

# BLASST: Dynamic BLocked Attention Sparsity via Softmax Thresholding

Jiayi Yuan<sup>\*1</sup>, **Cameron Shinn**<sup>\*2</sup>, Kai Xu<sup>3</sup>, Jingze Cui<sup>3</sup>, George Klimiashvili<sup>3</sup>, Guangxuan Xiao<sup>3</sup>, Perkz Zheng<sup>3</sup>, Bo Li<sup>3</sup>, Yuxin Zhou<sup>3</sup>, Zhouhai Ye<sup>3</sup>, Weijie You<sup>3</sup>, Tian Zheng<sup>3</sup>, Dominic Brown<sup>3</sup>, Pengbo Wang<sup>3</sup>, Markus Hoehnerbach<sup>4</sup>, Richard Cai<sup>3</sup>, Julien Demouth<sup>3</sup>, John D. Owens<sup>2</sup>, Xia Hu<sup>1</sup>, Song Han<sup>3</sup>, Timmy Liu<sup>3</sup>, Huizi Mao<sup>3</sup>

<sup>1</sup>Rice University, <sup>2</sup>UC Davis, <sup>3</sup>NVIDIA, <sup>4</sup>Meta



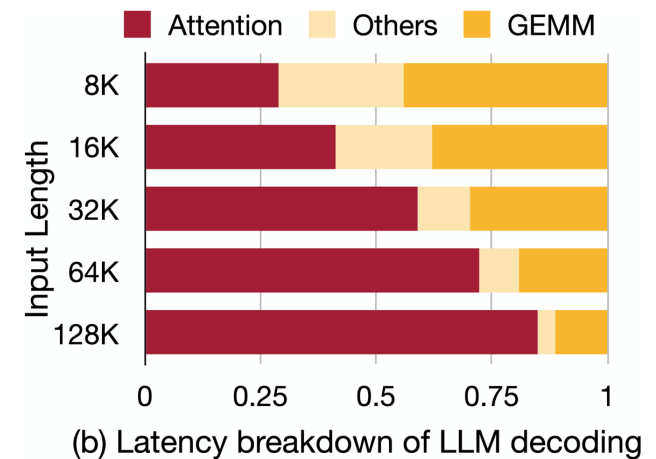
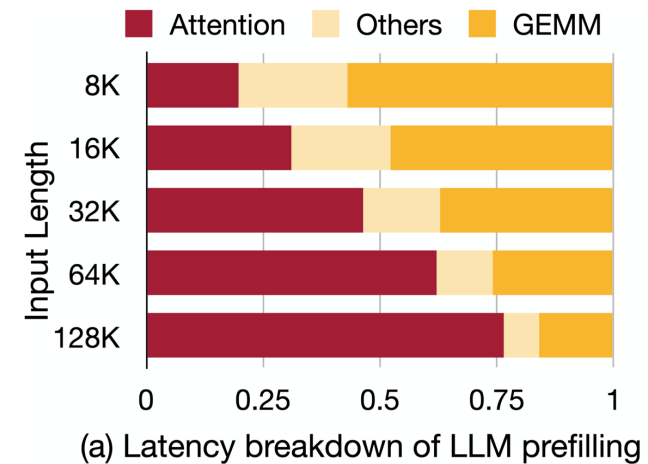
# Background – LLMs Continue to Scale

LLM inference context lengths are growing dramatically.

- System prompt tokens 📶
- Reasoning and chain-of-thought tokens 📶
- Content tokens (code, images, video) 📶
- VLM and diffusion models are rising in popularity and consume many tokens.

*Attention dominates prefill and decode time as sequence length increases. Data from Llama-3-8B on A100 GPU.*

*Figure credit: LServe, Yang et al., 2025*

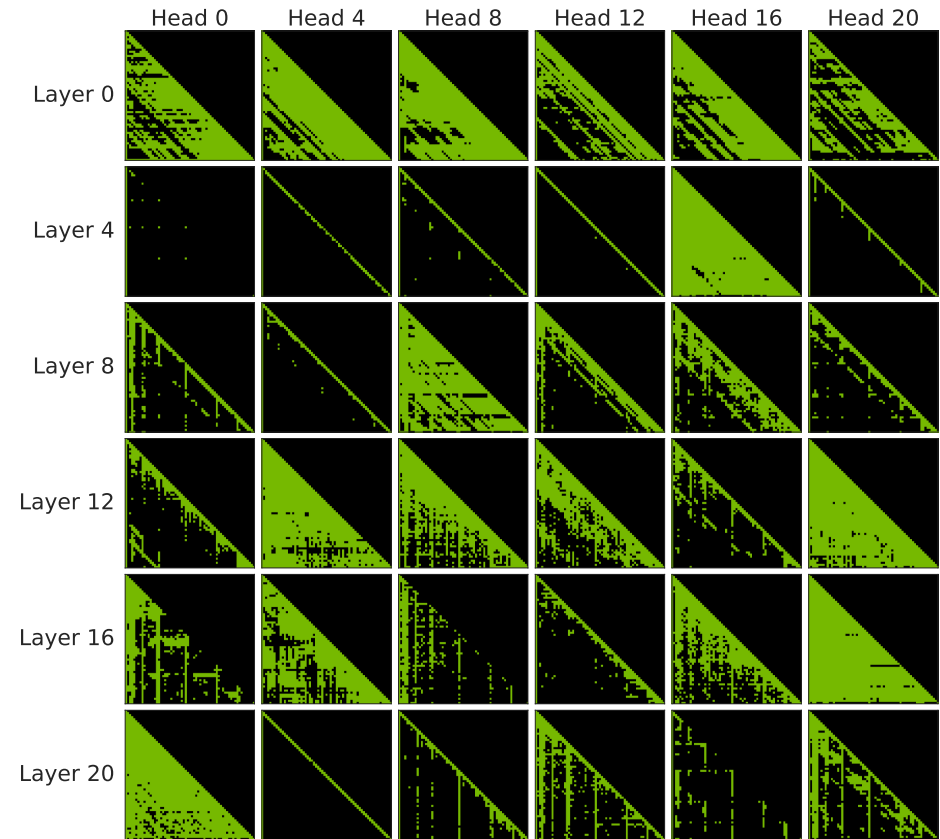


# Background – Attention Is Sparse

Attention is intrinsically sparse, with specific patterns being observed.

- Global, windowed and random sparsity (Big Bird, Zaheer et al., 2020).
- Attention sinks (StreamingLLM, Xiao et al., 2023).
- $\Lambda$ -shape, vertical-slash, and block-sparse patterns (MInference, Jiang et al., 2024).
- DeepSeek dynamic sparse variants: HCA/CSA/DSA

Large body of existing work on how to leverage attention sparsity.



Llama 3.1 8B sparsity patterns from RULER sample

# Barriers to Adoption

Existing works can be difficult to deploy for a few reasons...

# Barriers to Adoption

Existing works can be difficult to deploy for a few reasons...

1. Pre-computation to decide sparsity pattern reduces speedup.

# Barriers to Adoption

Existing works can be difficult to deploy for a few reasons...

1. Pre-computation to decide sparsity pattern reduces speedup.
2. Requiring training for added layers or a new architectures.

# Barriers to Adoption

Existing works can be difficult to deploy for a few reasons...

1. Pre-computation to decide sparsity pattern reduces speedup.
2. Requiring training for added layers or a new architectures.
3. Exclusive to either prefill or decode phase.

# Barriers to Adoption

Existing works can be difficult to deploy for a few reasons...

1. Pre-computation to decide sparsity pattern reduces speedup.
2. Requiring training for added layers or a new architectures.
3. Exclusive to either prefill or decode phase.
4. Not optimized or supported on modern GPUs.

# Barriers to Adoption

Existing works can be difficult to deploy for a few reasons...

1. Pre-computation to decide sparsity pattern reduces speedup.
2. Requiring training for added layers or a new architectures.
3. Exclusive to either prefill or decode phase.
4. Not optimized or supported on modern GPUs.
5. Intrusive modifications to attention APIs that are hard for frameworks to adopt.

