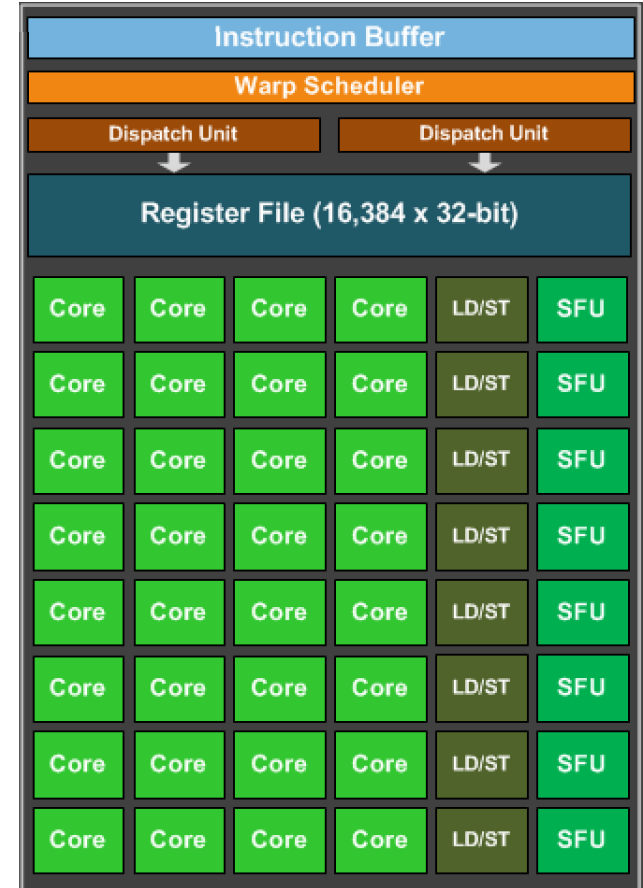


CSE113: Parallel Programming

March 10, 2023

- **Topics:**

- Homework 5
- Javascript webworkers
- GPU programming



Announcements

- HW 2 grades are out, let us know if there are any issues
 - Especially let us know if there are issues with throughput
- Work on Homework 4 (You have until tomorrow to turn it in)
- HW 5 is out
- Last week of class!

Previous Quiz

Previous Quiz

The C++ Parallel and Concurrent schedulers are the same, the only difference is that parallel is optimized to run on multiple cores, while concurrent is meant to timeshare on a single core.

True

False

Previous Quiz

Here are two statements:

(a) The car will never roll down a hill

(b) The car will eventually travel from UCSC to Natural Bridges

These statements are:

-
- Both are safety properties

 - Both are liveness properties

 - a is a liveness property and b is a safety property

 - a is a safety property and b is a liveness property

Previous Quiz

This is the last lecture of lecture 4: please provide any feedback you might have about the module: the material, lectures, slides, homework. Please let me know what you liked and what you didn't like so I can improve the course for future students!

Teaching GPU programming

- This is difficult!
- Nvidia GPUs have the most straightforward programming model (CUDA). They also have great PR.
- It is extremely difficult to get a class of 120 students access to Nvidia GPUs these days.
 - AWS? Expensive and often oversubscribed w.r.t. GPUs
 - Department? ML folks get priority and super computing clusters are painful

Homework 5 - first look

- It is the first time offering this homework, so feedback is very welcome and we will be generous with support.
- Thanks to Mingun Cho who basically did all the work setting up the assignment!



Homework 5 - first look

- Prerequisites
 - Google Chrome Canary
 - (if you have linux, Google Chrome Dev might work)
- Why do we need the Canary?
 - WebGPU is new and support is inconsistent on main (Although it is officially supported)
 - Perhaps more interesting is the shared array buffer.

Homework 5 - first look

- Javascript shared array buffer:
 - How javascript threads can actually share memory
 - Similar to memory in C++

Shared memory and high-resolution timers were effectively **disabled at the start of 2018** [↗](#) in light of **Spectre** [↗](#). In 2020, a new, secure approach has been standardized to re-enable shared memory.

With a few security measures, `postMessage()` will no longer throw for `SharedArrayBuffer` objects and shared memory across threads will be available:

As a baseline requirement, your document needs to be in a **secure context**.

Your application will be in a secure context (you are writing and running locally!)

Homework 5 - first look

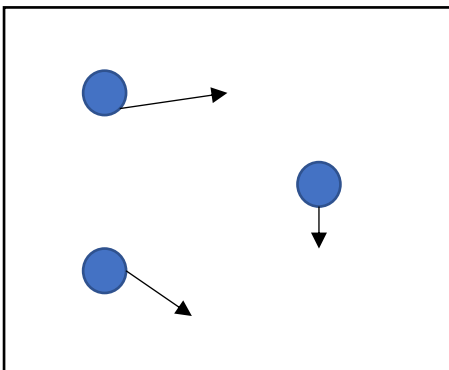
- You will also need Python3 to run a little server

Homework 5 - first look

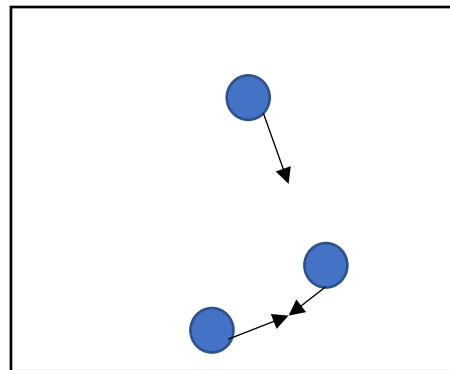
- Let's have a look!

Homework 5 - first look

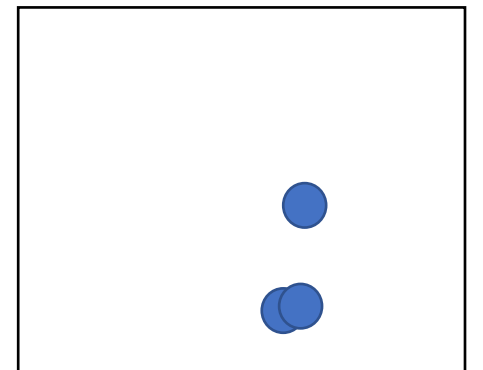
- Your assignment:
 - N-body simulation
- Each particle interacts with every other particle



time = 0



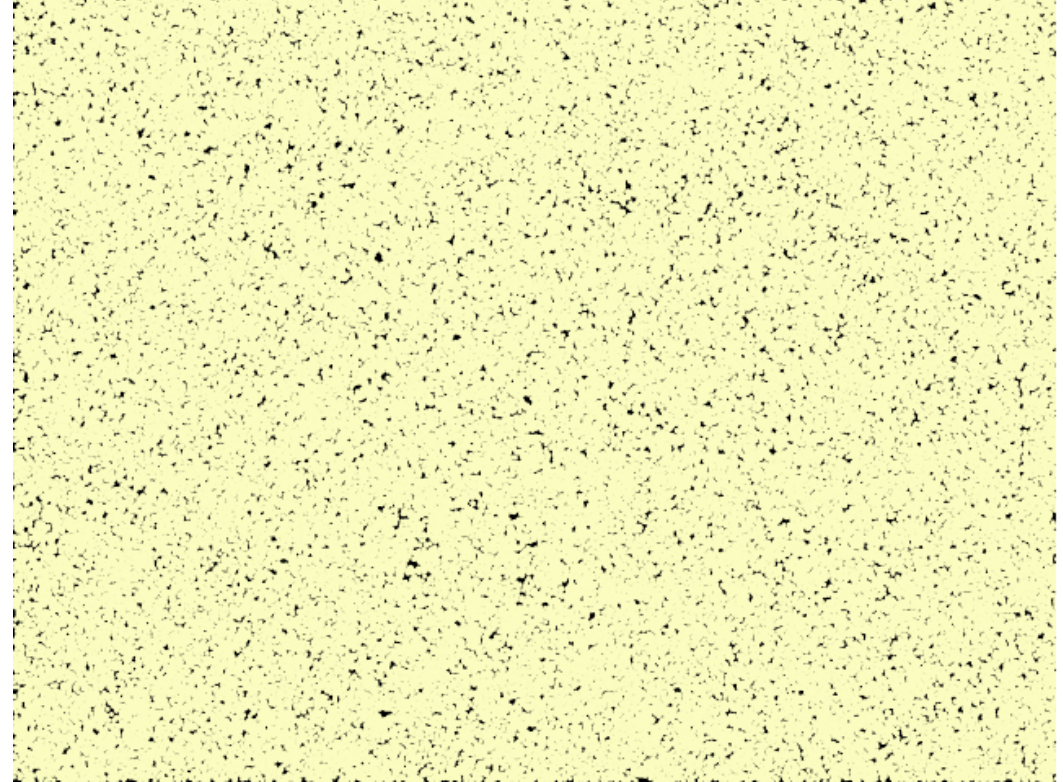
time = 1



time = 2

Examples

- Gravity:
- Boids:
 - <https://en.wikipedia.org/wiki/Boids>



Your homework

- Part 1 of your homework will do this on a single javascript thread
- Demo

Your homework

- Looks good, but with more particles, things start to go slower...

Your homework

- Looks good, but with more particles, things start to go slower...
- Part 2 of the homework is to implement with multiple CPU threads using javascript webworkers
 - Should get around a linear speedup
- Part 3 is to implement with webGPU
 - Should get a BIG speedup!
- You need to explore how many particles you can simulate while keeping a 60 FPS framerate.

Let's look at the code and see some javascript

- look at HTML
- how to print to the console (with interpolation).
 - Syntax errors
- how to interact with HTML
 - overwrite elements
 - modify elements

Shared Array Buffer

- Like Malloc, allocates a “pointer” to a contiguous array of bytes
- Can pass the “pointer” to different threads (webworkers)
- Need to instantiate a typed array to access the values
- Example

Web Workers

- How to do multi-threading in javascript
- Async
 - Concurrent (executes on the same thread)
 - Good for I/O and user interactions
- Web Workers will execute on multiple cores
 - Better for compute intensive applications
 - Better performance

How to use?

- Create a new worker with a file
 - Doesn't do anything yet
- File contains a function: "on message"
- Main file calls "post message" to start the thread along with arguments
- Worker sends a message back to the main file, it can catch the data

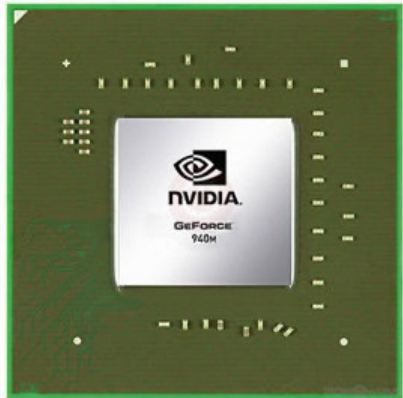
Web Workers

Example with data and arrays

Start on GPU lectures

Programming a GPU

Tiny GPU in an embedded system



Nvidia Jetson Nano (whole chip, CPU + GPU)

2 Billion transistors

10 TDP

Est. \$99

<https://www.techpowerup.com/gpu-specs/geforce-940m.c2648>

https://www.alibaba.com/product-detail/Intel-Core-i7-9700K-8-Cores_62512430487.html

<https://www.prolast.com/prolast-elevated-boxing-rings-22-x-22/>

Fight!



The CPU in my professor workstation



Intel i7-9700K

2.16 Billion transistors

95 TDP

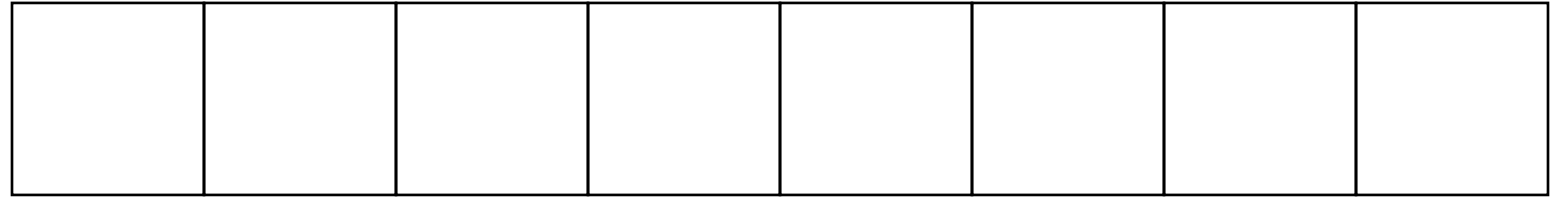
Est. \$316

Programming a GPU

- The problem: Vector addition

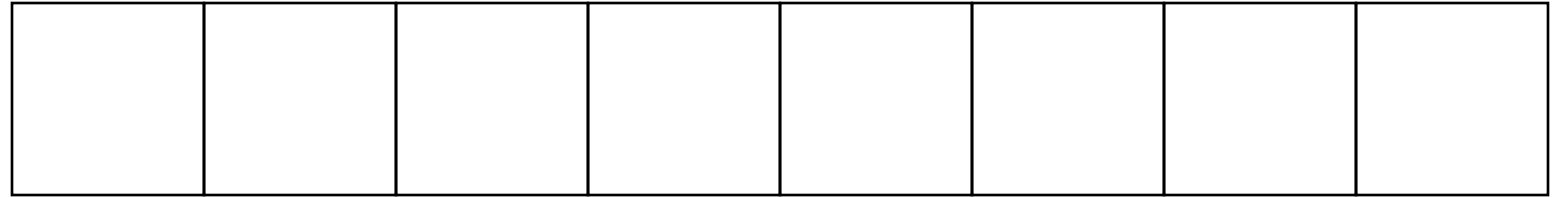
Embarrassingly parallel

array a



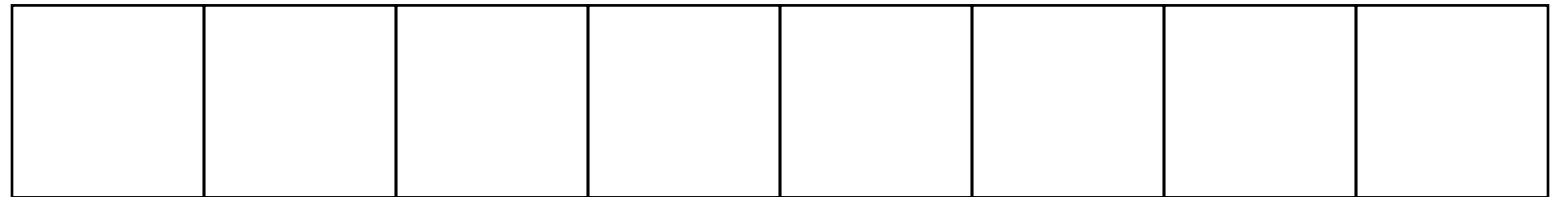
+ + + + + + + +

array b



= = = = = = = =

array c



Computation
can easily be
divided into
threads

- Thread 0 - Blue
- Thread 1 - Yellow
- Thread 2 - Green
- Thread 3 - Orange

Programming a GPU

- The problem: Vector addition
- Who can do it faster?

Lets set up the CPU

- CPU code

Now for the GPU

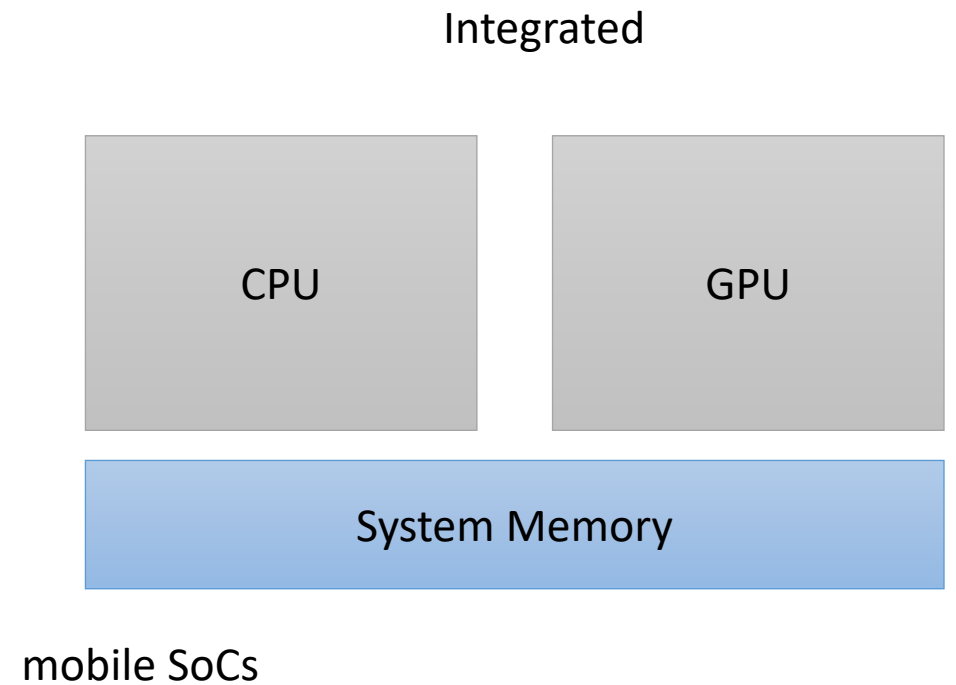
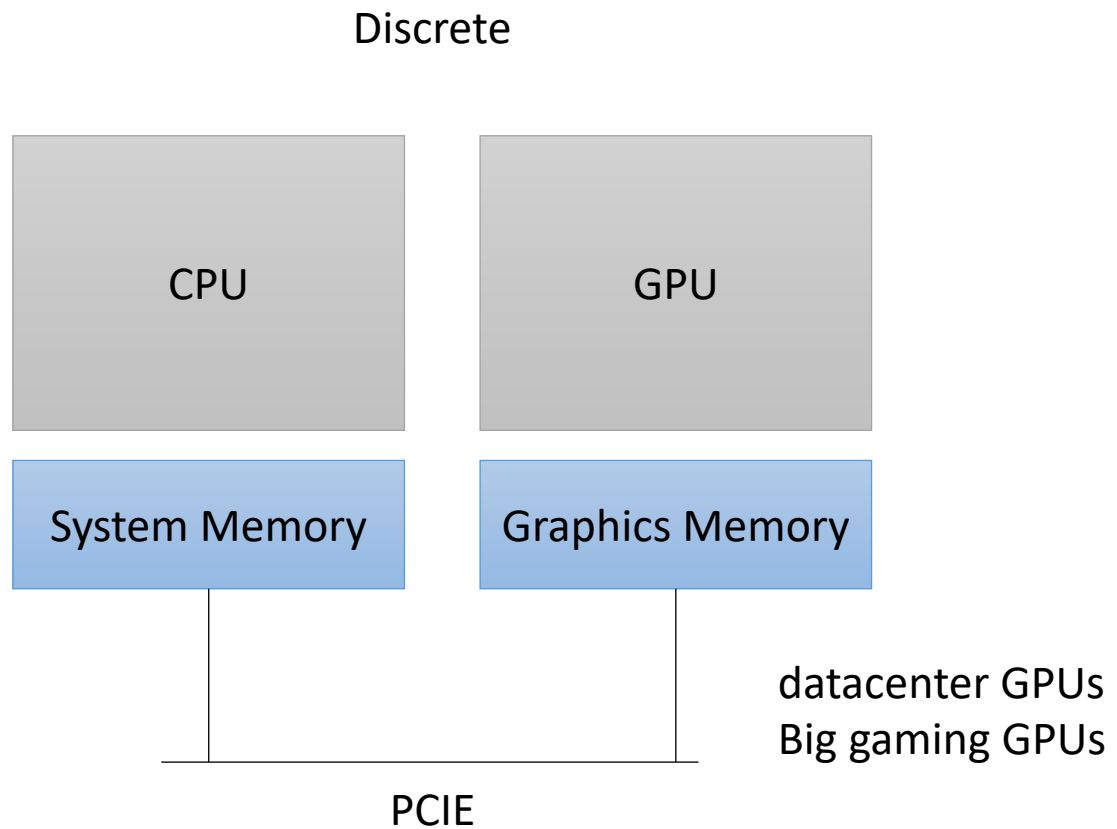
- Its going to take a bit of work....

GPU set up

- We need to allocate and initialize memory

GPU set up

- GPUs come in two flavors

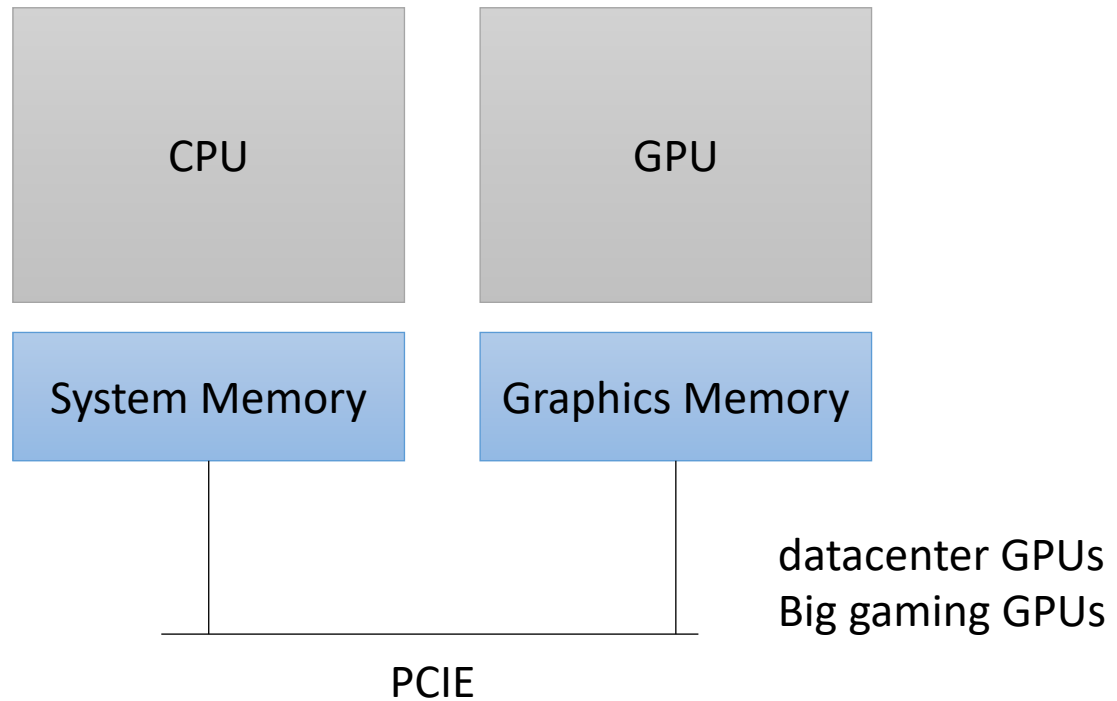


GPU set up

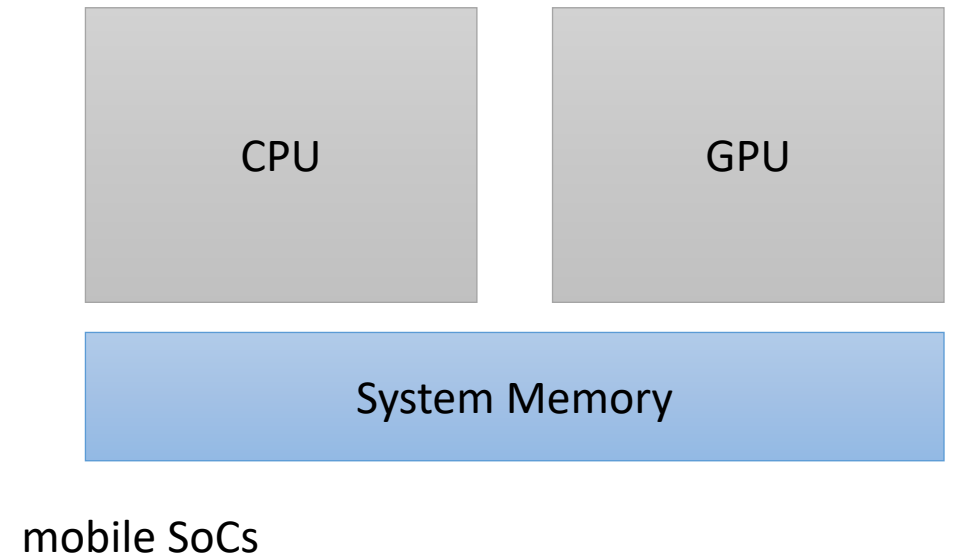
Pros and cons of each?

