

LESS IS MORE: TRAINING-FREE SPARSE ATTENTION WITH GLOBAL LOCALITY FOR EFFICIENT REASONING

Lijie Yang^{*§†} Zhihao Zhang^{*†} Arti Jain[†] Shijie Cao[‡] Baihong Yuan[†] Yiwei Chen[†]

Zhihao Jia[†] Ravi Netravali[§]

[§]Princeton University [†]Carnegie Mellon University [‡]Microsoft Research
ly3223@princeton.edu, zhihaoz3@cs.cmu.edu, shijiecao@microsoft.com
zhihao@cmu.edu, rnetravali@cs.princeton.edu

ABSTRACT

Large reasoning models achieve strong performance through test-time scaling but incur substantial computational overhead, particularly from excessive token generation when processing short input prompts. While sparse attention mechanisms can reduce latency and memory usage, existing approaches suffer from significant accuracy degradation due to accumulated errors during long-generation reasoning. These methods generally require either high token retention rates or expensive retraining. We introduce LessIsMore, a training-free sparse attention mechanism for reasoning tasks, which leverages global attention patterns rather than relying on traditional head-specific local optimizations. LessIsMore aggregates token selections from local attention heads with recent contextual information, enabling unified cross-head token ranking for future decoding layers. This unified selection improves generalization and efficiency by avoiding the need to maintain separate token subsets per head. Evaluation across diverse reasoning tasks and benchmarks shows that LessIsMore preserves—and in some cases improves—accuracy while achieving a $1.1\times$ average decoding speed-up compared to full attention. Moreover, LessIsMore attends to $2\times$ fewer tokens without accuracy loss, achieving a $1.13\times$ end-to-end speed-up compared to existing sparse attention methods.¹

1 INTRODUCTION

Recent advancements in large reasoning models (LRMs) have significantly enhanced the capabilities of large language models (LLMs) for complex reasoning tasks. Models such as DeepSeek-R1 (DeepSeek-AI, 2025), Gemini-2.5-pro (DeepMind, 2025), OpenAI-o3 (OpenAI, 2025), Qwen3 (Team, 2025), and gpt-oss (OpenAI, 2025) leverage test-time scaling—generating large numbers of tokens—to enhance accuracy on challenging reasoning benchmarks (Wei et al., 2023; AoPS, 2025; Rein et al., 2023).

Unlike traditional language processing tasks, which involve long inputs and short outputs, reasoning tasks exhibit a different computational profile: they generate extensive multi-step derivations—often spanning tens of thousands of output tokens (Research, 2024)—from relatively concise problem statements. This decode-heavy nature incurs substantial computational overhead (Liu et al., 2025). For example, using full attention in the HuggingFace framework, DeepSeek-R1-Distill-Llama-8B consumes more than 20 minutes on a single NVIDIA RTX A5000 GPU to generate 32,768 tokens for one AIME problem.

This computational profile creates a unique optimization opportunity: while input processing benefits from full attention for accurate context understanding, the lengthy generation phase is well-suited for sparse attention mechanisms (Cai et al., 2025; Gao et al., 2025). Sparse attention selectively attends to a subset of critical tokens, offering a promising approach to reduce computational complexity and generation latency. Current techniques can be categorized into selection-based (Yang et al., 2024; Tang et al., 2024; Hao et al., 2025; Liu et al., 2024; Gao et al., 2025; Yuan et al., 2025) and

^{*}Equal contribution

¹Code available at <https://github.com/DerrickYLJ/LessIsMore>

eviction-based methods (Li et al., 2024; Xiao et al., 2023; Zhang et al., 2023; Adnan et al., 2024; Cai et al., 2025). Both identify important tokens using predefined criteria; eviction-based approaches permanently discard unselected tokens, while selection-based approaches maintain the full key-value (KV) cache.

However, existing sparse attention approaches suffer from significant accuracy degradation on reasoning tasks due to the accumulation of selection errors over long generation sequences (Gao et al., 2025). While standard generation tasks tolerate moderate information loss, step-by-step reasoning requires that crucial contextual information to be preserved throughout the entire process to maintain logical consistency (Lee & Hockenmaier, 2025). For instance, TidalDecode (Yang et al., 2024) achieves over 99.9% sparsity with no accuracy loss on retrieval task, but must reduce sparsity below 50% to preserve accuracy on AIME-24 reasoning tasks. In these settings, even small selection inaccuracies compound over thousands of generated tokens, leading to attention recall degradation and cascading accuracy drops. Moreover, prior work has shown that inaccurate sparse attention in reasoning models can also lengthen the generation process, which further harms the model’s inference efficiency (Gao et al., 2025).

These limitations motivated us to investigate the intrinsic attention distributions of reasoning models and tasks, with the goal of identifying patterns that enable more accurate and efficient token selection. Our token-level analysis across the reasoning process reveals two key observations on attention localities that fundamentally challenge the selection principles used in existing sparse attention methods.

First, reasoning tasks exhibit prominent *spatial locality* across attention heads, particularly within the Grouped Query Attention (GQA) frameworks prevalent in open-source LLMs (Touvron et al., 2023; AI, 2024; Team, 2025). Contrary to conventional wisdom that different attention heads perform specialized roles that require distinct token subsets (Yang et al., 2024; Xiao et al., 2024; Tang et al., 2024), we observe substantial overlap in token-importance rankings across heads within the same decoding layer. This overlap reveals that per-head top- k selection yields only a local optimum, overfitting to head-specific query patterns while potentially missing globally important tokens that could enhance performance in future decoding layers.

Second, we observe a *recency locality* pattern across decoding steps: tokens that receive high attention in one decoding step tend to continue attracting substantial attention over multiple subsequent steps. Notably, the ratio between the size of this “recency window” and the total number of selected tokens remains relatively constant throughout decoding, reflecting the intuition that each logical step in reasoning builds directly on the conclusions of preceding steps (Lee & Hockenmaier, 2025).

Building on these insights, this paper presents LessIsMore, a novel *training-free* sparse attention approach that achieves *higher* accuracy on reasoning tasks with *lower* latency by attending to *fewer* tokens. LessIsMore aggregates head-specific local information into a global attention pattern that is both more robust and more accurate. In each selection layer, LessIsMore exploits the identified locality patterns through a unified token selection process: each attention head first identifies its approximate top- k tokens using tailored selection schemes; these tokens are then aggregated across heads, globally ranked, and pruned to satisfy a predefined token budget. To capture recency locality, LessIsMore reserves a fixed proportion of this for a stable recency window, ensuring that recently generated tokens—critical for step-by-step reasoning—are consistently attended to.

Evaluation on Qwen3-8B and Qwen3-4B across diverse reasoning tasks, including AIME-24/25, GPQA-Diamond, and MATH500, demonstrates that LessIsMore consistently and significantly outperforms existing sparse attention baselines, including reasoning-focused methods that require retraining. LessIsMore maintains full accuracy at substantially higher sparsity levels, achieving up to 87.5% sparsity and $1.1\times$ end-to-end inference speed-up on AIME-24 with lossless accuracy, all without increasing reasoning length. These gains are further enabled by our customized kernel support for GQA models.

In summary, our contributions are:

- We present the first detailed, token-level analysis of attention distributions in the reasoning process, revealing fundamental *spatial* and *recency* locality patterns that challenge the conventional assumptions of highly specialized, independent attention heads.

- We propose LessIsMore, a training-free sparse attention mechanism featuring: (1) Unified Attention Head Selection globally aggregates head-level top- k selections, and (2) Stable Recency Window reserves recent contextual information for reasoning coherence.
- We show that LessIsMore matches or improves accuracy on challenging reasoning benchmarks while achieving a $1.1\times$ average decoding speed-up compared with the full attention baseline. Compared to state-of-the-art sparse attention methods, LessIsMore attends to at least $2\times$ fewer tokens, achieves a $1.13\times$ end-to-end speed-up, and shortens generation length by 7% without sacrificing accuracy.

2 OBSERVATION

Attention mechanisms are central to the functionality of today’s transformer-based LLMs. For each attention head i , attention scores and outputs are computed using the scaled-dot product of the query (Q_i), key (K_i), and value (V_i) tensors:

$$W_i = \frac{Q_i K_i^T}{\sqrt{d}}, \quad O_i = \text{softmax}(W_i) V_i \quad (1)$$

Here, W_i represents the attention weights (scores) between tokens, and O_i is the output from the i -th attention head.

