

# MoBA: Mixture of Block Attention for Long-Context LLMs

Jiezhong Qiu

Hangzhou Institute of Medicine Chinese Academy of  
Sciences

# MoBA

- Paper: <https://arxiv.org/pdf/2502.13189>
- Code: <https://github.com/MoonshotAI/MoBA>

---

## MoBA: MIXTURE OF BLOCK ATTENTION FOR LONG-CONTEXT LLMs

---

TECHNICAL REPORT

**Enzhe Lu<sup>1</sup>   Zhejun Jiang<sup>1</sup>   Jingyuan Liu<sup>1</sup>   Yulun Du<sup>1</sup>   Tao Jiang<sup>1</sup>   Chao Hong<sup>1</sup>  
Shaowei Liu<sup>1</sup>   Weiran He<sup>1</sup>   Enming Yuan<sup>1</sup>   Yuzhi Wang<sup>1</sup>   Zhiqi Huang<sup>1</sup>   Huan Yuan<sup>1</sup>  
Suting Xu<sup>1</sup>   Xinran Xu<sup>1</sup>   Guokun Lai<sup>1</sup>   Yanru Chen<sup>1</sup>   Huabin Zheng<sup>1</sup>   Junjie Yan<sup>1</sup>  
Jianlin Su<sup>1</sup>   Yuxin Wu<sup>1</sup>   Neo Y. Zhang<sup>1</sup>   Zhilin Yang<sup>1</sup>  
Xinyu Zhou<sup>1,‡</sup>   Mingxing Zhang<sup>2,\*</sup>   Jiezhong Qiu<sup>3,‡</sup> \*†**

# MoBA

- Full Attention

$$\text{Attn}(\mathbf{q}, \mathbf{K}, \mathbf{V}) = \text{Softmax}(\mathbf{q}\mathbf{K}^\top) \mathbf{V}$$

- MoBA:

$$\text{MoBA}(\mathbf{q}, \mathbf{K}, \mathbf{V}) = \text{Softmax}(\mathbf{q}\mathbf{K}[I]^\top) \mathbf{V}[I]$$

$$I_i = [(i - 1) \times B + 1, i \times B]$$

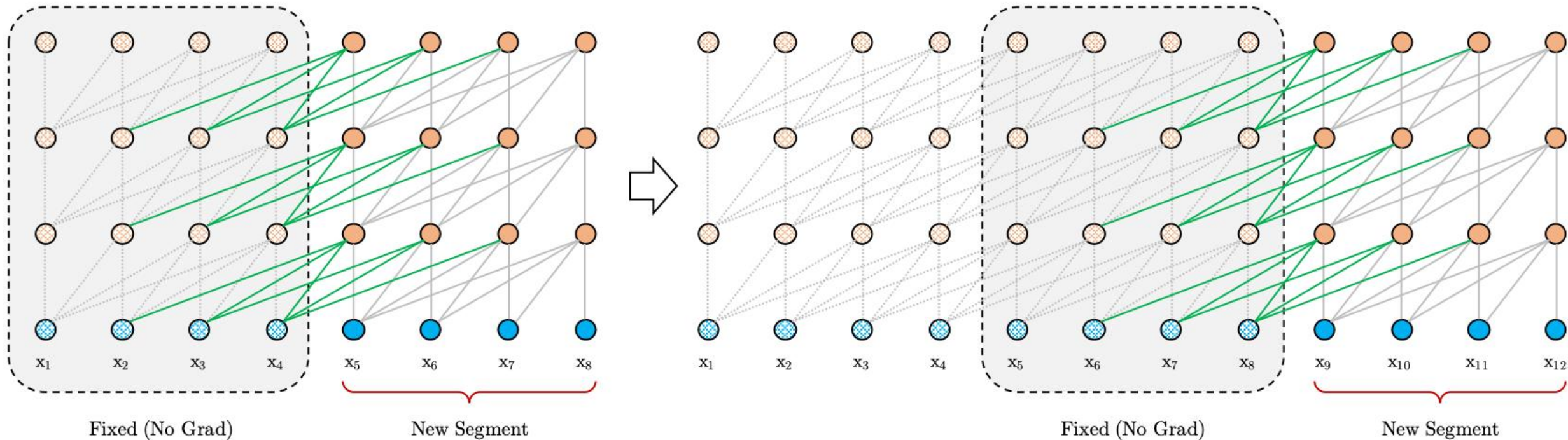
$$I = \bigcup_{g_i > 0} I_i.$$

$$g_i = \begin{cases} 1 & s_i \in \text{Topk}(\{s_j | j \in [n]\}, k) \\ 0 & \text{otherwise} \end{cases}$$

