

Every Token Counts: Generalizing 16M Ultra-Long Context in Large Language Models

Anonymous ACL submission

Abstract

This work explores efficient ultra-long context modeling. We posit that an effective solution requires three fundamental properties: **sparsity**, **random-access flexibility**, and **length generalization**. To achieve this, we leverage Hierarchical Sparse Attention (HSA), a novel attention mechanism that satisfies all three properties. We integrate HSA into the Transformer architecture to develop HSA-UltraLong, an 8B-parameter Mixture-of-Experts (MoE) model trained on over 8 trillion tokens. We rigorously evaluate the model across tasks with both in-domain and out-of-domain context lengths to validate its capabilities. Our model demonstrates comparable performance to full-attention baselines on in-domain sequence lengths. Crucially, it achieves over 90% accuracy on most in-context retrieval tasks with contexts up to 512 times the pre-training context length. This work reports our findings and remaining issues throughout the experiments, offering insights for future research in ultra-long context modeling.

1 Introduction

Despite the impressive capabilities of Large Language Models (LLMs) (Brown et al., 2020; Achiam et al., 2023; Touvron et al., 2023), their world knowledge is confined to static parameters, making it inflexible to update and impossible to learn dynamically from daily user interactions. This limitation motivates a fundamental question: how can we build machines that truly remember? Effective memory is critical for future AI agents, enabling each user to have a personalized agent that accumulates unique experiences over time. Human memory spans the entire context from birth to the present, suggesting that the problem of machine memory is closely related to ultralong context modeling. Imagine if Transformers could efficiently handle infinite-length contexts—encompassing all pre-trained tokens—so that most world knowledge

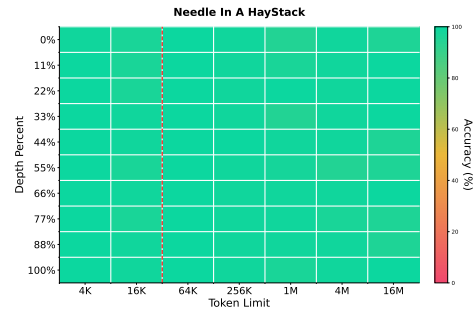


Figure 1: Despite being pre-trained with an 8K context window and mid-trained up to 32K, HSA-UltraLong achieves near-perfect accuracy on S-NIAH even at a 16M-token context length. The red dashed line at 32K marks the boundary between in-domain (left) and out-of-domain (right).

can be retrieved from context rather than compressed into model parameters. Furthermore, skills and the latest information could be acquired via in-context learning rather than through costly model retraining. Such advances would dramatically improve the online learning of knowledge and skills.

However, the Transformer (Vaswani et al., 2017) architecture, the backbone of modern LLMs, faces a fundamental efficiency challenge when processing ultra-long sequences, due to both poor length generalization and the quadratic computational complexity of full attention. Supporting longer contexts requires training models with extended context windows, yet simply scaling context length is computationally prohibitive. If we consider the extreme case of extending ultra-long context modeling to infinite context modeling, the following three points become necessary:

- **Sparsity:** Drawing inspiration from human long-term memory, which operates via selective activation and retrieval of relevant fragments (Cowan, 2008), full-attention for infinitely long contexts is clearly infeasible. Therefore, sparsity is a necessary prerequisite.
- **Random-Access Flexibility:** The utility of spar-

068	sity is predicated on the accurate retrieval of relevant past information. This necessitates designing an <i>intrinsic retrieval mechanism</i> within the model and optimizing it end-to-end under the guidance of an auto-regressive loss.	119
069		120
070		121
071		122
072		123
073	• Length Generalization: Pretraining with an infinite context is impossible. To achieve the goal, the path must involve generalizing retrieval ability from short to long contexts.	124
074		125
075		126
076		127
077	While several approaches show promising paths to achieve the goal, each presents notable shortcomings. Recurrent architectures, such as Mamba (Gu and Dao, 2023; Dao and Gu, 2024) and Linear Attentions (Katharopoulos et al., 2020; Yang et al., 2025b), compress past variable-length information into a fixed-dimensional state vector. This introduces an information bottleneck and sacrifices random access to distant tokens. Similarly, sliding-window attention (Beltagy et al., 2020) suffers from the same fundamental constraint on distant context accessibility. Sparse attention approaches like NSA (Yuan et al., 2025) and MoBA (Lu et al., 2025) improve training and inference efficiency over long sequences, but our empirical studies show they suffer from inaccurate chunk selection, which leads to both in-domain and out-of-domain performance degradation on in-context retrieval tasks.	128
078		129
079		130
080		131
081		132
082		133
083		134
084		135
085		136
086		137
087		138
088		139
089		140
090		141
091		142
092		143
093		144
094		145
095	A recent line of work that combines model-inherent retrieval (Mohtashami and Jaggi, 2023; Hu et al., 2025b) with chunk-wise sparse attention, such as Hierarchical Sparse Attention (HSA) (Hu et al., 2025a), has shown promising results in long-context modeling. Empirical studies (Leng et al., 2025) report that an HSA-based model pre-trained with a 4K context length can extrapolate to more than 10M context length while keeping high accuracy on the RULER (Hsieh et al., 2024) and BabiLong (Kuratov et al., 2024) benchmark, which simultaneously satisfies sparsity , random-access flexibility and length generalization . The method partitions text into fixed-length chunks with landmark representations; each token retrieves top-k relevant past chunks via these landmarks. The core innovation of HSA is to conduct attention with each chunk <i>separately</i> , and then <i>fuse the results weighted by the retrieval scores</i> . The overall process closely resembles the Mixture-of-Experts (MoE) (Shazeer et al., 2017), as illustrated in Figure 2. This design allows the retrieval scores to be integrated into the forward pass, enabling them to receive gradient updates during backpropagation.	146
096		147
097		148
098		149
099		150
100		151
101		152
102		153
103		154
104		155
105		156
106		157
107		158
108		159
109		160
110		161
111		162
112		163
113		164
114		165
115		166
116		167
117		
118		
	As a result, the model learns to assign higher retrieval scores to chunks that are more helpful for next token prediction. However, current work in this area is limited in scale and lacks results on data and parameter scaling.	
	We introduce HSA-UltraLong, an architecture combining sliding-window attention with HSA, and validate its effectiveness by training a 0.5B dense model and an 8BA1B MoE model from scratch on 8 trillion tokens. Through long-context extension and annealing, we verify that these models significantly enhance out-of-domain long-context capabilities while preserving in-domain performance without degradation. Our key findings include:	
	• Effective length generalization requires the combination of chunk-wise attention, retrieval score-based fusion, and NoPE (No Positional Encoding); all three are essential.	
	• Sliding-window attention and HSA interact in nontrivial ways. HSA’s long-range generalization arises from learning to retrieve over short contexts and transferring that ability to long contexts. However, an overly large sliding window can weaken HSA’s learning of short-range dependencies, degrading generalization.	
	• The effective context length in the training corpus strongly influences length extrapolation.	
	This work makes the first empirical demonstration of training-free length generalization—scaling from a 32K pretraining context to 16M tokens—in a model with billions of parameters trained on trillions of tokens. It details the training recipe and associated challenges, offering key insights for modeling ultra-long contexts.	
	2 Preliminary	
	2.1 Limitations of Chunk Selection in NSA	
	NSA is a highly inspiring contribution to sparse attention. However, our empirical study in Table 1 shows that its chunk selection mechanism does not always pick the most relevant chunk. On the RULER benchmark (Hsieh et al., 2024), NSA fails to achieve perfect accuracy even on in-domain context lengths such as Multi-Query NIAH. We trace this to a key limitation: the inaccurate chunk selection action. Further, our analysis of NSA’s extrapolation ability indicates that performance on in-context retrieval degrades rapidly as context length increases. Regarding positional encoding, we also	

