

TDT4200 Parallel Programming

Bart van Blokland

Lecture 9

&group

Last Time:

Im in ur CPU

Nao in ur GPU

Two important things:

- Differences between a CPU and GPU
 - Threads, Warps, and the Streaming Multiprocessor

The CPU can be faster than the GPU.

The CPU can be faster than the GPU.



Professional sprinter:

- + Covers small distances super quickly
- Takes a while to run 1000 km



City Marathon:

- Individual people are kind of slow
- + Combined distance of runners sums up to 1000km very quickly

The CPU can be faster than the GPU.



CPU:

- + Single thread performance is amazing
- Throughput could be higher



GPU:

- Single thread performance is terrible
- + High throughput

PCI Express 3.0 Host Interface

GigaThread Engine

Memory Controller

Memory Controller

Memory Controller

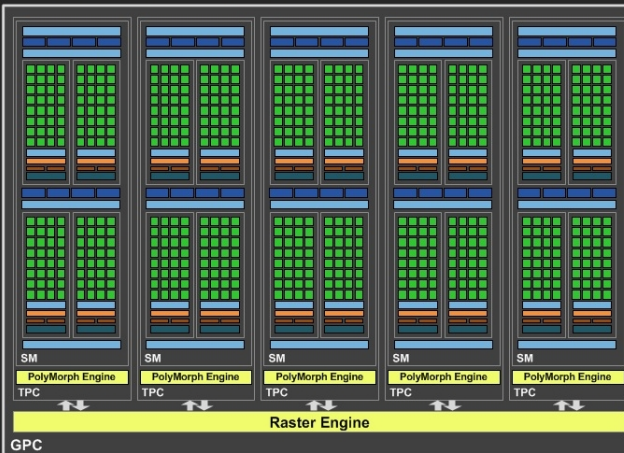
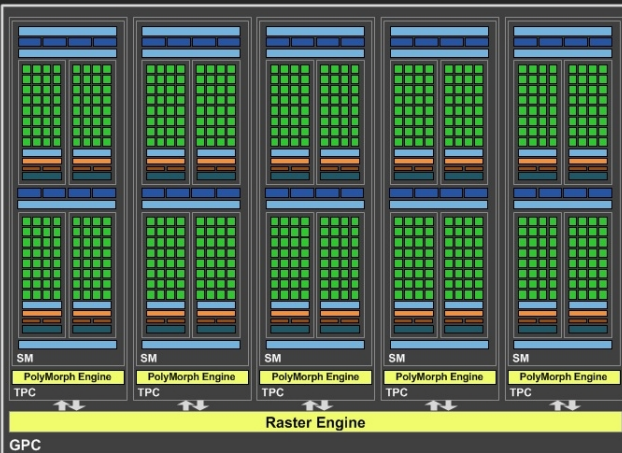
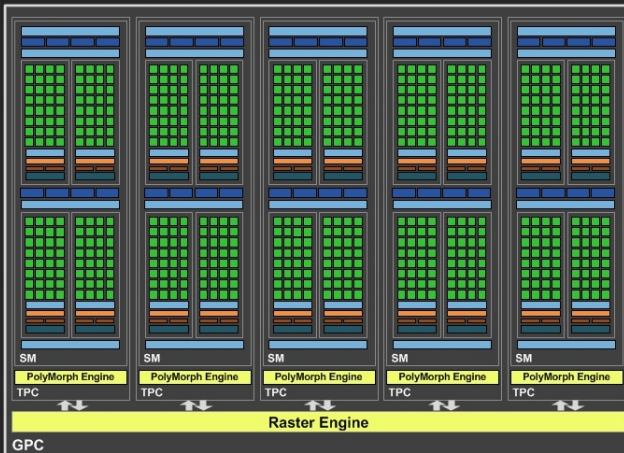
Memory Controller

Memory Controller

Memory Controller



L2 Cache



Memory Controller

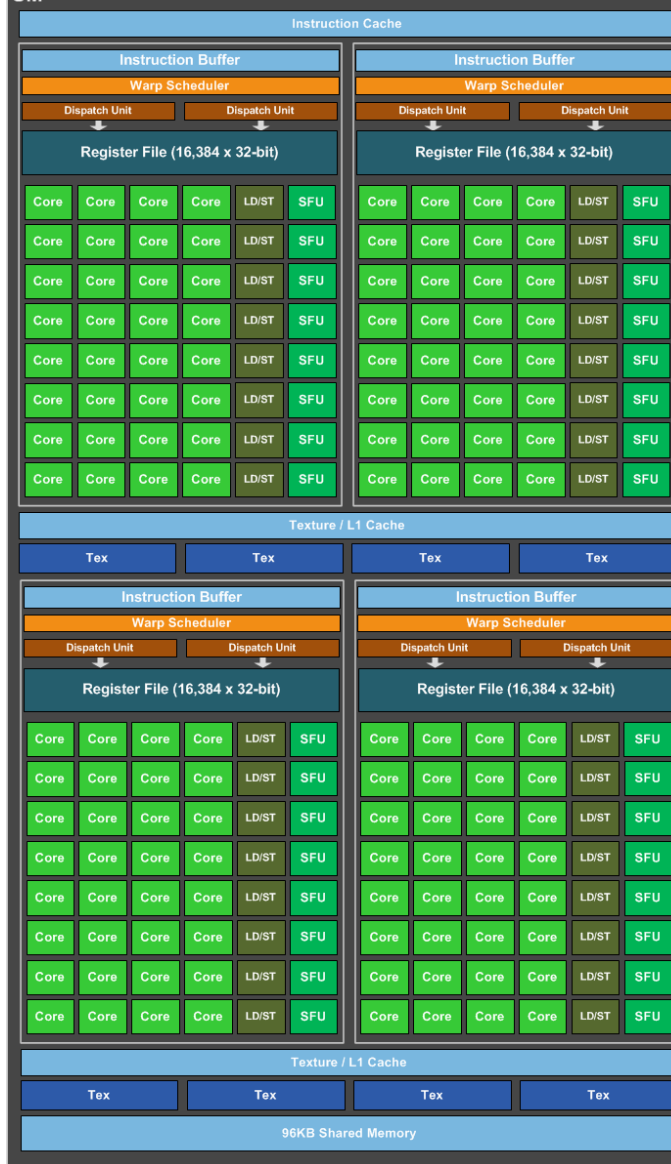
Memory Controller

Memory Controller

Memory Controller

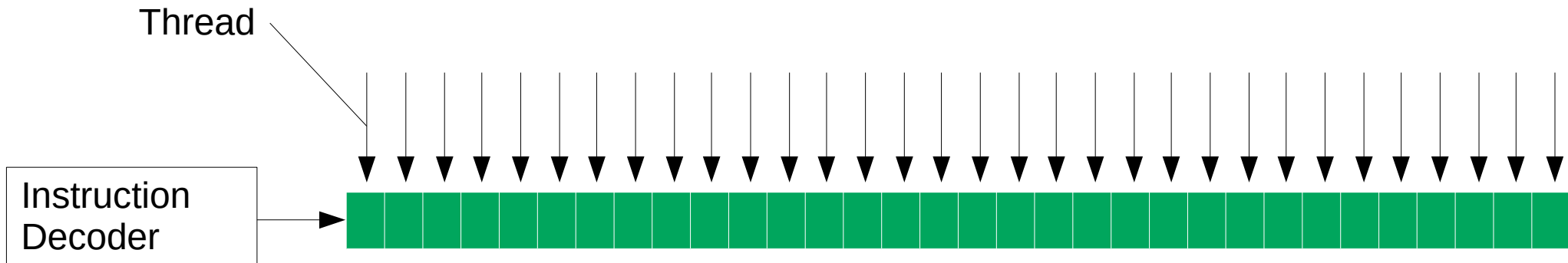
Memory Controller

Memory Controller

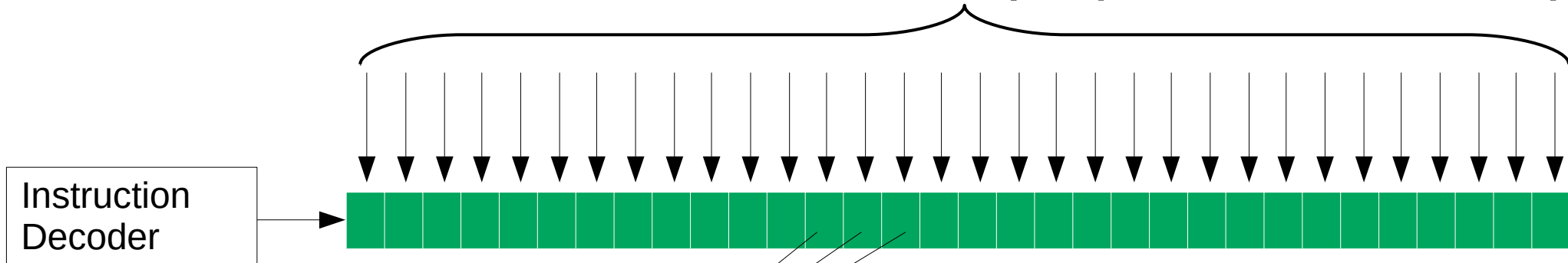


Roughly a single core
CPU with a bunch of
ALU's and FPU's!





“Warp” (= GPU “thread”)



FPU's and ALU's

Today:

Let's write a CUDA program!

```
#include <stdio.h>
```

```
__global__ void runMeOnGPU() {  
    printf("I'm thread %i!\n", threadIdx.x);  
}
```

```
int main() {  
    runMeOnGPU<<<1, 32>>>();  
    cudaDeviceSynchronize();  
    return 0;  
}
```

```
#include <stdio.h>
```

```
__global__ void runMeOnGPU() {  
    printf("I'm thread %i!\n", threadIdx.x);  
}
```

```
int main() {  
    runMeOnGPU<<<1, 32>>>();  
    cudaDeviceSynchronize();  
    return 0;  
}
```

“Kernel”



```
#include <stdio.h>
```

```
__global__ void runMeOnGPU() {  
    printf("I'm thread %i!\n", threadIdx.x);  
}
```

```
int main() {  
    runMeOnGPU<<<1, 32>>>();  
    cudaDeviceSynchronize();  
    return 0;  
}
```

```
#include <stdio.h>
```

```
__global__ void runMeOnGPU() {  
    printf("I'm thread %i!\n", threadIdx.x);  
}
```

```
int main() {  
    runMeOnGPU<<<1, 32>>>();  
    cudaDeviceSynchronize();  
    return 0;  
}
```

The CUDA API distinguishes
the CPU from the GPU:

Host: CPU side

Device: GPU side

GPU { `__device__ void doSomething();`
`__global__ void mainKernel();`

CPU { `__host__ void cpuFunction();`

Host and device functions cannot call each other.

They live a lonely life.


```
#include <stdio.h>
```

```
__global__ void runMeOnGPU() {  
    printf("I'm thread %i!\n", threadIdx.x);  
}
```

```
int main() {  
    runMeOnGPU<<<1, 32>>>();  
    cudaDeviceSynchronize();  
    return 0;  
}
```

Kernel launches are asynchronous

→ CPU can do other stuff after kicking off a job

.. But needs to wait if it requires the results.

```
#include <stdio.h>
```

```
__global__ void runMeOnGPU() {  
    printf("I'm thread %i!\n", threadIdx.x);  
}
```

```
int main() {  
    runMeOnGPU<<<1, 32>>>>();  
    cudaDeviceSynchronize();  
    return 0;  
}
```

Thread Hierarchy

Thread Hierarchy

```
for(int index = 0; index < 1000; index++) {  
    bigArray[index] = computeSomething();  
}
```

Thread Hierarchy

```
for(int index = 0; index < 1000; index++) {  
    bigArray[index] = computeSomething();  
}
```



```
int index = threadIndex;  
bigArray[index] = computeSomething();
```

X 1000

Done, right?

Note quite.



There is no time to run a scheduling algorithm.

→ All scheduling runs in hardware



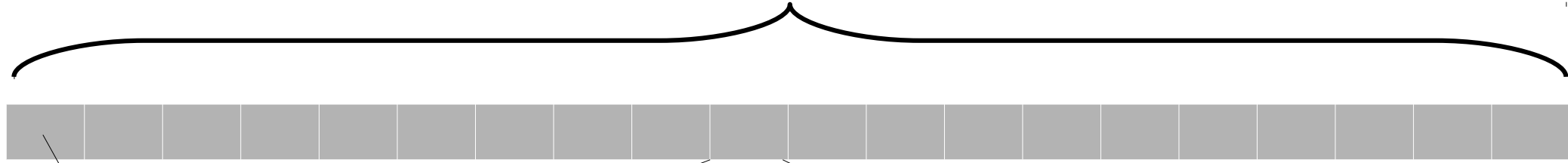
There is no time to run a scheduling algorithm.

→ All scheduling runs in hardware

Compromise: threads are grouped in “Blocks”.