# Barriers to Adoption

Existing works can be difficult to deploy for a few reasons...

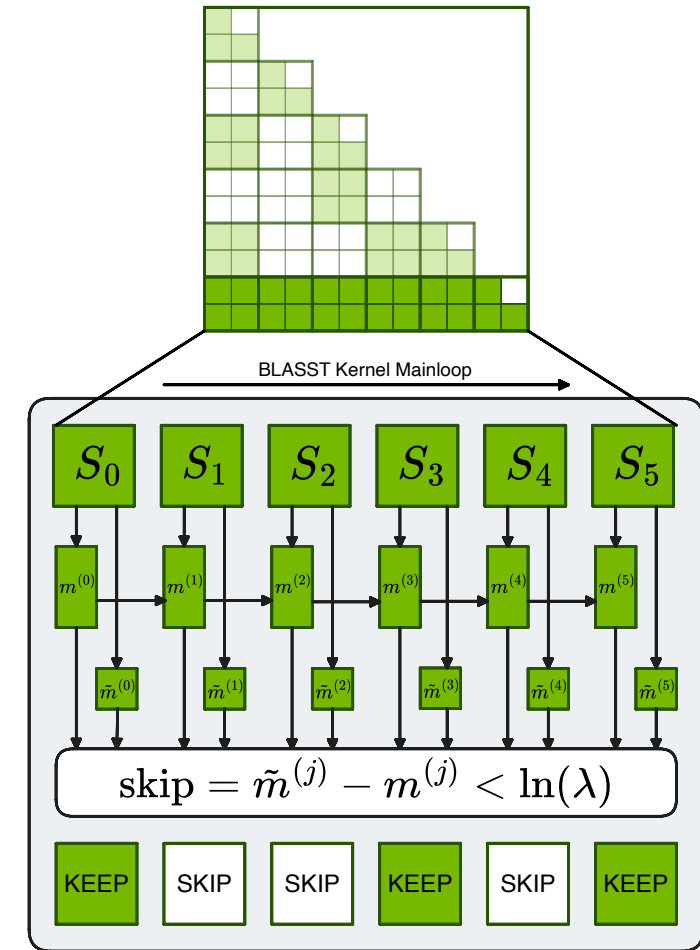
1. Pre-computation to decide sparsity pattern reduces speedup.
2. Requiring training for added layers or a new architectures.
3. Exclusive to either prefill or decode phase.
4. Not optimized or supported on modern GPUs.
5. Intrusive modifications to attention APIs that are hard for frameworks to adopt.

How can we avoid these and make sparse attention easier to use?

# Our Solution: BLASST

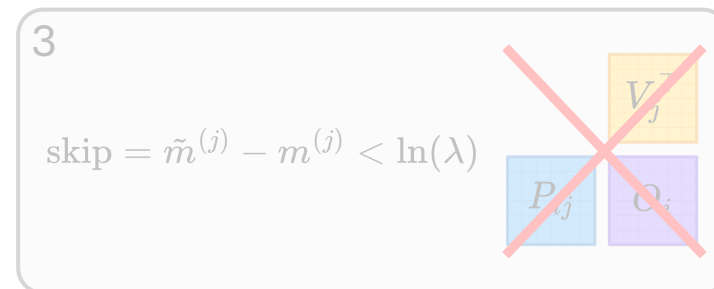
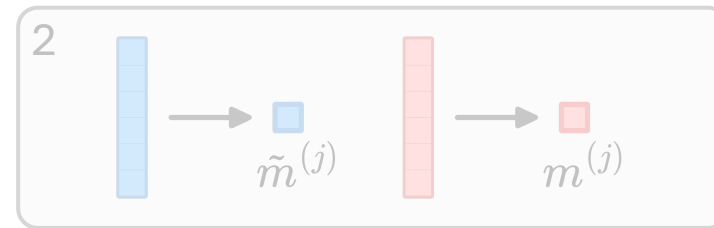
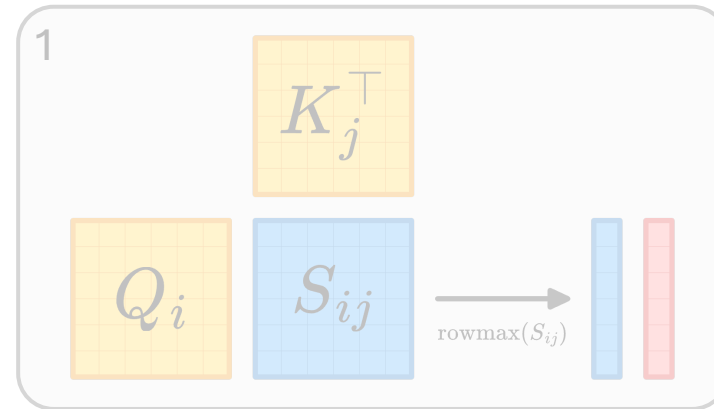
Dynamic block sparse attention using softmax thresholding

- Minimal impact on accuracy.
- Accelerates both **prefill** and **decode** stages with **~no overhead**.
- **Training-free**.
- **Blackwell and Hopper** kernel support.
- **Seamless integration** existing frameworks, with a simple kernel backend swap.



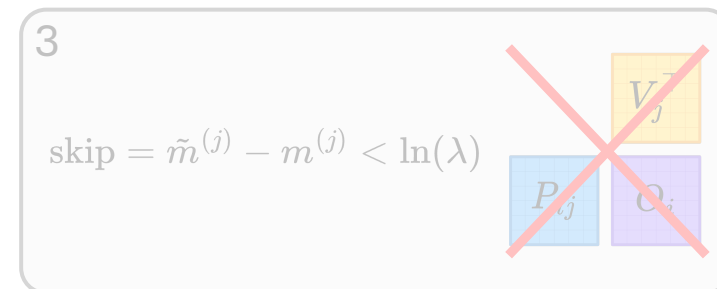
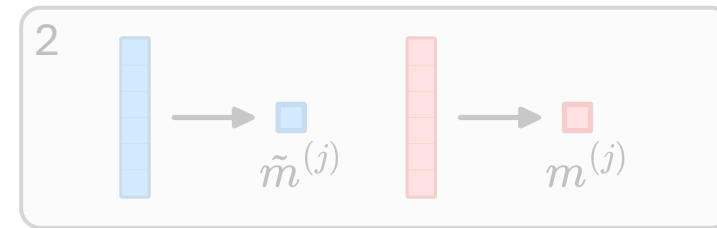
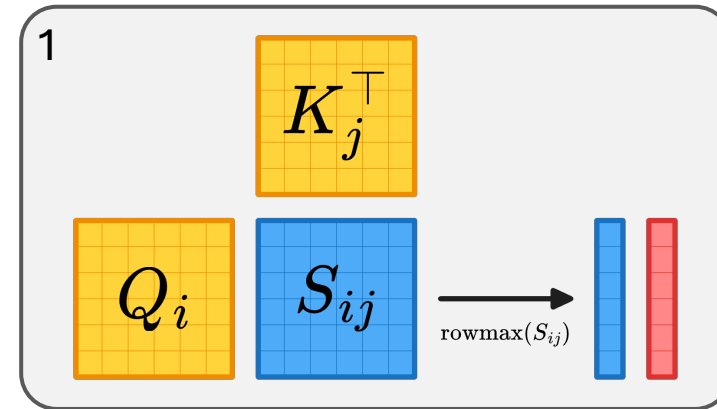
# BLASST Algorithm

1. After  $QK^T$ , online softmax already computes **local row max** of  $S$ , and updates **running row max**.
2. Take max across rows for both.
3. If **local block max** is below **running block max** by more than  $\ln(\lambda)$ : **Skip softmax, V load and  $PV^T$  matmul.**



