

## Flosh Attention

Safe Softmax

Making Softmax Safe

$$S = QK^T \in \mathbb{R}^{N \times N}, P = \text{Softmax}(S) \in \mathbb{R}^{N \times N}, O = PV \in \mathbb{R}^{N \times d}$$

$Q \rightarrow$  Seg token  $(N, d)$

$K \rightarrow (N, d) \rightarrow (d, N)$

$V \rightarrow (N, d)$

$$QK^T = (N, N)$$

$$(N, d)$$

$q_1^T K_1$	$q_1^T K_2$	$q_1^T K_3$	$q_1^T K_4$	$q_1^T K_5$
.	.	.	.	.
.	.	.	.	.
.	.	.	.	.
$q_5^T K_1$	$q_5^T K_2$	$q_5^T K_3$	$q_5^T K_4$	$q_5^T K_5$

Softmax  $\rightarrow$

$(5, 5)$

0.1	0.05	0.5	0.15	0.2
0.3	0.1	0.35	0.2	0.05

$$\Rightarrow \sum = 1$$

$(5, 5)$

$$x = [3, 2, 5, 1]$$

$$1) X_{\max} = 5$$

$$2) e^{3-5} + e^{2-5} + e^{5-5} + e^{1-5}$$

$$e^{-2} + e^{-3} + e^0 + e^{-4} = 1$$

$$3) x_1 = \frac{e^{3-5}}{\lambda} \quad x_2 = \frac{e^{2-5}}{\lambda} \quad x_3 = \frac{e^{5-5}}{\lambda} \quad x_4 = \frac{e^{1-5}}{\lambda}$$

Sequential call "Slow"

$$m_0 = -\infty$$

$$m_1 = \max(-\infty, 3) = 3$$

$$m_2 = \max(3, 2) = 3$$

$$m_3 = \max(3, 5) = 5$$

$$m_4 = \max(5, 1) = 5$$

$$x_1 = \frac{e^{3-5}}{\lambda_4} \quad x_3 = \frac{e^{5-5}}{\lambda_4}$$

$$x_2 = \frac{e^{2-5}}{\lambda_4} \quad x_4 = \frac{e^{1-5}}{\lambda_4}$$

$$x = [3, 2, 5, 1]$$

Fusion:

Step 1:

$$\max_1 = 3$$

$$\lambda_1 = e^{3-3} = e^0$$

$\lambda_3$  computed is wrong!

$$\lambda_3 = e^{3-3} + e^{2-3} + e^{5-5}$$

Wrong as global max is 5

Step 2:

$$\max_2 = \max(3, 2)$$

$$\lambda_2 = \lambda_1 + e^{2-3}$$

correction factor

Fix!

$$\lambda_3 = \lambda_2 \cdot e^{3-5} + e^{5-5}$$

$$= (e^{3-3} + e^{2-3}) e^{3-5} + e^{5-5}$$

$$= e^{+3-3+3-5}$$

Step 3:

$$\max_3 = \max(3, 5) = 5$$

$$\lambda_3 = \lambda_2 + e^{5-5} = e^{3-3} + e^{2-3} + e^{5-5}$$

$$= e^{3-3+3-5} + e^{2-3+3-5} + e^{5-5} = e^{3-5} + e^{2-5} + e^{5-5}$$

current max

correct!

$$l_3 = l_2 \cdot e^{3-5} + e^{5-5} = (e^{3-3} + e^{2-3}) e^{3-5} + e^{5-5}$$

New Pseudocode

$$m_0 = -\infty$$

$$l_0 = 0$$

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  $N$

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

we want to prove  
that at the end of this loop

$$m_N = \max(x_i) = x_{\max}$$

$$l_N = \sum_{j=1}^N e^{x_j - x_{\max}}$$

Proof by Induction

1. Size  $N=1$

$$m_1 = \max(-\infty, x_1) = x_1 = \max(x_i) = x_{\max}$$

$$l_1 = 0 \times e^{-\infty} + e^{x_1-x_1} = \sum_{j=1}^N e^{x_j - x_{\max}}$$

*erroneous factor*

2. Size  $N+1$

$$m_{N+1} = \max(m_N, x_{N+1}) = \max(x_i)$$

$$l_{N+1} = l_N e^{m_N - m_{N+1}} + e^{x_{N+1} - m_{N+1}} =$$

$$= \left( \sum_{j=1}^N e^{x_j} \right)$$

2) If we assume it holds for a vector of size  $N$ , does it hold for a vector of size  $N+1$ ?

$$m_{N+1} = \max(m_N, x_{N+1}) = \max_i (x_i)$$

$$\begin{aligned} l_{N+1} &= l_N e^{m_N - m_{N+1}} + e^{\delta l_{N+1} - m_{N+1}} \\ &= \left( \sum_{j=1}^N e^{x_j - m_N} \right) e^{m_N - m_{N+1}} + e^{x_{N+1} - m_{N+1}} \\ &\quad \text{local max iteration factor} \qquad \text{current element} \qquad \text{current set of max (local max)} \\ &= \sum_{j=1}^N e^{x_j - m_{N+1}} + e^{x_{N+1} - m_{N+1}} \\ &= \sum_{j=1}^{N+1} e^{x_j - m_{N+1}} \end{aligned}$$

### BLOCK MATRIX MULTIPLICATION

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

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

Original (8, 9)  
Block (2, 2)

Original (4, 8)  
Block (2, 4)

$A_{11} B_{11}$	$A_{11} B_{12}$	$A_{11} B_{13}$	$A_{11} B_{14}$
$+ A_{12} B_{21}$	$+ A_{12} B_{21}$	$+ A_{12} B_{23}$	$+ A_{12} B_{24}$
$A_{21} B_{11}$	$A_{21} B_{12}$	$A_{21} B_{13}$	$A_{21} B_{14}$
$+ A_{22} B_{21}$	$+ A_{21} B_{22}$	$+ A_{22} B_{23}$	$+ A_{22} B_{24}$

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

$$Q, K, V \in R^{N \times d}$$

Q

$Q_1 \{$	1 - - - 128
$Q_2 \{$	
$Q_3 \{$	
$Q_4 \{$	1 - - - 128

K

1 - - - 128
(1 - - - 128)
1 - - - 128

(n, M)

(n, N)

(N, M)

128	816	818	113
128	816	818	113

(3, 8) Interpol.

(3, 8) divide

8A	12A
8A	12A

(3, 8) Interpol.  
(3, 8) divide

$$S_{11} \quad (2, 128) \times (128, 2) = (2, 2)$$

$Q_1 K_1^T$	$Q_1 K_2^T$	$Q_1 K_3^T$	$Q_1 K_4^T$
$Q_2 K_1^T$	$Q_2 K_2^T$	$Q_2 K_3^T$	$Q_2 K_4^T$
$Q_3 K_1^T$	$Q_3 K_2^T$	$Q_3 K_3^T$	$Q_3 K_4^T$
$Q_4 K_1^T$	$Q_4 K_2^T$	$Q_4 K_3^T$	$Q_4 K_4^T$

Softmax\*  
=>

$$\text{Original} = (8, 8)$$

$$\text{Block} = (4, 4)$$

$$\text{Softmax}(x_i) = \frac{\exp(x_i - \bar{x}_{\text{max}})}{\sum_{j=1}^n \exp(x_j - \bar{x}_{\text{max}})}$$

$$\text{Softmax}^*(x_i) = \exp(x_i - \bar{x}_{\text{max}})$$

$$S_{11} = \begin{array}{|c|c|} \hline \textcircled{a} & b \\ \hline c & \textcircled{d} \\ \hline \end{array} \xrightarrow{\text{Softmax}^*} \begin{array}{|c|c|} \hline \exp(a-d) & \exp(b-d) \\ \hline \exp(c-d) & \exp(d-d) \\ \hline \end{array}$$

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

$P_{11}$	$P_{12}$	$P_{13}$	$P_{14}$
$P_{21}$	$P_{22}$	$P_{23}$	$P_{24}$
$P_{31}$	$P_{32}$	$P_{33}$	$P_{34}$
$P_{41}$	$P_{42}$	$P_{43}$	$P_{44}$

$$\text{Original} = (8, 8)$$

$$\text{BLOCK} = (4, 4)$$

$$\text{SOFTMAX}^*(S_{ij}) = \exp \left[ S_{ij} - \text{row\_max}(S_{ij}) \right]$$

$$P_{11} = \text{Softmax}^*(Q_1 K_1^T)$$

$$P_{12} = \text{Softmax}^*(Q_1 K_2^T)$$

etc

(2, 2)

Let's multiply by V

$$2P_{11} = (2 \times 2) \quad m_{11} = (2 \times 1)$$

$P_{11}$	$P_{12}$	$P_{13}$	$P_{14}$
$P_{21}$	$P_{22}$	$P_{23}$	$P_{24}$
$P_{31}$	$P_{32}$	$P_{33}$	$P_{34}$
$P_{41}$	$P_{42}$	$P_{43}$	$P_{44}$

$\frac{1}{1}$	$\frac{2}{1}$	$\frac{3}{1}$	$\frac{2}{1}$
$1 \dots 128$	$\frac{2}{1}$	$\frac{3}{1}$	$\frac{2}{1}$
$1 \dots 128$	$\frac{2}{1}$	$\frac{3}{1}$	$\frac{2}{1}$
$1 \dots 128$	$\frac{2}{1}$	$\frac{3}{1}$	$\frac{2}{1}$
$1 \dots 128$	$\frac{2}{1}$	$\frac{3}{1}$	$\frac{2}{1}$

$$\text{Orginal} = (8, 8)$$

$$\text{BLOCK} = (4, 4)$$

$$(4, 2) \times (2, 128) = (2, 128)$$

$$P_{11}V_1 + P_{12}V_2 + P_{13}V_3 + P_{14}V_4$$

$$\theta =$$

$(0, 1)$	$(0, 1)$	$(0, 1)$	$(0, 1)$
$(0, 1)$	$(0, 1)$	$(0, 1)$	$(0, 1)$
$(0, 1)$	$(0, 1)$	$(0, 1)$	$(0, 1)$
$(0, 1)$	$(0, 1)$	$(0, 1)$	$(0, 1)$

## Normal Attention:

$Q K^T$

$P_{ii} \rightarrow$  Local maximum  
for softmax $^*(\text{line})$  and next  $\rightarrow$   $\alpha_i$

$Q_1^T K_1$	$Q_1^T K_2$	$Q_1^T K_3$	$Q_1^T K_4$	$Q_1^T K_5$
$Q_2^T K_1$	$Q_2^T K_2$	$Q_2^T K_3$	$Q_2^T K_4$	$Q_2^T K_5$
$Q_3^T K_1$	$Q_3^T K_2$	$Q_3^T K_3$	$Q_3^T K_4$	$Q_3^T K_5$
$Q_4^T K_1$	$Q_4^T K_2$	$Q_4^T K_3$	$Q_4^T K_4$	$Q_4^T K_5$
$Q_5^T K_1$	$Q_5^T K_2$	$Q_5^T K_3$	$Q_5^T K_4$	$Q_5^T K_5$

## PSEUDO CODE

FOR EACH BLOCK

$$O_i = \text{ZEROS}(2, 128)$$

FOR EACH BLOCK  $K_j$

$$P_{ij} = \text{softmax}^*(Q_i K_j^T)$$

$$O_i \leftarrow O_i + P_{ij} Y_j$$

END FOR

END FOR

## Fix! The Online Softmax

$$m_0 = -\infty$$

$$l_0 = 0$$

for  $i=1$  to  $N$

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

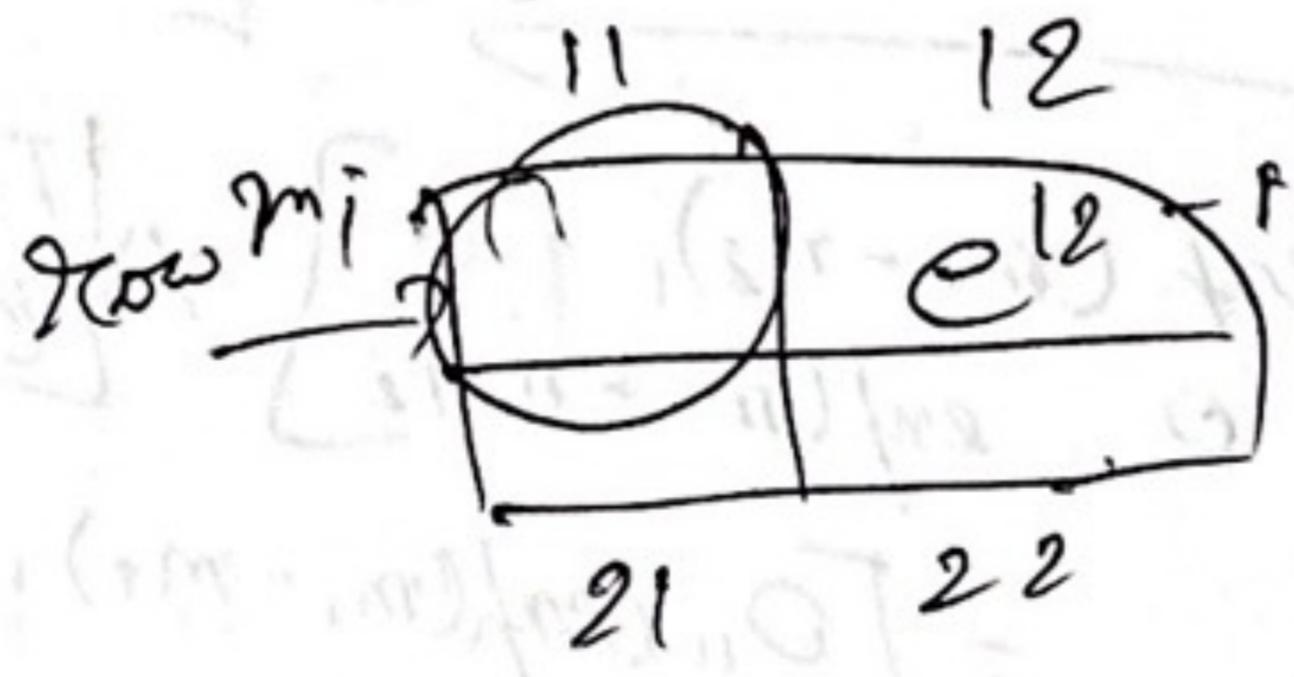
$$d_i = d_{i-1} \cdot e^{m_{i-1} - m_i} + e$$

for  $k=1$  to  $N$

$$x_k \leftarrow \frac{e^{x_k - m_N}}{\ln d_N}$$

## 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:

$$m_1 = \max(\text{row max}(Q_1 K_1^T), m_0)$$

$$S_1 = Q_1 K_1^T$$

$$\lambda_1 = \text{row sum} [\exp(S_1 - m_1)] + d_0 \cdot \exp(m_0 - m_1)$$

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

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



Step 1:

$$m_1 = \max(\text{row max}(Q_1 K_1^T), m_0)$$

$\downarrow$   
 $S_{11}$

$$S_1 = Q_1 K_1^T$$

$$\lambda_1 = \text{row sum} [\exp(S_1 - m_1)] + d_0 \cdot \exp(m_0 - m_1)$$

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

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

Step 2:

$$m_2 = \max(\text{row max}(Q_1 K_2^T), m_1)$$

$\downarrow$   
 $S_{12}$

$$S_2 = Q_1 K_2^T$$

$$\lambda_2 = \text{row sum}$$

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

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

$$O_1 = \begin{bmatrix} O_{11} & O_{12} & 0 & \dots & O_{1,128} \\ O_{21} & O_{22} & \dots & \dots & O_{2,128} \end{bmatrix}$$

