# ECSE 420
# Programming Models

Zeljko Zilic

McConnell Engineering Building

Room 546

# Reminder: Grading Scheme

○ 40% homeworks (4)

○ 30% exam

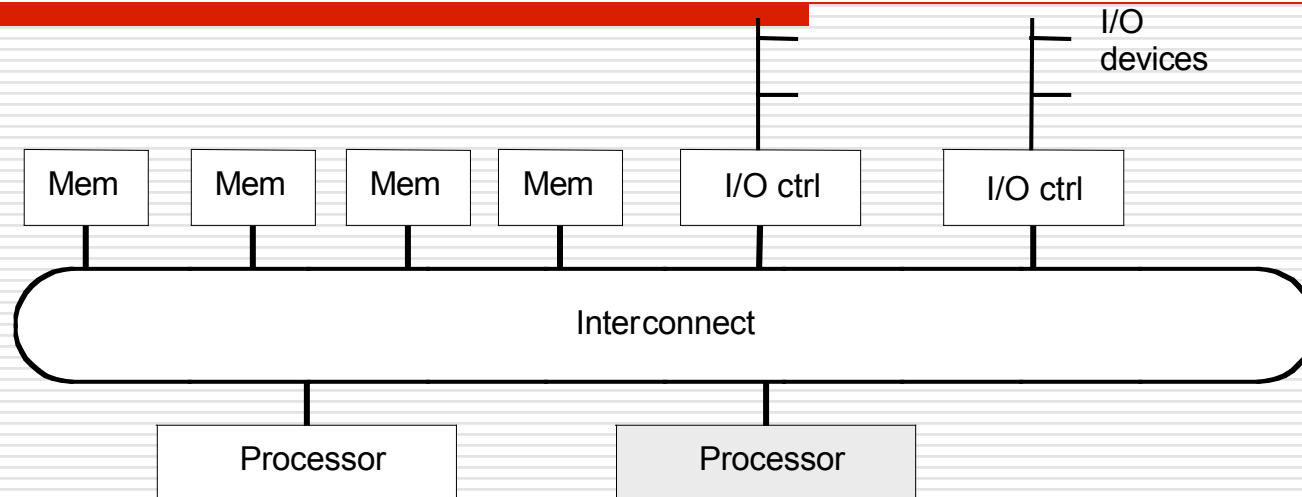○ 30% project (teams of 1-2)

# Programming Models

○ *Conceptualization of the machine that programmer uses in coding applications*

 ■ How parts cooperate and coordinate their activities
 ■ Specifies communication and synchronization operations

○ Multiprogramming

 ■ no communication or synch. at program level

○ *Shared address space*

 ■ like bulletin board

○ *Message passing*

 ■ like letters or phone calls, explicit point to point

○ *Data parallel*:

 ■ more regimented, global actions on data
 ■ Implemented with shared address space or message passing

# Shared Memory  (Shared Address Space)

- ○ Bottom-up engineering factors
- ○ Programming concepts
- ○ Why its attactive

ECSE 420/520
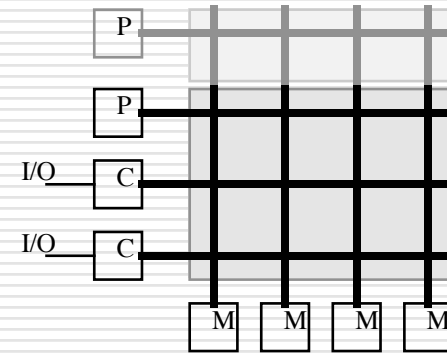Parallel Computing

# Adding Processing Capacity

```
                                                    ┌─ I/O
                                              │     │  devices
                                              │     │
  ┌─────┐ ┌─────┐ ┌─────┐ ┌─────┐ ┌─────────┐ ┌─────────┐
  │ Mem │ │ Mem │ │ Mem │ │ Mem │ │ I/O ctrl│ │ I/O ctrl│
  └──┬──┘ └──┬──┘ └──┬──┘ └──┬──┘ └────┬────┘ └────┬────┘
  ╭──┴───────┴───────┴───────┴──────────┴───────────┴────╮
  │                   Interconnect                        │
  ╰──────────────┬─────────────────────┬─────────────────╯
          ┌──────┴──────┐       ┌──────┴──────┐
          │  Processor  │       │  Processor  │
          └─────────────┘       └─────────────┘
```

