

Native Sparse Attention: Hardware-Aligned and Natively Trainable Sparse Attention

Amit Kumar

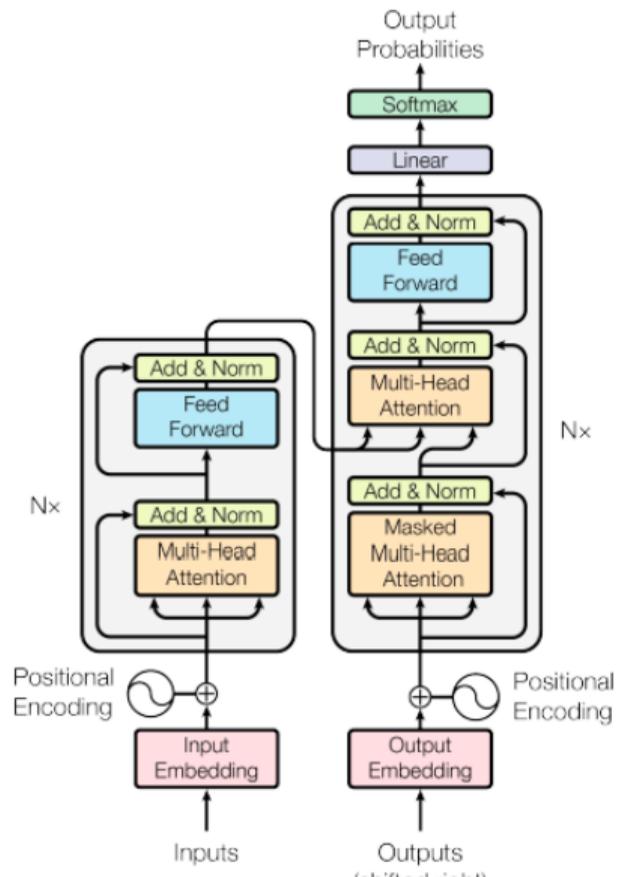
AI/NLP Engineer at E42.ai

Table of contents

1. Introduction
2. Native Sparse Attention (NSA)
3. Token Compression Attention
4. Token selection attention
5. Sliding window attention
6. Results
7. Conclusion

- Long-context modeling is critical for next-gen language models.
- Standard attention has quadratic complexity: $O(N^2)$ for sequence length N .
- This leads to huge computational and memory costs for long sequences (e.g., 64K tokens).
- Existing attention methods:
 1. Often hardware-inefficient
 2. Lack end-to-end trainability
 3. Use fixed heuristics limiting adaptability

