

Kinetics: Rethinking Test-Time Scaling Laws

Anonymous Authors¹

Abstract

We rethink test-time scaling laws from a *practical efficiency* perspective, revealing that the effectiveness of smaller models is significantly overestimated. Prior work, grounded in compute-optimality, overlooks critical memory access bottlenecks introduced by inference-time strategies (e.g., Best-of- N , long CoTs). Our holistic analysis, spanning models from 0.6B to 32B parameters, reveals a new *Kinetics Scaling Law* that better guides resource allocation by incorporating both computation and memory access costs. *Kinetics Scaling Law* suggests that test-time compute is more effective when used on models above a threshold than smaller ones. A key reason is that in TTS, attention, rather than parameter count, emerges as the dominant cost factor. Motivated by this, we propose a new scaling paradigm centered on *sparse attention*, which lowers per-token cost and enables longer generations and more parallel samples within the same resource budget. Empirically, we show that sparse attention models consistently outperform dense counterparts, achieving over **60 point** gains in low-cost regimes and over **5 point** gains in high-cost regimes for problem-solving accuracy on *AIME* and *LiveCodeBench*. These results suggest that sparse attention is essential for realizing the full potential of test-time scaling because, unlike training, where parameter scaling saturates, test-time accuracy continues to improve through increased generation.

1 Introduction

Test-time scaling (TTS) has recently emerged as a powerful strategy (e.g., Best-of- N , Long-CoT (Wei et al., 2022)) for enhancing the reasoning capabilities of large language models (LLMs) (Guo et al., 2025; Jaech et al., 2024; Team,

2025b), particularly in scenarios where agents interact with complex environments, e.g., writing code, browsing the web (Nakano et al., 2021; Yao et al., 2023b) or reinforcement learning (RL) with LLMs-in-the-loop (Huang et al., 2022; Driess et al., 2023; Chen et al., 2025a). These capabilities, however, introduce substantial inference-time costs, making it critical to understand performance scaling in this new paradigm. Existing scaling law studies (Brown et al., 2024; Snell et al., 2024; Wu et al., 2024) focus on floating-point operations (FLOPs) while ignoring memory access costs, which are often the dominant factor in determining wall-clock latency in TTS regimes. As shown in Figure 1a, this gap can lead to sub-optimal deployment decisions.

In Section 3, we introduce the *Kinetics Scaling Law* for TTS, derived from a novel cost model that explicitly incorporates memory access costs. This new perspective reveals markedly different conclusions about Pareto-optimal strategies for allocating test-time compute (Figure 1a). Specifically, we find that: (1) prior scaling laws consistently **overestimate** the effectiveness of small models enhanced with inference-time strategies; and (2) computational resources are best spent first on increasing model size up to a critical threshold (empirically around 14B), before investing in test-time strategies, such as Best-of- N sampling or long CoTs. Guided by the Kinetics Scaling Law, our approach yields up to a **3 \times** throughput improvement on B200 hardware.

Our roofline analysis across a suite of state-of-the-art reasoning models reveals that the shift in optimal test-time compute strategies arises because test-time strategies (e.g., long CoTs, Best-of- N) disproportionately increase attention costs rather than parameter costs (Figure 2a). Our Iso-cost analysis shows that the quadratic growth of attention with generation length, combined with the disproportionate scaling of KV memory relative to model parameters, drives a preference for scaling up model size over generations. This imbalance is further exacerbated by MoE architectures (Shazeer et al., 2017; Du et al., 2021; Fedus et al., 2022; AI@Meta, 2025; Dai et al., 2024; Jiang et al., 2024), which reduce active parameter count without alleviating attention overhead.

Building on this analysis, in Section 4 we introduce a new scaling paradigm, centered on **sparse** attention, which fundamentally reshapes the scaling law and significantly en-

¹Anonymous Institution, Anonymous City, Anonymous Region, Anonymous Country. Correspondence to: Anonymous Author <anon.email@domain.com>.

Preliminary work. Under review by the International Conference on Machine Learning (ICML). Do not distribute.

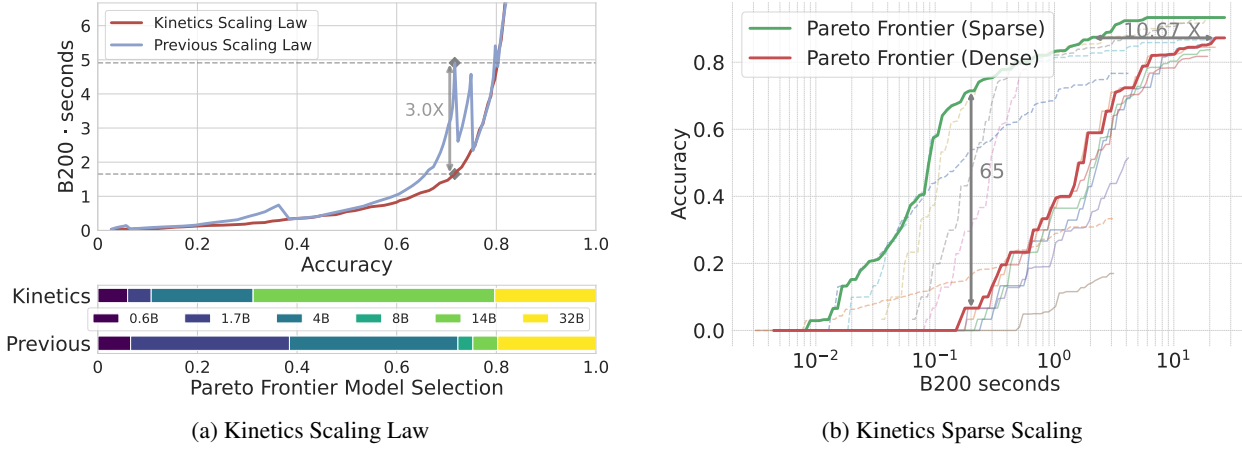


Figure 1. (a): Pareto Frontier for Qwen3 series on AIME24. Previous test-time scaling laws (Brown et al., 2024; Snell et al., 2024; Wu et al., 2024) focus solely on compute optimality, neglecting the significant bottleneck of memory access in long-sequence generation. This leads to suboptimal resource utilization. By incorporating memory access, the *Kinetics Scaling Law* reduces resource demands by up to 3 \times to achieve the same accuracy. (b): Inspired by the Kinetics Scaling Law, we show that *sparse attention models* scale significantly better than dense models, achieving over 50-point improvements in AIME24 in the low-cost regime and consistently outperforming dense models in the high-cost regime, in addition to substantial efficiency gains. B200 second represents the amount of work performed by a single B200 at full utilization for one second.

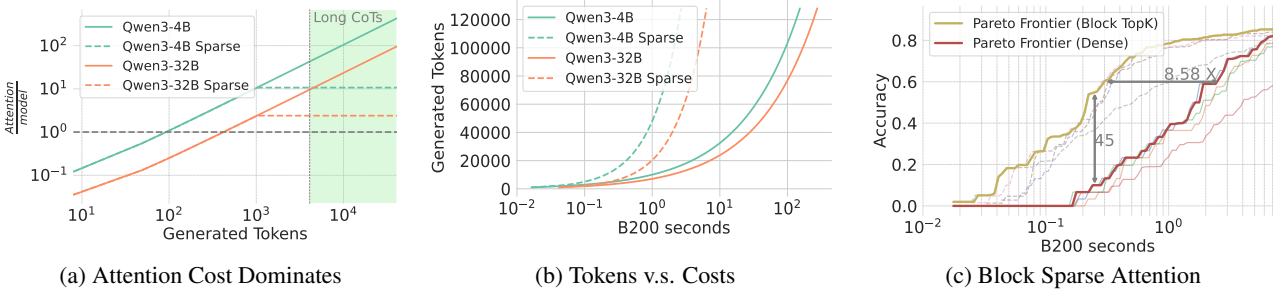


Figure 2. (a) Attention dominates inference cost, exceeding parameter computation by 100 \sim 1000 \times ; sparse attention alleviates this bottleneck. (b) Under equal resource constraints, sparse attention enables significantly more token generation, enhancing test-time scaling. (c) Simple block-sparse attention delivers major gains—improving accuracy by 45 points in low-cost settings and matching dense accuracy with 8.58 \times fewer resources.

hances the scalability of TTS (Figure 1b). According to our *Kinetics Sparse Scaling Law*, computational resources are best allocated to test-time strategies rather than reducing sparsity. As more computing is invested at test time, lower sparsity becomes increasingly critical to fully leveraging the benefits of these strategies. Guided by this principle, it increases problem-solving rates by up to 60 points in the low-cost regime and over 5 points in the high-cost regime on AIME24 and LiveCodeBench, through massive generated tokens, which is unaffordable for dense counterparts.

While sparsity has traditionally been employed either for regularization in small models (Tibshirani, 1996; Molchanov et al., 2017) or to reduce computation in over-parameterized networks (Mishra et al., 2021; Chen et al., 2021; Hoefler et al., 2021; Dao et al., 2021; Frantar & Alishtarh, 2023; Liu et al., 2023), our work introduces a funda-

mentally different perspective: sparsity as a central enabler of efficient and scalable test-time inference. In contrast to pretraining – where scaling laws often exhibit diminishing returns (Ilya) – TTS continues to benefit from increased token generation and more optimized inference paths. We hope this study can guide and encourage future co-design of model architectures, inference-time strategies, and hardware to fully unlock the next wave of scaling at deployment.

2 Cost Model and eFLOPs

We propose a cost model that captures both compute and memory access overhead during inference, focusing on realistic deployment settings (batch size $\gg 1$, model parallelism, and shared prompt cache). Notation is in Table 1.

Computation and Memory. Following (Brown et al., 2024), the compute cost combines linear layer operations and self-attention:

$$C_{\text{comp}} = 2PL_{\text{out}} + r(2L_{\text{in}} + L_{\text{out}})L_{\text{out}}D$$

Memory access includes both parameter loading and KV cache reads:

$$C_{\text{mem}} = 2PL_{\text{out}} + 2L_{\text{in}}L_{\text{out}}D + L_{\text{out}}^2D$$

In practice, parameter loading is amortized across large batches (DeepSeek-AI, 2025), so we omit that term and share prompt KV cache across N trials. The final per-task compute and memory cost becomes:

$$C_{\text{comp}}(N) = 2PNL_{\text{out}} + 2rNL_{\text{in}}L_{\text{out}}D + rNL_{\text{out}}^2D \quad (1)$$

$$C_{\text{mem}}(N) = 2L_{\text{in}}L_{\text{out}}D + NL_{\text{out}}^2D \quad (2)$$

eFLOPs. We define eFLOPs (equivalent FLOPs) as a linear combination of compute and memory cost, scaled by hardware intensity I to capture the memory and computational operations under the same scale:

$$\text{eFLOPs} = C_{\text{comp}} + C_{\text{mem}} \cdot I$$

We use $I = 562.5$ FLOPs-s/GB from NVIDIA B200.

2.1 Takeaway: Attention Bottleneck

Our eFLOPs calculation reveals that **attention dominates inference cost for long generations**. The ratio of attention to parameter cost, $\Phi = \frac{2rL_{\text{in}}D + (rD + ID)L_{\text{out}}}{2P}$, can exceed $100\times$ for $L_{\text{out}} > 4\text{k}$ (Figure 2a). This effect is magnified in MoEs (AI@Meta, 2025; Dai et al., 2024), which reduce linear FLOPs and further shift the bottleneck to attention.

Scalability Implication. Given long-CoT usage, where $L_{\text{out}} \gg L_{\text{in}}$, inference cost is increasingly governed by the quadratic term L_{out}^2D , motivating our *Kinetics Scaling Law*, akin to kinetic energy: $E_k = \frac{1}{2}mv^2$.

More details are in Appendix A.

2.2 Experimental Setup

We evaluate LLMs under our cost model on three challenging reasoning tasks: AIME24 (MAA, 2024), AIME25 (MAA, 2025), and LiveCodeBench (Jain et al., 2024), using the Qwen3 (Yang et al., 2025) model family. Our theoretical estimates assume NVIDIA B200 hardware. Our evaluations are focused on two representative inference strategies, *Long-CoTs* and *Best-of-N*.

3 Rethinking Test-time Scaling Law

In Section 3.1, we first introduce the Kinetics Scaling Law, derived from empirical investigations across the Qwen3

model series. Then, we explore the underlying reasons for the divergence between Kinetics and prior scaling laws through an Iso-Cost analysis in Section 3.2.

3.1 Kinetics Scaling Law

we study the scaling behavior of Qwen3 (Yang et al., 2025) considering the following problem:

Given a fixed inference budget (eFLOPs per question), what is the Pareto frontier of achievable accuracy across different LLM configurations?

In the *Long CoTs* setting (single trial per question, $N_T = 1$), we vary generation length n_T to evaluate performance across cost levels. Results in Figure 3 reveal two key findings of our Kinetics Scaling Law.

- **Efficiency of small models is overestimated.** As shown in Figures 2b and 3 (a, c), smaller models like 4B and 8B are outperformed by the 14B model even at low accuracy levels (e.g., below 40%). The 0.6B model appears on the Pareto frontier only when accuracy is negligible. In contrast to prior scaling laws, which gave smaller models more prominence, our results show they are often suboptimal in practice.
- **CoT length more effective than parameter size only beyond a critical model scale (empirically, 14B).** The Kinetics Scaling Law shows that, under limited compute, scaling up the model yields greater benefits than extending CoT length. As seen in Figure 3 (b, d), only the 14B and 32B models gain from CoTs longer than 10K tokens. For smaller models (e.g., 1.7B and 4B), switching to a larger model is more effective when $L_{\text{out}} < 5\text{K}$. This suggests compute should primarily be allocated to increasing model size, not generation length (Figure 3 (d)). In contrast, previous scaling laws assumed longer CoTs consistently improved performance across all model sizes and only favored model scaling once those gains plateaued.