$$s_i = \langle \mathbf{q}, \text{mean\_pool}(\mathbf{K}[I_i]) \rangle$$

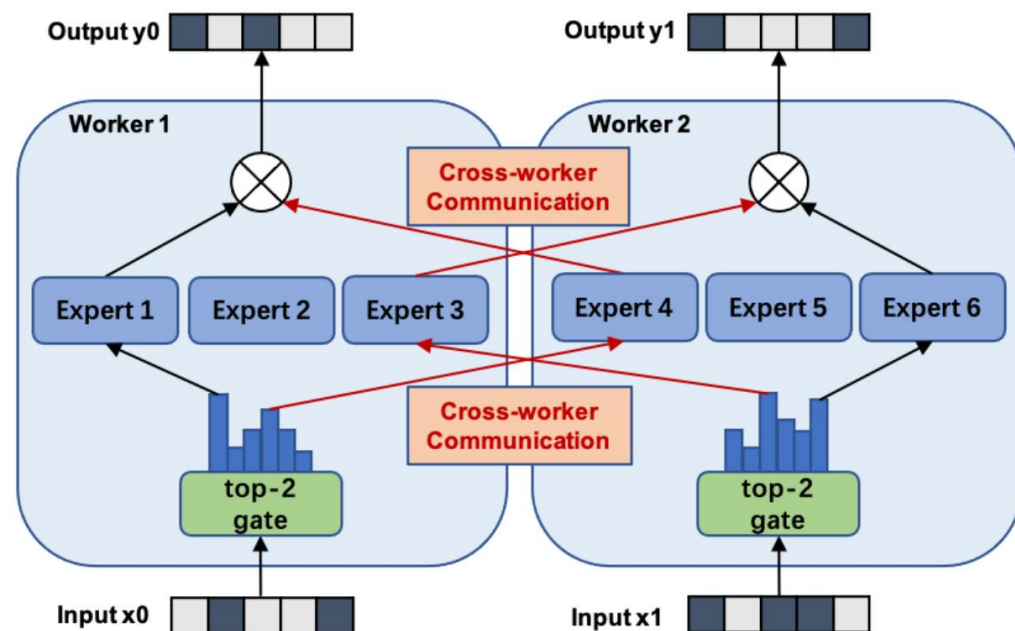
# Transformer-XL

- Transformer-XL can be viewed as special case of MoBA where
  - a pre-defined gating function that only attends to the most recent block
  - now grad to historical blocks



# FastMoE

- Expert-parallel
- Easy-to-use in Megatron



```
model = ...  
  
from fmoe.megatron import fmoeify  
model = fmoeify(model, fmoe_num_experts=<number of experts per worker>)  
  
train(model, ...)
```

He, Jiaao, **Jiezhong Qiu**, Aohan Zeng, Zhilin Yang, Jidong Zhai, and Jie Tang. "Fastmoe: A fast mixture-of-expert training system." arXiv preprint arXiv:2103.13262 (2021).  
<https://github.com/laekov/fastmoe>

# Implementation

---

**Algorithm 1** MoBA (Mixture of Block Attention) Implementation

---

**Require:** Query, key and value matrices  $\mathbf{Q}, \mathbf{K}, \mathbf{V} \in \mathbb{R}^{N \times h \times d}$ ; MoBA hyperparameters (block size  $B$  and top- $k$ );  $h$  and  $d$  denote the number of attention heads and head dimension. Also denote  $n = N/B$  to be the number of blocks.

- 1: // Split KV into blocks
  - 2:  $\{\tilde{\mathbf{K}}_i, \tilde{\mathbf{V}}_i\} = \text{split\_blocks}(\mathbf{K}, \mathbf{V}, B)$ , where  $\tilde{\mathbf{K}}_i, \tilde{\mathbf{V}}_i \in \mathbb{R}^{B \times h \times d}, i \in [n]$
  - 3: // Compute gating scores for dynamic block selection
  - 4:  $\bar{\mathbf{K}} = \text{mean\_pool}(\mathbf{K}, B) \in \mathbb{R}^{n \times h \times d}$
  - 5:  $\mathbf{S} = \mathbf{Q}\bar{\mathbf{K}}^\top \in \mathbb{R}^{N \times h \times n}$
  - 6: // Select blocks with causal constraint (no attention to future blocks)
  - 7:  $\mathbf{M} = \text{create\_causal\_mask}(N, n)$
  - 8:  $\mathbf{G} = \text{topk}(\mathbf{S} + \mathbf{M}, k)$
  - 9: // Organize attention patterns for computation efficiency
  - 10:  $\mathbf{Q}^s, \tilde{\mathbf{K}}^s, \tilde{\mathbf{V}}^s = \text{get\_self\_attn\_block}(\mathbf{Q}, \tilde{\mathbf{K}}, \tilde{\mathbf{V}})$
  - 11:  $\mathbf{Q}^m, \tilde{\mathbf{K}}^m, \tilde{\mathbf{V}}^m = \text{index\_select\_moba\_attn\_block}(\mathbf{Q}, \tilde{\mathbf{K}}, \tilde{\mathbf{V}}, \mathbf{G})$
  - 12: // Compute attentions separately
  - 13:  $\mathbf{O}^s = \text{flash\_attention\_varlen}(\mathbf{Q}^s, \tilde{\mathbf{K}}^s, \tilde{\mathbf{V}}^s, \text{causal}=\text{True})$
  - 14:  $\mathbf{O}^m = \text{flash\_attention\_varlen}(\mathbf{Q}^m, \tilde{\mathbf{K}}^m, \tilde{\mathbf{V}}^m, \text{causal}=\text{False})$
  - 15: // Combine results with online softmax
  - 16:  $\mathbf{O} = \text{combine\_with\_online\_softmax}(\mathbf{O}^s, \mathbf{O}^m)$
  - 17: **return**  $\mathbf{O}$
-