X

$$\begin{bmatrix} \exp(m_1 - m_2), 0 \\ 0 & \exp(m_1 - m_2) \end{bmatrix}$$

X

$$\begin{bmatrix} \exp(m_1 - m_2), 0 \\ 0 & \exp(m_1 - m_2) \end{bmatrix} \times \begin{bmatrix} O_{11} & O_{12} & \dots & O_{1,128} \\ O_{21} & O_{22} & \dots & O_{2,128} \end{bmatrix}$$

$$= \begin{bmatrix} O_{11} \exp(m_1 - m_2), \\ \vdots \end{bmatrix}$$

Step 5:

$$O_5 = [\text{diag}(\lambda_4)]^{-1} O_4$$

$$\boxed{O_1, \exp(m_1 - m_2), + O(O_2)}$$

$$l_4 = \begin{bmatrix} l_4^{(1)} \\ l_4^{(2)} \end{bmatrix} \quad Q_A = \begin{bmatrix} C_1 \dots & 128 \\ C_2 \dots & 128 \end{bmatrix}$$

$$\left( \begin{bmatrix} l_4^{(1)} & 0 \\ 0 & l_4^{(2)} \end{bmatrix} \right)^{-1} = \begin{bmatrix} \frac{1}{l_4^{(1)}} & 0 \\ 0 & \frac{1}{l_4^{(2)}} \end{bmatrix} \times 0 \begin{bmatrix} C_1 \dots & 128 \\ C_2 \dots & 128 \end{bmatrix}_{(2, 128)}$$

$$= \boxed{\cancel{\frac{1}{l_4^{(1)}} O_{11} + O_{21} \frac{1}{l_4^{(2)}}}}$$

$$= \begin{bmatrix} \frac{1}{l_4^{(1)}} O_{11} + 0 O_{21}, & \frac{1}{l_4^{(2)}} \\ 0 \times O_{11} + \frac{1}{l_4^{(2)}} O_{21}, & \end{bmatrix}_{(2, 128)}$$

### "Flash Attention - 2 Feed Forward"

$Q, K, V \in \mathbb{R}^{N \times d}$  by HBM block, sizes  $B_C, B_R$

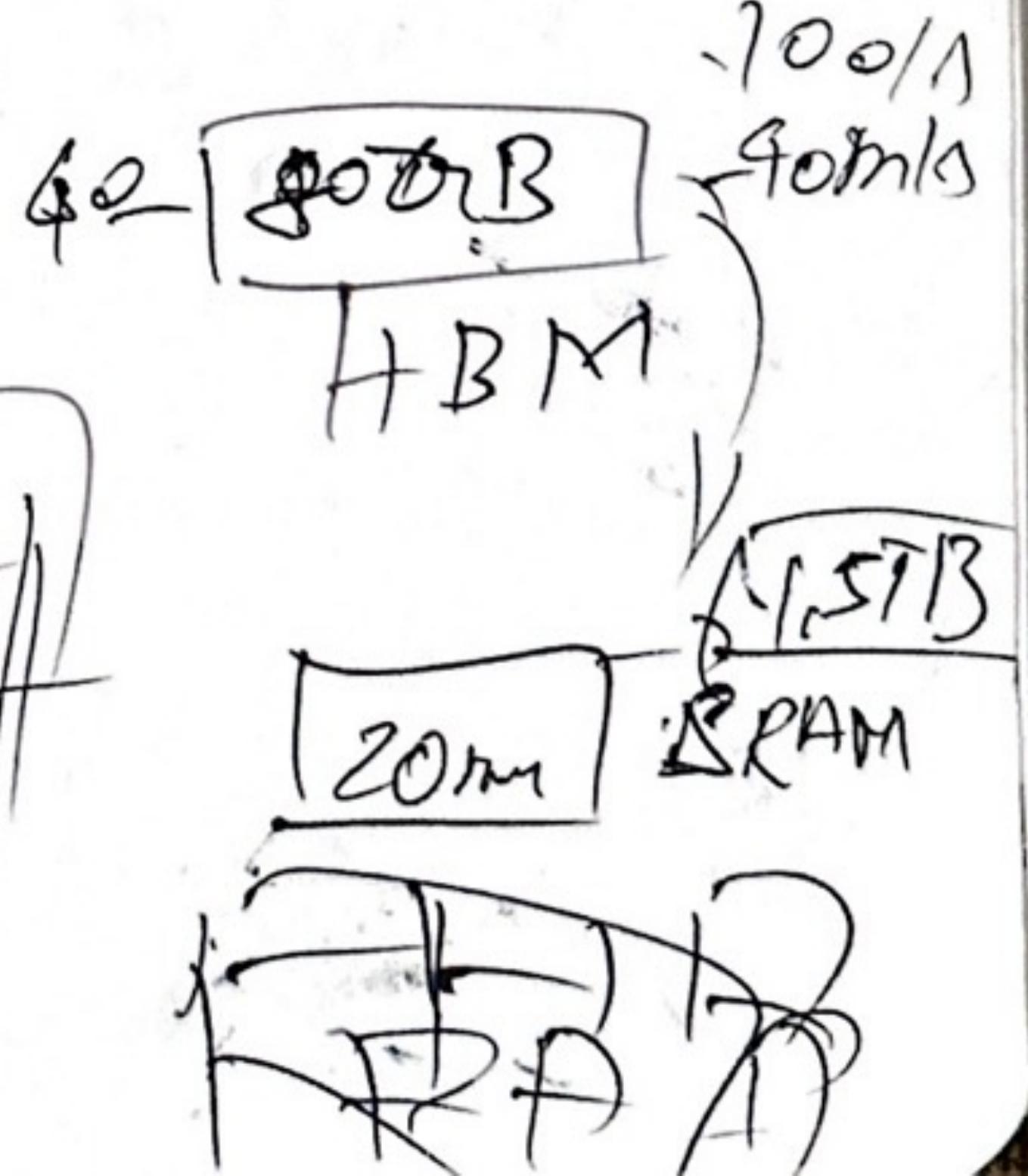
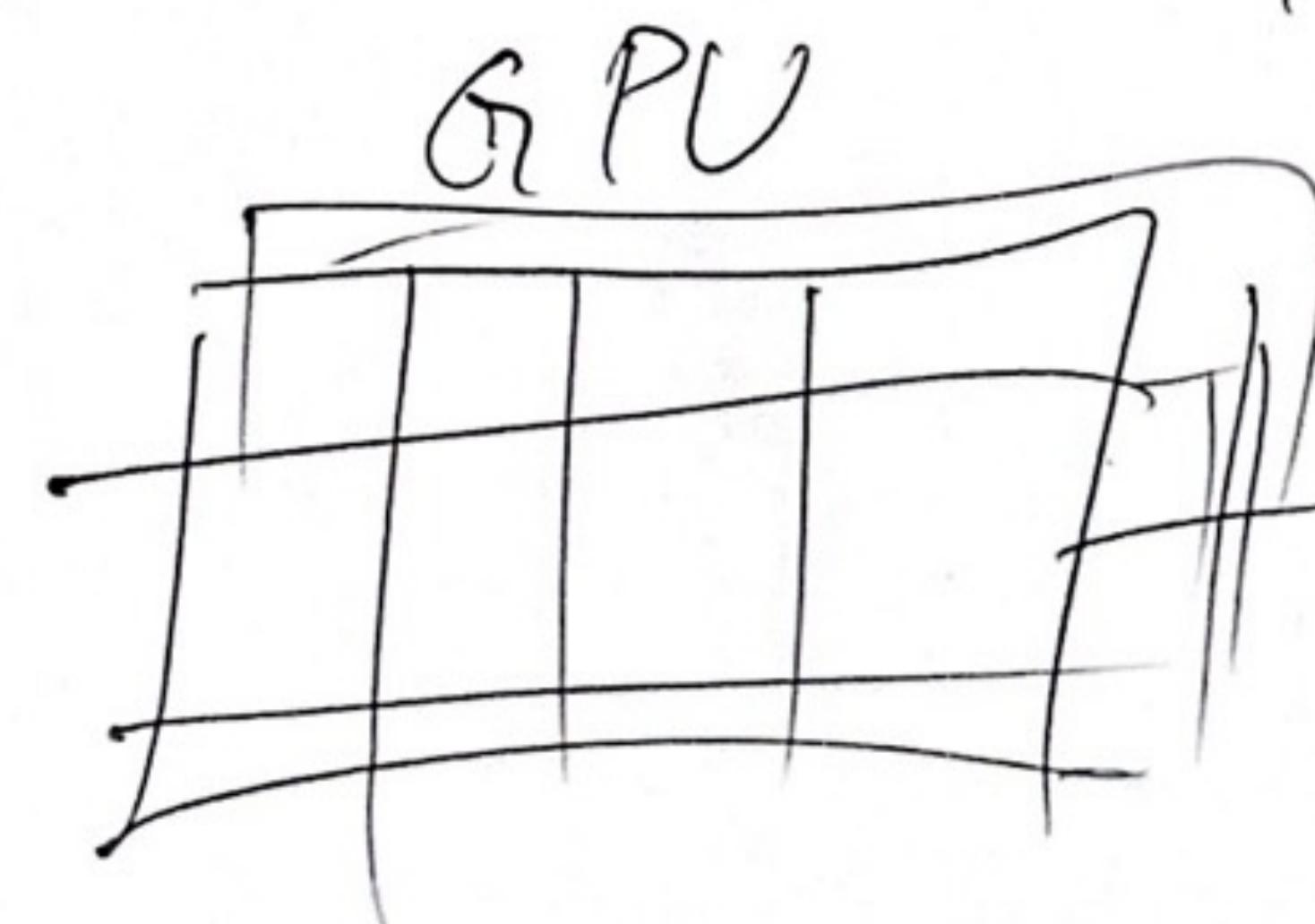
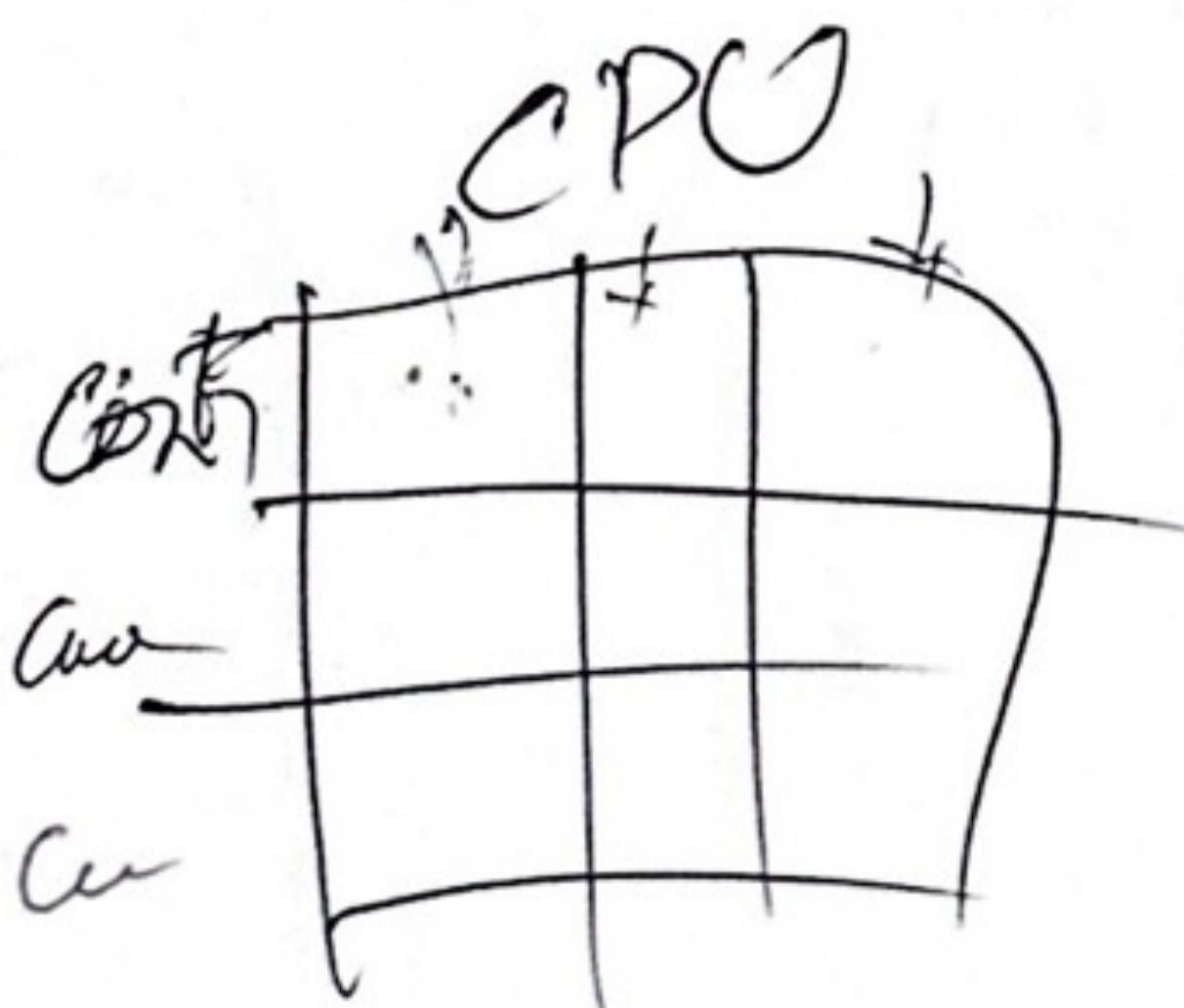
D,W  $Q$  into  $T_R = \left[ \frac{N_R}{B_R} \right]$  blocks  $Q_1, \dots, Q_{T_R}$  of size  $B_R \times d$

$K, V$  into  $T_C = \left[ \frac{N_R}{B_C} \right]$  blocks  $K_1, \dots, K_{T_C}$  and  $V_1, \dots, V_{T_C}, \dots$

B

A100

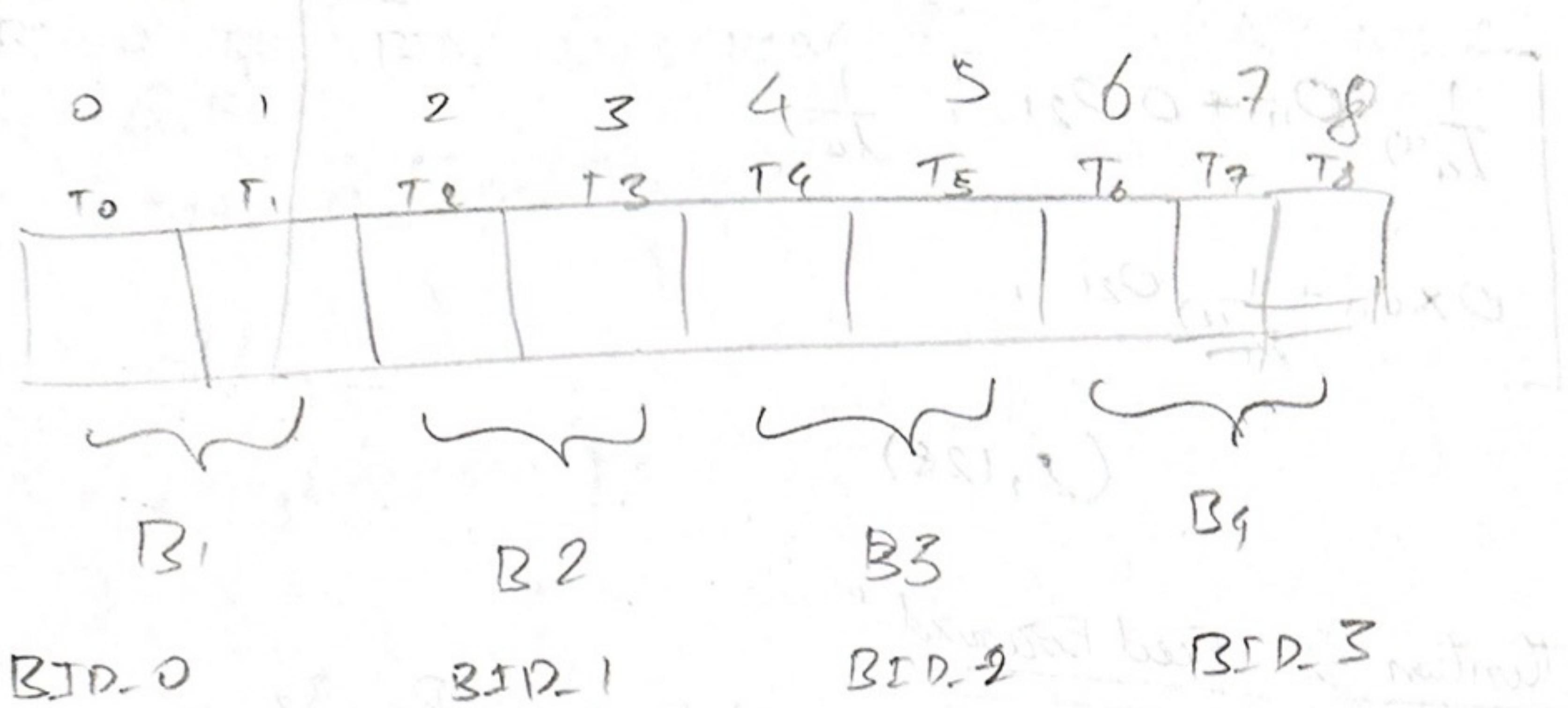
for  $1 \leq i \leq T_R : d_w$



## Vector Addition

	$T_1$	$T_2$	$\dots$	$T_7$	$T_8$							
$A =$	1	4	1	3	1	5	1	6	1	7	1	2
$B =$	3	1	2	6	1	2	5	6	9			
$C =$	4	6	9	10	7	16	18	11				

