

CUDA: Memory

Mark Greenstreet

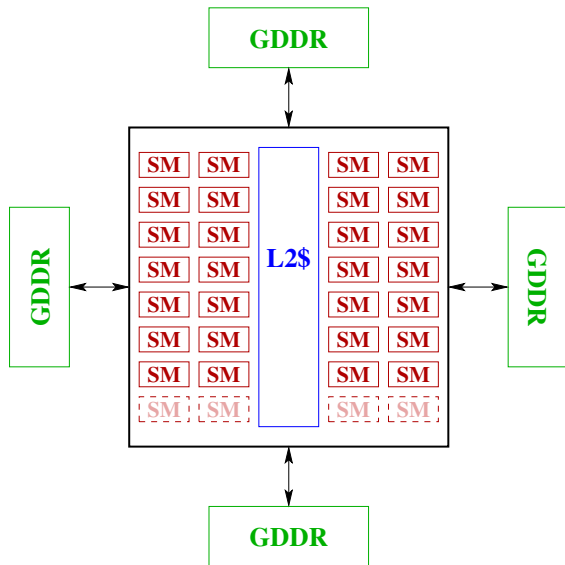
CpSc 418 – March 13 & 15, 2017

- [Architecture Snapshot](#)
- [Registers](#)
- [Shared Memory](#)
- [Global Memory](#)
- [Other Memory: texture memory, constant memory, caches](#)
- [Summary, preview, review, tide-chart](#)



Unless otherwise noted or cited, these slides are copyright 2017 by Mark Greenstreet and are made available under the terms of the Creative Commons Attribution 4.0 International license <http://creativecommons.org/licenses/by/4.0/>

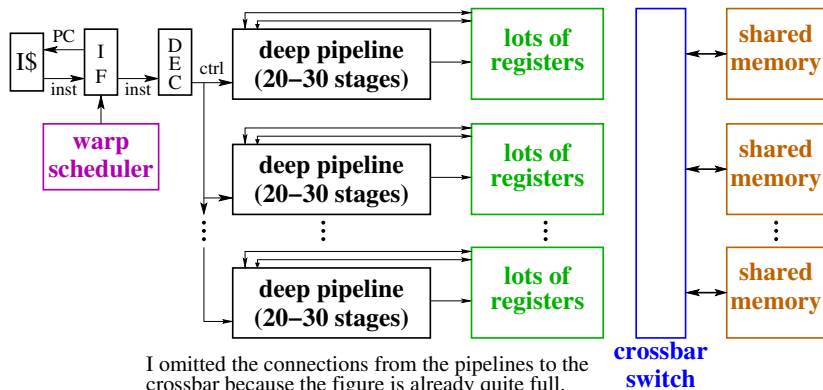
First, GPU Architecture Review



A “typical”, high-end GPU.

- 28 SMs:
 - ▶ 128 SPs/SM.
 - ▶ That’s 3584 SPs on the chip.
 - ▶ Each SM can schedule 4 warps in a single cycle.
- ~ 1.6GHz clock frequency.
- 11 GBytes of GDDR memory,
~ 484GBytes/sec. memory bandwidth.

A Streaming Multiprocessor (SM)



- Each of the pipelines is an SP (streaming processor)
- Lots of deep pipelines.
- Lots of threads: when we encounter an architectural challenge:
 - ▶ Raising throughput is **easy**, lowering latency is **hard**.
 - ▶ Solve problems by increasing latency and adding threads.
 - ▶ Make the programmer deal with it.

Why do we need a memory hierarchy

```
__global__ void saxpy(uint n, float a, float *x, float *y) {  
    uint myId = blockDim.x*blockIdx.x + threadIdx.x;  
    if(myId < n)  
        y[myId] = a*x[myId] + y[myId];  
}
```

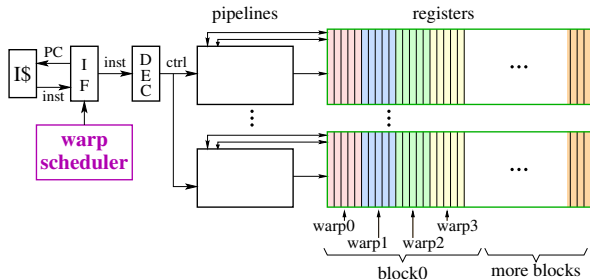
- A GPU with 3584 SPs, and a 1.6GHz clock rate (see [slide 2](#) can perform over 5700 single-precision GFlops.
 - ▶ With a main memory bandwidth of 484 GBytes/sec., and 4 bytes per `float`, a CUDA kernel needs to perform $\frac{3584 * 1.6 * 4}{484} \approx 48$ floating point operations per memory read or write.
 - ▶ Otherwise, memory bandwidth becomes the bottleneck.
- Registers and shared memory let us use a value many times without going to the off-chip, GDDR memory.
 - ▶ But, we need to program carefully to make this work.
- Is `saxpy` a good candidate for GPU execution?

