# CSE524 Parallel Computation

Lawrence Snyder

www.cs.washington.edu/CSEp524

3 April 2007

---

# Administrivia

o Grading: 20% HW, 70% Proj, 10% Talk

o Burn your book!

o Consumer software + CMPs

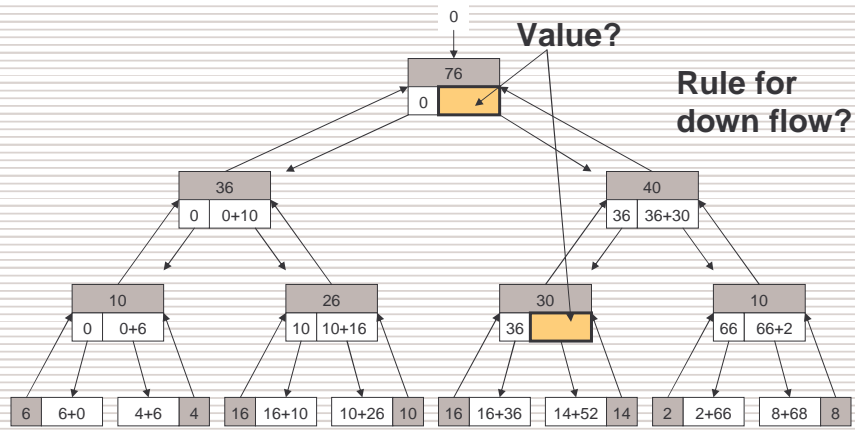   n **MozoDojo,** a graphic mosaic construction tool

> **MultiThreading :** To make the best use possible of multi-core and multi-processor machines, MozoDojo is heavily multithreaded. All computations are dispatched on available processors. The difference is huge between a single G4 and a CoreDuo processor.

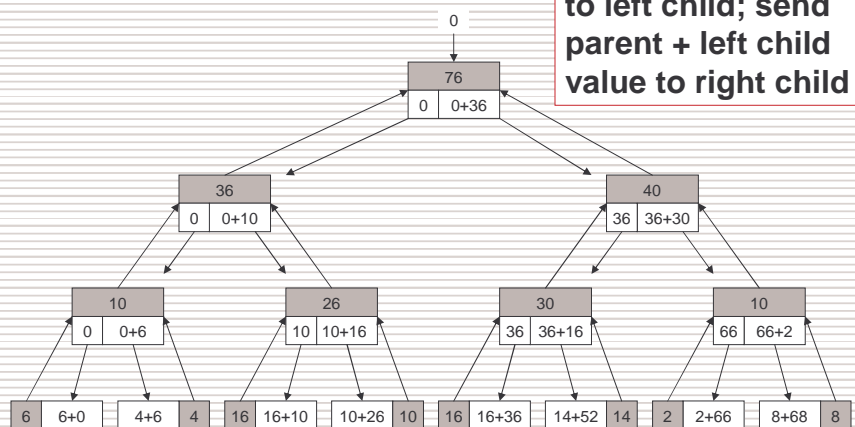o There a HW assigned at end today

Sign up for the email list

# Review



**Value?**

**Rule for down flow?**

# Review Solution

**Send parent value to left child; send parent + left child value to right child**

## Plan for today

o Parallel Hardware
  - n Shared: Multicore, SMP
  - n Distributed: Cluster, HPC Fire breather
o Models of Computation
  - n RAM
  - n PRAM
  - n CTA
  - n Reflection
o Communication modes
  - n Shared, Message Passing, One Sided

5

## Flynn's Taxonomy

o Michael Flynn had an early way to classify machines, two forms of which are still used:

**Flynn's Taxonomy**

Single
Multiple
} Instruction Stream, Single
Multiple
} Data Stream

- n SIMD -- single instruction, multiple data
- n MIMD -- multiple instruction, multiple data

**Our interest is exclusively with MIMD**

6

# Digression on SIMD

o Applying one instruction to multiple values…
  n Was important when memory was expensive
  n Powerful for tight, "inner loop" crunching
  n Performs in rigid lock-step w/apply|not protocol
  n SIMD implements most parallelism **in**efficiently

