# FlashAttention in Triton
Tuesday, March 18, 2025   5:17 PM

Flash Attention 1:
- Computations are on GPU
- Accessing HBM (GPU Global memory) or CPU DRAM is MUCH slower than accessing GPU Shared Memory
- Problem: Calculating Attention is I/O Bound -> as it keeps needing to access GPU Global Memory
- Solution: Computing Attention INSIDE Shared Memory itself, avoid repeated Global memory accesses
- Consequence: Computation of Attention is MUCH CLOSER to GPU cores, Bandwidth bottleneck addressed
- How we'll do it:
  ○ Split Attention input into smaller blocks that can FIT INSIDE Shared Memory using Tiling/Block Matmul
  ○ Compute Attention using K, Q, V for EACH Block (of the entire Matrix) in the Shared Memory (One Block at a time, Shared Memory can fit only 1 Block), then COPY this back to Global Memory

FlashAttention1 vs FlashAttention2:
- Structure of 2 'for' loops in FA1:
  ○ 1st loop: K (Keys)
  ○ 2nd loop: Q (Queries)
- Structure of 2 'for' loops in FA2:
  ○ 1st loop: Q (Queries)
  ○ 2nd loop: K (Keys)

Online Softmax algorithm (by hand):



$m_0 = -\infty$
$l_0 = 0$
for $i = 1$ to $N$
$\quad m_i = \max(m_{i-1}, x_i)$
$\quad l_i = l_{i-1} \cdot e^{m_{i-1} - m_i} + e^{x_i - m_i}$
for $u = 1$ to $N$
$\quad x_u \leftarrow \dfrac{e^{x_u - m_N}}{l_N}$

We want to prove that at the end of this loop:
$m_N = \max(x_i) = x_{max}$
$l_N = \sum_{i=1}^{N} e^{x_i - x_{max}}$

Correction Factor

If:
- Current Max > Previous Max: Correction Factor corrects
- If Current Max = Previous Max: Correction Factor becomes 1, no correction takes place

Online Softmax algorithm (implemented in code):



STEP 2.

$m_2 = \max(\text{rowmax}(Q_1 k_2^T), m_1)$
$S_2 = Q_1 k_2^T$
$l_2 = \text{row sum}[\exp(S_2 - m_2)] + l_1 \cdot \exp(m_1 - m_2)$
$P_{12} = \exp(S_2 - m_2)$
$O_2 = \text{diag}(\exp(m_1 - m_2))O_1 + P_{12}V_2$

new max
previous max
Current row max
Accumulate Normalization Factor terms at each step
We divide by the Accumulated Normalization Factor at the end to get the final Output

Update Output at each step. For 1st step, output = PV (only PV)

Applying Softmax_Star() on each Q.K^T
- Q and K are matrix BLOCKS (tiles/blocks from tiling/blocking)

Note on Local Memory and Global Memory:
- 'Local', 'Global' in scope and present in Device Memory (not actually physically partitioned in Device Memory)

CUDA Mem Model maps to the Nvidia GPU Mem Hierarchy

Nvidia GPU Memory Hierarchy:     On-Chip
Registers [Fastest]
L1 Data Cache + Shared Memory
L2 Cache                         Off-Chip
Device Memory (HBM) [Slowest]

CUDA Memory Model (how CUDA perceives GPU memory):
Registers & Local Memory (Private to its corresponding Thread)
                    V
Shared Memory (Every Thread in the Thread Block has access)
                    V
Global Memory (Every Thread in all Grids have access)

On-chip Shared Memory (is user-controlled L1 Cache, only Threads within the SAME Thread Block can access)

Off-chip Global Memory (shared across multiple Thread Blocks)

goes in
goes in



2. Data moved to compute units & SRAM for computation
1. inputs start out in HBM (GPU memory)
3. Output written back to HBM

Streaming Multiprocessors
Slow Data Transfer
HBM

SRAM: 19 TB/s (20 MB)
GPU HBM: 1.5 TB/s (40 GB)
Main Memory (CPU DRAM): 12.8 GB/s (>1 TB)

Memory Hierarchy with Bandwidth & Memory Size

The biggest cost is in moving data
Standard implementation requires repeated R/W from slow GPU memory

FlashAttention2 (using the Online Softmax algo above)