Matrix Multiplication and Memory

```
for(int i = 0; i < M; i++)
    for(int j = 0; j < N; j++)
        for(int k = 0; k < L; k++)
            c[i,j] += a[i,k]*b[k,j];
```

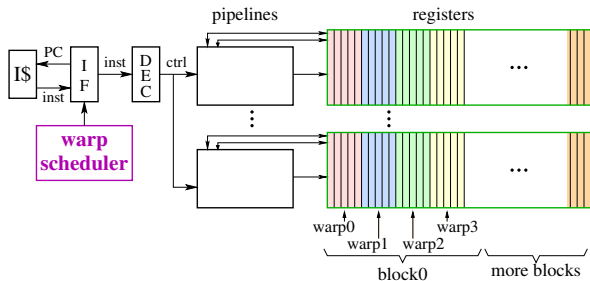
- Focus on the innermost loop: `for(k ...)`
 - ▶ Why?
- How many floating point operations per iteration?
- How many memory reads?
- How many memory writes?
- What is the “Compute-to-Global-Memory-Access” ratio (CGMA)?

Registers



- Each SP has its own register file.
- The register file is partitioned between threads executing on the SP.
- Local variables are placed in registers.
 - ▶ The compiler in-lines functions when it can
 - ★ A kernel with recursive functions or deeply nested calls can cause register spills to main memory – this is **slow**.
 - ▶ Local array variables are mapped to global memory – **watch out**.

More Registers



- In recent versions of CUDA, threads in the same warp can swap registers.
 - ▶ Provides very efficient intra-warp communication.
 - ▶ But not available in CUDA 2.1. ☹️
- Performance trade-offs
 - ▶ A thread can avoid slow, global memory accesses by keeping data in registers.
 - ▶ But, using too many registers reduces the number of threads that can run at the same time.

Registers and Memory Bandwidth

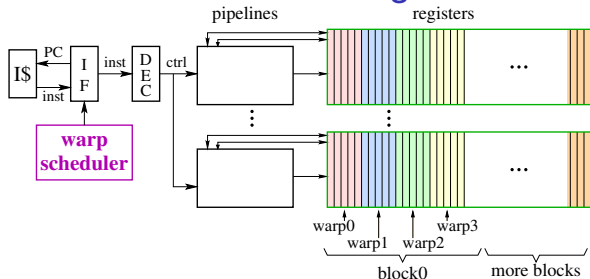
The GPU on [slide 2](#) has 28 SMs, each with 128 SPs.

- Each SP has access to a register file.
- I'll guess two register reads and one write per clock cycle, per SP.
- I'll assume 4-byte registers.
- We get

$$28\text{SM} * 128 \frac{\text{SP}}{\text{SM}} * 3 \frac{\text{RW}}{\text{SP} * \text{cycle}} * 1.6 \times 10^9 \frac{\text{cycle}}{\text{sec.}} * 4 \frac{\text{Byte}}{\text{RW}} = 68813 \frac{\text{GByte}}{\text{sec.}}$$

- 142 times faster than main memory bandwidth!
- ymmv: the GPUs in the linXX box are older.

Registers and Thread Scheduling



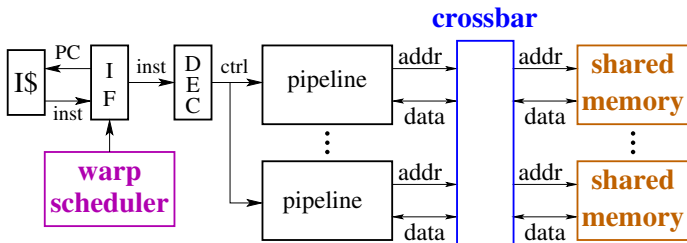
- Each SM has 256K registers, and 64 active warps, with 32 threads/warp.
 - ▶ That's 32 4-byte registers per thread.
- If a thread uses more registers
 - ▶ The SM cannot fully use its warp scheduler, or
 - ▶ Registers will spill to main memory – **slow**
- The numbers are smaller for older GPUs.
 - ▶ The GTX 550 Ti GPUs in the linXX boxes support 21 registers/thread.

