

Ananth Grama, Computing Research Institute and Computer Sciences, Purdue University.

Limitations of Memory System Performance

Often, the primary bottleneck to performance is not the CPU, rather, it is the memory system. Typical computations only run at 10-50% of peak CPU utilization because of memory bottlenecks. The key question here is how to connect a 50 ns latency memory to a processor that runs a 0.5 ns clock!

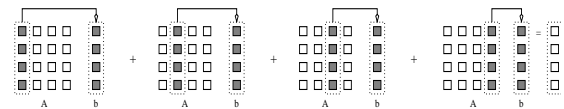
- Improving Effective Memory Latency Using Caches: A hierarchy of small, fast stores bridge the gap between processor and memory. These stores rely on repeated accesses to data to deliver higher aggregate performance.
- Improving Memory System Performance by Threading: Find something else to do while you are waiting for data to arrive from memory.

Ananth Grama, Computing Research Institute and Computer Sciences, Purdue University.

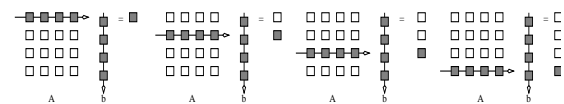
Impact of Strided Access on Program Performance

```
// Fragment 1: Summing columns of a matrix.
for (i = 0; i < 1000; i++)
    column_sum[i] = 0.0;
    for (j = 0; j < 1000; j++)
        column_sum[i] += b[j][i];
```

```
// Fragment 2: Fragment 1, rewritten.
for (i = 0; i < 1000; i++)
    column_sum[i] = 0.0;
for (j = 0; j < 1000; j++)
    for (i = 0; i < 1000; i++)
        column_sum[i] += b[j][i];
```



(a) Column major data access



(b) Row major data access.

Ananth Grama, Computing Research Institute and Computer Sciences, Purdue University.

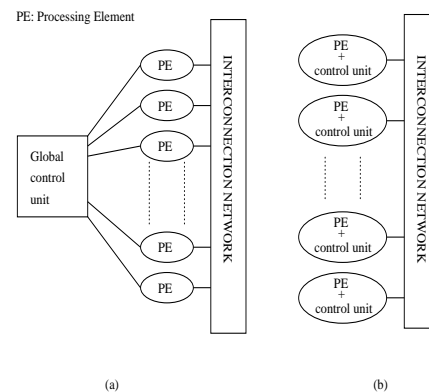
Dichotomy of Parallel Computing Platforms

Control Structure of Parallel Platforms: What is the nature of concurrent tasks?

Communication Model of Parallel Platforms: How do multiple tasks cooperate with each other?

- Message Passing Platforms
- Shared-Address-Space Platforms

Control Structure of Parallel Machines:



(a)

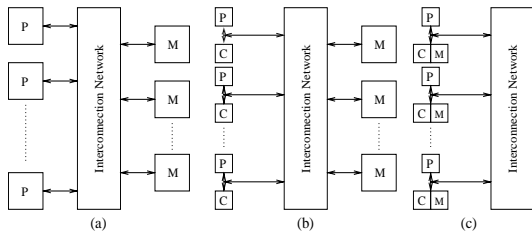
(b)

In a Single Instruction Multiple Data (SIMD) paradigm (a), all processing elements execute the same instruction. In a Multiple Instruction Multiple Data paradigm (b), all processing elements execute possibly different instructions, independently. An intermediate paradigm, called Single Program Multiple Data (SPMD) is the most popular programming paradigm.

Ananth Grama, Computing Research Institute and Computer Sciences, Purdue University.

Ananth Grama, Computing Research Institute and Computer Sciences, Purdue University.

