

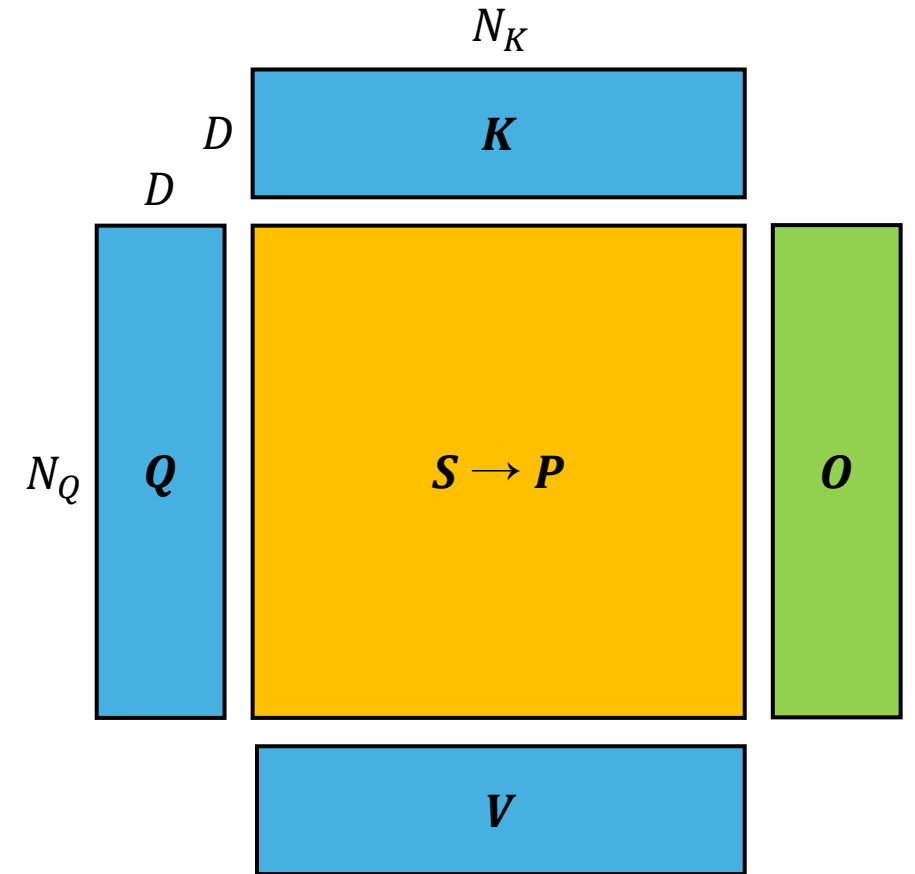
MSRA SH Triton Study Group

5. Flash-Attention (Algorithm)

2025/08/15

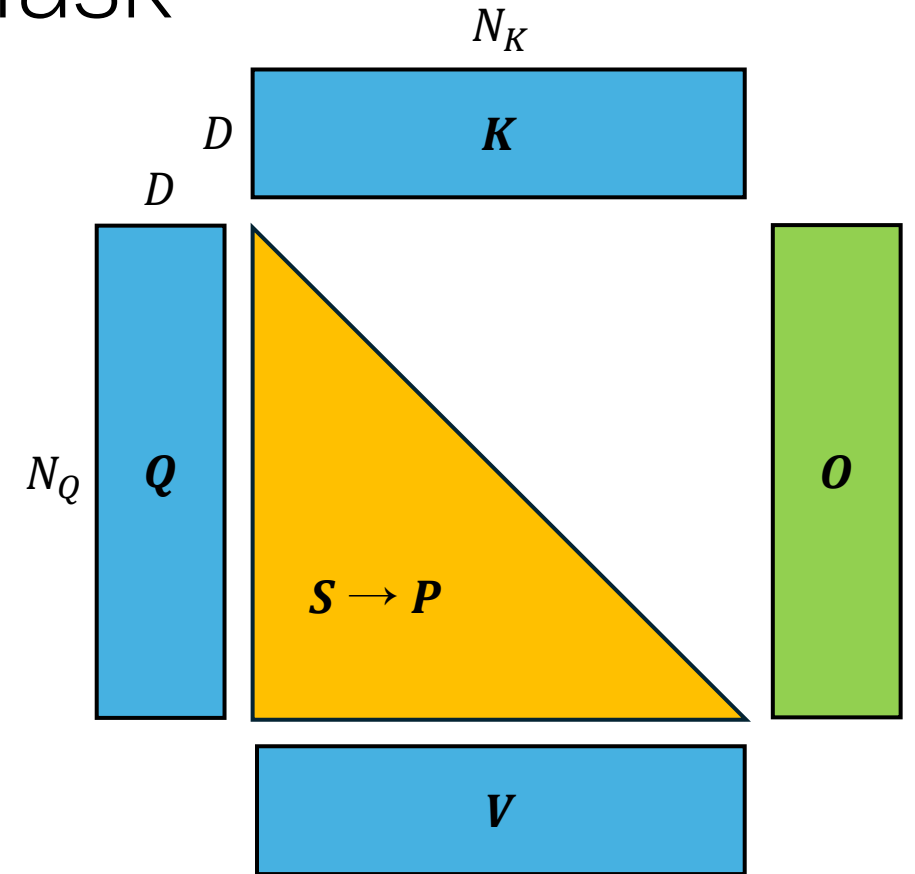
Attention

- $\mathbf{Q} \in \mathbb{R}^{N_Q \times D}$; $\mathbf{K}, \mathbf{V} \in \mathbb{R}^{N_K \times D}$
- $\mathbf{S} = \frac{\mathbf{Q}\mathbf{K}^T}{\sqrt{D}} \in \mathbb{R}^{N_Q \times N_K}$
- $\mathbf{P}_{ij} = \text{softmax}(\mathbf{S}_{ij}) = \frac{\exp(\mathbf{S}_{ij} - \max_j(\mathbf{S}_{ij}))}{\sum_j \exp(\mathbf{S}_{ij} - \max_j(\mathbf{S}_{ij}))}$
- $\mathbf{O} = \mathbf{P}\mathbf{V} \in \mathbb{R}^{N_Q \times D}$

