



Introduction to Systems / HPC

Abhinav Bhatele, Daniel Nichols



UNIVERSITY OF
MARYLAND

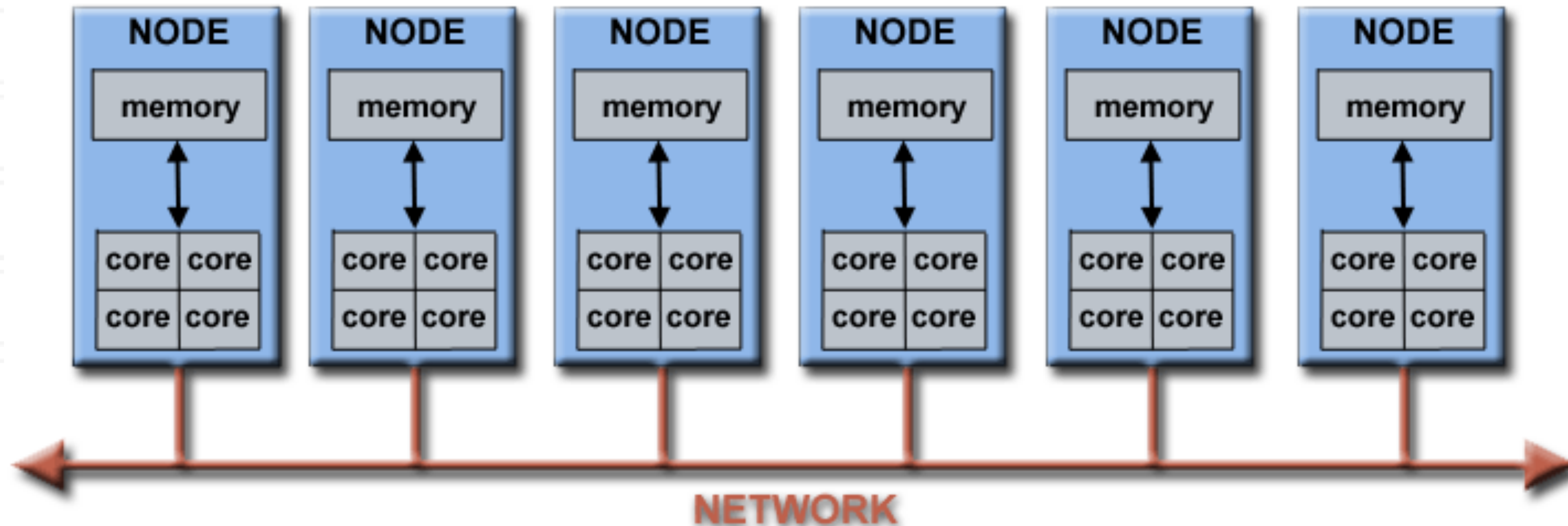
Getting started with zaratan

- Over 360 nodes with AMD Milan processors (128 cores/node, 512 GB memory/node)
- 20 nodes with four NVIDIA A100 GPUs (40 GB per GPU)
- 8 nodes with four NVIDIA H100 GPUs (80 GB per GPU)

```
ssh username@login.zaratan.umd.edu
```

Data center / HPC cluster

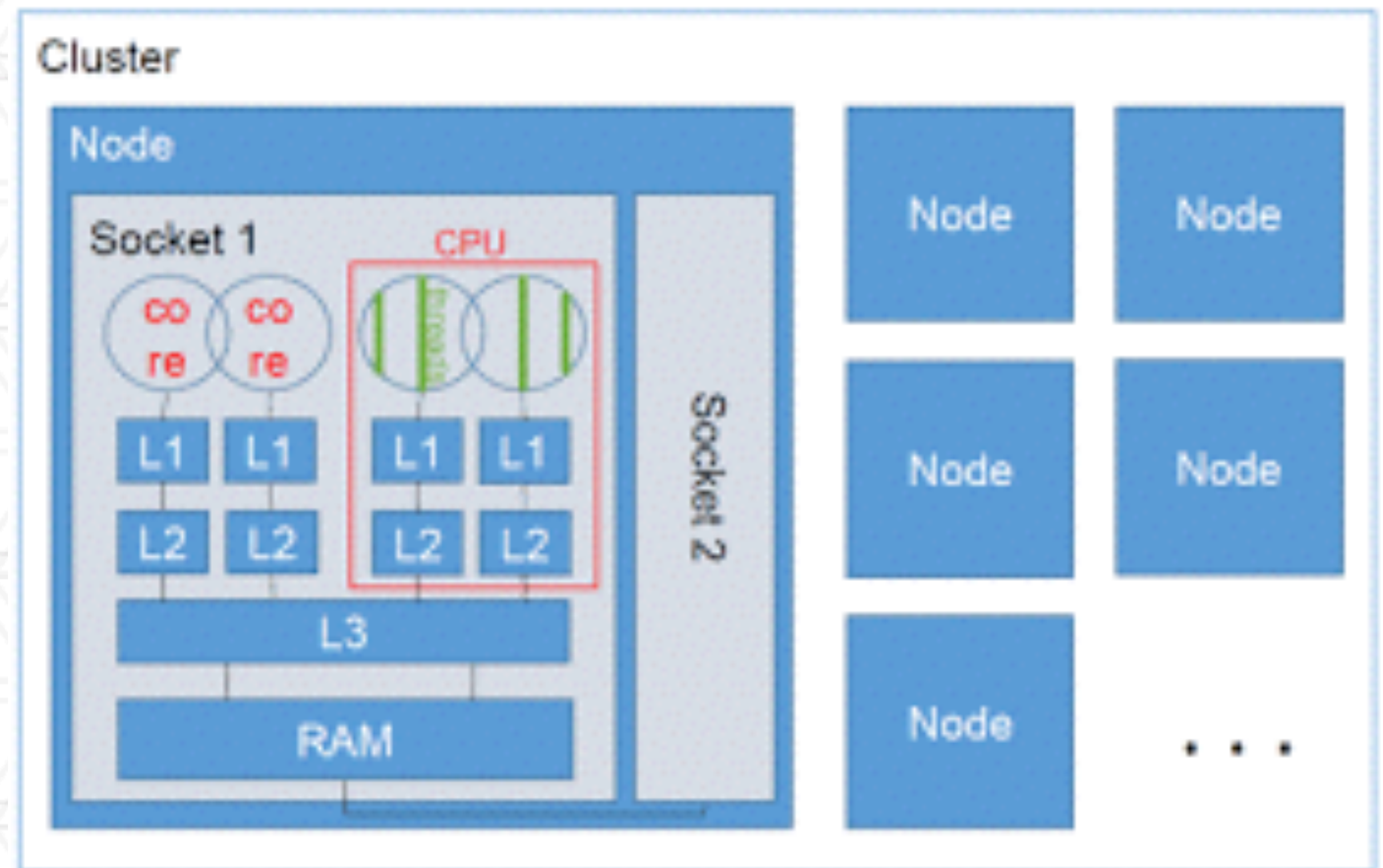
- A set of nodes or processing elements connected by a network.
- Compute node: A shared-memory unit (optionally has GPUs)



https://computing.llnl.gov/tutorials/parallel_comp

Cores, sockets, nodes

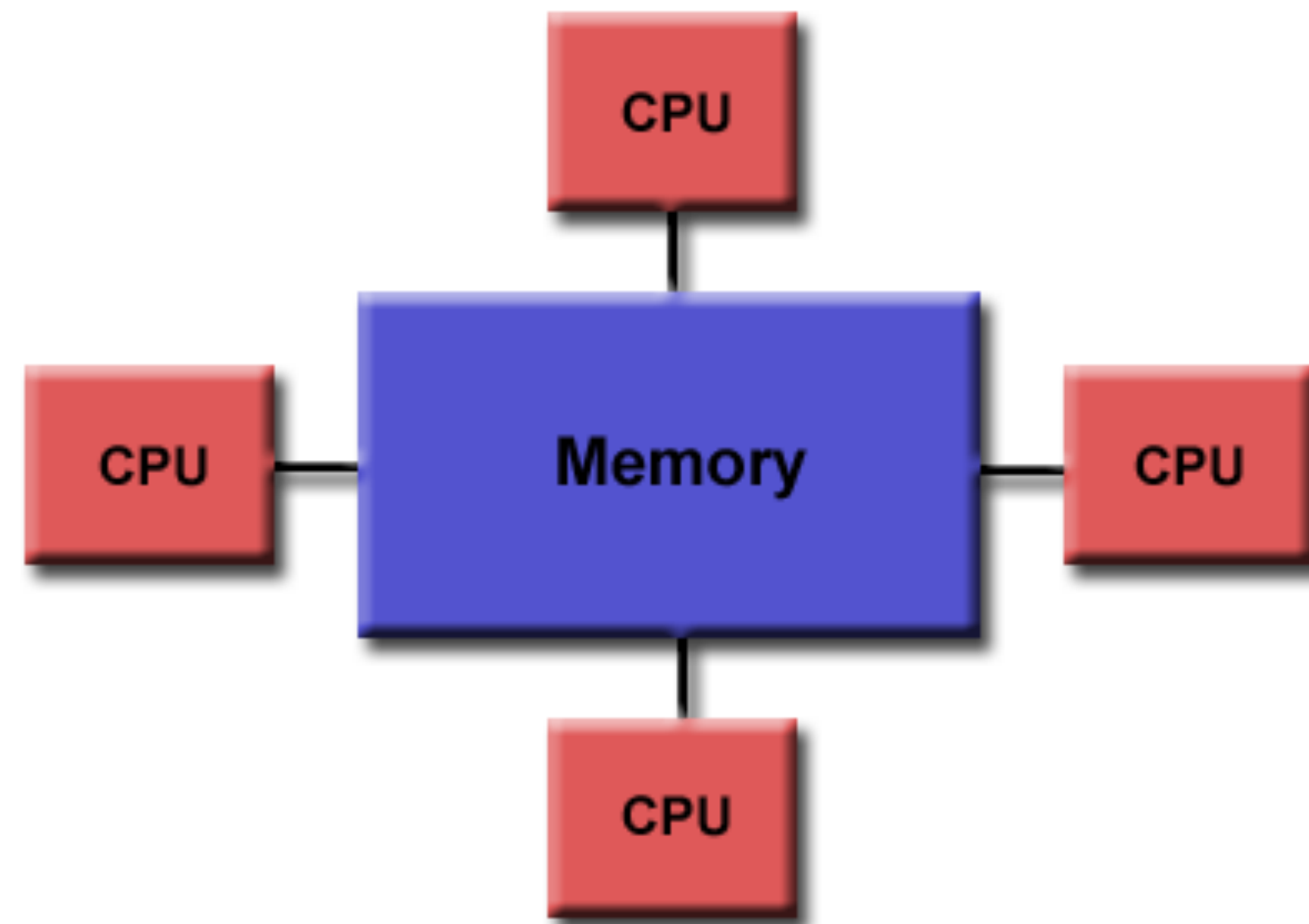
- Core: a single execution unit that has a private L1 cache and can execute instructions independently
- Processor: several cores on a single Integrated Circuit (IC) or chip are called a multi-core processor
- Socket: physical connector into which an IC/chip or processor is inserted.
- Node: a packaging of sockets — motherboard or printed circuit board (PCB) that has multiple sockets



