

# HT-SPARSE: TRAINING-FREE QUERY-GUIDED HEAD-TOKEN SPARSIFICATION FOR LONG-VIDEO MULTIMODAL INFERENCE

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Long-video multimodal inference is limited by the quadratic cost of dense attention, cumulative KV-cache growth during decoding, and cross-modal interference, while retraining sparsity-aware variants is often impractical. We present **HT-Sparse**, a *training-free, query-guided* hierarchical sparsification that performs *joint head-token computation* to reduce both latency and memory without parameter updates. The method comprises two components executed adaptively across layers: (i) *query-conditioned head sparsification*, which ranks attention heads via analytically stable saliency statistics to retain the most informative subspaces for the current query; (ii) *cross-modal token sparsification*, which selects salient visual tokens by query-vision attention, enabling efficient computation and persistent KV-cache savings. We further introduce *joint head-token routing* in selected layers: top-ranked heads attend to the *full* visual token set, whereas secondary heads operate on the *reduced* (selected) set, preserving semantics while amortizing compute and cache. Across long-video benchmarks, HT-Sparse delivers faster inference with reduced end-to-end latency and lower KV-cache memory, while achieving equal or higher accuracy, all on the same pretrained model with no fine-tuning. The approach is model-agnostic and plug-in deployable, offering a flexible route to scalable long-video reasoning.

## 1 INTRODUCTION

Foundation models with vision-language capability have rapidly advanced multimodal reasoning (Anil et al., 2023; Liu et al., 2024a; 2025b), yet their inference on long-form video remains prohibitive due to the quadratic complexity of dense attention (Vaswani et al., 2017), cumulative KV-cache growth across decoding (Liu et al., 2024b; Ge et al., 2024), and cross-modal interference that disperses evidence over thousands of visual tokens (Li et al., 2023a). In production settings, retraining or fine-tuning sparsity-aware variants is often impractical due to data governance, latency-to-market, and the risk of distribution shift (Xiao et al., 2024a; Han et al., 2024). Consequently, there is a pressing need for *training-free, input-adaptive* mechanisms that reduce computation and memory while preserving task performance (Gao et al., 2025; Wu et al., 2025).

**Limitations of existing approaches.** Fixed or hand-crafted sparse patterns constrain the model to a single layout of query-key interactions, which under-utilizes content adaptivity and can degrade accuracy on heterogeneous inputs (Lee et al., 2025; Liu et al., 2021). Heuristic token dropping improves throughput but is brittle for long video streams where informative regions are temporally sparse and semantically entangled with distractors (Fu et al., 2024; Tao et al., 2025; Zhang et al., 2024; 2025a; Liu et al., 2025a). Methods that operate at a single granularity (e.g., only head pruning or only token selection) leave substantial efficiency on the table: attention heads specialize different subspaces, while visual tokens carry complementary, partially redundant information (Fu et al., 2025b). Finally, mechanisms that rely on parameter updates or task-specific retraining limit deployability across models and domains (Pan et al., 2024; Li et al., 2023b).

---

Code availability: we will release code and evaluation scripts upon acceptance.

054  
055  
056  
057  
058  
059  
060  
061  
062  
063  
064  
065  
066  
067  
068  
069  
070  
071  
072  
073  
074  
075  
076  
077  
078  
079  
080  
081  
082  
083  
084  
085  
086  
087  
088  
089  
090  
091  
092  
093  
094  
095  
096  
097  
098  
099  
100  
101  
102  
103  
104  
105  
106  
107

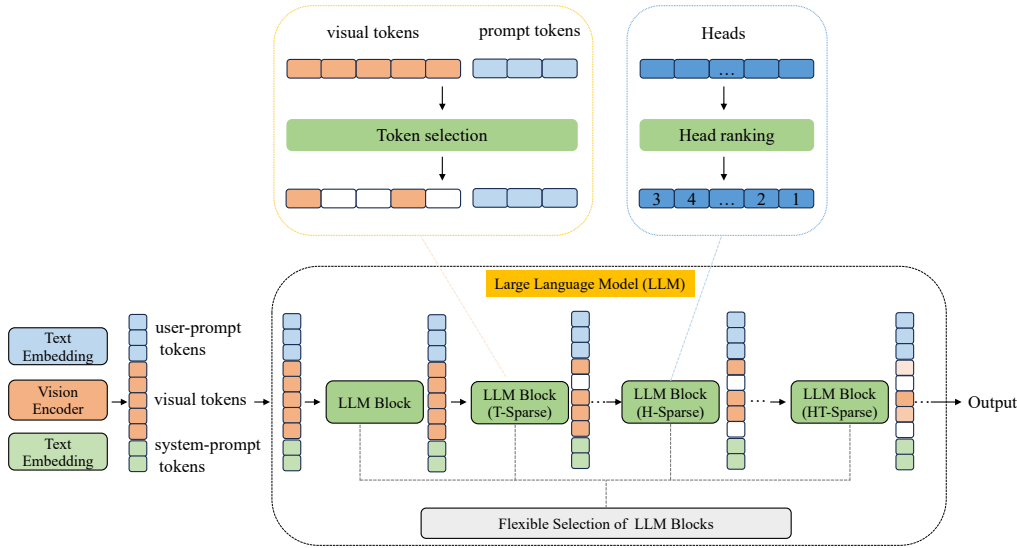


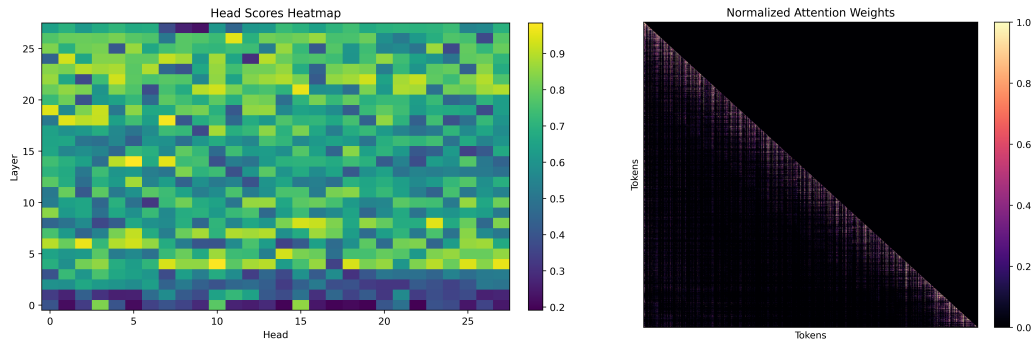
Figure 1: **HT-Sparse selection mechanism.** Given two input streams—a question and a long video—the multimodal model routes computation across layers and, at each layer, selectively enables one of three training-free strategies: token selection, head ranking, or joint head–token selection.

**Our approach in brief.** We propose **HT-Sparse**, a *training-free, query-guided* hierarchical sparsification for multimodal long-video inference. HT-Sparse performs *joint head–token computation* that adapts per input and per layer: (i) *query-conditioned head sparsification* ranks attention heads using analytically stable saliency statistics to retain the most informative subspaces for the current query; (ii) *cross-modal token sparsification*, which selects salient visual tokens via query–vision attention, yielding persistent savings in computation and KV-cache usage. When resource envelopes are tight, an optional *in-attention low-rank projection* further contracts dimensionality while preserving fidelity. The mechanism plugs into existing VLMs without parameter updates, making it suitable for latency- and memory-critical deployments. Crucially, HT-Sparse performs *joint head–token routing* in selected layers: a small set of top heads attends to all visual tokens to guard semantic coverage, while secondary heads attend to the *selected tokens* to yield persistent compute and KV-cache savings.

