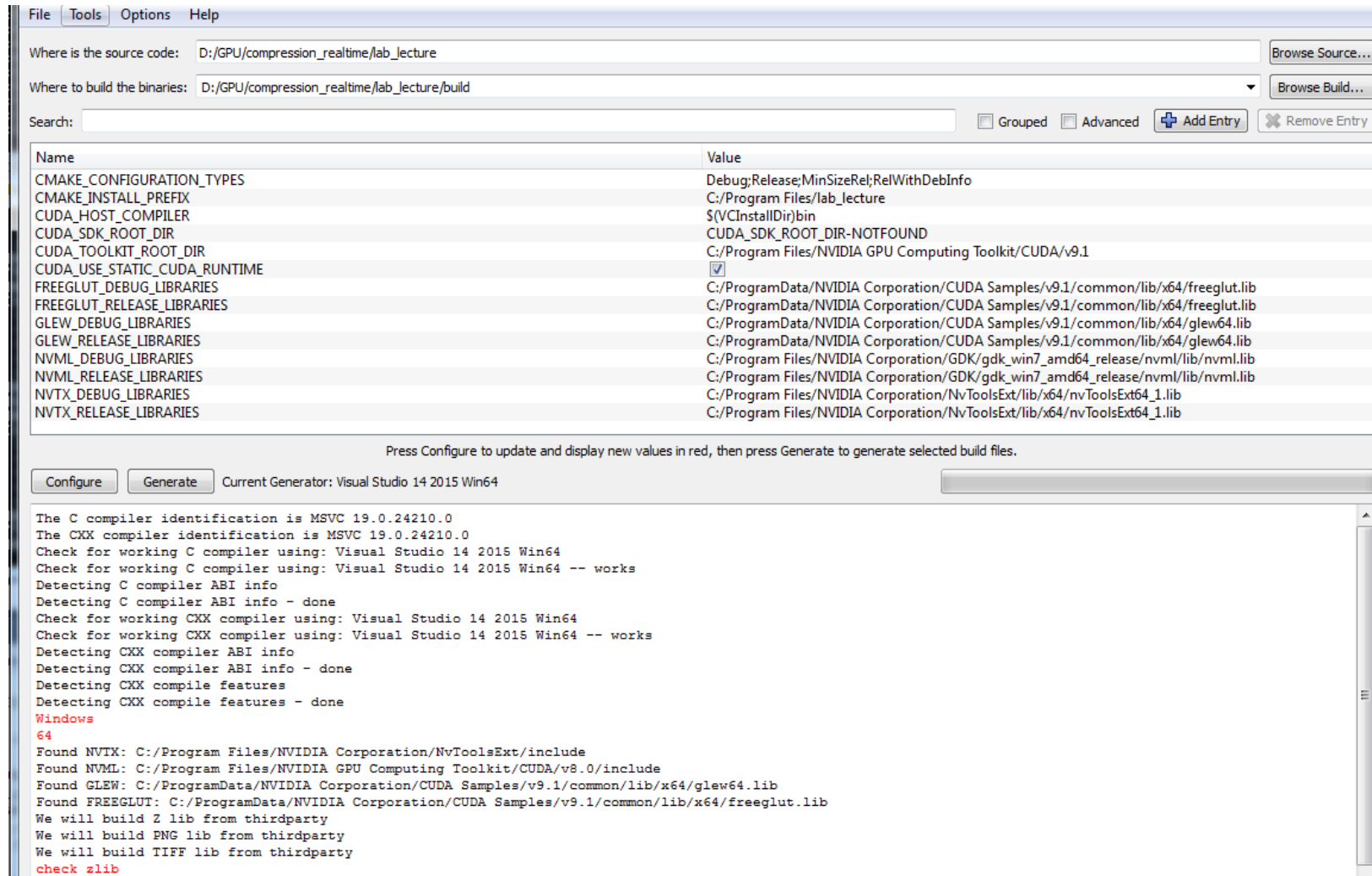# Lab Lecture

ECE 277
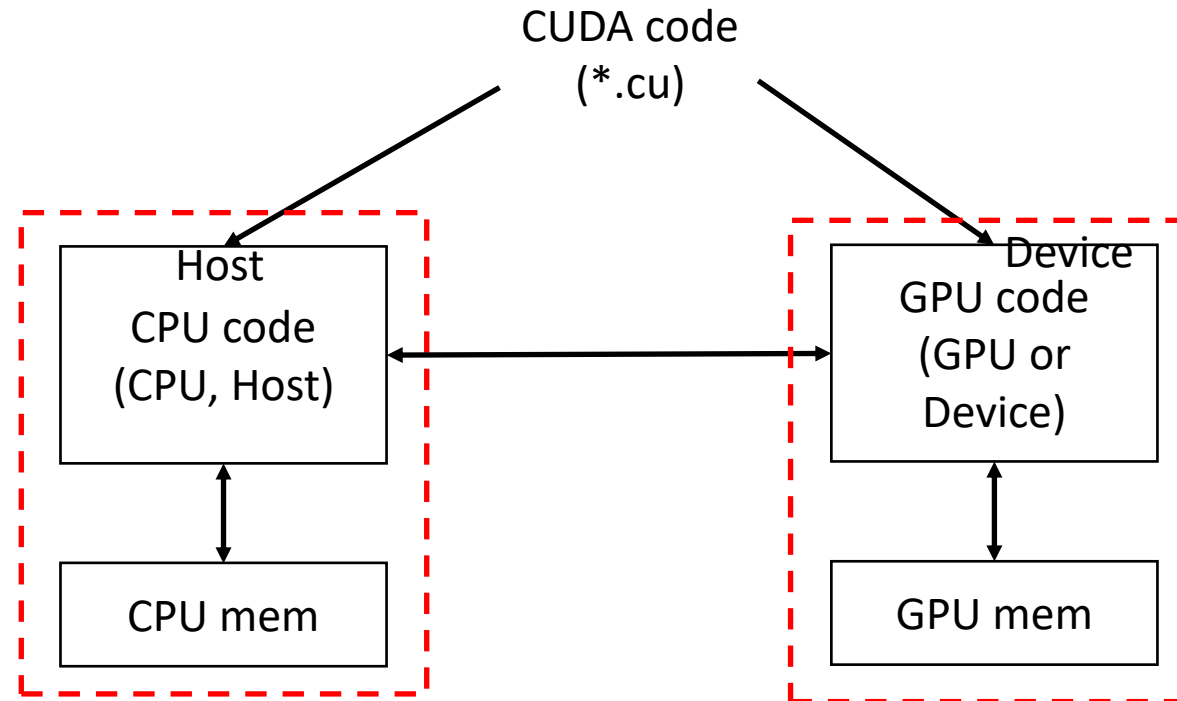
Cheolhong An

# CMake setup: Create a VS project solution
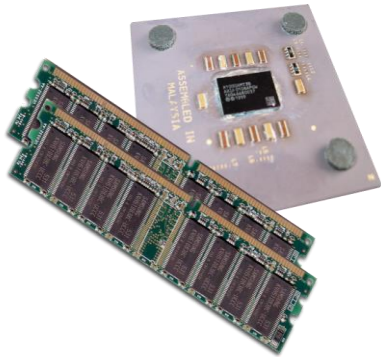
# Heterogeneous Computing

- A single code can run both CPU and GPU
- CUDA (*.cu) is a heterogenous programming language

# Heterogeneous Computing

- Terminology:
  - *Host*    The CPU and its memory (host memory)
  - *Device*  The GPU and its memory (device memory)

Host (CPU)                    Device (GPU)

