

A TYPICAL GPU PROGRAM

- ① CPU ALLOCATES STORAGE ON GPU *cudaMalloc*
- ② CPU COPIES INPUT DATA FROM CPU → GPU *cudaMemcpy*
- ③ CPU LAUNCHES KERNEL(S) ON GPU TO PROCESS THE DATA
kernel launch
- ④ CPU COPIES RESULTS BACK TO CPU FROM GPU
cudaMemcpy

CONFIGURING THE KERNEL LAUNCH

SQUARE <<< 1, 64 >>> (d_out, d_in)

NUMBER OF
BLOCKS

THREADS PER
BLOCK

128 THREADS? SQUARE <<< 1, 128 >>> (...)

DEFINING THE GPU COMPUTATION

BIG IDEA

THIS IS
IMPORTANT

KERNELS LOOK LIKE SERIAL PROGRAMS

WRITE YOUR PROGRAM AS IF IT WILL RUN ON ONE THREAD

THE GPU WILL RUN THAT PROGRAM ON MANY THREADS

MAKE SURE YOU UNDERSTAND THIS

WHAT IS THE GPU GOOD AT ?

- ☐ LAUNCHING A SMALL NUMBER OF THREADS EFFICIENTLY
- ☒ LAUNCHING A LARGE NUMBER OF THREADS EFFICIENTLY
- ☐ RUNNING ONE THREAD VERY QUICKLY
- ☐ RUNNING ONE THREAD THAT DOES LOTS OF WORK IN PARALLEL
- ☒ RUNNING A LARGE NUMBER OF THREADS IN PARALLEL

WHAT IS THE GPU GOOD AT?

① EFFICIENTLY LAUNCHING LOTS OF THREADS

② RUNNING LOTS OF THREADS IN PARALLEL

You may be used to other programming environments where launching threads is an expensive process.

SIMPLE EXAMPLE:

IN: FLOAT ARRAY $[0 \ 1 \ 2 \ \dots \ 63]$

OUT: FLOAT ARRAY $[0 \ 1^2 \ 2^2 \ \dots \ 63^2]$
 $[0 \ 1 \ 4 \ 9 \ \dots \]$

KERNEL: SQUARE

- ① CPU
- ② GPU (theory, no code)
- ③ GPU (code)

We're going to start by looking at how we'd run this code on the CPU.

CPU CODE: SQUARE EACH ELEMENT OF AN ARRAY

```
for (i=0; i<64; i++) {  
    out[i] = in[i] * in[i];  
}
```

CPU CODE: SQUARE EACH ELEMENT OF AN ARRAY

```
for (i=0; i < 64; i++) {  
    out[i] = in[i] * in[i];  
}
```

(1) ONLY ONE THREAD
OF EXECUTION

("thread" = "one
independent path of
execution through the
code")

(2) NO EXPLICIT PARALLELISM

CPU CODE: SQUARE EACH ELEMENT OF AN ARRAY

```
for (i=0; i < 64; i++) {  
    out[i] = in[i] * in[i];  
}
```

QUIZ:

HOW MANY MULTIPLICATIONS?

1 * TAKES 2 ns.

HOW LONG TO EXECUTE?

64

128

(1) ONLY ONE THREAD
OF EXECUTION

("thread" = "one
independent path of
execution through the
code")

(2) NO EXPLICIT PARALLELISM

GPU CODE: A HIGH-LEVEL VIEW

CPU

ALLOCATE MEMORY

COPY DATA TO/FROM GPU

LAUNCH KERNEL

↑
SPECIFIES DEGREE
OF PARALLELISM

GPU

Express $OUT = IN \cdot IN$

↑
SAYS NOTHING
ABOUT THE DEGREE
OF PARALLELISM.

CPU code: squareKernel <<< 64 >>> (outArray, inArray)

BUT HOW DOES IT WORK IF I LAUNCH 64 INSTANCES OF THE SAME PROGRAM?

CPU LAUNCHES 64 THREADS:

64



I KNOW I'M
THREAD 0!

...

I KNOW I'M
THREAD 63!

"WORK ON
ITEM N
OF
THE ARRAY!"

BUT HOW DOES IT WORK IF I LAUNCH 64 INSTANCES OF THE SAME PROGRAM?

CPU LAUNCHES 64 THREADS:

64

I KNOW I'M
THREAD 0!

I KNOW I'M
THREAD 63!

"WORK ON
ITEM N
OF
THE ARRAY!"

QUIZ:

HOW MANY MULTIPLICATIONS?

64

IF EACH MULT. TAKES
10 NS, HOW LONG
FOR THE ENTIRE
COMPUTATION?

