

Introduction to GPU Computing

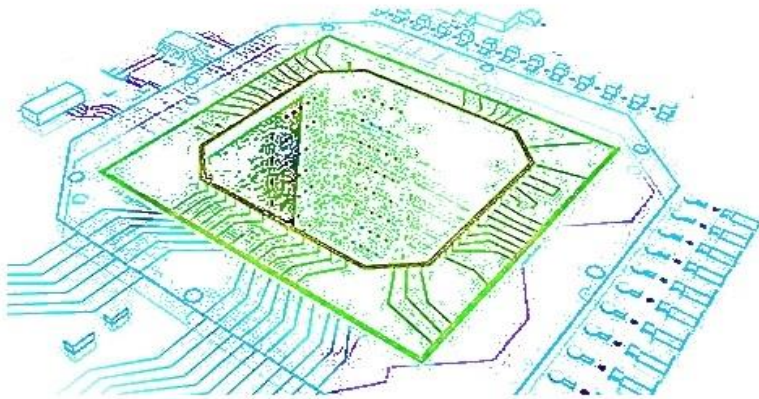
Adriana Cavada

June, 2016



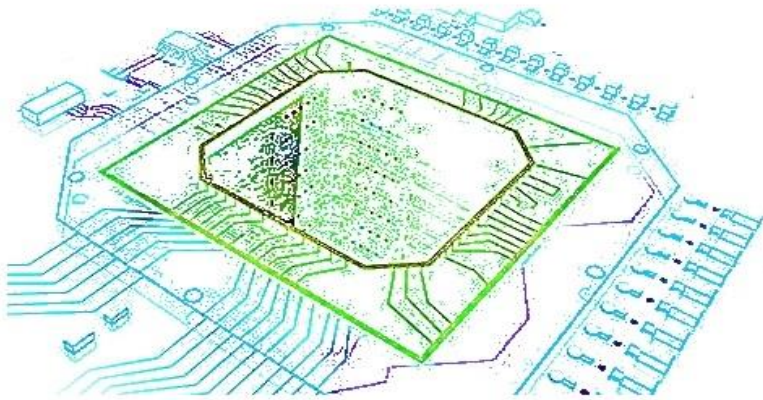
What is GPU?

GPU = Graphics Processor Unit



Computer chip mainly designed for image processing and real-time rendering

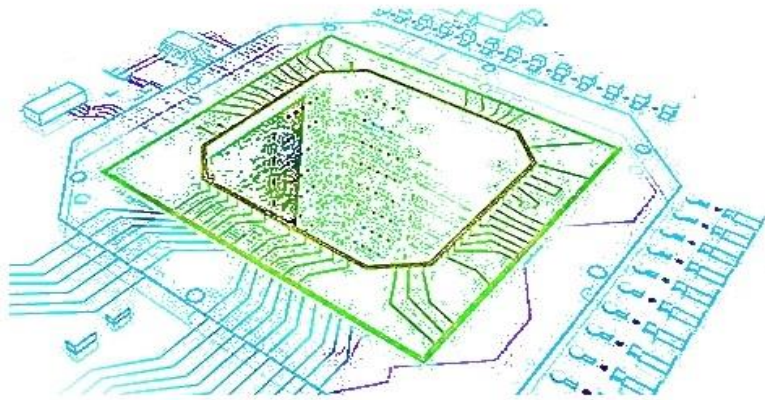
GPU = Graphics Processor Unit



Parallel processor mainly designed for image processing and real-time rendering



GPU = Graphics Processor Unit



Parallel processor mainly designed for image processing and real-time rendering

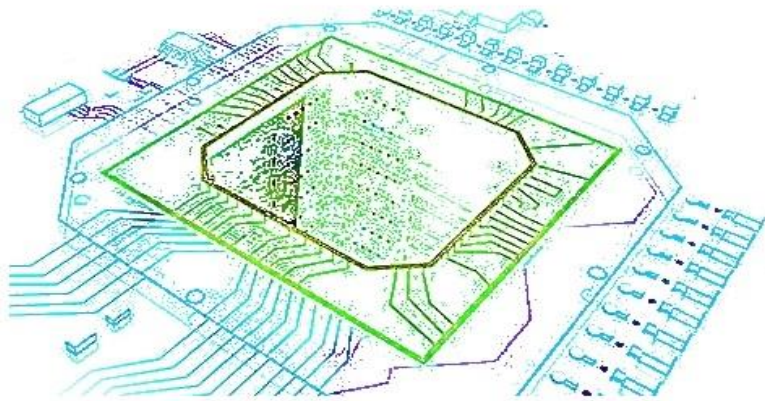
“THE WORLD'S FIRST GPU”



NVIDIA GeForce 256
(1999)



GPU = Graphics Processor Unit



Parallel processor mainly designed for image processing and real-time rendering

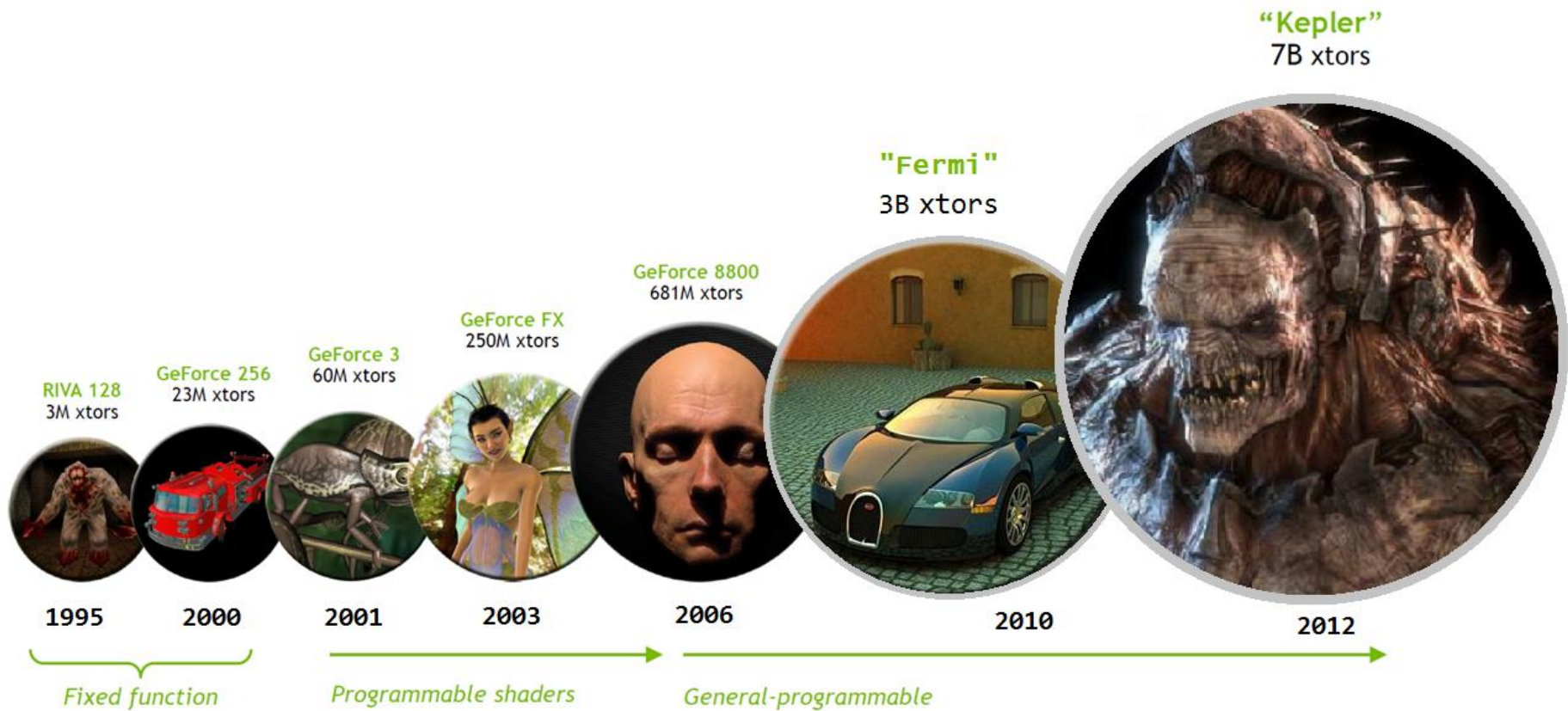
“THE WORLD'S FIRST GPU”



NVIDIA GeForce 256
(1999)