**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 $\mathbf{K}, \mathbf{V}$ in to $T_c = \left\lceil \frac{N}{B_c} \right\rceil$ 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** $1 \le i \le T_r$ **do**  ← FOR EACH $Q_i$ BLOCK
4: $\quad$ Load $\mathbf{Q}_i$ from HBM to on-chip SRAM.
5: $\quad$ 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: $\quad$ **for** $1 \le j \le T_c$ **do**  ← FOR EACH $K_j$ BLOCK
7: $\quad\quad$ Load $\mathbf{K}_j, \mathbf{V}_j$ from HBM to on-chip SRAM.
8: $\quad\quad$ On chip, compute $\mathbf{S}_i^{(j)} = \mathbf{Q}_i \mathbf{K}_j^T \in \mathbb{R}^{B_r \times B_c}$.
9: $\quad\quad$ 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: $\quad\quad$ 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: $\quad$ **end for**
12: $\quad$ On chip, compute $\mathbf{O}_i = \text{diag}(\ell_i^{(T_c)})^{-1} \mathbf{O}_i^{(T_c)}$.
13: $\quad$ On chip, compute $L_i = m_i^{(T_c)} + \log(\ell_i^{(T_c)})$.
14: $\quad$ Write $\mathbf{O}_i$ to HBM as the $i$-th block of $\mathbf{O}$.
15: $\quad$ Write $L_i$ to HBM as the $i$-th block of $L$.
16: **end for**
17: Return the output $\mathbf{O}$ and the logsumexp $L$.

CUDA:
- We have to tell each Thread what to do, ie software engineer determines mapping between Thread and its Data Elements (CUDA knows only # of Blocks & # of Threads [in each Block])
- When running 'N'-threads in parallel (done on each core), when we ask CUDA to run the N-threads in parallel, it will:
  ○ Allocate N threads
  ○ Assign a unique 'Thread Identifier' to each thread
  ○ Run the threads according to the Thread Identifier assigned to it



Why we have the 'if (i < N)' in cuda_vector_add_simple() method:
  ○ Threads are launched by CUDA in multiples of 32 (ie, launched as a warp, 1 warp has 32 threads)
    Eg: if no of threads intended by software engineer is 34 (ie, 34 elements in vector, 32+2), CUDA launches 64 threads
So, supposing we want to launch 8 threads (as we have 8 elements in vector), we have to tell the remaining 24 threads to NOT do anything.
- The N threads being launched have the same Control Unit. This means: they have run the same program, but on different data elements (SIMT programming model)
- Control Divergence:
  ○ In earlier case, 8 threads were launched, meaning remaining 24 threads were also launched, but remain IDLE as they don't do any work
  ○ This is called Control Divergence and causes lesser throughput, which software engineer should aim to minimize
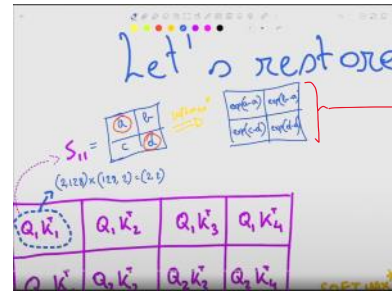


Eg: N = 8, Block Size = 2, GPU has 4 cores
So, # of Blocks = 8 / 2 = 4

Expression to access each element in N:
element_id = B_id * Block_size + thread_id

Num - Blocks = $\dfrac{N}{2 \to \text{Block-size}}$ = 4

Takes 0 or 1
Block ID: 0, 1, 2 or 3

General formula to calculate # of Blocks:
# of Blocks = Total # of Data Elements (N) / Block Size

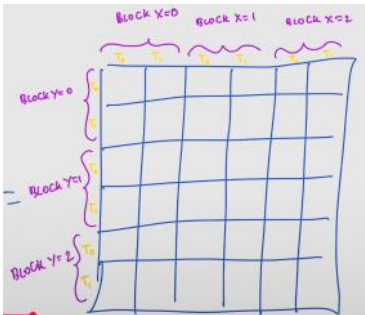# of Threads per block (considering 1 Thread per Data Element)

DOUBTS
1. Why does each Block have 2 Threads here?