# Combine with online softmax

$$o_1 = \text{Attn}(q, K[I_1], V[I_1]) = \sum_{i=1}^a \frac{\exp(x_i)}{\exp(lse_1)} v_i$$

$$o_2 = \text{Attn}(q, K[I_2], V[I_2]) = \sum_{i=a+1}^{a+b} \frac{\exp(x_i)}{\exp(lse_2)} v_i$$

$$lse_1 = \log \sum_{i=1}^a \exp(x_i)$$

$$lse_2 = \log \sum_{i=a+1}^{a+b} \exp(x_i)$$

$$o = \text{Attn}(q, K[I_1||I_2], V[I_1||I_2]) = \sum_{i=1}^{a+b} \frac{\exp(x_i)}{\exp(lse)} v_i$$

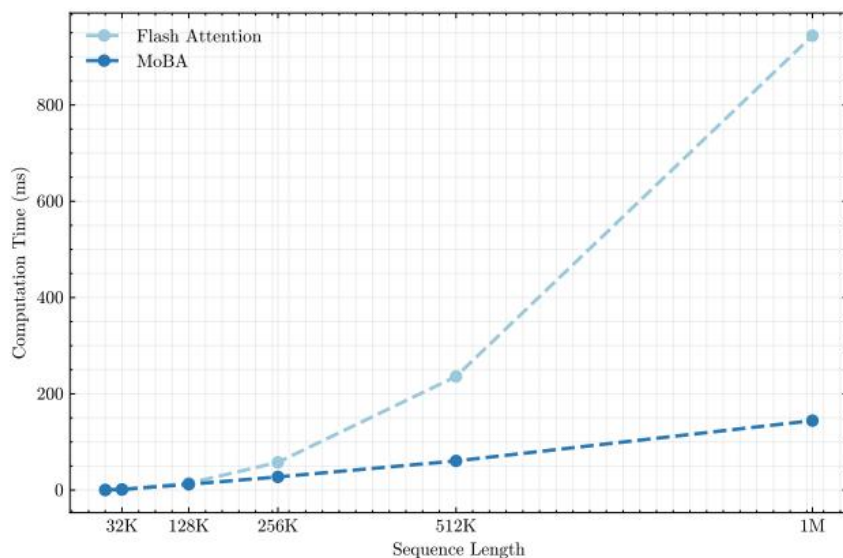
$$m = \max(lse_1, lse_2)$$

$$lse - m = \log [\exp(lse_1 - m) + \exp(lse_2 - m)]$$

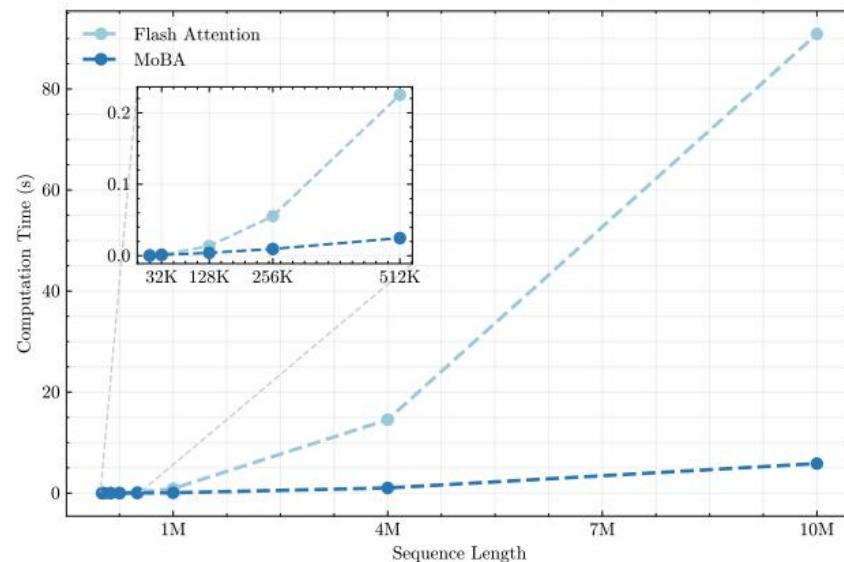
$$o = \exp[(lse_1 - m) - (lse - m)] o_1 + \exp[(lse_2 - m) - (lse - m)] o_2$$



# Speedup



(a)



(b)

Figure 2: Efficiency of MoBA vs. full attention (implemented with Flash Attention). **(a)** 1M Model speedup evaluation: Computation time scaling of MoBA versus Flash Attention on 1M model with increasing sequence lengths (8K-1M). **(b)** Fixed Sparsity Ratio scaling: Computation time scaling comparison between MoBA and Flash Attention across increasing sequence lengths (8K-10M), maintaining a constant sparsity ratio of 95.31% (fixed 64 MoBA blocks with variance block size and fixed top-k=3).