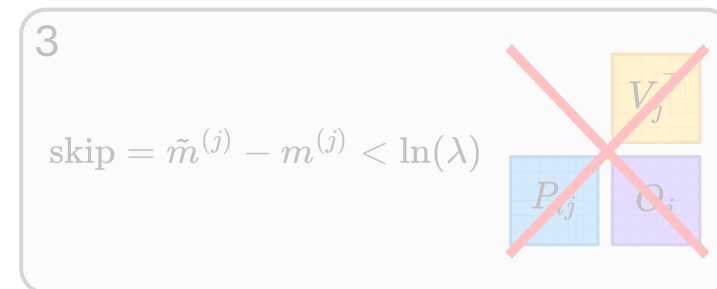
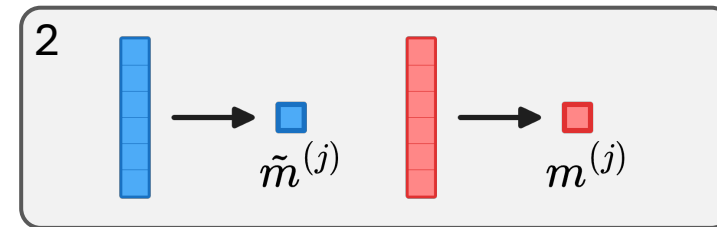
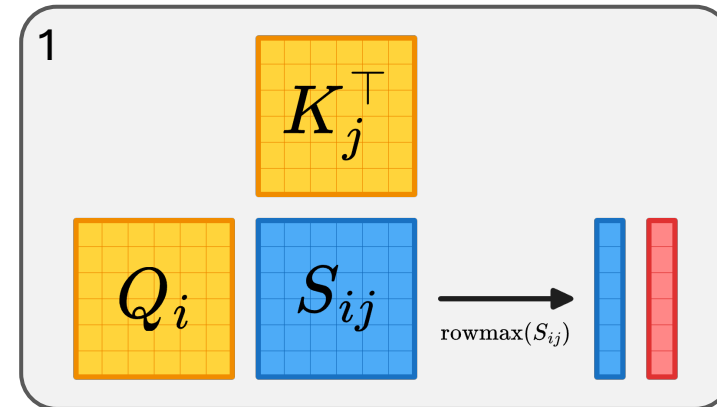
# BLASST Algorithm

1. After  $QK^T$ , online softmax already computes **local row max** of  $S$ , and updates **running row max**.
2. Take max across rows for both.
3. If **local block max** is below **running block max** by more than  $\ln(\lambda)$ : **Skip softmax, V load and  $PV^T$  matmul.**



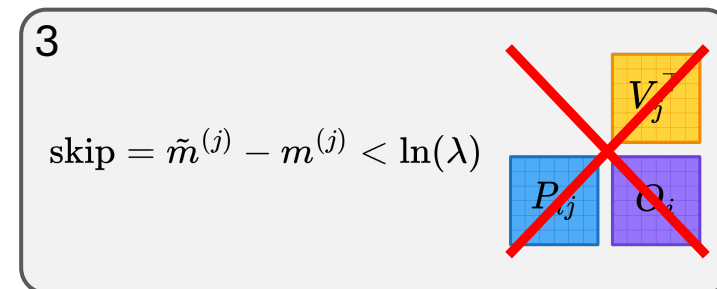
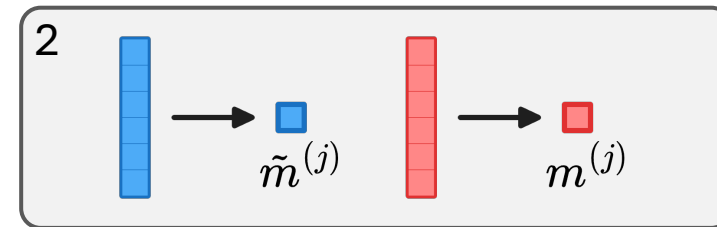
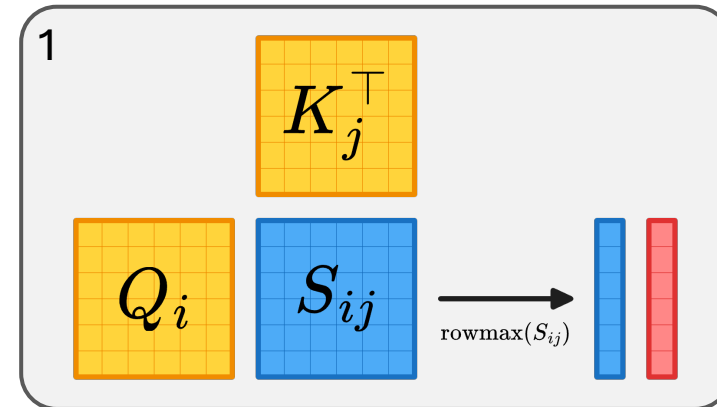
# BLASST Algorithm

1. After  $QK^T$ , online softmax already computes **local row max** of  $S$ , and updates **running row max**.
2. Take max across rows for both.
3. If **local block max** is below **running block max** by more than  $\ln(\lambda)$ : **Skip softmax, V load and  $PV^T$  matmul.**



# BLASST Algorithm

1. After  $QK^T$ , online softmax already computes **local row max** of  $S$ , and updates **running row max**.
2. Take max across rows for both.
3. If **local block max** is below **running block max** by more than  $\ln(\lambda)$ : **Skip softmax, V load and  $PV^T$  matmul.**



# Minimal Changes to FlashAttention Algorithm

- Very minor algorithmic changes to support skipping (denoted in **green**).
- Aim to introduce as few extra instructions as possible. Re-use rowmax operation.
- Same block size as existing FA, maps directly onto tiled computation.

---

## Algorithm 1 FlashAttention with BLASST

---

**Require:** Query blocks  $\{Q_i\}_{i=1}^{T_r}$ , Key blocks  $\{K_j\}_{j=1}^{T_c}$ ,  
Value blocks  $\{V_j\}_{j=1}^{T_c}$ , **threshold  $\lambda$**