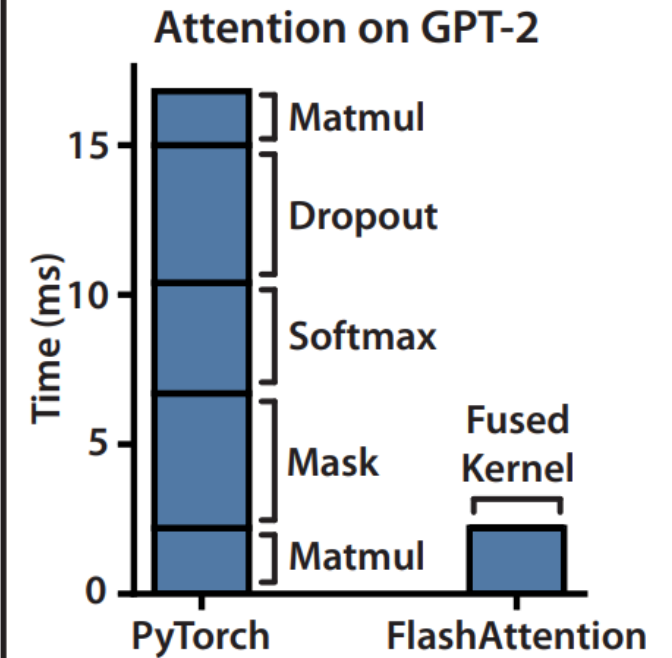
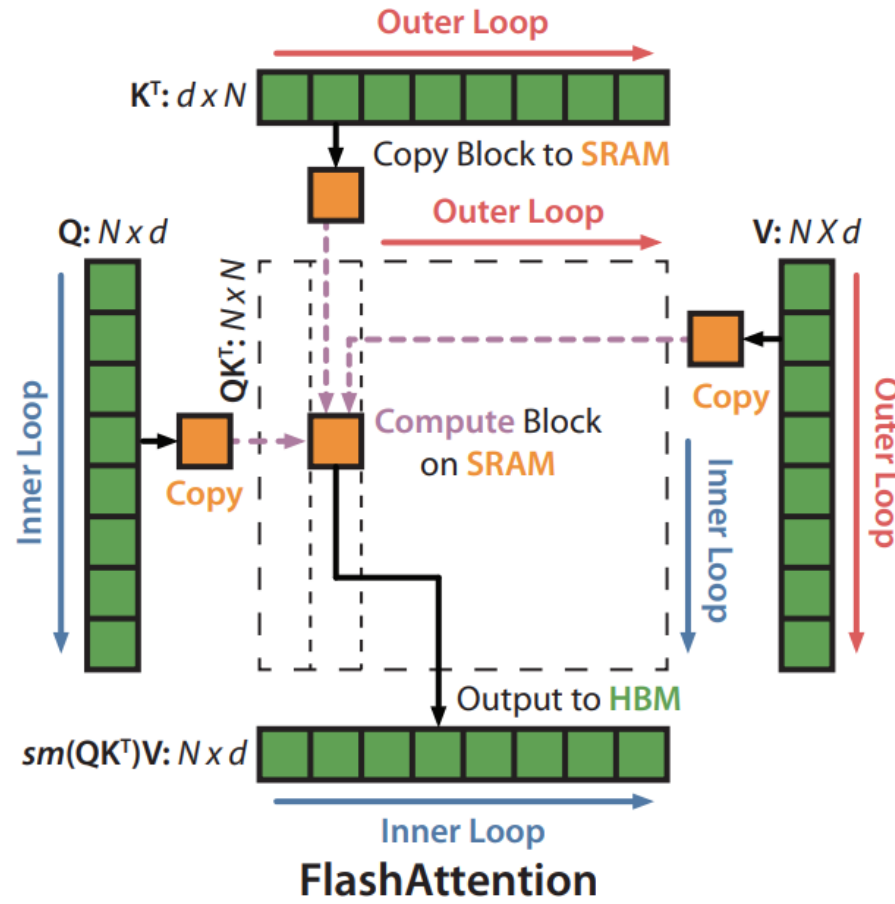
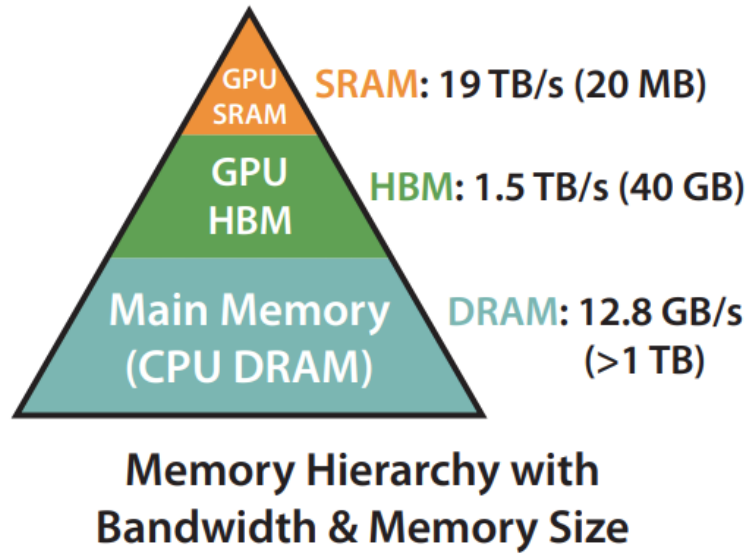


Self-Attention with Causal Mask

- Definition:
 - $N_Q = N_K$
 - $\mathbf{P}_{ij} \leftarrow \begin{cases} \mathbf{P}_{ij}, i \leq j \\ -\infty, i > j \end{cases}$
- Challenges:
 - The big $N_Q \times N_K$ tensor
 - `softmax()` is inefficient
 - Redundant calculation on masked \mathbf{P}_{ij}



Flash Attention



Fused Attention: GeMM-Lever Efficiency

- Assume 100% L2 cache hit ratio

	Total HBM Access	Total Calculation	Arithmetic Intensity
$\mathbf{S} = \mathbf{QK}^T / \sqrt{D}$	$N_Q D + N_K D + N_Q N_K$	$N_Q N_K D$	Memory-bound
$\mathbf{P}_{ij} = \text{softmax}(\mathbf{S}_{ij})$	$2N_Q N_K$	$C N_Q N_K$	Memory-bound
$\mathbf{O} = \mathbf{PV}$	$N_Q N_K + N_K D + N_Q D$	$N_Q N_K D$	Memory-bound
Fused attn()	$2N_Q D + 2N_K D$	$(2D + C)N_Q N_K$	Balanced