Blocks contain a number of warps, and are scheduled on

Grid



Block



Thread

Blocks and grids can be launched in multiple dimensions:

`blockIdx.x`

`blockIdx.y`

`blockIdx.z`

`threadIdx.x`

`threadIdx.y`

`threadIdx.z`

Blocks and grids can be launched in multiple dimensions:

`blockIdx.x`

`blockIdx.y`

`blockIdx.z`

`threadIdx.x`

`threadIdx.y`

`threadIdx.z`

You can request the block and grid sizes at runtime:

`gridDim.x`

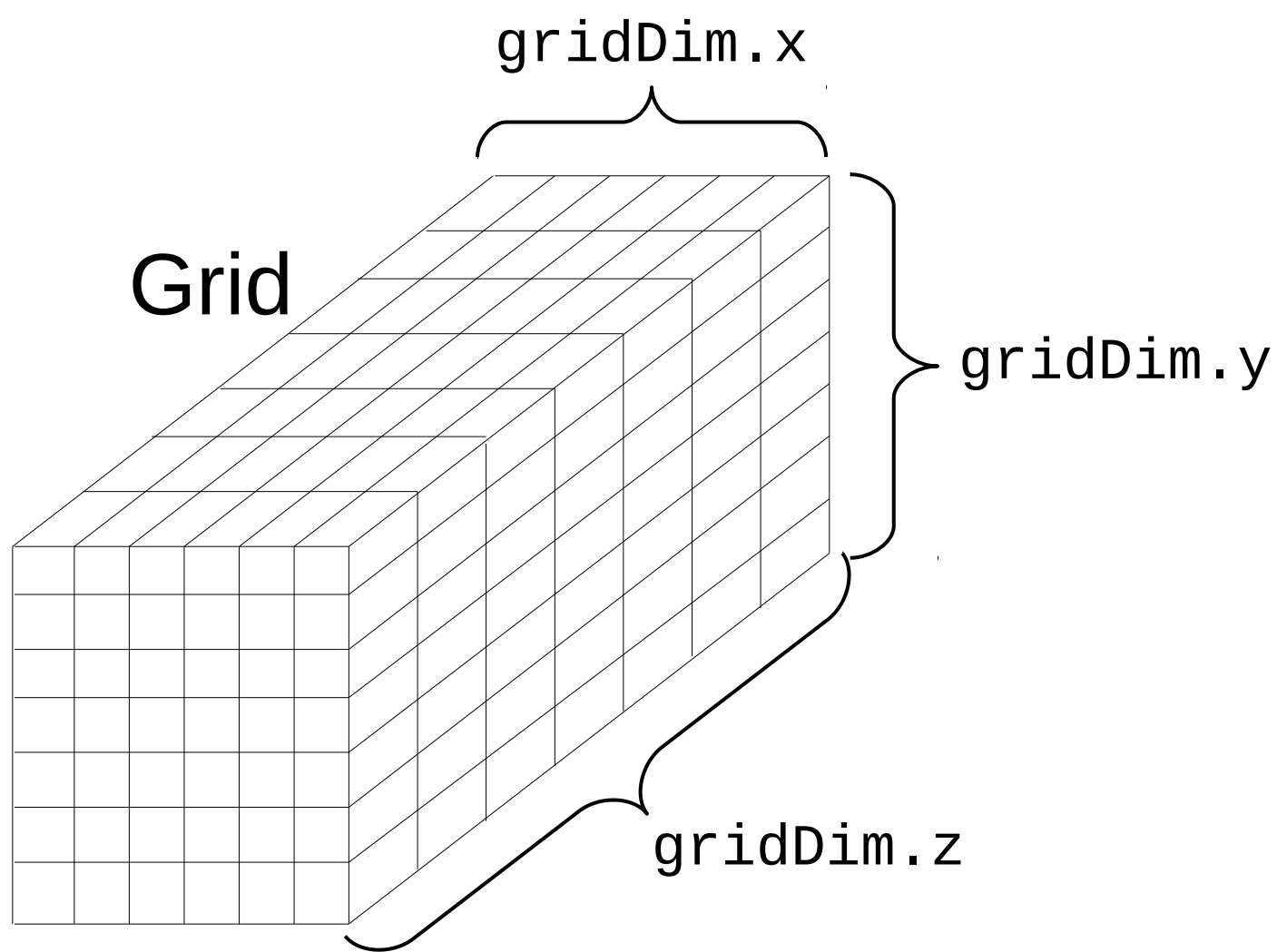
`gridDim.y`

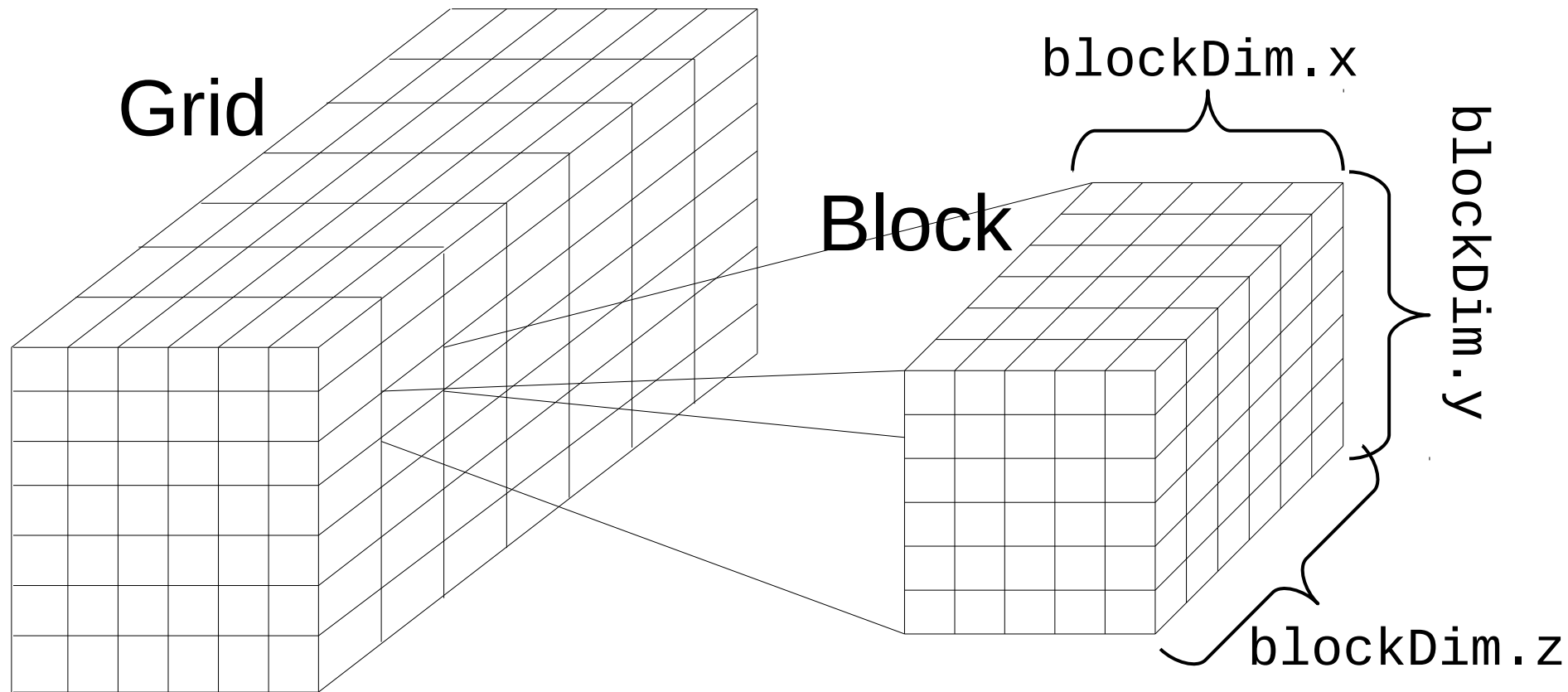
`gridDim.z`

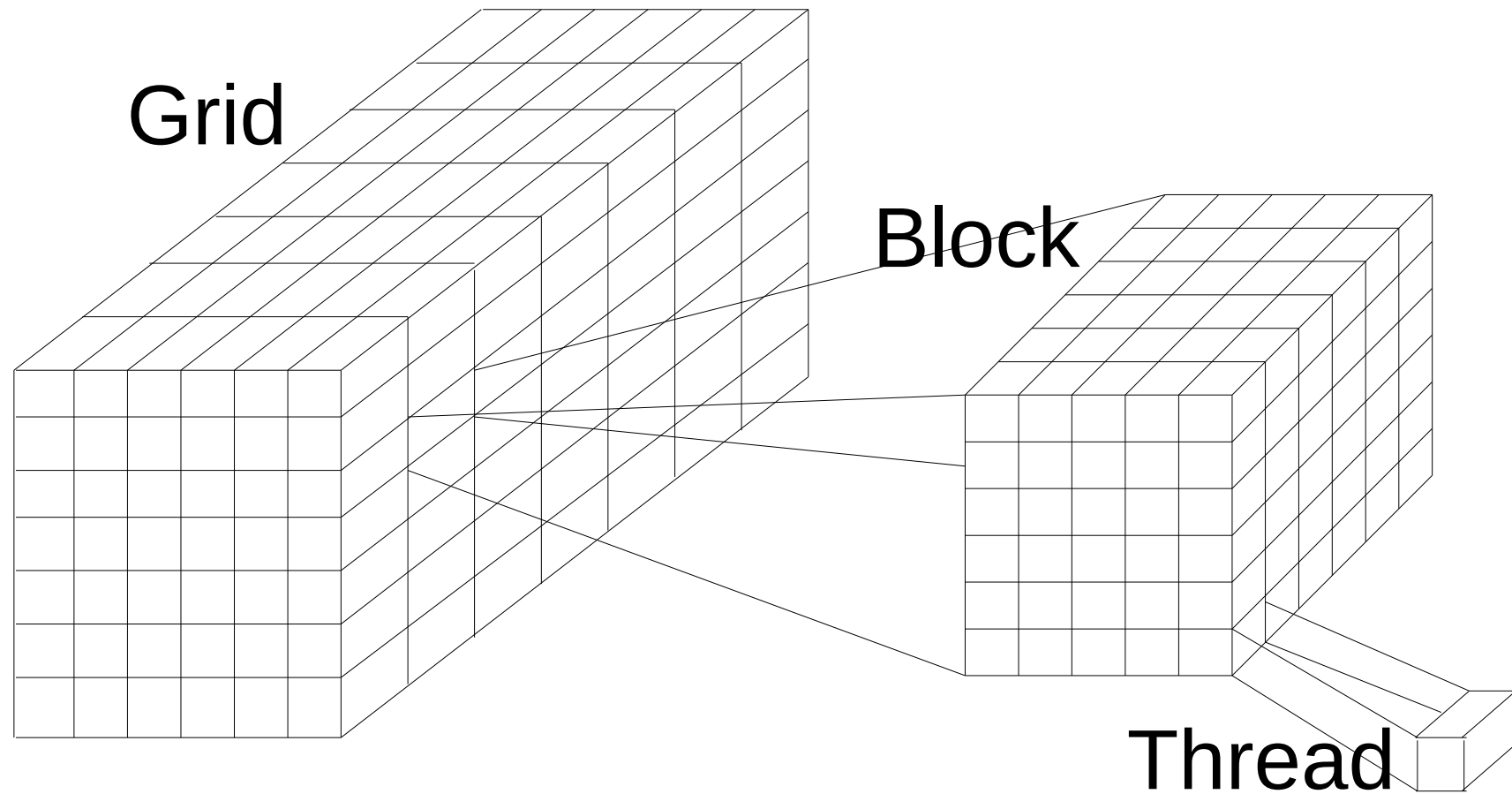
`blockDim.x`

`blockDim.y`

`blockDim.z`



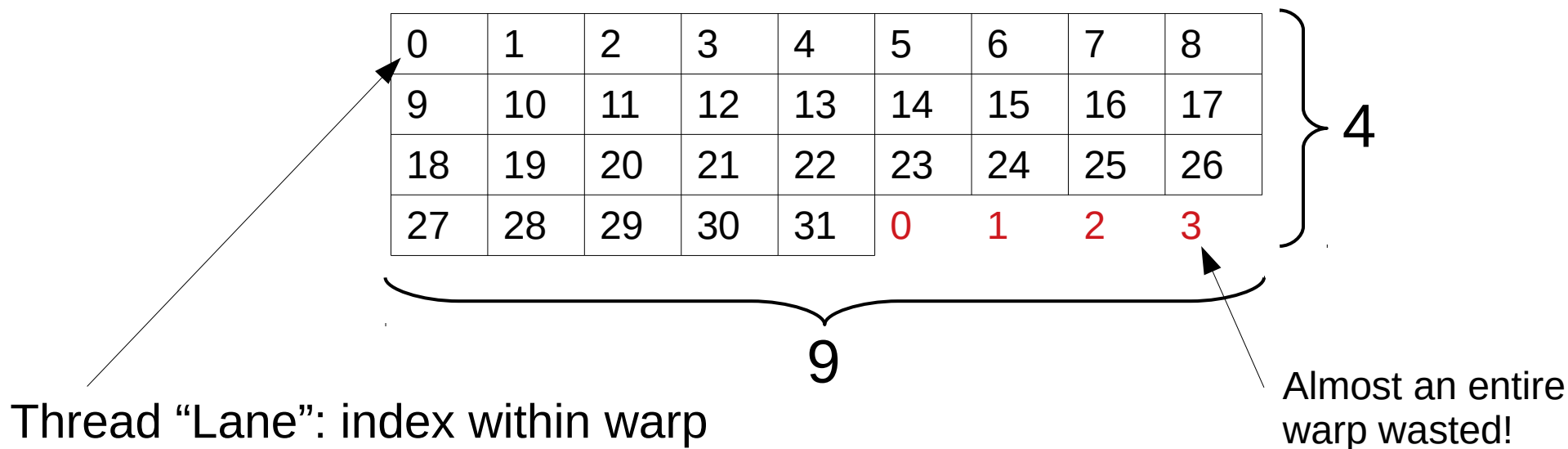




Caveats

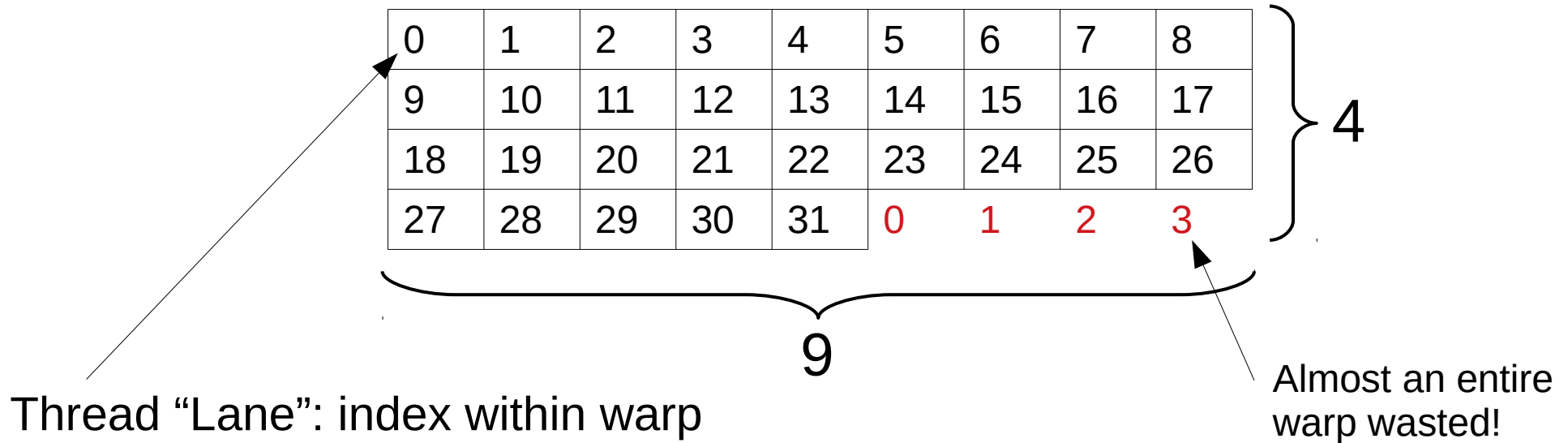
Caveats

Threads in a block are “flattened”, and chopped up into warps.



Caveats

Threads in a block are “flattened”, and chopped up into warps.



$\text{blockDim.x} * \text{blockDim.y} * \text{blockDim.z}$
should normally be a multiple of 32!

Caveats

Limits

For recent cards:

- Max grid size: $[2^{31} - 1, 65535, 65535]$
- Max block size: $[1024, 1024, 64]$
- Max threads per block: 1024
- Max number of warps per SM: 64
- Max number of threads per SM: 2048

Launching a kernel launches a grid of blocks

Each block contains threads,
rounded up to the nearest multiple of 32.

Blocks are assigned to SM's.

Warps execute on SM's.

```
#include <stdio.h>
```

```
__global__ void runMeOnGPU() {  
    printf("I'm thread %i!\n", threadIdx.x);  
}
```

```
int main() {  
    runMeOnGPU<<<1, 32>>>();  
    cudaDeviceSynchronize();  
    return 0;  
}
```

```
#include <stdio.h>
```

```
__global__ void runMeOnGPU() {  
    printf("I'm thread %i!\n", threadIdx.x);  
}
```

```
int main() {  
    dim3 grid(2, 1, 1);  
    dim3 block(3, 1, 1);  
    runMeOnGPU<<<grid, block>>>();  
    cudaDeviceSynchronize();  
    return 0;  
}
```

```
#include <stdio.h>
```

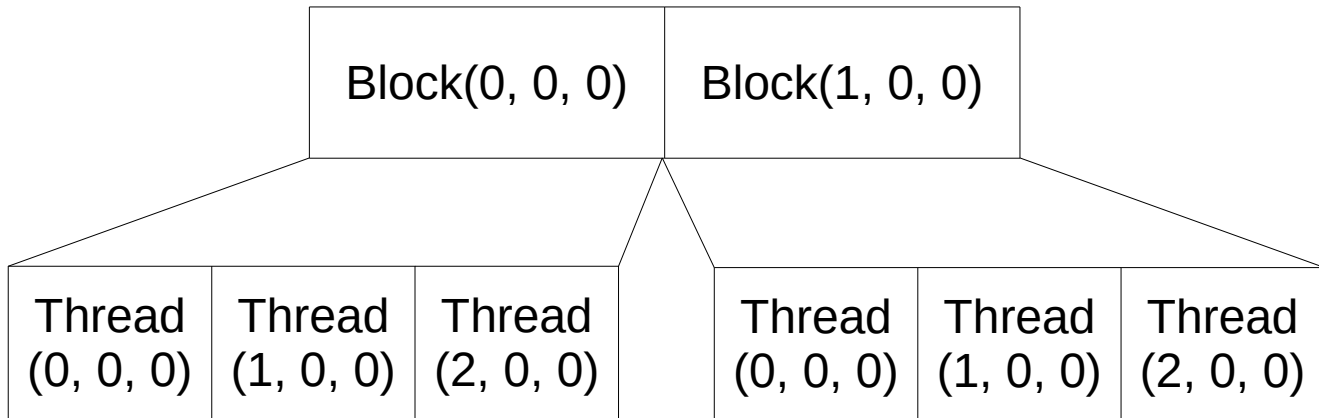
```
__global__ void runMeOnGPU() {  
    printf("I'm thread %i!\n", threadIdx.x);  
}
```

```
int main() {  
    dim3 grid(2, 1, 1);  
    dim3 block(3, 1, 1);  
    runMeOnGPU<<<grid, block>>>();  
    cudaDeviceSynchronize();  
    return 0;  
}
```

Small problem..

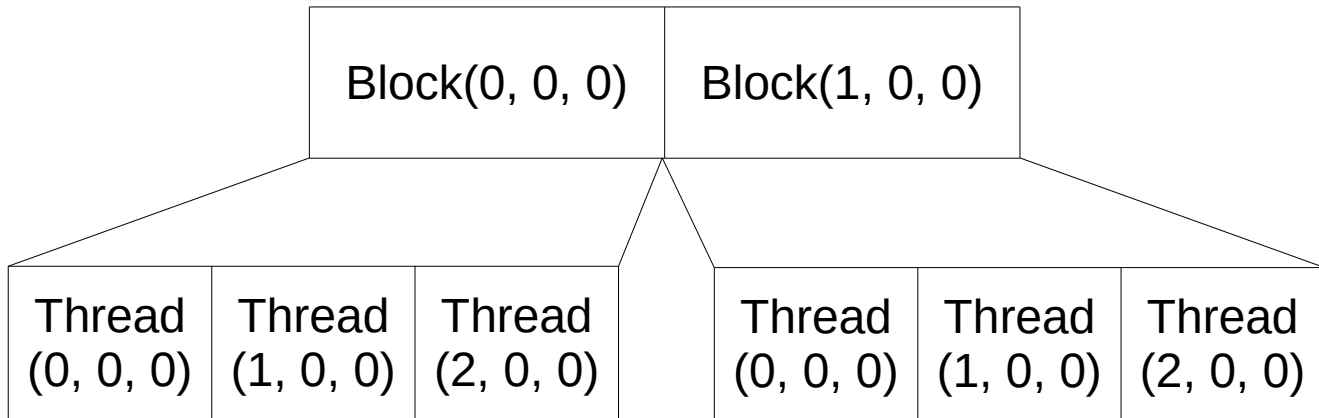
```
I'm thread 0!  
I'm thread 1!  
I'm thread 2!  
I'm thread 0!  
I'm thread 1!  
I'm thread 2!
```

Fortunately, all blocks have equal sizes



```
dim3 grid(2, 1, 1);  
dim3 block(3, 1, 1);
```


Fortunately, all blocks have equal sizes



```
int threadIndex = blockIdx.x * blockDim.x + threadIdx.x
```

```
dim3 grid(2, 1, 1);  
dim3 block(3, 1, 1);
```

Let's practice this a bit!

