$$\text{Attention}(Q, k, V) = \text{Softmax}\left(\frac{Q k^T}{\sqrt{d_k}}\right) V$$

$$\text{head}_i = \left(\text{Attention}(Q W_i^Q, k W_i^k, V W_i^v)\right)$$

GPU

more I/o bound on HBM.

$\checkmark$

HBM $\searrow$ shared

alt : compute attention in shared memory

(which is smaller)

(Global)

Blocked computation.

⊙ numerical instability of softmax due to $e^x$

**Safe softmax**

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

"break in" some constant.

$k = \max_i (x_i)$

## Pseudocode

$$m_0 = -\infty$$

for $i = 1$ to $N$

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

$l_0 = 0$

for $j = 1$ to $N$

$$l_j = l_{j-1} + e^{x_j - m_N} \quad \longrightarrow \text{②}$$

for $k = 1$ to $N$,

$$x_k \leftarrow \frac{e^{x_k - m_N}}{l_N} \quad \longrightarrow \text{③}$$

① accessing vector three times.

(Three sequential operations)

→ for given query, fill its row in this way.

→ To fuse the computation ①, and ②, use local maxima instead of global maxima $(m_N)$ and fix it on fly using correction factor.

(m_j)

**new pseudocode**

$$m_0 = -\infty, \quad l_0 = 0 \qquad \underline{\text{Online softmax}}$$

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 $N$

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

## Proof by Induction

① Prove it holds for $N = 1$

② Given it holds for $N$, does it hold for $N+1$?

fixed vector

# Block Matrix Multiplication

$$
Q_1 \begin{bmatrix} Q_1 \\ Q_2 \\ Q_3 \\ Q_4 \end{bmatrix}
\quad \times \quad
\begin{bmatrix} & k_1^T & k_2^T & \cdots & k_4^T \\ & K_1^T & K_2^T & - & K_4^T \end{bmatrix}
$$

$$Q \qquad\qquad\qquad\qquad K^T$$

Original $(8, 128)$ $\qquad$ Original $(128, 8)$

$\downarrow$ to $\qquad\qquad\qquad \downarrow$ to

4 blocks of size $(2, 128)$ $\qquad$ 4 blocks of size $(128, 2)$

$\downarrow$ $\qquad\qquad\qquad\qquad \downarrow$

$(4, 128)$ $\qquad\qquad\qquad\qquad (128, 4)$ $\qquad\qquad \checkmark$

$$
S = \begin{bmatrix}
Q_1 k_1^T & Q_1 k_2^T & - & - & Q_8 k_4^T \\
Q_2 k_1^T & & | & & | \\
| & & | & & | \\
| & - & - & - & Q_4 k_4^T
\end{bmatrix}
\qquad
\begin{matrix} V_1 \\ V_2 \\ V_4 \end{matrix}
\begin{bmatrix} V_1 \\ V_2 \\ V_3 \\ V_4 \end{bmatrix}
$$

Org : $(8, 128)$
block : $(4, 128)$

each block
is made up of
two acess.

shape $(8, 8)$ $\longleftarrow$ not scalars
but $2 \times 2$ matrix.

$(S \cdot V)$

$$
\text{Output} = \begin{bmatrix}
\underbrace{(Q_1 k_1^T) V_1}_{S_1} + \underbrace{(Q_1 k_2^T) V_2}_{S_2} + \underbrace{(Q_1 k_3^T) V_3}_{S_3} + \underbrace{(Q_1 k_4^T) V_4}_{S_4} \\
| \qquad\qquad | \qquad\qquad | \\
|
\end{bmatrix}
$$

## Pseudo code

For each block $Q_i$:

$$O_i = zeroes(2, 128)$$

For each block $k_j$:

$$O_i \leftarrow O_i + (Q_i k_j^T) V_j$$

End for

End for

But with (normalization) softmax, (For block)

### Initialization

$$m_o = \begin{bmatrix} -\infty \\ -\infty \end{bmatrix}, \quad l_o = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

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

### Step 1

$$m_1 = \max(\text{row max}(Q_1 k_1^T), m_o)$$

$$S_1 = Q_1 k_1^T$$

$$l_1 = \text{row sum}\left[\exp(S_1 - m_1)\right] + l_o \cdot \exp(m_o - m_1)$$

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

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

<u>Step 2</u>

$$m_2 = \max\left(\text{row max}\left(Q_1 k_2^T\right),\ m_1\right)$$

$$S_2 = Q_1 k_2^T$$

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

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

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

$$\begin{bmatrix} \exp(m_1 - m_2)_1 & 0 \\ 0 & \exp(m_1 - m_2)_2 \end{bmatrix} \times \begin{bmatrix} O_{11} & O_{12} & O_{13} & - & - \\ O_{21} & O_{22} & O_{23} & - & - \end{bmatrix}$$

$$= \begin{bmatrix} O_{11}\, \exp(m_1 - m_2)_1 & - & - \\ O_{21}\, \exp(m_1 - m_2)_2 & - & - \end{bmatrix}$$

Correcting each row from previous block error.

And so on until the last step. Then at the end we apply "$l$" normalization factor.

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

because.

$$\left(\begin{bmatrix} l_4^{(1)} & 0 \\ 0 & l_4^{(2)} \end{bmatrix}\right)^{-1} = \begin{bmatrix} \frac{1}{l_4^{(1)}} & 0 \\ 0 & \frac{1}{l_4^{(2)}} \end{bmatrix}$$

GPU : Hardware
CUDA : Software for GPU

## Cuda Programming Example : "Vector Addition"

Allocate memory in cuda and copy A and B to those
from the host m/c.
Execute the kernel $\Rightarrow$ launches N threads.
                                     ( N parallel operations).

Each thread get its ID and we're supposed
to define what we want to run in each.

control
units are
expensive part

(launches near two-multiple)
Ex: 34 $\rightarrow$ 64

Group of threads run (by) same
instruction, but data (may) differ. (same control unit)
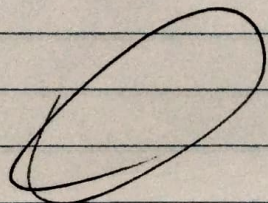    $\gg$ SINGLE INSTRUCTION MULTIPLE DATA
    $\gg$ " " " " " " " " THREADS

Blocks of threads .     $\hookrightarrow$ we'll have $\lceil \frac{N}{32} \rceil$ blocks.

element_id = B_id $\times$ Block_size + T_id    no. of threads in each block
                                         (or) Blocksize
$\hookrightarrow$ we are specifying this, not cuda.

- Access to DRAM is slow.

but access to shared memory is very fast.

CPU → GPU           Use shared
(global mem    ⤷    memory
which is slower)      copy to it
                      whenever needed

⇒ Each kernel works              (shared by all threads
with one "Query" block            in a same block)
and iterates through "key" blocks
⇒ Each "head" computes attention independently
⇒ Each Q blocks independently

SEQ_LEN / BLOCK_SIZE_Q

max parallel programs we can run. ⤷ BATCH_SIZE × NUM_HEADS × (~~BLOCK_SIZE_Q~~)

num of block we can divide our Query sequence into