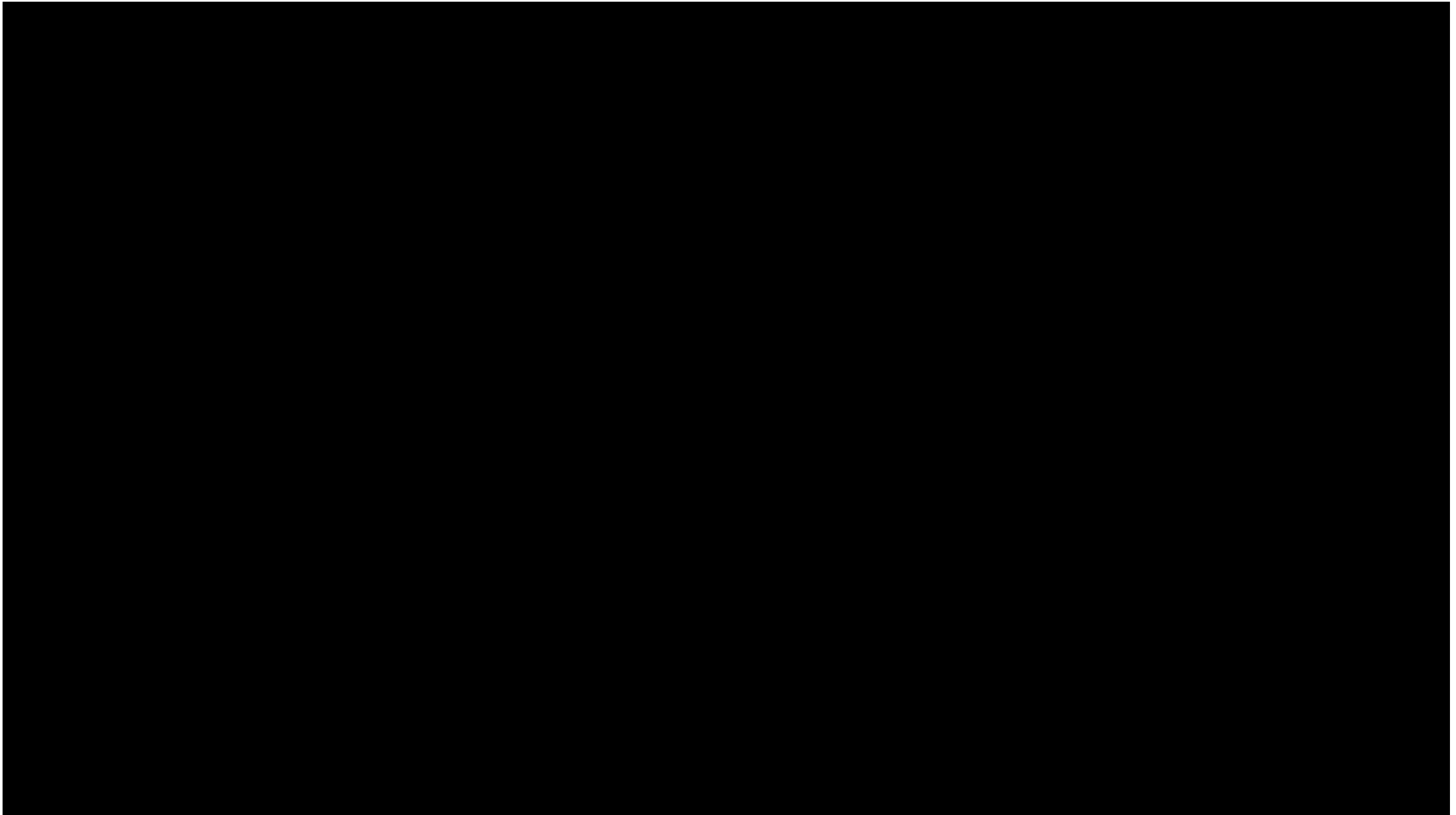
Evolution



Ref: NVIDIA Corporation, *GPU Hardware and Programming Models*

CPU vs GPU

CPU vs GPU



The Mythbusters, Adam Savage and Jamie Hyneman demonstrate the power of GPUs vs CPUs - NVIDIA Corporation

What's better?



What's better?



What's better?



Deliver a package as soon as possible?

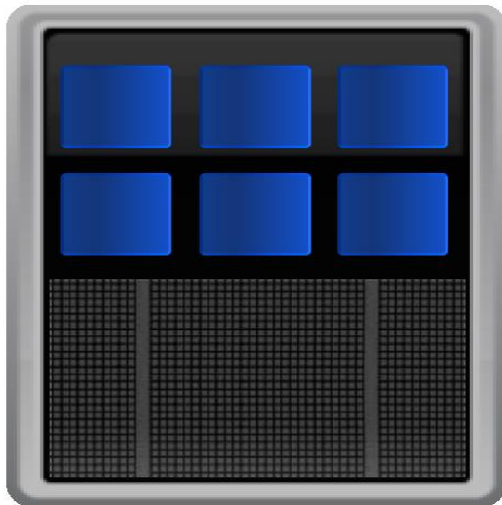


Or deliver many packages
within a reasonable
timescale?