**Dilemma**: How to get SIMD crunching power with the flexibility needed to control program logic efficiently? **Cell** is perhaps an answer.
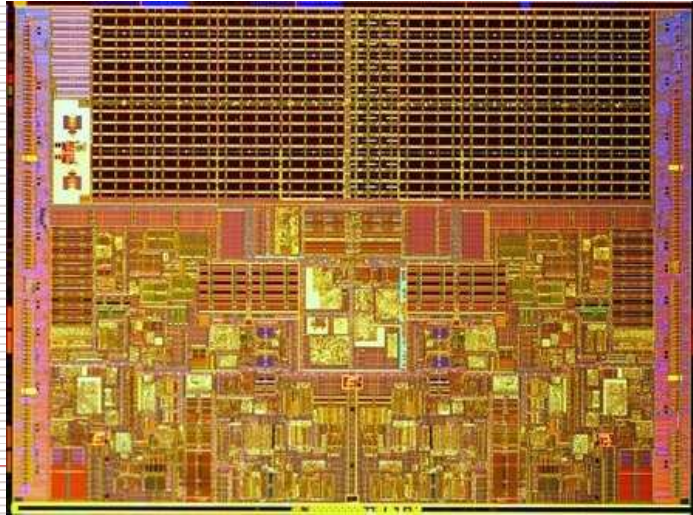
7

# The Problem w/Parallel Architectures

o The problem with parallel machines is
  n They are different from sequential machines
  n They are different from each other
o Both problems complicate programming

o Our solution: Adopt a machine model that abstracts performance critical features

But first, Let's look at some specific machines

8
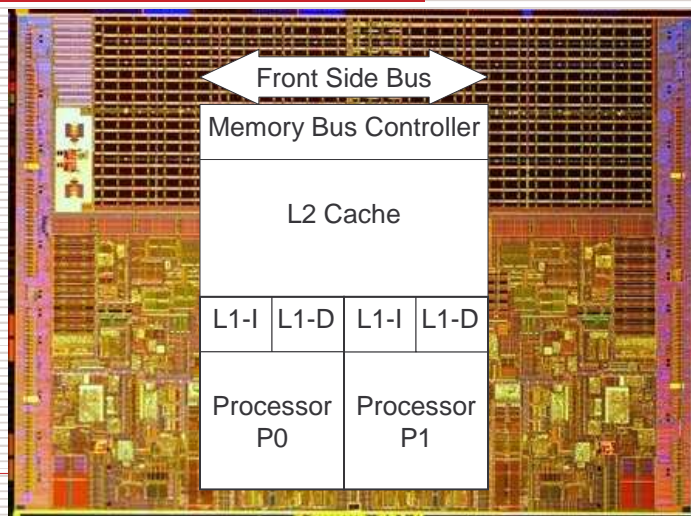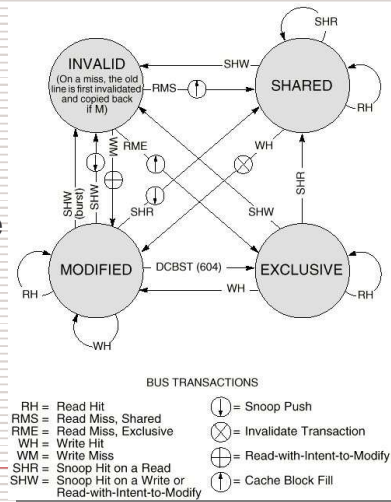
# Intel Core-Duo



9

# Intel Core-Duo



Front Side Bus

Memory Bus Controller

L2 Cache

| L1-I | L1-D | L1-I | L1-D |
|------|------|------|------|

| Processor P0 | Processor P1 |
|--------------|--------------|

10

# MESI Protocol

o Standard Protocol for cache - coherent shared memory
  n Mechanism for multiple caches to give single memory image
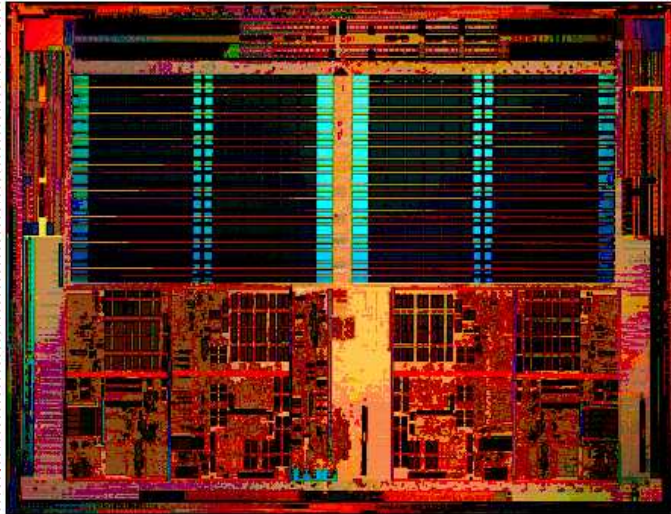  n We will not study it
  n 4 states can be amazingly rich



