

CUDA Programming model

CUDA

CUDA -Compute Unified DeviceArchitecture

- is a parallel computing platform and programming model created by NVIDIA
- Implemented by the GPUs
- CUDA gives developers access to the instruction set and memory of the parallel computational elements in CUDA GPUs.
- Using CUDA, GPUs become accessible for computation like CPUs.

CUDA

- functions for the **GPU(device)** and functions for the system processor(**host**),
- CUDA uses **_device_or_global_** for **GPU** and -
-host-for the **CPU**.
- CUDA variables declared in the device or global functions are allocated to the GPU Memory

CUDA Philosophy

SIMT philosophy

- Single Instruction Multiple Threads

Computationally intensive

- The time spent on computation significantly exceeds the time spent on transferring data to and from GPU memory

Massively parallel

- The computations can be broken down into hundreds or thousands of independent units of work

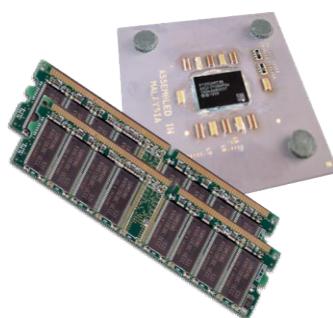
Definitions

- `dimGrid` - dimensions of the code (in blocks)
- `dimBlock`- dimensions of a block (in threads)
- `BlockIdx`- identifier for blocks
- `threadIdx`- identifier for threads per block
- `blockDim`- number of threads per block

Heterogeneous Computing

Host

- CPU and its memory
(host memory)



Device

- GPU and its memory
(device memory)



Heterogeneous Computing

```
#include <iostream>
#include <algorithm>

using namespace std;

#define N 1024
#define RADIUS 3
#define BLOCK_SIZE 16

_global void stencil_1d(int *in, int *out) {
    const int tempIndex = (int)(BLOCK_SIZE + 2 * RADIUS); int
    index = threadIdx.x * blockDim.x * blockDim.x; int
    index = threadIdx.x * blockDim.x + RADIUS;

    // Read input elements into shared memory
    temp[index] = in[index];
    if (threadIdx.x < RADIUS) {
        temp[(index - RADIUS)] = in[(index - RADIUS)];
        temp[(index + RADIUS)] = in[(index + RADIUS)];
    }

    // Synchronize (ensure all the data is available)
    __syncthreads();

    // Apply the stencil
    int result = 0;
    for (int offset = -RADIUS; offset <= RADIUS; offset++) {
        result += temp[(index + offset)];
    }

    // Store the result
    out[index] = result;
}

void fill_ints(int *x, int n) {
    fill_n(x, n, 1);
}

int main(void) {
    int *in, *out;           // host copies of a, b, c
    int *d_in, *d_out;       // device copies of a, b, c
    int size = (N + 2*RADIUS) * sizeof(int);

    // Alloc space for host copies and setup values
    in = (int *)malloc(size); fill_ints(in, N + 2*RADIUS); out
    = (int *)malloc(size); fill_ints(out, N + 2*RADIUS);

    // Alloc space for device copies
    cudaMalloc((void **)&d_in, size);
    cudaMalloc((void **)&d_out, size);

    // Copy to device
    cudaMemcpy(d_in, in, size, cudaMemcpyHostToDevice);
    cudaMemcpy(d_out, out, size, cudaMemcpyHostToDevice);

    // Launch stencil_1d kernel on GPU
    stencil_1d<<(NBLOCK_SIZE,BLOCK_SIZE>>>(d_in + RADIUS,
    d_out + RADIUS);

    // Copy result back to host
    cudaMemcpy(out, d_out, size, cudaMemcpyDeviceToHost);

    // Cleanup
    free(in); free(out);
    cudaFree(d_in); cudaFree(d_out);
    return 0;
}
```