## GPU Block and Thread Allocation

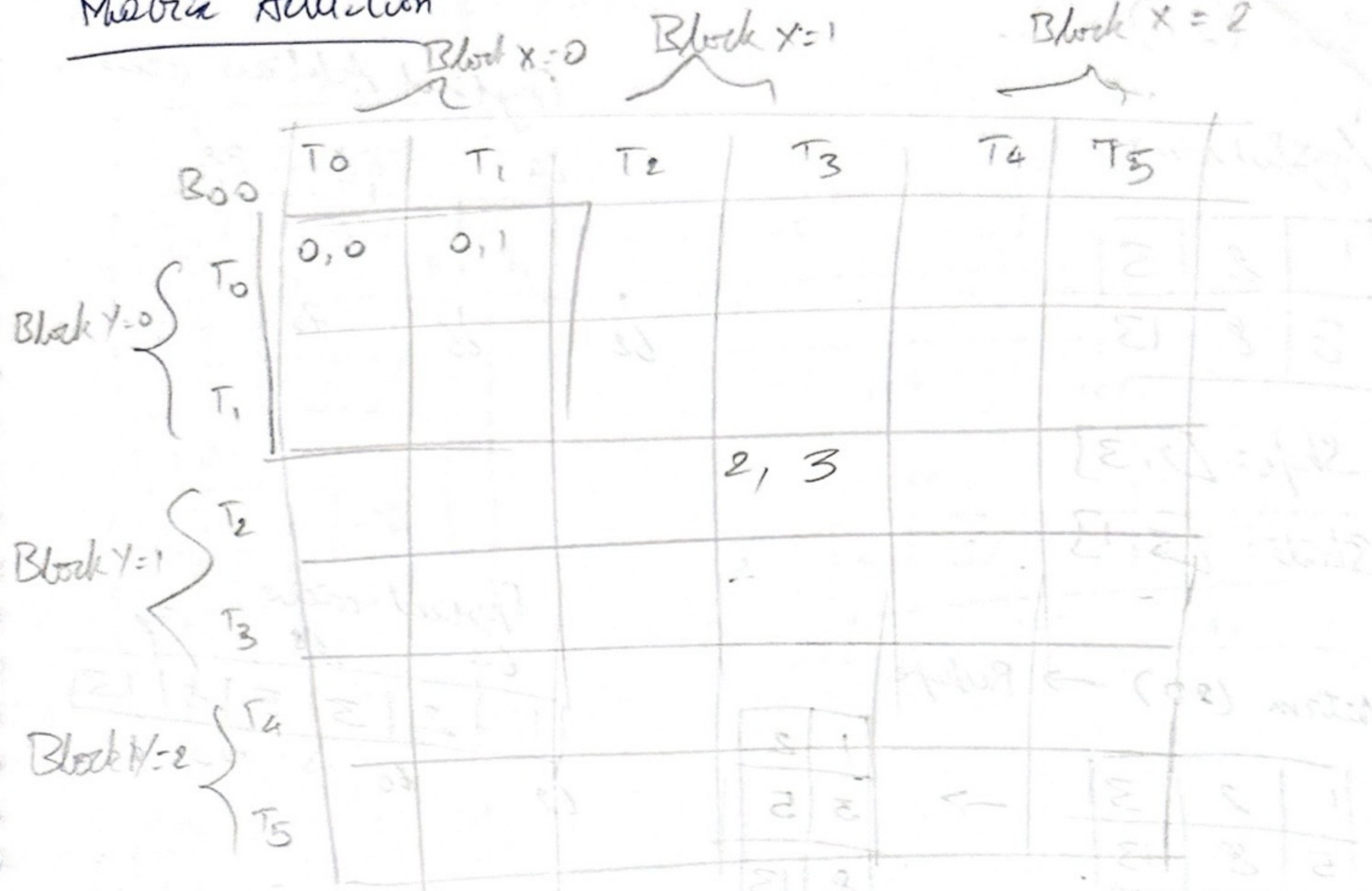


$$\text{Num. Blocks} = \frac{N}{\text{Block Size}} = 4$$

$$\text{Elem\_id} = i \left( B\text{-ID} \times \text{Block Size} + T\text{-id} \right)$$

$$N = 8 \quad \text{cols} = 4$$

## Matrix Addition



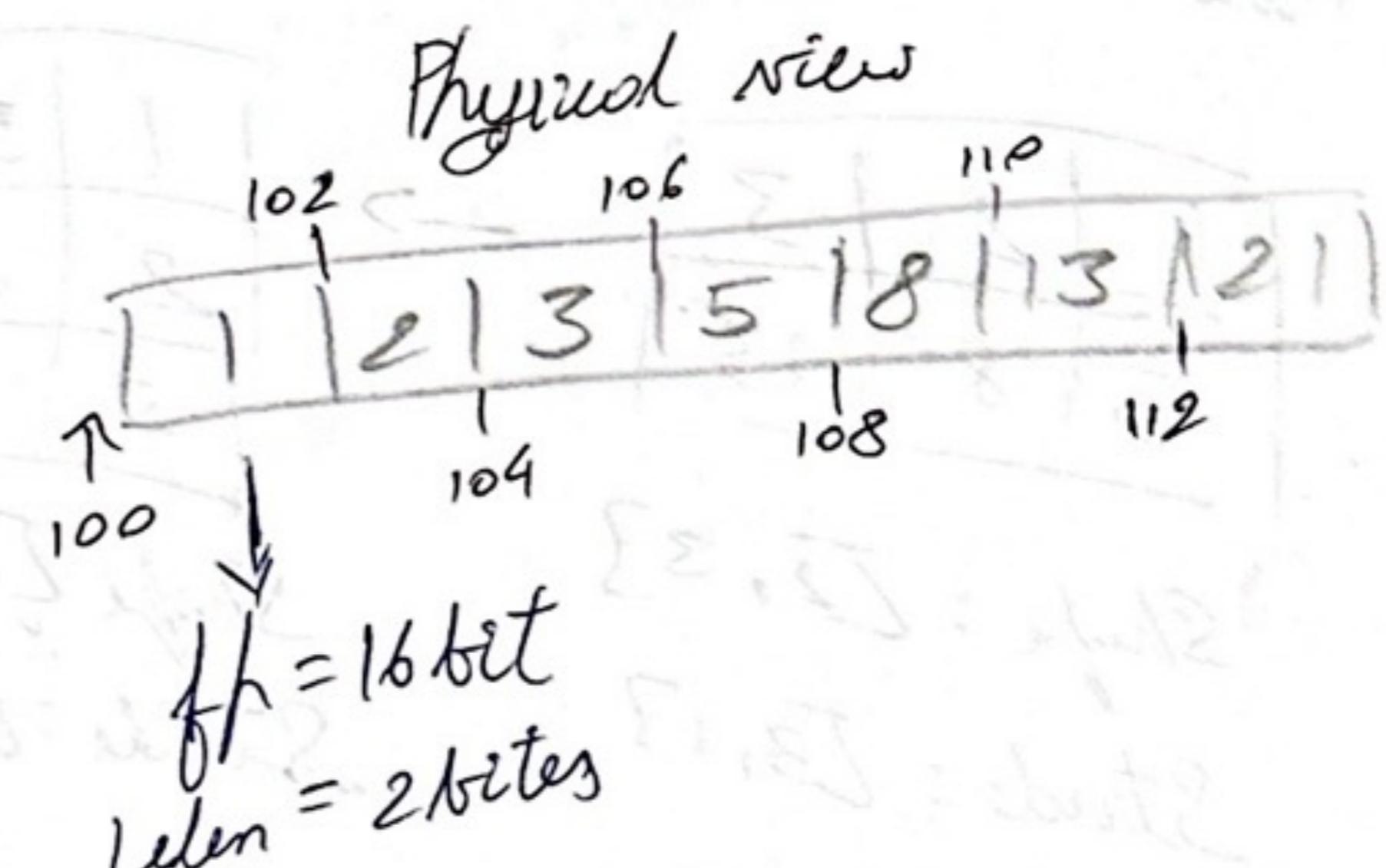
block\\_row\\_ident = bld\_y \* block\\_size + thread ID. Y  
 block\\_col\\_ident = bld\_x \* bl\\_size + thread ID. X

## Tensor layouts

array / vector (1D)

Logical view

1	2	3	4	5	6	7	8
---	---	---	---	---	---	---	---



Shape: [7]

Stride: [1]

memory block  $\rightarrow$  pointer

## Matrix (2D)

Logical view

1	2	3
5	8	13

Shape: [2, 3]

Stride: [3, 1]

Matrix (2D) → Reshape

1	2	3
5	8	13

Shape: [2, 3]

Stride: [3, 1]

1	2
3	5
8	13

Shape: [3, 2]  
Stride: [2, 1]

Physical Address view

64	68	72
1	2	3
5	8	13

Physical view

64	68	72
1	2	3
5	8	13

Matrix (2D) → Transpose

1	2	3
5	8	13

Shape: [2, 3]

Stride: [3, 1]

1	5
2	8
3	13

Shape: [3, 2]  
Stride: [1, 3]

$$S_{\text{stride}} =$$

$$\left\{ \begin{array}{l} S_{\text{stride}}[i] = \prod_{j=i+1}^{N-1} \text{shape}[j] \\ S_{\text{stride}}[N] = 1 \end{array} \right.$$

## Tensor (3D)

1	2	3
5	8	13
21	34	55
9	11	13

72	42	2
31	1	92
7	4	32
88	3	14

$\frac{1}{10^2}$	$\frac{1}{10^6}$
1	1

$\frac{1}{10^4}$	$\frac{1}{10^8}$
1	1

shape: [2, 4, 3]

stride: [12, 3, 1]

Algorithm: 1 Flash Attention - 2 forward pass

Matrices  $Q, K, V \in \mathbb{R}^{N \times d}$  in HBM, block size  $B_C, B_R$

Divide  $Q$  into  $T_Q = \left\lceil \frac{N}{B_R} \right\rceil$  blocks  $Q_1, \dots, Q_T$  of size  $B_R \times d$

each, and divide  $K, V$  into  $T_E = \left\lceil \frac{N}{B_C} \right\rceil$  blocks

$\text{BATCH\_SIZE} \times \text{NUM\_HEADS} \times \text{NUM\_BLOCKS}_n^Q$

$Q = [\text{B\_S}, \text{NumH}, \text{Seq\_len}, \text{Head}]$



$Q = [\text{Batch\_Size}, \text{Num\_Heads}, \text{Seq\_len}, \text{head\_dim}]$

Tensor :



Size : [Batch-Size, Num heads, Seq len, Head-Dim]  
[  
  Batch-Size  
  Num heads  
  Seq len  
  Head-Dim]  
  Seq len  
  Head-Dim  
]

Outer loop Q\_Block

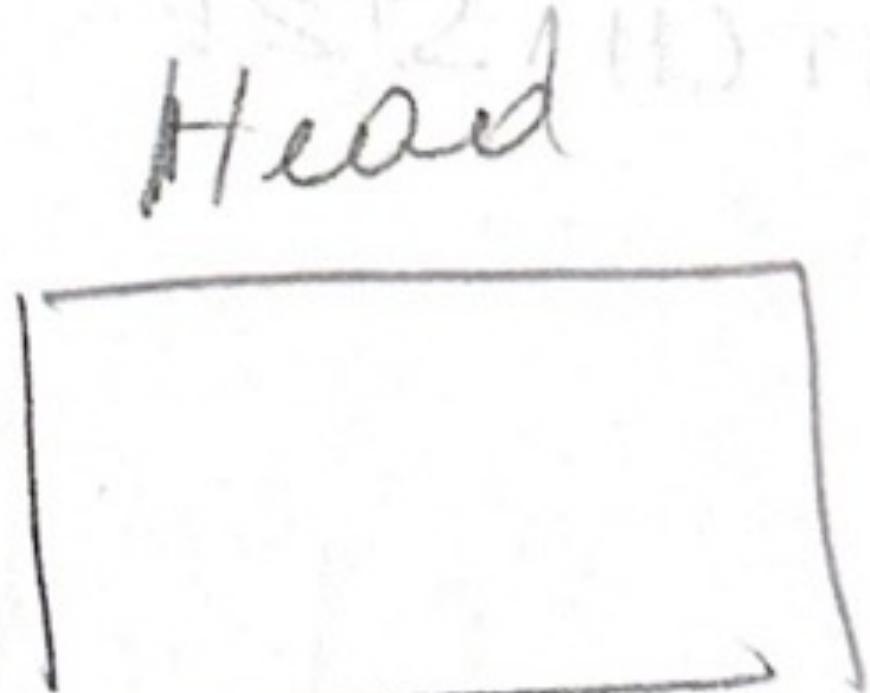
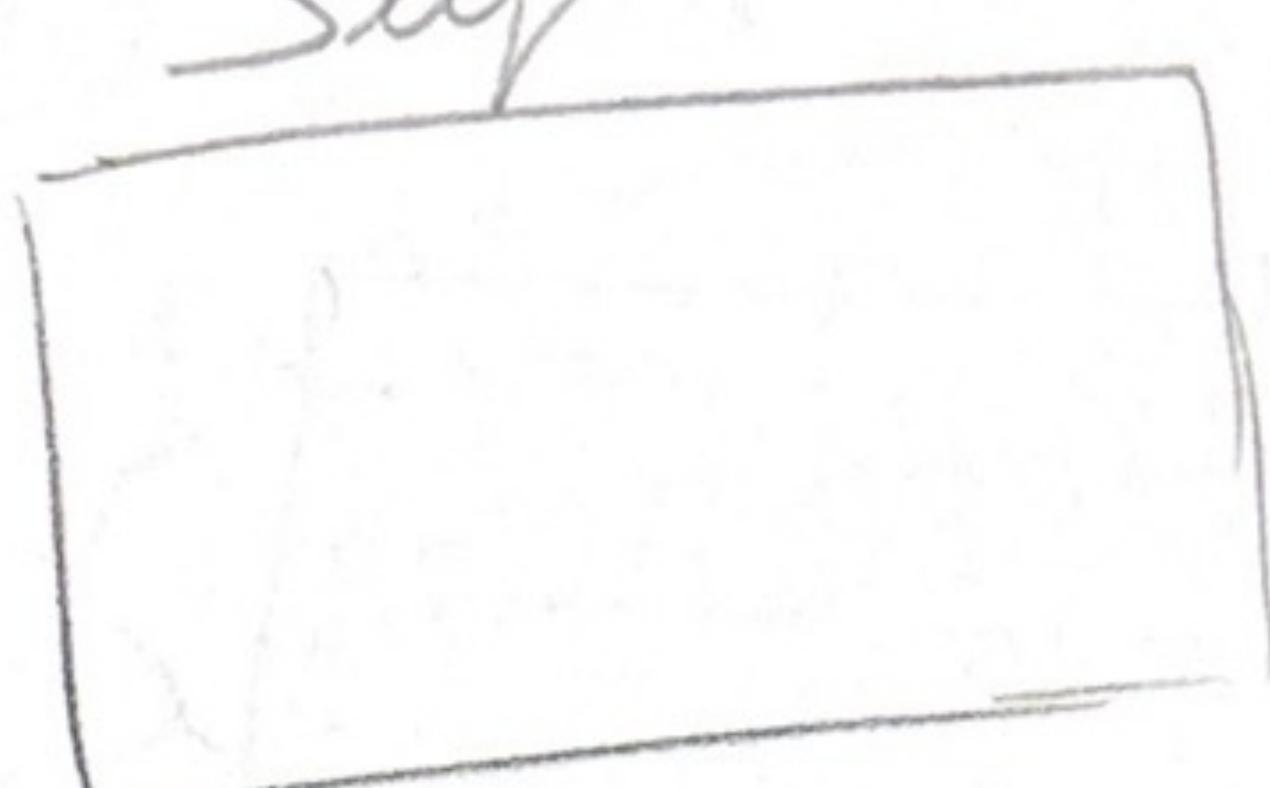
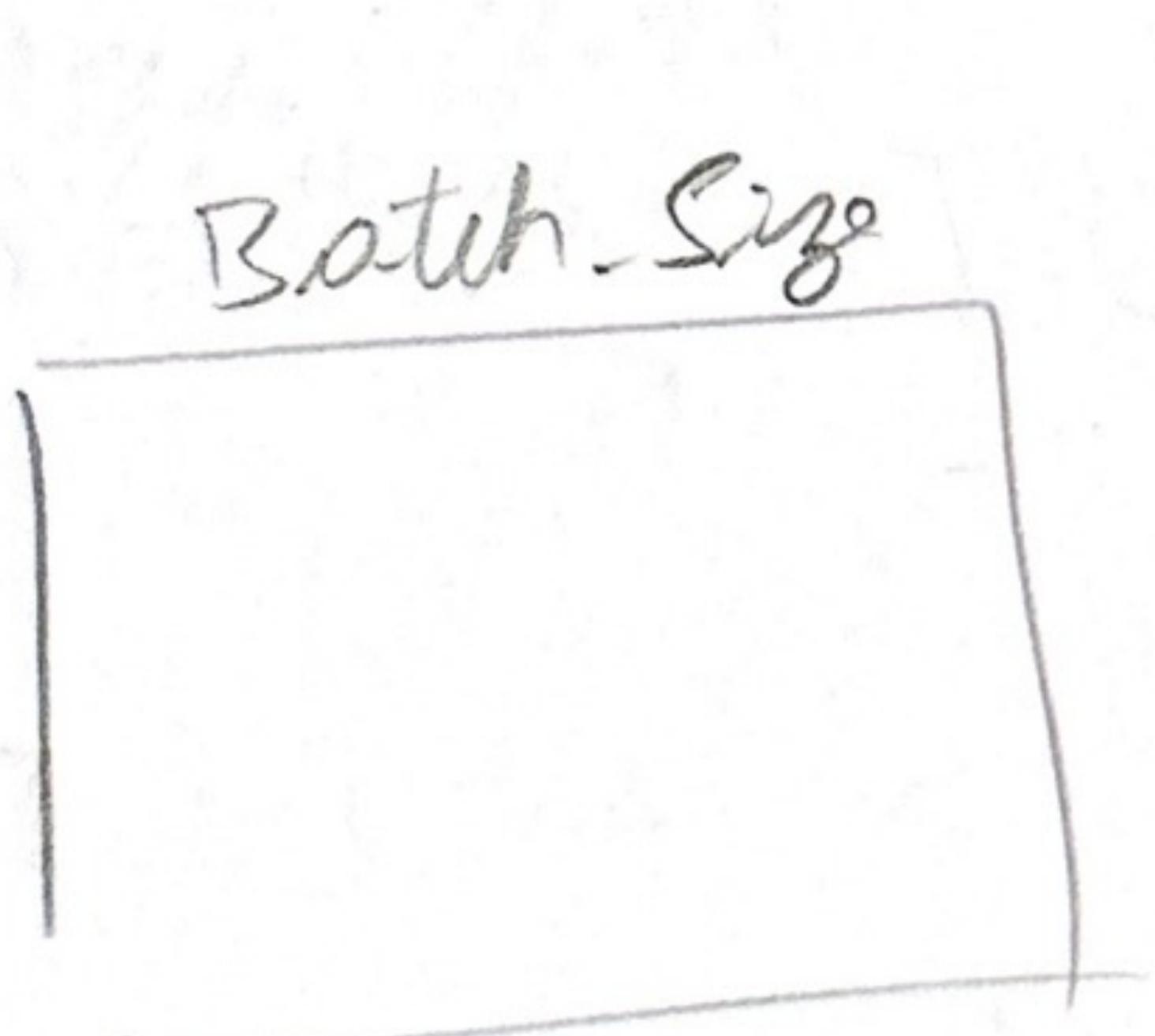
Inner loop all K\_Block

