



TECHNISCHE
UNIVERSITÄT
DARMSTADT

System and Parallel Programming

Prof. Dr. Felix Wolf

PARALLEL ARCHITECTURES

Outline



- Classification
- Memory architecture
- Interconnection networks
- Examples
 - Lichtenberg Cluster
- Cache coherence
- Memory consistency
- Synchronization

Taxonomies

- Number of instruction streams vs. number of data streams
- Memory architecture
- Network architecture
- Degree of heterogeneity
- Degree of customization

Flynn's classification [1966]

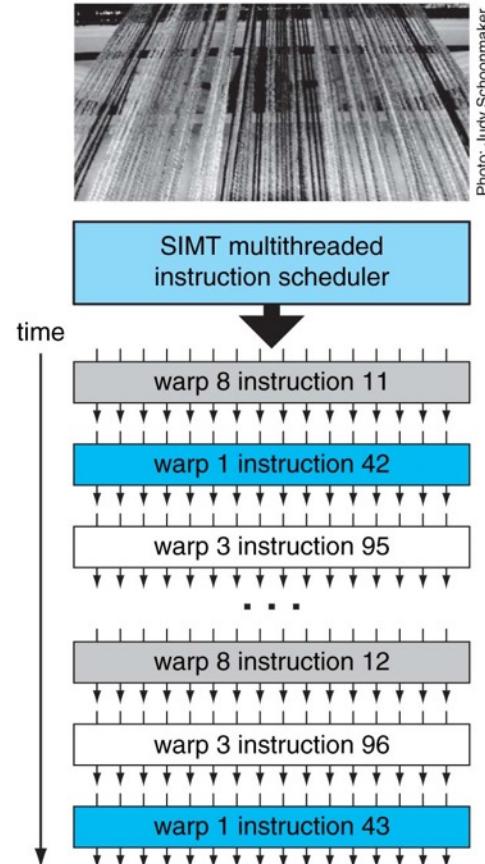
		#Instruction streams
		Single
#Data streams	Single	SISD <ul style="list-style-type: none">• Classical uniprocessor
	Multiple	MISD <ul style="list-style-type: none">• No commercial multiprocessor of this type ever built
#Data streams	Multiple	SIMD <ul style="list-style-type: none">• Same instruction is executed by multiple processors using different data streams• Data parallelism• Examples: SIMD extensions for multimedia, vector processors
	Multiple	MIMD <ul style="list-style-type: none">• Each processor fetches its own instructions and operates on its own data• Thread-level parallelism

Single-instruction multiple threads (SIMT)

- Used on GPUs -



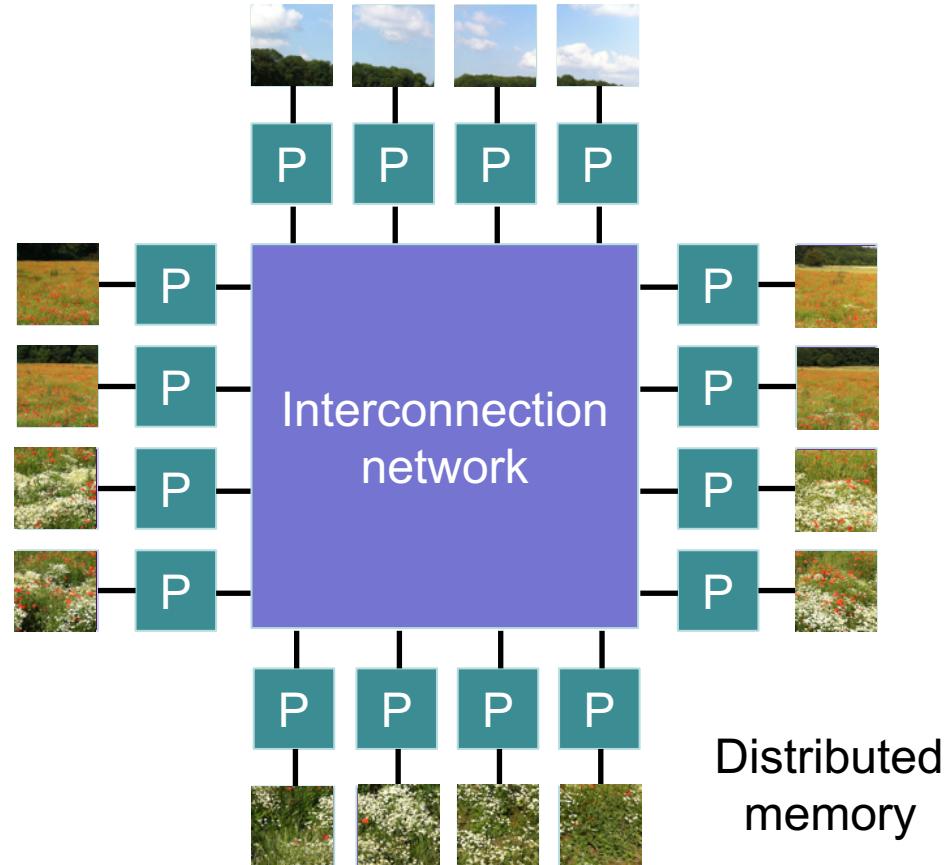
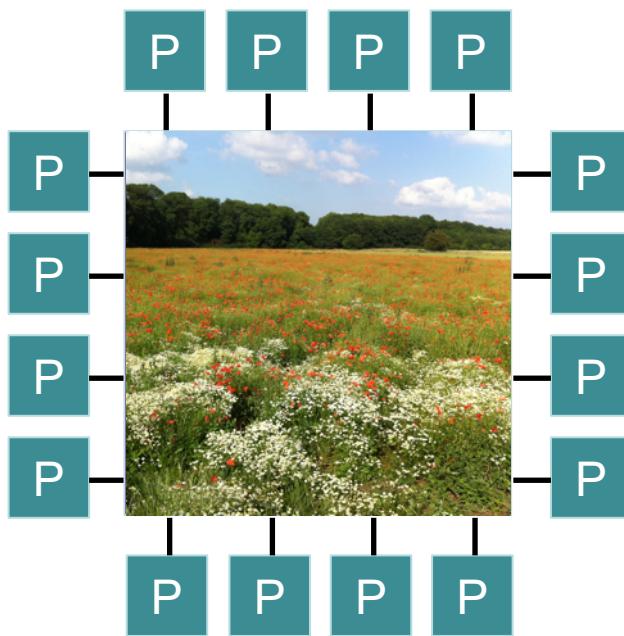
- Creates, manages, schedules, and executes threads in groups of parallel threads called **warps**
- At each instruction issue time, SIMT instruction unit
 - Selects warp that is ready to execute its next instruction
 - Broadcasts instruction to all active threads of that warp
- Individual threads may be inactive to do independent branching



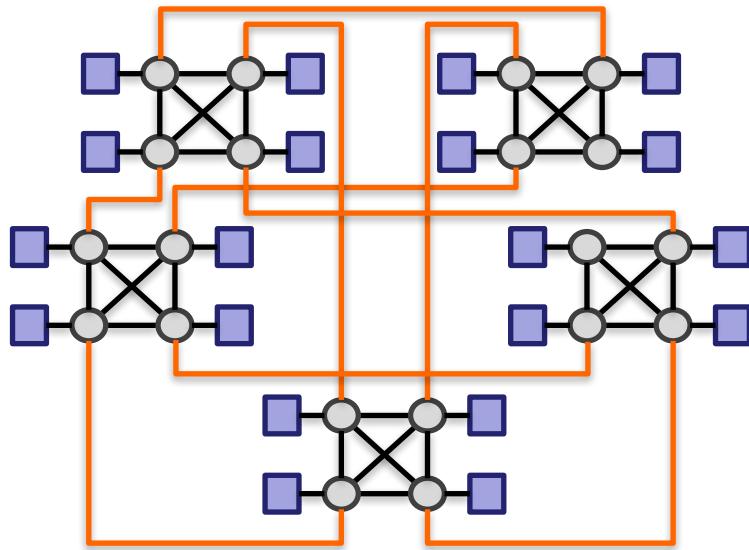


- Architecture of choice for general-purpose multiprocessors
- Offers high degree of flexibility
 - High performance for one application or multi-programmed multiprocessor
- Can take advantage of off-the-shelf processors
- Popular execution model - Single Program Multiple Data (SPMD)
 - The same program is executed in parallel with each instance having a potentially different control flow

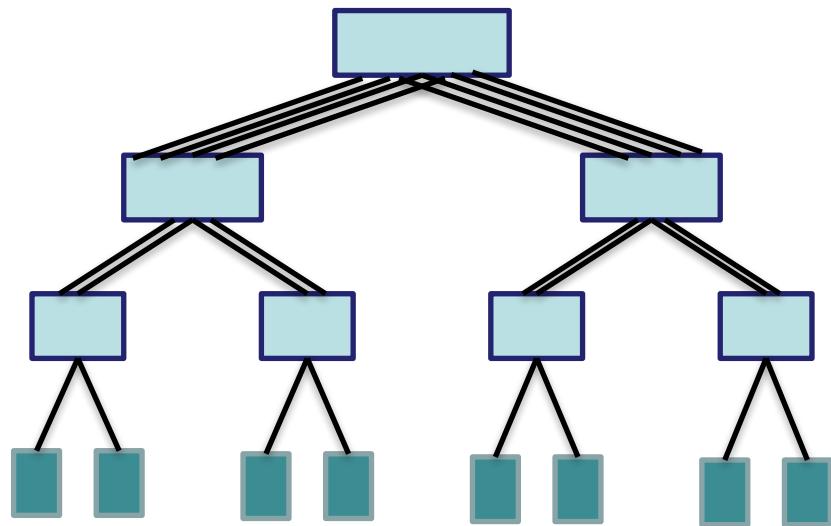
Memory architecture



Popular network architectures for distributed memory systems



Dragonfly
(distributed switched network)

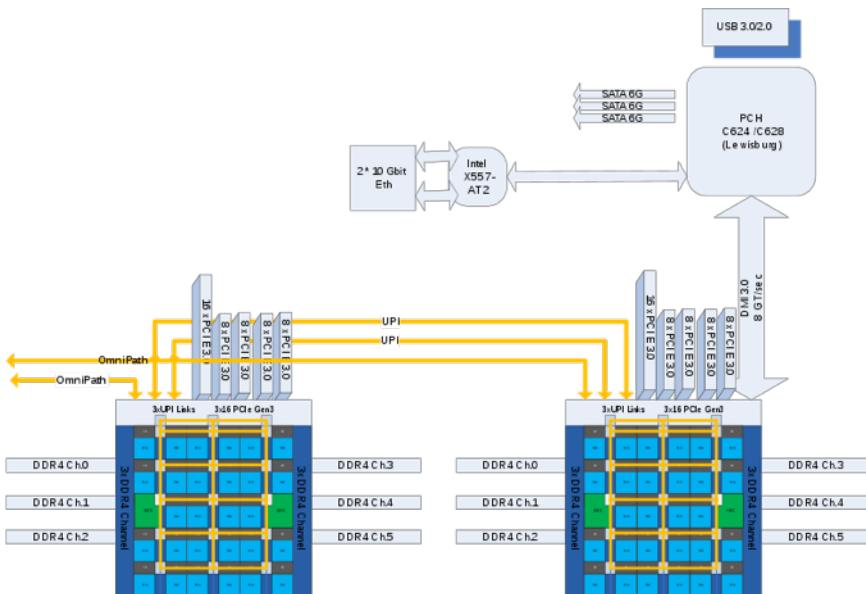


Fat tree
(centralized switched network)

Degree of heterogeneity



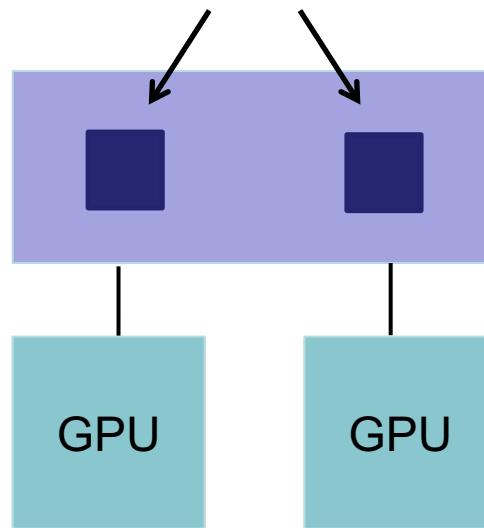
Homogeneous node architecture



Intel Skylake dual-socket system

Von Michael Wandinger - Eigenes Werk (Originaltext: selbst gezeichnet),
CC BY-SA 3.0 de, <https://commons.wikimedia.org/w/index.php?curid=63346901>

Classic server CPUs (e.g. Intel Xeon)



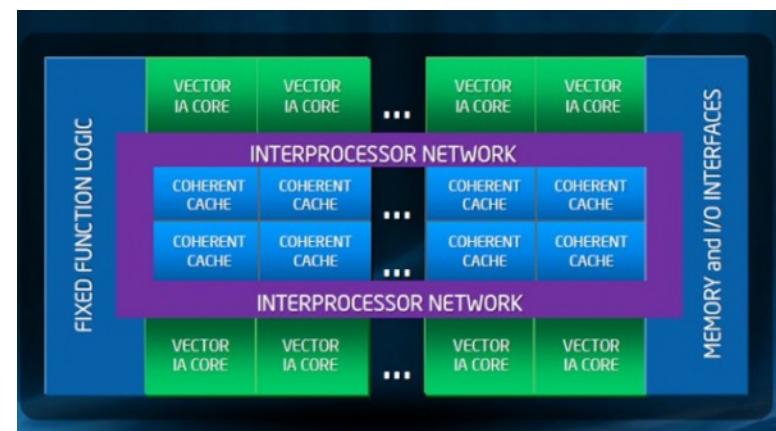
Accelerators

Heterogeneous node architecture

Accelerators



NVIDIA Volta GPU



Intel Xeon Phi

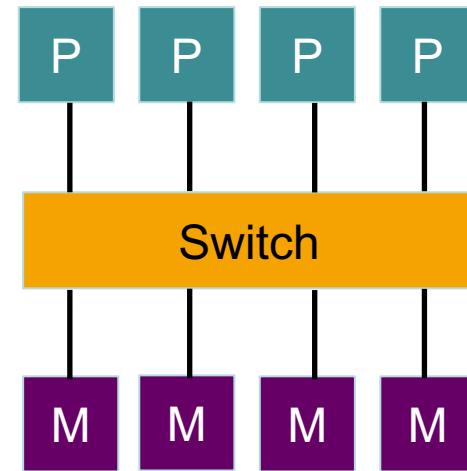
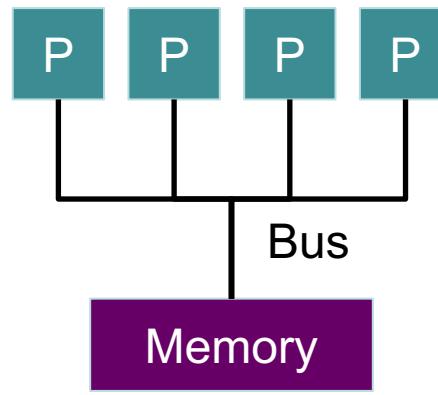
Degree of customization

- **Commodity clusters** – standard nodes and standard network
 - Focus on applications with small communication requirements
 - Example: Beowulf cluster
- **Custom clusters** – custom nodes and custom network
 - Also called massively parallel processors
 - Focus on applications that exploit large amounts of parallelism on single problem
 - Example: IBM Blue Gene/Q
- Above classes are extremes of a broad spectrum

Shared memory

UMA (Uniform memory access)

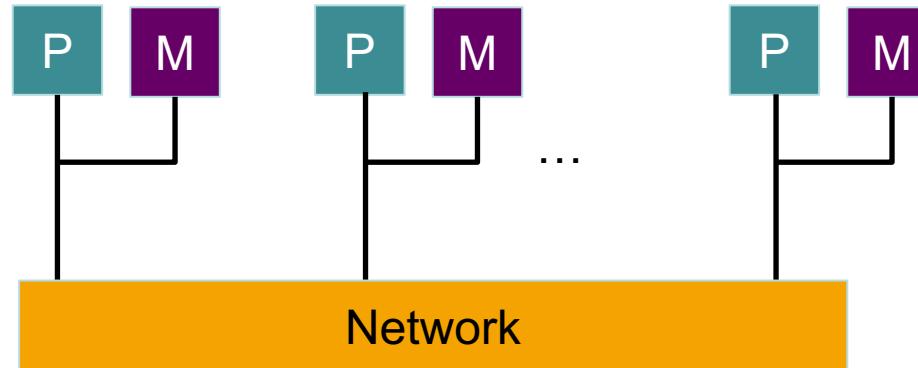
- Each CPU has same access time to each memory address
- Simple design but limited scalability (multicore or less)