x_max does over each row or over entire matrix? => over each row

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

| $Q_1 k_1^T$ | $Q_1 k_2^T$ | $Q_1 k_3^T$ | $Q_1 k_4^T$ |
| $Q_2 k_1^T$ | $Q_2 k_2^T$ | $Q_2 k_3^T$ | $Q_2 k_4^T$ |

Defining launch grid for 2D matrix (like on the left):
- Defining # of Block Rows and # of Block Columns

```
// Define the launch grid
int num_blocks_ROWS = (NUM_ROWS + ROWS_block_size - 1) / ROWS_block_size;
int num_blocks_COLS = (NUM_COLS + COLS_block_size - 1) / COLS_block_size;
```

- # of Blocks we have & # of Threads we have:
```
dim3 grid(num_blocks_COLS, num_blocks_ROWS, 1);
dim3 block(COLS_block_size, ROWS_block_size, 1);
```
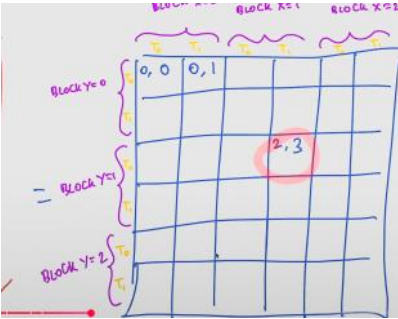
# of Element Columns (decided from total Elements 'N')
# of Element Columns we set per Column Block
Fancy way of writing ceil(NUM_COLS / COLS_block_size)
# of Blocks
# of Threads



```
void test_matrix_add(int NUM_ROWS, int NUM_COLS, int ROWS_block_size, int COLS_block_size)
    for (int i = 0; i < NUM_ROWS; i++)
        for (int j = 0; j < NUM_COLS; j++)
            size_t index = static_cast<size_t>(i) * NUM_COLS + j;
            A[index] = rand() % 100;
            B[index] = rand() % 100;

    // Allocate device memory for a
    CUDA_CHECK(cudaMalloc((void **)&d_A, sizeof(EL_TYPE) * NUM_ROWS * NUM_COLS));
```

4 threads per block or 2 threads per block here? Video says 2 threads per block

Relationship between warp and block (depends on threads per block)

- **Fewer than 32 threads per block:**
  Even if a block has fewer than 32 threads (say, 4 threads), the GPU still allocates a full warp (32 threads) for that block, but only the threads corresponding to your block are active while the remaining lanes remain inactive.
- **Exactly 32 threads per block:**
  The block perfectly fits into one warp.
- **More than 32 threads per block:**
  The block is divided into multiple warps. For example, if a block has 64 threads, it will be split into 2 warps.

Defines the grid

- Element ID expressions:

```
int row_index = blockIdx.y + blockDim.y + threadIdx.y;
int col_index = blockIdx.x + blockDim.x + threadIdx.x;
```



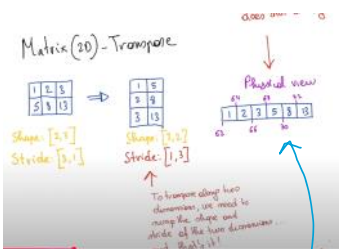Tensor Layout (in CPU and GPU memory):
- Row-major layout:



Stride: tells how many elements to be skipped (including present element) to:
  ○ Go to the next row dimension (eg: 3 here)
  ○ Go to the next column dimension (eg: 1 here)
In summary:
  ○ we skip 3 elements to go to next row
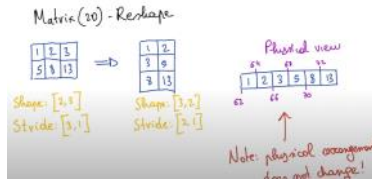  ○ we skip 1 element to go to next element in same row



Swapping Stride transposes matrix

Tensor Reshaping

In conclusion,
Stride is used for:
  ○ Indexing Tensor (pointer to starting element of physical location)
  ○ Allows reshaping Tensor for free
  ○ Allows Transposing: by swapping dims

Triton:
- Unlike CUDA (where we need to define how many Blocks and how many Threads per Block), Triton only needs # of Blocks (it decides # of Threads per Block on its own) We tell what each Thread should do.

