Yong Li, and Xiaoping Zhang. 2025. Plphp: Per-layer per-head vision token pruning for efficient large vision-language models . <i>Preprint</i> , arXiv:2502.14504.	1055
		1056
	A. Mishra, K. Alahari, and C. V. Jawahar. 2012. Scene text recognition using higher order language priors . In <i>BMVC</i> .	1057
		1058

1151	denotations: New similarity metrics for semantic inference over event descriptions. <i>Transactions of the Association for Computational Linguistics</i> , 2:67–78.	<i>IEEE/CVF International Conference on Computer Vision</i> , pages 20180–20192.	1207
1152			1208
1153			
1154	Weihao Yu, Zhengyuan Yang, Linjie Li, Jianfeng Wang, Kevin Lin, Zicheng Liu, Xinchao Wang, and Lijuan Wang. 2023. Mm-vet: Evaluating large multimodal models for integrated capabilities. <i>arXiv preprint arXiv:2308.02490</i> .	Junjie Zhou, Yan Shu, Bo Zhao, Boya Wu, Zhengyang Liang, Shitao Xiao, Minghao Qin, Xi Yang, Yongping Xiong, Bo Zhang, and 1 others. 2025. Mlvu: Benchmarking multi-task long video understanding. In <i>Proceedings of the Computer Vision and Pattern Recognition Conference</i> , pages 13691–13701.	1209
1155			1210
1156			1211
1157			1212
1158			1213
1159	Quan-Sheng Zeng, Yunheng Li, Qilong Wang, Peng-Tao Jiang, Zuxuan Wu, Ming-Ming Cheng, and Qibin Hou. 2025. A glimpse to compress: Dynamic visual token pruning for large vision-language models. <i>Preprint</i> , arXiv:2508.01548.	Xunyu Zhu, Jian Li, Yong Liu, Can Ma, and Weiping Wang. 2024. A survey on model compression for large language models. <i>Transactions of the Association for Computational Linguistics</i> , 12:1556–1577.	1215
1160			1216
1161			1217
1162			1218
1163			
1164	Xiaohua Zhai, Basil Mustafa, Alexander Kolesnikov, and Lucas Beyer. 2023. Sigmoid loss for language image pre-training. <i>Preprint</i> , arXiv:2303.15343.	Xianwei Zhuang, Zhihong Zhu, Yuxin Xie, Liming Liang, and Yuexian Zou. 2025. Vaspase: Towards efficient visual hallucination mitigation via visual-aware token sparsification. In <i>Proceedings of the Computer Vision and Pattern Recognition Conference</i> , pages 4189–4199.	1219
1165			1220
1166			1221
1167			1222
1168	Ce Zhang, Kaixin Ma, Tianqing Fang, Wenhao Yu, Hongming Zhang, Zhisong Zhang, Yaqi Xie, Katia Sycara, Haitao Mi, and Dong Yu. 2025a. Vscan: Rethinking visual token reduction for efficient large vision-language models. <i>Preprint</i> , arXiv:2505.22654.	Xin Zou, Di Lu, Yizhou Wang, Yibo Yan, Yuanhuiyi Lyu, Xu Zheng, Linfeng Zhang, and Xuming Hu. 2025. Don't just chase "highlighted tokens" in mllms: Revisiting visual holistic context retention. <i>Preprint</i> , arXiv:2510.02912.	1223
1169			1224
1170			1225
1171			1226
1172			1227
1173			1228
1174			1229
1175	Junyang Zhang, Mu Yuan, Ruiguang Zhong, Puhao Luo, Huiyou Zhan, Ningkan Zhang, Chengchen Hu, and Xiang-Yang Li. 2025b. A-vl: Adaptive attention for large vision-language models. In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , volume 39, pages 22461–22469.	A Appendix.	1230
1176			
1177			
1178			
1179	Qizhe Zhang, Aosong Cheng, Ming Lu, Renrui Zhang, Zhiyong Zhuo, Jiajun Cao, Shaobo Guo, Qi She, and Shanghang Zhang. 2025c. Beyond text-visual attention: Exploiting visual cues for effective token pruning in vlms. In <i>Proceedings of the IEEE/CVF International Conference on Computer Vision</i> , pages 20857–20867.	A.1 Tables.	1231
1180			
1181			
1182			
1183			
1184			
1185			
1186	Qizhe Zhang, Mengzhen Liu, Lichen Li, Ming Lu, Yuan Zhang, Junwen Pan, Qi She, and Shanghang Zhang. 2025d. Beyond attention or similarity: Maximizing conditional diversity for token pruning in mllms. <i>Preprint</i> , arXiv:2506.10967.	To organize the rapidly growing body of work on LVLM token pruning, we provide a systematic literature review structured as follows in this appendix. Table 2 provides a detailed comparison of vision token pruning methods, including their key innovations, training requirements, and pruning granularity. Table 3 complements this by surveying related compression techniques—such as token merging, knowledge distillation, and quantization—that address similar efficiency goals but through different mechanisms.	1232
1187			1233
1188			1234
1189			1235
1190			1236
1191	Yuan Zhang, Chun-Kai Fan, Junpeng Ma, Wenzhao Zheng, Tao Huang, Kuan Cheng, Denis Gudovskiy, Tomoyuki Okuno, Yohei Nakata, Kurt Keutzer, and Shanghang Zhang. 2025e. Sparsevlm: Visual token sparsification for efficient vision-language model inference. <i>Preprint</i> , arXiv:2410.04417.		1237
1192			1238
1193			1239
1194			1240
1195			1241
1196			1242
1197	Zhenyu Zhang, Ying Sheng, Tianyi Zhou, Tianlong Chen, Lianmin Zheng, Ruisi Cai, Zhao Song, Yundong Tian, Christopher Ré, Clark Barrett, and 1 others. 2023. H2o: Heavy-hitter oracle for efficient generative inference of large language models. <i>Advances in Neural Information Processing Systems</i> , 36:34661–34710.		
1198			
1199			
1200			
1201			
1202			
1203			
1204	Yiwu Zhong, Zhuoming Liu, Yin Li, and Liwei Wang. 2025. Aim: Adaptive inference of multi-modal llms via token merging and pruning. In <i>Proceedings of the</i>		
1205			
1206			

