



# **Fast and Memory-Efficient Exact Attention with IO-Awareness**

NeurIPS 2022

HUMANE Lab 김태균

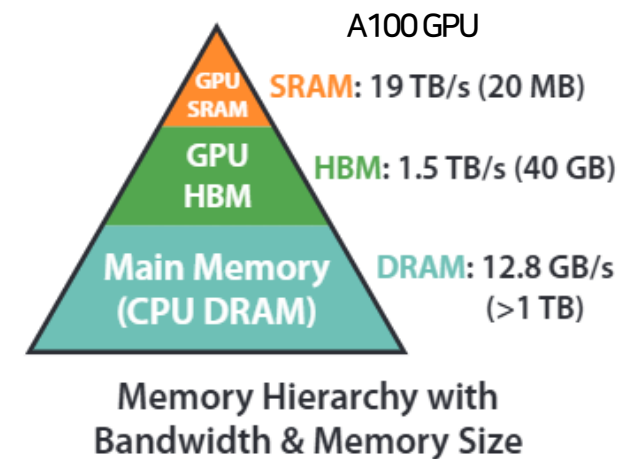
2025.03.14

# Introduction

---

- Sparse attention focuses on FLOP reduction but tends to ignore overheads from memory access (IO), resulting in lower wall-clock speed
- An IO-aware attention algorithm is needed to address this issue

=> Therefore, the goal is to reduce HBM access



# How to reduce HBM access?

---

1. Computing the softmax reduction
2. Not storing the large intermediate attention matrix for the backward pass

# How to reduce HBM access?

---

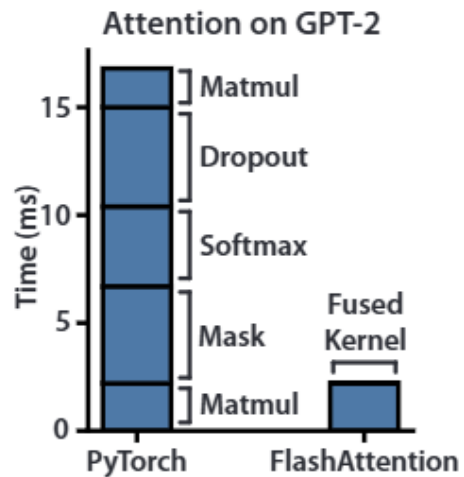
1. Computing the softmax reduction => **Tiling**
2. Not storing the large intermediate attention matrix for the backward pass => **Recomputation**

# Background

---

- Operations
  - compute-bound
  - memory-bound
- Kernel fusion

=> arithmetic intensity ( $\frac{FLOPs}{Byte}$ )



# FlashAttention

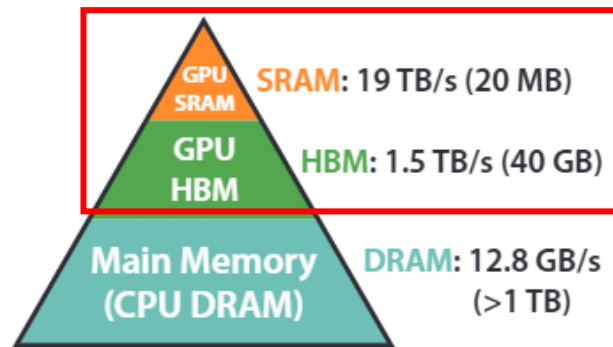
---

- Goal : avoid r/w the attention matrix to and from HBM
  - Tiling
  - Recomputation

$$\begin{aligned}S &= QK^T \\P &= \textit{softmax}(S) \\O &= PV\end{aligned}$$

# Tiling

- A technique that divides the entire matrix into blocks for computation
  - memory hierarchy
  - matrix multiplication



Memory Hierarchy with Bandwidth & Memory Size

A					B					C				
1	2	3	4	x	1	2	3	4	=	C <sub>11</sub>	C <sub>12</sub>			
5	6	7	8		5	6	7	8		C <sub>21</sub>	C <sub>22</sub>			
9	10	11	12		9	10	11	12						
13	14	15	16		13	14	15	16						

$$\textcircled{1} \quad \begin{aligned} C_{11} &= 1 \times 1 + 2 \times 5 + 3 \times 9 + 4 \times 13 \\ C_{12} &= \\ C_{21} &= \\ C_{22} &= \end{aligned}$$