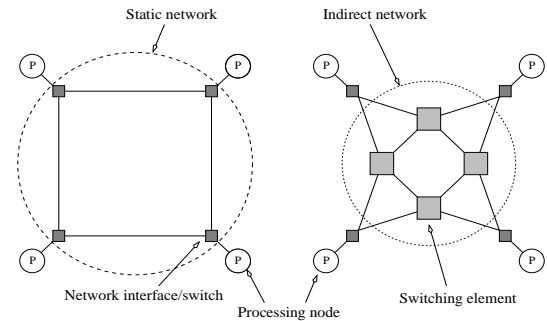
Communication Model of Parallel Platforms



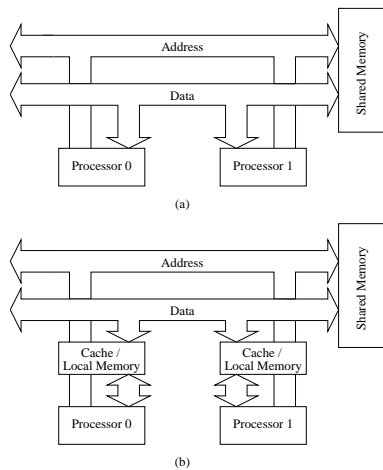
In a uniform memory access (UMA) shared memory model (a), all processors access memory through an interconnect. In (b), we show a UMA shared memory machine with caches, and in (c), we show a non-uniform memory access (NUMA) shared memory machine with private memories only. The logical name for all of these platforms is a shared address space machine.

Physical Organization of Parallel Platforms

The characterizing feature of parallel platforms is the underlying interconnection network. These networks can be static (a) or dynamic (b).

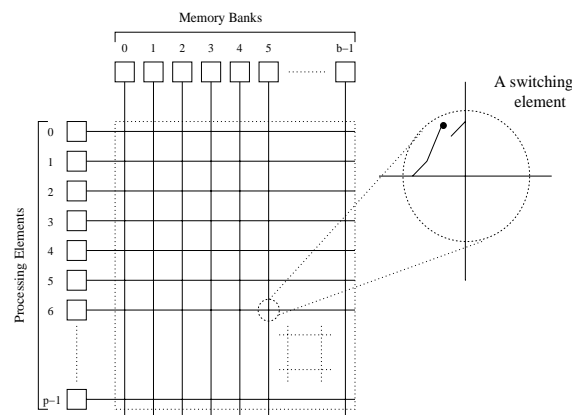


Direct Interconnection Networks:



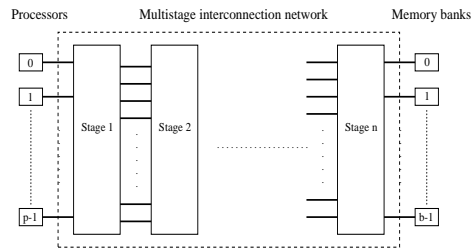
Bus-based interconnects (without (a) and with caches (b)) were the first networks in early commercially available platforms (Sequent Symmetry/Balance).

Direct Interconnection Networks:



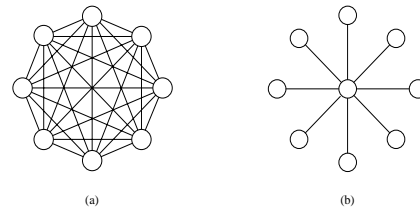
The other extreme in terms of performance and cost compared to buses, is the crossbar network.

Direct Interconnection Networks

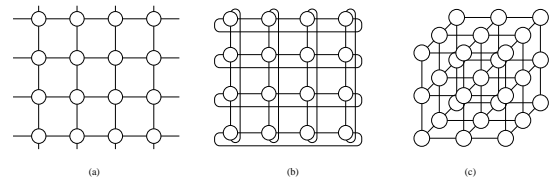


Multistage networks such as the Omega network fall between buses and crossbars in terms of cost and performance.