- ○ Memory capacity increased by adding modules
- ○ I/O by controllers and devices
- ○ Add processors for processing!
  - ■ For higher-throughput multiprogramming, or parallel programs
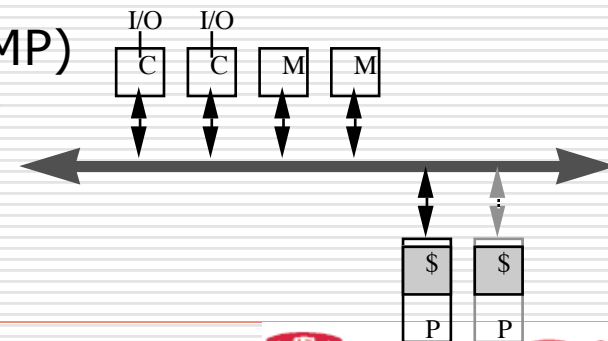
McGill

# Historical Development

○ "Mainframe" approach
- Motivated by multiprogramming
- Extends crossbar used for Mem and I/O
- Processor cost-limited => crossbar
- Bandwidth scales with $p$
- High incremental cost
  - use multistage instead



○ "Minicomputer" approach
- Almost all microprocessors have bus
- Motivated by multiprogramming, TP
- Used heavily for parallel computing
- Called symmetric multiprocessor (SMP)
- Latency larger than for uniprocessor
- Bus is bandwidth bottleneck
  - caching is key: _coherence_ problem
- Low incremental cost

ECSE 420/520
Parallel Computing

# Shared Physical Memory

- Any processor can directly reference any memory location

- Any I/O controller - any memory

- Operating system can run on any processor, or all.
    - OS uses shared memory to coordinate
- Communication occurs implicitly as result of loads and stores

- What about application processes?
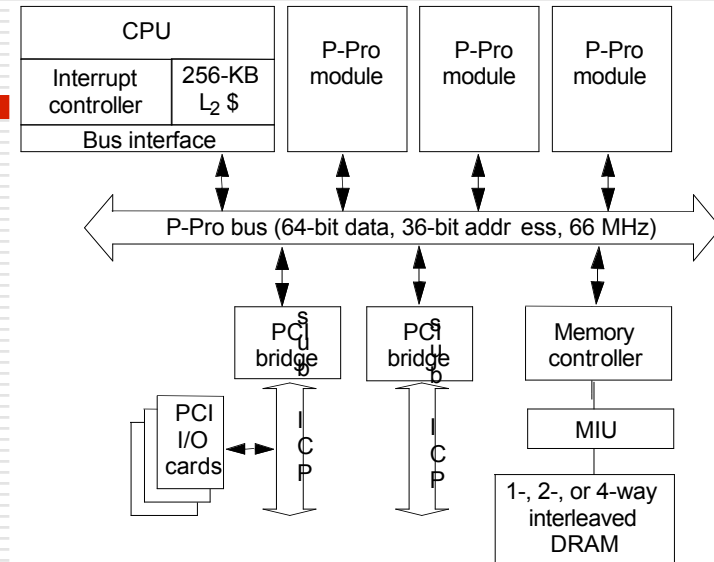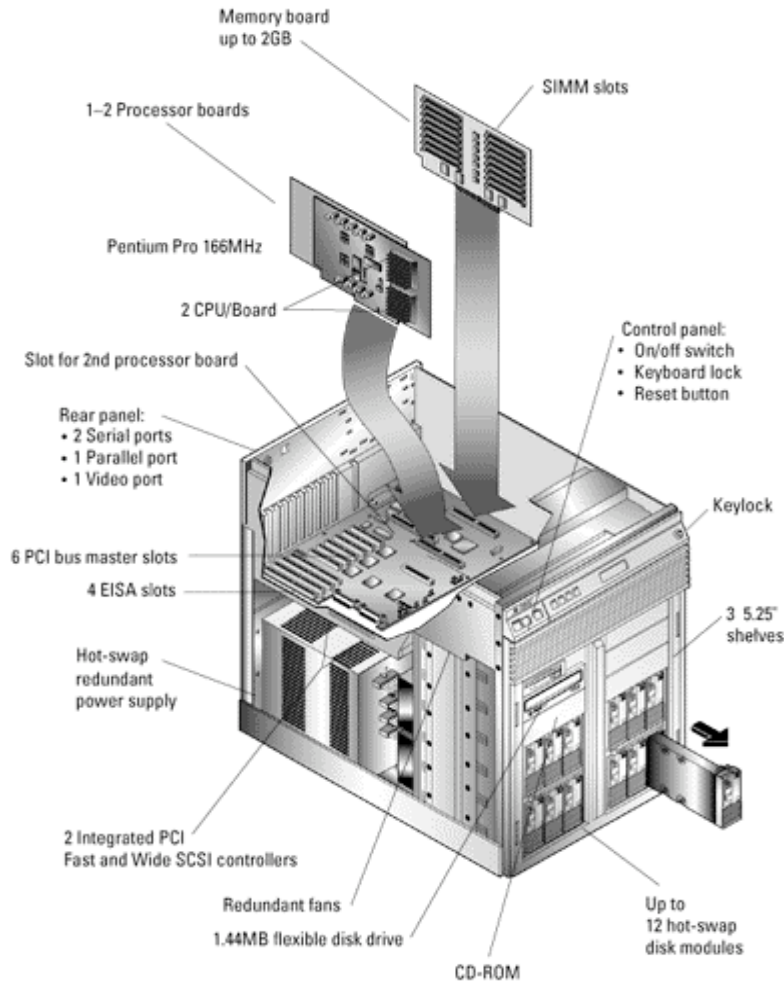
McGill