10


```
8
9 int main(int argc, char ** argv) {
10     const int ARRAY_SIZE = 64;
11     const int ARRAY_BYTES = ARRAY_SIZE * sizeof(float);
12
13     // generate the input array on the host
14     float h_in[ARRAY_SIZE];
15     for (int i = 0; i < ARRAY_SIZE; i++) {
16         h_in[i] = float(i);
17     }
18     float h_out[ARRAY_SIZE];
19
20     // declare GPU memory pointers
21     float * d_in;
22     float * d_out;
23
24     // allocate GPU memory
25     cudaMalloc((void **) &d_in, ARRAY_BYTES);
26     cudaMalloc((void **) &d_out, ARRAY_BYTES);
27
28     // transfer the array to the GPU
29     cudaMemcpy(d_in, h_in, ARRAY_BYTES, cudaMemcpy::host_to_device);
30
31     // launch the kernel
```

Data on the CPU, the host, starts with h underscore. Data on the GPU, the device, starts with d underscore.

```
27 // transfer the array to the GPU
28 cudaMemcpy(d_in, h_in, ARRAY_BYTES, ???);
29
30 // launch the kernel
31 cube<<<1, ARRAY_SIZE>>>(d_out, d_in);
32
33 // copy back the result array to the CPU
34 cudaMemcpy(h_out, d_out, ARRAY_BYTES, ???);
35
36
```

1

cudaMemcpyHostToDevice

2

cudaMemcpyDeviceToHost

CONFIGURING THE KERNEL LAUNCH

SQUARE<<<1, 64>>> (d_out, d_in)

NUMBER OF
BLOCKS

THREADS PER
BLOCK

- (1) CAN RUN MANY BLOCKS AT ONCE
- (2) MAXIMUM NUMBER OF THREADS/BLOCK <
 - 512 (OLDER GPUS)
 - 1024 (NEWER GPUS)

Two, each block has a maximum number of threads that it can support.

CONFIGURING THE KERNEL LAUNCH

SQUARE <<< 1, 64 >>> (d_out, d_in)

NUMBER OF
BLOCKS

THREADS PER
BLOCK

128 THREADS?

SQUARE <<< 1, 128 >>> (...)

1280 THREADS?

SQUARE <<< 10, 128 >>> (...)

SQUARE <<< 5, 256 >>> (...)

~~SQUARE <<< 1, 1280 >>> (...)~~

But we can't call square 1,1280 because that's too many threads per block.

CONFIGURING THE KERNEL LAUNCH

SQUARE <<< 1, 64 >>> (d.out, d.in)

NUMBER OF
BLOCKS

THREADS PER
BLOCK

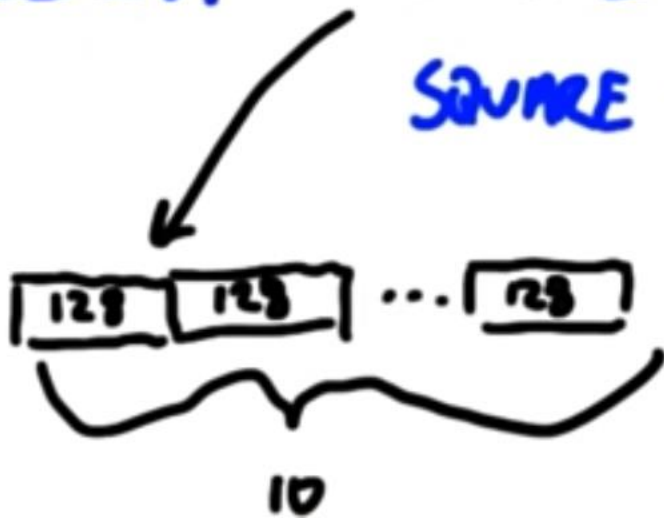
128 THREADS?

SQUARE <<< 1, 128 >>> (...)