For the next problems, determine:

- the launch parameters (block and grid size)
- how each thread computes its index.

You'd like to process a 1D array

What grid and block sizes do you use?

How do you compute the pixel index?

Limits

- Max grid size: $[2^{31} - 1, 65535, 65535]$
- Max block size: $[1024, 1024, 64]$
- Max threads per block: 1024
- Max number of warps per SM: 64
- Max number of threads per SM: 2048

Additional info

- Each thread processes exactly one element in the array

<code>gridDim.x</code>	<code>blockIdx.x</code>
<code>gridDim.y</code>	<code>blockIdx.y</code>
<code>gridDim.z</code>	<code>blockIdx.z</code>

<code>blockDim.x</code>	<code>threadIdx.x</code>
<code>blockDim.y</code>	<code>threadIdx.y</code>
<code>blockDim.z</code>	<code>threadIdx.z</code>

You'd like to process a 2D image of an arbitrary size.

What grid and block sizes do you use?

How do you compute the pixel index?

Limits

- Max grid size: $[2^{31} - 1, 65535, 65535]$
- Max block size: $[1024, 1024, 64]$
- Max threads per block: 1024
- Max number of warps per SM: 64
- Max number of threads per SM: 2048

Additional info

- Each thread processes exactly one pixel

<code>gridDim.x</code>	<code>blockIdx.x</code>
<code>gridDim.y</code>	<code>blockIdx.y</code>
<code>gridDim.z</code>	<code>blockIdx.z</code>

<code>blockDim.x</code>	<code>threadIdx.x</code>
<code>blockDim.y</code>	<code>threadIdx.y</code>
<code>blockDim.z</code>	<code>threadIdx.z</code>

You'd like to process a sequence
of 16x16 pixel images

What grid and block sizes do you use?

How do you compute the image index and pixel coord?

Limits

- Max grid size: $[2^{31} - 1, 65535, 65535]$
- Max block size: $[1024, 1024, 64]$
- Max threads per block: 1024
- Max number of warps per SM: 64
- Max number of threads per SM: 2048

Additional info

- Each thread processes exactly one pixel
- Images are stored sequentially on a row by row basis.

<code>gridDim.x</code>	<code>blockIdx.x</code>
<code>gridDim.y</code>	<code>blockIdx.y</code>
<code>gridDim.z</code>	<code>blockIdx.z</code>

<code>blockDim.x</code>	<code>threadIdx.x</code>
<code>blockDim.y</code>	<code>threadIdx.y</code>
<code>blockDim.z</code>	<code>threadIdx.z</code>

You'd like to process a sequence
of 14x14 pixel images

What grid and block sizes do you use?

How do you compute the image index and pixel coord?

Limits

- Max grid size: $[2^{31} - 1, 65535, 65535]$
- Max block size: $[1024, 1024, 64]$
- Max threads per block: 1024
- Max number of warps per SM: 64
- Max number of threads per SM: 2048

Additional info

- Each thread processes exactly one pixel
- Images are stored sequentially on a row by row basis.

