

# CME213/ME339

## Lecture 9

Eric Darve  
Erich Elsen  
Austin Gibbons

Department of Mechanical Engineering  
Institute for Computational and Mathematical Engineering  
Stanford University

Spring 2012



# Occupancy calculator

---

- Google “occupancy calculator”
- Input required:
  - 1 Compute capability: 2.0
  - 2 Threads per block
  - 3 Shared memory
- Get these numbers by compiling your code with the option:  
`--ptxas-options=-v`



# Matrix-matrix product

---

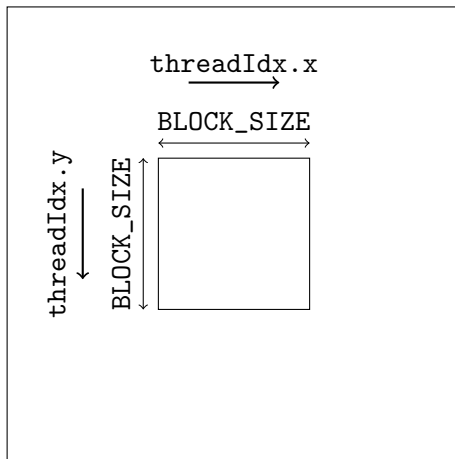
Let's illustrate these concepts using a matrix-matrix product:

$$C_{ij} = \sum_k A_{ik} B_{kj}$$



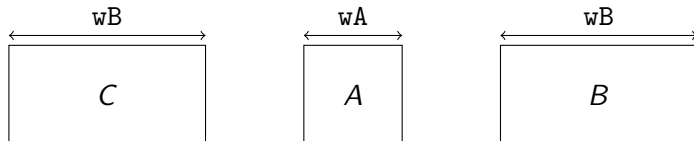
# Layout

---



# Matrix sizes

---



```

1  __global__ void
2  matrixMul_slow( float* C, float* A, float* B, int wA, int wB)
3  {
4      // Block index
5      int bx = blockIdx.x;
6      int by = blockIdx.y;
7
8      // Thread index
9      int tx = threadIdx.x;
10     int ty = threadIdx.y;
11
12     // Index of the first sub-matrix of A processed by the block
13     int aBegin = wA * blockDim.y * by;
14
15     // Index of the last sub-matrix of A processed by the block
16     int aEnd   = aBegin + wA - 1;
17
18     // Index of the first sub-matrix of B processed by the block
19     int bBegin = blockDim.x * bx;
20     ...

```



```

1  ...
2  // Csub is used to store the element of the block sub-matrix
3  // that is computed by the thread
4  float Csub = 0.0f;
5
6  // Loop over all the sub-matrices of A and B
7  // required to compute the block sub-matrix
8  for (int a = aBegin, b = bBegin; a <= aEnd;
9      ++a, b += wB) {
10
11     // Multiply the two matrices together;
12     // each thread computes one element
13     // of the block sub-matrix
14     Csub += A[a + wA * ty] * B[b + tx];
15 }
16
17 // Write the block sub-matrix to device memory
18 int c = wB * blockDim.y * by + blockDim.x * bx;
19 C[c + wB * ty + tx] = Csub;
20 }

```



# Output from compiler

---

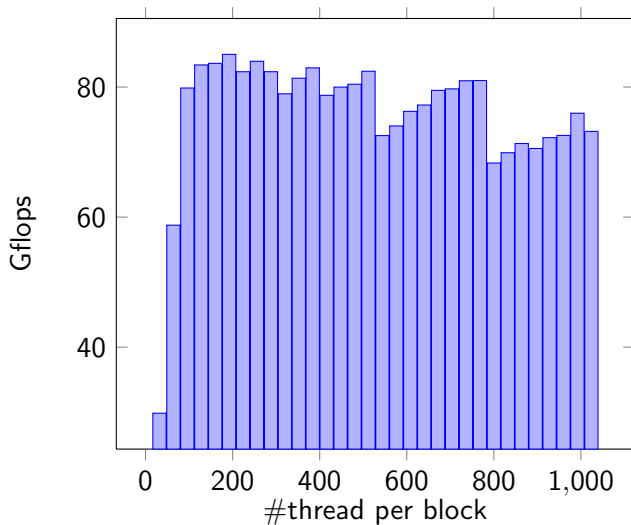
```
ptxas info      : Compiling entry function
                  '_Z12matrixMul_v2PfS_S_ii' for 'sm_20'
ptxas info      : Function properties for _Z12matrixMul_v2PfS_S_ii
                  0 bytes stack frame, 0 bytes spill stores, 0 bytes spill loads
ptxas info      : Used 18 registers, 64 bytes cmem[0]
```

18 registers, 64 bytes constant memory, no shared memory.





# Threads per block



Performance follows the occupancy.

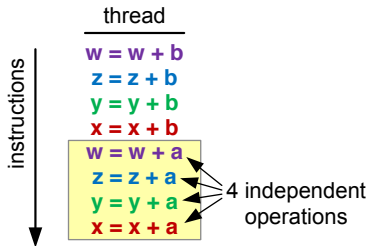
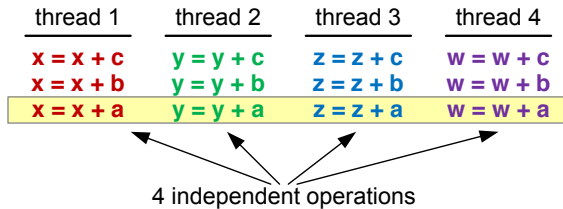


# Generating concurrency

---

- In some cases running enough threads is not possible to reach peak performance.
- For example, if the latency is 24 cycles, we need as many as 24 warps to hide this latency, which is 50% occupancy.
- Parallelism can also be generated at the level of a single thread. That is the hardware will take advantage of multiple instructions if there are no dependencies between them.