- $-N_Q N_K$ GPU HBM cost, $-4N_Q N_K$ GPU HBM access

Fused Attention: GeMM-Lever Efficiency

- Assume no L2 cache; $T_Q = T_K = D$; “best intensity” $E := \frac{T_Q T_K}{T_Q + T_K}$

	Total HBM Access	Total Calculation	Arithmetic Intensity	Latency
$\mathbf{S} = \frac{\mathbf{QK}^T}{\sqrt{D}}$	$\frac{N_K}{T_K} N_Q D + \frac{N_Q}{T_Q} N_K D + N_Q N_K$	$N_Q N_K D$	$\frac{T_Q T_K}{T_Q + T_K + \frac{T_Q T_K}{D}} < E$	$1 := \text{mem}(3N_Q N_K)$
$\mathbf{P}_{ij} = \text{softmax}(\mathbf{S}_{ij})$	$2N_Q N_K$	$C N_Q N_K$	$\frac{C}{2} \ll E$	$\frac{2}{3}$
$\mathbf{O} = \mathbf{PV}$	$N_Q N_K + \frac{N_Q}{T_Q} N_K D + N_Q D$	$N_Q N_K D$	$\frac{T_Q T_K}{\frac{T_Q T_K}{D} + T_K + \frac{T_K}{N_K} T_Q} \approx E$	$\approx \frac{2}{3}$
Fused attn()	$2 \frac{N_K}{T_K} N_Q D + 2 \frac{N_Q}{T_Q} N_K D$	$(2D + C) N_Q N_K$	$\frac{2D + C}{2D} \frac{T_Q T_K}{T_Q + T_K} \approx E$	$\frac{4}{3}$

Online Softmax

$$\bullet \begin{cases} \mathbf{P}_{ij}^0 = \frac{\exp(\mathbf{s}_{ij}^0 - \max_{0 \leq l < T}(\mathbf{s}_{il}^0))}{\sum_{l=0}^T \exp(\mathbf{s}_{ij}^0 - \max_{0 \leq l < T}(\mathbf{s}_{il}^0))} = \frac{\exp(\mathbf{s}_{ij}^0 - M_i^0)}{L_i^0} \\ \mathbf{P}_{ij}^1 = \frac{\exp(\mathbf{s}_{ij}^1 - \max_{T \leq l < 2T}(\mathbf{s}_{il}^1))}{\sum_{l=T}^{2T} \exp(\mathbf{s}_{ij}^1 - \max_{T \leq l < 2T}(\mathbf{s}_{il}^1))} = \frac{\exp(\mathbf{s}_{ij}^1 - M_i^1)}{L_i^1} \end{cases}$$

$$\bullet M_i^{[0,1]} = \max(M_i^0, M_i^1)$$

$$\bullet L_i^{[0,1]} = \exp(M_i^0 - M_i^{[0,1]}) L_i^0 + \exp(M_i^1 - M_i^{[0,1]}) L_i^1$$

$$\bullet \mathbf{P}_{ij}^{[0,1]} = \frac{\exp(\mathbf{s}_{ij}^{[0,1]} - \max_{0 \leq l < 2T}(\mathbf{s}_{il}^{[0,1]}))}{\sum_{l=0}^{2T} \exp(\mathbf{s}_{ij}^{[0,1]} - \max_{0 \leq l < 2T}(\mathbf{s}_{il}^{[0,1]}))} = \frac{\exp(\mathbf{s}_{ij}^{[0,1]} - M_i^{[0,1]})}{L_i^{[0,1]}}$$

- For $1 \leq j \leq \frac{N}{T}$
 - Load \mathbf{s}_{ij}
 - $M_i^{\text{local}} \leftarrow \max_{-1}(\mathbf{s}_{ij})$
 - $L_i^{\text{local}} \leftarrow \text{sum}_{-1}(\exp(\mathbf{s}_{ij} - M_i^{\text{local}}))$
 - $\mathbf{P}_{ij}^{\text{local}} \leftarrow \frac{\exp(\mathbf{s}_{ij} - M_i^{\text{local}})}{L_i^{\text{local}}}$
 - Update (M_i, L_i)

$$\bullet \text{ For } 1 \leq j \leq \frac{N}{T}$$

$$\bullet \mathbf{P}_{ij} \leftarrow \frac{\exp(\mathbf{s}_{ij} - M_i)}{L_i}$$

$$= \exp(M_i^{\text{local}} - M_i) \frac{L_i^{\text{local}}}{L_i} \mathbf{P}_{ij}^{\text{local}}$$