Attention Definition and Intuition



$$O = f(QK^T)V$$

$$Q \in \mathbb{R}^{N,d}$$

$$K \in \mathbb{R}^{M,d}$$

$$V \in \mathbb{R}^{M,D}$$

$$O \in \mathbb{R}^{N,D}$$

Attention formulation

$$O = f(QK^T)V$$

$$Q \in \mathbb{R}^{B,H,N,d}$$

$$K \in \mathbb{R}^{B,H,M,d}$$

$$V \in \mathbb{R}^{B,H,M,D}$$

$$O \in \mathbb{R}^{B,H,N,D}$$

Attention formulation with batch and head dims

Attention Equation Visualization

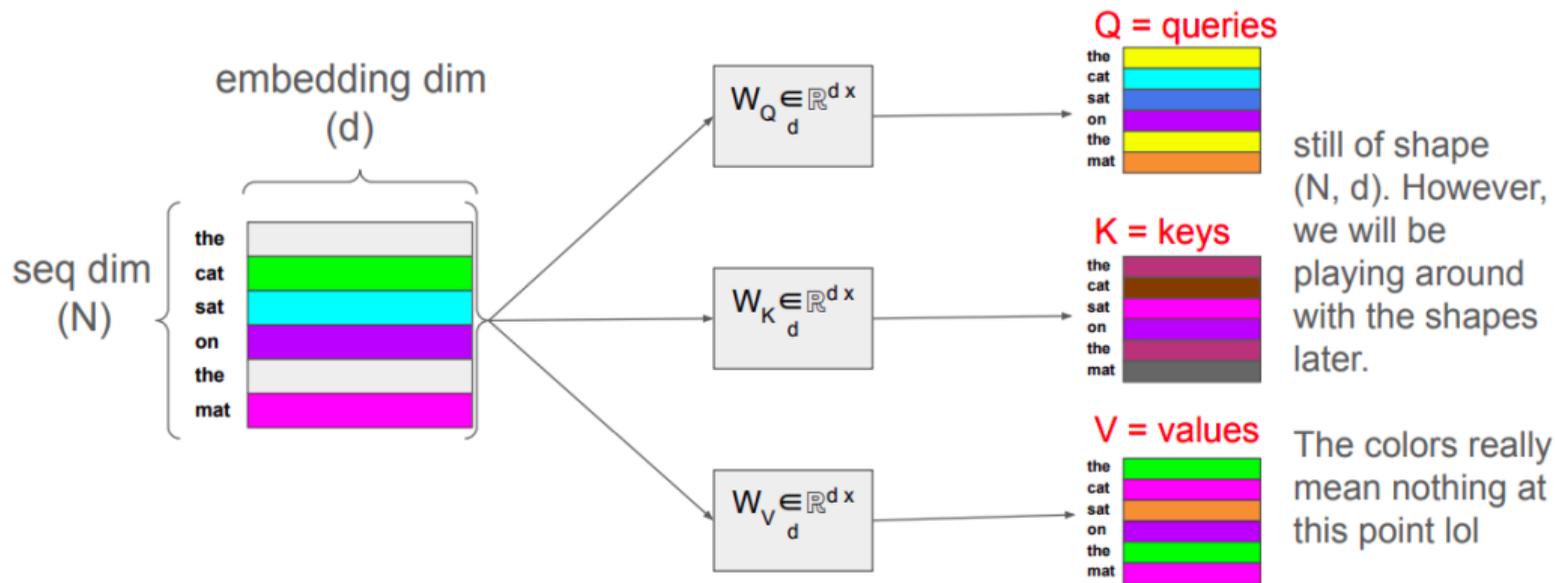
$$O = \text{Softmax}(QK^T)V$$

(N, d) (d, M) (M, d)

Literally just a lucky choice of function that is fantastic

$$O_t = \frac{\sum_{s=1}^M e^{\overbrace{Q_t K_s^T}^{(d)}} V_s}{\sum_{s=1}^M \underbrace{e^{\overbrace{Q_t K_s^T}^{(d)}}}_{(1) \text{ scalar}}}$$

$$\begin{aligned} Q &\in \mathbb{R}^{N,d} \\ K &\in \mathbb{R}^{M,d} \\ V &\in \mathbb{R}^{M,D} \\ O &\in \mathbb{R}^{N,D} \end{aligned}$$