# Scaling Law

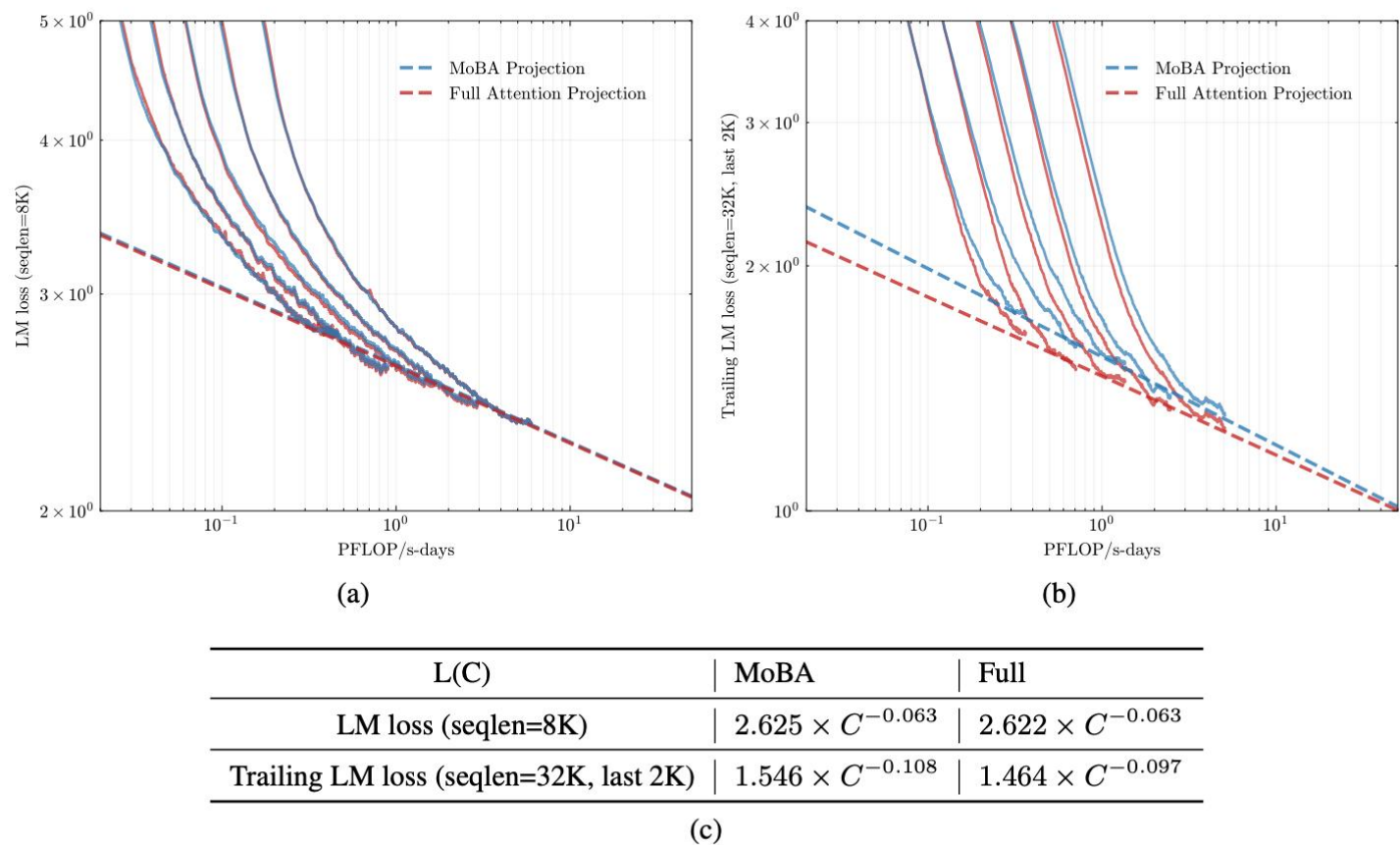
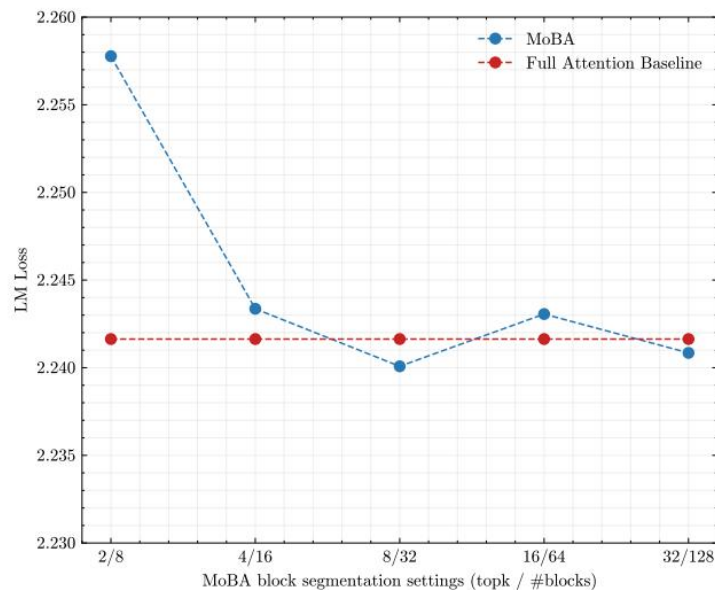


Figure 3: Scaling law comparison between MoBA and full attention. (a) LM loss on validation set (seqlen=8K); (b) trailing LM loss on validation set (seqlen=32K, last 1K tokens); (c) fitted scaling law curve.

