

SADAS: Enabling Efficient On-Device Reasoning via State-Aware Dynamic Attention Scheduling

Anonymous ACL submission

Abstract

The deployment of reasoning-intensive LLMs on edge devices is severely hampered by the memory bandwidth bottlenecks inherent in processing long Chain-of-Thought sequences. To address this, we propose State-Aware Dynamic Attention Scheduling (SADAS), a novel framework that aligns computational complexity with cognitive phases. SADAS leverages control tokens to autonomously toggle between bandwidth-efficient Sliding Window Attention for intermediate reasoning steps and high-fidelity Full Attention for final answer integration. This mechanism reduces memory pressure without sacrificing global context. Experimental results on commodity CPUs demonstrate that SADAS transforms unresponsive long-chain reasoning into a real-time capability, achieving up to $3.88\times$ end-to-end speedups. Crucially, SADAS-1.7B enjoys superior performance that matches or exceeds the results of leading 8B-scale reasoning models on challenging reasoning benchmarks like AIME24, while maintaining robust agentic instruction-following capabilities. To the best of our knowledge, SADAS is the first architecture to implement token-level dynamic attention mixing, establishing a new paradigm for responsive on-device intelligence.

1 Introduction

The paradigm of Natural Language Processing (NLP) is being profoundly reshaped by Large Language Models (LLMs). Currently, a major research direction clearly points towards endowing models with complex reasoning capabilities, a cutting-edge field exemplified by works such as OpenAI’s o1 (Jaech et al., 2024), Alibaba’s Qwen3 (Yang et al., 2025), and DeepSeek’s DeepSeek-R1 (Guo et al., 2025). The core mechanism of these models involves simulating and executing complex cognitive tasks by generating explicit Chain-of-Thought (CoT) sequences (Wei et al., 2022), leading to

breakthrough progress in various logic-intensive benchmarks. Crucially, this deep reasoning capability serves as the foundational engine for a wide range of downstream applications, from solving intricate mathematical problems to empowering autonomous agents on edge devices that must plan and execute complex user instructions (Xi et al., 2025).

However, this transition to generative reasoning fundamentally shifts the computational bottleneck. The demand for deep reasoning renders LLMs increasingly decoding-oriented. Before generating a final answer, models must produce extensive internal thought processes token-by-token. For a CoT of length N , the decoding process necessitates N forward passes, each involving self-attention with $O(N^2)$ complexity (Chou et al., 2024). While manageable on cloud GPUs, this quadratic cost becomes prohibitive when deploying reasoning models on resource-constrained edge devices (e.g., smartphones, laptops). On commodity CPUs with limited memory bandwidth, generating long reasoning chains results in unacceptable latency, making on-device reasoning capabilities practically unresponsive. Analyzing its output structure more deeply, this process functionally bifurcates the model’s generation task into two implicit stages: an exploratory intermediate thought stage, which focuses on rapid generation and path exploration; and an integrative final answer stage, which emphasizes reviewing and precisely consolidating global information. This raises a core architectural design question: Is a single homogeneous attention architecture optimal for these two functionally distinct computational stages?

The Linear Attention mechanism (Shen et al., 2021; Han et al., 2024; Ma et al., 2021) offers a highly promising direction to address this challenge, with its $O(N)$ computational complexity significantly superior to the $O(N^2)$ of standard attention mechanisms. However, linear attention

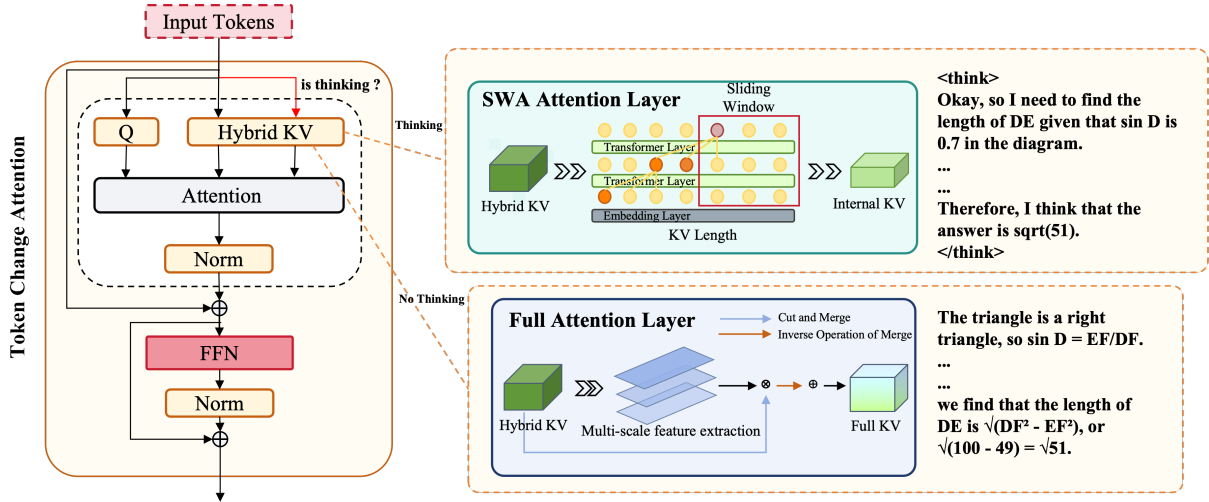


Figure 1: Architectural Overview of SADAS

width on edge devices, dominating decoding latency. Various mechanisms aim to alleviate this.

Linear Attention (e.g., RWKV (Peng et al., 2023), Mamba (Gu and Dao, 2023), Lightning Attention-2 (Qin et al., 2024)) reduces complexity to $\mathcal{O}(N)$. While efficient for mobile deployment, these methods often suffer from forgetfulness in high-fidelity recall tasks (Sun et al., 2025), posing risks for agents executing complex instructions. Similarly, Sparse Attention variants like SWA (Fu et al., 2025) restrict attention to local windows. Although cache-friendly for edge CPUs, SWA’s limited receptive field degrades performance in global reasoning tasks (Fan et al., 2025).

Hybrid approaches such as MoBA (Lu et al., 2025) and MoM (Du et al., 2025) combine sparse and dense attention. However, these typically employ spatial mixing and do not take advantage of the stage-specific structure of reasoning. They do not dynamically minimize resources during the transient thinking phase. SADAS addresses this via temporal mixing, adapting attention modes to the model’s real-time cognitive state to maximize edge hardware potential.

3 Method

To achieve an optimal balance between the exploratory and integrative stages of reasoning on resource-constrained devices, we propose a State-Aware Dynamic Attention Mechanism. The core idea of this framework is to enable the model to reconfigure its core attention computation paradigm in real-time at the token level during autoregressive generation, based on its current cognitive task.

3.1 Architectural Overview

The overall architectural design of SADAS is based on a key insight: complex reasoning processes are naturally divided into two cognitive stages, exploratory thinking akin to agentic planning and integrative answering corresponding to execution. To map this cognitive model onto a computational architecture, we designed a dynamic switching framework, as depicted in Fig. 1. This framework seamlessly integrates two computation modes: bandwidth-efficient SWA for the exploratory phase of quickly generating a Chain-of-Thought, and information-fidelity preserving Full Attention for the final answer stage, which requires integrating global information.

To enable autonomous switching between modes, we introduce a set of predefined control tokens: `<think>` and `</think>`. These tokens serve as explicit signals for the model’s internal state. When the model generates `<think>`, it automatically enters the hardware-friendly SWA mode; upon generating `</think>`, it switches back to the high-fidelity Full Attention mode. This design empowers the model with autonomous control over its computational flow.

At the implementation level, SADAS utilizes a unified and append-only KV cache (Zhou et al., 2024). In SWA mode, attention computation only accesses the most recent portion of the cache; whereas in Full Attention mode, it can access the entire historical data within the cache. This design is particularly advantageous for CPU-based inference as it eliminates the costly memory copy operations typically associated with cache eviction

in standard sliding window implementations. It significantly simplifies architectural complexity, ensuring a high degree of continuity and parallelism in the computational process.