Q = queries

sat 

K = keys

the 
cat 
sat 
on 
the 
mat 

V = values

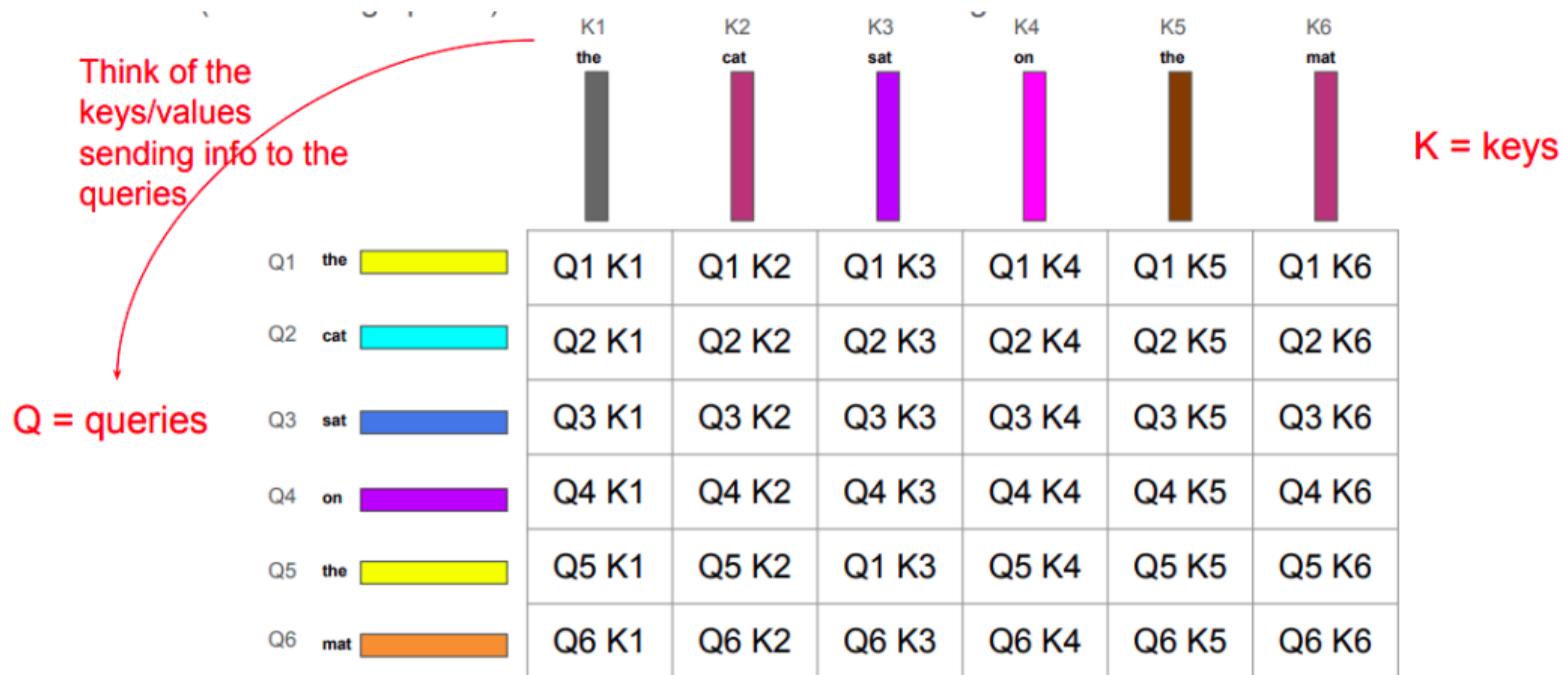
the 
cat 
sat 
on 
the 
mat 

$$O_{\text{sat}} = \frac{e^{Q_{\text{sat}}K_{\text{the}}^T}V_{\text{the}} + e^{Q_{\text{sat}}K_{\text{cat}}^T}V_{\text{cat}} + e^{Q_{\text{sat}}K_{\text{sat}}^T}V_{\text{sat}} + e^{Q_{\text{sat}}K_{\text{on}}^T}V_{\text{on}} + e^{Q_{\text{sat}}K_{\text{the}}^T}V_{\text{the}} + e^{Q_{\text{sat}}K_{\text{mat}}^T}V_{\text{mat}}}{e^{Q_{\text{sat}}K_{\text{the}}^T} + e^{Q_{\text{sat}}K_{\text{cat}}^T} + e^{Q_{\text{sat}}K_{\text{sat}}^T} + e^{Q_{\text{sat}}K_{\text{on}}^T} + e^{Q_{\text{sat}}K_{\text{the}}^T} + e^{Q_{\text{sat}}K_{\text{mat}}^T}}$$

Notice how the query determines the word we are modeling in the output

We do this for all words in the sequence!

Attention matrix



What exactly does this represent?