Sparse attention methods address the computational overhead associated with attending to all tokens by selectively attending to a limited subset. Existing approaches aim to retain tokens most likely to yield high attention scores, constrained by a fixed token budget k , through different approximation functions (Tang et al., 2024; Cai et al., 2025; Yang et al., 2024):

$$\arg \max_{\rho} f(Q_i, K_i[\rho], V_i[\rho], k), \quad |\rho| = k \quad (2)$$

where ρ denotes the selected subset of tokens in the KV cache, and the approximation function f is usually an efficient estimation of the underlying ground-truth attention score to obtain ρ . The primary goal of this approximation is to maximize attention recall, defined as the proportion of the ground-truth attention scores that the selected subset covers:

$$R_i = \frac{\sum(\text{softmax}(W_i)[\rho])}{\sum(\text{softmax}(W_i))} \quad (3)$$

A higher attention recall indicates a better coverage of the attention mass with the selected token, thereby improving overall accuracy.

2.1 LIMITATIONS OF SPARSE ATTENTION IN REASONING

Despite the critical role of attention recall in sparse attention mechanisms, existing methods demonstrate significant limitations when applied to reasoning tasks. Current approaches either overestimate token importance (Tang et al., 2024; Xiao et al., 2023) or focus on locally optimal selections without adequately capturing global attention patterns across layers or decoding steps (Yang et al., 2024).

As illustrated in Figure 1a, both selection- and eviction-based sparse attention methods exhibit substantial attention recall degradation on the AIME-24 task, with degradation becoming more pronounced as decoding length increases. Specifically, selection-based approaches reach only 75% attention recall, while eviction-based approaches fall below 65%. Additionally, as shown in Figure 1b, although TidalDecode achieves comparable attention recall on both AIME-24 and simpler retrieval tasks (e.g., needle-in-the-haystack), the reasoning tasks inherently involve much longer generation sequences. Consequently, inaccuracies from token selection accumulate over prolonged generation, significantly degrading the reasoning quality in sparse attention. These observations underscore the necessity for designing a selection approach capable of capturing critical tokens globally.

2.2 LOCALITIES IN REASONING

Previous studies on traditional tasks have identified the locality property of attention patterns—different decoding layers can share a similar set of critical tokens (Yang et al., 2024). It

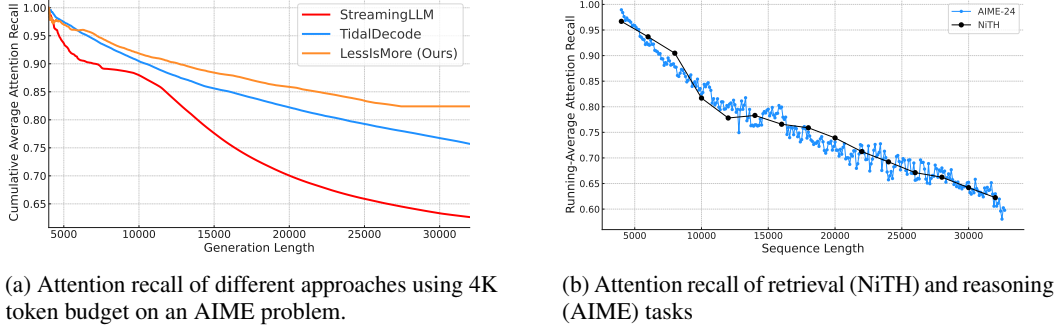


Figure 1: Analysis of attention recall degradation for sparse attention methods on reasoning tasks. Figure 1a compares cumulative average attention recall between eviction-based StreamingLLM (Xiao et al., 2023) and selection-based TidalDecode (Yang et al., 2024) on an AIME-24 reasoning task, using a token budget of 4K and generation length up to 32K tokens on Qwen3-8B. Figure 1b contrasts running-average attention recall of TidalDecode between the reasoning-intensive AIME-24 task and the simpler needle-in-the-haystack retrieval task under the same token budget on Qwen3-8B.

has also been suggested that different attention heads have distinct functional roles, thereby benefiting from specialized token subsets (Xiao et al., 2024). In Figure 2 and Figure 3, we analyze the distribution of the top-4K tokens across all 32 attention heads at different decoding steps using Qwen3-8B, a 36-layer model with Group Query Attention (GQA). Our analysis extends these findings by demonstrating that attention localities in reasoning tasks manifest both spatially (within the same key-value group and even across attention heads) and with recency (across decoding steps).

2.2.1 SPATIAL LOCALITY ACROSS ATTENTION HEADS

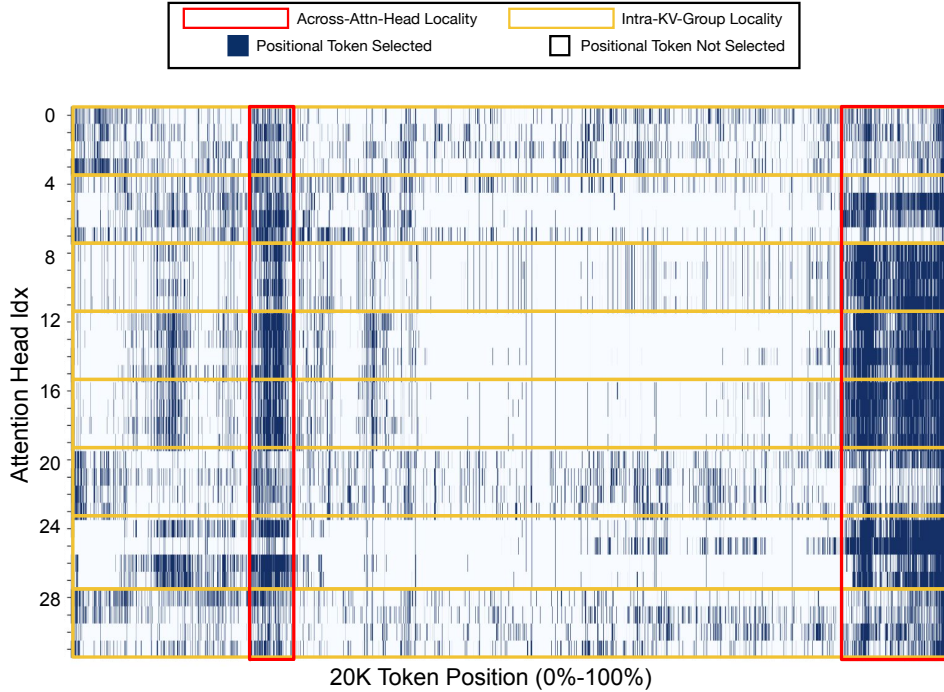


Figure 2: The distribution of the ground-truth top-4K tokens across all attention heads at 20K-th decoding step at Layer 4 on AIME-24 with Qwen3-8B. The dark blue positions stand for tokens included in ground-truth top-4K budget. We enclose the highly overlapped area of attention heads within the same kv group and across all heads with different colors

The visualization of Figure 2 highlights high overlap among selected tokens across consecutive groups of four attention heads within the same key-value group (yellow regions). Additionally, broader overlaps spanning all attention heads (red regions) include frequently attended recent tokens. This observation contradicts the common belief that each attention head serves specialized functions with distinct attention score distributions, thereby requiring head-wise token selection with different token subsets for optimal performance. Instead, our findings suggest that reasoning tasks exhibit remarkable consistency in token importance across attention heads, indicating that a unified global selection strategy may be more effective than maintaining separate token subsets per head in the reasoning process.