# Fine-grained expert $\rightarrow$ fine-grained block

**Ablation Study on Fine-Grained Block Segmentation.** We further ablate the block granularity of MoBA. We carry out a series of experiments using a 1.5B parameter model with a 32K context length. The hyperparameters of block size and top-k are adjusted to maintain a consistent level of attention sparsity. Specifically, we divide the 32K context into 8, 16, 32, 64, and 128 blocks, and correspondingly select 2, 4, 8, 16, and 32 blocks, ensuring an attention sparsity of 75% across these configurations. As shown in Figure 4, MoBA's performance is significantly affected by block granularity. Specifically, there is a performance difference of  $1e-2$  between the coarsest-grained setting (selecting 2 blocks from 8) and the settings with finer granularity. These findings suggest that fine-grained segmentation appears to be a general technique for enhancing the performance of models within the MoE family, including MoBA.



# Hybrid Strategy of MoBA

- Hybrid Training: 90% tokens with MoBA + 10% with Full Attention
- Layer Hybrid: switch the last several layers from MoBA to full attention

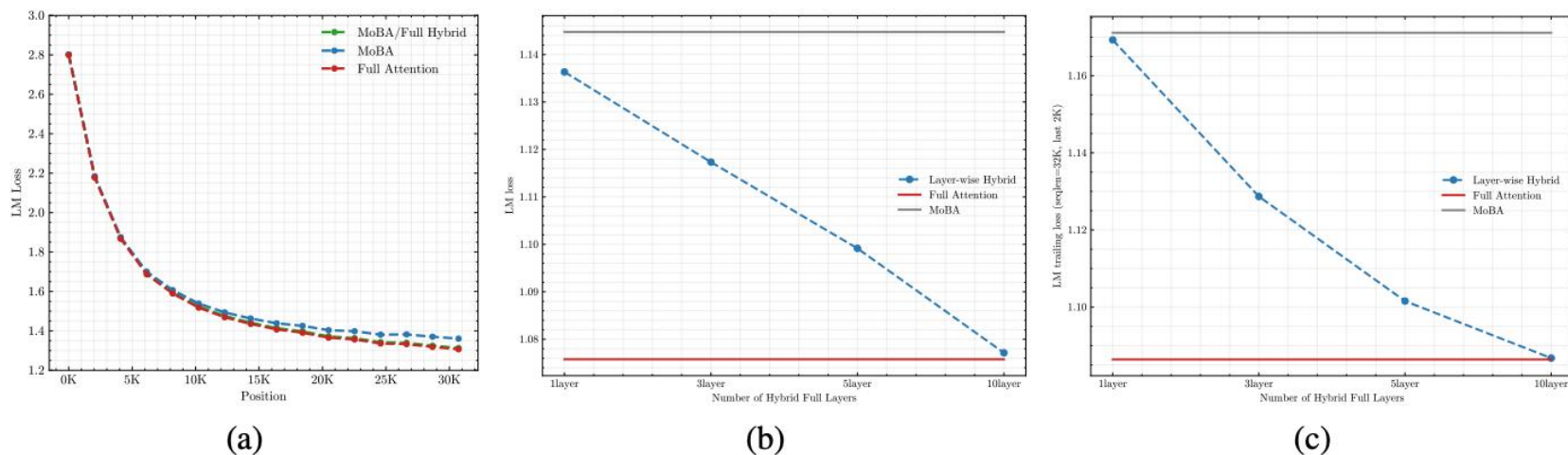


Figure 5: Hybrid of MoBA and full attention. (a) position-wise LM loss for MoBA, full attention, and MoBA/full hybrid training; (b) SFT LM loss w.r.t the number of full attention layers in layer-wise hybrid; (c) SFT trailing LM loss (seqLen=32K, last 2K) w.r.t the number of full attention layers in layer-wise hybrid.