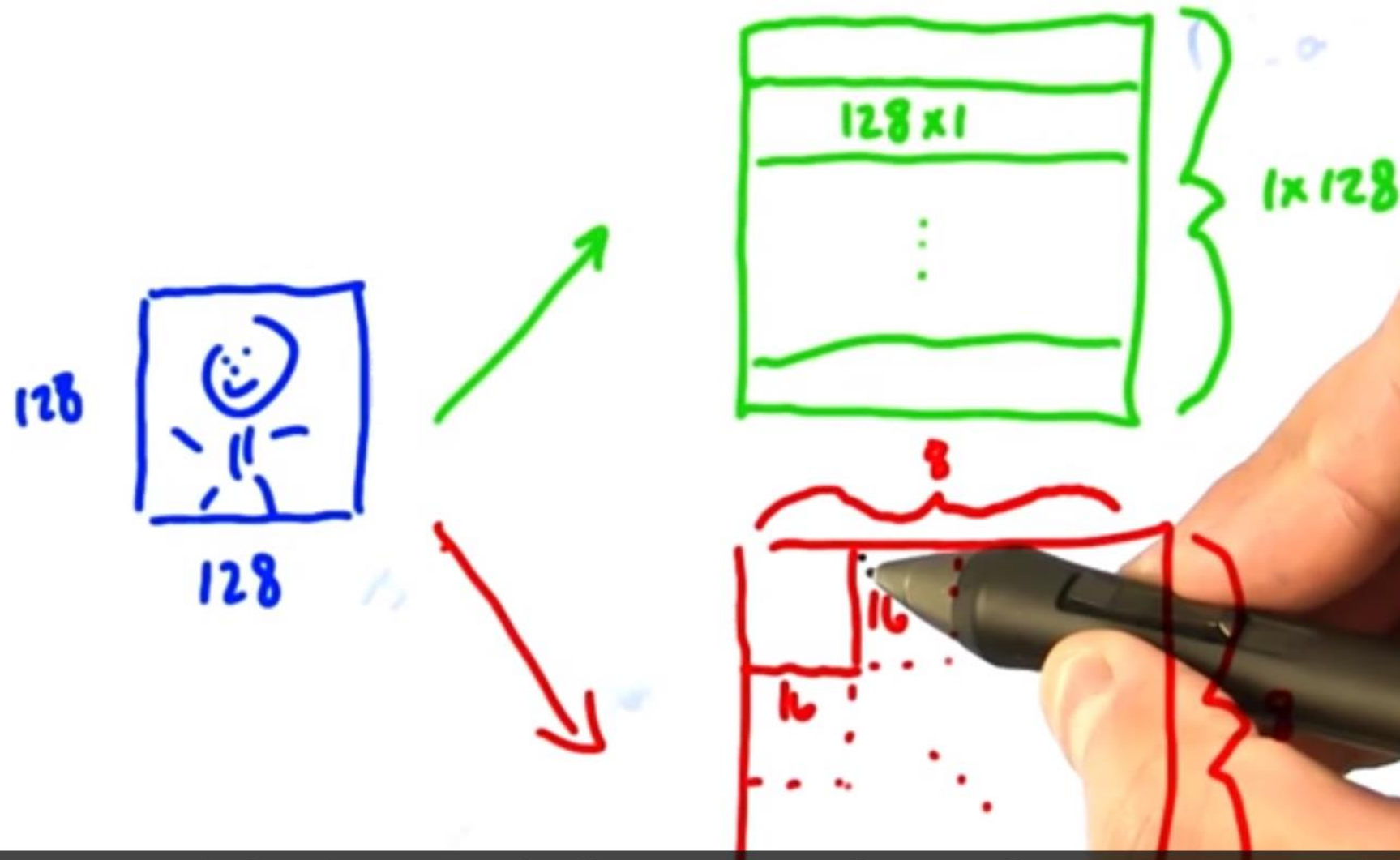
1280 THREADS?

SQUARE <<< 10, 128 >>> (...)

SQUARE <<< 5, 256 >>> (...)



CONFIGURING THE KERNEL LAUNCH



Or we might instead choose to launch an 8 x 8 grid of blocks where each block is 16 threads by 16 threads.

CONFIGURING THE KERNEL LAUNCH

KERNEL <<< GRID OF BLOCKS, BLOCK OF THREADS >>> (...)



1, 2, or 3 D

So CUDA supports 1, 2, or 3 dimensional thread blocks.

CONFIGURING THE KERNEL LAUNCH

KERNEL <<< GRID OF BLOCKS, BLOCK OF THREADS >>> (...)

↓
1, 2, or 3D

↓
1, 2, or 3D

We can also arrange thread blocks into 1, 2, or 3 dimensional grids.

CONFIGURING THE KERNEL LAUNCH

KERNEL <<< GRID OF BLOCKS, BLOCK OF THREADS >>> (...)

↓
1, 2, or 3D

↓
1, 2, or 3D

$\text{dim3}(x, y, z)$

$\text{dim3}(w, 1, 1) == \text{dim3}(w) == w$

CONFIGURING THE KERNEL LAUNCH

KERNEL <<< GRID OF BLOCKS, BLOCK OF THREADS >>> (...)

square <<< dim3(bx, by, bz), dim3(tx, ty, tz), shmem >>> (...)

↓
grid of blocks
bx · by · bz

↓
block of threads
tx · ty · tz

↓
shared
memory
per
block in
bytes

The most general kernel launch we can do looks like thi:, square of 3 parameters.

