

FlashAttention

cool paper

Attention

- Typical Attention Computing Flow:

Algorithm 0 Standard Attention Implementation

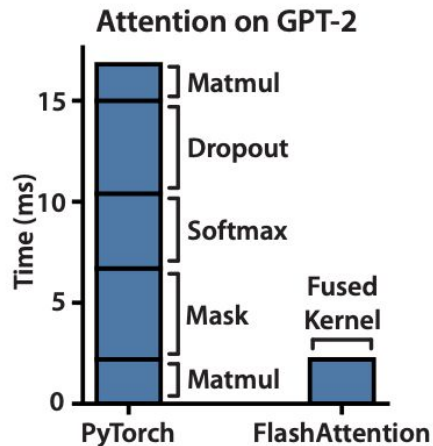
Require: Matrices $\mathbf{Q}, \mathbf{K}, \mathbf{V} \in \mathbb{R}^{N \times d}$ in HBM.

- 1: Load \mathbf{Q}, \mathbf{K} by blocks from HBM, compute $\mathbf{S} = \mathbf{Q}\mathbf{K}^\top$, write \mathbf{S} to HBM.
 - 2: Read \mathbf{S} from HBM, compute $\mathbf{P} = \text{softmax}(\mathbf{S})$, write \mathbf{P} to HBM.
 - 3: Load \mathbf{P} and \mathbf{V} by blocks from HBM, compute $\mathbf{O} = \mathbf{P}\mathbf{V}$, write \mathbf{O} to HBM.
 - 4: Return \mathbf{O} .
-

- TimeComplexity $O(N^2 * d)$
- Memory $O(N^2)$
- Usually $N \gg d$ (eg. $N=4000, d=128$)

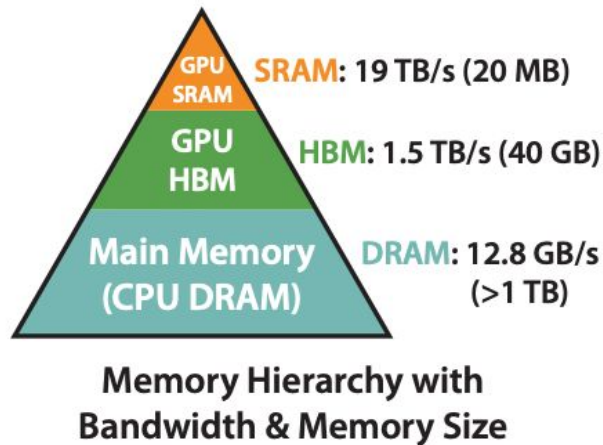
FlashAttention

- Fast and Memory Efficient **Exact** Attention with **IO-Awareness**
- Contribution
 - Exact: No approximation. (But it can be modified...)
 - Faster: BERT-large (15% quicker, E2E)
 - Use Less Memory: linear in sequence length



Background: GPU Memory Hierarchy

- A100
 - HBM(40-80G), 1.5TB/s bandwidth
 - GPU SRAM (192 KB per 108 streaming multiprocessors), 19TB/s
- Compute Bound vs Memory Bound
 - Compute Bound:
 - Matrix multiply with large inner dimension, and convolution with large number of channels
 - Memory Bound
 - elementwise (e.g., activation, dropout), and reduction (e.g., sum, softmax, batch norm, layer norm)



Algorithm - Challenges

Algorithm 0 Standard Attention Implementation

Require: Matrices $\mathbf{Q}, \mathbf{K}, \mathbf{V} \in \mathbb{R}^{N \times d}$ in HBM.

- 1: Load \mathbf{Q}, \mathbf{K} by blocks from HBM, compute $\mathbf{S} = \mathbf{Q}\mathbf{K}^\top$, write \mathbf{S} to HBM.
 - 2: Read \mathbf{S} from HBM, compute $\mathbf{P} = \text{softmax}(\mathbf{S})$, write \mathbf{P} to HBM.
 - 3: Load \mathbf{P} and \mathbf{V} by blocks from HBM, compute $\mathbf{O} = \mathbf{P}\mathbf{V}$, write \mathbf{O} to HBM.
 - 4: Return \mathbf{O} .
-

