## GP-GPU and High Performances Computing

Lecture 07 – Matrix Multiplication

January 21, 2024

- Example of scan algorithm
- Processor hardware
  - Max threads per SM : 2048
  - Max threads per block : 1024
  - Max warps per SM : 64
- If 2 blocks are assigned to an SM and each block has 1024 threads, how many warps are there in an SM?
  - Each block is divided into 1024/32 = 32 warps
  - There are 32 * 2 = 64 Warps
- At any point in time, only 4 of the 64 warps will be selected for instruction fetch and execution.
  - One instruction is issued for 1 warp at every cycle (by design).
  - 16 cycles are needed to execute 1 instruction on all threads of the block.
- SM will interleaved warps to optimize execution.

Given two squares matrices in $\boldsymbol{R}^{Width \times Width}$, $M$ and $N$, we multiply $M$ by $N$ to compute a third square matrix in $\boldsymbol{R}^{Width \times Width}$, $P$.
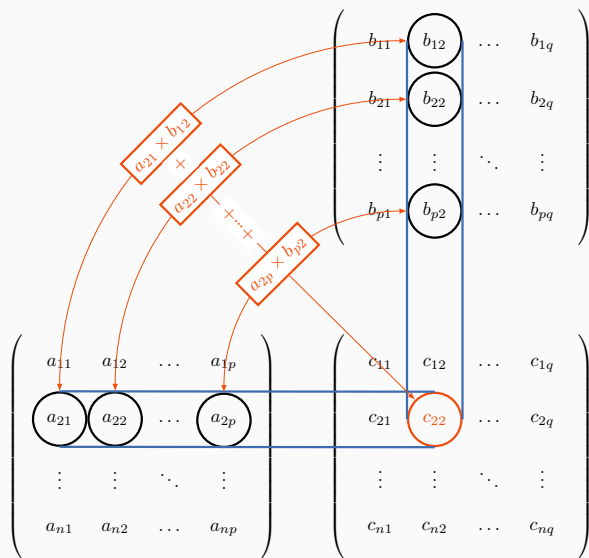
$$P = MN$$

In terms of the elements of P, matrix multiplication implies computing, for all $1 \le i, j \le Width$

$$P_{ij} = \sum_{k=1}^{Width} M_{ik} N_{kj}$$

The complexity for the naïve computation is $\mathcal{O}(Width^3)$.

```
1    void GMM_CPU(float* M, float* N, float* P, int Width)
2    {
3      for (int i = 0; i < Width; ++i)
4        for (int j = 0; j < Width; ++j) {
5          float sum = 0;
6          for (int k = 0; k <Width; ++k) {
7            float a = M[i * Width + k];
8            float b = N[k * Width + j];
9            sum += a * b;
10         }
11         P[i * Width + j] = sum;
12       }
13   }
```

- 3 possibles choices: $M$, $N$ or $P$.
- The outer loops are all independent computations.
- We will focus on the computation of the elements of $P$.
- The inner loop is a scalar product between a row of $M$ and a column of $N$. It can be parallelize using reduction.

$P$ is $2D$, one possible choice is to organize threads in $2D$ as well:

- Split the output $P$ into square tiles of size $TILE\_WIDTH \times TILE\_WIDTH$ (a preprocessor user defined constant).
- Each thread block produces one tile of $[TILE\_WIDTH]^2$ elements.
- Create $[ceil(Width/TILE\_WIDTH)]^2$ thread blocks to cover the output matrix.