$\Rightarrow 32$  memory access

$$\textcircled{2} \quad \begin{bmatrix} C_{11} & C_{12} \\ C_{21} & C_{22} \end{bmatrix} = \begin{bmatrix} 1 & 2 \\ 5 & 6 \end{bmatrix} \begin{bmatrix} 1 & 2 \\ 5 & 6 \end{bmatrix} + \begin{bmatrix} 3 & 4 \\ 7 & 8 \end{bmatrix} \begin{bmatrix} 9 & 10 \\ 13 & 14 \end{bmatrix}$$

$\Rightarrow 16$  memory access

# Tiling

---

- A technique that divides the entire matrix into blocks for computation
  - memory hierarchy
  - matrix multiplication

$$\begin{aligned} S &= QK^T \\ P &= \text{softmax}(S) \\ O &= PV \end{aligned}$$



# Tiling

- Safe softmax
  - subtract the maximum value to scale down
  - for numerical stability (to prevent overflow)

$$y_i = \frac{e^{x_i}}{\sum_{j=1}^V e^{x_j}}$$

---

**Algorithm 1** Naive softmax

---

```
1:  $d_0 \leftarrow 0$ 
2: for  $j \leftarrow 1, V$  do
3:    $d_j \leftarrow d_{j-1} + e^{x_j}$ 
4: end for
5: for  $i \leftarrow 1, V$  do
6:    $y_i \leftarrow \frac{e^{x_i}}{d_V}$ 
7: end for
```

---

memory access : 3V

$$y_i = \frac{e^{x_i - \max_{k=1}^V x_k}}{\sum_{j=1}^V e^{x_j - \max_{k=1}^V x_k}}$$

---

**Algorithm 2** Safe softmax

---

```
1:  $m_0 \leftarrow -\infty$ 
2: for  $k \leftarrow 1, V$  do
3:    $m_k \leftarrow \max(m_{k-1}, x_k)$ 
4: end for
5:  $d_0 \leftarrow 0$ 
6: for  $j \leftarrow 1, V$  do
7:    $d_j \leftarrow d_{j-1} + e^{x_j - m_V}$ 
8: end for
9: for  $i \leftarrow 1, V$  do
10:   $y_i \leftarrow \frac{e^{x_i - m_V}}{d_V}$ 
11: end for
```

---

memory access : 4V

Online normalizer calculation for softmax, arXiv preprint 2018

# Tiling

- Online safe softmax
  - compute m online
  - reduce the number of memory accesses by one

---

**Algorithm 3** Safe softmax with online normalizer calculation

---

```
1:  $m_0 \leftarrow -\infty$ 
2:  $d_0 \leftarrow 0$ 
3: for  $j \leftarrow 1, V$  do
4:    $m_j \leftarrow \max(m_{j-1}, x_j)$ 
5:    $d_j \leftarrow d_{j-1} \times e^{m_{j-1}-m_j} + e^{x_j-m_j}$ 
6: end for
7: for  $i \leftarrow 1, V$  do
8:    $y_i \leftarrow \frac{e^{x_i-m_V}}{d_V}$ 
9: end for
```

---

memory access :  $3V$

$$\begin{aligned} d_1 &= e^{\lambda_1 - m_1} \\ d_2 &= e^{\lambda_1 - m_1} e^{m_1 - m_2} + e^{\lambda_2 - m_2} \\ &= e^{\lambda_1 - m_2} + e^{\lambda_2 - m_2} \\ d_3 &= (e^{\lambda_1 - m_2} + e^{\lambda_2 - m_2}) e^{m_2 - m_3} + e^{\lambda_3 - m_3} \\ &= e^{\lambda_1 - m_3} + e^{\lambda_2 - m_3} + e^{\lambda_3 - m_3} \\ &\vdots \end{aligned}$$

Can be proven by mathematical induction

Online normalizer calculation for softmax, arXiv preprint 2018

# Tiling

- Online safe softmax
  - compute m online
  - reduce the number of memory accesses by one

