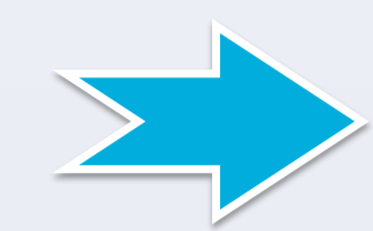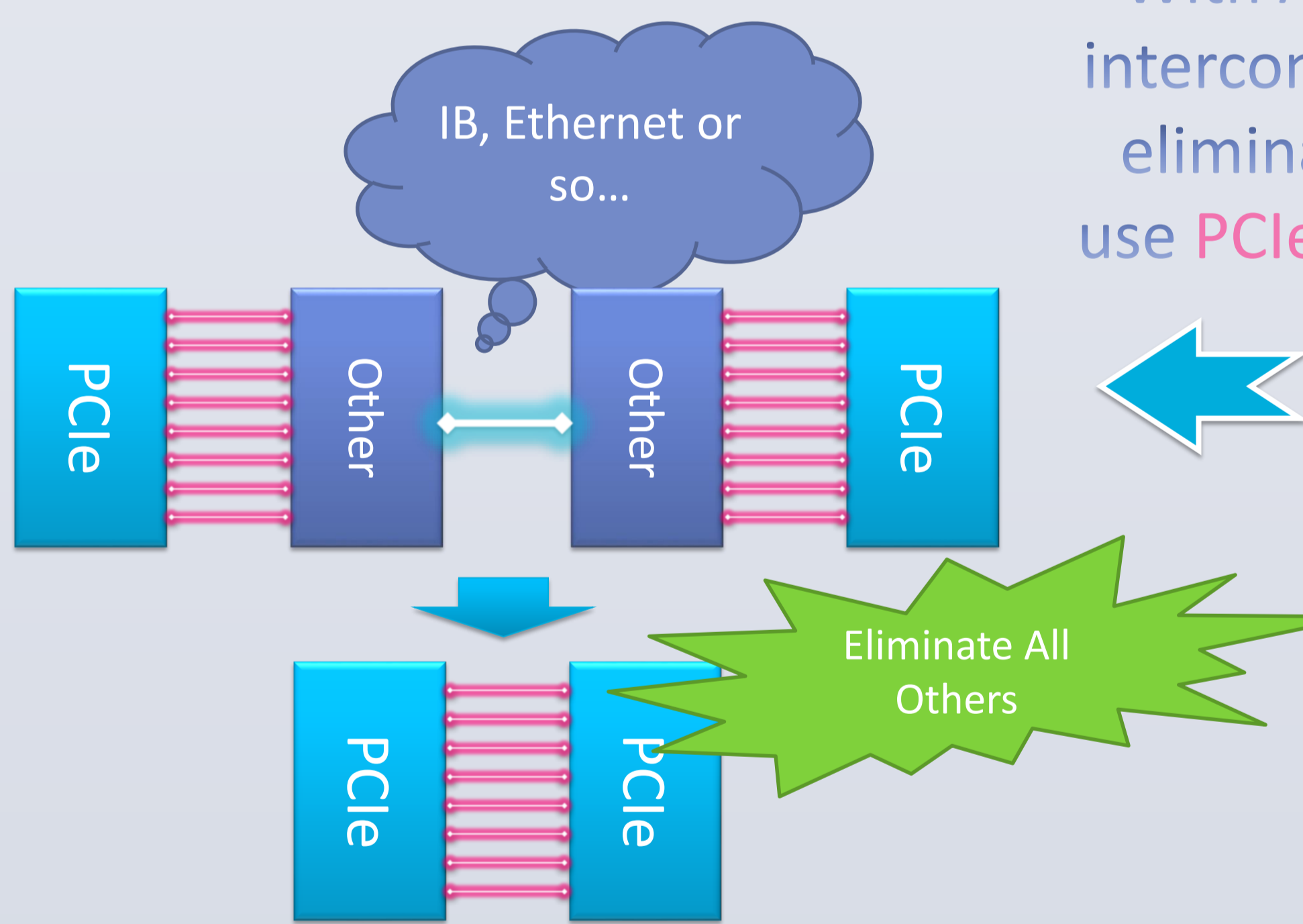# V-OpenCL: A method to use remote GPGPU