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| $P_{0,0}$ | $P_{0,1}$ | $P_{0,2}$ | $P_{0,3}$ | $P_{0,4}$ | $P_{0,5}$ | $P_{0,6}$ | $P_{0,7}$ |
| $P_{1,0}$ | $P_{1,1}$ | $P_{1,2}$ | $P_{1,3}$ | $P_{1,4}$ | $P_{1,5}$ | $P_{1,6}$ | $P_{1,7}$ |
| $P_{2,0}$ | $P_{2,1}$ | $P_{2,2}$ | $P_{2,3}$ | $P_{2,4}$ | $P_{2,5}$ | $P_{2,6}$ | $P_{2,7}$ |
| $P_{3,0}$ | $P_{3,1}$ | $P_{3,2}$ | $P_{3,3}$ | $P_{3,4}$ | $P_{3,5}$ | $P_{3,6}$ | $P_{3,7}$ |
| $P_{4,0}$ | $P_{4,1}$ | $P_{4,2}$ | $P_{4,3}$ | $P_{4,4}$ | $P_{4,5}$ | $P_{4,6}$ | $P_{4,7}$ |
| $P_{5,0}$ | $P_{5,1}$ | $P_{5,2}$ | $P_{5,3}$ | $P_{5,4}$ | $P_{5,5}$ | $P_{5,6}$ | $P_{5,7}$ |
| $P_{6,0}$ | $P_{6,1}$ | $P_{6,2}$ | $P_{6,3}$ | $P_{6,4}$ | $P_{6,5}$ | $P_{6,6}$ | $P_{6,7}$ |
| $P_{7,0}$ | $P_{7,1}$ | $P_{7,2}$ | $P_{7,3}$ | $P_{7,4}$ | $P_{7,5}$ | $P_{7,6}$ | $P_{7,7}$ |

Block size : each block as $2 \times 2 = 4$ threads.

Number of blocks : $\dfrac{Width}{TILE\_WIDTH} \implies 16$ blocks.

| $P_{0,0}$ | $P_{0,1}$ | $P_{0,2}$ | $P_{0,3}$ | $P_{0,4}$ | $P_{0,5}$ | $P_{0,6}$ | $P_{0,7}$ |
|---|---|---|---|---|---|---|---|
| $P_{1,0}$ | $P_{1,1}$ | $P_{1,2}$ | $P_{1,3}$ | $P_{1,4}$ | $P_{1,5}$ | $P_{1,6}$ | $P_{1,7}$ |
| $P_{2,0}$ | $P_{2,1}$ | $P_{2,2}$ | $P_{2,3}$ | $P_{2,4}$ | $P_{2,5}$ | $P_{2,6}$ | $P_{2,7}$ |
| $P_{3,0}$ | $P_{3,1}$ | $P_{3,2}$ | $P_{3,3}$ | $P_{3,4}$ | $P_{3,5}$ | $P_{3,6}$ | $P_{3,7}$ |
| $P_{4,0}$ | $P_{4,1}$ | $P_{4,2}$ | $P_{4,3}$ | $P_{4,4}$ | $P_{4,5}$ | $P_{4,6}$ | $P_{4,7}$ |
| $P_{5,0}$ | $P_{5,1}$ | $P_{5,2}$ | $P_{5,3}$ | $P_{5,4}$ | $P_{5,5}$ | $P_{5,6}$ | $P_{5,7}$ |
| $P_{6,0}$ | $P_{6,1}$ | $P_{6,2}$ | $P_{6,3}$ | $P_{6,4}$ | $P_{6,5}$ | $P_{6,6}$ | $P_{6,7}$ |
| $P_{7,0}$ | $P_{7,1}$ | $P_{7,2}$ | $P_{7,3}$ | $P_{7,4}$ | $P_{7,5}$ | $P_{7,6}$ | $P_{7,7}$ |

Block size : each block as $4 \times 4 = 16$ threads.

Number of blocks : $\dfrac{n}{TILE\_WIDTH} \implies 4$ blocks.

```
1    // TILE_WIDTH is a #define constant
2    dim3 dimGrid(ceil((1.0*Width)/TILE_WIDTH), ceil((1.0*Width)/TILE_WIDTH), 1);
3    dim3 dimBlock(TILE_WIDTH, TILE_WIDTH, 1);
4    // Launch the device computation threads!
5    MatrixMulKernel<<<dimGrid, dimBlock>>> (Md, Nd, Pd, Width);
```

```
1       // Matrix multiplication kernel - per thread code
2       __global__ void MatrixMulKernel(float* d_M, float* d_N,
3                                       float* d_P, int Width) {
4
5       // Pvalue is used to store the element of the matrix
6       // that is computed by the thread
7         float Pvalue = 0;
```

# Simple multiplication kernel

```
1    __global__ void MatrixMulKernel(float* d_M, float* d_N, float* d_P, int Widt
2      // Calculate the row index of the d_P element and d_M
3      int Row = blockIdx.y*blockDim.y+threadIdx.y;
4      // Calculate the column index of d_P and d_N
5      int Col = blockIdx.x*blockDim.x+threadIdx.x;
6      if ((Row < Width) && (Col < Width)) {
7        float Pvalue = 0;
8        // each thread computes one element of the block sub-matrix
9        for (int k = 0; k < Width; ++k)
10          Pvalue += d_M[Row*Width+k] * d_N[k*Width+Col];
11        d_P[Row*Width+Col] = Pvalue;
12      }
13    }
```

That's a simple implementation:

➤ GPU kernel is the CPU code with the outer loops replaced
➤ with per-thread index calculations!