Program 2d =

AXIS 0 = [SEQ-LEN / BLOCK-SIZE-Q]  
2 4

AXIS 1 = [BLOCK-SIZE \* NUM-HEADS]

2 6

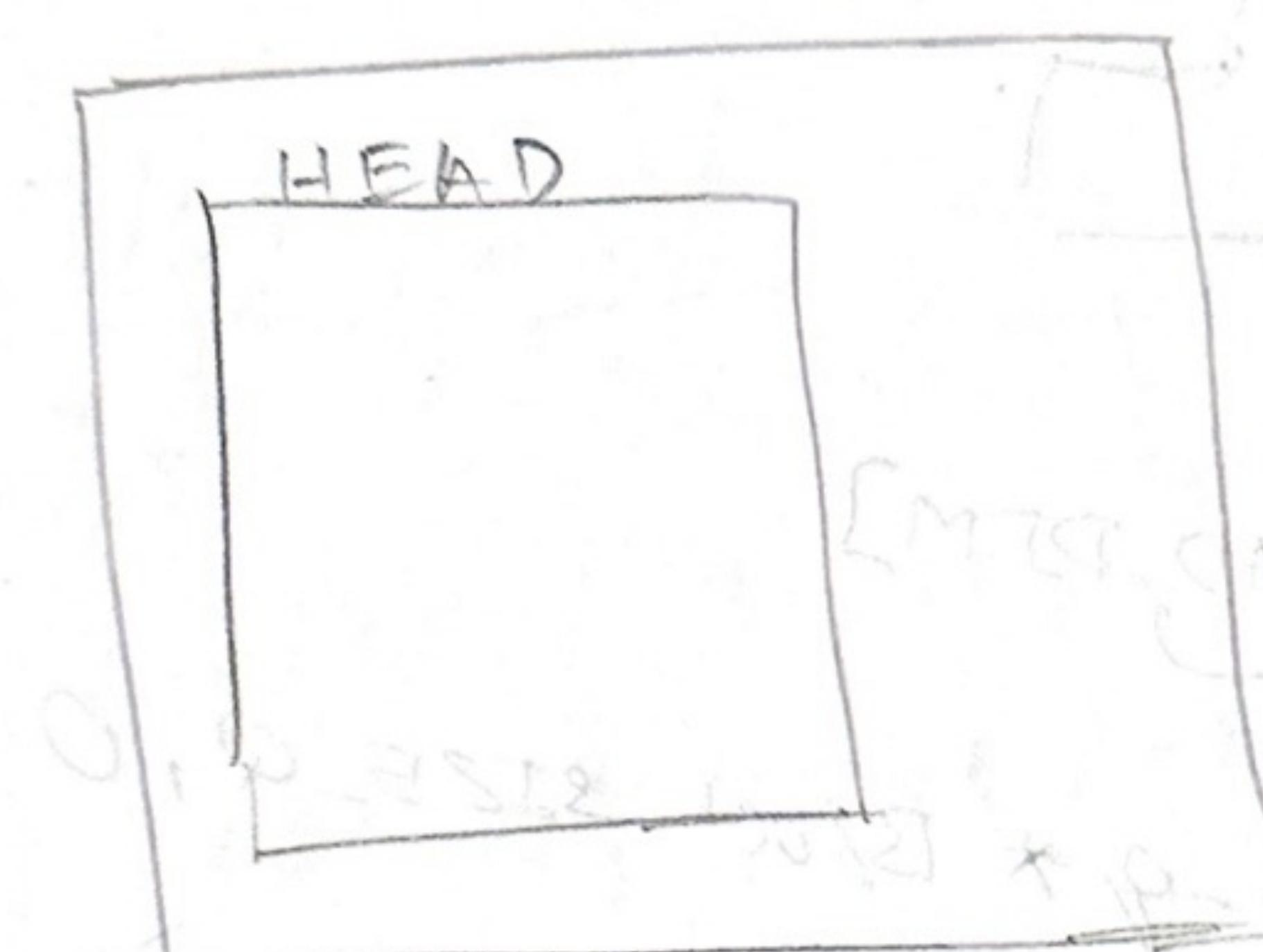


Batched

ProgramId = 0

$$Axis\_0 = \frac{\text{Seq\_len}}{\text{BlockSize}}$$

$$Axis\_1 = \text{Batch\_size} \times \text{num\_head}$$



$$Q = [B, \text{num\_heads}, \text{Seq\_len}, \text{Head\_dim}]$$

↑  
↑  
→  
pre-ded

↓  
→  
→

$$Q \rightarrow \text{offset}(\text{Block\_size}, * \\ \text{Block\_size} \cdot Q, 0)$$

↓  
↓  
↓

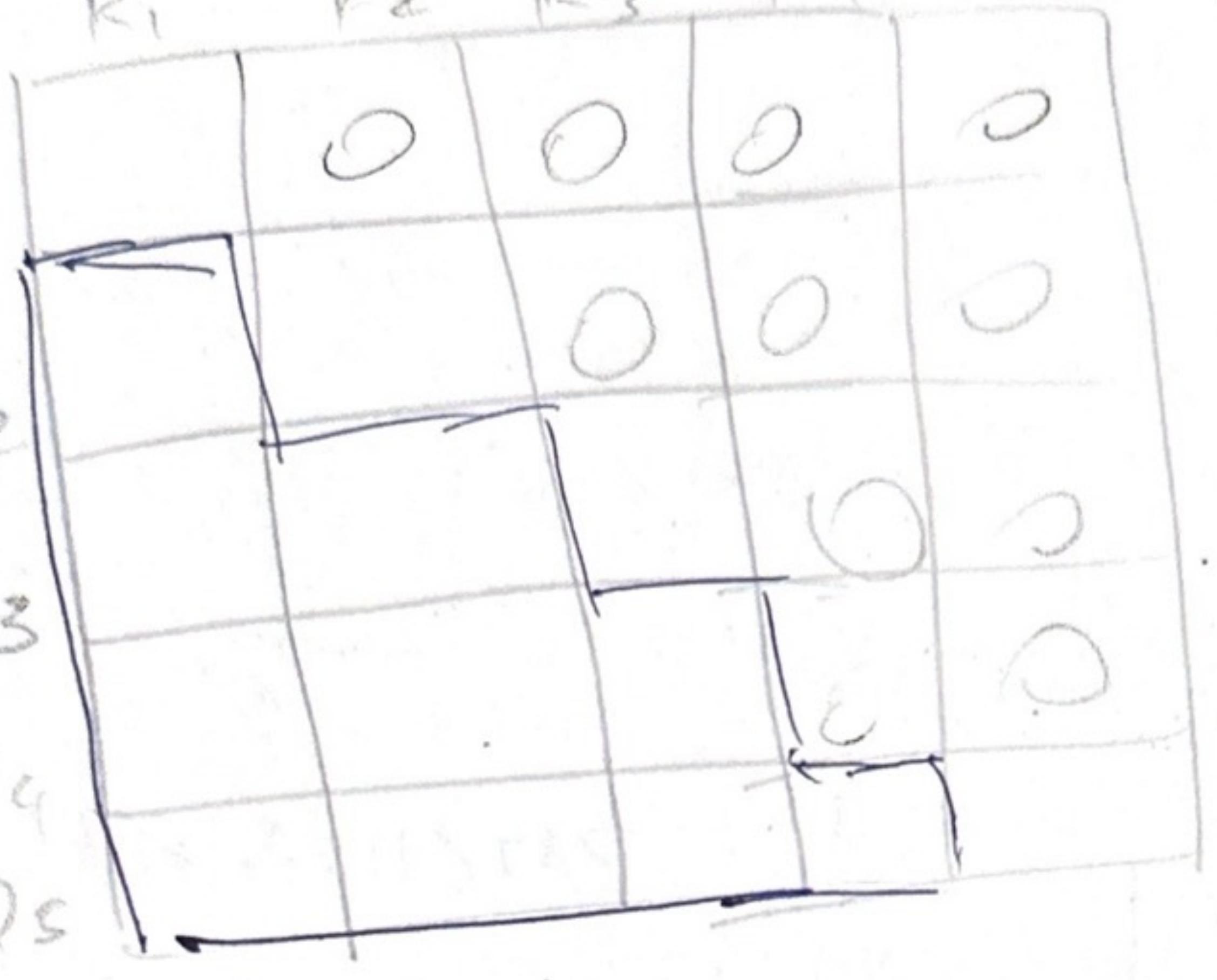
QBlock

↓

Seq-Q

order()

$0, 1, 2 \quad 3, 4, 5 \quad 6, 7, 8 \dots$   
 $k_1 \quad k_2 \quad k_3 \quad k_4 \quad 155$



$$D_{2y} = \text{blk\_idx\_q} * \text{blk\_size\_Q}_1 + (\text{blk\_idx\_q} + 1) * \text{blk\_size\_Q}$$

$$Q = [B, \underbrace{\text{NUM\_H}, \text{SEQ}, \text{HEAD\_DIM}}_{\downarrow \text{Block\_idx\_q} * \text{Block\_size\_Q}}, 0]$$

~~Inner loop~~

Inner loop  
 $\oplus k, v \rightarrow (Q \& k)$   
 $\downarrow$   
 max each row  $\rightarrow$  softmax  $\star + l$

$O =$

Kernel Attention = Non Correlate

~~exists~~  
 $\text{index } k \leq \text{index } Q$

## From derivatives to Jacobians

Derivative: scalar input, scalar output

$$f: \mathbb{R} \rightarrow \mathbb{R}$$

$$f'(x) = \lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h} \quad \begin{matrix} \leftarrow \text{how much the "output" changes} \\ \leftarrow \text{how much the "input" changes} \end{matrix}$$

$$f'(x) = \frac{\partial f(x)}{\partial x} = \frac{\partial y}{\partial x} = \lim_{h \rightarrow 0}$$

$$\frac{f(x+h) - f(x)}{h}$$

$$f(x+h) \cong f'(x)h + f(x)$$

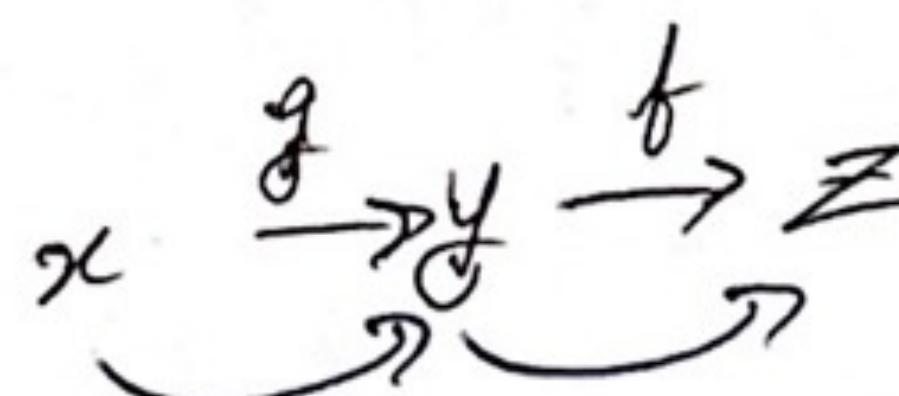
$$f(x+\Delta x) \cong f'(x)\Delta x + f(x) \quad y^{\text{NEW}} \cong \frac{\partial y}{\partial x} \Delta x + y^{\text{OLD}}$$

$$f(x+\Delta x) \cong \frac{\partial y}{\partial x} \Delta x + f(x)$$

$$x^{\text{new}} \rightarrow x^{\text{old}} + \Delta x \Rightarrow y^{\text{new}} \cong y^{\text{old}} + \frac{\partial y}{\partial x} \Delta x$$

Chain Rule

$$z = f(g(x))$$



$$\Delta y$$

$$x^{\text{new}} \rightarrow x^{\text{old}} + \Delta x \Rightarrow y^{\text{new}} \cong y^{\text{old}} + \frac{\partial y}{\partial x} \Delta x$$

$$y^{\text{NEW}} \rightarrow y^{\text{old}} + \Delta y \Rightarrow z^{\text{new}} \cong z^{\text{old}} + \frac{\partial z}{\partial y} \Delta y$$

$$z^{\text{new}} \rightarrow z^{\text{old}} + \frac{\partial z}{\partial y} \cdot \frac{\partial y}{\partial x} \Delta x$$

$$\frac{\partial z}{\partial x} = \frac{\partial z}{\partial y} \cdot \frac{\partial y}{\partial x}$$