How to Fuse: Forward

- $$\begin{cases} \mathbf{s}_{ij}^0 = D^{-1/2} \sum_{k=0}^D \mathbf{Q}_{ik} \mathbf{K}_{jk}^0 \\ \mathbf{s}_{ij}^1 = D^{-1/2} \sum_{k=0}^D \mathbf{Q}_{ik} \mathbf{K}_{jk}^1 \end{cases}$$

- $$\begin{cases} \mathbf{P}_{ij}^0 = \frac{\exp(\mathbf{s}_{ij}^0 - \max_{0 \leq l < T}(\mathbf{s}_{il}^0))}{\sum_{l=0}^T \exp(\mathbf{s}_{ij}^0 - \max_{0 \leq l < T}(\mathbf{s}_{il}^0))} = \frac{\exp(\mathbf{s}_{ij}^0 - M_i^0)}{L_i^0} \\ \mathbf{P}_{ij}^1 = \frac{\exp(\mathbf{s}_{ij}^1 - \max_{T \leq l < 2T}(\mathbf{s}_{il}^1))}{\sum_{l=T}^{2T} \exp(\mathbf{s}_{ij}^1 - \max_{T \leq l < 2T}(\mathbf{s}_{il}^1))} = \frac{\exp(\mathbf{s}_{ij}^1 - M_i^1)}{L_i^1} \end{cases}$$

- $$M_i^{[0,1]} = \max(M_i^0, M_i^1)$$

- $$L_i^{[0,1]} = \exp(M_i^0 - M_i^{[0,1]}) L_i^0 + \exp(M_i^1 - M_i^{[0,1]}) L_i^1$$

- $$\mathbf{P}_{ij}^{[0,1]} = \frac{\exp(\mathbf{s}_{ij}^{[0,1]} - \max_{0 \leq l < 2T}(\mathbf{s}_{il}^{[0,1]}))}{\sum_{l=0}^{2T} \exp(\mathbf{s}_{ij}^{[0,1]} - \max_{0 \leq l < 2T}(\mathbf{s}_{il}^{[0,1]}))} = \frac{\exp(\mathbf{s}_{ij}^{[0,1]} - M_i^{[0,1]})}{L_i^{[0,1]}}$$

- $$\begin{cases} \mathbf{O}_{ik}^0 = \sum_{j=0}^T \mathbf{P}_{ij}^0 \mathbf{V}_{jk}^0 \\ \mathbf{O}_{ik}^1 = \sum_{j=T}^{2T} \mathbf{P}_{ij}^1 \mathbf{V}_{jk}^1 \end{cases}$$

- $$\mathbf{O}_{ik}^{[0,1]} = \sum_{j=0}^{2T} \mathbf{P}_{ij}^{[0,1]} \mathbf{V}_{jk}^{[0,1]} = \sum_{j=0}^T \mathbf{P}_{ij}^{[0,1]} \mathbf{V}_{jk}^0 + \sum_{j=T}^{2T} \mathbf{P}_{ij}^{[0,1]} \mathbf{V}_{jk}^1$$

- $$= \sum_{j=0}^T \frac{\exp(\mathbf{s}_{ij}^0 - M_i^{[0,1]})}{L_i^{[0,1]}} \mathbf{V}_{jk}^0 + \sum_{j=T}^{2T} \frac{\exp(\mathbf{s}_{ij}^1 - M_i^{[0,1]})}{L_i^{[0,1]}} \mathbf{V}_{jk}^1$$

- $$= \frac{\exp(M_i^{[0,1]})}{\exp(M_i^0)} \frac{L_i^0}{L_i^{[0,1]}} \sum_{n=0}^T \frac{\exp(\mathbf{s}_{ij}^0 - M_i^0)}{L_i^0} \mathbf{V}_{jk}^0 +$$

$$\frac{\exp(M_i^{[0,1]})}{\exp(M_i^1)} \frac{L_i^1}{L_i^{[0,1]}} \sum_{n=T}^{2T} \frac{\exp(\mathbf{s}_{ij}^1 - M_i^1)}{L_i^1} \mathbf{V}_{jk}^1$$

