

BLOCK MATRIX MULTIPLICATION

$$\begin{array}{ccc}
 \boxed{A} & \times & \boxed{B} = \boxed{C} \\
 (M, K) & & (K, N) \quad (M, N)
 \end{array}$$

$$\begin{array}{ccc}
 \begin{array}{|c|c|} \hline A_{11} & A_{12} \\ \hline A_{21} & A_{22} \\ \hline \end{array} & \times & \begin{array}{|c|c|c|c|} \hline B_{11} & B_{12} & B_{13} & B_{14} \\ \hline B_{21} & B_{22} & B_{23} & B_{24} \\ \hline \end{array} \\
 \text{ORIGINAL}(8, 4) & & \text{ORIGINAL}(4, 8) \\
 \text{BLOCK}(2, 2) & & \text{BLOCK}(2, 4)
 \end{array}
 =
 \begin{array}{|c|c|c|c|} \hline A_{11}B_{11} + A_{12}B_{21} & A_{11}B_{12} + A_{12}B_{22} & A_{11}B_{13} + A_{12}B_{23} & A_{11}B_{14} + A_{12}B_{24} \\ \hline A_{21}B_{11} + A_{22}B_{21} & A_{21}B_{12} + A_{22}B_{22} & A_{21}B_{13} + A_{22}B_{23} & A_{21}B_{14} + A_{22}B_{24} \\ \hline \end{array}$$

$(4,2) \times (2,2) = (4,2)$
 $A_{11} \times B_{12}$

$\text{ORIGINAL}(8, 8)$
 $\text{BLOCK}(2, 4)$

$$\begin{array}{|c|c|} \hline A_{11} & A_{12} \\ \hline A_{21} & A_{22} \\ \hline \end{array}
 \times
 \begin{array}{|c|c|c|c|} \hline B_{11} & B_{12} & B_{13} & B_{14} \\ \hline B_{21} & B_{22} & B_{23} & B_{24} \\ \hline \end{array}
 =
 \begin{array}{|c|c|c|c|} \hline A_{11}B_{11} + A_{12}B_{21} & & & \\ \hline & & & \\ \hline & & & \\ \hline & & & \\ \hline \end{array}$$

\blacksquare = BLOCK MATRIX
 \blacksquare = ORIGINAL MATRIX

WHY SHOULD WE CARE?

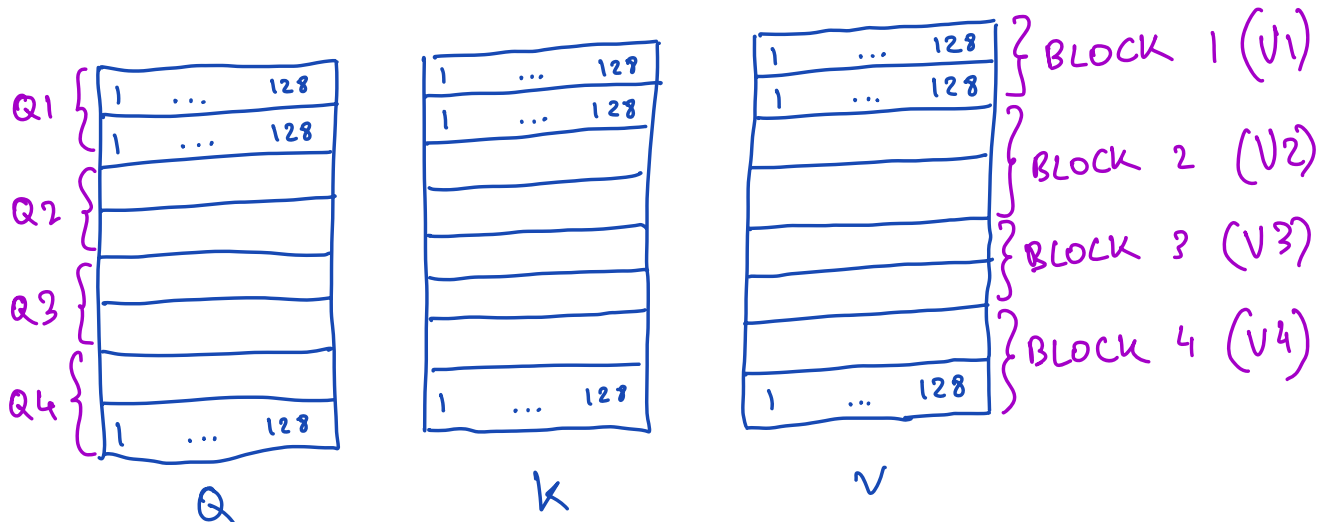
$$S = QK^T \in \mathbb{R}^{N \times N}, \quad \text{Let's ignore the softmax for a moment...} \quad O = \frac{S}{R}V \in \mathbb{R}^{N \times d}$$