2.2.2 RECENCY LOCALITY OF RECENT TOKENS

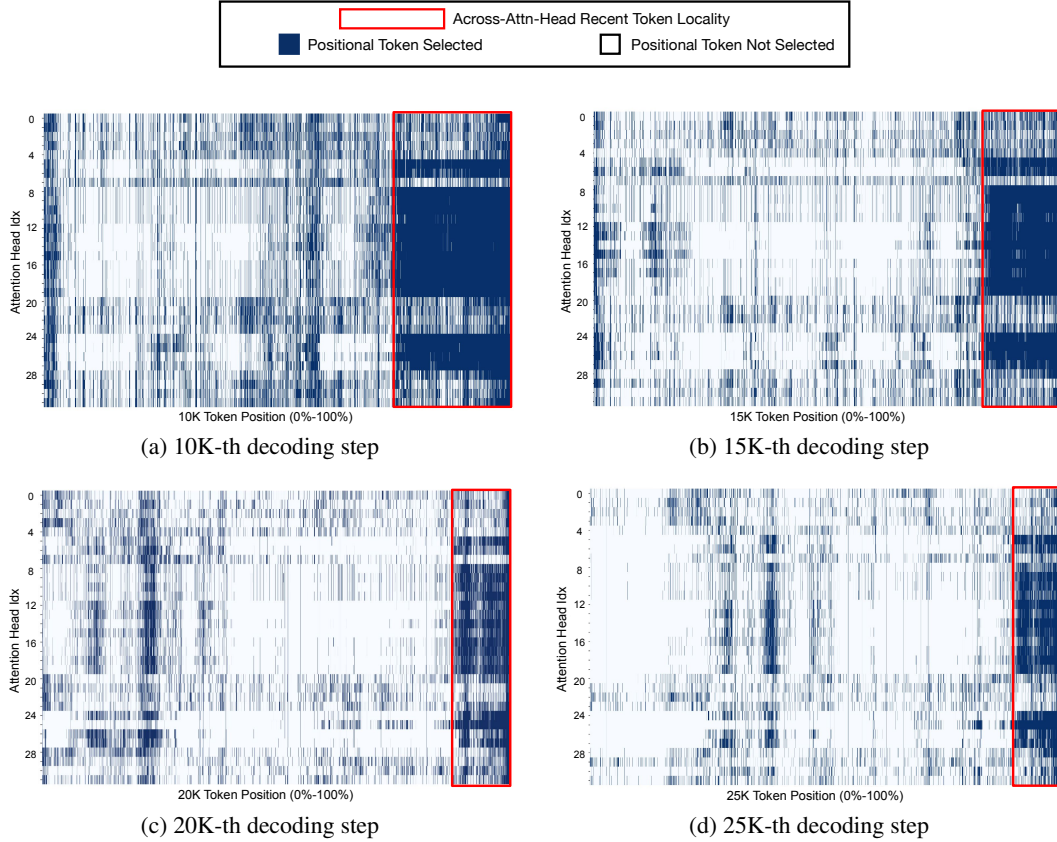


Figure 3: The distribution of the ground-truth top-4K tokens across all attention heads at 10K-, 15K-, 20K-, and 25K-th decoding step at Layer 4 on AIME-24 with Qwen3-8B. We enclose the highly overlapped area of attention heads within the same KV group with red, which forms a most recent window across all decoding steps

Figure 3 illustrates the distribution of top-4K tokens across multiple decoding lengths. As shown, the most recently generated tokens consistently receive high attention scores in subsequent steps. Furthermore, the size of this “recency window” remains relatively constant throughout the decoding process, as shown in all subfigures of Figure 3.

This recency locality directly reflects the nature of step-by-step reasoning, where each new logical step maintains coherence with immediately preceding conclusions (Lee & Hockenmaier, 2025). Prior work like StreamingLLM (Xiao et al., 2023) has recognized the importance of recent tokens by maintaining a fixed sliding window alongside attention sink tokens. Building on this, our analysis reveals that the ratio between critical recency window size and the number of selected tokens remains stable throughout the reasoning process. This observation supports the design of adaptive token

selection mechanisms that allocate a fixed proportional budget to recent tokens to maintain reasoning accuracy efficiently.

3 METHODOLOGY

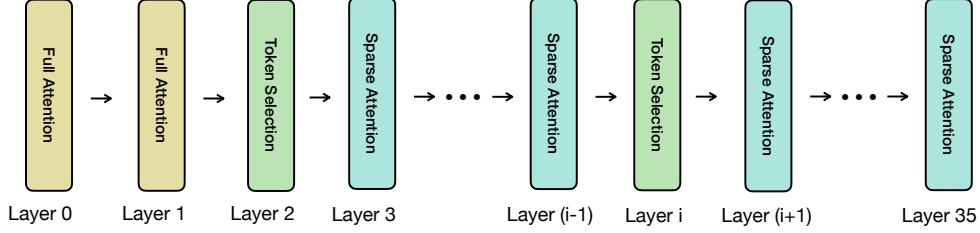


Figure 4: Walkthrough of a single decoding step of the selection-based approach TidalDecode (Yang et al., 2024), which performs full attention for the first two layers, full attention with token selection for the third layer and a middle layer, and sparse attention for the other layers.

This section introduces LessIsMore, an efficient and effective sparse attention system explicitly designed for reasoning tasks. LessIsMore exploits the unified global attention patterns identified through the analysis of spatial and recency localities presented in Section 2.2.1 and Section 2.2.2 by integrating two key techniques: Unified Attention Head Selection and Stable Recency Window. In this paper, LessIsMore adopts TidalDecode (Yang et al., 2024), one of the best-performing sparse attention methods, as a backbone; specifically, it starts with two full attention layers, includes two dedicated token selection layers, and employs sparse attention in the remaining layers. However, as discussed in Section 1 and shown in Section 5.2, LessIsMore’s underlying techniques can be incorporated into any sparse attention using approximation algorithms.

3.1 LESSISMORE

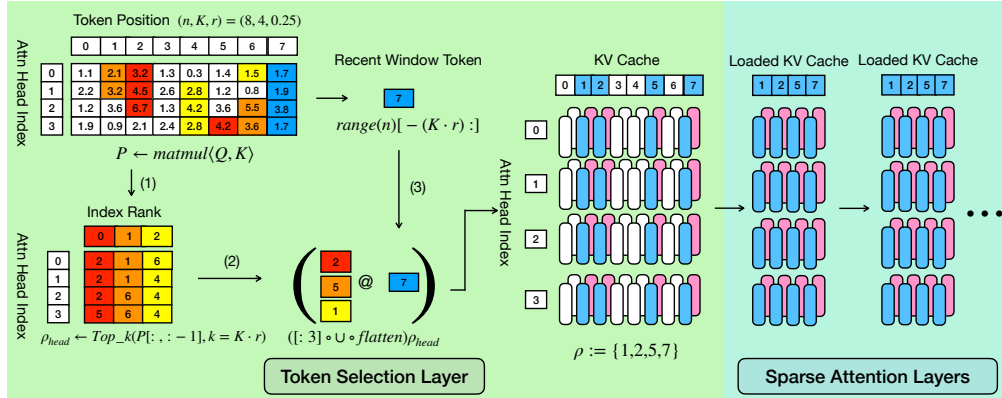


Figure 5: The selection process of LessIsMore is three-fold: (1) under token budget $K = 4$, $r = 0.25$, compute attention score matrix W and extract the top- k ($k = K \cdot (1 - r)$, r is the ratio of the most recent tokens that will be, by default, reserved as recency window) token indices for each attention head as ρ_{head} ; (2) flatten and union the selected indices for all heads, keeping the first k indices; (3) concatenate them with the most recent tokens, resulting in final token set ρ . The sparse attention layers only load the tensors of tokens in ρ from KV cache for all attention heads until the next selection layer or the end of the decoding step.

The detailed mechanism of token selection and sparse attention layers in LessIsMore is illustrated in Figure 5 and formalized in Algorithm 1. As shown in lines 4-19 of Algorithm 1, LessIsMore

Algorithm 1 LessIsMore Decoding Pipeline

```

1: Input: Current hidden state  $h$ , KV cache  $\mathcal{C}$ , token budget  $K$ , static ratio  $r$ 
2: Output: Logits
3: Initialize:  $\rho = []$  ▷ Token buffer for selected indices
4: for each decoder layer  $i$  do
5:    $q, k, v = f(W_{qkv}, h)$ 
6:    $\mathcal{C}.append(k, v)$ 
7:   if  $i$  is Full Attention Layer then
8:      $o = \text{FullAttention}(q, \mathcal{C}[:])$ 
9:   else if  $i$  is Token Selection Layer then
10:     $o = \text{FullAttention}(q, \mathcal{C}[:]), P := q \cdot \mathcal{C}.K^\top$  ▷ Full attention and Extract QK product
11:     $\rho_{\text{head}} = \text{TopKIndices}(P[:, :-(K \cdot r)], k = K \cdot (1 - r))$  ▷ Top-k for each head
12:     $\rho_{\text{unified}} = \text{UnionFlatten}(\rho_{\text{head}})$ 
13:     $\rho_{\text{recent}} = \text{Recent}(K \cdot r)$ 
14:     $\rho = \rho_{\text{unified}}[:, K \cdot (1 - r)] \cup \rho_{\text{recent}}$ 
15:   else
16:      $o = \text{SparseAttention}(q, \mathcal{C}[\rho])$  ▷ Use selected token indices
17:   end if
18:    $h = \text{FFN}(o)$ 
19: end for
20: return  $\text{lm\_head}(h)$ 

```

processes each decoding step with three distinct layer types. For full attention layers (lines 7-8), standard attention computation is performed. For token selection layers (lines 9-15), the overall token budget K is partitioned into two subsets: top-k tokens and the most recent tokens. During token selection, full attention is first used to calculate the QK-product $P = q \cdot \mathcal{C}.K^\top$ (line 10).