In the *Best-of-N* setting, we fix the maximum number of generated tokens at n_T , and vary the number of reasoning trials N to evaluate the problem-solving rate (i.e., the probability that at least one trial produces a correct answer). We have similar observations in Figures 4a to 4c. Under the previous scaling laws (Figure 4b), the most cost-effective strategy to achieve high accuracy is to apply repeated sampling using smaller models. Kinetics Scaling Law Figure 4a reveals that deploying a 14B model with fewer reasoning trials is more efficient. We also observe a critical size of 14B. For models smaller than 14B, increasing compute is best allocated toward model scaling rather than additional trials. For models at or above 14B, however, further computation is more effectively spent on increasing the number of reasoning trials, up to diminishing returns.

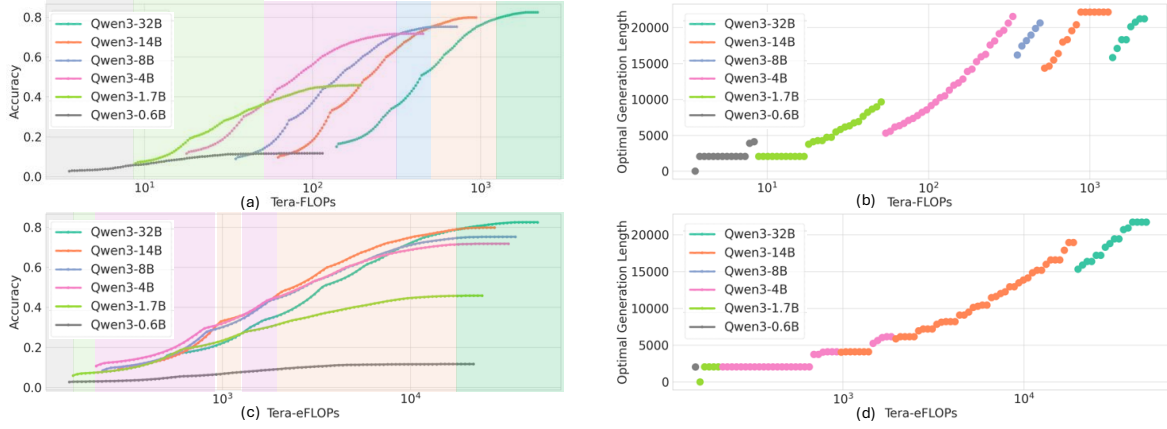


Figure 3. AIME Pareto Frontier (Long-CoTs). We first launch evaluations for Qwen3 series models. By controlling the maximum generation lengths, we control the inference cost in eFLOPs (ab) for our scaling law) or FLOPs (cd for previous scaling law) and measure the accuracy (Pass@1) in AIME24. The optimal model is marked with different colors in (ac). The optimal generation length is in (bd).

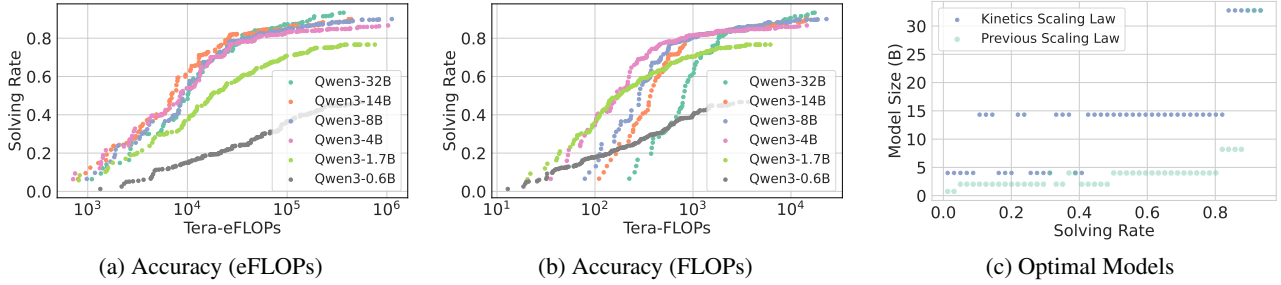


Figure 4. AIME24 Score Curve Envelope (Best-of-N). We control the incurred inference cost in eFLOPs (a) or FLOPs (b) and measure the solving rate (Coverage) in AIME24 for various models by varying the maximum allowed number of reasoning trials. By taking the curve envelopes, we can project the optimal models in (c).

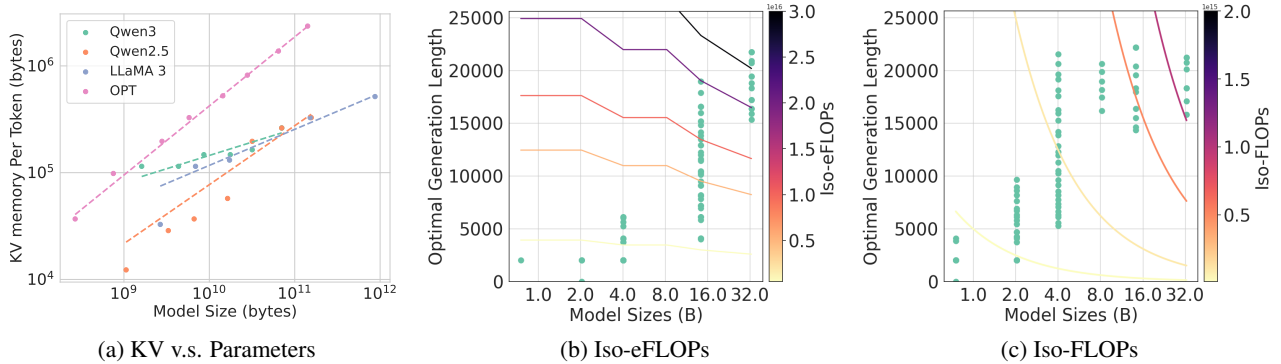


Figure 5. Explanation of the New Scaling Law. (a) Analysis across four LLM families reveals a trend of disproportionately slower KV memory growth relative to model size. (bc) Iso-Cost landscapes under the proposed cost model (b) and the traditional model (c).

3.2 Iso-Cost Study

We attribute the above divergence between Kinetics and previous scaling laws to two reasons.

Disproportionation between KV memory size D and model parameters P . Smaller models tend to require significantly more KV cache relative to their parameter size. For example, Qwen3-0.6B demands 3.5GB of KV cache to store 32K tokens, despite the model itself occupying only 1.2GB. In contrast, Qwen3-32B uses just 8GB of KV cache for the same sequence length. Empirically, doubling model parameters results in only a $1.18\times$ increase in KV cache size. As shown in Figure 5a, this phenomenon is consistently observed across model families such as OPT (Zhang et al., 2022) ($1.55\times$), Qwen2.5 (Yang et al., 2024) ($1.46\times$), and LLaMA3 (Grattafiori et al., 2024) ($1.27\times$).

Shift from linear to quadratic cost model. Under this revised model, increasing generation length incurs a substantially higher cost than scaling model size; consequently, the tradeoff between model capacity and token budget shifts meaningfully. For instance, under the linear LP model, the cost of generating 8K tokens with a 14B model (which is usually insufficient to solve complex tasks) is treated as equivalent to generating 24K tokens with a 4B model (sufficient to complete most tasks). However, under the L^2D model, the same 14B@8K generation is only comparable in cost to a 4B@9K generation. This tighter bound makes it much harder for smaller models to compensate for their limited capacity through extended generation alone. Thus, only if the gap in model capacities is small enough (e.g., 32B only improves the accuracy by 3% on AIME24 compared to 14B), the benefits of extending generation length might be more effective than directly enlarging model parameters.

Figures 5b and 5c show an Iso-Cost analysis comparing two cost models. Under Kinetics Scaling Law, the cost grows quadratically with L_{out} , while the KV cache scales sublinearly with model parameters P . As a result, when total budget is low, the Iso-eFLOPs contours tend to stretch horizontally, favoring larger model sizes over longer generation lengths. This implies that increasing model size is a more efficient use of resources than generating longer outputs. In contrast, the traditional FLOPs-based model leads to steeply vertical contours, encouraging longer generation before increasing model size. More details are in Appendix B.

4 Sparse Test-time Scaling Law

Based on our findings in Section 3, we propose a new scaling paradigm centered on sparse attention. Sparse attention fundamentally reshapes the Kinetics Scaling Law in Section 3 and enhances the scalability of TTS.

Sparse attention significantly enhances problem-solving

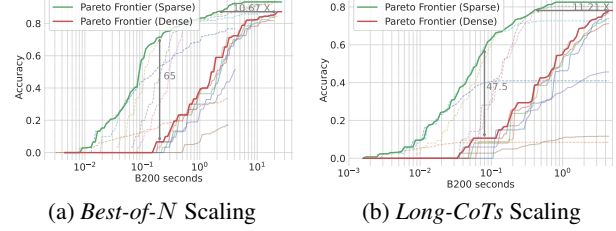


Figure 6. **Sparse Attention Boosts Test-Time Scaling.** In (a) and (c), we show that sparse attention models significantly improve the cost-accuracy trade-off under both inference strategies.

performance. As shown in Figures 6a and 6b, compared to dense baselines, for both of the inference strategies and models of various sizes, sparse attention models improve problem-solving rates by up to 60 points in the low-cost regime and over 5 points in the high-cost regime. From an efficiency perspective, dense models require over $10\times$ more eFLOPs to match the same solving rate. These findings underscore sparse attention as a key enabler for unlocking the full benefits of test-time scaling.

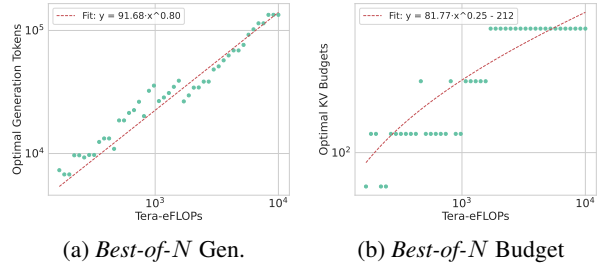


Figure 7. **Tradeoff Between Generated Tokens and KV Budget.** We empirically investigate how to balance the tradeoff between generating more tokens and allocating a larger KV cache budget, which may yield more accurate but potentially shorter outputs. Using Qwen3-8B as a representative model, we fit curves to characterize this tradeoff. For *Best-of-N*, we find that for every doubling of the total compute cost, the optimal KV budget increases by a factor of $1.18\times$, while the total number of generated tokens increases by $1.74\times$.

Sparse attention becomes increasingly valuable in high-cost scenarios. We investigate the tradeoff between KV budget B and reasoning trials (N). Our analysis reveals a consistent trend: allocating additional compute toward generating more tokens is generally more effective than expanding the KV cache. In *Best-of-N* frontier, doubling the cost leads to only a $1.18\times$ increase in KV budget, compared to a $1.74\times$ increase in total generated tokens.

Sparse attention reshapes the Kinetics Scaling Law. As shown in Section 4, applying sparse attention significantly improves the efficiency of smaller models (0.6B, 1.7B, 4B), allowing them to re-emerge on the Pareto frontier across a

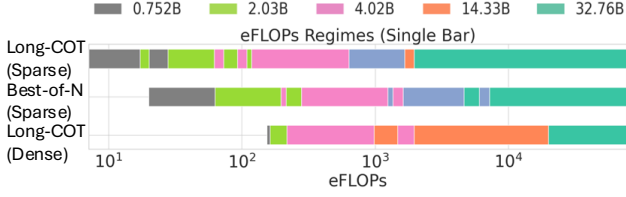


Figure 8. Compared to the scaling law for the dense models (a), small models (0.6B, 1.7B, 4B) are more effective with sparse attention. They occupy more space in the Pareto Frontier (Figure 6a).

broader range. Sparse attention reduces attention memory access from a quadratic cost term (L^2D) to a linear one (LBD), making it negligible or comparable when compared to the cost of computing with model parameters (LP).

More details are in Appendices C and D.

5 Experimental Validation

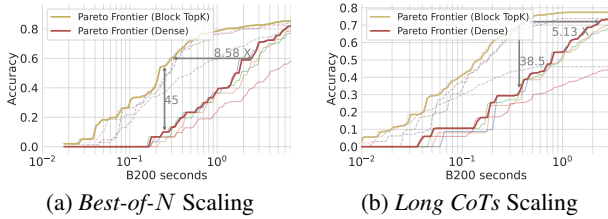


Figure 9. **Sparse Attention Algorithms.** We illustrate the optimality of block top- k sparse attention in terms of TTS on AIME24 dataset. Although upper bounded by the oracle top- k attention performance, block top- k achieves a good trade-off between effectiveness and tractability.

In this section, we demonstrate the practicality of our sparse scaling law through block top- k attention. We report empirical improvements in task throughput (number of tasks performed per unit time) using our block-sparse implementation and conduct ablation studies with alternative sparsification strategies, such as local attention, to highlight the importance of the KV selection mechanism.

5.1 Block Top- k Attention

While top- k attention offers attractive theoretical scaling, it is computationally intractable in practice. Instead, we adopt block top- k attention for two key reasons. *First*, it exploits temporal locality in attention patterns (Sun et al., 2024a) to retrieve semantically related key-value (KV) blocks. *Second*, its localized retrieval is hardware-efficient and integrates seamlessly with paged attention (Kwon et al., 2023), enabling high-throughput decoding. In practice, we compute a representative vector for each KV block by averaging its key vectors, and use these to score the relevance of blocks to each query. Importance scores are shared across query heads within a group, following the Grouped Query Attention (GQA) scheme. As shown in Figures 9a and 9b, block

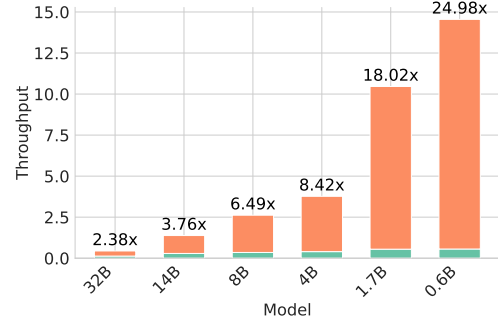


Figure 10. Task throughput improvement with block top- k .

top- k achieves a good trade-off between effectiveness and tractability, scaling far beyond dense counterparts.

