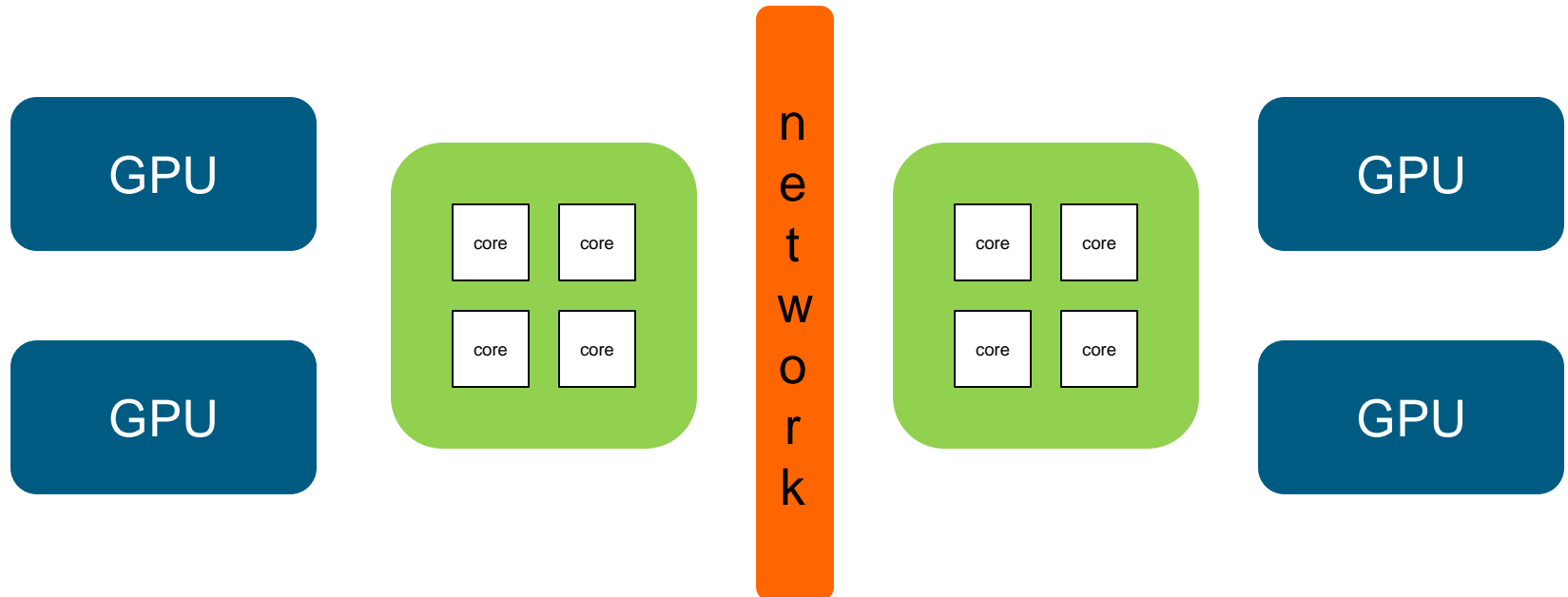# MultiGPU programming

Presenter: Jiri Kraus (NVIDIA)
Suraj Prabhakaran  |  April 9, 2014

German Research School for Simulation Sciences GmbH
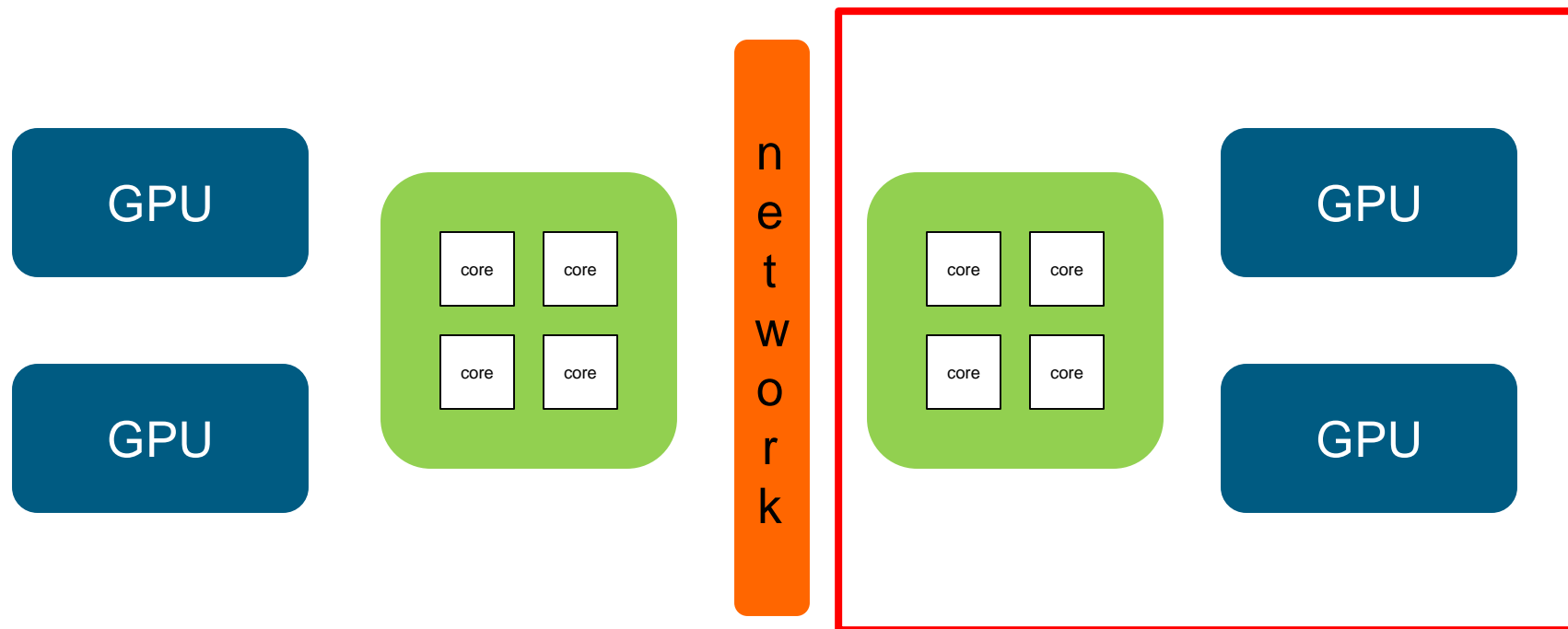Laboratory for Parallel Programming

# Using Multi GPUs

- Further speedup computations
- Single GPU memory not sufficient
- Increases performance/W

- Intra-node Multi-GPU
  - Easy-to-use, directly use the CUDA API

- Inter-node Multi-GPU
  - Network communication with MPI

# Application scenario

# Application scenario

GPU

GPU

core core

core core

network

core core

core core

GPU

GPU

# Intra-node Multi-GPU

- Single CPU thread access Multiple GPUs
- CUDA calls issued to _current_ GPU
- cudaSetDevice(x) sets the current GPU.
- Example

```
cudaSetDevice(0);
cudaMalloc(dst_0,…);
cudaMemcpy(dst_0, …);
cudaSetDevice(1);
cudaMalloc(dst_1,…);
cudaMemcpy(dst_1, …);
```

# Intra-node Multi-GPU

- Current GPU can be changed even when async calls (kernels, async memcopies) are running
- Example

```
cudaSetDevice(0);
kernel<<<…>>>(…);
cudaSetDevice(1);
cudaMemcpyAsync(…);
```

