
AdaptInfer: Adaptive Token Pruning for Vision–Language Model Inference with Dynamical Text Guidance

Weichen Zhang

Global Innovation Exchange
Tsinghua University
Beijing, China
weic_zhang@mails.tsinghua.edu.cn

Zhui Zhu

Department of Automation
Tsinghua University
Beijing, China
z-zhu22@mails.tsinghua.edu.cn

Kebin Liu *

Global Innovation Exchange
Tsinghua University
Beijing, China
kebinliu2021@tsinghua.edu.cn

Yunhao Liu

Global Innovation Exchange
Tsinghua University
Beijing, China
yunhao@tsinghua.edu.cn

Abstract

Vision–language models (VLMs) have achieved impressive performance on multi-modal reasoning tasks such as visual question answering, image captioning and so on, but their inference cost remains a significant challenge due to the large number of vision tokens processed during the prefill stage. Existing pruning methods often rely on directly using the attention patterns or static text prompt guidance, failing to exploit the dynamic internal signals generated during inference. To address these issues, we propose AdaptInfer, a plug-and-play framework for adaptive vision token pruning in VLMs. First, we introduce a fine-grained, dynamic text-guided pruning mechanism that reuses layer-wise text-to-text attention maps to construct soft priors over text-token importance, allowing more informed scoring of vision tokens at each stage. Second, we perform an offline analysis of cross-modal attention shifts and identify consistent inflection locations in inference, which inspire us to propose a more principled and efficient pruning schedule. Our method is lightweight and plug-and-play, also generalizable across multi-modal tasks. Experimental results have verified the effectiveness of the proposed method. For example, it reduces CUDA latency by 61.3% while maintaining an average accuracy of 93.1% on vanilla LLaVA-1.5-7B. Under the same token budget, AdaptInfer surpasses SOTA in accuracy.

1 Introduction

In recent years, building on the success of LLMs [Bommasani et al., 2021, Touvron et al., 2023, Brown et al., 2020], vision–language models (VLMs) have emerged to tackle multimodal reasoning by combining visual encoders [Liu et al., 2022, Dosovitskiy et al., 2021], with LLMs’ text decoders [Du et al., 2022]. This integration enables impressive performance on tasks such as captioning [Lin et al., 2014], image retrieval [Faghri et al., 2018], and visual question answering (VQA) [Antol et al., 2015], but it also introduces a new computational challenge: the sheer number of vision tokens.

*Corresponding author.

During the inference process of VLMs, the number of vision tokens is often much larger than that of textual tokens, sometimes by an order of magnitude or more. For instance, an image of size 672×672 processed by a visual encoder with a patch size of 14×14 typically results in 2304 vision tokens [Radford et al., 2021, Zhang et al., 2025], whereas the corresponding text prompt may contain fewer than 100 tokens [Hudson and Manning, 2019, Fu et al., 2023]. Also, much previous research suggests that the vision tokens are more redundant and semantically repetitive [Zhang et al., 2024, Tong et al., 2025].

As a result, explorations in VLM acceleration primarily focus on the efficient pruning, compression or sparsification of vision tokens. This paradigm aims to reduce computational overhead by retaining only the most valuable vision tokens. Among them, some works introduce sparsity strategies within the visual encoders to generate less but useful enough vision tokens [Li et al., 2024, Chen et al., 2025b], while others further prune tokens based on either self-attention or cross-attention patterns during the prefill stage [Lin et al., 2025, Xing et al., 2024, Chen et al., 2025a, Li et al., 2025]. Nevertheless, not all the attention logits should be involved in the vision token ranking for the dispersion of the full attention patterns [Zhang et al., 2024]. Granting voting rights only to the most salient tokens sharpens guidance, enabling more aggressive yet accurate vision-token pruning.

To address this, SparseVLM [Zhang et al., 2025] introduces the concept of text prompt-guided pruning, selecting the most salient text tokens offline before the prefill pass. While this approach acknowledges the importance of textual cues, it does not fully address the underlying challenge: **the dynamic nature of text token importance during inference**. In practice, the informativeness of text tokens evolves across layers as the model progressively refines its internal representations [Tenney et al., 2019, Clark et al., 2019]. Our observations in Figure 1a also indicate that the most prominent text tokens vary significantly across layers, making any static selection inherently suboptimal.

Therefore, to truly harness the benefits of text-guided sparsification, it is crucial to develop a pruning strategy matched with dynamic fluidity of information, reflecting the effective cross-modal interaction throughout the inference process. In this work, we propose to reconstruct the dynamic importance ranking of text tokens at each layer by utilizing the text-to-text (t2t) attention maps. These attention maps provide a natural, layer-specific prior distribution for text-token importance, which we then use to reweight text-to-vision (t2v) attention scores for vision token pruning. Importantly, since the t2t attention maps can be directly extracted from the model’s attention computations, our method does not introduce additional computational overhead.

Moreover, current methods determine pruning hyperparameters (e.g. the pruning locations) primarily through either empirical rule-of-thumb or extensive hyperparameter optimization experiments [Xing et al., 2024, Chen et al., 2025a, Zhang et al., 2024]. However, we argue that **relying on manual tuning or grid search not only imposes substantial offline computational overhead, but also leads to task- or dataset-specific heuristics**. In this work, we take the first step in providing a principled pruning schedule. We provide our insights based on systematically analyzing the distributional characteristics of attention shifts to vision tokens during VLM inference. Specifically, we identify consistent attention inflection points at layer 1, 10, and 20 in LLava-1.5-7B [Liu et al., 2023b], suggesting that aggressive pruning immediately after these layers is a more effective and computationally efficient strategy on LLava.

The solutions we propose effectively address the limitations of existing works in the field of vision token sparsification for VLM acceleration. Our main contribution involves:

- We propose AdaptInfer, an adaptive vision token sparsification framework in which VLM dynamically determines text token guidance during inference. AdaptInfer is a plug-and-play solution.
- We introduce a novel observation of the distributional characteristics of attention shifts, and gain insights in a more effective and reasonable pruning schedule.
- We implement and evaluate our proposed solution, AdaptInfer, across multiple benchmarks and different vision token budget settings. Within the same token budget, our AdaptInfer outperforms state-of-the-art (SOTA) methods on the metric of the accuracy.

2 Related Work

In this section, we will briefly introduce the previous works that are correlated with ours.