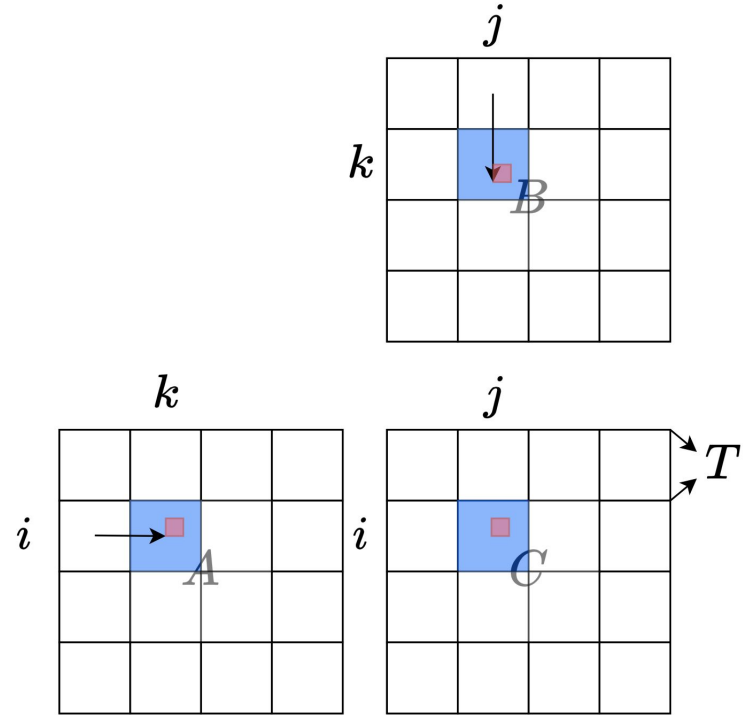
1. Matrix Multiplication - (Memory)

2. Softmax - (Memory, I/O)

3. Backward - (Memory)

Algorithm- Matrix Multiplication

- Tiling
- Only need $3T^2$ elements stored on chips



Algorithm- Online Softmax

- Safe Softmax (3-pass):

$$\frac{e^{x_i}}{\sum_{j=1}^N e^{x_j}} = \frac{e^{x_i - m}}{\sum_{j=1}^N e^{x_j - m}}$$

- m_N must be computed after one pass.
 - ... or not?

Algorithm 3-pass safe softmax

NOTATIONS

$\{m_i\}$: $\max_{j=1}^i \{x_j\}$, with initial value $m_0 = -\infty$.

$\{d_i\}$: $\sum_{j=1}^i e^{x_j - m_N}$, with initial value $d_0 = 0$, d_N is the denominator of safe softmax.

$\{a_i\}$: the final softmax value.

BODY

for $i \leftarrow 1, N$ do

$$m_i \leftarrow \max(m_{i-1}, x_i) \quad (7)$$

end

for $i \leftarrow 1, N$ do

$$d_i \leftarrow d_{i-1} + e^{x_i - m_N} \quad (8)$$

end

for $i \leftarrow 1, N$ do

$$a_i \leftarrow \frac{e^{x_i - m_N}}{d_N} \quad (9)$$

end

Algorithm- Online Softmax

- Safe Softmax (2-pass):
- Can we do 1 pass then?
 - Unfortunately, no.
 - ... But we are not computing $\text{Softmax}(Q^* K^T)$ only, we want $\text{Softmax}(Q^* K^T) * V$

$$\begin{aligned}d'_i &= \sum_{j=1}^i e^{x_j - m_i} \\&= \left(\sum_{j=1}^{i-1} e^{x_j - m_i} \right) + e^{x_i - m_i} \\&= \left(\sum_{j=1}^{i-1} e^{x_j - m_{i-1}} \right) e^{m_{i-1} - m_i} + e^{x_i - m_i} \\&= d'_{i-1} e^{m_{i-1} - m_i} + e^{x_i - m_i}\end{aligned}$$

Algorithm- Online Softmax

- Magically, we can save one more pass.

$$\mathbf{o}_i := \sum_{j=1}^i \left(\frac{e^{x_j - m_N}}{d'_N} V[j, :] \right) \quad (13)$$

This still depends on m_N and d_N which cannot be determined until the previous loop completes. But we can play the “surrogate” trick in section 3 again, by creating a surrogate sequence \mathbf{o}' :

$$\mathbf{o}'_i := \left(\sum_{j=1}^i \frac{e^{x_j - m_i}}{d'_i} V[j, :] \right)$$