CPU vs GPU

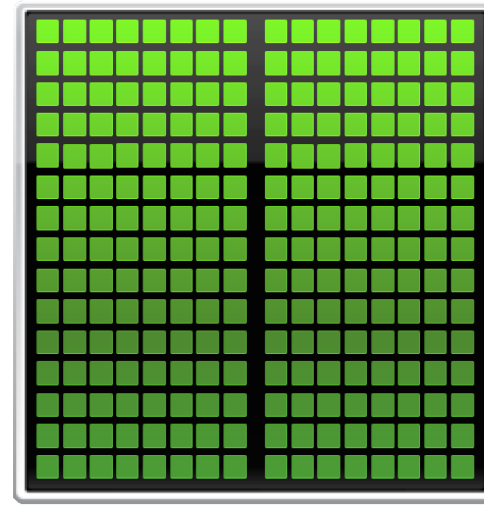
Compute a job as fast
as possible

CPU



Compute many jobs within
a reasonable timeframe

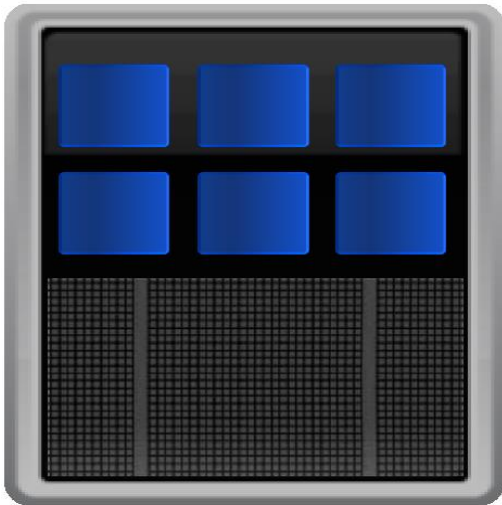
GPU



CPU vs GPU

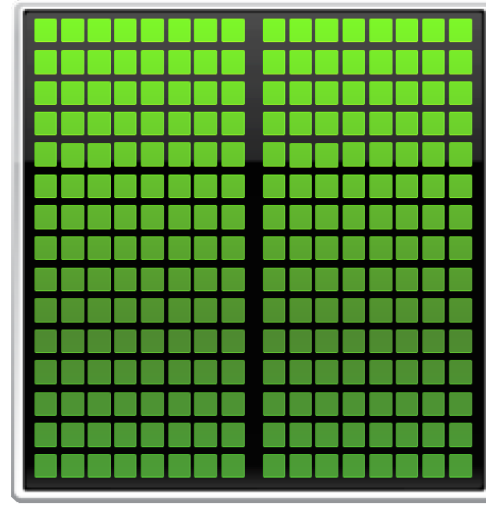
Low Latency

CPU



High Throughput

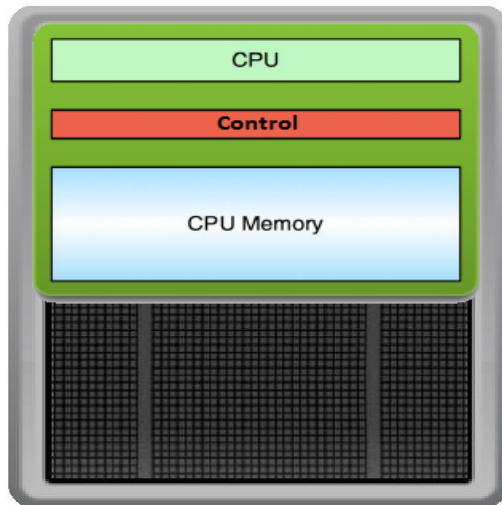
GPU



CPU vs GPU

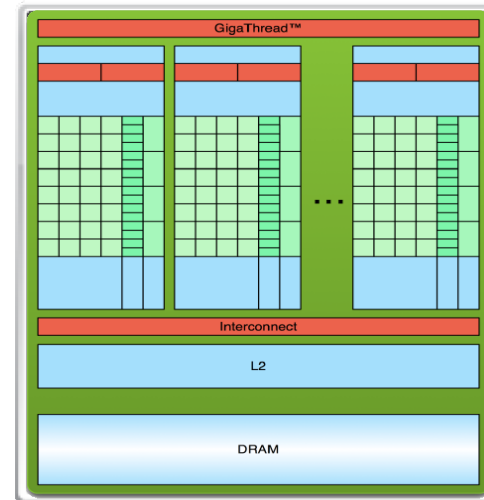
Low Latency

CPU

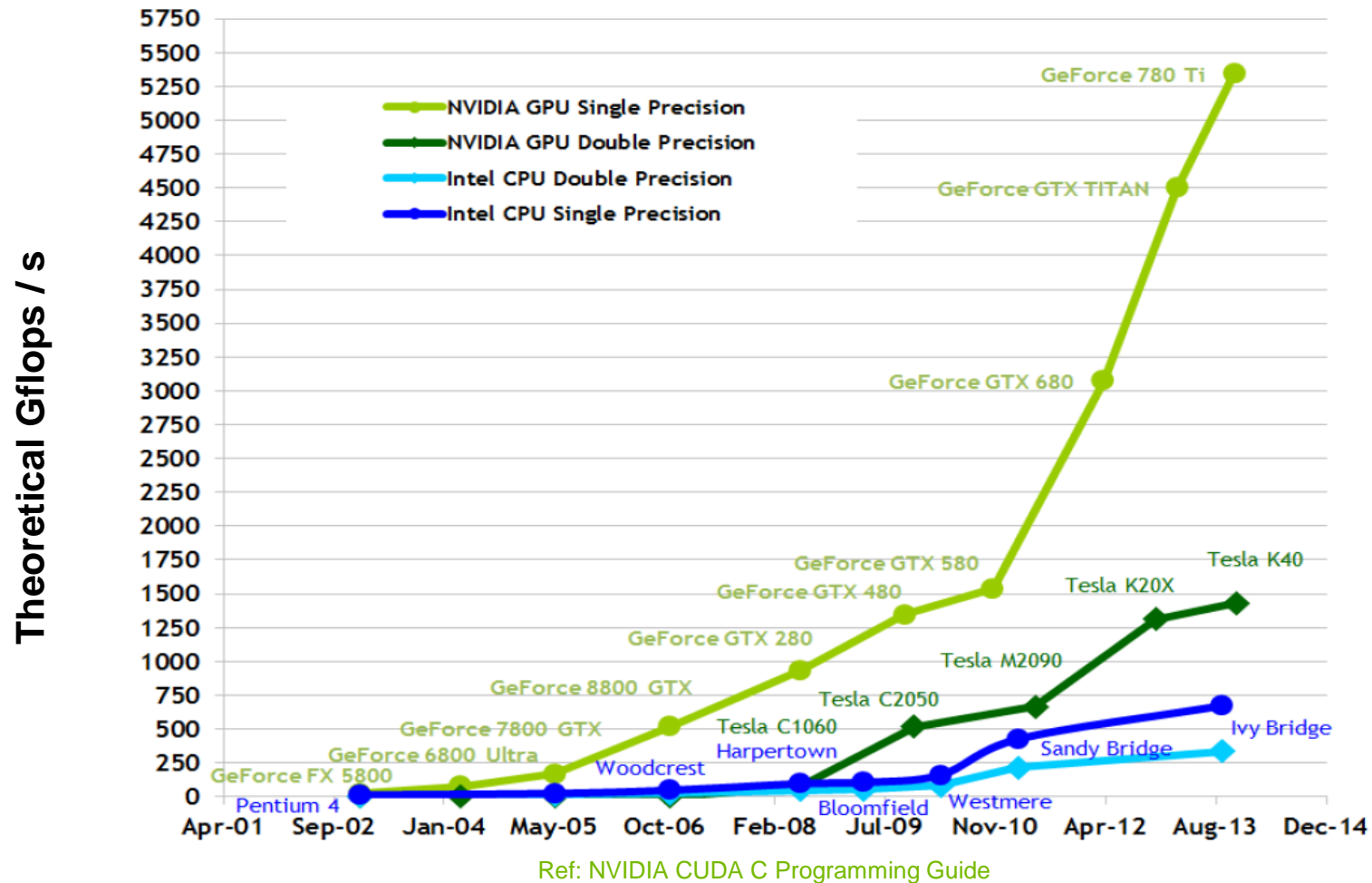


High Throughput

GPU

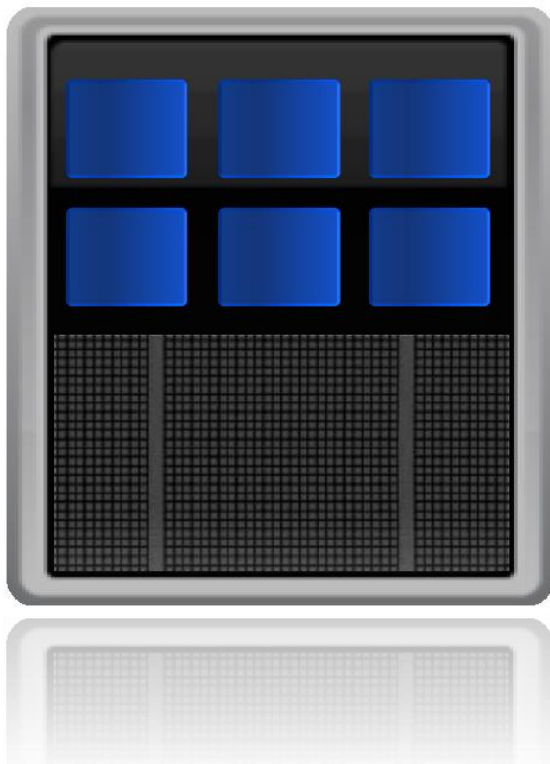


Floating Point Operations

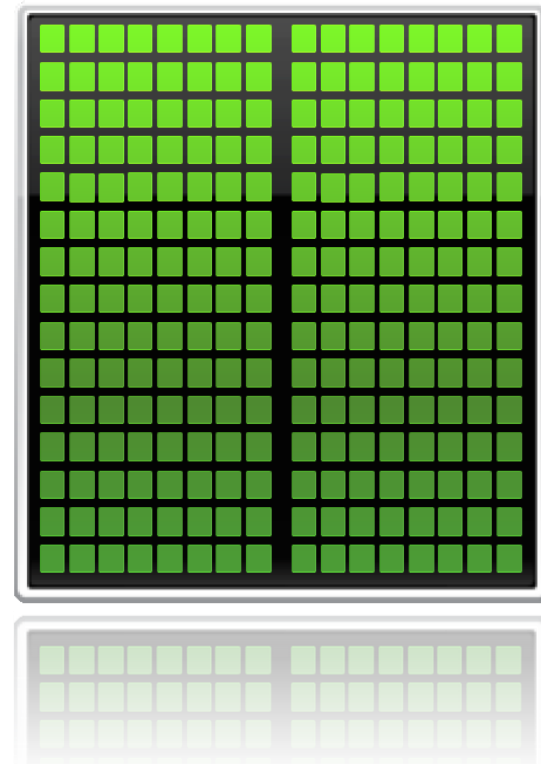


Better adding tan choosing

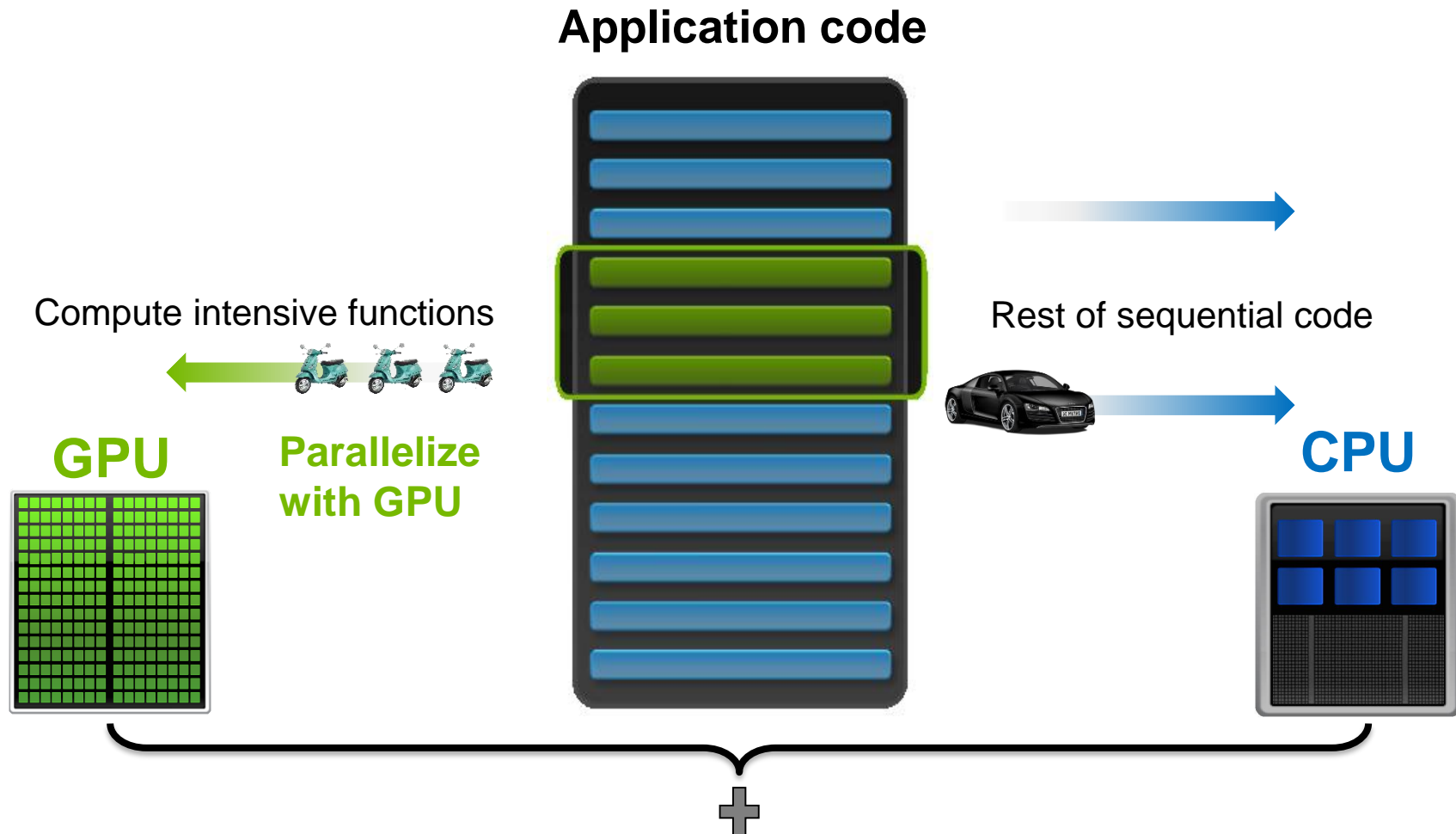
CPU



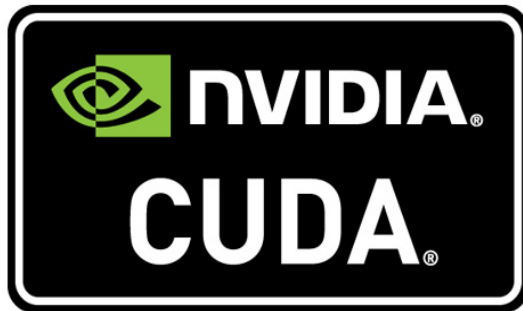
GPU



Adding GPUs

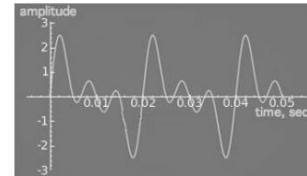


How starting?

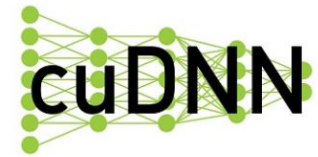


What is CUDA?

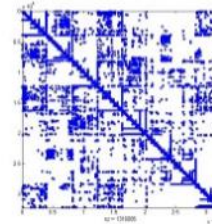
- Parallel computing platform
- Programming model
- CUDA uses GPUs for general purpose
- C, C++, Fortran...
- A lot of libraries



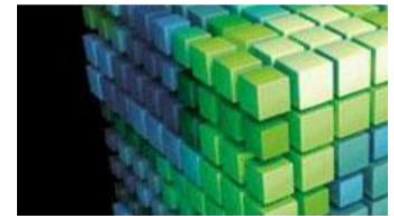
cuFFT



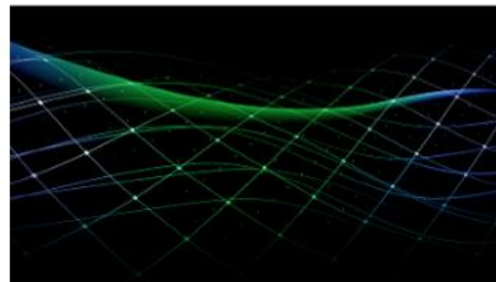
cuDNN



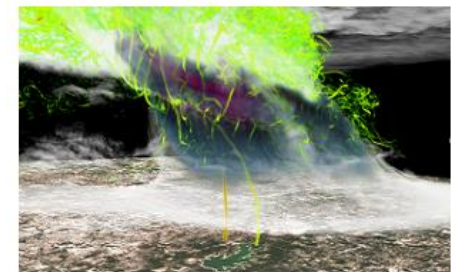
cuSPARSE



cuBLAS



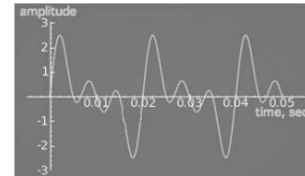
cuSOLVER



Index Framework

What is CUDA?