LessIsMore then applies Unified Attention Head Selection through the unified token selection process (lines 11-14), excluding the most recent tokens to focus the selection process on historical context. Each attention head independently selects its top-k token indices based on attention scores (line 11), followed by global aggregation and sorting across all heads (line 12). The top-ranked indices are then combined with the recent token indices determined by Stable Recency Window (lines 13-14) to form the final selected token set ρ . LessIsMore then fetches the corresponding key-value tensors from the KV cache for the selected indices. The indices are subsequently shared by all sparse attention layers (lines 16) until the next selection layer or the end of the current decoding step.

3.2 UNIFIED ATTENTION HEAD SELECTION

Unified Attention Head Selection aims to take advantage of the spatial locality observed with the token attention of Section 2.2.1 to improve efficiency and effectiveness. The core mechanism is implemented following lines 11-12 of Algorithm 1, where each attention head independently computes attention scores and identifies the top-k tokens most relevant to its query through $\text{TopKIndices}(P[:, :-(K \cdot r)], k = K \cdot (1 - r))$. Instead of maintaining separate sets of tokens per head, which increases the selection overhead and complexity of KV cache access, Unified Attention Head Selection aggregates the independently selected tokens from all attention heads.

The aggregation process, formalized as $\text{UnionFlatten}(\rho_{\text{head}})$ in line 12, flattens the top-k token indices selected by each head into a single unified set. This combined set is then globally sorted according to the tokens' rank within its attention head. Subsequently, only the globally highest-ranked tokens, limited by the predefined token budget $K \cdot (1 - r)$ (line 14), are retained. This unified selection strategy not only improves the attention recall shown in Section 4.4 by using the observed spatial location, where tokens frequently overlap in importance across heads, but also significantly simplifies token retrieval, enhancing computational efficiency during sparse attention computation.

3.3 STABLE RECENCY WINDOW

Stable Recency Window addresses the consistent recency locality observed in reasoning tasks, where recently generated tokens are repeatedly and consistently attended by all attention heads in Figure 3.

To exploit this pattern, Stable Recency Window dedicates a fixed ratio of the total token budget K exclusively to the most recently generated tokens, forming a “stable recency window.” This mechanism is implemented in lines 13-14 of Algorithm 1.

Prior sparse attention training approaches maintain a constant number of tokens as the sliding window size regardless of token budgets (Yuan et al., 2025). In contrast, the stable recency window in LessIsMore is determined by a predefined ratio r , typically a small fraction of K , through $\text{Recent}(K \cdot r)$ (line 13). The final set of tokens used for sparse attention computation is formed by the union operation $\rho = \rho_{\text{unified}}[: K \cdot (1 - r)] \cup \rho_{\text{recent}}$ (line 14), consisting of the globally unified top- k tokens selected through Unified Attention Head Selection and the most recent tokens determined by Stable Recency Window. This design directly reflects the empirical observation that recently generated tokens inherently possess critical contextual information necessary for accurate and coherent step-by-step reasoning. By explicitly allocating resources to recent tokens through this algorithmic approach, Stable Recency Window effectively ensures high attention recall and improved reasoning quality while maintaining computational efficiency, as shown in Section 4.4.

4 EXPERIMENTS

4.1 EXPERIMENT SETUP

We conduct extensive experiments to evaluate the accuracy and efficiency of LessIsMore. Our experiments consider two widely-used reasoning models, Qwen3-8B and Qwen3-4B (Team, 2025) backed up with GQA. Both models are specifically trained for reasoning tasks and perform the effective thinking process by generating extensive tokens. Further, we evaluate on multiple mainstream reasoning tasks, including AIME-24/25, GPQA-Diamond, and MATH500.

In Section 4.2, we compare LessIsMore with both training-free (TidalDecode, Quest) and training-required (SeerAttention-r) sparse attention baselines. For training-free approaches, we keep the first two full attention layers for TidalDecode. To ensure a fair comparison, we set the same selection layer for LessIsMore and TidalDecode - Layer 12 for Qwen3-8B and Layer 20 for Qwen3-4B.² The static ratio r of LessIsMore is set to 0.25 and 4 tokens are always reserved for attention sink in this section. We maintain the same experiment configuration of Quest and SeerAttention-r in (Gao et al., 2025), where the block size is set to 64 and all layers perform sparse attention. To guarantee the consistency of evaluation results, the maximum generation length of Qwen3-8B and 4B models are set to 32,768 across all reasoning tasks. For LessIsMore and TidalDecode, we calculate average pass@1 accuracy over 16 samples for AIME-24/25, 8 samples for MATH500 and GPQA-Diamond, with Full Attention, Quest, and SeerAttention-r results using the reported sampling configurations in (Gao et al., 2025).

In Section 4.3, we compare the per-token decoding latency of our system implemented using customized kernels with the selection-based model (TidalDecode) and full attention using FlashInfer (Ye et al., 2024). Finally, in Section 4.4, we show that LessIsMore’s global selection techniques generalize to other sparse attention approaches, analyze the impact of sparse attention on reasoning length, and examine the effect of the ratio of recent window size on LessIsMore’s attention recall throughout the generation in reasoning.

4.2 EVALUATION ON REASONING TASKS

Figure 6 presents the accuracy comparison among Full Attention baseline, LessIsMore, and other sparse attention methods across mainstream reasoning benchmarks, including the complex mathematical contests AIME-24/25, and the simpler reasoning tasks MATH500 and GPQA-Diamond, evaluated on the reasoning-focused language models Qwen3-8B and Qwen3-4B. For the challenging AIME tasks, experiments span token budgets of 2K, 4K, and 6K tokens, corresponding to overall generation lengths averaging between 14K and 20K tokens. In contrast, the MATH500 and GPQA-Diamond tasks involve shorter reasoning sequences averaging between 4K and 9K tokens. LessIsMore consistently achieves the highest accuracy across all evaluated tasks and token budgets, closely matching the performance of full attention, even under stringent token constraints. Specifically, for Qwen3-8B on the AIME-24 task at the smallest token budget (2K tokens), LessIsMore attains a nearly lossless

²the choice is justified in Appendix A.1

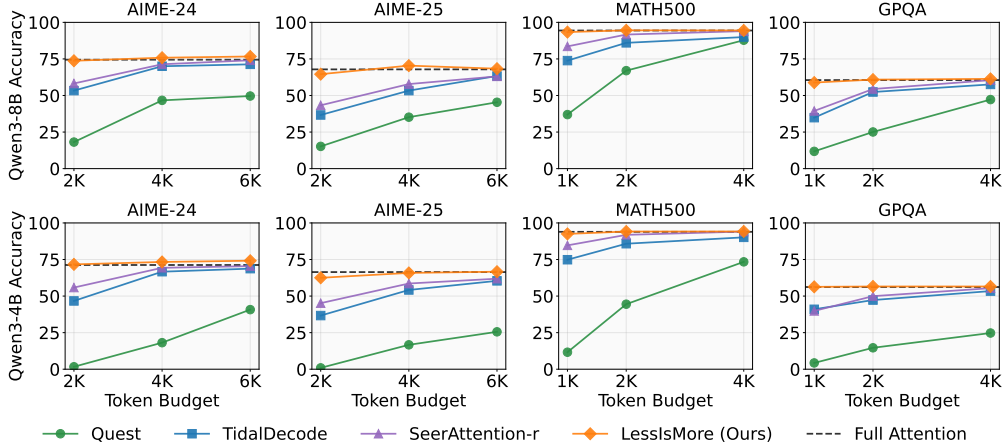


Figure 6: Accuracy results of LessIsMore (ours), Quest, TidalDecode, SeerAttention-r, and Full Attention across for multiple main-stream reasoning tasks. Across all evaluated tasks, LessIsMore consistently achieves the lossless accuracy with small token budgets (1K or 2K), always outperforming all other baselines.

performance, markedly surpassing training-free methods such as Quest and TidalDecode, as well as the training-required SeerAttention-r, all of which suffer notable accuracy degradation at reduced token budgets. A similar trend is observed across all tasks, underscoring the capability of LessIsMore to effectively retain critical contextual information and facilitate accurate reasoning even with limited computational resources.