The complete workflow of SADAS is as follows: when the model needs to perform complex reasoning, it first generates the <think> token, which triggers the framework to switch its computational core to the SWA Attention Layer. In this mode, the model efficiently generates a series of intermediate thought steps, minimizing memory bandwidth consumption. When the thinking process concludes, the model generates the </think> token, marking the completion of the exploratory phase. At this point, the framework seamlessly switches the computational core back to the Full Attention Layer, and the model begins to integrate all contextual information, including the thought chain, to finally generate an accurate answer.

3.2 State-Aware Dynamic Attention Mechanism

To realize the architecture described, we designed a dynamic scheduling mechanism that relies on precise tracking of each sequence’s cognitive state and is formalized by dynamic attention computation graph reconfiguration. For the i -th sequence in a batch, at each generation timestep t , we define its cognitive state $S_t^{(i)}$ and a termination state latch $F_t^{(i)}$, where $S_t^{(i)}, F_t^{(i)} \in \{0, 1\}$. $S_t^{(i)}$ indicates whether the model is currently in thinking mode and $F_t^{(i)}$ marks whether the entire thinking process has concluded.

3.2.1 State Transition Dynamics

The model’s state evolution is driven by the previously generated token $y_{t-1}^{(i)}$. Let y_{think} and $y_{\text{end_think}}$ denote the IDs of the control tokens, respectively. The cognitive state $S_t^{(i)}$ represents whether the model is in the thinking mode. Its transition logic is defined as follows:

$$S_t^{(i)} = \begin{cases} 1 & \text{if } y_{t-1}^{(i)} = y_{\text{think}} \\ 0 & \text{if } y_{t-1}^{(i)} = y_{\text{end_think}} \\ S_{t-1}^{(i)} & \text{otherwise} \end{cases} \quad (1)$$

where $S_t^{(i)} = 1$ corresponds to the thinking mode, and $S_t^{(i)} = 0$ to the integrative mode. To ensure that once reasoning is complete, the model stably enters the high-fidelity integrative mode, we introduce the termination state latch $F_t^{(i)}$. This is a

unidirectional trigger that, once activated, remains persistent. Its update rule is as follows:

$$F_t^{(i)} = F_{t-1}^{(i)} \vee (y_{t-1}^{(i)} = y_{\text{end_think}}) \quad (2)$$

where \vee denotes the logical OR operation, with initial state $F_0^{(i)} = 0$. This latch ensures that after a complete reasoning chain, the model will be locked into the integrative mode and will not switch back to SWA.

3.2.2 State-Dependent Attention Selection

The computation of the standard self-attention mechanism can be abstractly represented as:

$$\text{Attention}(Q, K, V) = \text{softmax} \left(\frac{QK^T}{\sqrt{d_k}} \right) V \quad (3)$$

The core innovation of our mechanism lies in transforming the key (K) and value (V) matrices provided to this function from static, undifferentiated historical records into a dynamically configured context set based on the cognitive state $F_t^{(i)}$. This process can be viewed as a real-time reconfiguration of the attention computation graph at each decoding step.

For the complete key-value history of sequence i at timestep t , we define it as $\mathcal{H}_t^{(i)} = \{(k_j^{(i)}, v_j^{(i)})\}_{j=1}^t$. When the model enters the high-fidelity integrative mode, $F_t^{(i)} = 1$, it must have access to the global context for accurate recall and integration. In this computational path, the attention mechanism is granted unrestricted access to the complete history $\mathcal{H}_t^{(i)}$. At this point, the effective key-value pairs ($\mathbf{K}_{\text{Full}}^{(i)}, \mathbf{V}_{\text{Full}}^{(i)}$) are formed by concatenating all key and value tensors from $\mathcal{H}_t^{(i)}$ along the sequence dimension.

Conversely, when the model is in the computationally more efficient exploratory thinking stage, $F_t^{(i)} = 0$, to accelerate thinking, its context access is actively restricted to a local window of a predefined size W_c . In this path, the effective key-value pairs ($\mathbf{K}_{\text{SWA}}^{(i)}, \mathbf{V}_{\text{SWA}}^{(i)}$) consist only of the most recent subset of the history $\mathcal{H}_t^{(i)}$. This dynamic constraint on historical records is crucial for achieving a trade-off between computational efficiency and model performance.

Therefore, the final key-value set ($\mathbf{K}_{\text{eff}}^{(i)}, \mathbf{V}_{\text{eff}}^{(i)}$) used by the effective attention computation $\text{Attn}_{\text{eff}}^{(l,i,t)}$ of the l decoding layer can be formally

described as a conditional selection function controlled by the state $F_t^{(i)}$:

$$(\mathbf{K}_{\text{eff}}^{(i)}, \mathbf{V}_{\text{eff}}^{(i)}) = \begin{cases} (\mathbf{K}_{\text{SWA}}^{(i)}, \mathbf{V}_{\text{SWA}}^{(i)}) & \text{if } F_t^{(i)} = 0 \\ (\mathbf{K}_{\text{Full}}^{(i)}, \mathbf{V}_{\text{Full}}^{(i)}) & \text{if } F_t^{(i)} = 1 \end{cases} \quad (4)$$

Finally, the output of this layer $h_t^{(l,i)}$ is computed based on this dynamically selected context set:

$$h_t^{(l,i)} = \text{Attention}(q_t^{(l,i)}, \mathbf{K}_{\text{eff}}^{(i)}, \mathbf{V}_{\text{eff}}^{(i)}) \quad (5)$$

In terms of implementation, this high-level stateful logic is elegantly mapped onto low-level computational optimizations. The state variable F_t is propagated across model layers via a boolean flag, ultimately passing a concrete window size to the underlying attention implementation to execute the logic of Equation 4. This achieves a direct mapping from abstract semantic states to concrete hardware-accelerated computational paths, thereby enabling dynamic control over the model’s computation mode without introducing significant architectural complexity.

4 Experiments

To empirically evaluate the effectiveness and efficiency of the SADAS framework, this section presents a series of experiments. We first detail the experimental setup, including model initialization, fine-tuning scheme, and evaluation benchmarks. Subsequently, we present a performance comparison of SADAS with current mainstream reasoning models and architectural variants. Finally, through a series of ablation studies, we deeply analyze the impact of the framework’s key components and hyperparameters.

4.1 Experimental Setup

4.1.1 Model Initialization and Baselines

SADAS was initialized by distilling and mapping weights from the pre-trained Qwen3 model (Yang et al., 2025). For comprehensive assessment, we compared SADAS against a representative set of baseline models, including top-tier closed-source models like GPT-4o (Hurst et al., 2024) and Claude-3.5-Sonnet. Open-source baselines included advanced reasoning models such as DeepSeek-R1-Distill-Llama-8B (Guo et al., 2025), Qwen3 (Yang et al., 2025), its pure Sliding Window Attention variant SWAQwen, and the latest hybrid architecture reasoning model M1-3B (Wang et al., 2025).

4.1.2 Experimental Setup

Lightweight Adaptation Fine-tuning. To enable the SADAS framework, we performed a lightweight fine-tuning stage (Team, 2025). We emphasize that this process is designed strictly for mechanism adaptation rather than knowledge injection. Its core objectives are twofold: (1) Cognitive Rhythm Alignment: enabling the model to learn the timing for generating control tokens; and (2) Computation-Mode Adaptation: allowing the model weights to adapt to the receptive field shifts between Sliding Window Attention and Full Attention. Detailed training parameters and settings are provided in the Appendix B.

