# Introduction to Multi-GPU Computing

Sebastian Kuckuk

Zentrum für Nationales Hochleistungsrechnen Erlangen (NHR@FAU)

# Enroll in the Course

- Go to

  **courses.nvidia.com/dli-event**
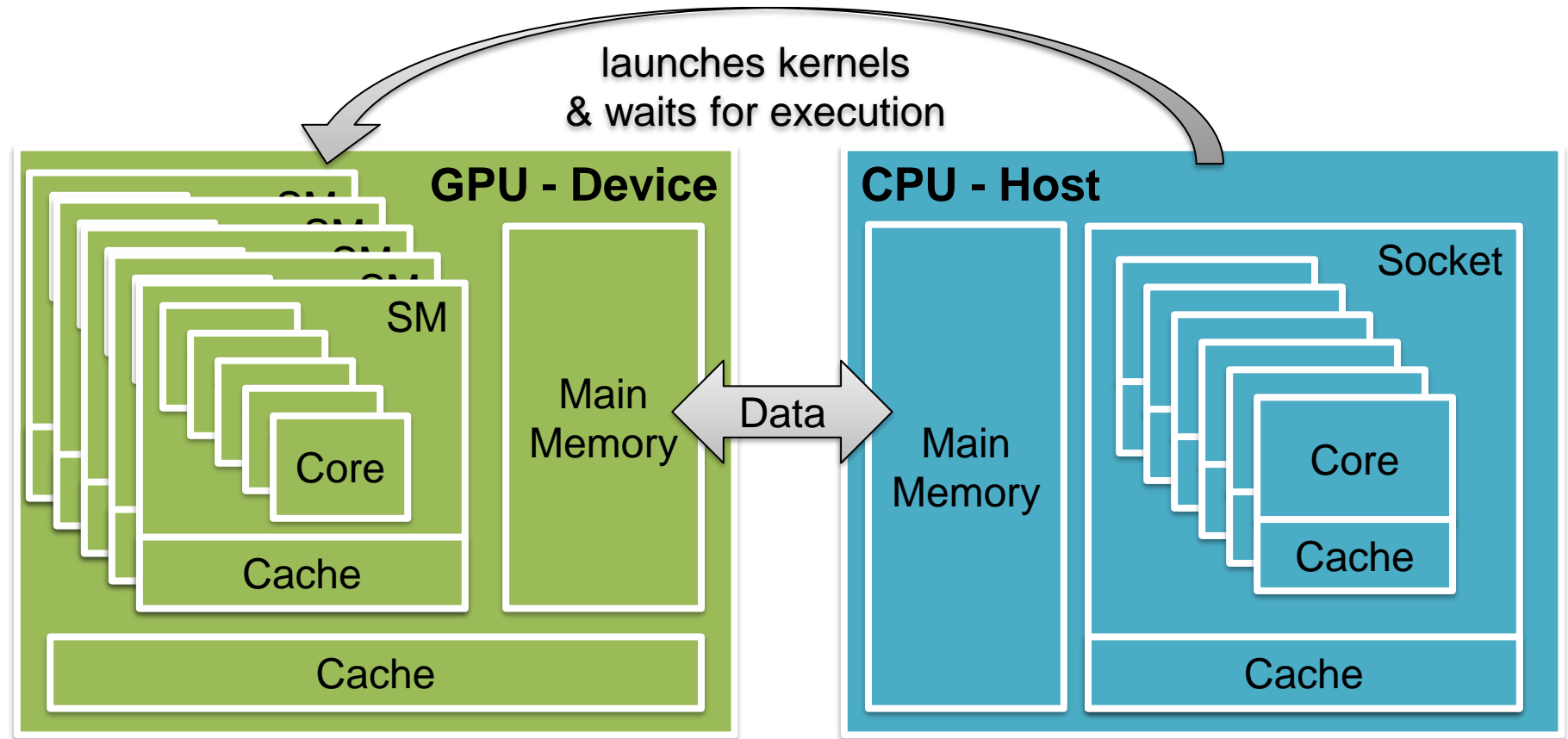
- Enter the code

  **[event code]**

- All your courses are available at

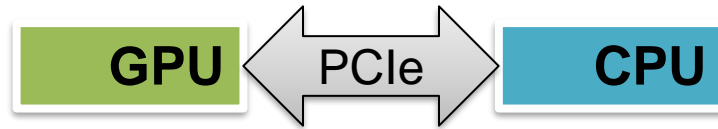  **courses.nvidia.com/dashboard**

# Enroll in the Course

- Set your Zoom name as [first name] [last name] ([affiliation])