4.3 EFFICIENCY EVALUATION

To evaluate the practical efficiency gains of LessIsMore, we implement customized kernels optimized for GQA-based models and conduct average decoding latency measurements on reasoning tasks. Our evaluation uses LLaMA-3.1-8B (AI, 2024) (as a reference to the evaluated Qwen3-8B with the same parameter size and attention mechanism GQA) deployed on a single NVIDIA RTX A5000 GPU, comparing the efficiency-accuracy trade-offs among LessIsMore, TidalDecode, and Full Attention within a unified system built on FlashInfer (Ye et al., 2024). The token index aggregation stage illustrated in Figure 5 is implemented using PyTorch primitives.

As demonstrated in Figure 7, even if TidalDecode is an efficient selection-based approach that only performs a two-time top-k selection for each decoding step, LessIsMore consistently outperforms TidalDecode in all token budget configurations, achieving substantially higher accuracy while maintaining faster average decoding speeds. Notably, LessIsMore delivers near-lossless performance (73.75 vs. 74.48 full attention baseline) using only a 2K token budget with a $1.10\times$ speed-up over full attention. In contrast, TidalDecode achieves merely 53.33 accuracy under identical constraints. In Table 1, even with a token budget of 6K, TidalDecode obtains a lower accuracy and generates 15.9K tokens. Meanwhile, LessIsMore achieves a $1.06\times$ average decoding speedup and 7% shorter generation length, which contributes to a $1.13\times$ end-to-end speedup compared to TidalDecode. This demonstrates LessIsMore’s superior ability to maintain reasoning quality while delivering meaningful computational savings.

4.4 ABLATION STUDY

4.4.1 EFFECTIVENESS OF LESSISMORE’S AGGREGATION ON GQA

Applying sparse attention to GQA-based models often requires aggregating tokens from the query heads to the KV group (Yuan et al., 2025). To demonstrate the generalizability of our unified selection in LessIsMore beyond the TidalDecode pipeline, we evaluate different token aggregation schemes on Qwen3-8B with GQA, where each KV head is shared across multiple query (or attention)

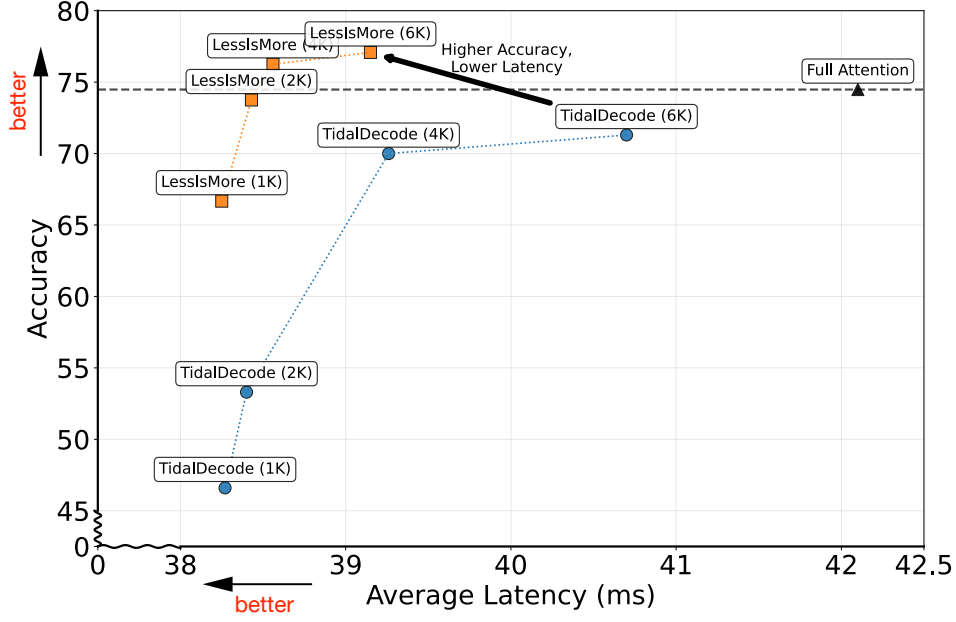


Figure 7: Efficiency-accuracy tradeoff comparison on AIME-24 using LLama-3.1-8B. Each point represents the average decoding latency across the corresponding average generation length in Table 1. LessIsMore (orange squares) consistently achieves higher accuracy than TidalDecode (blue circles) while maintaining lower latency across all token budgets (1K, 2K, 4K, 6K). The closer to the top-left corner, the better the method performs. Full Attention baseline (triangle) provides the accuracy upper bound but with higher computational cost.

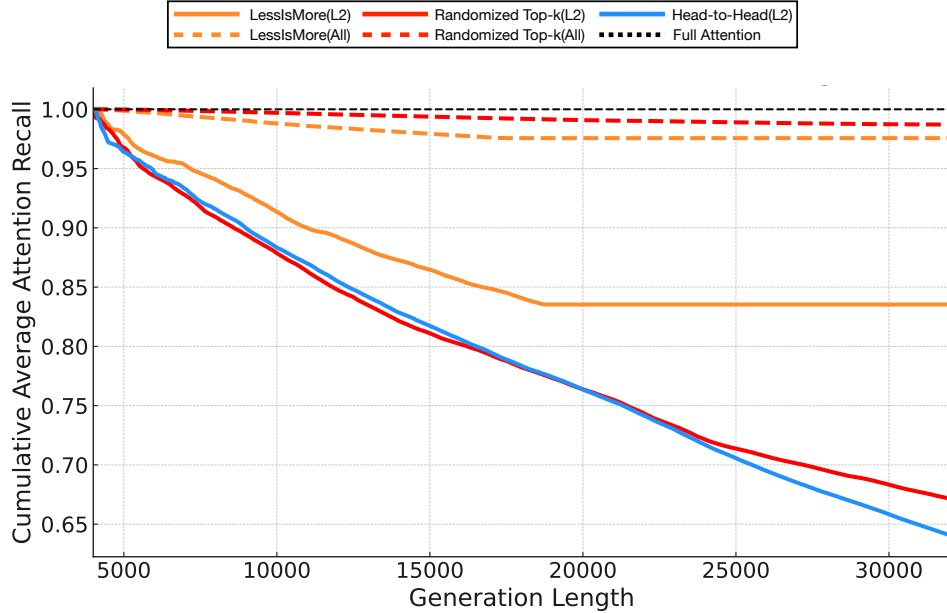


Figure 8: The Top-4K attention recall of different selection schemes applied only on Layer 2(L2) or all decoding layers(All). (1) **LessIsMore**: our unified top-k selection across all attention heads with 25% tokens for recency window, (2) **Randomized Top-k**: random application of one query head's top-k tokens to the entire KV group, and (3) **Head-to-Head**: direct utilization of top-k tokens for each individual attention head

heads. We compare three selection strategies LessIsMore, Randomized Top-k, and Head-to-Head in Figure 8.

The results reveal a critical distinction between local and global optimization strategies. When selection is performed on all decoding layers (All), locally optimal schemes like Randomized Top-k can achieve competitive performance by overfitting to layer-specific patterns. However, when selection frequency is reduced to only Layer 2 (L2)—a more general scenario that reduces computational overhead—our unified approach significantly outperforms alternative aggregation schemes.