Unfortunately, performance is quite bad. Why? With the given approach,

➤ global memory bandwidth
➤ can't supply enough data
➤ to keep the SMs busy!

```
1    __global__ void MatrixMulKernel(float* d_M, float* d_N,
2                                    float* d_P, int Width) {
3      // Calculate the row index of d_P  and d_M
4      int Row = blockIdx.y*blockDim.y+threadIdx.y;
5      // Calculate the column index of d_P and d_N
6      int Col = blockIdx.x*blockDim.x+threadIdx.x;
7      if ((Row < Width) && (Col < Width)) {
8        float Pvalue = 0;
9        // each thread computes one element of the block sub-matrix
10       for (int k = 0; k < Width; ++k)
11         Pvalue += d_M[Row*Width+k] * d_N[k*Width+Col];
12       d_P[Row*Width+Col] = Pvalue;
13     }
14   }
```

- Each threads access global memory
  - 4B each, or 8B per pair.
  - (And once TOTAL to P per thread—ignore it.)

  –for elements of M and N:

- With each pair of elements, a thread does a single multiply-add, –2 FLOP—floating-point operations.

- So for every FLOP, a thread needs 4B from memory: **4B / FLOP**.

- One generation of GPUs: 1,000 GFLOP/s of compute power, and 150 GB/s of memory bandwidth.
- Dividing bandwidth by memory requirements: 150 GB/S, Host 4B.Flop = 37.5 GFLOP/S which limits computation!

# Reuse Memory Accesses!

- ➤ 37.5 GFLOPs is a limit.
- ➤ In an actual execution, memory is not busy all the time, and the code runs at about 25 GFLOPs.
- ➤ To get closer to 1,000 GFLOPs, we need to drastically cut down accesses to global memory.

# A common programming strategy

- ➤ The dilemma:
  - ➤ Matrices M and N are large.
  - ➤ They fit easily in global memory, but that's slow.
  - ➤ Shared memory is fast, but M and N don't fit.
- ➤ The solution:
  - ➤ Break M and N into tiles (called blocks in the much older CPU literature).
  - ➤ Read a tile into shared memory.
  - ➤ Use the tile from shared memory.
  - ➤ Repeat until done.

# A common programming strategy

- In a GPU, only threads in a block can use shared memory.
- Thus, each block operates on separate tiles: -
  - Read tile(s) into shared memory using multiple threads to exploit memory-level parallelism.
  - Compute based on shared memory tiles.
  - Repeat.
  - Write results back to global memory.

```
1
2      __global__ void MatrixMulKernel(float* M, float* N, float* P, int Width)
3      {
4        __shared__ float subTileM[TILE_WIDTH][TILE_WIDTH];
5        __shared__ float subTileN[TILE_WIDTH][TILE_WIDTH];
```

Identify a tile of global data that are accessed by multiple threads

➤ Load the tile from global memory into on-chip memory

➤ Have the multiple threads to access their data from the on-chip memory

➤ Move on to the next block/tile

## Idea: Place global memory data into Shared Memory for reuse

- ► Each input element is used to calculate WIDTH elements of $P$.
- ► Load each element into Shared Memory
- ► have several threads use the local version to reduce memory bandwidth.

Break up the execution of the kernel into phases so that the data accesses in each phase are focused on one subset (tile) of M and N

➤ All threads in a block participate
➤ Each thread loads
  ➤ one M element (corresponding to the global index of the thread)
  ➤ one N element (corresponding to the global index of the thread)
  ➤ in basic tiling code: Assign the loaded element to each thread such that the accesses within each warp is coalesced (more later)

But …

➤ How can a thread know that another thread has finished its part of the tile?

➤ Or that another thread has finished using the previous tile?