**Design principles.** HT-Sparse is built on three principles:

1. **Hierarchical adaptivity.** Sparsification operates at complementary levels (heads and tokens) and adapts across layers, enabling the method to align with the evolving representational needs of the query as it propagates through the network.
2. **Cross-modal selectivity.** Visual token selection is driven by query–vision interactions, reducing cardinality and cache footprint while retaining the most informative tokens.
3. **Joint routing for fidelity and efficiency.** Top heads operate on full tokens to prevent semantic loss, while secondary heads consume reduced tokens, coupling head- and token-level sparsification within the same layer.
4. **Deployment practicality.** The procedure is training-free, model-agnostic, and requires only lightweight additions around standard attention, facilitating immediate integration into production inference stacks.

**Technical overview.** Let  $H$  denote the number of heads,  $L_v$  the number of visual tokens, and  $d$  the head dimension. Dense cross-modal attention scales as  $\mathcal{O}(H L_v^2 d)$  in the prefill stage and induces large KV-cache costs in decoding. HT-Sparse reduces both factors: head sparsification effectively replaces  $H$  by  $\tilde{H} \ll H$  through query-conditioned ranking; token sparsification reduces  $L_v$  to  $\tilde{L}_v$ ,



(a) Head saliency is layer- and task-dependent rather than uniform, revealing stable non-uniform importance patterns.

(b) Visual tokens are selectively attended by the query, revealing non-uniform, task-relevant patterns instead of diffuse attention.

Figure 2: **Structured sparsity in multimodal attention.** (a) Head saliency is layer- and task-dependent rather than uniform, exhibiting stable non-uniform importance across layers; (b) query-vision attention concentrates on a selective subset of visual tokens instead of diffusing broadly. Together, these observations reveal inherent, input-adaptive sparsity that motivates joint head-token sparsification at inference.

retaining only the most salient tokens for computation and cache. When enabled, low-rank projection replaces  $d$  by  $\tilde{d}$  for queries and keys within attention, while values remain full-dimensional to preserve representation fidelity. The resulting attention score cost scales as  $\mathcal{O}(\tilde{H} \tilde{L}_v^2 \tilde{d})$ , with value aggregation still computed in the original  $d$ -dimensional space, leading to proportional savings in KV-cache during decoding.

### Contributions.

- **A training-free, query-guided hierarchical sparsification for multimodal long-video inference.** HT-Sparse jointly performs head and token sparsification under a unified, input-adaptive procedure that requires no parameter updates.
- **Cross-modal token sparsification.** A content-aware selection of salient visual tokens driven by query-vision attention, reducing compute and cache persistence while retaining the most informative tokens.
- **Optional low-rank contraction within attention.** A lightweight projection inside attention further reduces dimensionality when resources are tight, complementing head-token sparsification.
- **Model-agnostic deployment and empirical validation.** The method integrates into standard VLMs without fine-tuning and consistently reduces end-to-end latency and KV-cache memory on long-video benchmarks while maintaining accuracy.

**Scope and implications.** HT-Sparse targets inference-time efficiency for vision-language models operating on extended video streams and remains applicable to other multimodal settings that exhibit long-context behavior (e.g., document understanding). By eliminating retraining, it decouples efficiency gains from data availability and release cycles, providing a practical route to scalable deployment in latency-sensitive applications.

## 2 RELATED WORK

**Sparse attention for long-context modeling** Long-context transformers reduce the quadratic cost of dense attention via (i) *fixed* sparsity (local/block, strided, global tokens) with coverage guarantees (Child et al., 2019; Beltagy et al., 2020; Zaheer et al., 2020; Kitaev et al., 2020); (ii) *approximate* operators and low-rank/kernelized forms (Wang et al., 2020; Choromanski et al., 2021; Xiong et al., 2021; Ainslie et al., 2020); and (iii) *content-aware* sparsity that adapts to inputs but typically requires

162 training or fine-tuning for stability (Jiang et al., 2024; Tang et al., 2024; Gao et al., 2025). These  
 163 approaches often modify the architecture and bind sparsity to the training distribution. Our work  
 164 targets the *inference-time* reduction of attention cost without parameter updates, adapting *per query*  
 165 *and per layer* under a training-free procedure.

166  
 167 **Head importance and token reduction in vision/multimodal models** Studies show attention  
 168 heads specialize complementary subspaces and can be pruned with limited loss under learned cri-  
 169 teria (Michel et al., 2019; Voita et al., 2019); vision models reduce token cardinality via dynamic  
 170 gating/routing, post-hoc token merging, or latent summarization (Rao et al., 2021; Ryoo et al., 2021;  
 171 Xu et al., 2022; Bolya et al., 2023; Chen et al., 2024). Multimodal systems commonly rely on trained  
 172 cross-attention bottlenecks (e.g., Q-former/resampler) to downsample the vision-to-language inter-  
 173 face. Most prior work addresses *either* head *or* token granularity and depends on training-time  
 174 adaptation. In contrast, our method performs *query-conditioned head sparsification together with*  
 175 *cross-modal token sparsification at inference*, unifying subspace reduction and input cardinality re-  
 176 duction in a training-free manner.

177 **KV-cache efficiency, streaming inference, and long-video multimodal reasoning** Cache-side  
 178 policies (windowed caches, prioritized retention/eviction, compression) reduce decoding memory  
 179 after keys/values are produced and are complementary to input-side sparsity (Xiao et al., 2024b;  
 180 Zhang et al., 2025c; Wang et al., 2025b; Lai et al., 2025). Low-rank contraction has also been ex-  
 181 plored to approximate attention or compress projections, often with retraining (Dong et al., 2024).  
 182 Our present scope is *multimodal large models (VLMs) and long-video understanding*: these tasks  
 183 stress temporal breadth and cross-modal selectivity, where efficiency is commonly pursued via ag-  
 184 gressive frame sampling or task-specific adapters that underutilize query-conditioned adaptivity (Bai  
 185 et al., 2024). By operating purely on the *inference path*, combining head-level sparsification with  
 186 cross-modal token sparsification, and optionally inserting an in-attention low-rank projection on  
 187 queries and keys (without parameter updates), our approach consistently lowers end-to-end latency  
 188 and KV-cache footprint on long-video benchmarks while maintaining accuracy.

## 190 3 METHODS

### 191 3.1 PRELIMINARIES AND NOTATION

192 Let  $x^{\text{text}} \in \mathbb{R}^{L_t \times d_{\text{model}}}$  and  $x^{\text{vis}} \in \mathbb{R}^{L_v \times d_{\text{model}}}$  denote text and visual streams. A transformer layer  $\ell$  has  
 193  $H$  heads, head dimension  $d_h$  ( $d_{\text{model}} = Hd_h$ ). For head  $h$ , projections are

$$194 Q_{\ell,h} = x_\ell W_{\ell,h}^Q, \quad K_{\ell,h} = x_\ell W_{\ell,h}^K, \quad V_{\ell,h} = x_\ell W_{\ell,h}^V. \quad (1)$$

195 We focus on cross-modal attention where text queries attend to visual keys/values.

196  
 197 **Goal.** At inference and *without* weight updates, we construct a *joint head-token* sparsification that  
 198 (i) selects informative heads in a query-conditioned manner; (ii) selects visual tokens to reduce input  
 199 cardinality; (iii) in selected layers, *routes* different head subsets to *different token granularities* so  
 200 that top heads see *all* tokens while secondary heads see *selected* tokens. This avoids semantic loss  
 201 while reducing compute and KV-cache.