Vector addition example



Defining triton grid:
- # of blocks = ceil(n_elements, block_size)

But here, block size is not # of Threads per block (unlike CUDA)
It is: # of Data Elements each block processes (in CUDA, assumption is: 1 thread per data element in the same cycle, that's not the assumption here)

- Unlike Pytorch, Triton does not take entire vectors just as they are (full matrices, broadcasting etc operations). It takes the POINTER to the 1st Data Element of that vector in memory. We have to compute the indices for the vectors from the Pointer (to the 1st data element).



Pointer to 1st element of vector 'x'

Pointer to 1st element of vector 'y'

Pointer to 1st element of vector 'output'

The Block_ID (CUDA) equivalent for triton is program_id

FlashAttention2 forward pass:



PyTorch operations:
- Always derived from torch.autograd.Function
- Operations like Softmax, ReLU, etc. are always implemented as torch.autograd.Function

Causal vs Non-causal Attention:
- Causal (Masked):
  ○ Token can ONLY attend to itself and tokens BEFORE it in the sequence, NOT any tokens AFTER it

**FlashAttention2 forward pass:**



Algorithm 1 FlashAttention-2 forward pass

Require: Matrices $Q, K, V \in \mathbb{R}^{N \times d}$ in HBM, block sizes $B_c$, $B_r$.
1: Divide $Q$ into $T_r = \lceil \frac{N}{B_r} \rceil$ blocks $Q_1, ..., Q_{T_r}$ of size $B_r \times d$ each, and divide $K, V$ in to $T_c = \lceil \frac{N}{B_c} \rceil$ blocks $K_1, ..., K_{T_c}$ and $V_1, ..., V_{T_c}$, of size $B_c \times d$ each.
2: Divide the output $O \in \mathbb{R}^{N \times d}$ into $T_r$ blocks $O_1, ..., O_{T_r}$ of size $B_r \times d$ each, and divide the logsumexp $L$ into $T_r$ blocks $L_1, ..., L_{T_r}$ of size $B_r$ each.
3: for $1 \le i \le T_r$ do $\leftarrow$ first, EACH $B_r$ BLOCK.
4:   Load $Q_i$ from HBM to on-chip SRAM.
5:   On chip, initialize $O_i^{(0)} = (0)_{B_r \times d}$, $\ell_i^{(0)} = (0)_{B_r}$, $m_i^{(0)} = (-\infty)_{B_r} \in \mathbb{R}^{B_r}$.
6:   for $1 \le j \le i$ do $\leftarrow$ THEN, EACH $B_c$ BLOCK.
7:     Load $K_j, V_j$ from HBM to on-chip SRAM.
8:     On chip, compute $S_i^{(j)} = Q_i K_j^T \in \mathbb{R}^{B_r \times B_c}$.
9:     On chip, compute $m_i^{(j)} = \max(m_i^{(j-1)}, \text{rowmax}(S_i^{(j)})) \in \mathbb{R}^{B_r}$, $\tilde{P}_i^{(j)} = \exp(S_i^{(j)} - m_i^{(j)}) \in \mathbb{R}^{B_r \times B_c}$ (pointwise), $\ell_i^{(j)} = e^{m_i^{(j-1)} - m_i^{(j)}} \ell_i^{(j-1)} + \text{rowsum}(\tilde{P}_i^{(j)}) \in \mathbb{R}^{B_r}$.
10:    On chip, compute $O_i^{(j)} = \text{diag}(e^{m_i^{(j-1)} - m_i^{(j)}})^{-1} O_i^{(j-1)} + \tilde{P}_i^{(j)} V_j$.
11:   end for
12:   On chip, compute $O_i = \text{diag}(\ell_i^{(T_r)})^{-1} O_i^{(T_r)}$.
13:   On chip, compute $L_i = m_i^{(T_r)} + \log(\ell_i^{(T_r)})$.
14:   Write $O_i$ to HBM as the $i$-th block of $O$.
15:   Write $L_i$ to HBM as the $i$-th block of $L$.
16: end for
17: Return the output $O$ and the logsumexp $L$.