- GPUs come in two flavors

Discrete



Integrated



GPU set up

- GPUs come in two flavors

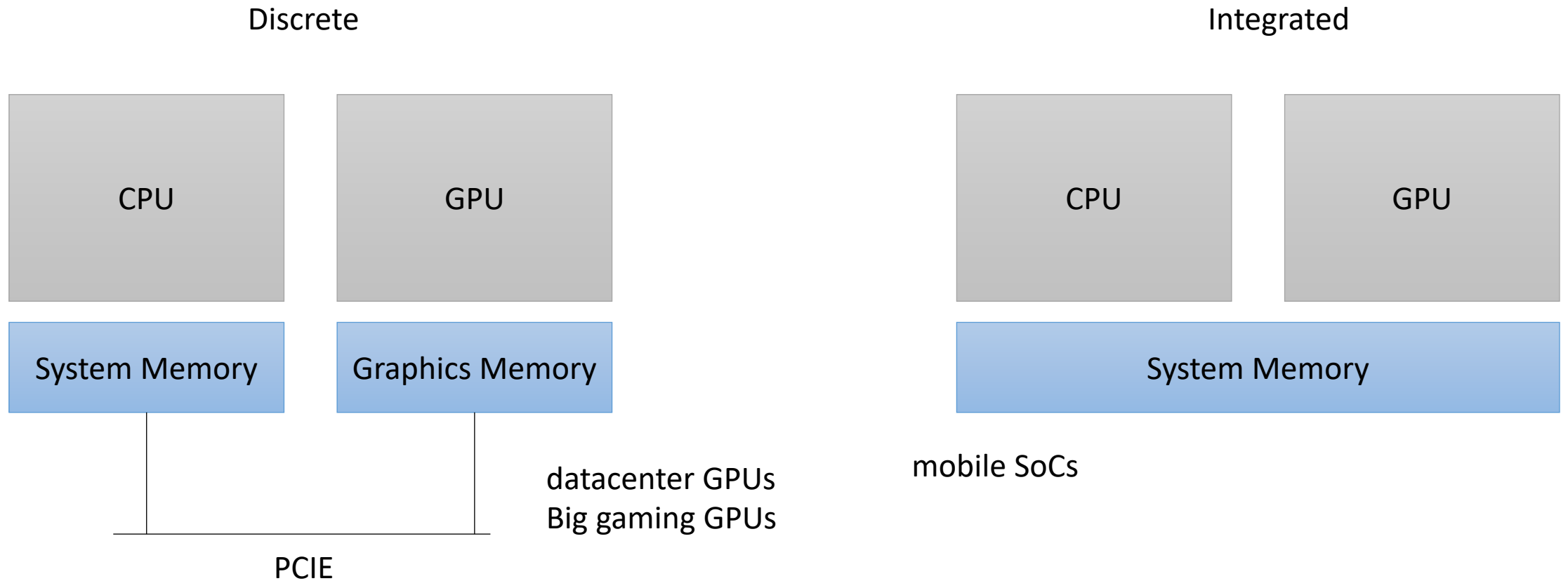
Pros and cons of each?

- *Different types of memory for discrete

- *Swappable for discrete

- *More energy efficient for integrated

- *Better memory utilization for integrated



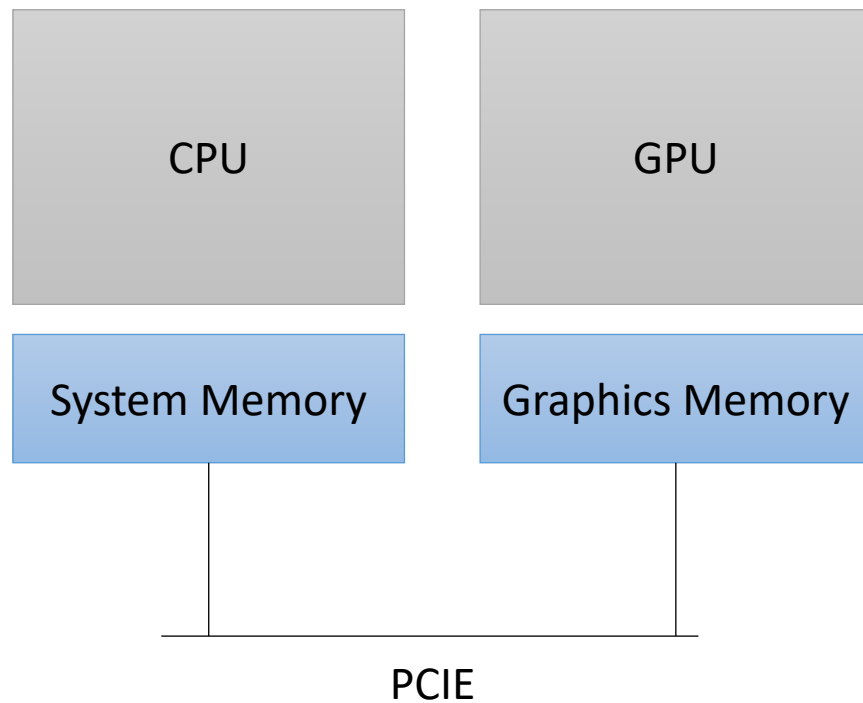
GPU set up

- GPUs come in two flavors

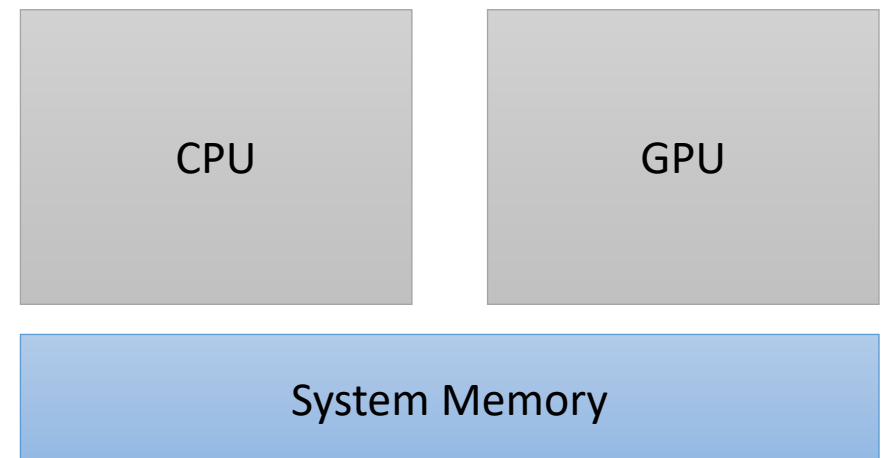
Although mobile GPUs share the system memory, Most still require you to program as if they didn't have shared memory.

Why?

Discrete



Integrated



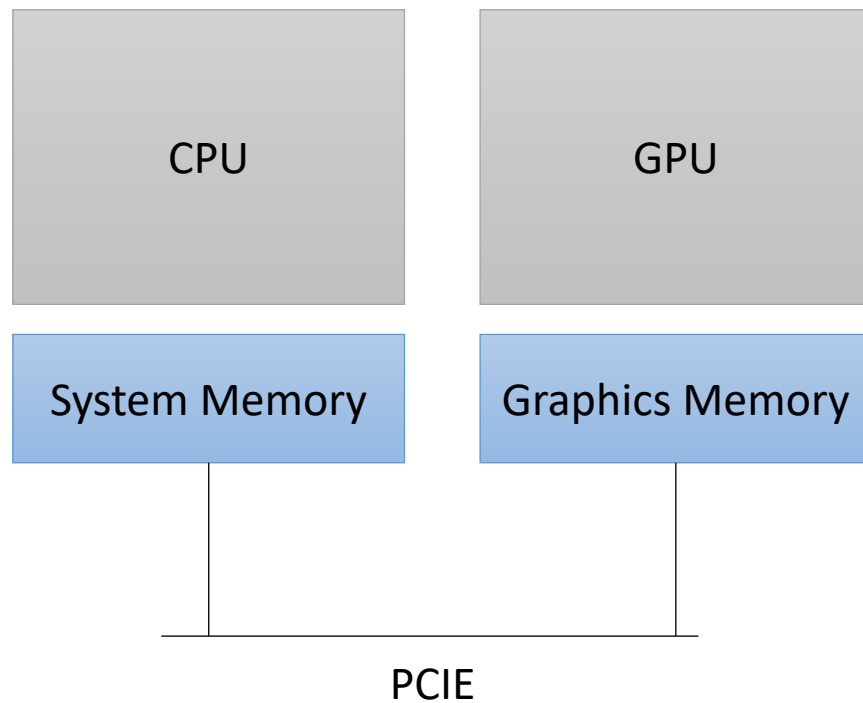
GPU set up

- GPUs come in two flavors

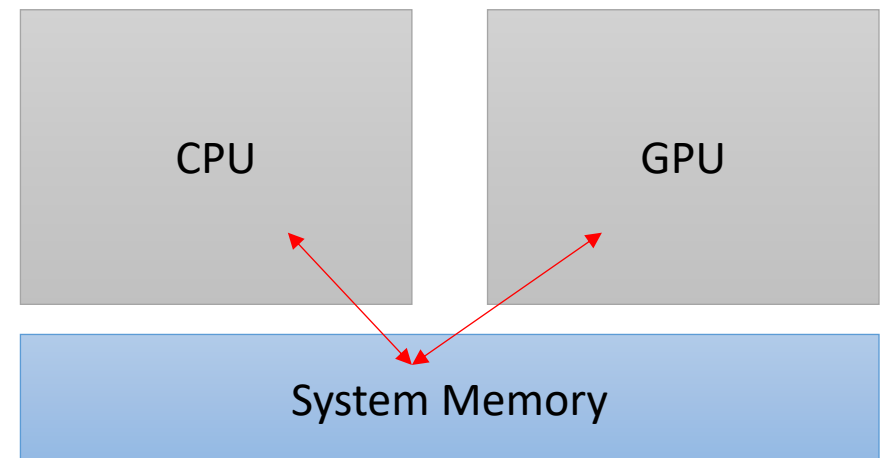
Although mobile GPUs share the system memory, Most still require you to program as if they didn't have shared memory.

Why?

Discrete



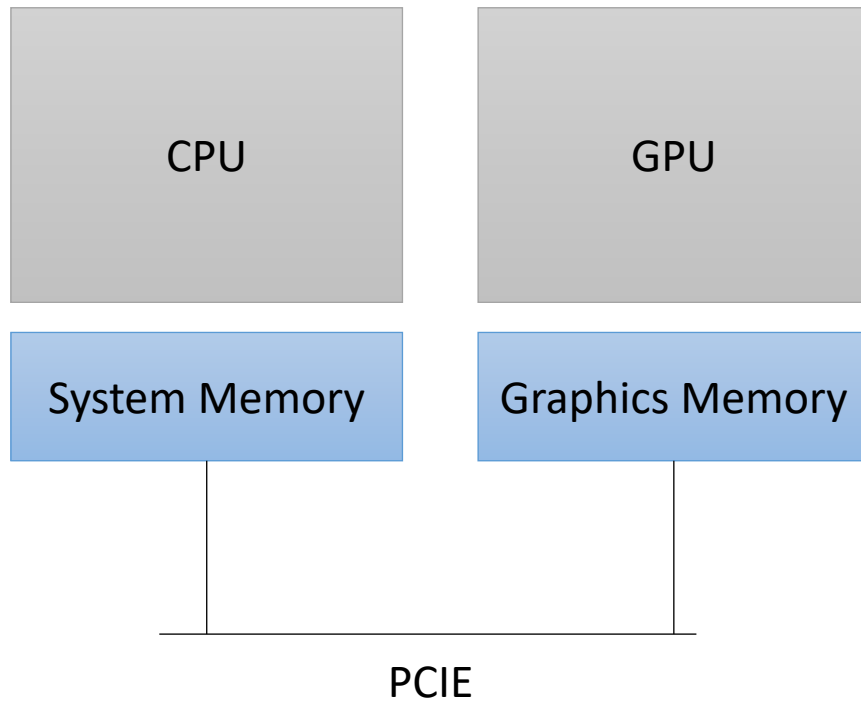
Integrated



GPU set up

- GPUs come in two flavors

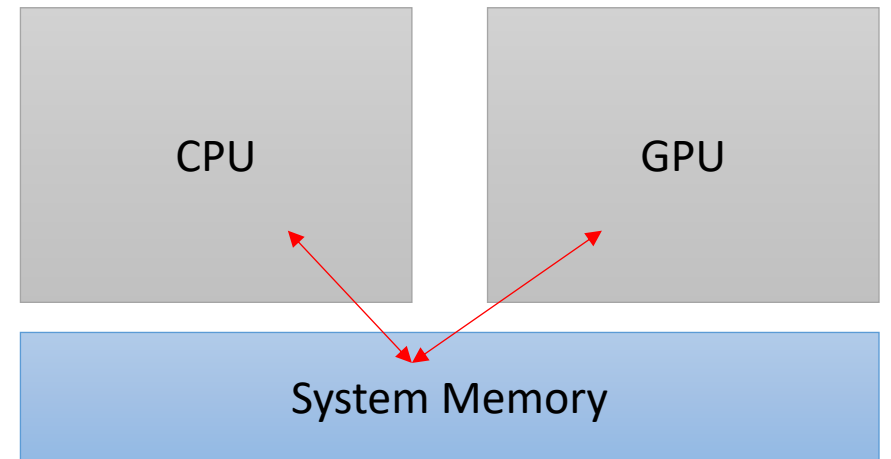
Discrete



Although mobile GPUs share the system memory, Most still require you to program as if they didn't have shared memory.

Why?

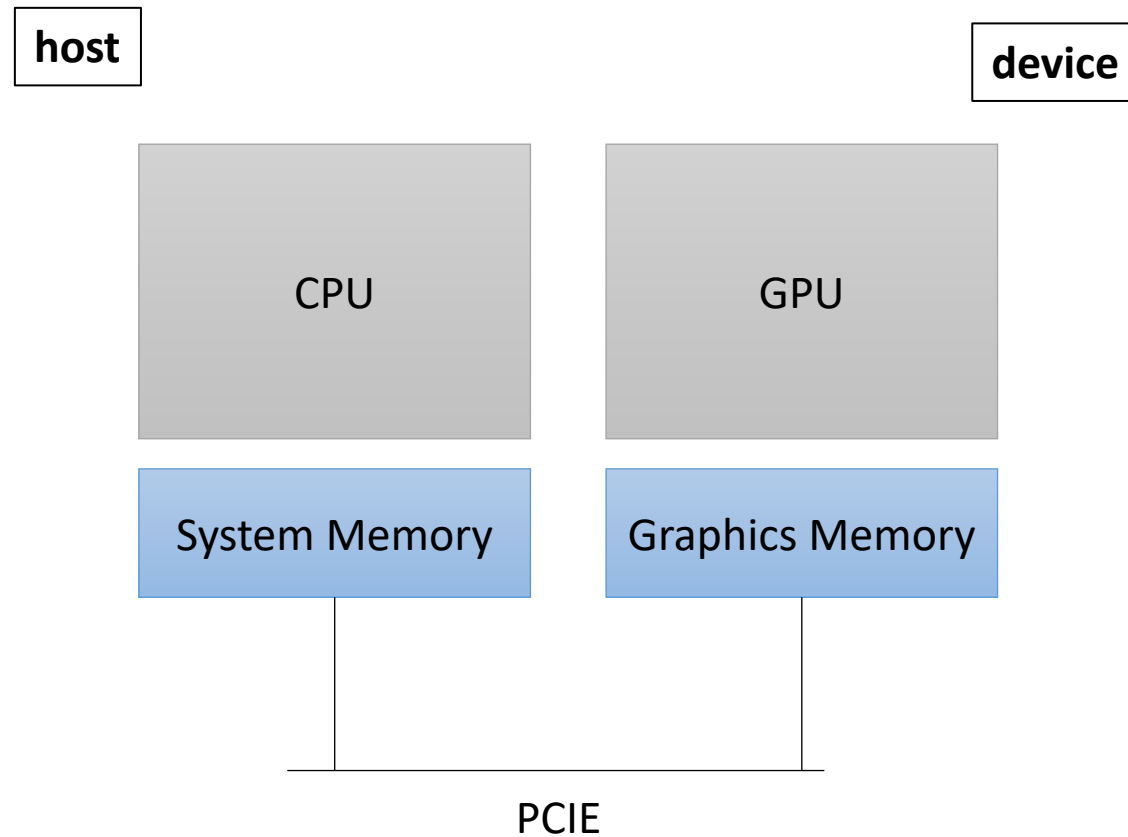
Integrated



In many cases, CPU-GPU communication is not fully supported coherence, fences, and RMWs might now be supported.

GPU set up

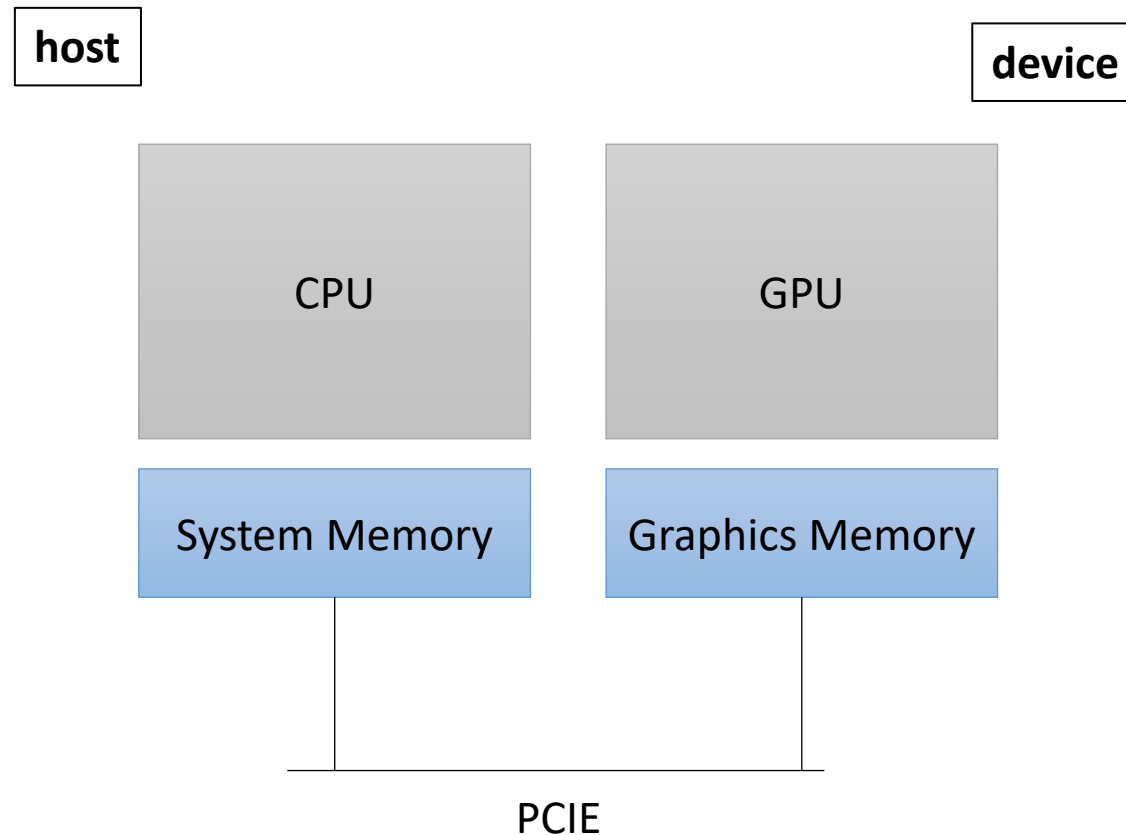
- Our heterogeneous, parallel, programming model



GPU set up

- Our heterogeneous, parallel, programming model

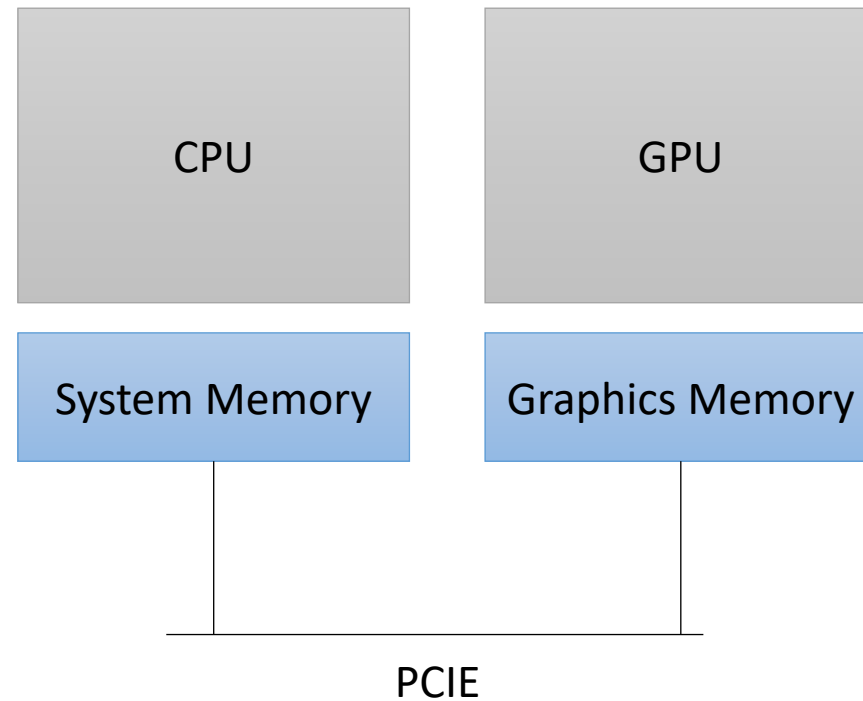
The host (CPU) will write a C++-like program that allocates and sets up memory on the GPU. The host will then call a GPU program called a kernel.



GPU set up

How do we allocate memory on a CPU?

- Our heterogeneous, parallel, programming model

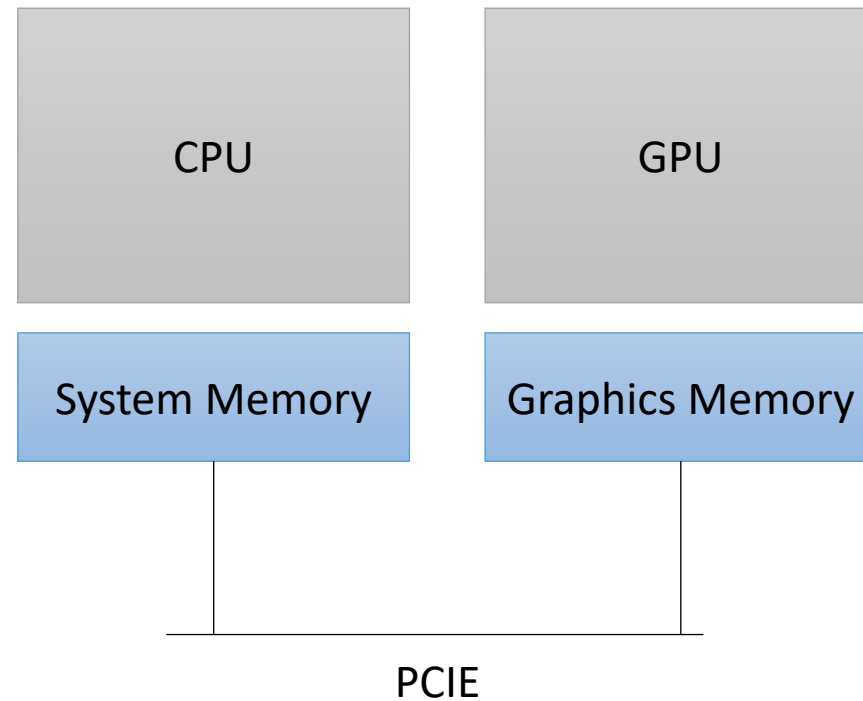


GPU set up

How do we allocate CPU memory on the host?