2.1 Vision-Language Models

Early multi-modal systems paired convolutional vision backbones with recurrent language decoders [Karpathy and Fei-Fei, 2017, Vinyals et al., 2015]. Modern VLMs instead follow the Transformer paradigm [Vaswani et al., 2017] by representing an image as a sequence of *visual tokens* that interact with textual tokens in a shared self-attention space [Liu et al., 2023b, Chen et al., 2024a,b]. BLIP-2 [Li et al., 2023] and MiniGPT-4 [Zhu et al., 2023] introduce lightweight linear adapters that project features from a frozen CLIP encoder into the hidden space of a large language model, enabling efficient training. These explorations bridge the gap between frozen encoders and LLMs. The LLaVA family [Liu et al., 2023a,b, 2024a] refines this recipe with stronger instruction tuning, while a few other efforts such as Flamingo [Alayrac et al., 2022], CogVLM [Wang et al., 2024] and GPT-4V [OpenAI, 2023] scale the approach to billions of parameters.

Recent explorations also contribute to resolution and multi-modal extensions. Models leveraging hierarchical perception (e.g., LLaVA-NeXT [Liu et al., 2024a]) and adaptive patching (e.g., Qwen-VL [Bai et al., 2023]) permit high-resolution inputs, while large-scale vision encoders are adopted to generate richer hidden states of the vision tokens [Chen et al., 2024a,b]. However, these gains come at the cost of dramatically more vision tokens, which is a key bottleneck addressed by our work.

2.2 Inference Acceleration of VLMs

Previous approaches mainly focus on vision token sparsification. This is because the number of vision tokens is often an order of magnitude (or more) larger than that of text tokens. In addition, visual embeddings are naturally much more sparse and repetitive than human-made texts [Marr, 2010]. In this field, there are two research directions including efficient vision encoders and vision token pruning in LLM networks.

For example, methods like LLaVA-PruMerge [Shang et al., 2025] and FlowCut [Tong et al., 2025] follow the first direction, which cuts the encoder outputs or uses a lightweight projector to reduce the number of vision tokens. Recoverable compression [Chen et al., 2025b] is also introduced to repair the information loss from token pruning within vision encoders. Solutions follow the second direction not only to drop the vision tokens [Chen et al., 2025a, Lin et al., 2025, Xing et al., 2024, Li et al., 2025, Luan et al., 2025, Ye et al., 2025] during prefill, but also merge them to compress the numbers of the vision tokens [Bolya et al., 2023] and recover them in certain inference stage [Wu et al., 2025]. SparseVLM [Zhang et al., 2025] tries to go deeper in exploration of static text prompt guidance but ignoring the evolving inherent of token information. Our approach contributes to the second paradigm.

3 Method

3.1 Observations

In this subsection, we will present a few observations and preliminary experiments that have illuminated our insights.

3.1.1 Text Token Importance

Given an input question prompt, we first tokenize it as $Prompt = [t_1, \dots, t_n]$. We then define the *importance* of a text token t_i at layer ℓ as the total attention weight it receives from all text tokens in the layer- ℓ text-to-text (t2t) attention map on average of attention heads:

$$\text{Imp}_\ell(t_i) = \frac{1}{H} \sum_{h=1}^H \sum_{j=1}^n \mathbf{A}_{\text{t2t}}^{(\ell,h)}[j, i], \quad (1)$$

where $\mathbf{A}_{\text{t2t}}^{(\ell,h)}[j, i]$ denotes the attention of head h directed from token t_j to t_i , H represents the number of the attention heads. Intuitively, tokens with high Imp_ℓ are the current key text tokens of the language stream and are therefore best suited to guide cross-modal pruning decisions at that layer.

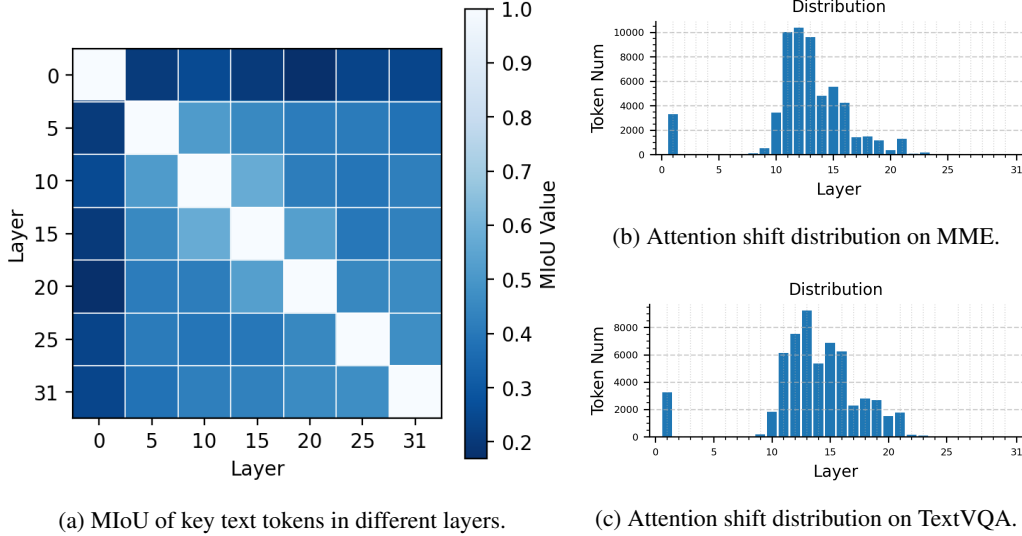


Figure 1: **Preliminary observations.** (a) The consistent low mIoU shows the key text token varies across layers. (b)(c) Layer-wise distribution of attention shifts on MME and TextVQA, which shows a highly consistent trend.

3.1.2 Dynamics of Text Token Importance

To demonstrate that the importance of text tokens evolves significantly across layers, we conduct a simple empirical study on LLaVa-1.5-7B [Liu et al., 2023b]. LLaVa-1.5-7B contains a 32-layer LLaMa [Touvron et al., 2023] as its language model. We extract t2t attention maps from 1,000 samples in the TextVQA dataset. On average, each sample contains approximately 100 text tokens. At each chosen layer, we select the top 20% text tokens as the key text tokens of that layer. We then compute the mean Intersection over Union (mIoU) between the indexes of the selected top tokens across layers, as reported in Figure 1a. The value in the i -th row and j -th column represents the mIoU of the key text token indexes of the i -th and j -th layers. Note that this figure is symmetrical, and the mIoU values along the diagonal are always 1 because they are comparisons within the same layer.