# Shared Virtual Address Space

- Process = address space plus thread of control
- Virtual-to-physical mapping can be established so that processes shared portions of address space.
  - User-kernel or multiple processes
- Multiple threads of control on one address space.
  - Popular approach to structuring OS's
  - Now standard application capability (ex: POSIX threads)
- Writes to shared address visible to other threads
  - Natural extension of uniprocessors model
  - conventional memory operations for communication
  - special atomic operations for synchronization
    - also load/stores
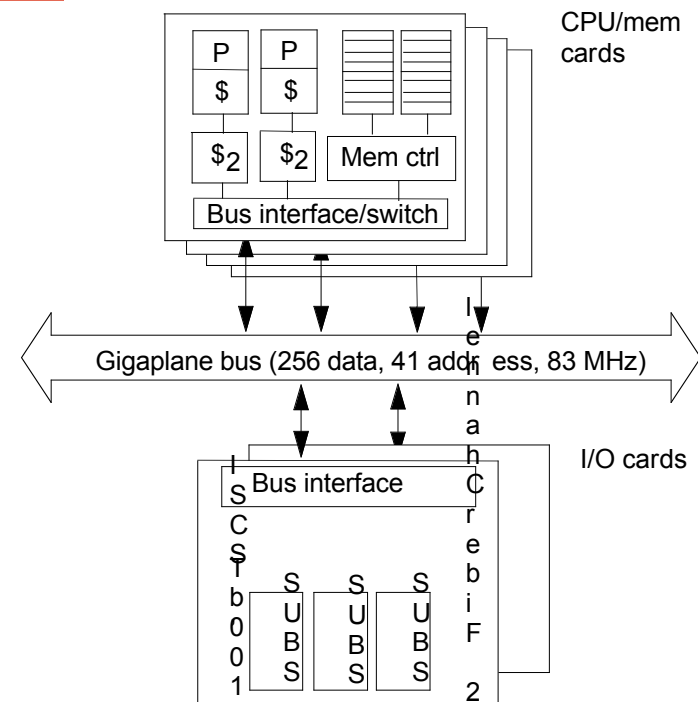
# Structured Shared Address Space

Virtual address spaces for a collection of processes communicating via shared addresses

Machine physical address space

Load

$P_n$

$P_2$

$P_1$

$P_0$

Store

Shared portion of address space

Private portion of address space

$P_n$ private

Common physical addresses

$P_2$ private

$P_1$ private

$P_0$ private

- Add hoc parallelism used in system code
- Most parallel applications have structured SAS
- Same program on each processor
  - shared variable X means the same thing to each thread

McGill

# Engineering: Intel Pentium Pro Quad



| CPU | | P-Pro module | P-Pro module | P-Pro module |
|---|---|---|---|---|
| Interrupt controller | 256-KB $L_2$ $ | | | |
| Bus interface | | | | |

P-Pro bus (64-bit data, 36-bit addr ess, 66 MHz)

| PCI bridge | PCI bridge | Memory controller |
|---|---|---|

PCI I/O cards

MIU

1-, 2-, or 4-way interleaved DRAM

- **All coherence and multiprocessing glue in processor module**
- **Highly integrated, targeted at high volume**
- **Low latency and bandwidth**

# Engineering: SUN Enterprise



CPU/mem cards

| P | P | | |
| $ | $ | | |
| $2 | $2 | Mem ctrl | |

Bus interface/switch

Gigaplane bus (256 data, 41 address, 83 MHz)

I/O cards

Bus interface

SUBS  SUBS  SUBS

○ Proc + mem card - I/O card

- 16 cards of either type
- All memory accessed over bus, so symmetric
- Higher bandwidth, higher latency bus