4.1.3 Benchmark

Benchmarks and Metrics. To provide a holistic assessment of SADAS as a foundation for on-device agents, we evaluated the model across two distinct capability dimensions:

Deep Reasoning: We utilized the challenging AIME24 and AIME25 datasets, employing the pass@ k metric (with $k = 8$) (Chen et al., 2021; Brown et al., 2024) to measure robust mathematical problem-solving. We also included Math-500 (Hendrycks et al., 2021), a curated collection of 500 diverse mathematical problems, to assess breadth.

General and Agentic Capabilities: To verify that the efficiency gains of SADAS do not compromise the general knowledge and strict instruction adherence required for agents, we incorporated MMLU (Hendrycks et al., 2020) for general world knowledge, BBH (Suzgun et al., 2023) for complex symbolic reasoning, and IFEval (Zhou et al., 2023) to specifically evaluate prompt compliance and instruction following fidelity.

All evaluations were conducted using the EvalScope framework, scoring mathematical problems based on strict equivalence. For consistency across accuracy benchmarks, the SWA layer window size was fixed at 4096 tokens, using standard decoding parameters (temperature 0.7, top-p 0.8) (Guo et al., 2025; Luo et al., 2025)

On-Device Inference Setup. To evaluate practical deployment efficiency, we conducted end-to-end latency tests on two representative consumer platforms: an Intel Core i5-12600KF desktop CPU and a Qualcomm Snapdragon 8 Gen 3 mobile SoC, utilizing the llama.cpp engine. All inference speed tests were conducted using Q4_0 quantization precision. To simulate realistic edge memory con-

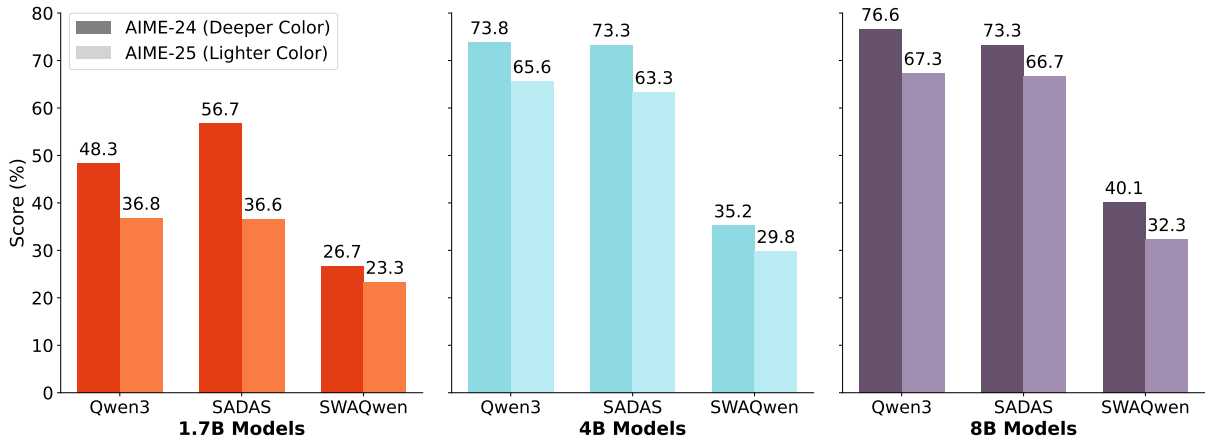


Figure 2: Comparison of SADAS and baseline Qwen3 variants on AIME24 and AIME25 datasets.

Table 1: Reasoning and General Evaluation Results for SADAS and Baseline Models

Model	Math500	AIME-24	AIME-25	MMLU	BBH	IFeval
GPT-4o-0513	74.6	9.3	11.6	87.2	–	84.3
Claude-3.5-Sonnet-1022	78.3	16.0	7.4	88.3	–	86.5
M1-3B	81.7	23.0	22.0	–	–	–
QWEN3-1.7b	89.0	48.3	36.8	55.7	48.2	33.6
DeepSeek-R1-Distill-Llama-8B	89.1	50.4	32.9	54.2	60.7	43.9
DeepSeek-R1-Distill-Qwen-1.5B	83.9	28.9	23.5	35.2	41.9	29.6
SADAS-1.7B	89.4	56.7	36.6	55.7	47.3	41.0
SADAS-4B	91.8	73.3	63.3	68.4	73.5	45.4

straints, the sliding window size was set to 1024 tokens. Tests were performed with a batch size of 1 and a fixed prompt length of 2048 tokens across varying output lengths from 1024 to 32768 tokens.

4.2 Main Results Comparing with Baselines

As shown in Table 1, SADAS shows excellent performance on all benchmarks. On relatively foundational benchmarks like Math500, SADAS performs comparably with other high-performance models, showing a slight improvement. However, SADAS’s advantage becomes particularly prominent when evaluated on benchmarks requiring deeper reasoning capabilities, such as AIME24 and AIME25. SADAS-4B achieved 73.3 on AIME24 and 63.3 on AIME25, significantly outperforming all comparable and even much larger models. Even at the 1.7B parameter scale, SADAS-1.7B’s performance surpassed that of other larger models, demonstrating improvements across all three datasets compared to the 8B-parameter DeepSeek-R1-Distill-Llama. Beyond mathematical reasoning, our results further validate the framework’s robustness in general and agentic capabilities. On MMLU,

SADAS-1.7B achieved 55.7, effectively matching the base model Qwen3-1.7B at 55.7, confirming that the efficiency optimization does not compromise general knowledge. Most notably, on the IFEval benchmark, SADAS-1.7B scored 41.0, marking a substantial improvement over the Qwen3-1.7B baseline of 33.6. This significant gain suggests that the dedicated thinking phase serves as an effective cognitive planning step, thereby enhancing the model’s ability to strictly adhere to complex instructions essential for on-device agents.

Table 2: Completion Token Lengths on AIME24

Model	Ave	Max	Min
Qwen3-1.7B	16737	28248	4771
SADAS-1.7B	18883	32768	5544

To further clarify the specific contributions of the SADAS framework, we conducted a direct comparison with the Full Attention model Qwen3 and its pure SWA variant, SWAQwen. As illustrated in Fig. 2, the performance of SADAS consistently positions it between Qwen3 and SWAQwen. No-