BUS TRANSACTIONS

RH = Read Hit
RMS = Read Miss, Shared
RME = Read Miss, Exclusive
WH = Write Hit
WM = Write Miss
SHR = Snoop Hit on a Read
SHW = Snoop Hit on a Write or Read-with-Intent-to-Modify

= Snoop Push
= Invalidate Transaction
= Read-with-Intent-to-Modify
= Cache Block Fill

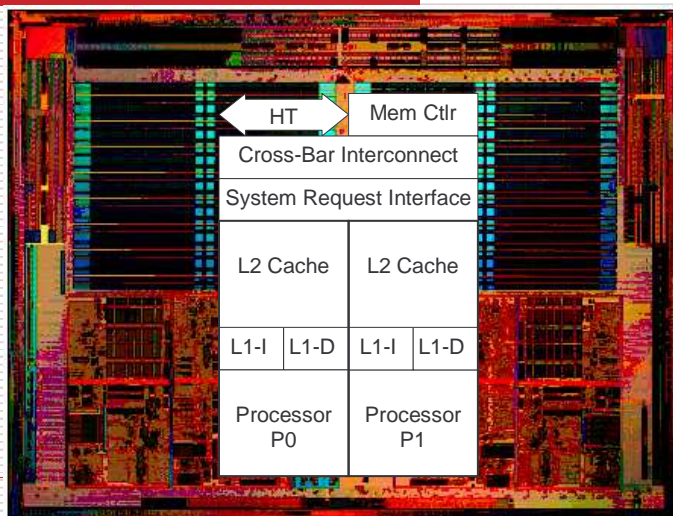Thanks: Slater & Tibrewala of CMU

---

# MESI, Intuitively

o Upon loading, a line is marked E, subsequent reads are OK; write marks M

o Seeing another load, mark as S

o A write to an S, sends I to all, marks as M

o Another's read to an M line, writes it back, marks it S

o Read/write to an I misses

o Related scheme: MOESI (used by AMD)
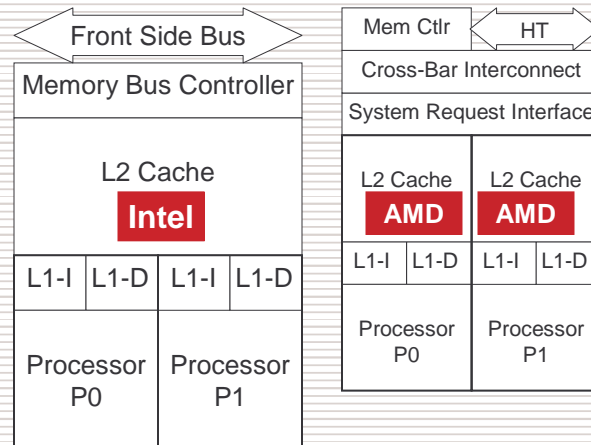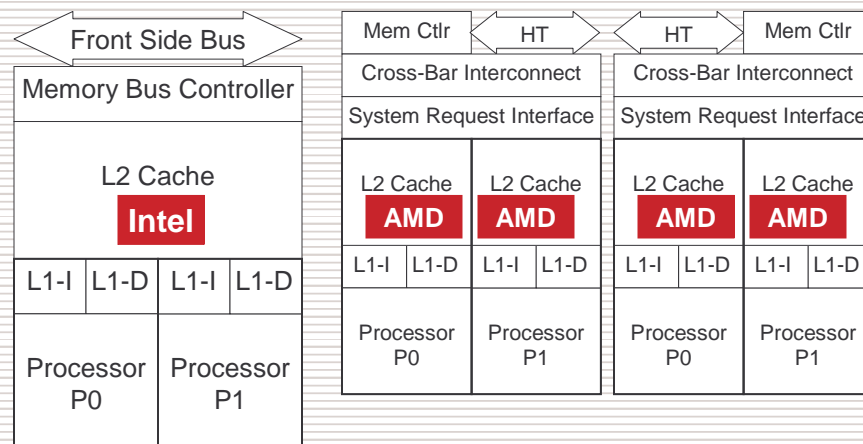
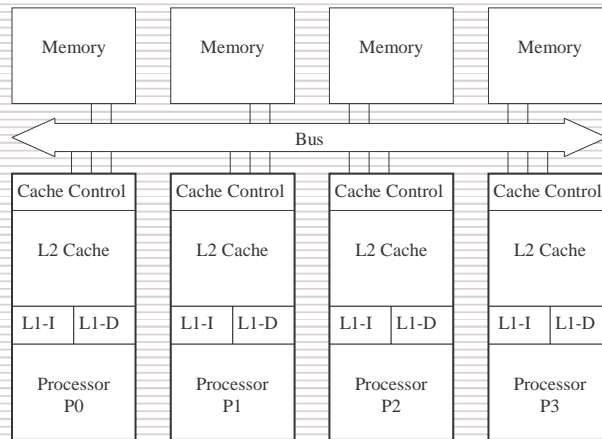# AMD Dual Core Operton



