ECS 289D: Datacenter Systems for LLMs

FlashAttention: Fast and Memory-Efficient Exact Attention with IO-Awareness

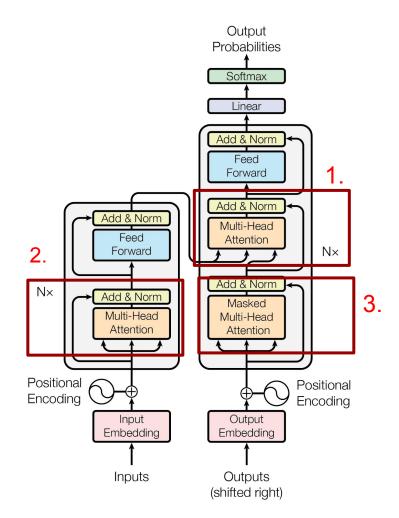
Authors: Tri Dao, Daniel Y. Fu, Stefano Ermon, Atri Rudra, and Christopher Ré

Presented by Ansha Prashanth

Introduction

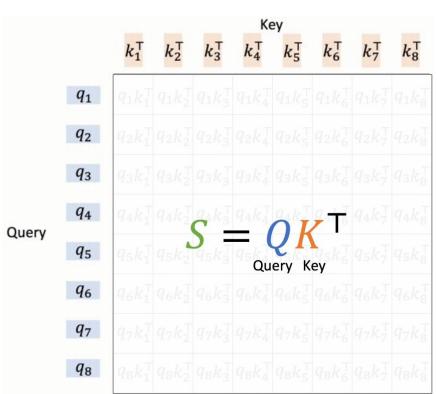
A transformer model contains many base attention layers:

- Cross-Attention layer: Decoder-encoder attention
- Global Self-attention layer: Encoder self-attention
- 3. Causal self-attention layer: Decoder self-attention



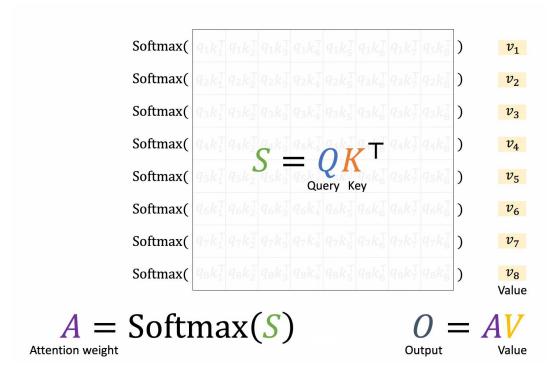
Self Attention

- Query (Q): Represents the current state or part of the sequence we are focusing on.
- Key (K): Represents all parts of the sequence we compare the query against.
- Value (V): Represents the actual information we want to use in the output.
- N is the sequence length and d is the head dimension
- Output (O)

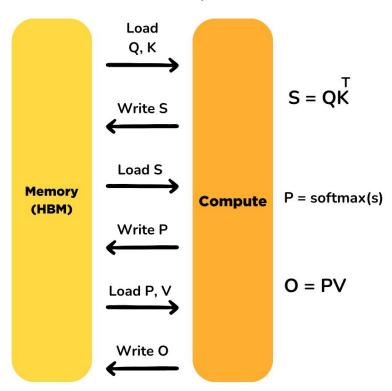


Self Attention

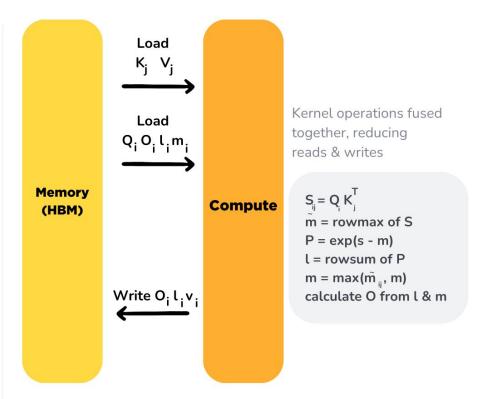
- Attention Score Matrix (S):
 Calculated via the dot product of Q and K
- Probability Matrix (P or A):
 Softmax is applied row-wise to
 S
- Final Output (O): Computed as the weighted sum of V
- The intermediate matrices, S and P, are N×N, requiring O(N²) memory, which is the root of the memory bottleneck.