# Scaling Up



"Dance hall"                    Distributed memory

- Problem is interconnect: cost (crossbar) or bandwidth (bus)
- Dance-hall:  bandwidth still scalable, but lower cost than crossbar
  - latencies to memory uniform, but uniformly large
- Distributed memory or non-uniform memory access (NUMA)
  - Construct shared address space  out of simple message transactions across a general-purpose network (e.g. read-request, read-response)
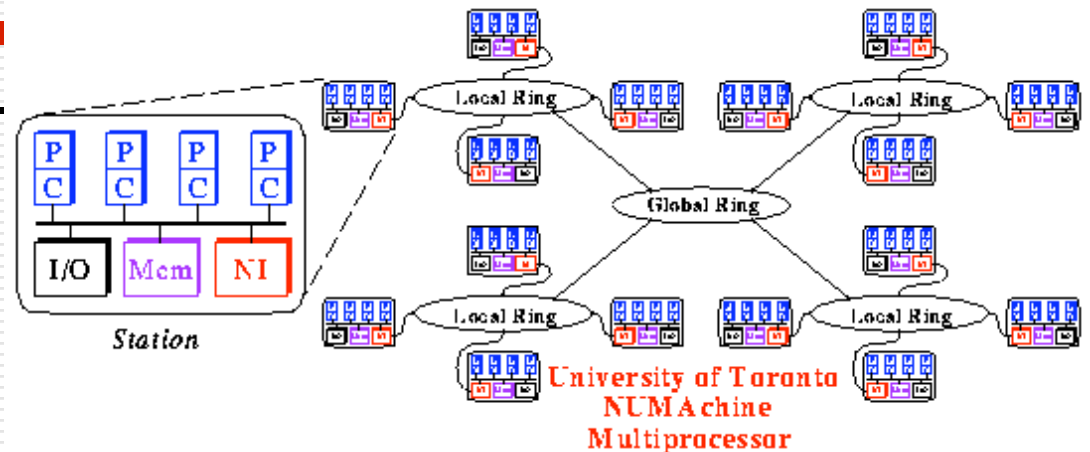- Caching shared (particularly nonlocal) data?

McGill

# Engineering: Cray T3E



- Scale up to 1024 processors, 480MB/s links
- Memory controller generates request message for non-local references
- No hardware mechanism for coherence
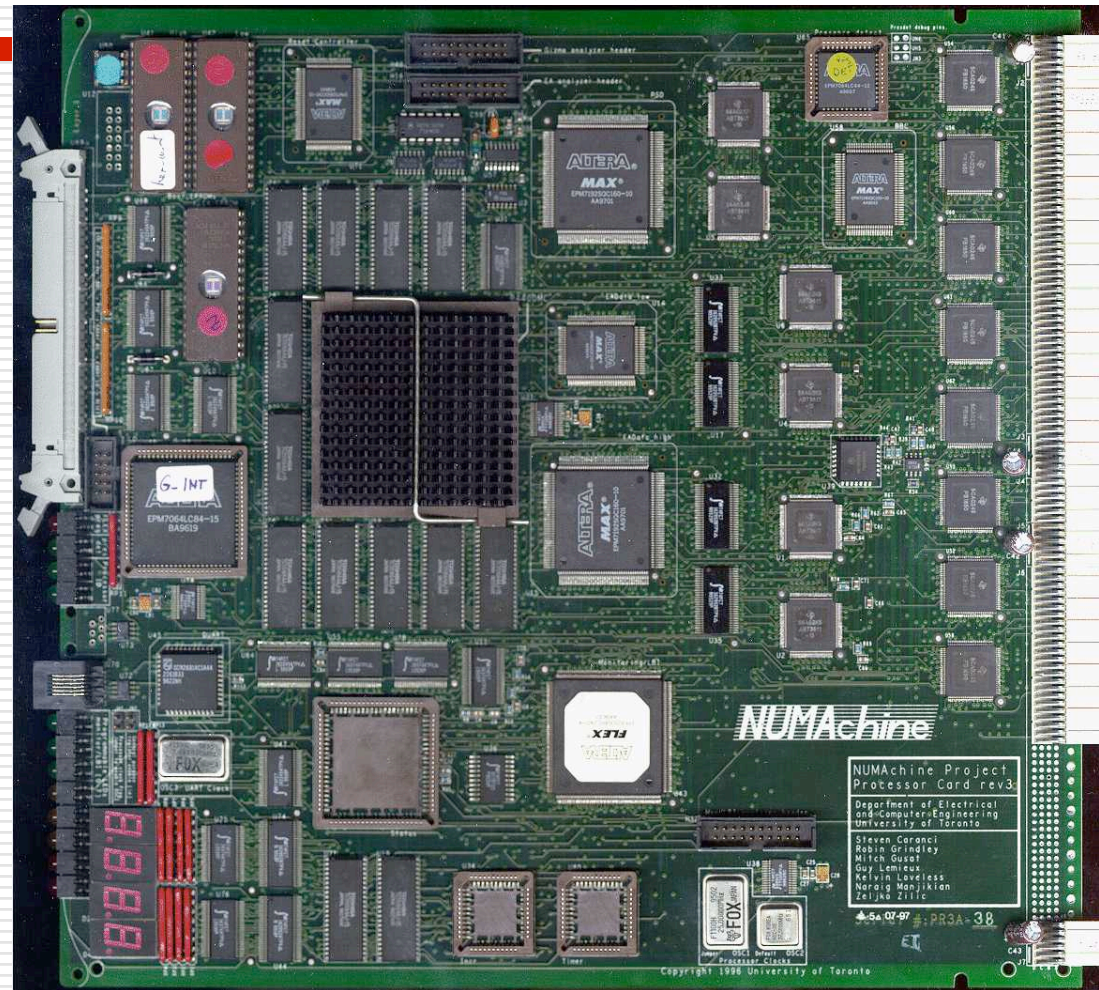  - SGI Origin etc. provide this