Suppose for a moment that we want to do the following operation

$$O = (QK^T)V \in \mathbb{R}^{N \times d}$$

We know that $Q, K, V \in \mathbb{R}^{N \times d}$

■ = BLOCK
■ = ORIGINAL MATRIX



ORIGINAL = (8, 128)

BLOCK = (4, 128)

$$Q = \begin{bmatrix} Q_1 \\ Q_2 \\ Q_3 \\ Q_4 \end{bmatrix} = \begin{bmatrix} 1 & \dots & 128 \\ 1 & \dots & 128 \\ & & \\ & & \\ & & \\ & & \\ 1 & \dots & 128 \end{bmatrix}$$

ORIGINAL = (8, 128)
BLOCK = (4, 1)

$$K^T = \begin{bmatrix} k_1^T & k_2^T & k_3^T & k_4^T \\ - & & & - \\ \vdots & & & \vdots \\ 128 & & & 128 \end{bmatrix}$$

ORIGINAL = (128, 8)
BLOCK = (1, 4)

$$Q = \begin{bmatrix} Q_1 \\ Q_2 \\ Q_3 \\ Q_4 \end{bmatrix} = (4, 1)$$

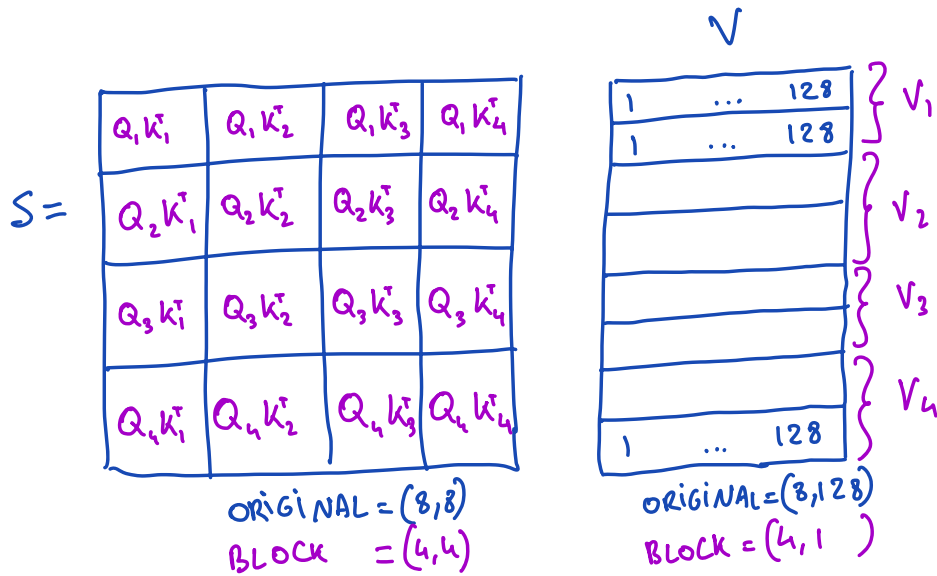
$$K^T = \begin{bmatrix} k_1^T & k_2^T & k_3^T & k_4^T \end{bmatrix} = (1, 4)$$

$(2, 128) \times (128, 2) = (2, 2)$

$$S = \begin{bmatrix} Q_1 k_1^T & Q_1 k_2^T & Q_1 k_3^T & Q_1 k_4^T \\ Q_2 k_1^T & Q_2 k_2^T & Q_2 k_3^T & Q_2 k_4^T \\ Q_3 k_1^T & Q_3 k_2^T & Q_3 k_3^T & Q_3 k_4^T \\ Q_4 k_1^T & Q_4 k_2^T & Q_4 k_3^T & Q_4 k_4^T \end{bmatrix}$$

ORIGINAL = (8, 8)
BLOCK = (4, 4)

Let's multiply by V



$O =$

$$\begin{aligned} & (Q_1 K_1^T) V_1 + (Q_1 K_2^T) V_2 + (Q_1 K_3^T) V_3 + (Q_1 K_4^T) V_4 \\ & (Q_2 K_1^T) V_1 + (Q_2 K_2^T) V_2 + (Q_2 K_3^T) V_3 + (Q_2 K_4^T) V_4 \\ & (Q_3 K_1^T) V_1 + (Q_3 K_2^T) V_2 + (Q_3 K_3^T) V_3 + (Q_3 K_4^T) V_4 \\ & (Q_4 K_1^T) V_1 + (Q_4 K_2^T) V_2 + (Q_4 K_3^T) V_3 + (Q_4 K_4^T) V_4 \end{aligned}$$

$\left[(2 \times 128) \times (128 \times 2) \right] \times (2 \times 128) = (2 \times 2) \times (2 \times 128) = (2 \times 128)$

EACH BLOCK OF THE OUTPUT MATRIX O IS ACTUALLY MADE UP OF TWO ROWS!