Standard Attention Implementation



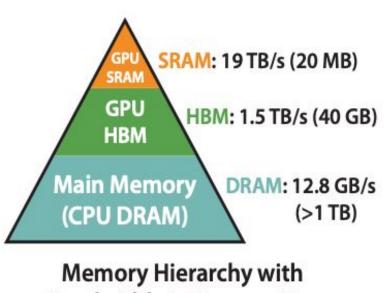
Flash Attention



Initialize O, I and m matrices with zeroes. m and I are used to calculate cumulative softmax. Divide Q, K, V into blocks (due to SRAM's memory limits) and iterate over them, for i is row & j is column.

GPU Memory Hierarchy

- **SRAM**: fast, on-chip, small
- **HBM (High Bandwidth Memory)**: slower than SRAM, large size. That's what we usually address as GPU memory.
- "On modern GPUs, compute speed has out-paced memory speed, and most operations in Transformers are bottlenecked by memory accesses". This is the main motivation for FlashAttention.



Bandwidth & Memory Size

FlashAttention

Tiling and recomputation are applied to overcome the technical challenge of computing exact attention with sub-quadratic HBM access

Tiling:

- Split Q, K, V into blocks and load them from HBM to SRAM
- Compute attention block by block without storing full S,P matrices
- Use online softmax statistics m(x), $\ell(x)$ to combine block results.

Recomputation:

- Only store the final output O and softmax statistics (m, ℓ)
- Recompute S, P on the fly during backpropagation using Q, K, V blocks

Tiling

$$C = A \times B$$

A

1	2	3	4
5	6	7	8
9	10	11	12
13	14	15	16

•	I	Т	٦
1	F	4	ď
J	L		

1	2	3	4
5	6	7	8
9	10	11	12
13	14	15	16

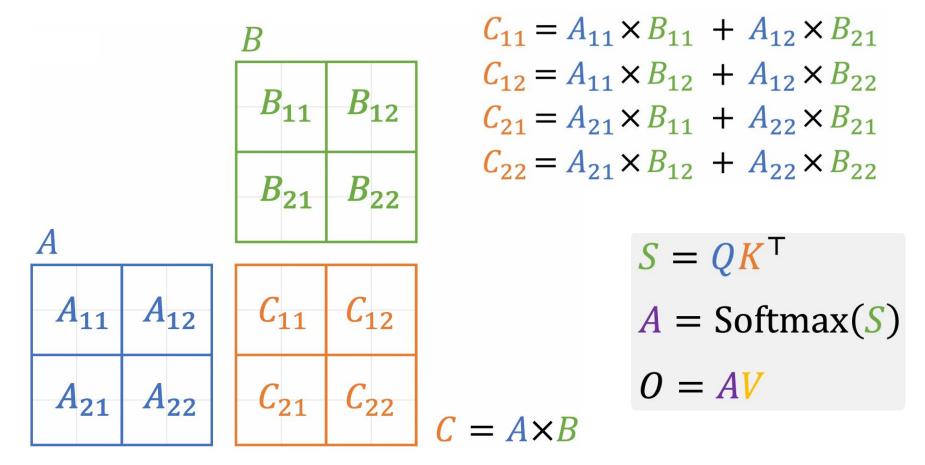
$$\begin{bmatrix} C_{1,1} & C_{1,2} \\ C_{2,1} & C_{2,2} \end{bmatrix} = \begin{bmatrix} 1 & 2 \\ 5 & 6 \end{bmatrix} \begin{bmatrix} 1 & 2 \\ 5 & 6 \end{bmatrix} + \begin{bmatrix} 3 & 4 \\ 7 & 8 \end{bmatrix} \begin{bmatrix} 9 & 10 \\ 13 & 14 \end{bmatrix}$$