---

**Algorithm 3** Safe softmax with online normalizer calculation

---

1:  $m_0 \leftarrow -\infty$

2:  $d_0 \leftarrow 0$

3: **for**  $j \leftarrow 1, V$  **do**

4:    $m_j \leftarrow \max(m_{j-1}, x_j)$

5:    $d_j \leftarrow d_{j-1} \times e^{m_{j-1}-m_j} + e^{x_j-m_j}$

6: **end for**

7: **for**  $i \leftarrow 1, V$  **do**

8:    $y_i \leftarrow \frac{e^{x_i-m_V}}{d_V}$

9: **end for**

---

memory access : 3V

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

Online normalizer calculation for softmax, arXiv preprint 2018

# Tiling

---

- Online safe softmax
  - compute m online
  - reduce the number of memory accesses by one

$$\begin{aligned} S &= QK^T \\ P &= \text{softmax}(S) \\ O &= PV \end{aligned}$$

# Tiling

for  $1 \leq i \leq N$

$$x_i = q k_i^T$$

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

$$d'_i = d'_{i-1} e^{m_{i-1} - m_i} + e^{x_i - m_i}$$

for  $1 \leq i \leq N$

$$a_i = e^{x_i - m_N} / d'_N$$

$$O_i = O_{i-1} + a_i V_i$$

return  $O_N$

where N is the  
number of blocks

trick

$$\begin{aligned}
 O_i &= \sum_{j=1}^i \frac{e^{x_j - m_N}}{d'_N} V_j \\
 O'_i &= \sum_{j=1}^i \frac{e^{x_j - m_i}}{d'_i} V_j \\
 &= \sum_{j=1}^{i-1} \frac{e^{x_j - m_i}}{d'_i} \cdot \left[ \frac{d'_{i-1}}{d'_{i-1}} \cdot \frac{e^{-m_{i-1}}}{e^{m_{i-1}}} \right] V_j + \frac{e^{x_i - m_i}}{d'_i} V_i \\
 &= \left( \sum_{j=1}^{i-1} \frac{e^{x_j - m_{i-1}}}{d'_{i-1}} V_j \right) \cdot \frac{d'_{i-1}}{d'_i} \cdot e^{m_{i-1} - m_i} + \frac{e^{x_i - m_i}}{d'_i} V_i \\
 &= O'_{i-1} \frac{d'_{i-1}}{d'_i} e^{m_{i-1} - m_i} + \frac{e^{x_i - m_i}}{d'_i} V_i
 \end{aligned}$$

remove  
dependency

# Tiling

for  $1 \leq i \leq N$

$$x_i = q k_i^T$$

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

$$d'_i = d'_{i-1} e^{m_{i-1} - m_i} + e^{x_i - m_i}$$

for  $1 \leq i \leq N$

$$a_i = e^{x_i - m_N} / d'_N$$

$$O_i = O_{i-1} + a_i V_i$$

return  $O_N$

fuse into a  
single loop



for  $1 \leq i \leq N$

$$x_i = q k_i^T$$

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

$$d'_i = d'_{i-1} e^{m_{i-1} - m_i} + e^{x_i - m_i}$$

$$O'_i = O'_{i-1} \frac{d'_{i-1}}{d'_i} e^{m_{i-1} - m_i} + \frac{e^{x_i - m_i}}{d'_i} V_i$$

return  $O'_N$

# Tiling

for  $1 \leq i \leq N$

$$x_i = qk_i^T$$

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

$$d'_i = d'_{i-1} e^{m_{i-1} - m_i} + e^{x_i - m_i}$$

for  $1 \leq i \leq N$

$$a_i = e^{x_i - m_N} / d'_N$$

$$O_i = O_{i-1} + a_i V_i$$

return  $O_N$

fuse into a  
single loop



~~$$S = QK^T$$~~

~~$$P = \text{softmax}(S)$$~~

~~$$O = PV$$~~

for  $1 \leq i \leq N$

$$x_i = qk_i^T$$

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

$$d'_i = d'_{i-1} e^{m_{i-1} - m_i} + e^{x_i - m_i}$$

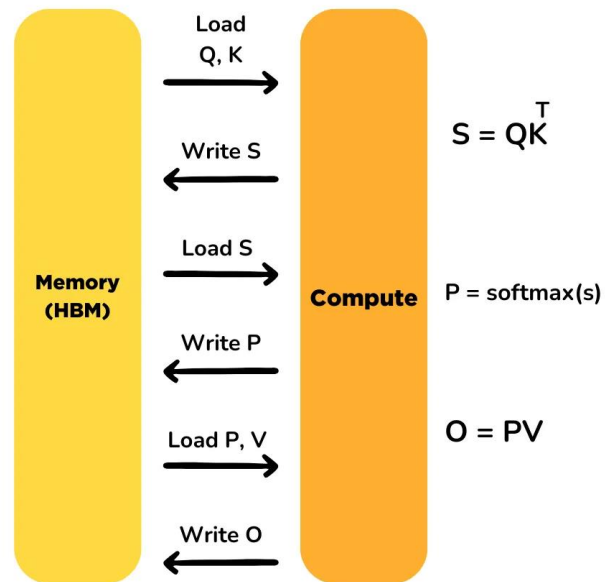
$$O'_i = O'_{i-1} \frac{d'_{i-1}}{d'_i} e^{m_{i-1} - m_i} + \frac{e^{x_i - m_i}}{d'_i} V_i$$