Static Interconnection Networks:



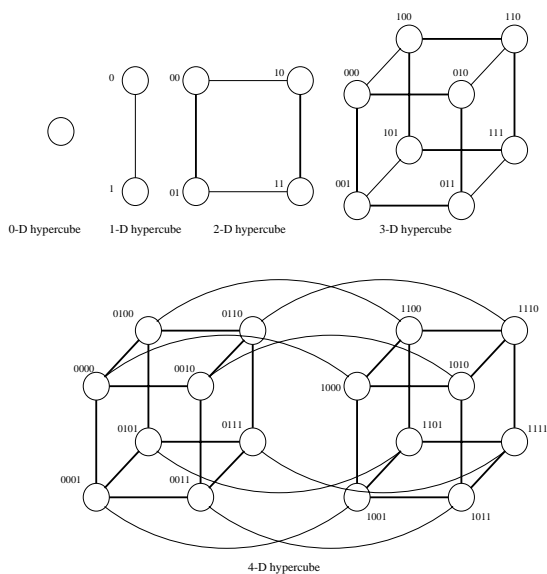
A completely connected network (a) is the strongest model for a static network. A star connected network (b) is the other extreme.



Meshes (2 and 3-D) are popular interconnects because of their desirable layout properties and performance for physical simulations.

Static Interconnection Networks:

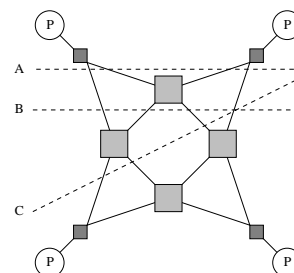
Hypercubes provide another popular interconnect.



Metrics for Interconnection Networks

Performance Metrics

- Link Bandwidth: how fat is each link?
- Arc Connectivity: how many links do I have to remove to separate a network.
- Bisection Bandwidth: how many pairs of people can have conversations at any given time, independently.

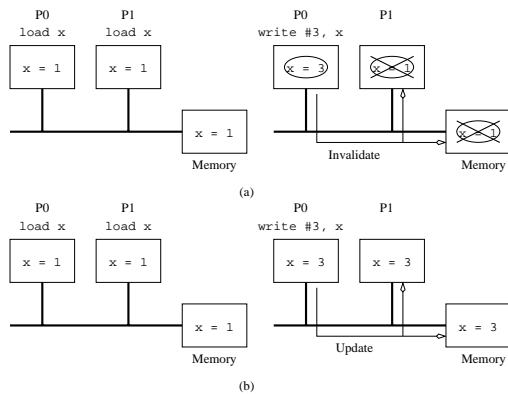


Cost Metrics

- Number of Links
- Layout Costs.

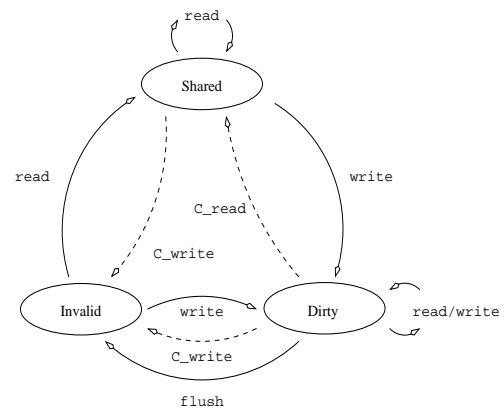
Cache Coherence in Multiprocessor Systems:

How do you deal with multiple copies of the same data item, being manipulated by different processors?



We can invalidate copies (a), or update them (b) when a processor changes a copy.

Coherence Protocols:



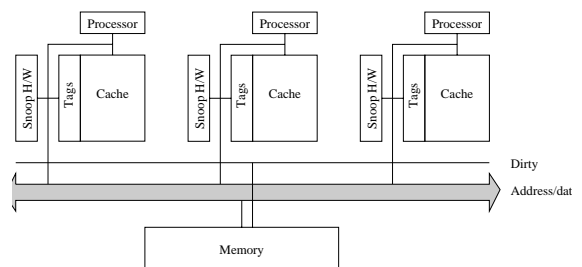
A simple three-state protocol for implementing cache-coherence.

Coherence Protocols:

An example of the coherence protocol.