Gradient: vector input, scalar output

$$f: \mathbb{R}^N \rightarrow \mathbb{R}$$

$$f \left( \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \right) = y$$

$$x \xrightarrow{\text{NEW}} x^{\text{OLD}} + \Delta x \Rightarrow y \xrightarrow{\text{NEW}} y^{\text{old}} + \frac{\partial y}{\partial x} \cdot \Delta x$$

vector  
sum

dot prod

$$\frac{\partial y}{\partial x} \leftarrow \text{gradient} = \left( \frac{\partial y}{\partial x_1}, \frac{\partial y}{\partial x_2}, \dots \right)$$

$$\frac{\partial y}{\partial x} \cdot \Delta x = \frac{\partial y}{\partial x_1} \cdot \Delta x_1 + \frac{\partial y}{\partial x_2} \cdot \Delta x_2 + \dots + \frac{\partial y}{\partial x_N} \cdot \Delta x_N$$

partial derivative

Cross Entropy loss:

$$\begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \Rightarrow \text{loss} = \text{scalar output}$$

Jacobian: vector input, vector output

$$f: \mathbb{R}^N \rightarrow \mathbb{R}^M$$

$$f \left( \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \right) = \begin{bmatrix} y_1 \\ y_2 \\ y_3 \end{bmatrix}$$

Jacobian =

$$\begin{bmatrix} \frac{\partial y_1}{\partial x_1} & \dots & \frac{\partial y_1}{\partial x_N} \\ \vdots & & \vdots \\ \frac{\partial y_M}{\partial x_1} & \dots & \frac{\partial y_M}{\partial x_N} \end{bmatrix}$$

$$x \xrightarrow{\text{NEW}} x^{\text{OLD}} + \Delta x \xrightarrow{\text{NEW}} y \xrightarrow{\text{NEW}} y^{\text{old}} + \frac{\partial y}{\partial x} \Delta x$$

Matrix-vector product

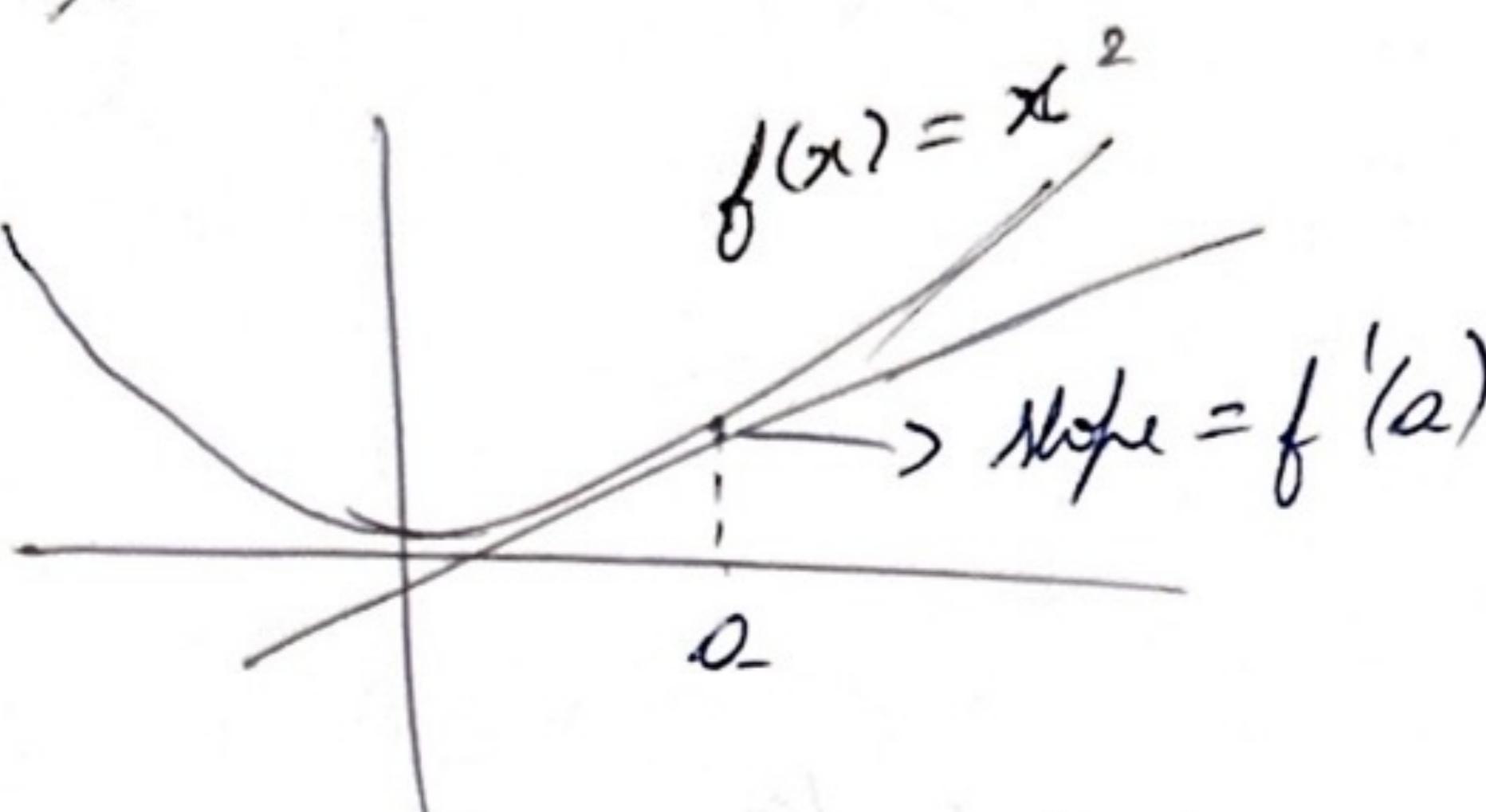
$$(M \times N) \times (N \times 1) = (M \times 1)$$

Derivatives:

$$f(x, y) = (x^2 - y^2, 2xy)$$

Derivatives = Slope

Derivatives ≠ Slopes



Linear Maps:

$$\begin{pmatrix} 1.33 \\ -0.73 \end{pmatrix} \begin{pmatrix} 1.17 \\ 0.75 \end{pmatrix} = (1.33 * 0.75) - (1.17 * -0.73)$$

Derivatives in 1D

$$f(x, y) = (x^2 - y^2, 3xy)$$

$$f(x, y) = (f_1(x, y), f_2(x, y))$$

$$f(x, y) = (x^2 - y^2, 3xy)$$

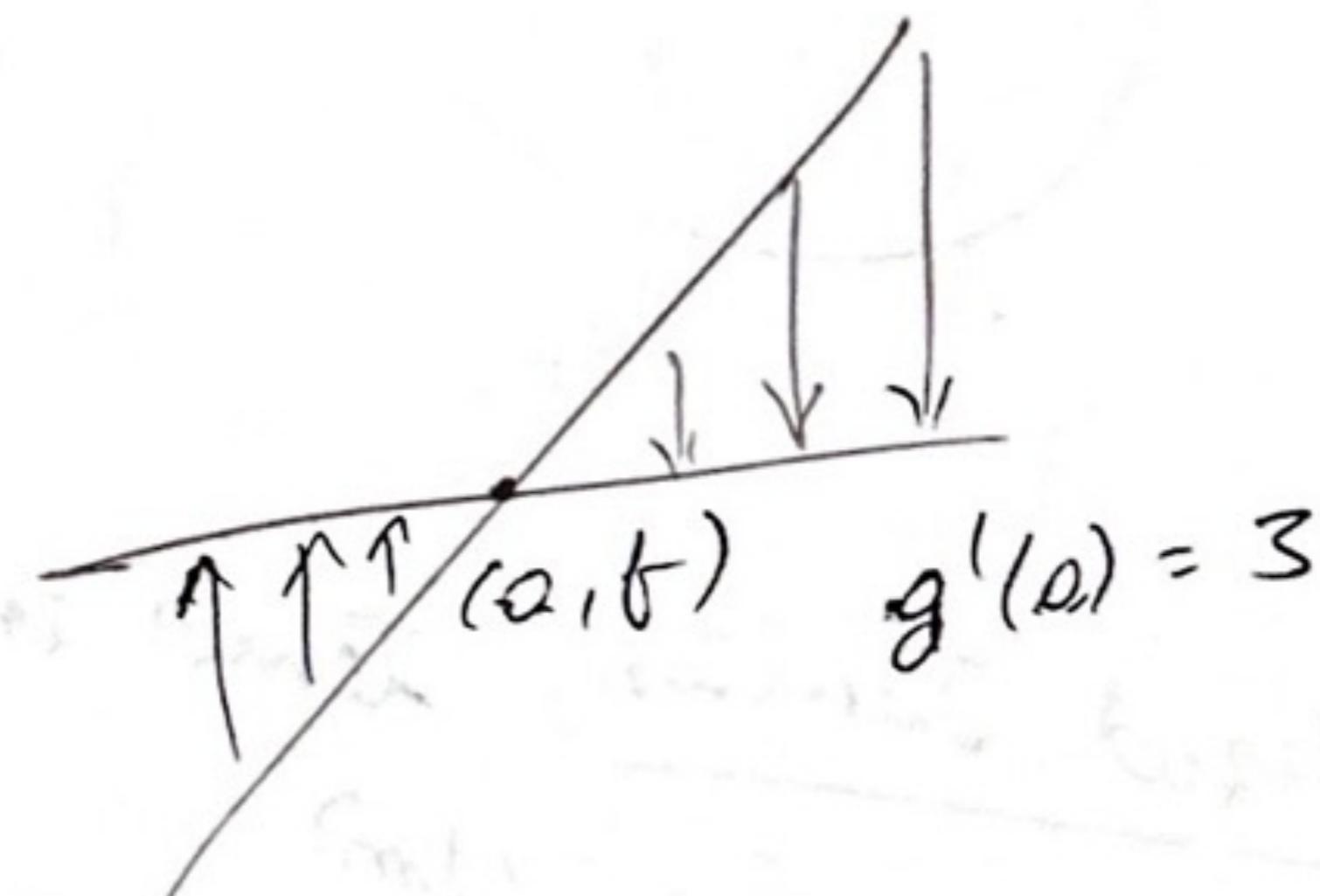
$$g(x) = f_1(x, b)$$

↓  
fixed

$$g'(a) = \frac{\partial f_1}{\partial x} \Big|_{(a, b)}$$

$$f(x, y) = (f_1(x, y), f_2(x, y))$$

$$f(x, y) = (x^2 - y^2, 3xy)$$



Jacobian Matrix

$$\begin{pmatrix} 1 \\ 1.5 \end{pmatrix}$$

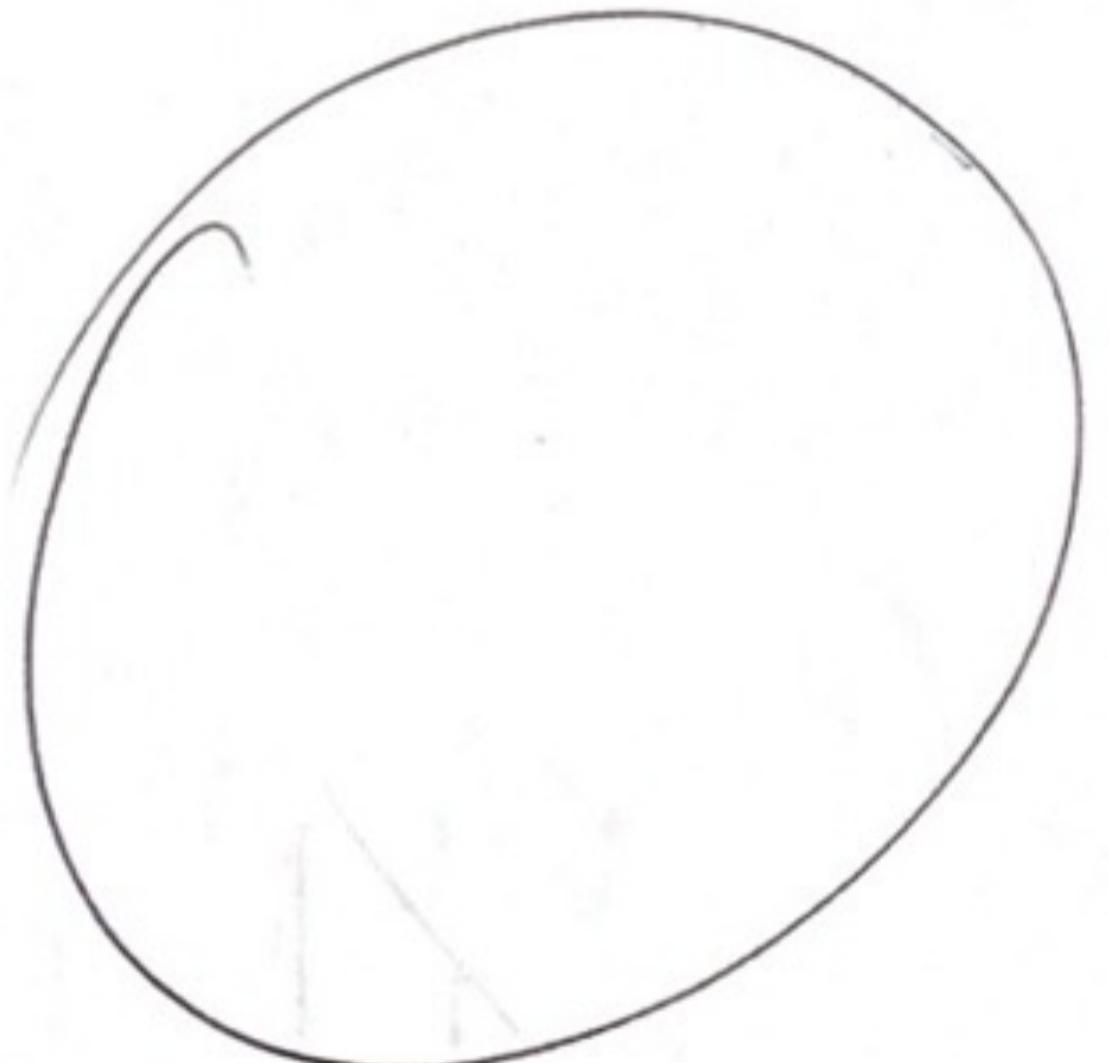
$$\begin{pmatrix} -1 \\ 1.5 \end{pmatrix} \Rightarrow f'(b)$$

$$\begin{pmatrix} \frac{\partial f_1}{\partial x} \Big|_{(a, b)} & \frac{\partial f_1}{\partial y} \Big|_{(a, b)} \\ \frac{\partial f_2}{\partial x} \Big|_{(a, b)} & \frac{\partial f_2}{\partial y} \Big|_{(a, b)} \end{pmatrix}$$