Mark Harris, **Introduction to CUDA C,** NVIDIA Corporation

# Heterogeneous Computing



parallel fn

serial code

parallel code

serial code

# Hello.cu

- Create your own project (Your_name)
    1) Create one directory (any name is fine, e.g. "test") under "Src" directory
    2) Add "add_subdirectory(your directory name)" in CmakeLists.txt under "Src" directory
        e.g. add_subdirectory(test)
    3) Copy CMakeLists.txt under "SimpleGL" into your directory ("test")
    4) Replace project name ("check_system") with your name, which is the project name in CMakeList.txt
    5) Create Hello.cu file in your directory ("test")
    6) Replace simpleGL.cu with Hello.cu in CMakeList.txt
    7) Edit Hello.cu with the next slide functions
    8) Configure CMake using CMake-gui  or rebuild solution
        If you encounter any error, you should use CMake-gui to check the error message

- Create multiple threads in a thread block
- Create multiple threadblocks in a grid

# Hello World! with Device Code

```
__global__ void mykernel(void) {
        // Print your name, site number
}


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

# Hello World! with Device Code

```
__global__ void mykernel(void) {

}
```

- CUDA C/C++ keyword __global__ indicates a function that:
  - Runs on the device
  - Is called from host code

- nvcc (CUDA compiler) separates source code into host and device components
  - Device functions (e.g. hello**()**) processed by NVIDIA compiler
  - Host functions (e.g. **main()**) processed by standard host compiler
    - **gcc, cl.exe**

# Common CUDA code structure

- Ex) vectorAdd

CPU memory allocation (malloc) <- Host
GPU memory allocation (cudaMalloc) <- Device

Transfer data from Host to Device
cudaMemcpy( d_a, a, size, cudaMemcpyHostToDevice)

Run GPU codes (Call kernel functions)
Kernel_fun<<<,>>> ()

Transfer processed data from Device to Host
cudaMemcpy(a, d_a, size, cudaMemcpyDeviceToHost)

# Transfer data from Host to Device



1. Copy input data from CPU memory to GPU memory

# Run GPU codes



CPU

Bridge

CPU Memory

PCI Bus

GigaThread™

Interconnect

L2

DRAM

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

# Transfer processed data from Device to Host

CPU

Bridge

PCI Bus

CPU Memory

GigaThread™

...

Interconnect

L2

DRAM

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

Mark Harris, **Introduction to CUDA C,** NVIDIA Corporation

# Pytorch High-level Interface

```python
device = torch.device("cuda" if use_cuda else "cpu")
model = Net().to(device)