- The DLI part of the course is composed of multiple modules (one IPython notebook each) and augmented with additional material available at https://github.com/SebastianKuckuk/accelerated-programming

- Pass the assessment(s) to get a certificate from NVIDIA

- You will have access to the course material at least six months
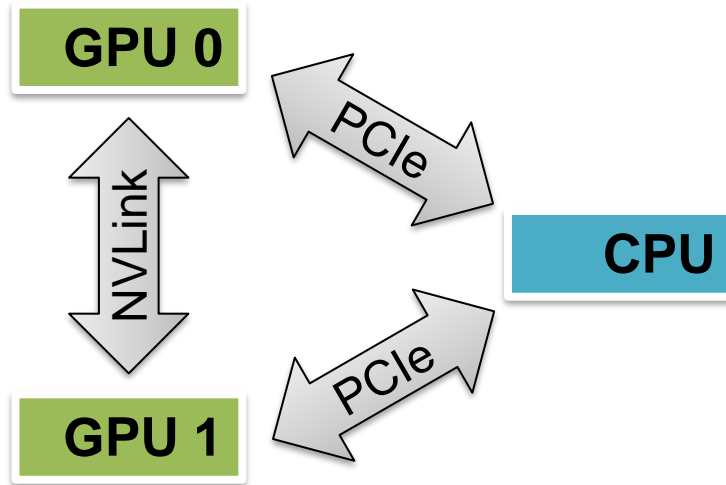
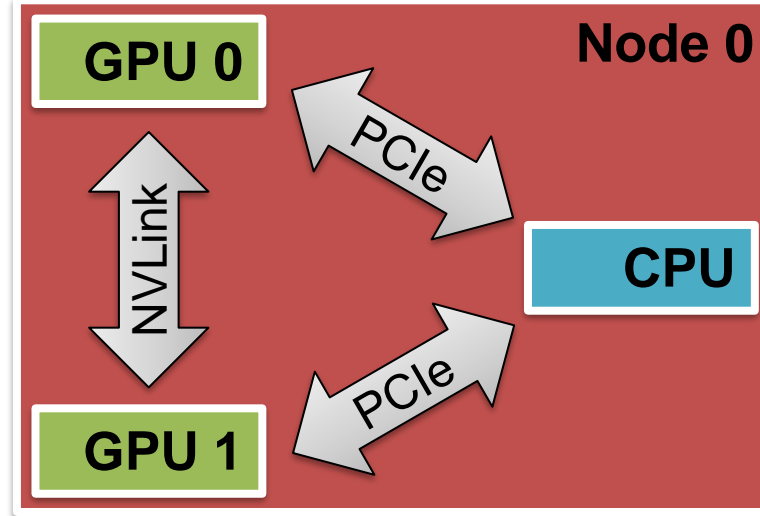- Feel free to interrupt and ask questions