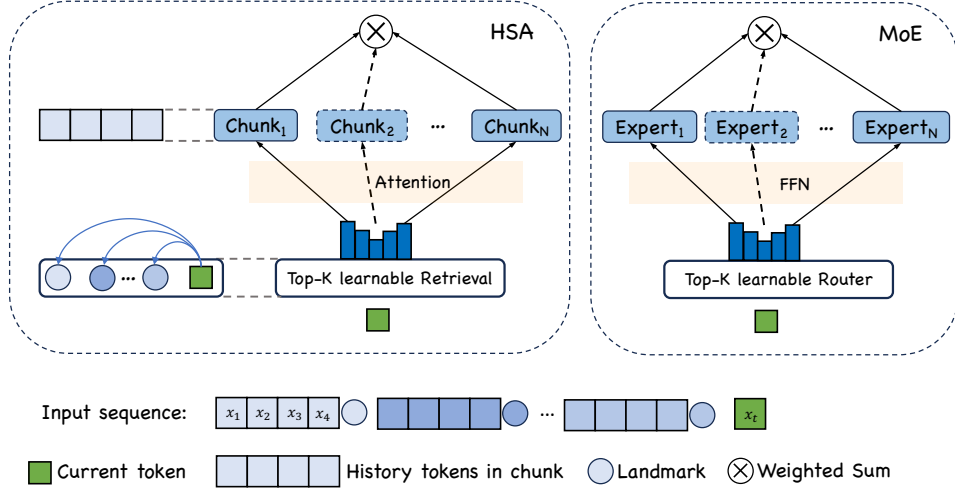


Figure 2: Hierarchical Sparse Attention (HSA) operates in a manner analogous to Mixture of Experts (MoE). First, the current token x_t computes dot products with the landmark representations of past chunks as retrieval scores, from which the top- k chunks are selected—similar to how MoE uses a router to select top- k experts. Subsequently, x_t performs attention with each of the k retrieved chunks **separately**, mirroring the process in MoE where x_t independently conducts Feedforward with k experts. Finally, the attention outputs from each chunk are weighted by the softmax-normalized retrieval scores and summed, which is functionally equivalent to MoE’s fusion of outputs from the selected FFNs.

Table 1: NSA ablation with 4K as the in-domain length. The higher scores are shown in **bold**.

Models	#params	Single-NIAH (ACC \uparrow)					MQ-NIAH (ACC \uparrow)				
		4K	8K	16K	32K	64K	4K	8K	16K	32K	64K
NSA(w/ RoPE)	370M	97.0	90.0	83.0	73.0	60.0	72.0	50.0	24.0	15.0	4.0
NSA(w/o RoPE)	370M	99.0	96.0	88.0	84.0	73.0	83.0	66.0	51.0	40.0	12.0

find that No Positional Encoding (NoPE) supports extrapolation better than RoPE (Su et al., 2024).

2.2 Attention with Intrinsic Chunk Retrieval

As we mentioned, the challenge of sparse attention lies in accurately retrieving the previous chunks. Hierarchical Sparse Attention (HSA) addresses the challenge by jointly learning chunk selection and attention in an end-to-end manner. Compared to NSA, HSA mainly makes two contributions:

- **Retrieval-oriented sparse attention.** Specifically, each token conducts attention with each past chunk **separately** and then fuses the attention results via retrieval scores.
- **RoPE for short, NoPE for long.** To mitigate the negative impact of RoPE on extrapolation, the sliding-window attention’s KV cache employs RoPE, while the HSA uses NoPE.

For an input sequence $\mathbf{S} = \{x_0, x_1, \dots, x_n\}$, where n is the length of the sequence, we denote the hidden states of tokens as $\mathbf{H} \in \mathbb{R}^{n \times d}$, where d is

the hidden dimension. The whole sequence is split into chunks according to a fixed length S , which is set to 64 by default to better align with hardware, thus we have $\frac{n}{S}$ chunks in total. We use indices with $[\cdot]$ to indicate that it is indexed by chunk rather than by token, e.g., $\mathbf{H}_{[i]} := \mathbf{H}_{iS:(i+1)S} \in \mathbb{R}^{S \times d}$. For each chunk, it has its own KV cache as $\mathbf{K}_{[i]}, \mathbf{V}_{[i]} \in \mathbb{R}^{S \times h \times d_h}$, with h as the number of heads satisfying $h \times d_h = d$, and its landmark representation as $\mathbf{L}_i \in \mathbb{R}^d$, which serves to summarize the content of the chunk. For each token, it uses $\mathbf{Q}_t^{slc} \in \mathbb{R}^d$ to retrieve chunks and $\mathbf{Q}_t^{attn} \in \mathbb{R}^{h \times d_h}$ to conduct attention with tokens inside chunks, both of which are derived from \mathbf{H}_i via linear transformations.

$$\begin{aligned} \mathbf{Q}^{slc} &= \mathbf{W}^{slc} \mathbf{H}, \quad \mathbf{Q}^{attn} = \mathbf{W}^{attn} \mathbf{H} \\ s_{t,i} &= \begin{cases} \mathbf{Q}_t^{slc \top} \mathbf{L}_i / \sqrt{d}, & i \leq \lfloor \frac{t}{S} \rfloor, \\ -\infty, & i > \lfloor \frac{t}{S} \rfloor, \end{cases} \quad (1) \\ \mathcal{I}_t &= \{i \mid \text{rank}(s_{t,i}) < K\}, \end{aligned}$$

where $\text{rank}(\cdot)$ denotes the ranking position in descending order, and \mathcal{I}_t is the indices of K chunks

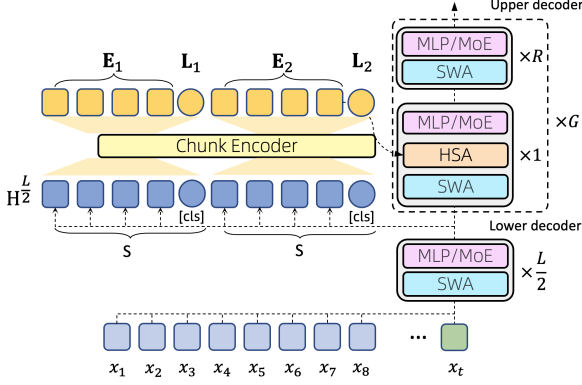


Figure 3: HSA-UltraLong model architecture.

with the highest relevance scores for x_t .

$$\begin{aligned} \bar{\mathbf{O}}_{t,i} &= \text{Attention}(\mathbf{Q}_t^{\text{attn}}, \mathbf{K}_{[i]}, \mathbf{V}_{[i]}) \\ &= \text{Softmax}\left(\underbrace{\frac{\text{norm}(\mathbf{Q}_t^{\text{attn}}) \text{norm}(\mathbf{K}_{[i]}^T)}{\sqrt{d_h}}}_{\text{intra-chunk attention}}\right) \mathbf{V}_{[i]}, \\ w_{t,i} &= \frac{\exp(s_{t,i})}{\sum_{k \in \mathcal{I}_t} \exp(s_{t,k})}, \quad \mathbf{O}_t = \underbrace{\sum_{k \in \mathcal{I}_t} w_{t,k} \bar{\mathbf{O}}_{t,k}}_{\text{inter-chunk fusion}}. \end{aligned} \quad (2)$$

norm is the Query-Key Normalization (Dehghani et al., 2023; Wortsman et al., 2023), which we find to be very important for the stability of HSA in practical trillion-token scale training.