Without tiling: 32 memory accesses

With tiling: 16 memory accesses

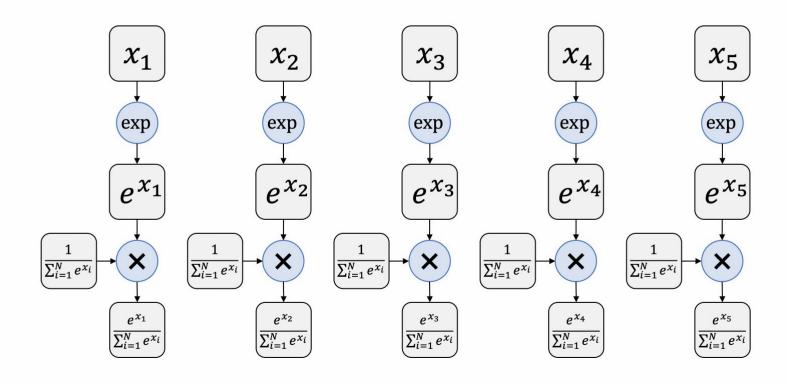
NxN block \rightarrow 1/N memory access

Tiling



$$A = Softmax(S)$$

$$S = \{x_1, x_2, \cdots, x_N\}$$



Safe Softmax in FlashAttention

- Standard softmax computation involves exponentiating values, which may cause numerical instability, especially in FP16 due to dynamic range.
 - Direct exponentiation of large numbers (eg. e¹² =162754) may exceed FP16 limits (65504) leading to numerical errors.
- Safe Softmax mitigates this by shifting all inputs by subtracting the maximum value:

$$m = \max(x_i), \quad \operatorname{softmax}(x_i) = rac{e^{x_i - m}}{\sum_j e^{x_j - m}}$$