This performance gap demonstrates that locally optimal selection methods, while effective for immediate layer optimization, fail to generalize robustly across future decoding layers. The superior attention recall of LessIsMore under sparse selection conditions indicates that our global aggregation strategy, combined with the stable recency window, provides more robust token importance estimation that generalizes effectively across the reasoning process. This finding validates the core principle that unified global selection, rather than head-specific local optimization, is essential for maintaining high attention recall in computation-constrained reasoning scenarios.

4.4.2 GENERATION LENGTH ANALYSIS UNDER SPARSE ATTENTION

Table 1: The AIME-24 accuracy followed by corresponding average reasoning length (in K) of different approaches on Qwen3-8B. The highest accuracy and the lowest generation length of each column are in bold, excluding the Full Attention row.

| Model (Task) | Method / Budget | K=2000 | K=4000 | K=6000 |
|------------------------|--------------------------|---------------------|---------------------|---------------------|
| Qwen-3-8B (AIME-24) | Quest | 18.15 (30.0) | 46.67 (22.9) | 49.63 (19.6) |
| | TidalDecode | 53.33 (17.4) | 70.00 (16.9) | 71.30 (15.9) |
| | SeerAttention-r | 58.23 (19.8) | 71.35 (16.3) | 74.06 (15.3) |
| | LessIsMore (Ours) | 73.75 (15.8) | 75.83 (14.8) | 76.67 (15.1) |
| | Full Attention | 74.48 (14.8) | 74.48 (14.8) | 74.48 (14.8) |

Sparse attention methods exhibit a concerning tendency that extends generation lengths on reasoning tasks, as demonstrated in Table 1 and corroborated by prior research (Gao et al., 2025). This phenomenon reflects the accumulation of selection errors discussed in Section 1, where imprecise token retention forces models into inefficient reasoning patterns that compromise both accuracy and computational efficiency.

Table 1 presents the average generation lengths of different approaches under various token budgets on AIME-24 using Qwen3-8B. Under restrictive token budgets (K=2000), existing methods generate substantially longer sequences compared to full attention: Quest, SeerAttention-r and TidalDecode each generate 30.0K, 19.8K, and 17.4K tokens, representing 103%, 34%, and 18% increases respectively over the full attention baseline of 14.8K tokens. These extended sequences indicate that sparse attention errors accumulate over time and may force models to engage in a redundant reasoning process. In contrast, LessIsMore maintains generation lengths closely aligned with full attention across all token budgets. At K=4000, LessIsMore generates the same number of tokens as full attention does while achieving better accuracy. Meanwhile, even with a token budget of 6K, TidalDecode obtains a significant lower accuracy and generates 15.9K tokens. Combining with the average decoding latency in Figure 7, LessIsMore achieves a $1.13\times$ end-to-end speedup compared to TidalDecode.

Since inaccurate token selection leads to extended generation lengths, attention recall serves as an indicator of both selection accuracy and computational efficiency. Therefore, evaluating attention recall dynamics throughout generation becomes more crucial for assessing sparse attention methods on reasoning tasks.

4.4.3 EFFECT OF RECENT WINDOW RATIO

To better understand the design choices in LessIsMore, we evaluate how varying the ratio of recent window size impacts attention recall and answer correctness in Figure 9. We run the AIME-24 reasoning task under a fixed 4K token budget and record the cumulative attention recall across

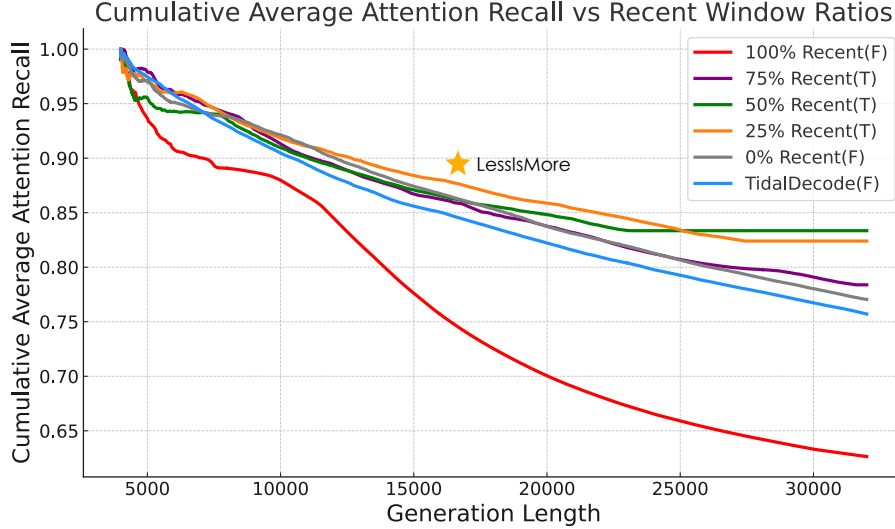


Figure 9: Ablation study on the impact of varying the recent window ratio in LessIsMore (★) on the AIME-24 reasoning task, using a token budget of 4K and generation length up to 32K tokens on Qwen3-8B. LessIsMore corresponds to the 25% recent setting combined with Unified Attention Head Selection, (labeled with (★)). We compare it against alternative recent window ratios, the 100% recent baseline (i.e., using only recent tokens), and TidalDecode. Curves are annotated with “(T)” or “(F)” to indicate whether the configuration yields the correct answer. Notably, only configurations that incorporate recent window with Unified Attention Head Selection (25%, 50%, 75%) succeed in solving the task.

different configurations. Using only recent tokens (100% recent) provides the lowest attention recall, as it discards distant but important contextual tokens. TidalDecode, which selectively retains tokens without explicitly accounting for reasoning-specific attention locality, significantly improves attention recall but still fails to produce the correct answer. Building upon TidalDecode, simply using Unified Attention Head Selection variant (0% recent) further improves attention recall by leveraging our selection scheme, yet it also fails to generate the correct answer. Incorporating a proportion of recent tokens consistently boosts attention recall. Specifically, configurations with 25%, 50%, and 75% recent windows all manage to generate the correct answer. Among them, 25% recent—which corresponds to the full design of LessIsMore—achieves the highest attention recall throughout the most generation process. This validates the design choice of allocating 25% of the token budget to the recent window in LessIsMore.

5 FUTURE WORK AND LIMITATIONS

5.1 ADAPTIVE SHORT-CONTEXT RATIO

While it is optimal for sparse attention mechanisms to dynamically determine the best token budget and parameter ratios at runtime, this remains challenging in practice (Yang et al., 2024; Gao et al., 2025; Tang et al., 2024; Cai et al., 2025). LessIsMore is no exception to this limitation. Currently, we employ a pre-defined token budget and a fixed 25% short-context ratio based on our observation of the consistent size of important recent context across reasoning tasks.

However, it would be optimal to determine this ratio adaptively during generation. One effective solution that has been explored by prior works is through top-p sampling-inspired approaches, where the ratio could be dynamically adjusted based on the attention score distribution at each decoding step (Chen et al., 2024; Lin et al., 2025).

5.2 GENERALIZATION OF LESSISMORE

Currently, LessIsMore is based upon the pipeline used by TidalDecode. Evaluations on GQA-based models Qwen3-4B and Qwen3-8B demonstrate superior reasoning performance compared to existing sparse attention approaches. More than just a model-specific optimization, LessIsMore offers a principled approach to designing more accurate and efficient sparse attention for reasoning tasks. Given the fundamental locality observations presented in Section 2.2, future work should prioritize extending LessIsMore’s locality-based principles to models beyond GQA, such as MLA or MoE architectures that employ different attention mechanisms.

5.3 FROM TOKEN SAVING TO MEMORY SAVING

System-level optimizations in LessIsMore remain unexplored that could translate token savings into substantial memory savings. Currently, several implementation optimizations limit the full potential of LessIsMore: (1) The union operation in selection is implemented using PyTorch primitives rather than optimized CUDA kernels. (2) The token selection process still requires maintaining the full KV cache in memory, limiting memory efficiency compared to eviction-based approaches.

5.4 BEYOND TRAINING-FREE: NATIVE TRAINING WITH HYBRID ATTENTION

The success of LessIsMore demonstrates that a training-free, hybrid attention approach which combines full attention and sparse attention can effectively address the efficiency challenges of reasoning models without sacrificing accuracy. This hybrid model architecture echoes the emerging trend in pretraining where the model architecture is moving towards combining different attention mechanisms [Dong et al.; OpenAI \(2025\)](#). Looking ahead, directly integrating the principles of LessIsMore into the pretraining process opens a compelling avenue for future research. In addition, LessIsMore employs an inter-layer hybrid approach, a more advanced step would be to explore intra-layer hybrid attention mechanisms.