<https://hpc-wiki.info/hpc/HPC-Dictionary>

Shared memory architecture

- All processors/cores can access all memory as a single address space

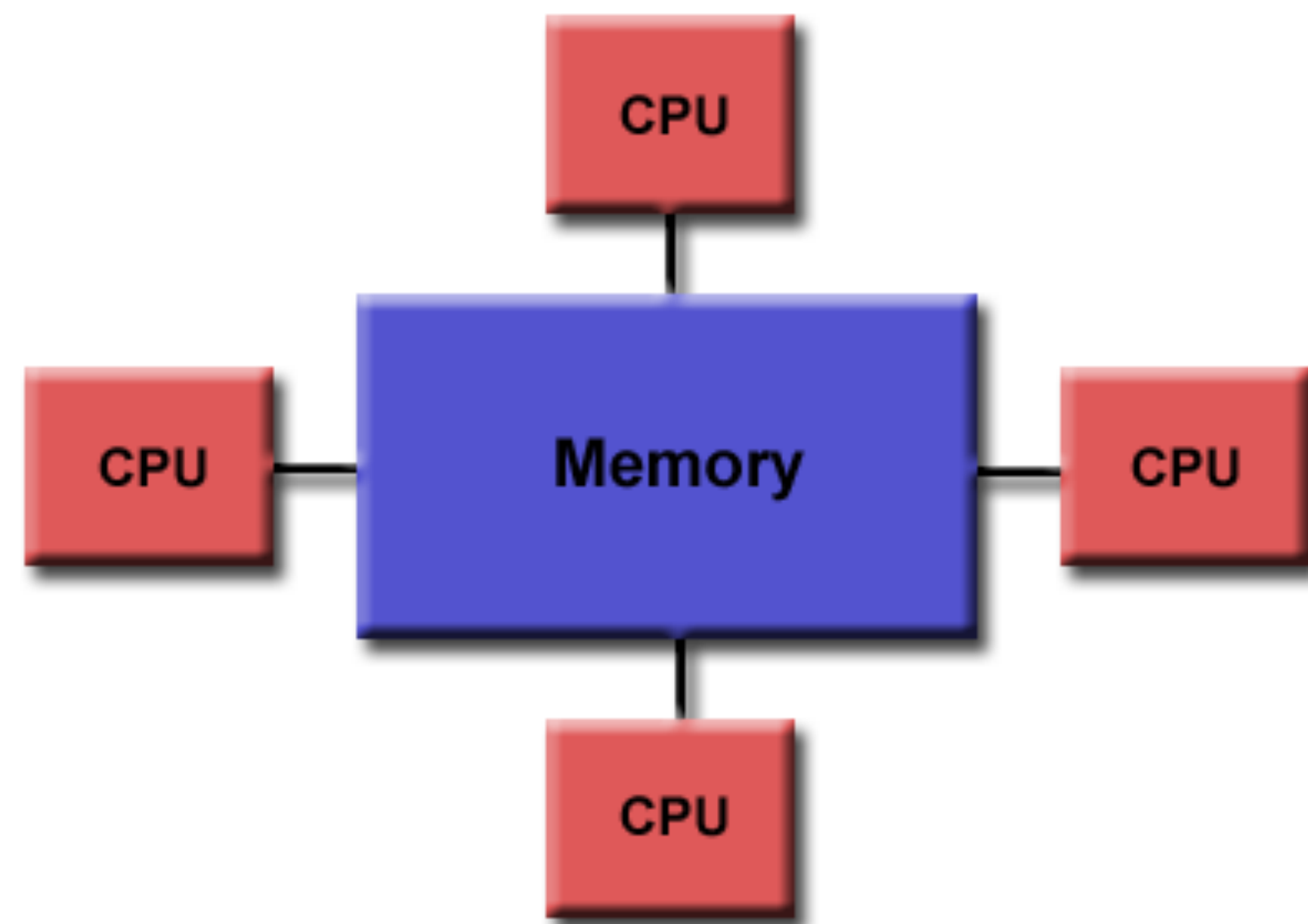


Uniform Memory Access

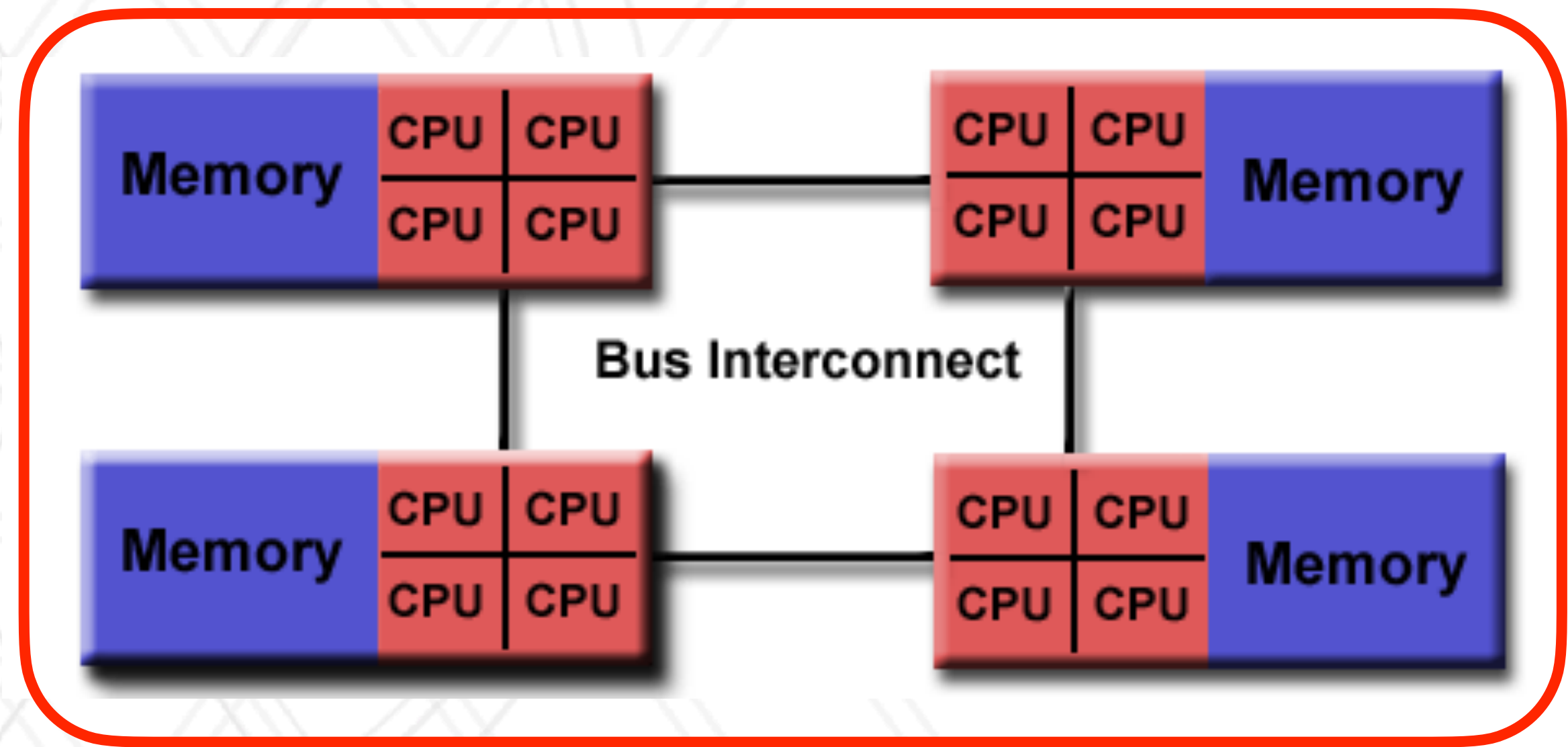
https://computing.llnl.gov/tutorials/parallel_comp/#SharedMemory

Shared memory architecture

- All processors/cores can access all memory as a single address space



Uniform Memory Access



Non-uniform Memory Access (NUMA)

https://computing.llnl.gov/tutorials/parallel_comp/#SharedMemory

Hopper H100 SM

- **CUDA Core**
 - Single serial execution unit
- **Each H100 Streaming Multiprocessor (SM) has:**
 - 128 FP32 cores
 - 64 INT32 cores
 - 64 FP64 cores
 - 84 Tensor cores
- **CUDA capable device or GPU**
 - Collection of SMs

Hopper H100 SM

- CUDA Core
 - Single serial execution unit
- Each H100 Streaming Multiprocessor (SM) has:
 - 128 FP32 cores
 - 64 INT32 cores
 - 64 FP64 cores
 - 84 Tensor cores
- CUDA capable device or GPU
 - Collection of SMs