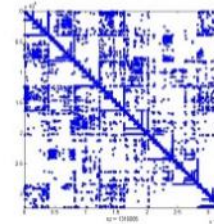
- Parallel computing platform
- Programming model
- CUDA uses GPUs for general purpose
- C, C++, Fortran...
- A lot of libraries



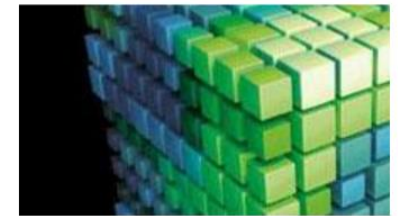
cuFFT



cuDNN



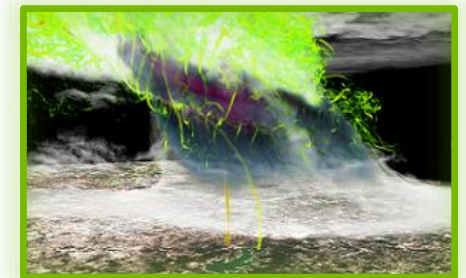
cuSPARSE



cuBLAS



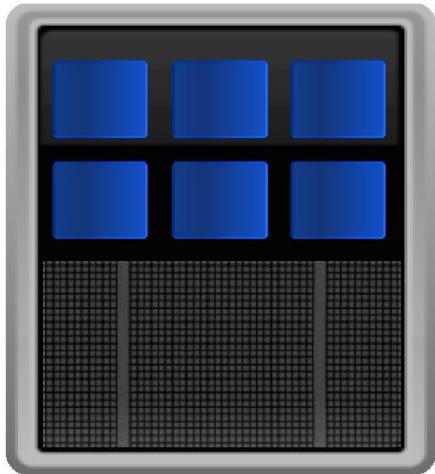
cuSOLVER



Index Framework

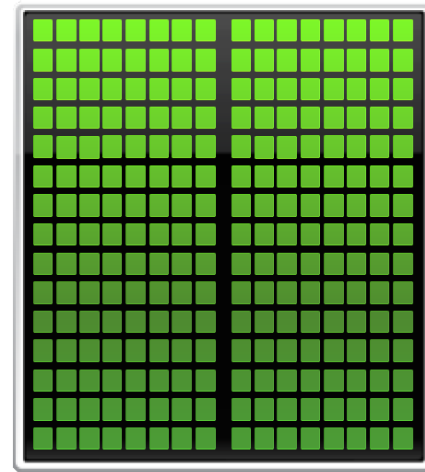
Terminology

Host



CPU + host memory

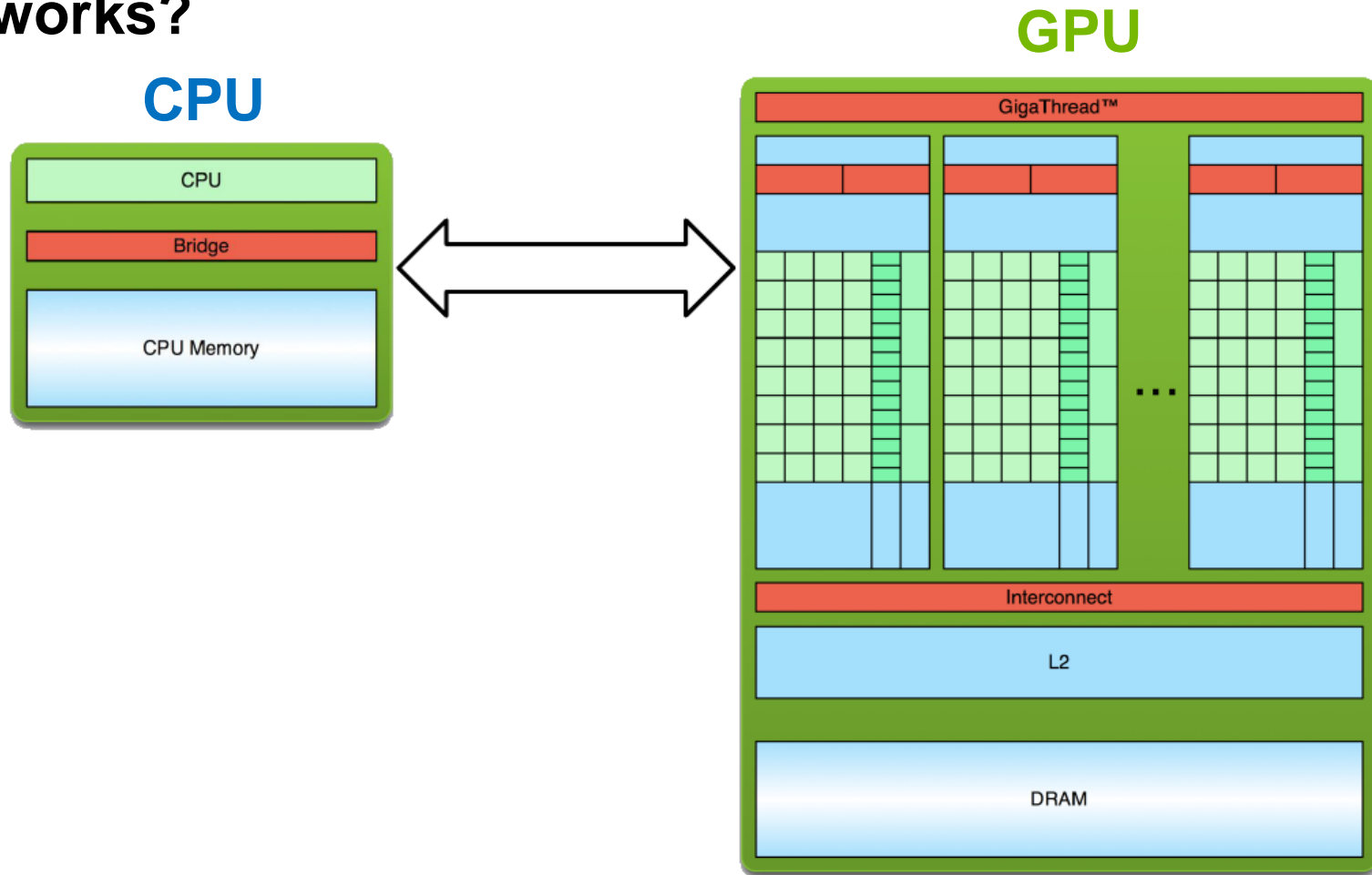
Device



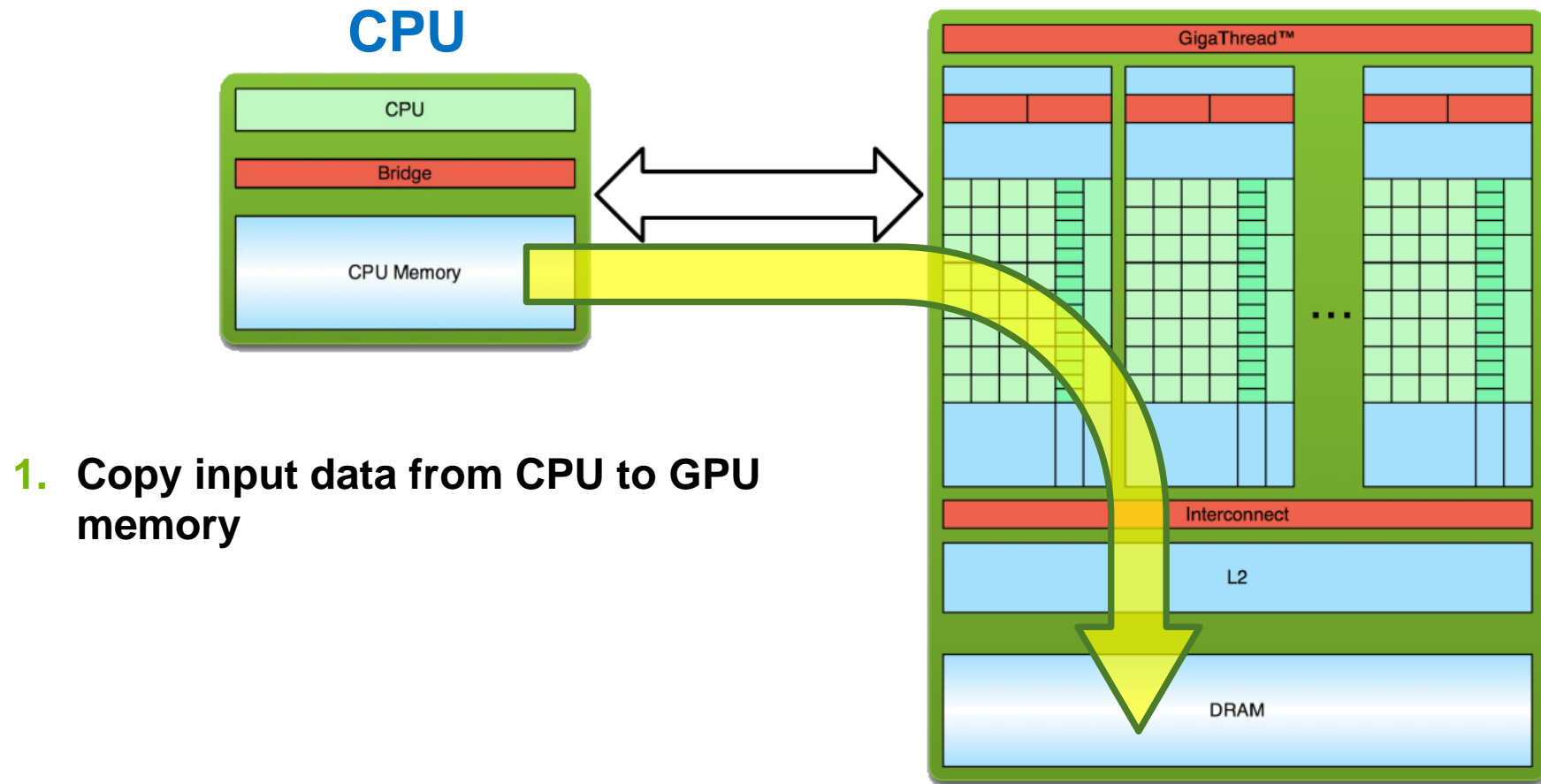
GPU + device memory

How it works?