<code>gridDim.x</code>	<code>blockIdx.x</code>
<code>gridDim.y</code>	<code>blockIdx.y</code>
<code>gridDim.z</code>	<code>blockIdx.z</code>

<code>blockDim.x</code>	<code>threadIdx.x</code>
<code>blockDim.y</code>	<code>threadIdx.y</code>
<code>blockDim.z</code>	<code>threadIdx.z</code>

```
#include <stdio.h>
```

```
__global__ void runMeOnGPU() {  
    printf("I'm thread %i!\n", threadIdx.x);  
}
```

```
int main() {  
    runMeOnGPU<<<1, 32>>>();  
    cudaDeviceSynchronize();  
    return 0;  
}
```


But how do we compile it?

```
#include <stdio.h>
```

```
__global__ void runMeOnGPU() {  
    printf("I'm thread %i!\n", threadIdx.x);  
}
```

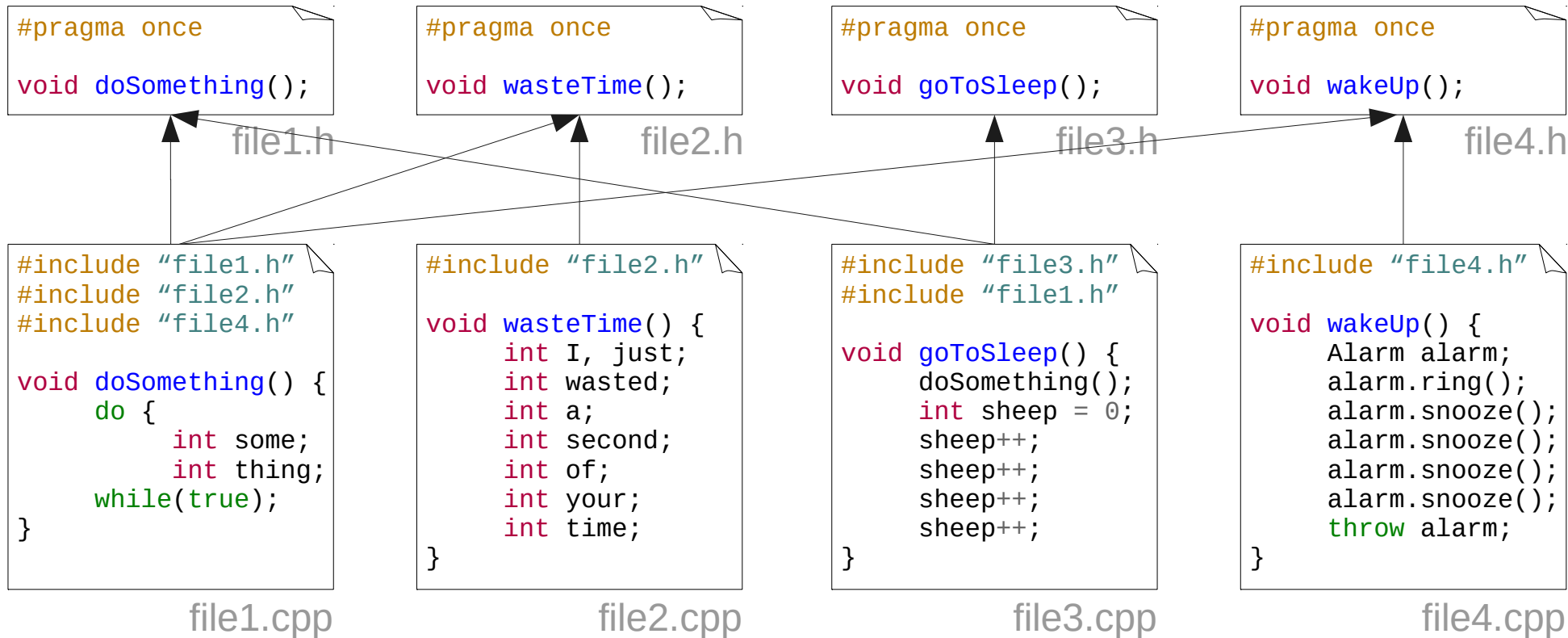
```
int main() {  
    runMeOnGPU<<<1, 32>>>();  
    cudaDeviceSynchronize();  
    return 0;  
}
```

NVCC (the CUDA Compiler)

```
nvcc sample.cu -o runongpu
```

```
nvcc sample.cu -o runongpu
```

What does GCC
have to do with that?



```
#pragma once
void doSomething();
```

file1.h

```
#pragma once
void wasteTime();
```

file2.h

```
#pragma once
void goToSleep();
```

file3.h

```
#pragma once
void wakeUp();
```

file4.h

```
#include "file1.h"
#include "file2.h"
#include "file4.h"

void doSomething() {
    do {
        int some;
        int thing;
        while(true);
    }
}
```

file1.cpp

```
#include "file2.h"

void wasteTime() {
    int I, just;
    int wasted;
    int a;
    int second;
    int of;
    int your;
    int time;
}
```

file2.cpp

```
#include "file3.h"
#include "file1.h"

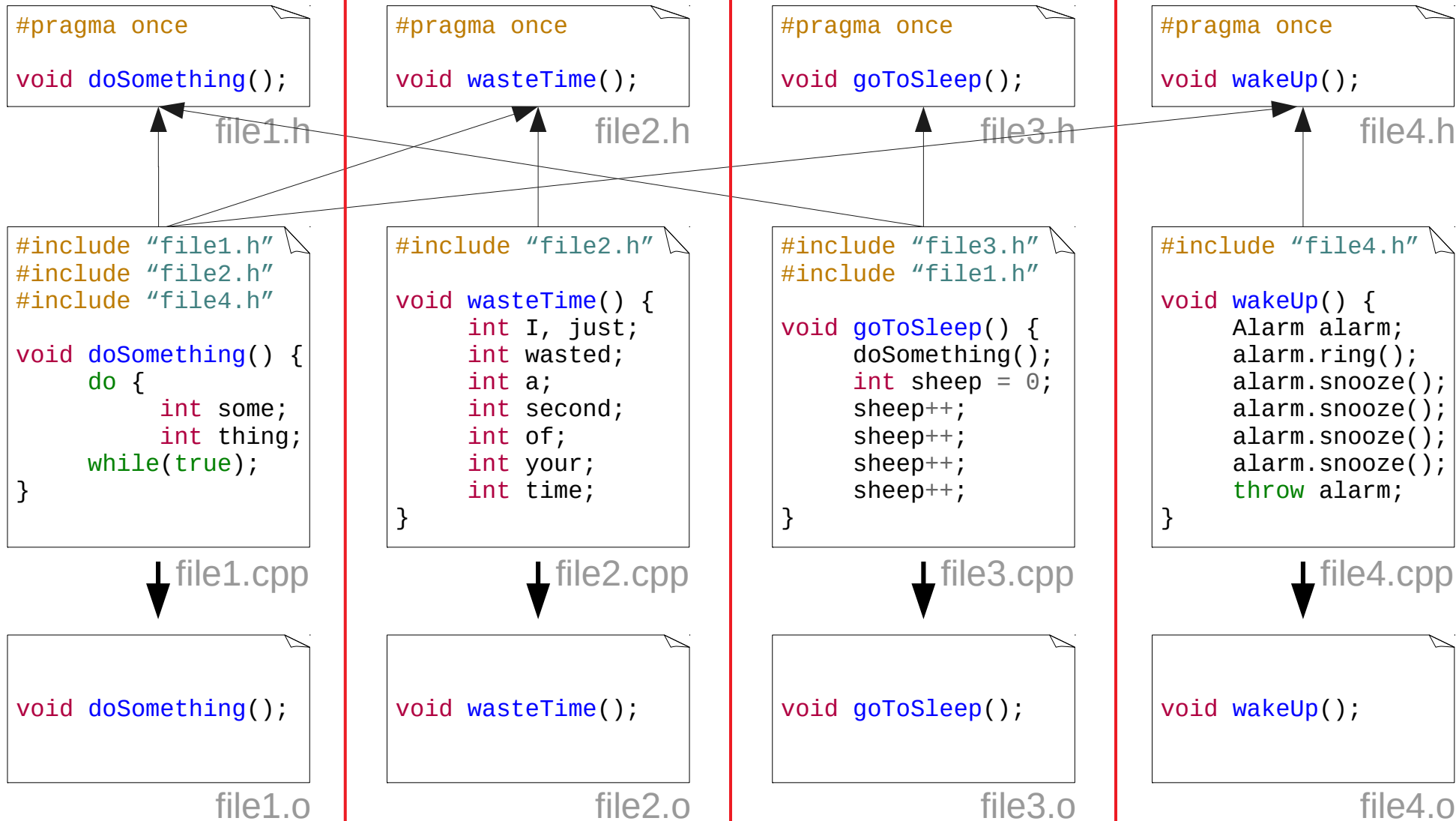
void goToSleep() {
    doSomething();
    int sheep = 0;
    sheep++;
    sheep++;
    sheep++;
    sheep++;
}
```

file3.cpp

```
#include "file4.h"

void wakeUp() {
    Alarm alarm;
    alarm.ring();
    alarm.snooze();
    alarm.snooze();
    alarm.snooze();
    alarm.snooze();
    throw alarm;
}
```

file4.cpp



```
void doSomething();
```

file1.o

```
void wasteTime();
```

file2.o

```
void goToSleep();
```

file3.o

```
void wakeUp();
```

file4.o

Executable



Linker



```
void doSomething();
```

file1.o

```
void wasteTime();
```

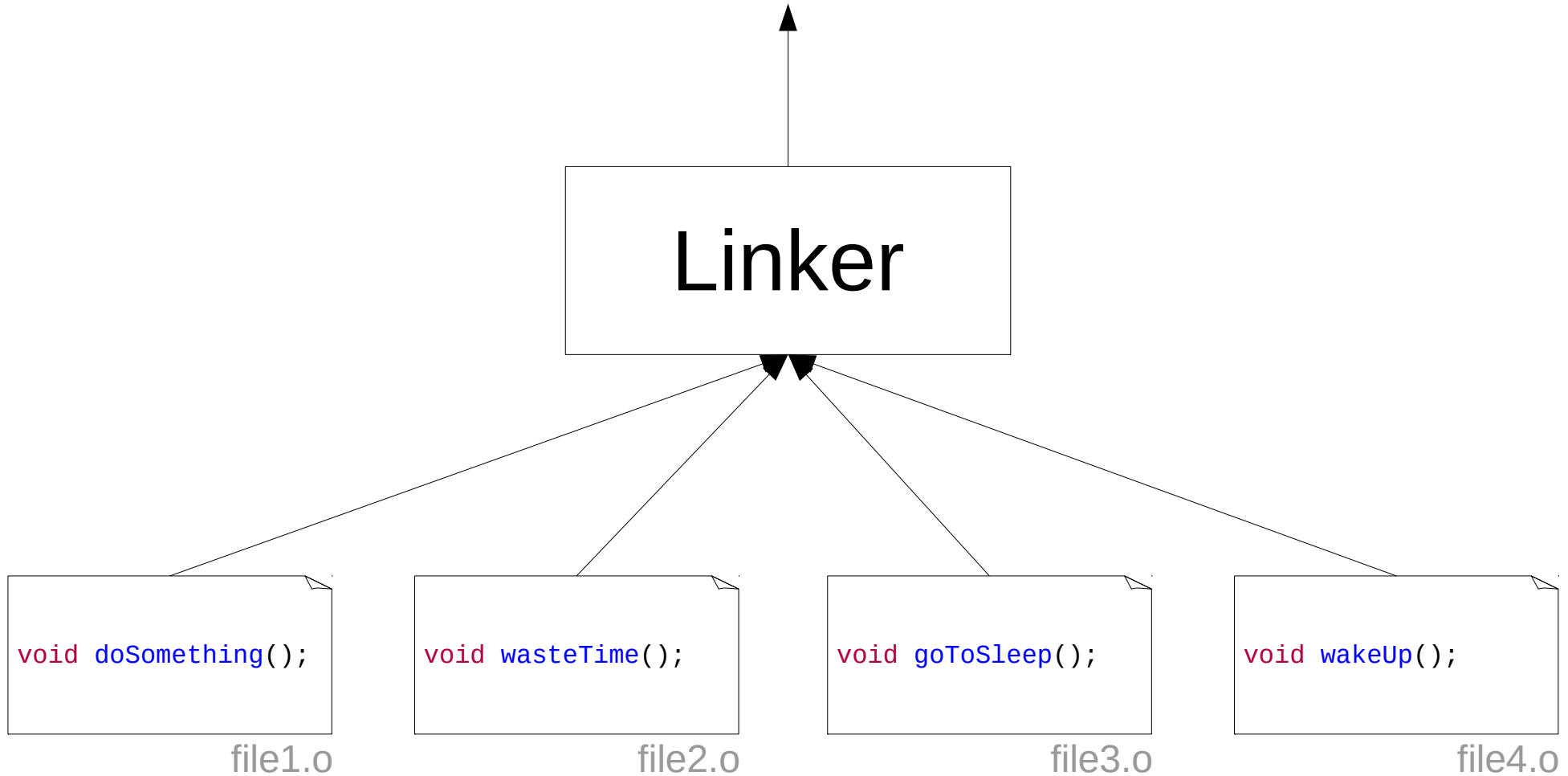
file2.o

```
void goToSleep();
```

file3.o

```
void wakeUp();
```

file4.o



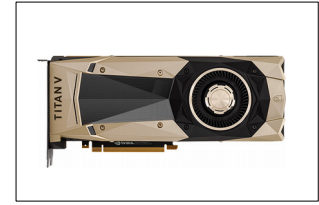
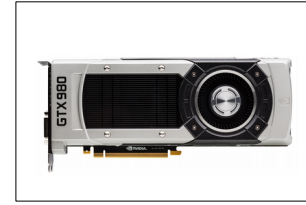
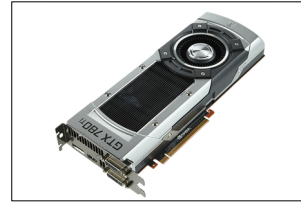
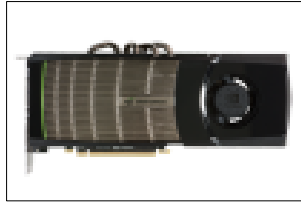
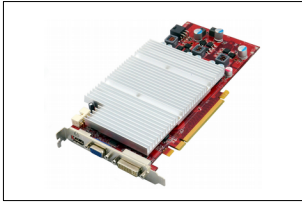
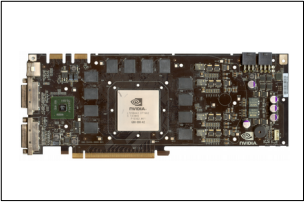
```
#include <stdio.h>
```

```
__global__ void runMeOnGPU() {  
    printf("I'm thread %i!\n", threadIdx.x);  
}
```