McGill

# U. Toronto NUMAchine



University of Toronto NUMAchine Multiprocessor

- Working state-of-the art cache coherent shared-memory multiprocessor

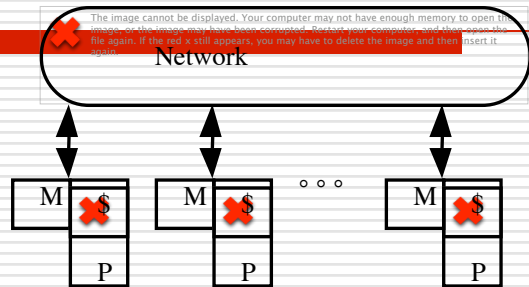  - Developed on a "shoebox" budget

- 64 processors (MIPS 4400)

McGill

# NUMAchine Processor Board

- Most complexity of the overall system

- Logic implemented completely with programmable logic

McGill

# Message Passing Approach

Network

M $ M $ ∘ ∘ ∘ M $

P P P

Systolic
Arrays

Generic
Architecture

SIMD

Dataflow

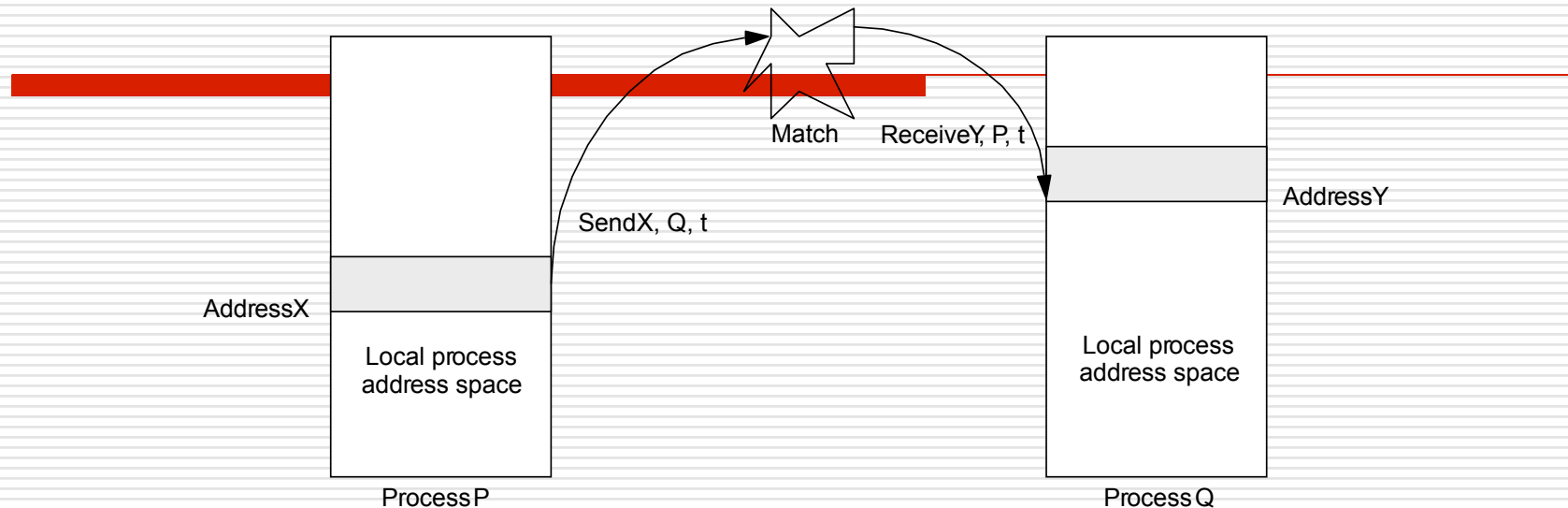Message Passing

Shared Memory

ECSE 420/520
Parallel Computing

McGill

# Message Passing Architectures

○ Complete computer as building block, including I/O
- Communication via explicit I/O operations