**Ensure:** Output blocks  $\{O_i\}_{i=1}^{T_r}$

```

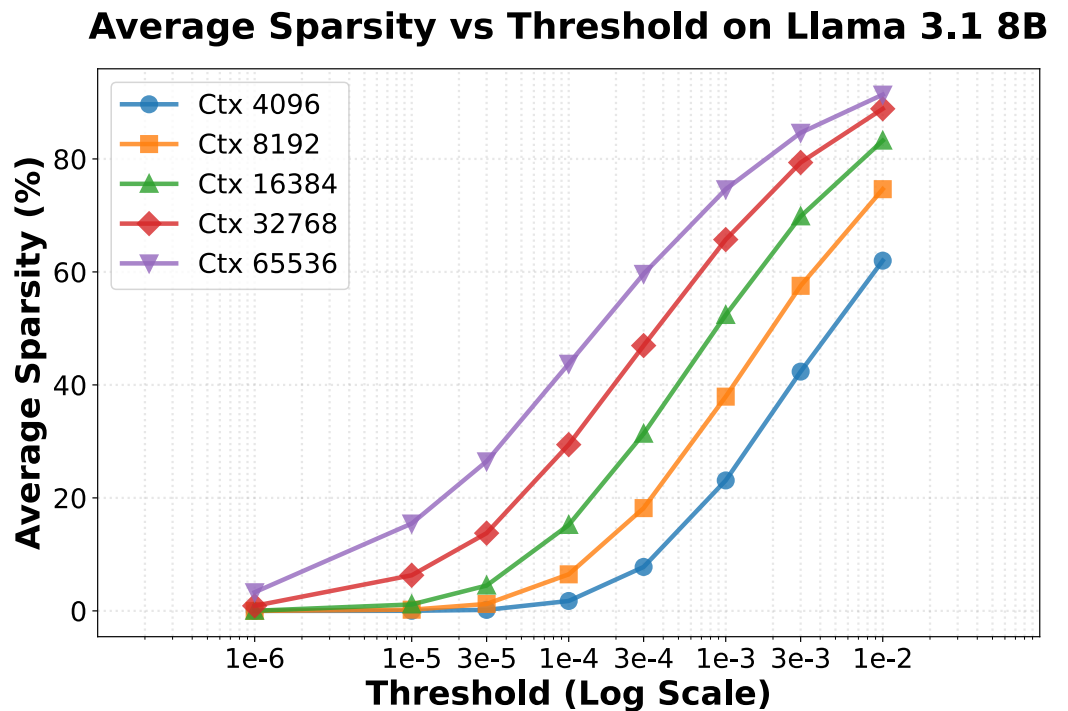
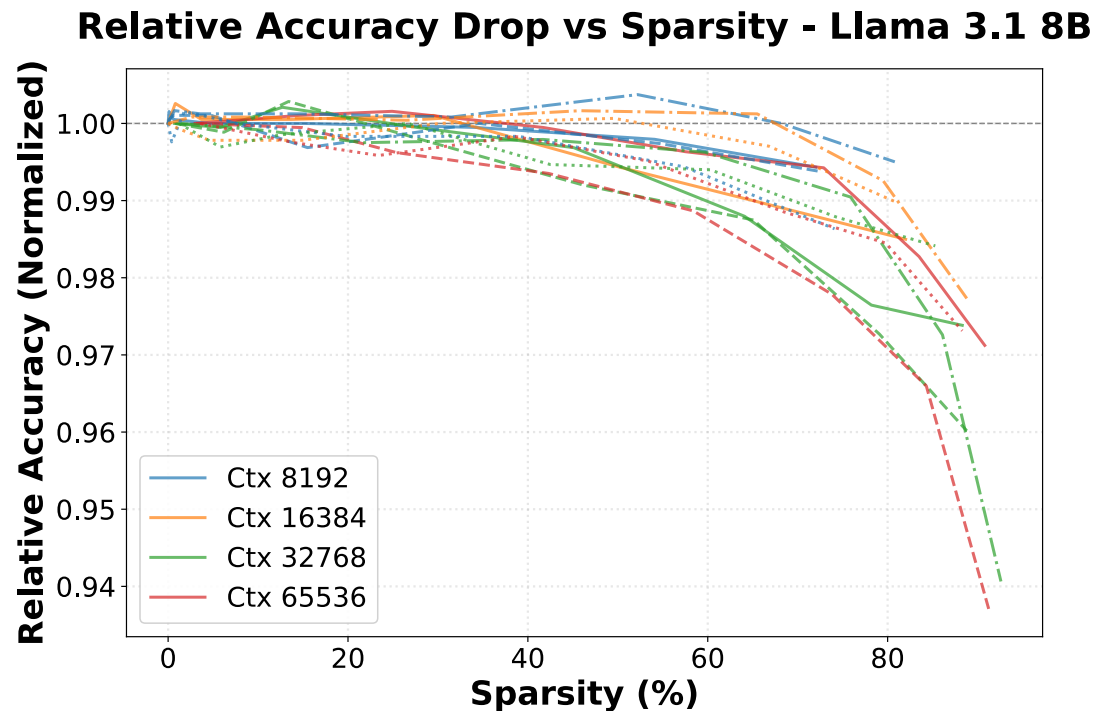
1: for  $i = 1$  to  $T_r$  do
2:   Initialize  $m_i^{(0)} = -\infty$ ,  $O_i^{(0)} = 0$ ,  $l_i^{(0)} = 0$ 
3:   for  $j = 1$  to  $T_c$  do
4:     Compute  $S_{ij} = Q_i K_j^\top$  ▷ Attention scores
5:      $\tilde{m}_i^{(j)} = \text{rowmax}(S_{ij})$  ▷ Local maximum
6:      $m_i^{(j)} = \max(m_i^{(j-1)}, \tilde{m}_i^{(j)})$  ▷ Running maximum
7:     if  $\tilde{m}_i^{(j)} - m_i^{(j)} < \ln(\lambda)$  then
8:       continue ▷ Skip this block
9:     end if
10:     $\tilde{P}_{ij} = \exp(S_{ij} - m_i^{(j)})$  ▷ Compute attn. weights
11:     $l_i^{(j)} = e^{m_i^{(j-1)} - m_i^{(j)}} l_i^{(j-1)} + \text{rowsum}(\tilde{P}_{ij})$ 
12:     $O_i^{(j)} = e^{m_i^{(j-1)} - m_i^{(j)}} O_i^{(j-1)} + \tilde{P}_{ij} V_j$ 
13:  end for
14:   $O_i = O_i^{(T_c)} / l_i^{(T_c)}$  ▷ Final normalization
15: end for
16: return  $\{O_i\}_{i=1}^{T_r}$ 

```

---

# How To Determine Threshold?

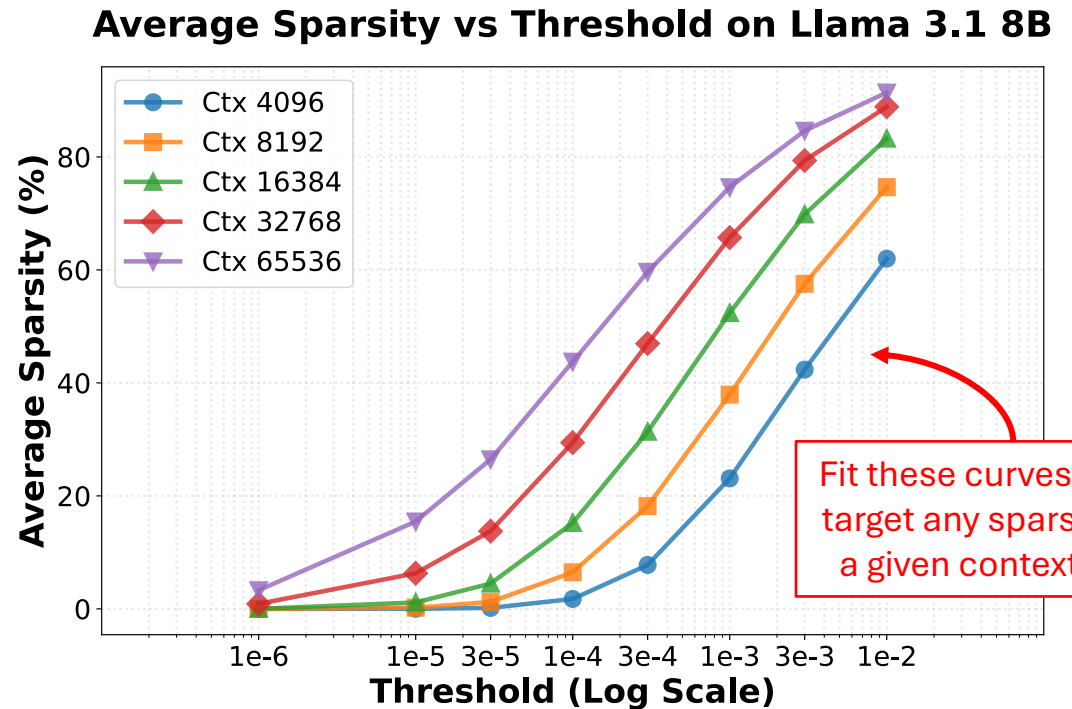
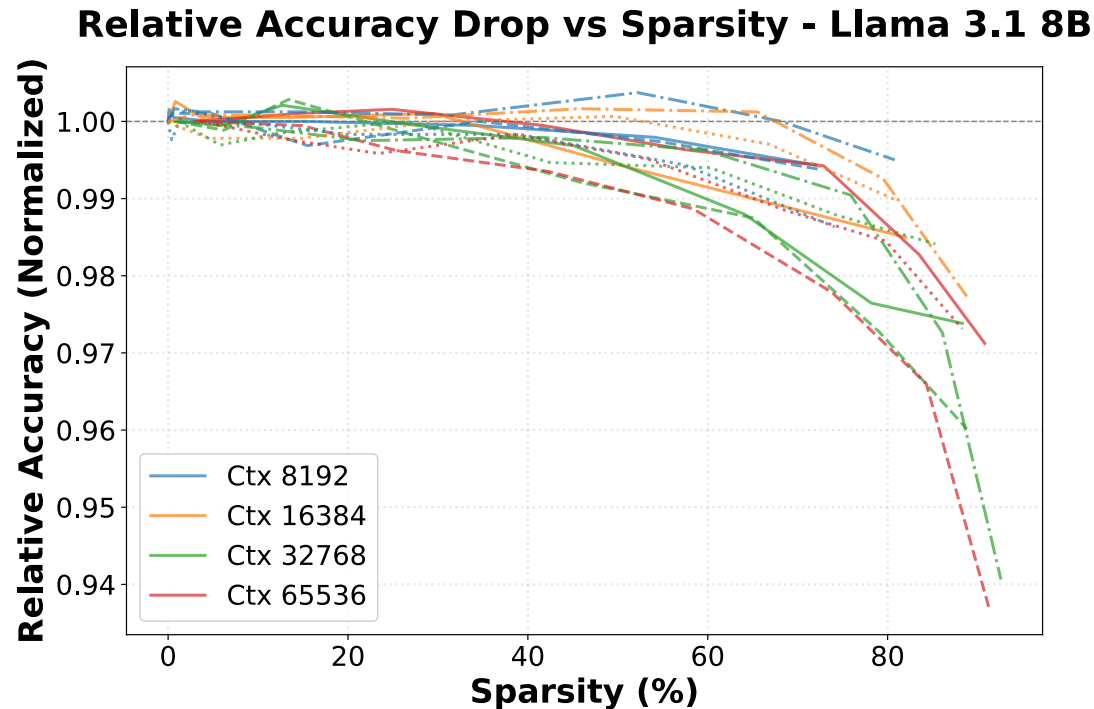
The right threshold depends on the **model**, **context length** and desired **accuracy**