parallel fn

serial code

parallel code

serial code



Heterogeneous Computing

```

#include <iostream>
#include <algorithm>

using namespace std;

#define N 1024
#define RADIUS 3
#define BLOCK_SIZE 16

__global void stencil_1d(int *in, int *out) {
    __shared__ int temp[BLOCK_SIZE + 2 * RADIUS]; int
    gridx = threadIdx.x + blockIdx.x * blockDim.x; int
    index = gridx * RADIUS; int

    // Read input elements into shared memory
    temp[threadIdx.x] = in[index];
    if (threadIdx.x == RADIUS) {
        temp[threadIdx.x - RADIUS] = in[index - RADIUS];
        temp[threadIdx.x + RADIUS] = in[index + RADIUS];
    }

    // Synchronize (ensure all the data is available)
    __syncthreads();

    // Apply the stencil
    int result = 0;
    for (int offset = -RADIUS; offset <= RADIUS; offset++)
        result += temp[index + offset];

    // Store the result
    out[index] = result;
}

void fill_ints(int *x, int n) {
    fill_n(x, n, 1);
}

int main(void) {
    int *in, *out;           // host copies of a, b, c
    int *d_in, *d_out;       // device copies of a, b, c
    int size = (N + 2 * RADIUS) * sizeof(int);

    // Alloc space for host copies and setup values
    in = (int *)malloc(size); fill_ints(in, N + 2 * RADIUS); out
    = (int *)malloc(size); fill_ints(out, N + 2 * RADIUS);

    // Alloc space for device copies
    cudaMalloc(&void **d_in, size);
    cudaMalloc(&void **d_out, size);

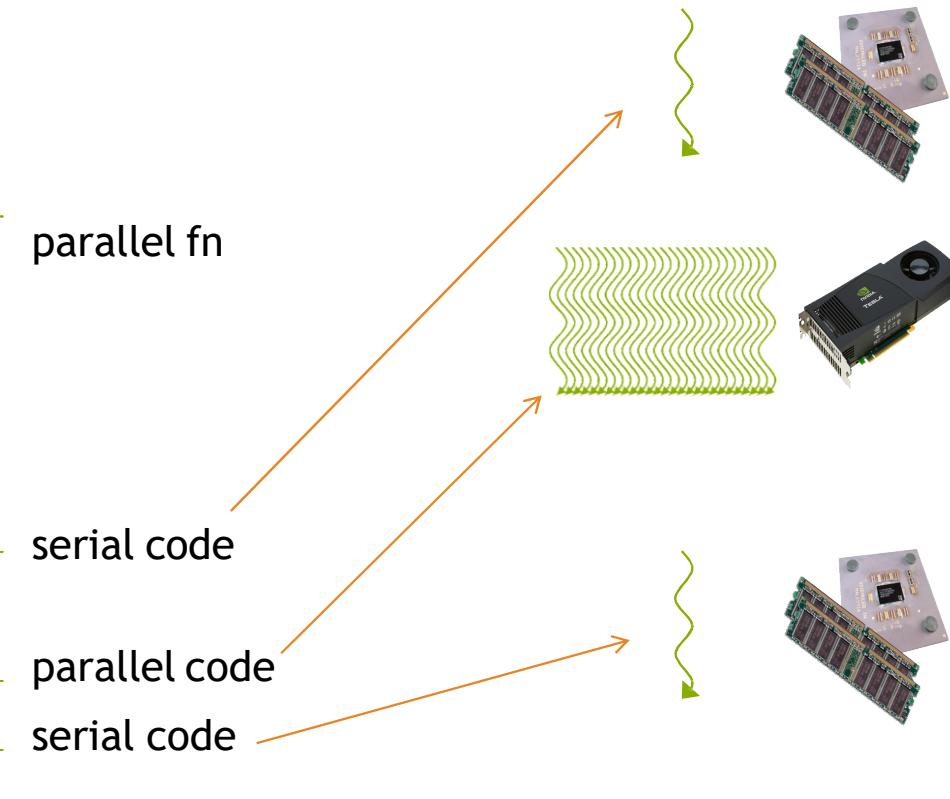
    // Copy to device
    cudaMemcpy(d_in, in, size, cudaMemcpyHostToDevice);
    cudaMemcpy(d_out, out, size, cudaMemcpyHostToDevice);

    // Launch stencil_1d() kernel on GPU
    stencil_1d<<N<<BLOCK_SIZE.BLOCK_SIZE>>(d_in + RADIUS,
    d_out + RADIUS);

    // Copy result back to host
    cudaMemcpy(out, d_out, size, cudaMemcpyDeviceToHost);

    // Cleanup
    free(in); free(out);
    cudaFree(d_in); cudaFree(d_out);
    return 0;
}