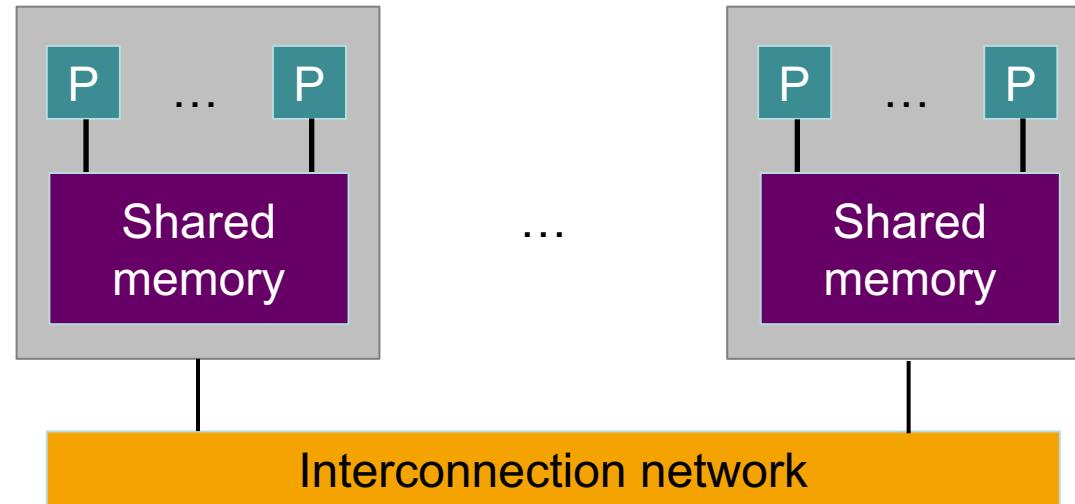
Shared memory (2)

NUMA (Non-uniform memory access)

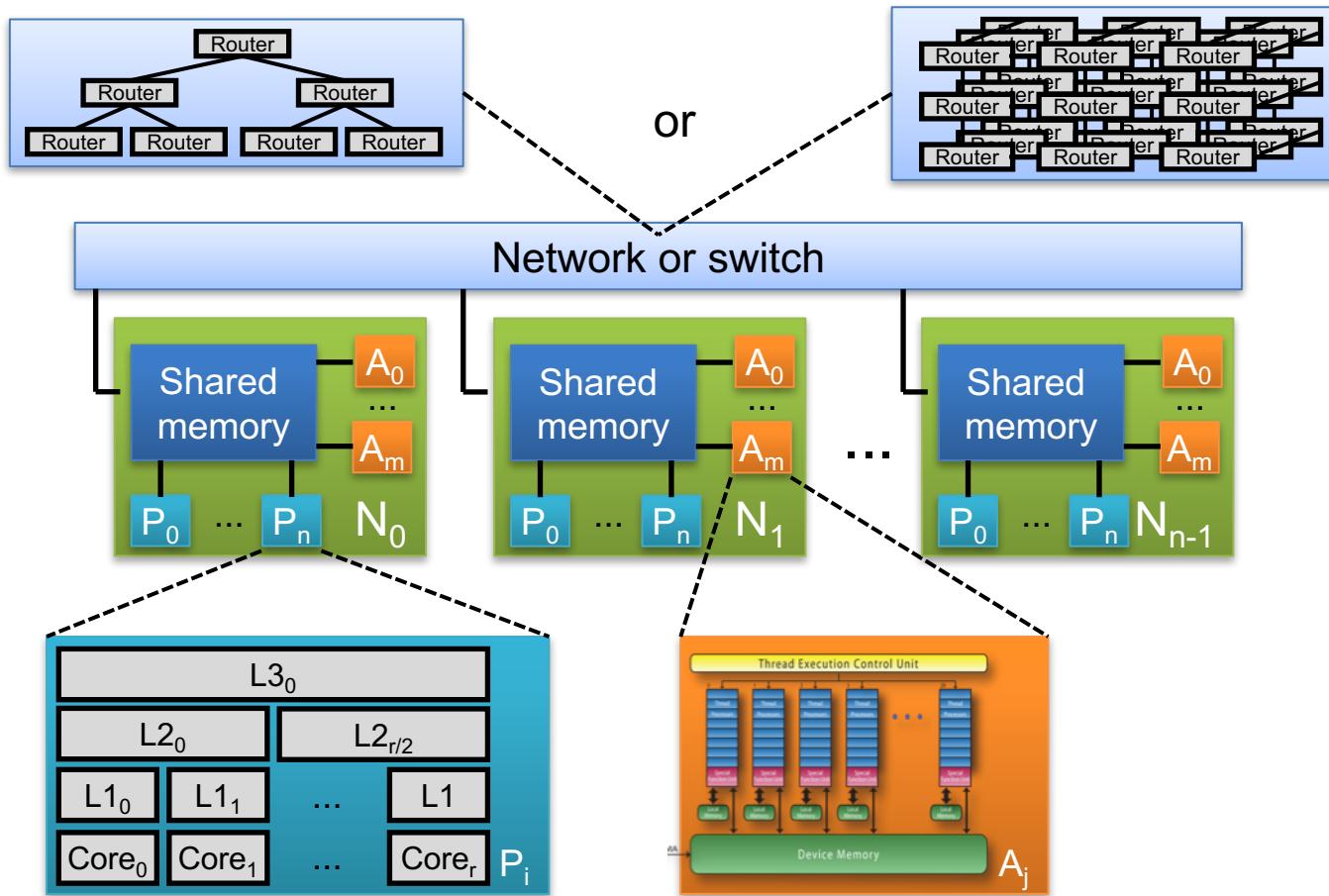
- Memory has affinity to a processor
- Access to local memory faster than to remote memory
- Harder to program but more scalable



Distributed memory (aka multicompiler)



Typical cluster architecture



Interconnection network



TECHNISCHE
UNIVERSITÄT
DARMSTADT

Physical link between components of a parallel system

- Between processors and memory
- Between nodes

Communication via exchange of messages

- Example: intermediate results, memory requests

Design elements

- **Topology** – determines geometric layout of links and switches
- **Routing technique** – determines paths of messages through network

Bandwidth

- Maximum rate at which information can be transferred
- Aggregate bandwidth – total data bandwidth supplied by network
- Effective bandwidth or throughput – fraction of aggregate bandwidth delivered to an application

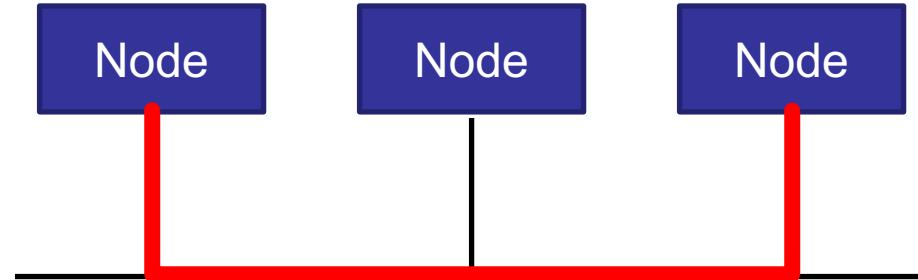
Latency

- Sending overhead + time of flight + receiving overhead



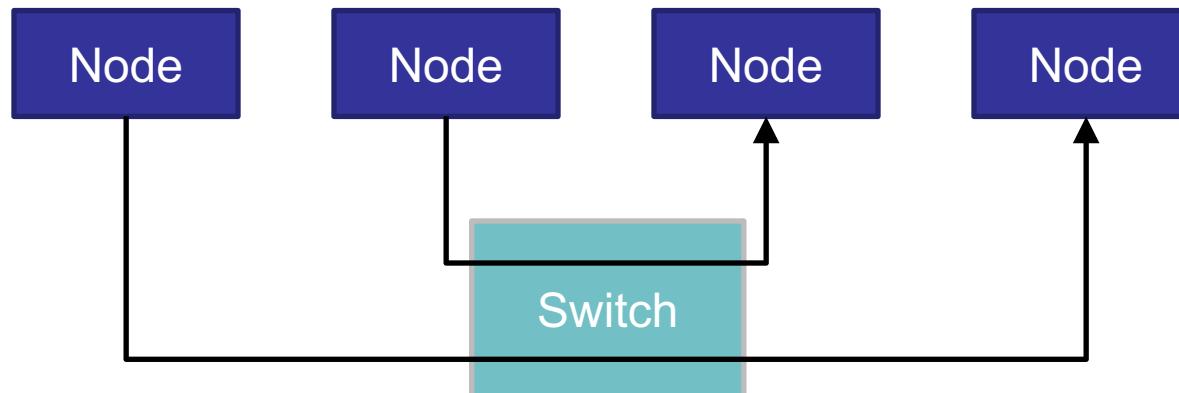
Shared-media networks

- Only one message at a time – processors broadcast their message over the medium
- Each processor “listens” to every message and receives the ones for which it is the destination
- Decentralized arbitration
 - Before sending a message, processors listen until medium is free
 - Message collision can degrade performance
- Low cost but not scalable
- Example – bus networks to connect processors to memory



Switched-media networks

- Support point-to-point messages between nodes
- Each node has its own communication path to the switch
- Advantages
 - Support concurrent transmission of multiple messages among different node pairs
 - Scale to very large numbers of nodes



Centralized switched networks

- Also called *indirect* or dynamic interconnection networks
- Connect processors / memory indirectly using several links and intermediate switches
- Examples: switching networks
- Used both for shared- and distributed-memory architectures

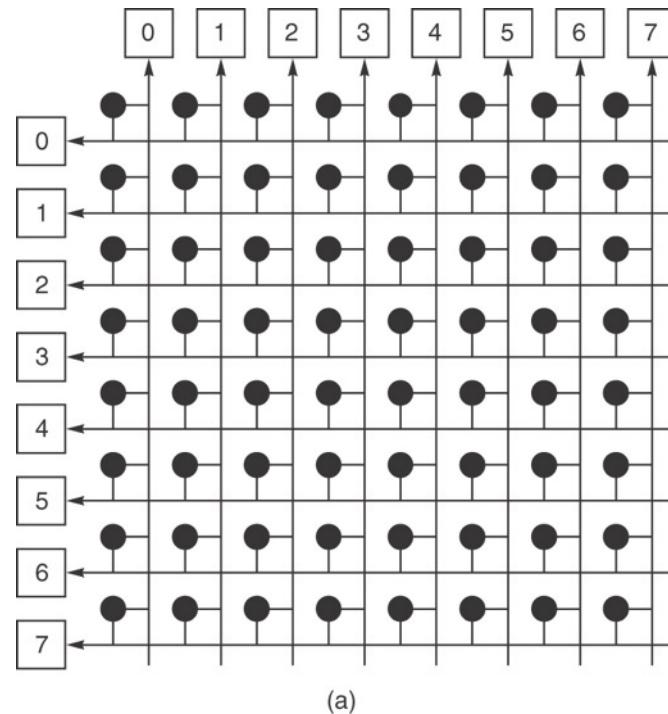
Crossbar switch

Non-blocking

- Links are not shared among paths to unique destinations

Requires N^2 crosspoint switches

- Limited scalability

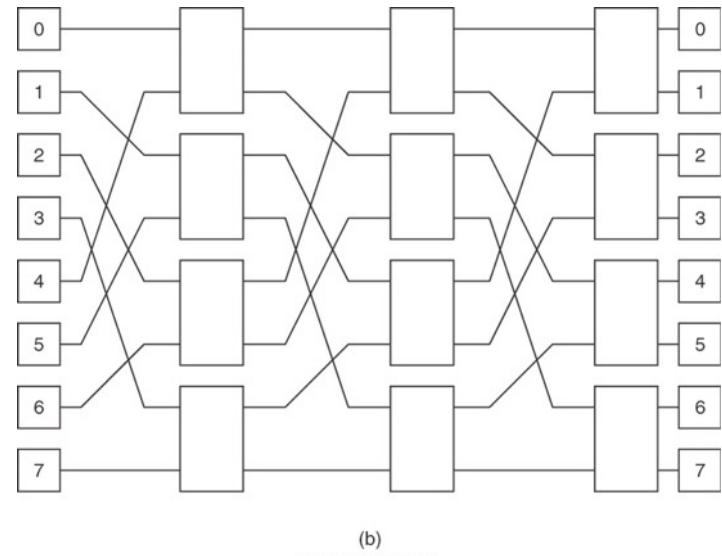


Source: Hennessy, Patterson: Computer Architecture, 4th edition, Morgan Kaufmann

Multistage interconnection network (MIN)

Example: Omega network

- Complexity $O(N \log N)$
- Perfect shuffle permutation at each stage
- Blocking due to paths between different sources and destinations simultaneously sharing network links
- Omega with $k \times k$ switches
 - $\log_k N$ stages ; $N/k \log_k N$ switches



(b)

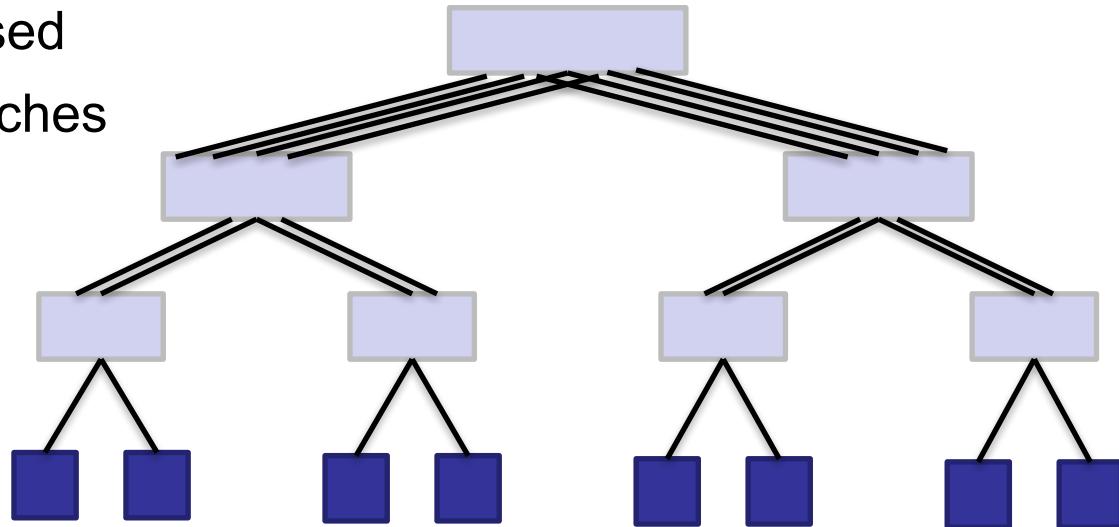
© 2007 Elsevier, Inc. All rights reserved.