3 Methodology

In terms of model design, we use SWA for local information access and HSA for global information access, fusing both local and global information through a stacking approach. A key challenge of long sequence inference is that the KV cache grows with the sequence length. Previous works (Wu and Tu, 2024; Rubin and Berant, 2024; Sun et al., 2024) have demonstrated that sharing the KV cache can significantly compress its size while maintaining comparable results. Inspired by these works, we share the intermediate layer KV cache among all HSA modules to serve as context memory.

3.1 Model Architecture

Overall, as shown in Figure 3, our architecture contains L layers, partitioned into an upper decoder and a lower decoder.

- The lower decoder is composed of $\frac{L}{2}$ standard Transformer layers, each utilizing Sliding Windowed Attention (SWA).
- The upper decoder is divided into G groups. Each group consists of one Transformer layer featuring both SWA and Head-based Sparse Attention

(HSA), followed by several layers employing SWA only.

We denote \mathbf{H}^l as the output hidden states of the l -th layer. To build compressed long-context memories, the intermediate layer output $\mathbf{H}^{\frac{L}{2}}$ is fed to a bidirectional encoder chunk by chunk. Specifically, each $\mathbf{H}_{[i]}^{\frac{L}{2}} \in \mathbb{R}^{S \times d}$ is followed by a [CLS] (acting as the landmark token) before being processed by the bi-directional encoder. This process yields the chunk embedding $\mathbf{E}_{[i]} \in \mathbb{R}^{S \times d}$ and the landmark embedding $\mathbf{L}_i \in \mathbb{R}^d$ used in Equation 1. Finally, the keys and values used in Equation 2 are obtained by applying a linear transformation to $\mathbf{E}_{[i]}$. Regarding MoE, we follow the design of Ling-2.0 (Ling-Team et al., 2025), where the first layer of the model adopts a dense MLP structure and all subsequent layers use MoE. Each MoE block has one shared expert, following the design in DeepSeek V3 (DeepSeek-AI et al., 2024). We use training-free balance strategy (Wang et al., 2024) as the expert balancing strategy.

3.2 Training

Previous work (Leng et al., 2025) demonstrates that with a 512-token sliding window, HSA-based models pre-trained on a 4K context length can generalize to 32M on the RULER task with high accuracy. However, such a small sliding window often sacrifices downstream performance. To address this, we increase the sliding window size to 4K tokens. However, when training models from scratch with a 4K sliding window, we observe that *they fail to generalize beyond the 4K context length*. We hypothesize that HSA’s length generalization ability arises from its inherent retrieval mechanism, which learns to generalize from short to long contexts. An oversized sliding window covers all short-range information, rendering short-range retrieval unnecessary for HSA. This inactivates HSA, preventing it from learning the crucial short-context retrieval mechanism required for generalization to longer contexts. To overcome this limitation, we introduce a warmup stage before the pre-training phase. The whole pre-training procedure is as follows:

- **Warm-up.** We use a short sliding window attention (SWA) of 512 tokens with HSA, whose *topk* is set large enough to cover all tokens, satisfying $\text{topk} \times \text{chunk_size} = \text{seq_len}$. Synthetic RULER tasks are randomly inserted into 1% of training samples. The warm-up phase is considered complete once the model achieves high

286 needle-in-a-haystack retrieval accuracy on con- 334
287 texts well beyond the 512-token window. At this 335
288 step, the context length is set to 16K. 336

- 289 • **Pre-training.** Following the warm-up, training 337
290 will resume from the checkpoint with the follow- 338
291 ing parameter adjustments: the SWA window 339
292 size is expanded to 4K, and the HSA’s *topk* is 340
293 decreased such that $topk \times chunk_size = 4K$. 341
294 The overall context length remains 16K. 342
- 295 • **Long-context mid-training.** Switch to corpora 343
296 with longer effective contexts and increase HSA 344
297 *topk* to cover all tokens. The context length is 345
298 expanded to 32K. 346
- 299 • **Annealing.** Perform annealing on high-quality 347
300 data while keeping a 32K context length. 348
- 301 • **Supervised fine-tuning.** Perform supervised 349
302 fine-tuning (SFT) with an 8K context length. 350

303 4 Experiments 351

304 4.1 Small-scale Preliminary Experiments 352

305 In this section, we detail the architecture of HSA- 353
306 UltraLong and validate its ability to balance in- 354
307 domain performance with length extrapolation ca- 355
308 pabilities through small-scale experiments. We 356
309 compared our model against Base Language Model 357
310 (BaseLM), which uses full attention across all lay- 358
311 ers. HSA-UltraLong employs SWA with a window 359
312 size of 4K across all layers, while strategically re- 360
313 placing two layers with HSA layers. The HSA 361
314 layers maintain a SWA window size of 512 and, 362
315 following our findings in Section 2.1, we removed 363
316 positional encoding information from these layers. 364
317 Each HSA layer includes an additional encoder sub- 365
318 layer, resulting in less than 5% parameter increase 366
319 compared to BaseLM. Within the HSA layers, we 367
320 set the chunk size to 64 and top-k to 64, estab- 368
321 lishing a fixed historical context window of 4,096 369
322 tokens. 370

323 Our initial experiments revealed limited length 371
324 extrapolation capabilities when training the model 372
325 directly with such configuration. We hypothe- 373
326 sized that this limitation stemmed from the pre- 374
327 dominance of pretraining data requiring only short- 375
328 range modeling capabilities within the 4K window 376
329 of SWA layers, leaving the HSA modules insuffi- 377
330 ciently trained. To address this issue, we explored 378
331 two warm-up strategies: 379

- 332 • **Self-copy warm-up:** We keep the model archi- 380
333 tecture unchanged and initialize training with a 381

334 self-copy objective. Given an input sequence 335
336 $\mathbf{S} = \{x_1, \dots, x_n\}$, we construct a target se- 337
338 quence $\mathbf{S}' = \{x_1, \dots, x_n, x_1, \dots, x_n\}$ by con- 339
340 catenating \mathbf{S} with itself. This objective encour- 341
342 ages the model to attend to and retrieve long- 343
344 range prefix information, enabling it to recon- 345
346 struct the second half of the sequence. 347

- 348 • **Full HSA + Short SWA warm-up:** Setting top- 349
350 k in HSA layers to 256 and sliding window size 351
352 to 512 during the initial training phase. 353

354 All experiments were conducted on a 0.5B pa- 355
356 rameter dense model trained on 100B tokens with 356
357 a pre-training context length of 16k. We incor- 357
358 porated 1% ruler-specific synthetic data into the 358
359 pre-training data to facilitate evaluation using ruler 359
360 benchmarks. Performance was evaluated based on 360
361 the perplexity of the last 4k tokens on the PG19 361
362 dataset and accuracy on the Multi-key NIAH (MK- 362
363 NIAH) task within the ruler benchmark. 363