```



Hello World with CUDA

```
#include <stdio.h>                                $ nvcc hello-world.cu
#include <cuda.h>

__global__ void hwkernel() {
    printf("Hello world!\n");
}

int main() {
    hwkernel<<<1, 1>>>();
}
```

\$./a.out

\$



Hello World with CUDA

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

```
#include <cuda.h>
```

```
$ nvcc hello-world.cu
```

```
_global__void hwkernel() {
```

```
$ ./a.out
```

```
printf("Hello world!\n");
```

```
Hello world!
```

```
}
```

```
$
```

```
int main() {
```

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

```
cudaDeviceSynchronize()
```

Program returns immediately after launching the kernel. To prevent program to finish before kernel is completed, we call `cudaDeviceSynchronize()`.

```
}
```

Hello World with CUDA

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

global __ void hwkernel() {
    printf("Hello world!\n");
}

int main()
{
    hwkernel<<<1, 32>>>();
    cudaThreadSynchronize();
}
```

```
$ nvcc hello-world.cu
$ ./a.out
Hello world!
Hello world!
Hello world!
Hello world!
Hello world!
...
...
$
```



How NVCC works?

- Nvcc is a driver program
 - Compiles and links all input files
 - Requires a general-purpose C/C++ host compiler
 - Uses gcc and g++ by default on Linux platforms
 - nvcc --version

<https://docs.nvidia.com/cuda/cuda-compiler-driver-nvcc/index.html>



CUDA Compilation Trajectory

- Conceptually, the flow is as follows
 - Input program is preprocessed for device compilation