- $$= \frac{\exp(M_i^{[0,1]})}{\exp(M_i^0)} \frac{L_i^0}{L_i^{[0,1]}} \mathbf{O}_{ik}^0 + \frac{\exp(M_i^{[0,1]})}{\exp(M_i^1)} \frac{L_i^1}{L_i^{[0,1]}} \mathbf{O}_{ik}^1$$

- $$:= \alpha^0 \beta^0 \mathbf{O}_{ik}^0 + \alpha^1 \beta^1 \mathbf{O}_{ik}^1$$

How to Fuse: Forward

- Tile \mathbf{Q}, \mathbf{O} into $\frac{N_Q}{T_Q}$ blocks $\mathbf{Q}_i, \mathbf{O}_i$; Tile \mathbf{K}, \mathbf{V} into $\frac{N_K}{T_K}$ blocks $\mathbf{K}_j, \mathbf{V}_j$
- Parallel for $1 \leq i \leq \frac{N_Q}{T_Q}$
 - in Registers: $\mathbf{Q}_i \leftarrow (0)_{\text{float32}}^{T_Q \times D}, M_i \leftarrow (-\infty)_{\text{float32}}^{T_Q}, L_i \leftarrow (0)_{\text{float32}}^{T_Q}$
 - For $1 \leq j \leq \frac{N_K}{T_K}$
 - $\mathbf{s}_{ij} \leftarrow \frac{\mathbf{Q}_i \mathbf{K}_j^T}{\sqrt{D}}, M_i^{\text{local}} \leftarrow \max_{-1}(\mathbf{s}_{ij}), L_i^{\text{local}} \leftarrow \text{sum}_{-1}(\exp(\mathbf{s}_{ij} - M_i^{\text{local}}))$
 - Update (M_i, L_i) and calc (α, β) given $(M_i, M_i^{\text{local}}, L_i, L_i^{\text{local}})$
 - $\mathbf{P}_{ij} \leftarrow \frac{\exp(\mathbf{s}_{ij} - M_i)}{L_i}, \mathbf{O}_i \leftarrow \alpha \beta \mathbf{O}_i + \mathbf{P}_{ij} \mathbf{V}_j$
 - Save \mathbf{O}_i to GPU HBM

How to Fuse: Backward

- $\frac{\partial \phi}{\partial \mathbf{V}} = \mathbf{P}^T \cdot \frac{\partial \phi}{\partial \mathbf{O}}, \quad \frac{\partial \phi}{\partial \mathbf{P}} = \frac{\partial \phi}{\partial \mathbf{O}} \cdot \mathbf{V}^T$
- $\frac{\partial \phi}{\partial \mathbf{S}_i} = \frac{\partial \phi}{\partial \mathbf{P}_i} (\text{diag}(\mathbf{P}_i) - \mathbf{P}_i^T \mathbf{P}_i) = \frac{\partial \phi}{\partial \mathbf{P}_i} \circ \mathbf{P}_i - \left(\frac{\partial \phi}{\partial \mathbf{P}_i} \cdot \mathbf{P}_i^T \right) \mathbf{P}_i$
- $\frac{\partial \phi}{\partial \mathbf{P}_i} \cdot \mathbf{P}_i^T = \left(\frac{\partial \phi}{\partial \mathbf{O}_i} \cdot \mathbf{V}^T \right) \cdot \mathbf{P}_i^T = \frac{\partial \phi}{\partial \mathbf{O}_i} \cdot \mathbf{O}_i^T = \sum_{k=0}^D \frac{\partial \phi}{\partial \mathbf{O}_{ik}} \mathbf{O}_{ik} := \Delta_i$
- $\frac{\partial \phi}{\partial \mathbf{S}} = \left(\frac{\partial \phi}{\partial \mathbf{P}} - \Delta \right) \circ \mathbf{P}$
- $\frac{\partial \phi}{\partial \mathbf{Q}} = \frac{\partial \phi}{\partial \mathbf{S}} \cdot \mathbf{K}, \quad \frac{\partial \phi}{\partial \mathbf{K}} = \frac{\partial \phi}{\partial \mathbf{S}^T} \cdot \mathbf{Q}$