# The post-training recipes (from short to long)

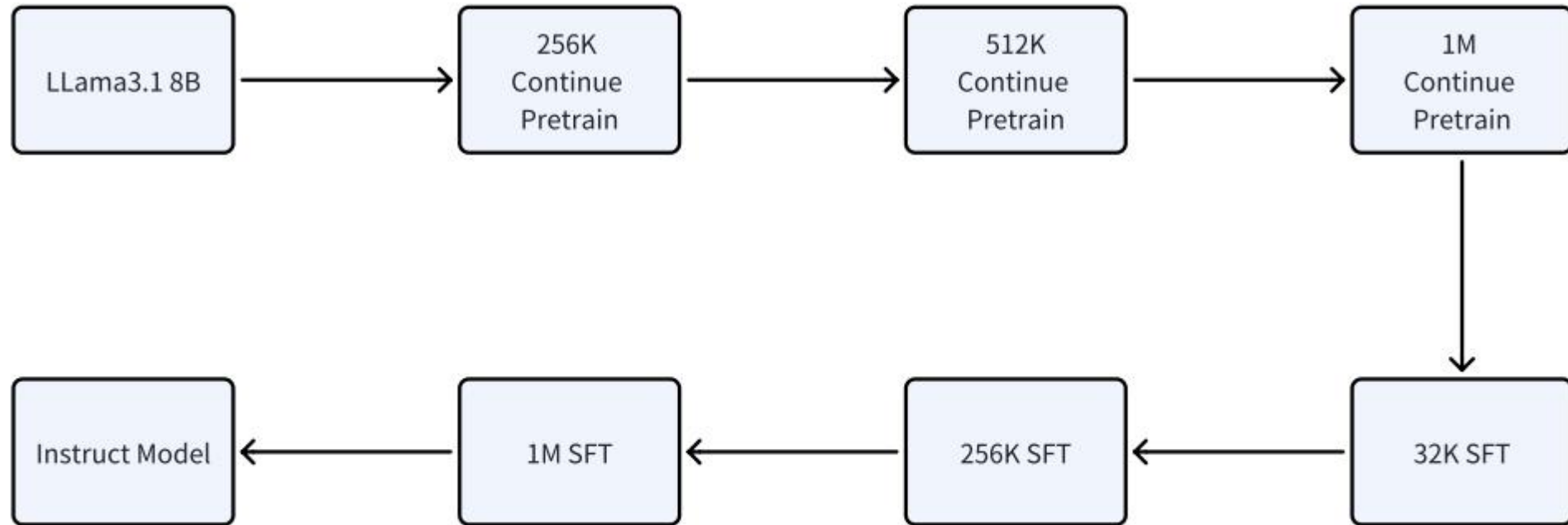


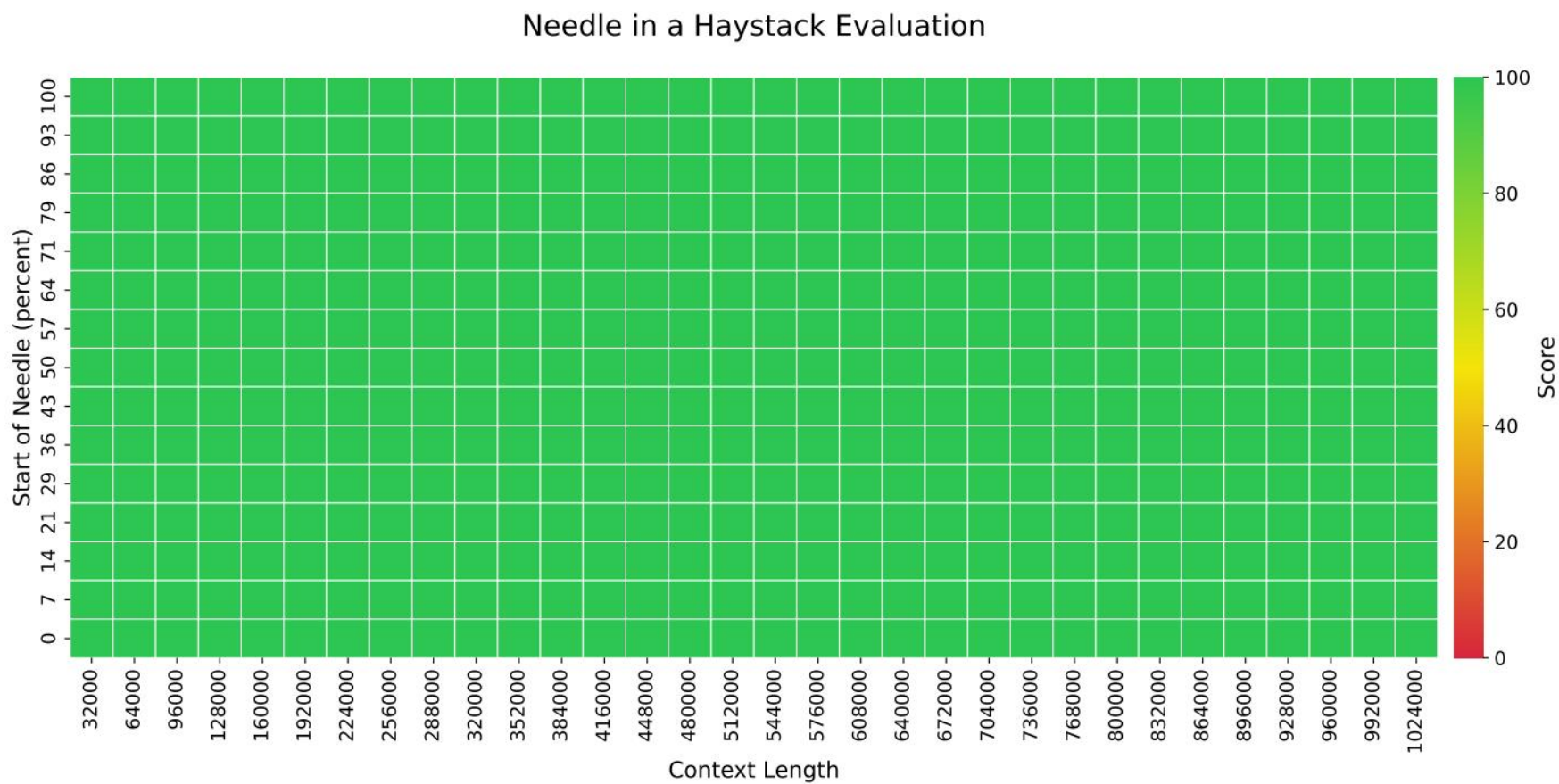
Figure 6: The continual pre-training and SFT recipes.