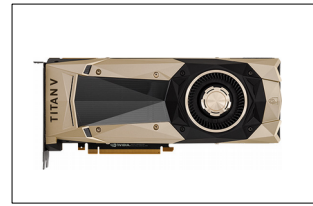
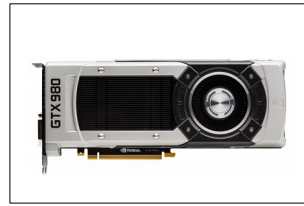
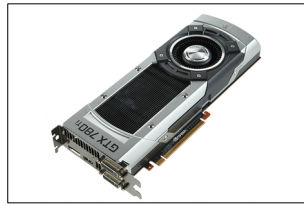
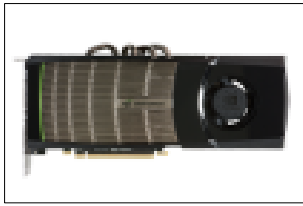
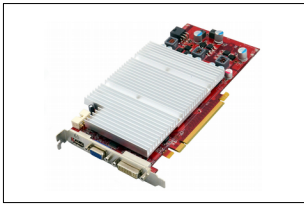
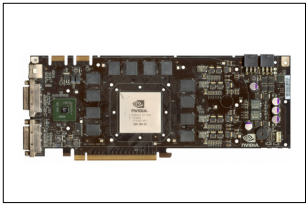
```
int main() {  
    runMeOnGPU<<<1, 32>>>();  
    cudaDeviceSynchronize();  
    return 0;  
}
```

Additional problems a GPU compiler has:

- Same executable may run on different generations of GPU hardware
- Exact hardware is only known at runtime
- Instruction set can vary wildly between generations.

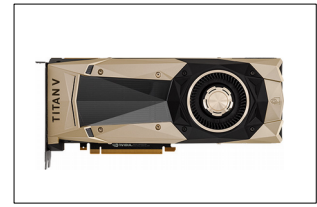
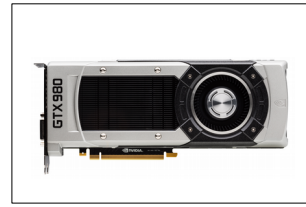
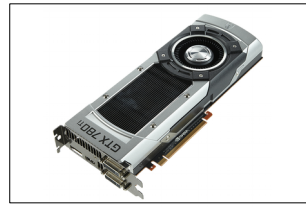
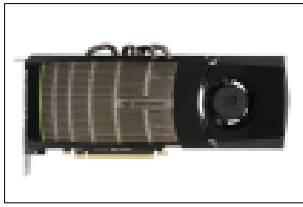
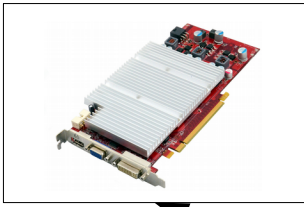
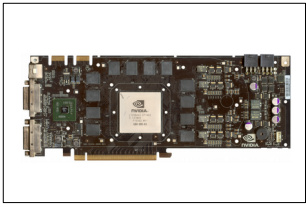


- Same executable may run on different generations of GPU hardware



?

- Exact hardware is only known at runtime



?

- Exact hardware is only known at runtime

Additional problems a GPU compiler has:

- Same executable may run on different generations of GPU hardware
- Exact hardware is only known at runtime
- **Instruction set can vary wildly between generations.**

```
#include <stdio.h>
```

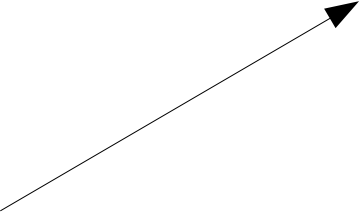
```
__global__ void runMeOnGPU() {  
    printf("I'm thread %i!\n", threadIdx.x);  
}
```

} GPU Kernel

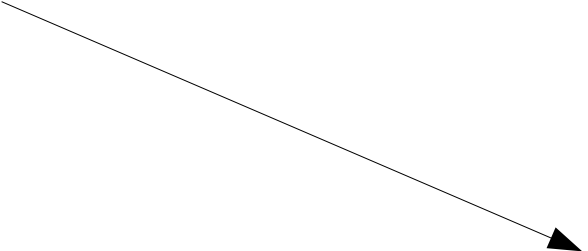
```
int main() {  
    runMeOnGPU<<<1, 32>>>();  
    cudaDeviceSynchronize();  
    return 0;  
}
```

} CPU Code

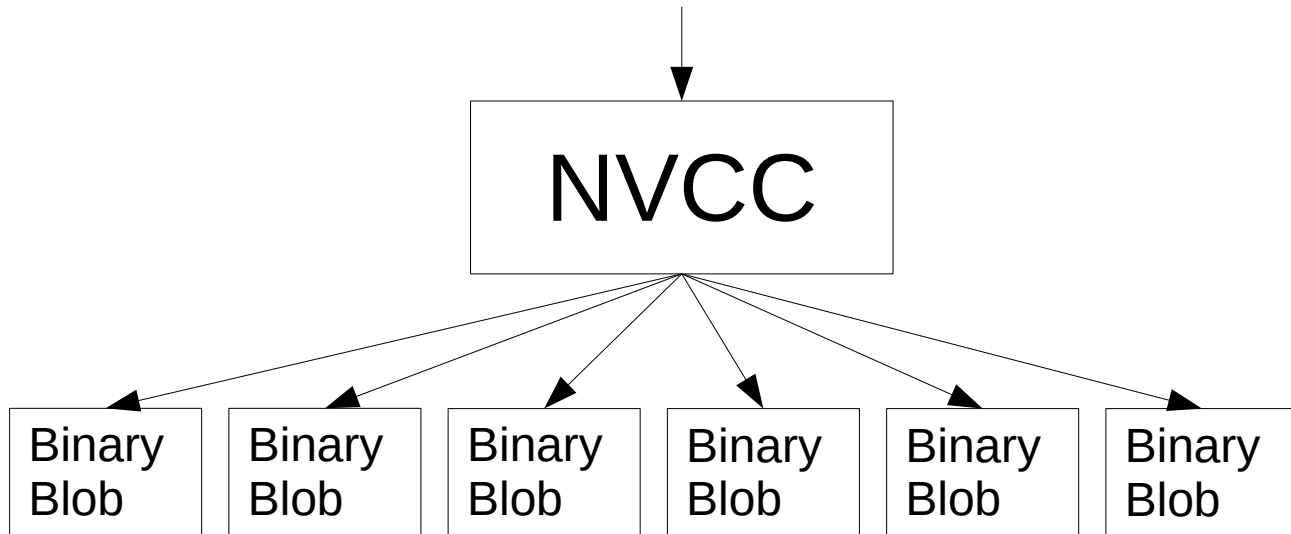

```
__global__ void runMeOnGPU() {  
    printf("I'm thread %i!\n", threadIdx.x);  
}
```

A black arrow originates from the left side of the image and points towards the closing curly brace of the runMeOnGPU function, indicating a call to this function.

```
int main() {  
    runMeOnGPU<<<1, 32>>>();  
    cudaDeviceSynchronize();  
    return 0;  
}
```

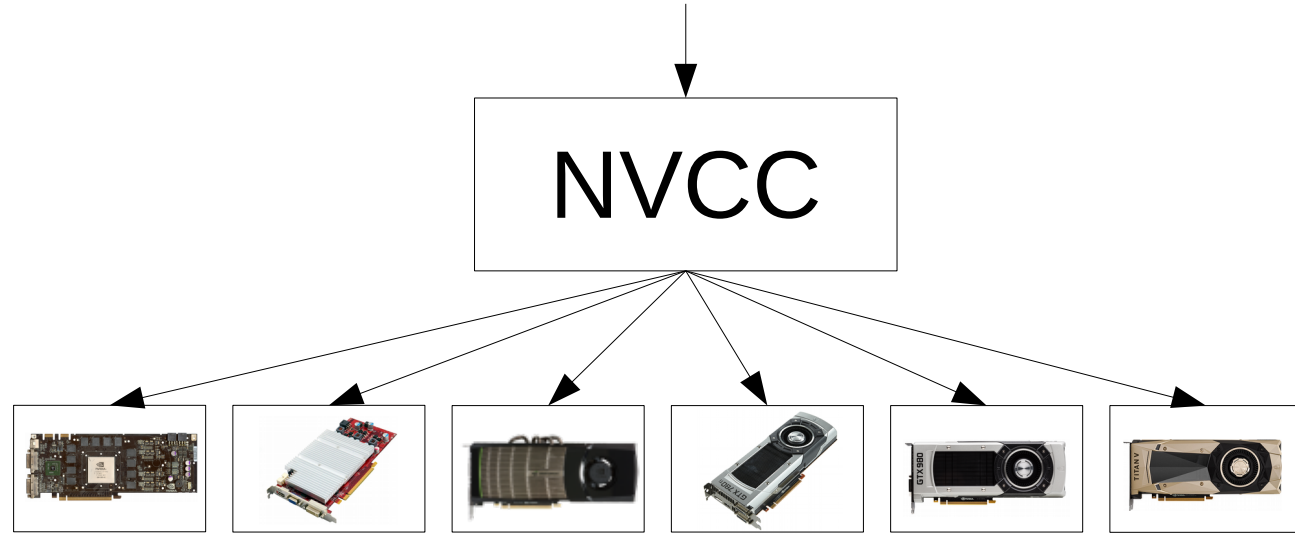
A black arrow originates from the left side of the image and points towards the opening curly brace of the main function, indicating a call to this function.

```
__global__ void runMeOnGPU() {  
    printf("I'm thread %i!\n", threadIdx.x);  
}
```



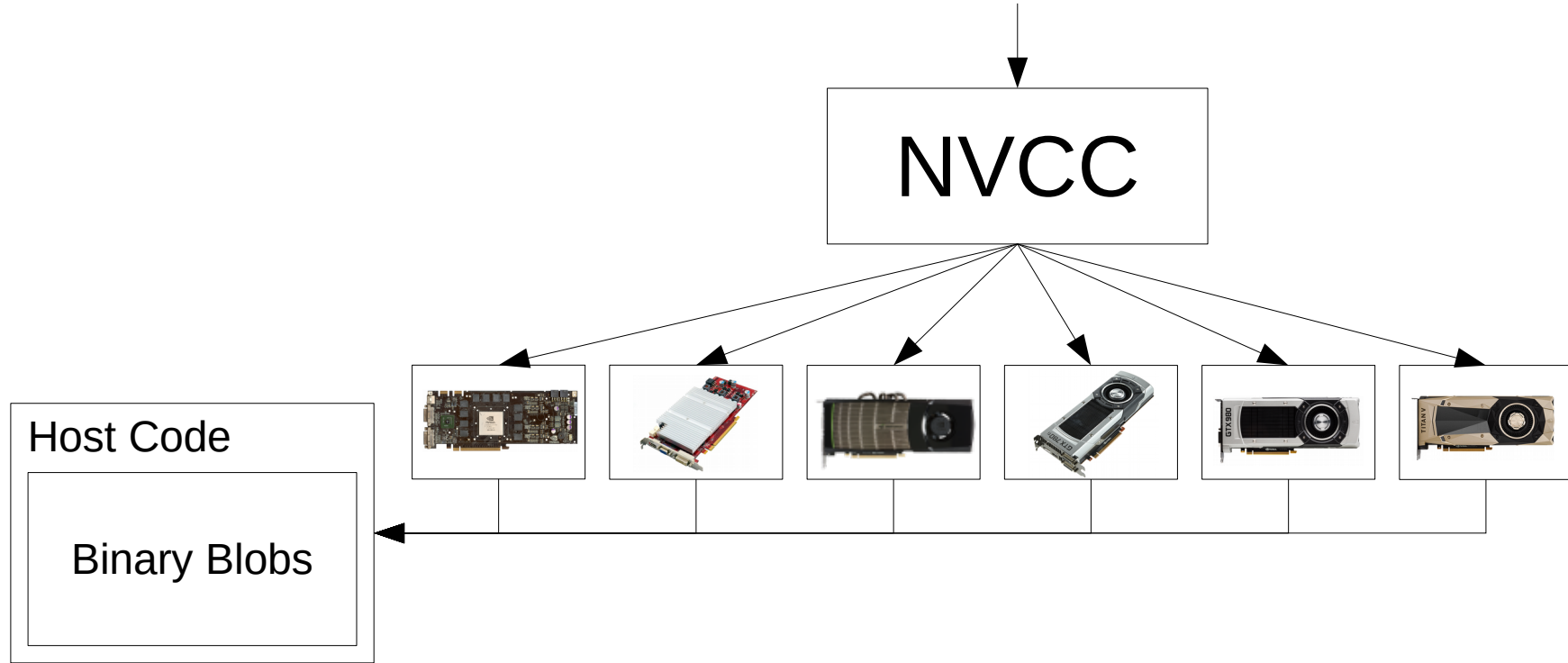
```
int main() {  
    runMeOnGPU<<<1, 32>>>();  
    cudaDeviceSynchronize();  
    return 0;  
}
```

```
__global__ void runMeOnGPU() {  
    printf("I'm thread %i!\n", threadIdx.x);  
}
```