# Multi-GPU Matrix Multiplication

$$C = A \times B$$

Split A and C into two sets of rows

$$C0 = A0 \times B \qquad \text{GPU 0}$$
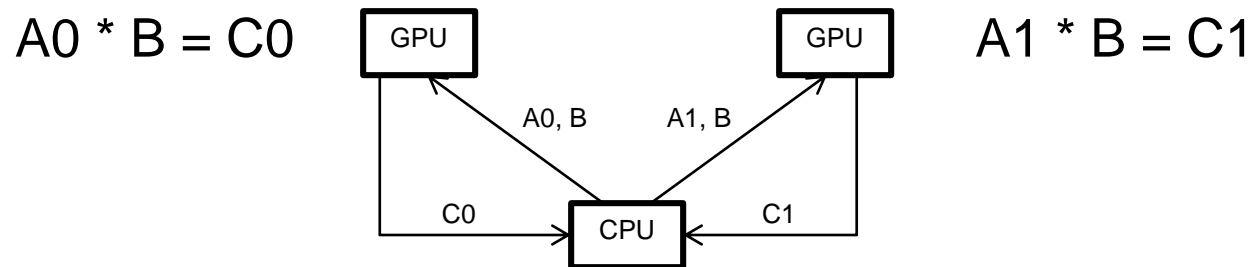
$$C1 = A1 \times B \qquad \text{GPU 1}$$

# Exercise: Multi-GPU Matrix Multiplication

- Use Multiple GPUs to speed up Simple Matrix Multiplication

- Split A into 2 set of rows

$A0 * B = C0$  [GPU]  ←A0, B    A1, B→  [GPU]  $A1 * B = C1$
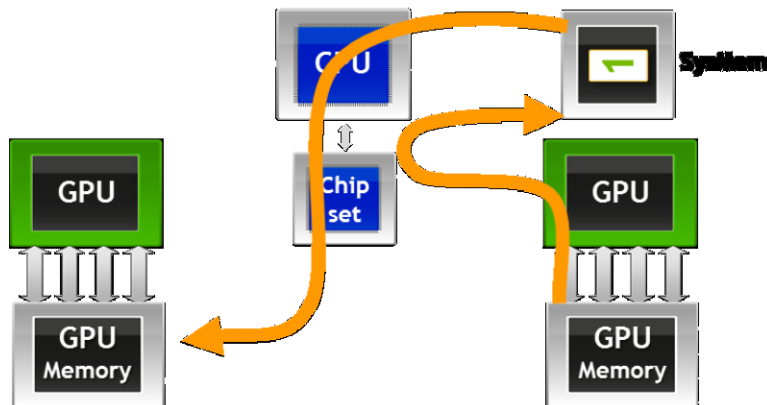
C0  [CPU]  C1

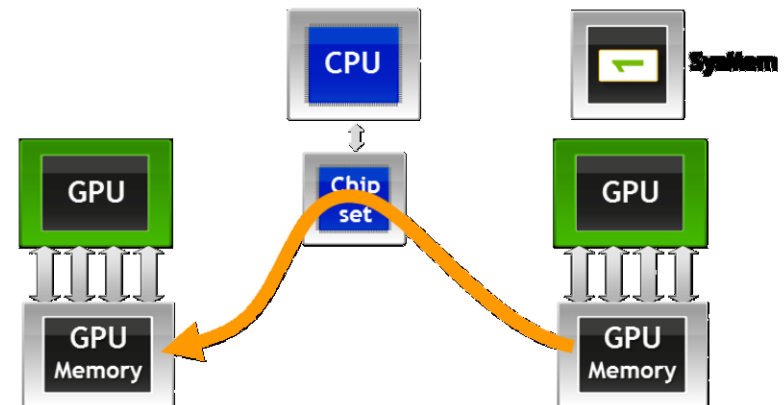- Verify with NVVP that two GPUs are used.

# Intra-node Multi-GPU Communication

- One GPU has to access data from another GPU
- Traditional method: Go about it through the CPU/Main Memory
- Due to UVA: Peer-to-peer memcopies (GPUDirect P2P)

**No GPUDirect P2P**

**GPUDirect P2P**

# Intra-node Multi-GPU Communication

- Check if the GPU can access Peer device

    cudaDeviceCanAccessPeer(&accessible, dev_x, dev_y);

- First enable Peer-to-peer communication

    cudaDeviceEnablePeerAccess(peer_device,0);

- Transfer data between two devices

    cudaMemcpy(dst, src, size, cudaMemcpyDeviceToDevice);

    - *Also works if peer access is not possible or not enabled (fall back with host memory staging)*

# **Exercise:**

- Compare memcopies between two devices using:
  - Manual staging through main memory
  - Using GPUDirect Peer to Peer