There is a need to synchronize

➤ Bulk synchronous execution: threads execute roughly in unison
  ➤ Do some work
  ➤ Wait for others to catch up
  ➤ Repeat
➤ Much easier programming model
  ➤ Threads only parallel within a section
  ➤ Debug lots of little programs
  ➤ Instead of one large one.
➤ Dominates high-performance applications

- ➤ How does it work?
- ➤ Use a barrier to wait for thread to 'catch up.'
- ➤ A barrier is a synchronization point:
    - ➤ each thread calls a function to enter barrier
    - ➤ threads block (sleep) in barrier function until all threads have called
    - ➤ after last thread calls function, all threads continue past the barrier.

- Use \_\_syncthreads for CUDA Blocks
- How does it work in CUDA?
- Only within thread blocks!
- The function: **void \_\_syncthreads(void);**
- All threads in block must enter (no subsets).
- All threads must enter the SAME static call (not the same as all threads calling function!).

- ➤ What exactly is guaranteed to have finished?
  - ➤ Are shared memory operations before a barrier (e.g., stores) guaranteed to have completed?
    - · What about global memory ops?
    - · What about atomic ops with no return values?
    - · What about I/O operations?
  - ➤ CUDA manual: all global and shared memory ops (which presumably includes atomic variants) have completed.
  - ➤ Avoid assumptions about I/O (such as `printf`).

# Tiled matrix multiplication

```
1    __global__ void MatrixMulKernel(float* M, float* N, float* P, int Width)
2    {
3      __shared__ float subTileM[TILE_WIDTH][TILE_WIDTH];
4      __shared__ float subTileN[TILE_WIDTH][TILE_WIDTH];
5      int bx = blockIdx.x; int by = blockIdx.y;
6      int tx = threadIdx.x; int ty = threadIdx.y;
7      // Identify the row and column of the P element to work on
8      int Row = by * TILE_WIDTH + ty; // note: blockDim.x == TILE_WIDTH
9      int Col = bx * TILE_WIDTH + tx; // blockDim.y == TILE_WIDTH
10     float Pvalue = 0;
11     // Loop over the M and N tiles required to compute the P element
12     // The code assumes that the Width is a multiple of TILE_WIDTH!
13     for (int m = 0; m < Width/TILE_WIDTH; ++m) {
14     // Collaborative loading of M and N tiles into shared memory
15       subTileM[ty][tx] = M[Row*Width + m*TILE_WIDTH+tx];
16       subTileN[ty][tx] = N[(m*TILE_WIDTH+ty)*Width+Col];
17       __syncthreads();
18       for (int k = 0; k < TILE_WIDTH; ++k)
19         Pvalue += subTileM[ty][k] * subTileN[k][tx];
20       __syncthreads();
21     }
22     P[Row*Width+Col] = Pvalue;
23   }
```

# Classical matrix multiply

```
1    {
2      // Calculate the row index of the P element and M
3      int Row = blockIdx.y * blockDim.y + threadIdx.y;
4      // Calculate the column index of P and N
5      int Col = blockIdx.x * blockDim.x + threadIdx.x;
6      if ((Row < Width) && (Col < Width)) {
7        float Pvalue = 0;
8        // each thread computes one element of the block sub-matrix
9        for (int k = 0; k < Width; ++k)
10          Pvalue += M[Row*Width+k] * N[k*Width+Col];
11       P[Row*Width+Col] = Pvalue;
12     }
13   }
```

➤ Recall our example GPU: 1,000 GFLOP/s, 150 GB/s
➤ 16x16 tiles use each operand for 16 operations
  ➤ reduce global memory accesses by a factor of 16
  ➤ 150GB/s bandwidth supports (150/4)*16 = 600 GFLOPs!
➤ 32x32 tiles use each operand for 32 operations
  ➤ reduce global memory accesses by a factor of 32
  ➤ 150 GB/s bandwidth supports (150/4)*32 = 1,200 GFLOPs!
  ➤ Memory bandwidth is no longer the bottleneck!