CONFIGURING THE KERNEL LAUNCH

KERNEL <<< GRID OF BLOCKS, BLOCK OF THREADS >>> (...)

square <<< dim3(bx, by, bz), dim3(tx, ty, tz), shmem >>> (...)

grid of blocks

$bx \cdot by \cdot bz$

block of threads

$tx \cdot ty \cdot tz$

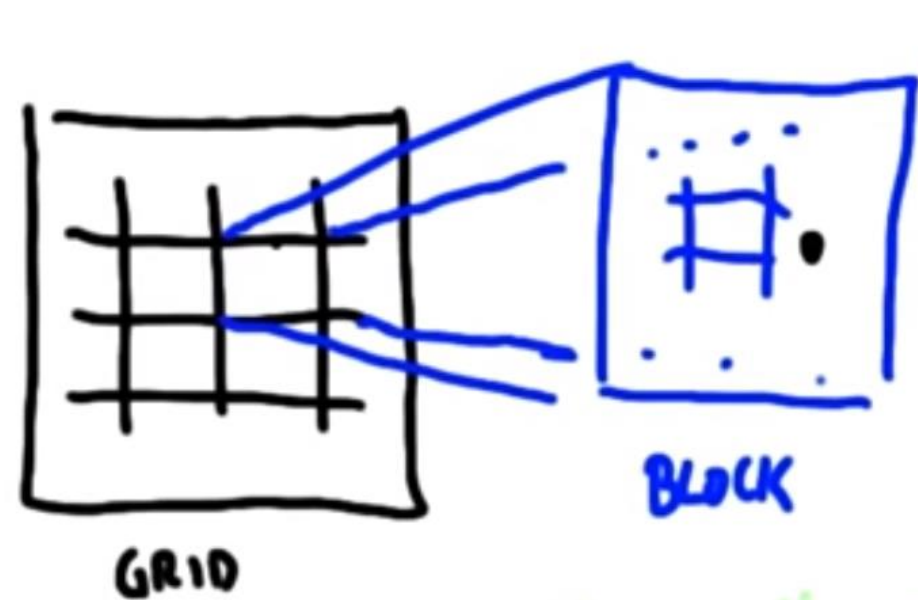
shared
memory

per
block in
bytes

CONFIGURING THE KERNEL LAUNCH

KERNEL <<< GRID OF BLOCKS, BLOCK OF THREADS >>> (...)

square <<< dim3(bx, by, bz), dim3(tx, ty, tz), stream >>> (...)



thread ldx : thread within block
thread ldx.x thread ldx.y

block Dim : size of a block

block ldx : block within grid

gridDim : size of grid

CONFIGURING THE KERNEL LAUNCH

KERNEL <<< GRID OF BLOCKS, BLOCK OF THREADS >>> (...)

square <<< dim3(bx, by, bz), dim3(tx, ty, tz), shmem >>> (...)

Quiz

kernel <<< dim3(8, 4, 2), dim3(16, 16) >>> (...)

How many blocks?

64

How many threads/block?

256

How many total threads?

16384

LESSONS FOR TODAY: WHAT WE KNOW

- WE WRITE A PROGRAM THAT LOOKS LIKE IT RUNS ON ONE THREAD
- WE CAN LAUNCH THAT PROGRAM ON ANY NUMBER OF THREADS
- EACH THREAD KNOWS ITS OWN INDEX IN THE BLOCK + THE GUI

MAP

- SET OF ELEMENTS TO PROCESS
- FUNCTION TO RUN ON EACH ELEMENT

[64 FLOATS]

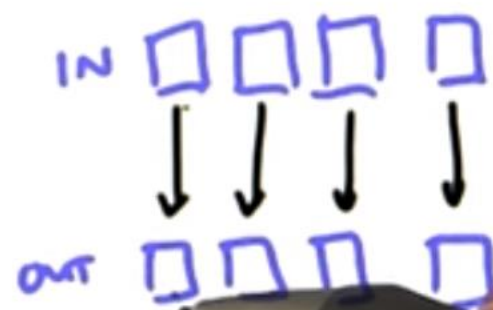
["SQUARE"]

MAP (ELEMENTS, FUNCTION)

GPUS ARE GOOD AT MAP

- GPUS HAVE MANY PARALLEL PROCESSORS
- GPUS OPTIMIZE FOR THROUGHPUT

MAP'S
COMMUNICATION
PATTERN



Quiz CHECK THE PROBLEMS THAT CAN BE SOLVED USING MAP.



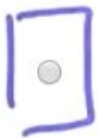
SORT AN INPUT ARRAY



ADD ONE TO EACH ELEMENT IN AN INPUT ARRAY



SUM UP ALL ELEMENTS IN AN INPUT ARRAY



COMPUTE THE AVERAGE OF AN INPUT ARRAY

WHAT WE LEARNED

- TECHNOLOGY TRENDS
 - THROUGHPUT VS LATENCY
 - GPU DESIGN GOALS
 - GPU PROGRAMMING MODEL
 - WITH EXAMPLE!
 - MAP
- 