



Flash Attention

Complexity of dot-production attention

• Original:

- Memory $O(N^2)$

- Computation complexity: $O(N^3)$

\rightarrow Quadratic

Difficulty in increasing the context length

• Approximate Attention methods: $O(N)$ \rightarrow Linear [Please REFER TO Sparse ATTN]

But often they are not able to achieve wall-clock speedup?

• Reducing FLOP operations, at the cost of memory access overhead

▷ key idea: Avoid reading and writing the full attention matrix to and from GPU's DRAM

- Using a tiling mechanism
- I/O-aware. Significantly fewer memory access
- Exact attention (not approximate)
- Wall-clock speedup at training.

challenges:

Computing the softmax step without accessing the whole QK^T

Solution - Tiling mechanism, gradually computation

challenges:

Backward-pass (gradient) into storing the attention matrix

Solution - Only store the normalization factor of softmax

- Recomputing attention on-chip (Faster than reading the attention matrix from memory)

Performance of computing Function $y = f(x)$

• Time for memory access: T_{mem}

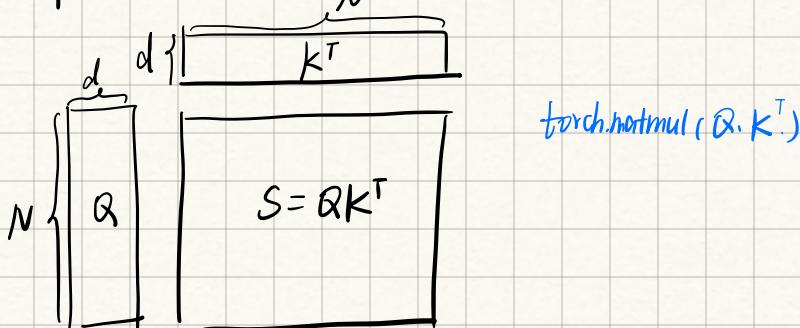
• Time for mathematical computations: T_{math}

\rightarrow Total time: $T_f = \max(T_{mem}, T_{math})$

| if $T_{math} > T_{mem}$ \rightarrow Compute-bound

| if $T_{math} < T_{mem}$ \rightarrow memory-bound

Dot-product Attention:



GPU Execution for Attention

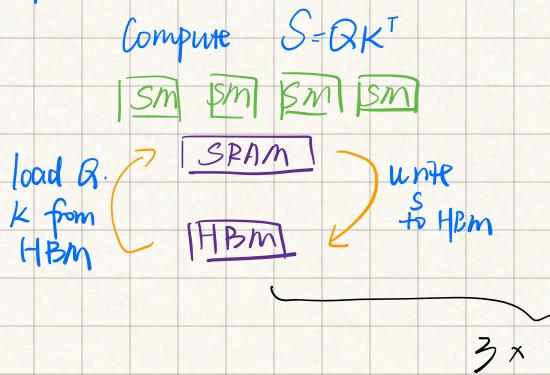
$S = \text{torch.matmul}(Q, K, \text{transpose}(-2, -1))$

$P = \text{torch.softmax}(S, \text{dim}=-1)$

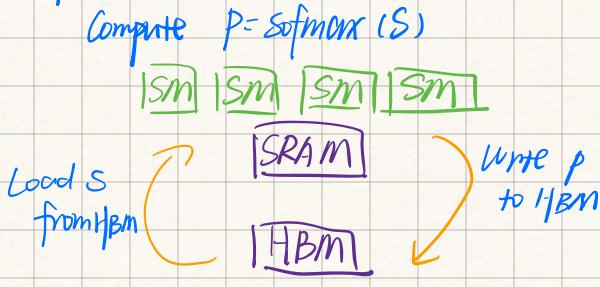
$\text{Output} = \text{torch.matmul}(P, V)$



Step 1:



Step 2:



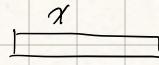
Flash Attention Algorithm

- Tiling (for the forward pass)
 - Loading blocks of Q, K, V and compute the output partially
- Re-computation (for the backward pass)
 - Using the stored normalization factors for each row

Reformulating Softmax Function

$$\text{softmax}(X)_i = \frac{\exp(X_i)}{\sum_k \exp(X_k)} \propto \frac{\exp(-m)}{\sum_k \exp(-m)}$$

To numerical stability, choose $m = \max(X)$. ensure all $X_i - m < 0, \therefore e^{X_i - m} \leq 1$



Re-formulating softmax:

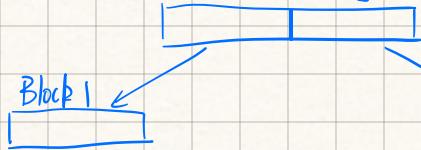
$$\vec{f}(x) = [e^{x_1-m}, e^{x_2-m}, \dots, e^{x_n-m}]$$

$$l(x) = \sum_j f(x)_j$$

So, how to compute softmax of entire vector x ?