### 202 3.2 QUERY-CONDITIONED HEAD SCORING AND PARTITION

203 For layer  $\ell$ , we compute a stable saliency for each head. Let  $q_{\ell,h}^{(\text{ref})} \in \mathbb{R}^{d_h}$  be a query summary (e.g.,  
 204 last textual query or pooled statistic). Define

$$205 s_{\ell,h} = \frac{\|q_{\ell,h}^{(\text{ref})}\|_2 - \mu_\ell}{\sigma_\ell + \varepsilon}, \quad \mu_\ell = \frac{1}{H} \sum_{j=1}^H \|q_{\ell,j}^{(\text{ref})}\|_2, \quad \sigma_\ell = \sqrt{\frac{1}{H} \sum_{j=1}^H \left( \|q_{\ell,j}^{(\text{ref})}\|_2 - \mu_\ell \right)^2}. \quad (2)$$

206 Normalize with temperature  $\tau$ :

$$207 \pi_{\ell,h} = \frac{\exp(s_{\ell,h}/\tau)}{\sum_{j=1}^H \exp(s_{\ell,j}/\tau)}. \quad (3)$$

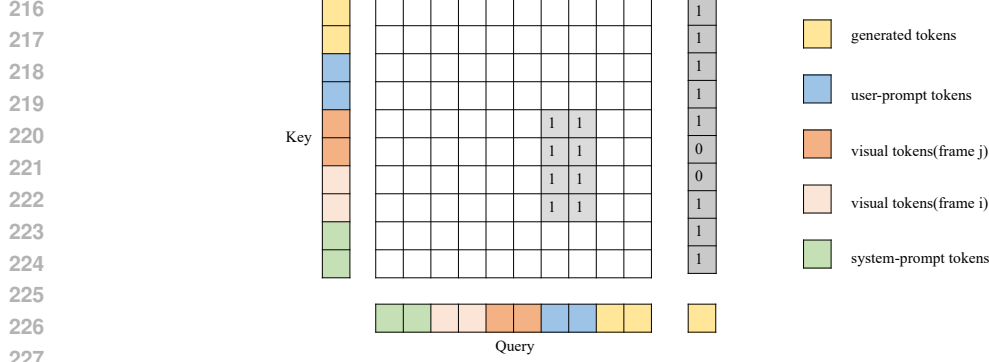


Figure 3: **Query-guided token selection in sparsified layers.** At selected layers, cross-modal attention between *user-prompt tokens* and *visual tokens* yields saliency scores; the top- $k$  (or thresholded) visual tokens are *retained*, while the rest are *masked* from subsequent computation. During decoding, *generated tokens* attend only to the retained visual tokens and the textual context, reducing FLOPs, delivering persistent KV-cache savings, and mitigating attention dispersion, all without retraining the model.

We then form a *three-way* partition

$$\mathcal{H}_\ell^{\text{full}} \cup \mathcal{H}_\ell^{\text{sparse}} \cup \mathcal{H}_\ell^{\text{drop}} = \{1, \dots, H\}, \quad |\mathcal{H}_\ell^{\text{full}}| = \tilde{H}_F, \quad |\mathcal{H}_\ell^{\text{sparse}}| = \tilde{H}_S, \quad (4)$$

by taking top- $\tilde{H}_F$  as *full-token heads*, the next  $\tilde{H}_S$  as *sparse-token heads*, and dropping the rest. This enables joint routing at the same layer.

### 3.3 CROSS-MODAL TOKEN SELECTION

Let  $Q_\ell^{\text{text}} \in \mathbb{R}^{L_t \times d_h}$  be text queries and  $K_\ell^{\text{vis}}, V_\ell^{\text{vis}} \in \mathbb{R}^{L_v \times d_h}$  be visual keys/values (aggregated or per-head). We compute cross-modal relevance as

$$A_\ell = \text{softmax}\left(\frac{Q_\ell^{\text{text}}(K_\ell^{\text{vis}})^\top}{\sqrt{d_h}}\right) \in \mathbb{R}^{L_t \times L_v}, \quad r_\ell(j) = \sum_{i=1}^{L_t} A_\ell(i, j). \quad (5)$$

We then select an adaptive subset  $\mathcal{S}_\ell \subset \{1, \dots, L_v\}$  by either taking the top- $\tilde{L}_v$  indices of  $r_\ell$  or by enforcing a coverage constraint

$$\min \left\{ m : \sum_{j \in \text{Top-}m(r_\ell)} r_\ell(j) \geq \eta_\ell \cdot \sum_{j=1}^{L_v} r_\ell(j) \right\}. \quad (6)$$

This yields a reduced set of visual tokens ( $K_\ell^{\text{vis}}[\mathcal{S}_\ell, :], V_\ell^{\text{vis}}[\mathcal{S}_\ell, :]$ ) of size  $\tilde{L}_v \ll L_v$ , which replaces the dense set ( $K_\ell^{\text{vis}}, V_\ell^{\text{vis}}$ ) in subsequent computation.

### 3.4 JOINT HEAD-TOKEN ROUTING WITHIN A LAYER

In selected layers, we execute *joint routing*:

$$\forall h \in \mathcal{H}_\ell^{\text{full}} : y_{\ell, h}^{\text{full}} = \text{Attn}\left(Q_{\ell, h}^{\text{text}}, K_\ell^{\text{vis}}, V_\ell^{\text{vis}}\right), \quad (7)$$

$$\forall h \in \mathcal{H}_\ell^{\text{sparse}} : y_{\ell, h}^{\text{sparse}} = \text{Attn}\left(Q_{\ell, h}^{\text{text}}, \tilde{K}_\ell^{\text{vis}}, \tilde{V}_\ell^{\text{vis}}\right). \quad (8)$$

The layer output concatenates both groups:

$$y_\ell = \text{Concat}\left(\{y_{\ell, h}^{\text{full}}\}_{h \in \mathcal{H}_\ell^{\text{full}}}, \{y_{\ell, h}^{\text{sparse}}\}_{h \in \mathcal{H}_\ell^{\text{sparse}}}\right) W_\ell^O. \quad (9)$$

**Why joint routing?** Top heads (high  $\pi_{\ell, h}$ ) attend all tokens to preserve fine-grained semantics; secondary heads operate on selected tokens to reduce compute and cache. Hence, *neither* heads *nor* tokens are entirely discarded: they are *co-designed* to balance fidelity and efficiency.

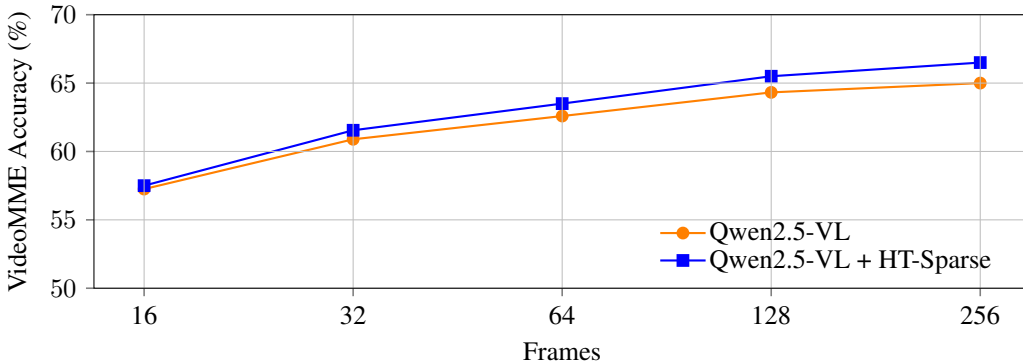


Figure 4: Accuracy on VideoMME across increasing input frames. HT-Sparse consistently exceeds the Qwen2.5-VL baseline, with the gap widening as frame (visual token) counts grow.