Source: Hennessy, Patterson: Computer Architecture, 4th edition, Morgan Kaufmann

MINs can be extended to **rearrangeably** non-blocking topologies

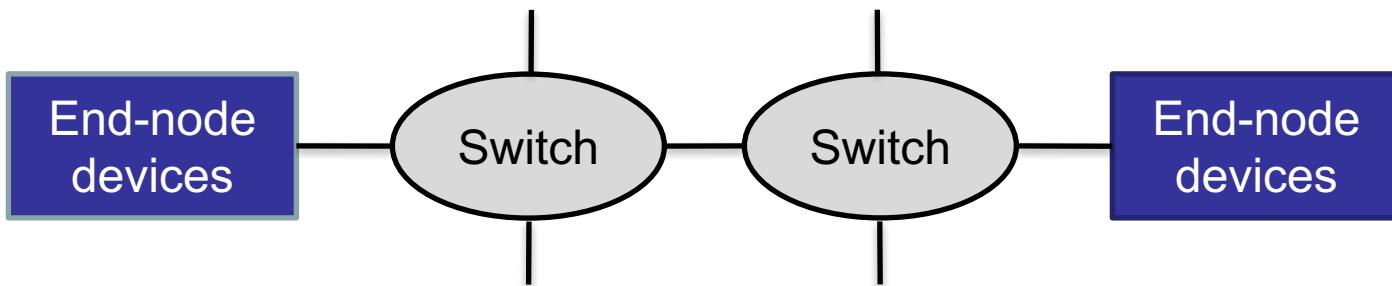
Fat tree

- Balanced tree where
 - Leaves = end node devices
 - Vertices = switches
- Total link bandwidth constant across all levels
- Switches often composed of multiple smaller switches
- Popular topology for cluster interconnects



Distributed switched networks

- Each network switch has one or more end node devices directly attached to it
- End node devices = processor(s) + memory
 - Directly connected to other nodes without going through external switches
 - Mostly used for distributed-memory architectures
- Also called *direct* or static interconnection networks
- Ratio of switches to nodes = 1:1



Evaluation criteria



Network degree

- Maximum node degree
- Node degree = number of adjacent nodes = (incoming + outgoing) edges

Diameter

- Largest minimum distance between two nodes

Bisection bandwidth

- Divide the network into two equal-sized partitions (in terms of #nodes) such that aggregate bandwidth between them is minimized

- Bisection bandwidth is aggregate bandwidth [bytes/s] between these two partitions

Edge / node connectivity

- Minimum number of edges / nodes that need to be removed to render network disconnected

Embedding

- Mapping of one network onto another

Requirements

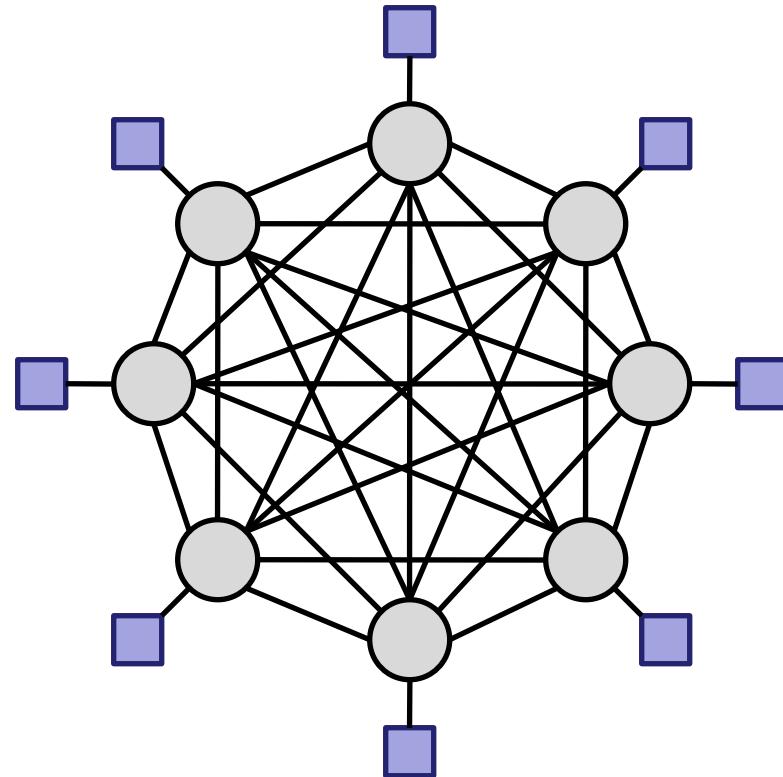


- Low network degree to reduce hardware costs
- Low diameter to ensure low distance (i.e., latency) for message transfer
- High bisection bandwidth to ensure high throughput
- High connectivity to ensure robustness
- Option to embed many other networks to ensure flexibility

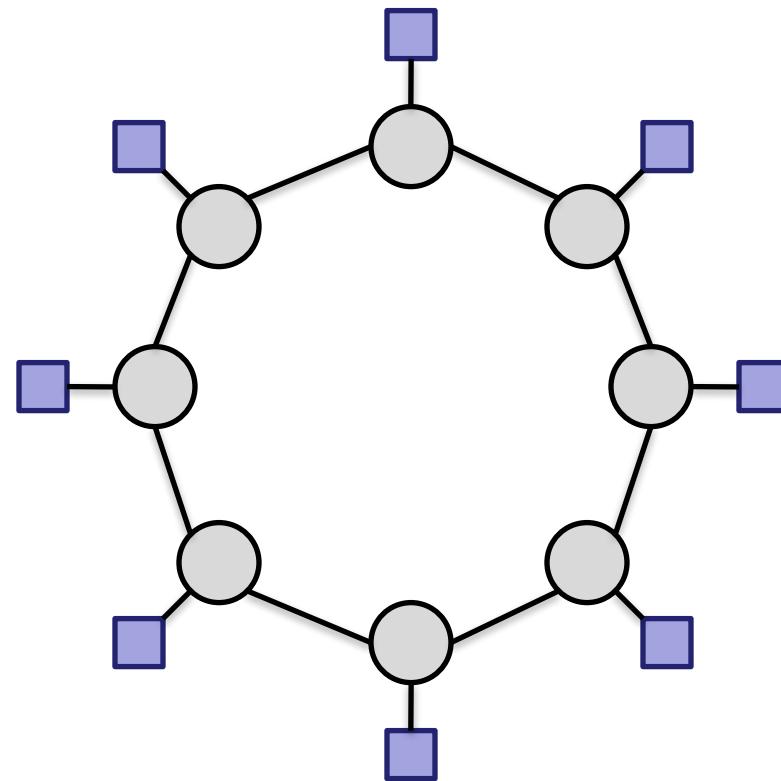
Often conflicting goals

Fully connected topology

- Each node is directly connected to every other node
- Expensive for large numbers of nodes
- Dedicated link between each pair of nodes
- Cheaper alternative: crossbar topology

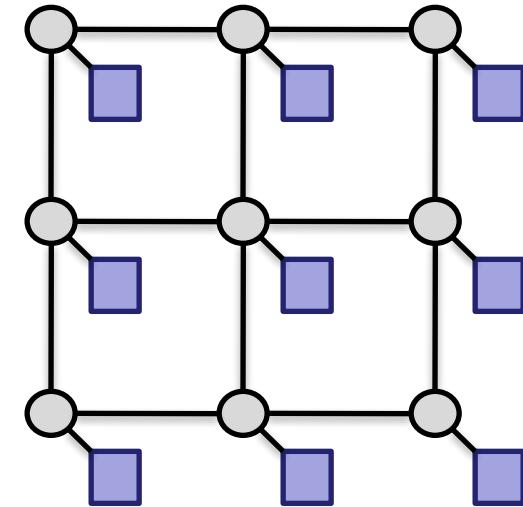


Ring topology



N-dimensional meshes

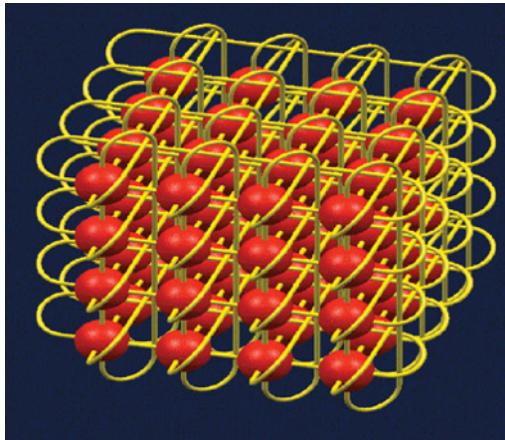
- Direct link to neighbors
- Each node has 1 or 2 neighbors per dimension
 - 2 in the center
 - Less for border or corner nodes
- Efficient nearest neighbor communication
- Suitable for large numbers of nodes



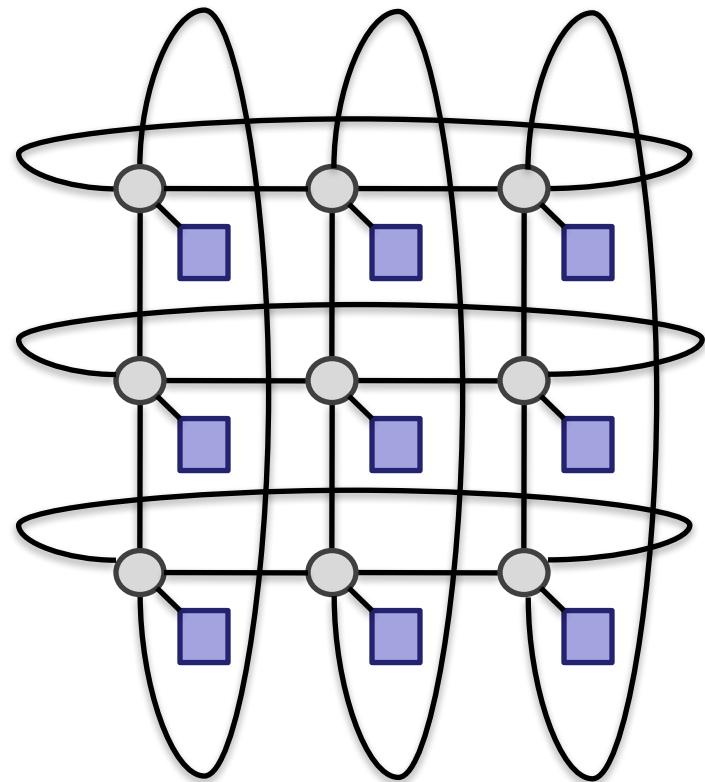
2D mesh

Torus

- Mesh with wrap-around connections
- Each node has exactly 2 neighbors per dimension
- Typically 3-6 dimensions



3D torus



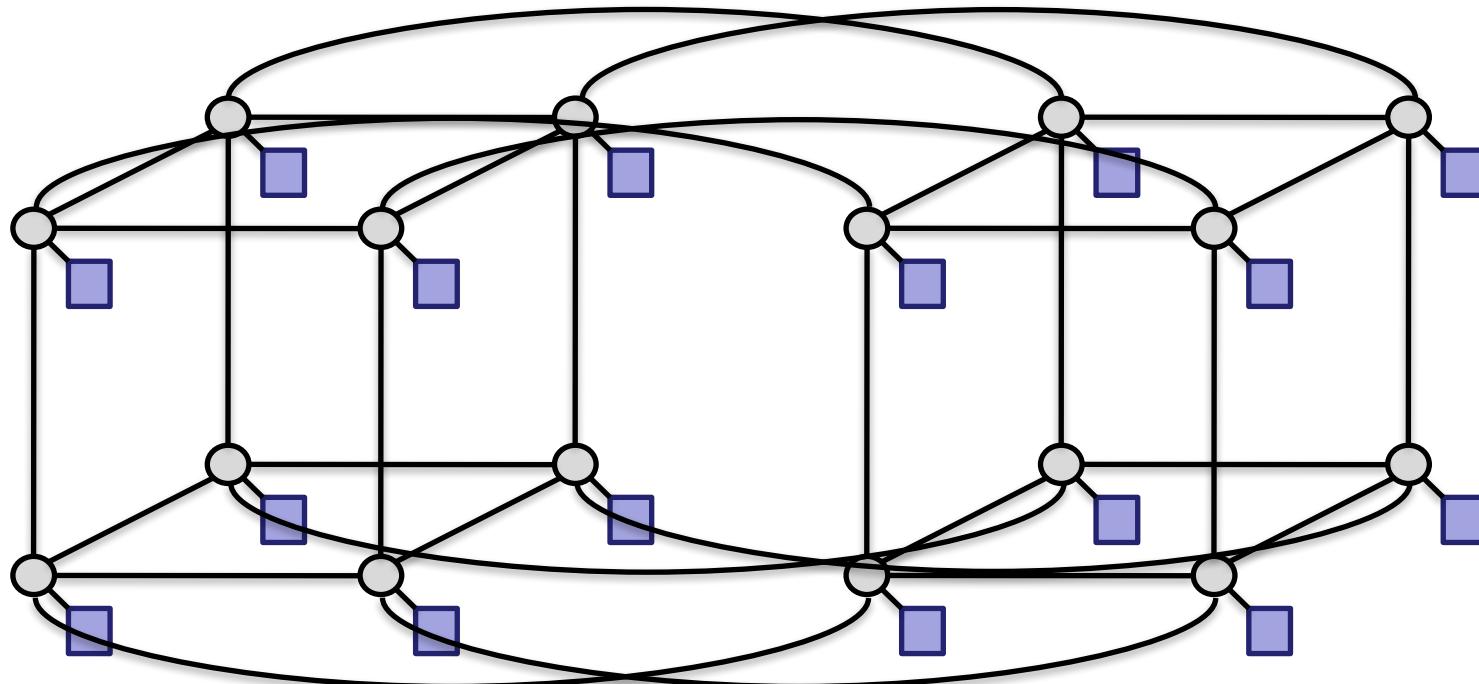
2D torus

Hypercube

16 nodes
 $(16 = 2^4 \text{ so } n = 4)$



TECHNISCHE
UNIVERSITÄT
DARMSTADT



Each node has one connection along each dimension ($n = \# \text{dimensions}$)

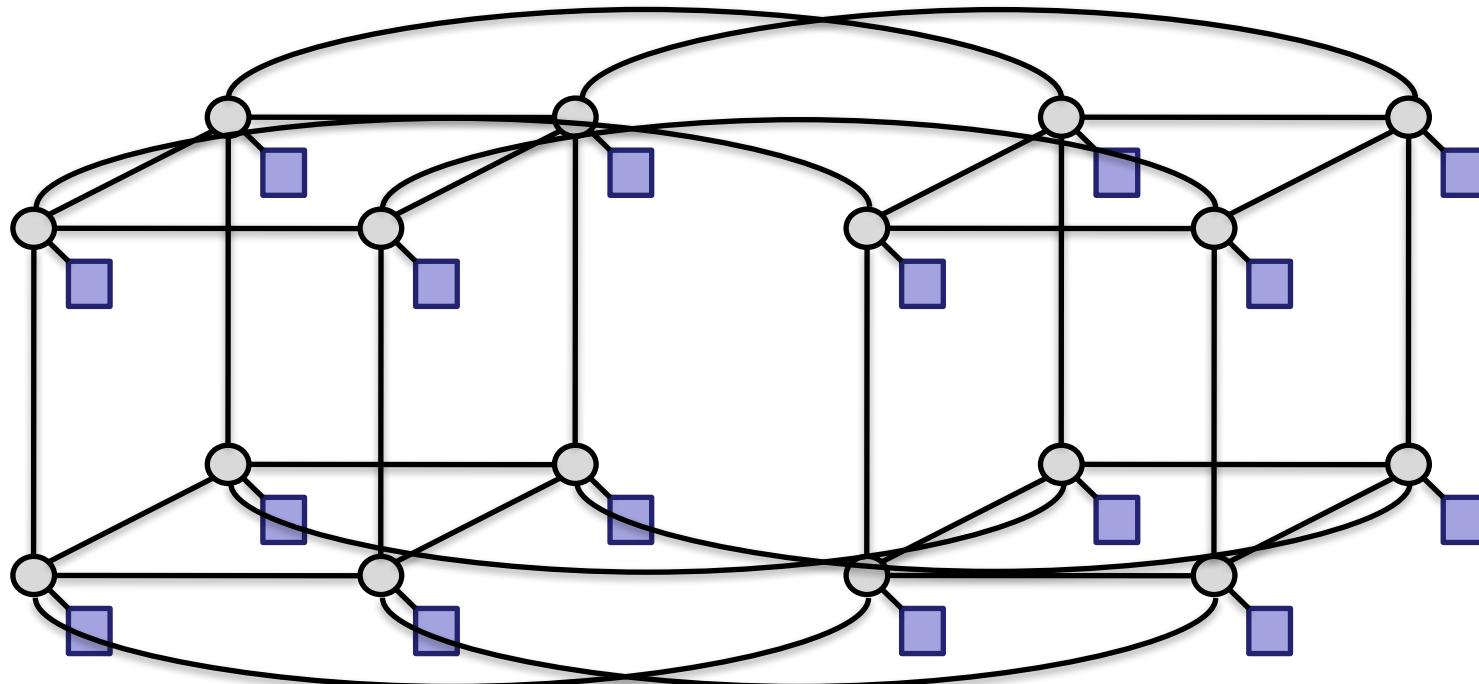
Usually better connectivity than tori at the expense of higher link and switch costs