 - It is compiled to a CUDA binary and/or PTX (Parallel Thread Execution) intermediate code which are encoded in a fatbinary
- Input program is processed for compilation of the host code
 - CUDA-specific C++ constructs are transformed to standard C++ code
 - Synthesized host code and the embedded fatbinary are linked together to generate the executable

Function Declarations in CUDA

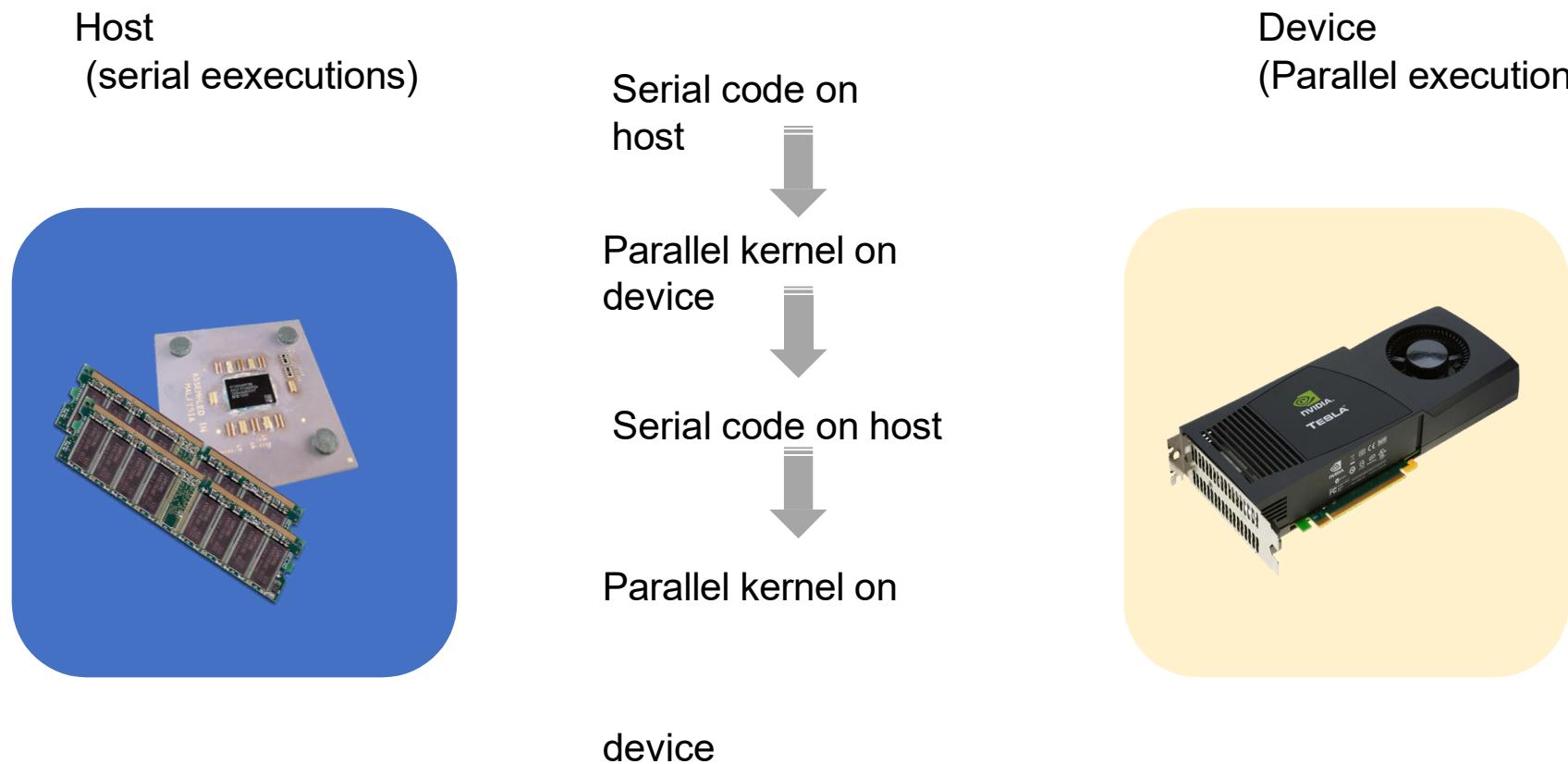
	Executed on	Callable from
<code>__device__ float deviceFunc()</code>	Device	Device
<code>__global__ void kernelFunc()</code>	Device	Host