Table 2: Survey of Vision Token Pruning Methods for LVLMs

Side	Paper	Category	Train?	Granularity	Key Innovation
Vision-side	LLaVA-PruMerge (Shang et al., 2025)	Similarity-based	×	Vision encoder output, before projector	[cls] attention sparsity + key-similarity clustering
	DivPrune (Alvar et al., 2025)	Similarity-based	×	After projector, before LLM	Diversity-driven selection by maximizing minimum pairwise distance
	CDPruner (Zhang et al., 2025d)	Semantic-aware	×	After projector, before LLM	Token similarity + textual relevance (text-guided diversity)
	MFPPruner (Luo et al., 2025)	Semantic-aware	×	Vision encoder output, before projector	Fuse [CLS] attention, token similarity, and instruction relevance via voting
	MMG-Vid (Ma et al., 2025b)	Semantic-aware	×	After projector, before LLM	Two-level marginal information gain maximization; adaptive temporal budget allocation
	ALTP (Bai et al., 2025)	Semantic-aware	×	After projector, before LLM	Region partitioning + density-adaptive token budget allocation
	FoPru (Jiang et al., 2024)	Attention-based	×	Vision encoder output, before projector	Variance-aware [CLS]-attention visual token pruning with global (rank) and local (row) selection strategies
	PTP (Liang et al., 2025b)	Attention-based	×	After projector, before LLM	region-level + token-level [CLS] attention
	HoloV (Zou et al., 2025)	Attention-based	×	Vision encoder output, before projector	Crop-wise variance + [CLS] attention importance
	VisPruner (Zhang et al., 2025d)	Attention-based	×	Vision encoder output, before projector	[CLS] attention importance + removes redundancy with token similarity
HiRED (Arif et al., 2025)	Attention-based	×	Vision encoder output, before projector	early [CLS]-attn for budget allocation; final [CLS]-attn for importance	
LLM-side	AdaptPrune (Luan et al., 2025)	One-shot	×	During decoding	Cross-attention scores + spatial distance + spatially diverse tokens
	DART (Wen et al., 2025)	One-shot	×	During decoding	Pivot tokens + duplicate dropping to remove redundancy
	Libra-Merging (Yang et al., 2025a)	One-shot	×	During decoding	pruning-merging trade-off
	GreedyPrune (Pei et al., 2025)	One-shot	×	between layer 1 and 2	Greedy subset selection for critical token retention
	FastV (Chen et al., 2024)	Attention-based	×	After layer 2	Drop low-attention visual tokens based on average attention scores
	PyramidDrop (Xing et al., 2025)	Attention-based	×	Stage-wise (keep shallow, drop deep)	Progressive pruning ratio scheduled by stages+ lightweight attention to the last instruction token
	FEATHER (Endo et al., 2025)	Attention-based	×	early-mid layer+ deeper layer	RoPE-free attention scores
	SparseVLM (Zhang et al., 2025e)	Attention-based	×	All layers (cross-attn guided)	Use text-visual cross-attention to induce sparsity
	PLPHP (Meng et al., 2025)	Attention-based	×	Per-layer, per-head	Layer/head-wise adaptive retention
	FitPrune (Ye et al., 2025a)	Attention-based	×	All layers	Minimizing divergence of self- and cross-attention distributions before/after pruning
	TransPrune (Li et al., 2025a)	Attention-based	×	Shallow-mid layers	Token Transition Variation of Attention for layer-wise selection
	ATP-LLaVA (Ye et al., 2025b)	Attention-based	✓	Can be inserted between any two LLM layers	Learnable adaptive token pruning
	V2Drop (Chen et al., 2025b)	Multi-factor	×	predefined layers	Variation-based token importance
Ctp2Fic (Lei et al., 2025)	Multi-factor	×	Layer 7 + Layer 22	Text-guided pruning at shallow layer + LSH-based clustering at deep layer	