Hypercube

16 nodes
 $(16 = 2^4 \text{ so } n = 4)$



TECHNISCHE
UNIVERSITÄT
DARMSTADT



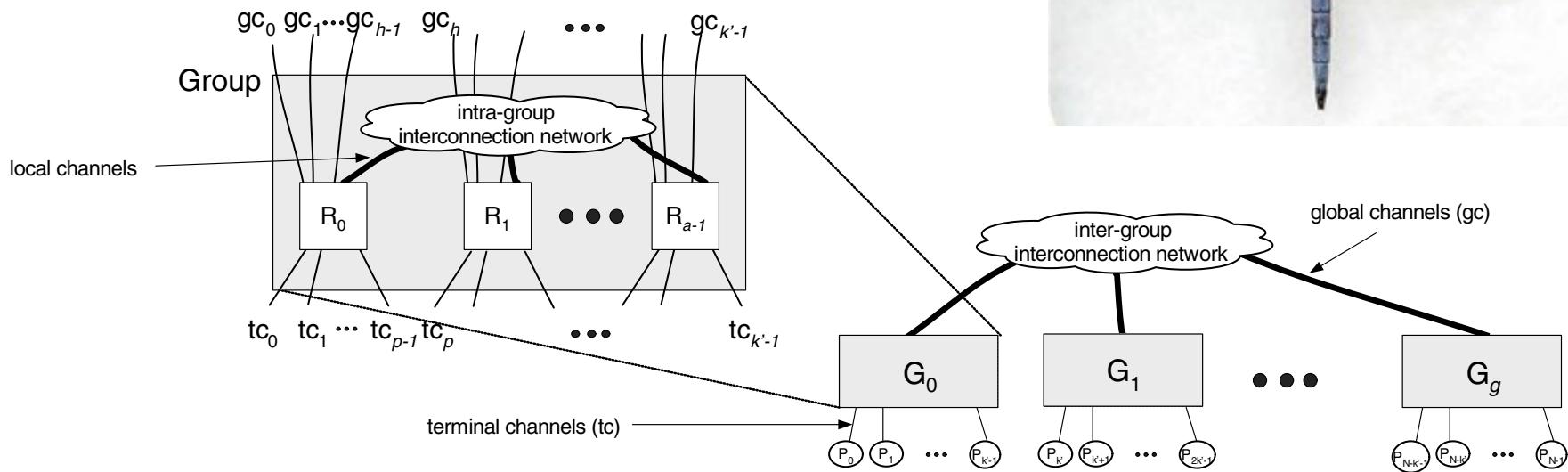
Each node has one connection along each dimension ($n = \# \text{dimensions}$)

Usually better connectivity than tori at the expense of higher link and switch costs

Dragonfly



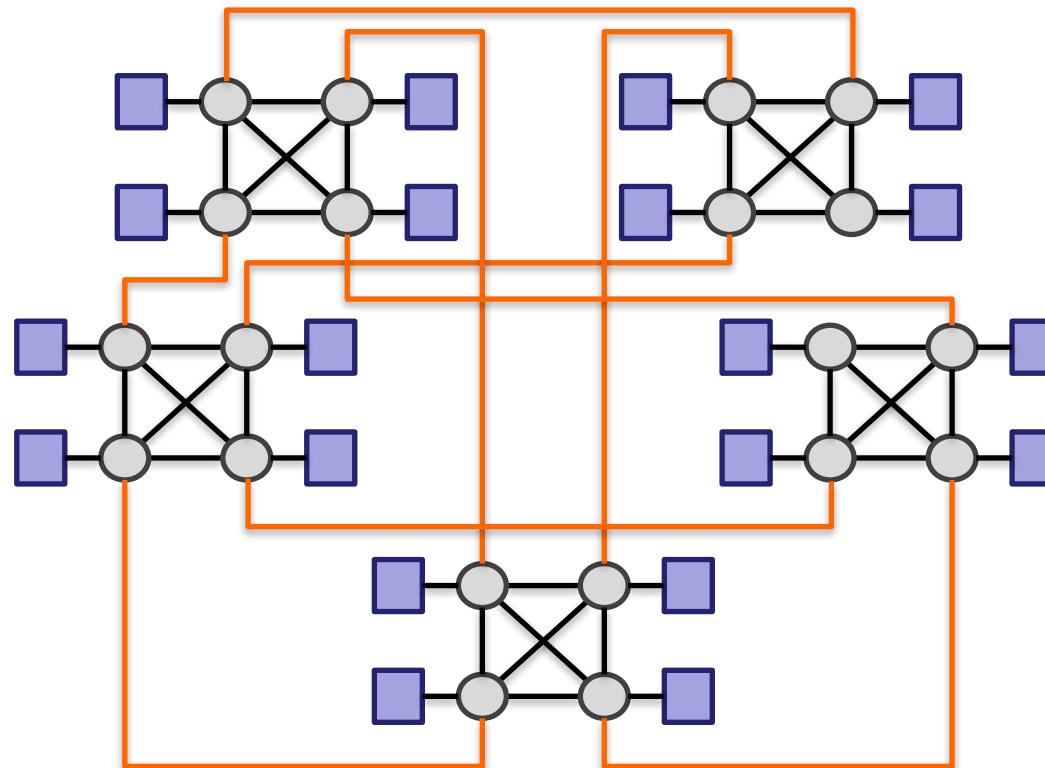
- Enabled by high-radix switches and long-distance optical signaling technology



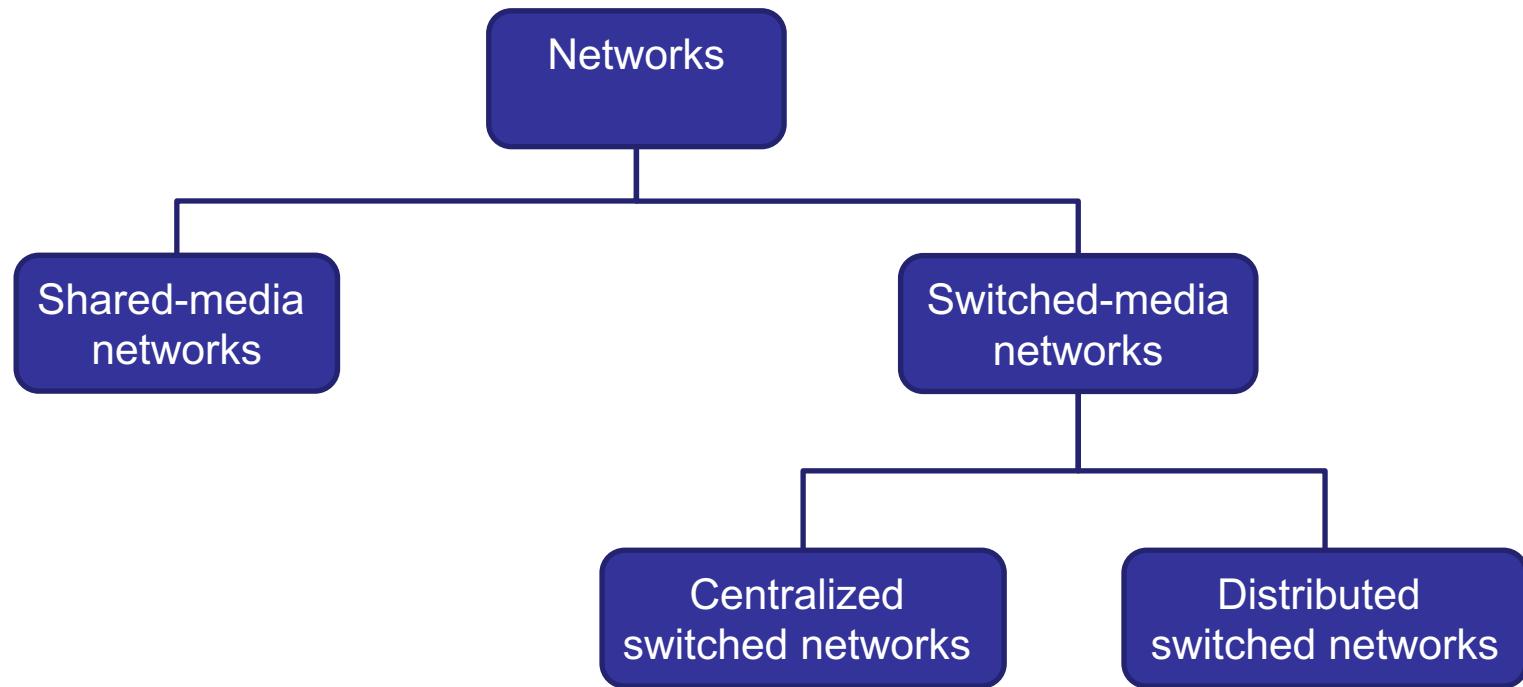
Source: J. Kim, W. Dally, S. Scott, and D. Abts, "Technology-driven, highly- scalable dragonfly topology". In Proc. of the 35th International Symposium on Computer Architecture (ISCA), June 2008, pp. 77–88.

Dragonfly (2)

- Arbitrary networks possible for intra-group and inter-group networks
 - Example: fully connected networks for both