PSEUDOCODE

```

FOR EACH BLOCK  $Q_i$ :
   $O_i = \text{ZEROS}(2, 128)$  // OUTPUT IS INITIALLY ZEROS
  FOR EACH BLOCK  $K_j$ :
     $O_i \leftarrow O_i + (Q_i K_j^T) V_j$ 
  END FOR
END FOR
  
```

WAIT... What happened to the
SOFT MAX?

Let's restore it... with a twist!

$$S = \begin{matrix} & \begin{matrix} Q_1 k_1^T & Q_1 k_2^T & Q_1 k_3^T & Q_1 k_4^T \end{matrix} \\ \begin{matrix} Q_2 k_1^T \\ Q_3 k_1^T \\ Q_4 k_1^T \end{matrix} & \begin{matrix} Q_2 k_2^T \\ Q_3 k_2^T \\ Q_4 k_2^T \end{matrix} & \begin{matrix} Q_2 k_3^T \\ Q_3 k_3^T \\ Q_4 k_3^T \end{matrix} & \begin{matrix} Q_2 k_4^T \\ Q_3 k_4^T \\ Q_4 k_4^T \end{matrix} \end{matrix}$$

ORIGINAL = (8,8)
BLOCK = (4,4)

SOFT MAX*
 \Rightarrow

$$P = \begin{matrix} & \begin{matrix} P_{11} & P_{12} & P_{13} & P_{14} \end{matrix} \\ \begin{matrix} P_{21} \\ P_{31} \\ P_{41} \end{matrix} & \begin{matrix} P_{22} \\ P_{32} \\ P_{42} \end{matrix} & \begin{matrix} P_{23} \\ P_{33} \\ P_{43} \end{matrix} & \begin{matrix} P_{24} \\ P_{34} \\ P_{44} \end{matrix} \end{matrix}$$

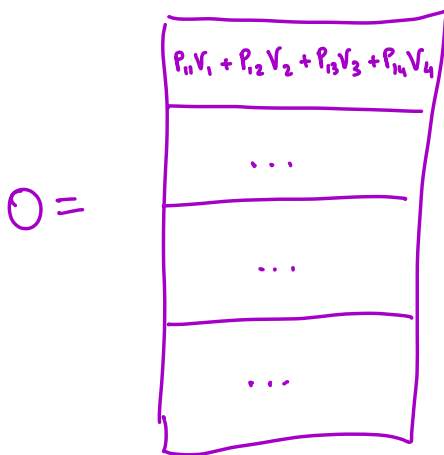
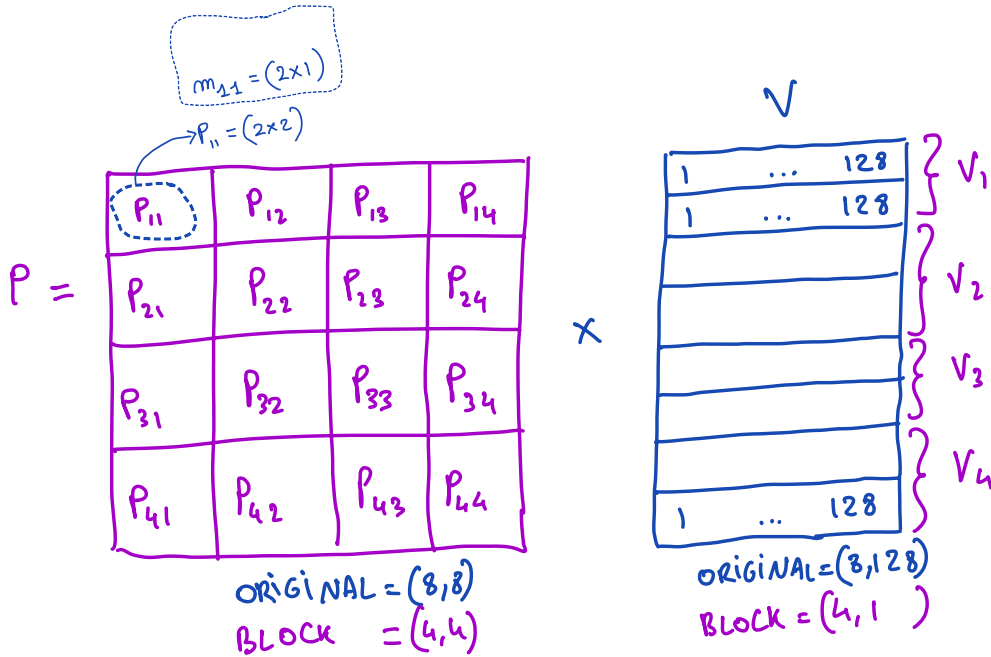
ORIGINAL = (8,8)
BLOCK = (4,4)

$$\text{SOFT MAX}^* (S_{ij}) = \exp [S_{ij} - \text{row max} (S_{ij})]$$

* Note: $P_{11} = \text{soft max}^* (Q_1 k_1^T)$
 $P_{12} = \text{soft max}^* (Q_1 k_2^T)$
 etc...