5.2 Empirical Results

We quantify TTS efficiency using *task throughput*, defined as the number of tasks completed per unit time. This metric is particularly relevant for reasoning tasks, where the utility of generation hinges entirely on the correctness of the final output—unlike tasks such as summarization or content creation, where partial outputs may still be useful. We illustrate the benefit of block top- k attention across different model sizes on $8 \times H200$ machines with an extremely large batch size of 4096. As shown in Figure 10, block top- k attention substantially improves task throughput, particularly for smaller models. Qwen3-0.6B model achieves a $12.6 \times$ to $25 \times$ increase in throughput as the generation length extends from 16k to 32k tokens. This improvement reflects the growing inefficiency of dense attention at longer contexts, which disproportionately affects smaller models. Thus, the use of sparse attention not only alleviates this bottleneck but also restores much of the practical utility of smaller models in resource-constrained settings by enabling them to use more test-time compute more cost-effectively.

6 Conclusion and Discussion

This work introduces the *Kinetics Scaling Law*, showing that attention costs—not parameter counts—dominate test-time inference. Sparse attention reshapes the scaling landscape, enabling longer generations and higher accuracy. We view Kinetics Scaling as a foundation for guiding LLM serving, agent systems, and RL environments, especially as progress slows in pretraining. Though our analysis focuses on NVIDIA GPUs, the core insight—that memory bandwidth is harder to scale than FLOPs—applies broadly. Ultimately, our findings call for co-designing models, inference algorithms, and hardware to enable the next generation of scalable LLMs. Test-time scaling can consume a substantial amount of energy, raising concerns about the environmental sustainability of widespread deployment. By promoting sparse attention, our work hopes to help to reduce the carbon footprint and energy consumption of inference systems and support the broader goal of sustainable AI.

References

- AI@Meta. Llama 4 model card. 2025. URL https://github.com/meta-llama/llama-models/blob/main/models/llama4/MODEL_CARD.md.
- Arora, D. and Zanette, A. Training language models to reason efficiently, 2025. URL <https://arxiv.org/abs/2502.04463>.
- Brown, B., Juravsky, J., Ehrlich, R., Clark, R., Le, Q. V., Ré, C., and Mirhoseini, A. Large language monkeys: Scaling inference compute with repeated sampling. *arXiv preprint arXiv:2407.21787*, 2024.
- Cai, R., Tian, Y., Wang, Z., and Chen, B. Lococo: Dropping in convolutions for long context compression. *arXiv preprint arXiv:2406.05317*, 2024.
- Chen, B., Dao, T., Winsor, E., Song, Z., Rudra, A., and Ré, C. Scatterbrain: Unifying sparse and low-rank attention. *Advances in Neural Information Processing Systems*, 34: 17413–17426, 2021.
- Chen, C., Borgeaud, S., Irving, G., Lespiau, J.-B., Sifre, L., and Jumper, J. Accelerating large language model decoding with speculative sampling. *arXiv preprint arXiv:2302.01318*, 2023.
- Chen, K., Cusumano-Towner, M., Huval, B., Petrenko, A., Hamburger, J., Koltun, V., and Krähenbühl, P. Reinforcement learning for long-horizon interactive llm agents. *arXiv preprint arXiv:2502.01600*, 2025a.
- Chen, M., Hui, B., Cui, Z., Yang, J., Liu, D., Sun, J., Lin, J., and Liu, Z. Parallel scaling law for language models. *arXiv preprint arXiv:2505.10475*, 2025b. URL <https://arxiv.org/abs/2505.10475>.
- Chen, Z., Sadhukhan, R., Ye, Z., Zhou, Y., Zhang, J., Nolte, N., Tian, Y., Douze, M., Bottou, L., Jia, Z., et al. Mag-icpig: Lsh sampling for efficient llm generation. *arXiv preprint arXiv:2410.16179*, 2024.
- Child, R., Gray, S., Radford, A., and Sutskever, I. Generating long sequences with sparse transformers. *arXiv preprint arXiv:1904.10509*, 2019.
- Choromanski, K., Likhoshesterov, V., Dohan, D., Song, X., Gane, A., Sarlos, T., Hawkins, P., Davis, J., Mohiuddin, A., Kaiser, L., et al. Rethinking attention with performers. *arXiv preprint arXiv:2009.14794*, 2020.
- Dai, D., Deng, C., Zhao, C., Xu, R., Gao, H., Chen, D., Li, J., Zeng, W., Yu, X., Wu, Y., et al. Deepseekmoe: Towards ultimate expert specialization in mixture-of-experts language models. *arXiv preprint arXiv:2401.06066*, 2024.
- Dao, T. Flashattention-2: Faster attention with better parallelism and work partitioning. *CoRR*, abs/2307.08691, 2023. doi: 10.48550/ARXIV.2307.08691. URL <https://doi.org/10.48550/arXiv.2307.08691>.
- Dao, T., Chen, B., Liang, K., Yang, J., Song, Z., Rudra, A., and Ré, C. Pixelated butterfly: Simple and efficient sparse training for neural network models. In *International Conference on Learning Representations (ICLR)*, 2021. URL <https://arxiv.org/abs/2112.00029>.
- Dao, T., Fu, D. Y., Ermon, S., Rudra, A., and Ré, C. Flashattention: Fast and memory-efficient exact attention with io-awareness. In Koyejo, S., Mohamed, S., Agarwal, A., Belgrave, D., Cho, K., and Oh, A. (eds.), *Advances in Neural Information Processing Systems 35: Annual Conference on Neural Information Processing Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November 28 - December 9, 2022*, 2022.
- DeepSeek-AI. Deepseek open infra index. 2025. URL https://github.com/deepseek-ai/open-infra-index/blob/main/202502OpenSourceWeek/day_6_one_more_things_deepseekV3R1_inference_system_overview.md.
- Dettmers, T., Lewis, M., Belkada, Y., and Zettlemoyer, L. Gpt3. int8 (): 8-bit matrix multiplication for transformers at scale. *Advances in neural information processing systems*, 35:30318–30332, 2022.
- Driess, D., Nguyen, M., Xia, F., et al. Palm-e: An embodied multimodal language model. *arXiv preprint arXiv:2303.03378*, 2023.
- Du, N., Huang, Y., Dai, A. M., Tong, S., Lepikhin, D., Xu, Y., Chen, D., Wu, Y., and Dean, J. Glam: Efficient scaling of language models with mixture-of-experts. *arXiv preprint arXiv:2112.06905*, 2021.
- Fedus, W., Zoph, B., and Shazeer, N. Switch transformers: Scaling to trillion parameter models with simple and efficient sparsity. *Journal of Machine Learning Research*, 23(1):5232–5270, 2022.
- Frantar, E. and Alistarh, D. Sparsegpt: Massive language models can be accurately pruned in one-shot. In *International Conference on Machine Learning*, pp. 10323–10337. PMLR, 2023.
- Frantar, E., Ashkboos, S., Hoefler, T., and Alistarh, D. Gptq: Accurate post-training quantization for generative pre-trained transformers. *arXiv preprint arXiv:2210.17323*, 2022.
- Fu, Y., Chen, J., Zhu, S., Fu, Z., Dai, Z., Qiao, A., and Zhang, H. Efficiently serving llm reasoning programs with certaintex. *arXiv preprint arXiv:2412.20993*, 2024.

- Grattafiori, A., Dubey, A., Jauhri, A., Pandey, A., Kadian, A., Al-Dahle, A., Letman, A., Mathur, A., Schelten, A., Vaughan, A., et al. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*, 2024.
- Gu, A. and Dao, T. Mamba: Linear-time sequence modeling with selective state spaces. *CoRR*, abs/2312.00752, 2023. doi: 10.48550/ARXIV.2312.00752. URL <https://doi.org/10.48550/arXiv.2312.00752>.
- Gu, A., Goel, K., and Ré, C. Efficiently modeling long sequences with structured state spaces. In *The Tenth International Conference on Learning Representations, ICLR 2022, Virtual Event, April 25-29, 2022*. OpenReview.net, 2022. URL <https://openreview.net/forum?id=uYLFozlVlAC>.
- Guo, D., Yang, D., Zhang, H., Song, J., Zhang, R., Xu, R., Zhu, Q., Ma, S., Wang, P., Bi, X., et al. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*, 2025.
- Hao, S., Sukhbaatar, S., Su, D., Li, X., Hu, Z., Weston, J., and Tian, Y. Training large language models to reason in a continuous latent space. *arXiv preprint arXiv:2412.06769*, 2024.
- Hoefler, T., Alistarh, D., Ben-Nun, T., Dryden, N., and Peste, A. Sparsity in deep learning: Pruning and growth for efficient inference and training in neural networks. *Journal of Machine Learning Research*, 22(241):1–124, 2021.
- Hooper, C., Kim, S., Mohammadzadeh, H., Mahoney, M. W., Shao, Y. S., Keutzer, K., and Gholami, A. Kvquant: Towards 10 million context length llm inference with kv cache quantization. In Globerson, A., Mackey, L., Belgrave, D., Fan, A., Paquet, U., Tomczak, J., and Zhang, C. (eds.), *Advances in Neural Information Processing Systems*, volume 37, pp. 1270–1303. Curran Associates, Inc., 2024. URL https://proceedings.neurips.cc/paper_files/paper/2024/file/028fcbcf85435d39a40c4d61b42c99a4-Paper-Conference.pdf.
- Huang, W., Fei, F., and Finn, C. Language models as zero-shot planners: Extracting actionable knowledge for embodied agents. *arXiv preprint arXiv:2201.07207*, 2022.
- Ilya, S. Ilya sutskever: “sequence to sequence learning with neural networks: what a decade”. URL <https://www.youtube.com/watch?v=lyvBqasHLZs>.
- Jaech, A., Kalai, A., Lerer, A., Richardson, A., El-Kishky, A., Low, A., Helyar, A., Madry, A., Beutel, A., Carney, A., et al. Openai o1 system card. *arXiv preprint arXiv:2412.16720*, 2024.
- Jain, N., Han, K., Gu, A., Li, W.-D., Yan, F., Zhang, T., Wang, S., Solar-Lezama, A., Sen, K., and Stoica, I. Livecodebench: Holistic and contamination free evaluation of large language models for code. *arXiv preprint arXiv:2403.07974*, 2024.
- Jiang, A. Q., Sablayrolles, A., Roux, A., Mensch, A., Savary, B., Bamford, C., Chaplot, D. S., Casas, D. d. l., Hanna, E. B., Bressand, F., et al. Mixtral of experts. *arXiv preprint arXiv:2401.04088*, 2024.
- Juravsky, J., Brown, B., Ehrlich, R., Fu, D. Y., Ré, C., and Mirhoseini, A. Hydragen: High-throughput llm inference with shared prefixes. *arXiv preprint arXiv:2402.05099*, 2024.
- Katharopoulos, A., Vyas, A., Pappas, N., and Fleuret, F. Transformers are rnns: Fast autoregressive transformers with linear attention. In *Proceedings of the 37th International Conference on Machine Learning, ICML 2020, 13-18 July 2020, Virtual Event*, volume 119 of *Proceedings of Machine Learning Research*, pp. 5156–5165. PMLR, 2020. URL <http://proceedings.mlr.press/v119/katharopoulos20a.html>.
- Kitaev, N., Kaiser, Ł., and Levskaya, A. Reformer: The efficient transformer. In *The International Conference on Machine Learning (ICML)*, 2020.
- Kwon, W., Li, Z., Zhuang, S., Sheng, Y., Zheng, L., Yu, C. H., Gonzalez, J. E., Zhang, H., and Stoica, I. Efficient memory management for large language model serving with pagedattention, 2023. URL <https://arxiv.org/abs/2309.06180>.
- Leviathan, Y., Kalman, M., and Matias, Y. Fast inference from transformers via speculative decoding. In *International Conference on Machine Learning*, pp. 19274–19286. PMLR, 2023.
- Li, Y., Huang, Y., Yang, B., Venkitesh, B., Locatelli, A., Ye, H., Cai, T., Lewis, P., and Chen, D. Snapkv: Llm knows what you are looking for before generation, 2024. URL <https://arxiv.org/abs/2404.14469>.
- Lin, C., Tang, J., Yang, S., Wang, H., Tang, T., Tian, B., Stoica, I., Han, S., and Gao, M. Twilight: Adaptive attention sparsity with hierarchical top- p pruning. *arXiv preprint arXiv:2502.02770*, 2025.
- Lin, J., Tang, J., Tang, H., Yang, S., Chen, W.-M., Wang, W.-C., Xiao, G., Dang, X., Gan, C., and Han, S. Awq: Activation-aware weight quantization for on-device llm compression and acceleration. *Proceedings of Machine Learning and Systems*, 6:87–100, 2024a.