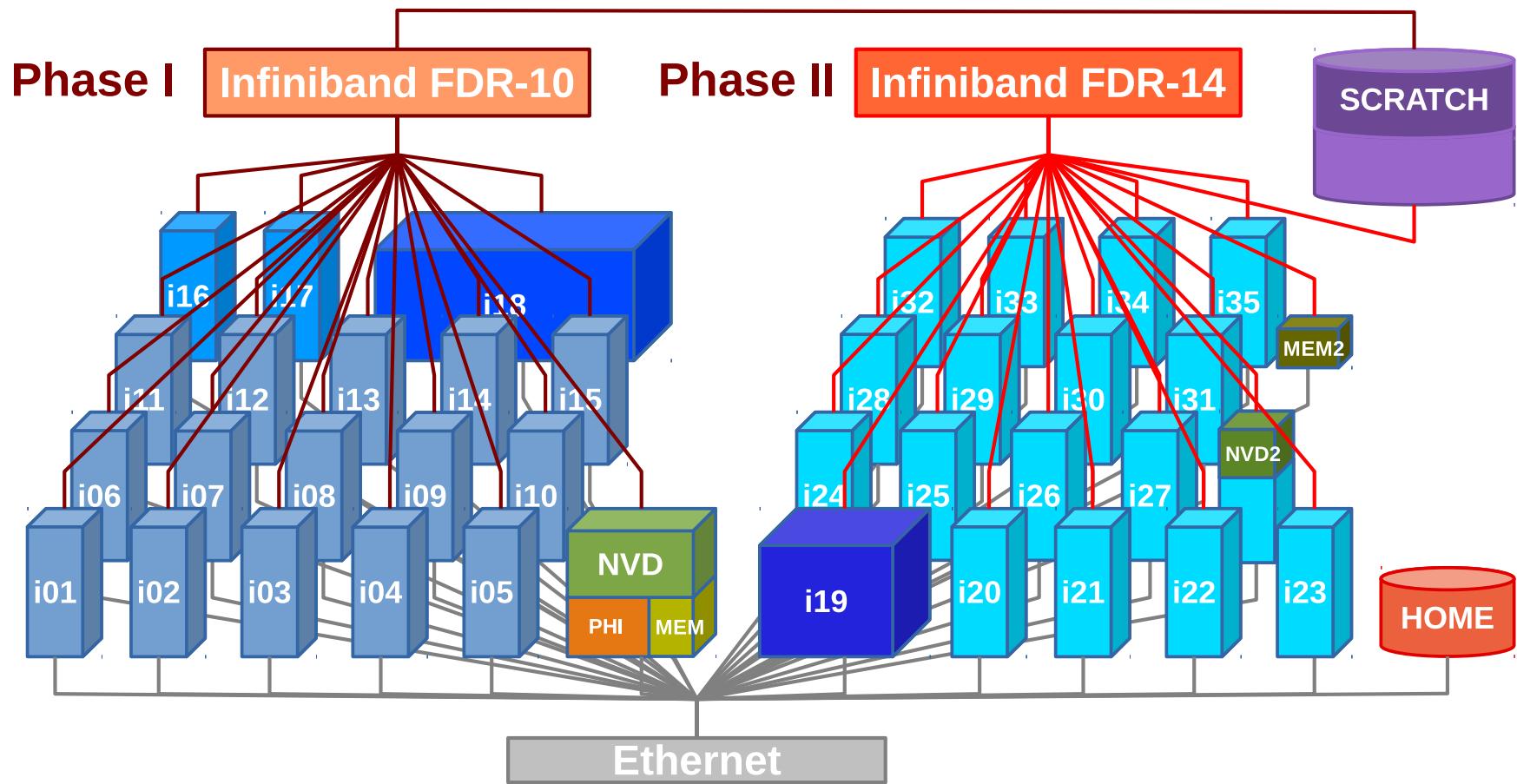
Network classes



Lichtenberg Cluster @ TU Darmstadt

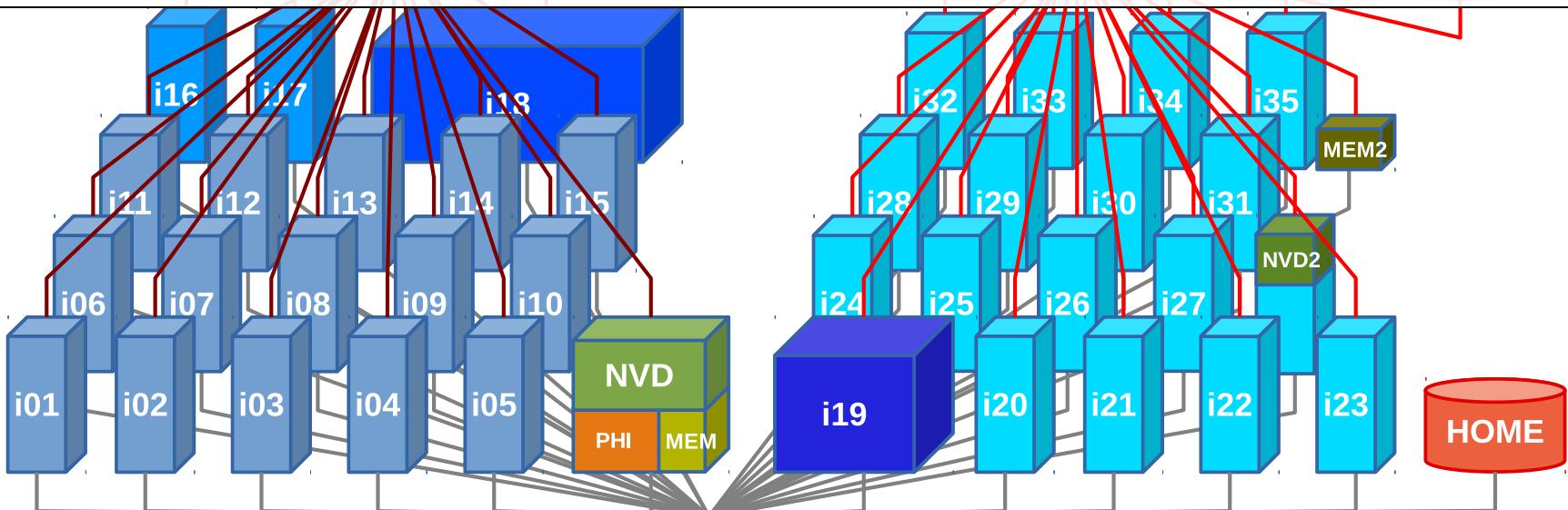


One cluster – multiple islands



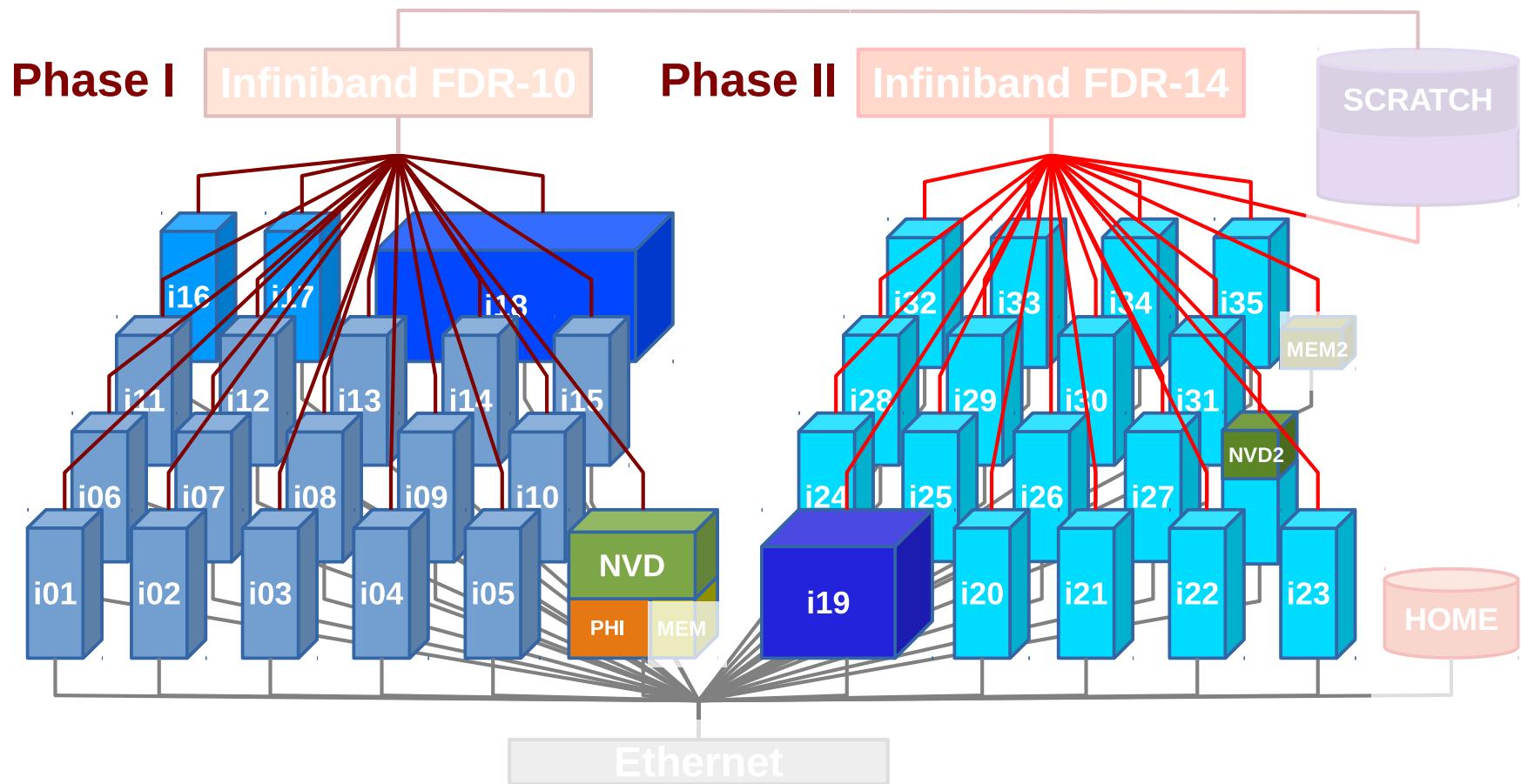
One cluster – multiple islands

- Cluster is divided into 2 phases
- Each phase is divided into several islands
- Rule of thumb: band FDR-10 Phase II Infiniband FDR-14
 $1 \text{ island} \hat{=} 32 \text{ compute nodes} \hat{=} 512 \text{ (ph. I) / } 768 \text{ (ph. 2) CPU cores}$
- For large computations, there are 2 islands with more than 2000 CPU cores



- Computation across more than one island is only possible on request, due to some technical limitations (across phases impossible).

Compute nodes (“mpi”, “nvd”, “phi”)



Compute nodes



Phase I (704+70 nodes):

Processors:

- 2 Intel Xeon E5-4650
(Sandy Bridge) processors
 $\hat{=} 2 \cdot 8 = 16 \text{ CPU cores}$

- **2.7 GHz**
(up to 3.3 GHz in turbo mode)

Main Memory:

- **32 GB RAM** (some have 64 GB)

Network:

- Gigabit Ethernet
- FDR-10 InfiniBand

Phase II (596+31 nodes):

Processors:

- 2 Intel Xeon E5-2680 v3
(Haswell) processors
 $\hat{=} 2 \cdot 12 = 24 \text{ CPU cores}$

- **2.5 GHz**
(up to 3.3 GHz in turbo mode)

Main Memory:

- **64 GB RAM**

Network:

- Gigabit Ethernet
- FDR-14 InfiniBand

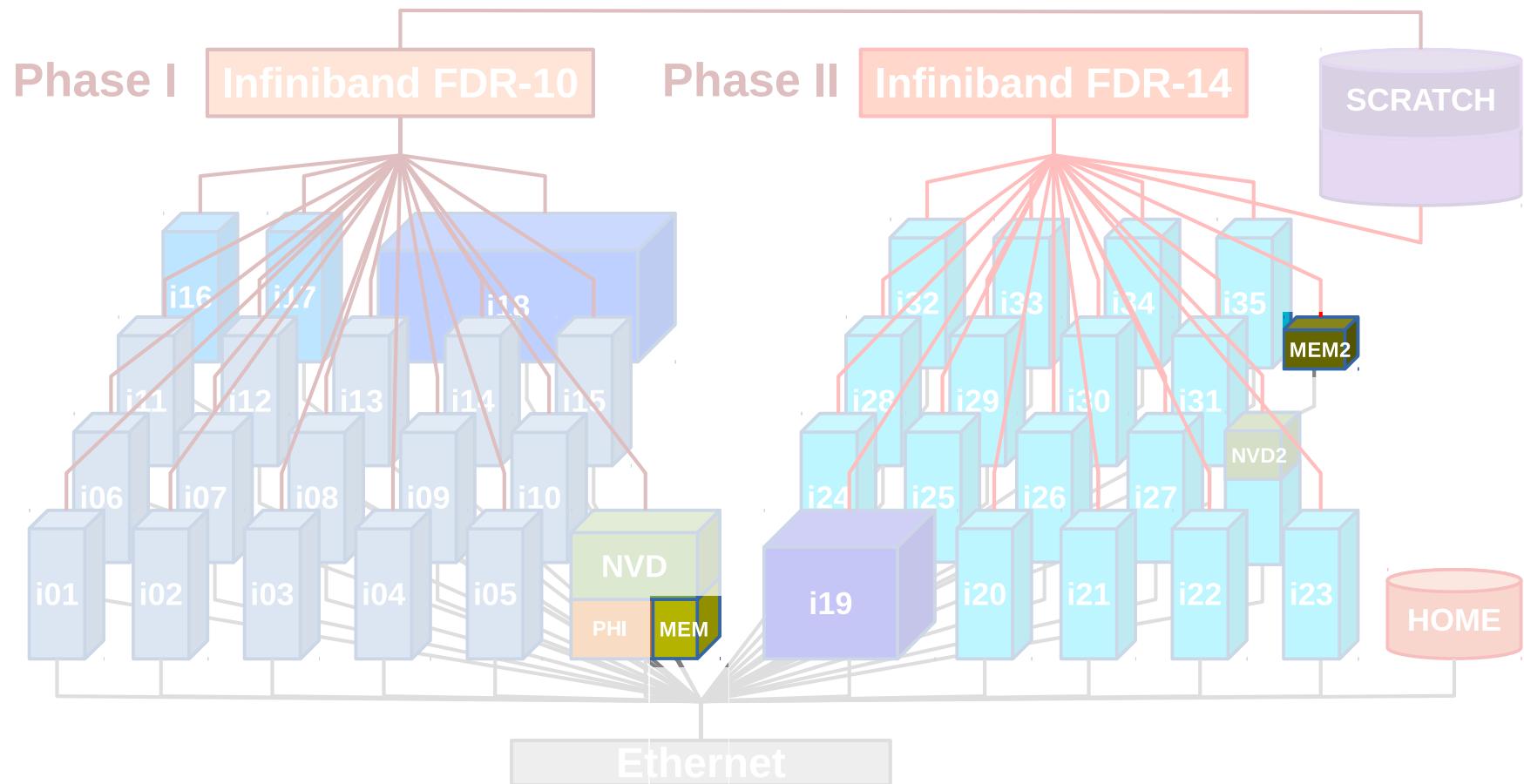
NVIDIA nodes

- 44 Sandy Bridge compute nodes have 2 **NVIDIA K20Xm** cards each
- 2 Haswell compute nodes have 2 **NVIDIA K40m** cards each
- 1 Haswell compute node has 2 **NVIDIA K80** cards

Xeon Phi Nodes

- 24 Sandy Bridge compute nodes have
2 Intel Xeon Phi 5110P cards each
- 2 Sandy Bridge compute nodes have
2 Intel Xeon Phi 7120P cards each

Big mem nodes (“mem”, “mem2”)



Mem nodes



Phase I (4 nodes):

Processors:

- 8 Intel Xeon E7-8837
(Westmere) processors
 $\hat{=} 8 \cdot 8 = \mathbf{64 \text{ CPU cores}}$

- **2.66 GHz**
(up to 2.8 GHz in turbo mode)

Main Memory:

- **1 TB (1024 GB) RAM**

Network:

- 10 Gigabit Ethernet
- $2 \cdot \text{FDR-10 InfiniBand}$

Phase II (4 nodes):

Processors:

- 4 Intel Xeon E7-4890 v2
(Ivy Bridge) processors
 $\hat{=} 4 \cdot 15 = \mathbf{60 \text{ CPU cores}}$

- **2.8 GHz**
(up to 3.4 GHz in turbo mode)

Main Memory:

- **1 TB (1024 GB) RAM**

Network:

- 10 Gigabit Ethernet
- $2 \cdot \text{FDR-14 InfiniBand}$

File systems



Mountpoint	/home	/work/scratch	/work/local
Size	> 300 TB	> 650 TB	> 100 GB per node
Access time	Normal (Ethernet)	Fast (InfiniBand)	Very fast (local HDD)
Accessibility	Global (cluster)	Global (cluster)	Local (node)
Data availability	permanent	≥ 1 month	Only during job runtime
Quota*	15 GB**	100 TB** 2 Mio. files**	none
Backup	Weekly + snapshots	none	none

* Use the command `cquota` to find out your current usage and quota.

** Can be increased on request.

Login nodes



4 nodes (hardware similar to Phase I):

Processors:

- 4 Intel Xeon E5-4650
(Sandy Bridge) processors
 $\hat{=} 4 \cdot 8 = 32$ CPU cores

- 2.7 GHz
(up to 3.3 GHz in turbo mode)

Main Memory:

- 128 GB RAM

Network:

- 2 · 10 Gigabit Ethernet
- 2 · FDR-10 InfiniBand

8 nodes (hardware similar to Phase II):

Processors:

- 2 Intel Xeon E5-2680 v3
(Haswell) processors
 $\hat{=} 2 \cdot 12 = 24$ CPU cores

- 2.5 GHz
(up to 3.3 GHz in turbo mode)

Main Memory:

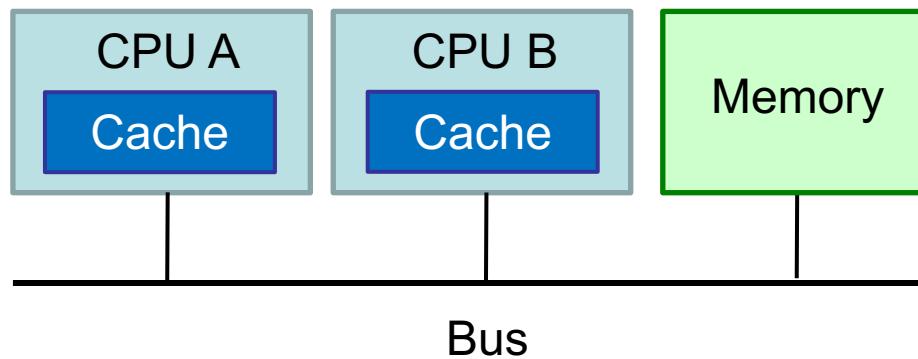
- 128 GB RAM

Network:

- 2 · 10 Gigabit Ethernet
- FDR-14 InfiniBand

Parallelism and memory hierarchy

A processor's view of the memory is through its cache



Cache coherence

- Problem – different processors may see different values
 - Without further precautions (!)

Time	Event	Cache CPU A	Cache CPU B	Memory
0				1
1	CPU A reads X	1		1
2	CPU B reads X	1	1	1
3	CPU A stores 0 in X	0	1	0

- Cache coherence – which value will be returned by a read?
 - Cache coherence protocols prevent different versions of the same cache line from appearing simultaneously in two or more caches

Coherence of a memory system

1. A read by a processor P to a location X that follows a write by P to X, with no writes of X by another processor occurring between the write and the read by P, always returns the value written by P
2. A read by a processor to location X that follows a write by another processor returns the written value if the read and the write are sufficiently separated in time and no other writes to X occur between the two accesses
3. Writes to the same location are serialized, that is, two writes to the same location by any two processors are seen in the same order by all processors.

Memory consistency

- Coherence refers to the behavior of the memory system when a single memory location is accessed by multiple threads
- Consistency refers to the ordering of accesses to different memory locations, observable from various threads in the system
 - When must a processor see a value that has been updated by another processor?
 - In what order does a processor observe the data writes of another processor?

Assumption for now: Sequential consistency



TECHNISCHE
UNIVERSITÄT
DARMSTADT

The result of any execution is the same as if the operations of all the processors were executed in some sequential order, and the operations of each individual processor appear in this sequence in the order specified by its program

[Lamport, 1979]

- Advantage – simple programming paradigm
- Disadvantage – potential performance degradation

Cache coherence protocols



Snooping

- Every cache that has a copy of a block of physical memory also has a copy of the sharing status of the block
- No centralized state is kept
- All caches are accessible via some broadcast medium (i.e., bus)
- All caches snoop on the medium to see whether they have a copy of a block that is being requested

Directory-based

- The sharing status of a block of physical memory is kept in one location (i.e., the directory)
- Two variants
 - Centralized directory (UMA)
 - Distributed directory (NUMA)
- Distributed directory needed if scalability is a concern



Write invalidate

- Remaining copies are invalidated on a write
- Most common for both snooping and directory-based protocols
- Guarantees exclusive access

Event	Bus activity	Cache CPU A	Cache CPU B	Memory
				0
CPU A reads X	Cache miss for X	0		0
CPU B reads X	Cache miss for X	0	0	0
CPU A writes 1 to X	Invalidation for X	1		0
CPU B reads X	Cache miss for X	1	1	1

Write update or write broadcast

- All copies are updated on a write
- Consumes more bandwidth – less common

Implementation of a write invalidate

- Invalidate is performed by broadcasting address to be invalidated on the bus
- All processors snoop on the bus and watch the addresses: if address is in their cache, data are invalidated
- Writes to a shared data item are serialized – only one processor can have access to the bus at a time

Finding data items on a miss

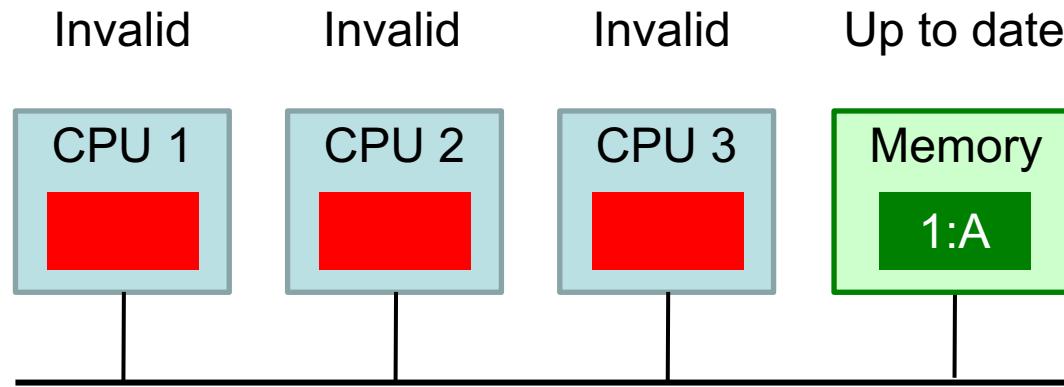
- Write-through caches
 - Data can be retrieved from main memory
- Write-back caches
 - Cache misses handled via snooping – if processor has dirty copy of a cache line, it provides the block in response to a request and causes memory access to be aborted
 - Complexity arises from transferring block to requesting processor – can take longer than getting it from main memory - especially if processors are on separate chips
 - Lower memory bandwidth demands & hence more scalable – often used at outermost cache levels