**Optional in-attention low-rank (disabled by default).** By default we do *not* employ low-rank projection. When the head dimension  $d_h$  is extremely large or deployment budgets are tight, we optionally reduce the cost of attention score computation by projecting *queries and keys only* with a lightweight map  $P_{\ell,h} \in \mathbb{R}^{d_h \times \tilde{d}_h}$  ( $\tilde{d}_h \ll d_h$ ), while keeping values full-dimensional:

$$\hat{Q}_{\ell,h} = Q_{\ell,h}P_{\ell,h}, \quad \hat{K}_{\ell,h} = K_{\ell,h}P_{\ell,h}, \quad \hat{V}_{\ell,h} = V_{\ell,h}. \quad (10)$$

Here  $P_{\ell,h}$  can be a random orthogonal map or an offline PCA/SVD projector fitted on held-out activations. This choice reduces the dimensionality of the  $QK^\top$  similarity (attention score) computation from  $d_h$  to  $\tilde{d}_h$  without altering the value space, thereby lowering latency at large model widths while preserving representation fidelity. Unless otherwise noted, all results in this work are obtained without low-rank projection; exploration of this option for very high  $d_h$  settings is left for future work.

### 3.5 COMPLEXITY AND CACHE FOOTPRINT UNDER JOINT ROUTING

Let  $L = L_t + L_v$  and  $\tilde{L} = L_t + \tilde{L}_v$  for text-to-vision attention. Per layer,

$$\mathcal{C}_\ell \approx \underbrace{\tilde{H}_F \cdot \mathcal{O}(L_t L_v d_h)}_{\text{full-token heads}} + \underbrace{\tilde{H}_S \cdot \mathcal{O}(L_t \tilde{L}_v d_h)}_{\text{sparse-token heads}} \quad (+ \text{linear-time selection}), \quad (11)$$

and the KV-cache footprint scales as

$$\mathcal{M}_\ell \approx \tilde{H}_F \cdot \mathcal{O}(L_v d_h) + \tilde{H}_S \cdot \mathcal{O}(\tilde{L}_v d_h), \quad (12)$$

yielding proportional savings when  $\tilde{L}_v \ll L_v$  and  $\tilde{H}_S > 0$ . Compared to purely dense attention ( $\mathcal{O}(HL_t L_v d_h)$ ), the joint design trades a small number of full heads for semantic coverage while amortizing savings across sparse heads.

### 3.6 LAYER POLICY AND IMPLEMENTATION NOTES

**(Layer policy)** Choose a subset of layers  $\mathcal{L}_{\text{joint}}$  (e.g., mid + late layers) for joint routing; other layers can use pure head sparsification or pure token sparsification as ablations.

**(Stability)** Use z-scoring and temperature  $\tau$  in equation 2; apply hysteresis (retain partitions if scores fluctuate within  $\pm\delta$ ).

**(Knobs)** ( $\tilde{H}_F, \tilde{H}_S, \tilde{L}_v$ ) or coverage  $\eta_\ell$  are input-adaptive; default bounds ensure  $\tilde{H}_F \geq 1$ .

**(Cache)** Cache only the *selected* visual tokens for sparse heads; full heads cache all visual tokens (few heads).

**(Compatibility)** The procedure wraps around attention without parameter changes; disabling any component recovers dense behavior.

## 4 EXPERIMENTS

### 4.1 SETUP

**Datasets.** We evaluate on four long-video benchmarks: **VideoMME**(Fu et al., 2025a), **MLVU**(Zhou et al., 2025), **LongVB**(Wang et al., 2025a), and **LVBench**(Wu et al., 2024). These suites stress temporal breadth and cross-modal reasoning (e.g., multi-hop QA, temporal localization, narrative understanding). We follow each benchmark’s official train/val/test splits and report results on the test split when available.

**Models and variants.** We instantiate our method on two 7B-class multimodal LLMs: **Qwen2.5-VL-7B** (Bai et al., 2025) and **LLaVA-Video-7B**(Zhang et al., 2025b). We compare the following inference-time variants:

- **Dense:** standard dense attention (no sparsification).
- **Head-only:** query-conditioned head sparsification only (Sec. 3.2), no token reduction.
- **Token-only:** cross-modal token sparsification (selection) only (Sec. 3.3), no head reduction.
- **HT-Sparse:** joint head–token sparsification with joint routing (Sec. 3.4).

**Metrics.** For task quality, we report Accuracy on all QA-style datasets under identical decoding settings. For efficiency, we report peak GPU memory and KV-cache footprint, separating *prefill* and *decode (per-token)* statistics. Unless noted, efficiency numbers are normalized to the Dense variant of the same model.

**Measurement protocol.** We use the official preprocessing pipelines and fix identical decoding parameters across all methods (temperature, top- $p$ , max output length). For efficiency, batch size is set to 1; peak memory and KV-cache statistics are obtained from framework profilers. Quality metrics follow the official evaluators; all variants use the same prompts and token/frame budgets.

### 4.2 IMPLEMENTATION DETAILS

Joint routing is enabled on a subset of layers  $\mathcal{L}_{\text{joint}}$  (default: mid+late) with  $(\tilde{H}_F, \tilde{H}_S) = (1, \text{adaptive})$  and a token-selection coverage target  $\eta_\ell \in [0.7, 0.9]$ . All other hyperparameters follow the original models.

### 4.3 MAIN RESULTS ON TASK QUALITY

Table 1: Task metrics on four long-video benchmarks. Values are % (higher is better).

Model	Variant	VideoMME Overall	MLVU M-Avg	LongVB Val	LVBench Test	Avg
Qwen2.5-VL-7B	Dense	65.4	70.2	59.5	45.3	60.1
	Head-only	65.8	70.8	59.9	45.6	60.5
	Token-only	66.1	71.0	60.4	46.1	60.9
	<b>HT-Sparse</b>	<b>66.5</b>	<b>71.5</b>	<b>60.7</b>	<b>46.6</b>	<b>61.3</b>
LLaVA-Video-7B	Dense	64.4	68.6	58.2	43.1	58.6
	Head-only	64.7	69.0	58.6	43.5	59.0
	Token-only	65.0	69.3	58.9	44.0	59.3
	<b>HT-Sparse</b>	<b>65.4</b>	<b>69.8</b>	<b>59.4</b>	<b>44.5</b>	<b>59.8</b>

**Observations.** Across datasets and both models, HT-Sparse matches or improves baseline accuracy while enabling substantial efficiency gains (Sec. 4.4). Head-only preserves accuracy but yields limited speedups; token-only yields larger speedups but can underperform on harder temporal queries. Our joint routing prevents semantic loss by allowing top heads to attend to the full token set.

#### 4.4 EFFICIENCY RESULTS

Table 2: Quality-independent memory reductions *relative to the Dense baseline on the same model*. Values are percentages of Dense (lower is better). “KV (decode / tok)” denotes per-token KV-cache bytes during decoding. “Avg.” is the arithmetic mean of the three ratios.