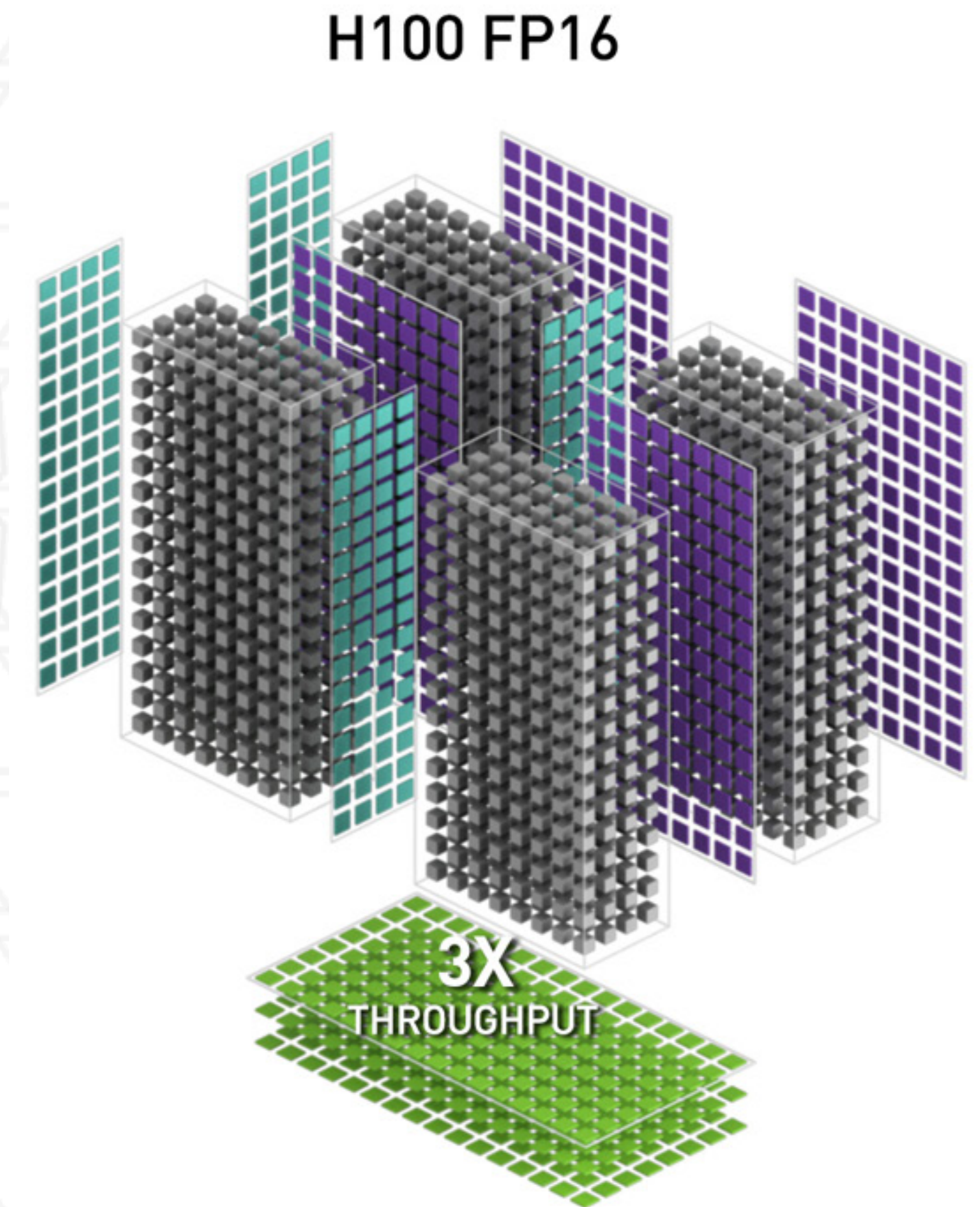


NVIDIA H100 chip



H100 tensor cores

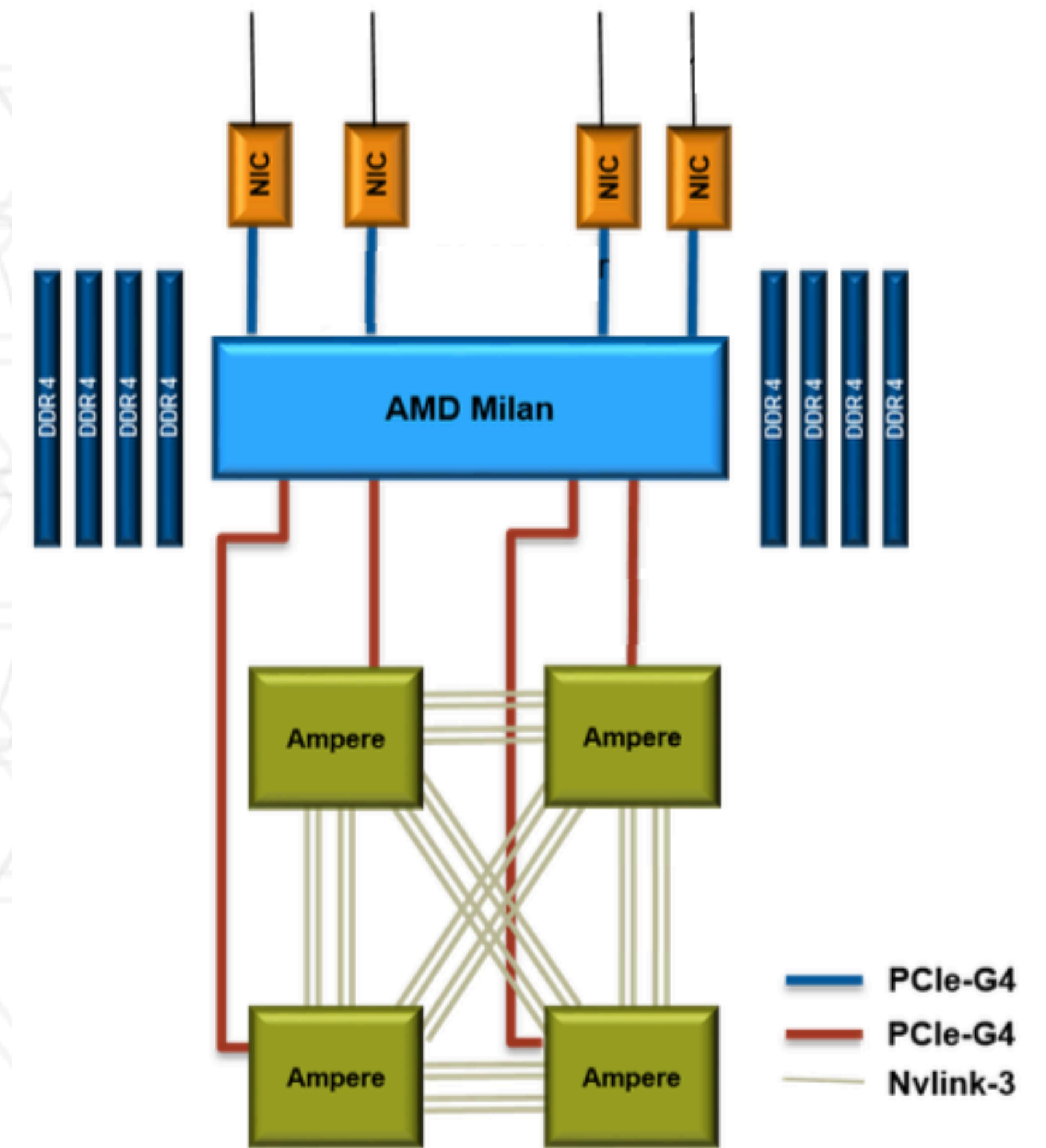
- Tensor cores are specialized cores for matrix multiply accumulate operations
- Operate in parallel across all SMs
- Multiply two 4×4 FP16 matrices and add to a 4×4 FP16 or FP32 matrix
- Mixed precision



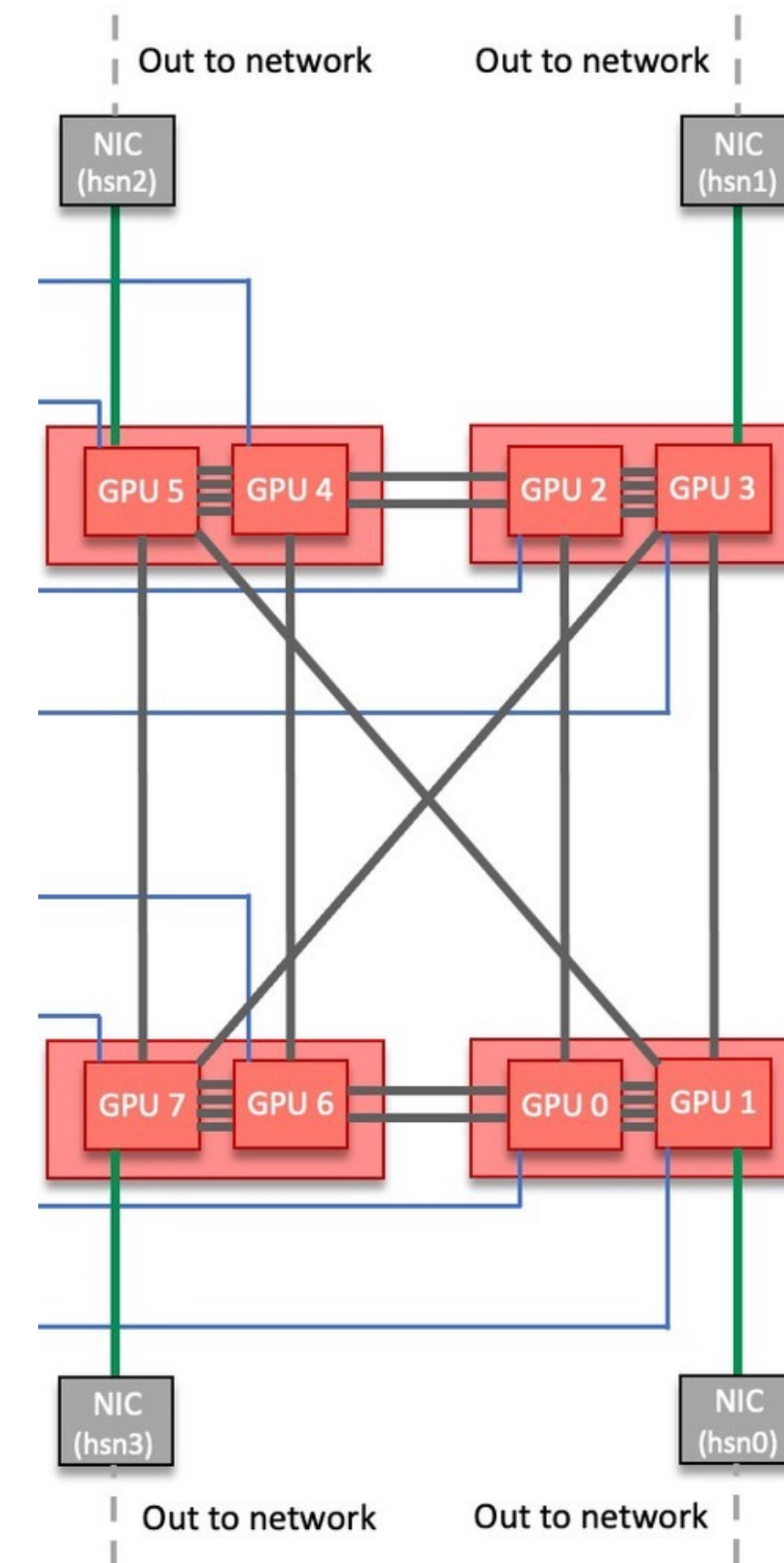
<https://resources.nvidia.com/en-us-tensor-core>