364 The results in Table 2 demonstrated that the self- 364
365 copy warm-up strategy yielded the best length ex- 365
366 trapolation performance, albeit with some negative 366
367 impact on in-domain performance. The full HSA 367
368 + short SWA warm-up approach achieved a bet- 368
369 ter balance, maintaining in-domain performance 369
370 while delivering reasonable length extrapolation 370
371 capabilities. 371

361 4.2 Large-scale experiments 361

362 4.2.1 Training & Evaluation Details 362

363 We detail the pre-training data ratios, hyperparam- 363
364 eters, and evaluation datasets in the Appendix A. 364

365 4.2.2 General Tasks Evaluation 365

366 For HSA-UltraLong, we developed two variants: 366
367 a 0.5B dense model and an 8B MoE model 367
368 with 1B activated parameters. We compared 368
369 the MoE variant against a standard Transformer- 369
370 based model (TRM-MoE) with similar parameter 370
371 count—trained on identical data with matching hy- 371
372 perparameters. The architectures are largely con- 372
373 sistent, with only one structural difference: HSA- 373
374 UltraLong modifies the MoE configuration from 374
375 32-expert/2-activated to 64-expert/4-activated with 375
376 halved expert dimensions. Additionally, HSA- 376
377 UltraLong uses 16K pretraining context compared 377
378 to TRM-MoE’s 4K. For evaluation, we used the 378
379 TRM-MoE’s 8T-token checkpoint (pre-annealing). 379

380 We benchmarked the dense variant against 380
381 Qwen 2.5-0.5B (Yang et al., 2024) and Qwen3- 381
382 0.6B (Yang et al., 2025a). These comparison mod- 382

Table 2: Preliminary experiments on HSA-UltraLong-Base with a training context length of 16k tokens. The highest and second-best scores are shown in **bold** and underlined, respectively.

Models	#params	Warmup	PG19 (PPL ↓)			MQ-NIAH(ACC ↑)			
			4K	8K	16K	4K	8K	64K	1M
BaseLM	519.6M	-	<u>18.61</u>	17.53	16.77	89.0	23.0	5.0	0.0
SWA+HSA	537.7M	self-copy	18.87	<u>17.44</u>	<u>16.50</u>	100.0	96.0	93.0	93.0
SWA+HSA	537.7M	short-swa,full-hsa	18.30	17.13	15.96	<u>99.0</u>	<u>95.0</u>	<u>90.0</u>	<u>66.0</u>

Table 3: Comparison among HSA-UltraLong-Base (HSA-UL-Base) and other baselines. All models were evaluated under a unified framework for fair comparison.

	Qwen2.5 Annealing	Qwen3 Annealing	HSA-UL Annealing	TRM-MoE Base	HSA-UL Base	HSA-UL Annealing
Architecture	Dense	Dense	Dense	MoE	MoE	MoE
# Total Params	0.5B	0.6B	0.5B	8B	8B	8B
# Activated Params	0.5B	0.6B	0.5B	1B	1B	1B
# Training Tokens	18T	36T	4T	8T	8T	8T
General Tasks						
BBH	32.27	41.28	18.15	50.34	51.70	60.11
ARC-C	55.25	66.10	46.10	72.20	67.80	71.53
AGIEval	30.01	33.58	29.29	38.64	36.52	44.08
HellaSwag	48.05	48.88	44.48	67.69	67.39	67.43
PIQA	70.46	71.33	70.29	77.48	78.84	80.69
MMLU	49.73	54.40	41.76	58.74	57.83	60.71
CMMLU	52.10	51.97	42.08	57.68	57.49	64.41
C-Eval	54.17	54.57	44.30	56.87	58.36	65.98
Math Tasks						
GSM8K	41.32	60.88	37.45	66.41	67.02	72.93
MATH	18.14	31.44	20.66	37.96	41.98	48.00
CMATH	52.09	66.67	60.75	74.59	74.13	82.88
Coding Tasks						
HumanEval+	24.39	26.83	29.27	48.17	50.61	61.59
MBPP+	32.80	38.36	20.63	50.26	55.82	62.17
CRUX-O	14.38	31.62	22.56	35.12	36.31	40.75
AVG	41.08	48.42	37.70	56.58	57.27	63.09

els have similar parameter counts but were trained on substantially larger datasets—4.5 times and 9 times our training data volume, respectively.

Our primary evaluation focused on assessing model performance within the pretraining context length across standard benchmarks. Results in Table 3 show that the HSA-UltraLong-MoE achieved parity with TRM-MoE in average performance scores, while the dense variant demonstrated only a 3.3-point deficit compared to Qwen 2.5-0.5B, despite having significantly less training data.

Additionally, we evaluated the MoE and Dense models after supervised fine-tuning. The results in Table 4 indicate that while a few general tasks showed no significant performance improvement after supervised fine-tuning, most tasks—particularly math and coding tasks—demonstrated substantial enhancements compared to the base models. Notably, our HSA-UltraLong-MoE achieved scores averaging

1.3 points higher than Qwen3-1.7B (Non-thinking), despite requiring fewer training flops. Similarly, our dense variant performed competitively, scoring only approximately 4 points below Qwen3-0.6B, despite being trained on a dataset merely one-ninth the size.

These findings demonstrate that HSA-UltraLong models maintain their capabilities within standard contexts while extending their effective context length to 16M tokens, further highlighting the superiority of our architectural approach.

4.2.3 Long-context Evaluation

During training, we randomly convert samples into RULER tasks with a 1% probability by inserting a “needle” in a long context and appending the Needle-in-a-Haystack (NIAH) prompt and answer at the end of the sample so the model can follow the NIAH instructions. This modification serves as a probe task to evaluate the model’s extrapolation ability while having minimal impact on training.

Table 4: Comparison among HSA-UltraLong-Inst (HSA-UL-Inst) and Qwen3 (Non-thinking) after supervised fine-tuning. All models were evaluated under a unified framework for fair comparison.

	Qwen3-Inst	HSA-UL-Inst	Qwen3-Inst	HSA-UL-Inst
Architecture	Dense	Dense	Dense	MoE
# Total Params	0.6B	0.5B	1.7B	8B
# Activated Params	0.6B	0.5B	1.7B	1B
# Training Tokens	36T	4T	36T	8T
General Tasks				
BBH	42.56	26.25	59.48	57.25
MMLU	45.87	42.24	63.05	61.34
CMMLU	41.64	43.33	60.84	64.06
C-Eval	43.81	45.41	62.70	62.86
Math Tasks				
GSM8K	55.65	55.42	79.00	82.94
MATH	45.26	40.76	64.32	61.56
MATH500	53.00	41.00	73.20	71.00
OlympiadBench	16.89	8.74	36.30	27.85
Coding Tasks				
HumanEval	40.24	39.63	65.24	71.95
MBPP	29.20	34.40	51.00	57.00
HumanEval+	35.37	37.20	61.59	70.73
MBPP+	34.39	39.95	59.52	65.87
CRUX-O	28.00	23.25	50.00	50.75
Alignment Tasks				
IFEval Strict Prompt	55.08	33.09	64.33	63.22
AVG	40.50	36.48	60.76	62.03