(Left) Relative accuracy drop across different datasets and context lengths shows consistent degradation patterns as observed sparsity increases. All curves are normalized to their initial accuracy. (Right) Relationship between threshold and observed sparsity levels across different sequence lengths, demonstrating the need for threshold calibration to maintain fixed sparsity across varying contexts.

# How To Determine Threshold?

The right threshold depends on the **model**, **context length** and desired **accuracy**



(Left) Relative accuracy drop across different datasets and context lengths shows consistent degradation patterns as observed sparsity increases. All curves are normalized to their initial accuracy. (Right) Relationship between threshold and observed sparsity levels across different sequence lengths, demonstrating the need for threshold calibration to maintain fixed sparsity across varying contexts.

# Threshold Calibration

One calibration pass  $\rightarrow$  fit  $\lambda \cdot L = a \cdot \exp(b \cdot S)$   $\rightarrow$  dial in any target sparsity at runtime

## Calibration Algorithm

1

### Sweep thresholds in one pass

Run the calibration set once, trying many candidate thresholds together.

2

### Record sparsity vs threshold

For each (threshold, sequence length), log the observed sparsity.

3

### Fit the curve

Fit  $\lambda \cdot L = a \cdot \exp(b \cdot S)$  — output is two scalars (a, b) per model.

4

### Deploy at any target sparsity

Pick S  $\rightarrow$  threshold is set automatically; rescales with context length L.

# Threshold Calibration

One calibration pass  $\rightarrow$  fit  $\lambda \cdot L = a \cdot \exp(b \cdot S)$   $\rightarrow$  dial in any target sparsity at runtime

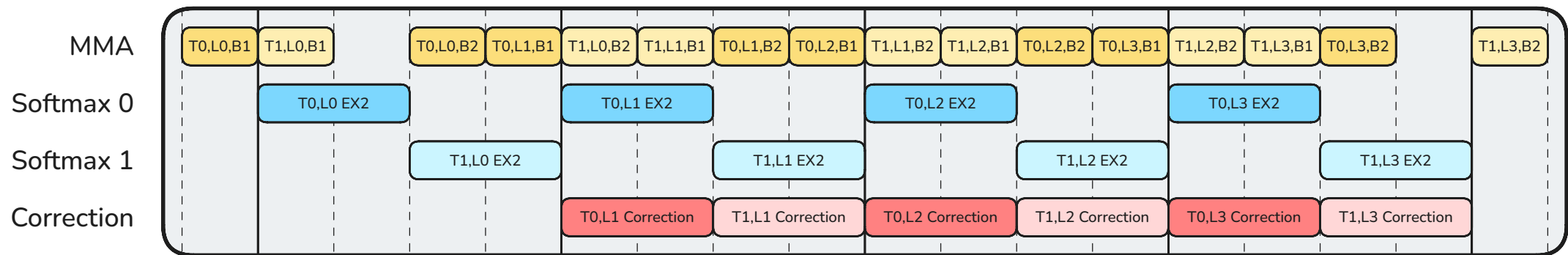
## Calibration Algorithm

- 1 Sweep thresholds in one pass**  
Run the calibration set once, trying many candidate thresholds together.
- 2 Record sparsity vs threshold**  
For each (threshold, sequence length), log the observed sparsity.
- 3 Fit the curve**  
Fit  $\lambda \cdot L = a \cdot \exp(b \cdot S)$  — output is two scalars (a, b) per model.
- 4 Deploy at any target sparsity**  
Pick S  $\rightarrow$  threshold is set automatically; rescales with context length L.

**One scalar per model.**  
 $\lambda$  rescales with L — no per-context recalibration, predictable speedup.

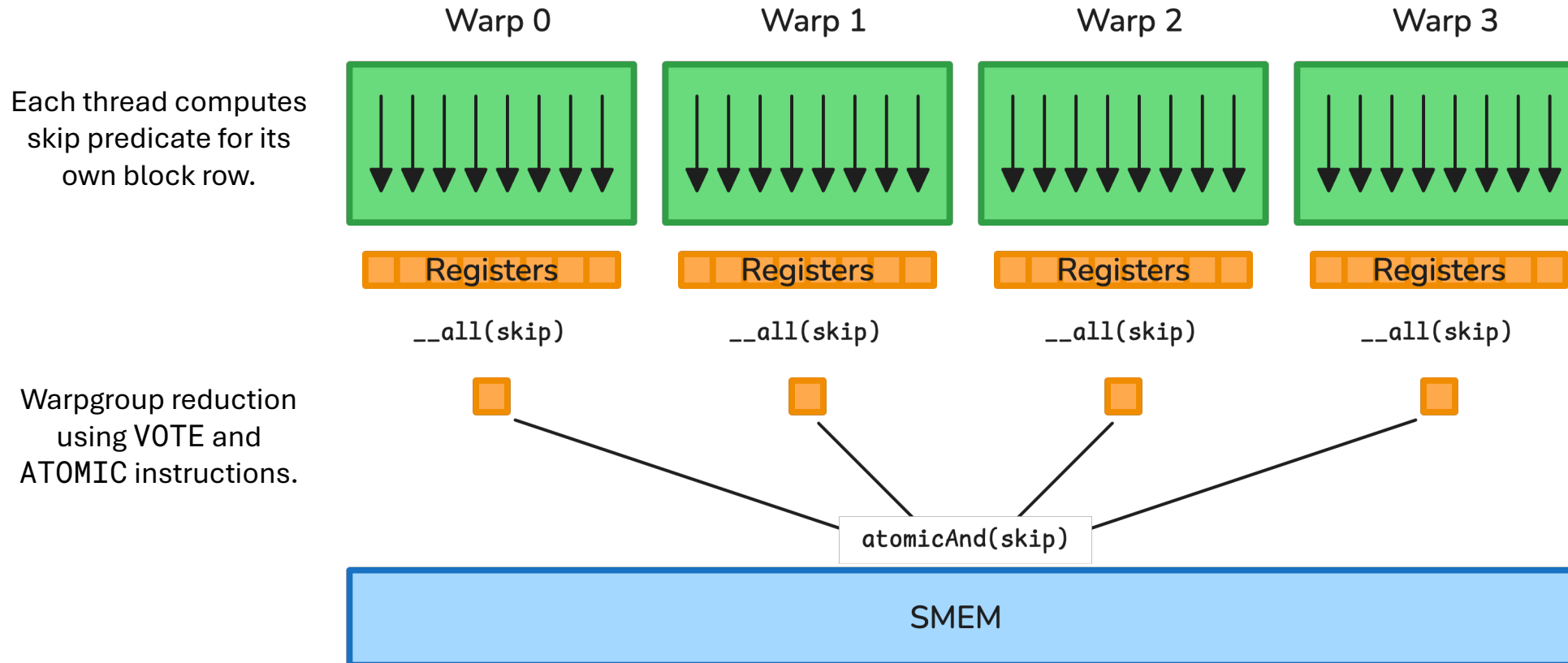
# Background – FlashAttention

- Avoid storing and loading  $N^2$  attention matrix in HBM by fusing surrounding operators,  $\text{softmax}(QK^T)V$ .
- Online softmax to avoid keeping entire attention row for rowmax.
- MMA and Softmax are on critical path for Blackwell: Specialize into their own warps (FlashAttention-4, Zadouri et al., 2026).