$$A = Softmax(S)$$

 $m = \max(\{x_1, x_2, \cdots, x_N\})$

 $S = \{x_1, x_2, \cdots, x_N\}$

$$x_{1}$$

$$x_{2}$$

$$x_{3}$$

$$x_{4}$$

$$x_{5}$$

$$x_{5}$$

$$x_{6}$$

$$x_{1}$$

$$x_{2}$$

$$x_{3}$$

$$x_{4}$$

$$x_{5}$$

$$x_{5}$$

$$x_{1}$$

$$x_{2}$$

$$x_{2}$$

$$x_{3}$$

$$x_{4}$$

$$x_{5}$$

$$x_{5}$$

$$x_{1}$$

$$x_{2}$$

$$x_{3}$$

$$x_{4}$$

$$x_{5}$$

$$x_{1}$$

$$x_{2}$$

$$x_{2}$$

$$x_{3}$$

$$x_{4}$$

$$x_{5}$$

$$x_{5}$$

$$x_{6}$$

$$x_{3}$$

$$x_{4}$$

$$x_{2}$$

$$x_{5}$$

$$x_{6}$$

$$x_{3}$$

$$x_{4}$$

$$x_{5}$$

$$x_{5}$$

$$x_{6}$$

$$x_{3}$$

$$x_{4}$$

$$x_{5}$$

$$x_{6}$$

$$x_{7}$$

$$x_{8}$$

$$x_{7}$$

$$x_{8}$$

$$x_{1}$$

$$x_{1}$$

$$x_{2}$$

$$x_{3}$$

$$x_{4}$$

$$x_{5}$$

$$x_{6}$$

$$x_{7}$$

$$x_{1}$$

$$x_{2}$$

$$x_{3}$$

$$x_{4}$$

$$x_{5}$$

$$x_{6}$$

$$x_{7}$$

$$x_{8}$$

$$x_{1}$$

$$x_{1}$$

$$x_{2}$$

$$x_{3}$$

$$x_{4}$$

$$x_{5}$$

$$x_{6}$$

$$x_{7}$$

$$x_{1}$$

$$x_{1}$$

$$x_{2}$$

$$x_{3}$$

$$x_{4}$$

$$x_{5}$$

$$x_{6}$$

$$x_{7}$$

$$x_{1}$$

$$x_{1}$$

$$x_{2}$$

$$x_{3}$$

$$x_{4}$$

$$x_{4}$$

$$x_{5}$$

$$x_{5}$$

$$x_{2}$$

$$x_{3}$$

$$x_{4}$$

$$x_{5}$$

$$x_{4}$$

$$x_{5}$$

$$x_{7}$$

$$x_{8}$$

$$x_{1}$$

$$x_{1}$$

$$x_{2}$$

$$x_{3}$$

$$x_{4}$$

$$x_{5}$$

$$x_{7}$$

$$x_{8}$$

$$x_{1}$$

$$x_{1}$$

$$x_{2}$$

$$x_{3}$$

$$x_{4}$$

$$x_{4}$$

$$x_{5}$$

$$x_{5}$$

$$x_{2}$$

$$x_{3}$$

$$x_{4}$$

$$x_{5}$$

$$x_{4}$$

$$x_{5}$$

$$x_{7}$$

$$x_{8}$$

$$x_{1}$$

$$x_{2}$$

$$x_{3}$$

$$x_{4}$$

$$x_{5}$$

$$x_{4}$$

$$x_{5}$$

$$x_{7}$$

$$x_{8}$$

$$x_{7}$$

$$x_{8}$$

$$x_{8}$$

$$x_{9}$$

$$x_{1}$$

$$x_{1}$$

$$x_{2}$$

$$x_{2}$$

$$x_{3}$$

$$x_{4}$$

$$x_{4}$$

$$x_{5}$$

$$x_{4}$$

$$x_{5}$$

$$x_{7}$$

$$x_{8}$$

$$x_{8}$$

$$x_{8}$$

$$x_{9}$$

$$x_{1}$$

$$x_{1}$$

$$x_{2}$$

$$x_{3}$$

$$x_{4}$$

$$x_{4}$$

$$x_{4}$$

$$x_{5}$$

$$x_{7}$$

$$x_{8}$$

$$x_{8}$$

$$x_{8}$$

$$x_{9}$$

$$x_{1}$$

$$x_{1}$$

$$x_{2}$$

$$x_{3}$$

$$x_{4}$$

$$x_{4}$$

$$x_{4}$$

$$x_{2}$$

$$x_{3}$$

$$x_{4}$$

$$x_{4}$$

$$x_{4}$$

$$x_{5}$$

$$x_{7}$$

$$x_{8}$$

$$x_{8}$$

$$x_{1}$$

$$x_{2}$$

$$x_{3}$$

$$x_{4}$$

$$x_{4}$$

$$x_{4}$$

$$x_{4}$$

$$x_{4}$$

$$x_{4}$$

$$x_{4}$$

$$x_{5}$$

$$x_{7}$$

$$x_{8}$$

$$x_{8}$$

$$x_{9}$$

$$x_{1}$$

$$x_{1}$$

$$x_{2}$$

$$x_{3}$$

$$x_{4}$$

$$x_{4}$$

$$x_{4}$$

$$x_{4}$$

$$x_{2}$$

$$x_{3}$$

$$x_{4}$$

$$x_{5}$$

$$x_{7}$$

$$x_{8}$$

$$x_{8}$$

$$x_{9}$$

$$x_{1}$$

$$x_{1}$$

$$x_{2}$$

$$x_{3$$

$$A = \operatorname{Softmax}(S) \qquad S = \{x_1, x_2, \dots, x_N\}$$

$$m_0 = -\infty \qquad \qquad x_1 \qquad x_2 \qquad x_{N-1} \qquad x_N$$

$$m_0 = -\infty \qquad \qquad x_1 \qquad x_2 \qquad x_{N-1} \qquad x_N$$

$$m_1 = \max(m_{i-1}, x_i) \qquad \qquad m_1 \qquad m_2 \qquad m_N$$

$$d_0 = 0 \qquad \qquad m_1 \qquad m_1 \qquad m_2 \qquad m_N$$

$$d_0 = 0 \qquad \qquad d_0 = 0$$

for $1 \le i \le N$ do

 $a_i = e^{x_i - m_N} / d_N$

 $A = \{a_1, a_2, \dots, a_N\}$

 $a_i = e^{x_i - m_N}/d_N$

$$A = Softmax(S)$$

$$S = \{x_1, x_2, \cdots, x_N\}$$

$$m_0 = -\infty$$

for $1 \le i \le N$ do
 $m_i = \max(m_{i-1}, x_i)$

$$d_i = \sum_{j=1}^i e^{x_j - m_N}$$
 $d'_i = \sum_{j=1}^i e^{x_j - m_i}$ $d'_N = d_N$