# MoBA v.s. Full

<b>Benchmark</b>	<b>Llama-8B-1M-MoBA</b>	<b>Llama-8B-1M-Full</b>
AGIEval [0-shot]	0.5144	<b>0.5146</b>
BBH [3-shot]	0.6573	<b>0.6589</b>
CEval [5-shot]	<b>0.6273</b>	0.6165
GSM8K [5-shot]	<b>0.7278</b>	0.7142
HellaSWAG [0-shot]	0.8262	<b>0.8279</b>
Loogle [0-shot]	<b>0.4209</b>	0.4016
Competition Math [0-shot]	0.4254	<b>0.4324</b>
MBPP [3-shot]	<b>0.5380</b>	0.5320
MBPP Sanitized [0-shot]	<b>0.6926</b>	0.6615
MMLU [0-shot]	0.4903	<b>0.4904</b>
MMLU Pro [5-shot][CoT]	0.4295	<b>0.4328</b>
OpenAI HumanEval [0-shot][pass@1]	0.6951	<b>0.7012</b>
SimpleQA [0-shot]	0.0465	<b>0.0492</b>
TriviaQA [0-shot]	<b>0.5673</b>	0.5667
LongBench @32K [0-shot]	<b>0.4828</b>	0.4821
RULER @128K [0-shot]	0.7818	<b>0.7849</b>

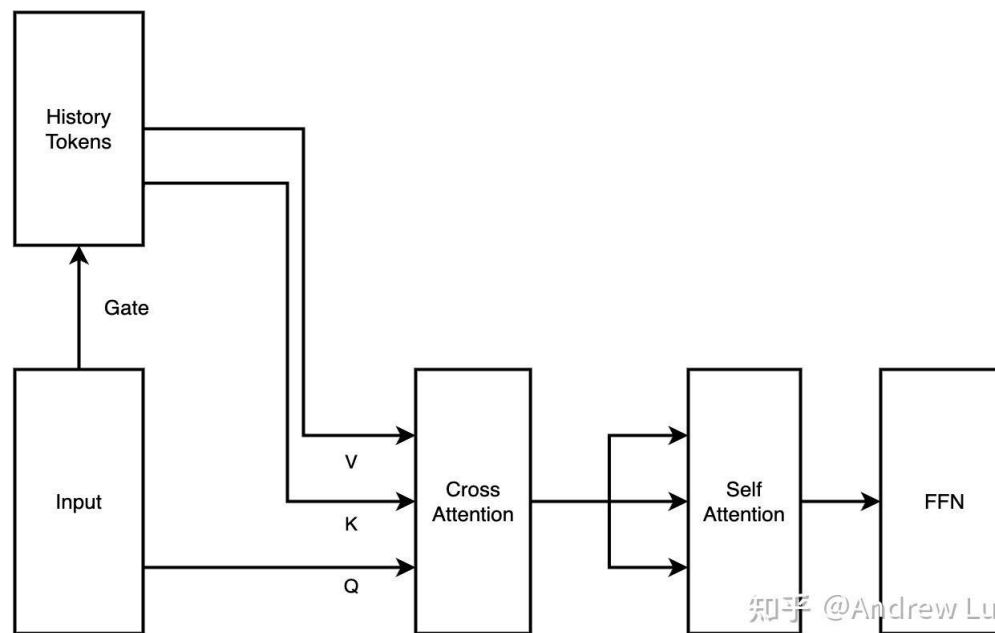
Table 2: Performance comparison between MoBA and full Attention across different evaluation benchmarks.