The consistently low mIoU values between different layers indicate that the set of the key text tokens changes substantially during inference. For example, the 0.169 mIoU between layer 0 and 24 refers that only 16.9% of key text tokens are overlapped while all others are different. This highlights the inherent dynamics of text token importance, where the VLM attends to different parts of the input question at different stages of reasoning. Consequently, any static text-prompt-guided pruning approach is likely to fail in capturing this evolving semantic alignment. These findings support the need for an adaptive, layer-wise text-guidance mechanism that can track and respond to the evolving attention distribution online.

3.1.3 Cross-Attention Shifts of VLM

We argue that a principled approach to setting pruning hyperparameters is not only necessary but preferable to complex, trial-and-error-based tuning. To support this claim, we performed an analysis to investigate the locations of cross-attention shifts during inference. Specifically, we calculate cumulative text-to-vision (t2v) attention scores for visual tokens at each layer using LLaVa-1.5-7B [Liu et al., 2023b]. For each sample, we first select the top 10% of vision tokens, approximately 58 tokens per image, which receive the highest total attentions in the prefill stage. Similar to the text tokens, these vision tokens are assumed to be the most critical ones for the corresponding image-based multi-modal task.

To understand when these important tokens become semantically salient, we analyze their cumulative attention trajectories across transformer layers. We apply change-point detection [Truong et al., 2020] on each curve to identify the layer where the model’s attention pattern changes significantly. The technical details of change-point detection are presented in the Appendix B. Intuitively, an attention

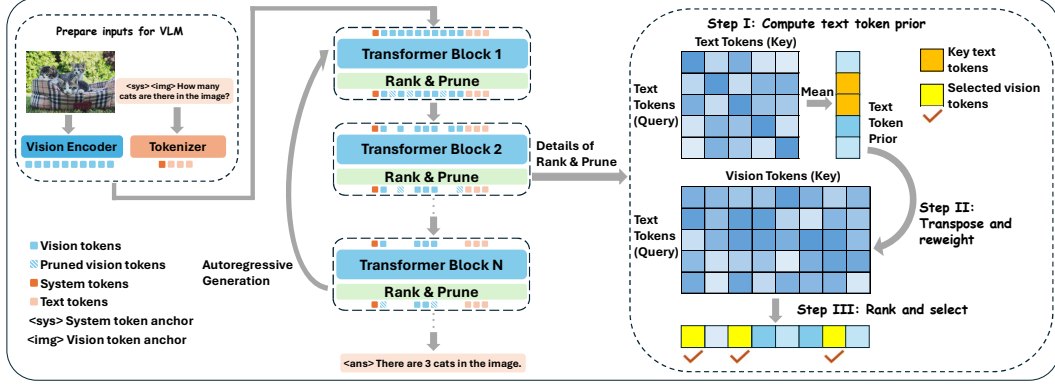


Figure 2: **The architecture of AdaptInfer.** Text token importance is computed and guides vision token selection at every pruning layer adaptively. The right panel illustrates the internal computation details of the Rank & Prune module.

shift point may indicate that either (1) the token begins to receive significantly more attention, becoming critical, or (2) the token’s informative content has already been fully extracted, becoming redundant. In both cases, these shifts mark the layer-wise transitions in how the model utilizes visual information.

Figure 1b and 1c shows the distribution of the attention shift locations aggregated from 1,000 samples each in the MME [Fu et al., 2023] and TextVQA [Singh et al., 2019] datasets. Despite data set differences, we observe a highly consistent trend: attention shifts cluster densely in layer 1 and round layers 10-20, while layers 2-8 and 20 + show a frequency of near-zero attention shifts. These findings provide a data-driven basis for pruning schedule design, where pruning locations are informed by the model’s own attention behavior, rather than by empirical intuition alone.

3.2 Adaptive Token Pruning

Based on the two key observations above, we propose AdaptInfer for VLM inference acceleration.

3.2.1 Dynamic Text Guidance Pruning

To guide visual token pruning in a more informed and adaptive way, we propose a dynamic text-guidance mechanism that leverages the model’s internal attention signals. Instead of statically selecting a fixed subset of text tokens before entering the language model [Zhang et al., 2025], we dynamically infer their relative importance during inference at each predefined pruning layer. We reuse the attention maps already computed by the model. The architecture of Our AdaptInfer is shown in Fig 2, and our dynamic text guidance mechanism can be divided into three steps.

Firstly, in each pruning layer, we extract the *text-to-text* attention matrix $\mathbf{A}_{t2t}^{(h)} \in \mathbb{R}^{T \times T}$, where T is the number of text tokens, h is the index of attention head with a total number of H . To estimate the importance of each text token, we aggregate attention scores across the query dimension:

$$\mathbf{w}^{(h)} = \sum_{i=1}^T \mathbf{A}_{t2t}^{(h)}[i, :] \in \mathbb{R}^T, \quad (2)$$

where \mathbf{w} serves as a soft prior distribution over the text tokens, averaged across all attention heads, indicating how much attention each token receives from the rest of the sequence.

Secondly, we then use this prior to reweight the *t2v* attention matrix $\mathbf{A}_{t2v}^{(h)} \in \mathbb{R}^{T \times V}$, where V is the number of visual tokens remained. $\mathbf{A}_{t2v}^{(h)}$ denotes the cross-attention matrix where text tokens are used as queries, and vision tokens serve as keys and values. The importance score of each visual token on average of all attention heads is computed as:

$$\mathbf{s} = \frac{1}{H} \sum_{h=1}^H \mathbf{w}^{(h)\top} \cdot \mathbf{A}_{t2v}^{(h)} \in \mathbb{R}^V. \quad (3)$$

Here, \mathbf{s}_j reflects the aggregated and weighted attention from all text tokens to visual token j .

Finally, based on these scores, we rank all visual tokens and retain the top- k for the current layer:

$$\mathcal{I}_k = \text{TopK}(\mathbf{s}, k). \quad (4)$$

Importantly, all text tokens participate in visual token scoring, but contribute in proportion to their dynamically inferred importance. Moreover, because both \mathbf{A}_{t2t} and \mathbf{A}_{t2v} are natively computed in standard forward passes, our method introduces little additional computational overhead. Note that this solution follows a training-free paradigm and can be seamlessly integrated as a plugin into existing VLMs.