- Our heterogeneous, parallel, programming model

```
int *x = (int*) malloc(sizeof(int)*SIZE);
```

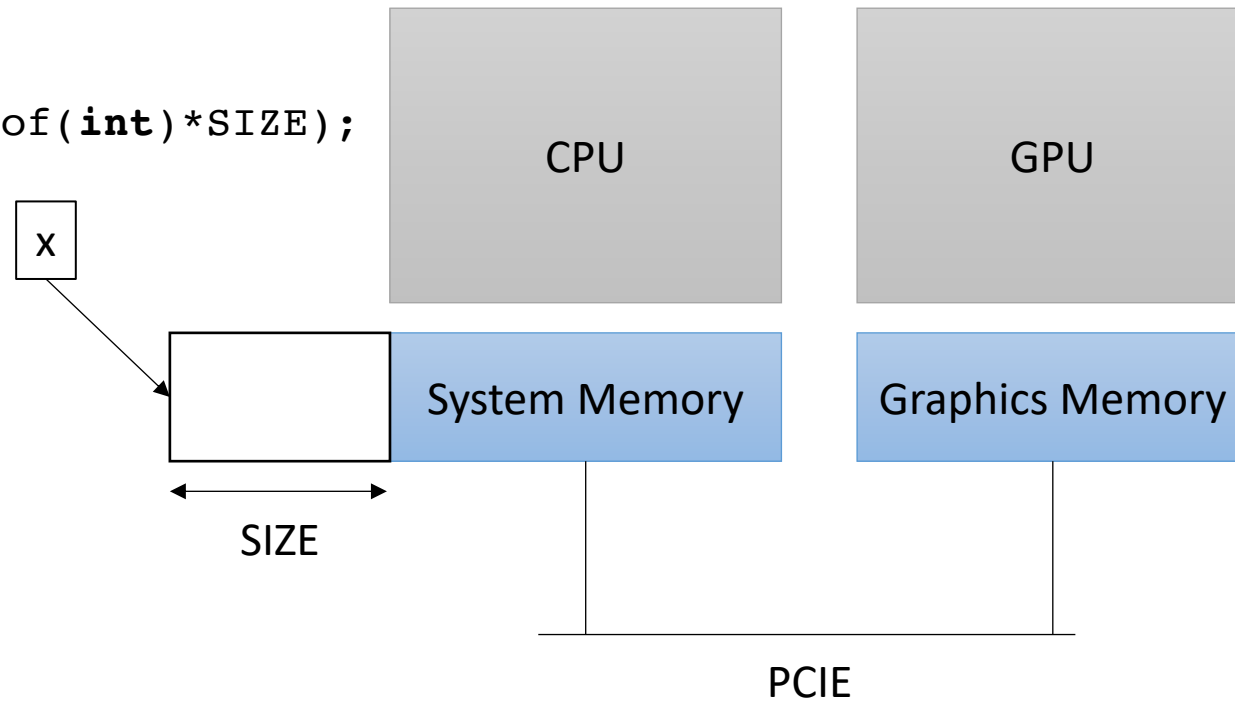


GPU set up

How do we allocate CPU memory on the host?

- Our heterogeneous, parallel, programming model

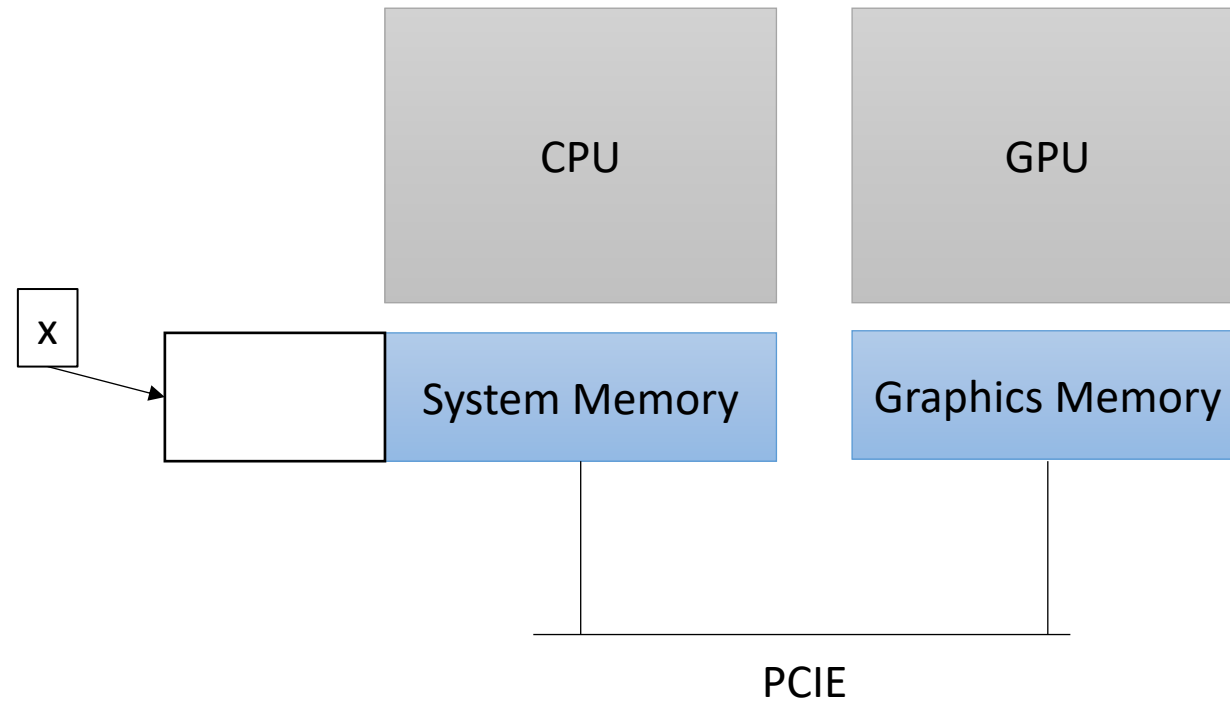
```
int *x = (int*) malloc(sizeof(int)*SIZE);
```



GPU set up

We need to allocate GPU memory on the host

- Our heterogeneous, parallel, programming model

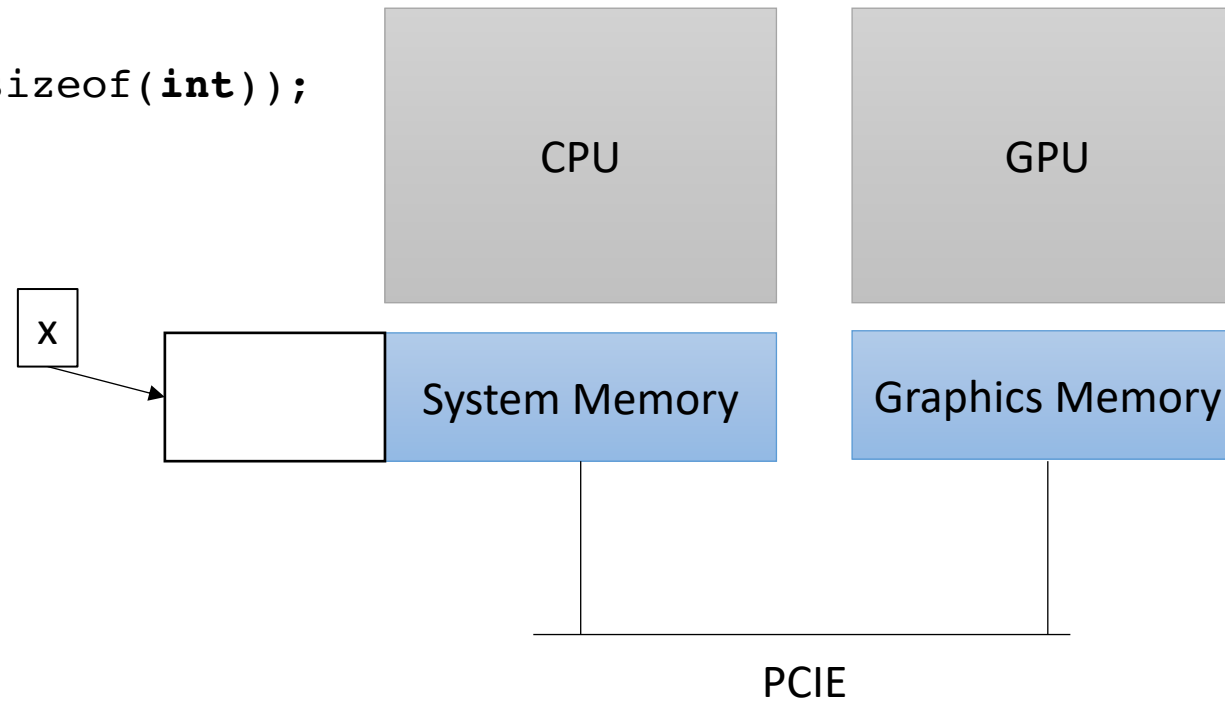


GPU set up

We need to allocate GPU memory on the host

- Our heterogeneous, parallel, programming model

```
int *d_x;  
cudaMalloc(&d_x, SIZE*sizeof(int));
```

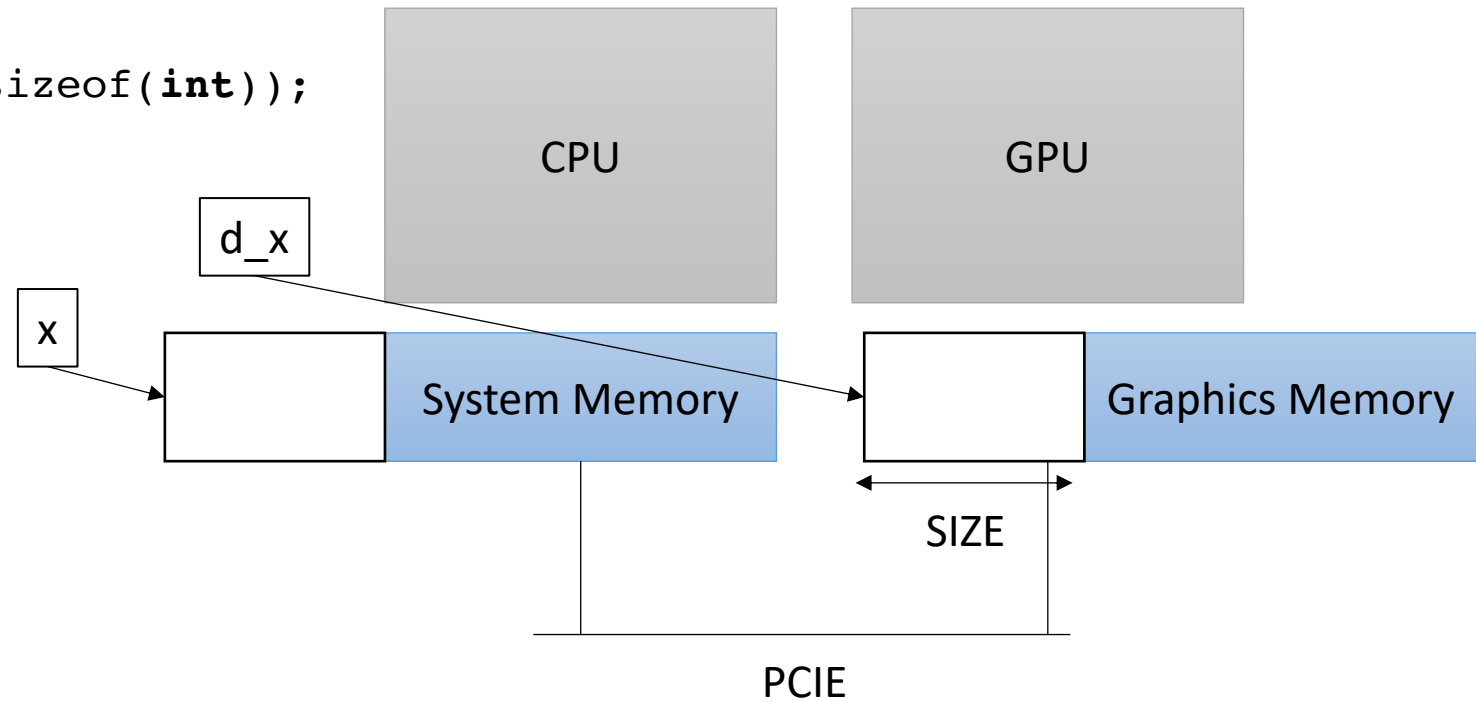


GPU set up

We need to allocate GPU memory on the host

- Our heterogeneous, parallel, programming model

```
int *d_x;  
cudaMalloc(&d_x, SIZE*sizeof(int));
```



GPU set up

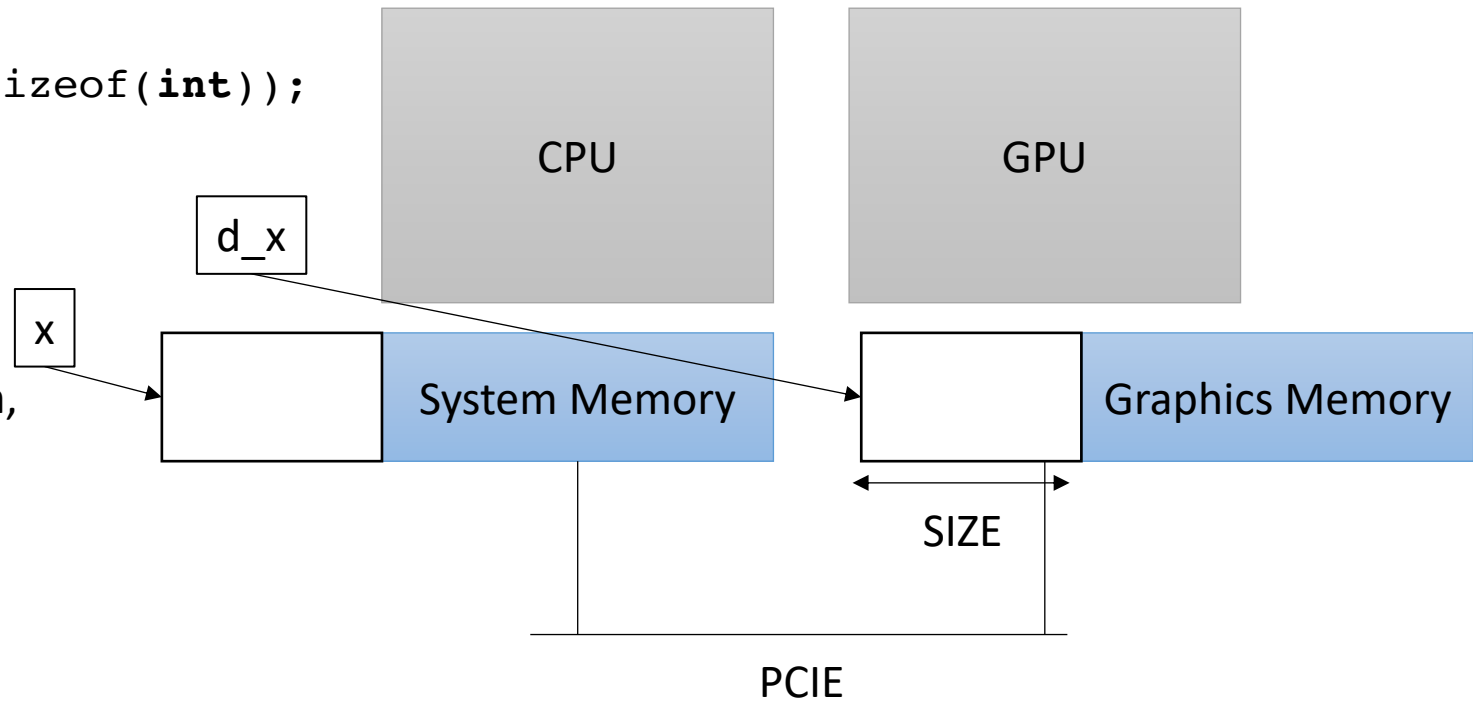
We need to allocate GPU memory on the host

- Our heterogeneous, parallel, programming model

```
int *d_x;  
cudaMalloc(&d_x, SIZE*sizeof(int));
```

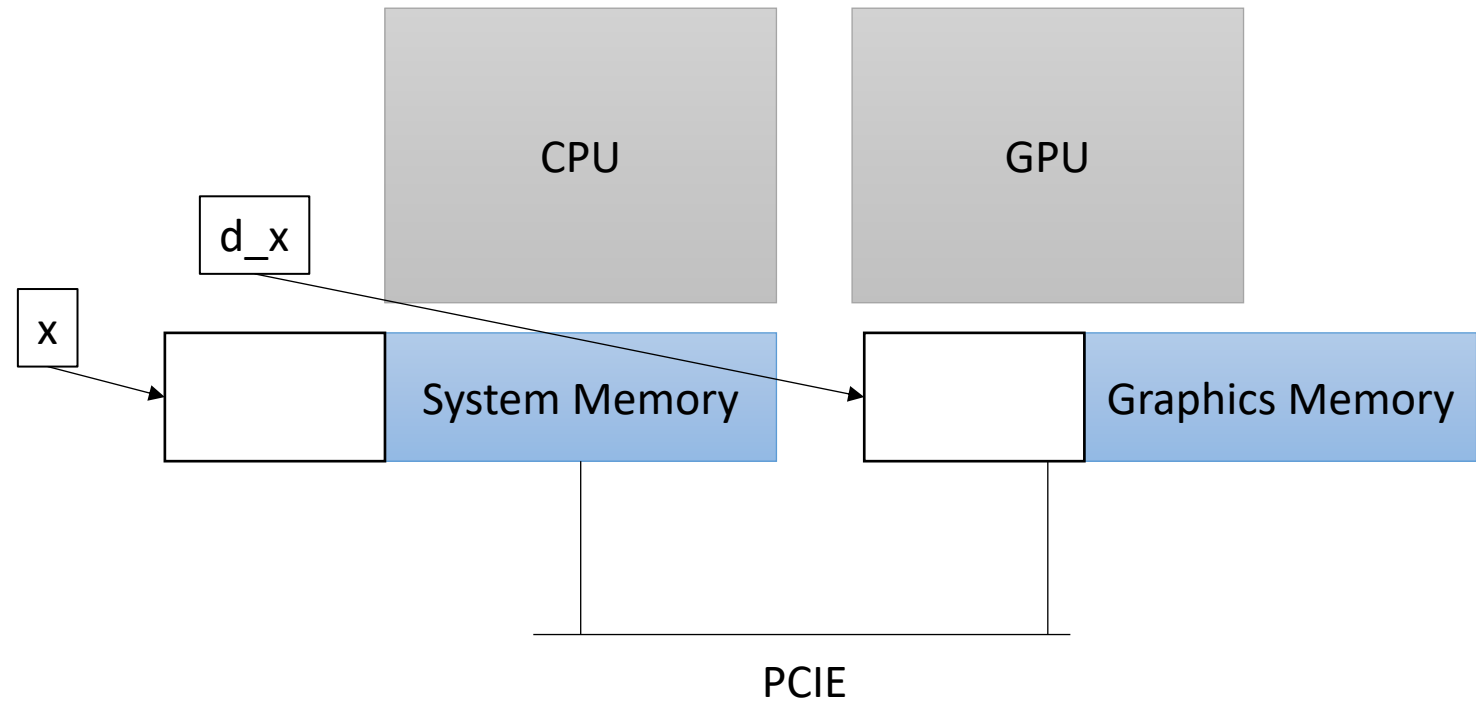
`d_x` is a pointer, in the CPU program, that points to memory on the GPU.

We can pass the pointer around, but the CPU cannot access the data
i.e. `d_x[0]` gives an error!



GPU set up

- Our heterogeneous, parallel, programming model

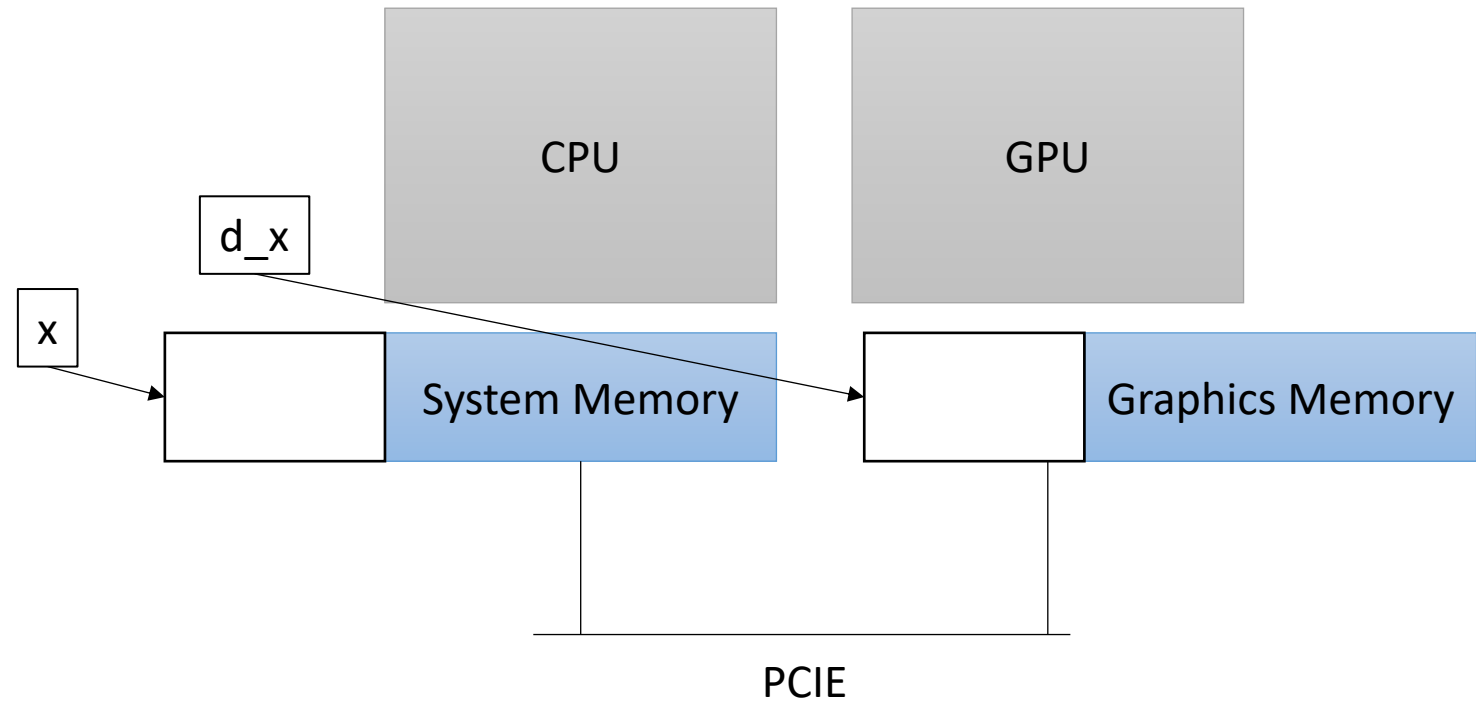


GPU set up

- Our heterogeneous, parallel, programming model

If we can't access d_x on the CPU, how do we initialize the memory?