- Parallelizations that can be made in this algo:
  1. Q (Queries) loop [1st for loop, outermost]:
     i. 1 kernel works on 1 Q iteration, and inside it works on all K iterations.
     ii. Next kernel works on next Q iteration, and inside it on all K iterations.
     iii. And so on
  2. Above algo is for a single Head and a single Sequence:
     i. Can parallelize for each Sequence
     ii. Can parallelize for each Head in each Sequence (as each Sequence can have multiple Heads)
- Effectively, what we're doing:
  Parallelize each Sequence (in a Batch of Sequences) ->
  inside each sequence, parallelize each Head ->
  inside each Head, parallelize each Query Block (Q)
- So, total # of triton programs (in triton, the name for kernel is program) running in parallel = BATCH_SIZE * NUM_OF_HEADS * (SEQ_LEN / BLOCK_SIZE_Q)

**PyTorch operations:**
- Always derived from torch.autograd.Function
- Operations like Softmax, ReLU, etc. are always Implemented as torch.autograd.Function
- The operations should provide two methods:
  1. forward(): for forward pass, should compute Outputs of the operation
  2. backward(): for backward pass, should compute gradients of Loss wrt Inputs of the operation

During backward pass, we need to softmax AGAIN, for that we need:
1. Normalization Factor
2. Max over rows term
As we don't want to compute these AGAIN, we save both terms during forward pass into L_i,
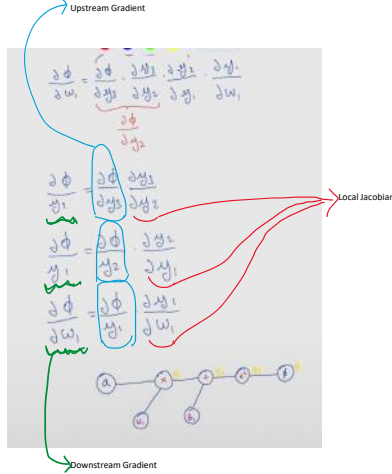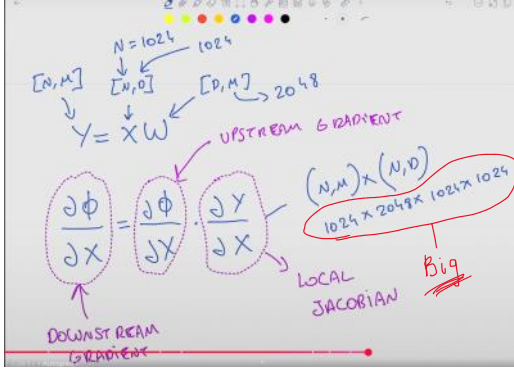And use them during backward pass

Go over Attention paper to clarify dimensions for all

**Umar's implementation:**
1. Not implemented in FP8, FP8 is faster actually on new GPUs
2. FlashAttention backward pass implemented for Causal AND Non-causal Attention, while Triton website only gives Causal
3. Explicit use of Softmax scale done where needed
4. Tips:
   a. If want to learn Triton optimizations: check Triton documentation

**Causal vs Non-causal Attention:**
- Causal (Masked):
  o Token can ONLY attend to itself and tokens BEFORE it in the sequence, NOT any tokens AFTER it
  o Prevents cheating in generative models (future tokens shouldn't inform current predictions)
  o Usage: Autoregressive settings (eg: language generation)
  o Math: For token at index 'i', attention to tokens at index 'j' > 'i' is zeroed out. Eg: if you have Q.K^T, you mask out the upper-right region (where j > i)
- Non-causal (Bidirectional):
  o Lets EVERY token attend to EVERY position in the sequence, including future ones
  o Usage: machine translation (encoder), BERT-style, not generating text one token at a time but learning bidirectional context
  o Math: no masking applied

---

**FlashAttention2 Backward pass:**
**What's needed:**
1. Normalization Factor, Row Max (we already have this from the Forward Pass) ''m_i + log(l_i)'' term
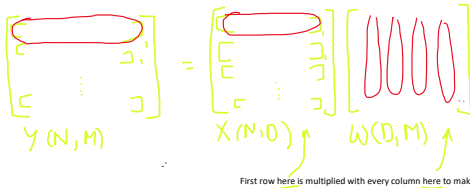2. PyTorch's Autograd
3. How to compute Gradients, Jacobians