3.2.2 Analysis of Computational Complexity

Assuming $n = T + V$ denotes the current sequence length, d denotes the VLM’s hidden state dimension, and m denotes the hidden size of projection layer in the FFN network. For each transformer layer in the prefill stage, the FLOPs can be estimated by

$$FLOPs^{\text{prefill}} = 4nd^2 + 2n^2d + 3ndm. \quad (5)$$

For each pruning layer, the additional FLOPs of are computed below:

$$FLOPs^{\text{prune}} = T^2 + 2TV. \quad (6)$$

Since both attention matrices are already computed during the forward pass, this additional cost is minimal relative to the main transformer computations. Then, during the decode stage, the FLOPs of each layer can be estimated by

$$FLOPs^{\text{decode}} = 4d^2 + 2nd + 3dm. \quad (7)$$

3.2.3 Layer-wise Pruning Schedule

Following the attention shift analysis described above, we design a pruning schedule that aligns with the attention dynamics observed in VLMs. Our goal is to perform aggressive pruning while maintaining the inference accuracy. To balance the trade-off, we suggest that pruning should avoid the regions with intense attention shifts. Because the probability of incorrect pruning is quite high, tokens that are not important in the past layers are likely to become important in the next layer. On the contrary, pruning just before or after the dense regions can reduce this risk to some extent.

Based on the consistent attention shift distributions across both MME and TextVQA datasets, we select layers 1, 10, and 20 as pruning locations in our framework. Pruning at layer 1 implies performing a forward pass through the first layer, and subsequently using its attention outputs to identify and remove redundant vision tokens based on their ranked importance. Since layer 1 exhibits the earliest dense attention shift, retaining the salient vision tokens at this stage allows them to effectively participate in subsequent information exchange. Layer 10 begins to mark a region of high attention volatility until layer 20. As a result, we preserve the vision tokens within this range, only compress the token numbers at the beginning and the end. At layer 20 and after, where shift activity is negligible, we suggest that the cross-modality information exchange is already sufficient, so we apply aggressive pruning at this layer by removing nearly all remaining visual tokens. This schedule balances caution and efficiency. Compared to heuristic or uniform pruning schemes, our approach is both data-driven and architecture-aware, requiring no expensive hyperparameter tuning while generalizing well across tasks and datasets.

3.3 Discussion

The chosen pruning hyperparameters come from empirical observations on LLaVA-1.5-7B, and thus are tailored specifically to VLMs built upon the LLaMA-7B backbone. Nevertheless, our proposed

Table 1: **Comparison of methods under different pruning budgets on LLava-1.5-7B.** We report results on five datasets, including average retained tokens, accuracy scores and average ratios. Results in bold present the highest accuracy under same pruning budgets.

Method	Avg. Tokens	MME	GQA	MMB	SQA	TVQA	Ratio (%)
Vanilla	576	1864	61.9	64.6	69.5	58.3	100
ToMe (ICLR23)	128	1343	52.4	53.3	59.6	49.1	81.8 (↓ 18.2)
FastV (ECCV24)	128	1490	49.6	56.1	68.6	52.5	87.1 (↓ 12.9)
PDrop (CVPR25)	128	1761	57.1	61.6	68.4	56.6	95.5 (↓ 4.5)
SparseVLM (ICML25)	128	1746	58.4	64.5	68.6	56.7	96.8 (↓ 3.2)
AdaptInfer (Ours)	128	1794	58.5	63.8	69.9	56.8	97.5 (↓ 2.5)
ToMe (ICLR23)	64	1138	48.6	43.7	50.0	45.3	71.4 (↓ 28.6)
FastV (ECCV24)	64	1255	46.1	47.2	68.7	45.9	78.5 (↓ 21.5)
PDrop (CVPR 25)	64	1561	47.5	58.8	69.0	50.6	87.5 (↓ 12.5)
SparseVLM (ICML 25)	64	1589	53.8	60.1	69.8	53.4	91.4 (↓ 8.6)
AdaptInfer (Ours)	64	1684	53.2	61.7	69.9	54.3	93.1 (↓ 6.9)

adaptive pruning schedule is generalizable for it can be easily transferred to other models with different parameter scales or architectures by performing a simple, offline attention shift analysis. Moreover, one of the datasets used in our observational studies is MME [Fu et al., 2023], a comprehensive multimodal benchmark comprising two major categories and fourteen subcategories. The consistent statistical patterns indicate that the attention dynamics are largely stable and transferable across different types of multimodal tasks.

4 Experiment

4.1 Experimental Settings

4.1.1 Datasets

For multimodal evaluation, we test our solutions on five widely used benchmarks, including MME [Fu et al., 2023], GQA [Hudson and Manning, 2019], MMBench (MMB) [Liu et al., 2024b], ScienceQA (SQA) [Lu et al., 2022] and TextVQA (TVQA) [Singh et al., 2019]. These datasets together provide a comprehensive evaluation, including vision-question-answering (VQA), optical character recognition (OCR), perception, reasoning, factual grounding and so on.

4.1.2 Baselines

We select four classic and latest vision token sparsification frameworks for VLM acceleration within the plug-and-play paradigm in LLM forward pass as baselines, including FastV [Chen et al., 2025a], ToMe [Bolya et al., 2023], Pyramid Drop (PDrop) [Xing et al., 2024] and SparseVLM [Zhang et al., 2025]. For PDrop, we only adopt the training-free version of PDrop to fit our requirements. Note that, methods performing token pruning or merging anywhere other than the prefill stage of language models, like vision encoders or the decode stage [Yang et al., 2024, Tong et al., 2025], are not included for comparison, since such approaches are orthogonal to ours and can be used together.

4.1.3 Implement Details

To ensure a fair and comprehensive comparison, we report the results in different average retained token budgets (128, 64, 48 and 32). AdaptInfer utilizes the pruning parameters described above, while other baselines keep their original settings for all core components. Our experiments are performed on two different types of VLMs, including LLava-1.5-7B [Liu et al., 2023b] and InternVL-chat-7B [Chen et al., 2024a]. Additional details are presented in Appendix C.

4.2 Main Results

In Table 1, we report the performance of AdaptInfer on LLava-1.5-7B. We provide comparisons with other baselines under average retained vision token numbers of 128 and 64. The accuracy scores have an upper bound of Vanilla LLava with all 576 tokens kept, so an average accuracy ratio is supported

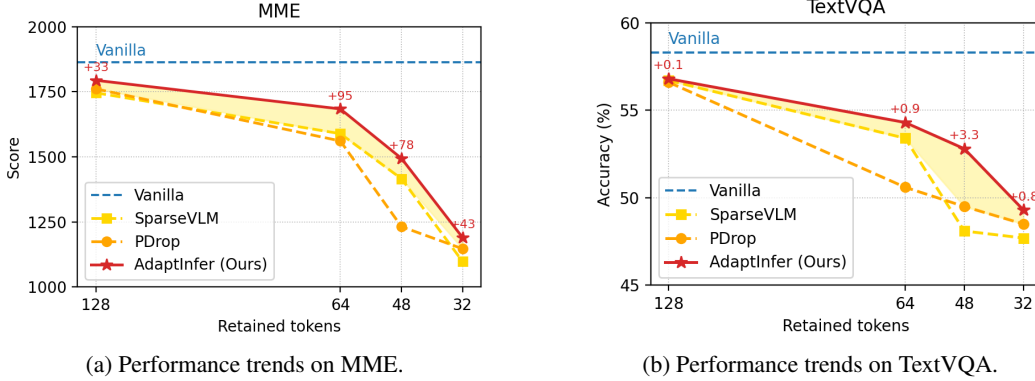


Figure 3: **Performance trends on two datasets.** AdaptInfer outperforms SparseVLM and PDrop across all retained token budgets.