Cong Wang, Tao Jiang, Rui Hou

Institute of Computing Technology, Chinese Academy of Sciences

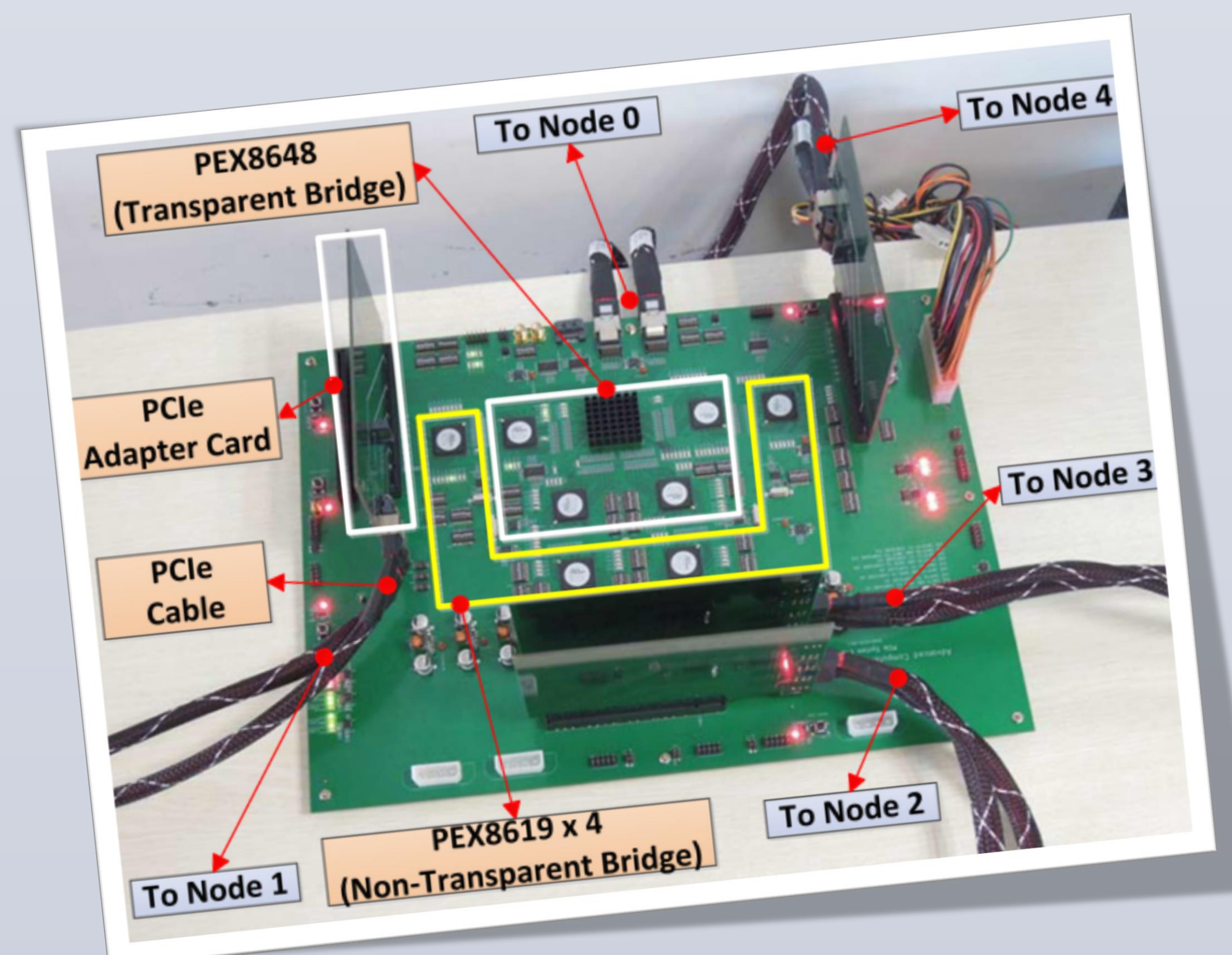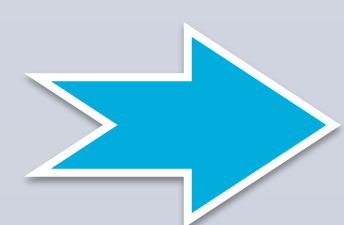In datacenter GPGPUs are useful but not needed all the time.