Continued on next page

Table 2 – continued from previous page

Side	Paper	Category	Train?	Granularity	Key Innovation
	PAR (Suo et al., 2025)	Multi-factor	✓	Adaptive	A meta-router trained via DPO
	FrameFusion (Fu et al., 2025d)	Multi-factor	×	Deep decoder layers	Shallow layers merging based on similarity + deep layers pruning based on importance
	LVLN_CSP (Chen et al., 2025a)	Multi-factor	×	Deep decoder layers	Coarse-to-fine pipeline: Clustering–Scattering–Pruning on the deep layers based on [SEG] attention
	MoB (Li et al., 2025c)	Multi-factor	×	A fixed early decoder layer	Balance prompt alignment + visual preservation
Hybrid	CATP (Li et al., 2025e)	Multi-stage pruning	×	pre-decoder + shallow LLM decoder layers	text–image semantic alignment + diversity (stage 1) +inter-layer attention changes (stage 2)
	MustDrop (Liu et al., 2024c)	Multi-stage pruning	×	Encode + prefill + decode + KV	Lifecycle-aware token dropping across multiple inference stages
	VisionDrop (Xu et al., 2025)	Multi-stage pruning + merging	×	End of vision encoder + multiple decoder stages	Visual-only token importance scoring + stage-wise dominant token selection
	IVTP (Huang et al., 2024)	Two-stage pruning	×	Vision encoder + early decoder layers	Group-wise Token Pruning + instruction-relevant visual tokens retaining
	AIM(Zhong et al., 2025)	merging+ pruning	×	Vision side merging + LLM decoder mid-layers pruning	similarity-based token merging before LLM + importance pruning inside LLM
	Glimpse(Zeng et al., 2025)	Attention-based	✓	single pruning after layer K	Learnable glimpse token + predictor
	STAR (Guo et al., 2025)	Multi-stage pruning	×	After vision encoder (self-attention pruning) + intermediate LLM layer	Early visual self-attention + later cross-modal guidance
	ToDRE (Li et al., 2025b)	Multi-stage pruning	×	After encoder + late LLM decoder layer	Diverse visual token selection with relevance awareness
	VScan (Zhang et al., 2025a)	merging + pruning	×	vision encoder merging + mid LLM decoder layers pruning	Global–local scanning; mid-layer pruning

Year	Paper	Category	Train?	Granularity	Key Innovation
2022	EfficientVLM(Wang et al., 2023)	Model compression	✓	Layer-/module-level	Distill-then-prune with modal-adaptive pruning that learns task-specific importance of vision vs. language encoders
2023	ToMe(Bolya et al., 2023)	Structural compression	×	Encoder merging	Training-free token merging for ViTs
2024	ECOFLAP(Sung et al., 2024)	Weight Pruning	×	Layer-wise (weight)	Coarse-to-fine layer-wise pruning with global importance estimated via zeroth-order gradients
2024	iLLaVA(Hu et al., 2024)	Visual token merging	×	Encoder + LLM	Attention-guided one-step token merging with information recycling across both image encoder and LLM
2024	VisionZip(Yang et al., 2025b)	Structural compression	✓	Encoder dominant-token + merge	Attention concentration + similarity merging
2025	ACCM(Fu et al., 2025c)	Visual token pruning	✓	Encoder + Decoder Pruning	Train a caption model to mitigate information loss / retain key information
2025	A-VL(Zhang et al., 2025b)	KV-cache optimization	×	Decoder-side	Modality-aware adaptive attention: hierarchical vision KV selection + multi-scale text cache
2025	EfficientLLaVA(Liang et al., 2025a)	Weight Pruning	✓	Layer-wise	Structural risk minimization: search layer-wise pruning ratios using few proxy samples and evolve the search space by optimizing the vision-language projector
2025	FiCoCo-V / -L (Han et al., 2025)	Token compression	×	Filter-Correlate-Compress	Three-phase design; V/L variants (training-free)
2025	Fwd2Bot(Bulat et al., 2025)	Visual token compression	✓	Decoder-side	Condense visual tokens into task-agnostic summary tokens, jointly optimized with autoregressive + contrastive losses
2025	short-LVLM (Ma et al., 2025a)	Layer pruning	×	Decoder-side layer-level	Token-aware layer localization (Token Importance Scores) + Subspace-Compensated Pruning
2025	LightVLM (Hu et al., 2025)	Token merging + KV cache compression	×	Encoder-side merging + Decoder-side KV cache	Pyramid token merging across LLM layers to hierarchically condense visual tokens + attention-guided KV cache compression
2025	UKMP(Wu et al., 2025)	Weight pruning	✓	MHA and FFN across vision + language	Balance modality- and block-wise importance + distillation
2025	DCP(Jiang et al., 2025)	Structured pruning	×	Dependency-aware channel pruning	Dependency-consistent pruning for efficiency
2025	Twilight (Lin et al., 2025)	Dynamic top- p	×	Decoder attention sparsity	Top- p instead of top- k for adaptive retention

Table 3: Other Efficiency Methods mentioned in this survey