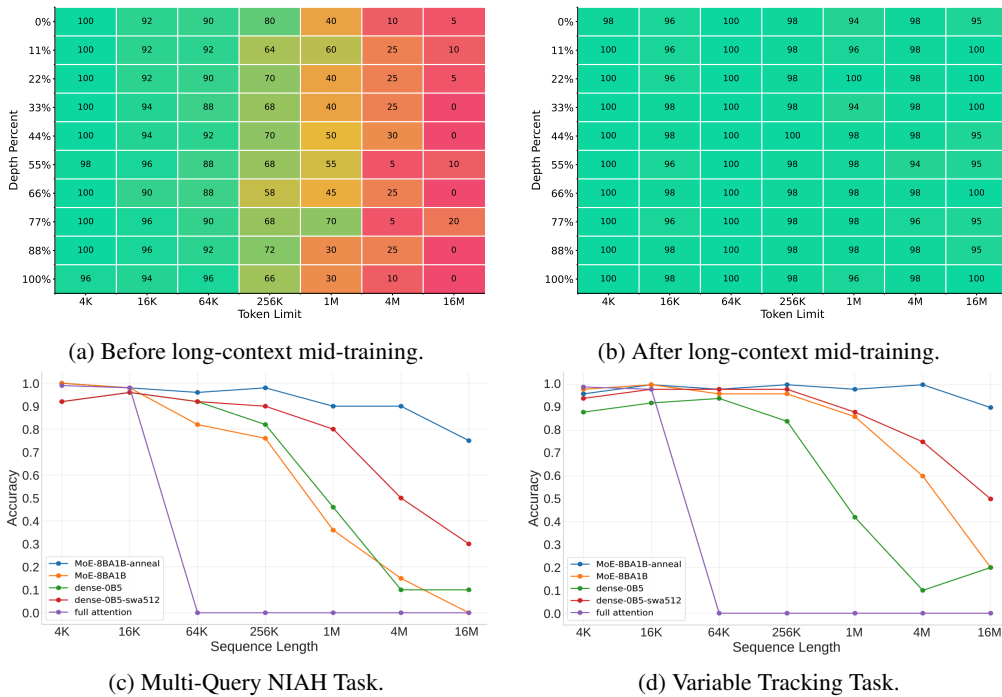


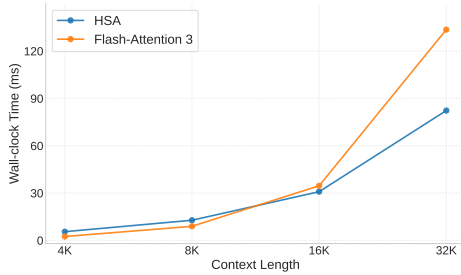
Figure 4: Evaluation of length generalization using the Needle-in-a-Haystack test. (a) and (b) present the results of the HSA-UltraLong-MoE before and after the long-context continued training phase on the Single-NIAH task at various depths. In (c) and (d), we evaluate the performance of different models on the Multi-Query NIAH Task (2 queries, 6 key-value pairs) and the Variable Tracking Task.

The results for RULER tasks are reported in Figure 4, we identify three key findings:

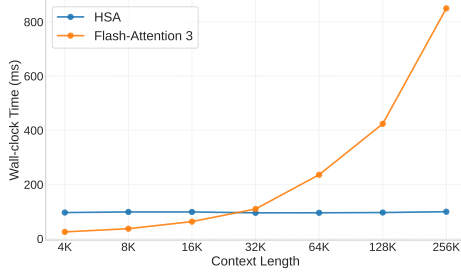
- **Effective context length of training data is critical for HSA extrapolation.** As shown in Figure 4(a), models pretrained on standard corpora

exhibit a progressive decline in retrieval accuracy with longer contexts. This occurs despite a 16K pretraining context window, because the effective context length of the data is often much shorter. In contrast, training on data with longer effective

428
429
430
431
432



(a) Training efficiency.



(b) Inference efficiency.

Figure 5: Comparison of training/inference efficiency between HSA kernel and Flash-attention 3.

contexts ($>32K$), as in Figure 4(b), yields substantially improved extrapolation. This principle underpins the trends in Figures 4(c) and (d).

- **A seesaw effect exists between HSA and Sliding Window Attention.** Figures 4(c) and (d) indicate that a smaller SWA window (512) during continued pretraining leads to better HSA extrapolation than a larger window (4K). Given that training from scratch with a 4K window fails to develop extrapolative HSA, we conclude that larger SWA windows impair HSA’s long-range generalization. We posit that HSA learns a form of retrieval-based extrapolation. Large SWA windows handle most short-range dependencies inherently, reducing the incentive for HSA to learn them and thus weakening its ability to generalize to longer sequences.
- **HSA capability scales with parameter size in reasoning-retrieval tasks.** While MoE-8B-A1B and Dense-0.5B exhibit comparable performance on the pure retrieval task (MQ-NIAH; Figure 4(c)), MoE-8B-A1B consistently outperforms Dense-0.5B in the variable-tracking task (Figure 4(d)), demonstrating that larger models better support joint reasoning and retrieval.

4.3 Training/Inference Efficiency

To further evaluate the efficiency of the sparse attention module, we benchmark the HSA operator against the FlashAttention-3 (Shah et al., 2024) op-

erator on H800 for both training and inference, with HSA implemented using TileLang (Wang et al., 2025). As shown in Figure 5, at shorter sequence lengths, FlashAttention-3 still leads in both training and inference, and HSA only gains an advantage with longer contexts. We attribute this to two factors: (1) the sparsity in HSA causes the kernel to incur more memory accesses compared to FlashAttention-3; and (2) FlashAttention-3 is implemented in CUDA, enabling it to better leverage the features of the Hopper architecture.

5 Related Works

Many works have explored sparse attention for efficient long-context modeling, which broadly fall into two lines: token-level sparsity (Lou et al., 2024; Gonçalves et al., 2025; DeepSeek-AI et al., 2025) and chunk-level sparsity. A central challenge in chunk-level sparse attention is how to accurately select relevant past chunks. NSA (Yuan et al., 2025) compresses the key-value pairs within each chunk into a single key-value pair and uses the resulting compressed attention to guide chunk selection. MoBA (Lu et al., 2025) heuristically selects chunks by summing unnormalized attention logits within each chunk. Seer Attention (Gao et al., 2025) learns sparse patterns via distillation from full attention. Our work is mainly based on HSA (Hu et al., 2025a), a chunk-level sparse attention mechanism. The original HSA framework is coupled with Mamba layers; however, in our scaling experiments, we observe that extrapolation degrades as model size increases. To stabilize extrapolation under scaling, we replace Mamba with SWA. Compared with NSA, MoBA, and Seer Attention, the key distinction of HSA is end-to-end learned chunk retrieval, without distillation or heuristic approximations, which enables efficient training while achieving higher in-domain accuracy and preserving extrapolation capability.

6 Conclusion

In this work, HSA-UltraLong presents a highly promising paradigm for long-context processing. The core insight of HSA is to *perform attention chunk by chunk and fuse the results via retrieval scores*, rather than selecting chunks and then concatenating them for attention. The experimental results provide a meaningful step toward effectively handling infinite-long context, advancing progress on long-term memory in machines.

