GPUs



_AXPY loop (aX + Y)

modern ILP

 \bullet

load x	load y	inc i		
mul	load x	load y	inc i	
add	mul	load x	load y	
store	add	mul	inc i	
branch	store	add		
branch	store			



DLP approaches



load x	load x	c load x	
mul	mul	mul	
load y	load y	load y	
add	add	add	
store	store	store	

require additions to ISA *and* hardware

SPMD approach

Single *program*, Multiple data – spawn many versions of same program Note that this is a *programming* model, not a hardware model!

How threads are scheduled is invisible to programmer



????

In what ways is SPMD a more flexible model than the SIMD processors we saw last week?



 \mathbf{O}

 \odot

Sparse matrices

An *nxn* matrix is *sparse* if the number *m* of nonzero entries is a small fraction of the total

CSR (compressed sparse row) representation allows for easy access to nonzero elts



CSR can be computed from original matrix. Can the conversion be parallelized?

Multiplication in CSR



can turn

this

(outer

loop!!)

into

threads!



???

How can we run a SPMD program using SIMD ideas?

SIMT

Single instruction, multiple threads

- Each thread uses the same program memory, PC (but its own registers, FUs, stack pointer)
- Warp (Nvidia term) of threads executes in lockstep





image source

(

S...P?I?...M?...D?T?!? What???

Recap:

Flynn's taxonomy: describes hardware computation model SISD: single instruction, single data (traditional uniprocessor) SIMD: single instruction, multiple data (ISA/hardware for DLP) *MIMD*: multiple instruction, multiple data (multiprocessor) SPMD: single program, multiple data (describes programming model) SIMT: single instruction, multiple threads; SIMD-style computation (hardware) to implement SPMD operation

Large data

```
What do we do if n = 8192??
for (int i = 0; i < n; i++) {</pre>
    y[i] = a * x[i] + y[i];
}
                                                  loop
void myThread(int n, float a, ...) {
    int i = tid; // thread ID
                                                  GPU
    if (i < n)
        y[i] = a * x[i] + y[i];
ζ
```

SISD approach: 「_(ツ)_/

Vector approach: vector inner loop, scalar outer

SPMD approach: get a ← aka: spawn 8192 of those guys on beefy hardware

What's a GPU?

Graphics Processing Unit

Accelerator that originally helped CPU render 3d graphics

Predecessors: arcade game circuits, VGA controllers

Huge parallelism, programmability: attractive for largescale computations



3D Graphics Rendering Pipeline: Output of one stage is fed as input of the next stage. A vertex has attributes such as (x, y, z) position, color (RGB or RGBA), vertex-normal (n_x, n_y, n_z) , and texture. A primitive is made up of one or more vertices. The rasterizer raster-scans each primitive to produce a set of grid-aligned fragments, by interpolating the vertices.

What's a GPU? (CS1952y answer)

A heck of a lot of SIMT processors, arranged in groups

Nvidia terminology: Streaming Multiprocessors (SMs)



SFU Processor 0 Processor 1 Processor 2 Processor 6 Processor 6	SFU Processor 1 Processor 2 Processor 2	SFU Processor 0 Processor 1 Processor 2 Processor 2 Processor 2 Processor 6	SFU Processor 0 Processor 1 Processor 2 Processor 5 Processor 6	SFU Processor 0 Processor 1 Processor 2 Processor 2 Processor 6	SFU Processor 1 Processor 1 Processor 2 Processor 5 Processor 6	SFU Processor 0 Processor 1 Processor 2 Processor 6	SFU Processor 0 Processor 1 Processor 5 Processor 6 Processor 6
Shared Memory	Processor 7 Shared Memory Instruction Fetch Unit	Processor 7 Processor 7 Shared Memory Instruction Fetch Unit	Processor 7 Processor 7 Shared Memory Instruction Fetch Unit	Processor 7 Processor 7 Shared Memory Instruction Fetch Unit	Processor 7 Processor 7 Shared Memory Instruction Fetch Unit	Processor 3 Processor 7 Shared Memory Instruction Fetch Unit	Shared Memory Instruction Fetch Unit
Instruction and Data Caches (L1)	Instruction and Data Caches (L1)	Instruction and Data Caches (L1)	Instruction and Data Caches (L1)	Instruction and Data Caches (L1)	Instruction and Data Caches (L1)	Instruction and Data Caches (L1)	Instruction and Data Caches (L1)
SFU Processor 0 Processor 1 Processor 2 Processor 7 Shared Memory Instruction Fetch Unit	SFU Processor 0 Processor 1 Processor 2 Srbu Instruction Fetch Unit	SFU Processor 0 Processor 1 Processor 2 Processor 3 SFU Processor 3 SFU Processor 6 Processor 6 Processor 6 Processor 6 Processor 7	SFU Processor 0 Processor 1 Processor 2 Processor 3 SFU Processor 3 SFU Processor 6 Processor 7 Processor 7 Proces	SFU Processor 2 Processor 2 Processor 3 SFU Processor 3 Stared Memory Instruction Fetch Unit	SFU Processor 0 Processor 1 Processor 2 Processor 3 SFU Processor 5 Processor 5 Processor 6 Processor 6 Processor 6 Processor 7 Processor 6 Processor 7 Processor 6 Processor 7	SFU Processor 1 Processor 2 Processor 2 Processor 3 SFU Processor 5 Processor 7 Shared Memory Instruction Fetch Unit	SFU Processor 1 Processor 2 Processor 3 Processor 3 Shared Memory Instruction Fetch Unit
Instruction and Data Caches (L1)	Instruction and Data Caches (L1)	Instruction and Data Caches (L1)	Instruction and Data Caches (L1)	Instruction and Data Caches (L1)	Instruction and Data Caches (L1)	Instruction and Data Caches (L1)	Instruction and Data Caches (L1)



???

Warp size is 32 (2⁵) threads We want to run 8192 (2¹³) threads Do we need 256 SMs?

Programming model + terminology

CUDA (Compute Unified Device Architecture): API for programming Nvidia GPUs at the thread level

Programmers don't worry about warps (controlled by hardware)

Thread block: threads running on the same SM (scheduled by hardware)

Exposed to programmer – can make assumptions about eg shared memory within a block

Grid: entirety of thread workload

so how are warps managed within a block?



Finegrained multithreading, again

P&H fig. B.4.2

(Truly more like multiwarping)

Why not just run all of warp 1, then 2, then 3?

warp scheduler time warp 8 instruction 11 warp 1 instruction 42 warp 3 instruction 95 warp 8 instruction 12 warp 3 instruction 96 warp 1 instruction 43

Recap

Programmer

Chooses block size, # threads

Writes one program to run on many, many threads (SPMD)

Hardware

Schedules each block to an SM

Threads inside blocks are split up into warps to enable lockstep execution



SAXPY example rewritten in CUDA

// CPU side to invoke 8192 threads to run 8000 computations
__host__
// 32 blocks, 256 threads per block

```
myThread<<<32, 256>>>(8000, 2.0, x, y);
```

```
// GPU side
__device__
void myThread(int n, float a, float* x, float* y) {
    int i = blockIdx.x * blockDim.x + threadIdx.x;
    if (i < n)
        y[i] = a*x[i] + y[i];
}
```