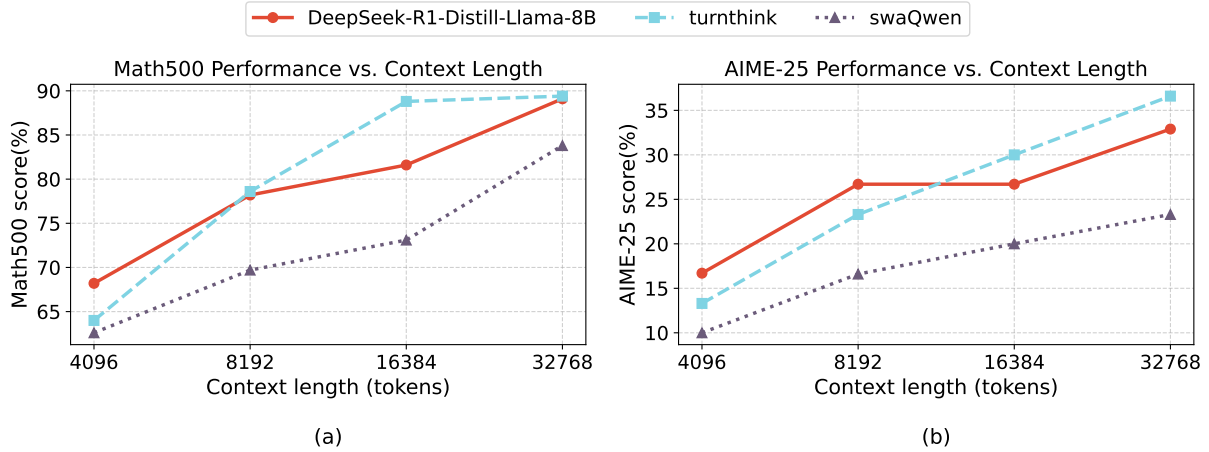


Figure 3: Relationship between Generation Length and Score Accuracy

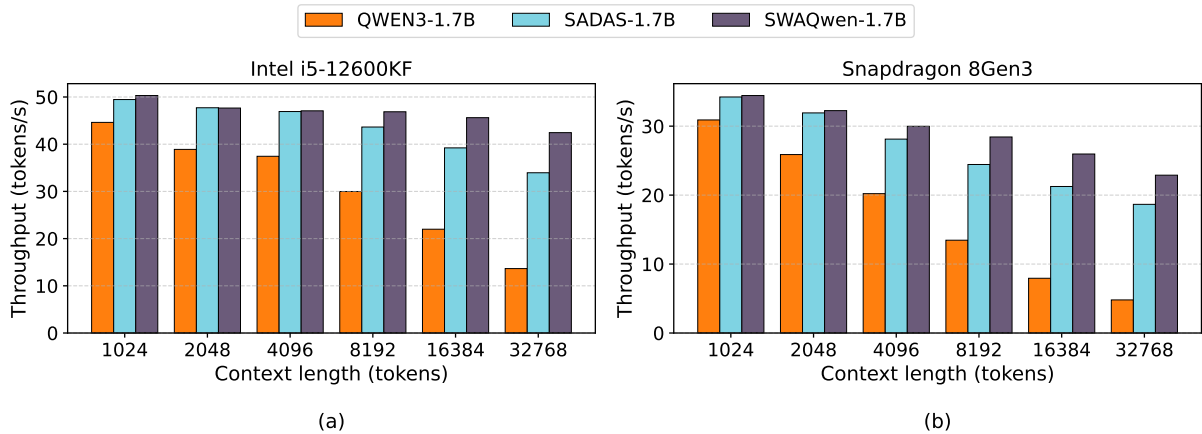


Figure 4: End-to-end inference latency of SADAS and Qwen3 variants under different output lengths (seconds)

tably, at the 1.7B scale, SADAS achieved a nearly 16% lead over Qwen3 on AIME24, representing a substantial improvement. As shown in Table 2, under the same 32k maximum output length constraint, SADAS demonstrated an approximate 12% improvement in the average length of generated answers compared to Qwen3, with maximum and minimum lengths increasing by nearly 16%. This confirms that SADAS successfully achieves an effective trade-off between computational efficiency and model fidelity, and to some extent, unlocks stronger reasoning potential. This is achieved by employing SWA during the thinking phase to accelerate reasoning and utilizing Full Attention for global information integration during the synthesis phase, resulting in only minor performance loss, far superior to pure SWA architectures that completely sacrifice global context.

Table 4 presents a detailed comparison of generation speeds for 1.7B-scale models on representative edge hardware (comprehensive results for

PC CPUs are provided in Appendix A). Under the severe memory bandwidth constraints of mobile devices, the speed of the standard Qwen3-1.7B model degrades drastically as the sequence length increases, plummeting to just 4.8 tokens/s at 32,768 tokens. By leveraging hardware-friendly SWA during the thinking phase, SADAS achieves faster end-to-end inference speeds than Qwen3 and its variants on both the Snapdragon 8 Gen 3 and Intel PC CPUs.

Specifically, on the Snapdragon 8 Gen 3 platform, when the output length reaches 32,768 tokens, SADAS maintains a speed of 18.66 tokens/s, achieving an approximate 3.88 \times speedup compared to Qwen3-1.7B. This effectively transforms long-chain reasoning from an unusable state to a real-time capability. On the Intel PC CPU, at 32,768 tokens, SADAS sustains 33.95 tokens/s, representing an approximate 2.5 \times speedup over Qwen3’s 13.65 tokens/s. Although SWAQwen is marginally faster than SADAS, its significant degra-

514 dation in accuracy renders it inferior to SADAS in
 515 the balance between usability and speed. This un-
 516 derlines the architectural superiority of SADAS:
 517 it utilizes an append-only mechanism to avoid the
 518 severe accuracy loss characteristic of pure SWA
 519 models, confirming its critical role in efficiently
 520 processing long-sequence reasoning on edge de-
 521 vices.

522 4.3 Ablation Study

523 4.3.1 Impact of Output Length

524 Figure 3 illustrates the impact of increased out-
 525 put length on complex reasoning. While all
 526 models benefit from extended reasoning depth,
 527 SADAS demonstrates the most significant effi-
 528 ciency gains. Quantitatively, extending the max-
 529 imum output length from 4096 to 32768 tokens,
 530 SADAS’s performance on Math500 improved
 531 by 39.7%, outperforming DeepSeek-R1-Distill-
 532 Llama-8B (30.6%) and SWAQwen (33.9%). On
 533 the more challenging AIME25, SADAS showed a
 534 remarkable 175.2% increase, significantly exceed-
 535 ing DeepSeek-R1-Distill-Llama-8B’s 97.0% and
 536 SWAQwen’s 133.0%.

537 This superior performance extension stems from
 538 SADAS’s dynamic switching, which efficiently
 539 generates long reasoning paths with SWA during
 540 the thinking phase and then leverages Full Atten-
 541 tion for lossless global recall and integration in the
 542 final synthesis. This allows SADAS to optimally
 543 translate reasoning depth into problem-solving ca-
 544 pabilities, particularly for complex tasks requiring
 545 a global perspective, where its performance curve
 546 exhibits the steepest growth.

547 4.3.2 Impact of Window Size

Table 3: Performance under Different Window Sizes

Model	Window Size	Math500	AIME24	AIME25
SADAS	1024	53.7	16.7	13.3
	2048	76.2	26.7	23.3
	4096	89.4	56.7	36.6
SWAQwen	1024	32.4	6.6	3.3
	2048	62.8	10.0	16.7
	4096	83.8	26.7	23.3

548 Table 3 presents performance results across dif-
 549 ferent window sizes, investigating how the width
 550 of local context in SADAS’s thinking phase affects
 551 model performance. The maximum output length
 552 was set to 32,768 tokens. Analysis reveals that in-
 553 creasing the window size significantly benefits both

554 SADAS and pure SWA architectures, underscoring
 555 the importance of a wider local receptive field dur-
 556 ing exploratory thinking. For instance, on AIME24,
 557 SADAS’s score surged by 239.5% (from 16.7 to
 558 56.7) when the window expanded from 1024 to
 559 4096.

560 Crucially, SADAS utilizes this expanded con-
 561 text much more efficiently than pure SWA. The
 562 performance gap between SADAS and SWAQwen
 563 systematically widens with increasing window size.
 564 On AIME24, SADAS’s lead over SWAQwen grew
 565 from 10.1 points (1024 window size) to 30.0 points
 566 (4096 window size), demonstrating SADAS’s per-
 567 formance more than doubling SWAQwen’s at the
 568 largest window size. This widening gap confirms
 569 SADAS’s hybrid mechanism maximizes the value
 570 of the SWA window. The SWA efficiently gener-
 571 ates high-quality thought content, which the subse-
 572 quent Full Attention layer, with its global access,
 573 meticulously processes for the final answer. This
 574 paradigm, which allocates computational resources
 575 to a stage that can be fully leveraged by a subse-
 576 quent high-fidelity module, is more effective than
 577 universally employing a single approximate atten-
 578 tion mechanism.

579 5 Conclusion

580 We propose SADAS, a dynamic inference frame-
 581 work designed to resolve the memory bandwidth
 582 bottleneck of long-chain reasoning on edge devices.
 583 By autonomously toggling between bandwidth-
 584 efficient sliding window attention and high-fidelity
 585 full attention, SADAS aligns computational modes
 586 with the model’s cognitive phases. Experimental
 587 results confirm that SADAS transforms unrespon-
 588 sive reasoning into a real-time mobile capability,
 589 achieving accuracy comparable to leading 8B-scale
 590 reasoning models. As the first architecture to im-
 591 plement token-level dynamic attention mixing, our
 592 work establishes a robust path for next-generation
 593 on-device agents that demand both speed and intel-
 594 ligence.