$$d_0 = 0$$
for $1 \le i \le N$ do
$$d_i = d_{i-1} + e^{x_i - m_N}$$

$$d'_{i} = \left(\sum_{j=1}^{i-1} e^{x_{j} - m_{i-1}}\right) e^{m_{i-1} - m_{i}} + e^{x_{i} - m_{i}}$$

$$d'_{i-1}$$

for
$$1 \le i \le N$$
 do
$$a_i = e^{x_i - m_N} / d_N$$

$$A = Softmax(S)$$

$$A = \text{Softmax}(S) \qquad S = \{x_1, x_2, \dots, x_N\}$$

$$m_0 = -\infty$$

for $1 \le i \le N$ do

 $m_i = \max(m_{i-1}, x_i)$
 $d_0 = 0$

for $1 \le i \le N$ do

 $d_i = d_{i-1} + e^{x_i - m_N}$

for $1 \le i \le N$ do

for
$$1 \le i \le N$$
 do
$$a_i = e^{x_i - m_N} / d_N$$

$$d'_{i} = d'_{i-1} e^{m_{i-1} - m_{i}} + e^{x_{i} - m_{i}}$$

$$m_{0} = -\infty$$

$$d_{0} = 0$$

$$for \ 1 \le i \le N \ do$$

$$m_{i} = \max(m_{i-1}, x_{i})$$

$$d'_{i} = d'_{i-1} e^{m_{i-1} - m_{i}} + e^{x_{i} - m_{i}}$$

for
$$1 \le i \le N$$
 do
$$a_i = e^{x_i - m_N} / d'_N$$

Online Softmax

- FlashAttention uses the mathematical properties of the softmax function to decompose the normalization.
- This allows us to compute 'local' softmax over each block and then rescale them to get the correct output incrementally.

$$S = QK^{T} \qquad A = \text{Softmax}(S) \qquad O = AV \qquad S = \{x_{1}, x_{2}, \dots, x_{N}\}$$

$$x_{i} = qk_{i}^{T}$$

$$d'_{i} = d'_{i-1} e^{m_{i-1} - m_{i}} + e^{x_{i} - m_{i}}$$

$$m_{0} = -\infty$$

$$d_{0} = 0$$

$$for \ 1 \le i \le N \ do$$

$$x_{i} = qk_{i}^{T}$$

$$m_{i} = \max(m_{i-1}, x_{i})$$

$$d'_{i} = d'_{i-1} e^{m_{i-1} - m_{i}} + e^{x_{i} - m_{i}}$$

$$o_{0} = 0$$

$$for \ 1 \le i \le N \ do$$

$$a_{i} = e^{x_{i} - m_{N}} / d'_{N}$$

$$o_{i} = o_{i-1} + a_{i} v_{i}$$
return o_{N}

$$S = QK^{T} \quad A = \text{Softmax}(S) \quad O = AV \quad S = \{x_{1}, x_{2}, \dots, x_{N}\}$$

$$x_{i} = qk_{i}^{T}$$

$$m_{0} = -\infty$$

$$d_{0} = 0 \quad o'_{i} = \sum_{j=1}^{i-1} \frac{e^{x_{j} - m_{i}}}{d'_{i}} \frac{d'_{i-1}}{d'_{i-1}} \frac{e^{-m_{i-1}}}{e^{-m_{i-1}}} v_{j} + \frac{e^{x_{i} - m_{i}}}{d'_{i}} v_{i}$$