return  $O'_N$

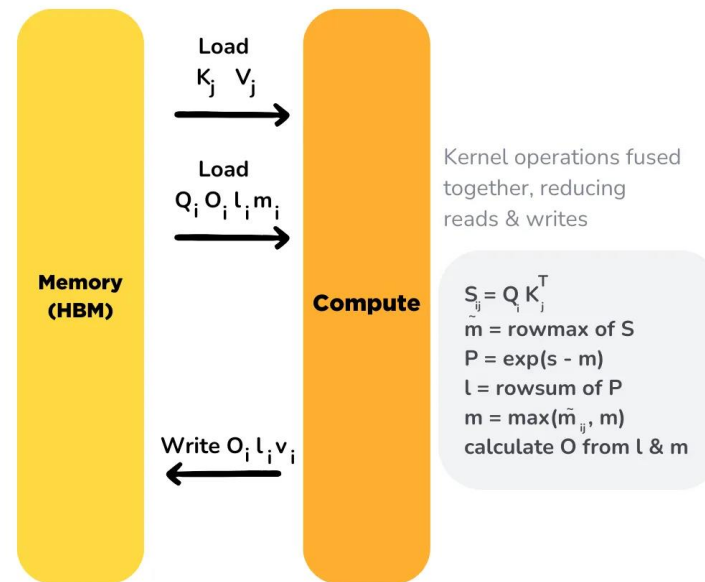
# Recomputation

- A technique of recompute intermediate values (S, P) during the backward pass instead of storing them

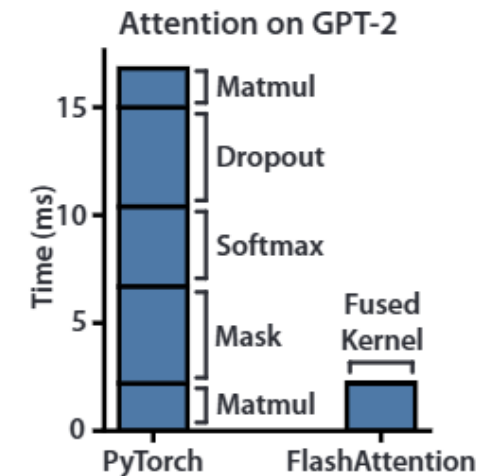
Standard Attention Implementation



Flash Attention



Initialize O, l and m matrices with zeroes. m and l 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.