MESI protocol



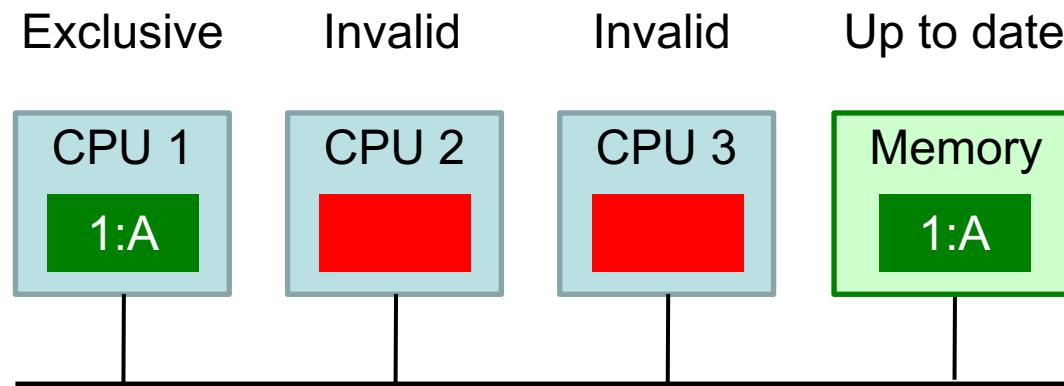
- Popular write-back cache coherence protocol
 - Used e.g. in Intel Core i7
- Each cache line can be in one of the following four states
 1. **Modified** = cache line is valid; memory is invalid; no other copies exist
 2. **Exclusive** = no other cache holds the line; memory up to date;
 3. **Shared** = multiple caches may hold the line; memory up to date
 4. **Invalid** = no valid data

MESI protocol - example



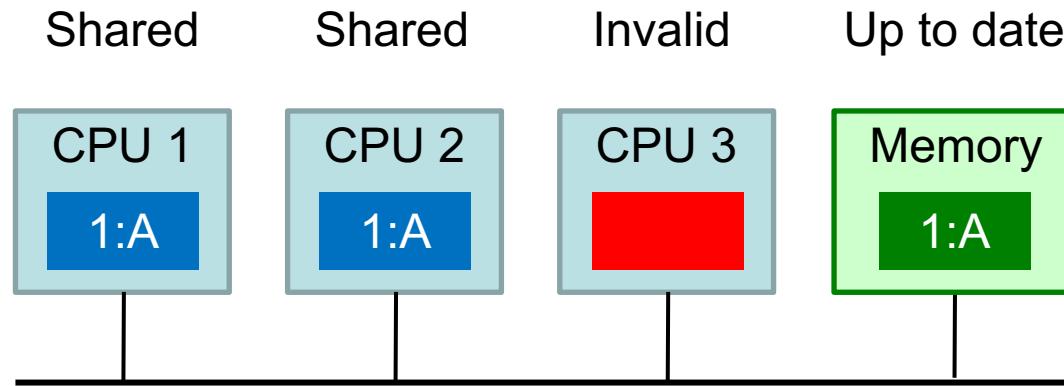
N:X = cache line address : contents version

MESI protocol - example



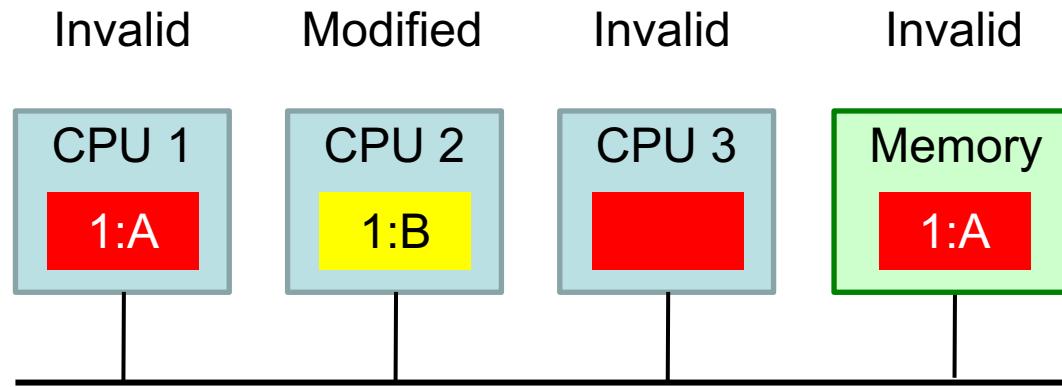
Initial state

MESI protocol - example



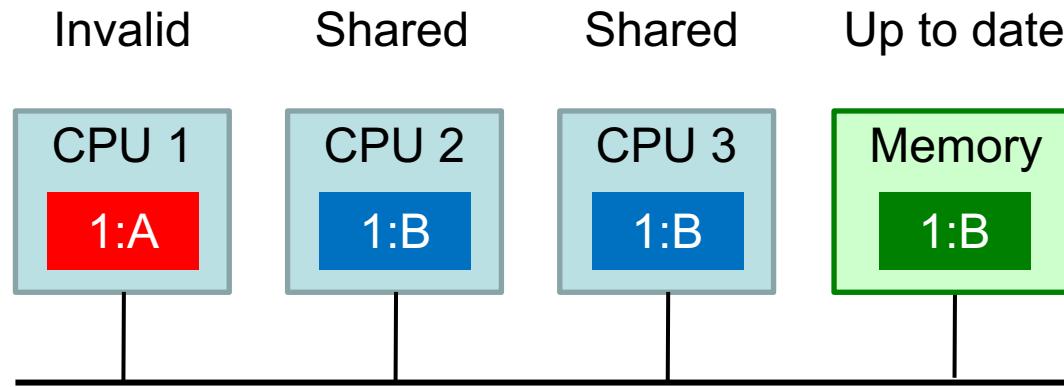
CPU2: read 1

MESI protocol - example



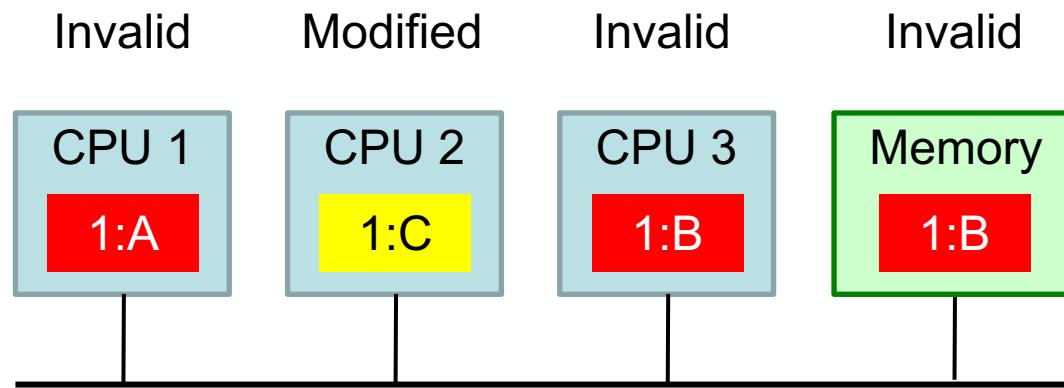
CPU2: write B to 1

MESI protocol - example



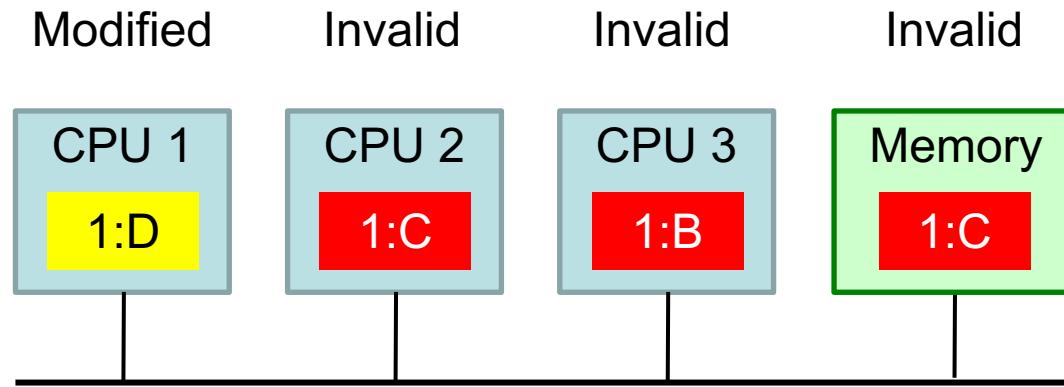
CPU3: read 1

MESI protocol - example



CPU2: write C to 1

MESI protocol - example



CPU1: write D to 1

Coherence misses



- Uniprocessor misses
 - Compulsory
 - Capacity
 - Conflict
- **True sharing miss** – cache miss occurring because a block was invalidated by another processor writing the same word
 - Arises from communication of data through coherence mechanism
 - Independent of cache-line size
- **False sharing miss** – cache miss occurring because a block was invalidated by another processor writing a different word
 - Miss that would not occur if block size were one word

Coherence misses (2)

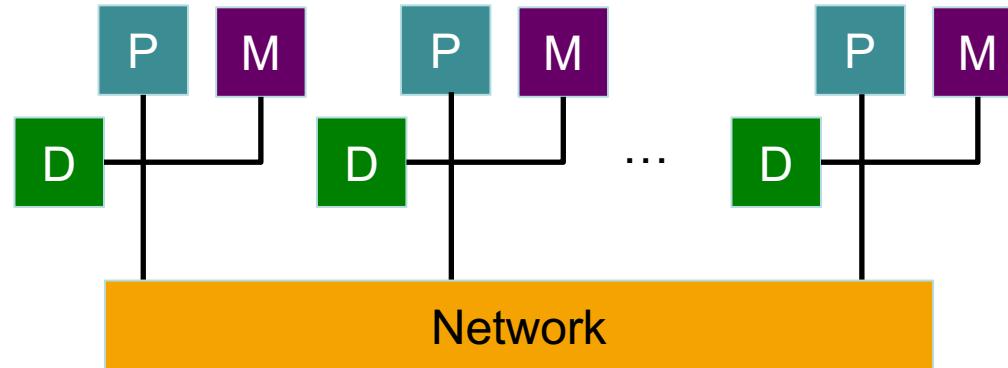
Assume X1 and X2 are in the same cache line, and initially in the caches of P1 and P2 in shared state

Time	P1	P2	Miss
1	Write X1		Invalidate X1 & X2 on P2
2		Read X2	False miss – X1 not used by P2
3	Write X1		Invalidate X1 & X2 on P2
4		Write X2	False miss – X1 not used by P2; Invalidate X1 & X2 on P1
5	Read X2		True miss – P2 has written X2 before

- Coherence misses most important for tightly-coupled applications that share significant amounts of user data

Cache-coherent NUMA systems

- Coherence of caches established via directory
 - Distributed database storing location and status cache of lines
 - Requires fast hardware because it must be queried on every memory reference





Example

- 256 nodes
- 1 CPU + 16 MB RAM per node
- Total memory 2^{32} bytes
- 2^{26} cache lines, 64 bytes each
- Memory statically allocated among nodes
 - 0-16M in node 0, 16-32M in node 1, etc.
- Directory of each node holds entries for the 2^{18} cache lines comprising its 2^{24} bytes of memory
- Assumption: a line can be held in at most one cache



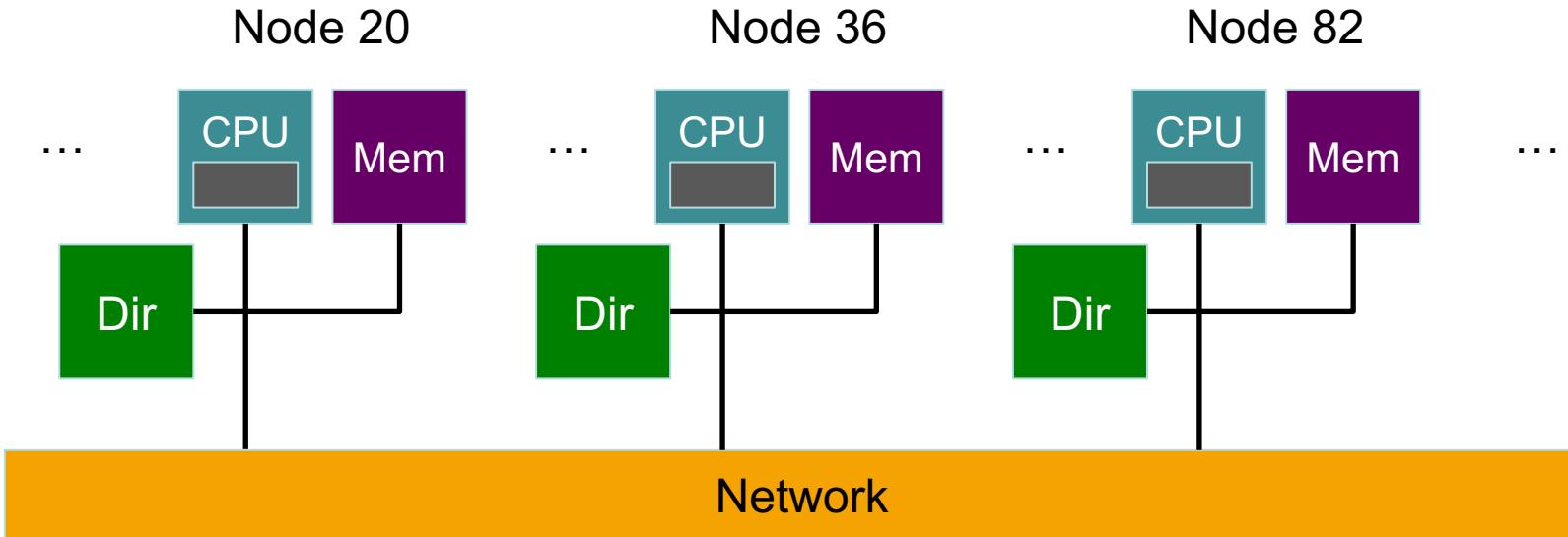
Example (2)

- Consider LOAD instruction from CPU 20 that references a cached line at address 0x24000108
- MMU splits address into three parts



- Address 0x24000108 = node 36 ; line 4 ; offset 8

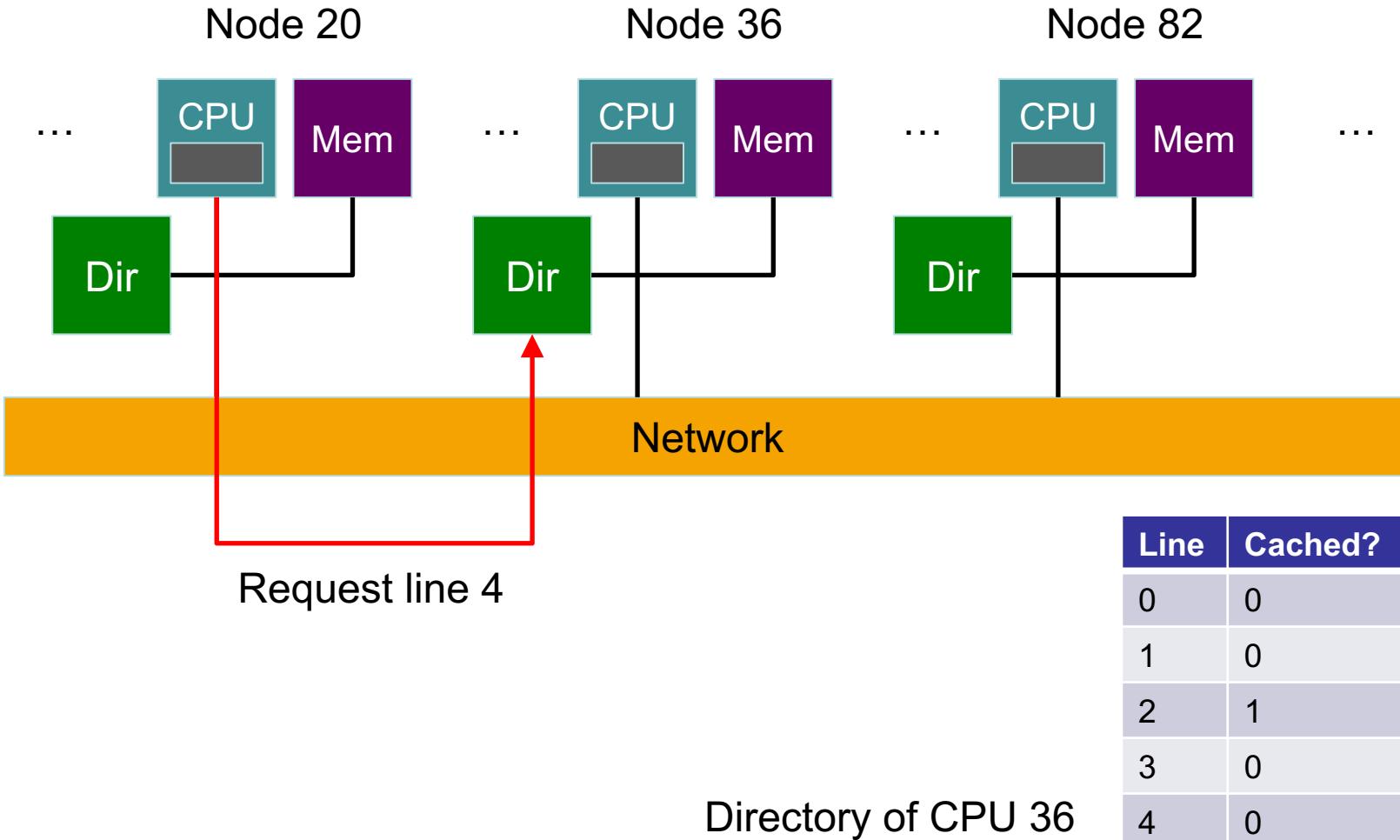
Example (3)