6 RELATED WORK

Efficient sparse attention with KV cache compression. Sparse attention mechanisms reduce the computational overhead and memory requirements of attention computation by selectively attending to only a subset of tokens, significantly improving inference efficiency for long-sequence tasks ([Yang et al., 2024](#); [Tang et al., 2024](#)). Current approaches can be broadly categorized into two main paradigms: eviction-based and selection-based methods. Eviction-based approaches ([Xiao et al., 2023](#); [Zhang et al., 2023](#); [Li et al., 2024](#); [Adnan et al., 2024](#)) permanently discard tokens from the KV cache based on predefined criteria, achieving better memory savings by maintaining a smaller cache size throughout generation. In contrast, selection-based methods ([Yang et al., 2024](#); [Tang et al., 2024](#); [Hao et al., 2025](#); [Liu et al., 2024](#)) retain the full KV cache but dynamically choose which tokens to attend to during computation, typically optimizing for locally maximal attention scores and choosing different tokens for each attention head. While both approaches demonstrate effectiveness on standard long-context tasks such as retrieval and summarization, they face significant challenges when applied to reasoning tasks due to the accumulation of selection errors over extended generation sequences.

Sparse attention in reasoning. Recent reasoning models leverage the principle that scaling test-time compute can be more effective than scaling model parameters ([Wei et al., 2023](#); [DeepSeek-AI, 2025](#)), generating extensive token sequences to enhance reasoning accuracy through deliberative processes. However, the lengthened generation nature of reasoning tasks poses unique challenges for applying existing sparse attention methods. When applied on reasoning tasks, prior works either suffer from significant accuracy degradation when using small token retention ratios ([Yang et al., 2024](#); [Tang et al., 2024](#); [Cai et al., 2025](#)) or require computationally expensive post-training procedures to mitigate accuracy loss accumulated during generation ([Gao et al., 2025](#)), both of which also significantly increase the generation length of reasoning tasks. In contrast, LessIsMore is proposed as a selection-based, training-free approach that leverages intrinsic spatial and recency attention patterns within the reasoning process to achieve high accuracy with substantially reduced token utilization ratios but without extending the generation length.

7 CONCLUSION

In this paper, we propose LessIsMore, a novel training-free sparse attention mechanism specifically designed to address the inaccurate selection limitation of prior approaches on reasoning tasks. Our approach fundamentally challenges the conventional paradigm that each attention head should independently select distinct subsets of tokens optimized locally. Instead, by leveraging observed spatial and recency locality patterns, LessIsMore globally aggregates token importance across all heads, significantly enhancing token selection accuracy. Notably, LessIsMore achieves near-lossless accuracy using extremely low token budgets across diverse reasoning benchmarks. For instance, it maintains full accuracy on the challenging AIME-24 task with merely a 2K token budget, outperforming existing methods that suffer significant accuracy drops under similar constraints. Furthermore, unlike prior sparse attention methods that inherently extend reasoning length due to accumulated selection errors, LessIsMore retains efficiency with generation lengths comparable to full attention. Consequently, comprehensive evaluations demonstrate that LessIsMore can preserve (and even improve) accuracy on challenging reasoning benchmarks while having $1.1\times$ average decoding speed-ups compared with the full attention baseline; moreover, it preserves the accuracy while attending to at least $2\times$ fewer tokens and achieves a $1.13\times$ end-to-end speedup stemming from 7% shorter generation length and $1.06\times$ average decoding speedup compared to other state-of-the-art sparse attention methods. Our results underscore the effectiveness and promise of exploiting global attention patterns in sparse attention mechanisms tailored for reasoning-intensive tasks.

REFERENCES

- Muhammad Adnan, Akhil Arunkumar, Gaurav Jain, Prashant J. Nair, Ilya Soloveychik, and Purushotham Kamath. Keyformer: Kv cache reduction through key tokens selection for efficient generative inference, 2024. URL <https://arxiv.org/abs/2403.09054>.
- Meta AI. Llama 3.1: Advanced open-source language model. <https://ai.meta.com/blog/meta-llama-3-1/>, 2024. Accessed: 2024-09-26.
- AoPS. Aime problems and solutions. https://artofproblemsolving.com/wiki/index.php/AIME_Problems_and_Solutions, 2025. Accessed: 2025-07-14.
- Zefan Cai, Wen Xiao, Hanshi Sun, Cheng Luo, Yikai Zhang, Ke Wan, Yucheng Li, Yeyang Zhou, Li-Wen Chang, Jiuxiang Gu, Zhen Dong, Anima Anandkumar, Abedelkadir Asi, and Junjie Hu. R-kv: Redundancy-aware kv cache compression for reasoning models, 2025. URL <https://arxiv.org/abs/2505.24133>.
- Zhuoming Chen, Ranajoy Sadhukhan, Zihao Ye, Yang Zhou, Jianyu Zhang, Niklas Nolte, Yuandong Tian, Matthijs Douze, Leon Bottou, Zhihao Jia, and Beidi Chen. Magicpig: Lsh sampling for efficient llm generation, 2024. URL <https://arxiv.org/abs/2410.16179>.
- Google DeepMind. Gemini 2.5: Our most intelligent ai model. <https://blog.google/technology/google-deepmind/gemini-model-thinking-updates-march-2025/>, March 2025. Accessed: 2025-07-14.
- DeepSeek-AI. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning, 2025. URL <https://arxiv.org/abs/2501.12948>.
- Xin Dong, Yonggan Fu, Shizhe Diao, Wonmin Byeon, Zijia Chen, Ameya Sunil Mahabaleshwarkar, Shih-Yang Liu, Matthijs Van Keirsbilck, Min-Hung Chen, Yoshi Suhara, et al. Hymba: A hybrid-head architecture for small language models, 2024. URL <https://arxiv.org/abs/2411.13676>.
- Yizhao Gao, Shuming Guo, Shijie Cao, Yuqing Xia, Yu Cheng, Lei Wang, Lingxiao Ma, Yutao Sun, Tianzhu Ye, Li Dong, Hayden Kwok-Hay So, Yu Hua, Ting Cao, Fan Yang, and Mao Yang. Seerattention-r: Sparse attention adaptation for long reasoning, 2025. URL <https://arxiv.org/abs/2506.08889>.