595 Limitations

596 While our proposed SADAS framework demon-
 597 strates significant improvements in balancing in-
 598 ference efficiency and reasoning accuracy, we ac-
 599 knowledge several limitations that offer avenues
 600 for future research.

- 601 1. **Static Window Size Configuration:** Cur-
 602 rently, SADAS employs a pre-defined fixed

603 window size (e.g., 1024 or 4096 tokens) for
 604 the sliding window attention phase. While ex-
 605 perimental results confirm this setting covers
 606 the majority of local context required for rea-
 607 soning steps, a static window may not be the-
 608 oretically optimal for every specific instance.
 609 Future work could explore adaptive mecha-
 610 nisms that dynamically adjust the window size
 611 based on real-time token entropy or seman-
 612 tic density, potentially further optimizing the
 613 trade-off between memory bandwidth usage
 614 and information recall.

615 **2. Memory Consumption of the KV Cache:**

616 Our implementation of SADAS prioritizes
 617 inference speed and architectural simplicity
 618 by utilizing a unified, append-only KV cache
 619 (DynamicCache). While this approach avoids
 620 the latency overhead associated with cache
 621 eviction in standard sliding window attention,
 622 it comes at the cost of increased memory con-
 623 sumption. The KV cache grows linearly with
 624 the length of the entire generated sequence.
 625 For extremely long reasoning chains (e.g., far
 626 exceeding 32k tokens), this could become a
 627 practical bottleneck on hardware with lim-
 628 ited VRAM, representing a direct trade-off
 629 between speed and memory footprint.

630 **3. Rigidity of the Bipartite Cognitive Model:**

631 The current SADAS framework is built upon
 632 a binary model of cognition, bifurcating the
 633 generation process into a single exploratory
 634 thinking phase (SWA) and a final integrative
 635 answering phase (Full Attention). However,
 636 complex reasoning may not always follow
 637 such a linear path. It could be iterative, re-
 638 quiring the model to switch back and forth
 639 between exploration and integration multiple
 640 times. Our current design, with its unidirec-
 641 tional termination latch ($F^{(i)}$), does not sup-
 642 port such complex, multi-turn cognitive state
 643 transitions, potentially limiting its effective-
 644 ness on tasks that require more sophisticated
 645 reasoning patterns.

646 **4. Generalizability to Other Long-Context**

647 **Tasks:** Our experiments have primarily fo-
 648 cused on mathematical and logical reason-
 649 ing benchmarks, where the Chain-of-Thought
 650 paradigm is well-established. The applicabil-
 651 ity and effectiveness of the SADAS frame-
 652 work for other long-context tasks, such as

653 long-form document summarization, narrative
 654 generation, or complex question-answering
 655 over large texts, remain to be thoroughly inves-
 656 tigated. The think-then-answer structure may
 657 not be as naturally suited or as beneficial for
 658 these domains, and adapting the framework
 659 might require designing new task-specific con-
 660 trol mechanisms.

661 **References**

662 Bradley Brown, Jordan Juravsky, Ryan Ehrlich, Ronald
 663 Clark, Quoc V Le, Christopher Ré, and Azalia Mirho-
 664 seini. 2024. Large language monkeys: Scaling infer-
 665 ence compute with repeated sampling. *arXiv preprint*
 666 *arXiv:2407.21787*.

667 Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan,
 668 Henrique Ponde De Oliveira Pinto, Jared Kaplan,
 669 Harri Edwards, Yuri Burda, Nicholas Joseph, Greg
 670 Brockman, and 1 others. 2021. Evaluating large
 671 language models trained on code. *arXiv preprint*
 672 *arXiv:2107.03374*.

673 Yuhong Chou, Man Yao, Kexin Wang, Yuqi Pan, Rui-
 674 Jie Zhu, Jibin Wu, Yiran Zhong, Yu Qiao, Bo Xu,
 675 and Guoqi Li. 2024. Metala: Unified optimal linear
 676 approximation to softmax attention map. *Advances*
 677 *in Neural Information Processing Systems*, 37:71034–
 678 71067.

679 Jusen Du, Weigao Sun, Disen Lan, Jiayi Hu, and
 680 Yu Cheng. 2025. Mom: Linear sequence mod-
 681 eling with mixture-of-memories. *arXiv preprint*
 682 *arXiv:2502.13685*.

683 Qihang Fan, Huaibo Huang, and Ran He. 2025. Break-
 684 ing the low-rank dilemma of linear attention. In *Pro-*
 685 *ceedings of the Computer Vision and Pattern Recog-*
 686 *nition Conference*, pages 25271–25280.

687 Zichuan Fu, Wentao Song, Yejing Wang, Xian Wu,
 688 Yefeng Zheng, Yingying Zhang, Derong Xu, Xue-
 689 tao Wei, Tong Xu, and Xiangyu Zhao. 2025. Sliding
 690 window attention training for efficient large language
 691 models. *arXiv preprint arXiv:2502.18845*.

692 Albert Gu and Tri Dao. 2023. Mamba: Linear-time
 693 sequence modeling with selective state spaces. *arXiv*
 694 *preprint arXiv:2312.00752*.

695 Daya Guo, Dejian Yang, Haowei Zhang, Junxiao
 696 Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shi-
 697 rong Ma, Peiyi Wang, Xiao Bi, and 1 others. 2025.
 698 Deepseek-r1: Incentivizing reasoning capability in
 699 llms via reinforcement learning. *arXiv preprint*
 700 *arXiv:2501.12948*.

701 Dongchen Han, Ziyi Wang, Zhuofan Xia, Yizeng Han,
 702 Yifan Pu, Chunjiang Ge, Jun Song, Shiji Song,
 703 Bo Zheng, and Gao Huang. 2024. Demystify mamba
 704 in vision: A linear attention perspective. *arXiv*
 705 *preprint arXiv:2405.16605*.