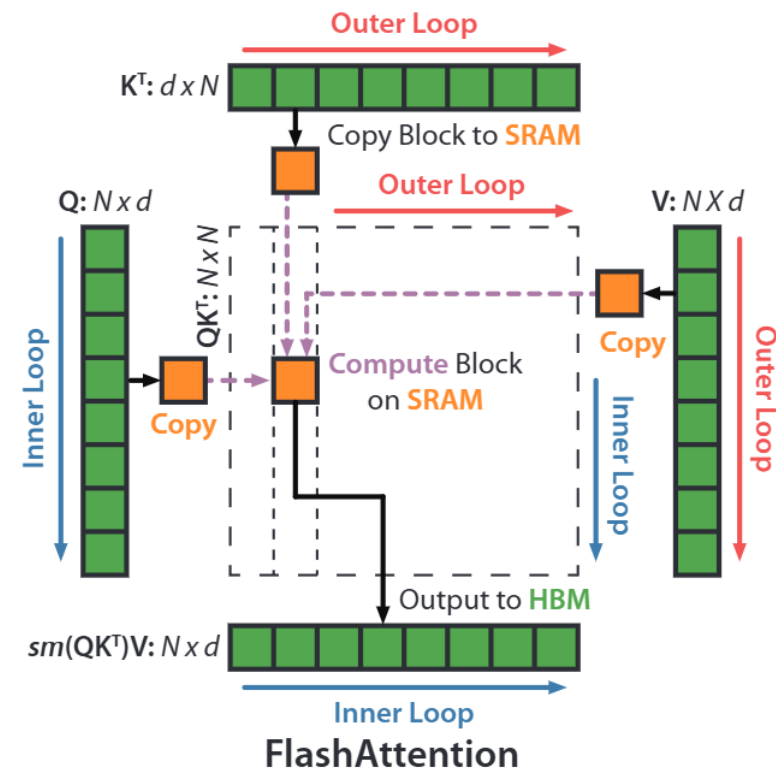


# Algorithm

## Algorithm 1 FLASHATTENTION

**Require:** Matrices  $\mathbf{Q}, \mathbf{K}, \mathbf{V} \in \mathbb{R}^{N \times d}$  in HBM, on-chip SRAM of size  $M$ .

- 1: Set block sizes  $B_c = \lceil \frac{M}{4d} \rceil, B_r = \min(\lceil \frac{M}{4d} \rceil, d)$ .
- 2: Initialize  $\mathbf{O} = (0)_{N \times d} \in \mathbb{R}^{N \times d}, \ell = (0)_N \in \mathbb{R}^N, m = (-\infty)_N \in \mathbb{R}^N$  in HBM.
- 3: Divide  $\mathbf{Q}$  into  $T_r = \lceil \frac{N}{B_r} \rceil$  blocks  $\mathbf{Q}_1, \dots, \mathbf{Q}_{T_r}$  of size  $B_r \times d$  each, and divide  $\mathbf{K}, \mathbf{V}$  in to  $T_c = \lceil \frac{N}{B_c} \rceil$  blocks  $\mathbf{K}_1, \dots, \mathbf{K}_{T_c}$  and  $\mathbf{V}_1, \dots, \mathbf{V}_{T_c}$ , of size  $B_c \times d$  each.
- 4: Divide  $\mathbf{O}$  into  $T_r$  blocks  $\mathbf{O}_i, \dots, \mathbf{O}_{T_r}$  of size  $B_r \times d$  each, divide  $\ell$  into  $T_r$  blocks  $\ell_i, \dots, \ell_{T_r}$  of size  $B_r$  each, divide  $m$  into  $T_r$  blocks  $m_1, \dots, m_{T_r}$  of size  $B_r$  each.
- 5: **for**  $1 \leq j \leq T_c$  **do**
- 6:   Load  $\mathbf{K}_j, \mathbf{V}_j$  from HBM to on-chip SRAM.
- 7:   **for**  $1 \leq i \leq T_r$  **do**
- 8:     Load  $\mathbf{Q}_i, \mathbf{O}_i, \ell_i, m_i$  from HBM to on-chip SRAM.
- 9:     On chip, compute  $\mathbf{S}_{ij} = \mathbf{Q}_i \mathbf{K}_j^T \in \mathbb{R}^{B_r \times B_c}$ .
- 10:    On chip, compute  $\tilde{m}_{ij} = \text{rowmax}(\mathbf{S}_{ij}) \in \mathbb{R}^{B_r}, \tilde{\mathbf{P}}_{ij} = \exp(\mathbf{S}_{ij} - \tilde{m}_{ij}) \in \mathbb{R}^{B_r \times B_c}$  (pointwise),  $\tilde{\ell}_{ij} = \text{rowsum}(\tilde{\mathbf{P}}_{ij}) \in \mathbb{R}^{B_r}$ .
- 11:    On chip, compute  $m_i^{\text{new}} = \max(m_i, \tilde{m}_{ij}) \in \mathbb{R}^{B_r}, \ell_i^{\text{new}} = e^{m_i - m_i^{\text{new}}} \ell_i + e^{\tilde{m}_{ij} - m_i^{\text{new}}} \tilde{\ell}_{ij} \in \mathbb{R}^{B_r}$ .
- 12:    Write  $\mathbf{O}_i \leftarrow \text{diag}(\ell_i^{\text{new}})^{-1} (\text{diag}(\ell_i) e^{m_i - m_i^{\text{new}}} \mathbf{O}_i + e^{\tilde{m}_{ij} - m_i^{\text{new}}} \tilde{\mathbf{P}}_{ij} \mathbf{V}_j)$  to HBM.
- 13:    Write  $\ell_i \leftarrow \ell_i^{\text{new}}, m_i \leftarrow m_i^{\text{new}}$  to HBM.
- 14:   **end for**
- 15: **end for**
- 16: Return  $\mathbf{O}$ .