Flash Attention: Tips

- Based on small head size (to store $\mathbf{Q}^i, \mathbf{K}^j, \mathbf{V}^j, \mathbf{O}^i$ in SRAM)
- The two for-loops are swapable
- Training: save M and L for backward
- Backward is much slower (recompute & limited by SRAM)
- Skip masked blocks

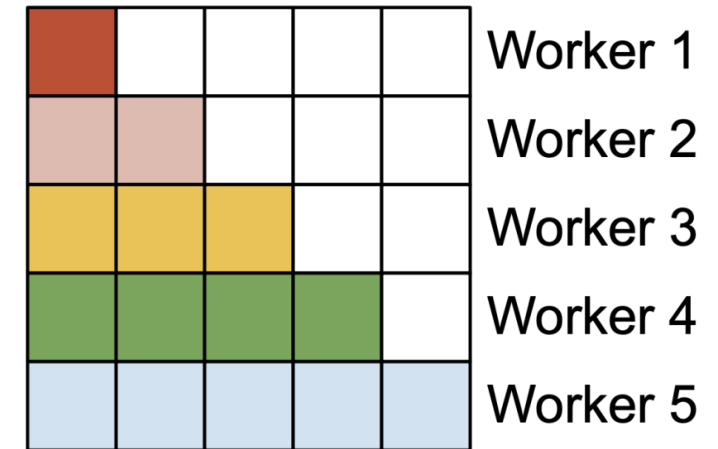
Flash Attention Implementations

- Official: [Dao-AI/flash-attention](#)
- Easy to modify: [openai/triton](#)
- For C++ template lovers: [nvidia/cutlass](#)
- Compatibility: [facebookresearch/xformers](#)

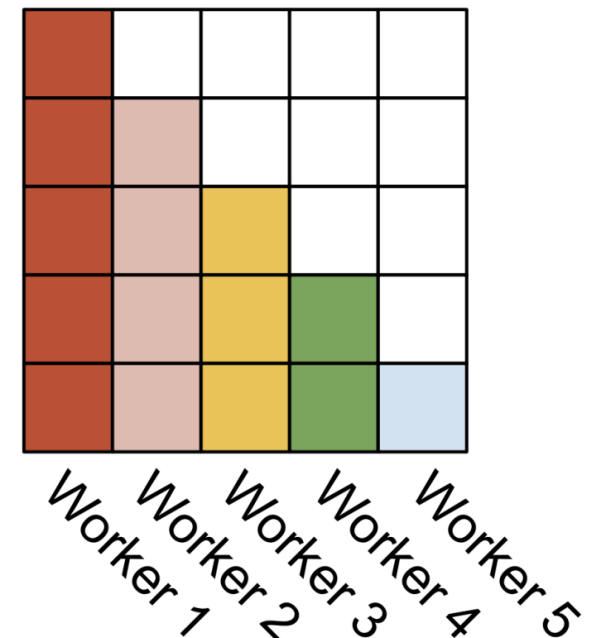
Flash Attention 2

- Forward:
 - $\mathbf{O}_i \leftarrow \alpha \mathbf{O}_i + \exp(\mathbf{S}_{ij} - M_i) \mathbf{V}_j; \mathbf{O}_i \leftarrow \frac{\mathbf{O}_i}{L_i}$
 - Save $L' = M + \log(L)$ for recompute
- Backward:
 - $\mathbf{P} = \frac{\exp(\mathbf{S} - M)}{L} = \exp(\mathbf{S} - L')$
- Long Context Generation:
 - Partition on K/V caches
 - Another kernel to combine results
- Parallelism Workload Balance

Forward pass



Backward pass



Reading Materials

- Flash Attention: [\[2205.14135\] FlashAttention: Fast and Memory-Efficient Exact Attention with IO-Awareness](#)
- Flash Attention 2: [\[2307.08691\] FlashAttention-2: Faster Attention with Better Parallelism and Work Partitioning](#)
- Flash Decoding: [Flash-Decoding for long-context inference – PyTorch](#)
- Flash Attention 3: [\[2407.08608\] FlashAttention-3: Fast and Accurate Attention with Asynchrony and Low-precision](#)