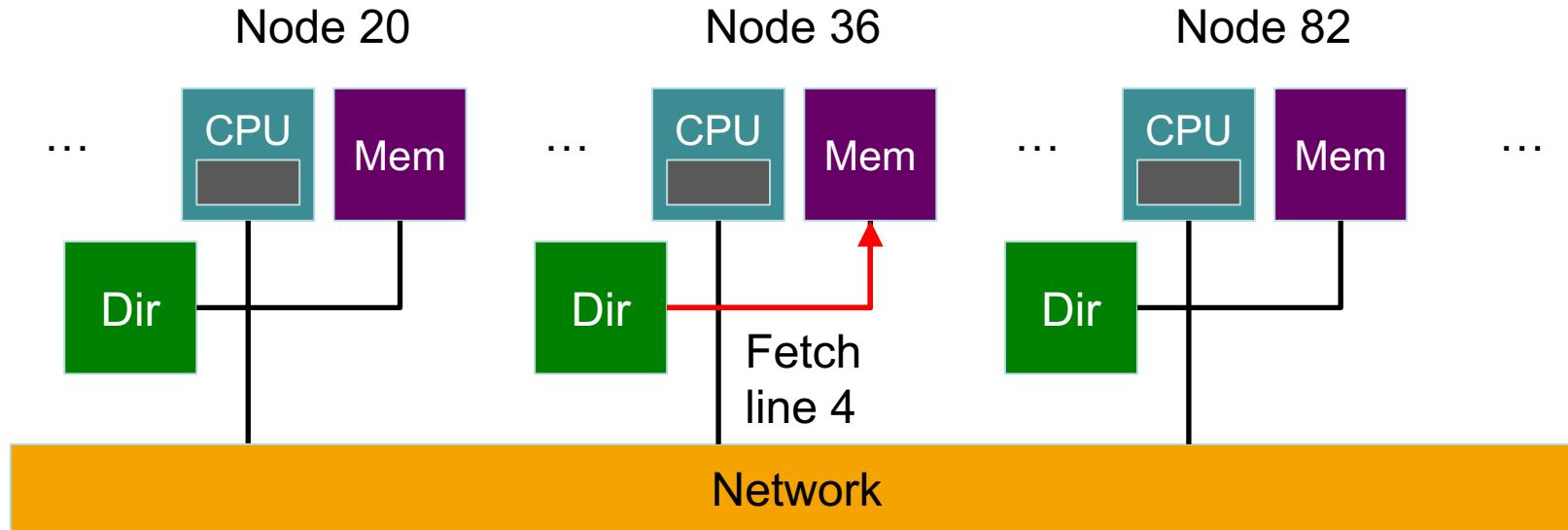
Directory of CPU 36

Line	Cached?	Where?
0	0	
1	0	
2	1	82
3	0	
4	0	

Example (3)



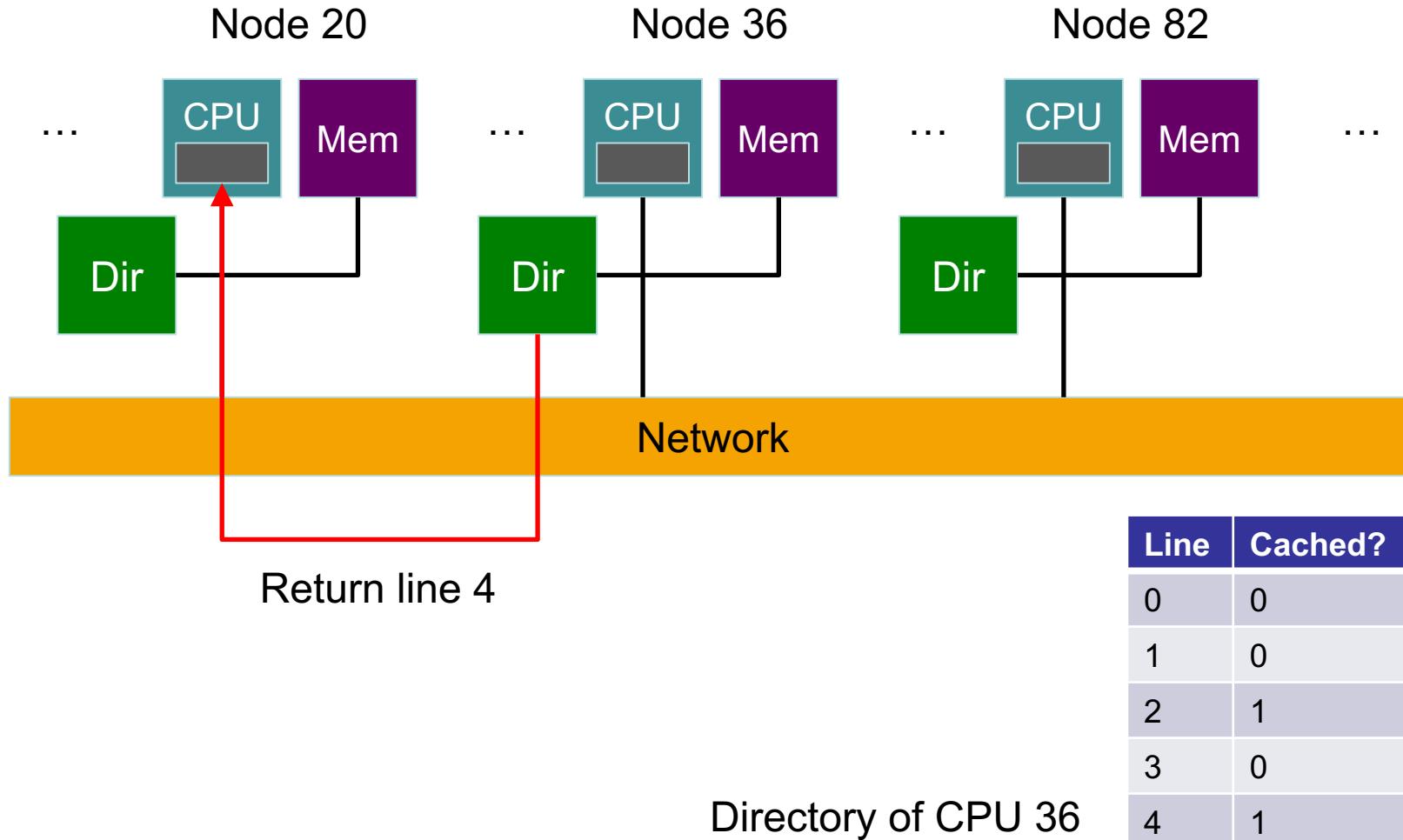
Example (3)



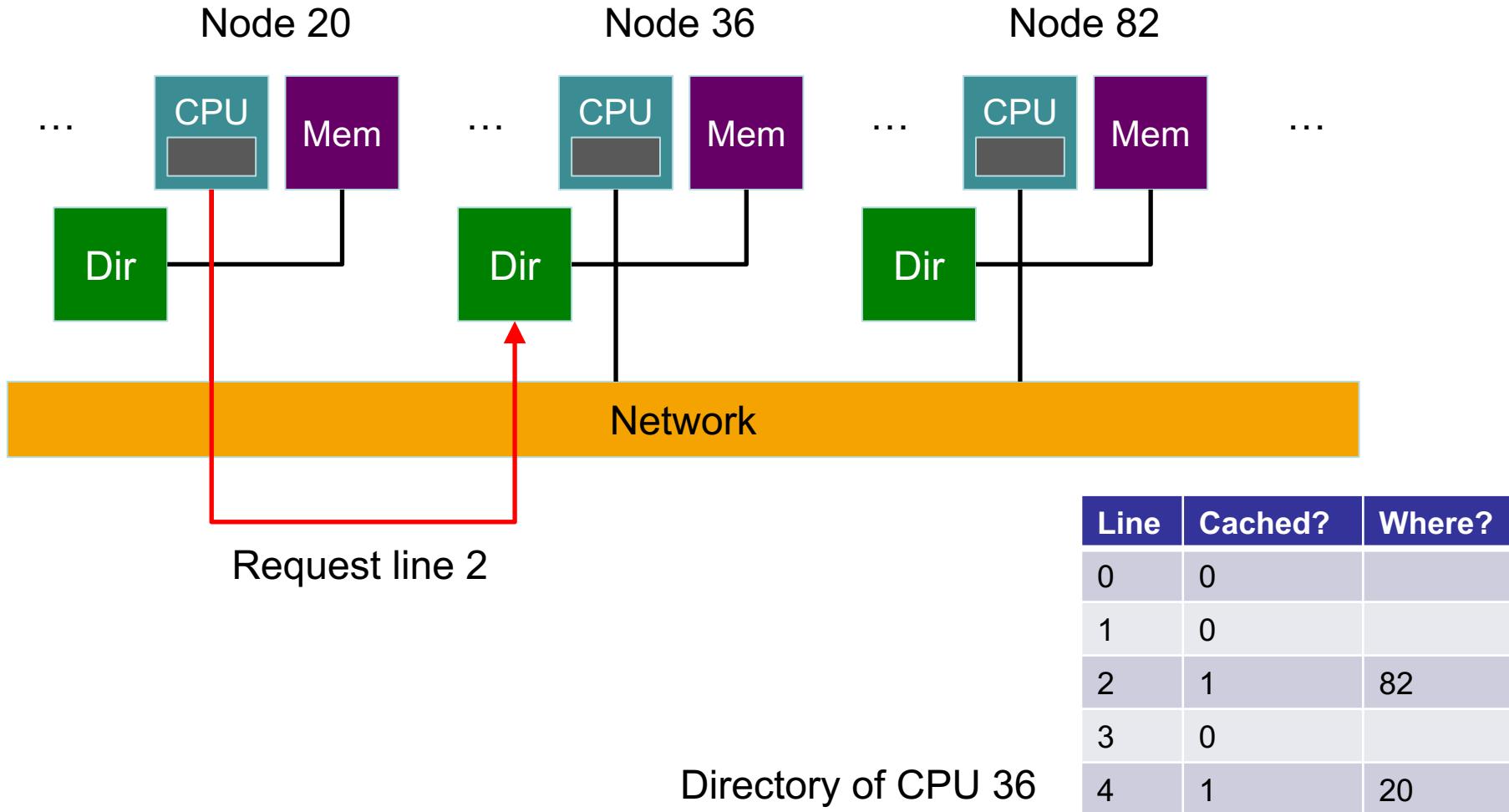
Directory of CPU 36

Line	Cached?	Where?
0	0	
1	0	
2	1	82
3	0	
4	0	

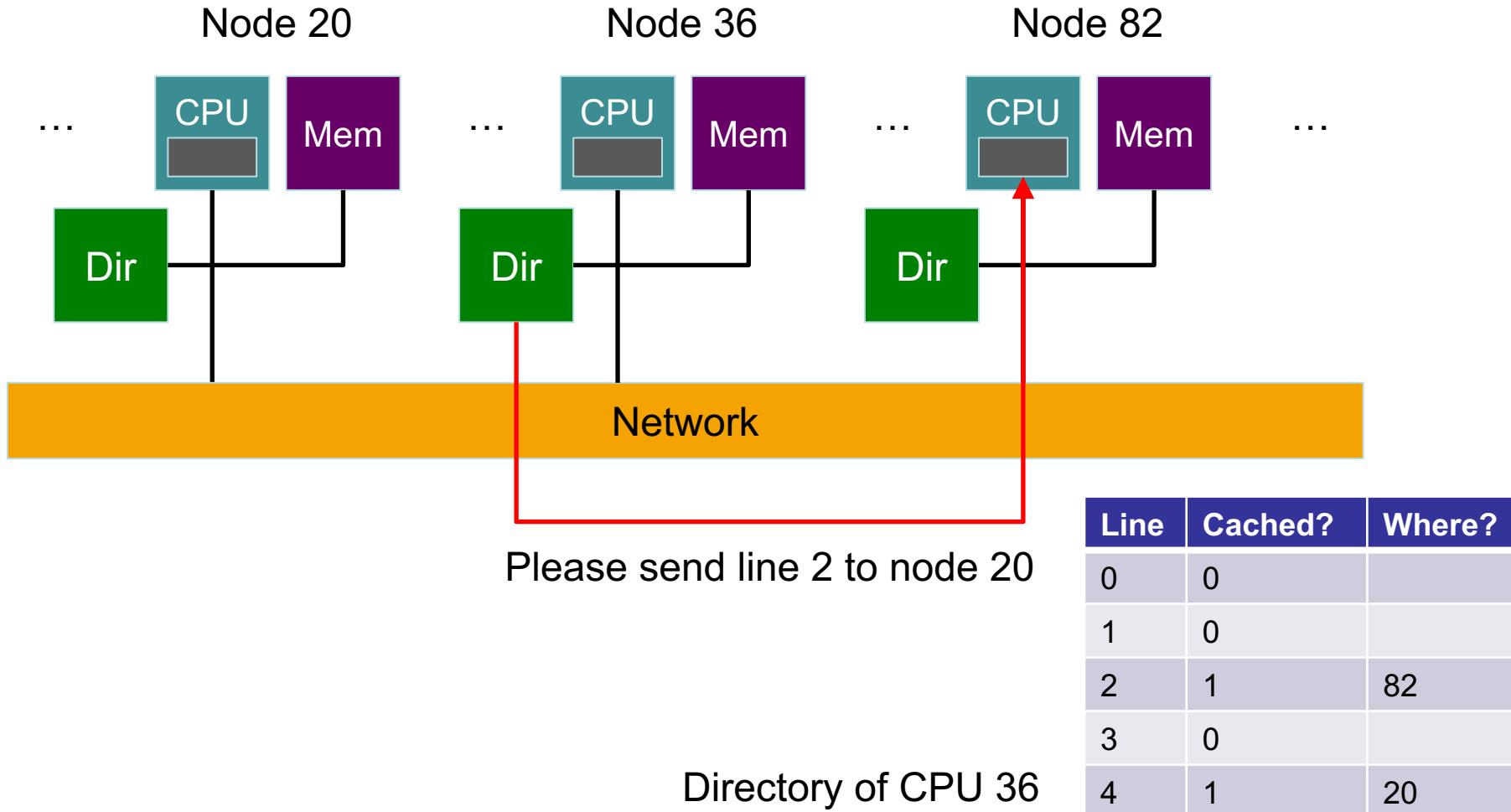
Example (3)