**Jacobians and Gradients**



Upstream Gradient

Local Jacobian

Downstream Gradient

Algorithm 2 FlashAttention-2 Backward Pass

Require: Matrices $Q, K, V, O, dO \in \mathbb{R}^{N \times d}$ in HBM, vector $L \in \mathbb{R}^N$ in HBM, block sizes $B_c$, $B_r$.
1: Divide $Q$ into $T_r = \lceil \frac{N}{B_r} \rceil$ blocks $Q_1, ..., Q_{T_r}$ of size $B_r \times d$ each, and divide $K, V$ in to $T_c = \lceil \frac{N}{B_c} \rceil$ blocks $K_1, ..., K_{T_c}$ and $V_1, ..., V_{T_c}$, of size $B_c \times d$ each.
2: Divide $O$ into $T_r$ blocks $O_1, ..., O_{T_r}$ of size $B_r \times d$ each, divide $dO$ into $T_r$ blocks $dO_1, ..., dO_{T_r}$ of size $B_r \times d$ each, and divide $L$ into $T_r$ blocks $L_i, ..., L_{T_r}$ of size $B_r$ each.
3: Initialize $dQ = (0)_{N \times d}$ in HBM and divide it into $T_r$ blocks $dQ_1, ..., dQ_{T_r}$ of size $B_r \times d$ each. Divide $dK, dV \in \mathbb{R}^{N \times d}$ in to $T_c$ blocks $dK_1, ..., dK_{T_c}$ and $dV_1, ..., dV_{T_c}$, of size $B_c \times d$ each.
4: Compute $D = \text{rowsum}(dO \circ O) \in \mathbb{R}^d$ (pointwise multiply), write $D$ to HBM and divide it into $T_r$ blocks $D_1, ..., D_{T_r}$ of size $B_r$ each.
5: for $1 \le j \le T_c$ do
6:   Load $K_j, V_j$ from HBM to on-chip SRAM.
7:   Initialize $dK_j = (0)_{B_c \times d}, dV_j = (0)_{B_c \times d}$ on SRAM.
8:   for $1 \le i \le T_r$ do
9:     Load $Q_i, O_i, dO_i, dQ_i, L_i, D_i$ from HBM to on-chip SRAM.
10:    On chip, compute $S_i^{(j)} = Q_i K_j^T \in \mathbb{R}^{B_r \times B_c}$.
11:    On chip, compute $P_i^{(j)} = \exp(S_{ij} - L_i) \in \mathbb{R}^{B_r \times B_c}$.
12:    On chip, compute $dV_j \leftarrow dV_j + (P_i^{(j)})^T dO_i \in \mathbb{R}^{B_c \times d}$.
13:    On chip, compute $dP_i^{(j)} = dO_i V_j^T \in \mathbb{R}^{B_r \times B_c}$.
14:    On chip, compute $dS_i^{(j)} = P_i^{(j)} \circ (dP_i^{(j)} - D_i) \in \mathbb{R}^{B_r \times B_c}$.
15:    Load $dQ_i$ from HBM to SRAM, then on chip, update $dQ_i \leftarrow dQ_i + dS_i^{(j)} K_j \in \mathbb{R}^{B_r \times d}$, and write back to HBM.
16:    On chip, compute $dK_j \leftarrow dK_j + dS_i^{(j)T} Q_i \in \mathbb{R}^{B_c \times d}$.
17:   end for
18:   Write $dK_j, dV_j$ to HBM.
19: end for
20: Return $dQ, dK, dV$.



$Y (N, M)$     $X (N, D)$     $W (D, M)$

First row here is multiplied with every column here to make: First row of Y

So, for Jacobian [dy/dx]:
- derivative of Y's First row with X' First row gives an output, but Y's First row with all other X's rows gives ZERO
- Hence, Jacobian turns out to be quite sparse
Also, Jacobian : can be very large, possible that it wont fit into GPU RAM

Can be shown (notebook) that:
Downstream Gradient wrt input = Upstream Gradient wrt output * Weight.T
Similarly,
Downstream Gradient wrt weight = Input.T * Upstream Gradient wrt output



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

In place of Jacobian

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



$A$ $(N, 3)$     $B$ $(3, 4)$     $(N, 4)$

$$O_i = a_{i1} b_{1i} + a_{i2} b_{2i} + a_{i3} b_{3i} = \sum_{j=1}^{3} a_{ij} b_{ji}$$

(2)

**FlashAttention1 Backward Pass algorithm**

Require: Matrices $Q, K, V, O, dO \in \mathbb{R}^{N \times d}$ in HBM, vectors $\ell, m \in \mathbb{R}^N$ in HBM, on-chip SRAM of size $M$, softmax scaling constant $\tau \in \mathbb{R}$, masking function MASK, dropout probability $p_{drop}$, pseudo-random number generator state $\mathcal{R}$ from the forward pass.