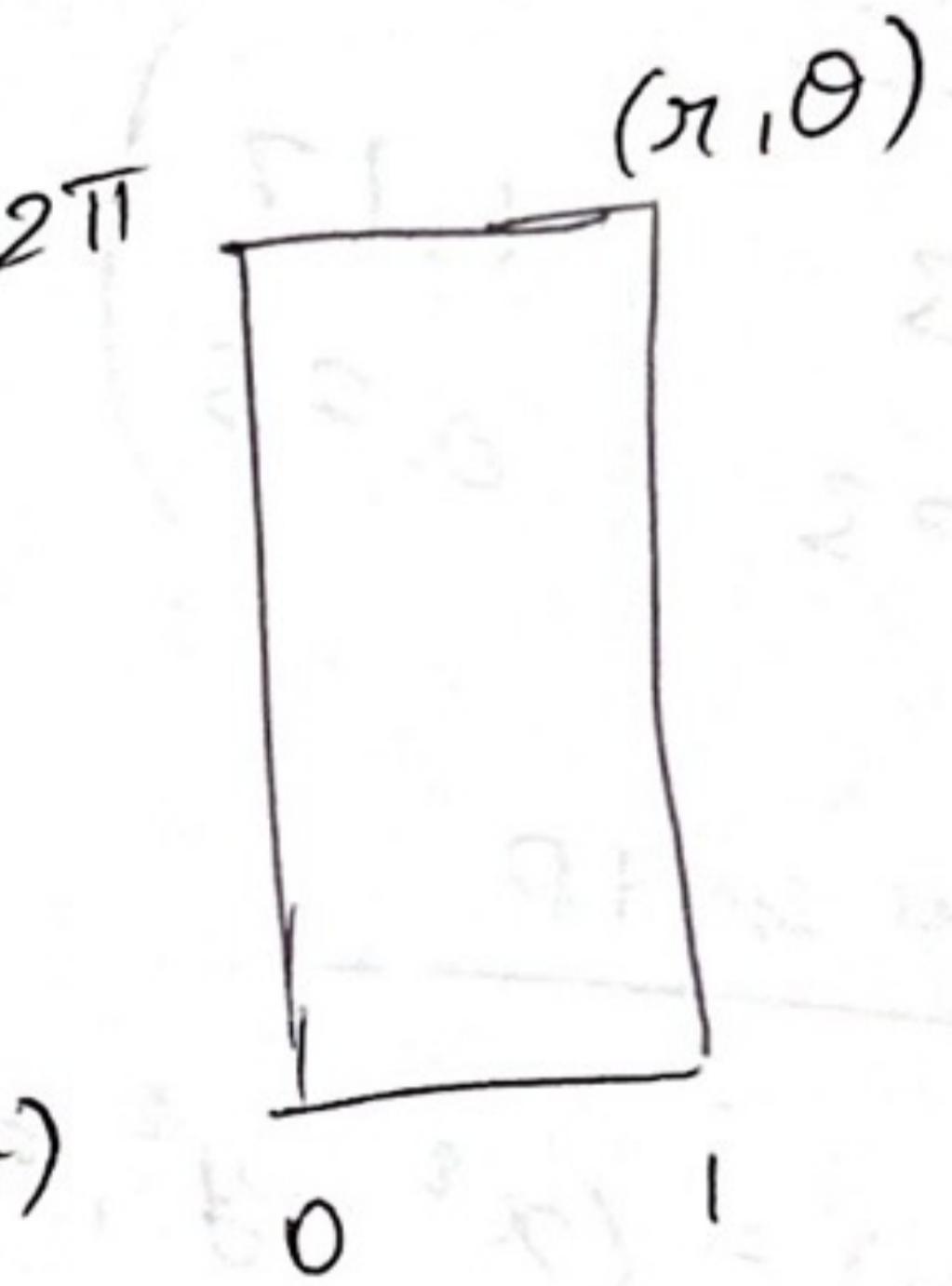
## Jacobien determinant

$$\begin{array}{|c|c|} \hline \frac{\partial f_1}{\partial x} & \left( \frac{\partial f_1}{\partial y}, (a, b) \right) \\ \hline \frac{\partial f_2}{\partial x} & \left( \frac{\partial f_2}{\partial y}, (a, b) \right) \\ \hline \end{array}$$

$$\iint_D f(x, y) dx dy$$



$$\iint_T f(g(r, \theta)) \cdot \text{abs}|\mathcal{J}| dr d\theta$$



Apply g

$$g(r, \theta) = (r \cos \theta, r \sin \theta)$$

Generalized Jacobian: tensor input, tensor output

$$f: \underbrace{\mathbb{R}^{N_1 \times \dots \times N_{D_X}}}_{\text{in tensor}} \xrightarrow{(\text{dim})} \underbrace{\mathbb{R}^{M_1 \times \dots \times M_{D_Y}}}_{\text{output tensor}}$$

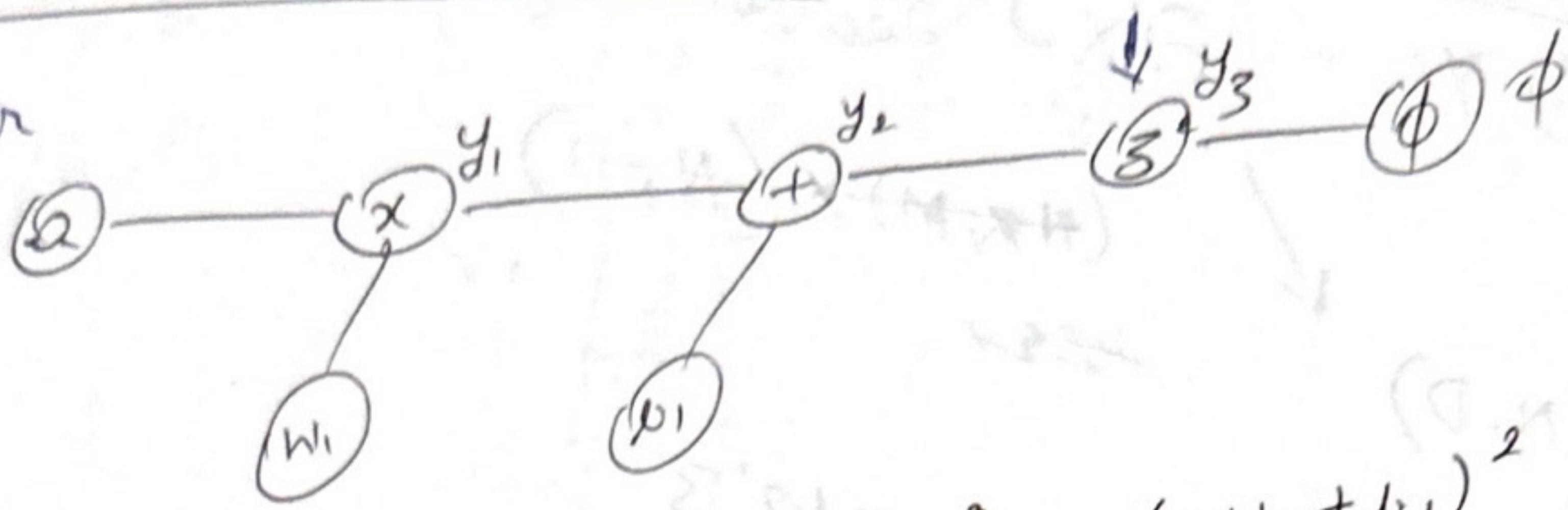
$$f(D_x - \text{dimension tensor}) = \underbrace{D_y - \text{dimension tensor}}_{\text{generalized Jacobian}}$$

$$x^{\text{NEW}} \rightarrow x^{\text{OLD}} + \Delta x \Rightarrow y^{\text{NEW}} \rightarrow y^{\text{OLD}} + \underbrace{\frac{\partial y}{\partial x} \Delta x}_{\text{Tensor product}}$$

$$(M_1 \times \dots \times M_{D_Y}) \times (N_1 \times \dots \times N_{D_X})$$

## Autograd with derivatives

Sector



$$\phi = y_3 = (y_2)^2 = (y_1 + b_1)^2 = (a w_1 + b_1)^2$$

$$\frac{\partial \phi}{\partial w_1} = 2(a w_1 + b_1)(a) = 2a(a w_1 + b_1)$$

$$\frac{\partial \phi}{\partial w_1} = \frac{\partial \phi}{\partial y_3} \cdot \frac{\partial y_3}{\partial y_2} \cdot \frac{\partial y_2}{\partial y_1} \cdot \frac{\partial y_1}{\partial w_1}$$

$$= 1 \cdot 2y_2 \cdot 1 \cdot a = 2a y_2 = 2a(a w_1 + b_1)$$

$$\frac{\partial \phi}{\partial w_1} = \underbrace{\frac{\partial \phi}{\partial y_3} \cdot \frac{\partial y_3}{\partial y_2}}_{\frac{\partial \phi}{\partial y_2}} \cdot \underbrace{\frac{\partial y_2}{\partial y_1} \cdot \frac{\partial y_1}{\partial w_1}}_{\begin{matrix} [N, d] \\ [d, M] \end{matrix}}$$

Pytorch

$$\frac{\partial \phi}{\partial y_2} = \frac{\partial \phi}{\partial y_3} \frac{\partial y_3}{\partial y_2}$$

$$\frac{\partial \phi}{\partial x} = \frac{\partial \phi}{\partial x} \cdot \frac{\partial x}{\partial x}$$

local variable

$$\frac{\partial \phi}{\partial y_1} = \frac{\partial \phi}{\partial y_2} \cdot \frac{\partial y_2}{\partial y_1}$$

Downstream  
gradient

$$\frac{\partial \phi}{\partial w_1} = \frac{\partial \phi}{\partial y_1} \cdot \frac{\partial y_1}{\partial w_1}$$

N=1024

d=1024

M=1024 2048

$$\frac{\partial \phi}{\partial x} = \frac{\partial \phi}{\partial y} \cdot \frac{\partial y}{\partial x} \quad \left. \begin{array}{l} \frac{\partial y}{\partial x} \\ \downarrow \end{array} \right\} \text{local Jacobian}$$

~~(N, M) × (N, D)~~

1024 × 2048 × 1024 × 1024 = 2.19 "B"

$$\begin{bmatrix} 1 & \dots & 1024 \end{bmatrix} \times \begin{bmatrix} 1 & \dots & 1024 \end{bmatrix} \times \begin{bmatrix} 1 & \dots & 1024 \end{bmatrix} \times \begin{bmatrix} 1 & \dots & 1024 \end{bmatrix} = (N, M)$$

(N, D)      (D, M)

Gradient of "Math Mat" operation

$$Y = XW \quad \frac{\partial \phi}{\partial Y} = \frac{\partial \phi}{\partial X} \frac{\partial \phi}{\partial W}$$

$$\frac{\partial \phi}{\partial X} = [N, D] \quad N = 1 \quad D = 3$$

$$W = [D, M] \quad M = 4$$

$$Y = [N, M] = [1, 4]$$

$$x = [N, D]$$

$$w = [D, M]$$

$$y = [N, M]$$

$$x = [1, 2, 3]$$

$$w = \begin{bmatrix} 1 & 2 & 3 & 4 \\ 1 & 2 & 3 & 4 \\ 1 & 2 & 3 & 4 \end{bmatrix}$$

$$y = \begin{bmatrix} 1 \\ 2 \\ 3 \\ 4 \end{bmatrix}$$

$$N = 1 \quad M = 4$$

$$D = 3$$

$$y = \begin{bmatrix} x_{11} & x_{12} & x_{13} \end{bmatrix} \times$$

X  
(1, 3)

$$\begin{bmatrix} w_{11} & w_{12} & w_{13} & w_{14} \\ w_{21} & w_{22} & w_{23} & w_{24} \\ w_{31} & w_{32} & w_{33} & w_{34} \end{bmatrix}$$

W  
(3, 4)

$$= \cancel{\left[ x_{11} w_{11} + x_{12} w_{21} + x_{13} w_{31} \right]}$$

$$= \left[ (x_{11} w_{11} + x_{12} w_{21} + x_{13} w_{31}) \quad (x_{11} w_{12} + x_{12} w_{22} + x_{13} w_{32}) \quad (x_{11} w_{13} + x_{12} w_{23} + x_{13} w_{33}) \quad (x_{11} w_{14} + x_{12} w_{24} + x_{13} w_{34}) \right] \quad y \in (1, 4)$$

$$\frac{\partial \phi}{\partial y} = \begin{bmatrix} dy_{11} & dy_{12} & dy_{13} & dy_{14} \end{bmatrix} - \underbrace{\begin{bmatrix} [N, M] & [M, D] \end{bmatrix}}_{\frac{\partial \phi}{\partial y} \cdot W^T = \begin{bmatrix} [N, D] \end{bmatrix}},$$

$$\frac{\partial \phi}{\partial x} = \frac{\partial \phi}{\partial y} \cdot \frac{\partial y}{\partial x}$$

$$\frac{\partial y}{\partial x} = \begin{bmatrix} w_{11} & w_{21} & w_{31} \\ w_{12} & w_{22} & w_{32} \\ w_{13} & w_{23} & w_{33} \\ w_{14} & w_{24} & w_{34} \end{bmatrix} \sim W^T$$

$$\frac{\partial \phi}{\partial w} = x^T \frac{\partial \phi}{\partial y}$$

$[D, N] \quad [N, M]$

$$N=2, M=4, D=3$$

$$X = \begin{bmatrix} N, D \\ 2, 3 \end{bmatrix} \quad X = \begin{bmatrix} x_{11} & x_{12} & x_{13} \\ x_{21} & x_{22} & x_{23} \end{bmatrix}$$

$$Y = \begin{bmatrix} H, M \\ 3, 4 \end{bmatrix}$$

$$\begin{bmatrix} w_{11} & w_{12} & w_{13} & w_{14} \\ w_{21} & w_{22} & w_{23} & w_{24} \\ w_{31} & w_{32} & w_{33} & w_{34} \end{bmatrix}$$

$$Y = \begin{bmatrix} N, M \\ 2, 4 \end{bmatrix}$$

$$(x_{11}w_{11} + x_{11}w_{21} + x_{11}w_{31}) \quad (x_{11}w_{12} + x_{11}w_{22} + x_{11}w_{32}) \quad (x_{11}w_{13} + \dots) \quad (x_{11}w_{14} + \dots)$$

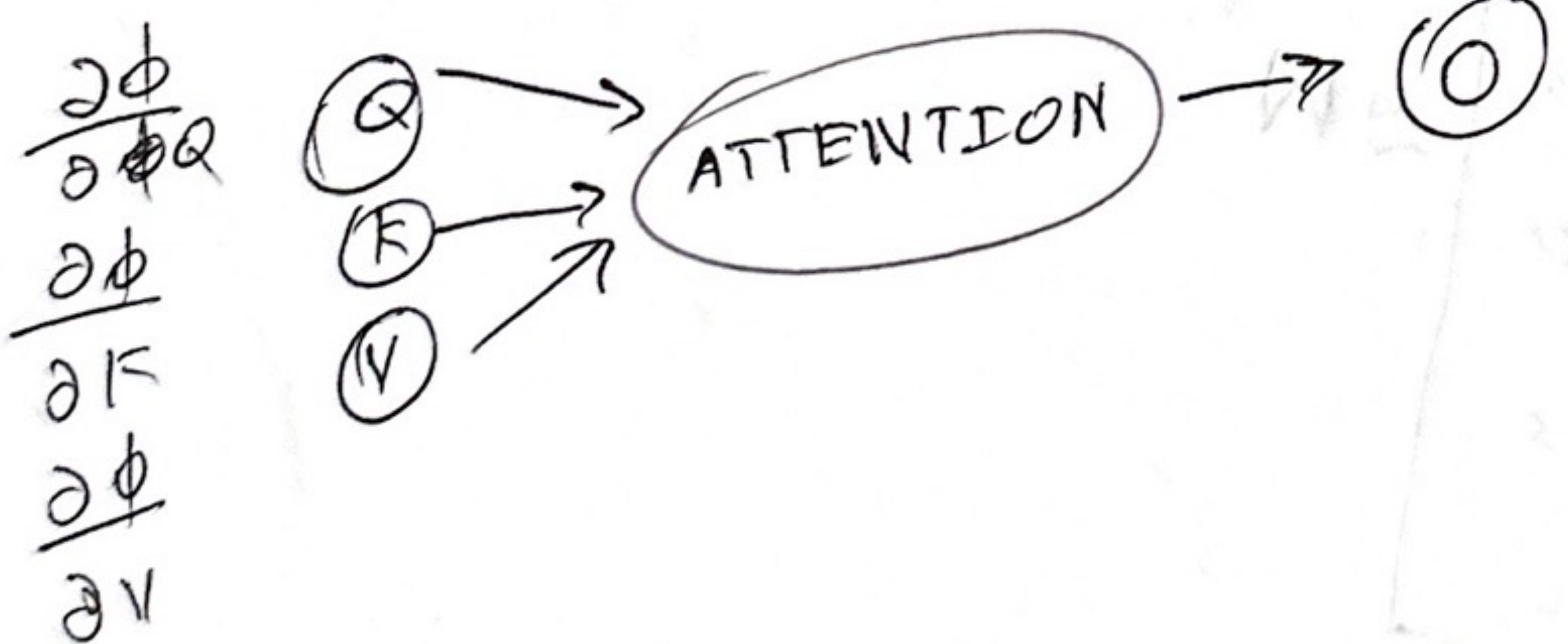
Gradient through the Softmax

$$S = QK^T$$

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

$$O = PV$$

$$\frac{\partial \phi}{\partial O}$$



$P = \text{Softmax Row}(S)$

$$S_i = S[i, :] \in \mathbb{R}^N$$

$$P_i = \text{Softmax}(S_i) \in \mathbb{R}^N$$

$$\text{Softmax}(P_{ij}) = \frac{e^{S_{ij}}}{\sum_{l=1}^N e^{S_{il}}} - S_{i\text{MAX}}$$