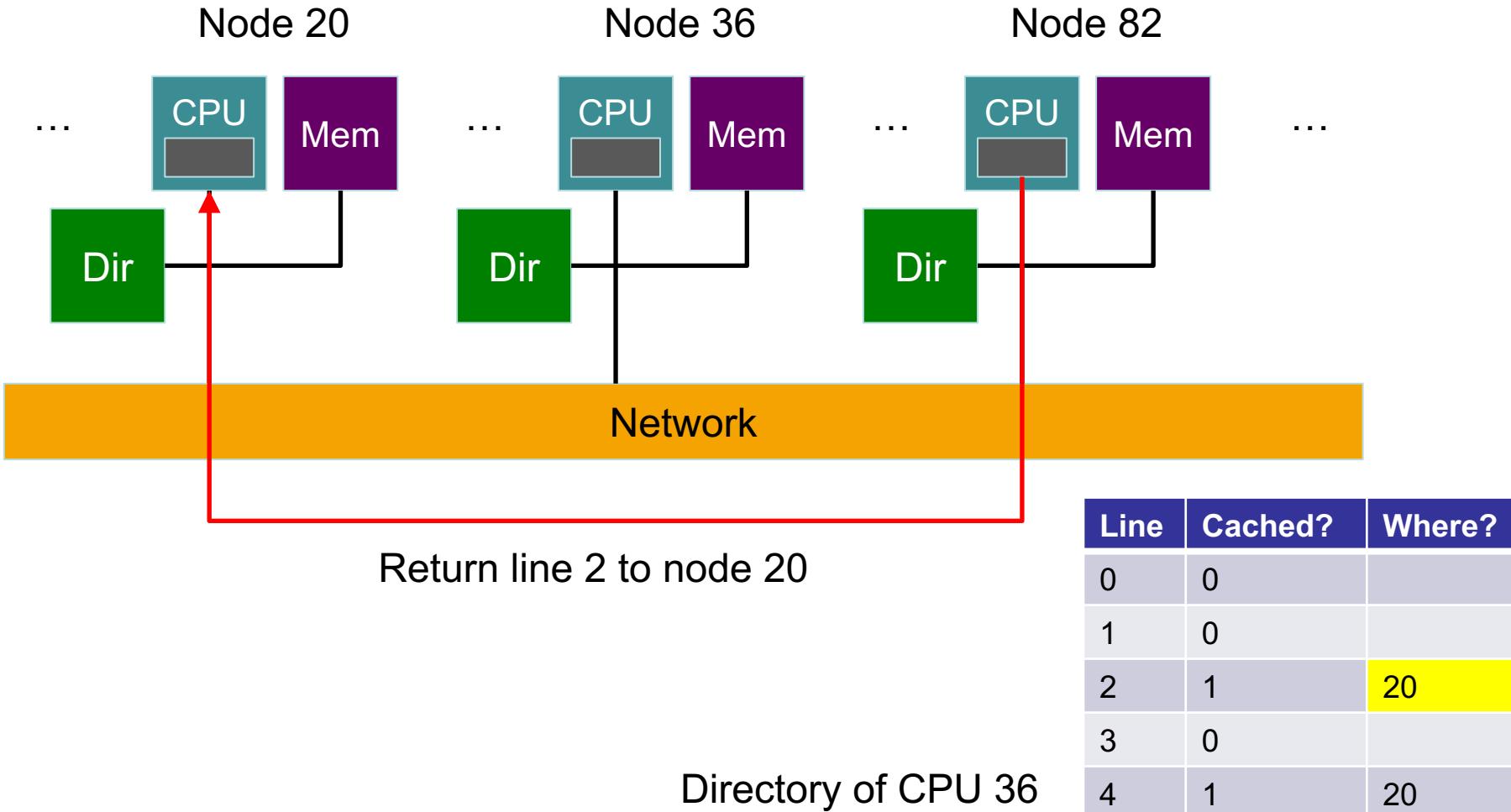
Example (3)



Example (3)



Example (3)





Number of copies

- Directory has 2^{18} 9-bit entries = 1.76% of total memory
- Limitation: a line can be cached at only one node
- Alternatives
 - k entries per line, allowing copies at up to k nodes
 - Bitmap with one bit per node (substantial increase in memory overhead)
 - 8-bit field as head of a linked list (requires extra storage for list pointers and time overhead for searching the list)

Status of a line

- Another optimization is to keep track of a line's status (dirty or clean)
- The home node can satisfy a read request for a clean line from its local memory without having to forward it to another node's cache
- A read request for a dirty line must still be forwarded to the node holding the copy
- No advantage if only one copy is allowed because any request requires invalidation of the previous copy
- Modification of a cached line requires home node to be informed and invalidation of all other copies
- Potentially significant coherence traffic

Memory semantics



- Shared memory = image of a **single shared address space**
 - Promises intuitive programming
- Implementation quite complex in reality
 - Many memory modules, each holding some portion of the physical memory
 - CPUs and memories often connected by complex interconnection network
 - Memory hierarchy with multiple levels (registers down to main memory)

Order of updates

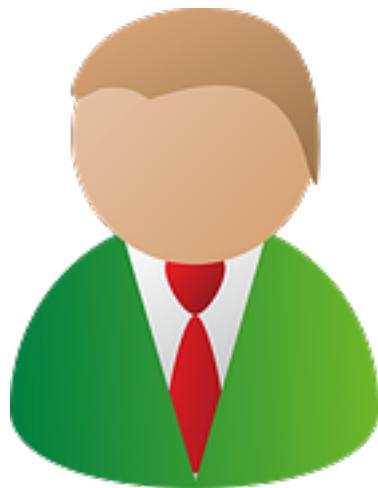
Can be influenced by two factors

- Order in which memory request “messages” arrive not necessarily the same as the one in which they were issued
 - A single thread may observe writes in an order different from the order another thread wrote them
 - Order may even differ among multiple readers
- Compiler may re-order instructions
(even possible on uni-processor systems)

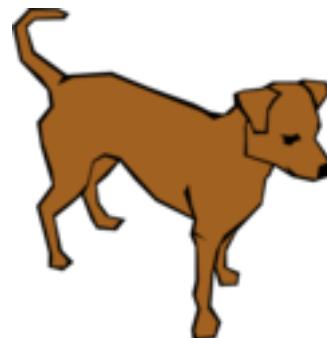
Piecewise update of reality



TECHNISCHE
UNIVERSITÄT
DARMSTADT



Bob

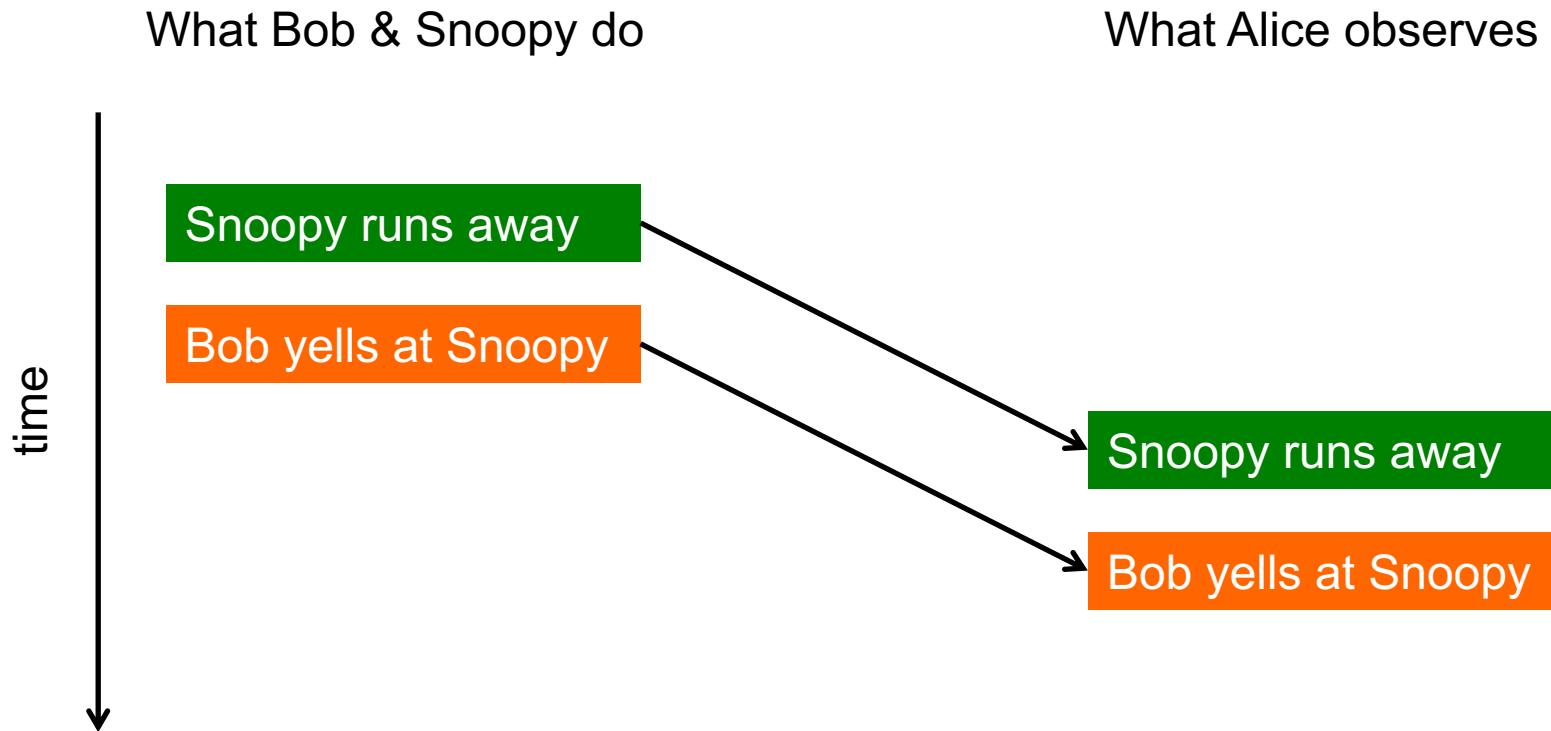


Snoopy
(Bob's dog)

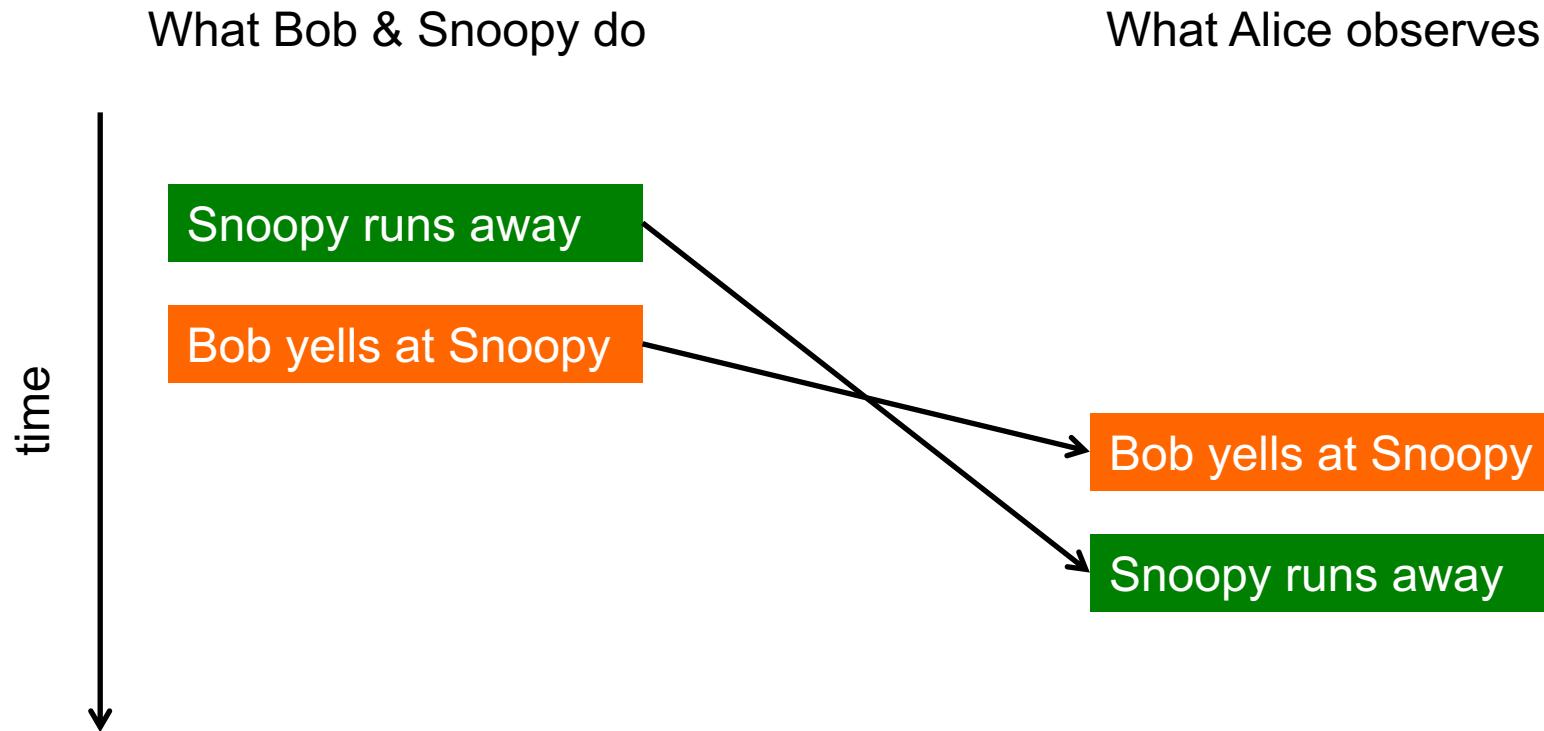


Alice
(observer)

Piecewise update of reality (2)



Perceived update order reversed



Instruction reordering

Thread 1

```
int x;
bool x_init;

void init()
{
    x = initialize();
    x_init = true;
    // ...
}
```

Thread 2

```
extern int x;
extern bool x_init;

void f2()
{
    int y;
    while (!x_init)
        usleep(10000);
    y = x;
    // ...
}
```

Instruction reordering (2)

- Compiler (or hardware instruction scheduler) may decide to execute `x_init = true` first
 - No thread-local dependence
 - Thread 2 may assign an uninitialized x to y
- Since thread 2 makes no assignments to `x_init`, an optimizer may decide to lift evaluation of `x_init` out of the loop
 - Thread 2 may sleep forever or not at all

Memory ordering

- Memory ordering specifies what a programmer can assume about the order in which a thread sees updates of memory locations
- Sequential consistency – simplest memory order
 - Every thread sees the effects of every operation in the same order
 - As if the operations of all the threads were executed in some sequential order (like by a single thread)
 - The operations of each individual thread appear in this sequence in the order specified by the program
 - Usually lowers concurrency
 - Still many possible sequentially consistent orders for a given set of threads

Relaxed consistency models

- Sequential consistency is convenient, but impractical because of potential performance degradation
- Idea of relaxed consistency models
 - Let reads and writes complete out of order
 - Use synchronization to enforce ordering where important
- Relaxed consistency models can be distinguished by the orderings they guarantee / relax

Processor consistency

1. Writes by any CPU are seen by all in the order they were issued
 - Assume CPU 1 issues writes to variable x with values 1A, 1B, 1C
 - No other processor will ever see 1B followed by 1A
2. For every memory word, all CPUs see writes to it in the same order
 - Requires every memory word to have unambiguous value at the end



Release consistency

- Finishing all pending writes after a synchronization expensive
- Instead, release consistency adopts model similar to critical sections
 - When leaving a critical section, a thread does not force all the writes to complete immediately
 - It only ensures that they are completed before any thread enters the critical sections again
- Synchronization split into two parts
 - **Acquire** – get exclusive access to shared data
 - **Release** – indicate that exclusive access is finished

Synchronization



- A program is synchronized if all accesses to shared data are ordered by synchronization operations
- Updates of a single location not ordered by synchronization are called **data races**
- Synchronization allows programs to behave as under sequential consistency even if the architecture implements a more relaxed consistency model
- Building synchronization mechanisms is hard

Synchronization (2)

- Synchronization mechanisms implemented using hardware supplied synchronization instructions
- Uninterruptible instruction or instruction sequence capable of **atomically reading and changing a value** together with the ability to tell whether read and write were performed atomically
- Synchronization can become a bottleneck in
 - Large-scale multiprocessors
 - High-contention situations

Summary