706	Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. 2020. Measuring massive multitask language understanding. <i>arXiv preprint arXiv:2009.03300</i> .	762
707		763
708		764
709		765
710	Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song, and Jacob Steinhardt. 2021. Measuring mathematical problem solving with the math dataset. <i>arXiv preprint arXiv:2103.03874</i> .	766
711		767
712		768
713		769
714		770
715	Namgyu Ho, Sangmin Bae, Taehyeon Kim, Hyunjik Jo, Yireun Kim, Tal Schuster, Adam Fisch, James Thorne, and Se-Young Yun. 2024. Block transformer: Global-to-local language modeling for fast inference. <i>Advances in Neural Information Processing Systems</i> , 37:48740–48783.	771
716		772
717		773
718		774
719		775
720		776
721	Aaron Hurst, Adam Lerer, Adam P Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark, AJ Ostrow, Akila Welihinda, Alan Hayes, Alec Radford, and 1 others. 2024. Gpt-4o system card. <i>arXiv preprint arXiv:2410.21276</i> .	777
722		778
723		779
724		780
725		
726	Aaron Jaech, Adam Kalai, Adam Lerer, Adam Richardson, Ahmed El-Kishky, Aiden Low, Alec Helyar, Aleksander Madry, Alex Beutel, Alex Carney, and 1 others. 2024. Openai o1 system card. <i>arXiv preprint arXiv:2412.16720</i> .	781
727		782
728		783
729		784
730		785
731	Lingjie Jiang, Xun Wu, Shaohan Huang, Qingxiu Dong, Zewen Chi, Li Dong, Xingxing Zhang, Tengchao Lv, Lei Cui, and Furu Wei. 2025. Think only when you need with large hybrid-reasoning models. <i>arXiv preprint arXiv:2505.14631</i> .	786
732		787
733		788
734		789
735		790
736	Ao Liu, Botong Zhou, Can Xu, Chayse Zhou, ChenChen Zhang, Chengcheng Xu, Chenhao Wang, Decheng Wu, Dengpeng Wu, Dian Jiao, and 1 others. 2025. Hunyuan-turbos: Advancing large language models through mamba-transformer synergy and adaptive chain-of-thought. <i>arXiv preprint arXiv:2505.15431</i> .	791
737		792
738		793
739		794
740		
741		795
742		796
743	Enzhe Lu, Zhejun Jiang, Jingyuan Liu, Yulun Du, Tao Jiang, Chao Hong, Shaowei Liu, Weiran He, Enming Yuan, Yuzhi Wang, and 1 others. 2025. Moba: Mixture of block attention for long-context llms. <i>arXiv preprint arXiv:2502.13189</i> .	797
744		798
745		799
746		800
747		
748	Michael Luo, Sijun Tan, Justin Wong, Xiaoxiang Shi, William Y. Tang, Manan Roongta, Colin Cai, Jeffrey Luo, Li Erran Li, Raluca Ada Popa, and Ion Stoica. 2025. Deepscaler: Surpassing o1-preview with a 1.5b model by scaling rl. Notion Blog.	801
749		802
750		803
751		804
752		805
753	Wenjie Ma, Jingxuan He, Charlie Snell, Tyler Griggs, Sewon Min, and Matei Zaharia. 2025. Reasoning models can be effective without thinking. <i>arXiv preprint arXiv:2504.09858</i> .	806
754		807
755		808
756		809
757		810
758		811
759		
760	Xuezhe Ma, Xiang Kong, Sinong Wang, Chunting Zhou, Jonathan May, Hao Ma, and Luke Zettlemoyer. 2021. Luna: Linear unified nested attention. <i>Advances in Neural Information Processing Systems</i> , 34:2441–2453.	812
761		813
		814
		815
		816
	Bo Peng, Eric Alcaide, Quentin Anthony, Alon Albalak, Samuel Arcadinho, Stella Biderman, Huanqi Cao, Xin Cheng, Michael Chung, Matteo Grella, and 1 others. 2023. Rwkv: Reinventing rns for the transformer era. <i>arXiv preprint arXiv:2305.13048</i> .	
	Zhen Qin, Weigao Sun, Dong Li, Xuyang Shen, Weixuan Sun, and Yiran Zhong. 2024. Lightning attention-2: A free lunch for handling unlimited sequence lengths in large language models. <i>arXiv preprint arXiv:2401.04658</i> .	
	Zhuoran Shen, Mingyuan Zhang, Haiyu Zhao, Shuai Yi, and Hongsheng Li. 2021. Efficient attention: Attention with linear complexities. In <i>Proceedings of the IEEE/CVF winter conference on applications of computer vision</i> , pages 3531–3539.	
	Weigao Sun, Disen Lan, Yiran Zhong, Xiaoye Qu, and Yu Cheng. 2025. Lasp-2: Rethinking sequence parallelism for linear attention and its hybrid. <i>arXiv preprint arXiv:2502.07563</i> .	
	Mirac Suzgun, Nathan Scales, Nathanael Schärli, Sebastian Gehrmann, Yi Tay, Hyung Won Chung, Aakanksha Chowdhery, Quoc Le, Ed Chi, Denny Zhou, and 1 others. 2023. Challenging big-bench tasks and whether chain-of-thought can solve them. In <i>Findings of the Association for Computational Linguistics: ACL 2023</i> , pages 13003–13051.	
	Ling Team. 2025. Ring-lite: Scalable reasoning via c3po-stabilized reinforcement learning for llms . <i>Preprint</i> , arXiv:2506.14731.	
	Junxiong Wang, Wen-Ding Li, Daniele Paliotta, Daniel Ritter, Alexander M Rush, and Tri Dao. 2025. M1: Towards scalable test-time compute with mamba reasoning models. <i>arXiv preprint arXiv:2504.10449</i> .	
	Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, and 1 others. 2022. Chain-of-thought prompting elicits reasoning in large language models. <i>Advances in neural information processing systems</i> , 35:24824–24837.	
	Zhiheng Xi, Wenxiang Chen, Xin Guo, Wei He, Yiwen Ding, Boyang Hong, Ming Zhang, Junzhe Wang, Senjie Jin, Enyu Zhou, and 1 others. 2025. The rise and potential of large language model based agents: A survey. <i>Science China Information Sciences</i> , 68(2):121101.	
	An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, and 1 others. 2025. Qwen3 technical report. <i>arXiv preprint arXiv:2505.09388</i> .	
	Jiarui Yao, Yifan Hao, Hanning Zhang, Hanze Dong, Wei Xiong, Nan Jiang, and Tong Zhang. 2025. Optimizing chain-of-thought reasoners via gradient variance minimization in rejection sampling and rl. <i>arXiv preprint arXiv:2505.02391</i> .	

817 Han Zhang, Ian Goodfellow, Dimitris Metaxas, and
818 Augustus Odena. 2019. Self-attention generative
819 adversarial networks. In *International conference on*
820 *machine learning*, pages 7354–7363. PMLR.

821 Libo Zhang, Zhaoning Zhang, Rui Li, Zhiliang Tian,
822 Songzhu Mei, Dongsheng Li, and 1 others. 2025a.
823 Dovetail: A cpu/gpu heterogeneous speculative de-
824 coding for llm inference. In *Proceedings of the 2025*
825 *Conference on Empirical Methods in Natural Lan-*
826 *guage Processing*, pages 17393–17406.

827 Renrui Zhang, Jiaming Han, Chris Liu, Peng Gao, Ao-
828 jun Zhou, Xiangfei Hu, Shilin Yan, Pan Lu, Hong-
829 sheng Li, and Yu Qiao. 2023. Llama-adapter: Effi-
830 cient fine-tuning of language models with zero-init
831 attention. *arXiv preprint arXiv:2303.16199*.

832 Shilong Zhang, Peize Sun, Shoufa Chen, Min Xiao,
833 Wenqi Shao, Wenwei Zhang, Yu Liu, Kai Chen,
834 and Ping Luo. 2025b. Gpt4roi: Instruction tuning
835 large language model on region-of-interest. In *Euro-*
836 *pean Conference on Computer Vision*, pages 52–70.
837 Springer.

838 Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang,
839 Xiaolei Wang, Yupeng Hou, Yingqian Min, Beichen
840 Zhang, Junjie Zhang, Zican Dong, and 1 others. 2023.
841 A survey of large language models. *arXiv preprint*
842 *arXiv:2303.18223*, 1(2).

843 Jeffrey Zhou, Tianjian Lu, Swaroop Mishra, Sid-
844 dhartha Brahma, Sujoy Basu, Yi Luan, Denny Zhou,
845 and Le Hou. 2023. Instruction-following evalua-
846 tion for large language models. *arXiv preprint*
847 *arXiv:2311.07911*.

848 Xiabin Zhou, Wenbin Wang, Minyan Zeng, Jiaxian Guo,
849 Xuebo Liu, Li Shen, Min Zhang, and Liang Ding.
850 2024. Dynamickv: Task-aware adaptive kv cache
851 compression for long context llms. *arXiv preprint*
852 *arXiv:2412.14838*.

853 A PC-CPU Inference Performance

854 Table 4 details the inference speed (tokens/s) across
855 1.7B, 4B, and 8B parameter scales on the Intel
856 Core i5-12600KF CPU. All inference speed tests
857 were conducted using Q4_0 quantization preci-
858 sion. The results demonstrate a consistent trend: as
859 the sequence length extends to 32,768 tokens, the
860 standard Qwen3 models suffer severe performance
861 degradation due to memory bandwidth saturation,
862 with the 4B and 8B models dropping to single-digit
863 speeds (6.84 and 5.13 tokens/s, respectively). In
864 contrast, SADAS maintains robust generation rates
865 across all scales. Notably, for the 4B model at
866 32k length, SADAS achieves 18.15 tokens/s, rep-
867 resenting a nearly 2.6× speedup over the baseline,
868 ensuring that even larger models remain responsive
869 on consumer-grade CPUs.