- Lin, Y., Tang, H., Yang, S., Zhang, Z., Xiao, G., Gan, C., and Han, S. Qserve: W4a8kv4 quantization and system co-design for efficient llm serving. *arXiv preprint arXiv:2405.04532*, 2024b.
- Liu, Z., Wang, J., Dao, T., Zhou, T., Yuan, B., Song, Z., Shrivastava, A., Zhang, C., Tian, Y., Re, C., et al. Dejavu: Contextual sparsity for efficient llms at inference time. In *International Conference on Machine Learning*, pp. 22137–22176. PMLR, 2023.
- Liu, Z., Yuan, J., Jin, H., Zhong, S., Xu, Z., Braverman, V., Chen, B., and Hu, X. Kivi: A tuning-free asymmetric 2bit quantization for kv cache. *arXiv preprint arXiv:2402.02750*, 2024.
- Ma, W., He, J., Snell, C., Griggs, T., Min, S., and Zaharia, M. Reasoning models can be effective without thinking. *arXiv preprint arXiv:2504.09858*, 2025a.
- Ma, X., Wan, G., Yu, R., Fang, G., and Wang, X. Cot-valve: Length-compressible chain-of-thought tuning. *arXiv preprint arXiv:2502.09601*, 2025b.
- MAA. American invitational mathematics examination 2024, 2024. URL https://artofproblemsolving.com/wiki/index.php/American_Invitational_Mathematics_Examination?srsltid=AfmBOoqiDCiaGTLQrsRTKsZui8RFnjOZqM4qIqY3yGn6tXt.
- MAA. American invitational mathematics examination 2025, 2025. URL https://artofproblemsolving.com/wiki/index.php/American_Invitational_Mathematics_Examination?srsltid=AfmBOoqiDCiaGTLQrsRTKsZui8RFnjOZqM4qIqY3yGn6tXt.
- Miao, X., Oliaro, G., Zhang, Z., Cheng, X., Wang, Z., Zhang, Z., Wong, R. Y. Y., Zhu, A., Yang, L., Shi, X., et al. Specinfer: Accelerating generative large language model serving with tree-based speculative inference and verification. *arXiv preprint arXiv:2305.09781*, 2023.
- Mishra, A., Latorre, J. A., Pool, J., Stosic, D., Stosic, D., Venkatesh, G., Yu, C., and Micikevicius, P. Accelerating sparse deep neural networks. *arXiv preprint arXiv:2104.08378*, 2021.
- Mohtashami, A., Pagliardini, M., and Jaggi, M. Cotformer: A chain-of-thought driven architecture with budget-adaptive computation cost at inference. *arXiv preprint arXiv:2310.10845*, 2023.
- Molchanov, D., Ashukha, A., and Vetrov, D. Variational dropout sparsifies deep neural networks. In *Proceedings of the 34th International Conference on Machine Learning*, pp. 2498–2507. PMLR, 2017.
- Muennighoff, N., Yang, Z., Shi, W., Li, X. L., Fei-Fei, L., Hajishirzi, H., Zettlemoyer, L., Liang, P., Candès, E., and Hashimoto, T. s1: Simple test-time scaling. *arXiv preprint arXiv:2501.19393*, 2025.
- Nakano, R., Hilton, J., Wu, J., et al. Webgpt: Browser-assisted question-answering with human feedback. *arXiv preprint arXiv:2112.09332*, 2021.
- Nawrot, P., Li, R., Huang, R., Ruder, S., Marchisio, K., and Ponti, E. M. The sparse frontier: Sparse attention trade-offs in transformer llms. *arXiv preprint arXiv:2504.17768*, 2025.
- Ning, X., Lin, Z., Zhou, Z., Wang, Z., Yang, H., and Wang, Y. Skeleton-of-thought: Prompting llms for efficient parallel generation. *arXiv preprint arXiv:2307.15337*, 2023.
- Pope, R., Douglas, S., Chowdhery, A., Devlin, J., Bradbury, J., Heek, J., Xiao, K., Agrawal, S., and Dean, J. Efficiently scaling transformer inference. *Proceedings of Machine Learning and Systems*, 5:606–624, 2023.
- Sadhu, R., Chen, J., Chen, Z., Tiwari, V., Lai, R., Shi, J., Yen, I. E.-H., May, A., Chen, T., and Chen, B. Magicdec: Breaking the latency-throughput tradeoff for long context generation with speculative decoding. *arXiv preprint arXiv:2408.11049*, 2024.
- Shazeer, N., Mirhoseini, A., Maziarz, K., Davis, A., Le, Q., Hinton, G., and Dean, J. Outrageously large neural networks: The sparsely-gated mixture-of-experts layer. *arXiv preprint arXiv:1701.06538*, 2017.
- Sheng, Y., Zheng, L., Yuan, B., Li, Z., Ryabinin, M., Chen, B., Liang, P., Ré, C., Stoica, I., and Zhang, C. Flexgen: High-throughput generative inference of large language models with a single gpu. In *International Conference on Machine Learning*, pp. 31094–31116. PMLR, 2023.
- Snell, C., Lee, J., Xu, K., and Kumar, A. Scaling llm test-time compute optimally can be more effective than scaling model parameters. *arXiv preprint arXiv:2408.03314*, 2024.
- Sun, H., Chen, Z., Yang, X., Tian, Y., and Chen, B. Triforce: Lossless acceleration of long sequence generation with hierarchical speculative decoding, 2024a. URL <https://arxiv.org/abs/2404.11912>.

- Sun, H., Haider, M., Zhang, R., Yang, H., Qiu, J., Yin, M., Wang, M., Bartlett, P., and Zanette, A. Fast best-of-n decoding via speculative rejection. *arXiv preprint arXiv:2410.20290*, 2024b.
- Svirschevski, R., May, A., Chen, Z., Chen, B., Jia, Z., and Ryabinin, M. Specexec: Massively parallel speculative decoding for interactive llm inference on consumer devices. *Advances in Neural Information Processing Systems*, 37:16342–16368, 2024.
- Tack, J., Lanchantin, J., Yu, J., Cohen, A., Kulikov, I., Lan, J., Hao, S., Tian, Y., Weston, J., and Li, X. Llm pretraining with continuous concepts. *arXiv preprint arXiv:2502.08524*, 2025.
- Tang, J., Zhao, Y., Zhu, K., Xiao, G., Kasikci, B., and Han, S. Quest: Query-aware sparsity for efficient long-context llm inference. *arXiv preprint arXiv:2406.10774*, 2024.
- Team, N. Sky-t1: Train your own ol preview model within \$450. <https://novasky-ai.github.io/posts/sky-t1>, 2025a. Accessed: 2025-01-09.
- Team, Q. Qwq-32b: Embracing the power of reinforcement learning, March 2025b. URL <https://qwenlm.github.io/blog/qwq-32b/>.
- Tibshirani, R. Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society: Series B (Methodological)*, 58(1):267–288, 1996.
- Tirumala, A. and Wong, R. Nvidia blackwell platform: Advancing generative ai and accelerated computing. In *2024 IEEE Hot Chips 36 Symposium (HCS)*, pp. 1–33. IEEE Computer Society, 2024.
- Wei, J., Wang, X., Schuurmans, D., Bosma, M., Xia, F., Chi, E., Le, Q. V., Zhou, D., et al. Chain-of-thought prompting elicits reasoning in large language models. *Advances in neural information processing systems*, 35:24824–24837, 2022.
- Wu, Y., Sun, Z., Li, S., Welleck, S., and Yang, Y. Inference scaling laws: An empirical analysis of compute-optimal inference for problem-solving with language models. *arXiv preprint arXiv:2408.00724*, 2024.
- Xiao, G., Tian, Y., Chen, B., Han, S., and Lewis, M. Efficient streaming language models with attention sinks, 2024. URL <https://arxiv.org/abs/2309.17453>.
- Yang, A., Yang, B., Zhang, B., Hui, B., Zheng, B., Yu, B., Li, C., Liu, D., Huang, F., Wei, H., Lin, H., Yang, J., Tu, J., Zhang, J., Yang, J., Yang, J., Zhou, J., Lin, J., Dang, K., Lu, K., Bao, K., Yang, K., Yu, L., Li, M., Xue, M., Zhang, P., Zhu, Q., Men, R., Lin, R., Li, T., Xia, T., Ren, X., Ren, X., Fan, Y., Su, Y., Zhang, Y., Wan, Y., Liu, Y., Cui, Z., Zhang, Z., and Qiu, Z. Qwen2.5 technical report. *arXiv preprint arXiv:2412.15115*, 2024.
- Yang, A., Li, A., Yang, B., Zhang, B., Hui, B., Zheng, B., Yu, B., Gao, C., Huang, C., Lv, C., Zheng, C., Liu, D., Zhou, F., Huang, F., Hu, F., Ge, H., Wei, H., Lin, H., Tang, J., Yang, J., Tu, J., Zhang, J., Yang, J., Yang, J., Zhou, J., Zhou, J., Lin, J., Dang, K., Bao, K., Yang, K., Yu, L., Deng, L., Li, M., Xue, M., Li, M., Zhang, P., Wang, P., Zhu, Q., Men, R., Gao, R., Liu, S., Luo, S., Li, T., Tang, T., Yin, W., Ren, X., Wang, X., Zhang, X., Ren, X., Fan, Y., Su, Y., Zhang, Y., Zhang, Y., Wan, Y., Liu, Y., Wang, Z., Cui, Z., Zhang, Z., Zhou, Z., and Qiu, Z. Qwen3 technical report, 2025. URL <https://arxiv.org/abs/2505.09388>.
- Yao, S., Yu, D., Zhao, J., Shafran, I., Griffiths, T., Cao, Y., and Narasimhan, K. Tree of thoughts: Deliberate problem solving with large language models. *Advances in neural information processing systems*, 36:11809–11822, 2023a.
- Yao, S., Zhao, J., Yu, D., et al. React: Synergizing reasoning and acting in language models. *Advances in Neural Information Processing Systems (NeurIPS)*, 2023b.
- Ye, Z., Chen, L., Lai, R., Lin, W., Zhang, Y., Wang, S., Chen, T., Kasikci, B., Grover, V., Krishnamurthy, A., and Ceze, L. Flashinfer: Efficient and customizable attention engine for llm inference serving. *arXiv preprint arXiv:2501.01005*, 2025. URL <https://arxiv.org/abs/2501.01005>.
- Yu, G.-I., Jeong, J. S., Kim, G.-W., Kim, S., and Chun, B.-G. Orca: A distributed serving system for Transformer-Based generative models. In *16th USENIX Symposium on Operating Systems Design and Implementation (OSDI 22)*, pp. 521–538, Carlsbad, CA, July 2022. USENIX Association. ISBN 978-1-939133-28-1. URL <https://www.usenix.org/conference/osdi22/presentation/yu>.
- Yuan, J., Gao, H., Dai, D., Luo, J., Zhao, L., Zhang, Z., Xie, Z., Wei, Y., Wang, L., Xiao, Z., et al. Native sparse attention: Hardware-aligned and natively trainable sparse attention. *arXiv preprint arXiv:2502.11089*, 2025.
- Yuan, Z., Shang, Y., Zhou, Y., Dong, Z., Zhou, Z., Xue, C., Wu, B., Li, Z., Gu, Q., Lee, Y. J., et al. Llm inference unveiled: Survey and roofline model insights. *arXiv preprint arXiv:2402.16363*, 2024.
- Zandieh, A., Han, I., Daliri, M., and Karbasi, A. Kdeformer: Accelerating transformers via kernel density estimation. In *International Conference on Machine Learning*, pp. 40605–40623. PMLR, 2023.

Zhang, S., Roller, S., Goyal, N., Artetxe, M., Chen, M.,
Chen, S., Dewan, C., Diab, M., Li, X., Lin, X. V.,
et al. Opt: Open pre-trained transformer language models.
arXiv preprint arXiv:2205.01068, 2022.

Zhang, Z., Sheng, Y., Zhou, T., Chen, T., Zheng, L., Cai,
R., Song, Z., Tian, Y., Ré, C., Barrett, C., et al. H2o:
Heavy-hitter oracle for efficient generative inference of
large language models. *Advances in Neural Information
Processing Systems*, 36:34661–34710, 2023.

Zheng, L., Yin, L., Xie, Z., Sun, C. L., Huang, J., Yu, C. H.,
Cao, S., Kozyrakis, C., Stoica, I., Gonzalez, J. E., et al.
Sglang: Efficient execution of structured language model
programs. *Advances in Neural Information Processing
Systems*, 37:62557–62583, 2024.

Appendix

Table of Contents

| | |
|--|----|
| – Cost Model | 12 |
| • Full Formulations of Cost Model | 12 |
| • Max Cost Model v.s. Additive Cost Model | 13 |
| • Details about Sparse Attention Cost Model | 14 |
| – Dense Scaling Law | 15 |
| • Additional Benchmarks | 15 |
| • Additional Reasoning Models | 15 |
| – Sparse Scaling Law | 17 |
| • Optimal Resource Allocation with Sparse Attention Models | 18 |
| • Greedy Algorithm for Optimal Resource Allocation | 18 |
| • Additional Benchmarks | 18 |
| • Additional Analysis | 19 |
| – Experimental Details | 20 |
| • Estimate Cost, Accuracy and Solving Rate | 20 |
| • Top- K Attention and Block Top- K Attention | 21 |
| – Related Work | 23 |
| – Limitations, Future Scope, and Broader Impact | 24 |

Table 1. Notation Used throughout the Paper.

| Symbol | Description | Symbol | Description |
|---------------------|-----------------|-----------|------------------|
| T, \mathcal{T} | Task (set) | L_{out} | # Gen tokens |
| M | Model | N, N_T | Reasoning trials |
| $C, C_{TTS}(\cdot)$ | Cost function | n, n_T | Max # tokens |
| \mathcal{A} | Algorithm | B, B_T | KV budget |
| L_{in} | Prompt length | P | Parameters |
| D | KV size / token | r | GQA ratio |

A Cost Model

In this section, we delve into the cost models used in the Kinetics Scaling Law. We show empirically that adopting a max cost model does not alter the scaling behavior and outline methods for calculating the cost of sparse attention models.