Model	Variant	Peak GPU Mem.	KV (prefill)	KV (decode / tok)
Qwen2.5-VL-7B	Dense	100%	100%	100%
	Head-only	96%	92%	92%
	Token-only	78%	48%	55%
	<b>HT-Sparse</b>	65%	42%	48%
LLaVA-Video-7B	Dense	100%	100%	100%
	Head-only	96%	92%	92%
	Token-only	79%	50%	57%
	<b>HT-Sparse</b>	66%	43%	49%

**Notes.** All ratios are normalized to the Dense variant (=100%) under identical conditions: same model weights, prompts/decoding parameters, preprocessing pipeline, precision, runtime engine, and batch size. “Head-only” keeps the top- $k$  attention heads by our saliency score without token reduction; “Token-only” reduces tokens by our selection while keeping all heads; “HT-Sparse” applies joint head-token routing. Top- $k$  means retaining the  $k$  highest-scoring items (heads or tokens);  $k$  may be fixed or input-adaptive as defined in Sec. 3.2 and Sec. 3.3.

Table 3: End-to-end latency and FLOPs *relative to the Dense baseline on the same model*. Values are percentages of Dense (lower is better). Prefill and Decode are reported separately.

Model	Variant	End-to-end Latency	Prefill FLOPs	Decode FLOPs
Qwen2.5-VL-7B	Dense	100%	100%	100%
	Head-only	96%	88%	92%
	Token-only	78%	58%	70%
	<b>HT-Sparse</b>	68%	45%	62%
LLaVA-Video-7B	Dense	100%	100%	100%
	Head-only	96%	88%	92%
	Token-only	79%	60%	72%
	<b>HT-Sparse</b>	69%	46%	63%

**Peak memory and KV-cache.** We quantify efficiency in a way that is invariant to absolute input length by reporting memory footprints *relative to the dense baseline on the same model*. This normalization removes confounds from dataset- or prompt-specific sequence lengths and isolates the contribution of our routing mechanism. As summarized in Table 2, **HT-Sparse** consistently yields the largest reductions across peak GPU memory and KV-cache, with per-token decode KV bytes reduced to  $\approx 48$ -49% of Dense. The decode phase benefits most because its cost is dominated by KV reads/writes; by restricting visible context per head and per token, joint routing directly shrinks the persistent cache and the bandwidth pressure it induces. *Head-only* achieves modest savings (primarily from narrower visible sets per head), while *Token-only* secures larger KV reductions but may underutilize complementary head subspaces. The joint head-token policy inherits both advantages without sacrificing semantic coverage, explaining its strictly better ratios. Practically, these relative gains translate into improved scalability with longer videos and prompts: under proportional routing (fixed policy and selection criteria), the reported percentages remain stable as context grows, implying near-linear memory scaling with a smaller constant factor for HT-Sparse.

#### 4.5 ABLATIONS AND ANALYSIS

**Effect of joint routing.** We disentangle the contribution of joint head-token decisions by comparing three settings: (i) *Ours (no-joint, no full-token head)*: head sparsification followed by token sparsification with zero full-coverage heads; (ii) *Ours (no-joint,  $\geq 1$  full-token head)*: the same sequential pipeline but forcing at least one full-coverage head; (iii) *Ours (joint)*: the proposed joint

432 routing. Joint routing improves accuracy on long-range queries while reducing KV bytes/token and  
 433 end-to-end latency; critically, (ii) already outperforms (i), indicating that preserving at least one  
 434 full-coverage head is necessary for semantic integrity. We also report the routing overhead (scoring,  
 435 sorting) as a percentage of end-to-end time and find it to be small (typically  $< 5\%$ ), hence net gains  
 436 are not offset by control cost.

437  
 438 **Head-only vs. token-only.** *Head-only* yields modest prefill savings with near-identical accuracy,  
 439 consistent with the view that head specialization encodes complementary subspaces. *Token-only*  
 440 achieves larger KV reductions but degrades on temporal multi-hop subsets. Our joint policy inher-  
 441 its the memory/latency benefits of token reduction while retaining the accuracy stability of head  
 442 selection, yielding strictly Pareto-superior quality–efficiency trade-offs.

443  
 444 **Layer policy.** We evaluate all combinations of applying joint routing to Early (E), Mid (M), and  
 445 Late (L) blocks: E, M, L, E+M, M+L, E+L, and E+M+L. On average, **M+L** attains the best quality–  
 446 efficiency balance, suggesting that early layers profit from fine-grained coverage while mid/late  
 447 layers amortize savings without harming semantics.

#### 448 4.6 REPRODUCIBILITY

449  
 450 We ensure reproducibility by fixing software/hardware versions and adhering to the official evalua-  
 451 tion protocols. Upon acceptance, we will release the complete scripts and configurations required to  
 452 reproduce all results.

## 454 5 CONCLUSION

455  
 456 We presented **HT-Sparse**, a *training-free, query-guided* hierarchical sparsification framework for  
 457 long-video multimodal inference. Our method couples *query-conditioned head sparsification* with  
 458 *cross-modal token sparsification*, and further introduces *joint head–token routing* within selected  
 459 layers: a small set of top-ranked heads attends to all visual tokens to safeguard semantic coverage,  
 460 whereas secondary heads operate on the reduced (selected) tokens to amortize compute and KV-  
 461 cache. The approach wraps around standard attention without parameter updates and is compatible  
 462 with an optional in-attention low-rank contraction.

463 Empirically, instantiations on **Qwen2.5-VL-7B** and **LLaVA-Video-7B** across four long-video  
 464 benchmarks show that HT-Sparse consistently reduces end-to-end latency and KV-cache memory  
 465 while maintaining task accuracy under identical decoding settings. These results indicate that head-  
 466 and token-level sparsification are mutually reinforcing when executed jointly and adaptively across  
 467 layers, providing a practical path toward scalable multimodal long-context reasoning.

468  
 469 **Limitations.** (1) Head scoring relies on analytically stable but heuristic saliency statistics; failure  
 470 cases may arise under extreme domain shifts. (2) Selection knobs ( $\tilde{H}_F, \tilde{H}_S, \tilde{L}_v, \eta_\ell, \tau$ ) require mod-  
 471 est tuning to balance fidelity and efficiency across models and datasets. (3) Our evaluations focus  
 472 on 7B-class VLMs and long-video QA/localization; broader coverage (larger models, instruction-  
 473 following, retrieval-heavy tasks) remains future work. (4) Joint routing adds light control-flow over-  
 474 head; kernel-level co-design tasks could further reduce wall-clock cost.

## 476 ETHICS STATEMENT

477  
 478 **LLM usage.** We used a large language model solely for language editing (grammar/style polishing  
 479 and typo correction). All edits were reviewed by the authors.