Nodes with GPUs

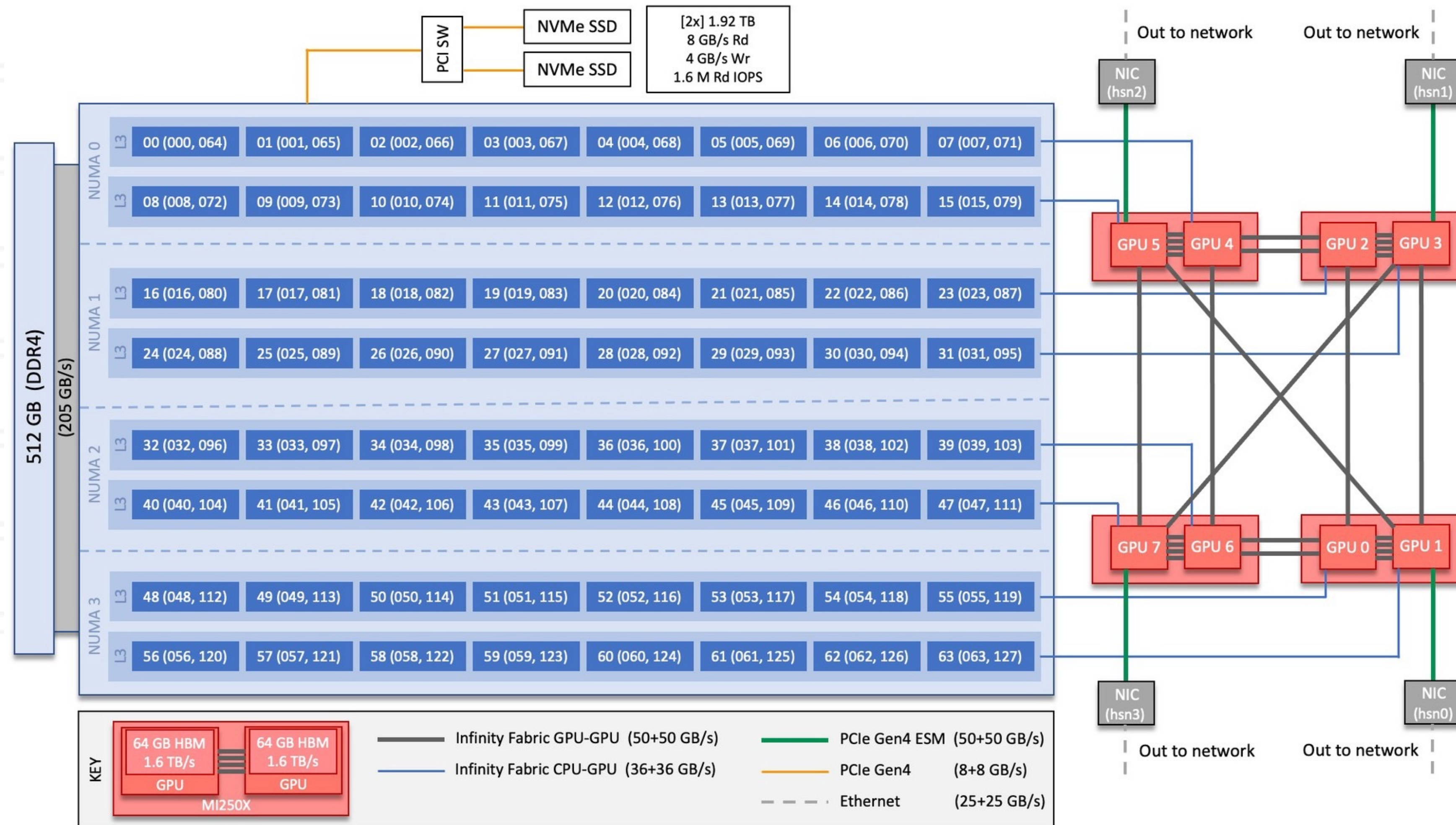
- NIC: Network interface card that connects the node to the network
- PCIe: high-speed interface often used to connect CPUs and GPUs
- NVLink: NVIDIA's high-speed interface often used between GPUs



Alternative node diagram

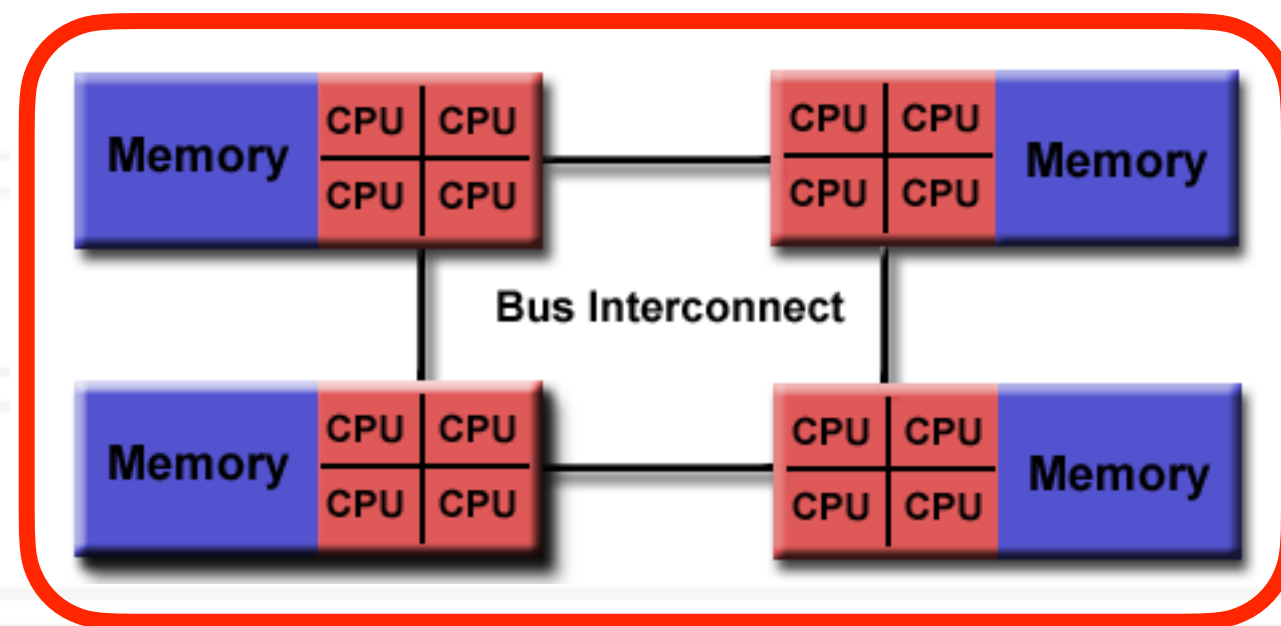


Alternative node diagram

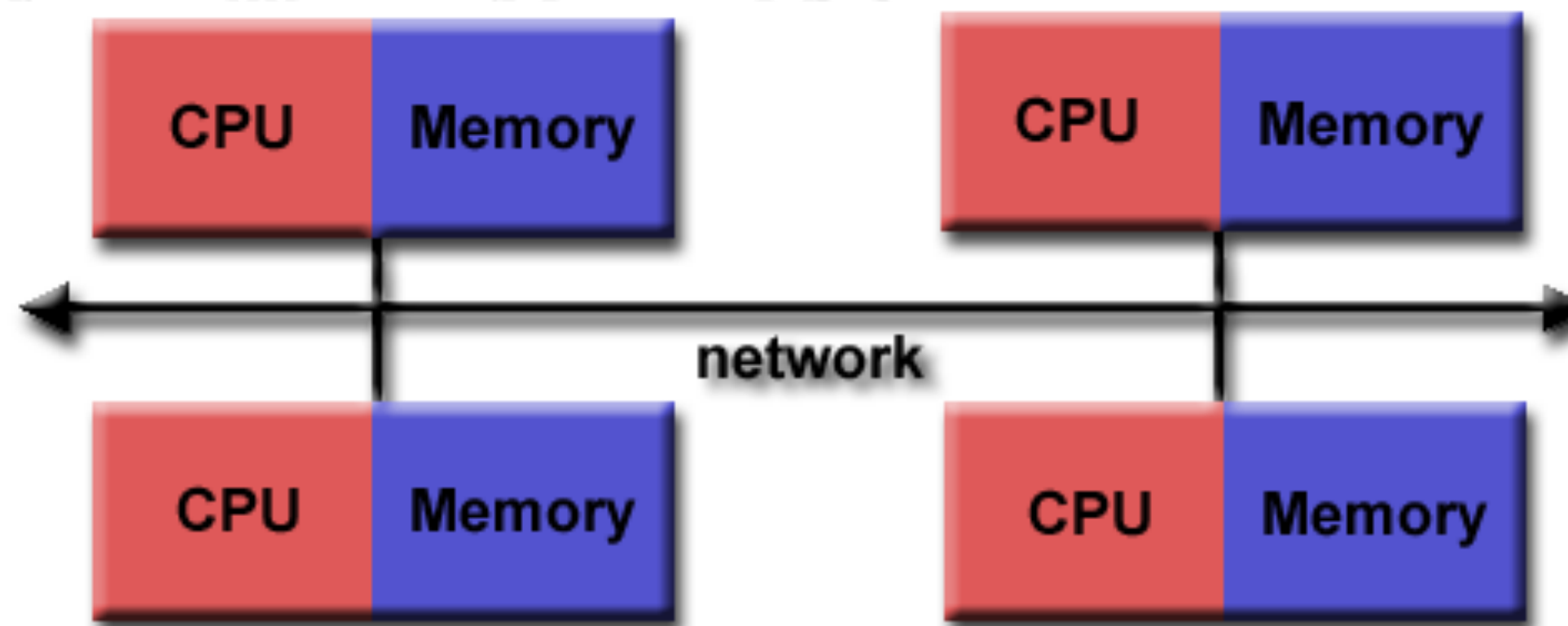


Distributed memory architecture

- Groups of processors/cores have access to their local memory
- Writes in one group's memory have no effect on another group's memory



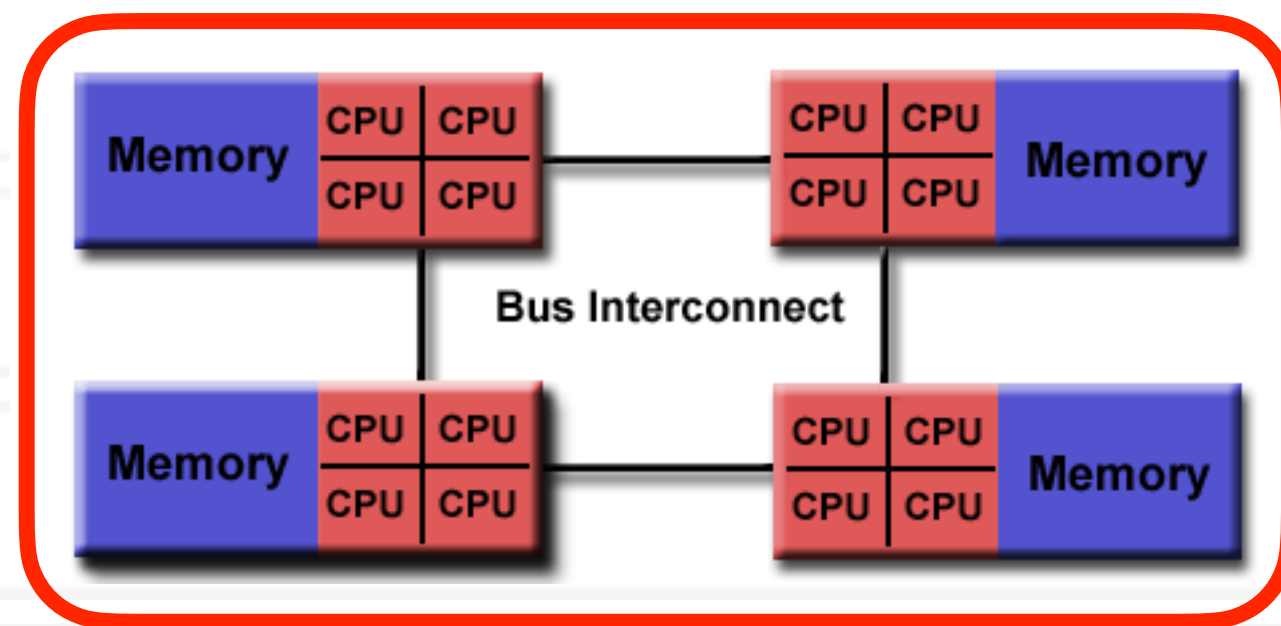
Shared memory (NUMA)