def train(args, model, device, train_loader, optimizer, epoch):
    model.train()
    for batch_idx, (data, target) in enumerate(train_loader):
        data = data.to(device)
        target = target.to(device)
        optimizer.zero_grad()
        output = model(data)
        loss = F.nll_loss(output, target)
        loss.backward()
        optimizer.step()
        if batch_idx % args.log_interval == 0:
            print('Train Epoch: {} [{}/{} ({:.0f}%)]\tLoss: {:.6f}'.format(
                epoch, batch_idx * len(data), len(train_loader.dataset),
                100. * batch_idx / len(train_loader), loss.item()))
```
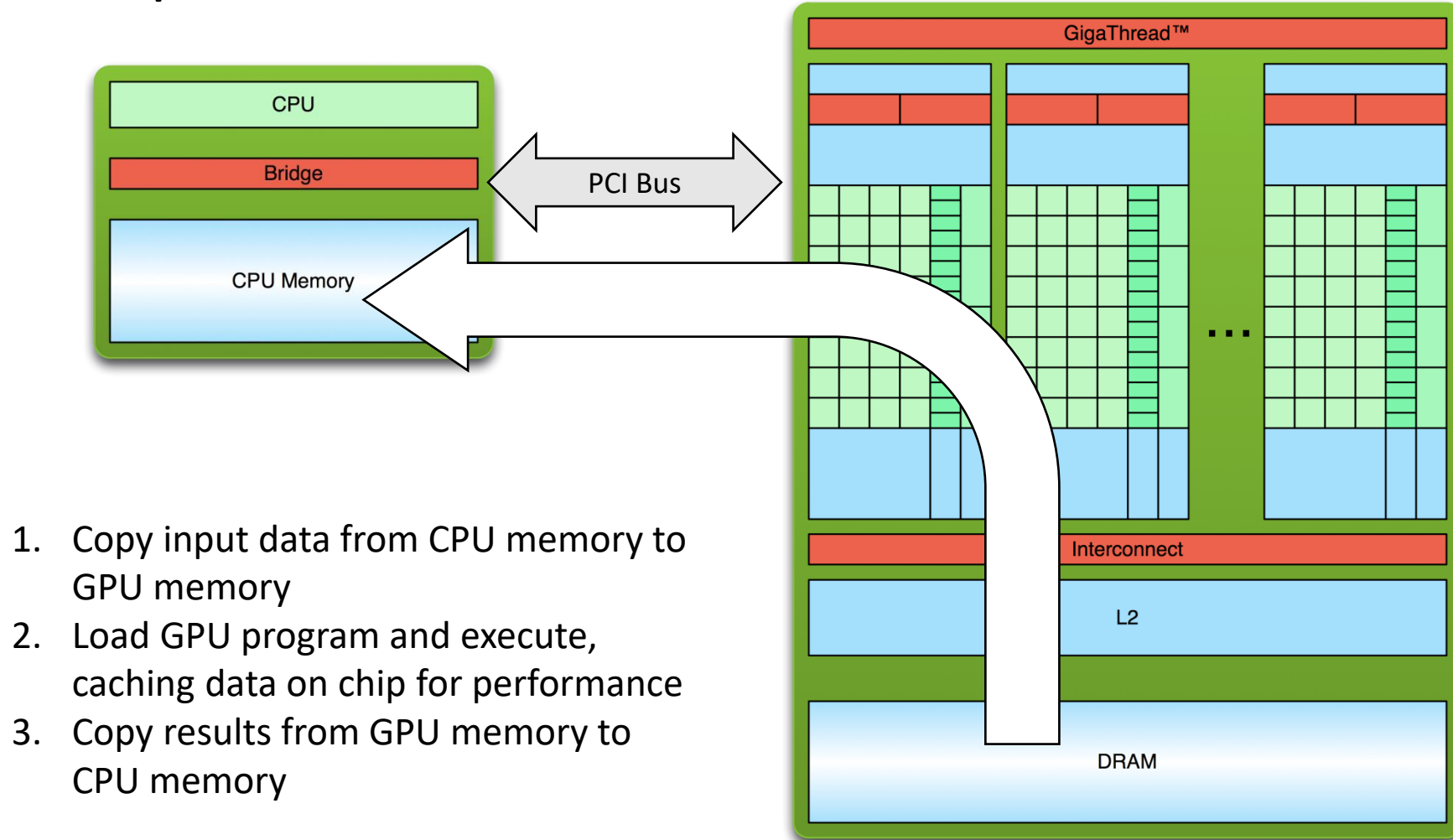
CPU to GPU transfer

CUDA operation

CUDA variables

# Pytorch framework with CUDA

Machine learning Applications

NLP

Medical Image Classification

Pytorch ML framework

Python Frontend

C++ Frontend

C++ Backend

Customized C++ (CUDA)

Tensors

CUDA

AD

AD: Automatic differentiation

# Check bandwidth

- Memory bandwidth
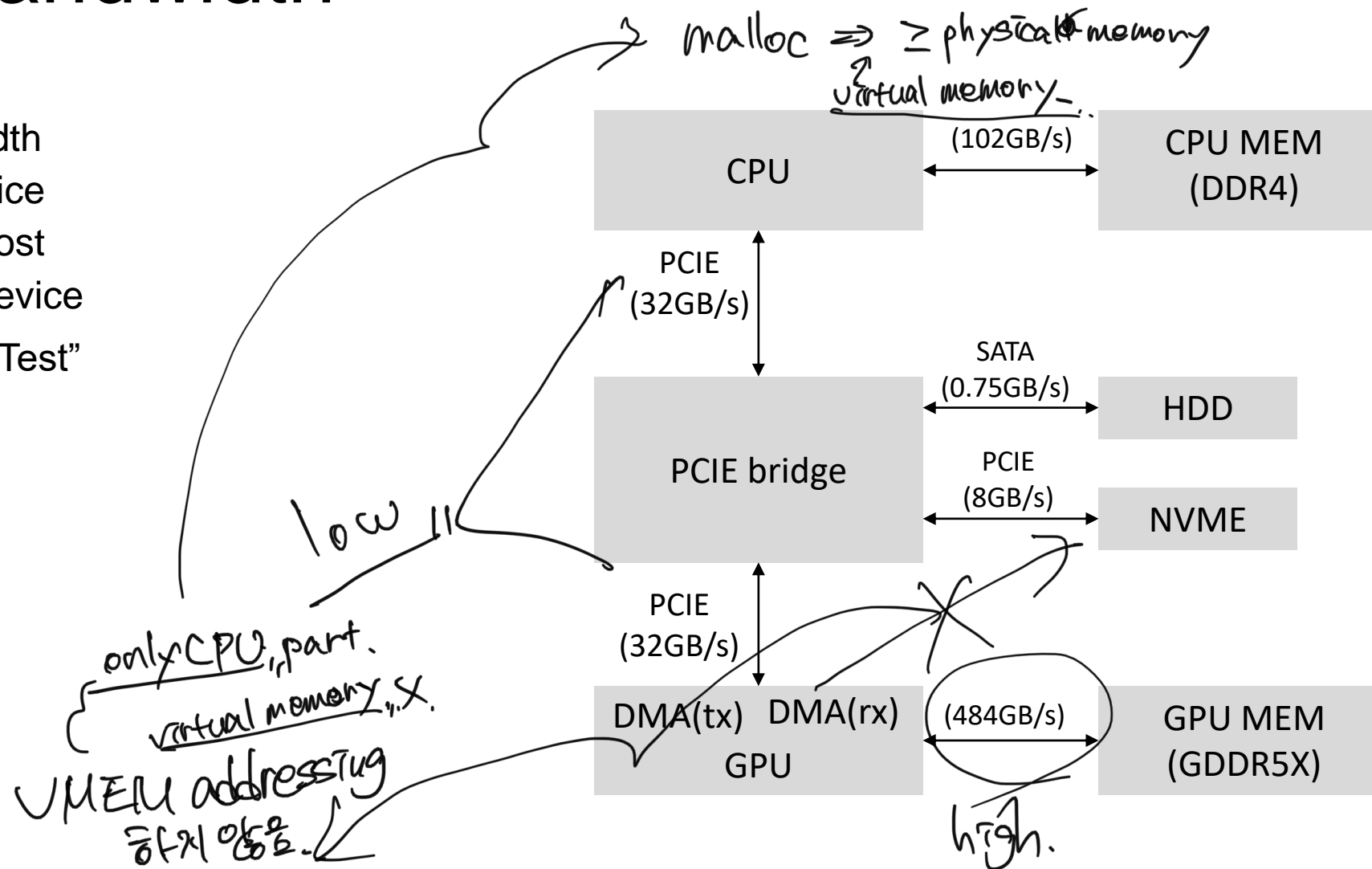  - Host to Device
  - Device to Host
  - Device to Device

- Run "bandwidth Test"

malloc ⇒ ≥ physical memory

virtual memory.

| CPU | (102GB/s) | CPU MEM (DDR4) |

PCIE (32GB/s)

| PCIE bridge | SATA (0.75GB/s) | HDD |
| | PCIE (8GB/s) | NVME |

PCIE (32GB/s)

low ll

only CPU, part.

virtual memory, X.

VMEM addressing 하지 않음.

| DMA(tx)  DMA(rx) GPU | (484GB/s) | GPU MEM (GDDR5X) |

high.

# Parallel Programming in CUDA C/C++

- But wait... GPU computing is about massive parallelism!

- We need a more interesting example...

- We'll start by adding two integers and build up to vector addition

a   b   c

# Addition on the Device

- A simple kernel to add two integers

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

- As before __global__ is a CUDA C/C++ keyword meaning
  - add() will execute on the device
  - add() will be called from the host

# Addition on the Device

- Note that we use pointers for the variables

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

- **add()**  runs on the device, so a, b and c must point to device memory

- We need to allocate memory on the GPU