○ Programming model
- direct access only to private address space (local memory),
- communication via explicit messages (send/receive)

○ High-level block diagram



- Communication integration?
  ○ Mem, I/O, LAN, Cluster
- Easier to build and scale than SAS

○ Programming model more removed from basic hardware operations
- Library or OS intervention
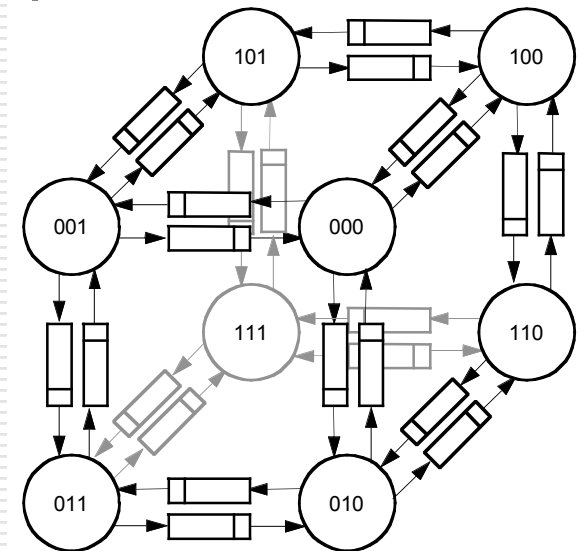
McGill

# Message-Passing Abstraction



- Send specifies buffer to be transmitted and receiving process
- Recv specifies sending process and application storage to receive to
- Memory to memory copy, but need to name processes
- Optional tag on send and matching rule on receive
- User process names local data and entities in process/tag space too
- In simplest form, send/recv match achieves pairwise synch event
  - Other variants too
- Many overheads: copying, buffer management, protection

# Evolution of Message-Passing Machines

○ Early machines: FIFO on each link

  ■ HW close to prog. Model;

  ■ synchronous ops

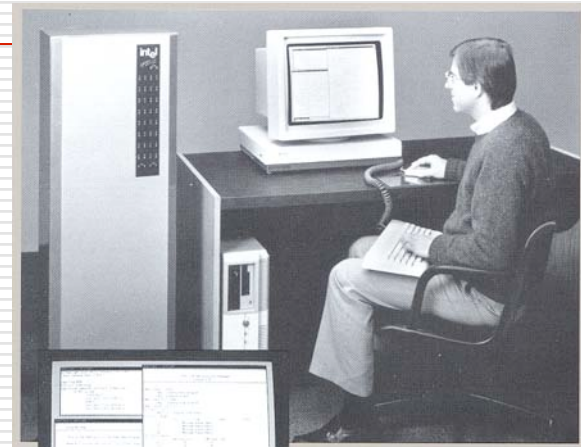  ■ topology central (hypercube algorithms)

**CalTech Cosmic Cube (Seitz, CACM Jan 95)**

ECSE 420/520
Parallel Computing

# Diminishing Role of Topology



○ Shift to general links
  ▪ DMA, enabling non-blocking ops
    ○ Buffered by system at destination until recv
  ▪ Store&forward routing
○ Diminishing role of topology
  ▪ Any-to-any pipelined routing