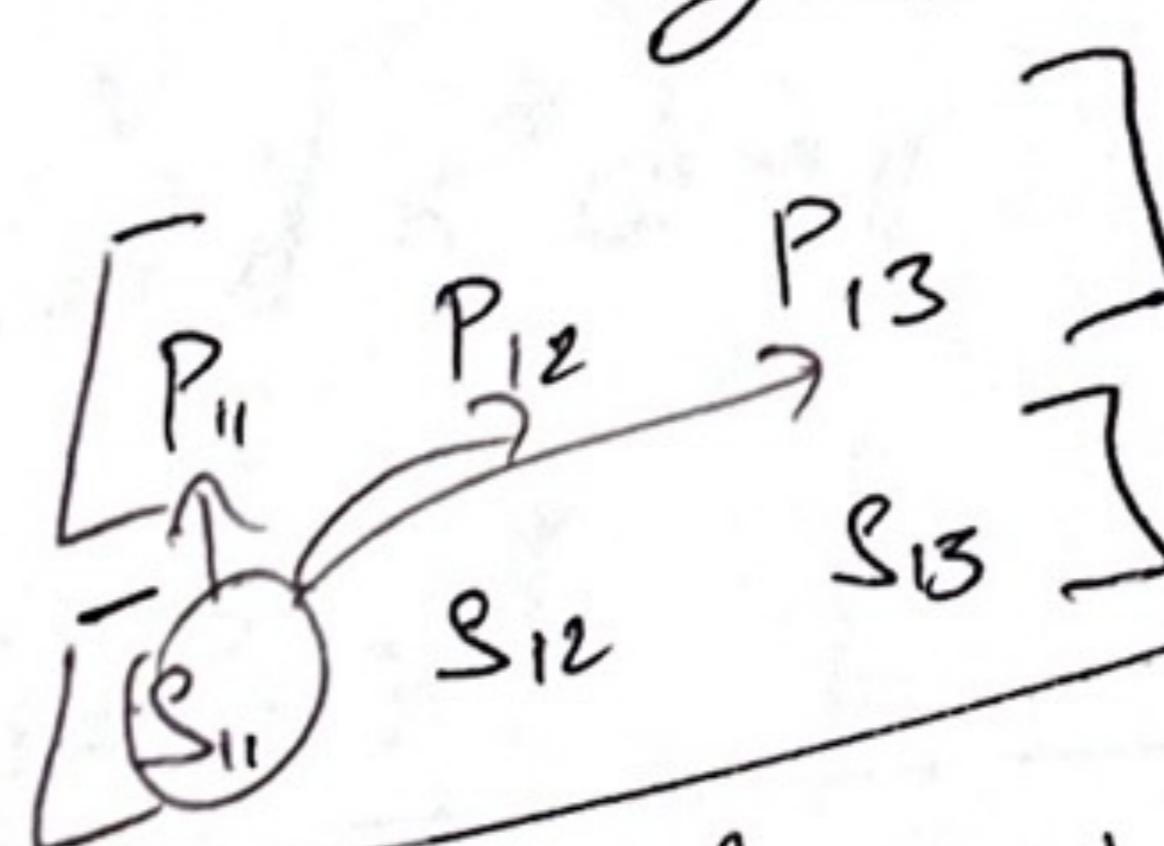
$$\cos^2(x) + \sin^2(x) = 1$$

$$y = \cos^2(x)$$

$$\frac{\partial y}{\partial x} = \frac{\partial(\cos^2(x))}{\partial x} \approx \frac{\partial(1 - \sin^2(x))}{\partial x}$$

$$\frac{\partial \phi}{\partial S_i} = \frac{\partial \phi}{\partial P_i} \cdot \frac{\partial P_i}{\partial S_i}$$

$$\frac{\partial P_{ij}}{\partial S_{ik}} = \frac{\partial}{\partial S_{ik}} \left[ \frac{e^{S_{ij}}}{\sum_{l=1}^N e^{S_{il}}} \right]$$



$$\frac{\partial \phi}{\partial S_i} = \frac{\partial \phi}{\partial P_i} \cdot \frac{\partial P_i}{\partial S_i}$$

$$\frac{\partial P_{ij}}{\partial S_{ik}} = \frac{\partial}{\partial S_{ik}} \left[ \frac{e^{S_{ij}}}{\sum_{l=1}^N e^{S_{il}}} \right]$$

$$\frac{\partial}{\partial S_{ik}}$$

$$P_{ii} = S_{ii}$$

=

$$\frac{e^{S_{ij}} \left( \sum_{l=1}^N e^{S_{il}} \right) - e^{S_{ik}} \cdot e^{S_{ij}}}{\left( \sum_{l=1}^N e^{S_{il}} \right)^2}$$

$$\frac{e^{S_{ij}} \left( \sum_{l=1}^N e^{S_{il}} - e^{S_{ik}} \right)}{\left( \sum_{l=1}^N e^{S_{il}} \right)^2}$$

$$P_{ij} \cdot (1 - P_{ik})$$

$$= \frac{e^{S_{ij}}}{\sum_{l=1}^N e^{S_{il}}} \times \frac{\sum_{l=1}^N e^{S_{il}} - e^{S_{ik}}}{\sum_{l=1}^N e^{S_{il}}}$$

$$\frac{\partial P_{ij}}{\partial s_{ik}} = \frac{\partial}{\partial s_{ik}} \left[ \frac{e^{s_{ij}}}{\sum_{l=1}^N e^{s_{il}}} \right]$$

$j \neq k$

$$\left[ \frac{f(x)}{g(x)} \right] = \frac{f'(x)g(x) - g'(x)f(x)}{[g(x)]^2}$$

$$\frac{0 * g - e^{s_{ik}} e^{s_{ij}}}{\left[ \sum_{l=1}^N e^{s_{il}} \right]^2}$$

~~$\frac{\partial P_{ii}}{\partial s_{ik}} = 0$~~

$$P = \begin{bmatrix} p_{11} & p_{12} & p_{13} \\ & & \\ & & \end{bmatrix}$$

$$S = \begin{bmatrix} s_{11} & s_{12} & s_{13} \\ & & \\ & & \end{bmatrix}$$

$$= \frac{-e^{s_{ik}}}{\left( \sum_{l=1}^N e^{s_{il}} \right)} \cdot \frac{e^{s_{ij}}}{\left( \sum_{l=1}^N e^{s_{il}} \right)}$$

$\downarrow$  Softm of  $k^{\text{th}}$  elem

$\downarrow$  Softm of  $j^{\text{th}}$  elem

$$= -P_{ik} P_{ij}$$

$$j=k = P_{ij} (1-P_{ik})$$

$$\frac{\partial P_{ij}}{\partial s_{ik}} = \frac{\partial}{\partial s_{ik}} \left[ -P_{ik} P_{ij} \right]$$

$$[N, N]$$

$$\frac{\partial P_{ij}}{\partial s_{ik}} =$$

$$\begin{bmatrix} P_{11} (1-P_{11}) & -P_{11} P_{12} & -P_{11} P_{13} & \dots & -P_{11} P_{1N} \\ -P_{12} P_{11} & P_{12} (1-P_{12}) & -P_{12} P_{13} & \dots & -P_{12} P_{1N} \\ -P_{13} P_{11} & -P_{13} P_{12} & P_{13} (1-P_{13}) & \dots & -P_{13} P_{1N} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ -P_{1N} P_{11} & -P_{1N} P_{12} & -P_{1N} P_{13} & \dots & P_{1N} (1-P_{1N}) \end{bmatrix}$$

$$P \cdot P^T$$

$$\{ \text{fading}(p_i) - p_i P_i^T \} \text{ As for } \text{Flesh Attn Paper}$$

$$PPT = \begin{bmatrix} P_1 \\ P_2 \\ P_3 \end{bmatrix} \begin{bmatrix} P_1 & P_2 & P_3 \end{bmatrix} = \begin{bmatrix} P_1 P_1 & P_1 P_2 & P_1 P_3 \\ P_2 P_1 & P_2 P_2 & P_2 P_3 \\ P_3 P_1 & P_3 P_2 & P_3 P_3 \end{bmatrix}$$

Formula as per Frost Attn Paper

$$y = \text{softmax}(x)$$

$$\text{Jacobian} = \text{diag}(y) - yy^T$$

B1. Memory efficient forward pass

$$o_i = p_i v = \sum_j p_{ij} v_j = z$$

$$a \begin{bmatrix} a_1 & a_2 & a_3 \end{bmatrix} \times \begin{bmatrix} b_{11} & b_{12} & b_{13} & b_{14} \\ b_{21} & b_{22} & b_{23} & b_{24} \\ b_{31} & b_{32} & b_{33} & b_{34} \end{bmatrix}$$

$$o_i = p_i v = \sum_j e$$

$$o_i = \begin{bmatrix} (a_{11}b_{11} + a_{12}b_{21} + a_{13}b_{31}) \\ (a_{11}b_{12} + a_{12}b_{22} + a_{13}b_{32}) \\ (a_{11}b_{13} + a_{12}b_{23} + a_{13}b_{33}) \\ (a_{11}b_{14} + a_{12}b_{24} + a_{13}b_{34}) \\ (a_{11}b) \end{bmatrix} \quad (N, 4)$$

$$\begin{bmatrix} a_1 & a_2 & a_3 \end{bmatrix} \times \begin{bmatrix} b_{11} & b_{12} & b_{13} & b_{14} \\ b_{21} & b_{22} & b_{23} & b_{24} \\ b_{31} & b_{32} & b_{33} & b_{34} \\ b_{41} & b_{42} & b_{43} & b_{44} \end{bmatrix} \quad (4, 4)$$

$$\begin{aligned} o_i &= a_{11} b_{11} + a_{12} b_{21} + a_{13} b_{31} \\ &\quad + a_{14} b_{41} \\ &= \sum_{j=1}^3 a_j \end{aligned}$$

$$A = \begin{bmatrix} a_{11} & a_{12} & a_{13} \end{bmatrix}$$

x

$$B = \begin{bmatrix} b_{11} & b_{12} & b_{13} & b_{14} \\ b_{21} & b_{22} & b_{23} & b_{24} \\ b_{31} & b_{32} & b_{33} & b_{34} \end{bmatrix}$$

(3, N)

(N, 3)

$$O_1 = \begin{bmatrix} (a_{11}b_{11} + a_{12}b_{21} + a_{13}b_{31}) \\ (a_{11}b_{12} + a_{12}b_{22} + a_{13}b_{32}) \\ (a_{11}b_{13} + a_{12}b_{23} + a_{13}b_{33}) \\ (a_{11}b_{14} + a_{12}b_{24} + a_{13}b_{34}) \end{bmatrix}$$

$$\begin{aligned} O_1 &= a_{11}b_{11} + a_{12}b_{21} + a_{13}b_{31} \\ &= \sum_{j=1}^3 a_{1j} b_{j1} \end{aligned}$$

$$O_i = P_{ij} Y_j = \sum_j e^{\frac{q_i^T k_j}{L_i}} Y_j$$

$$O = PV$$

$$\frac{\partial \phi}{\partial V} = P^T \frac{\partial \phi}{\partial O} =$$

$$\frac{\partial V}{\partial O} = P^T \phi \rightarrow P_{phi}$$

$$\frac{\partial \phi}{\partial P} = \frac{\partial \phi}{\partial O} V^T = (\frac{\partial P}{\partial O} V^T)$$

$$Y = XW$$

$$\frac{\partial \phi}{\partial X} = \frac{\partial \phi}{\partial Y} W^T$$

$$\frac{\partial \phi}{\partial W} = X^T \frac{\partial \phi}{\partial X}$$

$$O = PV$$

$$\frac{\partial \phi}{\partial V} = P^T \frac{\partial \phi}{\partial O}$$

$$dV = P^T dO$$

$$\frac{\partial \phi}{\partial P} = \frac{\partial \phi}{\partial O} V^T$$

$$dP = dO V^T$$

$$dV = F \frac{\partial \phi}{\partial O} dO$$

$$dV_j = \sum_i F_{ji} \downarrow dO_i =$$

Scalar  
Vector Multiplication

$$\frac{\partial \phi}{\partial x} = \frac{\partial \phi}{\partial y} \cdot \frac{\partial y}{\partial x}$$

S6257786

negnehyoorder@protonmail.com

Pytorch

$$\frac{\partial \phi}{\partial y} = [dy_{11} \ dy_{12} \ dy_{13} \ dy_{14}]$$

$$\frac{\partial \phi}{\partial x} = \frac{\partial \phi}{\partial y} \times \frac{\partial y}{\partial x}$$

$$x = \begin{bmatrix} x_1 & x_{12} & x_{13} \end{bmatrix}_{(1,3)} \quad \begin{bmatrix} W_{11} & W_{12} & W_{13} & W_{14} \\ W_{21} & W_{22} & W_{23} & W_{24} \\ W_{31} & W_{32} & W_{33} & W_{34} \end{bmatrix}_{(3,4)}$$

$$[(x_1 w_{11} + x_{12} w_{21} + x_{13} w_{31})]$$

Pytro

[N, M] [M, P]

$$\frac{\partial \phi}{\partial x} = \frac{\partial \phi}{\partial y} \cdot \frac{\partial y}{\partial x} = \frac{\partial \phi}{\partial y} \cdot W^T$$

$$\frac{\partial y}{\partial x} = \begin{bmatrix} W_{11} & W_{21} & W_{31} \\ W_{12} & W_{22} & W_{32} \\ W_{13} & W_{23} & W_{33} \\ W_{14} & W_{24} & W_{34} \end{bmatrix} = W^T$$

$$\frac{\partial \phi}{\partial w} = x^T \frac{\partial \phi}{\partial y}$$

$$\frac{\partial \phi}{\partial x} = \frac{\partial \phi}{\partial y} \cdot \frac{\partial y}{\partial x} = \frac{\partial \phi}{\partial y} \cdot W^T \quad [N, D]$$

$$\frac{\partial \phi}{\partial w} = x^T \frac{\partial \phi}{\partial y}$$

gradient through softmax

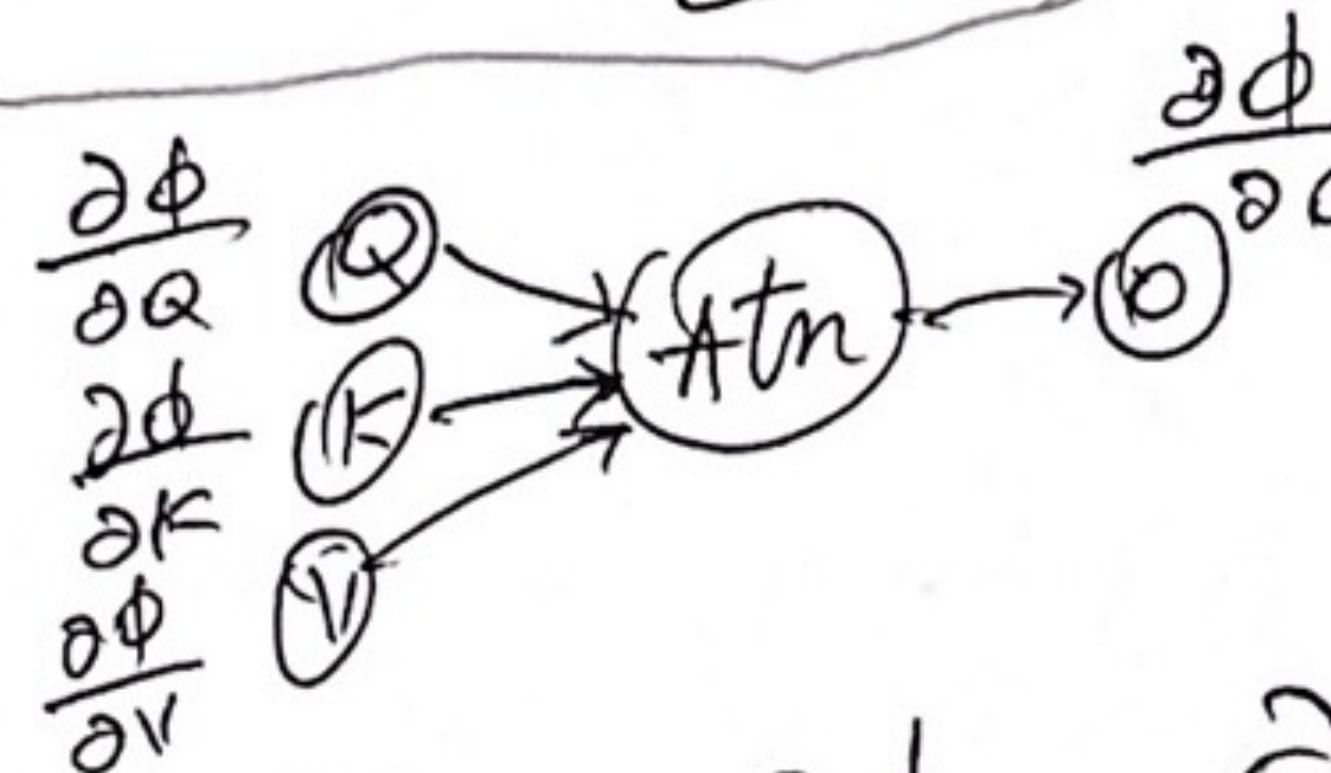
$$g = Q F^T \quad p = \text{Softmax}(s_i) \quad o = p v$$