The n -th element of \mathbf{o} and \mathbf{o}' are the identical: $\mathbf{o}'_N = \mathbf{o}_N$, and we can find a recurrence relation between \mathbf{o}'_i and \mathbf{o}'_{i-1} :

$$\begin{aligned} \mathbf{o}'_i &= \sum_{j=1}^i \frac{e^{x_j - m_i}}{d'_i} V[j, :] \\ &= \left(\sum_{j=1}^{i-1} \frac{e^{x_j - m_i}}{d'_i} V[j, :] \right) + \frac{e^{x_i - m_i}}{d'_i} V[i, :] \\ &= \left(\sum_{j=1}^{i-1} \frac{e^{x_j - m_{i-1}}}{d'_{i-1}} \frac{e^{x_j - m_i}}{e^{x_j - m_{i-1}}} \frac{d'_{i-1}}{d'_i} V[j, :] \right) + \frac{e^{x_i - m_i}}{d'_i} V[i, :] \\ &= \left(\sum_{j=1}^{i-1} \frac{e^{x_j - m_{i-1}}}{d'_{i-1}} V[j, :] \right) \frac{d'_{i-1}}{d'_i} e^{m_{i-1} - m_i} + \frac{e^{x_i - m_i}}{d'_i} V[i, :] \\ &= \mathbf{o}'_{i-1} \frac{d'_{i-1} e^{m_{i-1} - m_i}}{d'_i} + \frac{e^{x_i - m_i}}{d'_i} V[i, :] \end{aligned} \quad (14)$$

Algorithm - All Together

Algorithm 1 FLASHATTENTION

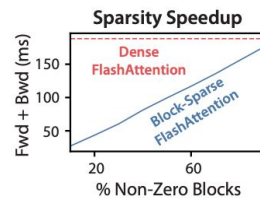
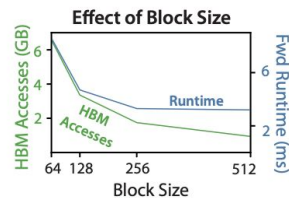
Require: Matrices $\mathbf{Q}, \mathbf{K}, \mathbf{V} \in \mathbb{R}^{N \times d}$ in HBM, on-chip SRAM of size M .

- 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} into $T_c = \lceil \frac{N}{B_c} \rceil$ 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 $\mathbf{O}_1, \dots, \mathbf{O}_{T_r}$ of size $B_r \times d$ each, divide ℓ into T_r blocks $\ell_1, \dots, \ell_{T_r}$ 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} .
-

Algorithm

- $O(N^2 \cdot d)$ FLOPs,
 $O(N)$ extra memory
- HBM access?
 - Standard: $O(Nd + N^2)$
 - Flashattention:
 $O(N^2 \cdot d^2 / M)$
 - In practice, $d = 64-128$,
 $M=100k$, $d^2 \ll M$

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



Algorithm - Backward?

- Recompute!
 - Q, K, V, O, I, m has already been stored.
 - For each block of dQ, dK and dV, recompute the block of S and P.

Algorithm 4 FLASHATTENTION Backward Pass

Require: Matrices $\mathbf{Q}, \mathbf{K}, \mathbf{V}, \mathbf{O}, \mathbf{dO} \in \mathbb{R}^{N \times d}$ in HBM, vectors $\ell, m \in \mathbb{R}^N$ in HBM, on-chip SRAM of size M , softmax scaling constant $\tau \in \mathbb{R}$, masking function MASK, dropout probability p_{drop} , pseudo-random number generator state \mathcal{R} from the forward pass.