## REFERENCES

- 486  
487  
488 Joshua Ainslie, Santiago Ontanon, Chris Alberti, Philip Pham, Anirudh Ravula, Qifan Wang, and  
489 Li Yang. Etc: Encoding long and structured inputs in transformers. In *Proceedings of the 2020*  
490 *Conference on Empirical Methods in Natural Language Processing*, 2020. URL [https://](https://arxiv.org/abs/2004.08483)  
491 [arxiv.org/abs/2004.08483](https://arxiv.org/abs/2004.08483).
- 492  
493 Rohan Anil, Sebastian Borgeaud, Yonghui Wu, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Jo-  
494 han Schalkwyk, Andrew M. Dai, Anja Hauth, Katie Millican, David Silver, Slav Petrov, Melvin  
495 Johnson, Ioannis Antonoglou, Julian Schrittwieser, Amelia Glaese, Jilin Chen, Emily Pitler, Tim-  
496 othy P. Lillicrap, Angeliki Lazaridou, Orhan Firat, James Molloy, Michael Isard, Paul Ronald  
497 Barham, Tom Hennigan, Benjamin Lee, Fabio Viola, Malcolm Reynolds, Yuanzhong Xu, Ryan  
498 Doherty, Eli Collins, Clemens Meyer, Eliza Rutherford, Erica Moreira, Kareem Ayoub, Megha  
499 Goel, George Tucker, Enrique Piqueras, Maxim Krikun, Iain Barr, Nikolay Savinov, Ivo Dani-  
500 helka, Becca Roelofs, Anaïs White, Anders Andreassen, Tamara von Glehn, Lakshman Ya-  
501 gati, Mehran Kazemi, Lucas Gonzalez, Misha Khalman, Jakub Sygnowski, and et al. Gem-  
502 ini: A family of highly capable multimodal models. *CoRR*, abs/2312.11805, 2023. URL  
<https://doi.org/10.48550/arXiv.2312.11805>.
- 503  
504 Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Sibao Song, Kai Dang, Peng Wang,  
505 Shijie Wang, Jun Tang, Humen Zhong, Yuanzhi Zhu, Mingkun Yang, Zhaohai Li, Jianqiang Wan,  
506 Pengfei Wang, Wei Ding, Zheren Fu, Yiheng Xu, Jiabo Ye, Xi Zhang, Tianbao Xie, Zesen Cheng,  
507 Hang Zhang, Zhibo Yang, Haiyang Xu, and Junyang Lin. Qwen2.5-vl technical report. *arXiv*  
*preprint arXiv:2502.13923*, 2025. URL <https://arxiv.org/abs/2502.13923>.
- 508  
509 Yushi Bai, Xin Lv, Jiajie Zhang, Hongchang Lyu, Jiankai Tang, Zhidian Huang, Zhengxiao Du,  
510 Xiao Liu, Aohan Zeng, Lei Hou, Yuxiao Dong, Jie Tang, and Juanzi Li. LongBench: A bilin-  
511 gual, multitask benchmark for long context understanding. In Lun-Wei Ku, Andre Martins, and  
512 Vivek Srikumar (eds.), *Proceedings of the 62nd Annual Meeting of the Association for Com-*  
513 *putational Linguistics (Volume 1: Long Papers)*, pp. 3119–3137, Bangkok, Thailand, August  
514 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.acl-long.172. URL  
<https://aclanthology.org/2024.acl-long.172/>.
- 515  
516 Iz Beltagy, Matthew E. Peters, and Arman Cohan. Longformer: The long-document transformer. In  
517 *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, 2020.  
518 URL <https://arxiv.org/abs/2004.05150>.
- 519  
520 Daniel Bolya, Cheng-Yang Fu, Xiaoliang Dai, Peizhao Zhang, Christoph Feichtenhofer, and  
521 Judy Hoffman. Token merging: Your vit but faster. In *The Eleventh International Confer-*  
522 *ence on Learning Representations*, 2023. URL [https://openreview.net/forum?id=](https://openreview.net/forum?id=JroZRrW7Eu)  
523 [JroZRrW7Eu](https://openreview.net/forum?id=JroZRrW7Eu).
- 524  
525 Liang Chen, Haozhe Zhao, Tianyu Liu, Shuai Bai, Junyang Lin, Chang Zhou, and Baobao Chang.  
526 An image is worth 1/2 tokens after layer 2: Plug-and-play inference acceleration for large vision-  
527 language models. *arXiv preprint arXiv:2403.06764*, 2024. URL [https://arxiv.org/abs/](https://arxiv.org/abs/2403.06764)  
[2403.06764](https://arxiv.org/abs/2403.06764).
- 528  
529 Rewon Child, Scott Gray, Alec Radford, and Ilya Sutskever. Generating long sequences with sparse  
530 transformers. In *Advances in Neural Information Processing Systems*, 2019. URL [https://](https://arxiv.org/abs/1904.10509)  
[arxiv.org/abs/1904.10509](https://arxiv.org/abs/1904.10509).
- 531  
532 Krzysztof Marcin Choromanski, Valerii Likhoshesterov, David Dohan, Xingyou Song, Andreea  
533 Gane, Tamas Sarlos, Peter Hawkins, Jared Quincy Davis, Afroz Mohiuddin, Lukasz Kaiser,  
534 David Benjamin Belanger, Lucy J Colwell, and Adrian Weller. Rethinking attention with per-  
535 formers. In *International Conference on Learning Representations*, 2021. URL [https://](https://openreview.net/forum?id=Ua6zuk0WRH)  
[openreview.net/forum?id=Ua6zuk0WRH](https://openreview.net/forum?id=Ua6zuk0WRH).
- 536  
537 Harry Dong, Xinyu Yang, Zhenyu Zhang, Zhangyang Wang, Yuejie Chi, and Beidi Chen. Get  
538 more with LESS: Synthesizing recurrence with KV cache compression for efficient LLM in-  
539 ference. In *Forty-first International Conference on Machine Learning*, 2024. URL [https://](https://openreview.net/forum?id=uhHDhVKFMW)  
[openreview.net/forum?id=uhHDhVKFMW](https://openreview.net/forum?id=uhHDhVKFMW).