GPU has no access to input devices e.g. disk



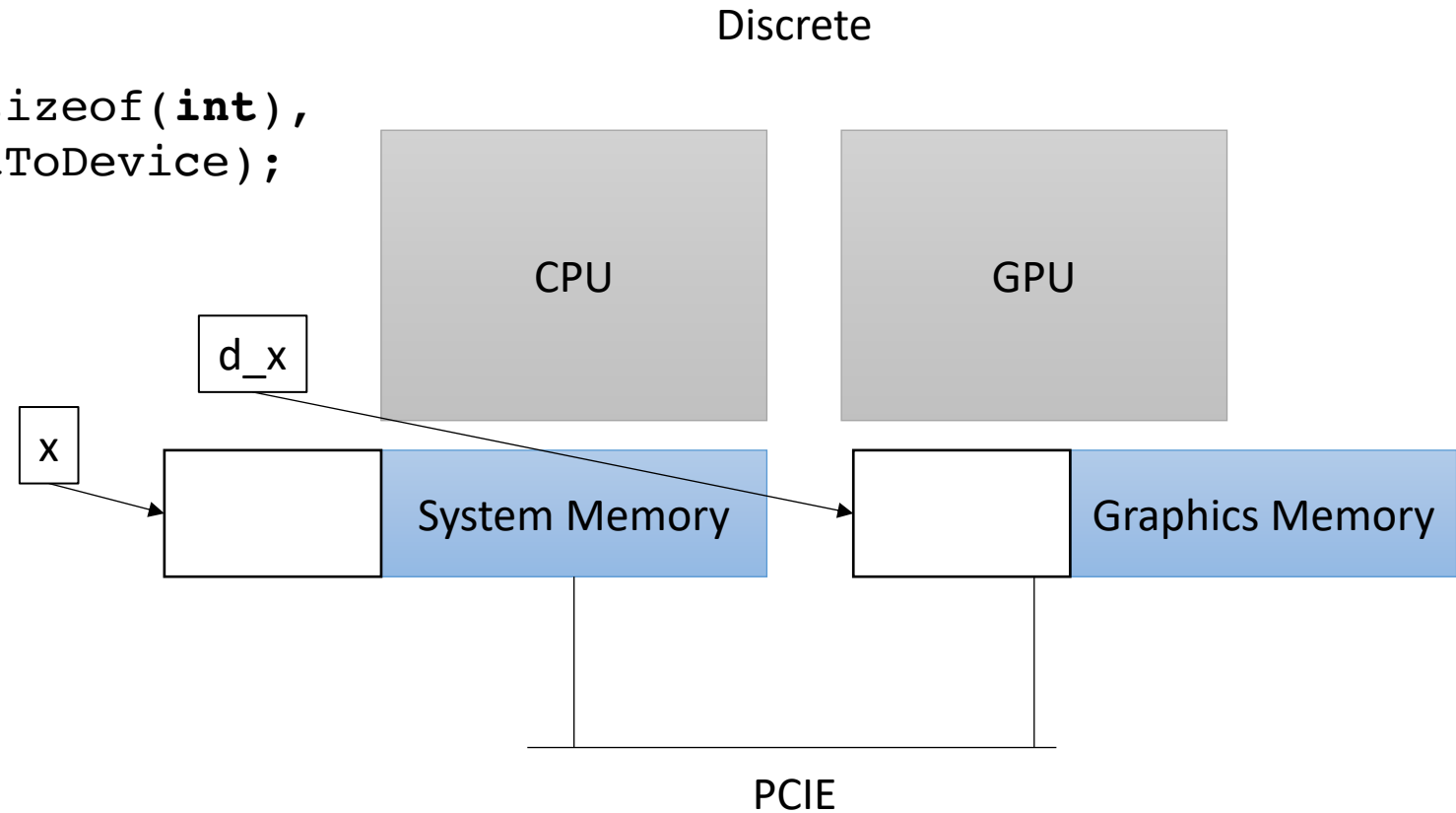
GPU set up

- Our heterogeneous, parallel, programming model

If we can't access `d_x` on the CPU, how do we initialize the memory?

GPU has no access to input devices e.g. disk

```
//initialize x on host  
cudaMemcpy(d_x, x, SIZE*sizeof(int),  
           cudaMemcpyHostToDevice);
```



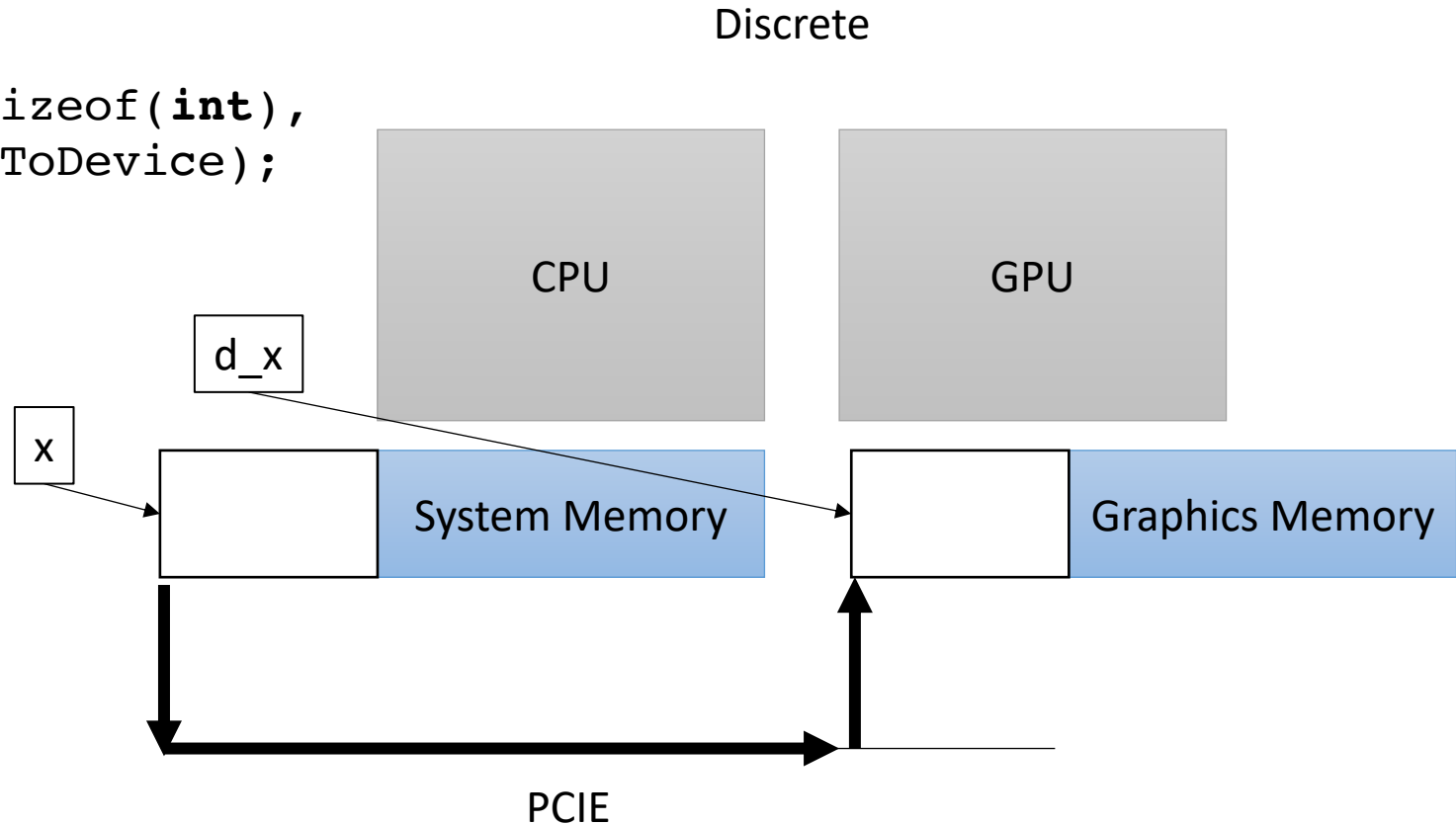
GPU set up

- Our heterogeneous, parallel, programming model

If we can't access `d_x` on the CPU, how do we initialize the memory?

GPU has no access to input devices e.g. disk

```
//initialize x on host  
cudaMemcpy(d_x, x, SIZE*sizeof(int),  
           cudaMemcpyHostToDevice);
```



How does this look in code?

How does this look in code?

Nothing too exciting yet.

The GPU Program

- Write a special function in your C++ code.
 - Called a Kernel
 - Use the new keyword `__global__`
 - Keywords in
 - OpenCL `__kernel`
 - Metal `kernel`
- Write it how you'd write any other function

The GPU Program

```
__global__ void vector_add(int * a, int * b, int * c, int size) {  
    for (int i = 0; i < size; i++) {  
        a[i] = b[i] + c[i];  
    }  
}
```

The GPU Program

```
__global__ void vector_add(int * a, int * b, int * c, int size) {  
    for (int i = 0; i < size; i++) {  
        a[i] = b[i] + c[i];  
    }  
}
```

calling the function

```
vector_add<<<1,1>>>(d_a, d_b, d_c, size);
```

The GPU Program

```
__global__ void vector_add(int * a, int * b, int * c, int size) {  
    for (int i = 0; i < size; i++) {  
        a[i] = b[i] + c[i];  
    }  
}
```

calling the function

What in the world?

special new CUDA syntax. We will talk more soon

```
vector_add<<<1,1>>>(d_a, d_b, d_c, size);
```

The GPU Program

```
__global__ void vector_add(int * d_a, int * d_b, int * d_c, int size) {  
    for (int i = 0; i < size; i++) {  
        d_a[i] = d_b[i] + d_c[i];  
    }  
}
```

Pass in pointers to memory on the device

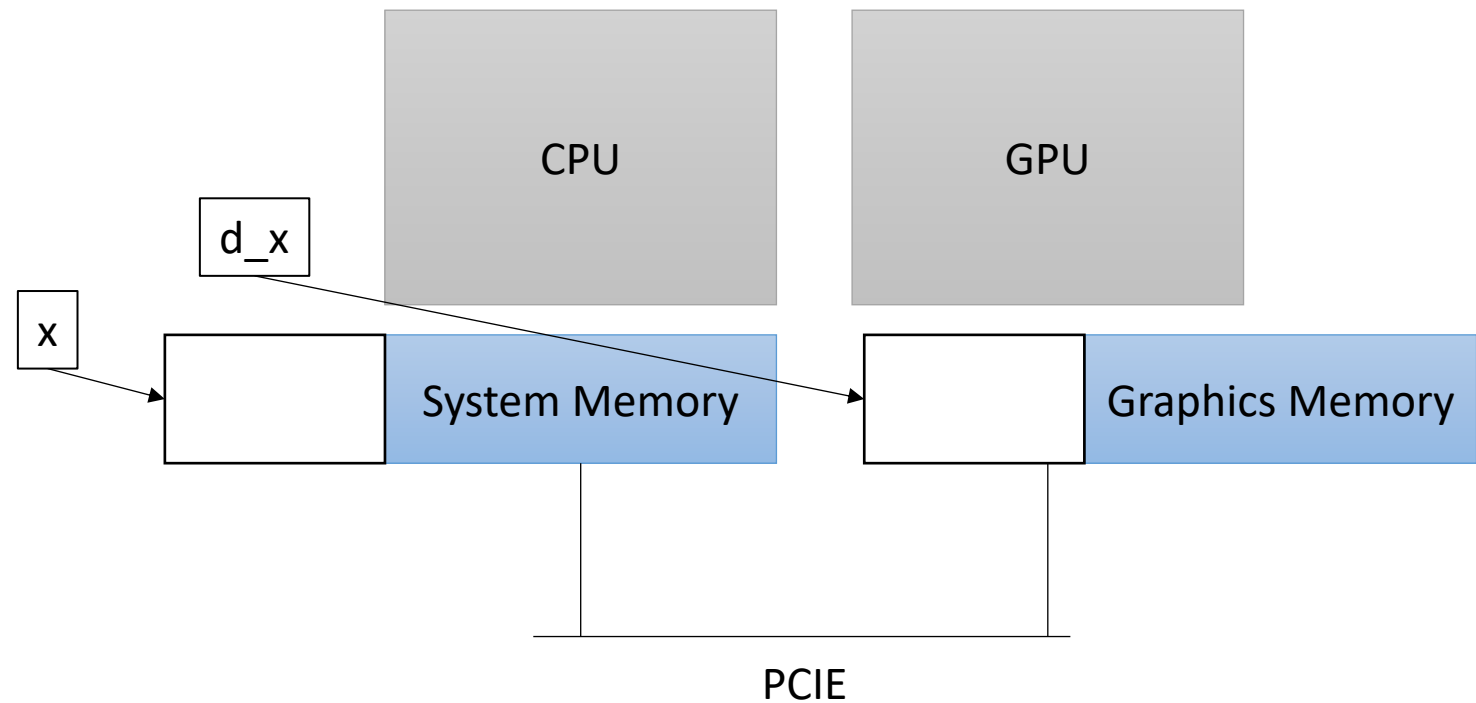
calling the function

```
vector_add<<<1,1>>>(d_a, d_b, d_c, size);
```


GPU set up

- Our heterogeneous, parallel, programming model

Remember, GPU needs to access its own memory



The GPU Program

```
__global__ void vector_add(int * d_a, int * d_b, int * d_c, int size) {  
    for (int i = 0; i < size; i++) {  
        d_a[i] = d_b[i] + d_c[i];  
    }  
}
```

Constants can be passed in regularly

calling the function

```
vector_add<<<1,1>>>(d_a, d_b, d_c, size);
```

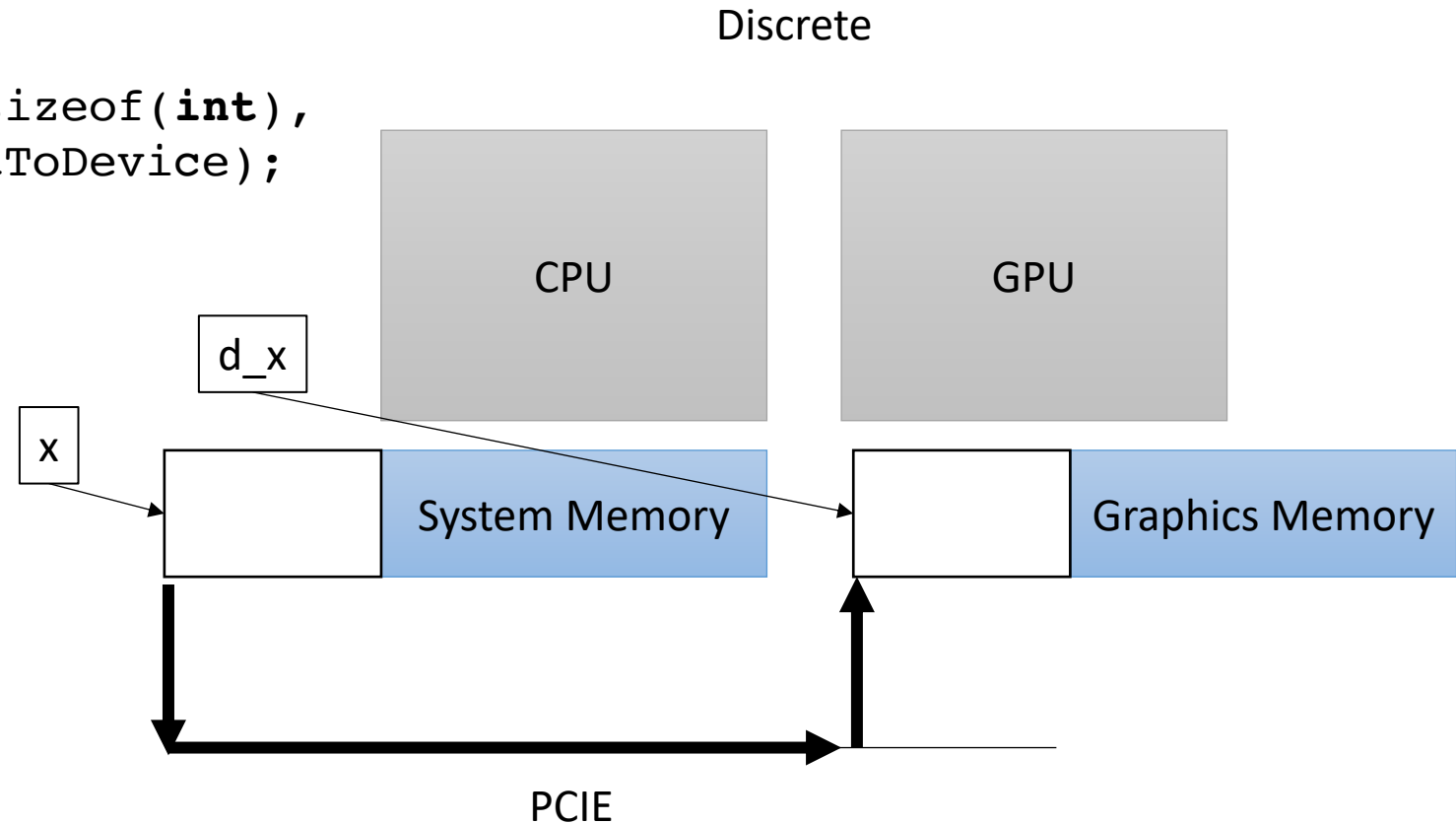
The GPU Program

Are we ready to run the program? What are we missing?

GPU set up

- Our heterogeneous, parallel, programming model

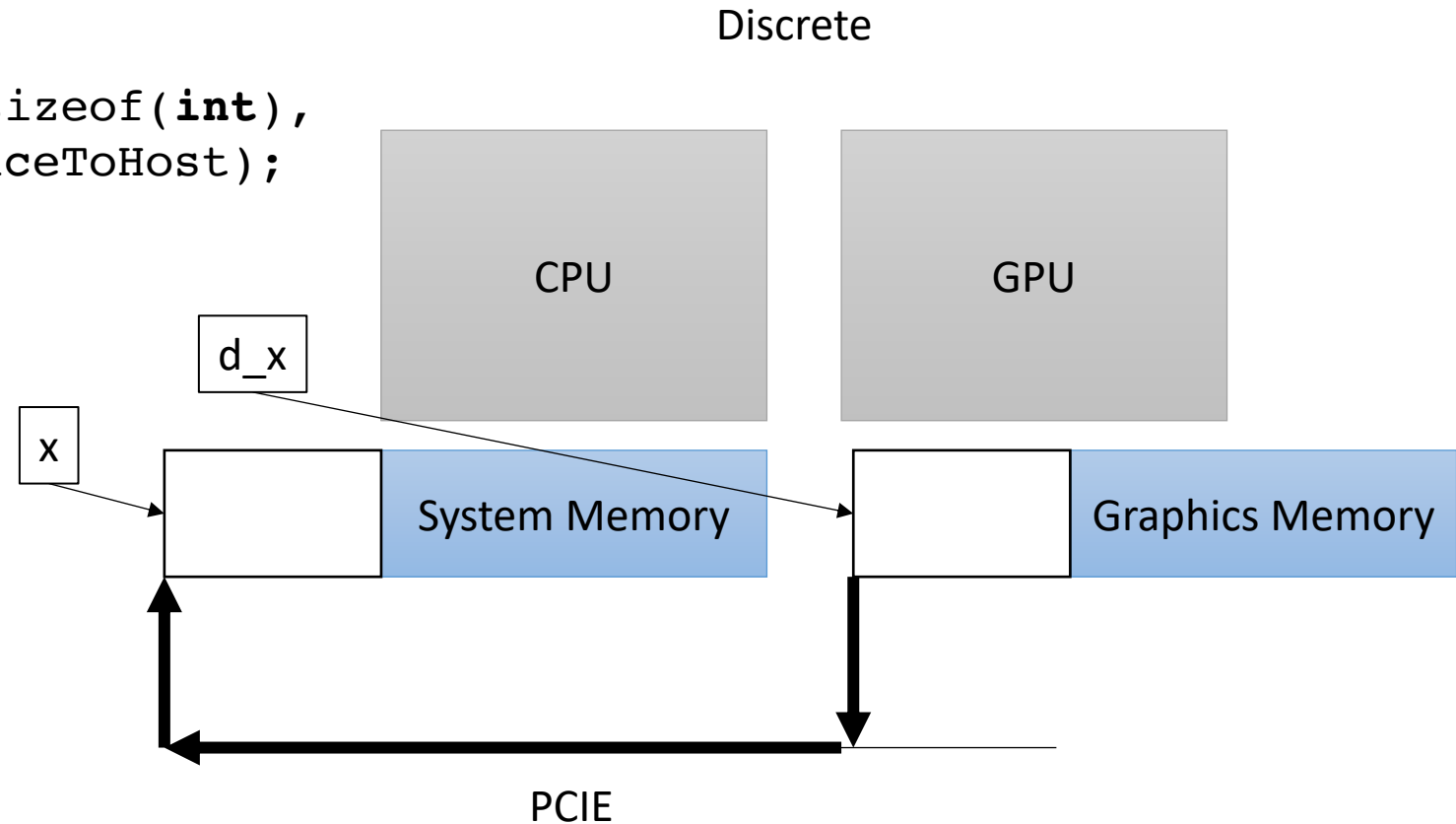
```
//initialize x on host  
cudaMemcpy(d_x, x, SIZE*sizeof(int),  
           cudaMemcpyHostToDevice);
```



GPU set up

- Our heterogeneous, parallel, programming model

```
//initialize x on host  
cudaMemcpy(x, d_x, SIZE*sizeof(int),  
           cudaMemcpyDeviceToHost);
```



The GPU Program

Finally, we can run the GPU program!

Lets see what all the hype is about

The GPU Program



It didn't do so well...

First parallelization attempt

- Lets look at some GPU documentation.
- The Maxwell whitepaper shows a diagram of one of the GPU cores



woah, 32 cores!

We should parallelize our application!



First parallelization attempt

```
__global__ void vector_add(int * d_a, int * d_b, int * d_c, int size) {  
    for (int i = 0; i < size; i++) {  
        d_a[i] = d_b[i] + d_c[i];  
    }  
}
```

calling the function

```
vector_add<<<1,1>>>(d_a, d_b, d_c, size);
```

First parallelization attempt

```
__global__ void vector_add(int * d_a, int * d_b, int * d_c, int size) {  
    for (int i = 0; i < size; i++) {  
        d_a[i] = d_b[i] + d_c[i];  
    }  
}
```

calling the function

```
vector_add<<<1, 32>>>(d_a, d_b, d_c, size);
```

number of threads to launch the program with

First parallelization attempt

```
__global__ void vector_add(int * d_a, int * d_b, int * d_c, int size) {  
    int chunk_size = size/blockDim.x;  
    int start = chunk_size * threadIdx.x;  
    int end = start + end;  
    for (int i = start; i < end; i++) {  
        d_a[i] = d_b[i] + d_c[i];  
    }  
}
```

calling the function

```
vector_add<<<1,32>>>(d_a, d_b, d_c, size);
```

number of threads

First parallelization attempt

```
__global__ void vector_add(int * d_a, int * d_b, int * d_c, int size) {  
    int chunk_size = size/blockDim.x;  
    int start = chunk_size * threadIdx.x;  
    int end = start + end;  
    for (int i = start; i < end; i++) {  
        d_a[i] = d_b[i] + d_c[i];  
    }  
}
```

calling the function

```
vector_add<<<1,32>>>(d_a, d_b, d_c, size);
```

number of threads
thread id

First parallelization attempt

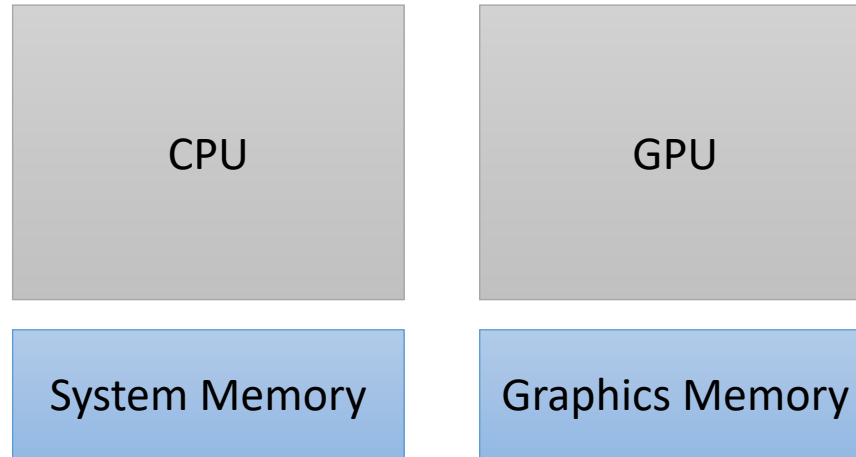
Lets try it! What do we think?