$$s_i = s[i, :] \in R^N$$

$$p_i = \text{Softmax}(s_i) \in R^N$$

$$\text{Softmax}(p_{ij}) = \frac{e^{s_{ij}}}{\sum_{k=1}^N e^{s_{ik}}}$$

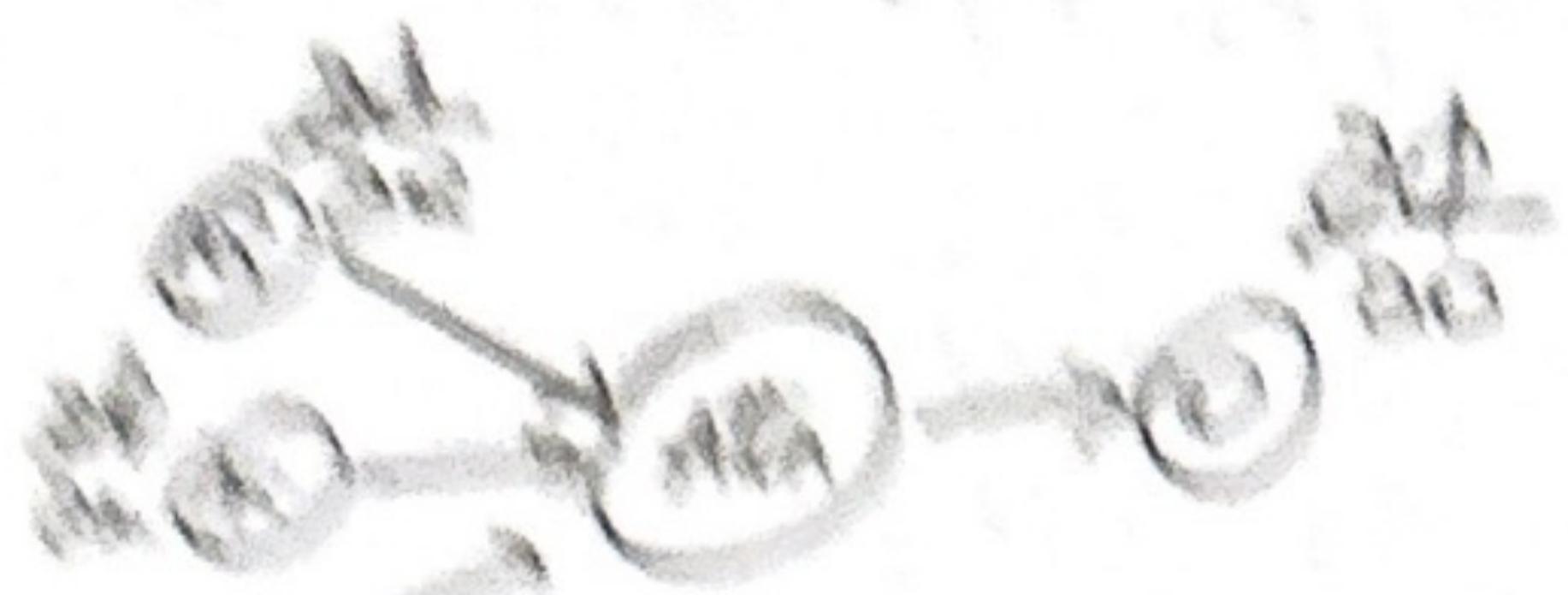
$$\left(\frac{f(x)}{g(x)}\right)' = \frac{f'(x)g(x) - g'(x)f(x)}{(g(x))^2}$$



$$\frac{\partial \phi}{\partial s_i} = \frac{\partial \phi}{\partial p_i} \cdot \frac{\partial p_i}{\partial s_i}$$

$$\frac{\partial p_i}{\partial s_{ik}} = \frac{\partial}{\partial s_{ik}} \left[ \frac{e^{s_{ij}}}{\sum_{k=1}^N e^{s_{ik}}} \right] =$$

gradient through softmax  
 $\hat{y} = \text{softmax}(x) \in \mathbb{R}^N$



$$\begin{aligned} \text{softmax function: } & y_i = \text{softmax}(x_i) \\ & y_i = \frac{e^{x_i}}{\sum_{j=1}^N e^{x_j}} \end{aligned}$$

$$\frac{\partial \hat{y}_i}{\partial x_j}$$

$$\frac{\partial \hat{y}_i}{\partial x_j} = \frac{\partial \text{softmax}(x_i)}{\partial x_j} = \frac{\partial}{\partial x_j} \left[ \frac{e^{x_i}}{\sum_{l=1}^N e^{x_l}} \right] = \frac{e^{x_i}}{\sum_{l=1}^N e^{x_l}} \cdot \frac{-\sum_{l=1}^N e^{x_l} \cdot e^{x_j}}{(\sum_{l=1}^N e^{x_l})^2} = \frac{e^{x_i}}{\sum_{l=1}^N e^{x_l}} \cdot \frac{-e^{x_j}}{\sum_{l=1}^N e^{x_l}}$$

$$\begin{aligned} g &= \mathbb{R}^{G, 13} \in \mathbb{R}^H \\ h &= \text{softmax}(g_i) \in \mathbb{R}^H \quad g_{i,j} = g_{\max} \\ \text{softmax}(p_{ij}) &= \frac{e^{g_{ij}}}{\sum_{l=1}^N e^{g_{il}}} = g_{\max} \end{aligned}$$

$$\frac{\partial \hat{y}_{ij}}{\partial g_{ik}} = \frac{\partial \phi}{\partial P} \cdot \frac{\partial P}{\partial g_{ik}} \quad \text{if } i=k \rightarrow e^{g_{ik}} \left( \sum_{l=1}^N e^{g_{il}} \right) = e^{g_{ik}} e^{g_{ik}}$$

$$e^{g_{ik}} \left( \sum_{l=1}^N e^{g_{il}} - e^{g_{ik}} \right)$$

$$\frac{e^{g_{ik}}}{\sum_{l=1}^N e^{g_{il}}} \times \frac{\sum_{l=1}^N e^{g_{il}} \cdot e^{g_{ik}}}{\sum_{l=1}^N e^{g_{il}}} \downarrow p_{ij} \quad (1-p_{ik})$$

$$\text{softmax}(x_i) = \frac{\exp(x_i - m_i)}{\sum_j \exp(x_j - m_j)}$$

$$\begin{aligned} \text{softmax}(x_i) &= \exp(x_i - m_i) \\ &= \frac{\exp(x_i - m_i)}{\exp(m_i)} \end{aligned}$$

$$\exp(a+b) = \exp(a) \cdot \exp(b)$$

$$\exp(a-b) = \frac{\exp(a)}{\exp(b)}$$

$$(9) \quad \exp$$

# THE ONLINE SOFTMAX

$$m_0 = \begin{bmatrix} -\infty \\ 0 \\ 0 \end{bmatrix}$$

**STEP 1**  
 $m_i = \max(\max(Q_i, k_i^T), m_0)$

$$l_0 = -\infty$$

$$l_0 = 0$$

for  $i = 1 \dots N$

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

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

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

STEP 2

$$\frac{m_2}{S_2} = \max(\max(Q_1, k_2^T), m_1) \quad O = \text{diag}(\exp(m_0 - m_1)) O_0 \rho_{k_2}^{m_2}$$

$$P_{12} = Q_1 k_2^T$$

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

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

$$O_1 = \begin{bmatrix} O_{11} & O_{12} & \cdots & O_{1N} \\ O_{21} & O_{22} & \cdots & O_{2N} \end{bmatrix} \times$$

$$\begin{bmatrix} k_4^{(1)} & & & \\ k_4^{(2)} & \ddots & & \\ & \ddots & \ddots & \\ & & \ddots & k_4^{(N)} \end{bmatrix} = \begin{bmatrix} C_1 & \cdots & C_N \\ C_1 & \cdots & C_N \\ \vdots & \ddots & \vdots \\ 0 & \ddots & 0 \end{bmatrix} \times \begin{bmatrix} k_4^{(1)} & & & \\ 0 & k_4^{(2)} & & \\ & 0 & \ddots & \\ & & \ddots & k_4^{(N)} \end{bmatrix}$$

$$= \begin{bmatrix} C_1 & \cdots & C_N \\ C_1 & \cdots & C_N \\ \vdots & \ddots & \vdots \\ 0 & \ddots & 0 \end{bmatrix} \times \begin{bmatrix} O_1 \exp(m_1 - m_2) + 0 & O_2 \exp(m_1 - m_2) + 0 & \cdots & O_N \exp(m_1 - m_2) + 0 \\ O_1 + O_2 \exp(m_1 - m_2) & O_2 + O_3 \exp(m_1 - m_2) & \cdots & O_N + O_{N-1} \exp(m_1 - m_2) \end{bmatrix}$$

$$O_5 = [\text{diag}(L_1)]^{-1} O_4$$

卷之三

卷之三

六

卷之二

(12814)

A rectangular blue ink stamp is positioned vertically on the left side of the page. The stamp contains the characters '卷之二' (Volume 2) in vertical columns at the top, followed by '三' (Three) at the bottom. Above the stamp, there is a faint, horizontal red seal impression.

826 = 8122

卷之三

11

	381	..
	381	.
	381	..

$$OPT_b = \text{C}_2 \times \text{C}_{128} = (2, 2)$$

QKT	QKT	QKT	QKT
QKT	QKT	QKT	QKT
QKT	QKT	QKT	QKT
QKT	QKT	QKT	QKT
QKT	QKT	QKT	QKT

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

$$S_{11} = (2 \times 128) \times (128 \times 2) = (2, 2)$$

$S =$

$Q_1 K_1^T$	$Q_1 K_2^T$	$Q_1 K_3^T$	$Q_1 K_4^T$
$Q_2 K_1^T$	$Q_2 K_2^T$	$Q_2 K_3^T$	$Q_2 K_4^T$
$Q_3 K_1^T$	$Q_3 K_2^T$	$Q_3 K_3^T$	$Q_3 K_4^T$
$Q_4 K_1^T$	$Q_4 K_2^T$	$Q_4 K_3^T$	$Q_4 K_4^T$

$$\xrightarrow{\text{SOFTMAX}} P =$$

$$= \exp(x_i - x_{\max})$$

$P_{1,1}$	$P_{1,2}$	$P_{1,3}$	$P_{1,4}$
$P_{2,1}$	$P_{2,2}$	$P_{2,3}$	$P_{2,4}$
$P_{3,1}$	$P_{3,2}$	$P_{3,3}$	$P_{3,4}$
$P_{4,1}$	$P_{4,2}$	$P_{4,3}$	$P_{4,4}$

ORIGINAL = (8, 2)

BLOCK = (4, 4)

$$\text{SOFTMAX} = \frac{\exp(x_i - x_{\max})}{\sum_{j=1}^n \exp(x_j - x_{\max})}$$

$S_{11}$

$$= \begin{bmatrix} Q_1 & K_1^T \\ Q_2 & K_2^T \\ Q_3 & K_3^T \\ Q_4 & K_4^T \end{bmatrix} \xrightarrow{\text{SOFTMAX}} \begin{bmatrix} \exp(a-a) & \exp(b-a) \\ \exp(c-a) & \exp(d-a) \end{bmatrix}$$

$$S_{11} = \exp \left[ s_{ij} - \text{row max}(s_{ij}) \right]$$

LOCAL MAX

PSEUDO CODE

\* Y

1 ... 128	128
1 ... 128	128

$\{v_1, v_2, v_3, v_4\}$

$$C_1(128) \times (128 \times 2) = C_2(2)$$

$$\rightarrow S_{11} \\ \cancel{X \otimes X \otimes X} \times (128 \times 2) = (2, 2)$$

$$Q_1 K_1^T \quad Q_1 K_2^T \quad Q_1 K_3^T \quad Q_1 K_4^T$$

$$S =$$

$$Q_2 K_1^T \quad Q_2 K_2^T \quad Q_2 K_3^T \quad Q_2 K_4^T$$

$$Q_3 K_1^T \quad Q_3 K_2^T \quad Q_3 K_3^T \quad Q_3 K_4^T$$

$$Q_4 K_1^T \quad Q_4 K_2^T \quad Q_4 K_3^T \quad Q_4 K_4^T$$

$$Q_1 K_1^T \quad Q_1 K_2^T \quad Q_1 K_3^T \quad Q_1 K_4^T$$

SOFTMAX\*

$$\frac{\exp(x_i - x_{\max})}{\sum_{j=1}^n \exp(x_j - x_{\max})}$$

$P_{11}$	$P_{12}$	$P_{13}$	$P_{14}$
$P_{21}$	$P_{22}$	$P_{23}$	$P_{24}$
$P_{31}$	$P_{32}$	$P_{33}$	$P_{34}$
$P_{41}$	$P_{42}$	$P_{43}$	$P_{44}$

$\Rightarrow P =$

$$\exp(x_i - x_{\max})$$

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

$$S_{11} = \begin{bmatrix} Q_1 & K_1^T \\ Q_2 & K_2^T \\ Q_3 & K_3^T \\ Q_4 & K_4^T \end{bmatrix} \rightarrow \begin{bmatrix} \exp(a-a) & \exp(b-a) \\ \exp(c-a) & \exp(d-a) \end{bmatrix}$$

LOGIC MAX

$S_{11}^*$  ( $S_{11}$ )

$$= \exp \left[ S_{11} - \text{LOGIC MAX}(S_{11}) \right]$$

PSUDO CODE

FOR EACH  $i$  &  $K_i$   
 $O_i = \text{ZERO}(2, n)$

FOR EACH  $j$  &  $K_j$

$$P_{ij} = S(Q_i K_j^T)$$

$$O_i \leftarrow O_i + P_{ij} V_j$$

$$(g, 128)$$

$$(q, 128)$$

1	...	128	
1	...	128	
3	$V_1$		
...	...	...	
1	...	128	

END FOR

END FOR

$$V = \begin{bmatrix} 1 & \dots & 1 & 2 & 3 \\ 1 & \dots & 1 & 2 & 8 \end{bmatrix}^T$$

$$\sum_{j=1}^4 V_j = \sum_{j=1}^4 \left( (Q_1 K_1^T) V_1 + (Q_2 K_2^T) V_2 + (Q_3 K_3^T) V_3 + (Q_4 K_4^T) V_4 \right)$$

$$O = (Q_1 K_1^T) V_1 + (Q_2 K_2^T) V_2 + (Q_3 K_3^T) V_3 + (Q_4 K_4^T) V_4$$

$$(Q_1 K_1^T) V_1 + (Q_3 K_3^T) V_2 + (Q_2 K_2^T) V_3 + (Q_4 K_4^T) V_4$$

$$(Q_1 K_1^T) V_1 + (Q_4 K_3^T) V_2 + (Q_3 K_2^T) V_3 + (Q_2 K_4^T) V_4$$

$$(Q_1 K_1^T) V_1 + (Q_4 K_2^T) V_2 + (Q_3 K_3^T) V_3 + (Q_2 K_1^T) V_4$$

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

EACH BLOCK OF BUFFER  
OUTPUT "O" IS MADE FOR  
TWO ROWS!

### PSEUDOCODE

FOR EACH BLOCK  $Q_i$

$O_i = \text{ZEROS}(2, 128)$  // OUTPUT INITIALLY MADE OF ZEROS

FOR EACH BLOCK  $K_j$

$O_i^{ij} = O_i + (Q_i K_j^T) V_j$

END FOR

END FOR