511 Limitations

512 Although HSA has shown promising extrapolation
513 capabilities, several challenges remain:

- 514 • The HSA/SWA seesaw problem. After training
515 on short SFT data, extrapolation can de-
516 grade. The main reason is that an exces-
517 sively long sliding-window attention reduces
518 the need for HSA to learn short-range depen-
519 dencies, which in turn hampers its ability to
520 extrapolate to long-range dependencies.
- 521 • The head ratio constraint. HSA currently re-
522 quires a 16:1 ratio of query heads to key–value
523 heads, creating a severe information bottle-
524 neck. Future work should pursue kernel-level
525 optimizations to alleviate this constraint.
- 526 • When sequences are short, training and
527 inference show no clear advantage over
528 FlashAttention-3; further kernel-level opti-
529 mizations are needed to improve efficiency.

530 Furthermore, there is a slight discrepancy in the
531 comparison between HSA-UL and TRM-MoE in
532 Table 3. While the total MoE parameters are iden-
533 tical, the baseline employs a top-2 routing strategy
534 among 32 experts (each with 1024 dimensions),
535 whereas HSA-UL employs a top-4 strategy among
536 64 experts (each with 512 dimensions). However,
537 this difference does not significantly impact the
538 overall conclusion: HSA-based models maintain
539 efficient extrapolation capabilities even as param-
540 eter counts and data volumes scale up.

541 References

- 542 Josh Achiam, Steven Adler, Sandhini Agarwal, Lama
543 Ahmad, Ilge Akkaya, Florencia Leoni Aleman,
544 Diogo Almeida, Janko Altenschmidt, Sam Altman,
545 Shyamal Anadkat, and 1 others. 2023. Gpt-4 techni-
546 cal report. *arXiv preprint arXiv:2303.08774*.
- 547 Iz Beltagy, Matthew E. Peters, and Arman Cohan.
548 2020. Longformer: The long-document transformer.
549 *arXiv:2004.05150*.
- 550 Sumithra Bhakthavatsalam, Daniel Khashabi, Tushar
551 Khot, Bhavana Dalvi Mishra, Kyle Richardson,
552 Ashish Sabharwal, Carissa Schoenick, Oyvind
553 Tafjord, and Peter Clark. 2021. *Think you have
554 solved direct-answer question answering? try arc-da,
555 the direct-answer AI2 reasoning challenge*. *CoRR*,
556 abs/2102.03315.
- 557 Yonatan Bisk, Rowan Zellers, Ronan Le Bras, Jianfeng
558 Gao, and Yejin Choi. 2020. *PIQA: reasoning about*

physical commonsense in natural language. In *The
559 Thirty-Fourth AAAI Conference on Artificial Intelli-
560 gence, AAAI 2020, The Thirty-Second Innovative Ap-
561 plications of Artificial Intelligence Conference, IAAI
562 2020, The Tenth AAAI Symposium on Educational
563 Advances in Artificial Intelligence, EAAI 2020, New
564 York, NY, USA, February 7-12, 2020*, pages 7432–
565 7439. AAAI Press. 566

Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie
Subbiah, and 1 others. 2020. *Language models are
567 few-shot learners*. In *Advances in Neural Information
568 Processing Systems 33: Annual Conference on Neu-
569 ral Information Processing Systems 2020, NeurIPS
570 2020, December 6-12, 2020, virtual*. 572

Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan,
Henrique Pondé de Oliveira Pinto, and 1 others. 2021.
Evaluating large language models trained on code.
CoRR, abs/2107.03374. 573-576

Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian,
Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias
Plappert, Jerry Tworek, Jacob Hilton, Reiichiro
Nakano, and 1 others. 2021. *Training Verifiers
577 to Solve Math Word Problems*. *arXiv preprint
578 arXiv:2110.14168*. 579-582

Nelson Cowan. 2008. *What are the differences be-
583 tween long-term, short-term, and working memory?*
584 *Progress in brain research*, 169:323–38. 585

Tri Dao and Albert Gu. 2024. *Transformers are ssms:
586 Generalized models and efficient algorithms through
587 structured state space duality*. In *Forty-first Interna-
588 tional Conference on Machine Learning, ICML 2024,
589 Vienna, Austria, July 21-27, 2024*. OpenReview.net.

DeepSeek-AI, Aixin Liu, Bei Feng, Bing Xue, Bingx-
uan Wang, and 1 others. 2024. *Deepseek-v3 techni-
591 cal report*. *Preprint*, arXiv:2412.19437. 592-593

DeepSeek-AI, Aixin Liu, Aoxue Mei, Bangcai Lin,
Bing Xue, Bingxuan Wang, Bingzheng Xu, Bochao
Wu, Bowei Zhang, Chaofan Lin, Chen Dong,
Chengda Lu, Chenggang Zhao, Chengqi Deng, Chen-
hao Xu, Chong Ruan, Damai Dai, Daya Guo, Dejian
Yang, and 245 others. 2025. *Deepseek-v3.2: Pushing
594 the frontier of open large language models*. *Preprint*,
595 arXiv:2512.02556. 596-601

Mostafa Dehghani, Josip Djolonga, Basil Mustafa, and
1 others. 2023. *Scaling vision transformers to 22
602 billion parameters*. *Preprint*, arXiv:2302.05442. 603-604

Yizhao Gao, Zhichen Zeng, Dayou Du, Shijie Cao,
Peiyuan Zhou, Jiaxing Qi, Junjie Lai, Hayden Kwok-
Hay So, Ting Cao, Fan Yang, and Mao Yang. 2025.
*Seerattention: Learning intrinsic sparse attention in
605 your llms*. *Preprint*, arXiv:2410.13276. 606-609

Nuno Gonçalves, Marcos V Treviso, and Andre Martins.
2025. *Adasplash: Adaptive sparse flash attention*. In
*Forty-second International Conference on Machine
610 Learning*. 611-613