First parallelization attempt



Getting better but we have a long ways to go!

GPU Memory



GPU Memory

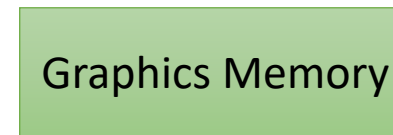
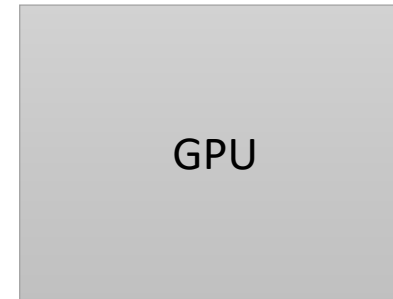
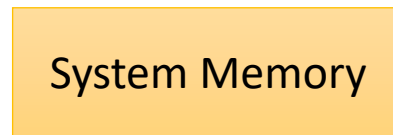
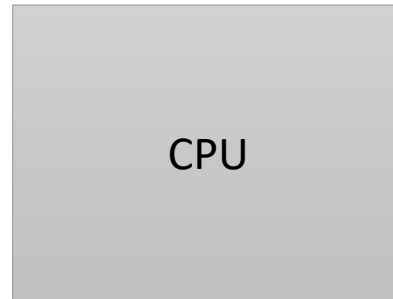
CPU Memory:

Fast: Low Latency

Easily saturated: Low Bandwidth

Scales well: up to 1 TB

DDR



GPU Memory:

slow: High Latency

hard to saturate: High Bandwidth

doesn't scale: 32 GB

GDDR, HBM

*2-lane straight highway
driven on by sports cars*

Different technologies

*16-lane highway on a windy
road driven by semi trucks*

GPU Memory

bandwidth:

~**700 GB/s** for GPU

~**50 GB/s** for CPUs

memory Latency:

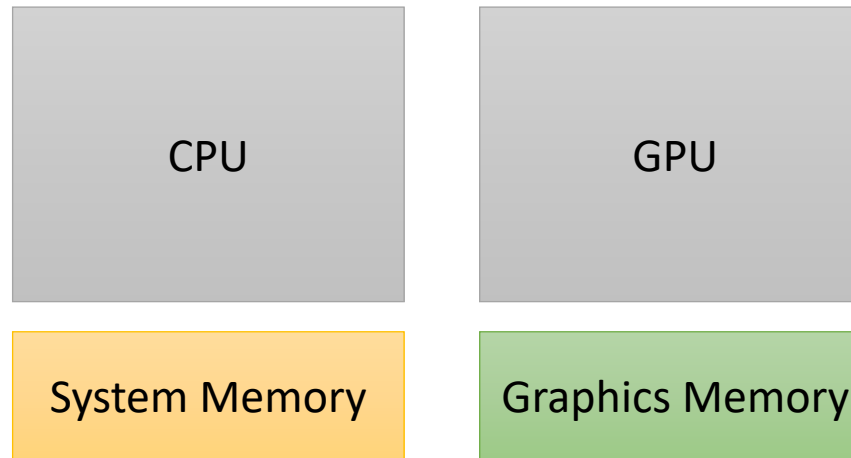
~**600** cycles for GPU memory

~**200** cycles for CPU memory

Cache Latency:

~**28** cycles for L1 hit for GPU

~**4** cycles for L1 hit on CPUs



Preemption and concurrency?

warp 0



GPU

Graphics Memory

Preemption and concurrency?

warp 0

all threads load from memory.

GPU

Graphics Memory

Preemption and concurrency?

warp 0

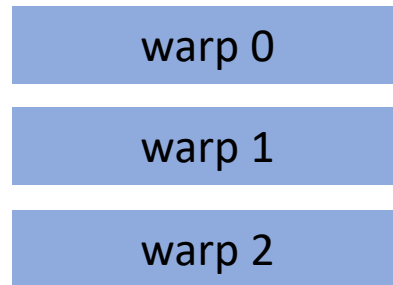
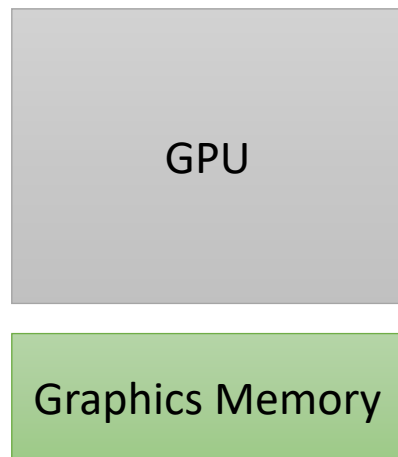
all threads load from memory.

GPU

600 cycles!

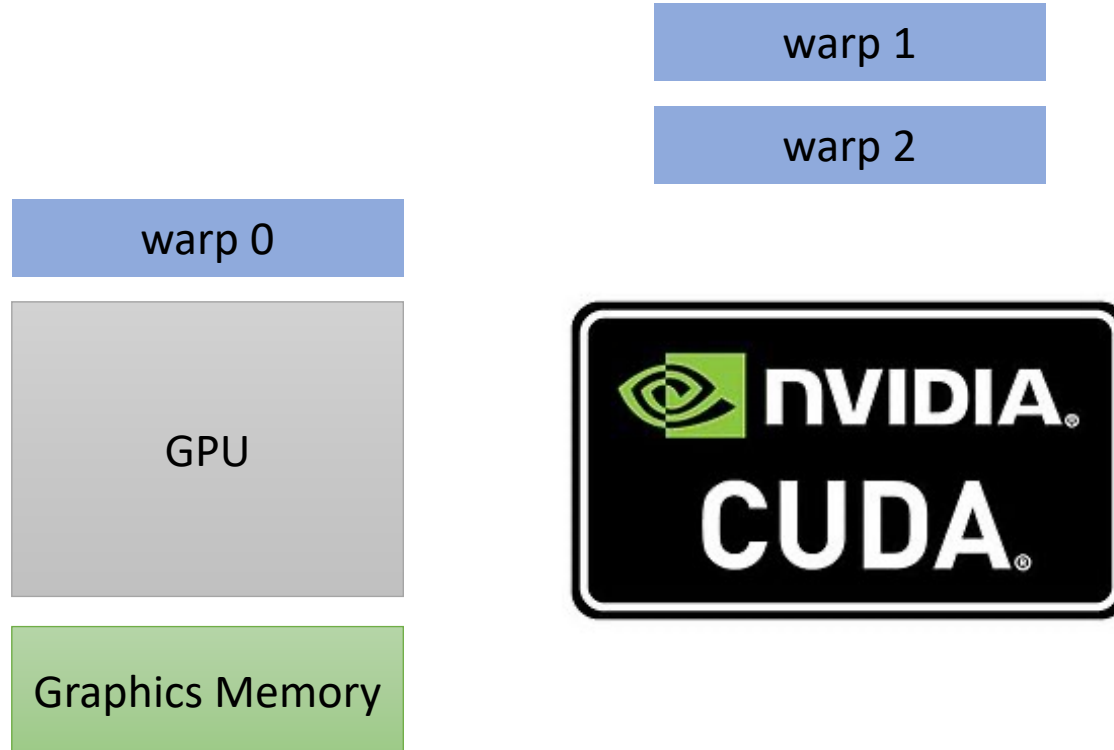
Graphics Memory

Preemption and concurrency?



We can hide latency through
preemption and concurrency!

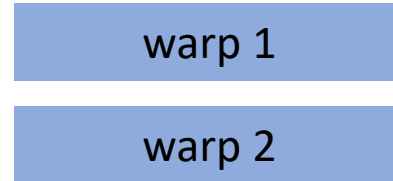
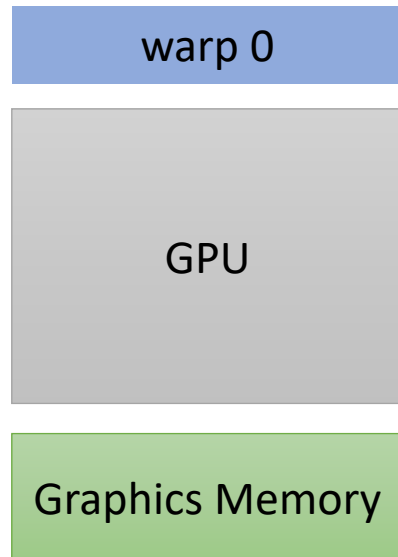
Preemption and concurrency?



We can hide latency through
preemption and concurrency!

Preemption and concurrency?

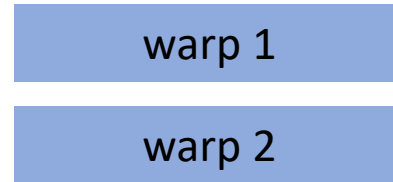
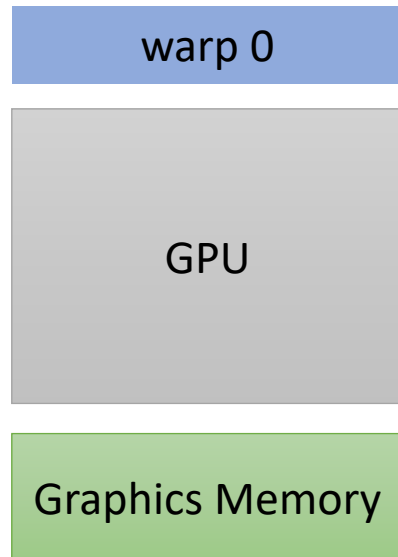
memory access
600 cycles



We can hide latency through
preemption and concurrency!

Preemption and concurrency?

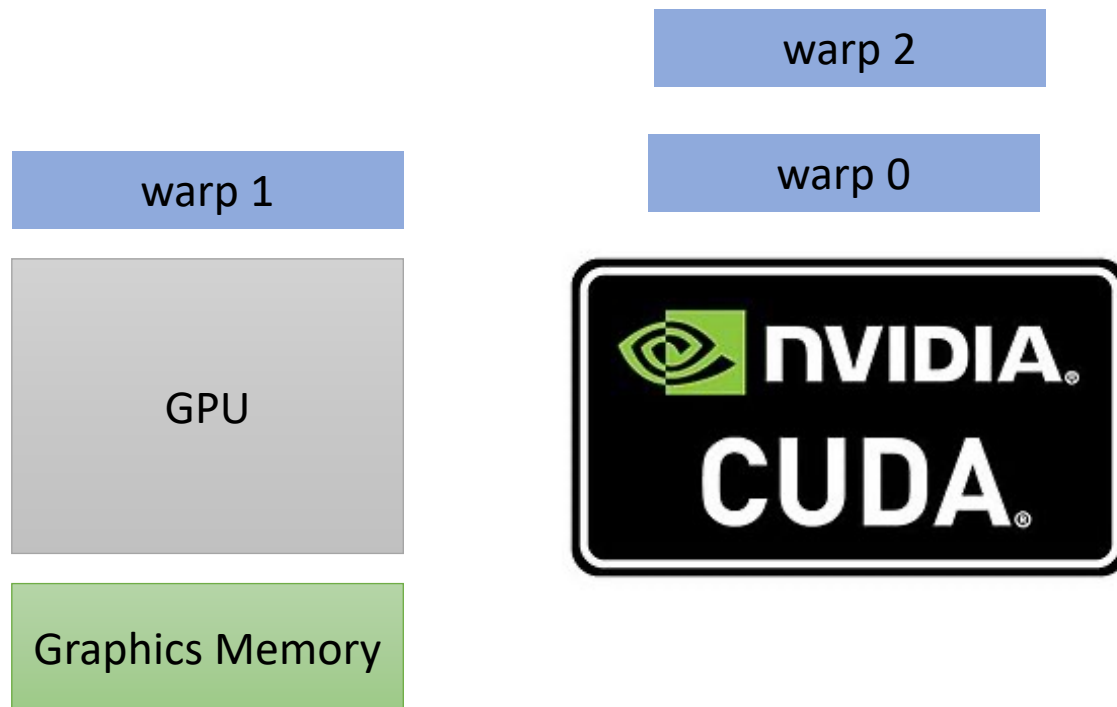
memory access
600 cycles



preempt warp 0
and put warp 1 on

We can hide latency through
preemption and concurrency!

Preemption and concurrency?



We can hide latency through
preemption and concurrency!

Preemption and concurrency?

memory access
600 cycles

warp 1

GPU

Graphics Memory

warp 2

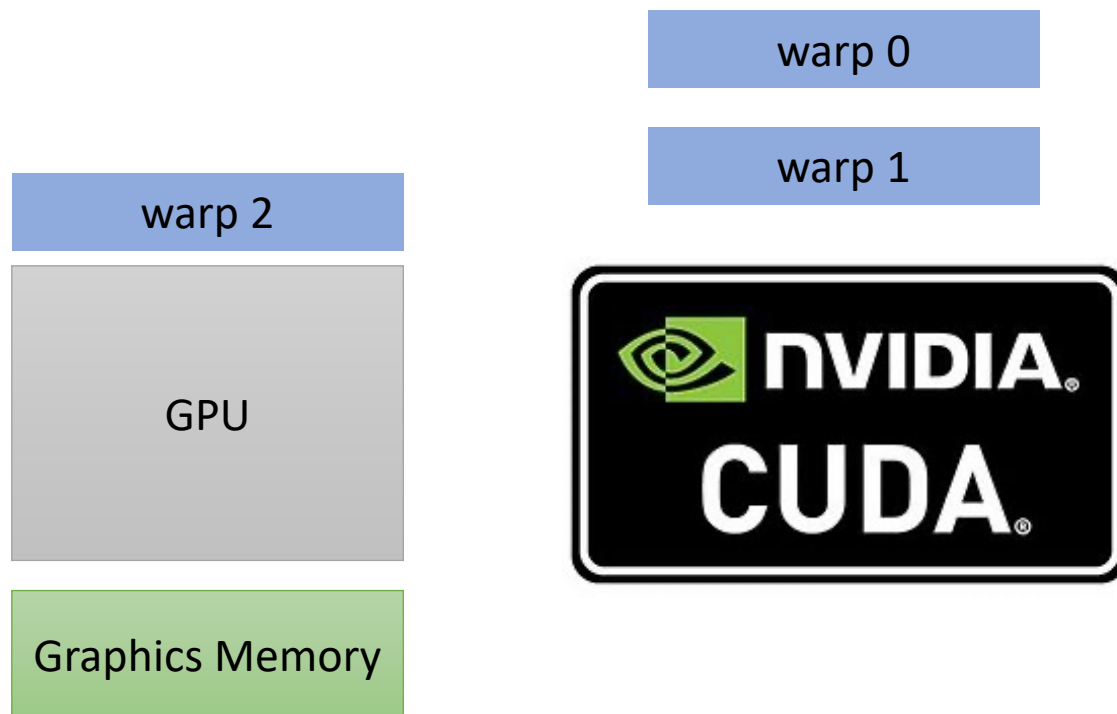
warp 0



preempt warp 1
and put warp 2 on

We can hide latency through
preemption and concurrency!

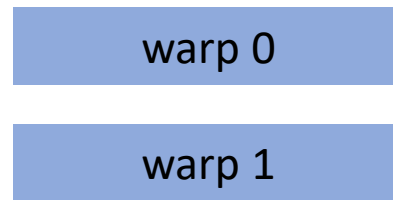
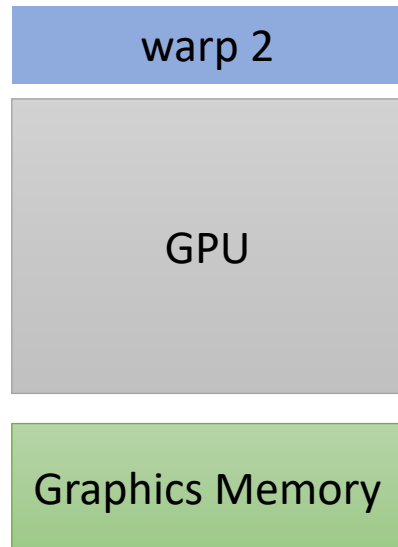
Preemption and concurrency?



We can hide latency through
preemption and concurrency!

Preemption and concurrency?

memory access
600 cycles

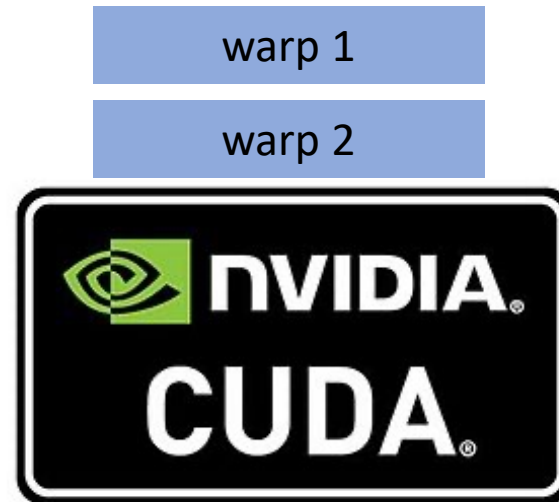
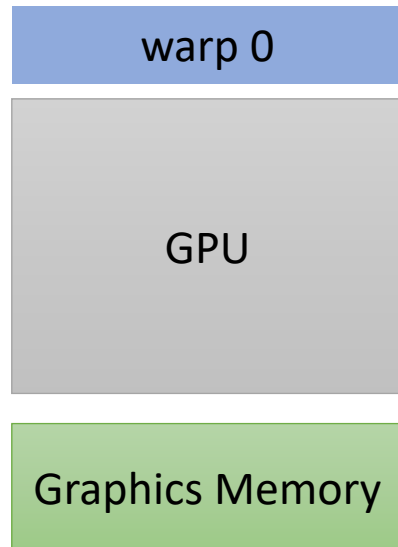


preempt warp 2
and put warp 0 on

We can hide latency through
preemption and concurrency!

Preemption and concurrency?

Hey, my memory has arrived!



preempt warp 2
and put warp 0 on

We can hide latency through
preemption and concurrency!

Preemption and concurrency?

But wait, I thought preemption was expensive?

Preemption and concurrency?



But wait, I thought preemption was expensive?

Registers all stay on chip

Preemption and concurrency?



But wait, I thought preemption was expensive?

dedicated scheduler logic

Preemption and concurrency?



But wait, I thought preemption was expensive?

bound on number of warps: 32

Go back to our program

```
__global__ void vector_add(int * d_a, int * d_b, int * d_c, int size) {  
    int chunk_size = size/blockDim.x;  
    int start = chunk_size * threadIdx.x;  
    int end = start + end;  
    for (int i = start; i < end; i++) {  
        d_a[i] = d_b[i] + d_c[i];  
    }  
}
```

calling the function

```
vector_add<<<1, 32>>>(d_a, d_b, d_c, size);
```

Lets launch with 32 warps

Go back to our program

```
__global__ void vector_add(int * d_a, int * d_b, int * d_c, int size) {  
    int chunk_size = size/blockDim.x;  
    int start = chunk_size * threadIdx.x;  
    int end = start + chunk_size;  
    for (int i = start; i < end; i++) {  
        d_a[i] = d_b[i] + d_c[i];  
    }  
}
```

calling the function

Lets launch with 32 warps

```
vector_add<<<1, 1024>>>(d_a, d_b, d_c, size);
```

Concurrent warps

Lets try it! What do we think?

Concurrent warps

Lets try it! What do we think?

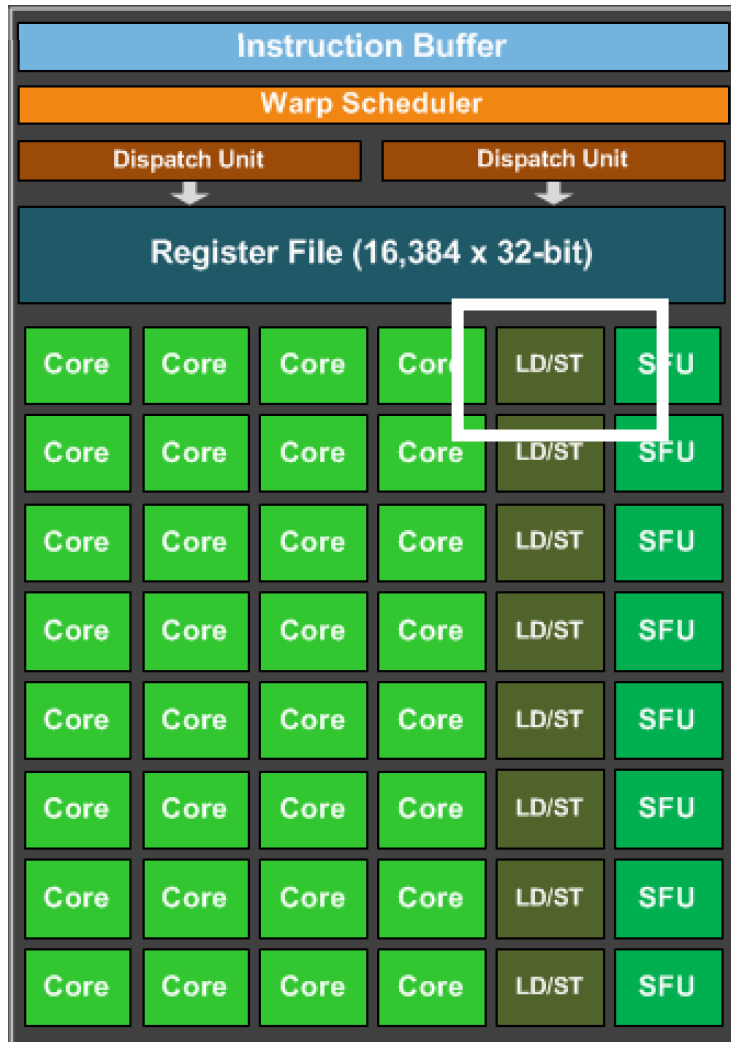


Getting better!

Optimizing memory accesses



Optimizing memory accesses



this is the load/store unit. The hardware component responsible for issuing loads and stores.

Why doesn't every core have one?

Optimizing memory accesses



This is the instruction cache... Why doesn't every core have a instruction buffer to keep track of its program?

this is the load/store unit. The hardware component responsible for issuing loads and stores.

Why doesn't every core have one?