$$for \quad 1 \le i \le N \quad do$$

$$x_{i} = qk_{i}^{T}$$

$$m_{i} = \max(m_{i-1}, x_{i})$$

$$d'_{i} = d'_{i-1}e^{m_{i-1} - m_{i}} + e^{x_{i} - m_{i}}$$

$$o_{0} = 0$$

$$for \quad 1 \le i \le N \quad do$$

$$a_{i} = e^{x_{i} - m_{N}} / d'_{N}$$

$$o_{i} = o_{i-1} + a_{i} v_{i}$$

$$return \quad o_{N}$$

$$S = QK^{T} \quad A = \text{Softmax}(S) \quad O = AV \quad S = \{x_{1}, x_{2}, \dots, x_{N}\}$$

$$x_{i} = qk_{i}^{T}$$

$$m_{0} = -\infty$$

$$d_{0} = 0 \quad o'_{i} = \left[\sum_{j=1}^{i-1} \frac{e^{x_{j} - m_{i-1}}}{d'_{i-1}} v_{j}\right] \frac{d'_{i-1}}{d'_{i}} e^{m_{i-1} - m_{i}} + \frac{e^{x_{i} - m_{i}}}{d'_{i}} v_{i}$$

$$for \quad 1 \le i \le N \quad do$$

$$x_{i} = qk_{i}^{T} \quad o'_{i-1}$$

$$m_{i} = \max(m_{i-1}, x_{i})$$

$$d'_{i} = d'_{i-1}e^{m_{i-1} - m_{i}} + e^{x_{i} - m_{i}}$$

$$o_{0} = 0$$

$$for \quad 1 \le i \le N \quad do$$

$$a_{i} = e^{x_{i} - m_{N}} / d'_{N}$$

$$o_{i} = o_{i-1} + a_{i} v_{i}$$