Table 2: **Performance of AdaptInfer on InternVL.** Our method still maintain high accuracy.

Tokens	MME	SQA	TVQA	Ratio (%)
Origin	1849	69.1	56.9	100
128	1758	69.2	55.7	97.7
64	1528	69.0	53.0	91.9
32	1183	68.1	46.7	81.5

by each baseline. As shown in this Table, AdaptInfer achieves the highest overall accuracy scores under both 128 and 64 vision token budgets. Under 128 token budgets, our average ratio reaches 97.5%, which is 0.7% higher than the second-best method, SparseVLM. While this may seem like a small margin, it is in fact approaching the upper bound of Vanilla LLaVA. In practice, retaining only 64 tokens on average reduces the prefill token load by 88.9%, enabling significantly more efficient inference. Despite this aggressive pruning, our framework still preserves 93.1% of the original inference accuracy and outperforms SparseVLM, which scores 91.4%, with a clear improvement of 1.7%. In Figure 3, we present the performance trends of AdaptInfer with the latest baseline methods SparseVLM and PDrop. AdaptInfer outperforms others with clear margins.

We further evaluate the proposed AdaptInfer on InternVL-Chat-7B [Chen et al., 2024a] in Table 2. InternVL adopts a more powerful 6B parameter ViT encoder, which produces finer-grained vision tokens with richer hidden states. Experiments confirm that our token-sparsification strategy preserves the original semantic richness of these vision features with high inference accuracy. Under token budgets of 128 and 64, AdaptInfer maintains a high accuracy ratio of 97.7% and 91.9% respectively. However, in the case of extreme pruning, where only 32 average vision tokens are retained, the accuracy has a significant drop to only 81.5%. This is because the inevitable removal of a large number of informative tokens affects the performance.

4.3 Latency Test

In addition, we conduct a comparative analysis of CUDA latency time, and FLOPs on LLaVA-1.5-7B, in order to show the real acceleration performance. The results are shown in Table 3. This experiment is carried out on a single NVIDIA RTX 4090 24G GPU. All results are the average values per sample.

We replicate the performance of PDrop [Xing et al., 2024] and SparseVLM [Zhang et al., 2025] under 64-token budget, and compare them with that of AdaptInfer. The comparisons are based on four commonly used benchmarks, including MME, GQA, MMBench (MMB) and TextVQA (TVQA). According to the table, while theoretical FLOPs estimations are close, our implementation on AdaptInfer reaches a lower average cuda latency of 33.0 ms per sample than 34.5 ms and 36.7 ms by PDrop and SparseVLM. After all, our method rarely introduces additional computational load, such as any extra attention computation or offline feature matching step.

Table 3: **Latency test of AdaptInfer.** While the FLOPs estimations are close, our method reduces more real CUDA latency.

Method	Tokens	Metrics	MME	GQA	MMB	TVQA	Average
Vanilla	576	FLOPs (T)	4.268	4.250	4.623	4.611	4.438
		Latency (ms)	82.0	76.5	91.3	91.2	85.3
PDrop	64	FLOPs (T)	0.975	0.958	1.316	1.305	1.138 (↓ 74.3%)
		Latency (ms)	33.0	32.0	36.4	36.6	34.5 (↓ 59.6%)
SparseVLM	64	FLOPs (T)	0.974	0.958	1.316	1.305	1.138 (↓ 74.3%)
		Latency (ms)	34.4	34.7	38.1	39.5	36.7 (↓ 57.0%)
AdaptInfer (Ours)	64	FLOPs (T)	0.975	0.959	1.317	1.305	1.139 (↓ 74.3%)
		Latency (ms)	31.0	30.9	34.5	35.5	33.0 (↓ 61.3%)

Table 4: **Performance w/ and w/o dynamic text-guidance.** Comparisons are between static (SparseVLM) and dynamic text-guidance (Ours) methods under extreme pruning on LLaVa-1.5-7B.

Method	Tokens	MME	TVQA	Ratio (%)
Vanilla	576	1864	58.3	100
SparseVLM	48	1416	48.1	79.2
AdaptInfer	48	1494	52.8	85.4
SparseVLM	32	1098	47.7	70.4
AdaptInfer	32	1190	49.3	74.2

4.4 Ablation Study

Firstly, in order to illustrate the necessity of dynamic text guidance, we provide an experiment to compare results with a static text guidance work SparseVLM in Table 4. Also, we conduct this study as an exploration of extreme pruning. In the experiment, AdaptInfer consistently outperforms SparseVLM under extreme low pruning budgets of only 48 and 32 vision tokens retained averagely. Specifically, AdaptInfer can keep a relatively high overall performance of 85.4% and 74.2% respectively. These results prove the effectiveness of our dynamic text guidance design.

Furthermore, we conduct an intuitive study to assess the effectiveness of our pruning hyperparameters. We compare our observation-driven pruning strategy with three baselines: uniform pruning, single-layer pruning, and random-layer pruning which averages results over five random configurations. The results shown in Table 5 confirm that the chosen hyperparameters are the best suited for AdaptInfer on the LLaMa-7B backbone.

Table 5: **Performance of different pruning locations on LLaVa-1.5-7B with 128 tokens.** The chosen hyperparameters, inspired by our observations, yield the highest scores.

Method	Pruning Loc.	MME	MMB	TVQA
Ours	1,10,20	1794	63.8	56.8
Uniform	0,9,18,27	1788	62.8	56.4
Uniform	2,12,22	1758	63.2	56.5
Single	1	1668	60.8	56.3
Random	-	1679	62.1	55.8

5 Conclusion

This paper proposes a novel plug-and-play solution AdaptInfer for VLM acceleration via dynamic text-guided pruning. We also provide an offline analysis of cross-attention shifts, which motivates a principled pruning schedule. Our VLM acceleration plugin introduces minimal additional computational overhead while maintaining high accuracy. In particular, AdaptInfer achieves SOTA accuracy under all the token budgets. For instance, AdaptInfer reduces CUDA latency by 61.3% and retains an accuracy of 93.1% on LLaVA-1.5-7B, with only 64 vision tokens preserved per layer on average.