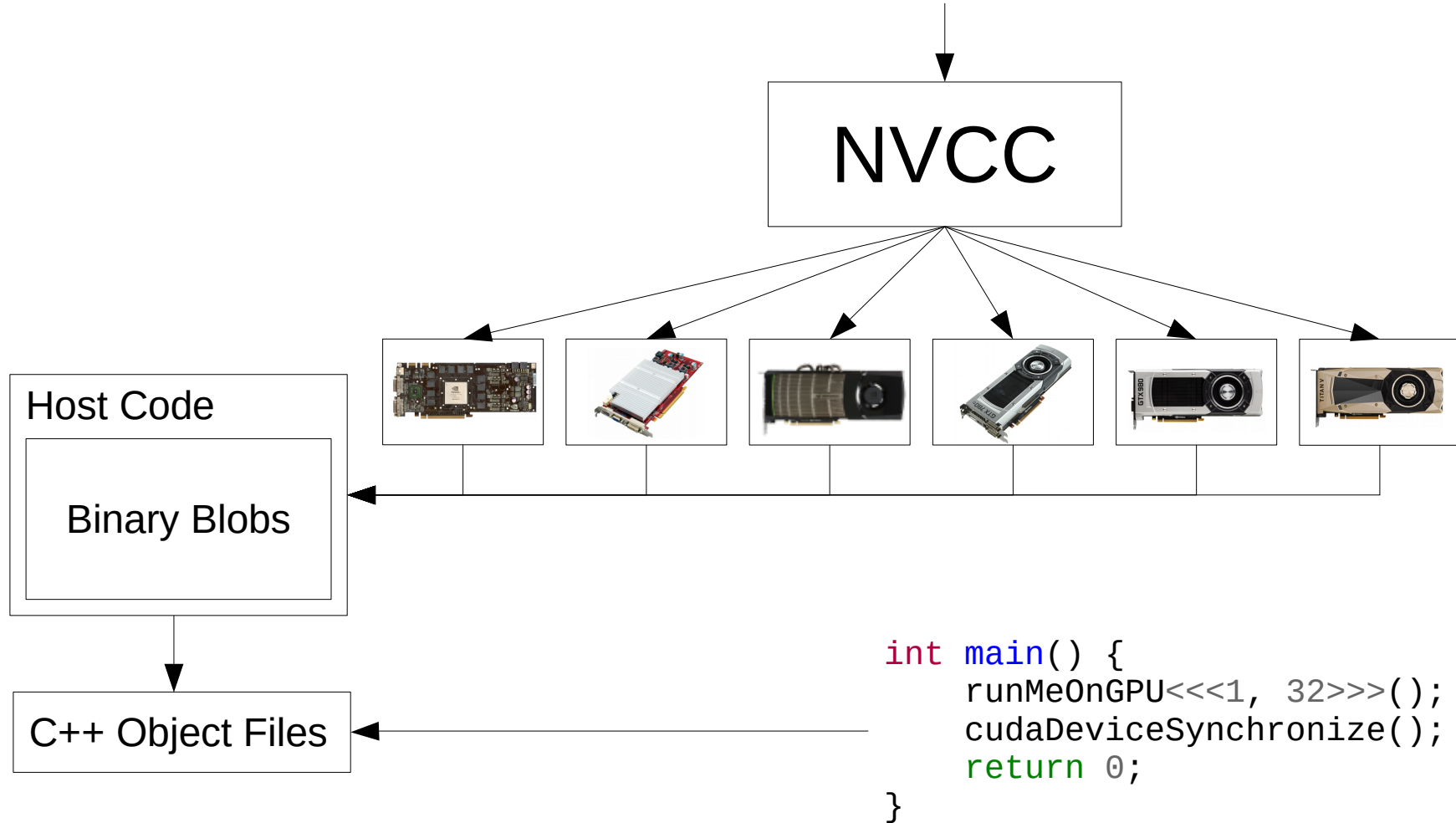
```
int main() {  
    runMeOnGPU<<<1, 32>>>();  
    cudaDeviceSynchronize();  
    return 0;  
}
```

```
__global__ void runMeOnGPU() {  
    printf("I'm thread %i!\n", threadIdx.x);  
}
```

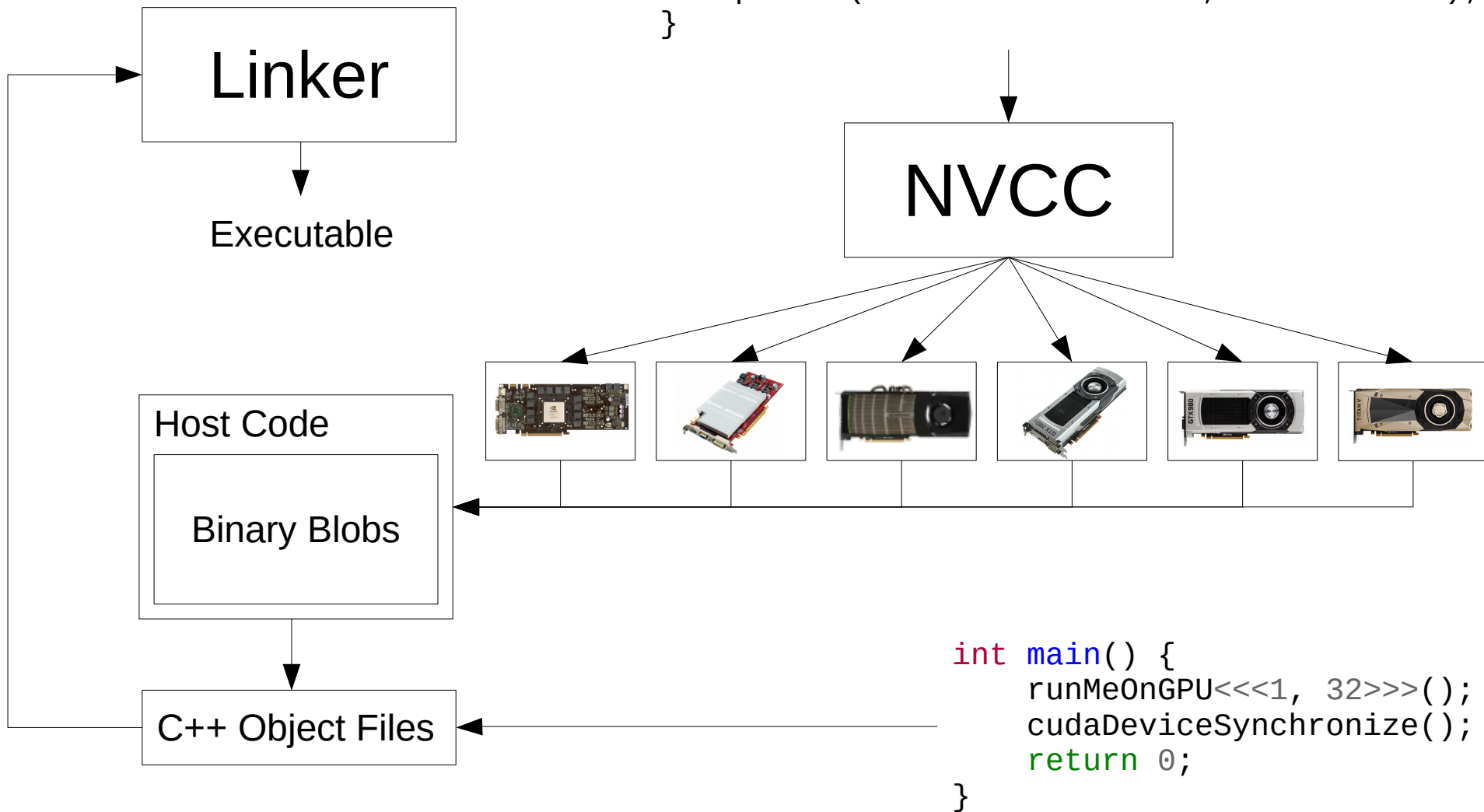


```
int main() {  
    runMeOnGPU<<<1, 32>>>();  
    cudaDeviceSynchronize();  
    return 0;  
}
```

```
__global__ void runMeOnGPU() {  
    printf("I'm thread %i!\n", threadIdx.x);  
}
```



```
__global__ void runMeOnGPU() {  
    printf("I'm thread %i!\n", threadIdx.x);  
}
```



GCC, MSVC, etc

Linker

Executable

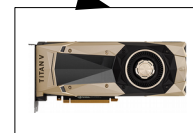
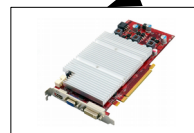
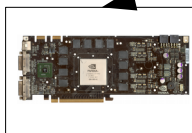
Host Code

Binary Blobs

C++ Object Files

```
__global__ void runMeOnGPU() {  
    printf("I'm thread %i!\n", threadIdx.x);  
}
```

NVCC



```
int main() {  
    runMeOnGPU<<<1, 32>>>();  
    cudaDeviceSynchronize();  
    return 0;  
}
```

Why this matters:

- CUDA is not C++
- NVCC has to support a particular compiler to compile the “host” code
- You may need to specify any architecture(s) you’d like to compile your kernels for.


```
#include <stdio.h>
```

```
__global__ void runMeOnGPU() {  
    printf("I'm thread %i!\n", threadIdx.x);  
}
```

```
int main() {  
    runMeOnGPU<<<1, 32>>>();  
    cudaDeviceSynchronize();  
    return 0;  
}
```

What about larger arrays?

```
#include <stdio.h>
```

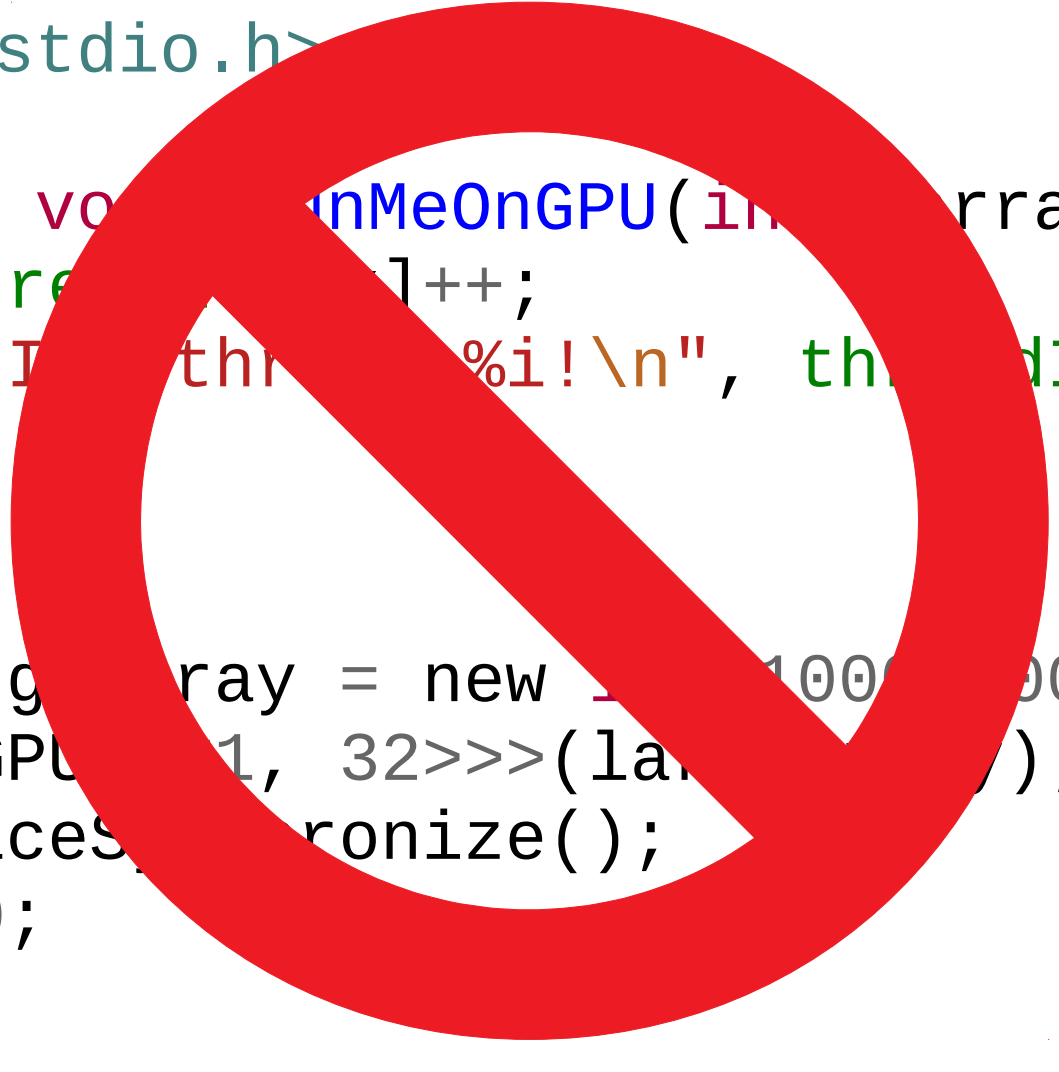
```
__global__ void runMeOnGPU(int* array) {  
    array[threadIdx.x]++;  
    printf("I'm thread %i!\n", threadIdx.x);  
}
```