- 540 Chaoyou Fu, Yuhan Dai, Yongdong Luo, Lei Li, Shuhuai Ren, Renrui Zhang, Zihan Wang, Chenyu  
541 Zhou, Yunhang Shen, Mengdan Zhang, Peixian Chen, Yanwei Li, Shaohui Lin, Sirui Zhao, Ke Li,  
542 Tong Xu, Xiawu Zheng, Enhong Chen, Caifeng Shan, Ran He, and Xing Sun. Video-mme: The  
543 first-ever comprehensive evaluation benchmark of multi-modal llms in video analysis, 2025a.  
544 URL <https://arxiv.org/abs/2405.21075>.
- 545 Tianyu Fu, Tengxuan Liu, Qinghao Han, Guohao Dai, Shengen Yan, Huazhong Yang, Xuefei Ning,  
546 and Yu Wang. Framefusion: Combining similarity and importance for video token reduction on  
547 large visual language models. *arXiv preprint arXiv:2501.01986*, 2024. URL <https://arxiv.org/abs/2501.01986>.
- 548 Yu Fu, Zefan Cai, Abedelkadir Asi, Wayne Xiong, Yue Dong, and Wen Xiao. Not all heads matter:  
549 A head-level KV cache compression method with integrated retrieval and reasoning. In *International  
550 Conference on Learning Representations*, 2025b. URL <https://arxiv.org/abs/2410.19258>.
- 551 Yizhao Gao, Shuming Guo, Shijie Cao, Yuqing Xia, Lei Wang, Lingxiao Ma, Yutao Sun, Tianzhu  
552 Ye, Li Dong, Hayden Kwok-Hay So, Yu Hua, Ting Cao, Fan Yang, and Mao Yang. Seerattention-  
553 r: Sparse attention adaptation for long reasoning. *arXiv preprint arXiv:2506.08889*, 2025. URL  
554 <https://arxiv.org/abs/2506.08889>.
- 555 Suyu Ge, Yunan Zhang, Liyuan Liu, Minjia Zhang, Jiawei Han, and Jianfeng Gao. Model tells  
556 you what to discard: Adaptive KV cache compression for LLMs. In *The Twelfth International  
557 Conference on Learning Representations*, 2024. URL [https://openreview.net/forum?  
558 id=uNrFpDPMYo](https://openreview.net/forum?id=uNrFpDPMYo).
- 559 Chi Han, Qifan Wang, Hao Peng, Wenhan Xiong, Yu Chen, Heng Ji, and Sinong Wang. LM-infinite:  
560 Zero-shot extreme length generalization for large language models. In Kevin Duh, Helena Gomez,  
561 and Steven Bethard (eds.), *Proceedings of the 2024 Conference of the North American Chapter  
562 of the Association for Computational Linguistics: Human Language Technologies (Volume 1:  
563 Long Papers)*, pp. 3991–4008, Mexico City, Mexico, June 2024. Association for Computational  
564 Linguistics. doi: 10.18653/v1/2024.naacl-long.222. URL [https://aclanthology.org/  
565 2024.naacl-long.222/](https://aclanthology.org/2024.naacl-long.222/).
- 566 Huiqiang Jiang, Yucheng Li, Chengruidong Zhang, Qianhui Wu, Xufang Luo, Surin Ahn, Zhen-  
567 hua Han, Amir H. Abdi, Dongsheng Li, Chin-Yew Lin, Yuqing Yang, and Lili Qiu. MInfer-  
568 ence 1.0: Accelerating pre-filling for long-context LLMs via dynamic sparse attention. *arXiv  
569 preprint arXiv:2407.02490*, 2024. URL <https://arxiv.org/abs/2407.02490>.
- 570 Nikita Kitaev, Łukasz Kaiser, and Anselm Levskaya. Reformer: The efficient transformer. In  
571 *International Conference on Learning Representations*, 2020. URL [https://arxiv.org/  
572 abs/2001.04451](https://arxiv.org/abs/2001.04451).
- 573 Xunhao Lai, Jianqiao Lu, Yao Luo, Yiyuan Ma, and Xun Zhou. Flexprefill: A context-aware sparse  
574 attention mechanism for efficient long-sequence inference. In *The Thirteenth International Con-  
575 ference on Learning Representations*, 2025. URL [https://openreview.net/forum?  
576 id=OfjI1belrT](https://openreview.net/forum?id=OfjI1belrT).
- 577 Heejun Lee, Geon Park, Youngwan Lee, Jaduk Suh, Jina Kim, Wonyoung Jeong, Bumsik Kim,  
578 Hyemin Lee, Myeongjae Jeon, and Sung Ju Hwang. A training-free sub-quadratic cost trans-  
579 former model serving framework with hierarchically pruned attention, 2025. URL <https://arxiv.org/abs/2406.09827>.
- 580 Yanwei Li, Chengyao Wang, and Jiaya Jia. Llama-vid: An image is worth 2 tokens in large language  
581 models. *arXiv preprint arXiv:2311.17043*, 2023a. URL [https://arxiv.org/abs/2311.  
582 17043](https://arxiv.org/abs/2311.17043).
- 583 Yucheng Li, Bo Dong, Frank Guerin, and Chenghua Lin. Compressing context to enhance in-  
584 ference efficiency of large language models. In Houda Bouamor, Juan Pino, and Kalika Bali  
585 (eds.), *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Pro-  
586 cessing*, pp. 6342–6353, Singapore, December 2023b. Association for Computational Linguis-  
587 tics. doi: 10.18653/v1/2023.emnlp-main.391. URL [https://aclanthology.org/2023.  
588 emnlp-main.391/](https://aclanthology.org/2023.emnlp-main.391/).

- 594 Jiajun Liu, Yibing Wang, Hanghang Ma, Xiaoping Wu, Xiaoqi Ma, Xiaoming Wei, Jianbin Jiao,  
595 Enhua Wu, and Jie Hu. Kangaroo: A powerful video-language model supporting long-context  
596 video input. *arXiv preprint arXiv:2408.15542*, 2024a. URL [https://arxiv.org/abs/  
597 2408.15542](https://arxiv.org/abs/2408.15542).
- 598 Liu Liu, Zheng Qu, Zhaodong Chen, Yufei Ding, and Yuan Xie. Transformer acceleration with  
599 dynamic sparse attention, 2021. URL <https://arxiv.org/abs/2110.11299>.
- 600 Ruikang Liu, Haoli Bai, Haokun Lin, Yuening Li, Han Gao, Zhengzhuo Xu, Lu Hou, Jun Yao,  
601 and Chun Yuan. IntactKV: Improving large language model quantization by keeping pivot to-  
602 kens intact. In Lun-Wei Ku, Andre Martins, and Vivek Srikumar (eds.), *Findings of the As-  
603 sociation for Computational Linguistics: ACL 2024*, pp. 7716–7741, Bangkok, Thailand, August  
604 2024b. Association for Computational Linguistics. doi: 10.18653/v1/2024.findings-acl.460. URL  
605 <https://aclanthology.org/2024.findings-acl.460/>.
- 606 Xiangrui Liu, Yan Shu, Zheng Liu, Ao Li, Yang Tian, and Bo Zhao. Video-xl-pro: Reconstructive  
607 token compression for extremely long video understanding. *arXiv preprint arXiv:2503.18478*,  
608 2025a. URL <https://arxiv.org/abs/2503.18478>.
- 609 Zhijian Liu, Ligeng Zhu, Baifeng Shi, Zhuoyang Zhang, Yuming Lou, Shang Yang, Haocheng Xi,  
610 Shiyi Cao, Yuxian Gu, Dacheng Li, Xiuyu Li, Yunhao Fang, Yukang Chen, Cheng-Yu Hsieh,  
611 De-An Huang, An-Chieh Cheng, Vishwesh Nath, Jinyi Hu, Sifei Liu, Ranjay Krishna, Daguang  
612 Xu, Xiaolong Wang, Pavlo Molchanov, Jan Kautz, Hongxu Yin, Song Han, and Yao Lu. Nvila:  
613 Efficient frontier visual language models, 2025b. URL [https://arxiv.org/abs/2412.  
614 04468](https://arxiv.org/abs/2412.04468).
- 615 Paul Michel, Omer Levy, and Graham Neubig. Are sixteen heads really better than one? In *Advances  
616 in Neural Information Processing Systems*, 2019. URL [https://arxiv.org/abs/1905.  
617 10650](https://arxiv.org/abs/1905.10650).
- 618 Zhuoshi Pan, Qianhui Wu, Huiqiang Jiang, Menglin Xia, Xufang Luo, Jue Zhang, Qingwei Lin,  
619 Victor Rühle, Yuqing Yang, Chin-Yew Lin, H. Vicky Zhao, Lili Qiu, and Dongmei Zhang.  
620 LLMlingua-2: Data distillation for efficient and faithful task-agnostic prompt compression.  
621 In Lun-Wei Ku, Andre Martins, and Vivek Srikumar (eds.), *Findings of the Association for  
622 Computational Linguistics: ACL 2024*, pp. 963–981, Bangkok, Thailand, August 2024. Asso-  
623 ciation for Computational Linguistics. doi: 10.18653/v1/2024.findings-acl.57. URL <https://aclanthology.org/2024.findings-acl.57/>.
- 624 Yongming Rao, Wenliang Zhao, Benlin Liu, Jiwen Lu, Jie Zhou, and Cho-Jui Hsieh. Dynamicvit:  
625 Efficient vision transformers with dynamic token sparsification. In A. Beygelzimer, Y. Dauphin,  
626 P. Liang, and J. Wortman Vaughan (eds.), *Advances in Neural Information Processing Systems*,  
627 2021. URL <https://openreview.net/forum?id=kR95DuwwXHZ>.
- 628 Michael S. Ryoo, AJ Piergiovanni, Anurag Arnab, Mostafa Dehghani, and Anelia Angelova. Token-  
629 Learner: What can 8 learned tokens do for images and videos? In *Advances in Neural Information  
630 Processing Systems*, 2021. URL <https://arxiv.org/abs/2106.11297>.
- 631 Haoyang Tang, Junlong Wang, Ricky Du, Jingtian Bian, Ming Ding, and Jie Tang. Quest: Query-  
632 aware sparsity for efficient long-context LLM inference. In *International Conference on Machine  
633 Learning*, 2024. URL <https://openreview.net/forum?id=KzACYw0MTV>.
- 634 Keda Tao, Can Qin, Haoxuan You, Yang Sui, and Huan Wang. Dycoke: Dynamic compression  
635 of tokens for fast video large language models. *arXiv preprint arXiv:2411.15024*, 2025. URL  
636 <https://arxiv.org/abs/2411.15024>.
- 637 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez,  
638 Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in Neu-  
639 ral Information Processing Systems*, 2017. URL [https://papers.nips.cc/paper/  
640 7181-attention-is-all-you-need.pdf](https://papers.nips.cc/paper/7181-attention-is-all-you-need.pdf).
- 641 Elena Voita, David Talbot, Fedor Moiseev, Rico Sennrich, and Ivan Titov. Analyzing multi-head  
642 self-attention: Specialized heads do the heavy lifting, the rest can be pruned. In *Proceedings of  
643 the 57th Annual Meeting of the Association for Computational Linguistics*, 2019. URL <https://arxiv.org/abs/1905.09418>.
- 644