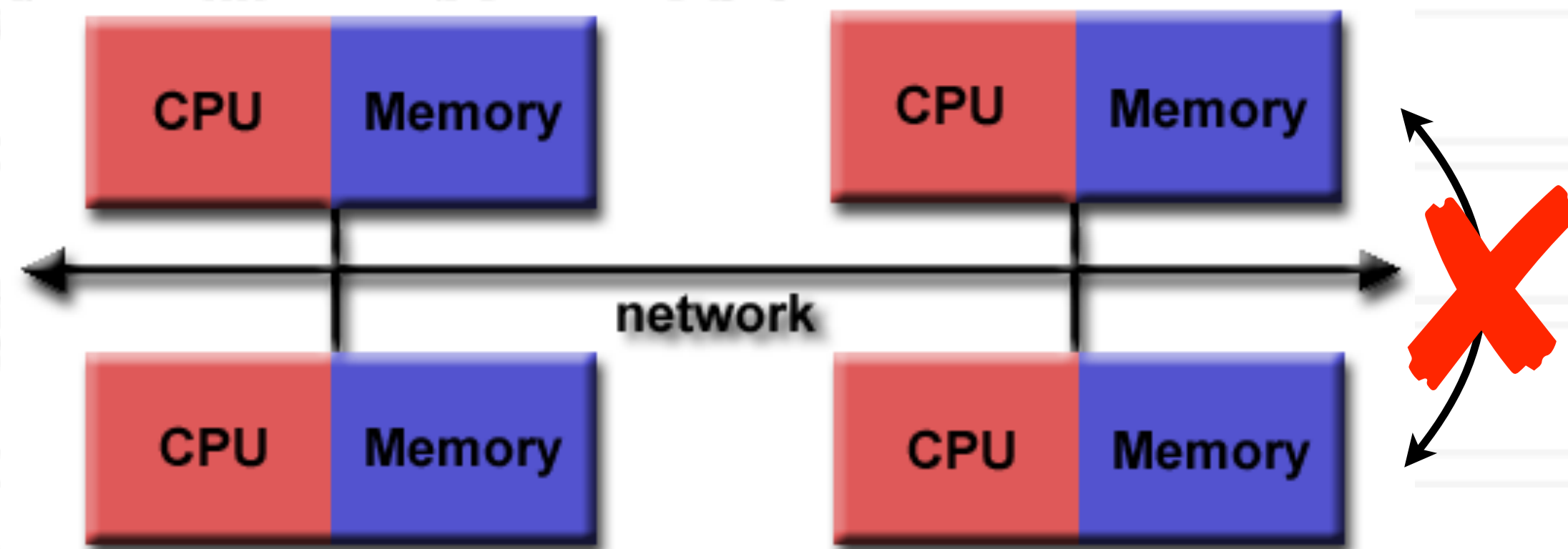
Distributed memory

Distributed memory architecture

- Groups of processors/cores have access to their local memory
- Writes in one group's memory have no effect on another group's memory



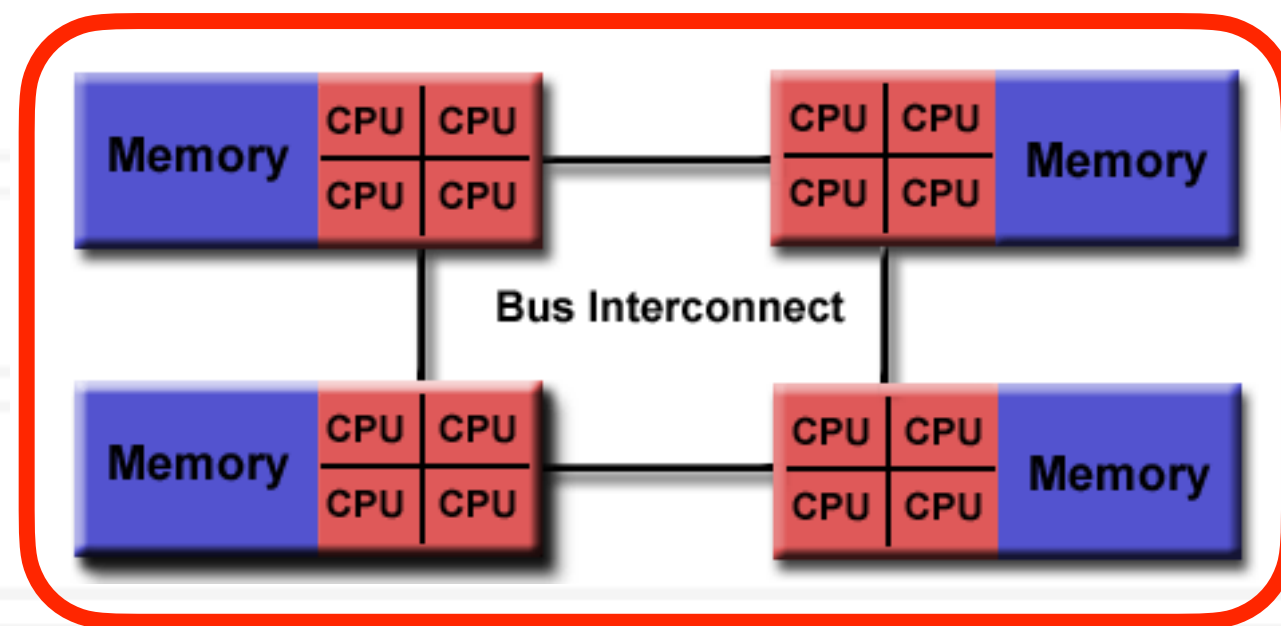
Shared memory (NUMA)



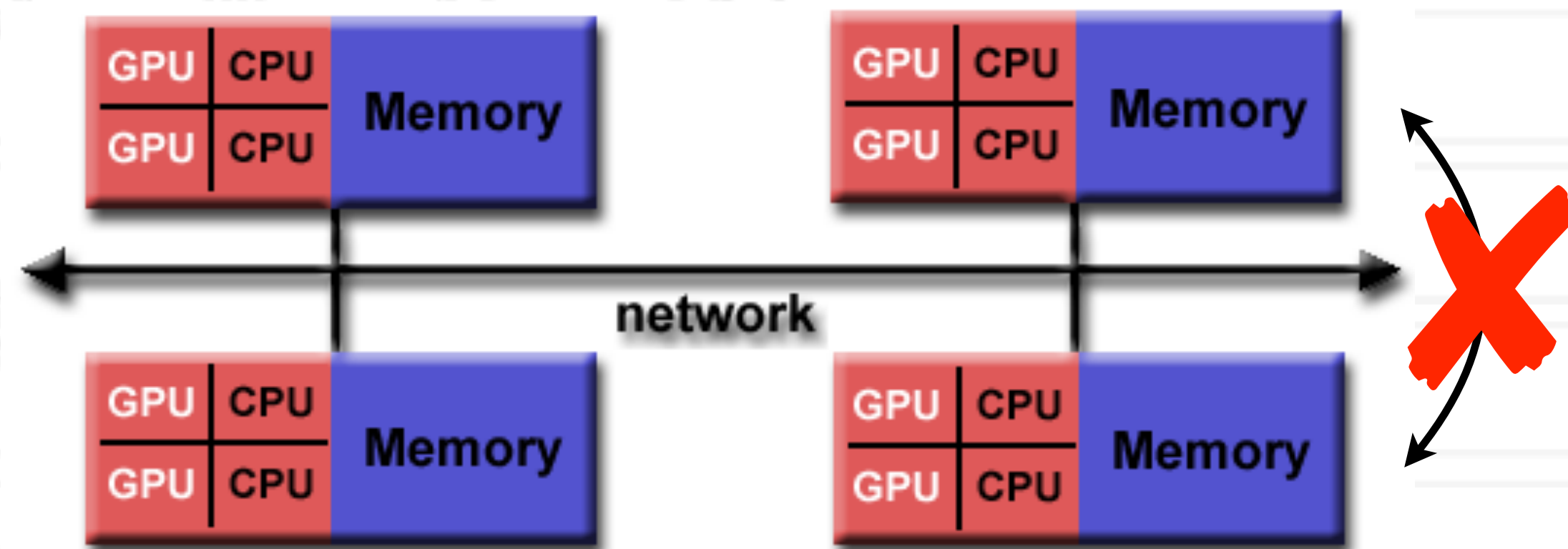
Distributed memory

Distributed memory architecture

- Groups of processors/cores have access to their local memory
- Writes in one group's memory have no effect on another group's memory

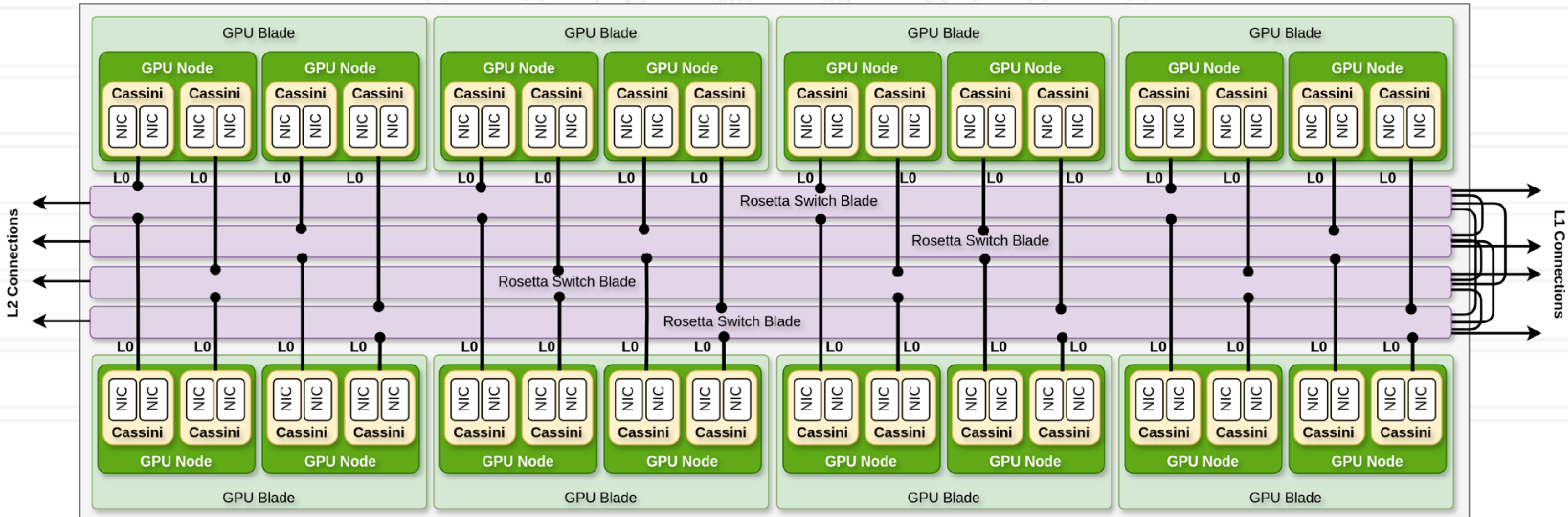


Shared memory (NUMA)