- 1: Set the pseudo-random number generator state to \mathcal{R} .
- 2: Set block sizes $B_c = \lceil \frac{M}{4d} \rceil$, $B_r = \min(\lceil \frac{M}{4d} \rceil, d)$.
- 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 $\mathbf{O}_1, \dots, \mathbf{O}_{T_r}$ of size $B_r \times d$ each, divide \mathbf{dO} into T_r blocks $\mathbf{dO}_1, \dots, \mathbf{dO}_{T_r}$ of size $B_r \times d$ each, divide ℓ into T_r blocks $\ell_1, \dots, \ell_{T_r}$ of size B_r each, divide m into T_r blocks m_1, \dots, m_{T_r} of size B_r each.
- 5: Initialize $\mathbf{dQ} = (0)_{N \times d}$ in HBM and divide it into T_r blocks $\mathbf{dQ}_1, \dots, \mathbf{dQ}_{T_r}$ of size $B_r \times d$ each. Initialize $\mathbf{dK} = (0)_{N \times d}$, $\mathbf{dV} = (0)_{N \times d}$ in HBM and divide \mathbf{dK}, \mathbf{dV} in to T_c blocks $\mathbf{dK}_1, \dots, \mathbf{dK}_{T_c}$ and $\mathbf{dV}_1, \dots, \mathbf{dV}_{T_c}$, of size $B_c \times d$ each.
- 6: **for** $1 \leq j \leq T_c$ **do**
- 7: Load $\mathbf{K}_j, \mathbf{V}_j$ from HBM to on-chip SRAM.
- 8: Initialize $\mathbf{dK}_j = (0)_{B_c \times d}$, $\mathbf{dV}_j = (0)_{B_c \times d}$ on SRAM.
- 9: **for** $1 \leq i \leq T_r$ **do**
- 10: Load $\mathbf{Q}_i, \mathbf{O}_i, \mathbf{dO}_i, \mathbf{dQ}_i, \ell_i, m_i$ from HBM to on-chip SRAM.
- 11: On chip, compute $\mathbf{S}_{ij} = \tau \mathbf{Q}_i \mathbf{K}_j^T \in \mathbb{R}^{B_r \times B_c}$.
- 12: On chip, compute $\mathbf{S}_{ij}^{\text{masked}} = \text{MASK}(\mathbf{S}_{ij})$.
- 13: On chip, compute $\mathbf{P}_{ij} = \text{diag}(\ell_i)^{-1} \exp(\mathbf{S}_{ij}^{\text{masked}} - m_i) \in \mathbb{R}^{B_r \times B_c}$.
- 14: On chip, compute dropout mask $\mathbf{Z}_{ij} \in \mathbb{R}^{B_r \times B_c}$ where each entry has value $\frac{1}{1-p_{\text{drop}}}$ with probability $1 - p_{\text{drop}}$ and value 0 with probability p_{drop} .
- 15: On chip, compute $\mathbf{P}_{ij}^{\text{dropped}} = \mathbf{P}_{ij} \circ \mathbf{Z}_{ij}$ (pointwise multiply).
- 16: On chip, compute $\mathbf{dV}_j \leftarrow \mathbf{dV}_j + (\mathbf{P}_{ij}^{\text{dropped}})^T \mathbf{dO}_i \in \mathbb{R}^{B_c \times d}$.
- 17: On chip, compute $\mathbf{dP}_{ij}^{\text{dropped}} = \mathbf{dO}_i \mathbf{V}_j^T \in \mathbb{R}^{B_r \times B_c}$.
- 18: On chip, compute $\mathbf{dP}_{ij} = \mathbf{dP}_{ij}^{\text{dropped}} \circ \mathbf{Z}_{ij}$ (pointwise multiply).
- 19: On chip, compute $\mathbf{D}_i = \text{rowsum}(\mathbf{dO}_i \circ \mathbf{O}_i) \in \mathbb{R}^{B_r}$.
- 20: On chip, compute $\mathbf{dS}_{ij} = \mathbf{P}_{ij} \circ (\mathbf{dP}_{ij} - \mathbf{D}_i) \in \mathbb{R}^{B_r \times B_c}$.
- 21: Write $\mathbf{dQ}_i \leftarrow \mathbf{dQ}_i + \mathbf{dS}_{ij} \mathbf{K}_j \in \mathbb{R}^{B_r \times d}$ to HBM.
- 22: On chip, compute $\mathbf{dK}_j \leftarrow \mathbf{dK}_j + \tau \mathbf{dS}_{ij} \mathbf{Q}_i \in \mathbb{R}^{B_c \times d}$.
- 23: **end for**
- 24: Write $\mathbf{dK}_j \leftarrow \mathbf{dK}_j, \mathbf{dV}_j \leftarrow \mathbf{dV}_j$ to HBM.
- 25: **end for**
- 26: Return $\mathbf{dQ}, \mathbf{dK}, \mathbf{dV}$.

Block-Sparse FlashAttention

- Approximate attention
 - Masking specific block for S to reduce FLOPs
 - $\Theta(ND + N^2 d^2 * s/M)$ HBM accesses
 - s is the fraction of nonzero blocks in the block-sparsity mask, s can be \sqrt{N} ...

$$\mathbf{P} = \text{softmax}(\mathbf{S} \odot \mathbf{1}_{\tilde{\mathbf{M}}}) \in \mathbb{R}^{N \times N},$$

Experiment

Table 2: GPT-2 small and medium using FLASHATTENTION achieve up to 3× speed up compared to Huggingface implementation and up to 1.7× compared to Megatron-LM. Training time reported on 8×A100s GPUs.