Warp execution



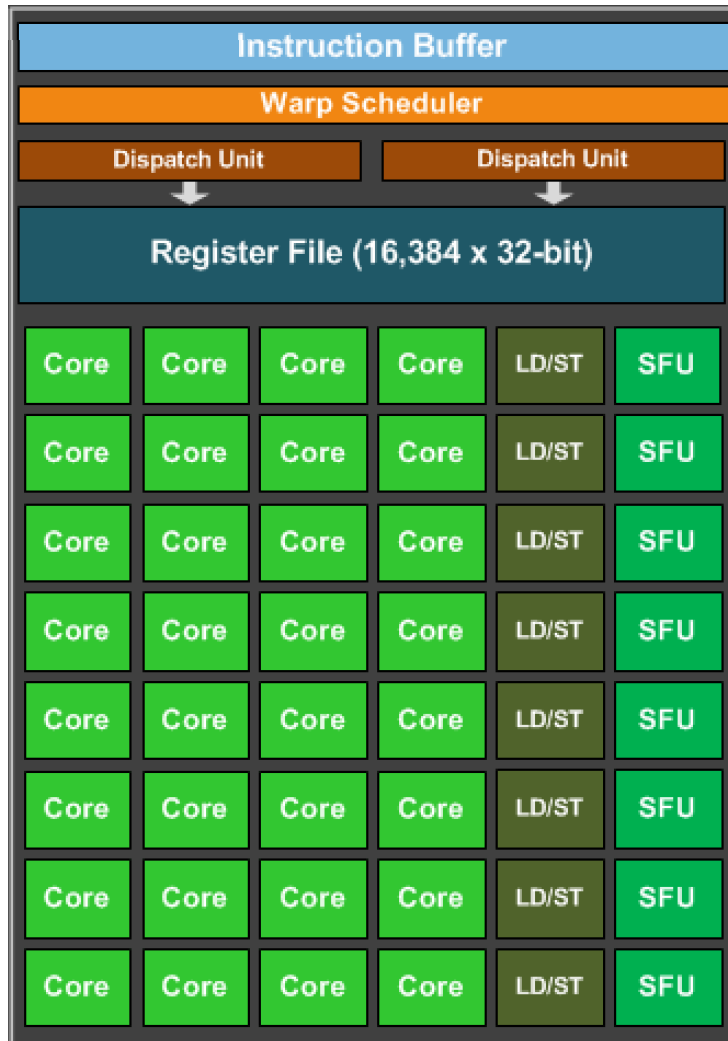
Groups of 32 threads are called a “warp”

They are executed in lock-step, i.e. they all execute the same instruction at the same time

Warp execution

Groups of 32 threads are called a “warp”

They are executed in lock-step, i.e. they all execute the same instruction at the same time



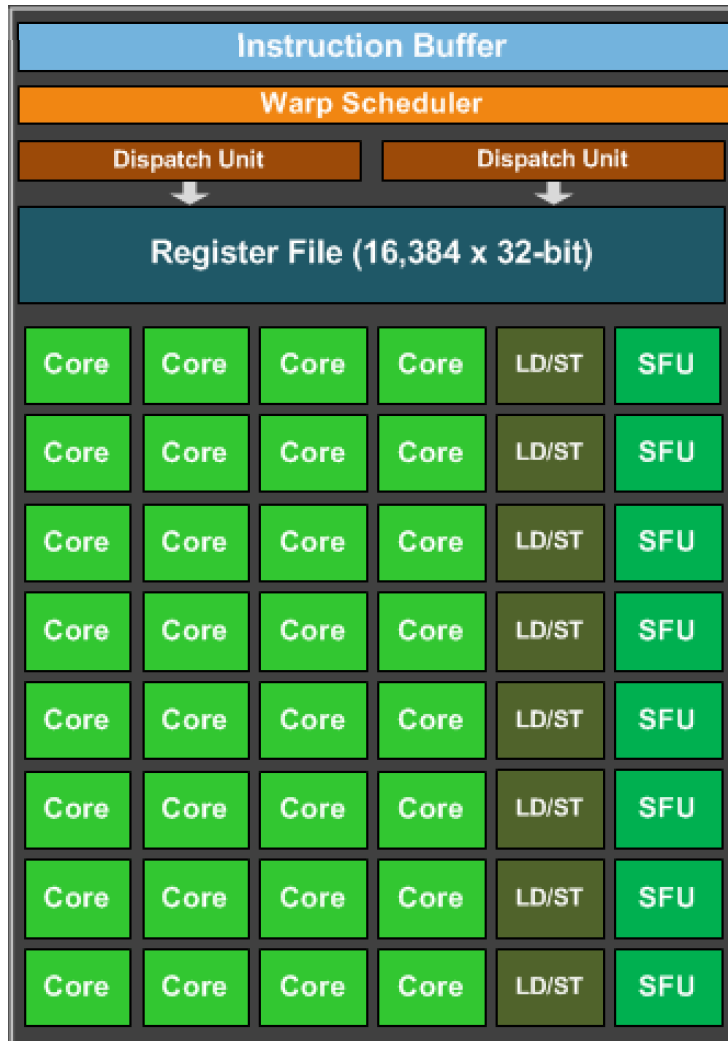
Program:

```
int variable1 = b[0];  
int variable2 = c[0];  
int variable3 = variable1 + variable2;  
a[0] = variable3;
```

Warp execution

Groups of 32 threads are called a “warp”

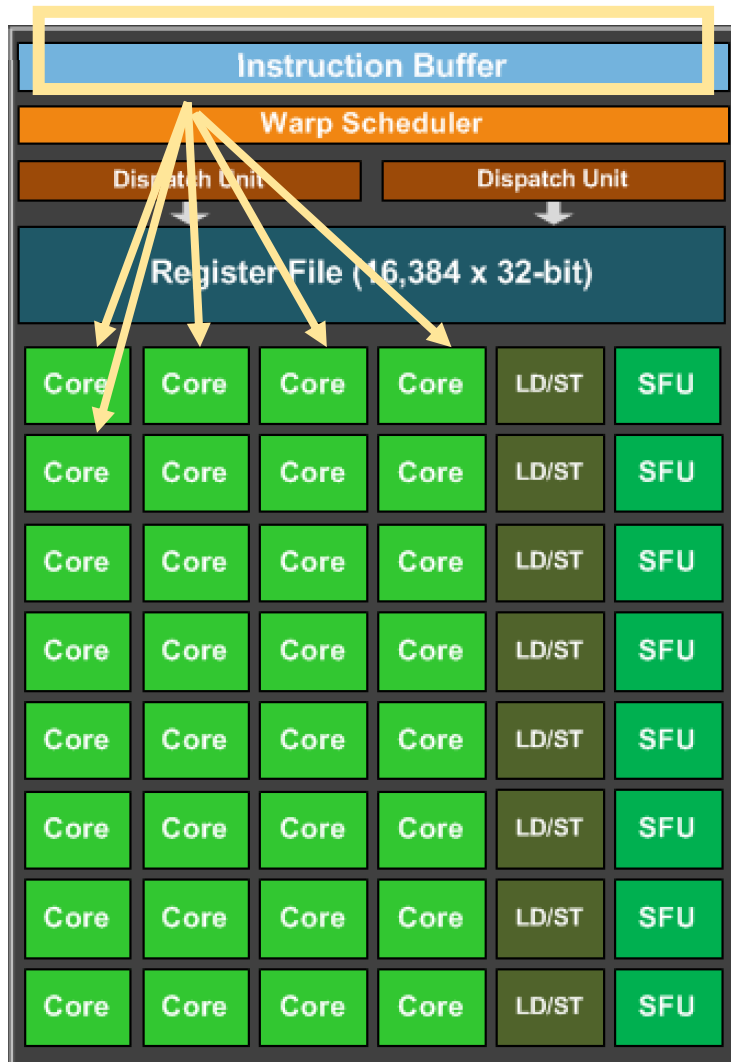
They are executed in lock-step, i.e. they all execute the same instruction at the same time



Program:

```
int variable1 = b[0];  
int variable2 = c[0];  
int variable3 = variable1 + variable2;  
a[0] = variable3;
```

Warp execution



Groups of 32 threads are called a “warp”

They are executed in lock-step, i.e. they all execute the same instruction at the same time

instruction is fetched from the buffer and distributed to all the cores.

Program:

```
int variable1 = b[0];  
int variable2 = c[0];  
int variable3 = variable1 + variable2;  
a[0] = variable3;
```

Warp execution

Groups of 32 threads are called a “warp”

They are executed in lock-step, i.e. they all execute the same instruction at the same time



Cores can a large register file
they share expensive HW units (load/store and special functions)

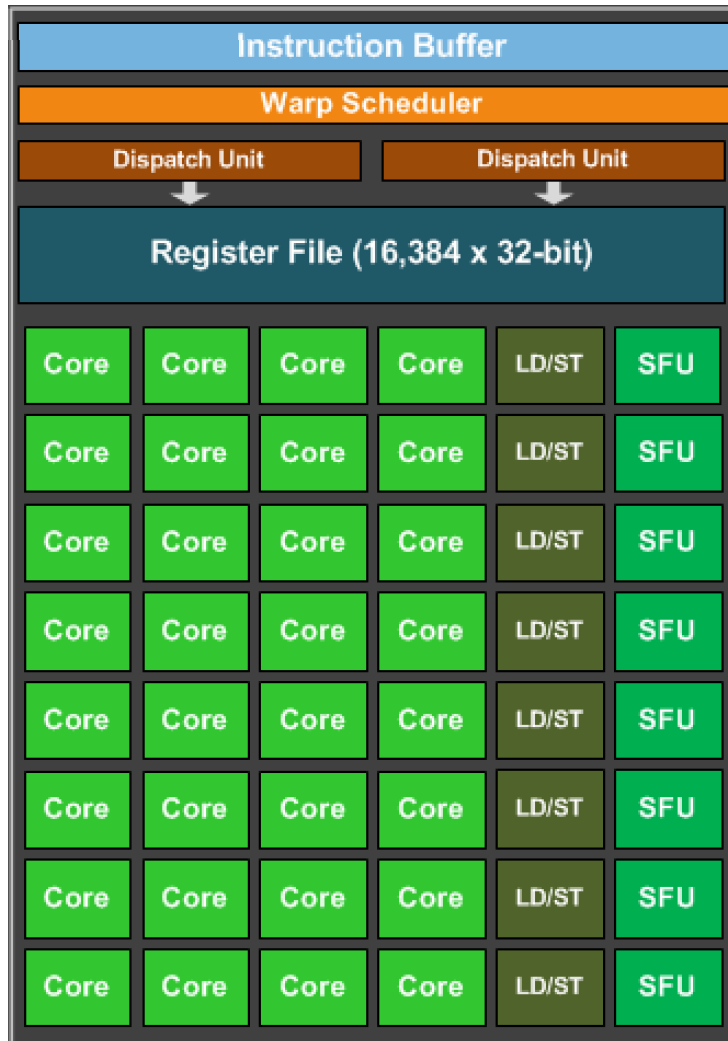
Program:

```
int variable1 = b[0];  
int variable2 = c[0];  
int variable3 = variable1 + variable2;  
a[0] = variable3;
```

Warp execution

Groups of 32 threads are called a “warp”

They are executed in lock-step, i.e. they all execute the same instruction at the same time

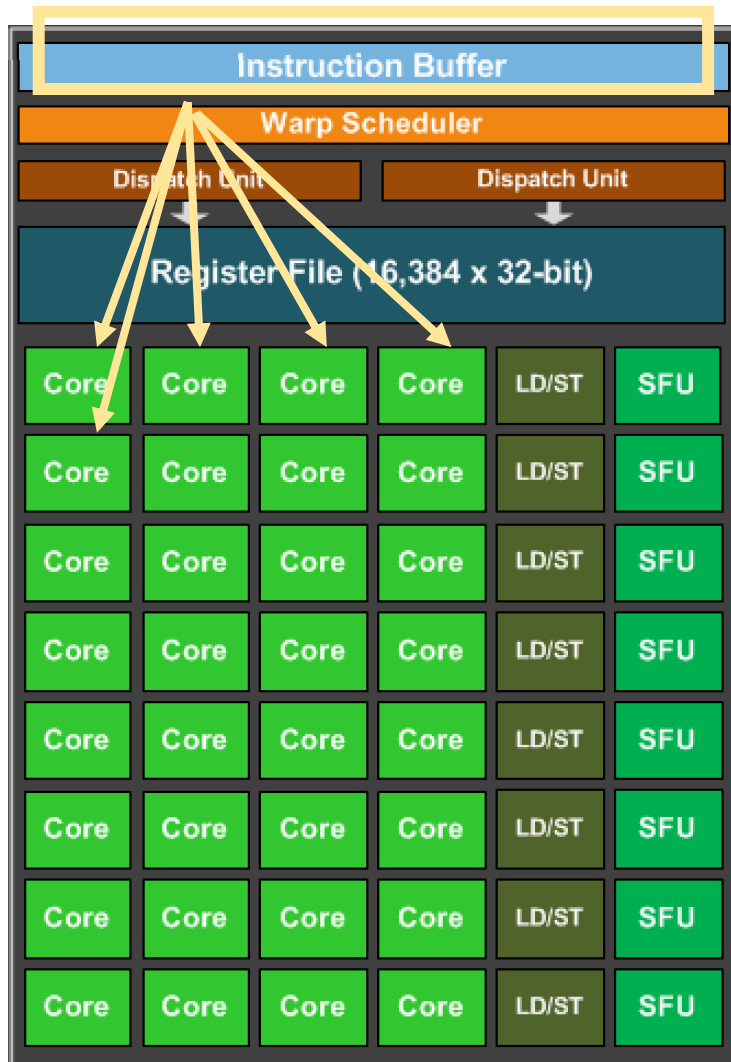


All cores need to wait until all cores finish the first instruction

Program:

```
int variable1 = b[0];  
int variable2 = c[0];  
int variable3 = variable1 + variable2;  
a[0] = variable3;
```


Warp execution



Start the next instruction.

Groups of 32 threads are called a “warp”

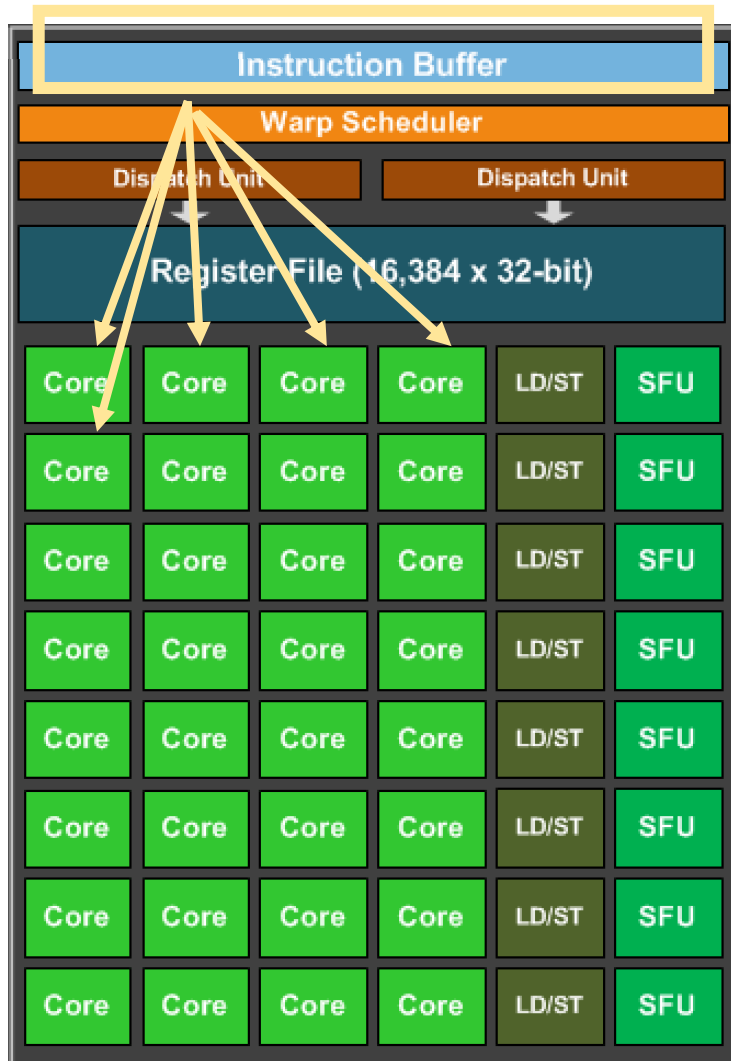
They are executed in lock-step, i.e. they all execute the same instruction at the same time

Program:

```
int variable1 = b[0];  
int variable2 = c[0];  
int variable3 = variable1 + variable2;  
a[0] = variable3;
```

Why would we have a programming model like this?

Warp execution



Groups of 32 threads are called a “warp”

They are executed in lock-step, i.e. they all execute the same instruction at the same time

Start the next instruction.

Program:

```
int variable1 = b[0];  
int variable2 = c[0];  
int variable3 = variable1 + variable2;  
a[0] = variable3;
```

Why would we have a programming model like this?

More cores (share program counters)

Can be efficient to share other hardware resources

Warp execution

Lets look closer at memory

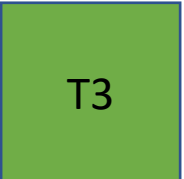
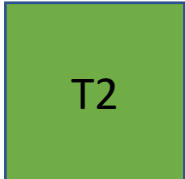
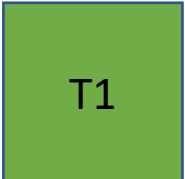
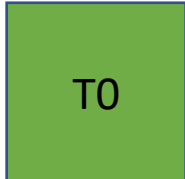
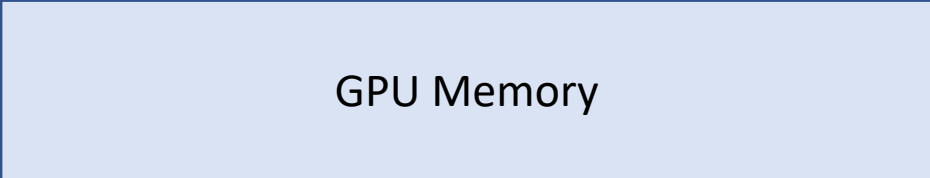


Program:

```
int variable1 = b[0];  
int variable2 = c[0];  
int variable3 = variable1 + variable2;  
a[0] = variable3;
```

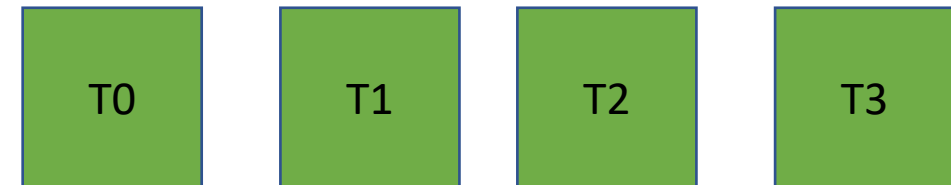
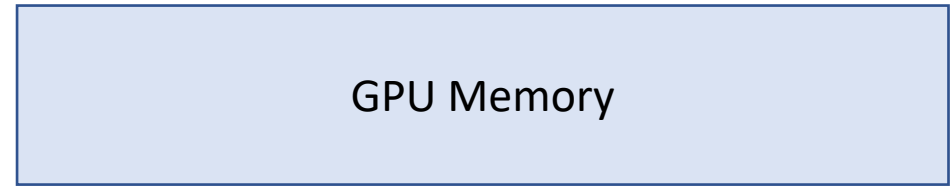
4 cores are accessing memory. what happens if they access the same value?

4 cores are accessing memory. What can happen



4 cores are accessing memory. What can happen

All read the same value



a[0]

a[0]

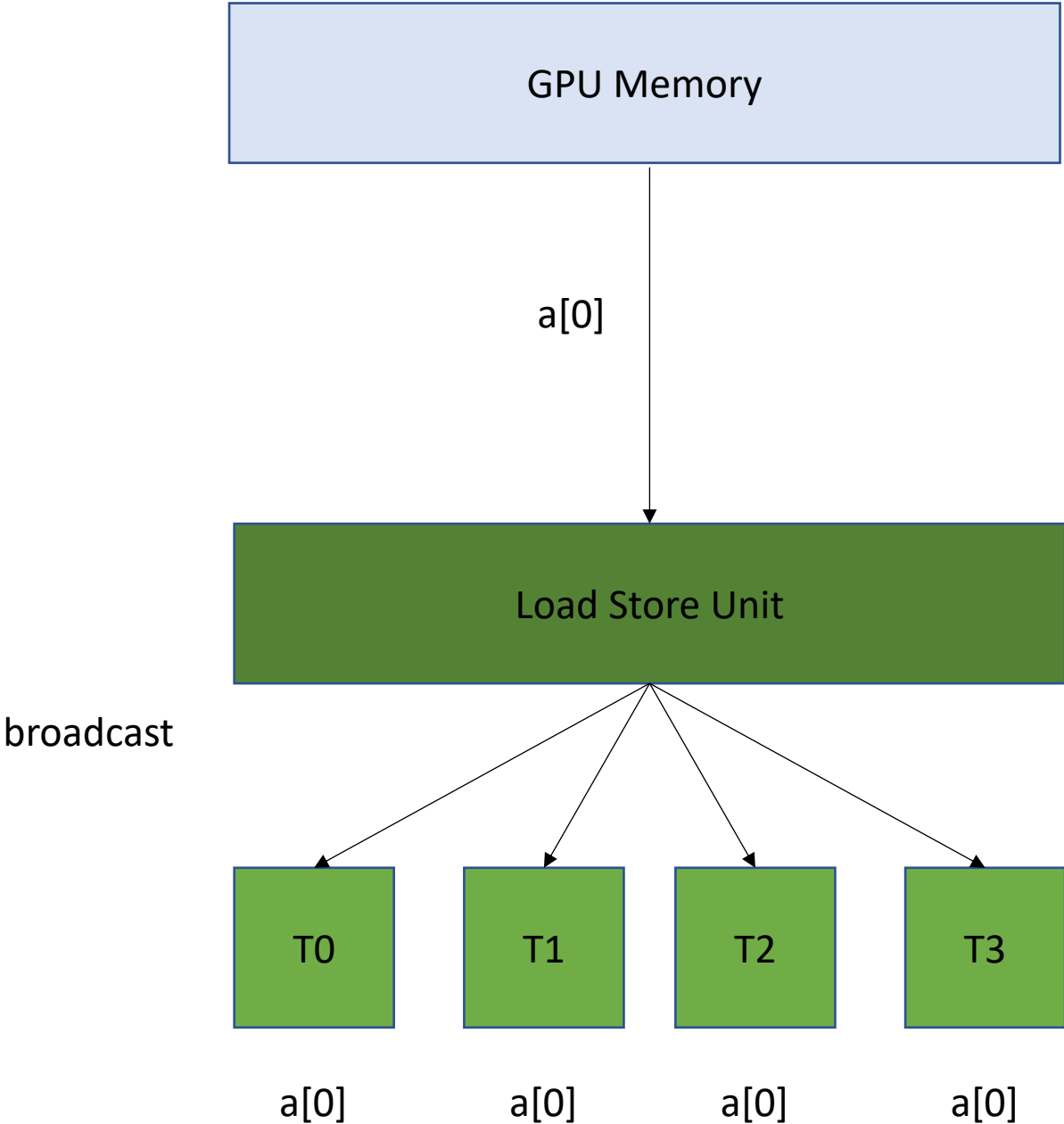
a[0]

a[0]

4 cores are accessing memory. What can happen

All read the same value

This is efficient: the load store unit can ask for the value and then broadcast it to all cores.