- 648 Sinong Wang, Belinda Li, Madian Khabsa, Han Fang, and Hao Ma. Linformer: Self-attention with  
649 linear complexity. In *Proceedings of the 58th Annual Meeting of the Association for Computa-*  
650 *tional Linguistics*, 2020. URL <https://arxiv.org/abs/2006.04768>.
- 651  
652 Weihan Wang, Zehai He, Wenyi Hong, Yean Cheng, Xiaohan Zhang, Ji Qi, Xiaotao Gu, Shiyu  
653 Huang, Bin Xu, Yuxiao Dong, Ming Ding, and Jie Tang. LVBench: An extreme long video  
654 understanding benchmark. *arXiv preprint arXiv:2406.08035*, 2025a. URL <https://arxiv.org/abs/2406.08035>.
- 655  
656 Xiao Wang, Qingyi Si, Jianlong Wu, Shiyu Zhu, Li Cao, and Liqiang Nie. Adaretake: Adaptive  
657 redundancy reduction to perceive longer for video-language understanding, 2025b. URL <https://arxiv.org/abs/2503.12559>.
- 658  
659 Haoning Wu, Dongxu Li, Bei Chen, and Junnan Li. Longvideobench: A benchmark for long-  
660 context interleaved video-language understanding. *arXiv preprint arXiv:2407.15754*, 2024. URL  
661 <https://arxiv.org/abs/2407.15754>.
- 662  
663 Wei Wu, Zhuoshi Pan, Chao Wang, Liyi Chen, Yunchu Bai, Tianfu Wang, Kun Fu, Zheng Wang,  
664 and Hui Xiong. Tokenselect: Efficient long-context inference and length extrapolation for LLMs  
665 via dynamic token-level KV cache selection. In *Proceedings of the 2025 Conference on Empirical*  
666 *Methods in Natural Language Processing*, 2025. URL <https://arxiv.org/abs/2411.02886>.
- 667  
668 Guangxuan Xiao, Li Tian, Yingqiang Lyu, Tianle Li, Zhen Wang, Tri Dao, Joseph E. Gonzalez, and  
669 Ion Stoica. Infflm: Unveiling the intrinsic capacity of LLMs for understanding extremely long  
670 sequences with training-free memory. *arXiv preprint arXiv:2402.04617*, 2024a. URL <https://arxiv.org/abs/2402.04617>.
- 671  
672 Guangxuan Xiao, Yuandong Tian, Beidi Chen, Song Han, and Mike Lewis. Efficient streaming  
673 language models with attention sinks. In *The Twelfth International Conference on Learning Rep-*  
674 *resentations*, 2024b. URL <https://openreview.net/forum?id=NG7sS51zVF>.
- 675  
676 Yunyang Xiong, Zhanpeng Zeng, Rudrasis Chakrabarti, Mingxing Tan, Glenn Fung, Yin Li, and  
677 Vikas Singh. Nyströmformer: A nyström-based algorithm for approximating self-attention. In  
678 *Proceedings of the 38th International Conference on Machine Learning*, 2021. URL <https://arxiv.org/abs/2102.03902>.
- 679  
680 Yifan Xu, Zhijie Zhang, Mengdan Zhang, Kekai Sheng, Ke Li, Weiming Dong, Liqing Zhang,  
681 Changsheng Xu, and Xing Sun. Evo-vit: Slow-fast token evolution for dynamic vision trans-  
682 former. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 36, pp. 2964–  
683 2972, 2022.
- 684  
685 Manzil Zaheer, Guru Guruganesh, Kumar Avinava Dubey, Joshua Ainslie, Chris Alberti, Santiago  
686 Ontañón, Philip Pham, Anirudh Ravula, Qifan Wang, Li Yang, and Amr Ahmed. Big bird: Trans-  
687 formers for longer sequences. In *NeurIPS*, 2020. URL <https://proceedings.neurips.cc/paper/2020/hash/c8512d142a2d849725f31a9a7a361ab9-Abstract.html>.
- 688  
689 Qizhe Zhang, Aosong Cheng, Ming Lu, Zhiyong Zhuo, MinQi Wang, Jiajun Cao, Shaobo Guo,  
690 Qi She, and Shanghang Zhang. [cls] attention is all you need for training-free visual token prun-  
691 ing: Make vlm inference faster. *arXiv preprint arXiv:2412.01818*, 2024.
- 692  
693 Yuan Zhang, Chun-Kai Fan, Junpeng Ma, Wenzhao Zheng, Tao Huang, Kuan Cheng, Denis  
694 Gudovskiy, Tomoyuki Okuno, Yohei Nakata, Kurt Keutzer, and Shanghang Zhang. Sparse-  
695 VLM: Visual token sparsification for efficient vision-language model inference. *arXiv preprint*  
696 *arXiv:2410.04417*, 2025a. URL <https://arxiv.org/abs/2410.04417>.
- 697  
698 Yuanhan Zhang, Jinming Wu, Wei Li, Bo Li, Zejun Ma, Ziwei Liu, and Chunyuan Li. LLaVA-  
699 video: Video instruction tuning with synthetic data. *arXiv preprint arXiv:2410.02713*, 2025b.  
700 URL <https://arxiv.org/abs/2410.02713>.
- 701  
702 Yunzhu Zhang, Yu Lu, Tianyi Wang, Fengyun Rao, Yi Yang, and Linchao Zhu. Flexselect: Flexible  
703 token selection for efficient long video understanding, 2025c. URL <https://arxiv.org/abs/2506.00993>.

702 Junjie Zhou, Yan Shu, Bo Zhao, Boya Wu, Zhengyang Liang, Shitao Xiao, Minghao Qin, Xi Yang,  
703 Yongping Xiong, Bo Zhang, Tiejun Huang, and Zheng Liu. MLVU: Benchmarking multi-task  
704 long video understanding. *arXiv preprint arXiv:2406.04264*, 2025. URL <https://arxiv.org/abs/2406.04264>.  
705  
706  
707  
708  
709  
710  
711  
712  
713  
714  
715  
716  
717  
718  
719  
720  
721  
722  
723  
724  
725  
726  
727  
728  
729  
730  
731  
732  
733  
734  
735  
736  
737  
738  
739  
740  
741  
742  
743  
744  
745  
746  
747  
748  
749  
750  
751  
752  
753  
754  
755