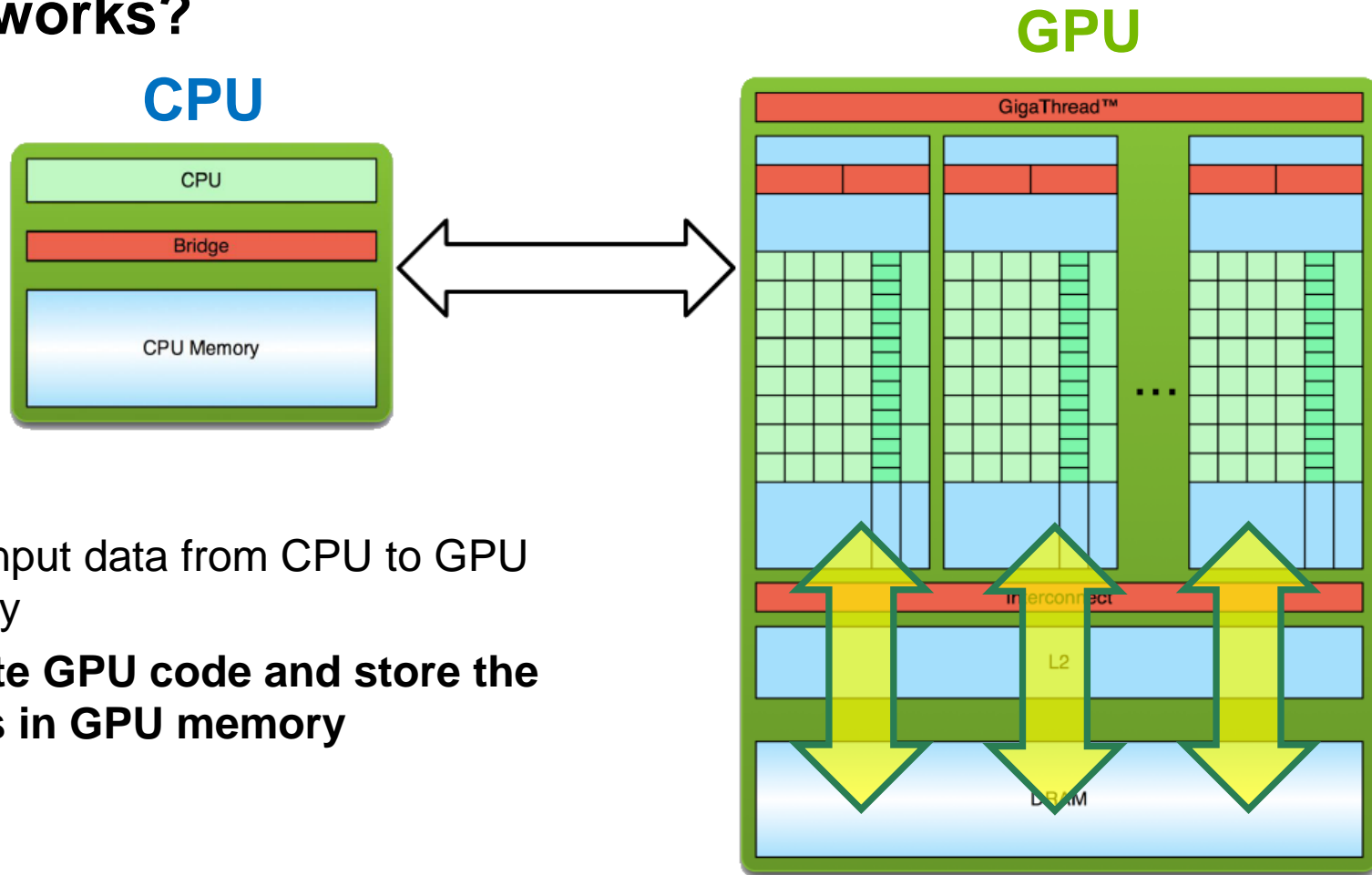
How it works?



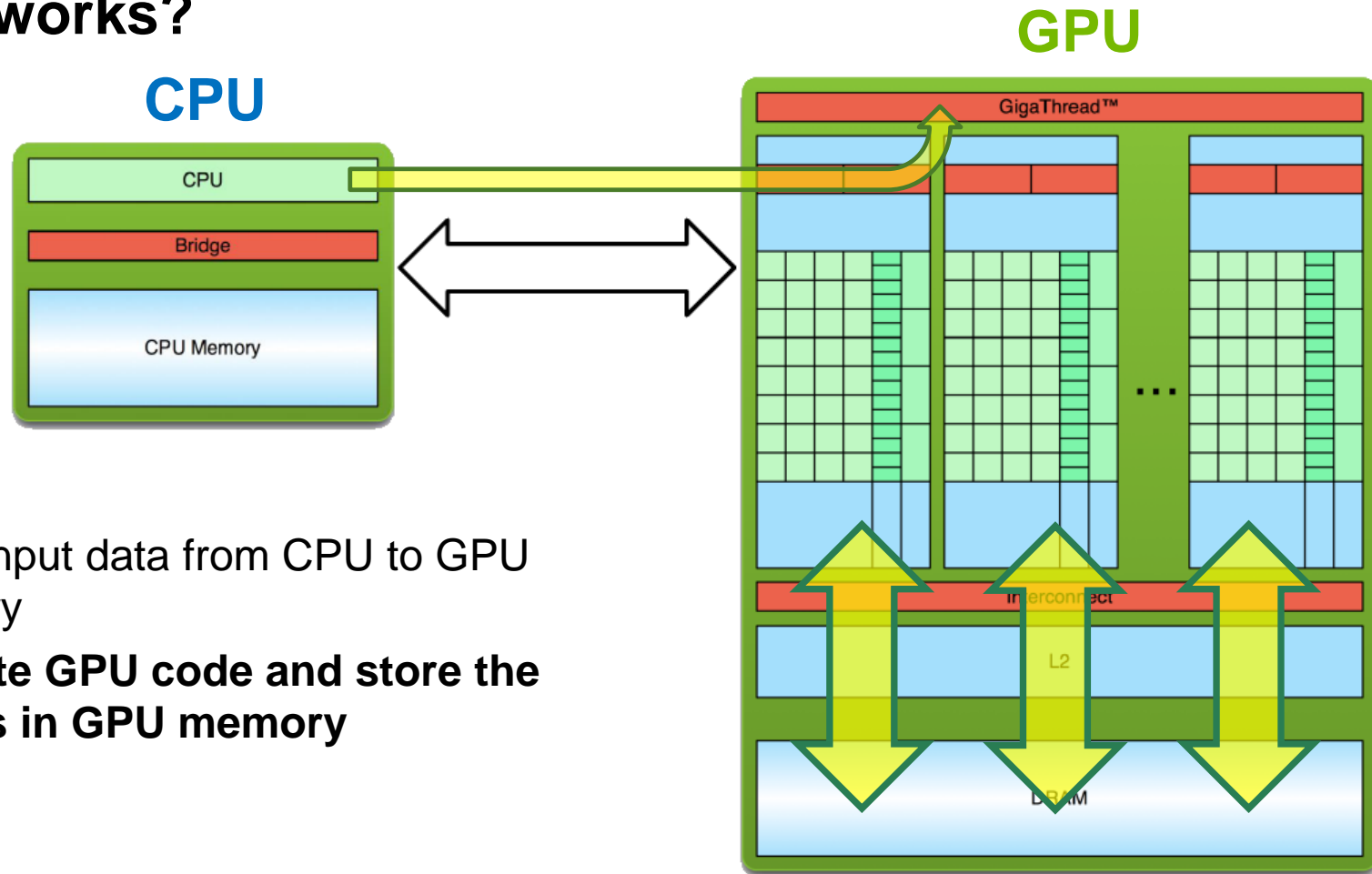
How it works?



How it works?