<code>__host__ float hostFunc()</code>	Host	Host

- `__global__` define a kernel function, must return `void`
- `__device__` functions can have return values
- `__host__` is default, and can be omitted
- Prepending `__host__` `__device__` causes the system to compile separate host and device versions of the function

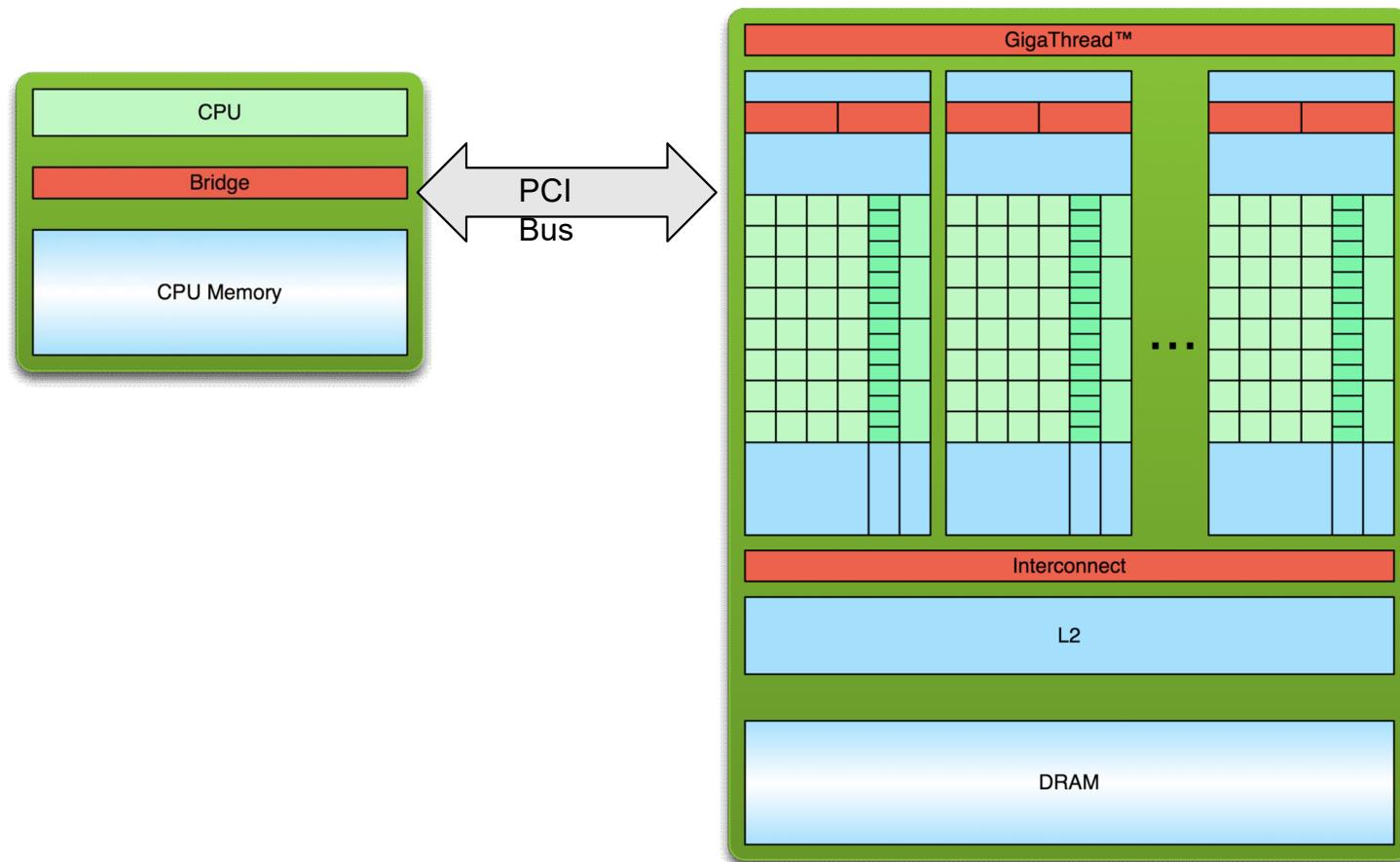
Dynamic Parallelism

- It is possible to launch kernels from other kernels
- Calling `_global` functions from the device is referred to as dynamic parallelism