870 B Training Details and Hyperparameters

871 We utilized the inclusionAI/Ring-lite-distill-
872 preview-sft-data. To demonstrate the data
873 efficiency of our approach and ensure that perfor-
874 mance gains are attributable to the architecture
875 rather than extensive training, we restricted the
876 training data to only the first 5,000 samples.
877 Fine-tuning the SADAS-1.7B model required
878 only 3.5 hours on a single NVIDIA A800 GPU.
879 This minimal training overhead confirms that the
880 model’s reasoning capabilities are preserved from
881 the base model, while the efficiency gains are
882 derived from the SADAS architecture.

883 Training was conducted for one epoch. To
884 accommodate the lengthy Chain-of-Thought se-
885 quences typical in reasoning tasks, we set the
886 maximum sequence length to 8192 tokens. We
887 employed a conservative initial learning rate of
888 7×10^{-6} with a cosine decay schedule to minimize
889 the disruption to the model’s pre-trained knowledge
890 representations.

891 C Analysis of the Necessity for 892 Fine-tuning

893 In our research, to activate and stabilize SADAS’s
894 dynamic switching capability, we performed con-
895 tinuous full-parameter fine-tuning on the base
896 model. We did not adopt parameter-efficient fine-
897 tuning (PEFT) methods such as LoRA. Our primary
898 considerations were to provide an unconstrained
899 upper bound for evaluating the potential of the
900 SADAS architecture and to ensure all model param-
901 eters adapted to the new framework, thereby avoid-
902 ing potential performance bottlenecks that PEFT
903 methods might introduce. To validate the necessity
904 of fine-tuning, we designed a set of comparative
905 experiments, directly testing performance under
906 an initial setup where only model weights were
907 imported and no fine-tuning was performed. The
908 model parameters were set to 1.7B, and training
909 settings referred to Appendix B.

910 From the data in Table 5, it is clearly evident
911 that fine-tuning is a prerequisite for the successful
912 operation of the SADAS framework. The non-fine-
913 tuned SADAS model exhibited significant perfor-
914 mance degradation across all key metrics, which
915 can be attributed to the following reasons:

- 916 • **Cognitive Mismatch:** The internal weights
917 of the non-fine-tuned model have not learned
918 the "think-answer" rhythm. Forcing attention

Table 4: End-to-end inference throughput (tokens/s) of turnthink and Qwen3 under different output lengths.

Model	1024	2048	4096	8192	16384	32768
QWEN3-1.7B	44.62	38.90	37.44	29.98	21.99	13.65
turnthink	49.46	47.72	46.92	43.64	39.22	33.95
swaQwen	50.31	47.66	47.07	48.87	46.62	42.45
QWEN3-4B	22.16	20.70	18.49	14.97	10.73	6.84
turnthink	24.02	23.11	21.64	21.14	20.94	18.15
swaQwen	24.12	23.59	23.38	23.10	22.35	21.92
QWEN3-8B	13.16	12.72	11.80	10.30	8.26	5.13
turnthink	13.25	12.93	12.67	12.33	11.68	10.56
swaQwen	13.31	13.12	12.98	12.78	12.11	11.53

Table 5: Performance Comparison of SADAS under Fine-tuned and Non-Fine-tuned Settings

Scheme	AIME-24	AIME-25	COT Stability
Fine-tuned	56.7	36.6	High
Non-Fine-tuned	12.3	8.7	Low

mode switching disrupts its inherent generation logic, leading to incoherent model outputs and even complete disorientation in complex reasoning tasks, as evidenced by the substantial drop in AIME-24 scores.

- Generation Instability:** As the model does not understand the semantics of control tokens and has not adapted to the dynamic changes in the attention calculation range, its generation process becomes highly unstable. We observed numerous instances of repetition, logical interruptions, or premature generation termination, which explains its "low" rating in CoT generation stability.
- Reduced Computational Efficiency:** Although SWA is theoretically faster, the actual inference speed of the non-fine-tuned model was lower than that of the end-to-end optimized fine-tuned model. This is due to uncoordinated internal states during mode switching, which leads to frequent interruptions in the GPU computation flow.

In summary, a targeted fine-tuning phase is crucial for SADAS. It not only teaches the model how to autonomously utilize control tokens but, more importantly, reshapes the model’s internal computation flow, enabling it to transition smoothly and efficiently between the two attention modes, thereby

truly converting SADAS’s architectural advantages into practical performance gains.

D Exploration and Trade-offs of Alternative Cache Management Schemes

The core of the SADAS framework lies in its dynamic nature, and KV cache management is the technical cornerstone for achieving this dynamism. The globally unified DynamicCache scheme adopted in the main text is the optimal solution we derived after comprehensive consideration of speed, accuracy, and implementation complexity. To more comprehensively illustrate our design decision process, this section will detail two other alternative schemes we explored, analyzing their advantages and limitations with experimental data.

Similarly, all experiments in this section were conducted with 1.7B model parameters, and all throughput tests used a batch size of 16, performing inference on a single A100 GPU.

D.1 Scheme 1: Hybrid Cache

This scheme aims to minimize GPU memory footprint. Its core mechanism involves using a fixed-size SWA Cache during the thinking phase. Upon transitioning to the answering phase, the contents of the SWA Cache are migrated to a new, infinitely growing Dynamic Cache via a one-time cache copy operation.

As shown in Table 6, the Hybrid Cache scheme demonstrates a significant advantage in peak GPU memory usage. For example, at an output length of 32768, it reduced memory demand by approximately 94.2% compared to SADAS. However, this advantage comes at the cost of sacrificing critical

Table 6: Performance Evaluation of the Hybrid Cache Scheme

Scheme		Output Length						
		512	1024	2048	4096	8192	16384	32768
Hybrid Cache	Speed (token/s)	17.09	33.65	71.93	169.41	453.71	1421	4853
	Peak Memory (GB)	0.889	1.763	1.763	1.764	1.766	1.768	3.221
SADAS	Speed (token/s)	18.9	37.84	75.31	163.97	417.68	1271.93	4316.15
	Peak Memory (GB)	0.889	1.765	3.515	7.016	14.019	28.02	55.808

Table 7: Performance Evaluation of the Truncated Dynamic Cache Scheme

Scheme		Output Length						
		512	1024	2048	4096	8192	16384	32768
Truncated Cache	Speed (token/s)	17.08	33.34	71.54	168.37	451.59	1418	4833
	Peak Memory (GB)	0.889	1.763	1.763	1.764	1.766	1.768	3.221
SADAS	Speed (token/s)	18.9	37.84	75.31	163.97	417.68	1271.93	4316.15
	Peak Memory (GB)	0.889	1.765	3.515	7.016	14.019	28.02	55.808

Table 8: Performance Evaluation of Alternative Schemes Integrating Sink Attention

Scheme		Output Length						
		512	1024	2048	4096	8192	16384	32768
Truncated Cache	Speed (token/s)	16.76	33.22	71.20	168.97	452.74	1416.22	4814
	Peak Memory (GB)	0.889	1.763	1.763	1.764	1.766	1.768	3.221
Hybrid Cache	Speed (token/s)	17.02	34.08	72.13	167.74	535.11	1418.22	4821
	Peak Memory (GB)	0.889	1.763	1.763	1.764	1.766	1.768	3.221
Truncated Cache + Sink	Speed (token/s)	17.08	33.34	71.54	168.37	451.59	1418	4833
	Peak Memory (GB)	0.889	1.764	1.764	1.765	1.766	1.768	3.223
Hybrid Cache + Sink	Speed (token/s)	17.09	33.65	71.93	169.41	453.71	1421	4853
	Peak Memory (GB)	0.889	1.764	1.764	1.765	1.766	1.768	3.223

Table 9: Reasoning Accuracy Evaluation of the Hybrid Cache Scheme

Scheme	AIME-24	AIME-25
Hybrid Cache	36.6	26.7
SADAS	56.7	36.6

performance. Its speed significantly decreased, primarily due to the substantial latency introduced by the cache copying operation, which represents a difficult-to-optimize serial bottleneck on the GPU. More critically, as shown in Table 9, its accuracy also suffered significantly, because the context for the final answering phase was limited to only

the latter part of the thinking process, leading the model to lose a large amount of crucial early information and fail to complete complex reasoning tasks requiring long-range recall.