# Number of grid blocks and number of threads per block — heuristics

---

- We need approx. 18–27 warps = 576–864 threads per SM for good performance (less if you can issue independent instructions).
- The number of blocks must be  $>$  than the number of SMs, so each SM has at least one block.
- Typically,  $\# \text{ of blocks} / \# \text{ of multiprocessors} > 2$ :
  - Multiple blocks can run concurrently in a multiprocessor.
  - Blocks that aren't waiting at a `__syncthreads()` keep the hardware busy.
  - Subject to resource availability – registers, shared memory.



# Performance of matrix-matrix product

---

- The number of Gflops we achieved with the matrix-matrix product is disappointing. It peaks at about 80 Gflops.
- Consider the memory access pattern:

```
1  int aBegin = wA * blockDim.y * by;
2  int bBegin = blockDim.x * bx;
3  float Csub = 0.f;
4
5  for (int a = aBegin, b = bBegin; a <= aEnd;
6      ++a, b += wB)
7      Csub += A[a + wA * ty] * B[b + tx];
8
9  int c = wB * blockDim.y * by + blockDim.x * bx;
10 C[c + wB * ty + tx] = Csub;
```



# Memory access

---

- $B[b + tx]$ : excellent access if  $\text{blockDim.x} == 32$
- $A[a + w_A * ty]$ : since we have  $++a$  between iterations, we can use the cache but chances are data get evicted before we finish the block.
- This suggests using shared memory.
- Load a block of  $A$  and a block of  $B$  in shared memory and then perform multiplication.



# Block multiplication

```
1      // Loop over all the sub-matrices of A and B
2      // required to compute the block sub-matrix
3      for (int a = aBegin, b = bBegin; a <= aEnd;
4           a += aStep, b += bStep) {
5
6          // Declaration of the shared memory arrays used to
7          // store the sub-matrices for A and B
8          __shared__ float As[BLOCK_SIZE][BLOCK_SIZE];
9          __shared__ float Bs[BLOCK_SIZE][BLOCK_SIZE];
10
11         // Load the matrices from device memory
12         // to shared memory
13         As[ty][tx] = A[a + wA * ty + tx];
14         Bs[ty][tx] = B[b + wB * ty + tx];
15
16         __syncthreads();
17     ...
```



```
1  ...
2      // Multiply the two blocks together
3  #pragma unroll
4      for (int k = 0; k < BLOCK_SIZE; ++k)
5          Csub += As[ty][k] * Bs[k][tx];
6
7      // Synchronize to make sure that the preceding
8      // computation is done before loading two new
9      // sub-matrices of A and B in the next iteration
10     __syncthreads();
11 }
```





# Memory access pattern

---

- $A_s[ty][tx] = A[a + w_A * ty + tx]$ : excellent; coalesced; 128 bytes if `BLOCK_SIZE == 32`
- $B_s[ty][tx] = B[b + w_B * ty + tx]$ : excellent.
- $B_s[k][tx]$ : conflict free access.
- $A_s[ty][k]$ : broadcast access; conflict free.
- Bandwidth from memory: GeForce GTX 480, bandwidth = 177.4 GB/s; peak flops = 1344.96 GFlops.
- To hide memory access we need at least  $\sim 30$  flops per word read from memory.



# Flop count

---

- Number of words read by a thread: 2
- Number of flops:  $2 \text{ BLOCK\_SIZE}$
- If  $\text{BLOCK\_SIZE} == 32$ , the ratio is 32. Enough to hide memory access!
- So we should get peak performance:

CUBLAS                      784.5752 GFlop/s, Time = 0.00100 s

CUDA matrixMul 235.8730 GFlop/s, Time = 0.00333 s

CUBLAS: library for linear algebra that is part of the CUDA SDK.



# Shared memory access

---

- We getting closer to CUBLAS but there is still room for improvement.
- We are currently slowed down by the access to the shared memory. It takes 4 cycles to read  $Bs[k][tx]$  and  $As[ty][k]$  but only one cycle to perform the operation.



# Instruction throughput

---

Guidelines to improve instruction throughput:

- Use single precision whenever possible.
- Use bitwise operations to calculate division or modulo by powers of 2

$i/2^k \rightarrow i \gg k$

$i\%2^k \rightarrow i \& (2^k-1)$



# Reciprocal square roots

---

The reciprocal square root should always be invoked explicitly as `rsqrtf()` for single precision and `rsqrt()` for double precision.

This is a common operation in graphics and therefore the hardware is highly optimized for this operation.



# Math functions and constants

---

- When doing floating point arithmetic use single-precision floating-point constants, defined with an `f` suffix such as `3.141592653589793f`, `1.0f`, `0.5f`.
- Use the appropriate math function

`rsqrtf(x)`, `sqrtf(x)`

`expf(x)`, `logf(x)`

`sinf(x)`, `cosf(x)`, `tanf(x)`

`rsqrt(x)`, `sqrt(x)`

`exp(x)`, `log(x)`

`sin(x)`, `cos(x)`, `tan(x)`



# Fast math

---

Some functions have a very fast implementation in hardware but with lower accuracy (?). Their names are prepended by underscores.

```
__fdividef(x,y) [x/y],  
__sinf(x,y), __cosf(x,y), __tanf(x,y),  
__expf(x,y), __logf(x,y)
```

See the programming guides for a full list.