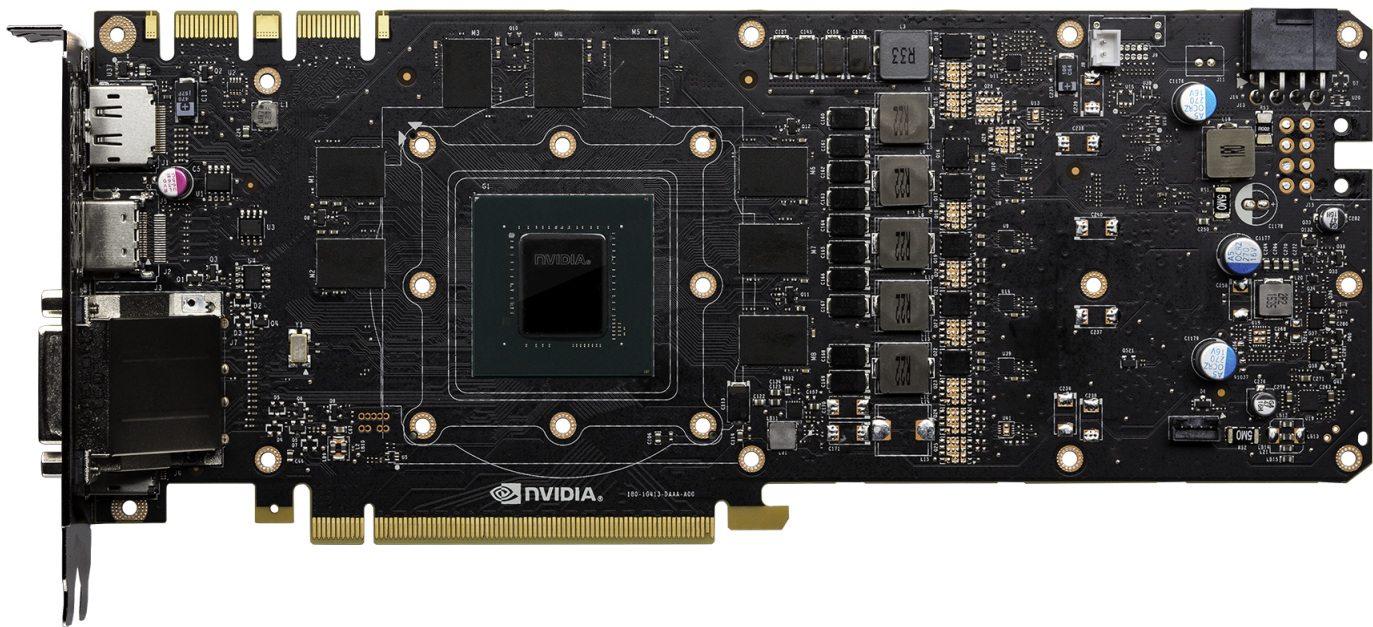
```
int main() {  
    int* largeArray = new int[10000000000];  
    runMeOnGPU<<<1, 32>>>(largeArray);  
    cudaDeviceSynchronize();  
    return 0;  
}
```

```
#include <stdio.h>
```

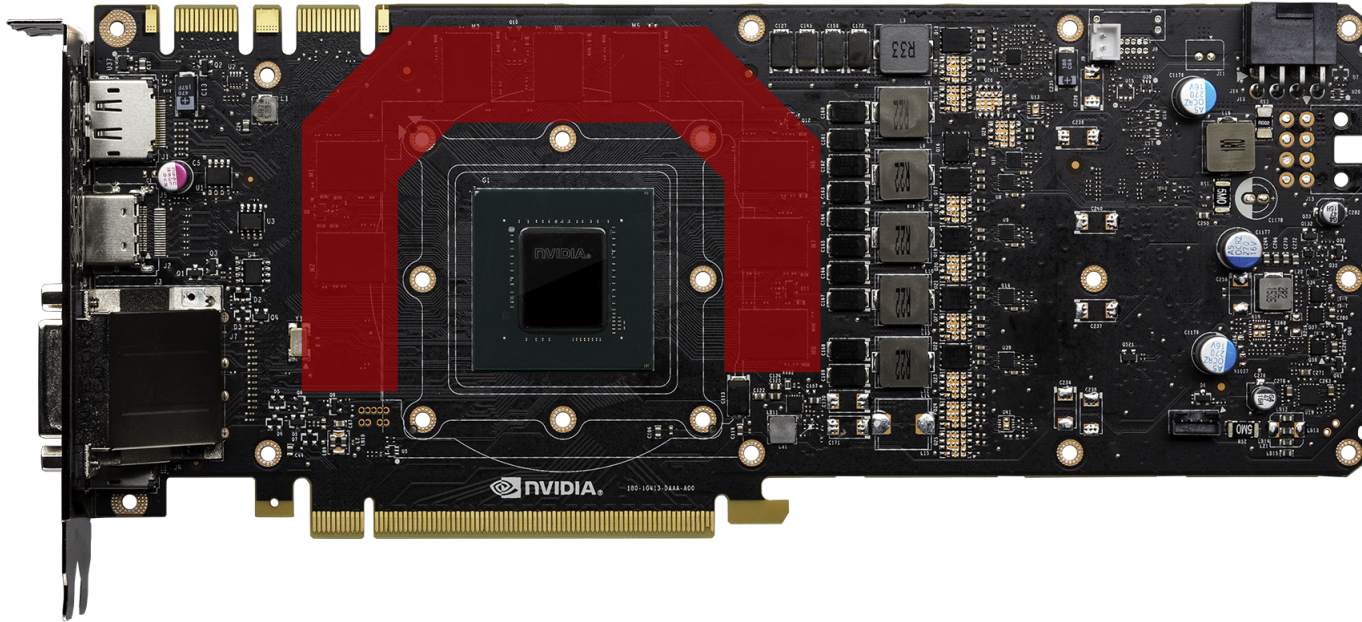
```
__global__ void runMeOnGPU(int array) {  
    array[threadIdx.x]++;  
    printf("Thread %i!\n", threadIdx.x);  
}
```

```
int main()  
{  
    int* largeArray = new int[100000000];  
    runMeOnGPU(1, 32>>>(largeArray));  
    cudaDeviceSynchronize();  
    return 0;  
}
```





The GPU has its own memory banks!



GPU Memory (VRAM)

0

3735913279

Address Range

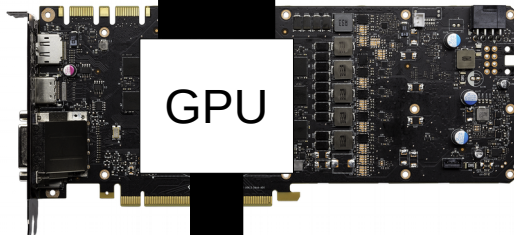
GPU

PCIe bus

0

3126770193

Main Memory (RAM)



GPU Memory (VRAM)

0

3735913279

Address Range

GPU

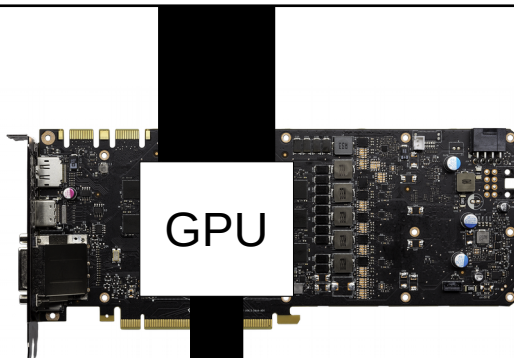
PCIe bus

```
int* largeArray = new int[1000000000];
```

0

3126770193

Main Memory (RAM)

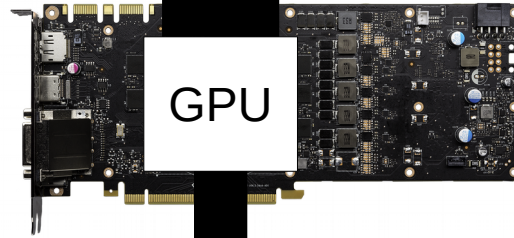


GPU Memory (VRAM)

0

Address Range

3735913279



```
__global__ void runMeOnGPU(int* array) {  
    array[threadIdx.x]++;  
    printf("I'm thread %i!\n", threadIdx.x);  
}
```

PCIe bus

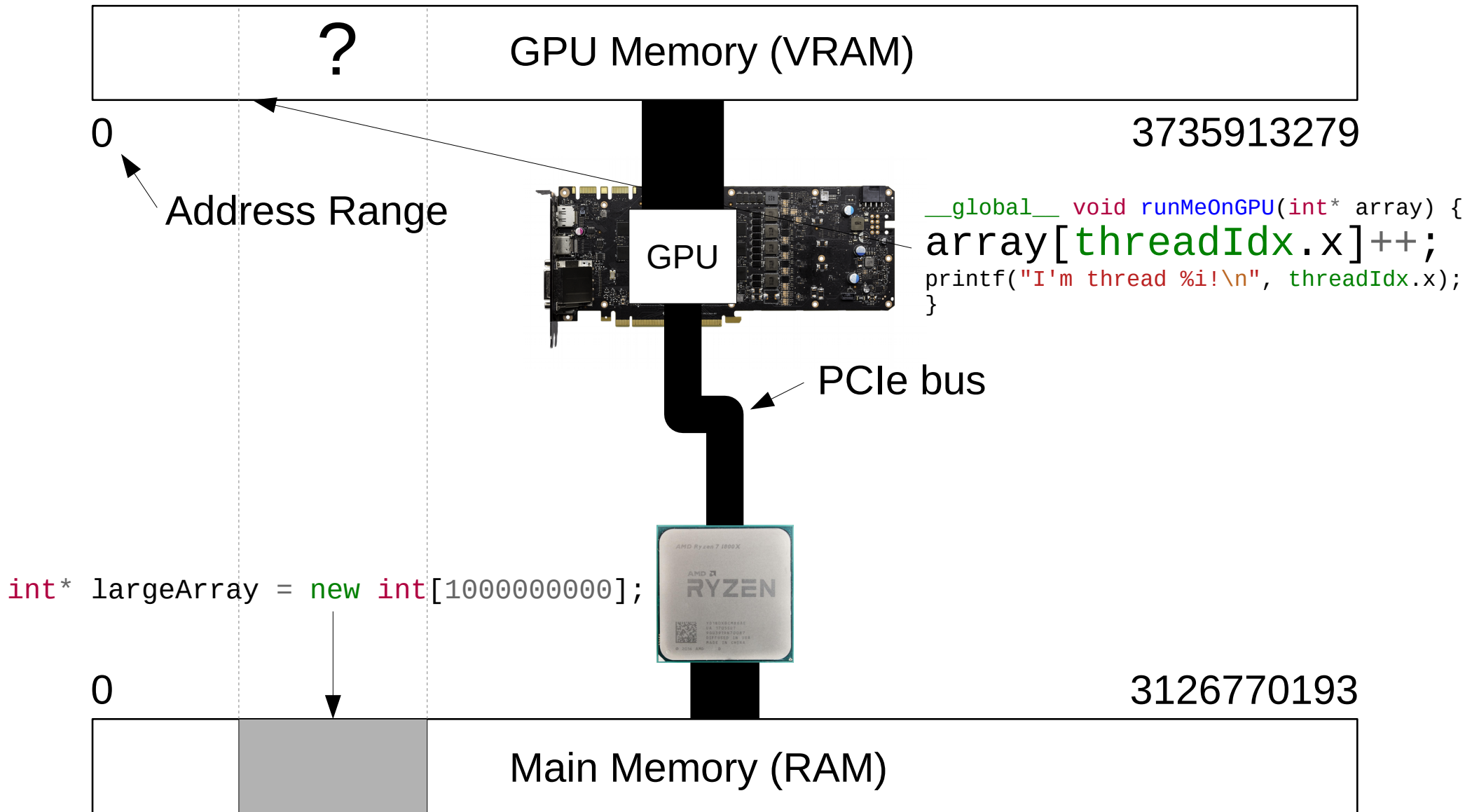
```
int* largeArray = new int[1000000000];
```