614	Albert Gu and Tri Dao. 2023. Mamba: Linear-time sequence modeling with selective state spaces . <i>CoRR</i> , abs/2312.00752.	670
615		671
616		
617	Alex Gu, Baptiste Rozière, Hugh Leather, Armando Solar-Lezama, Gabriel Synnaeve, and Sida I Wang. 2024. CruxEval: A Benchmark for Code Reasoning, Understanding and Execution . <i>arXiv preprint arXiv:2401.03065</i> .	672
618		673
619		674
620		675
621		
622	Chaoqun He, Renjie Luo, Yuzhuo Bai, Shengding Hu, Zhen Leng Thai, Junhao Shen, Jinyi Hu, Xu Han, Yujie Huang, Yuxiang Zhang, and 1 others. 2024. OlympiadBench: A Challenging Benchmark for Promoting AGI with Olympiad-Level Bilingual Multimodal Scientific Problems . <i>arXiv preprint arXiv:2402.14008</i> .	676
623		677
624		678
625		679
626		680
627		
628		
629	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> .	681
630		682
631		683
632		684
633	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> .	685
634		686
635		
636		
637		
638	Cheng-Ping Hsieh, Simeng Sun, Samuel Kriman, Shantanu Acharya, Dima Rekeshe, Fei Jia, and Boris Ginsburg. 2024. RULER: What’s the real context size of your long-context language models? In <i>First Conference on Language Modeling</i> .	687
639		688
640		689
641		690
642		691
643	Xiang Hu, Jiaqi Leng, Jun Zhao, Kewei Tu, and Wei Wu. 2025a. Hardware-aligned hierarchical sparse attention for efficient long-term memory access . In <i>The Thirty-ninth Annual Conference on Neural Information Processing Systems</i> .	692
644		693
645		694
646		695
647		696
648	Xiang Hu, Zhihao Teng, Jun Zhao, Wei Wu, and Kewei Tu. 2025b. Efficient length-generalizable attention via causal retrieval for long-context language modeling . In <i>Forty-second International Conference on Machine Learning</i> .	697
649		
650		
651		
652		
653	Yuzhen Huang, Yuzhuo Bai, Zhihao Zhu, Junlei Zhang, Jinghan Zhang, Tangjun Su, Junteng Liu, Chuancheng Lv, and 1 others. 2023. C-Eval: A Multi-Level Multi-Discipline Chinese Evaluation Suite for Foundation Models . <i>Advances in Neural Information Processing Systems</i> , 36:62991–63010.	698
654		699
655		700
656		701
657		
658		
659	Angelos Katharopoulos, Apoorv Vyas, Nikolaos Pappas, and François Fleuret. 2020. Transformers are rnns: Fast autoregressive transformers with linear attention . In <i>Proceedings of the 37th International Conference on Machine Learning, ICML 2020, 13-18 July 2020, Virtual Event</i> , volume 119 of <i>Proceedings of Machine Learning Research</i> , pages 5156–5165. PMLR.	702
660		703
661		704
662		705
663		706
664		707
665		708
666		709
667	Yuri Kuratov, Aydar Bulatov, Petr Anokhin, Ivan Rodkin, Dmitry Sorokin, Artyom Sorokin, and Mikhail Burtsev. 2024. Babilong: Testing the limits of llms with long context reasoning-in-a-haystack . <i>Preprint</i> , arXiv:2406.10149.	710
668		711
669		712
		713
	Jiaqi Leng, Xiang Hu, Junxiong Wang, Jianguo Li, Wei Wu, and Yucheng Lu. 2025. Understanding and improving length generalization in hierarchical sparse attention models . <i>Preprint</i> , arXiv:2510.17196.	714
		715
		716
		717
	Haonan Li, Yixuan Zhang, Fajri Koto, Yifei Yang, Hai Zhao, Yeyun Gong, Nan Duan, and Timothy Baldwin. 2023. CMMLU: Measuring Massive Multitask Language Understanding in Chinese . <i>arXiv preprint arXiv:2306.09212</i> .	718
		719
		720
		721
		722
		723
	Hunter Lightman, Vineet Kosaraju, Yuri Burda, Harrison Edwards, Bowen Baker, Teddy Lee, Jan Leike, John Schulman, Ilya Sutskever, and Karl Cobbe. 2023. Let’s Verify Step by Step . In <i>The Twelfth International Conference on Learning Representations</i> .	
	Ling-Team, Ang Li, Ben Liu, Binbin Hu, Bing Li, Bingwei Zeng, Borui Ye, and 1 others. 2025. Every activation boosted: Scaling general reasoner to 1 trillion open language foundation . <i>Preprint</i> , arXiv:2510.22115.	
	Jiawei Liu, Chunqiu Steven Xia, Yuyao Wang, and Lingming Zhang. 2023. Is Your Code Generated by ChatGPT Really Correct? Rigorous Evaluation of Large Language Models for Code Generation . <i>Advances in Neural Information Processing Systems</i> , 36:21558–21572.	
	Chao Lou, Zixia Jia, Zilong Zheng, and Kewei Tu. 2024. Sparser is faster and less is more: Efficient sparse attention for long-range transformers . <i>Preprint</i> , arXiv:2406.16747.	
	Enzhe Lu, Zhejun Jiang, Jingyuan Liu, Yulun Du, Tao Jiang, Chao Hong, Shaowei Liu, Weiran He, Enming Yuan, Yuzhi Wang, Zhiqi Huang, Huan Yuan, Suting Xu, Xinran Xu, Guokun Lai, Yanru Chen, Huabin Zheng, Junjie Yan, Jianlin Su, and 6 others. 2025. MoBA: Mixture of block attention for long-context LLMs . In <i>The Thirty-ninth Annual Conference on Neural Information Processing Systems</i> .	
	Amirkeivan Mohtashami and Martin Jaggi. 2023. Random-access infinite context length for transformers . In <i>Thirty-seventh Conference on Neural Information Processing Systems</i> .	
	Ohad Rubin and Jonathan Berant. 2024. Retrieval-pretrained transformer: Long-range language modeling with self-retrieval . <i>Transactions of the Association for Computational Linguistics</i> , 12:1197–1213.	
	Jay Shah, Ganesh Bikshandi, Ying Zhang, Vijay Thakkar, Pradeep Ramani, and Tri Dao. 2024. Flashattention-3: Fast and accurate attention with asynchrony and low-precision . In <i>The Thirty-eighth Annual Conference on Neural Information Processing Systems</i> .	