13

# AMD Dual Core Operton



| HT | Mem Ctlr |
|---|---|
| Cross-Bar Interconnect | |
| System Request Interface | |
| L2 Cache | L2 Cache |
| L1-I | L1-D | L1-I | L1-D |
| Processor P0 | Processor P1 |

14

7

# Comparing Core Duo/Dual Core

| Front Side Bus | | Mem Ctlr | HT |
|---|---|---|---|
| Memory Bus Controller | | Cross-Bar Interconnect | |
| | | System Request Interface | |
| L2 Cache **Intel** | | L2 Cache **AMD** | L2 Cache **AMD** |
| L1-I | L1-D | L1-I | L1-D |
| Processor P0 | Processor P1 | Processor P0 | Processor P1 |

15

# Comparing Core Duo/Dual Core

| Front Side Bus | | Mem Ctlr | HT | HT | Mem Ctlr |
|---|---|---|---|---|---|
| Memory Bus Controller | | Cross-Bar Interconnect | | Cross-Bar Interconnect | |
| | | System Request Interface | | System Request Interface | |
| L2 Cache **Intel** | | L2 Cache **AMD** | L2 Cache **AMD** | L2 Cache **AMD** | L2 Cache **AMD** |
| L1-I | L1-D | L1-I | L1-D | L1-I | L1-D | L1-I | L1-D |
| Processor P0 | Processor P1 | Processor P0 | Processor P1 | Processor P0 | Processor P1 |

16

# SMP on a Bus

| Memory | Memory | Memory | Memory |
|--------|--------|--------|--------|

Bus

| Cache Control | Cache Control | Cache Control | Cache Control |
|---------------|---------------|---------------|---------------|
| L2 Cache | L2 Cache | L2 Cache | L2 Cache |
| L1-I \| L1-D | L1-I \| L1-D | L1-I \| L1-D | L1-I \| L1-D |
| Processor P0 | Processor P1 | Processor P2 | Processor P3 |

17

# SMP on a Bus

o The bus is a point that serializes references
o A serializing point is a shared mem enabler

| Memory | Memory | Memory | Memory |
|--------|--------|--------|--------|

Bus

| Cache Control | Cache Control | Cache Control | Cache Control |
|---------------|---------------|---------------|---------------|
| L2 Cache | L2 Cache | L2 Cache | L2 Cache |
| L1-I \| L1-D | L1-I \| L1-D | L1-I \| L1-D | L1-I \| L1-D |
| Processor P0 | Processor P1 | Processor P2 | Processor P3 |

18

# Sun Fire E25K



Diagram labels:
- 18 Sun Fire E25K system expander boards
- Diagram shown: Sun Fire E25K system
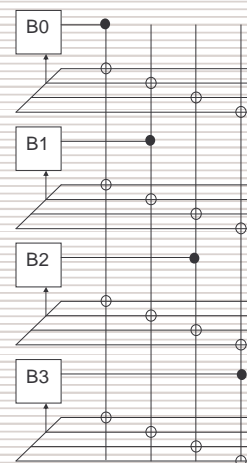- 18 x 18 address and response crossbars
- M, C, PCI, I components