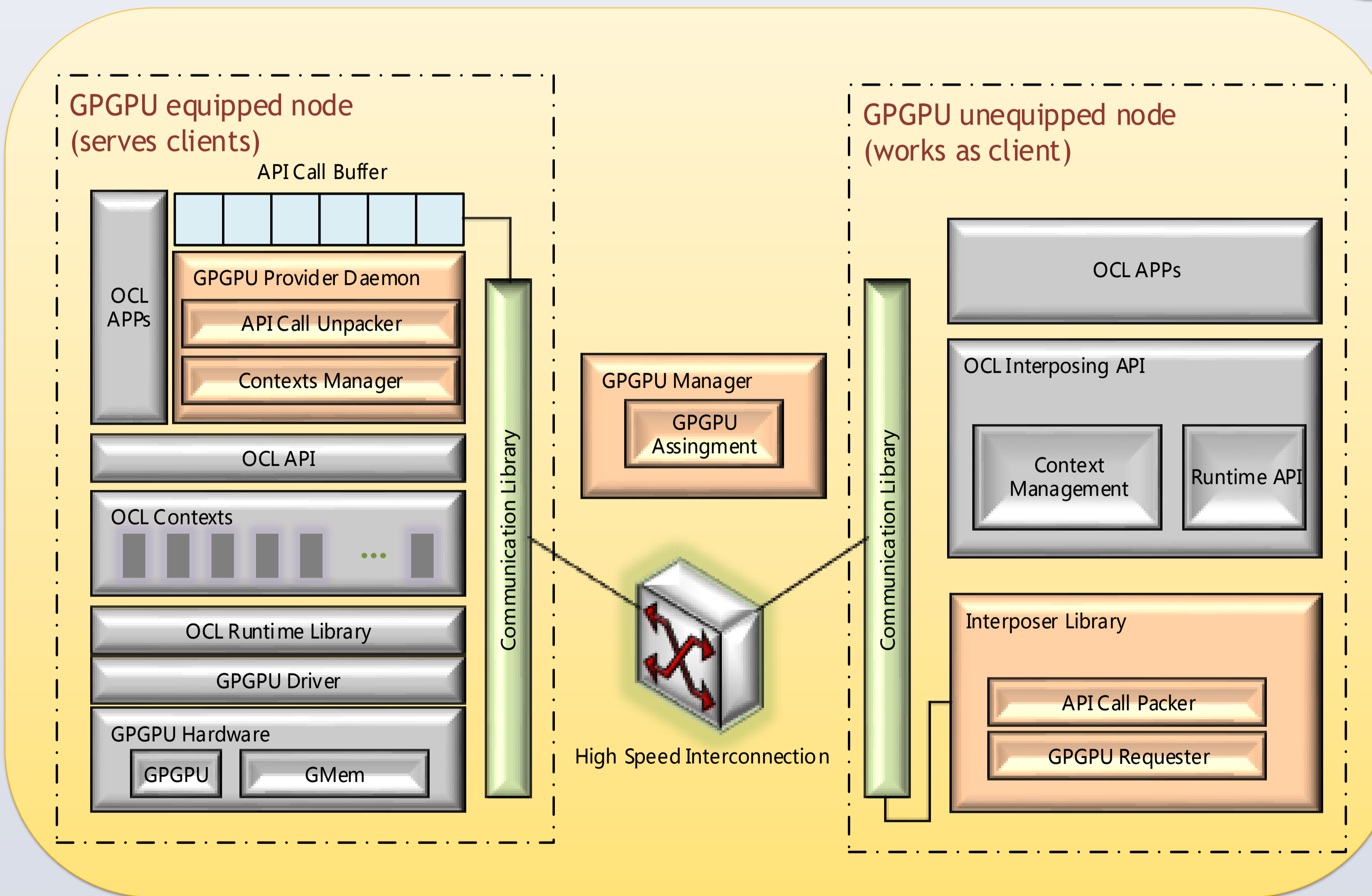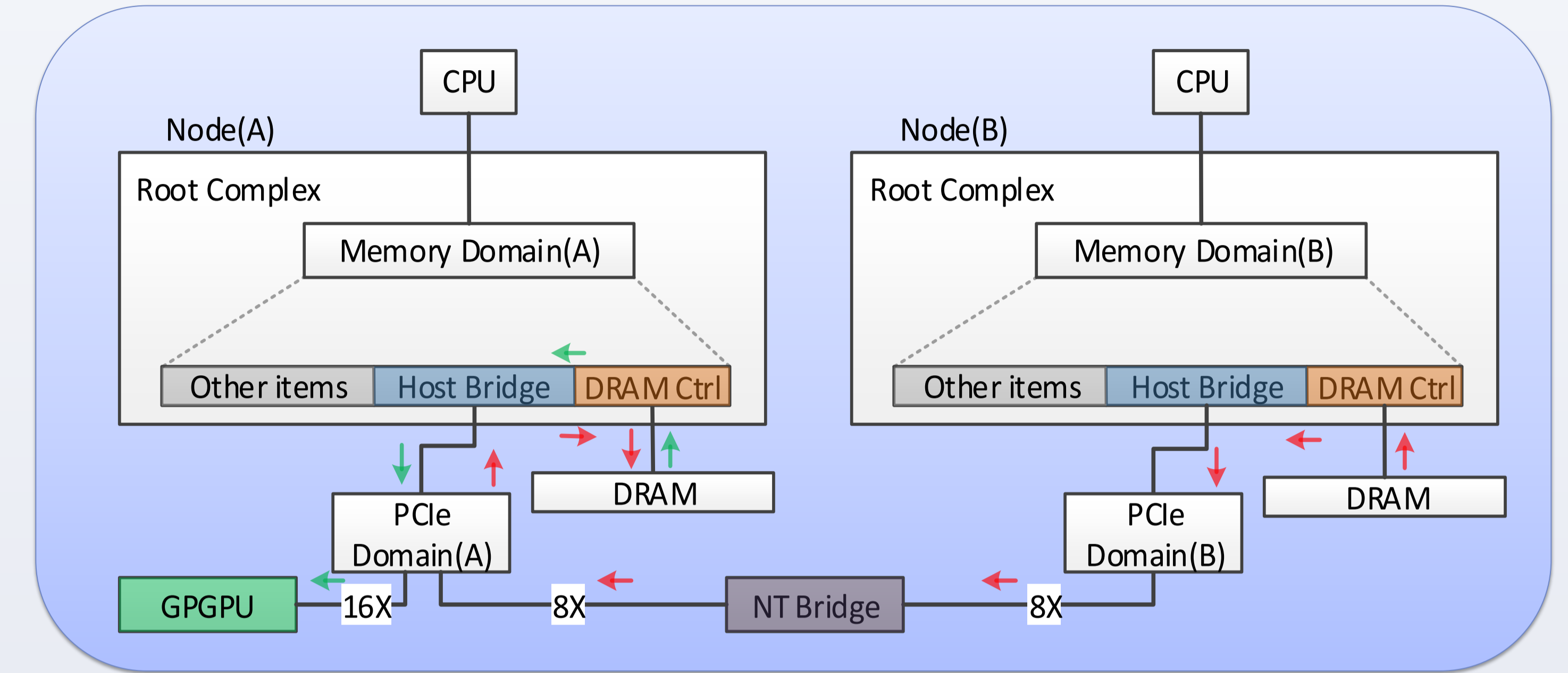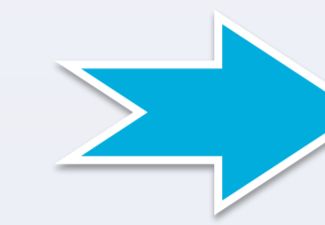Partially equip nodes with GPGPU and share them between nodes.