A.1 Full Formulations of Cost Model

We first calculate the inference cost for the cases where the batch size is 1, and then extend to a more general case in TTS. Finally, we propose our cost model using equivalent FLOPs.

Computation. As discussed in (Brown et al., 2024), the computation consists of two parts: linear modules and self-attention, which is (we assume the model is served in BFloat16.)

$$C_{\text{comp}} = \underbrace{2PL_{\text{out}}}_{\text{model parameters computation}} + \underbrace{r(2L_{\text{in}} + L_{\text{out}})L_{\text{out}}D}_{\text{self-attention}}$$

Memory Access. Memory access also consists of two parts: model parameters and KV cache.

$$C_{\text{mem}} = \underbrace{2PL_{\text{out}}}_{\text{model parameter access}} + \underbrace{2L_{\text{in}}L_{\text{out}}D}_{\text{prompt KV cache}} + \underbrace{L_{\text{out}}^2D}_{\text{decoding KV cache}}$$

In real serving scenarios, a large batch size will be used (DeepSeek-AI, 2025) with growing GPU VRAM (Tirumala & Wong, 2024) and model parallelism (Pope et al., 2023). The access to the model parameter will be amortized across requests in a batch shows parameter access time is negligible when the batch size is large). Thus, we only consider the second term (i.e., KV cache loading) in our cost function. Furthermore, in the cases that we have N reasoning trials, the prompt cache access (Juravsky et al., 2024; Zheng et al., 2024) is also shared across these N trials. Thus,

$$C_{\text{comp}}(N) = 2PNL_{\text{out}} + 2rNL_{\text{in}}L_{\text{out}}D + rNL_{\text{out}}^2D \quad (3)$$

$$C_{\text{mem}}(N) = 2L_{\text{in}}L_{\text{out}}D + NL_{\text{out}}^2D \quad (4)$$

eFLOPs. We propose eFLOPs (equivalent FLOPs) to capture both compute and memory access cost,

$$\text{eFLOPs} = C_{\text{comp}} + C_{\text{mem}} \times I \quad (5)$$

where I is the arithmetic intensity of hardware, which reflects that modern accelerators usually have a much larger computation capacity over memory bandwidth, and the gap is growing over the years (Sadhukhan et al., 2024). In this work, we use $I = 562.5$ (unit: FLOPs \times s / GB) from NVIDIA B200 (Tirumala & Wong, 2024).

With Equations (3) to (5), we obtain the final cost model.

$$C_{\text{TTS}} = \underbrace{2NPL_{\text{out}}}_{\text{linear modules computation}} + \underbrace{2rNL_{\text{in}}DL_{\text{out}} + rNDL_{\text{out}}^2}_{\text{self-attention computation}} + \underbrace{2IL_{\text{in}}DL_{\text{out}} + INDL_{\text{out}}^2}_{\text{KV access}} \quad (6)$$

where P, r, D are hyper-parameters determined by model M^1 .

A.2 Max Cost Model v.s. Additive Cost Model

Max cost model is widely used in performance modeling (Yuan et al., 2024). It assumes that computation and memory operations can be fully overlapped with each other and only considers the bottleneck operation for cost measurement.

$$C_{\text{max-cost}} = \max(C_{\text{comp}}, C_{\text{mem}} \times I)$$

where C_{comp} denotes the compute cost, C_{mem} the memory cost per access, and I the memory intensity.

In this section, we analyze the Kinetics Scaling Law using the max cost model. For clarity, we refer to the cost model $C_{\text{comp}} + C_{\text{mem}} \times I$, which is used in the main paper, as **the additive cost model**.

We draw two conclusions from empirical results **under the max cost model**:

- **Kinetics scaling law for dense models still holds.** We re-plot Figure 3(a)(b) and Figure 4a under the measurement of max cost models in Figures 11 and 12. We find except that in Long-CoTs scenarios, large models become slightly more effective in low-cost regime (with accuracy \sim 0.3), the overall trends are very close to the plots with additive cost models.
- **Sparse attention solves problems more cost-effectively.** We re-plot Figures 6a and 6b in Figures 13a and 13b. Under the max cost models, in Long-CoTs, the accuracy and efficiency gaps increase from 47.5 points and $11.21\times$ to 52.8 points and $15.71\times$, respectively. In Best-of- N , the gaps widen from 65 points and $10.67\times$ to 69.4 points and $19.64\times$. These results indicate that under the max cost model, our claim that sparse attention can enhance problem-solving performance is strengthened. Compared to dense attention models, sparse attention models tend to have more balanced memory and compute costs. Thus omitting one of them via a max cost model will favor sparse attention models.

¹Since L_{out} might differ across reasoning trials, we take the expectation for $\mathbb{E}[L_{\text{out}}]$ and $\mathbb{E}[L_{\text{out}}^2]$.

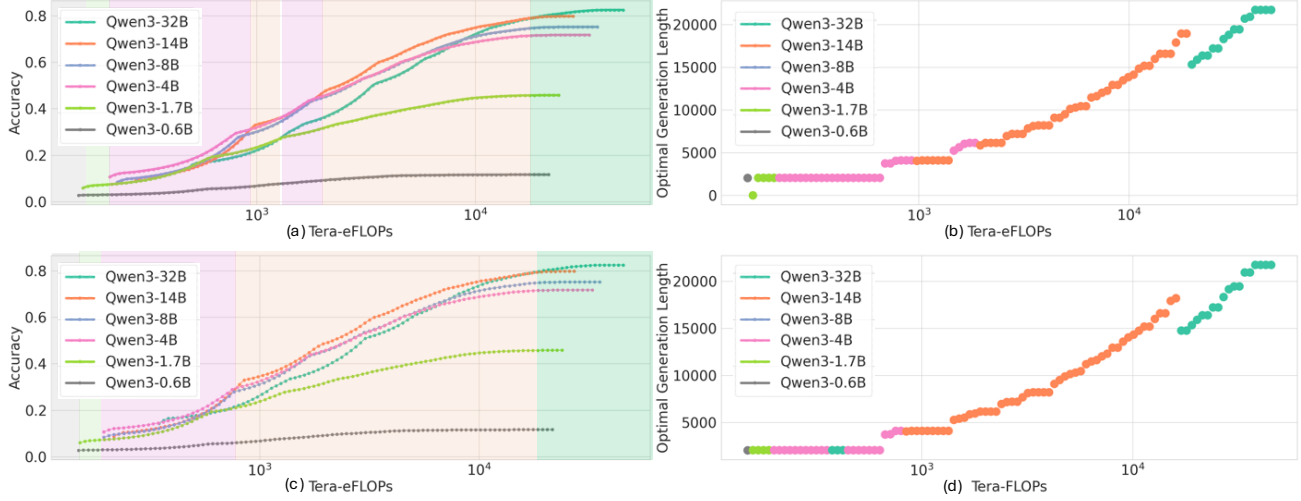


Figure 11. **AIME Pareto Frontier (Long-CoTs) with Max Cost Models.** (a)(b) is the original plot with the additive cost model. (c)(d) is the corresponding plot using max cost models. Compared to the original plots, the overall trend is similar except that larger models span a slightly broader region on the Pareto frontier. For example, the 14B model now consistently outperforms the 4B model with a noticeable gap around accuracy 0.3 and maintains dominance thereafter. In contrast, under the additive cost model in Figure 3(a), the two models alternate in performance until accuracy exceeds 0.4. This suggests that, when evaluated using a max cost model, larger models appear slightly more efficient relative to their performance under additive cost models.

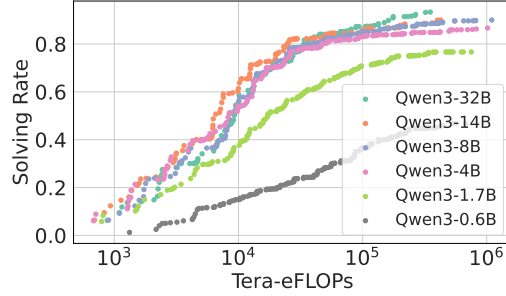


Figure 12. **AIME Pareto Frontier (Best-of-N) with Max Cost Models.** We re-plot Figure 4a using max cost models. The Pareto Frontier is very similar under different cost models.

A.3 Details about Sparse Attention Cost Model

Sparse attention models follow different cost functions due to the sparsification of KV memory access. In this paper, we focus on algorithms that impose a uniform KV budget (denoted as B) per attention head for each decoded token. We consider $L_{in} \geq B$ for the sake of simplicity. Under this setting, the cost model for sparse attention is given by:

$$C_{\text{sparse}} = \underbrace{2NPL_{\text{out}} + 2rNDBL_{\text{out}}}_{\text{compute}} + \underbrace{2INDBL_{\text{out}}}_{\text{memory}}. \quad (7)$$

In practical implementations, we must also account for the overhead associated with retrieving or searching KV memory, denoted as C_{search} , which depends on the specific sparse attention algorithm \mathcal{A} . For example, in block top- k selection, the search cost is:

$$C_{\text{search}} = \underbrace{\frac{2NL_{in}DL_{\text{out}} + rNDL_{\text{out}}^2}{2\text{Block-Size}}}_{\text{compute}} + \underbrace{\frac{2IL_{in}DL_{\text{out}} + INDL_{\text{out}}^2}{2\text{Block-Size}}}_{\text{memory}}. \quad (8)$$

In our work, we choose the Block-Size in such a way that C_{sparse} and C_{search} are roughly balanced, so that the sparse attention

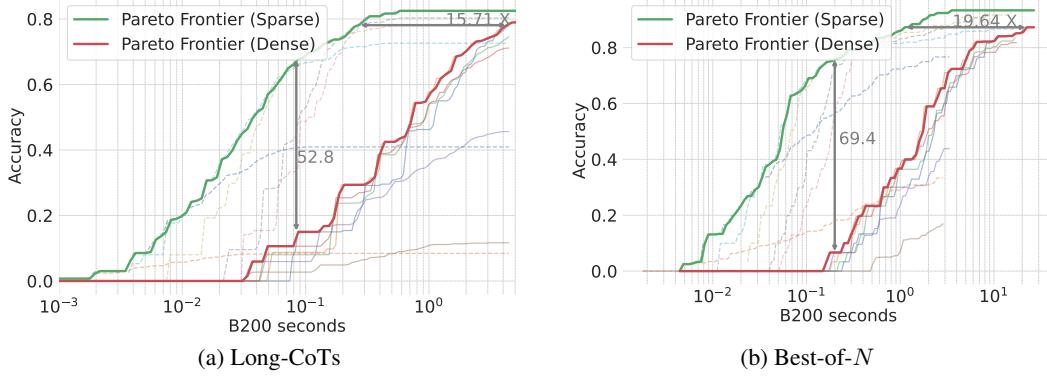


Figure 13. **Sparse attention scales significantly better under max cost models.** We re-plot Figures 6a and 6b using max cost models. Compared to the original plots, the performance and efficiency gaps between sparse attention models and dense models become more pronounced. In Long-CoTs, the accuracy and efficiency gaps increase from 47.5 points and $11.21\times$ to 52.8 points and $15.71\times$, respectively. In Best-of- N , the gaps widen from 65 points and $10.67\times$ to 69.4 points and $19.64\times$.

cost increases sub-linearly with generation length.

For local attention and oracle top- k attention, we assume no search overhead, i.e., $C_{\text{search}} = 0$.

Many sparse attention algorithms skip the first layer (Tang et al., 2024; Chen et al., 2024; Zhang et al., 2023), resulting in only a minor increase in total cost. For the Qwen3 series, this additional overhead is bounded by 3.57% for the 0.6B model and by 1.56% for the 32B model.

B Dense Scaling Law

In this section, we further verify Kinetics Scaling Law for dense models proposed in Section 3 with Iso-Cost analysis and extended experimental results of different benchmarks and model series.

B.1 Additional Benchmarks

We evaluate on AIME25 in Figures 14 and 15a to 15c and LiveCodeBench² in Figures 16 and 17a to 17c (excluding the 0.6B model), following the setting described in Section 3. The empirical results support the Kinetics Scaling Law: across both benchmarks, the 0.6B and 1.7B models are consistently less effective, and the Pareto frontier is almost always dominated by the 14B models.

B.2 Additional Reasoning Models

In Figures 18 and 19a to 19c, we evaluate DeepSeek-R1 Distilled Qwen models (abbreviated as DS models) (Guo et al., 2025) on AIME24. The DeepSeek series models further demonstrate that previous scaling laws—those based on FLOPs—significantly overestimate the effectiveness of the 1.5B model. As predicted by the Kinetics Scaling Law, increasing the number of generated tokens for the 1.5B model is less effective than scaling up the model size, such as using the 7B or larger variants.

Interestingly, we observe a shift in the emerging model size: unlike Qwen3, where the 14B model dominates, the 7B model becomes the dominant choice in the DeepSeek series. In Figures 18, 19a and 19c, the 7B model spans most of the Pareto frontier, and Figure 18 shows that 7B models with long CoTs are more efficient and effective than 14B models with short generations. We attribute this to an architectural outlier in the DeepSeek-R1 (Qwen2.5) model series. As shown in Table 2, the DeepSeek-R1 7B model is significantly more KV memory-efficient than the Qwen3-8B model. Unlike most model series illustrated in Figure 5a, where KV cache size typically grows sublinearly with respect to model parameters, DeepSeek-R1 shows a deviation from this trend: the 14B model has approximately $3.4\times$ more KV memory than the 7B model, while

²For LiveCodeBench dataset, we have sampled 50 examples from the v5 subset consisting 167 examples. Our subset comprises 24 hard, 16 medium and 10 easy examples respectively.