**Intel iPSC/1 -> iPSC/2 -> iPSC/860**

  ▪ Node-network interface dominates communication time

$$H \times (T_0 + n/B)$$
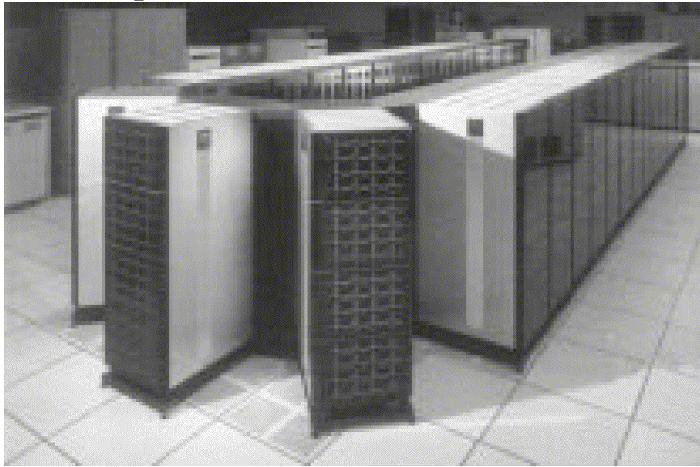$$\text{vs}$$
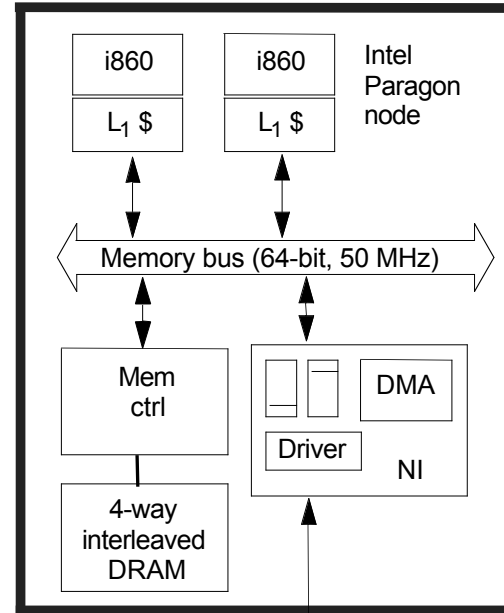$$T0 + H\Delta + n/B$$



  ▪ Simplifies programming
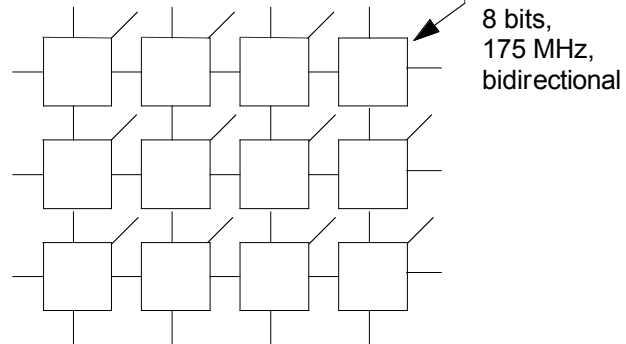  ▪ Allows richer design space
    ○ grids vs hypercubes

McGill

# Example Intel Paragon
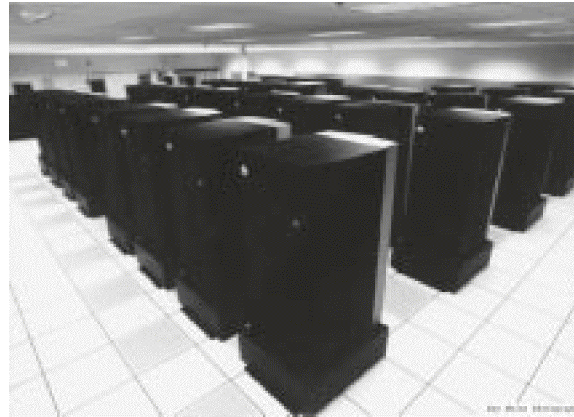


Sandia's Intel Paragon XP/S-based Supercomputer

| i860 | i860 | Intel Paragon node |
| L$_1$ $ | L$_1$ $ | |

Memory bus (64-bit, 50 MHz)

Mem ctrl

DMA

Driver

NI

4-way interleaved DRAM

8 bits, 175 MHz, bidirectional

2D grid network with processing node attached to every switch

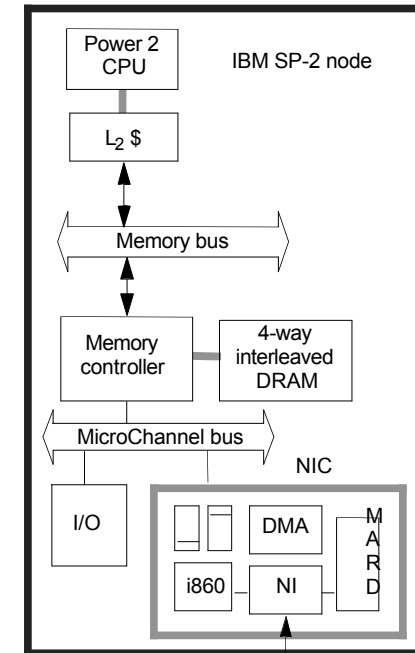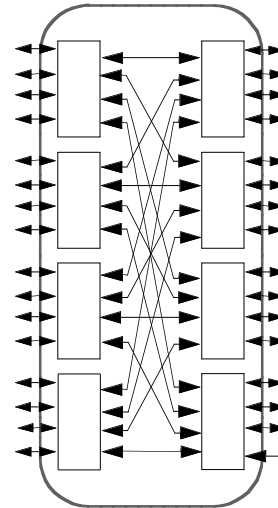ECSE 420/520
Parallel Computing