This is a 2x2 matrix

Let's multiply by V



WRONG!

EACH OF THE P_{ij} BLOCK HAS BEEN INDEPENDENTLY CALCULATED, SO THE MAX ELEMENT FOR EACH ROW IN EACH BLOCK IS NOT THE GLOBAL FOR EACH ROW, BUT THE ONE LOCAL TO EACH BLOCK

PSEUDOCODE

```

FOR EACH BLOCK  $Q_i$ :
   $O_i = \text{ZEROS}(2, 128)$  // OUTPUT IS INITIALLY ZEROS
  FOR EACH BLOCK  $V_j$ :
     $P_{ij} = \text{softmax}(Q_i \cdot V_j)$ 
     $O_i \leftarrow O_i + P_{ij} \cdot V_j$ 
  END FOR
END FOR

```

HOW CAN WE FIX THE PREVIOUS ITERATION'S OUTPUT?

IF ONLY WE HAD A WAY TO FIX THE SOFTMAX...

THE ONLINE SOFTMAX

$$m_0 = -\infty$$

$$l_0 = 0$$

for $i = 1$ to N

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

$$l_i = l_{i-1} \cdot e^{m_{i-1} - m_i} + e^{x_i - m_i}$$

for $k = 1$ to N

$$x_k \leftarrow \frac{e^{x_k - m_N}}{l_N}$$

THE IDEA

IF WE CAN "FIX" THE SOFTMAX WHILE ITERATING ON A ROW, WE CAN ALSO FIX BLOCKS OF ROWS, SINCE THE SOFTMAX IS APPLIED INDEPENDENTLY TO EACH ROW

Let's see how to apply the online softmax here...

$$O_i = P_{i1}V_1 + P_{i2}V_2 + P_{i3}V_3 + P_{i4}V_4 \quad \Leftarrow \text{WE NEED TO FIX THESE}$$



PSEUDO CODE

```
FOR EACH BLOCK Q:  
  O_i = ZEROS(2, 128) // OUTPUT IS INITIALLY ZEROS  
  FOR EACH BLOCK KJ  
    P_ij = softmax*(Q_i * K_j)  
    O_i ← O_i + P_ij * V_j  
  END FOR  
END FOR
```

INITIALIZATION

$$m_0 = \begin{bmatrix} -\infty \\ -\infty \end{bmatrix}$$

$$l_0 = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

$$O_0 = \begin{bmatrix} 0 & 0 & 0 & \dots & 0 & 0 \\ 0 & 0 & 0 & \dots & 0 & 0 \end{bmatrix} \quad 2 \times 128 \text{ matrix}$$

STEP 1

$$m_1 = \max(\text{rowmax}(Q_1 k_1^T), m_0)$$

$$S_1 = Q_1 k_1^T$$

$$l_1 = \text{rowsum}[\exp(S_1 - m_1)] + l_0 \cdot \exp(m_0 - m_1)$$

$$p_{11} = \exp(S_1 - m_1)$$

$$O_1 = \text{diag}(\exp(m_0 - m_1)) O_0 + p_{11} V_1$$

As you can see, we never normalize the "softmax" values. We will do it at the end.

STEP 2

$$m_2 = \max(\text{rowmax}(Q_2 k_2^T), m_1)$$

$$S_2 = Q_2 k_2^T$$

$$l_2 = \text{rowsum}[\exp(S_2 - m_2)] + l_1 \cdot \exp(m_1 - m_2)$$

$$p_{12} = \exp(S_2 - m_2)$$

$$O_2 = \text{diag}(\exp(m_1 - m_2)) O_1 + p_{12} V_2$$

AND SO ON UNTIL THE LAST STEP. THEN, WE APPLY THE "l" NORMALIZATION FACTOR.

STEP 5

$$O_5 = [\text{diag}(l_4)]^{-1} O_4$$

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** ← FOR EACH \mathbf{Q}_i BLOCK
 - 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** ← FOR EACH \mathbf{K}_j BLOCK
 - 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 .
-