Acknowledgment

This research was supported by National Natural Science Foundation of China (Grants No. 62472248, Grants No. 62432008).

References

- J.-B. Alayrac, J. Donahue, P. Luc, A. Miech, I. Barr, Y. Hasson, K. Lenc, A. Mensch, K. Millicah, M. Reynolds, R. Ring, E. Rutherford, S. Cabi, T. Han, Z. Gong, S. Samangooei, M. Monteiro, J. Menick, S. Borgeaud, A. Brock, A. Nematzadeh, S. Sharifzadeh, M. Binkowski, R. Barreira, O. Vinyals, A. Zisserman, and K. Simonyan. Flamingo: a visual language model for few-shot learning. In *Proceedings of the 36th International Conference on Neural Information Processing Systems*, NIPS '22, Red Hook, NY, USA, 2022. Curran Associates Inc. ISBN 9781713871088.
- S. Antol, A. Agrawal, J. Lu, M. Mitchell, D. Batra, C. L. Zitnick, and D. Parikh. Vqa: Visual question answering. In *2015 IEEE International Conference on Computer Vision (ICCV)*, pages 2425–2433, 2015. doi: 10.1109/ICCV.2015.279.
- J. Bai, S. Bai, S. Yang, S. Wang, S. Tan, P. Wang, J. Lin, C. Zhou, and J. Zhou. Qwen-vl: A versatile vision-language model for understanding, localization, text reading, and beyond. *arXiv preprint arXiv:2308.12966*, 2023.
- D. Bolya, C.-Y. Fu, X. Dai, P. Zhang, C. Feichtenhofer, and J. Hoffman. Token merging: Your ViT but faster. In *International Conference on Learning Representations*, 2023.
- R. Bommasani, D. A. Hudson, E. Adeli, R. B. Altman, S. Arora, S. von Arx, M. S. Bernstein, J. Bohg, A. Bosselut, E. Brunskill, E. Brynjolfsson, S. Buch, D. Card, R. Castellon, N. S. Chatterji, A. S. Chen, K. Creel, J. Q. Davis, D. Demszky, C. Donahue, M. Doumbouya, E. Durmus, S. Ermon, J. Etchemendy, K. Ethayarajh, L. Fei-Fei, C. Finn, T. Gale, L. E. Gillespie, K. Goel, N. D. Goodman, S. Grossman, N. Guha, T. Hashimoto, P. Henderson, J. Hewitt, D. E. Ho, J. Hong, K. Hsu, J. Huang, T. Icard, S. Jain, D. Jurafsky, P. Kalluri, S. Karamcheti, G. Keeling, F. Khani, O. Khattab, P. W. Koh, M. S. Krass, R. Krishna, R. Kuditipudi, and et al. On the opportunities and risks of foundation models. *CoRR*, abs/2108.07258, 2021. URL <https://arxiv.org/abs/2108.07258>.
- T. B. Brown, B. Mann, N. Ryder, M. Subbiah, J. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell, S. Agarwal, A. Herbert-Voss, G. Krueger, T. Henighan, R. Child, A. Ramesh, D. M. Ziegler, J. Wu, C. Winter, C. Hesse, M. Chen, E. Sigler, M. Litwin, S. Gray, B. Chess, J. Clark, C. Berner, S. McCandlish, A. Radford, I. Sutskever, and D. Amodei. Language models are few-shot learners. In *Proceedings of the 34th International Conference on Neural Information Processing Systems*, NIPS '20, Red Hook, NY, USA, 2020. Curran Associates Inc. ISBN 9781713829546.
- L. Chen, H. Zhao, T. Liu, S. Bai, J. Lin, C. Zhou, and B. Chang. An image is worth 1/2 tokens after layer 2: Plug-and-play inference acceleration for large vision-language models. In A. Leonardis, E. Ricci, S. Roth, O. Russakovsky, T. Sattler, and G. Varol, editors, *Computer Vision – ECCV 2024*, pages 19–35, Cham, 2025a. Springer Nature Switzerland. ISBN 978-3-031-73004-7.
- Y. Chen, J. Xu, X.-Y. Zhang, W.-Z. Liu, Y.-Y. Liu, and C.-L. Liu. Recoverable compression: A multimodal vision token recovery mechanism guided by text information. *Proceedings of the AAAI Conference on Artificial Intelligence*, 39(2):2293–2301, Apr. 2025b. doi: 10.1609/aaai.v39i2.32229. URL <https://ojs.aaai.org/index.php/AAAI/article/view/32229>.
- Z. Chen, W. Wang, H. Tian, S. Ye, Z. Gao, E. Cui, W. Tong, K. Hu, J. Luo, Z. 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, 2024a.
- Z. Chen, J. Wu, W. Wang, W. Su, G. Chen, S. Xing, M. Zhong, Q. Zhang, X. Zhu, L. Lu, et al. Intervl: 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, 2024b.
- K. Clark, U. Khandelwal, O. Levy, and C. D. Manning. What does BERT look at? an analysis of bert’s attention. *CoRR*, abs/1906.04341, 2019. URL <http://arxiv.org/abs/1906.04341>.