- Shared memory size
  - implementation dependent
  - 164kB per SM in Ampere (163kB max per block)
- Given `TILE_WIDTH` of 16 (256 threads / block),
  - each thread block uses $2 \times 256 \times 4B = 2kB$ of shared memory, which limits active blocks to 82;
  - maximum of 2048 threads per SM, which limits blocks to 8.
  - Thus up to $8 \times 512 = 4,096$ pending loads (2 per thread, 256 threads per block)

➤ Given `TILE_WIDTH` of 32 (1 024 threads / block),
  ➤ each thread block uses $2 \times 1024 \times 4B = 8kB$ of shared memory, which limits active blocks to 20;
  ➤ maximum of 2,048 threads per SM, which limits blocks to 2.
  ➤ Thus up to $2 \times 2,048 = 4,096$ pending loads (2 per thread, 1,024 threads per block) (same memory parallelism exposed)

- Number of devices in the system

```
1    int dev_count;
2    cudaGetDeviceCount( &dev_count);
```

- Capability of devices

```
1    cudaDeviceProp dev_prop;
2    for (i = 0; i < dev_count; i++) {
3      cudaGetDeviceProperties( &dev_prop, i);
4      // decide if device has sufficient resources and capabilities
5    }
```

- cudaDeviceProp is a built-in C structure type

- dev_prop.dev_prop.maxThreadsPerBlock

- dev_prop.sharedMemoryPerBlock

- How to Handle Matrices of Other Sizes? Use tiles :
  assumed integral number of tiles (thread blocks) in all matrix
  dimensions.
- How can we avoid this assumption?
  One answer: add padding, but not easy to reformat data, and adds
  transfer time. Other ideas?

- Threads that calculate valid P elements but can step outside valid input
  : Second tile of Block(0,0), all threads when k is 1
- Threads that do not calculate valid P elements
  - Block(1,1), Thread(1,0), non-existent row
  - Block(1,1), Thread(0,1), non-existing column
  - Block(1,1), Thread(1,1), non-existing row/column

- Test during tile load: is target within input matrix?
  - If yes, proceed to load;
  - otherwise, just write 0 to shared memory.
- The benefit?
  - No specialization during tile use!
  - Multiplying by 0 guarantees that unwanted terms do not contribute to the inner product.

If a thread is not within P,

➤ All terms in sum are 0.
➤ No harm in performing FLOPs.
➤ No harm in writing to registers.
➤ Must not be allowed to write to global memory!

So: Threads outside of P calculate 0, but store nothing.

➤ `for (int m = 0; m < Width/TILE_WIDTH; ++m)`
The bound for m implicitly assumes that Width is a multiple of `TILE_WIDTH`. We need to round up.

➤
`for (int m = 0; m < (Width \item 1)/TILE_WIDTH + 1; ++m)`

  ➤ For non-multiples of `TILE_WIDTH`:
    · quotient is unchanged;
    · add one to round up.
  ➤ For multiples of `TILE_WIDTH`:
    · quotient is now one smaller,
    · but we add 1.

We had …

```
1    // Collaborative loading of M and N tiles into shared memory
2    subTileM[ty][tx] = M[Row*Width + m*TILE_WIDTH+tx];
3    subTileN[ty][tx] = N[(m*TILE_WIDTH+ty)*Width+Col];
```

Note: the tests for M and N tiles are NOT the same.

```
1    if (Row < Width && m*TILE_WIDTH+tx < Width) {
2    // as before
3      subTileM[ty][tx] = M[Row*Width + m*TILE_WIDTH+tx];
4    } else {
5      subTileM[ty][tx] = 0;
6    }
```

We had …

```
1    for (int k = 0; k < TILE_WIDTH; ++k)
2      Pvalue += subTileM[ty][k] * subTileN[k][tx];
```

Note: no changes are needed, but we might save a little energy (fewer floating-point ops)?

```
1    if (Row < Width && Col < Width) {
2      // as before
3      for (int k = 0; k < TILE_WIDTH; ++k)
4        Pvalue += subTileM[ty][k] * subTileN[k][tx];
5    }
```

We had

```
1    P[Row*Width+Col] = Pvalue;
```

We must test for threads outside of P:

```
1    if (Row < Width && Col < Width) {
2    // as before
3    P[Row*Width+Col] = Pvalue;
4    }
```

- ➤ For each thread, conditions are different for
  - ➤ Loading M element
  - ➤ Loading N element
  - ➤ Calculation/storing output elements
- ➤ Branch divergence
  - ➤ affects only blocks on boundaries,
  - ➤ should be small for large matrices.
- ➤ What about rectangular matrices?

## Conclusion

- Themes of this class
  - Organization of a computation with respect to data and architecture
  - Usage of shared memory
- Computations should carefully adapt to the architecture and maximize the usage of the ressources