4 cores are accessing memory. What can happen

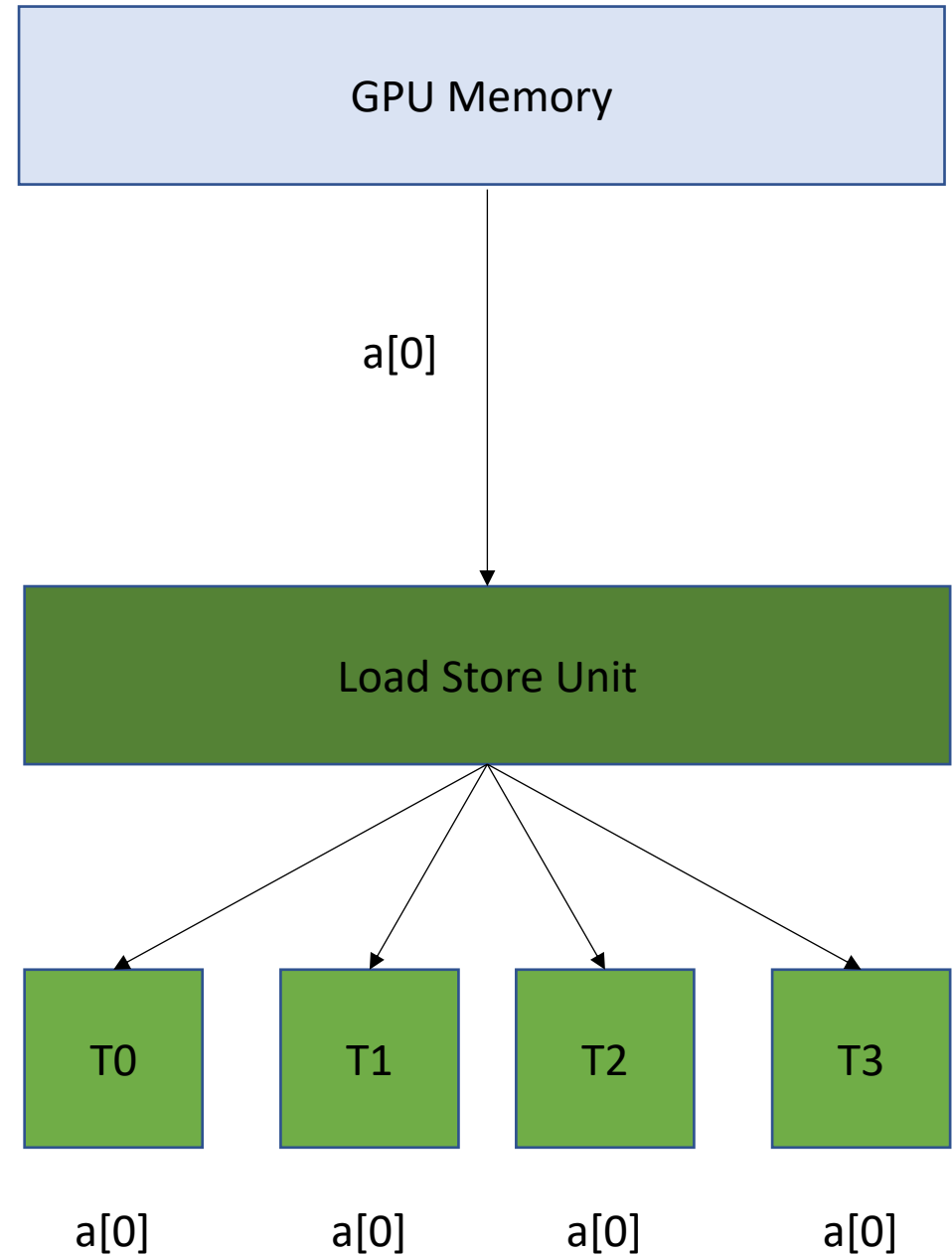
All read the same value

This is efficient: the load store unit can ask for the value and then broadcast it to all cores.

1 request to GPU memory

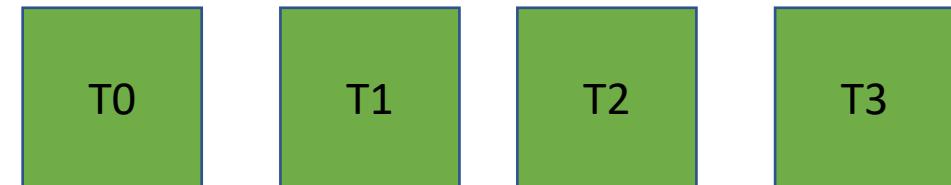
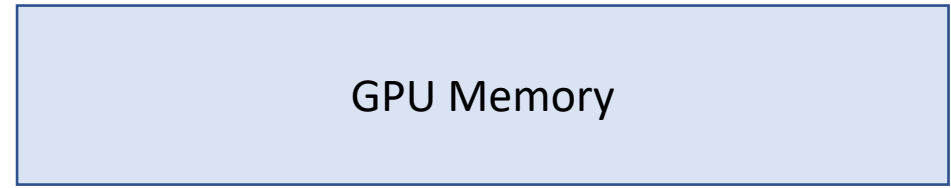
Efficient, but probably not too common.

broadcast



4 cores are accessing memory. What can happen

Read contiguous values



a[0]

a[1]

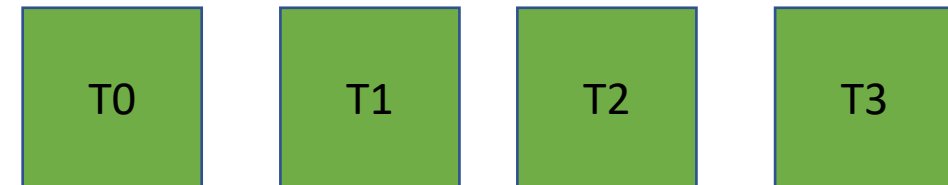
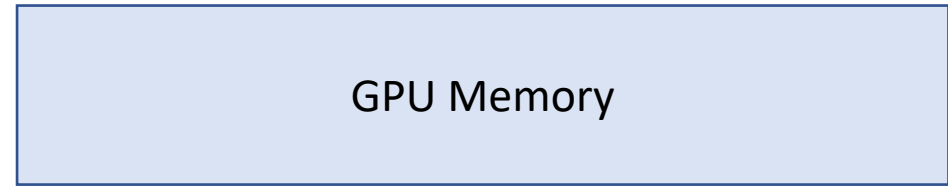
a[2]

a[3]

4 cores are accessing memory. What can happen

Read contiguous values

Like the CPU cache, the Load/Store Unit reads in memory in chunks. 16 bytes



a[0]

a[1]

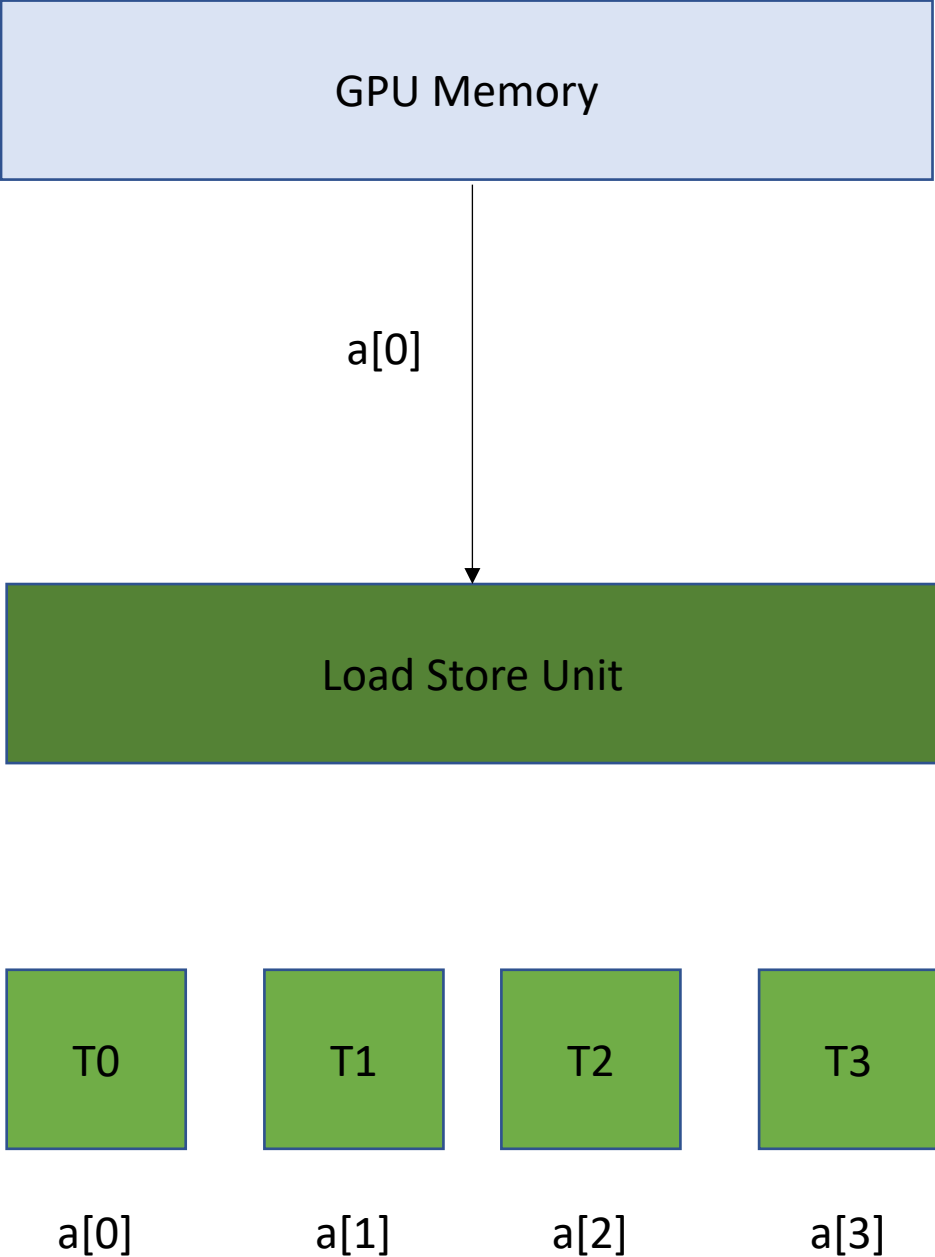
a[2]

a[3]

4 cores are accessing memory. What can happen

Read contiguous values

Like the CPU cache, the Load/Store Unit reads in memory in chunks. 16 bytes

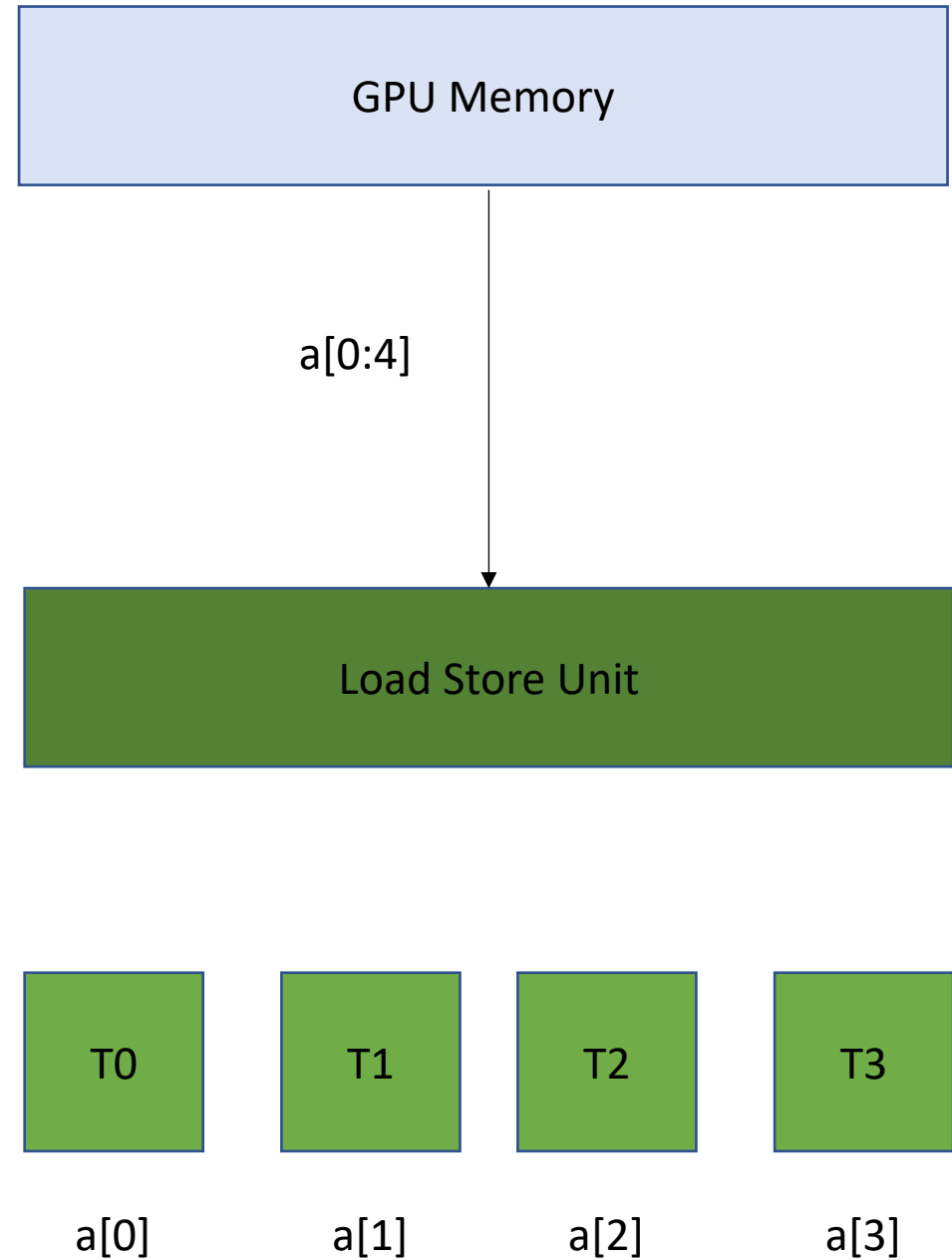


4 cores are accessing memory. What can happen

Read contiguous values

Like the CPU cache, the Load/Store Unit reads in memory in chunks. 16 bytes

Can easily distribute the values to the threads



4 cores are accessing memory. What can happen

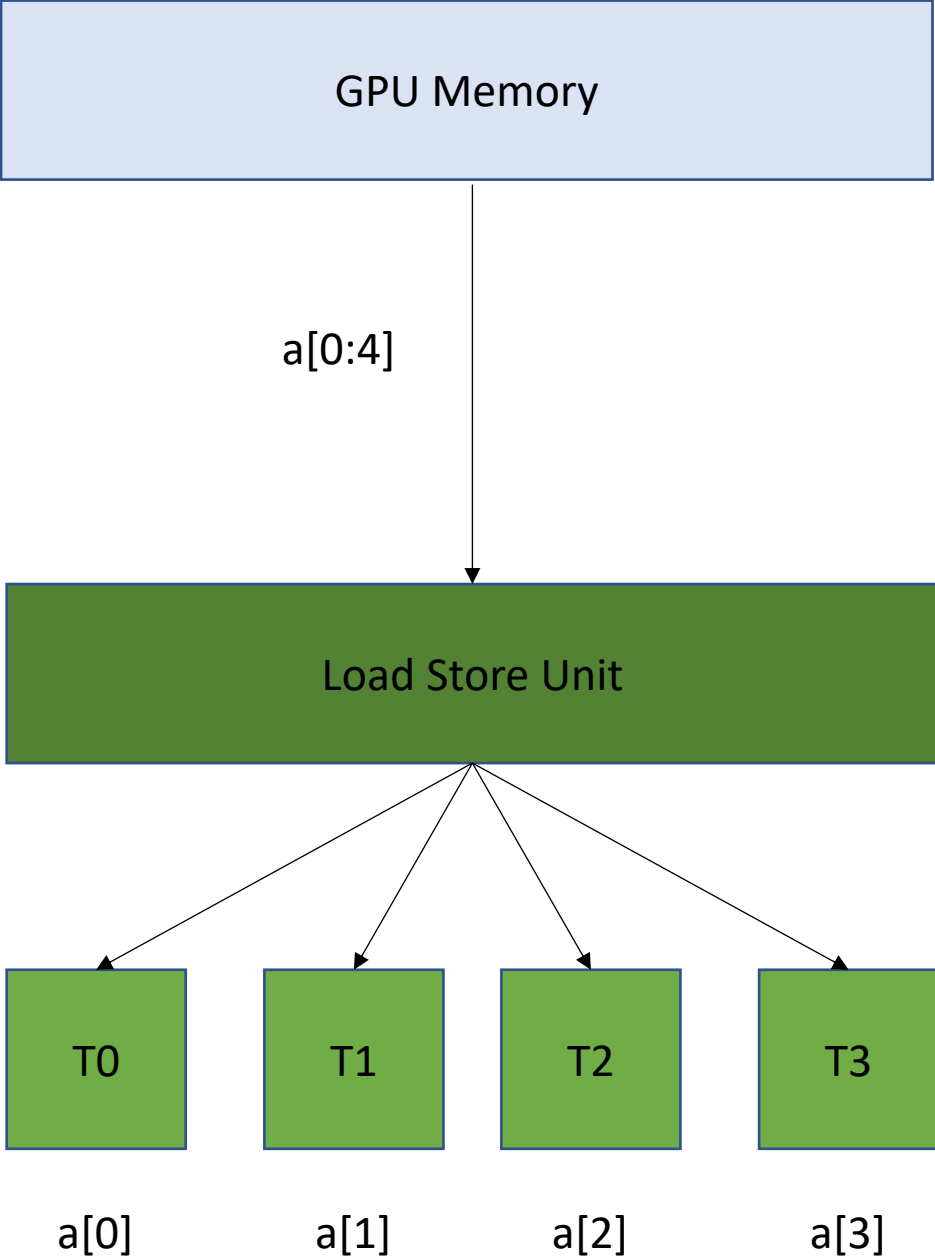
Read contiguous values

Like the CPU cache, the Load/Store Unit reads in memory in chunks. 16 bytes

Can easily distribute the values to the threads

1 request to GPU memory

stream



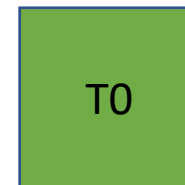
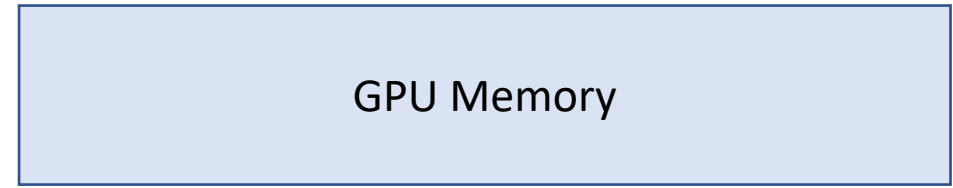
4 cores are accessing memory. What can happen

Read non-contiguous values

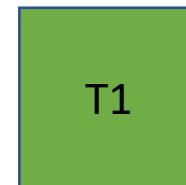
Not good!

Accesses are Serialized.

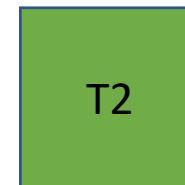
You need 4 requests to GPU memory



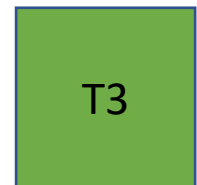
a[x]



a[y]



a[z]



a[w]

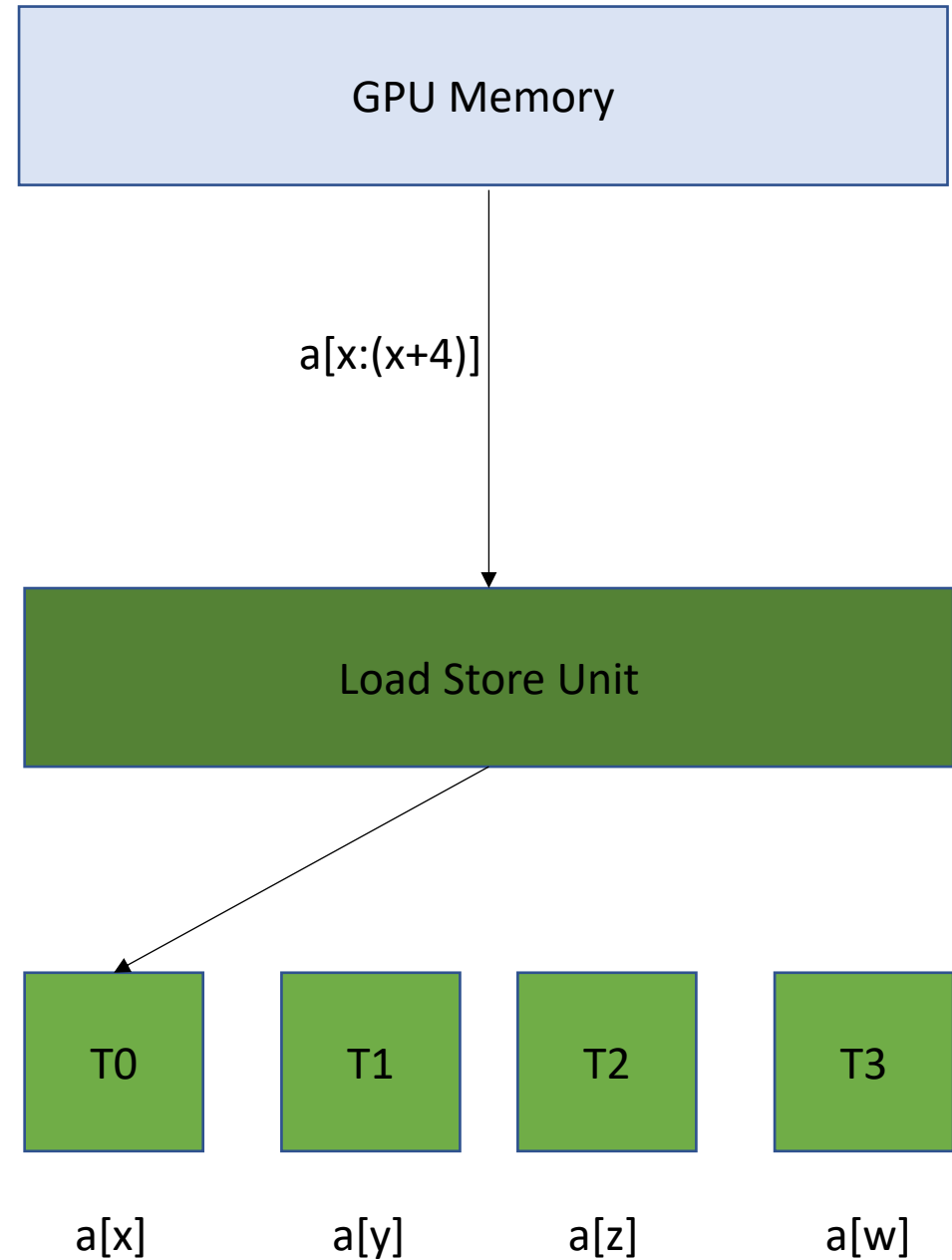
4 cores are accessing memory. What can happen

Read non-contiguous values

Not good!

Accesses are Serialized.

You need 4 requests to GPU memory



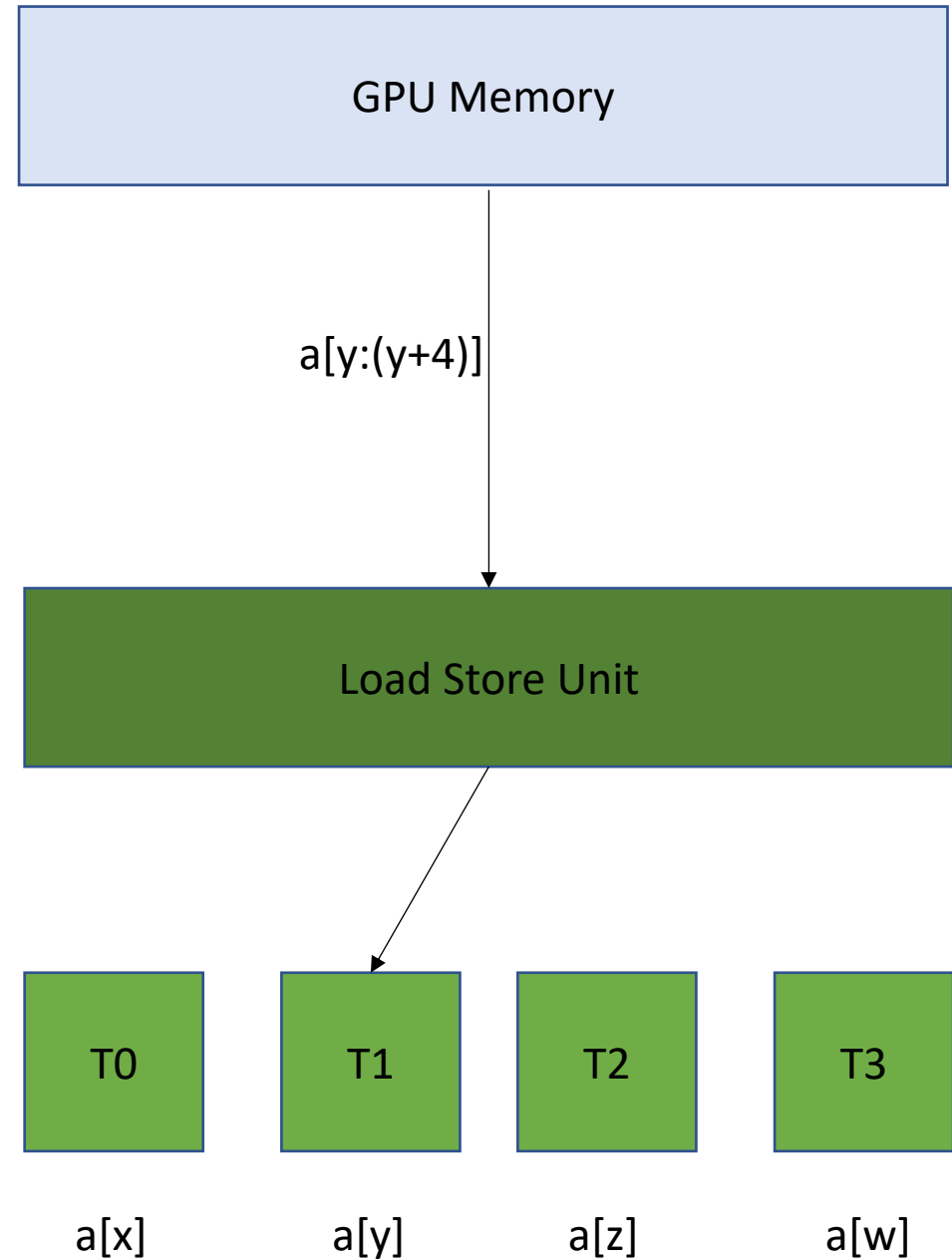
4 cores are accessing memory. What can happen

Read non-contiguous values

Not good!

