## **Execution on GPUs**



#### **GPU overview**

#### Programmer

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

Chooses block size, # threads

#### Hardware

Schedules each block to a Streaming Multiprocessor

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



## A note about nVidia

They're not paying me

These slides use nVidia/CUDA terminology for practicality (matches dominant platform)

See H&P Fig. 4.25 for translations to their terms and OpenCL

source

Is this good for the consumer/society? I'm not an economist but I'm happy to hear your takes



Total AIB share and units



source

## CUDA programming model

API for programming nVidia GPUs at the thread level

Allows for specifying host (CPU) code (to set up threads) + device (GPU) code (parallel **kernel**)

Host and device have separate memories

Note: this is not a graphics course, and not a GPU programming course (also GPUs are being used, developed, and marketed for more than graphics) We're studying the details of highly parallel architectures that use the SPMD programming model



### nVidia PTX <del>ISA</del>

PTX = parallel thread execution

Full summary in P&H fig. B.4.3, documentation here

Suffixes may define operand data type, memory space

Arithmetic: add, sub, mul, etc

Special function: sqrt, sin, cos, etc

Logical: and, or, xor, etc

Memory: Id, st, tex (texture lookup), atom

Control: branch, call, ret, sync, exit

#### Why isn't it a "true" ISA?

- Actual machine instructions are proprietary/may differ by device – device has to translate PTX to its own instructions
- Virtual registers

## **CSR** multiplication in CUDA

```
host
// set up matrix in CSR format here: ...
// 8 blocks, 256 threads per block to do multiplication
csrMult<<<8, 256>>>(2048, Rp, C, V, x, y);
// GPU side
device
void csrMult(int n, int* Rp, int* C, float* V, float* x, float* y) {
                                                        Instead of loop, use these
    int r = blockIdx.x * blockDim.x + threadIdx.x;
                                                         provided variables as
    if (r < n) {
                                                       per-thread "coordinates"
        int rBeg = Rp[r];
                                                          for data access
        int rSize = Rp[r + 1] - Rp[r];
        y[r] = multRow(rSize, C + rBeg, V + rBeg, x);
```

## $\mathbf{O}$

# ???

What do we need to watch out for when programming using CUDA?

#### Hardware-aware SPMD

SPMD promises us the power of writing scalar programs and running them with high parallelism

It's (relatively) easy to write a SPMD program that runs, BUT:





#### **Kernel execution**



### Active threads

```
if (threadIdx.x < 4) {
    A;
    B;
} else {
    X;
    Y;
}
Z;</pre>
```

<u>image source</u>

X; Y; A; B; Z;

How to keep track of which threads are active? How to keep track of when to reconverge? Time

#### **Execution** mask

Use single PC and keep track of which threads are using that PC



#### Execution mask & loops





# ???

Where might execution mask usage become complicated?

 $\odot$ 

Allows switching between execution paths ...how?

ī.

#### nVidia Volta



<u>image source</u>

#### **Tensor cores**

An accelerator on the accelerator (see also Ray Tracing cores)

Circuits on GPUs optimized for AI (matrix operation D = A \* B + C)

Uses mixed-precision to speed up math, reduce memory demands

Claim huge (8x-32x per generation) performance gains over SM matrix multiply



#### <u>image source</u>

#### **Sparse acceleration**

 $\odot$ 

Sparse Tensor Input activations Core Select ш Dot-product = zero entry TTT Fine-grained Compress structured pruning 2:4 sparsity: 2 nonzero out of 4 entries **Output activations Dense trained** Non-zero Indices data values weights Fine-tuned sparse and **Fine-tune weights** compressed weights

1

7

<u>source</u>

source

#### **GA100** SM in Ampere