Registers and Matrix Multiply

```
for(int i = 0; i < M; i +=2) {
    for(int j = 0; j < N; j +=2) {
        sum00 = sum01 = sum10 = sum11 = 0.0;
        for(int k = 0; k < L; k += 2) {
            a00 = a[i,k]; a01 = a[i,k+1];
            a10 = a[i+1,k]; a11 = a[i+1,k+1];
            b00 = b[k,j]; b01 = b[k,j+1];
            b10 = b[k+1,j]; b11 = b[k+1,j+1];
            sum00 += a00*b00 + a01*b10;
            sum01 += a00*b01 + a01*b11;
            sum10 += a10*b00 + a11*b10;
            sum11 += a10*b01 + a11*b11;
        }
        c[i,j] = sum00; c[i,j+1] = sum01;
        c[i+1,j] = sum10; c[i+1,j+1] = sum11;
    } }
}
```

- use a register to accumulate $c[i, j]$
- hold blocks of each matrix in main memory.
- can use each value loaded from $a[i, k]$ or $b[k, j]$ three or four times.
- What is the CGMA for the example above?

Shared Memory



- On-chip, one bank per SP.
- Banks are interleaved by:
 - ▶ Early CUDA GPUs: 4-byte word
 - ▶ Later GPUs: programmer configurable 4-byte or 8-byte words
 - ▶ Why?
- Shared memory is a limited resource: 48KBytes to 96Kbytes/SM.
 - ▶ Each SM has more registers than shared-memory.
 - ▶ Shared memory demands limit how many blocks can execute concurrently on a SM.

Shared Memory Example: Reduce

```
--shared-- float v[1024];

--device-- float f(float x) {
    return((5/2)*(x*x*x - x));
}

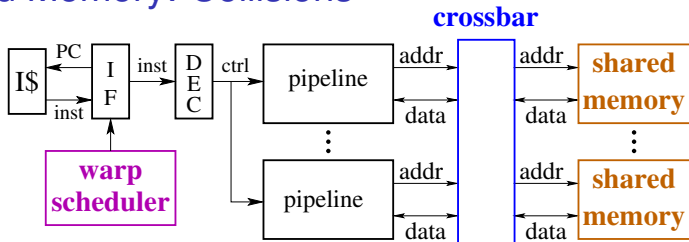
--device-- void compute_and_reduce(uint n, uint m, float *x) {
    uint myId = threadIdx.x;
    if(myId < n) {
        float y = x[myId];
        for(uint i = 0; i < m; i++)
            y = f(y); %  $y = f^{i+1}(x[\text{myId}])$ 
        v[myId] = y; % shared memory is much faster than global memory
        for(uint j = n >> 1; m > 0; m = n >> 1) { % reduce
            n -= j;
            __syncthreads();
            if(myId < j)
                v[myId] += v[myId+n];
        }
        x[myId] = v[myId]; % move result to global memory
    }
}
```

See notes on the next slide.

Notes on Reduce

- We calculate $f^m(x[\text{myId}])$ locally using the register variable y .
- For the reduce, we use the **shared** array v
 - ▶ This avoids the penalty of off-chip, global memory access for each step of the reduce.
 - ▶ All threads in a warp can access shared memory on the same cycle.
 - ▶ We can have multiple blocks running on multiple SMs
 - ★ Each SM has its own shared memory.
 - ★ Blocks running on different SMs in parallel can **all** access their shared memories in parallel.
 - ★ **But**, threads in one block do not share shared-memory with threads in other blocks.
 - ▶ To perform a reduce across blocks:
 - ★ Each block writes its subtotal to the **global** memory.
 - ★ The results from the blocks are combined on the host CPU or by launching a new kernel.
- At the end, we copy our value from shared memory to the global memory so the CPU or a subsequent kernel can access it.