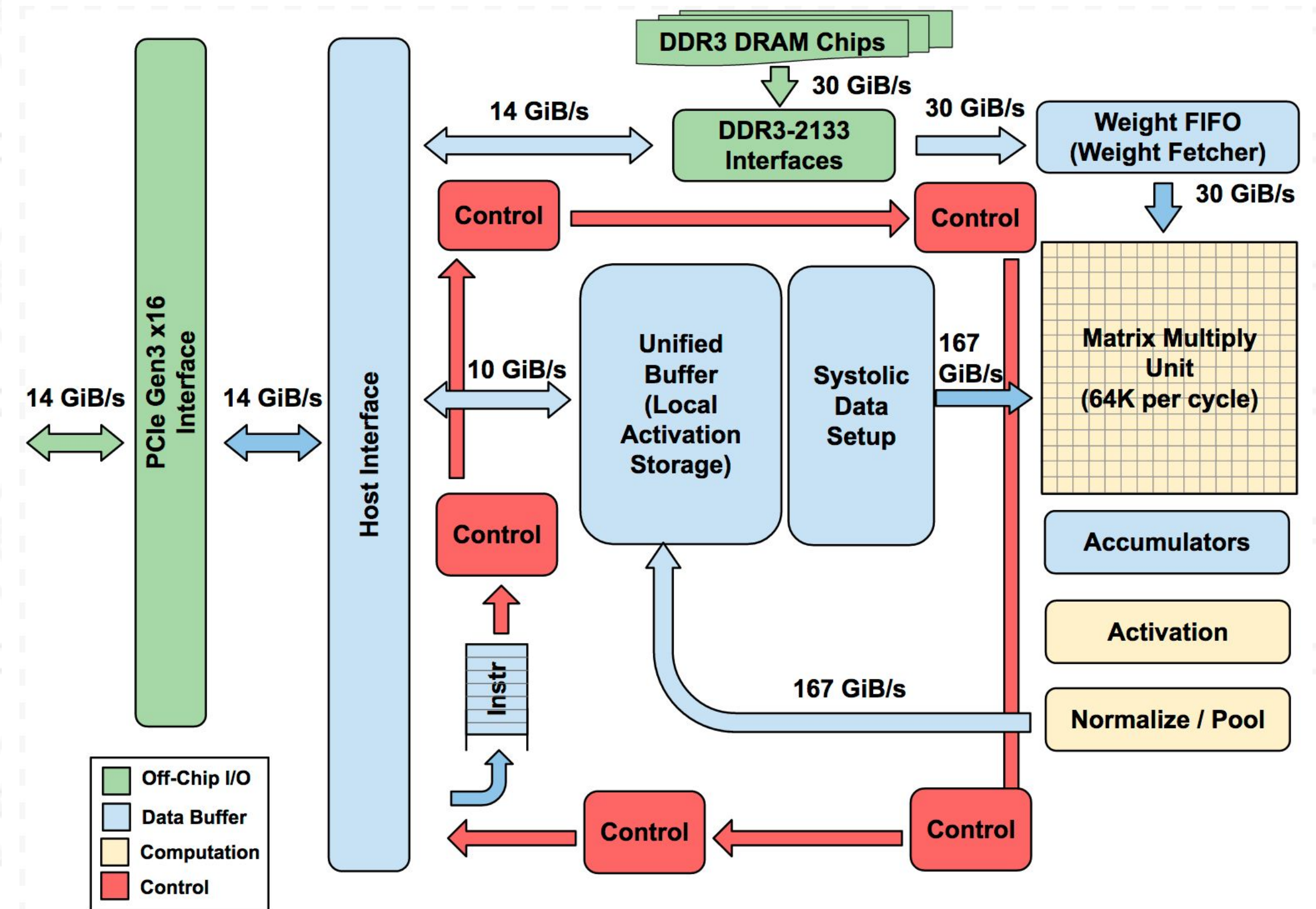
Distributed memory

A realistic cluster



Google's Tensor Processing Unit

- TPU is an ASIC (Application-specific Integrated Circuit)
- Co-processor just like GPUs
- Each TPU can have one or multiple MMUs
- TPU Pod is a collection of TPUs



Network components

- Network interface controller or card
- Router or switch
- Network cables: copper or optical