D.2 Scheme 2: Truncated Dynamic Cache

To address the speed issues of the Hybrid Cache scheme, we designed the Truncated Dynamic Cache. This scheme uses a single global cache, but during the thinking phase, it logically truncates access to older cache entries outside the window through attention masks.

As can be seen from the data in Table 7, this scheme avoids data copying, and its speed performance is close to that of our final adopted

Table 10: Reasoning Accuracy Evaluation of the Truncated Dynamic Cache Scheme

Scheme	AIME-24	AIME-25
Truncated Cache	36.6	26.7
SADAS	56.7	36.6

scheme. Its memory footprint is the same as the final scheme, as it still retains all KV pairs at the underlying level. However, as shown in Table 10, during the thinking phase, the model is similarly unable to recall early thinking steps outside the window, resulting in its accuracy still being lower than the final scheme. While this involves less information loss than the Hybrid Cache scheme, this limitation remains critical in complex reasoning chains that require repeated backtracking and verification.

Through a rigorous evaluation of the three schemes discussed, we concluded that while the alternative schemes offer attractive benefits in terms of memory savings, they all compromise the model’s peak reasoning ability by introducing some form of permanent information loss during the thinking phase. The globally unified Dynamic-Cache scheme we ultimately adopted has its core advantage in ensuring information completeness and flexible access. It allows the model to recall the entire history at any time by modifying attention masks, a capability that proved crucial for achieving efficient inference without sacrificing accuracy. Although it demands higher GPU memory, it provides the best overall performance in terms of both speed and accuracy, which is fully consistent with the original design philosophy of SADAS.

E GPU Inference Performance Analysis

E.1 Performance Observation

To quantify computational efficiency on server-grade hardware, we conducted end-to-end latency tests on a single NVIDIA A800 GPU. With a fixed prompt length of 2048 tokens and a batch size of 16, models generated output sequences of varying lengths (from 512 to 32,768 tokens). As detailed in Table 11, while SADAS still outperforms baselines at extreme lengths, the acceleration is less pronounced compared to edge CPUs.

Specifically, at 32,768 tokens, SADAS-1.7B reduces total generation time from 4937.93s (Qwen3) to 4316.15s, a 12.6% speedup. However, at inter-

mediate lengths (e.g., 4096 tokens), the pure SWA variant actually exhibits higher latency (192.53s) compared to the full-attention baseline (183.34s). This counter-intuitive phenomenon—where a theoretically lighter mechanism runs slower—provides critical insight into the hardware-specific bottlenecks of linear attention.

E.2 Mechanism Analysis

Why does a pure SWA mechanism sometimes result in slower inference speeds? Our investigation revealed that its core bottleneck lies in the Sliding-WindowCache’s management. To maintain a fixed window size, it necessitates frequent eviction of the oldest KV pairs at each generation step. This non-trivial, often unparallelizable overhead—especially against highly optimized attention computations on GPUs—can negate the theoretical gains from reduced attention scope, leading to a theoretically fast, practically slow outcome.

In contrast, SADAS bypasses this issue by utilizing a globally unified DynamicCache. This cache is append-only, meaning KV pairs are simply added without costly eviction or rolling operations. This design simplifies data flow, eliminates cache management overhead, and ensures high computational continuity and parallelism. Consequently, the inherent speed benefits of linear attention methods like SWA are fully realized in SADAS, leading to robust acceleration (e.g., outperforming pure SWA by avoiding cache rolling) even on powerful accelerators, despite potentially higher memory consumption for very long sequences compared to fixed-window caches.

F Future Work: Exploration of Integrating Sink Attention

Building upon the success of the SADAS framework, our future work will focus on exploring token-level dynamic inference models with improved performance, higher efficiency, and lower memory footprint. A promising direction is to combine the cache optimization schemes we explored in Appendix B with cutting-edge long-sequence inference techniques, particularly the Attention Sink concept proposed in StreamingLLM.

The Attention Sink mechanism posits that in autoregressive models, the initial few tokens are crucial for maintaining attention stability and integrating global information, even when they fall outside the attention window. Retaining these ini-

Table 11: End-to-end inference latency of SADAS and Qwen3 variants under different output lengths (seconds)

Model	512	1024	2048	4096	8192	16384	32768
Qwen3-1.7B	17.16	34.60	72.24	183.34	595.11	1427.36	4937.93
SADAS-1.7B	18.90	37.84	75.31	163.97	417.68	1271.93	4316.15
SWAQwen-1.7B	16.87	33.65	71.73	192.53	508.92	1419.33	4617.54
Qwen3-4B	22.14	43.70	94.77	224.73	604.05	1862.90	–
SADAS-4B	25.86	47.99	97.64	214.26	547.98	1644.21	–
SWAQwen-4B	27.52	55.44	107.70	216.82	567.08	1800.73	–
Qwen3-8B	23.76	47.47	100.73	237.90	628.27	2097.96	–
SADAS-8B	26.28	52.30	106.53	248.72	616.95	1813.29	–
SWAQwen-8B	30.30	60.10	125.63	281.58	699.16	1963.14	–

tial tokens can effectively mitigate performance degradation caused by window sliding in long sequences. Inspired by this, we improved the two alternative schemes discussed in Appendix B:

1. **Hybrid Cache + Sink:** Building upon the original Hybrid Cache scheme, we permanently retained the initial few sink tokens within the SWA Cache. During cache copying, these sink tokens, along with the tokens within the sliding window, were copied to the new dynamic cache.
2. **Truncated Dynamic Cache + Sink:** In the Truncated Dynamic Cache scheme, we modified the attention mask to allow it to permanently attend to the initial sink tokens during the thinking phase, in addition to the tokens within the sliding window.

We conducted preliminary experiments on these two improved schemes, also using 1.7B model parameters, and all models underwent the same fine-tuning procedure.

Table 12: Reasoning Accuracy Evaluation of Alternative Schemes Integrating Sink Attention

Scheme	AIME-24	AIME-25
SADAS	56.7	36.6
Hybrid Cache	36.6	26.7
Hybrid Cache+Sink	46.6	32.9
Truncated Cache	36.6	26.7
Truncated Cache+Sink	46.6	32.9

From the results in Table 12, we can observe that by introducing the Attention Sink mechanism, the accuracy of both alternative schemes significantly

improved, showing an increase of approximately 27.3% (relative to their non-Sink counterparts) on AIME-24. This demonstrates the importance of retaining initial global information for maintaining long-range reasoning capabilities, even under cache-constrained conditions. Despite the improved accuracy, these enhanced schemes did not show an advantage in inference speed. As shown in Table 8, versions integrating Sink attention even exhibited slightly slower inference speeds. The speed bottleneck for the Hybrid Cache scheme remains the costly cache copying operation, while the Truncated Dynamic Cache scheme experienced a slight increase in computational complexity after introducing additional attention to sink tokens, leading to a marginal decrease in speed.

This preliminary exploration points us towards a clear direction for future work. We have demonstrated that Attention Sink can effectively compensate for the accuracy shortcomings of cache-optimized schemes. Therefore, our core future research will focus on fundamentally addressing the speed bottleneck while preserving Sink information and optimizing memory usage. Possible exploration paths include designing more efficient, copy-free cache update mechanisms, or leveraging hardware-aware algorithms to optimize access to non-contiguous caches (sliding window + Sink). The ultimate goal is to build a next-generation dynamic inference architecture that achieves state-of-the-art performance in accuracy, speed, and memory efficiency.