# Simplified Architecture

# Simplified Architecture
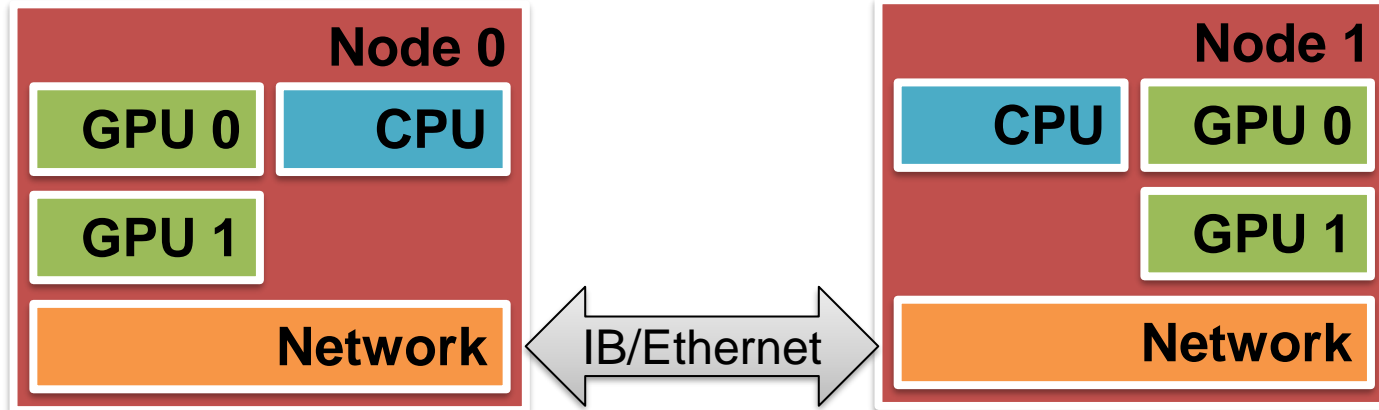


GPU ⟷ PCIe ⟷ CPU

# Simplified Architecture

# Simplified Architecture

# Simplified Architecture

# Workflow (single GPU)

1. Initialize data on CPU

2. Copy data from CPU to GPU

3. Launch GPU kernels

4. Do independent work on CPU (optional)

5. Synchronize GPU

6. Copy data from GPU to CPU

7. Post-process data on CPU

# Workflow Example

- Goal: repetition of basic CUDA C++ programming elements

- Sample application: copy array and increase each element by 1

- Full code available at
  - https://github.com/SebastianKuckuk/accelerated-programming/blob/master/

# Workflow

- 0. Allocate Data        (managed)    (explicit)

```
int main(int argc, char *argv[]) {
    size_t nx   = atoi(argv[1]);
    size_t size = sizeof(double) * nx;

    double *src, *dest;                    double *src, *dest;
    cudaMallocManaged(&src,  size);        cudaMallocHost(&src,  size);
    cudaMallocManaged(&dest, size);        cudaMallocHost(&dest, size);

                                           double *d_src, *d_dest;
                                           cudaMalloc(&d_src,  size);
                                           cudaMalloc(&d_dest, size);

    // ...
}
```

# Workflow

- 1. Initialize data on CPU

```
void initOnCPU(double *src, size_t nx) {
    for (size_t i = 0; i < nx; ++i)
        src[i] = 1337.;
}


int main(/* ... */) {
    // allocate

    initOnCPU(src, nx);

    // ...
}
```
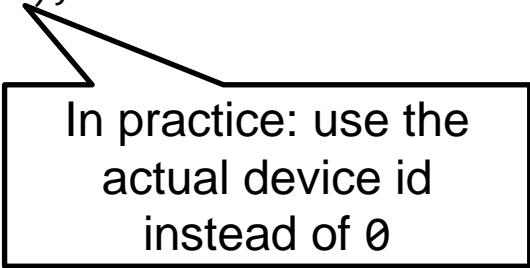
# Workflow

- 2. Copy data from CPU to GPU

```
// allocate & init

cudaMemPrefetchAsync(src,  size, 0);        cudaMemcpy(d_src, src, size,
                                                       cudaMemcpyHostToDevice);


cudaMemPrefetchAsync(dest, size, 0);

// ...
```

In practice: use the actual device id instead of 0

# Workflow

- 3. Launch GPU kernels

```cuda
__global__ void copyOnGPU(double *src, double *dest, size_t nx) {
    size_t i = blockIdx.x * blockDim.x + threadIdx.x;


    if (i < nx)
        dest[i] = src[i] + 1;
}

// ... in main
copyOnGPU<<<(nx + 255) / 256, 256>>>(  src,   dest, nx);
// ... for managed, or for explicit
copyOnGPU<<<(nx + 255) / 256, 256>>>(d_src, d_dest, nx);
```