McGill

# Building on the mainstream: IBM SP-2

- Made out of essentially complete RS6000 workstations

- Network interface integrated in I/O bus (bw limited by I/O bus)

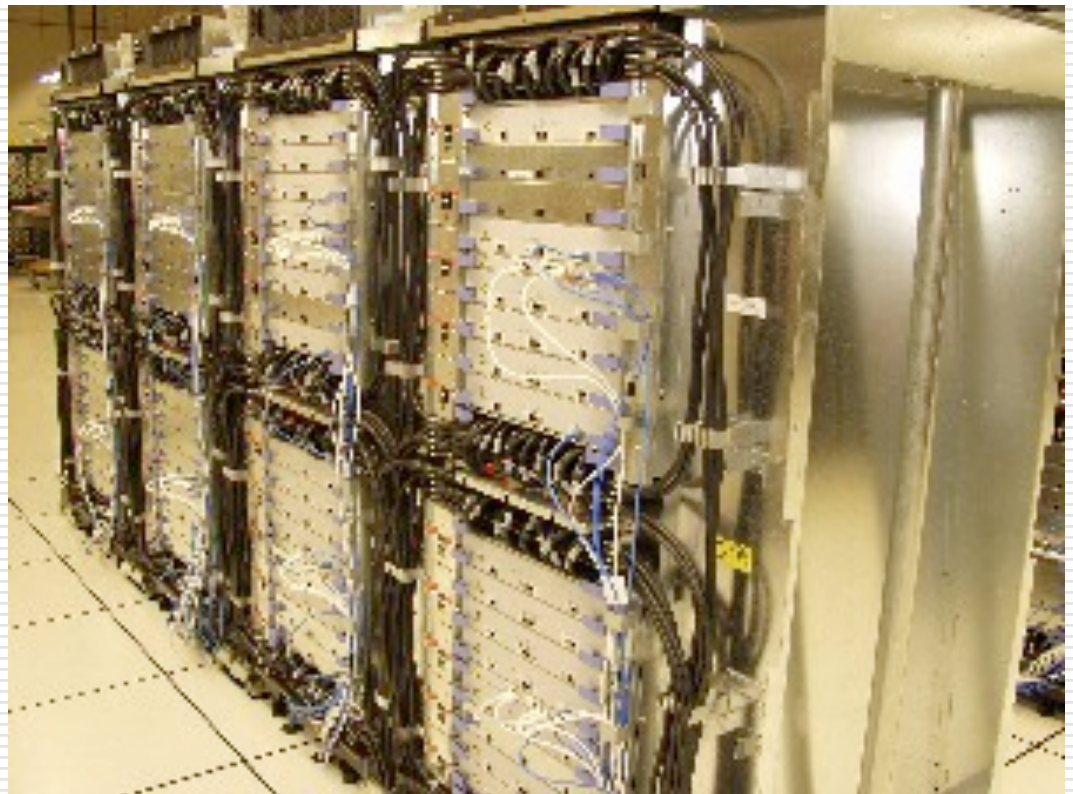General inter connection network formed fr om 8-port switches

Power 2 CPU

IBM SP-2 node

$L_2$ $

Memory bus

Memory controller

4-way interleaved DRAM

MicroChannel bus

NIC

I/O

DMA

i860

NI

M A R D

McGill

# Berkeley NOW



- 100 Sun Ultra2 workstations
- Inteligent network interface
  - proc + mem
- Myrinet Network
  - 160 MB/s per link
  - 300 ns per hop

McGill

# IBM Blue Gene /L

○ Currently, occupies few top spots in top500

○ Lots of embedded processors - PowerPC

McGill

# Toward Architectural Convergence

- Evolution and role of software have blurred boundary
  - Send/recv supported on SAS machines via buffers
  - Can construct global address space on MP    (GA -> P | LA)
  - Page-based (or finer-grained) shared virtual memory
- Hardware organization converging too
  - Tighter NI integration even for MP (low-latency, high-bandwidth)
  - Hardware SAS passes messages
- Even clusters of workstations/SMPs are parallel systems
  - Emergence of fast system area networks (SAN)
- Programming models distinct, but organizations converging
  - Nodes connected by general network and communication assists
  - Implementations also converging, at least in high-end machines

McGill

# Acknowledgments

○ D. Koester, MITRE

○ NUMAchine group

○ Authors of recommended textboks

McGill