- A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly, J. Uszkoreit, and N. Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. In *International Conference on Learning Representations*, 2021. URL <https://openreview.net/forum?id=YicbFdNTTy>.
- Y. Du, Z. Liu, J. Li, and W. X. Zhao. A survey of vision-language pre-trained models. In L. D. Raedt, editor, *Proceedings of the Thirty-First International Joint Conference on Artificial Intelligence, IJCAI-22*, pages 5436–5443. International Joint Conferences on Artificial Intelligence Organization, 7 2022. doi: 10.24963/ijcai.2022/762. URL <https://doi.org/10.24963/ijcai.2022/762>. Survey Track.
- F. Faghri, D. J. Fleet, J. R. Kiros, and S. Fidler. Vse++: Improving visual-semantic embeddings with hard negatives. In *Proceedings of BMVC*, 2018. URL <https://github.com/fartashf/vsepp>.
- C. Fu, P. Chen, Y. Shen, Y. Qin, M. Zhang, X. Lin, J. Yang, X. Zheng, K. Li, X. Sun, et al. Mme: A comprehensive evaluation benchmark for multimodal large language models. *arXiv preprint arXiv:2306.13394*, 2023.
- D. A. Hudson and C. 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 (CVPR)*, June 2019.
- A. Karpathy and L. Fei-Fei. Deep visual-semantic alignments for generating image descriptions. *IEEE Trans. Pattern Anal. Mach. Intell.*, 39(4):664–676, Apr. 2017. ISSN 0162-8828. doi: 10.1109/TPAMI.2016.2598339. URL <https://doi.org/10.1109/TPAMI.2016.2598339>.
- J. Li, D. Li, S. Savarese, and S. Hoi. Blip-2: bootstrapping language-image pre-training with frozen image encoders and large language models. In *Proceedings of the 40th International Conference on Machine Learning, ICML’23*. JMLR.org, 2023.
- Y. Li, C. Wang, and J. Jia. Llama-vid: An image is worth 2 tokens in large language models. In *Computer Vision – ECCV 2023*, 2024.
- Y. Li, H. Jiang, C. Zhang, Q. Wu, X. Luo, S. Ahn, A. H. Abdi, D. Li, J. Gao, Y. Yang, and L. Qiu. Mminference: Accelerating pre-filling for long-context vlms via modality-aware permutation sparse attention, 2025. URL <https://arxiv.org/abs/2504.16083>.
- T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick. Microsoft coco: Common objects in context. In D. Fleet, T. Pajdla, B. Schiele, and T. Tuytelaars, editors, *Computer Vision – ECCV 2014*, pages 740–755, Cham, 2014. Springer International Publishing. ISBN 978-3-319-10602-1.
- Z. Lin, M. Lin, L. Lin, and R. Ji. Boosting multimodal large language models with visual tokens withdrawal for rapid inference. *Proceedings of the AAAI Conference on Artificial Intelligence*, 39(5):5334–5342, Apr. 2025. doi: 10.1609/aaai.v39i5.32567. URL <https://ojs.aaai.org/index.php/AAAI/article/view/32567>.
- H. Liu, C. Li, Y. Li, and Y. J. Lee. Improved baselines with visual instruction tuning, 2023a.
- H. Liu, C. Li, Q. Wu, and Y. J. Lee. Visual instruction tuning, 2023b.
- H. Liu, C. Li, Y. Li, B. Li, Y. Zhang, S. Shen, and Y. J. Lee. Llava-next: Improved reasoning, ocr, and world knowledge, January 2024a. URL <https://llava-vl.github.io/blog/2024-01-30-llava-next/>.
- Y. Liu, H. Duan, Y. Zhang, B. Li, S. Zhang, W. Zhao, Y. Yuan, J. Wang, C. He, Z. Liu, et al. Mmbench: Is your multi-modal model an all-around player? In *Computer Vision – ECCV 2023*, pages 216–233. Springer, 2024b.
- Z. Liu, H. Mao, C.-Y. Wu, C. Feichtenhofer, T. Darrell, and S. Xie. A convnet for the 2020s. In *2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 11966–11976, 2022. doi: 10.1109/CVPR52688.2022.01167.

- P. Lu, S. Mishra, T. Xia, L. Qiu, K.-W. Chang, S.-C. Zhu, O. Tafjord, P. Clark, and A. Kalyan. Learn to explain: Multimodal reasoning via thought chains for science question answering. In *The 36th Conference on Neural Information Processing Systems (NeurIPS)*, 2022.
- B. Luan, W. Zhou, H. Feng, Z. Wang, X. Li, and H. Li. Multi-cue adaptive visual token pruning for large vision-language models, 2025. URL <https://arxiv.org/abs/2503.08019>.
- D. Marr. *Vision: A Computational Investigation into the Human Representation and Processing of Visual Information*. The MIT Press, 07 2010. ISBN 9780262514620. doi: 10.7551/mitpress/9780262514620.001.0001. URL <https://doi.org/10.7551/mitpress/9780262514620.001.0001>.
- OpenAI. Gpt-4v system card. Online, September 2023. Available at <https://cdn.openai.com/papers/GPTVsystemcard.pdf>.
- A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell, P. Mishkin, J. Clark, G. Krueger, and I. Sutskever. Learning transferable visual models from natural language supervision. In M. Meila and T. Zhang, editors, *Proceedings of the 38th International Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*, pages 8748–8763. PMLR, 18–24 Jul 2021. URL <https://proceedings.mlr.press/v139/radford21a.html>.
- Y. Shang, M. Cai, B. Xu, Y. J. Lee, and Y. Yan. Llava-prumerge: Adaptive token reduction for efficient large multimodal models. In *ICCV*, 2025.
- A. Singh, V. Natarjan, M. Shah, Y. Jiang, X. Chen, D. Parikh, and M. Rohrbach. Towards vqa models that can read. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 8317–8326, 2019.
- I. Tenney, D. Das, and E. Pavlick. BERT rediscovers the classical NLP pipeline. In A. Korhonen, D. Traum, and L. Màrquez, editors, *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 4593–4601, Florence, Italy, July 2019. Association for Computational Linguistics. doi: 10.18653/v1/P19-1452. URL <https://aclanthology.org/P19-1452/>.
- J. Tong, W. Jin, P. Qin, A. Li, Y. Zou, Y. Li, Y. Li, and R. Li. Flowcut: Rethinking redundancy via information flow for efficient vision-language models, 2025. URL <https://arxiv.org/abs/2505.19536>.
- H. Touvron, T. Lavril, G. Izacard, X. Martinet, M.-A. Lachaux, T. Lacroix, B. Rozière, N. Goyal, E. Hambro, F. Azhar, A. Rodriguez, A. Joulin, E. Grave, and G. Lample. Llama: Open and efficient foundation language models, 2023. URL <https://arxiv.org/abs/2302.13971>.
- C. Truong, L. Oudre, and N. Vayatis. Selective review of offline change point detection methods. *Signal Processing*, 167:107299, Feb. 2020. doi: 10.1016/j.sigpro.2019.107299. URL <https://hal.science/hal-02442692>.
- A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin. Attention is all you need. In *Proceedings of the 31st International Conference on Neural Information Processing Systems, NIPS’17*, page 6000–6010, Red Hook, NY, USA, 2017. Curran Associates Inc. ISBN 9781510860964.
- O. Vinyals, A. Toshev, S. Bengio, and D. Erhan. Show and tell: A neural image caption generator. In *2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 3156–3164, 2015. doi: 10.1109/CVPR.2015.7298935.
- W. Wang, Q. Lv, W. Yu, W. Hong, J. Qi, Y. Wang, J. Ji, Z. Yang, L. Zhao, X. Song, J. Xu, B. Xu, J. Li, Y. Dong, M. Ding, and J. Tang. Cogvlm: Visual expert for pretrained language models, 2024. URL <https://arxiv.org/abs/2311.03079>.
- Z. Wu, J. Chen, and Y. Wang. Unified knowledge maintenance pruning and progressive recovery with weight recalling for large vision-language models. *Proceedings of the AAAI Conference on Artificial Intelligence*, 39(8):8550–8558, Apr. 2025. doi: 10.1609/aaai.v39i8.32923. URL <https://ojs.aaai.org/index.php/AAAI/article/view/32923>.