Life-cycle of a message

Source

Source

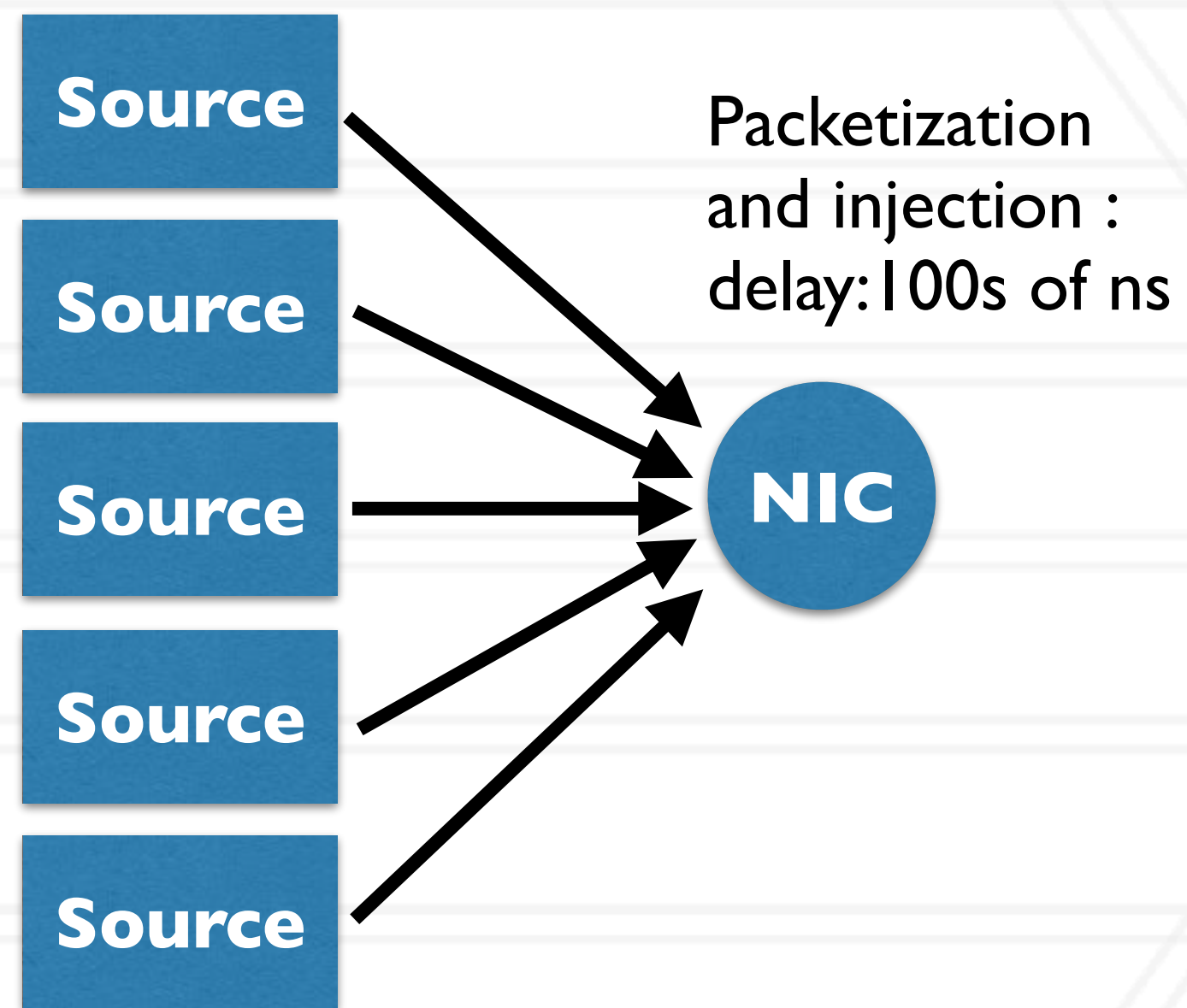
Source

Source

Source

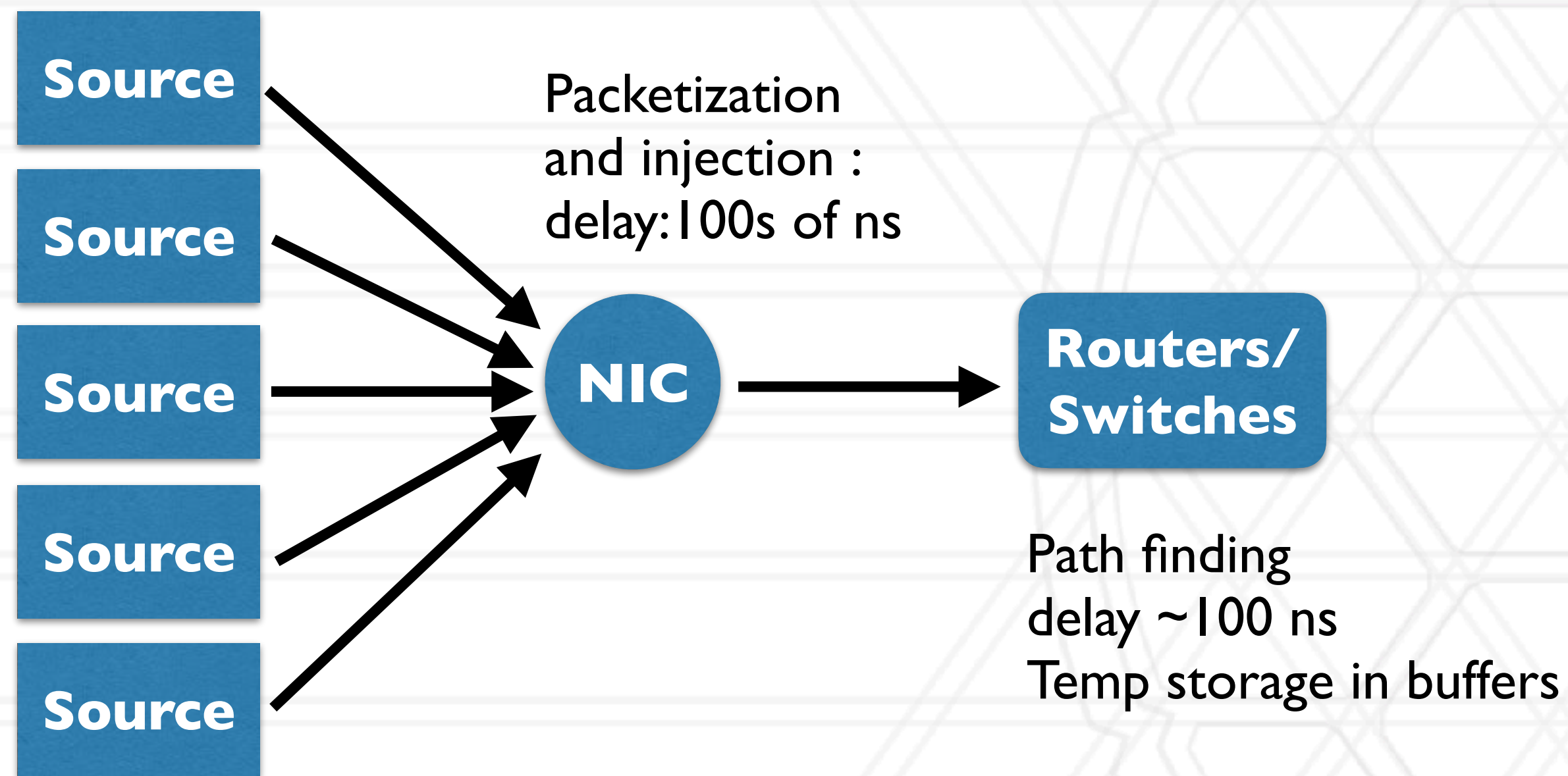
Message origin points :
destination, frequency,
size, etc. determined
by application
1 micro sec - 10s of sec

Life-cycle of a message



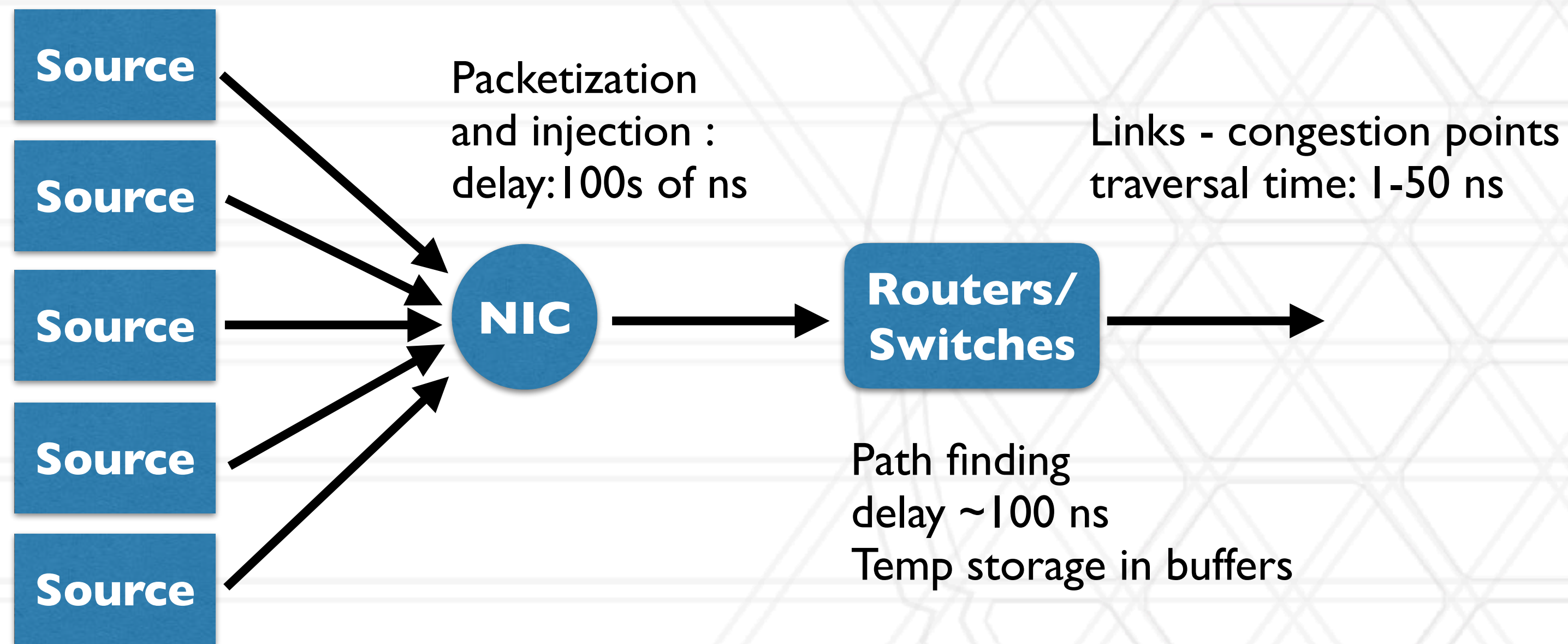
Message origin points :
destination, frequency,
size, etc. determined
by application
1 micro sec - 10s of sec

Life-cycle of a message



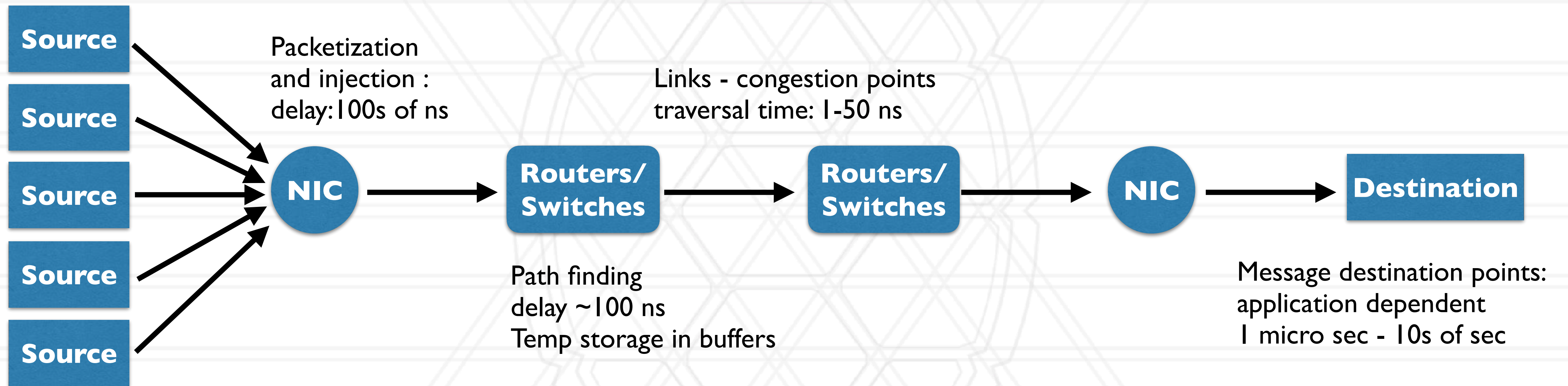
Message origin points :
destination, frequency,
size, etc. determined
by application
1 micro sec - 10s of sec

Life-cycle of a message



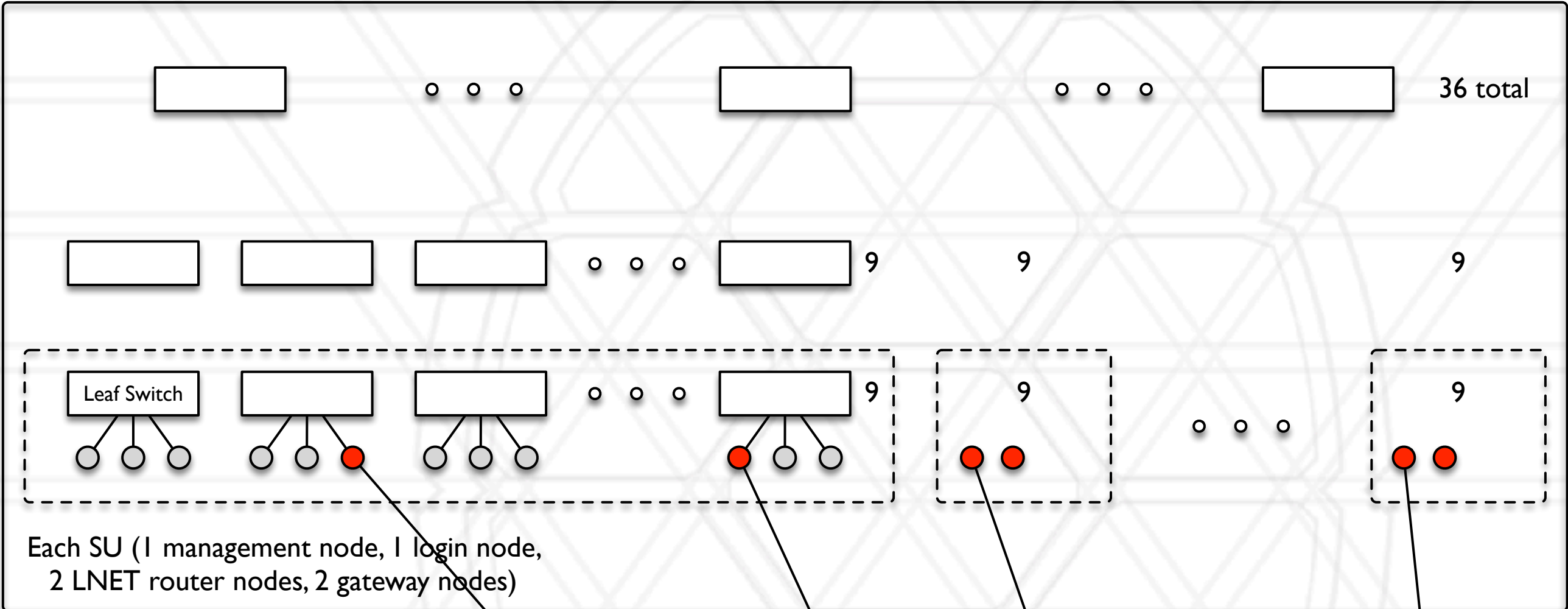
Message origin points :
destination, frequency,
size, etc. determined
by application
1 micro sec - 10s of sec