- Jitai Hao, Yuke Zhu, Tian Wang, Jun Yu, Xin Xin, Bo Zheng, Zhaochun Ren, and Sheng Guo. OmniKV: Dynamic context selection for efficient long-context LLMs. In *The Thirteenth International Conference on Learning Representations*, 2025. URL <https://openreview.net/forum?id=ulCAPXYXfa>.
- Jinu Lee and Julia Hockenmaier. Evaluating step-by-step reasoning traces: A survey, 2025. URL <https://arxiv.org/abs/2502.12289>.
- Yuhong Li, Yingbing Huang, Bowen Yang, Bharat Venkitesh, Acyr Locatelli, Hanchen Ye, Tianle Cai, Patrick Lewis, and Deming Chen. Snapkv: Llm knows what you are looking for before generation, 2024. URL <https://arxiv.org/abs/2404.14469>.
- Chaofan Lin, Jiaming Tang, Shuo Yang, Hanshuo Wang, Tian Tang, Boyu Tian, Ion Stoica, Song Han, and Mingyu Gao. Twilight: Adaptive attention sparsity with hierarchical top- p pruning, 2025. URL <https://arxiv.org/abs/2502.02770>.
- Di Liu, Meng Chen, Baotong Lu, Huiqiang Jiang, Zhenhua Han, Qianxi Zhang, Qi Chen, Chengruidong Zhang, Bailu Ding, Kai Zhang, Chen Chen, Fan Yang, Yuqing Yang, and Lili Qiu. Retrievalattention: Accelerating long-context llm inference via vector retrieval, 2024. URL <https://arxiv.org/abs/2409.10516>.
- Yue Liu, Jiaying Wu, Yufei He, Hongcheng Gao, Hongyu Chen, Baolong Bi, Ruihan Gong, Jiaheng Zhang, Zhiqi Huang, and Bryan Hooi. Efficient inference for large reasoning models: A survey, 2025. URL <https://arxiv.org/abs/2503.23077>.
- OpenAI. Introducing gpt-oss, 2025. URL <https://openai.com/index/introducing-gpt-oss/>.
- OpenAI. Introducing openai o3 and o4-mini. <https://openai.com/index/introducing-o3-and-o4-mini/>, April 2025. Accessed: 2025-07-14.
- David Rein, Betty Li Hou, Asa Cooper Stickland, Jackson Petty, Richard Yuanzhe Pang, Julien Dirani, Julian Michael, and Samuel R. Bowman. Gpqa: A graduate-level google-proof q&a benchmark, 2023. URL <https://arxiv.org/abs/2311.12022>.
- Nous Research. Open reasoning tasks: Llm reasoning tasks collection, 2024. URL <https://github.com/NousResearch/Open-Reasoning-Tasks>.
- Jiaming Tang, Yilong Zhao, Kan Zhu, Guangxuan Xiao, Baris Kasikci, and Song Han. Quest: Query-aware sparsity for efficient long-context llm inference, 2024. URL <https://arxiv.org/abs/2406.10774>.
- Qwen Team. Qwen3 technical report, 2025. URL <https://arxiv.org/abs/2505.09388>.
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. Llama: Open and efficient foundation language models, 2023.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed Chi, Quoc Le, and Denny Zhou. Chain-of-thought prompting elicits reasoning in large language models, 2023. URL <https://arxiv.org/abs/2201.11903>.
- Guangxuan Xiao, Yuandong Tian, Beidi Chen, Song Han, and Mike Lewis. Efficient streaming language models with attention sinks. *arXiv*, 2023.
- Guangxuan Xiao, Jiaming Tang, Jingwei Zuo, Junxian Guo, Shang Yang, Haotian Tang, Yao Fu, and Song Han. Duoattention: Efficient long-context llm inference with retrieval and streaming heads, 2024. URL <https://arxiv.org/abs/2410.10819>.
- Lijie Yang, Zhihao Zhang, Zhuofu Chen, Zikun Li, and Zhihao Jia. Tidaldecode: Fast and accurate llm decoding with position persistent sparse attention, 2024. URL <https://arxiv.org/abs/2410.05076>.

Zihao Ye, Ruihang Lai, Roy Lu, Chien-Yu Lin, Size Zheng, Lequn Chen, Tianqi Chen, and Luis Ceze. Cascade inference: Memory bandwidth efficient shared prefix batch decoding. <https://flashinfer.ai/2024/01/08/cascade-inference.html>, Jan 2024. URL <https://flashinfer.ai/2024/01/08/cascade-inference.html>. Accessed on 2024-02-01.

Jingyang Yuan, Huazuo Gao, Damai Dai, Junyu Luo, Liang Zhao, Zhengyan Zhang, Zhenda Xie, Y. X. Wei, Lean Wang, Zhiping Xiao, Yuqing Wang, Chong Ruan, Ming Zhang, Wenfeng Liang, and Wangding Zeng. Native sparse attention: Hardware-aligned and natively trainable sparse attention, 2025. URL <https://arxiv.org/abs/2502.11089>.

Zhenyu Zhang, Ying Sheng, Tianyi Zhou, Tianlong Chen, Lianmin Zheng, Ruisi Cai, Zhao Song, Yuandong Tian, Christopher Ré, Clark Barrett, Zhangyang Wang, and Beidi Chen. H₂O: Heavy-hitter oracle for efficient generative inference of large language models, 2023.

A APPENDIX

A.1 CHOICE OF TIDALDECODE OPTIMAL RE-SELECTION ON QWEN3

Following the procedure of choosing optimal re-selection layer of TidalDecode (Yang et al., 2024), we conduct a simple 5K-context-length needle-in-the-haystack test on TidalDecode with Qwen3-8B and Qwen3-4B. With a token budget of 256, Layer 12 achieves the best accuracy of 86% on Qwen3-8B; Layer 12 and 20 achieve the best accuracies of 86% and 84% on Qwen3-4B, respectively. Moreover, prior work has found that in the same model family, the optimal re-selection layer is similar. For Qwen3, we validate that Layer 12 is an important layer. To demonstrate the generalization of our approach on different models, we choose different re-selection layers for different Qwen3 models. In this paper’s experiments Section 4, we apply the same re-selection layer index on TidalDecode and LessIsMore for a fair comparison - Layer 12 and Layer 20 for Qwen3-8B and Qwen3-4B, respectively.

A.2 REASONING EVALUATION RESULTS

Table 2: Results of 2K-, 4K-, and 6K-token-budget in Figure 6 evaluated on Qwen3-8B and Qwen3-4B for AIME24 and AIME25 benchmarks.

| Model (Task) | Method / Budget | K=2000 | K=4000 | K=6000 |
|------------------------|--------------------------|--------------|--------------|--------------|
| Qwen-3-8B (AIME-24) | Quest | 18.15 | 46.67 | 49.63 |
| | TidalDecode | 53.33 | 70.00 | 71.30 |
| | SeerAttention-r | 58.23 | 71.35 | 74.06 |
| | LessIsMore (Ours) | 73.75 | 75.83 | 76.67 |
| | Full Attention | 74.48 | 74.48 | 74.48 |
| Qwen-3-8B (AIME-25) | Quest | 15.2 | 35.18 | 45.37 |
| | TidalDecode | 36.67 | 53.33 | 63.33 |
| | SeerAttention-r | 43.30 | 57.81 | 63.07 |
| | LessIsMore (Ours) | 64.58 | 70.42 | 68.33 |
| | Full Attention | 67.86 | 67.86 | 67.86 |
| Qwen-3-4B (AIME-24) | Quest | 1.67 | 18.14 | 40.74 |
| | TidalDecode | 46.67 | 66.67 | 68.75 |
| | SeerAttention-r | 55.83 | 69.32 | 70.47 |
| | LessIsMore (Ours) | 71.67 | 73.33 | 74.17 |
| | Full Attention | 71.25 | 71.25 | 71.25 |
| Qwen-3-4B (AIME-25) | Quest | 0.92 | 16.67 | 25.56 |
| | TidalDecode | 38.12 | 54.17 | 60.41 |
| | SeerAttention-r | 45.16 | 58.59 | 61.88 |
| | LessIsMore (Ours) | 62.50 | 65.83 | 66.67 |
| | Full Attention | 66.41 | 66.41 | 66.41 |

Table 3: Results of 1K-, 2K-, and 4K-token-budget in Figure 6 evaluated on Qwen3-8B and Qwen3-4B for MATH500 and GPQA benchmarks.

| Model (Task) | Method / Budget | K=1000 | K=2000 | K=4000 |
|------------------------|--------------------------|--------------|--------------|--------------|
| Qwen-3-8B (MATH500) | Quest | 36.95 | 66.98 | 87.80 |
| | TidalDecode | 73.85 | 86.00 | 89.95 |
| | SeerAttention-r | 83.57 | 91.67 | 94.00 |
| | LessIsMore (Ours) | 93.35 | 94.55 | 94.45 |
| | Full Attention | 94.43 | 94.43 | 94.43 |
| Qwen-3-8B (GPQA) | Quest | 11.80 | 25.06 | 47.22 |
| | TidalDecode | 34.84 | 52.39 | 57.57 |
| | SeerAttention-r | 39.43 | 54.41 | 60.48 |
| | LessIsMore (Ours) | 58.84 | 60.86 | 61.36 |
| | Full Attention | 60.54 | 60.54 | 60.54 |
| Qwen-3-4B (MATH500) | Quest | 11.62 | 44.45 | 73.40 |
| | TidalDecode | 74.85 | 85.80 | 90.15 |
| | SeerAttention-r | 84.67 | 91.85 | 94.10 |
| | LessIsMore (Ours) | 92.50 | 94.12 | 94.16 |
| | Full Attention | 93.93 | 93.93 | 93.93 |
| Qwen-3-4B (GPQA) | Quest | 4.29 | 14.64 | 24.74 |
| | TidalDecode | 40.97 | 47.28 | 53.41 |
| | SeerAttention-r | 39.84 | 49.94 | 55.40 |
| | LessIsMore (Ours) | 56.31 | 56.56 | 56.56 |
| | Full Attention | 56.19 | 56.19 | 56.19 |