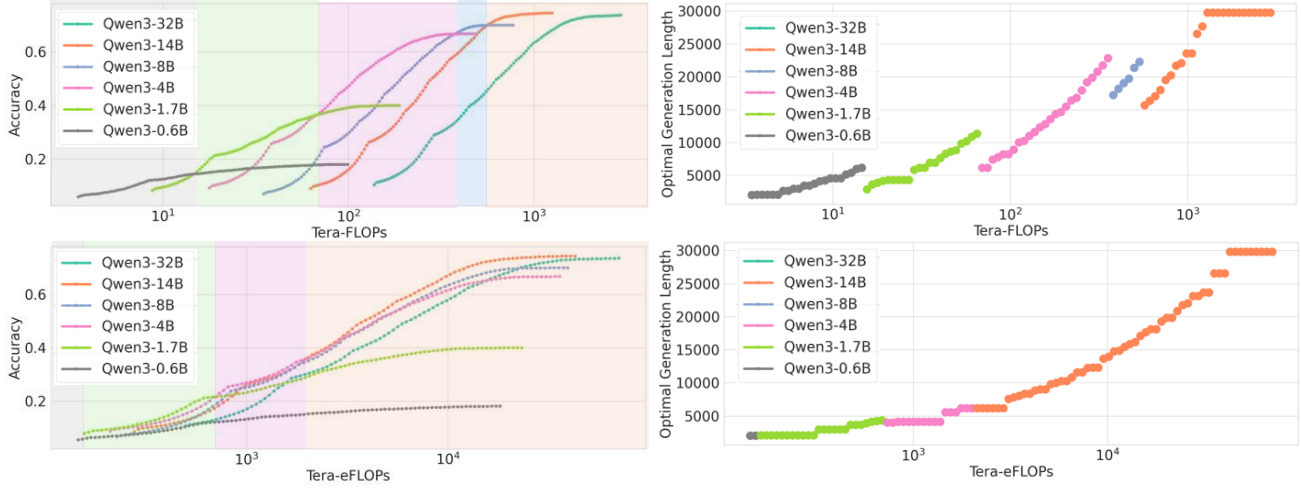


Figure 14. AIME25 Pareto Frontier (Long-CoTs). We conduct the same experiments as Figure 3.

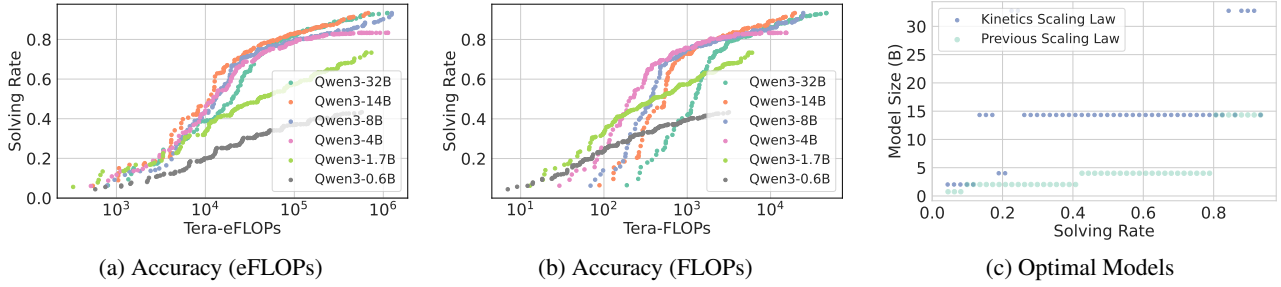


Figure 15. AIME25 Score Curve (Best-of- N). We conduct the same experiments as Figures 4a to 4c.

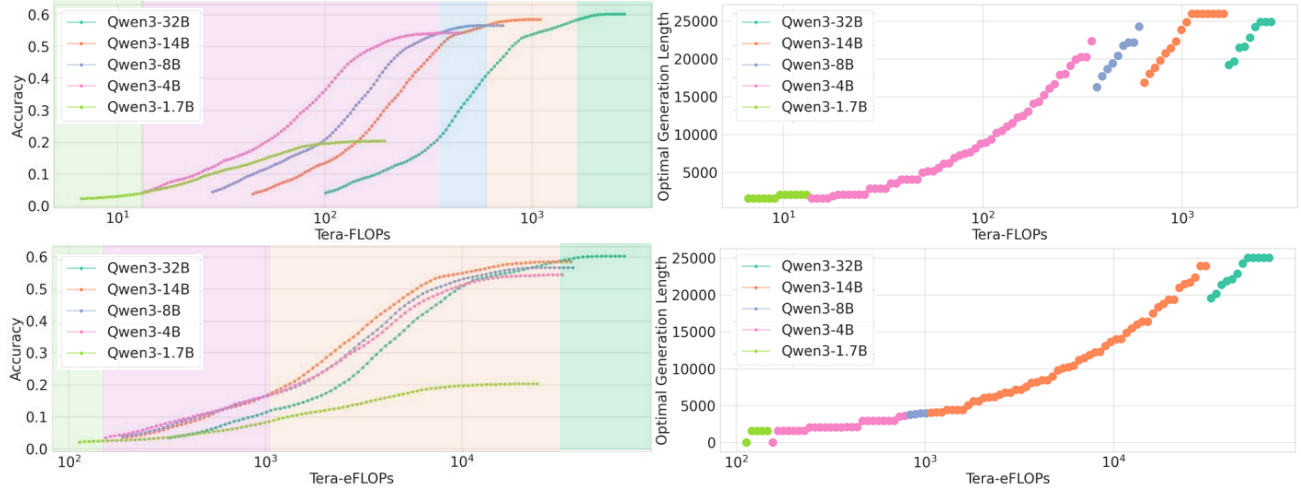


Figure 16. LiveCodeBench Pareto Frontier (Long-CoTs). We conduct the same experiments as Figure 3.

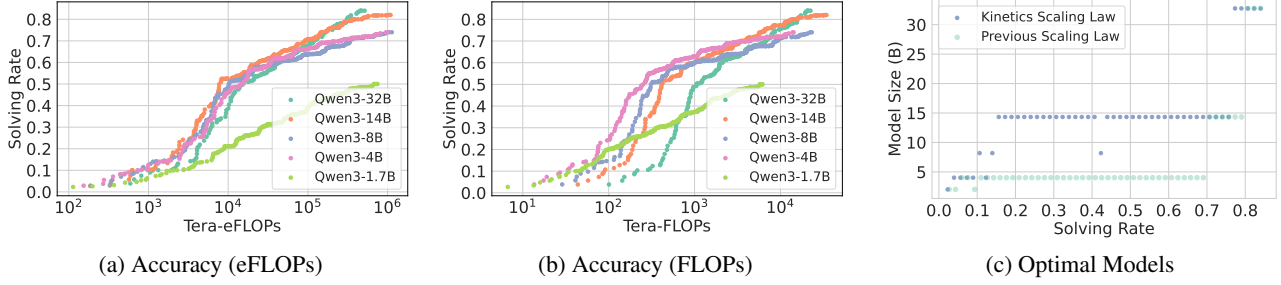


Figure 17. LiveCodeBench Score Curve (Best-of- N). We conduct the same experiments as Figures 4a to 4c.

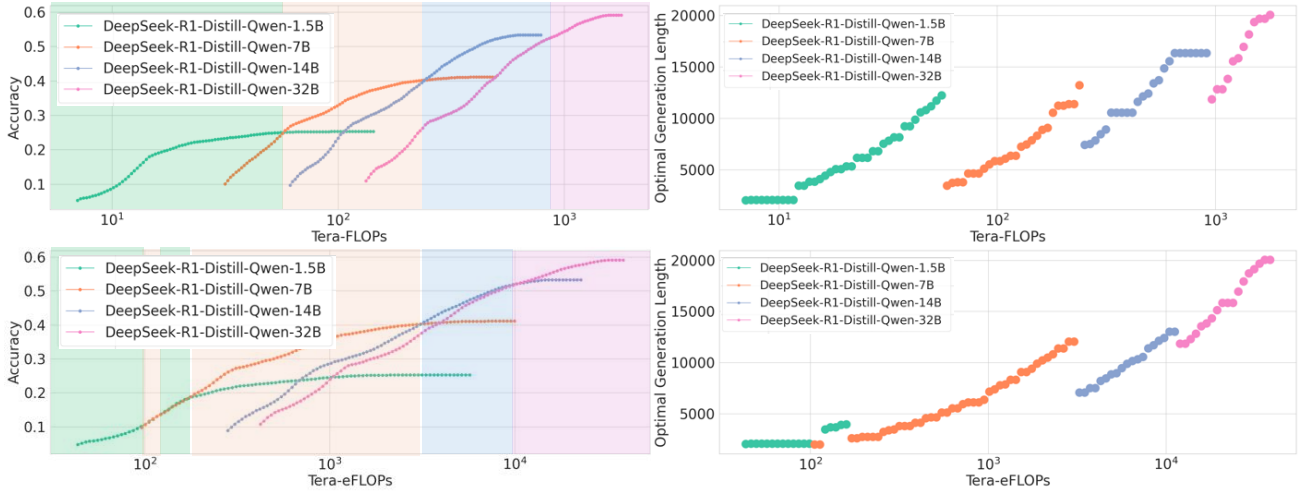


Figure 18. AIME24 Pareto Frontier (Long-CoTs). We conduct the same experiments as Figure 3 on DeepSeek Distilled Qwen series.

having only $2\times$ more parameters.

Table 2. KV memory Size for Qwen3 and DeepSeek-R1 Distilled models (per 32K tokens, unit: GB).

| Qwen3 | Qwen3-1.7B | Qwen3-8B | Qwen3-14B | Qwen3-32B |
|----------|------------|----------|-----------|-----------|
| | 3.5 | 4.5 | 6 | 8 |
| DeepSeek | DS-1.5B | DS-7B | DS-14B | DS-32B |
| | 0.875 | 1.75 | 6 | 8 |

This finding highlights the importance of concrete model architecture design, rather than focusing solely on the number of model parameters. Whether KV memory size is directly related to reasoning performance remains an open question, which we leave for future investigation.

C Sparse Scaling Law

We present how we find the Pareto frontier of sparse attention models through an optimal resource allocation, which demonstrates the upper bound of scalability of a certain sparse attention algorithms. Then we present additional results supporting the kinetics sparse scaling law across multiple tasks and demonstrate how these insights enable scalable test-time scaling with sparse attention.

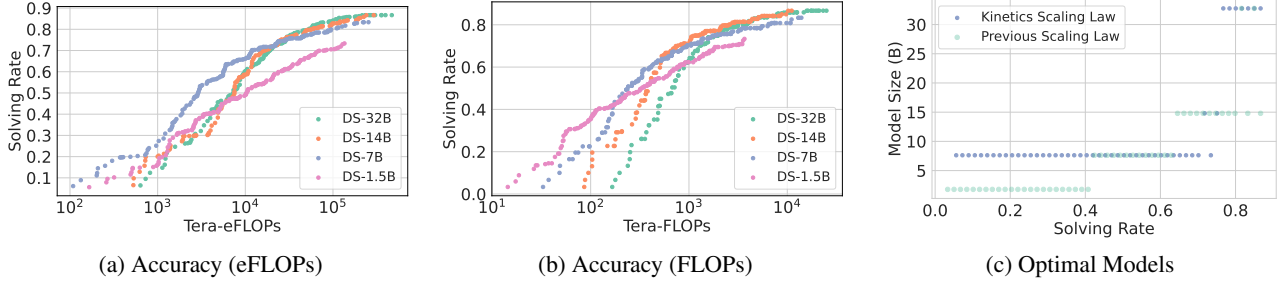


Figure 19. AIME24 Score Curve Envelope (Best-of-N). We conduct the same experiments as Figures 4a to 4c on DeepSeek Distilled Qwen series.

C.1 Optimal Resource Allocation with Sparse Attention Models

Problem statement. Let \mathcal{A} denote the corresponding sparsity patterns (e.g., top- k , block sparse and local). Our goal is to explore the optimal tradeoff among three factors: model M , KV budget B , and number of trials, and the maximum generation length (N, n) . Specifically,

$$\begin{aligned} (N, n)_*, M_*, B_* &= \arg \max_{(N, n), M, B} \text{Acc}(N, n, B, \mathcal{A}, M; T) \\ \text{s.t. } C_{\text{TTS}}(N, n, B, \mathcal{A}, M; T) &\leq C \end{aligned} \quad (9)$$

C.2 Greedy Algorithm for Optimal Resource Allocation

We present a method to optimally schedule generation parameters (N, n) and the KV budget B for each task, establishing an upper bound on achievable performance and enabling analysis of the core tradeoff between TTS strategies and sparsity. We begin by solving the subproblem for each individual task T^3 :

$$\max \text{Acc}(N_T, n_T, B_T, \mathcal{A}, M; T) \quad \text{s.t. } C_{\text{TTS}}(N_T, n_T, B_T, \mathcal{A}, M; T) \leq C \quad (10)$$

Empirically, we discretize the searching space. For instance, in *Best-of-N*, we discretize the space of N and B by producing a search grid:

$$G = \{N_0, N_1, \dots, N_i\} \otimes \{B_0, B_1, \dots, B_j\}$$

For each pair $(N_a, B_b) \in G$, we compute the corresponding cost $C_{T, (a, b)}$ and accuracy $\text{Acc}_{T, (a, b)}$. We use $(N_T, B_T) \in G$ which maximizes the accuracy under the cost constraint C as an approximation for Equation (10). By combining the optimal configurations (N_T, B_T) for all tasks T , we obtain a solution to the overall problem in Equation (9). Similar discretizations also applies for *Long-CoTs*. Thus we find the optimal resource allocation.

We describe the procedure for identifying optimal resource allocations and establishing the Pareto frontier for sparse attention models in Algorithms 1 and 2, as a supplement to Appendix C.1. Given a fixed cost constraint C , we perform a grid search over key parameters: KV budgets and either reasoning trials or maximum generation lengths.