Life-cycle of a message

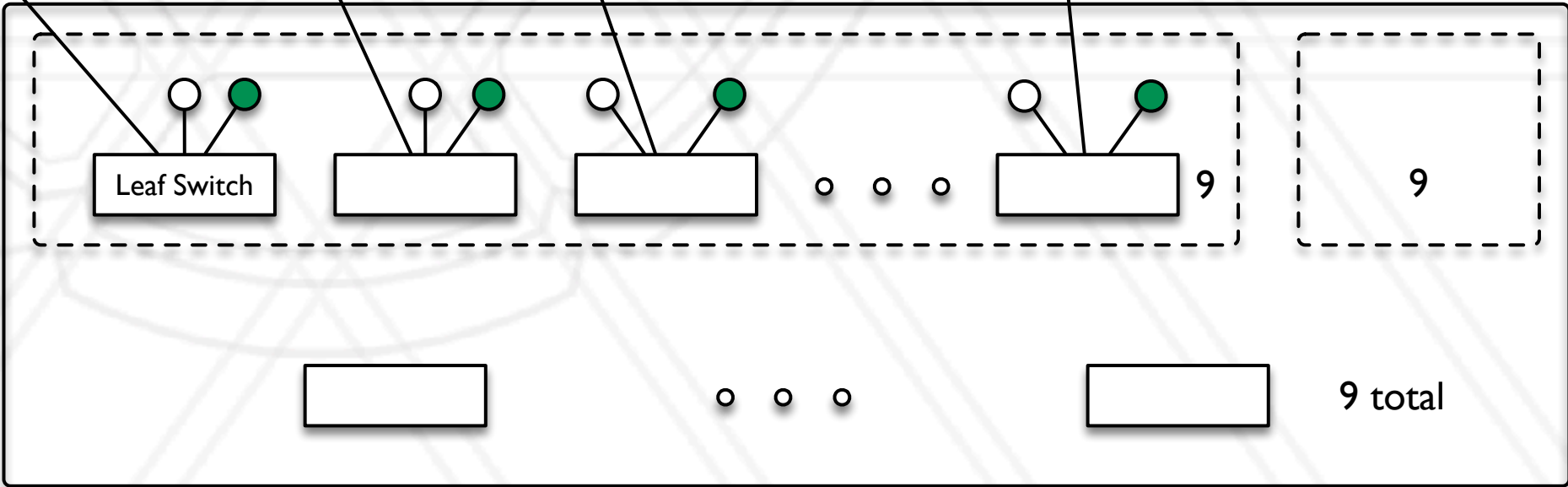


Message origin points :
destination, frequency,
size, etc. determined
by application
1 micro sec - 10s of sec

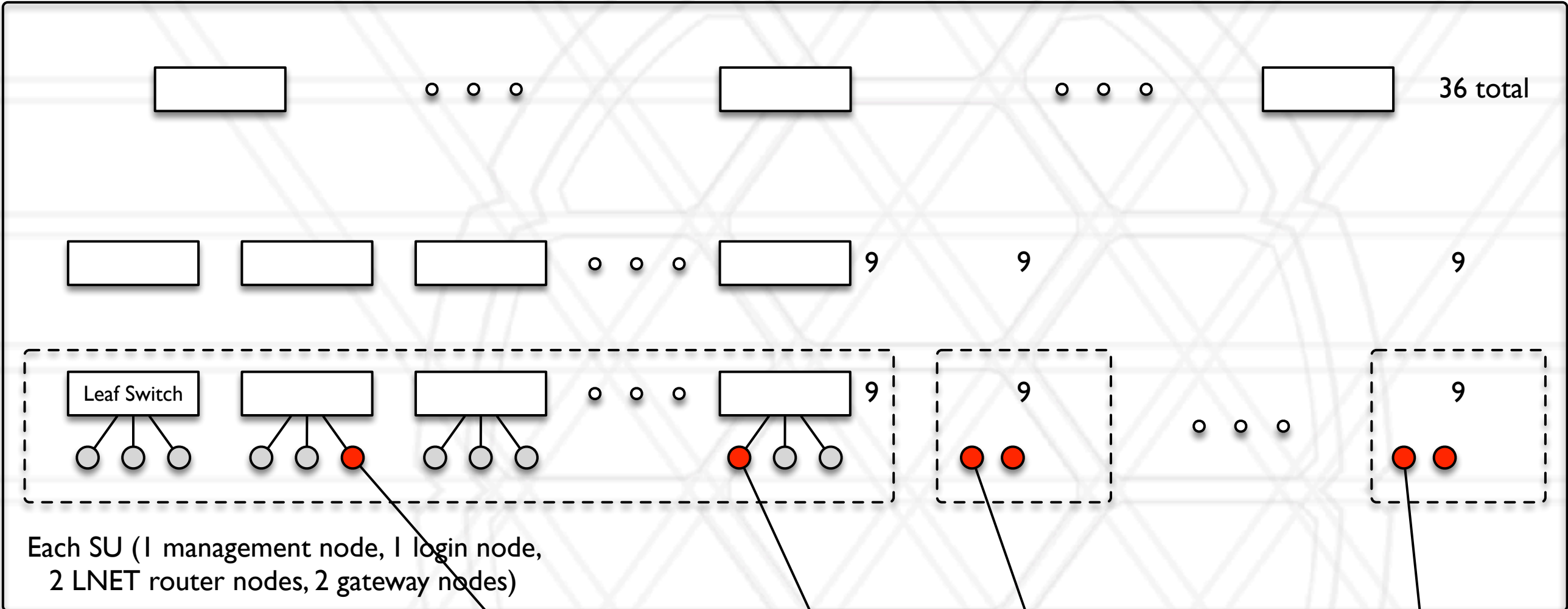
Parallel file system or I/O sub-system



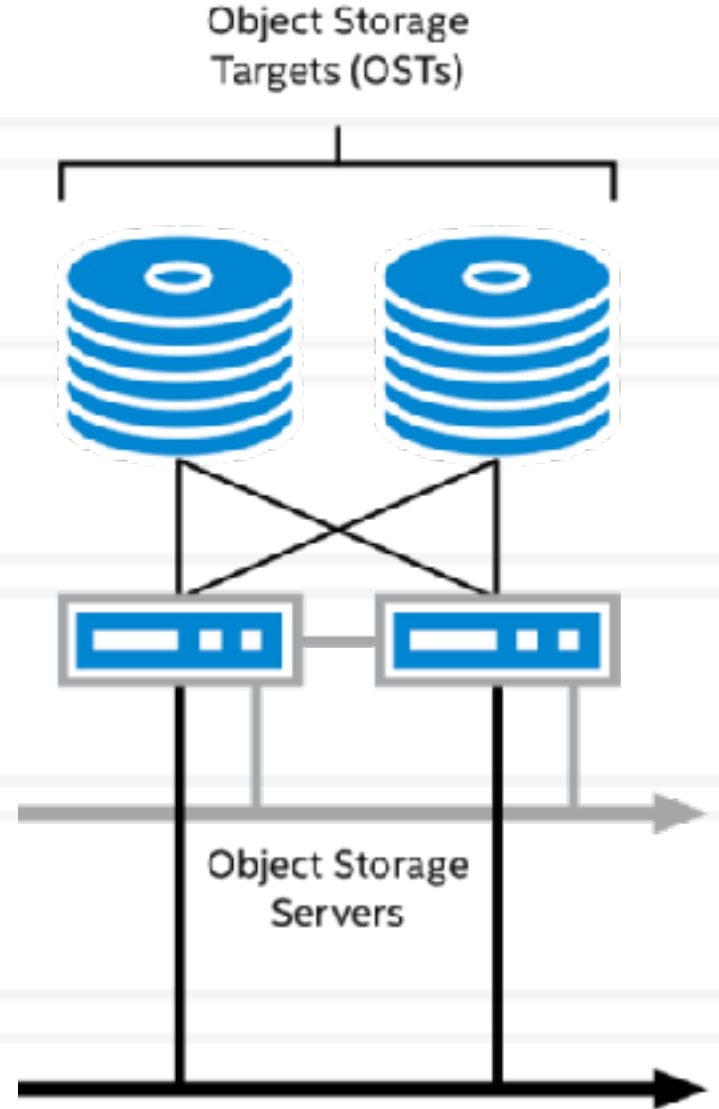
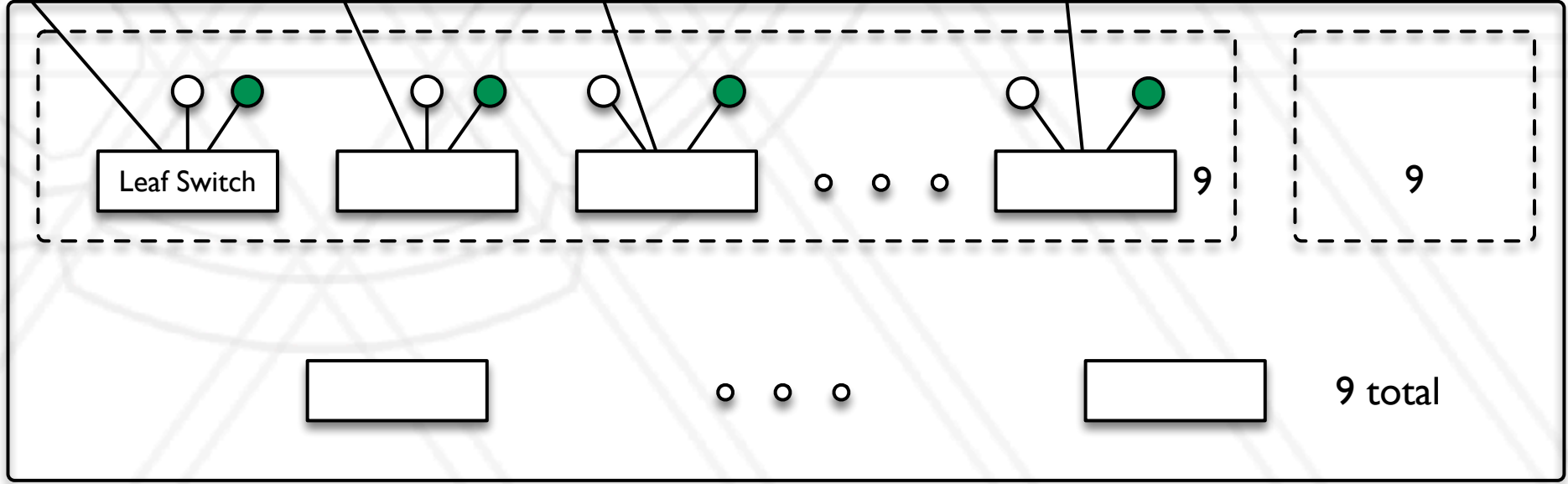
- Compute node
- LNET router node
- Object storage server (OSS)



Parallel file system or I/O sub-system



- Compute node
- LNET router node
- Object storage server (OSS)



Group Projects

- Self form into groups of 2-3
- Project will be ideally at the intersection of systems + ML
 - Using parallel systems to optimize an ML workload
- Timeline (all deadlines are midnight):
 - Group formation and project proposal: March 4
 - Interim report: April 17
 - Final presentation: May 6-13
 - Final report and code: May 15



UNIVERSITY OF
MARYLAND