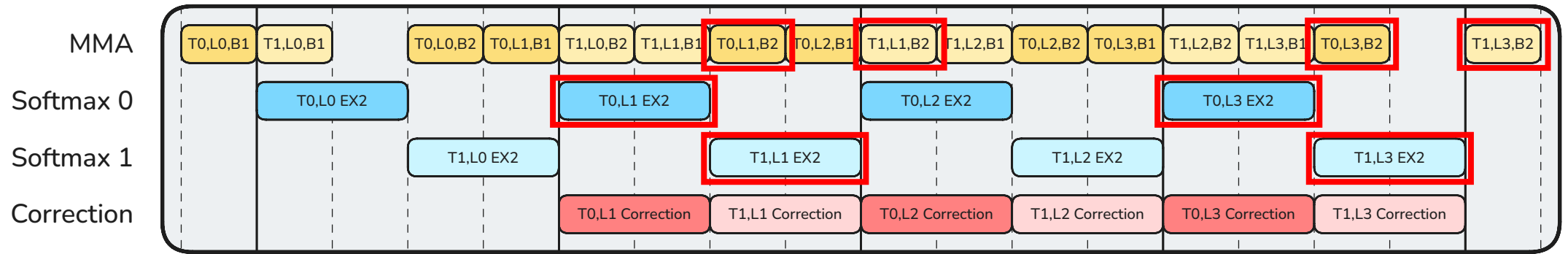
# Skip Predicate Reduction

Compute block max off the critical path (in correction warpgroup)

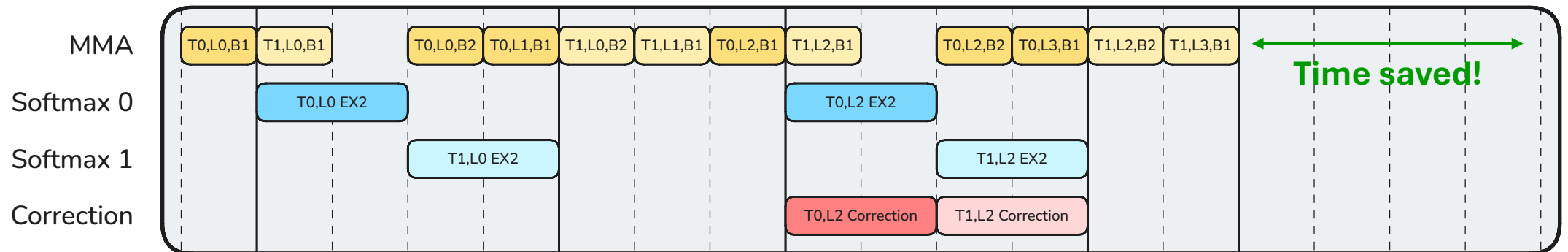


# Prefill Pipeline Comparison

Compute bound → Skip matmul and softmax

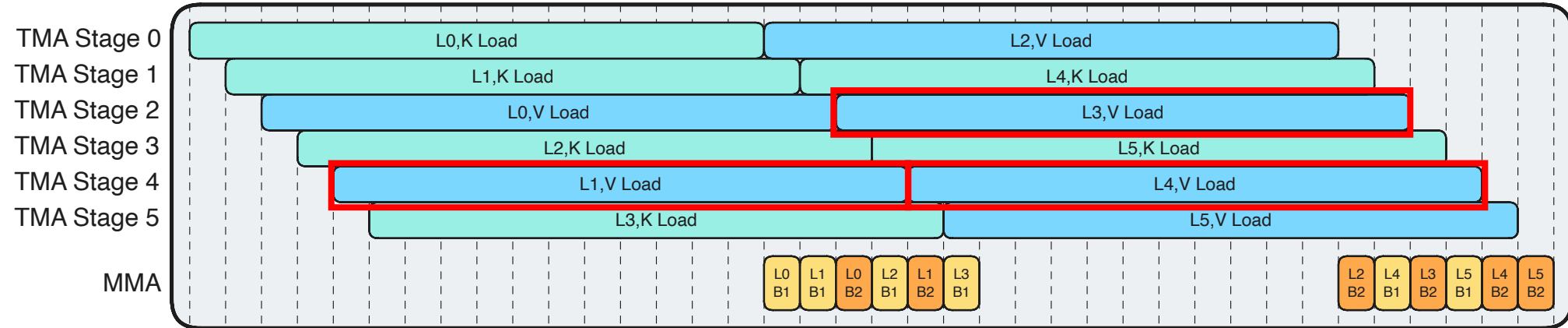


Skip loop 1 and loop 3 for example

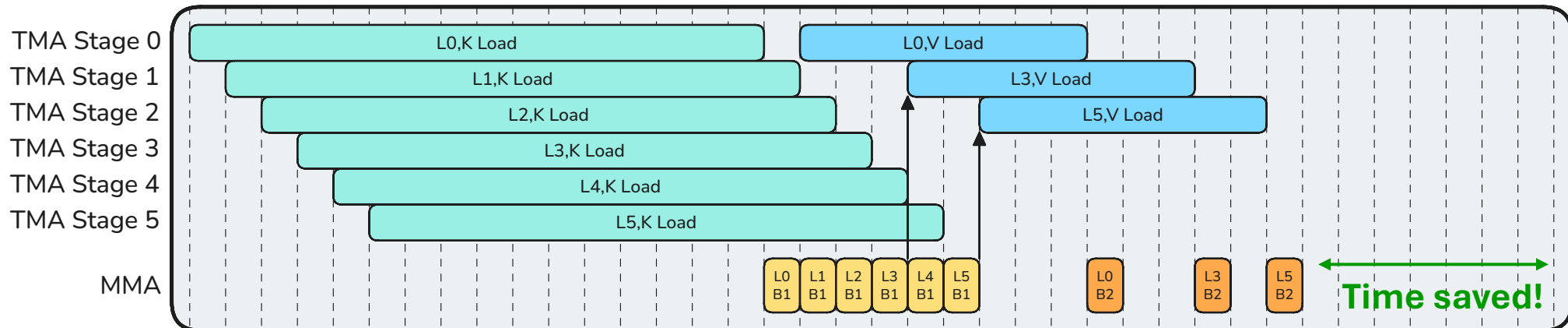


# Decode Pipeline Comparison

Memory bound → Skip V loads



Skip loop 1, 3 and 4 for example



Enable BLASST

\*For compute bound decoding, e.g. MLA, we can also skip matmul and softmax

# BLASST Preserves Accuracy

| Model        | Target Sparsity | Prefill Phase |           | Decode Phase |          |       | Prefill + Decode Phase |           |
|--------------|-----------------|---------------|-----------|--------------|----------|-------|------------------------|-----------|
|              |                 | RULER-32K     | LongBench | MATH500      | AIME2024 | GPQA  | RULER-32K              | LongBench |
| Llama-3.1-8B | Dense           | 92.33         | 31.40     | 73.40        | 46.66    | 46.71 | 92.33                  | 31.40     |
|              | 50%             | 91.81         | 31.80     | 73.71        | 46.15    | 46.31 | 91.79                  | 32.40     |
|              | 75%             | 91.67         | 31.80     | 73.89        | 46.01    | 45.95 | 91.67                  | 31.80     |
| Qwen3-8B     | Dense           | 91.90         | 33.60     | 95.87        | 75.00    | 61.21 | 91.90                  | 33.60     |
|              | 50%             | 92.08         | 35.10     | 96.23        | 76.50    | 61.56 | 92.07                  | 33.30     |
|              | 75%             | 92.11         | 34.40     | 96.07        | 75.33    | 61.51 | 91.74                  | 33.10     |

Performance of BLASST at different sparsity levels across all models and benchmarks. We evaluate on Llama-3.1-8B and Qwen3-8B across three deployment scenarios: prefill-only optimization (long-context tasks: RULER, LongBench); decode-only optimization (reasoning tasks: MATH500, AIME 2024, GPQA); and combined prefill+decode optimization. Results show minimal accuracy degradation even at ~75% sparsity, with occasional improvements over the dense baseline.