# From MoE to MoBA



# MoBA v0.5

- Cross Attention + Self Attention + FFN
- Bad news: introduce new layers/parameters

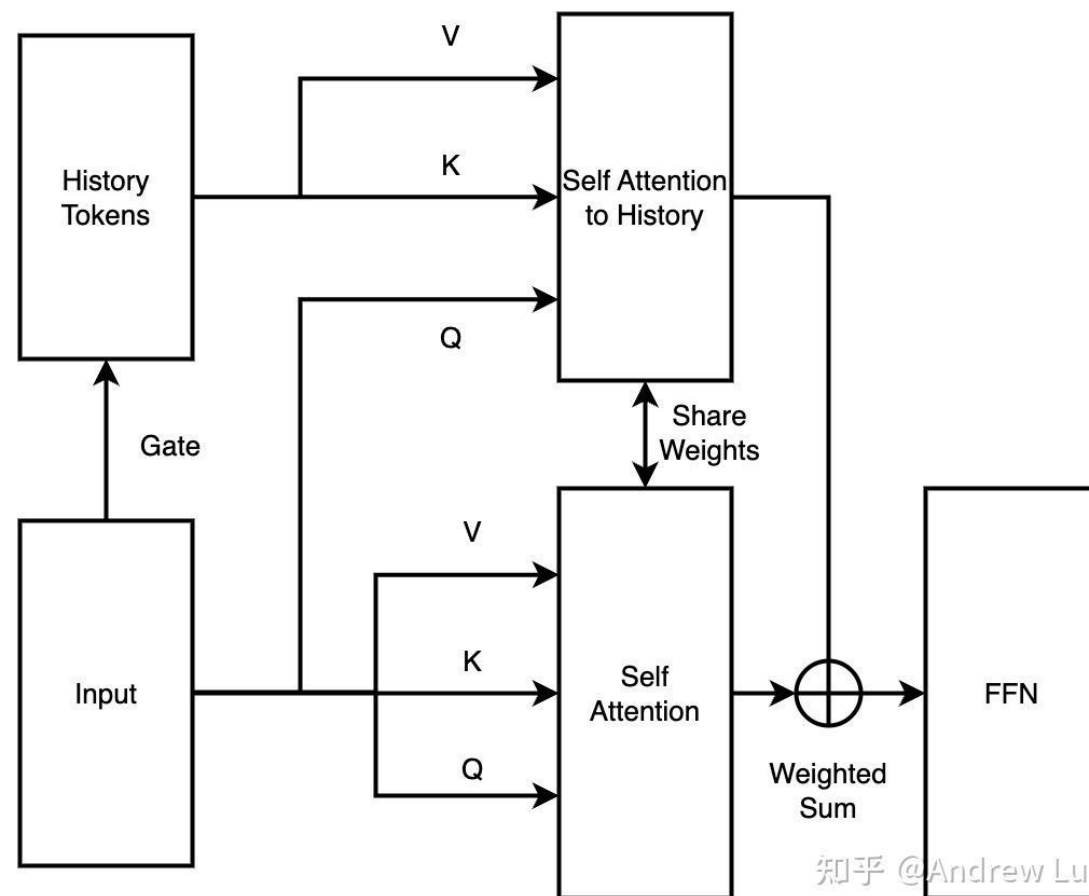




# MoBA v1.0

- MoBA but with a MoE-style weighted sum

$$o = \sum_i g_i \text{Attn}(q, K[I_i], V[I_i])$$



# MoBA with Context Parallel

- MoE has expert parallel to support a large number of experts
- Place KV blocks distributedly ...
- Bad news: load balance...

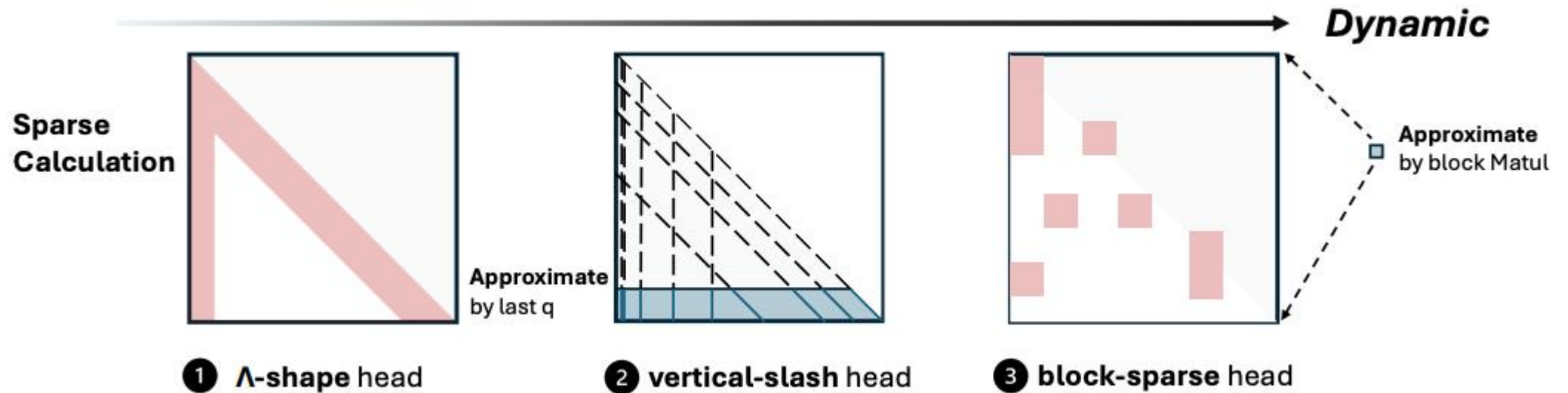


Figure 4: The three sparse methods in MInference.

Thank You!