		K1 / V1	K2 / V2	K3 / V3	K4 / V4	K5 / V5	K6 / V6	
		the	cat	sat	on	the	mat	
Q1	the	1.0	0.0	0.0	0.0	0.0	0.0	O1
Q2	cat	0.0	0.8	0.1	0.0	0.0	0.1	O2
Q3	sat	0.0	0.3	0.6	0.0	0.0	0.1	O3
Q4	on	0.0	0.0	0.3	0.7	0.0	0.0	O4
Q5	the	0.0	0.0	0.0	0.0	1.0	0.0	O5
Q6	mat	0.0	0.0	0.0	0.1	0.0	0.9	O6

$$O_3 = e^{Q_3 K_1} V_1 + e^{Q_3 K_2} V_2 + e^{Q_3 K_3} V_3 + e^{Q_3 K_4} V_4 + e^{Q_3 K_5} V_5 + e^{Q_3 K_6} V_6$$

$$O_3 = 0.0 * V_1 + 0.3 * V_2 + 0.6 * V_3 + 0.0 * V_4 + 0.0 * V_5 + 0.1 * V_6$$

$$O_{\text{sat}} = 0.0 * V_{\text{the}} + 0.3 * V_{\text{cat}} + 0.6 * V_{\text{sat}} + 0.0 * V_{\text{on}} + 0.0 * V_{\text{the}} + 0.1 * V_{\text{mat}}$$

Causal attention

K = keys

Q = queries

	K1	K2	K3	K4	K5	K6	
Q1	Q1 K1	$-\infty$	$-\infty$	$-\infty$	$-\infty$	$-\infty$	→ output 1
Q2	Q2 K1	Q2 K2	$-\infty$	$-\infty$	$-\infty$	$-\infty$	→ output 2
Q3	Q3 K1	Q3 K2	Q3 K3	$-\infty$	$-\infty$	$-\infty$	→ output 3
Q4	Q4 K1	Q4 K2	Q4 K3	Q4 K4	$-\infty$	$-\infty$	→ output 4
Q5	Q5 K1	Q5 K2	Q5 K3	Q5 K4	Q5 K5	$-\infty$	→ output 5
Q6	Q6 K1	Q6 K2	Q6 K3	Q6 K4	Q6 K5	Q6 K6	→ output 6