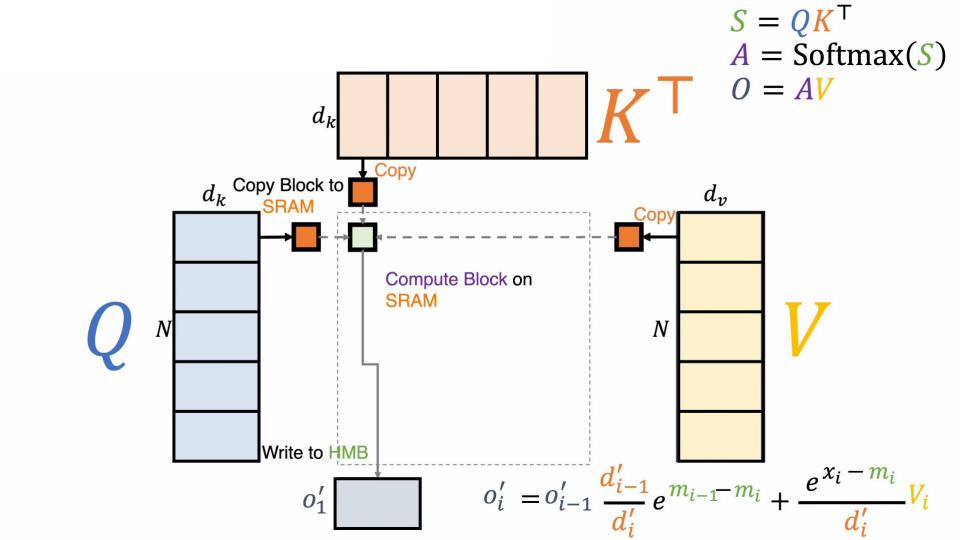
$$return \quad o_{N}$$

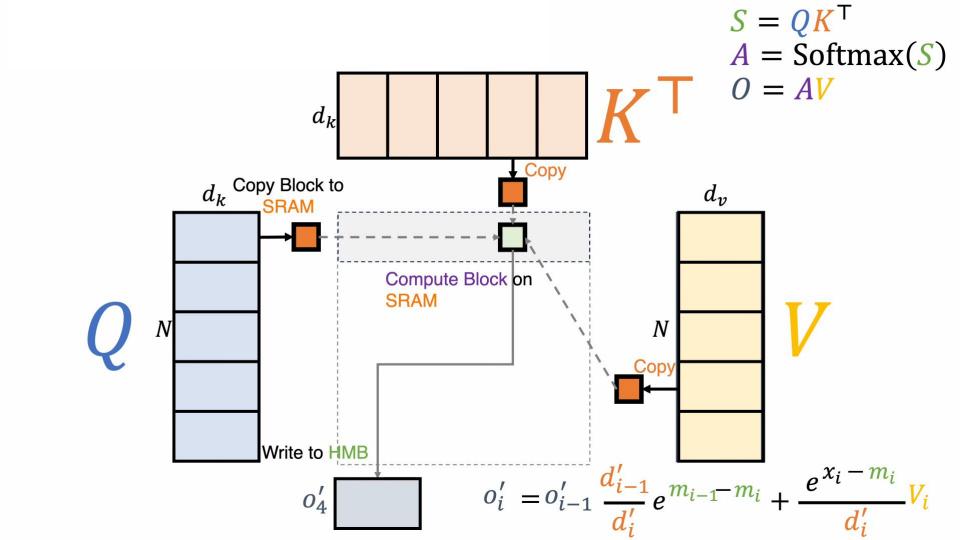
$$S = QK^{T}$$
 $A = \text{Softmax}(S)$ $O = AV$ $S = \{x_1, x_2, \dots, x_N\}$ $x_i = qk_i^{T}$ for $1 \le i \le N$ do $x_i = qk_i^{T}$

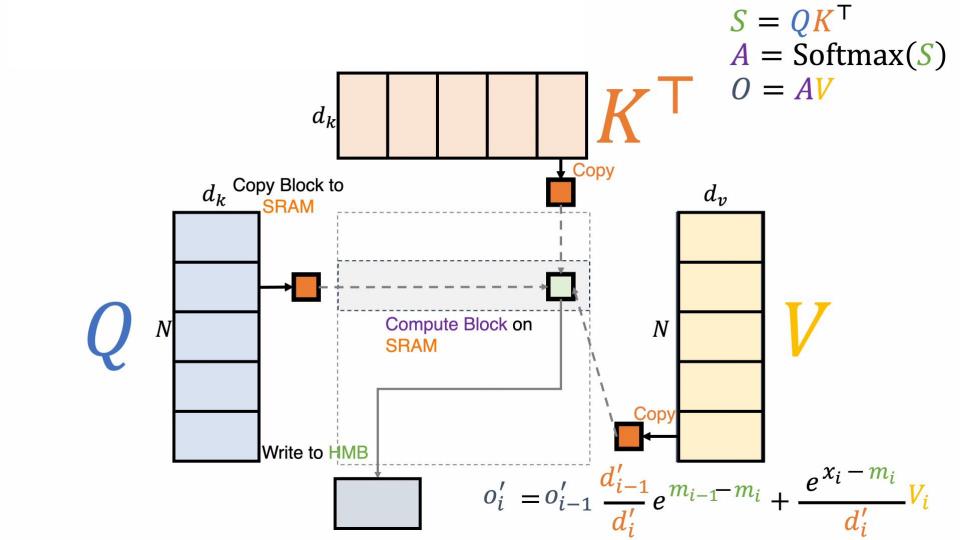
$$m_{i} = \max(m_{i-1}, x_{i})$$

$$d'_{i} = d'_{i-1}e^{m_{i-1}-m_{i}} + e^{x_{i}-m_{i}}$$

$$o'_{i} = o'_{i-1}\frac{d'_{i-1}}{d'_{i}}e^{m_{i-1}-m_{i}} + \frac{e^{x_{i}-m_{i}}}{d'_{i}}v_{i}$$
return o'_{N}







Backward Pass and Recomputation

Standard Attention:

 The backward pass requires the original N×N attention score matrix (S) and probability matrix (P) to compute the gradients dQ, dK, and dV using the chain rule.

FlashAttention Solution: Recomputation

- This strategy is a form of selective gradient checkpointing.
- We only store the output O and the compact, linear-sized softmax normalization statistics (m and l) from the forward pass.