Time ↓	Instruction at Processor 0	Instruction at Processor 1	Variables and their states at Processor 0	Variables and their states at Processor 1	Variables and their states in Global mem.
					x = 5, D y = 12, D
	read x		x = 5, S		x = 5, S
		read y		y = 12, S	y = 12, S
	x = x + 1		x = 6, D		x = 5, I
		y = y + 1		y = 13, D	y = 12, I
	read y		y = 13, S		y = 13, S
		read x	x = 6, S	x = 6, S	x = 6, S
	x = x + y		x = 19, D	x = 6, I	x = 6, I
		y = x + y	y = 13, I	y = 19, D	y = 13, I
	x = x + 1		x = 20, D		x = 6, I
		y = y + 1		y = 20, D	y = 13, I

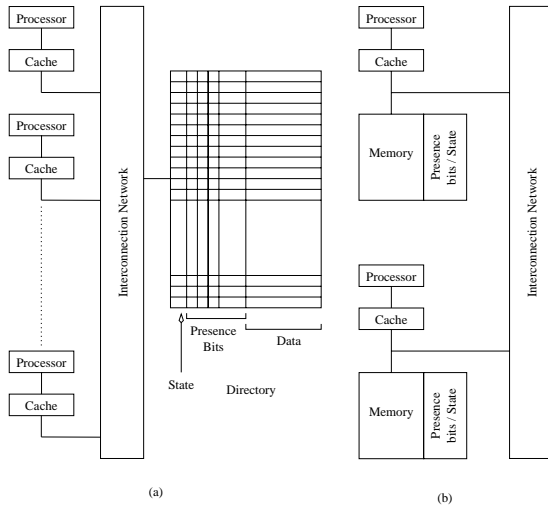
Implementing Coherence Protocols:



Using a snoopy bus to implement the coherence protocol. Snoopy buses do not scale to large configurations.

Implementing Coherence Protocols:

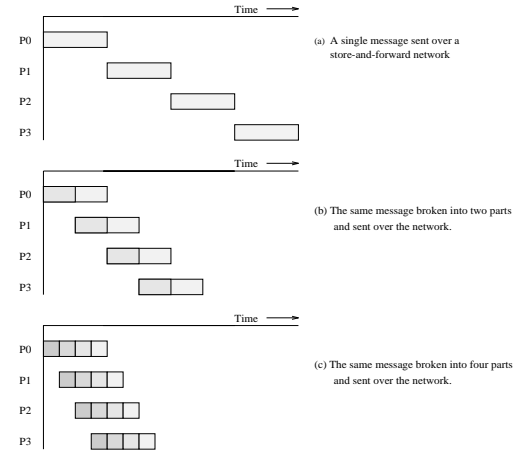
Directories provide a more scalable solution than snoopy buses.



Ananth Grama, Computing Research Institute and Computer Sciences, Purdue University.

Routing Messages in Parallel Platforms:

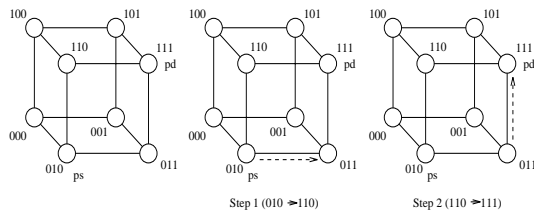
- Store-and-forward routing
- Packet Routing
- Cut-Through routing



Ananth Grama, Computing Research Institute and Computer Sciences, Purdue University.

Routing in Parallel Platforms:

- Dimension ordered / E-cube routing



- Hot-potato routing
- Randomized Routing

Ananth Grama, Computing Research Institute and Computer Sciences, Purdue University.