# Prefill-Only Accuracy

| Method          | RULER        |              |              |              |              |              | LongBench   |             |             |             |             |             |
|-----------------|--------------|--------------|--------------|--------------|--------------|--------------|-------------|-------------|-------------|-------------|-------------|-------------|
|                 | 4K           | 8K           | 16K          | 32K          | 64K          | Average      | Easy        | Hard        | Short       | Medium      | Long        | Overall     |
| Dense Attention | 96.16        | 95.07        | 94.80        | 92.33        | 87.69        | 93.21        | 29.7        | 32.5        | 38.3        | 28.8        | 25.0        | 31.4        |
| FlexPrefill     | 95.99        | 93.67        | 92.73        | 88.14        | 81.14        | 87.72        | 28.8        | 23.8        | 24.4        | 26.5        | 26.2        | 25.7        |
| MInference      | <b>96.54</b> | 94.06        | 91.37        | 85.79        | 83.03        | 84.15        | 28.6        | <b>32.8</b> | 36.7        | 30.2        | 24.1        | 31.2        |
| XAttention      | 96.37        | 94.47        | 94.48        | <b>91.91</b> | 85.01        | 92.44        | 29.2        | 31.5        | <b>38.3</b> | 26.0        | <b>26.9</b> | 30.6        |
| BLASST (~50%)   | 96.17        | <b>94.70</b> | <b>94.61</b> | 91.81        | <b>87.06</b> | <b>92.87</b> | <b>30.7</b> | 32.5        | <b>38.3</b> | <b>29.8</b> | 25.0        | <b>31.8</b> |

*Prefill phase comparison on Llama-3.1-8B-Instruct across RULER and LongBench. Best (non-dense) score in each column is denoted in bold. Targeting 50% sparsity, BLASST achieves the best accuracy among all sparse attention methods, closely matching dense attention in addition to being the easiest to use.*

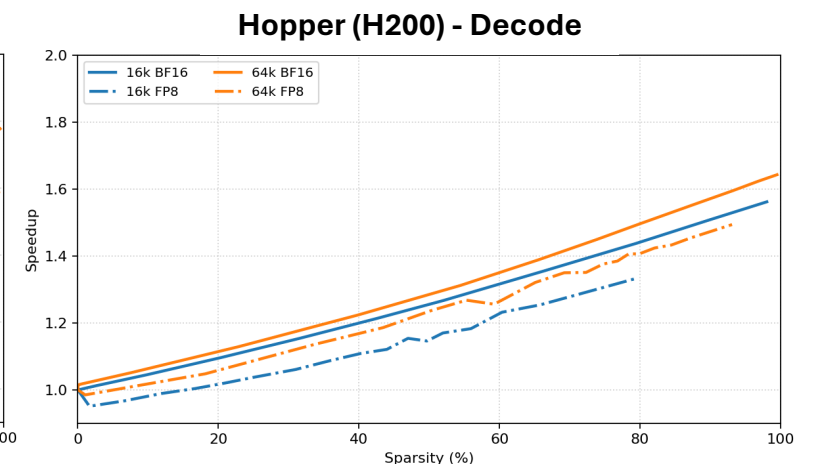
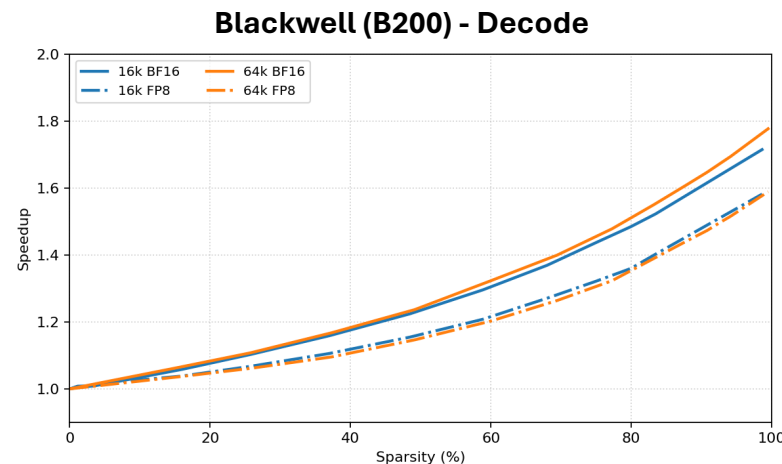
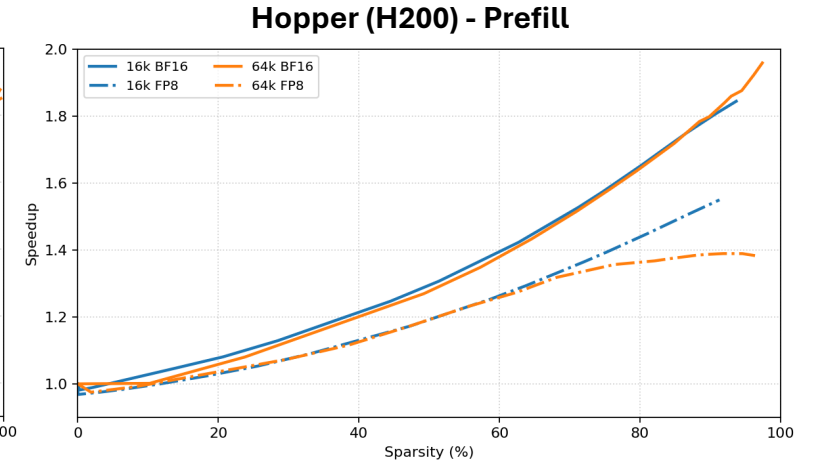
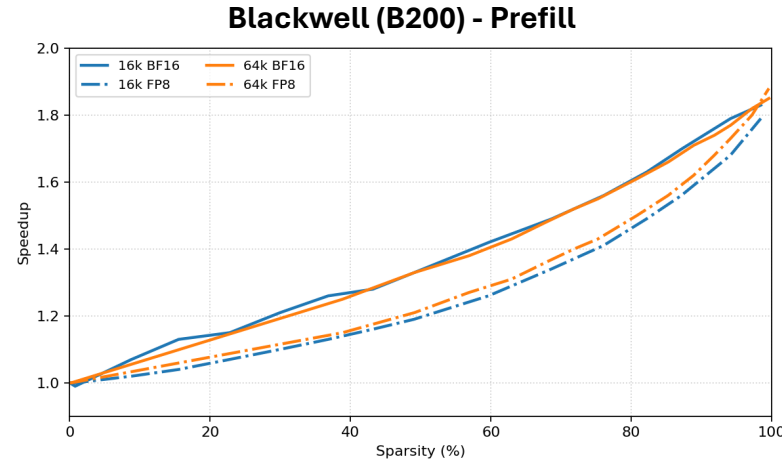
# Decode-Only Accuracy

| Method          | RULER-32K    | LongBench    | MATH500      | AIME 2024    | LiveCodeBench | GPQA         | Average      |
|-----------------|--------------|--------------|--------------|--------------|---------------|--------------|--------------|
| Dense Attention | 91.90        | 33.60        | 95.87        | 75.00        | 53.83         | 61.21        | 68.57        |
| Quest           | 56.23        | 30.30        | 94.18        | 71.50        | 52.17         | 60.12        | 60.75        |
| RocketKV        | 87.89        | 30.60        | 95.88        | 73.54        | 53.10         | 60.50        | 66.91        |
| BLASST ~50%     | <b>91.55</b> | <b>33.90</b> | <b>96.23</b> | <b>76.50</b> | <b>54.15</b>  | <b>61.51</b> | <b>68.97</b> |