1. Set the pseudo-random number generator state to $\mathcal{R}$.
2. Set block sizes $B_c = \lceil \frac{M}{4d} \rceil$, $B_r = \min(\lceil \frac{M}{4d} \rceil, d)$.
3. Divide $Q$ into $T_r = \lceil \frac{N}{B_r} \rceil$ blocks $Q_1, \ldots, Q_{T_r}$ of size $B_r \times d$ each, and divide $K, V$ in to $T_c = \lceil \frac{N}{B_c} \rceil$ blocks $K_1, \ldots, K_{T_c}$ and $V_1, \ldots, V_{T_c}$, of size $B_c \times d$ each.
4. Divide $O$ into $T_r$ blocks $O_1, \ldots, O_{T_r}$ of size $B_r \times d$ each, divide $dO$ into $T_r$ blocks $dO_1, \ldots, dO_{T_r}$ of size $B_r \times d$ each, divide $\ell$ into $T_r$ blocks $\ell_1, \ldots, \ell_{T_r}$ of size $B_r$ each, divide $m$ into $T_r$ blocks $m_1, \ldots, m_{T_r}$ of size $B_r$ each.
5. Initialize $dQ = (0)_{N \times d}$ in HBM and divide it into $T_r$ blocks $dQ_1, \ldots, dQ_{T_r}$ of size $B_r \times d$ each. Initialize $dK = (0)_{N \times d}, dV = (0)_{N \times d}$ in HBM and divide $dK, dV$ in to $T_c$ blocks $dK_1, \ldots, dK_{T_c}$ and $dV_1, \ldots, dV_{T_c}$, of size $B_c \times d$ each.
6. for $1 \le j \le T_c$ do
7.   Load $K_j, V_j$ from HBM to on-chip SRAM.   ← Outer Loop thru all K and V blocks
8.   Initialize $dK_j = (0)_{B_c \times d}, dV_j = (0)_{B_c \times d}$ on SRAM.
9.   for $1 \le i \le T_r$ do
10.    Load $Q_i, O_i, dO_i, dQ_i, \ell_i, m_i$ from HBM to on-chip SRAM.   ← Inner loop thru all Q blocks
11.    On chip, compute $S_{ij} = \tau Q_i K_j^T \in \mathbb{R}^{B_r \times B_c}$.
12.    On chip, compute $S_{ij}^{masked} = \text{MASK}(S_{ij})$.
13.    On chip, compute $P_{ij} = \text{diag}(\ell_i)^{-1} \exp(S_{ij}^{masked} - m_i) \in \mathbb{R}^{B_r \times B_c}$.
14.    On chip, compute dropout mask $Z_{ij} \in \mathbb{R}^{B_r \times B_c}$ where each entry has value $\frac{1}{1 - p_{drop}}$ with probability $1 - p_{drop}$ and value 0 with probability $p_{drop}$.
15.    On chip, compute $P_{ij}^{dropped} = P_{ij} \circ Z_{ij}$ (pointwise multiply).
16.    On chip, compute $dV_j \leftarrow dV_j + (P_{ij}^{dropped})^T dO_i \in \mathbb{R}^{B_c \times d}$.
17.    On chip, compute $dP_{ij}^{dropped} = dO_i V_j^T \in \mathbb{R}^{B_r \times B_c}$.
18.    On chip, compute $dP_{ij} = dP_{ij}^{dropped} \circ Z_{ij}$ (pointwise multiply).
19.    On chip, compute $D_i = \text{rowsum}(dO_i \circ O_i) \in \mathbb{R}^{B_r}$.
20.    On chip, compute $dS_{ij} = P_{ij} \circ (dP_{ij} - D_i) \in \mathbb{R}^{B_r \times B_c}$.
21.    Write $dQ_i \leftarrow dQ_i + \tau dS_{ij} K_j \in \mathbb{R}^{B_r \times d}$ to HBM.
22.    On chip, compute $dK_j \leftarrow dK_j + \tau dS_{ij}^T Q_i \in \mathbb{R}^{B_c \times d}$.
23.  end for
24.  Write $dK_j \leftarrow dK_j, dV_j \leftarrow dV_j$ to HBM.
25. end for
26. Return dQ, dK, dV