Empirically, we sweep over KV budgets $\{32, 64, 128, 256, 512, 1024\}$; reasoning trials $\{1, 2, 4, 8, 16, 32\}$ (with a reduced upper limit for the 14B and 32B models to save computation time); and generation lengths $\{2k, 4k, 6k, 8k, 10k, 12k, 14k, 16k, 18k, 20k, 22k, 24k, 26k, 28k, 30k, 32k\}$.

It is important to note that we do not consider inter-request resource scheduling strategies, such as early stopping or dynamic reallocation across requests (Fu et al., 2024), since we aim to ensure fairness across all inputs. Instead, the cost constraint C is interpreted as the maximum allowable cost per request (not the average), even if some requests achieve saturated accuracy below that threshold.

C.3 Additional Benchmarks

Beyond AIME24, we evaluate our approach on LiveCodeBench (Jain et al., 2024) and AIME25 (MAA, 2025). LiveCodeBench features complex programming problems from recent coding contests, while AIME25 consists of challenging

³For fairness, we do not schedule resources across tasks, but consider a resource upper bound for all the tasks.

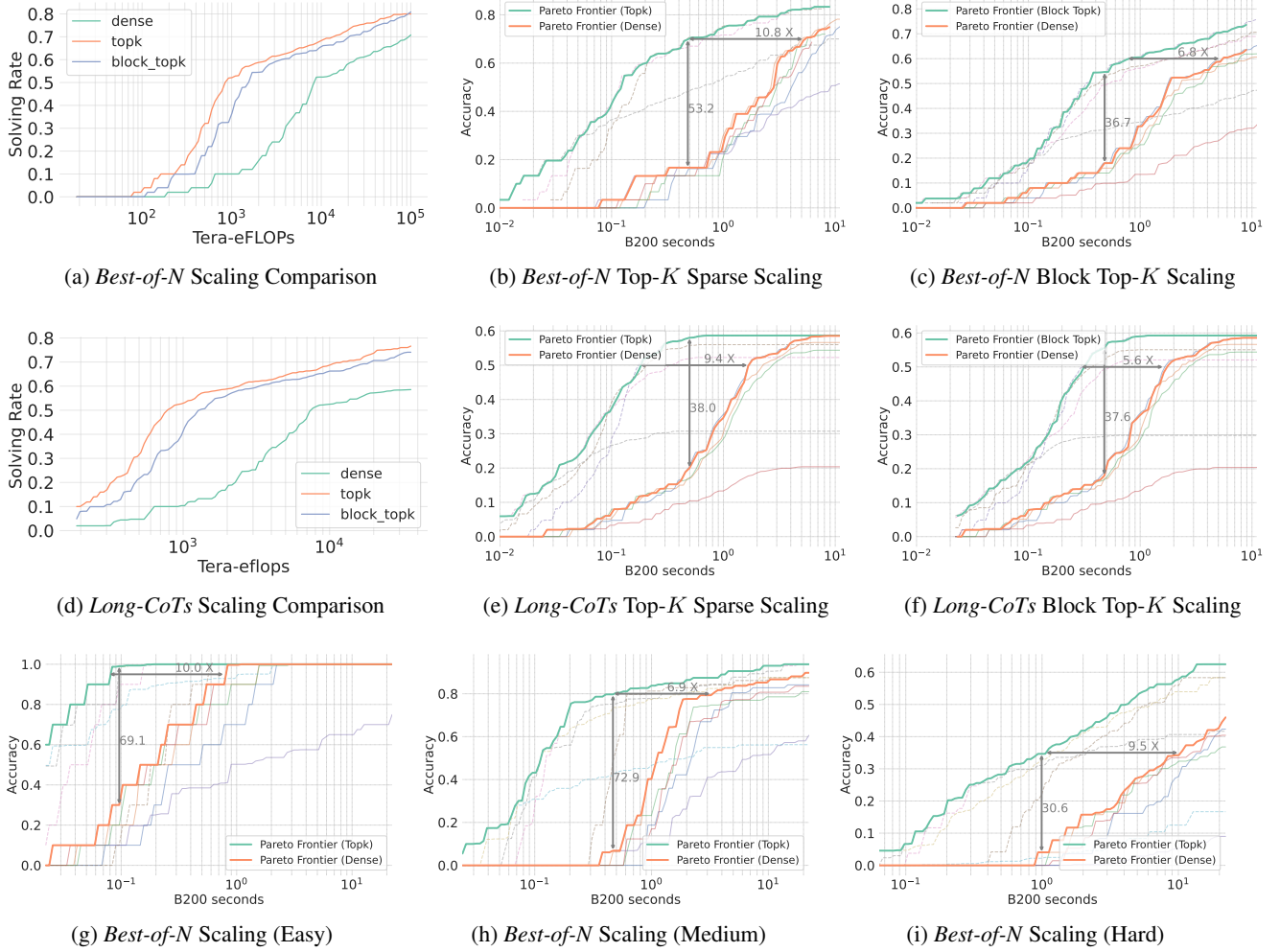


Figure 20. LiveCodeBench Sparse Scaling. We evaluate sparse scaling laws for Qwen3-14B model using oracle top- k and block-top- k attention on the LiveCodeBench dataset. (a)(d) compare block-top- k and oracle top- k with dense scaling under *Best-of-N* and *long-CoT* TTS settings. (b)(e) show cost-accuracy trade-offs for top- k attention. (c)(f) show trade-offs for block-top- k attention. (g)(h)(i) compare the oracle top- k scaling for easy, medium and hard difficulty questions.

math problems. In both cases, sparse attention—particularly oracle top- k —consistently outperforms dense attention. Block top- k attention, a tractable alternative, closely matches the performance of the oracle.

For LiveCodeBench, we sample 50 problems from the v5 subset (24 hard, 16 medium, 10 easy). As shown in Figure 20, oracle top- k attention can achieve $\sim 10\times$ speedup in high-accuracy regimes and improves coverage by 40–50% in low-cost regimes. Conversely, the tractable alternative, Block top- k yields 5–6 \times speedup and 30–40% coverage gains. We further show how the benefits of sparse attention scale with problem difficulty (Figures 20g to 20i).

Figure 21 confirms similar trends for AIME25, with substantial gains in both accuracy and efficiency under sparse attention.

C.4 Additional Analysis

Fixing a model (e.g., Qwen3-8B), we investigate the tradeoff between generating more tokens through *Best-of-N* and increasing the KV budget in Figures 22a to 22d. As the figures suggest, on AIME25, each doubling of total compute cost increases the optimal KV budget by 1.13 \times , while generated tokens grow by 1.67 \times ; on LiveCodeBench, these factors are 1.14 \times and 1.89 \times , respectively. We find that although the concrete numbers depend on the types of tasks, the overall results confirm our suggestions in the main paper that allocating compute toward generating more responses is generally more

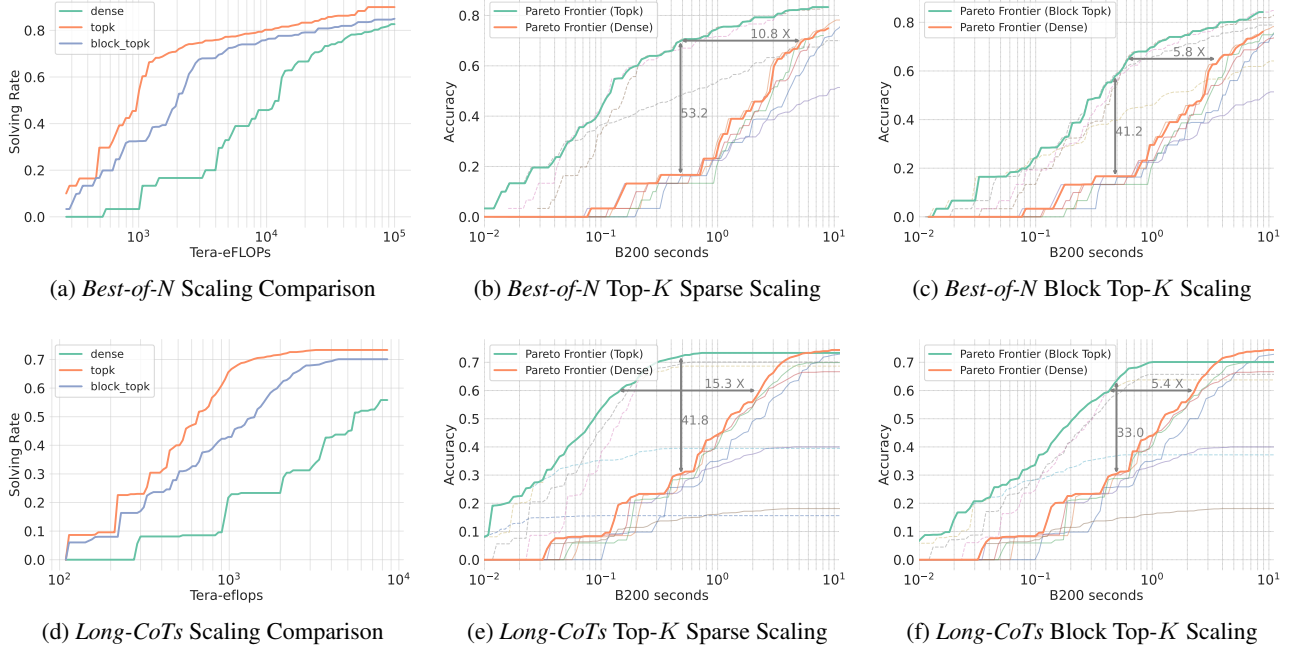


Figure 21. **AIME25 Sparse Scaling.** We evaluate sparse scaling laws for Qwen3-14B model using oracle top- k and block-top- k attention on the AIME25 dataset. (a)(d) compare block-top- k and oracle top- k with dense scaling under *Best-of-N* and *long-CoT* settings. (b)(e) show cost-accuracy trade-offs for oracle top- k attention. (c)(f) show trade-offs for block-top- k attention.

effective than expanding KV budget, highlighting the scalability of sparse attention.

D Experimental Details

In this section, we explain the details about our experiments.

D.1 Estimate Cost, Accuracy and Solving Rate

When empirically measuring cost, one major challenge is the difficulty of controlling the actual generation length. Although it is possible to set an upper bound on the number of generated tokens, there is no guarantee that the model will utilize the full budget. For instance, in our *Best-of-N* experiments, we set the maximum number of generated tokens to 32,768, yet the average generation length was only 14K–16K tokens.

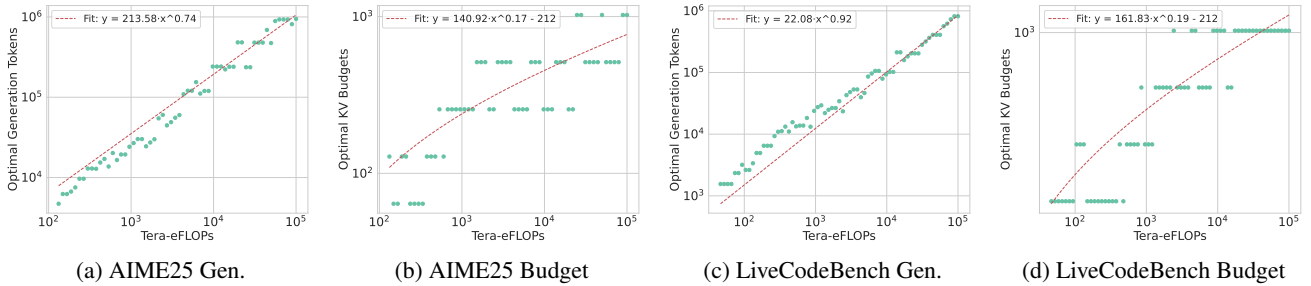


Figure 22. **Tradeoff Between Generated Tokens and KV Budget.** We empirically characterize the tradeoff between increasing generation length and allocating a larger KV cache budget using Qwen3-8B. For AIME25 ((a)(b)) and LiveCodeBench ((c)(d)), we identify the optimal KV budget and generated tokens (defined as number of reasoning trials times the average generated tokens per trial) to achieve the highest problem-solving rate under every cost constraint C .

Algorithm 1 Best-of- N optimal resource allocation under cost C

Data: Tasks \mathcal{T} , KV budgets $\{B_1, \dots, B_j\}$, trial counts $\{N_1, \dots, N_i\}$, cost limit C

Result: Average of maximum accuracy per task under cost C

```

1 AccumBestAcc  $\leftarrow$  0 Count  $\leftarrow$  0 for task  $T$  in  $\mathcal{T}$  do
2   for KV budget  $B_b$  do
3     Generate  $S \geq \max\{N_1, \dots, N_i\}$  responses using  $B_b$  for task  $T$  for trial count  $N_a$  do
4       compute cost  $c_{b,a}^{(T)}$  if  $c_{b,a}^{(T)} \leq C$  then
5         Compute accuracy  $\text{Acc}_{b,a}^{(T)} = \text{Pass}@N_a$ ;
6         if  $\text{Acc}_{b,a}^{(T)} > \text{BestAcc}$  then
7           BestAcc  $\leftarrow$   $\text{Acc}_{b,a}^{(T)}$ ;
8         end if
9       end if
10    end for
11  end for
12  AccumBestAcc  $+=$  BestAcc; Count  $+=$  1;
13 end for
14 AvgBestAcc = AccumBestAcc/Count return AvgBestAcc
    
```

Furthermore, it is important to model the relationship between actual inference cost and performance metrics, such as accuracy in Long-CoTs or solving rate in Best-of- N . Relying solely on the maximum allowed generation length to estimate cost can substantially underestimate the efficiency of models that solve problems with much shorter responses—an ability that **may** reflect higher capability.