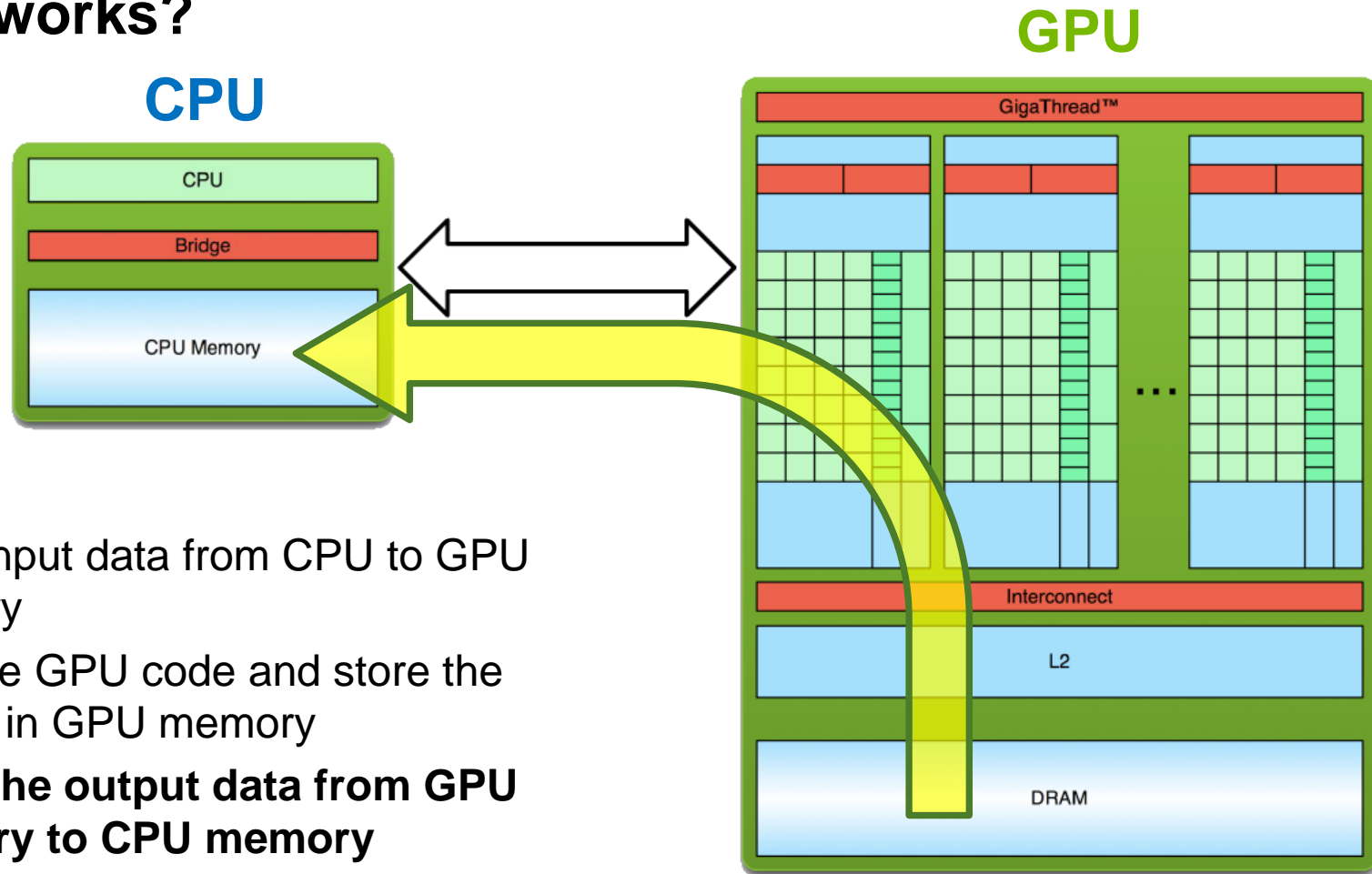


How it works?



1. Copy input data from CPU to GPU memory
2. **Execute GPU code and store the results in GPU memory**

How it works?



Hello World!

Hello world!

Standard C code

```
int main(void)
{
    printf("Hello World!\n");
    return 0;
}
```

Run on the host

NVIDIA compiler (nvcc) can compile programs with no device

Hello world!

Standard C code

```
int main(void)
{
    printf("Hello World!\n");
    return 0;
}
```

C with CUDA

```
__global__ void mykernel(void)
{
    int main(void)
    {
        mykernel<<<1,1>>>>();
        printf("Hello World!\n");
        return 0;
    }
}
```

Two new elements...

Hello world!

C with CUDA

```
__global__ void mykernel(void)
{
    int main(void)
    {
        mykernel<<<1,1>>>();
        printf("Hello World!\n");
        return 0;
    }
}
```

Two new elements...

Hello world!

- Keyword `__global__` indicates that the function is a device function. It is called from host code
- NVIDIA compiler separates code into host and device:
 - **Device functions** are processed by `nvcc`
 - **Host functions** are processed by standard compiler like `gcc`

C with CUDA

```
__global__ void mykernel(void)
{
    int main(void)
    {
        mykernel<<<1,1>>>();
        printf("Hello World!\n");
        return 0;
    }
}
```

Two new elements...

Hello world!

C with CUDA

```
__global__ void mykernel(void)
{
    mykernel(<math>d</math>);
    {
        mykernel<<<1,1>>>>();
        printf("Hello World!\n");
        return 0;
    }
}
```

Two new elements...

Hello world!

- Triple angle brackets are a call from host code to device code
- Also called a “**kernel launch**”
- The kernel *mykernel()* does nothing, but a device function is needed

C with CUDA

```
__global__ void mykernel(void)
{
    // ...
}

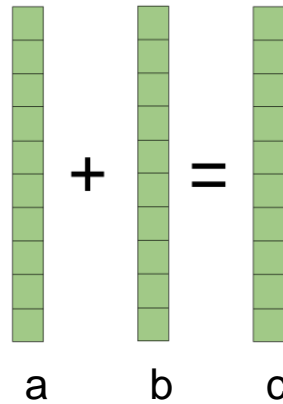
int main()
{
    mykernel<<<1,1>>>>();
    printf("Hello World!\n");
    return 0;
}
```

Two new elements...

but massive parallelism?

We need a more interesting example...

Vector addition



Integer addition

- The kernel to add two integers

```
__global__ void add(int *a, int *b, int *c)
{
    *c = *a + *b;
}
```

Integer addition

- The kernel to add two integers

```
__global__ void add(int *a, int *b, int *c)
{
    *c = *a + *b;
}
```

- As before the keyword `__global__` means that
 - `add()` is a device function
 - `add()` will be called from the host

Integer addition

- The kernel to add two integers

```
__global__ void add(int *a, int *b, int *c)
{
    *c = *a + *b;
}
```

- As before the keyword `__global__` means that

- `add()` is a device function
- `add()` will be called from the host

- We use device pointers `*a`, `*b`, `*c` point to GPU memory

**Pointers
reminder?**

Integer addition

- The kernel to add two integers

```
__global__ void add(int *a, int *b, int *c)
{
    *c = *a + *b;
}
```

- As before the keyword `__global__` means that

- `add()` is a device function
- `add()` will be called from the host

- We use device pointers `*a`, `*b`, `*c` point to GPU memory

- CUDA functions to handle device memory


- `cudaMalloc()`, `cudaFree()`, `cudaMemcpy()`

Pointers
reminder?

Integer addition

- Let's look at the main() function

Host copies of a, b, c



```
int main(void)
{
    int a, b, c;
    int *d_a, *d_b, *d_c;
    int size = sizeof(int);

    cudaMalloc((void **)&d_a, size);
    cudaMalloc((void **)&d_b, size);
    cudaMalloc((void **)&d_c, size);

    a = 2;
    b = 7;

    cudaMemcpy(d_a, &a, size, cudaMemcpyHostToDevice);
    cudaMemcpy(d_b, &b, size, cudaMemcpyHostToDevice);

    add<<<1,1>>>(d_a, d_b, d_c);

    cudaMemcpy(&c, d_c, size, cudaMemcpyDeviceToHost);


    cudaFree(d_a); cudaFree(d_b); cudaFree(d_c);
    return 0;
}
```

Integer addition

■ Let's look at the main() function

Host copies of a, b, c

Device copies of a, b, c



```
int main(void)
{
    int a, b, c;
    int *d_a, *d_b, *d_c;
    int size = sizeof(int);

    cudaMalloc((void **)&d_a, size);
    cudaMalloc((void **)&d_b, size);
    cudaMalloc((void **)&d_c, size);

    a = 2;
    b = 7;

    cudaMemcpy(d_a, &a, size, cudaMemcpyHostToDevice);
    cudaMemcpy(d_b, &b, size, cudaMemcpyHostToDevice);

    add<<<1,1>>>>(d_a, d_b, d_c);

    cudaMemcpy(&c, d_c, size, cudaMemcpyDeviceToHost);

    cudaFree(d_a); cudaFree(d_b); cudaFree(d_c);
    return 0;
}
```


Integer addition

■ Let's look at the main() function

Host copies of a, b, c

Device copies of a, b, c

Allocate space in
device for a, b, c



```
int main(void)
{
    int a, b, c;
    int *d_a, *d_b, *d_c;
    int size = sizeof(int);

    cudaMalloc((void **)&d_a, size);
    cudaMalloc((void **)&d_b, size);
    cudaMalloc((void **)&d_c, size);

    a = 2;
    b = 7;

    cudaMemcpy(d_a, &a, size, cudaMemcpyHostToDevice);
    cudaMemcpy(d_b, &b, size, cudaMemcpyHostToDevice);

    add<<<1,1>>>(d_a, d_b, d_c);

    cudaMemcpy(&c, d_c, size, cudaMemcpyDeviceToHost);

    cudaFree(d_a); cudaFree(d_b); cudaFree(d_c);
    return 0;
}
```

Integer addition

■ Let's look at the main() function

Host copies of a, b, c

Device copies of a, b, c

Allocate space in
device for a, b, c

Inputs

```
int main(void)
{
    int a, b, c;
    int *d_a, *d_b, *d_c;
    int size = sizeof(int);

    cudaMalloc((void **)&d_a, size);
    cudaMalloc((void **)&d_b, size);
    cudaMalloc((void **)&d_c, size);

    a = 2;
    b = 7;

    cudaMemcpy(d_a, &a, size, cudaMemcpyHostToDevice);
    cudaMemcpy(d_b, &b, size, cudaMemcpyHostToDevice);

    add<<<1,1>>>(d_a, d_b, d_c);

    cudaMemcpy(&c, d_c, size, cudaMemcpyDeviceToHost);

    cudaFree(d_a); cudaFree(d_b); cudaFree(d_c);
    return 0;
}
```

Integer addition

■ Let's look at the main() function

Host copies of a, b, c

Device copies of a, b, c

Allocate space in
device for a, b, c

Inputs

Copy inputs to device

```
int main(void)
{
    int a, b, c;
    int *d_a, *d_b, *d_c;
    int size = sizeof(int);

    cudaMalloc((void **)&d_a, size);
    cudaMalloc((void **)&d_b, size);
    cudaMalloc((void **)&d_c, size);

    a = 2;
    b = 7;

    cudaMemcpy(d_a, &a, size, cudaMemcpyHostToDevice);
    cudaMemcpy(d_b, &b, size, cudaMemcpyHostToDevice);

    add<<<1,1>>>(d_a, d_b, d_c);

    cudaMemcpy(&c, d_c, size, cudaMemcpyDeviceToHost);

    cudaFree(d_a); cudaFree(d_b); cudaFree(d_c);
    return 0;
}
```