# Workflow

- 3. Launch GPU kernels (grid-stride loop)

```
__global__ void copyOnGPU(double *src, double *dest, size_t nx) {
    size_t start  = blockIdx.x * blockDim.x + threadIdx.x;
    size_t stride = gridDim.x * blockDim.x;

    for (size_t i = start; i < nx; i += stride)
        dest[i] = src[i] + 1;
}

// ... in main
copyOnGPU<<<1280, 256>>>(  src,   dest, nx);
// for managed or for explicit
copyOnGPU<<<1280, 256>>>(d_src, d_dest, nx);
```

# Workflow

- 5. Synchronize GPU

```
cudaDeviceSynchronize();
```

# Workflow

- 6. Copy data from GPU to CPU

```
// computation

cudaMemPrefetchAsync(dest, size,          cudaMemcpy(dest, d_dest, size,
                     cudaCpuDeviceId);                cudaMemcpyDeviceToHost);


// ...
```

# Workflow

- 7. Post-process data on CPU

```
void checkOnCPU(double *dest, size_t nx) {
    for (size_t i = 0; i < nx; ++i)
        assert(1338. == dest[i]);
}
```
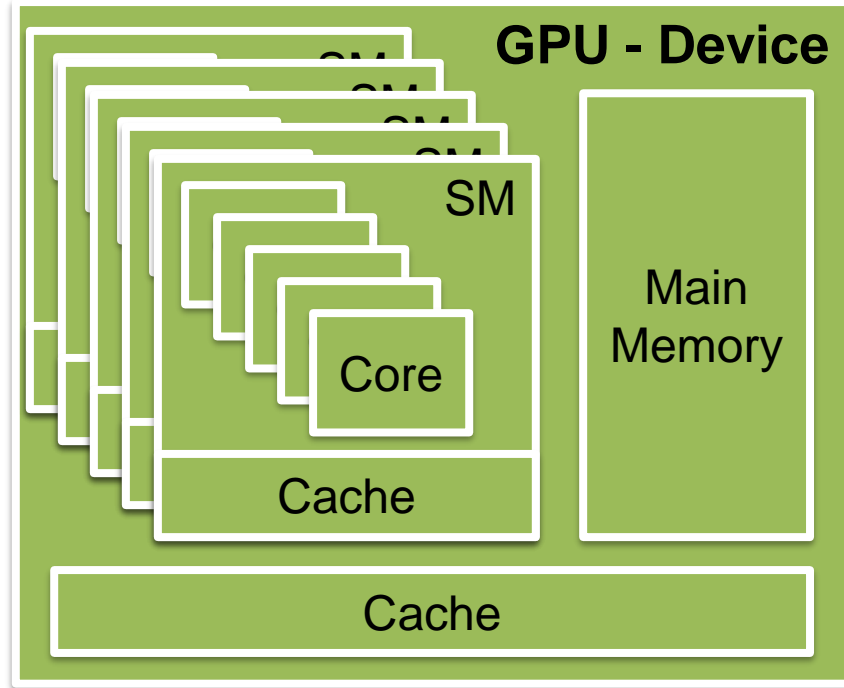
# Workflow

- 8. De-Allocate Data

```
// post-processing

cudaFree(src);                          cudaFree(d_src);
cudaFree(dest);                         cudaFree(d_dest);

                                        cudaFreeHost(src);
                                        cudaFreeHost(dest);
```

# CUDA Mapping



- **Grids** are mapped to **devices**

- **Blocks** are mapped to **SMs**

- **Threads** are mapped to **cores**

- **Threads of a block** are executed in **warps** (groups of 32 threads)

# Welcome to *[course name]*

# We will start at 9:00

# Enjoy your coffee break!

# We will reconvene at 10:30

# Enjoy your lunch break!

# We will reconvene at 13:30

# Enjoy your break!

# We will reconvene at 10:30

# Please remember to start the next lab