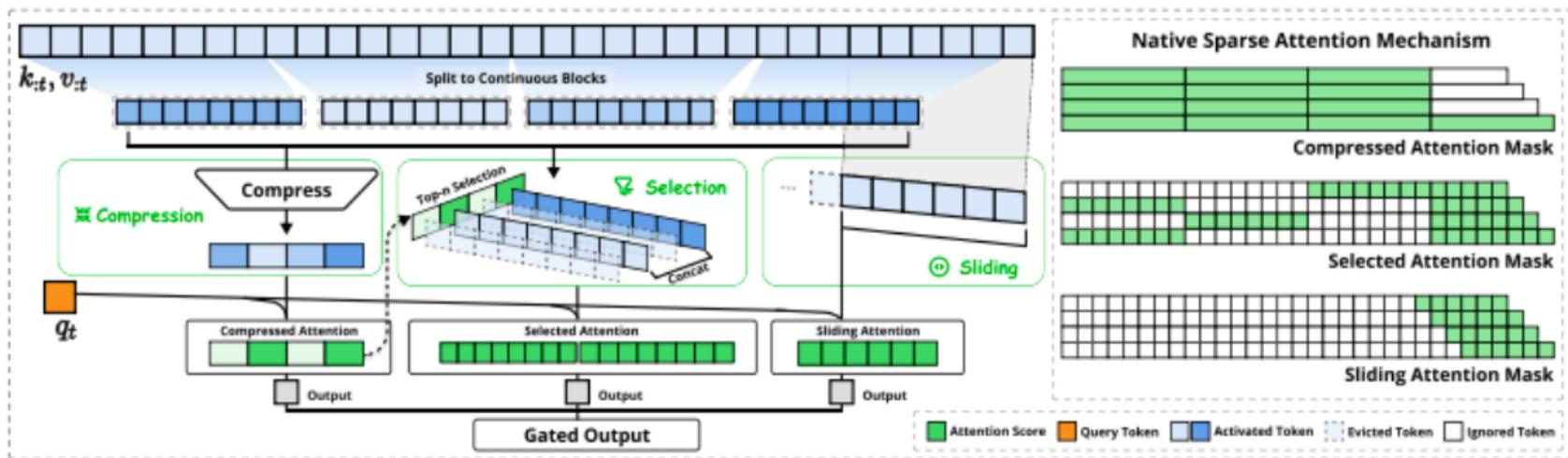
Causal Padding Mask



What is Native Sparse Attention (NSA)?

- A natively trainable sparse attention mechanism designed for:
 - Hardware alignment (efficient on GPUs/TPUs)
 - End-to-end differentiability (learns sparse patterns during training)
 - Efficient long-context modeling
- Combines algorithmic innovations with hardware-aware optimizations
- Input sequence split into blocks (e.g., 32 tokens/block)
- Three complementary pathways:
 - Compression: Learnable MLP compresses blocks into summary vectors
 - Selection: Differentiable top-k selects important blocks globally
 - Sliding Window: Local attention on recent tokens (e.g., 512 tokens)
- Balances global context awareness and local precision

NSA architecture



Compression Pathway in Native Sparse Attention

- The input sequence of tokens is partitioned into **blocks** of size B .
- Each block of tokens is passed through a **learnable MLP** (Multi-Layer Perceptron) with positional encoding.
- The MLP compresses the block into a **compact summary vector** representing the block.
- This reduces the global attention complexity from $O(N^2)$ to approximately $O(N \times \frac{N}{B})$, which is near-linear.

Simple example:

- Suppose sequence length $N = 16$, block size $B = 4$.
- The sequence is divided into $M = \frac{N}{B} = 4$ blocks.
- Each block of 4 tokens is compressed into 1 summary vector.
- Global attention attends to these 4 summary vectors instead of all 16 tokens.

Block-Level Compressed Key Representation (Simplified)

Idea: Instead of using every single key or value, we group them into blocks. Each block is turned into a single, compressed key (or value) that summarizes the whole block. This makes processing faster and helps capture the main information.

- Take blocks of length l from the sequence of keys.
- Move forward by d steps to get the next block (stride).
- Use a function φ (MLP) to turn each block into one compressed key.

Equation:

$$\tilde{\mathbf{K}}_t^{\text{cmp}} = f_K^{\text{cmp}}(\mathbf{k}_t) = \left\{ \varphi(\mathbf{k}_{id+1:id+l}) \mid 0 \leq i \leq \left\lfloor \frac{t-l}{d} \right\rfloor \right\}$$

Notes:

- l : block size
- d : step size between blocks
- φ : function that compresses a block into a single key

Leveraging Attention Scores for Block Importance

- Computing block importance scores can introduce significant overhead.
- Attention computation of compression tokens produces intermediate attention scores.
- These scores can be leveraged to induce selection block importance scores:

$$p_{cmp}^t = \text{Softmax} \left(q_t^T \tilde{K}_{cmp}^t \right),$$

where $p_{cmp}^t \in \mathbb{R}^{\lfloor \frac{t-l}{d} \rfloor + 1}$ represents attention scores between q_t and compression keys \tilde{K}_{cmp}^t .

- Retain tokens within top- n sparse blocks ranked by importance scores.

$$I_t = \{i \mid \text{rank}(p_{slc}^{t'}[i]) \leq n\}$$

$$\tilde{K}_{slc}^t = \text{Cat}\{k_{il'+1:(i+1)l'} \mid i \in I_t\}$$

- $\text{rank}(\cdot)$ ranks scores in descending order (rank 1 = highest).
- $\tilde{K}_{slc}^t \in \mathbb{R}^{d_k \times nl'}$ is the tensor of selected compression keys.
- Selected keys and values then participate in attention computation with q_t .