724	Noam Shazeer, *Azalia Mirhoseini, *Krzysztof	Haoyi Wu and Kewei Tu. 2024. Layer-condensed KV cache for efficient inference of large language models.	779
725	Maziarz, Andy Davis, Quoc Le, Geoffrey Hinton,	In <i>Proceedings of the 62nd Annual Meeting of the</i>	780
726	and Jeff Dean. 2017. Outrageously large neural networks: The sparsely-gated mixture-of-experts layer.	<i>Association for Computational Linguistics (Volume 1: Long Papers)</i> , ACL 2024, Bangkok, Thailand, August	781
727	In <i>International Conference on Learning Representations.</i>	11-16, 2024, pages 11175–11188. Association for	782
728		Computational Linguistics.	783
729			784
730	Jianlin Su, Murtadha H. M. Ahmed, Yu Lu, Shengfeng	Haoyuan Wu, Haoxing Chen, Xiaodong Chen, Zhan-	786
731	Pan, Wen Bo, and Yunfeng Liu. 2024. Roformer: Enhanced transformer with rotary position embedding.	chao Zhou, Tiejuan Chen, Yihong Zhuang, Guoshan	787
732	<i>Neurocomputing</i> , 568:127063.	Lu, Zenan Huang, Junbo Zhao, Lin Liu, and 1 oth-	788
733		ers. 2025. Grove moe: Towards efficient and supe-	789
734	Yutao Sun, Li Dong, Yi Zhu, Shaohan Huang, Wen-	rior moe llms with adjugate experts. <i>arXiv preprint</i>	790
735	hui Wang, Shuming Ma, Quanlu Zhang, Jianyong	<i>arXiv:2508.07785.</i>	791
736	Wang, and Furu Wei. 2024. You only cache once: Decoder-decoder architectures for language models.	An Yang, Anfeng Li, Baosong Yang, Beichen Zhang,	792
737	In <i>The Thirty-eighth Annual Conference on Neural</i>	Binyuan Hui, Bo Zheng, Bowen Yu, Chang	793
738	<i>Information Processing Systems.</i>	Gao, Chengen Huang, Chenxu Lv, and 1 others.	794
739		2025a. Qwen3 Technical Report. <i>arXiv preprint</i>	795
740	Mirac Suzgun, Nathan Scales, Nathanael Schärli, Se-	<i>arXiv:2505.09388.</i>	796
741	bastian Gehrmann, Yi Tay, Hyung Won Chung,	An Yang, Baosong Yang, Beichen Zhang, Binyuan	797
742	Aakanksha Chowdhery, Quoc V Le, Ed H Chi, Denny	Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Day-	798
743	Zhou, and 1 others. 2022. Challenging Big-Bench	iheng Liu, Fei Huang, Haoran Wei, and 1 others.	799
744	Tasks and Whether Chain-of-Thought Can Solve	2024. Qwen2.5 technical report. <i>arXiv preprint</i>	800
745	Them. <i>arXiv preprint arXiv:2210.09261.</i>	<i>arXiv:2412.15115.</i>	801
746	Ning Tao, Anthony Ventresque, Vivek Nallur, and Tak-	Songlin Yang, Jan Kautz, and Ali Hatamizadeh. 2025b.	802
747	farinas Saber. 2024. Enhancing program synthesis with large language models using many-objective grammar-guided genetic programming.	Gated delta networks: Improving mamba2 with delta rule.	803
748	<i>Algorithms</i> , 17(7):287.	In <i>The Thirteenth International Conference on</i>	804
749		<i>Learning Representations, ICLR 2025, Singapore,</i>	805
750		<i>April 24-28, 2025.</i> OpenReview.net.	806
751	Hugo Touvron, Louis Martin, Kevin Stone, Peter Al-	Jingyang Yuan, Huazuo Gao, Damai Dai, Junyu Luo,	807
752	bert, Amjad Almahairi, Yasmine Babaei, Nikolay	Liang Zhao, Zhengyan Zhang, Zhenda Xie, Yux-	808
753	Bashlykov, Soumya Batra, Prajwal Bhargava, Shruti	ing Wei, Lean Wang, Zhiping Xiao, and 1 others.	809
754	Bhosale, and 1 others. 2023. Llama 2: Open founda-	2025. Native sparse attention: Hardware-aligned and	810
755	tion and fine-tuned chat models. <i>arXiv preprint</i>	natively trainable sparse attention. In <i>Proceedings</i>	811
756	<i>arXiv:2307.09288.</i>	<i>of the 63rd Annual Meeting of the Association for</i>	812
757	Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob	<i>Computational Linguistics (Volume 1: Long Papers)</i> ,	813
758	Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz	pages 23078–23097.	814
759	Kaiser, and Illia Polosukhin. 2017. Attention is all	Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali	815
760	you need. <i>Advances in neural information processing</i>	Farhadi, and Yejin Choi. 2019. Hellaswag: Can a machine really finish your sentence?	816
761	<i>systems</i> , 30.	In <i>Proceedings</i>	817
762	Lean Wang, Huazuo Gao, Chenggang Zhao, Xu Sun,	<i>of the 57th Conference of the Association for Compu-</i>	818
763	and Damai Dai. 2024. Auxiliary-loss-free load balancing strategy for mixture-of-experts.	<i>tational Linguistics, ACL 2019, Florence, Italy, July</i>	819
764	<i>CoRR</i> , abs/2408.15664.	28- August 2, 2019, Volume 1: Long Papers, pages	820
765		4791–4800. Association for Computational Linguis-	821
766	Lei Wang, Yu Cheng, Yining Shi, Zhengju Tang, Zhi-	tics.	822
767	wen Mo, Wenhao Xie, Lingxiao Ma, Yuqing Xia,	Wanjun Zhong, Ruixiang Cui, Yiduo Guo, Yaobo Liang,	823
768	Jilong Xue, Fan Yang, and Zhi Yang. 2025. Tile-	Shuai Lu, Yanlin Wang, Amin Saied, Weizhu Chen,	824
769	lang: A composable tiled programming model for ai	and Nan Duan. 2024. Agieval: A human-centric	825
770	systems. <i>Preprint</i> , arXiv:2504.17577.	benchmark for evaluating foundation models.	826
771	Tianwen Wei, Jian Luan, Wei Liu, Shuang Dong, and	In <i>Findings of the Association for Computational Lin-</i>	827
772	Bin Wang. 2023. CMATH: can your language model	<i>guistics: NAACL 2024, Mexico City, Mexico, June</i>	828
773	pass chinese elementary school math test?	16-21, 2024, pages 2299–2314. Association for Com-	829
774	<i>CoRR</i> , abs/2306.16636.	putational Linguistics.	830
775	Mitchell Wortsman, Peter J. Liu, Lechao Xiao, Katie	Jeffrey Zhou, Tianjian Lu, Swaroop Mishra, Sid-	831
776	Everett, and 1 others. 2023. Small-scale prox-	dhartha Brahma, Sujoy Basu, Yi Luan, Denny Zhou,	832
777	ies for large-scale transformer training instabilities.	and Le Hou. 2023. Instruction-Following Evalua-	833
778	<i>Preprint</i> , arXiv:2309.14322.	tion for Large Language Models. <i>arXiv preprint</i>	834
		<i>arXiv:2311.07911.</i>	835

A Pre-training setups

Training Data In the first phase of general pre-training, we utilized a large-scale deduplicated, multi-domain dataset totaling 10T tokens, with differential sampling ratios across various sub-datasets from different domains. This distribution comprised predominantly Web content (50%), followed by Code (14.4%), Math (12.0%), Code-nlp (5.6%), Reason (5%), Multilingual (4.0%), Books (2.0%), Wikipedia (1.5%), and Others (5.5%). During this phase, the MoE model processed 8T tokens while the dense model was trained on 4T tokens. The second phase uses a dataset of 32K-length long-text sequences totaling 175B tokens and the third phase consisted of 400B tokens with a high proportion of reasoning data. During the Supervised Fine-tuning phase, we utilized the same dataset as described in (Wu et al., 2025).

Hyperparameters All models are trained using AdamW optimizer with a weight decay of 0.01, $\beta_1 = 0.9$, $\beta_2 = 0.95$, and gradient clipping norm of 1.0. We use FSDP2 for distributed training. For the MoE model, we employ a learning rate of $3.87e-4$, sequence length of 16,384, and batch size of 16.8M tokens. The dense model is trained with a learning rate of $4.96e-4$ and batch size of 5.2M tokens. The learning rate schedule begins with a linear warmup phase followed by a constant learning rate maintained until training completion. For the Supervised Fine-tuning stage, we adopt a cosine decay learning rate schedule. The dense model uses a learning rate of $5.5e-5$ and was trained for up to 5 epochs, while the MoE model uses a learning rate of $3.87e-4$ and was trained for up to 3 epochs. In both cases, we select the checkpoint from the epoch that yields the best performance.

Evaluation Benchmarks To conduct a comprehensive evaluation of the model, we selected a diverse range of assessment tasks, encompassing four major categories: general tasks, mathematical tasks, coding tasks, and alignment tasks:

- **General Tasks:** MMLU (Hendrycks et al., 2020), CMMLU (Li et al., 2023), C-Eval (Huang et al., 2023), ARC (Bhaktavatsalam et al., 2021), AGIEval (Zhong et al., 2024), PIQA (Bisk et al., 2020), HellaSwag (Zellers et al., 2019) and BBH (Suzgun et al., 2022).
- **Math Tasks:** GSM8K (Cobbe et al.,

2021), MATH (Hendrycks et al., 2021), CMATH (Wei et al., 2023), MATH-500 (Lightman et al., 2023) and Olympiad-Bench (He et al., 2024).

- **Coding Tasks:** HumanEval (Chen et al., 2021), HumanEval+ (Liu et al., 2023), MBPP (Tao et al., 2024), MBPP+ (Liu et al., 2023), and CRUX-O (Gu et al., 2024).
- **Alignment Tasks:** We report the average prompt-level strict accuracy of IFEval (Zhou et al., 2023)