- L. Xing, Q. Huang, X. Dong, J. Lu, P. Zhang, Y. Zang, Y. Cao, C. He, J. Wang, F. Wu, and D. Lin. Pyramiddrop: Accelerating your large vision-language models via pyramid visual redundancy reduction. *CoRR*, abs/2410.17247, 2024. URL <https://doi.org/10.48550/arXiv.2410.17247>.
- S. Yang, Y. Chen, Z. Tian, C. Wang, J. Li, B. Yu, and J. Jia. Visionzip: Longer is better but not necessary in vision language models, 2024. URL <https://arxiv.org/abs/2412.04467>.
- W. Ye, Q. Wu, W. Lin, and Y. Zhou. Fit and prune: Fast and training-free visual token pruning for multi-modal large language models. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 39, pages 22128–22136, 2025.
- Q. Zhang, A. Cheng, M. Lu, Z. Zhuo, M. Wang, J. Cao, S. Guo, Q. She, and S. Zhang. [cls] attention is all you need for training-free visual token pruning: Make vlm inference faster. *arXiv preprint arXiv:2412.01818*, 2024.
- Y. Zhang, C.-K. Fan, J. Ma, W. Zheng, T. Huang, K. Cheng, D. Gudovskiy, T. Okuno, Y. Nakata, K. Keutzer, et al. Sparsevlm: Visual token sparsification for efficient vision-language model inference. In *International Conference on Machine Learning*, 2025.
- D. Zhu, J. Chen, X. Shen, X. Li, and M. Elhoseiny. Minigpt-4: Enhancing vision-language understanding with advanced large language models, 2023. URL <https://arxiv.org/abs/2304.10592>.

A Appendix and LLM Usage statement

The chapters below are all technical appendix.

The authors declare that they use LLM tools for grammar refinement and language polishing only. The authors take full responsibility for the entire content of this manuscript.

B Change-Point Detection

To identify shifts in cross-attention across transformer layers, we employ change-point detection [Truong et al., 2020] on cumulative cross-attention curves. The intuition is that a significant change in attention behavior corresponds to a noticeable inflection in the cumulative attention distribution over layers.

Given a token’s cross-attention curve across L layers, denoted as:

$$\mathbf{a} = [a_1, a_2, \dots, a_L], \quad \text{where } a_\ell \in \mathbb{R} \quad (8)$$

we first compute the cumulative attention:

$$\mathbf{c} = \left[\sum_{i=1}^1 a_i, \sum_{i=1}^2 a_i, \dots, \sum_{i=1}^L a_i \right] \quad (9)$$

This sequence \mathbf{c} captures the aggregate build-up of attention, which we treat as a univariate time series. We then apply change-point detection to identify the most likely layer index where a shift in trend occurs.

We adopt the ruptures Python library² for offline detection. Specifically, we use the Binary Segmentation (Binseg) algorithm with an ℓ_2 cost model. This algorithm aims to find a segmentation point that minimizes the within-segment variance. Formally, given the cumulative curve $\mathbf{c} = [y_1, y_2, \dots, y_L]$, we solve:

$$\min_b \left(\sum_{t=1}^b (y_t - \mu_1)^2 + \sum_{t=b+1}^L (y_t - \mu_2)^2 \right), \quad (10)$$

²<https://github.com/deepcharles/ruptures>

where μ_1 and μ_2 are the mean values of the first and second segments, respectively. This corresponds to the ℓ_2 cost model used in ruptures, which measures the homogeneity of each segment by its squared deviation from the mean.

We restrict the number of change points to 1, making the detection both efficient and robust. The returned breakpoint $b \in \{1, \dots, L\}$ indicates the first significant shift in attention accumulation for the given token. This layer index is treated as the token’s *cross-attention shift point*, which guides our downstream pruning strategy.

C Experiment Setup

In this appendix section, we will briefly introduce the details of our experiment setups.

C.1 Hardware Environment

To facilitate reproducibility, we ran every experiment on a single workstation that hosts an Intel Xeon Platinum 8358P processor together with eight NVIDIA RTX 4090 GPUs, each furnished with 24 GB of VRAM. The system operates under Ubuntu Linux 24.04, and all code was compiled using GCC 13.2.0 with Binutils 2.42 as the linker.

C.2 Software Configuration

To improve reproducibility, we provide the main software configuration employed in our experiments. All runs were performed with the package versions listed below. We introduce the core deep-learning frameworks first, follow with performance-oriented extensions, and conclude with supporting libraries.

- **Python:** 3.10.18
- **PyTorch:** 2.1.2
- **Transformers:** 4.37.0
- **Accelerate:** 0.21.0
- **Flash-Attn:** 2.3.3
- **LLaVA:** 1.7.0.dev0
- **Tokenizers:** 0.15.1
- **TorchVision:** 0.16.2
- **ruptures:** 1.1.9

C.3 Implement Details

In our evaluation of AdaptInfer, the pruning ratios at each pruning location are the hyperparameters need to be selected as well. For these hyperparameters, we follow the same sparsification-level adaptation method which SparseVLM [Zhang et al., 2025] introduced.

In the paper, we evaluated three vision–language models, including LLaVA-1.5-7B, LLaVA-1.5-13B, and InternVL-Chat-7B [Chen et al., 2024a] in total. To stay within GPU-memory limits, we ran the LLaVA-7B variant on one GPU, the LLaVA-13B variant on three GPUs, and InternVL-Chat-7B on two GPUs (InternVL-Chat-7B contains a 6B-parameter vision encoder). All training and inference are performed in torch.float16 precision to further conserve memory usage. All the results of AdaptInfer reported in the context are the average of five runs.

C.4 Randomness Report

Table 5 in the main paper contrasts our proposed *layer-wise pruning schedule* with three baselines: *uniform-layer pruning*, *single-layer pruning*, and *random-layer pruning*. For the random baseline, we executed five independent runs, each time drawing a different set of transformer layers to prune while keeping the token budget fixed. The sampled pruning locations are $\{3\}$, $\{2,15\}$, $\{2,8,16\}$, $\{2,4,8,16\}$, $\{3,6,23\}$.