Sliding Window attention

- Local patterns adapt faster and can dominate learning, hindering compression and selection token learning.
- Introduce a *sliding window* dedicated to local context, allowing other branches to focus on their features without shortcutting by local patterns.
- **Mechanism:**
 - Maintain recent tokens in window w :

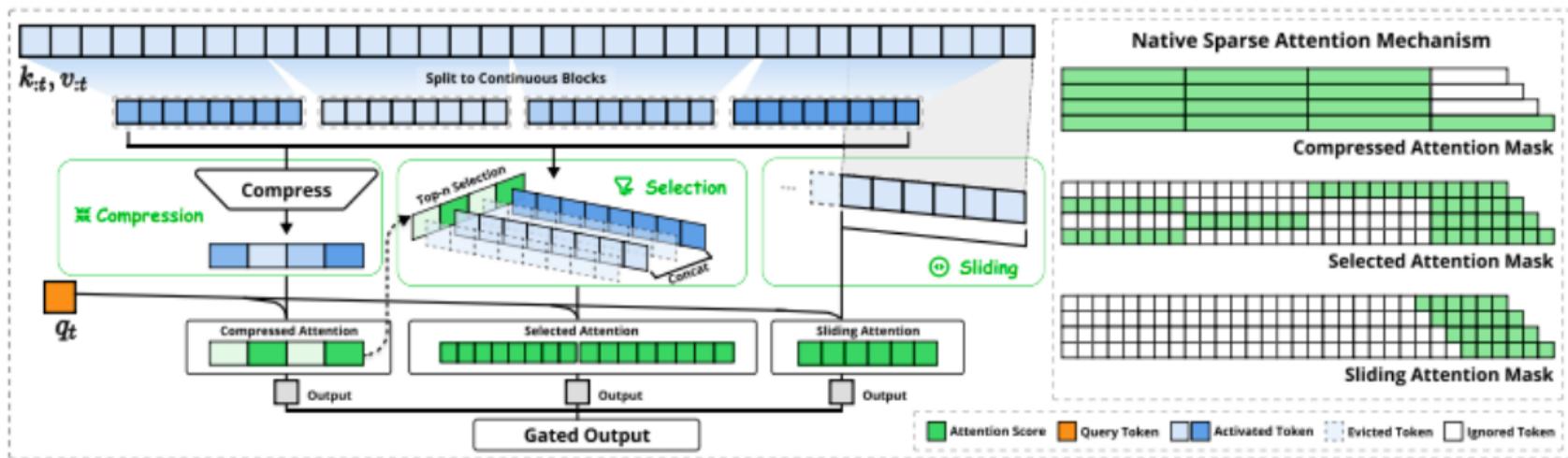
$$\tilde{K}_t^{win} = k_{t-w:t}, \quad \tilde{V}_t^{win} = v_{t-w:t}$$

- Separate attention computations into three branches:

Compression: $(\tilde{K}_t^{cmp}, \tilde{V}_t^{cmp})$, Selection: $(\tilde{K}_t^{slc}, \tilde{V}_t^{slc})$, Sliding Window: $(\tilde{K}_t^{win}, \tilde{V}_t^{win})$