Shared Memory: Collisions



- When one thread in a warp accesses shared memory, **all** active threads in the warp access shared memory.
- If each thread accesses a different bank, then all accesses are performed in a single cycle.
 - ▶ Otherwise, the load or store can take multiple cycles.
 - ▶ Multiple accesses to the same bank are called **collisions**.
 - ▶ The **worst-case** occurs when **all threads access the same bank**.
- The programmer needs to think about the index calculations to avoid collisions.
 - ▶ When programming GPUs, the programmer needs to think about index calculations a lot.

Shared Memory Example: Matrix Multiply

Running example from the textbook: $C = AB$

- Each thread-block loads a 16×16 block from A and B .
 - ▶ The threads to these loads “cooperatively”:
 - ▶ Read $A_{I,K}$ and $B_{K,J}$ from global memory with “coalesced” loads.
 - ▶ Write these blocks to shared-memory in a way that avoids bank conflicts.
- Compute: $C_{I,J} += A_{I,K}B_{K,J}$.
 - ▶ This takes $16^3 = 4096$ fused multiply-adds.
 - ▶ Loading $A_{I,K}$ fetches $16^2 = 256$ floats from global memory.
 - ▶ Likewise for $B_{K,J}$. Total of 512 floats fetched.
 - ▶ $CGMA = 4096/512 = 8$.
- Note: the L2 cache may help here: A and B are read-only.
 - ▶ Need to try more experiments.

Global Memory

- Off-chip DRAM
 - ▶ GDDR supports higher-bandwidth than than regular DDR.
 - ▶ A GPU can have multiple memory interfaces.
 - ▶ Total bandwidth 80 to 484+ GBytes/sec
- Memory accesses can be a big bottleneck.
 - ▶ CGMA: compute to global memory access ratio

Now for a word about DRAM

The memory that you plug into your computer is mounted on DIMMs (dual-inline memory modules).

- A DIMM typically has 16 or 18 chips
- E.g. each chip of an 8Gbyte DIMM holds 512MBytes = 4Gbits.
- Each chip consists of many “tiles”,
 - ▶ a typical chip has 1Mbit/tile
 - ▶ that's 4096 tiles for a 4Gbit chip.
- Each tile is an array of capacitors.
 - ▶ each capacitor holds 1 bit.
 - ▶ a typical tile could have 1024 rows and 1024 columns.

Writing and reading DRAM

- Writing: easy
 - ▶ drive all 1024 column-lines to the values you want to write.
 - ▶ open up all the valves for one row.
 - ▶ the drinking cups for each column in that row get filled or emptied.
 - ▶ note: you end up writing **every** column in the row; so writes are often preceded by reads.
- Reading: hard
 - ▶ drive all 1024 column-lines to “half-way”, and let them “float”.
 - ▶ open up all the valves for one row.
 - ▶ if the level in the pipe goes up a tiny amount, that cup held a 1.
 - ▶ if the level in the pipe goes down a tiny amount, that cup held a 0.
 - ▶ it's a delicate measurement – it takes time to set it up.
 - ▶ This is why DRAM is slow.
- But: we just read 1024 bits, from each chip of the DIMM.
 - ▶ That's 16Kbits = 2Kbytes total.
 - ▶ Conclusion: DRAM has awful latency, but we can get very good bandwidth.
 - ★ The bandwidth bottleneck is the wires from the DIMM to the CPU or GPU.
 - ★ But I'm pretty sure that Ian won't let me give a lecture on transmission lines, phase-locked loops, equalizers, and all the other cool stuff in the DDR (or GDDR) interface.

GPUs meet DRAM

- DRAM summary: terrible latency (60-200ns or more), fairly high bandwidth.
- The GPU lets the program take advantage of high bandwidth.
 - ▶ If the 32 loads from a warp access 32 consecutive memory location,
 - ★ The GPU does **one** GDDR access,
 - ★ and it transfers a large block of data.
 - ▶ The same optimization is applied to stores, and to loads from the on-chip caches.
- In CUDA-speak, if the loads from a warp access consecutive locations, we say that the memory accesses are **coalesced**.
- It's a big deal to make sure that your memory accesses are coalesced.
 - ▶ Note that the memory optimizations are exposed to the programmer.
 - ▶ You can get the performance by considering the memory model.
 - ▶ But, it's not automatic.