To address this challenge, we first sample S independent reasoning traces r_1, r_2, \dots, r_S from model M on task T , with the maximum allowed number of tokens set to n . We slightly generalize Equation (6) as:

$$\begin{aligned}
 C_{\text{TTS}} &= 2NP\mathbb{E}[L_{\text{out}}] + 2rNL_{\text{in}}D\mathbb{E}[L_{\text{out}}] + rND\mathbb{E}[L_{\text{out}}^2] \\
 &\quad + 2IL_{\text{in}}D\mathbb{E}[L_{\text{out}}] + IND\mathbb{E}[L_{\text{out}}^2] \\
 &= a\mathbb{E}[L_{\text{out}}] + b\mathbb{E}[L_{\text{out}}^2] + c,
 \end{aligned} \tag{11}$$

where a , b , and c are constants determined by the model architecture and test-time strategies (e.g., the value of n). The expectations are estimated from the sampled traces, whose distribution is influenced by the model M , the token limit n , and the task T .

For Long-CoTs, we fix $N = 1$ in Equation (11) and vary n . From the sampled traces, we estimate the accuracy (Pass@1), and compute the corresponding cost by substituting the empirical values of $\mathbb{E}[L_{\text{out}}]$ and $\mathbb{E}[L_{\text{out}}^2]$ measured under each n .

For Best-of- N , we fix $n = 32,768$, and estimate the solving rate (Pass@ K) following the methodology of Brown et al. (2024). The corresponding cost is then computed by substituting $N = K$ into Equation (11).

Similarly, we can estimate the cost for sparse attention models using Equations (7) and (8).

Advanced control of generation lengths is an active area of research (Yang et al., 2025; Muennighoff et al., 2025; Ma et al., 2025a), but it is beyond the scope of this paper.

D.2 Top- K Attention and Block Top- K Attention

In this section, we explain the sparse attention algorithms discussed in the main paper, namely *Top- K Attention* and *Block Top- K Attention*.

During the decoding phase of a large language model (LLM), the self-attention mechanism computes a weighted average of past values as follows:

$$o = \text{Softmax}\left(\frac{qK^\top}{\sqrt{d}}\right)V = wV, \quad q \in \mathbb{R}^{1 \times d}, \quad K, V \in \mathbb{R}^{n \times d}, \quad w \in \mathbb{R}^{1 \times n}, \tag{12}$$

Algorithm 2 Long-CoTs optimal resource allocation under cost C

Data: Tasks \mathcal{T} , KV budgets $\{B_1, \dots, B_j\}$, gen. lengths $\{n_1, \dots, n_i\}$, samples S , cost limit C
Result: Average of maximum accuracy per task under cost C

```

1155 Algorithm 2 Long-CoTs optimal resource allocation under cost  $C$ 
1156 Data: Tasks  $\mathcal{T}$ , KV budgets  $\{B_1, \dots, B_j\}$ , gen. lengths  $\{n_1, \dots, n_i\}$ , samples  $S$ , cost limit  $C$ 
1157 Result: Average of maximum accuracy per task under cost  $C$ 
1158 14 AccumBestAcc  $\leftarrow 0$  Count  $\leftarrow 0$  for task  $T$  in  $\mathcal{T}$  do
1159   15 BestAcc  $\leftarrow 0$  for gen. length  $n_a$  do
1160     16 for KV budget  $B_b$  do
1161       17 Generate  $S$  responses using  $(B_b, n_a)$ ; compute cost  $c_{b,a}^{(T)}$  if  $c_{b,a}^{(T)} \leq C$  then
1162       18 Compute accuracy  $\text{Acc}_{b,a}^{(T)} = \text{Pass}@1$ ;
1163       19 if  $\text{Acc}_{b,a}^{(T)} > \text{BestAcc}$  then
1164         20 BestAcc  $\leftarrow \text{Acc}_{b,a}^{(T)}$ ;
1165       21 end if
1166     22 end for
1167   23 end for
1168   AccumBestAcc  $+=$  BestAcc; Count  $+= 1$ ;
1169 24 end for
1170 AvgBestAcc = AccumBestAcc/Count return AvgBestAcc
1171

```

where d is the head dimension and n is the context length. The key and value matrices are given by $K = [k_1, k_2, \dots, k_n]$, $V = [v_1, v_2, \dots, v_n]$, where each $k_i, v_i \in \mathbb{R}^{1 \times d}$ are cached from previous decoding steps.

Top- K Attention. Top- K Attention is a sparsification method where only the K most relevant tokens (i.e., those with the highest attention scores) are selected to compute the output. Formally, instead of computing the full softmax, we define a sparse attention weight vector:

$$w_i = \begin{cases} \frac{\exp(s_i)}{\sum_{j \in \mathcal{I}_K} \exp(s_j)} & \text{if } i \in \mathcal{I}_K, \\ 0 & \text{otherwise,} \end{cases} \quad \text{where } s_i = \frac{qk_i^\top}{\sqrt{d}}, \quad \mathcal{I}_K = \text{TopK}_K(s), \quad (13)$$

Here, \mathcal{I}_K denotes the indices of the top K attention scores s_i . By masking out the less important positions, this approach reduces the computational and memory cost of attention from $\mathcal{O}(n)$ to $\mathcal{O}(K)$, where $K \ll n$.

Block Top- K . Block Top- K Attention is a block-level sparse attention mechanism. Instead of selecting individual tokens based on attention scores, this method selects entire blocks of tokens, thereby reducing the number of attention computations.

Specifically, assume the full sequence of n keys is divided into $m = \frac{n}{\text{BLOCK_SIZE}}$ consecutive blocks, each of size BLOCK_SIZE:

$$K = [k_1, \dots, k_n] \rightarrow \{K_1, K_2, \dots, K_m\}, \quad K_i \in \mathbb{R}^{\text{BLOCK_SIZE} \times d}$$

For each block K_i , we first compute the average key vector:

$$\bar{k}_i = \frac{1}{\text{BLOCK_SIZE}} \sum_{j=1}^{\text{BLOCK_SIZE}} k_{i,j}$$

Next, we compute the attention score between the query q and each block's average key:

$$s_i = \frac{q\bar{k}_i^\top}{\sqrt{d}}, \quad \text{for } i = 1, 2, \dots, m$$

We then select the top $K' = \frac{K}{\text{BLOCK_SIZE}}$ blocks based on the scores s_i , denoted by the index set $\mathcal{J}_{K'} = \text{TopK}_{K'}(s)$. Attention is computed only over the tokens within the selected blocks. The sparse attention weights are defined as:

$$w_i = \begin{cases} \frac{\exp(s_i)}{\sum_{j \in \mathcal{I}_K} \exp(s_j)} & \text{if } i \in \mathcal{I}_K \subseteq \text{tokens in selected blocks,} \\ 0 & \text{otherwise} \end{cases}$$

For both algorithms, K is the KV budget. For GQA, we conduct an average pooling across all the query heads in a group, ensuring that the total number of retrieved key-value vectors does not exceed the allocated KV budget.

Implementation. Here we provide details of our block top- k attention implementation. We build our inference backend on Flashinfer (Ye et al., 2025), incorporating support for paged attention (Kwon et al., 2023) and continuous batching (Yu et al., 2022). Alongside the paged KV cache, we introduce an auxiliary data structure to store block-level average key vectors. The KV block size is chosen such that the memory load from the block-average vectors and the selected top- k KV blocks remains balanced. This design enables sub-quadratic KV loading cost as the number of reasoning tokens increases.

E Related Work

Efficient Attention. Sparse attention (Kitaev et al., 2020; Zandieh et al., 2023; Chen et al., 2021; 2024; Zhang et al., 2023; Xiao et al., 2024; Yuan et al., 2025; Nawrot et al., 2025; Child et al., 2019; Li et al., 2024; Cai et al., 2024) has been comprehensively studied to reduce the attention cost when processing long sequences. In parallel, approaches like FlashAttention (Dao et al., 2022; Dao, 2023) accelerate attention by maximizing hardware efficiency. To address the quadratic complexity of standard attention, researchers have also explored linear attention architectures (Gu & Dao, 2023; Gu et al., 2022; Katharopoulos et al., 2020; Choromanski et al., 2020). Additionally, quantization and low-precision methods (Liu et al., 2024; Hooper et al., 2024; Lin et al., 2024b) have been broadly applied for improving inference efficiency.

Efficient Inference. Orca (Yu et al., 2022), vLLM (Kwon et al., 2023), and SGLang (Zheng et al., 2024) are widely adopted to enhance the efficiency of LLM serving. Our analysis builds on the practical designs and implementations of these systems. In parallel, speculative decoding (Leviathan et al., 2023; Chen et al., 2023; Miao et al., 2023; Sadhukhan et al., 2024) has been proposed to mitigate the memory-bandwidth bottleneck during LLM decoding. Additionally, model compression and offloading (Dettmers et al., 2022; Lin et al., 2024a; Svirschevski et al., 2024; Sheng et al., 2023; Frantar et al., 2022) techniques are playing a crucial role in democratizing LLM deployment.

Efficient Test-time Strategies. Optimizing reasoning models to generate fewer tokens has been shown to directly reduce inference-time cost (Team, 2025a; Arora & Zanette; Ma et al., 2025b). Recent work such as CoCoNut (Hao et al., 2024) and CoCoMix (Tack et al., 2025) explores conducting reasoning in a latent space, thereby reducing decoding time. Methods like ParScale (Chen et al., 2025b), Tree-of-Thoughts (Yao et al., 2023a), and Skeleton-of-Thoughts (Ning et al., 2023) aim to improve efficiency by enabling parallel reasoning. Architectural innovations such as CoTFormer (Mohtashami et al., 2023) further enhance efficiency by adaptively allocating computational resources across tokens. Efficient reward-model-based (Wu et al., 2024; Snell et al., 2024; Sun et al., 2024b) test-time scaling algorithms are also comprehensively studied.

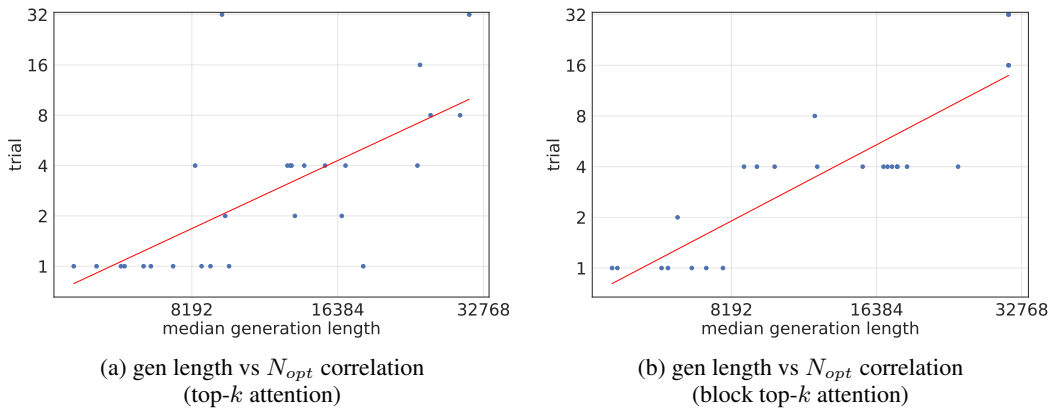


Figure 23. **Correlation between Generation Length and Number of Trials.** Longer generations correlate strongly with the optimal number of trials (N_{opt}), serving as a proxy for problem difficulty. (a) shows this trend for top- k and block top- k attention on the AIME24 dataset using the Qwen3-8B model.

F Limitations, Future Scope, and Broader Impact

Limitations. Our experiments primarily focus on **Qwen3** (Yang et al., 2025) and **DeepSeek-R1-Distilled-Qwen** (Guo et al., 2025), two state-of-the-art pretrained reasoning model series, evaluated from the inference perspective. However, the effects of training and post-training strategies are not fully explored and may influence the performance gaps and robustness to sparse attention mechanisms. In addition, our cost analysis assumes a cloud-based serving environment, where computational resources are typically sufficient and large batch sizes are feasible. In contrast, local deployment scenarios, such as those using Ollama⁴, often face limited VRAM where access to model parameters can dominate inference costs. Smaller models may be more appropriate in such settings, and our findings may not fully extend to these use cases.

Future Scope. Our sparse scaling law offers valuable insights for enriching the applications of sparse attention algorithms and the design space of test-time scaling strategies. On one hand, except for top- k , currently we only discuss a simple variant, i.e., block top- k , and have already demonstrated strong scalability. More advanced sparse attention algorithms (Tang et al., 2024; Chen et al., 2024; Yuan et al., 2025; Lin et al., 2025) are emerging these days. We do believe they can eventually push the scalability of test-time scaling to a much higher boundary. On the other hand, test-time scaling algorithms are proposed to adaptively allocate computation to tasks, or even to tokens (Arora & Zanette; Mohtashami et al., 2023; Ma et al., 2025b;a). Extending them towards to new resource allocation problems in sparse attention is critical to reach the limit of Kinetics sparse scaling law. For instance, since generation length strongly correlates with the optimal number of trials under sparse attention (as shown in Figure 23), it can be used as a dynamic signal to adjust the number of trials and KV budget. Moreover, sparse attention drastically reduces inference cost, enabling more reasoning trials and longer generations. This unlocks greater flexibility in configuring TTS strategies within a fixed resource budget.

Broader Impact. This work aims to contribute to the understanding of efficiency and scalability challenges in the test-time scaling era, spanning model architecture, system-level implementation, and hardware design. We highlight the central role of sparsity in addressing these challenges. Our study is algorithmic in nature and does not target specific applications. While large language models can be misused in harmful ways, this work does not introduce new capabilities or risks beyond those already present in existing systems. Test-time scaling can consume a substantial amount of energy, raising concerns about the environmental sustainability of widespread deployment. By promoting sparse attention, our work hopes to help to reduce the carbon footprint and energy consumption of inference systems and support the broader goal of sustainable AI.

⁴<https://github.com/ollama/ollama>