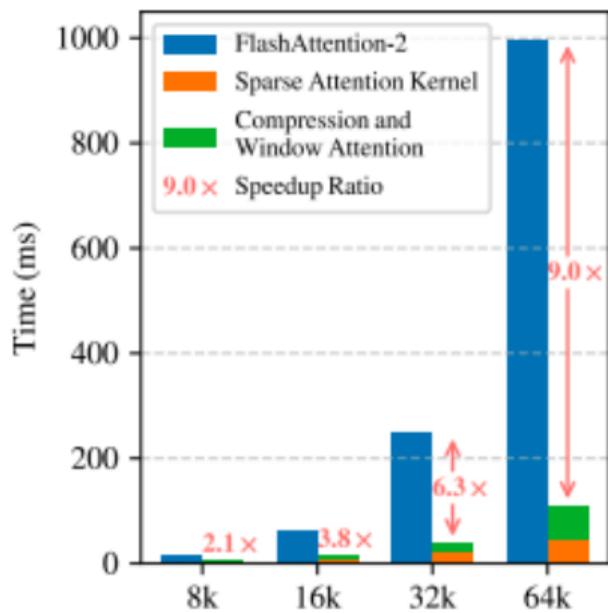
- Use independent keys and values per branch to prevent shortcut learning and gradient interference.
- Aggregate branch outputs via a learned gating mechanism.

NSA architecture

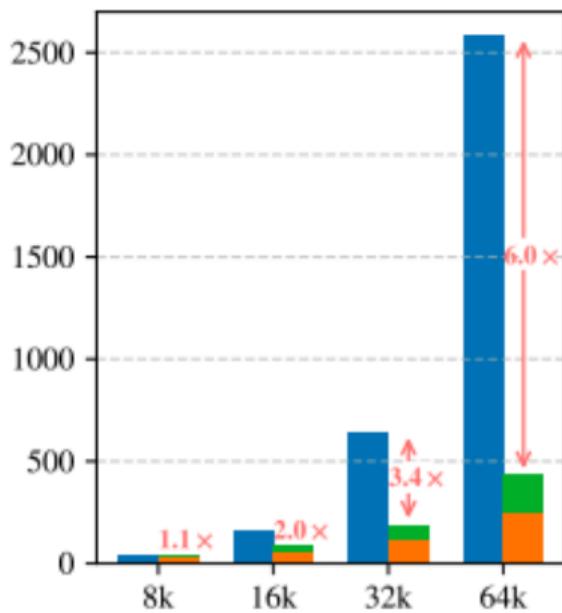


Comparison of NSA kernel with FlashAttention-2

Forward Time Comparison



Backward Time Comparison



Context Length	8192	16384	32768	65536
Full Attention	8192	16384	32768	65536
NSA	2048	2560	3584	5632
Expected Speedup	4×	6.4×	9.1×	11.6×

NSA achieving lower loss values

Model	SQA			MQA				Synthetic		Code	Avg.
	MFQA-en	MFQA-zh	Qasper	HPQ	2Wiki	GovRpt	Dur	PassR-en	PassR-zh	LCC	
H2O	0.428	0.429	0.308	0.112	0.101	0.231	0.208	0.704	0.421	0.092	0.303
InfLLM	0.474	0.517	0.356	0.306	0.250	0.277	0.257	0.766	0.486	0.143	0.383
Quest	0.495	0.561	0.365	0.295	0.245	0.293	0.257	0.792	0.478	0.135	0.392
Exact-Top	0.502	0.605	0.397	0.321	0.288	<u>0.316</u>	0.291	0.810	0.548	0.156	0.423
Full Attn	0.512	<u>0.623</u>	<u>0.409</u>	<u>0.350</u>	<u>0.305</u>	0.324	<u>0.294</u>	<u>0.830</u>	0.560	<u>0.163</u>	<u>0.437</u>
NSA	<u>0.503</u>	0.624	0.432	0.437	0.356	0.307	0.341	0.905	<u>0.550</u>	0.232	0.469

- **Efficient Long-Context Modeling:** NSA uses hierarchical token compression and blockwise token selection to handle very long input sequences efficiently.
- **Hardware-Optimized Design:** It is specially designed to work well with modern GPUs, reducing memory access and speeding up training and inference.
- **Maintains High Accuracy:** NSA matches or exceeds the performance of full attention models on benchmarks and long-context tasks.
- **Significant Speedup:** It achieves much faster computation (up to $9\times$ faster in some cases) while keeping the model's reasoning and accuracy intact.

Thank You