Example: Matrix Multiplication

- In C, matrices are usually stored in row-major order.
 - ▶ $A[i, k]$ and $A[i, k+1]$ are at adjacent locations, but
 - ▶ $B[k, j]$ and $B[k+1, j]$ are N words apart (for $N \times N$ matrices).
- For matrix multiplication, accesses to A are naturally coalesced, but accesses to B .
- The optimized code loads a block of B into shared memory.
 - ▶ This allows accesses to be coalesced.
 - ▶ But we need to be careful about how we store the data in the shared memory to avoid bank conflicts.

Other Memory

- Constant memory: cached, read-only access of global memory.
- Texture memory: global memory with special access operations.
- L1 and L2 caches: only for memory reads.

Summary

- GPUs can have thousands of execution units, but only a few off-chip memory interfaces.
 - ▶ This means that the GPU can perform 10-50 floating point operations for every memory read or write.
 - ▶ Arithmetic operations are very cheap compared with memory operations
- To mitigate the off-chip memory bottleneck
 - ▶ GPUs have, limited on-chip memory
 - ▶ Registers and the per-block, shared-memory will be our main concerns in this class.
- Moving data between different kinds of storage is the programmer's responsibility.
 - ▶ The programmer explicitly declares variables to be stored in shared memory.
 - ▶ The programmer needs to be aware of the per-thread register usage to achieve good SM utilization.
 - ▶ The only way to communicate between thread blocks is to write to global memory, end the kernel, and start a new kernel (ouch!)

Preview

March 15: CUDA Memory: examples

March 17: CUDA Performance

Reading [Kirk & Hwu](#) 3rd ed., Ch. 5 (Ch. 6 in 2nd ed.)

Mini Assignment Mini Assignment 5 due at 10am.

March 20: Matrix multiplication, Part 1

March 22: Matrix multiplication, Part 2

March 24: Complete CUDA

March 27 – April 3: this may change

March 27: Using Parallel Libraries

March 29 – April 3: Verification of/and Parallel Programs

April 5: Party: 50th Anniversary of Amdahl's Law

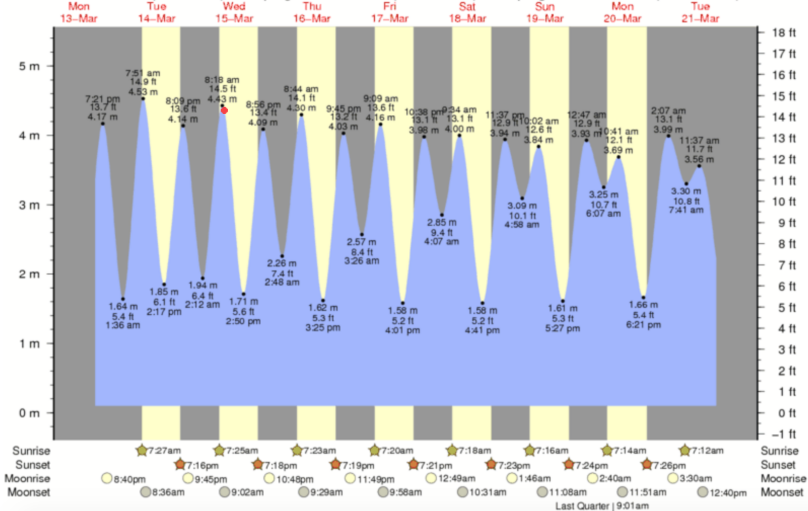
Review

- What is CGMA?
- On [slide 15](#) we computed the CGMA for matrix-multiplication using 16×16 blocks of the A , B , and C matrices.
 - ▶ How many such thread-blocks can execute concurrently on an SM with 48KBytes of memory?
 - ▶ How does the CGMA change if we use 32×32 blocks?
 - ▶ If we use the larger matrix-blocks, how many thread blocks can execute concurrently on an SM with 48Kbytes of memory?
 - ▶ If we use the larger matrix-blocks, how many thread blocks can execute concurrently on an SM with 96Kbytes of memory?
- What are bank conflicts?
- How can increasing the number of registers used by a thread improve performance?
- How can increasing the number of registers used by a thread degrade performance?
- What is a “coalesced memory access”?

Beware the Tides of March

Vancouver, British Columbia (max. tidal range 4.97m 16.3ft)

Times are PDT (UTC -7.0hrs). Last Spring Tide on Fri 10 Mar (h=4.64m 15.2ft). Next Spring Tide on Wed 29 Mar (h=4.56m 15.0ft)



From <https://www.tide-forecast.com/locations/Vancouver-British-Columbia/tides/latest>