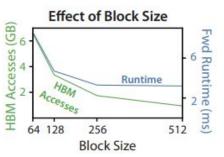
How it works:

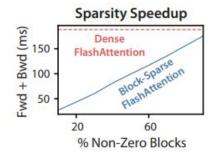
- We load blocks of Q, K, V, dO, O, m, and ℓ into SRAM
- We recompute the attention matrix block (S_{ij}) and probability matrix block (P_{ij})
 on-chip using the stored m and ℓ and the original input Q and K
- We then use this recomputed P_{ij} along with the output gradient (dO) and the original O to compute the input gradients dQ, dK, and dV locally.

How impactful is FlashAttention?

Attention	Standard	FLASHATTENTION
GFLOPs	66.6	75.2
HBM R/W (GB)	40.3	4.4
Runtime (ms)	41.7	7.3

Forward + backward runtime for GPT-2 on A100 GPU





BERT Implementation	Training time (minutes)
Nvidia MLPerf 1.1 [58]	20.0 ± 1.5
FLASHATTENTION (ours)	17.4 ± 1.4

Models	ListOps	Text	Retrieval	Image	Pathfinder	Avg	Speedup
Transformer	36.0	63.6	81.6	42.3	72.7	59.3	-
FLASHATTENTION	37.6	63.9	81.4	43.5	72.7	59.8	2.4×
Block-sparse FlashAttention	37.0	63.0	81.3	43.6	73.3	59.6	2.8×

Model implementations	OpenWebText (ppl)	Training time (speedu
GPT-2 small - Huggingface [87]	18.2	$9.5 \text{ days } (1.0 \times)$
GPT-2 small - Megatron-LM [77]	18.2	$4.7 \text{ days } (2.0\times)$
GPT-2 small - FlashAttention	18.2	$2.7 \text{ days } (3.5 \times)$
GPT-2 medium - Huggingface [87]	14.2	$21.0 \text{ days } (1.0 \times)$
GPT-2 medium - Megatron-LM [77]	14.3	11.5 days $(1.8\times)$
GPT-2 medium - FLASHATTENTION	14.3	$6.9 ext{ days } (3.0 \times)$

15% wall-clock speed-up on BERT-large (seq. Length 512);

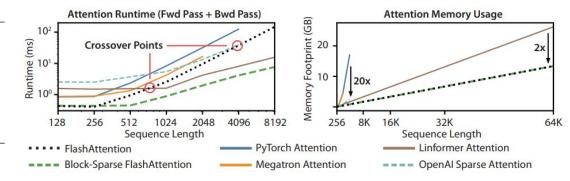
3x on GPT-2 (seq. Length 1K);

2.4x on long-range arena (seq. length 1K-4K)

How impactful is FlashAttention?

Model implementations	Context length	OpenWebText (ppl)	Training time (speedup)
GPT-2 small - Megatron-LM	1k	18.2	4.7 days (1.0×)
GPT-2 small - FLASHATTENTION	1k	18.2	2.7 days $(1.7 \times)$
GPT-2 small - FLASHATTENTION	2k	17.6	3.0 days $(1.6 \times)$
GPT-2 small - FLASHATTENTION	4k	17.5	$3.6 \text{ days } (1.3\times)$

\mathbf{Model}	Path-X	Path-256
Transformer	X	X
Linformer [84]	X	×
Linear Attention [50]	X	×
Performer [12]	X	×
Local Attention [80]	X	×
Reformer [51]	X	×
SMYRF [19]	X	×
FLASHATTENTION	61.4	X
Block-sparse FlashAttention	56.0	63.1



Summary

- FlashAttention: IO-aware exact attention
- Reduces HBM access. Avoids storage of QK^T for both forward and backward passes.
- Softmax is an algebraic operation
- Huge gains across GPT-2, BERT, and Transformer

Thank You

References:

- Dao, T., Fu, D., Ermon, S., Rudra, A. and Ré, C., 2022. Flashattention: Fast and memory-efficient exact attention with io-awareness. Advances in neural information processing systems, 35, pp.16344-16359.
- 2. Dao, T., 2023. Flashattention-2: Faster attention with better parallelism and work partitioning. arXiv preprint arXiv:2307.08691.
- 3. https://www.youtube.com/@jbhuang0604