Accesses are Serialized.

You need 4 requests to GPU memory



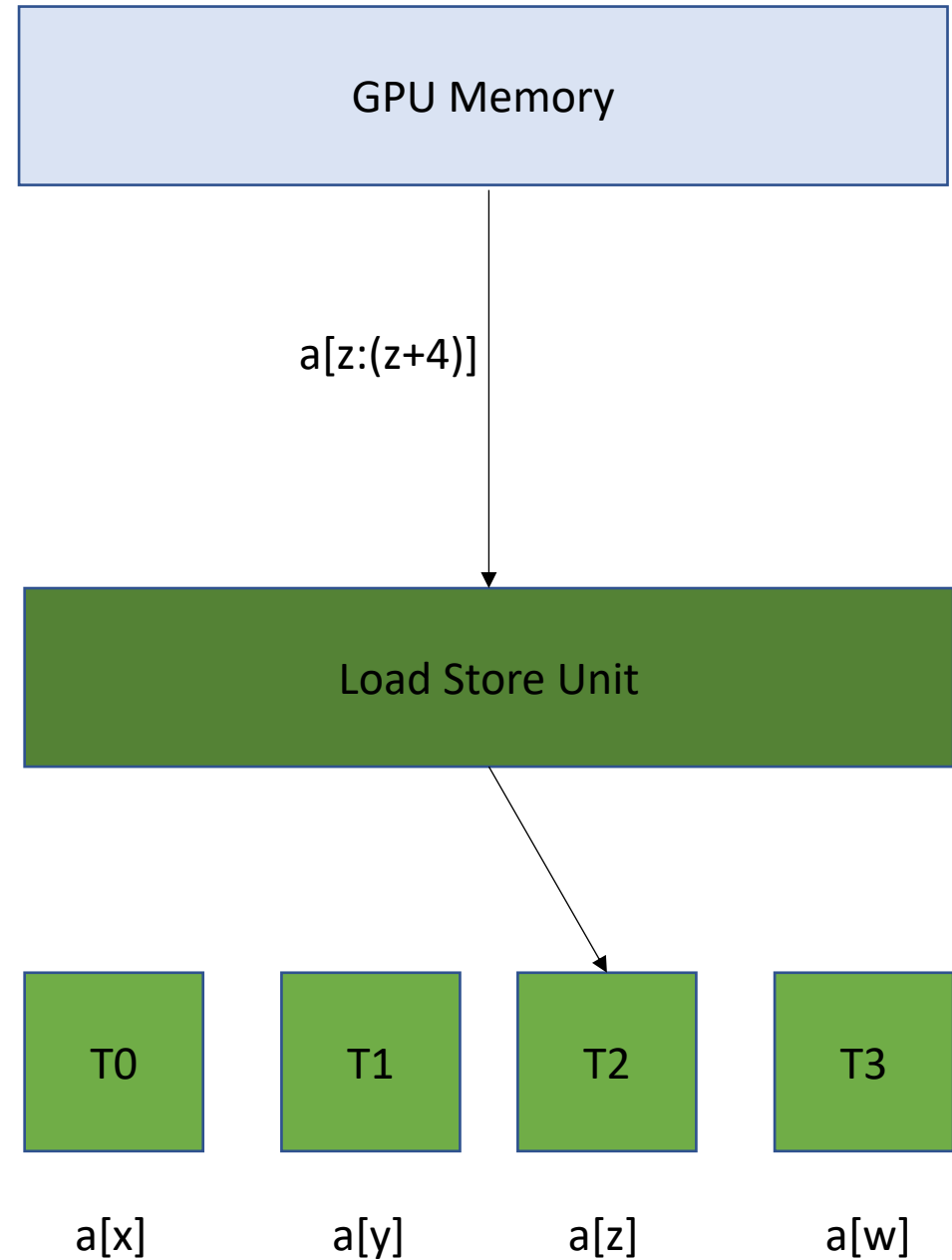
4 cores are accessing memory. What can happen

Read non-contiguous values

Not good!

Accesses are Serialized.

You need 4 requests to GPU memory



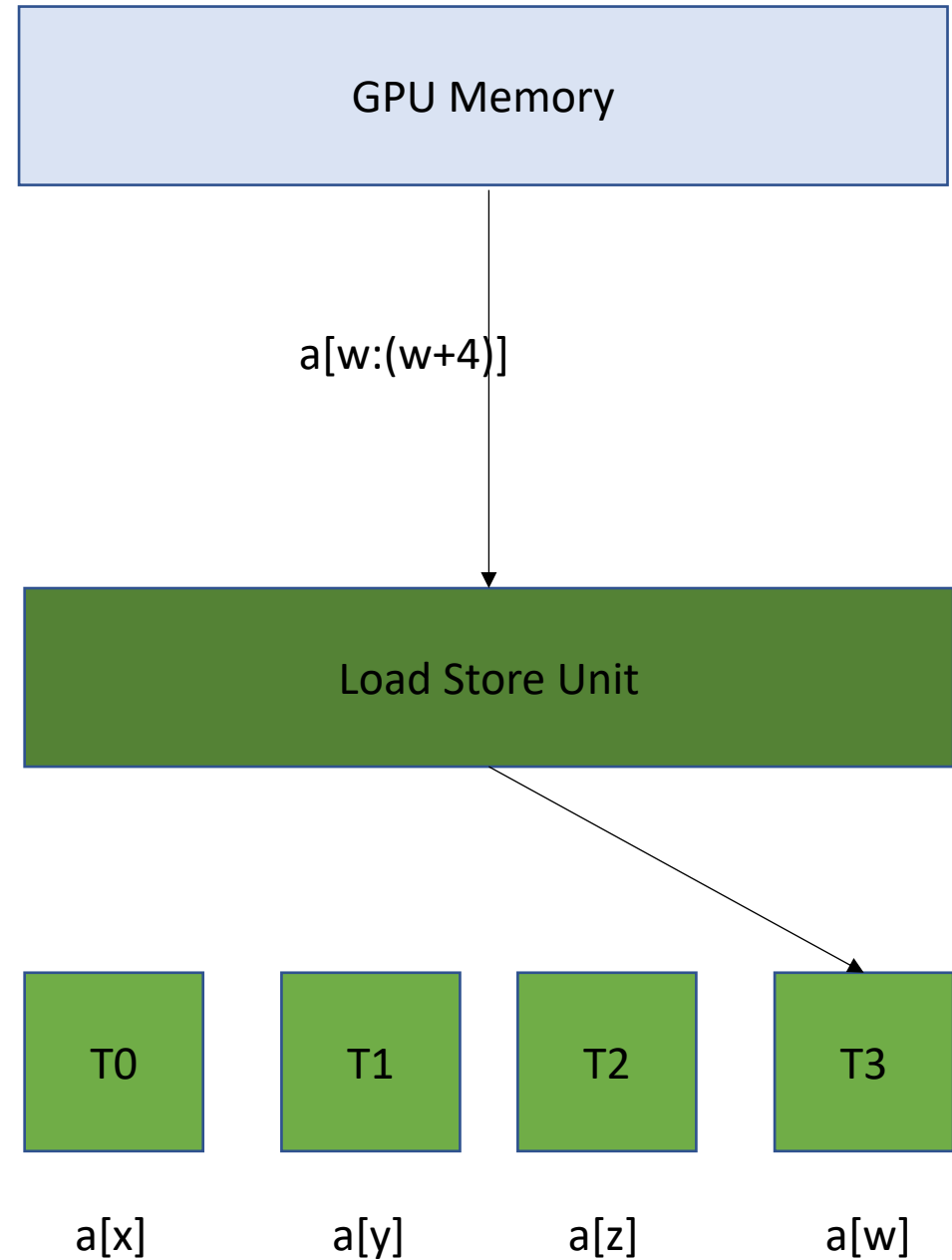
4 cores are accessing memory. What can happen

Read non-contiguous values

Not good!

Accesses are Serialized.

You need 4 requests to GPU memory



Go back to our program

```
__global__ void vector_add(int * d_a, int * d_b, int * d_c, int size) {  
    int chunk_size = size/blockDim.x;  
    int start = chunk_size * threadIdx.x;  
    int end = start + end;  
    for (int i = start; i < end; i++) {  
        d_a[i] = d_b[i] + d_c[i];  
    }  
}
```

calling the function

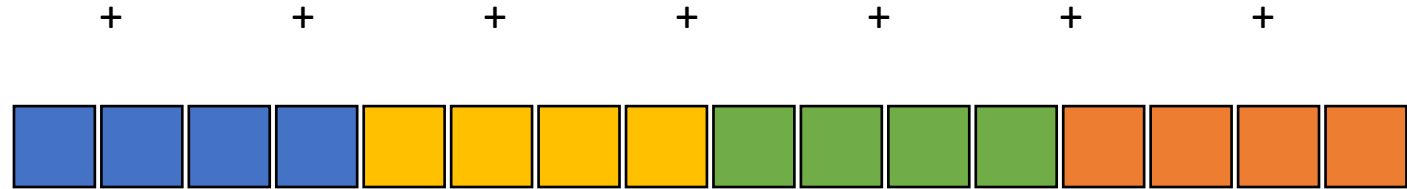
```
vector_add<<<1,32>>>(d_a, d_b, d_c, size);
```

Chunked Pattern

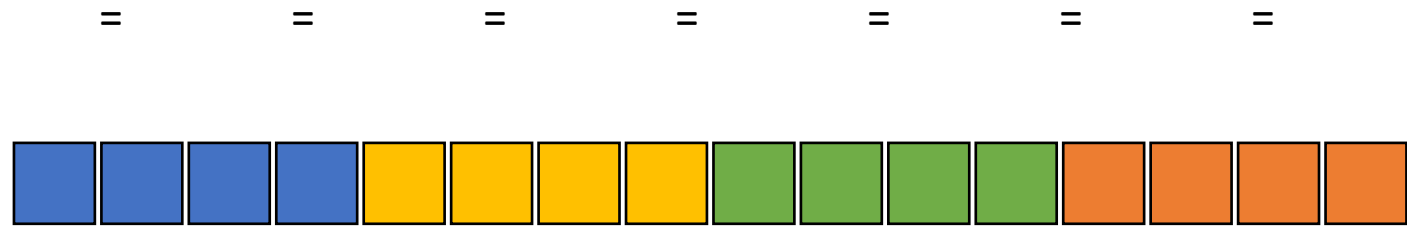
array a



array b



array c



Computation
can easily be
divided into
threads

- Thread 0 - Blue
- Thread 1 - Yellow
- Thread 2 - Green
- Thread 3 - Orange

Chunked Pattern

the first element accessed by the 4 threads sharing a load store unit. What sort of access is this?

array a



+ + + + + + +

array b



= = = = = = =

array c



Computation can easily be divided into threads

- Thread 0 - Blue
- Thread 1 - Yellow
- Thread 2 - Green
- Thread 3 - Orange

Chunked Pattern

the first element accessed by the 4 threads sharing a load store unit. What sort of access is this?

array a



+ + + + + + +

array b



= = = = = = =

array c



Computation can easily be divided into threads

- Thread 0 - Blue
- Thread 1 - Yellow
- Thread 2 - Green
- Thread 3 - Orange

How can we fix this

Stride Pattern

array a



+ + + + + + +

array b



= = = = = = =

array c



Computation
can easily be
divided into
threads

- Thread 0 - Blue
- Thread 1 - Yellow
- Thread 2 - Green
- Thread 3 - Orange

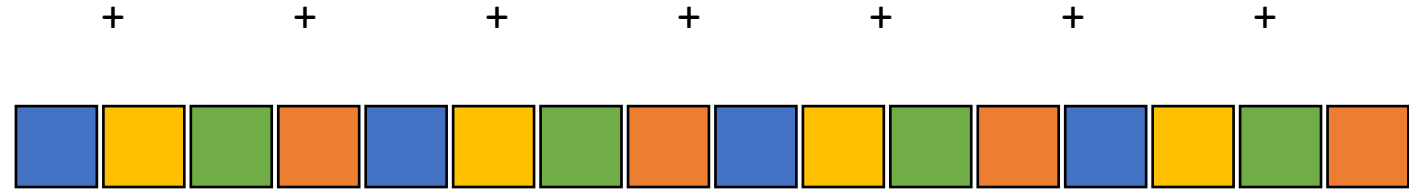
Stride Pattern

What sort of pattern is this?

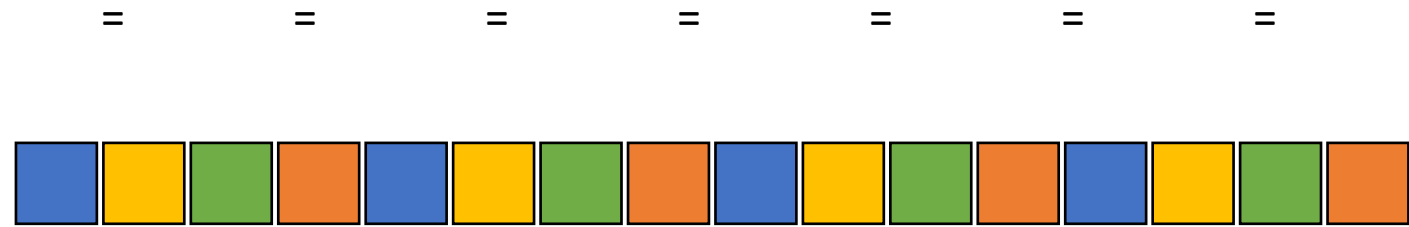
array a



array b



array c



Computation
can easily be
divided into
threads

- Thread 0 - Blue
- Thread 1 - Yellow
- Thread 2 - Green
- Thread 3 - Orange

Go back to our program

```
__global__ void vector_add(int * d_a, int * d_b, int * d_c, int size) {  
    int chunk_size = size/blockDim.x;  
    int start = chunk_size * threadIdx.x;  
    int end = start + end;  
    for (int i = start; i < end; i++) {  
        d_a[i] = d_b[i] + d_c[i];  
    }  
}
```

calling the function

```
vector_add<<<1,1024>>>(d_a, d_b, d_c, size);
```

Lets change this to a stride pattern

Go back to our program

```
__global__ void vector_add(int * d_a, int * d_b, int * d_c, int size) {  
    for (int i = threadIdx.x; i < size; i+=blockDim.x) {  
        d_a[i] = d_b[i] + d_c[i];  
    }  
}
```

calling the function

```
vector_add<<<1,1024>>>(d_a, d_b, d_c, size);
```

Coalesced memory accesses

Lets try it! What do we think?

Coalesced memory accesses

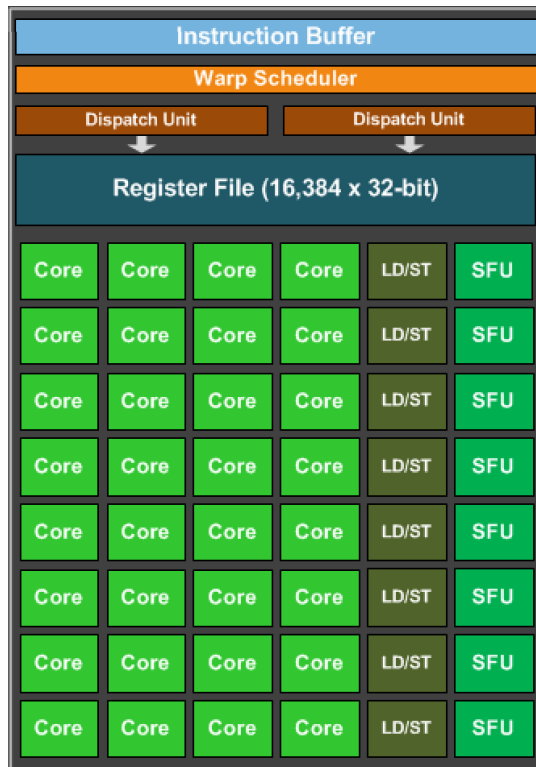
Lets try it! What do we think?



What else can we do?

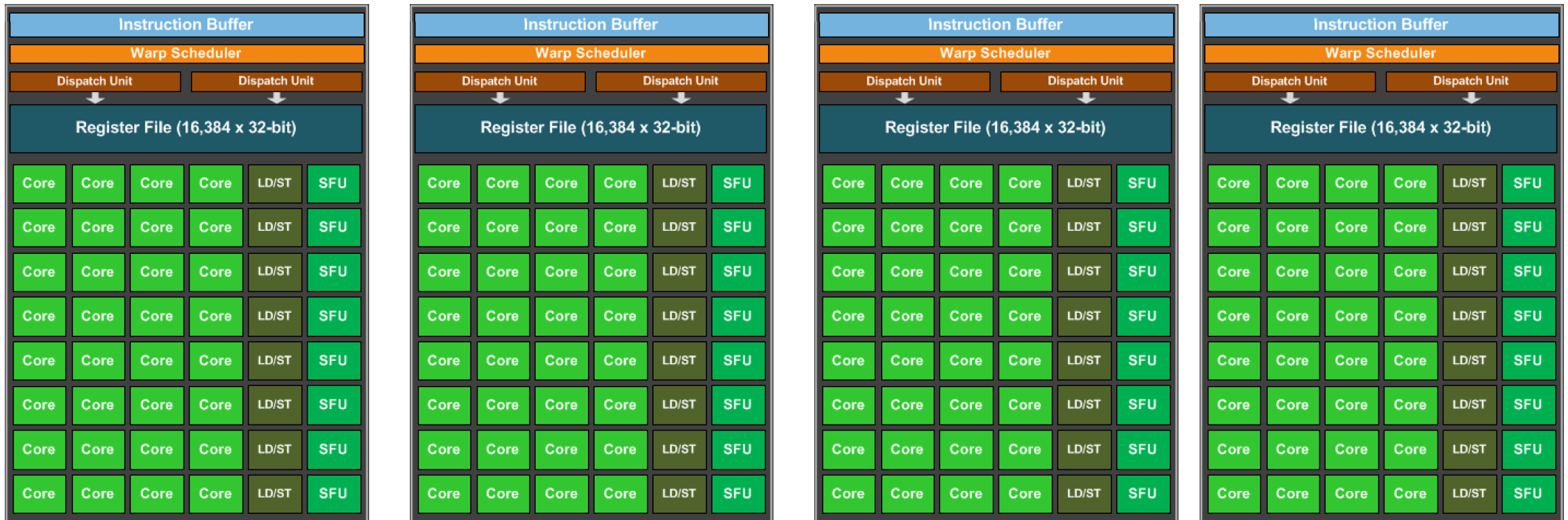
Multiple streaming multiprocessors

*We've been talking only about 1 streaming multiprocessor, most GPUs have multiple SMs
big ML GPUs have 32. My GPU has 4*



Multiple streaming multiprocessors

*We've been talking only about 1 streaming multiprocessor, most GPUs have multiple SMs
big ML GPUs have 32. This little GPU has 1*

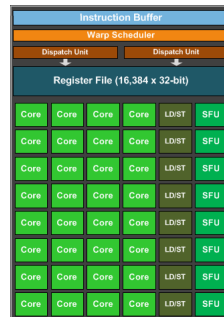
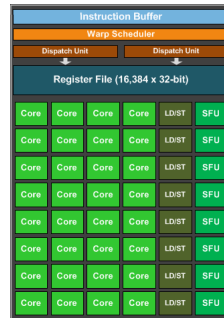
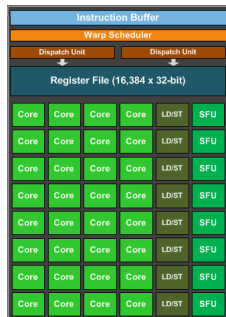
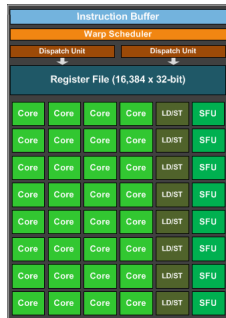
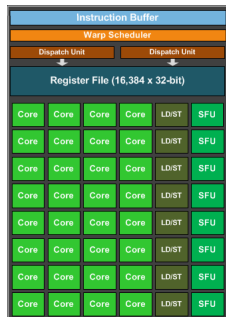


Multiple streaming multiprocessors

CUDA provides virtual streaming multiprocessors called **blocks**

Very efficient at launching and joining **blocks**.

No limit on blocks: launch as many as you need to map 1 thread to 1 data element



Go back to our program

```
__global__ void vector_add(int * d_a, int * d_b, int * d_c, int size) {  
    for (int i = threadIdx.x; i < size; i+=blockDim.x) {  
        d_a[i] = d_b[i] + d_c[i];  
    }  
}
```

calling the function

```
vector_add<<<1,1024>>>(d_a, d_b, d_c, size);
```

Launch with many thread blocks

Go back to our program

```
__global__ void vector_add(int * d_a, int * d_b, int * d_c, int size) {  
    int i = blockIdx.x * blockDim.x + threadIdx.x;  
    d_a[i] = d_b[i] + d_c[i];  
}
```

calling the function

```
vector_add<<<1024,1024>>>(d_a, d_b, d_c, size);
```

```
#define SIZE (1024*1024)
```

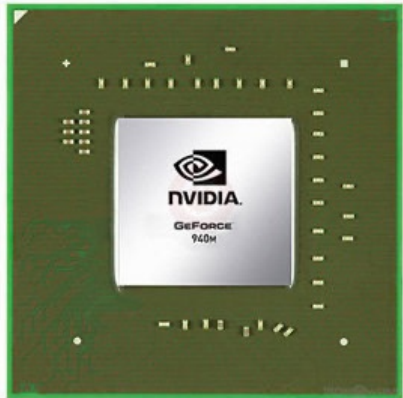
Need to recalculate some thread ids.

Launch with many thread blocks

Now we have 1 thread for each element

Final Round

Tiny GPU in an
embedded system



Nvidia Jetson Nano (whole chip, CPU + GPU)

2 Billion transistors

10 TDP

Est. \$99

<https://www.techpowerup.com/gpu-specs/geforce-940m.c2648>

https://www.alibaba.com/product-detail/Intel-Core-i7-9700K-8-Cores_62512430487.html

<https://www.prolast.com/prolast-elevated-boxing-rings-22-x-22/>

Fight!



The CPU in
my professor
workstation



Intel i7-9700K

2.16 Billion transistors

95 TDP

Est. \$316

See you on Wednesday

- Turn in HW 4 if you haven't already
- Working on GPU programming!