# Block-sparse FlashAttention

---

- Better IO complexity than FlashAttention by a factor proportional to the sparsity ratio
  - Given a predefined block sparsity mask  $M$

$$\mathbf{S} = \mathbf{Q}\mathbf{K}^\top \in \mathbb{R}^{N \times N}, \quad \mathbf{P} = \text{softmax}(\mathbf{S} \odot \mathbf{1}_{\tilde{\mathbf{M}}}) \in \mathbb{R}^{N \times N}, \quad \mathbf{O} = \mathbf{P}\mathbf{V} \in \mathbb{R}^{N \times d}$$

# Experiments

---

- Training time & model accuracy

Model implementations	OpenWebText (ppl)	Training time (speedup)
GPT-2 small - Huggingface [87]	18.2	9.5 days (1.0×)
GPT-2 small - Megatron-LM [77]	18.2	4.7 days (2.0×)
GPT-2 small - FLASHATTENTION	18.2	<b>2.7 days (3.5×)</b>
GPT-2 medium - Huggingface [87]	14.2	21.0 days (1.0×)
GPT-2 medium - Megatron-LM [77]	14.3	11.5 days (1.8×)
GPT-2 medium - FLASHATTENTION	14.3	<b>6.9 days (3.0×)</b>

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	<b>2.7 days (1.7×)</b>
GPT-2 small - FLASHATTENTION	2k	17.6	3.0 days (1.6×)
GPT-2 small - FLASHATTENTION	4k	<b>17.5</b>	3.6 days (1.3×)

# Experiments

- Attention runtime & memory benchmarks

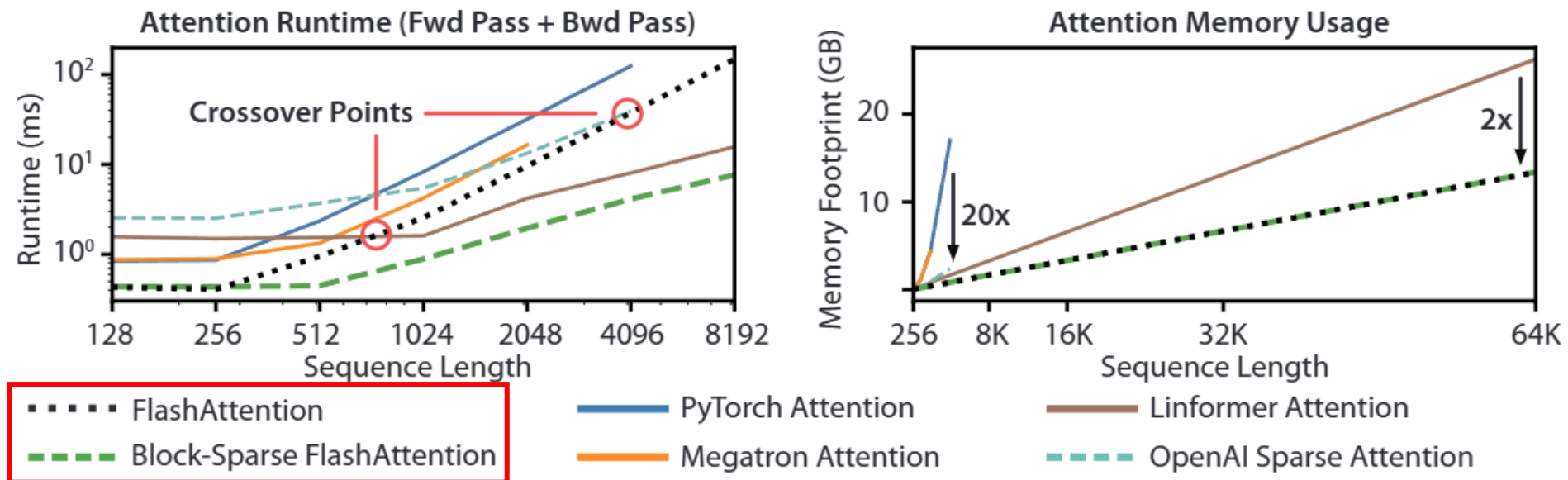


Figure 3: **Left:** runtime of forward pass + backward pass. **Right:** attention memory usage.

# Conclusion

---

- Showing significant reduction in HBM accesses compared to standard attention