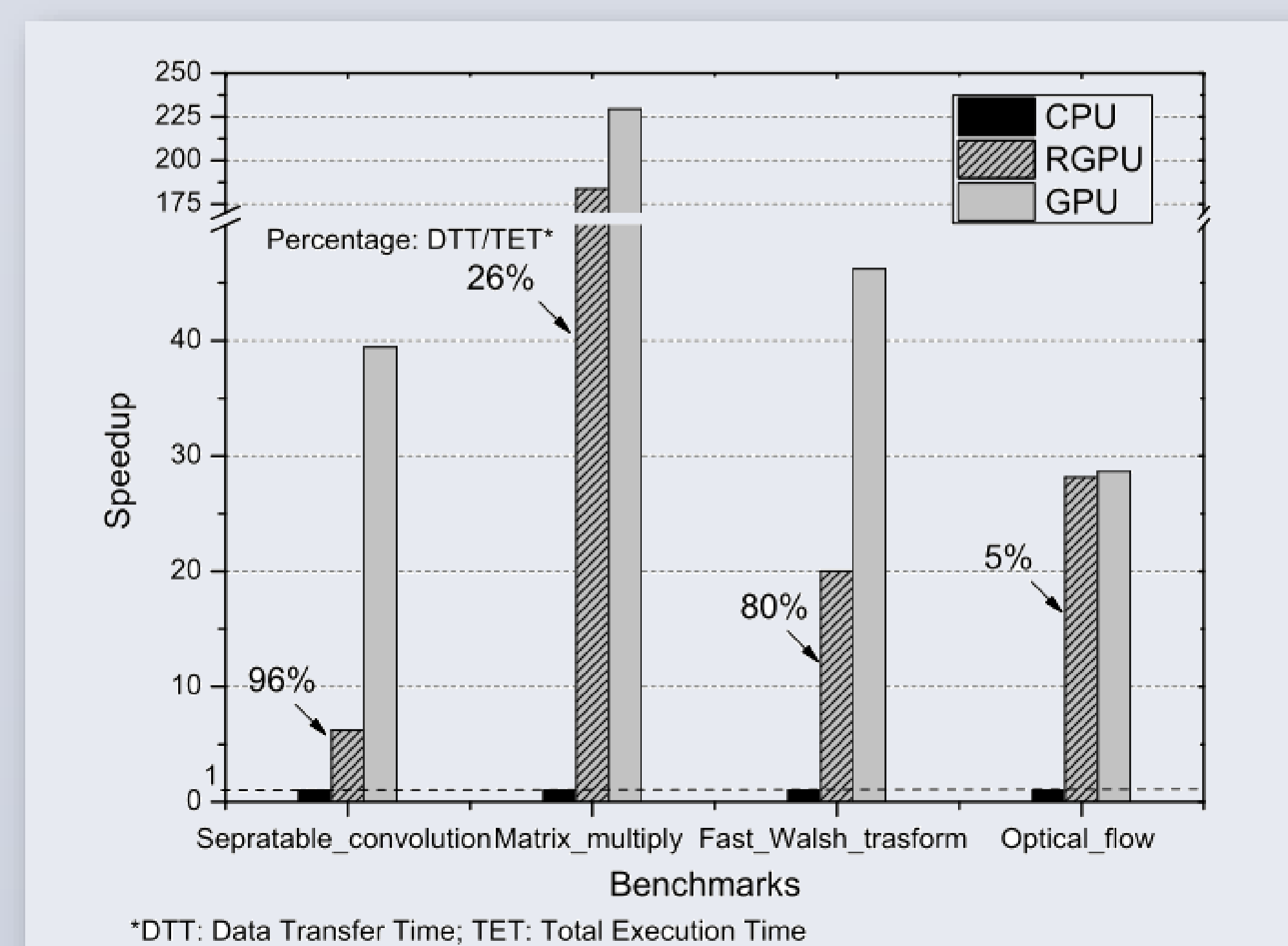With All other interconnections eliminated we use PCIe Directly.

This is our primitive PCIe switch board. Server nodes connect to it with PCIe cables.

With PCIe, memory can be mapped to different nodes. Data and binary are first copied from DRAM of Node B to Node A and then sent to GPGPU



- OpenCL API calls are packed up and redirected to GPGPU equipped node(s).
- An API call buffer is set for each GPGPU equipped node to hold the incoming call packets and sort them by context ID.
- A GPGPU manager software is set to track the states of GPGPUs and help nodes to find GPGPU(s).
- This mechanism is suitable for common interconnections, not only PCIe.

Four sample programs are used for evaluation. All of them run much faster on remote GPGPU than on CPU although data transfer takes time. Of the over all run time, transfer time takes large portion in two of the programs while takes very little portion in the other two.