Solve:

$$\vec{x} = [\vec{x}^{111}, \vec{x}^{122}]$$



$$m_1 = \max(\vec{x}^{111})$$

$$\vec{f}_1 = e^{\vec{x}^{111}-m_1}$$

$$l_1 = \sum_j e^{x_j^{111}-m_1}$$

$$m_2 = \max(\vec{x}^{122})$$

$$\vec{f}_2 = e^{\vec{x}^{122}-m_2}$$

$$l_2 = \sum_j e^{x_j^{122}-m_2}$$

softmax of TWO concatenated vectors

$$\vec{x} = [\vec{x}^{111}, \vec{x}^{122}]$$

$$\begin{cases} m_1 = \max(\vec{x}^{111}) \\ l_1 = \sum_j e^{x_j^{111}-m_1} \\ \vec{f}_1 = e^{\vec{x}^{111}-m_1} \end{cases}$$

$$\begin{cases} m_2 = \max(\vec{x}^{122}) \\ l_2 = \sum_j e^{x_j^{122}-m_2} \\ \vec{f}_2 = e^{\vec{x}^{122}-m_2} \end{cases}$$

Overall max : $m(x) = \max(m_1, m_2)$

Adjusted \vec{f}_1 and \vec{f}_2 :

$$\vec{f}(x) = [e^{m_1-m(x)} f_1, e^{m_2-m(x)} f_2]$$

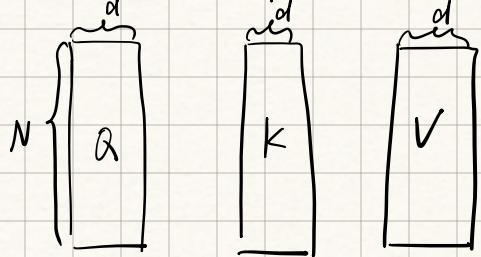
Adjusted normalization factor :

$$l(x) = e^{m_1-m(x)} l_1 + e^{m_2-m(x)} l_2$$

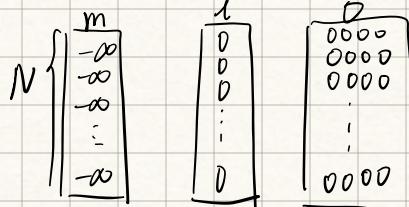
$$\Rightarrow \text{softmax}(\vec{x}) = \frac{\vec{f}(x)}{l(x)}$$



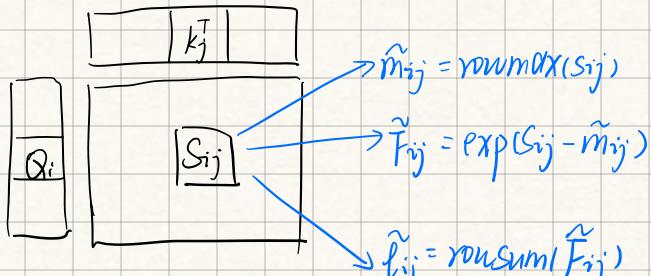
⚡ FlashAttention - Input and Initialization



initialize intermediate vectors and the outputs:

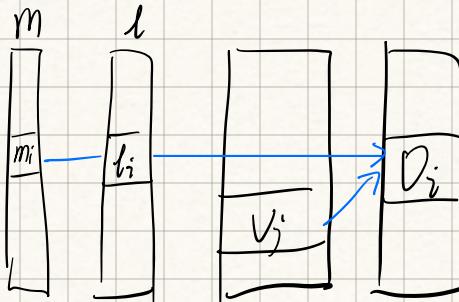


⚡ FlashAttention - Partial Updates

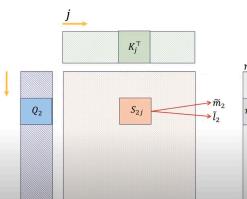


Partial update for block matrix D_i :
 $D_i^{\text{new}} = \text{diag}(l_i^{\text{new}})^{-1} (\text{diag}(l_i) e^{\hat{m}_i - m_i^{\text{new}}} D_i + e^{\hat{m}_i - m_i^{\text{new}}} \hat{F}_{ij} \hat{V}_j)$

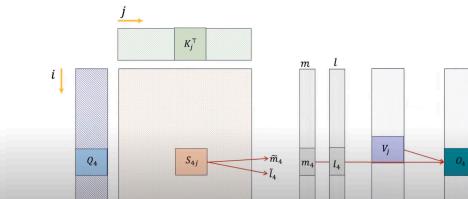
$D_i = \text{softmax}(Q_i K_j^T) V$, $S_{ij} = Q_i K_j^T$, $\hat{m}_{ij} = \text{softmax}(S_{ij})$
 $F_{ij} = e^{S_{ij} - \hat{m}_{ij}}$)



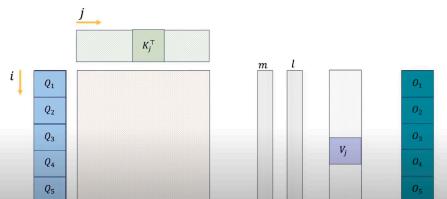
FlashAttention - Partial Updates



FlashAttention - Partial Updates



FlashAttention - Partial Updates



⚡ HBM Access Comparison

Naive Algorithm
 $O(Nd^2 + N^3)$

\Rightarrow

Flash Attn
 $O(Nd^2 M^{-1})$
M: Size of SRAM

\Rightarrow Flash Attn reduces number of HBM access

Large bm slow

High Bandwidth Memory (HBM): External GPU memory - large capacity, but slow access
Static RAM (SRAM)/on-chip Memory: Fast on-chip cache (shared memory) located inside the GPU core, but small but fast

N : sequence length eg 1024 tokens

d : Hidden dimension

Reference: Machine Learning Studio - Flash Attn: Accelerate LLM training