Integer addition

■ Let's look at the main() function

Host copies of a, b, c

Device copies of a, b, c

Allocate space in
device for a, b, c

Inputs

Copy inputs to device

Launch add() kernel on GPU

```
int main(void)
{
    int a, b, c;
    int *d_a, *d_b, *d_c;
    int size = sizeof(int);

    cudaMalloc((void **)&d_a, size);
    cudaMalloc((void **)&d_b, size);
    cudaMalloc((void **)&d_c, size);

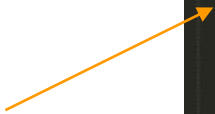
    a = 2;
    b = 7;

    cudaMemcpy(d_a, &a, size, cudaMemcpyHostToDevice);
    cudaMemcpy(d_b, &b, size, cudaMemcpyHostToDevice);

    add<<<1,1>>>>(d_a, d_b, d_c);

    cudaMemcpy(&c, d_c, size, cudaMemcpyDeviceToHost);

    cudaFree(d_a); cudaFree(d_b); cudaFree(d_c);
    return 0;
}
```



Integer addition

■ Let's look at the main() function

Host copies of a, b, c

Device copies of a, b, c

Allocate space in
device for a, b, c

Inputs

Copy inputs to device

Launch add() kernel on GPU

Copy result to host

```
int main(void)
{
    int a, b, c;
    int *d_a, *d_b, *d_c;
    int size = sizeof(int);

    cudaMalloc((void **)&d_a, size);
    cudaMalloc((void **)&d_b, size);
    cudaMalloc((void **)&d_c, size);

    a = 2;
    b = 7;

    cudaMemcpy(d_a, &a, size, cudaMemcpyHostToDevice);
    cudaMemcpy(d_b, &b, size, cudaMemcpyHostToDevice);

    add<<<1,1>>>(d_a, d_b, d_c);

    cudaMemcpy(&c, d_c, size, cudaMemcpyDeviceToHost);

    cudaFree(d_a); cudaFree(d_b); cudaFree(d_c);
    return 0;
}
```

**but massive parallelism?
vector addition?**

Let's running in parallel...

Moving to parallel

- Instead of executing `add()` once, execute N times in parallel

`add<<<1,1>>>();`

```
int main(void)
{
    int a, b, c;
    int *d_a, *d_b, *d_c;
    int size = sizeof(int);

    cudaMalloc((void **)&d_a, size);
    cudaMalloc((void **)&d_b, size);
    cudaMalloc((void **)&d_c, size);

    a = 2;
    b = 7;

    cudaMemcpy(d_a, &a, size, cudaMemcpyHostToDevice);
    cudaMemcpy(d_b, &b, size, cudaMemcpyHostToDevice);

    add<<<1,1>>>(d_a, d_b, d_c);

    cudaMemcpy(&c, d_c, size, cudaMemcpyDeviceToHost);

    cudaFree(d_a); cudaFree(d_b); cudaFree(d_c);
    return 0;
}
```

Moving to parallel

- Instead of executing `add()` once, execute `N` times in parallel

`add<<<1,1>>>();`



`add<<<N,1>>>();`

```
int main(void)
{
    int a, b, c;
    int *d_a, *d_b, *d_c;
    int size = sizeof(int);

    cudaMalloc((void **)&d_a, size);
    cudaMalloc((void **)&d_b, size);
    cudaMalloc((void **)&d_c, size);

    a = 2;
    b = 7;

    cudaMemcpy(d_a, &a, size, cudaMemcpyHostToDevice);
    cudaMemcpy(d_b, &b, size, cudaMemcpyHostToDevice);

    add<<<1,1>>>(d_a, d_b, d_c);

    cudaMemcpy(&c, d_c, size, cudaMemcpyDeviceToHost);

    cudaFree(d_a); cudaFree(d_b); cudaFree(d_c);
    return 0;
}
```

Moving to parallel

- Instead of executing `add()` once, execute `N` times in parallel

```
add<<<1,1>>>();
```



```
add<<<N,1>>>();
```

- Each parallel invocation of `add()` is referred to as a **block**

```
int main(void)
{
    int a, b, c;
    int *d_a, *d_b, *d_c;
    int size = sizeof(int);

    cudaMalloc((void **)&d_a, size);
    cudaMalloc((void **)&d_b, size);
    cudaMalloc((void **)&d_c, size);

    a = 2;
    b = 7;

    cudaMemcpy(d_a, &a, size, cudaMemcpyHostToDevice);
    cudaMemcpy(d_b, &b, size, cudaMemcpyHostToDevice);

    add<<<1,1>>>(d_a, d_b, d_c);

    cudaMemcpy(&c, d_c, size, cudaMemcpyDeviceToHost);

    cudaFree(d_a); cudaFree(d_b); cudaFree(d_c);
    return 0;
}
```

Moving to parallel

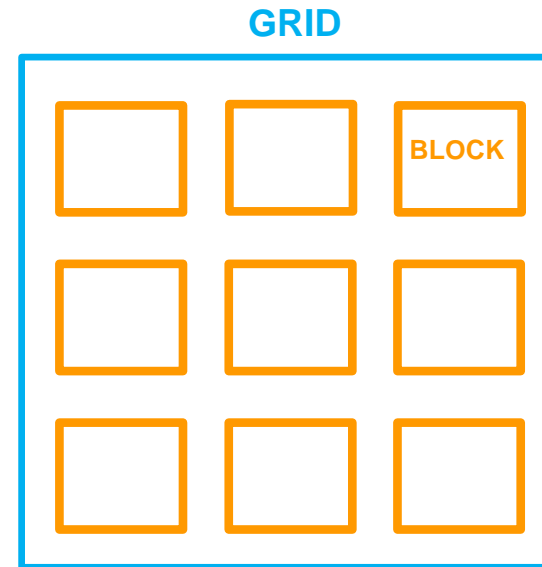
- Instead of executing `add()` once, execute `N` times in parallel

```
add<<<1,1>>>();
```



```
add<<<N,1>>>();
```

- Each parallel invocation of `add()` is referred to as a **block**



Moving to parallel

- Instead of executing `add()` once, execute `N` times in parallel

```
add<<<1,1>>>();
```

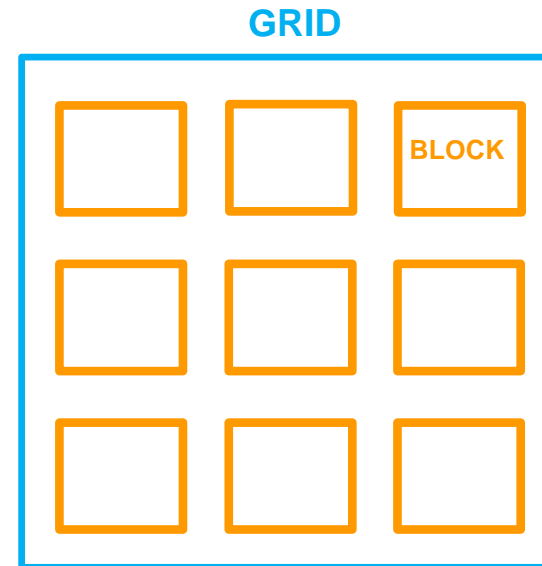


```
add<<<N,1>>>();
```

- Each parallel invocation of `add()` is referred to as a **block**

- block indexing using `blockIdx.x`

```
__global__ void add(int *a, int *b, int *c)
{
    c[blockIdx.x] = a[blockIdx.x] + b[blockIdx.x];
}
```



Block 0


```
c[0] = a[0] + b[0]
```

Block 1

```
c[1] = a[1] + b[1]
```

Vector addition

Number of blocks



```
#define N 512
int main(void)
{
    int *a, *b, *c;
    int *d_a, *d_b, *d_c;
    int size = N * sizeof(int);

    cudaMalloc((void **)&d_a, size);
    cudaMalloc((void **)&d_b, size);
    cudaMalloc((void **)&d_c, size);

    a = (int *)malloc(size); random_ints(a,N);
    b = (int *)malloc(size); random_ints(b,N);
    c = (int *)malloc(size);

    cudaMemcpy(d_a, &a, size, cudaMemcpyHostToDevice);
    cudaMemcpy(d_b, &b, size, cudaMemcpyHostToDevice);

    add<<<N,1>>>>(d_a, d_b, d_c);

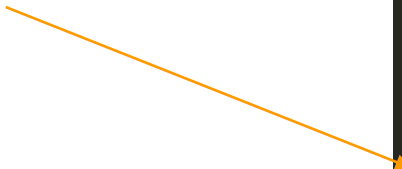
    cudaMemcpy(&c, d_c, size, cudaMemcpyDeviceToHost);

    free(a); free(b); free(c);
    cudaFree(d_a); cudaFree(d_b); cudaFree(d_c);
    return 0;
}
```


Vector addition

Number of blocks

New setup of inputs



```
#define N 512
int main(void)
{
    int *a, *b, *c;
    int *d_a, *d_b, *d_c;
    int size = N * sizeof(int);

    cudaMalloc((void **)&d_a, size);
    cudaMalloc((void **)&d_b, size);
    cudaMalloc((void **)&d_c, size);

    a = (int *)malloc(size); random_ints(a,N);
    b = (int *)malloc(size); random_ints(b,N);
    c = (int *)malloc(size);

    cudaMemcpy(d_a, &a, size, cudaMemcpyHostToDevice);
    cudaMemcpy(d_b, &b, size, cudaMemcpyHostToDevice);

    add<<<N,1>>>>(d_a, d_b, d_c);

    cudaMemcpy(&c, d_c, size, cudaMemcpyDeviceToHost);

    free(a); free(b); free(c);
    cudaFree(d_a); cudaFree(d_b); cudaFree(d_c);
    return 0;
}
```