*Decode phase comparison on Qwen3-8B across diverse reasoning and generation tasks. Best (non-dense) score in each column is denoted in bold. Targeting 50% sparsity, BLASST matches or exceeds dense baseline on all benchmarks, including mathematical reasoning (MATH500, AIME 2024), graduate-level science (GPQA), and code generation (LiveCodeBench), while maintaining long-context performance (RULER, LongBench).*

# Kernel Performance

- Compared against SoL FlashAttention kernels in TensorRT LLM.
- We see **~1.5x BF16 speedup at 75% sparsity**, where accuracy remains good.
- **Minimal to no overhead at 0% sparsity** – skip computation costs are hidden.



# Enable With a Few Lines of Code

# Enable With a Few Lines of Code

## 1. ModelOpt finds threshold for your model

```
from modelopt.torch.sparsity.attention_sparsity.config \
import SKIP_SOFTMAX_CALIB

# Load your model
model = AutoModelForCausalLM.from_pretrained(
    "Qwen/Qwen3-30B-A3B-Instruct-2507",
    attention_implementation="eager",
    torch_dtype=torch.bfloat16
)

# Apply sparse attention
model = mtsa.sparsify(model, config=SKIP_SOFTMAX_CALIB)
```

# Enable With a Few Lines of Code

## 1. ModelOpt finds threshold for your model

```
from modelopt.torch.sparsity.attention_sparsity.config \
import SKIP_SOFTMAX_CALIB

# Load your model
model = AutoModelForCausalLM.from_pretrained(
    "Qwen/Qwen3-30B-A3B-Instruct-2507",
    attention_implementation="eager",
    torch_dtype=torch.bfloat16
)

# Apply sparse attention
model = mtsa.sparsify(model, config=SKIP_SOFTMAX_CALIB)
```

## 2a. Deploy on TensorRT LLM server

```
from tensorrt_llm import LLM
from tensorrt_llm.llmapi import SkipSoftmaxAttentionConfig

sparse_attn_config = SkipSoftmaxAttentionConfig(
    threshold_scale_factor=1000.0
)

# Or separate threshold for prefill/decode
sparse_attn_config = SkipSoftmaxAttentionConfig(
    threshold_scale_factor={
        "prefill": 1000.0, "decode": 500.0
    }
)

llm = LLM(
    model="Qwen/Qwen3-30B-A3B-Instruct-2507",
    sparse_attention_config=sparse_attn_config,
    # Other arguments...
)
```

# Enable With a Few Lines of Code

## 1. ModelOpt finds threshold for your model

```
from modelopt.torch.sparsity.attention_sparsity.config \
import SKIP_SOFTMAX_CALIB

# Load your model
model = AutoModelForCausalLM.from_pretrained(
    "Qwen/Qwen3-30B-A3B-Instruct-2507",
    attention_implementation="eager",
    torch_dtype=torch.bfloat16
)

# Apply sparse attention
model = mtsa.sparsify(model, config=SKIP_SOFTMAX_CALIB)
```

## 2b. Deploy on SGLang server

```
SGLANG_SKIP_SOFTMAX_PREFILL_THRESHOLD_SCALE_FACTOR=1000.0 \
SGLANG_SKIP_SOFTMAX_DECODE_THRESHOLD_SCALE_FACTOR=500.0 \
python3 -m sglang.launch_server \
    --model-path "Qwen/Qwen3-30B-A3B-Instruct-2507" \
    # Other arguments...
```

## 2a. Deploy on TensorRT LLM server

```
from tensorrt_llm import LLM
from tensorrt_llm.llmapi import SkipSoftmaxAttentionConfig

sparse_attn_config = SkipSoftmaxAttentionConfig(
    threshold_scale_factor=1000.0
)

# Or separate threshold for prefill/decode
sparse_attn_config = SkipSoftmaxAttentionConfig(
    threshold_scale_factor={
        "prefill": 1000.0, "decode": 500.0
    }
)

llm = LLM(
    model="Qwen/Qwen3-30B-A3B-Instruct-2507",
    sparse_attention_config=sparse_attn_config,
    # Other arguments...
)
```

# Framework Support

Kernels available in TensorRT LLM and FlashInfer.



[Link](#)



[Link](#)

Threshold calibration tool and recipes available in Model Optimizer.



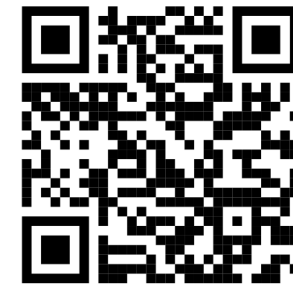
[Link](#)

# Framework Support



[Link](#)

Coming soon



[Link](#)

# References

Yang, S., Guo, J., Tang, H., Hu, Q., Xiao, G., Tang, J., Lin, Y., Liu, Z., Lu, Y. and Han, S., 2025. **Lserve: Efficient long-sequence llm serving with unified sparse attention**. In *Proceedings of Machine Learning and Systems*, 7.

Zaheer, M., Guruganesh, G., Dubey, K. A., Ainslie, J., Alberti, C., Ontanon, S., Pham, P., Ravula, A., Wang, Q., Yang, L., et al. **Big Bird: Transformers for longer sequences**. In *Advances in Neural Information Processing Systems*, volume 33, pp. 17283–17297, December 2020. doi: 10.48550/arXiv.2007.14062.

Xiao, G., Tian, Y., Chen, B., Han, S., and Lewis, M. **Efficient streaming language models with attention sinks**. In *The Twelfth International Conference on Learning Representations*, May 2024b. doi: 10.48550/arXiv.2309.17453

Jiang, H., Li, Y., Zhang, C., Wu, Q., Luo, X., Ahn, S., Han, Z., Abdi, A. H., Li, D., Lin, C.-Y., Yang, Y., and Qiu, L. **MIInference 1.0: Accelerating prefilling for long-context LLMs via dynamic sparse attention**. In *Advances in Neural Information Processing Systems*, volume 37, pp. 52481–52515, December 2024. doi: 10.48550/arXiv.2407.02490. URL <https://openreview.net/forum?id=fPBACAbqSN>.

Zadouri, T., Hoehnerbach, M., Shah, J., Liu, T., Thakkar, V., and Dao, T. **FlashAttention-4: Algorithm and kernel pipelining co-design for asymmetric hardware scaling**. *CoRR*, abs/2603.05451, March 2026. doi: 10.48550/arXiv.2603.05451.

# BLASST Is Robust and Easy to Use

| Method          | Accelerates Prefill | Accelerates Decode | No Training | No Pre-Computation |
|-----------------|---------------------|--------------------|-------------|--------------------|
| H2O             | ✗                   | ✓                  | ✓           | ✓                  |
| SnapKV          | ✗                   | ✓                  | ✓           | ✓                  |
| RocketKV        | ✗                   | ✓                  | ✓           | ✗                  |
| Quest           | ✗                   | ✓                  | ✓           | ✗                  |
| DuoAttention    | ✓                   | ✓                  | ✗           | ✓                  |
| DSA             | ✓                   | ✓                  | ✗           | ✓                  |
| MInference      | ✓                   | ✗                  | ✓           | ✗                  |
| SpargeAttention | ✓                   | ✗                  | ✓           | ✗                  |
| XAttention      | ✓                   | ✗                  | ✓           | ✗                  |
| <b>BLASST</b>   | ✓                   | ✓                  | ✓           | ✓                  |