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	2.7 days (3.5×)
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	6.9 days (3.0×)

FlashAttention 2

- [Paper](#)
- Improvement
 - Better Iteration
 - Better usage of MatMul Acceleration
 - Better Parallelism

FlashAttention 2 - Fix

- O_i is the temporary rows of O for Q_i
- Why not lifting Q_i to outer loops?

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} into $T_c = \lceil \frac{N}{B_c} \rceil$ 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 $\mathbf{O}_1, \dots, \mathbf{O}_{T_r}$ of size $B_r \times d$ each, divide ℓ into T_r blocks $\ell_1, \dots, \ell_{T_r}$ 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 $\mathbf{O}_i \leftarrow \text{diag}(\ell_i^{\text{new}})^{-1} (\text{diag}(\ell_i) e^{m_i - m_i^{\text{new}}} \mathbf{O}_i + e^{m_{ij} - m_i^{\text{new}}} \tilde{\mathbf{P}}_{ij} \mathbf{V}_j)$ 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} .
-

FlashAttention 2 - Online Softmax optimization

- Do not need to divide by L_i each time, only do it in the last time.

Algorithm 1 FLASHATTENTION-2 forward pass

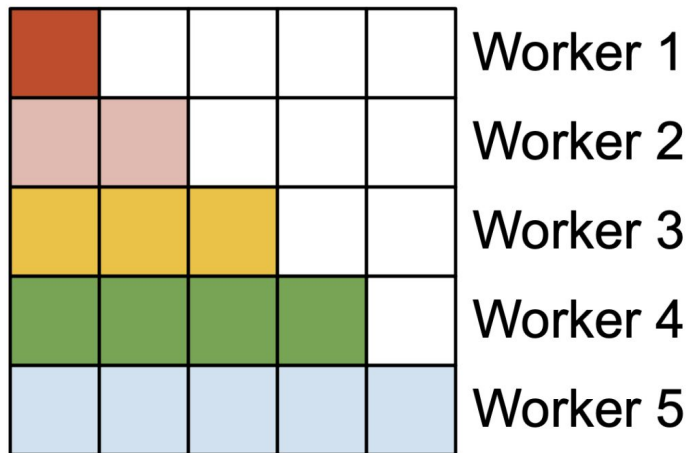
Require: Matrices $\mathbf{Q}, \mathbf{K}, \mathbf{V} \in \mathbb{R}^{N \times d}$ in HBM, block sizes B_c, B_r .

- 1: Divide \mathbf{Q} into $T_r = \left\lceil \frac{N}{B_r} \right\rceil$ blocks $\mathbf{Q}_1, \dots, \mathbf{Q}_{T_r}$ of size $B_r \times d$ each, and divide \mathbf{K}, \mathbf{V} into $T_c = \left\lceil \frac{N}{B_c} \right\rceil$ 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.
 - 2: Divide the output $\mathbf{O} \in \mathbb{R}^{N \times d}$ into T_r blocks $\mathbf{O}_1, \dots, \mathbf{O}_{T_r}$ of size $B_r \times d$ each, and divide the logsumexp L into T_r blocks L_1, \dots, L_{T_r} of size B_r each.
 - 3: **for** $1 \leq i \leq T_r$ **do**
 - 4: Load \mathbf{Q}_i from HBM to on-chip SRAM.
 - 5: On chip, initialize $\mathbf{O}_i^{(0)} = (0)_{B_r \times d} \in \mathbb{R}^{B_r \times d}, \ell_i^{(0)} = (0)_{B_r} \in \mathbb{R}^{B_r}, m_i^{(0)} = (-\infty)_{B_r} \in \mathbb{R}^{B_r}$.
 - 6: **for** $1 \leq j \leq T_c$ **do**
 - 7: Load $\mathbf{K}_j, \mathbf{V}_j$ from HBM to on-chip SRAM.
 - 8: On chip, compute $\mathbf{S}_i^{(j)} = \mathbf{Q}_i \mathbf{K}_j^T \in \mathbb{R}^{B_r \times B_c}$.
 - 9: On chip, compute $m_i^{(j)} = \max(m_i^{(j-1)}, \text{rowmax}(\mathbf{S}_i^{(j)})) \in \mathbb{R}^{B_r}, \tilde{\mathbf{P}}_i^{(j)} = \exp(\mathbf{S}_i^{(j)} - m_i^{(j)}) \in \mathbb{R}^{B_r \times B_c}$
(pointwise), $\ell_i^{(j)} = e^{m_i^{(j-1)} - m_i^{(j)}} \ell_i^{(j-1)} + \text{rowsum}(\tilde{\mathbf{P}}_i^{(j)}) \in \mathbb{R}^{B_r}$.
 - 10: On chip, compute $\mathbf{O}_i^{(j)} = \text{diag}(e^{m_i^{(j-1)} - m_i^{(j)}})^{-1} \mathbf{O}_i^{(j-1)} + \tilde{\mathbf{P}}_i^{(j)} \mathbf{V}_j$.
 - 11: **end for**
 - 12: On chip, compute $\mathbf{O}_i = \text{diag}(\ell_i^{(T_c)})^{-1} \mathbf{O}_i^{(T_c)}$.
 - 13: On chip, compute $L_i = m_i^{(T_c)} + \log(\ell_i^{(T_c)})$.
 - 14: Write \mathbf{O}_i to HBM as the i -th block of \mathbf{O} .
 - 15: Write L_i to HBM as the i -th block of L .
 - 16: **end for**
 - 17: Return the output \mathbf{O} and the logsumexp L .
-