 - Requires CUDA devices of compute capability 3.5 and CUDA 5.0 or higher

Execution Model

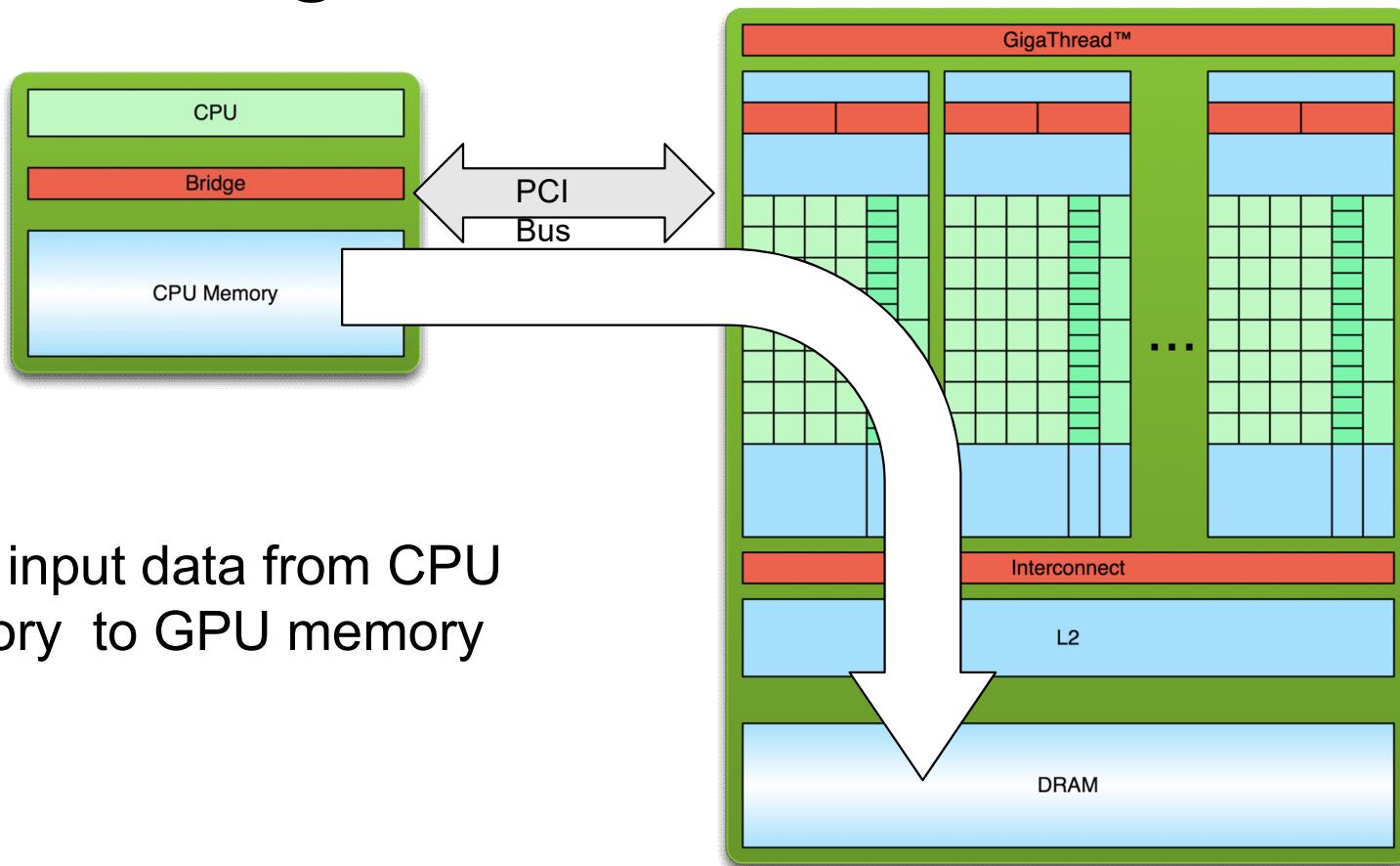


Simple Processing Flow

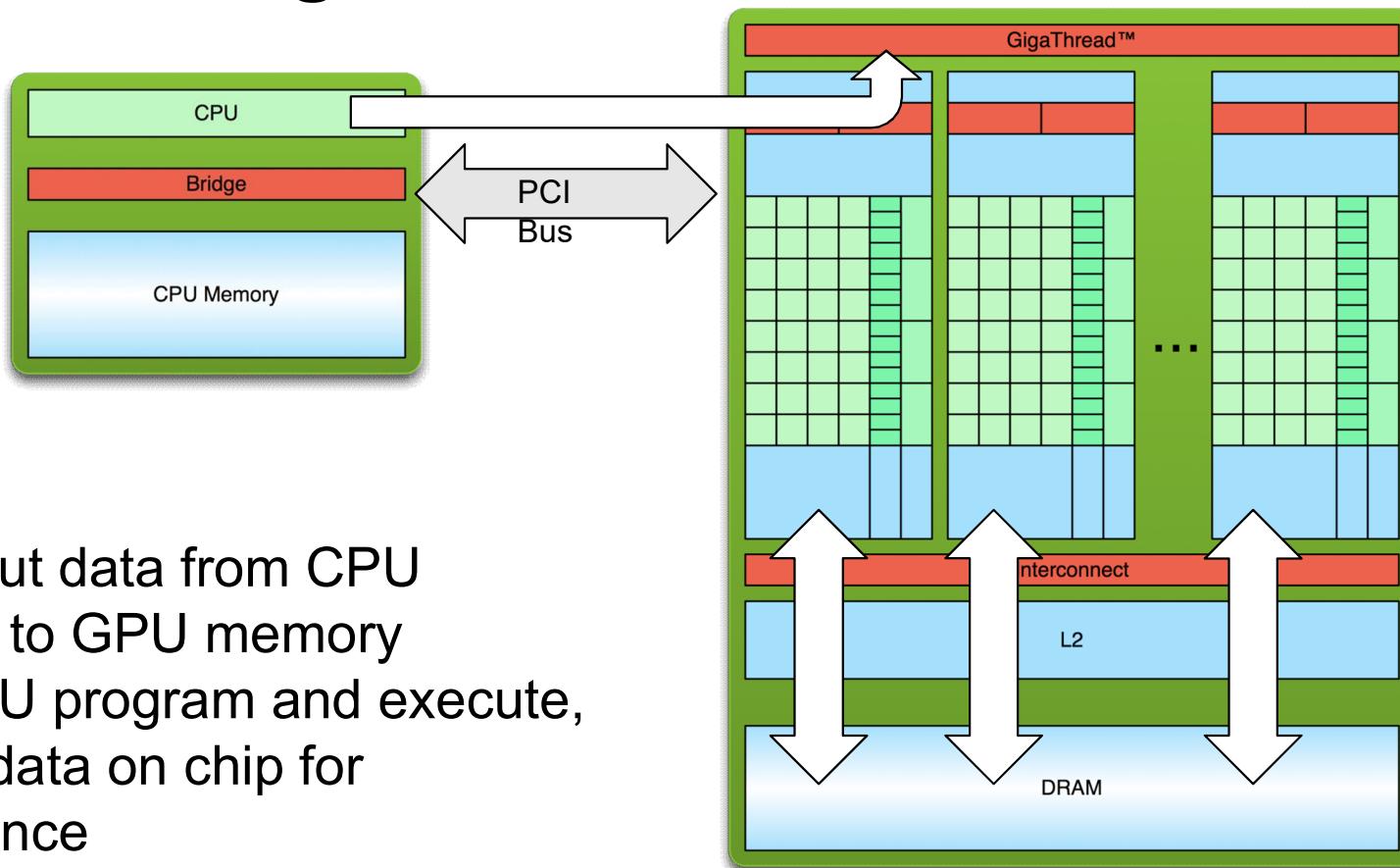


Simple Processing Flow

1. Copy input data from CPU memory to GPU memory

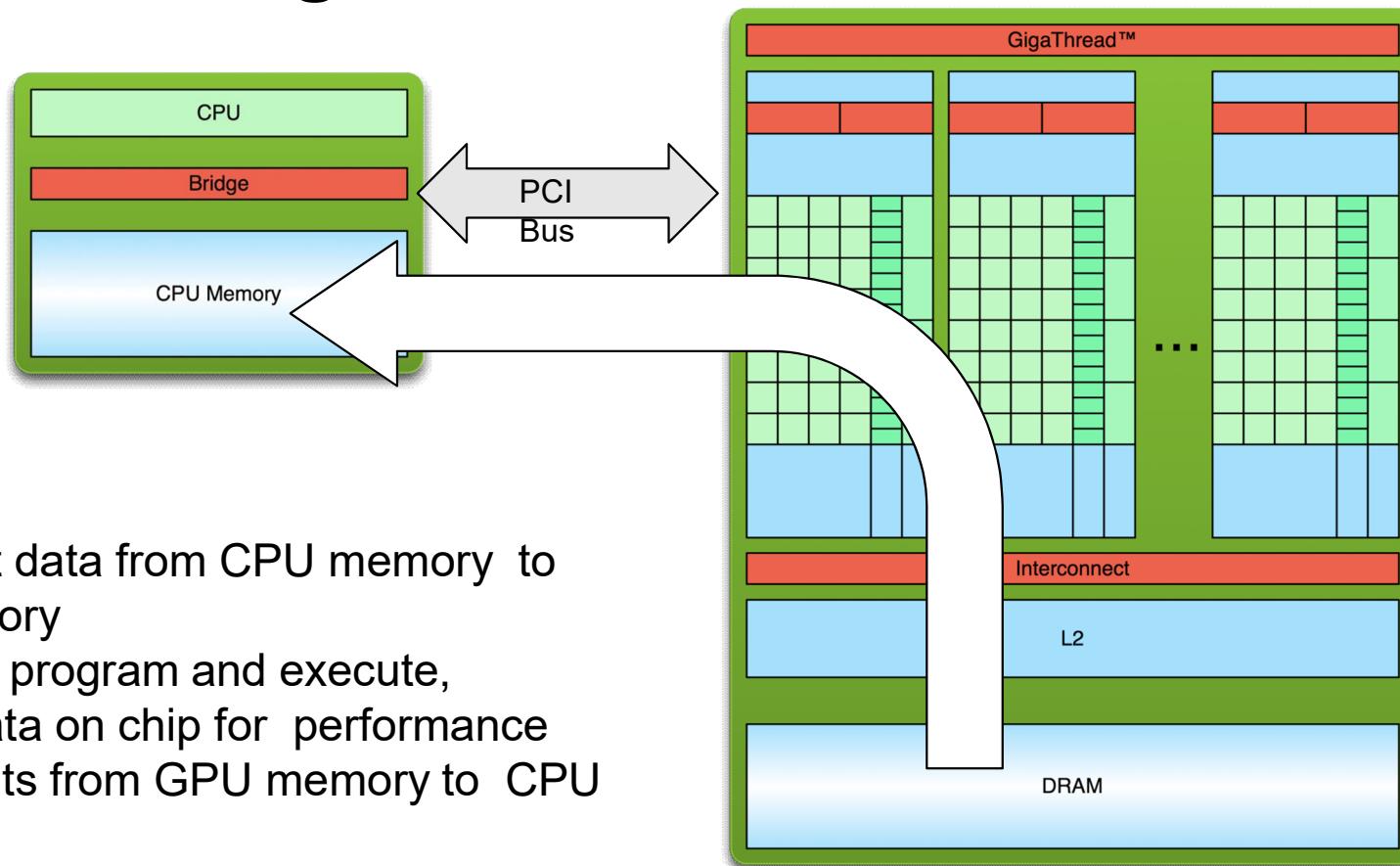


Simple Processing Flow



1. Copy input data from CPU memory to GPU memory
2. Load GPU program and execute, caching data on chip for performance

Simple Processing Flow



ssn