Embeddings and Overhead:

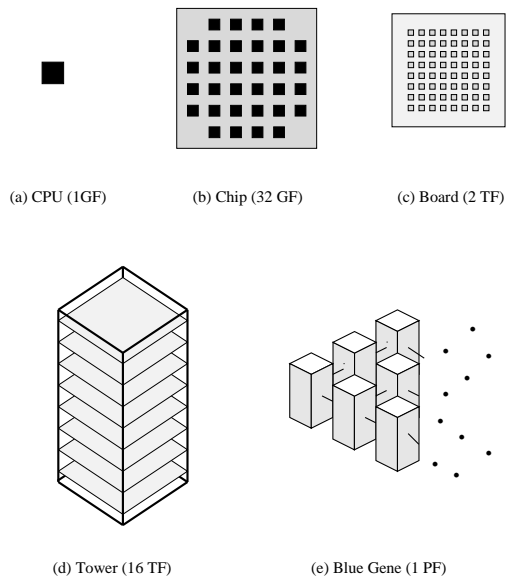
It is often possible to map a weaker architecture on to a stronger one with no performance overhead. Conversely, mapping a stronger architecture on to a weaker one results in performance penalties, depending on the program.

Algorithms for mapping between structured networks are well studied (see handout).

Ananth Grama, Computing Research Institute and Computer Sciences, Purdue University.

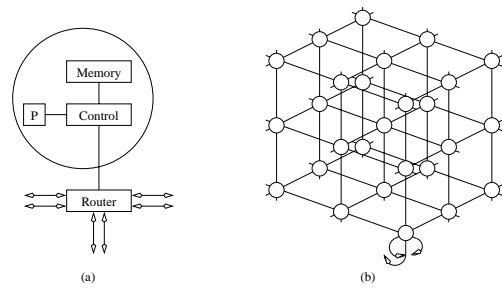
Case Studies:

The IBM Blue Gene.



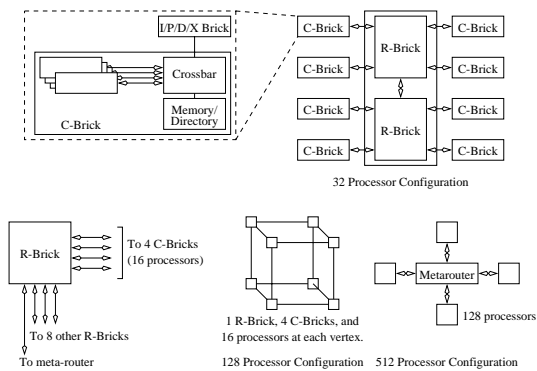
Ananth Grama, Computing Research Institute and Computer Sciences, Purdue University.

The Cray T3E.



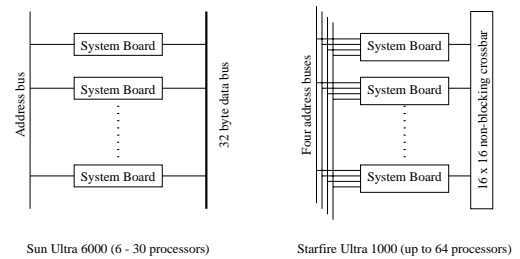
Ananth Grama, Computing Research Institute and Computer Sciences, Purdue University.

SGI Origin 3000



Ananth Grama, Computing Research Institute and Computer Sciences, Purdue University.

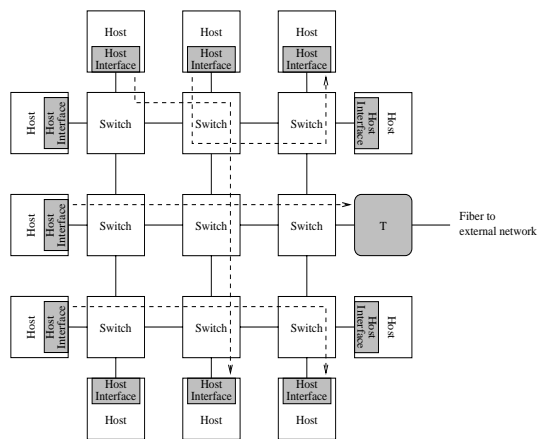
Sun Enterprise Servers



Ananth Grama, Computing Research Institute and Computer Sciences, Purdue University.

Commonly used networks:

Myrinet:



Other networks include Gigabit Ethernet, FiberChannel, HiPPI, etc.