FlashAttention 2 - Parallelism

- FlashAttention 1: parallelizes over batch size and number of heads
- FlashAttention 2: Also over the length of input text.

Forward pass



FlashAttention 2 - Work Partitioning Between Warps

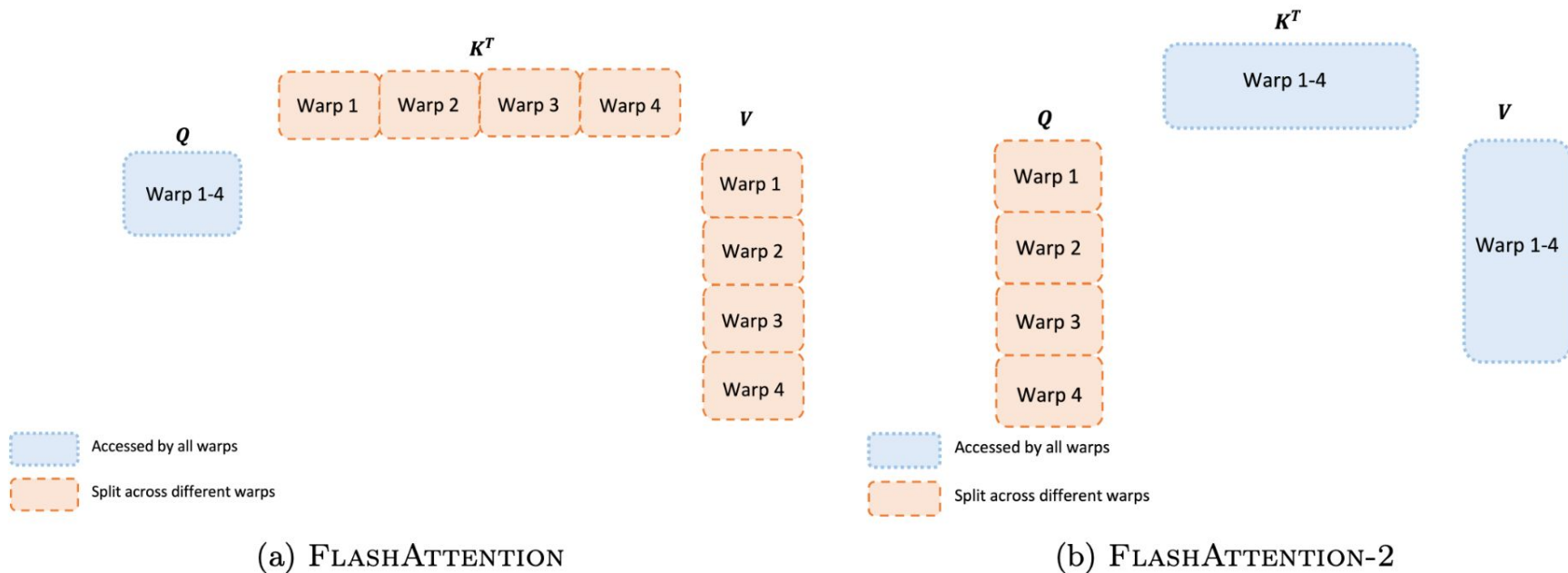
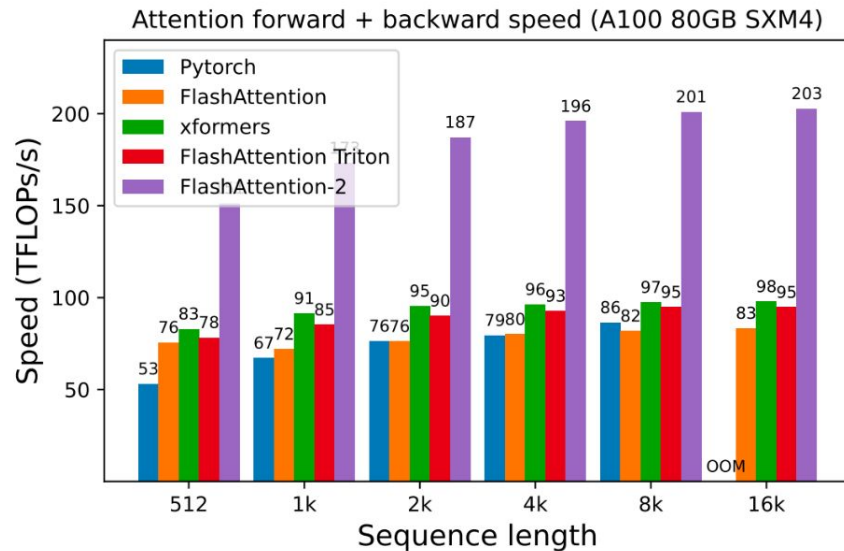
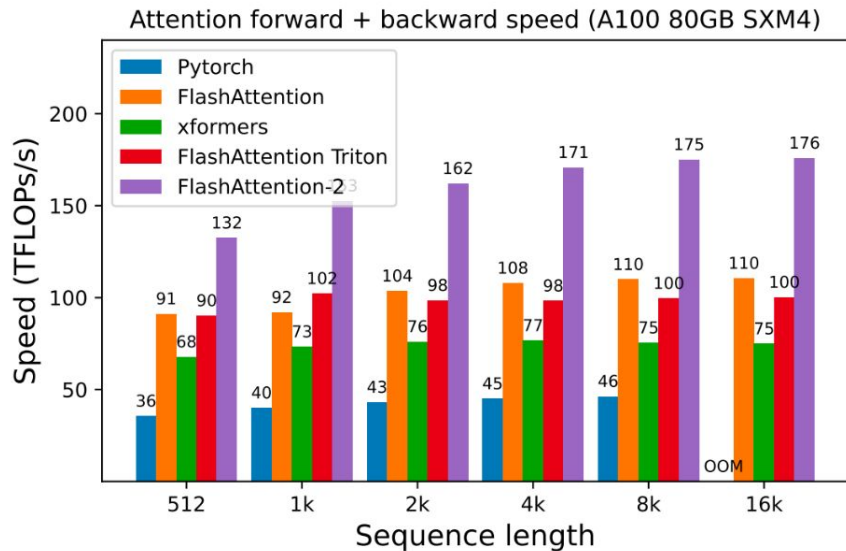


Figure 3: Work partitioning between different warps in the forward pass

FlashAttention 2 - Experiment



FlashAttention 3 (Flash-Decoding)

- Optimization on decoding/inference
- Why?
 - During inference, the query length is typically 1. So the partition of query length doesn't work.
- Solution:
 - First, we split the keys/values in smaller chunks
 - We compute the attention of the query with each of these splits in parallel using FlashAttention. We also write 1 extra scalar per row and per split: the log-sum-exp of the attention values.
 - Finally, we compute the actual output by reducing over all the splits, using the log-sum-exp to scale the contribution of each split.
- Experiment
 - The up to 8x speedup end-to-end measured earlier is made possible because the attention itself is up to 50x faster than FlashAttention.

Reference

- [FlashAttention](#)
- [FlashAttention V2](#)
- <https://medium.com/@sthanikamsanthosh1994/introduction-to-flash-attention-a-breakthrough-in-efficient-attention-mechanism-3eb47e8962c3>
- https://huggingface.co/docs/text-generation-inference/en/conceptual/flash_attention
- [From Online Softmax to FlashAttention](#)
- [FlashAttention V3](#)