0

3126770193



Main Memory (RAM)

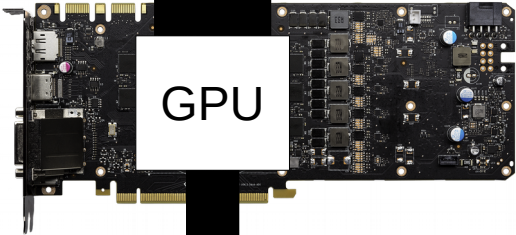


Memory needs to be allocated and
copied to the GPU explicitly

GPU Memory (VRAM)

0

3735913279



0

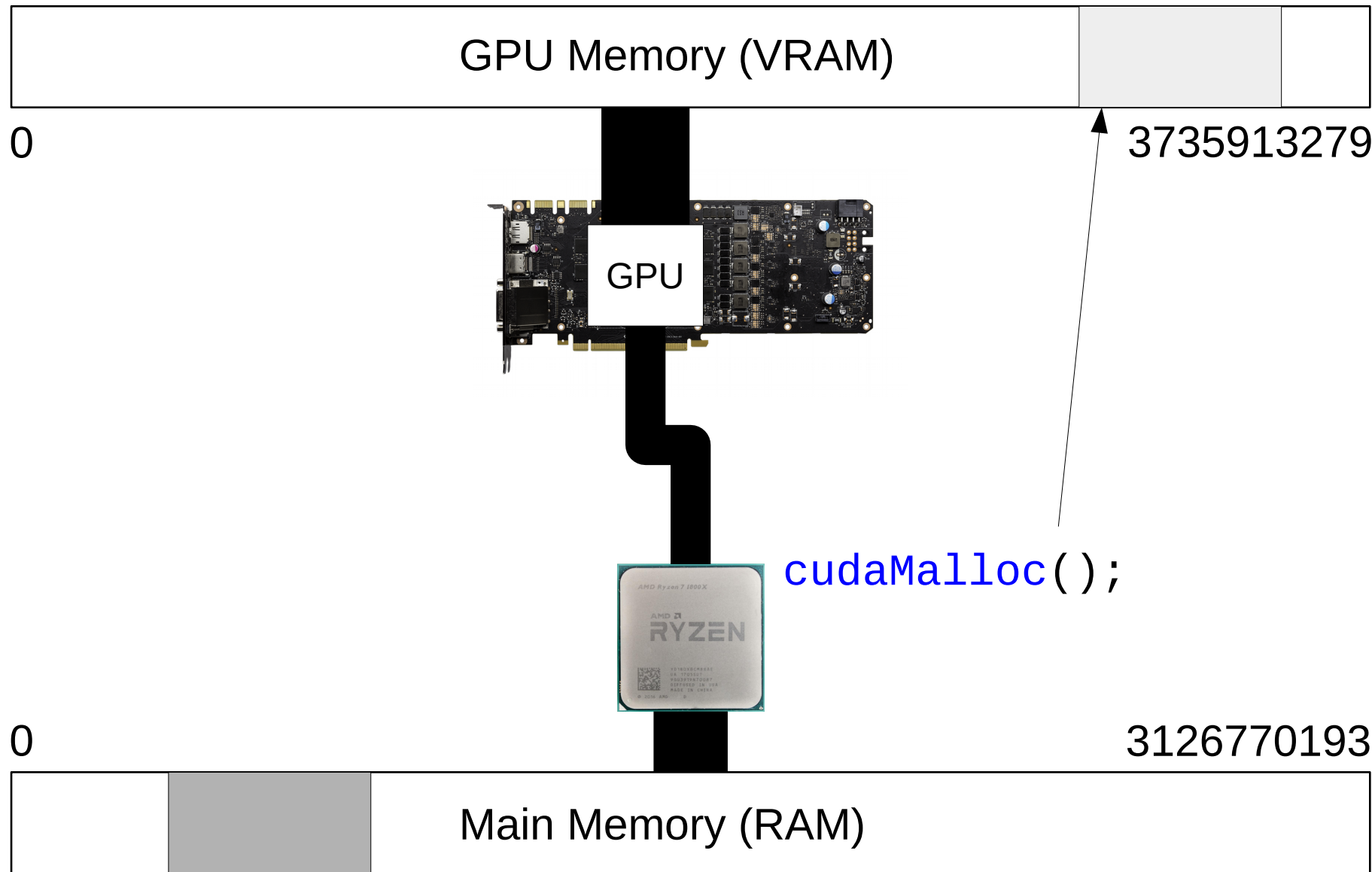
3126770193

Main Memory (RAM)

```
cudaError_t cudaMalloc(void** devPtr, size_t size);
```

```
cudaError_t cudaMalloc(void** devPtr, size_t size);
```

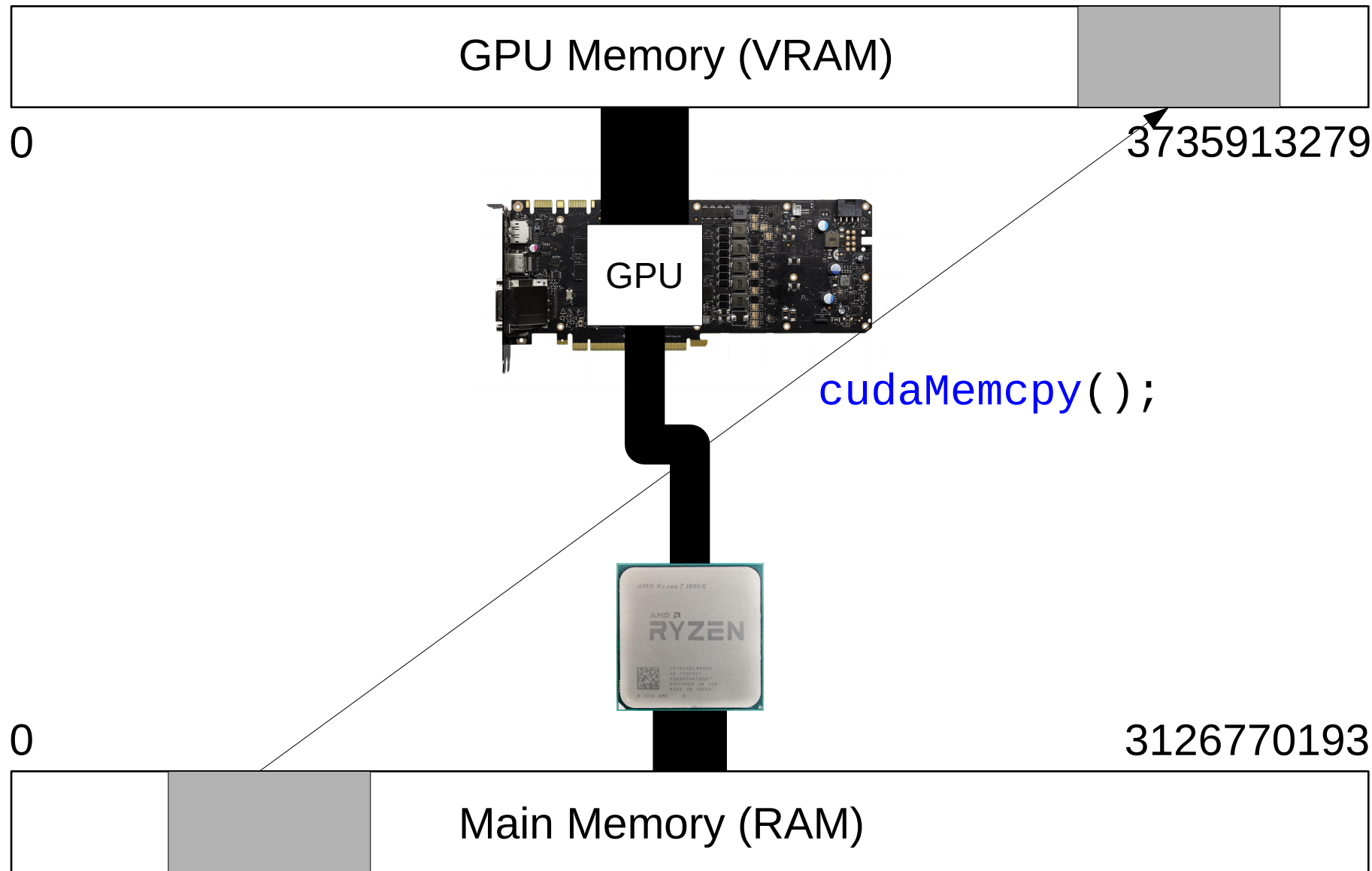
```
float* device_bigArray;  
cudaMalloc(&device_bigArray, sizeof(float) * length);
```




```
cudaError_t cudaMemcpy(void* dest, const void* src,  
                        size_t count, cudaMemcpyKind kind);
```

```
cudaMemcpy(device_bigArray, bigArray,  
            sizeof(float) * length, cudaMemcpyHostToDevice);
```

```
cudaMemcpy(bigArray, device_bigArray,  
            sizeof(float) * length, cudaMemcpyDeviceToHost);
```



Remember that `std::vector` is a two-part structure!

```
cudaError_t cudaFree(void** devicePointer);
```

That's how we copy memory back and forth.

But using memory inside a kernel is slow.

What else can we do?



Shared Memory



Shared Memory

- Blocks can share memory
 - Works as a “programmable cache”
 - Allows communication and cooperation between warps
- ~10x faster than main memory
~10x slower than registers

Last topics:

- What about race conditions?
 - Thread divergence
- Debugging CUDA code
 - Practical notes

What about race conditions?



Bad news:

They still exist

But we have atomic operations!

```
int atomicSub(int* address, int val);
int atomicExch(int* address, int val);
int atomicMin(int* address, int val);
int atomicMax(int* address, int val);
unsigned int atomicInc(..);
unsigned int atomicDec(..);
int atomicCAS(int* address, int compare, int val);
int atomicAnd(int* address, int val);
int atomicOr(int* address, int val);
int atomicXor(int* address, int val);
```

But no mutexes!

(can create one using atomicCAS, but it's prone to deadlocks
on anything but brand new hardware)

Debugging CUDA code

- Can do stepwise debugging, though may need alternate GPU to draw the screen

Debugging CUDA code

- Can do stepwise debugging, though may need alternate GPU to draw the screen
- In my experience:
 - Use assertions
 - Use printf() to find anomalous values
 - Worst case: write debug values into a buffer, dump buffer into excel sheet (lot of work).

Practical notes:

CUDA is not C++

Practical notes:

Windows and the WDDM

Practical notes:

Nvtop: <https://github.com/Syllo/nvtop>

```
bart@MECHANINJA: ~  
File Edit View Search Terminal Help  
Device 0 [GeForce GTX 980 Ti] PCIe GEN 3@16x RX: 0.000 kB/s TX: 0.000 kB/s  
GPU 999MHz MEM 3505MHz TEMP 63M- FAN 32% POW 79 / 250 W  
GPU-Util[ | 5%] MEM-Util[ | 0.9G/6.4G] Encoder[ 0%] Decoder[ 0%]  
  
Device 1 [Quadro P5000] PCIe GEN 1@ 8x RX: 0.000 kB/s TX: 0.000 kB/s  
GPU 139MHz MEM 405MHz TEMP 33M- FAN 26% POW 6 / 180 W  
GPU-Util[ 0%] MEM-Util[ 0.0G/17.1G] Encoder[ 0%] Decoder[ 0%]  
  
PID USER GPU TYPE MEM Command  
4892 bart 0 Graphic 3MB 0.1% /usr/lib/firefox/firefox  
7926 bart 0 Graphic 3MB 0.1% gnome-control-center  
3206 bart 0 Graphic 12MB 0.2% share/jetbrains-toolbox/jetbrains-toolbox  
1516 root 0 Graphic 17MB 0.3% /usr/lib/xorg/Xorg  
6540 bart 0 Graphic 31MB 0.5% /snap/discord/79/usr/share/discord/Discord --  
3416 bart 0 Graphic 42MB 0.7% /usr/lib/slack/slack --type=gpu-process --no-  
2335 gdm 0 Graphic 53MB 0.8% /usr/bin/gnome-shell  
2867 bart 0 Graphic 296MB 4.7% /usr/bin/gnome-shell  
2709 root 0 Graphic 452MB 7.1% /usr/lib/xorg/Xorg  
  
F1Kill F2Sort F3Quit
```

Next week