Umar's Implementation:
Works with both Causal as well as Non-causal Attention

1. Inner for Loop split into 2 parts, as otherwise HBM needs to be written to at every Inner Loop iteration -> too costly
   Split like:
   - Each dq depends upon Loops over Ks:
     Fix Q block, loop over all KV blocks
   - Each dk depends upon Loops over Qs:
     Fix K block, loop over all Q blocks
2. To compute dq and dk vectors, need 'D_i':

$$dq_i = \sum_j dS_{ij} k_j = \sum_j P_{ij}(dP_{ij} - D_i)k_j = \sum_j \frac{e^{q_i^T k_j}}{L_i}(do_i^T v_j - D_i)k_j.$$

$$dk_j = \sum_i dS_{ij} q_i = \sum_i P_{ij}(dP_{ij} - D_i)q_i = \sum_i \frac{e^{q_i^T k_j}}{L_i}(do_i^T v_j - D_i)q_i.$$

where D_i is:

$$D_i = P_i^T dP_i = \sum_j \left(\frac{e^{q_i^T k_j}}{L_i}\right) do_i^T v_j = do_i^T \sum_j \frac{e^{q_i^T k_j}}{L_i} v_j = do_i^T o_i$$

So effectively, we:
- Calculate D_i
- Calculate dq
- Calculate dk
- Fix Q, loop over all KV
- Fix K, loop over all Q

Triton Autotuning:
- Triton does Thread Coarsening for us, unlike CUDA, based on:
  - Block size
  - Number of Warps
- But, we need to give Triton like:

```
@triton.autotune(
    [
        triton.Config(
            {"BLOCK_SIZE_Q": BLOCK_SIZE_Q, "BLOCK_SIZE_KV": BLOCK_SIZE_KV},
            num_stages=num_stages,
            num_warps=num_warps,
        )
        for BLOCK_SIZE_Q in [64, 128]
        for BLOCK_SIZE_KV in [32, 64]
        for num_stages in [3, 4, 7]
        for num_warps in [2, 4]
    ],
    key=["SEQ_LEN", "HEAD_DIM"],
)
```

Have to try various configs for BLOCK_SIZE based on what works best (heuristic) based on timing

Triton runs a config for each pair of SEQ_LEN and HEAD_DIM, choosing the best throughput running in least amount of time

Software Pipelining:
- Piece of code that can be parallelized
- 
```
Imagine you have a for loop

for i = 1 to N
    A = LOAD (...)          → spawn
    B = LOAD (...)          → spawn another

    C = A * B               → Check if both spawns are done
    STORE (...)               before doing matmul
```

- If done sequentially, above code does not make optimal use of GPU
- We parallelize like:

Prologue | All units working | Epilogue

- Conditions:
  - Parallelization can be done using 'async' operations
  - More memory needed for this as SRAM needs to hold more memory:
    - When Compute Unit does MM1, at the same time step (3rd time step):
      - Reads for two matrices (RDA2 and RDB2) are already done
      - Reads for two matrices are half done already (RDA3)
  - Num of Iterations of for loop have to be MUCH GREATER THAN num of Stages of software pipeline
    (4 stages in example above: [RDA, RDB, MM, WR])
    # Iterations in loop >>> # stages in software pipeline, for pipelining to work well

===============================================================================

Credit:
1. hkproj/triton-flash-attention: https://github.com/hkproj/triton-flash-attention
2. Tri Dao et. al (2022) FlashAttention: Fast and Memory-Efficient Exact Attention with IO-Awareness. arXiv:2205.14135
3. Tri Dao (2023) FlashAttention-2: Faster Attention with Better Parallelism and Work Partitioning. arXiv:2307.08691
4. Basics on NVIDIA GPU Hardware Architecture, NASA : https://www.nas.nasa.gov/hecc/support/kb/basics-on-nvidia-gpu-hardware-architecture_704.html#:~:text=Memory%20access%20latency%20is%20lower,227%20KB%20of%20shared%20memory, accessed 10/22
5. USC EE508 Hardware Foundations of ML Notes

===============================================================================