Vector addition

Number of blocks

New setup of inputs

Launch kernel in N blocks

```
#define N 512
int main(void)
{
    int *a, *b, *c;
    int *d_a, *d_b, *d_c;
    int size = N * sizeof(int);

    cudaMalloc((void **)&d_a, size);
    cudaMalloc((void **)&d_b, size);
    cudaMalloc((void **)&d_c, size);

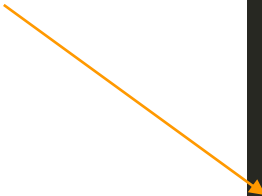
    a = (int *)malloc(size); random_ints(a,N);
    b = (int *)malloc(size); random_ints(b,N);
    c = (int *)malloc(size);

    cudaMemcpy(d_a, &a, size, cudaMemcpyHostToDevice);
    cudaMemcpy(d_b, &b, size, cudaMemcpyHostToDevice);

    add<<<N,1>>>>(d_a, d_b, d_c);

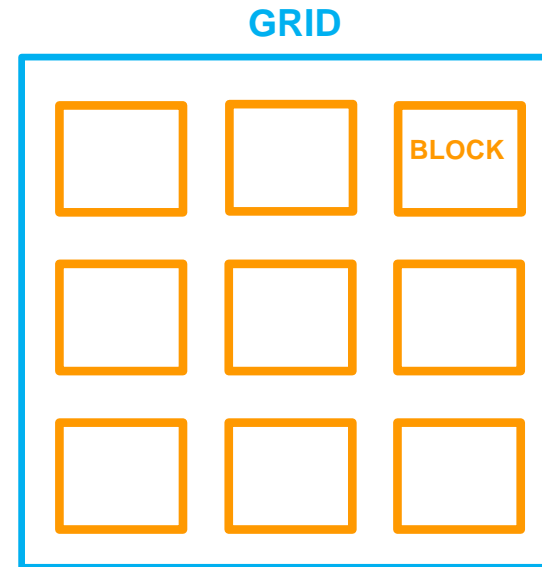
    cudaMemcpy(&c, d_c, size, cudaMemcpyDeviceToHost);

    free(a); free(b); free(c);
    cudaFree(d_a); cudaFree(d_b); cudaFree(d_c);
    return 0;
}
```

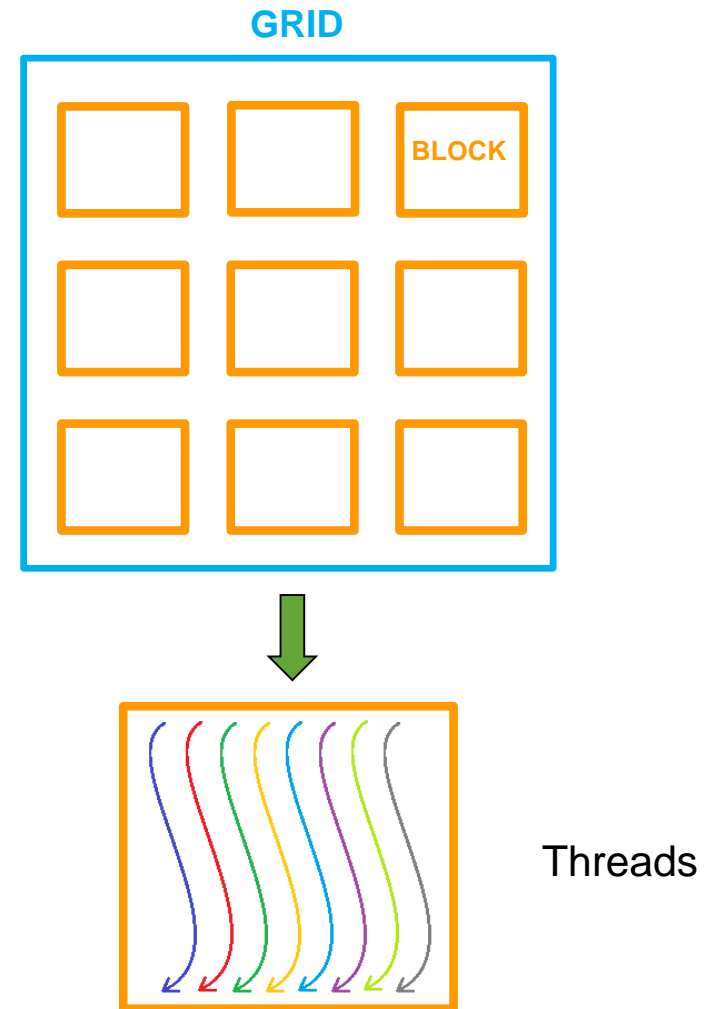


Introducing threads

Threads



Threads

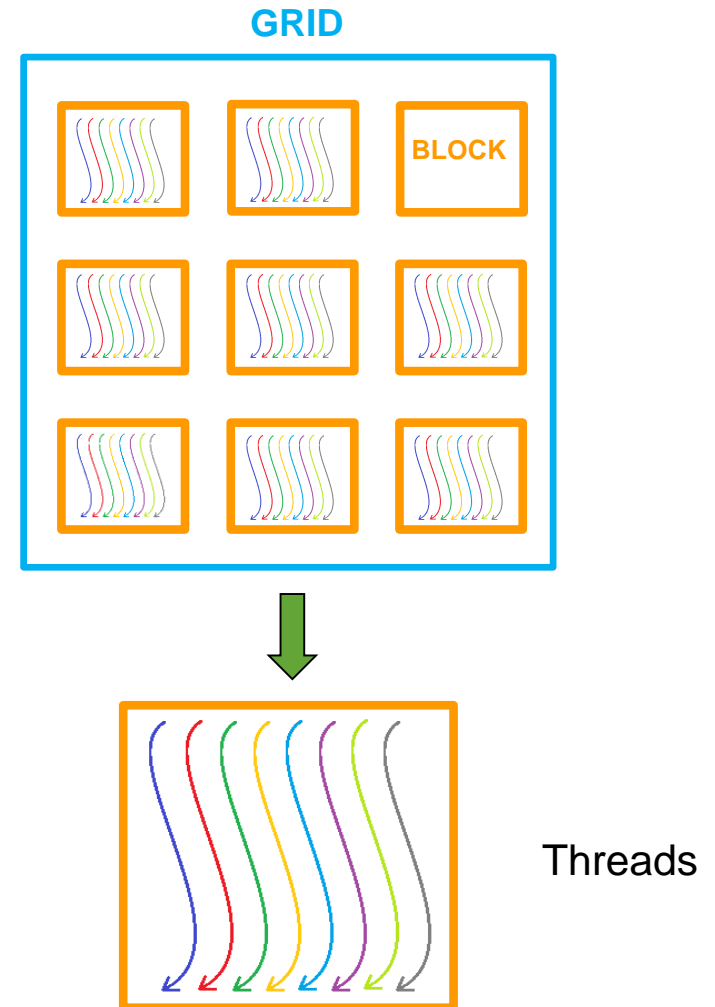


Threads

- Change `add()` to use parallel threads

`add<<<N, 1>>>();` Blocks

↓
`add<<<1, N>>>();` Threads



Threads

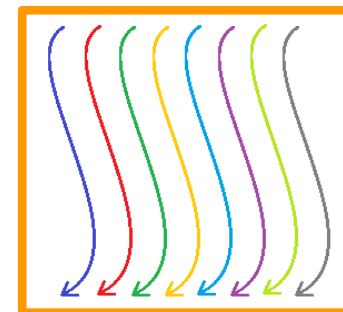
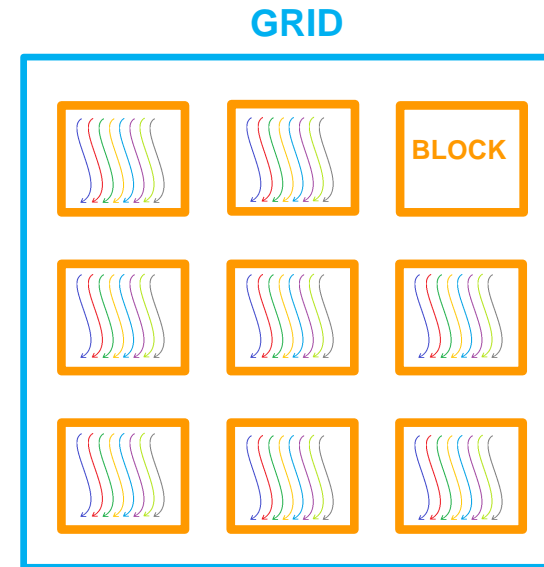
- Change add() to use parallel threads

add<<<N,1>>>(); Blocks

add<<<1,N>>>(); Threads

```
__global__ void add(int *a, int *b, int *c)
{
    c[blockIdx.x] = a[blockIdx.x] + b[blockIdx.x];
}
```

```
__global__ void add(int *a, int *b, int *c)
{
    c[threadIdx.x] = a[threadIdx.x] + b[threadIdx.x];
}
```



Threads

Review

- Using `__global__` to declare a function as device code
 - Executes on the device
 - Called from the host
- Basic device memory management
 - `cudaMalloc()`
 - `cudaMemcpy()`
 - `cudaFree()`
- Launching parallel kernels
 - With blocks `<<<N,1>>>`
 - With threads `<<<1,N>>>`

Review

■ Using `__global__` to declare a function as device code

- Executes on the device
- Called from the host

■ Basic device memory management

- `cudaMalloc()`
- `cudaMemcpy()`
- `cudaFree()`

■ Launching parallel kernels

- With blocks `<<<N,1>>>`
- With threads `<<<1,N>>>`

Steps in CUDA

1. Declare and allocate host and device memory.
2. Initialize host data.
3. Transfer data from the host to the device.
4. Execute one or more kernels.
5. Transfer results from the device to the host.

3 Ways to accelerate applications

Applications

Libraries

“Drop-in”
Acceleration

OpenACC
Directives

Easily
Accelerate
Applications

Programming
Languages

Maximum
Flexibility

Applications

Libraries

“Drop-in”
Acceleration

OpenACC
Directives

Easily
Accelerate
Applications

Programming
Languages

Maximum
Flexibility

Applications

Libraries

“Drop-in”
Acceleration

OpenACC
Directives

Easily
Accelerate
Applications

Programming
Languages

Maximum
Flexibility

B/S/H/

The background of the slide is an abstract, pixelated image. It features a dense grid of small squares in various shades of green, from bright lime green to dark forest green, set against a black background. The pattern is somewhat irregular, with some squares appearing brighter or more saturated than others, creating a sense of depth and movement. The overall effect is reminiscent of a digital or data-themed background.

Questions?