# Cross-Bar Switch

o A crossbar is a network connecting each processor to every other processor

o Used in CMU's 1971 C.MMP, 16 proc PDP-11s

o Crossbars grow as $n^2$ making them impractical for large $n$



B0, B1, B2, B3

# Sun Fire E25K

o X-bar gives low latency for snoops allowing for shared memory

o 18 x 18 X-bar is basically the limit

o Raising the number of processors per node will, on average, increase congestion
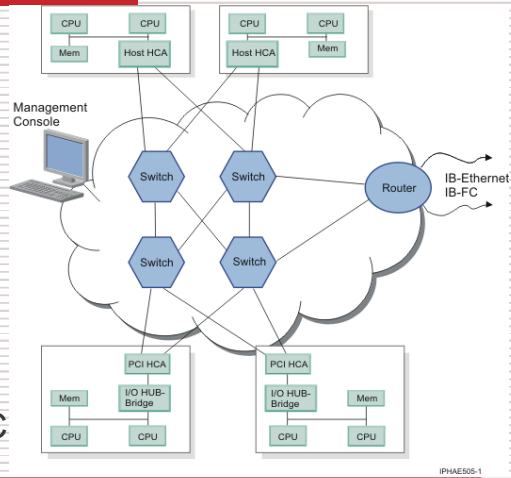
o How could we make a larger machine?

# Co-Processor Architectures

o A powerful parallel design is to add 1 or more subordinate processors to std design

   n Floating point instructions once implemented this way

   n Graphics Processing Units - deep pipelining

   n Cell Processor - multiple SIMD units

   n Attached FPGA chip(s) - compile to a circuit

o These architectures will be discussed later

# Clusters

o **Interconnecting with InfiniBand**

o **Switch-based technology**

  n Host channel adapters (HCA)

  n Peripheral computer interconnect (PC



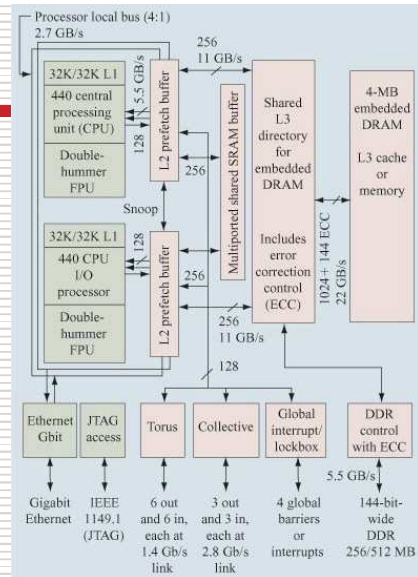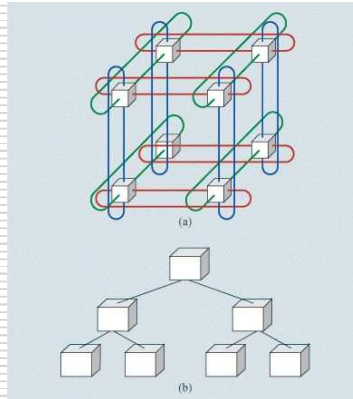Thanks: IBM's Clustering sytems using InfiniBand Hardware   23

---

# Clusters

o Cheap to build using commodity technologies

o Effective when interconnect is "switched"

o Easy to extend, usually in increments of 1

o Processors often have disks "nearby"

o No shared memory

o Latencies are usually large

o Programming uses message passing

24

# Supercomputer

o **BlueGene/L**

# BlueGene/L Specs

o A 64x32x32 torus = 65K 2-core processors

o Cut-through routing gives a worst-case latency of 6.4 μs

o Processor nodes are dual PPC-440 with "double hummer" FPUs

o Collective network performs global reduce for the "usual" functions

o #1 on November's Top 500 at 280 TF

# Summarizing Architectures

o Two main classes
  n Complete connection: CMPs, SMPs, X-bar
    o Preserve single memory image
    o Complete connection limits scaling to …
    o Available to everyone
  n Sparse connection: Clusters, Supercomputers, Networked computers used for parallelism (Grid)
    o Separate memory images
    o Can grow "arbitrarily" large
    o Available to everyone with air conditioning
o Differences are significant; world views diverge

# The Parallel Programming Problem

o Some computations can be platform specific

o Most should be platform independent

o Parallel Software Development Problem: How do we neutralize the machine differences given that
  n Some knowledge of execution behavior is needed to write programs that perform
  n Programs must port across platforms effortlessly, meaning, by at most recompilation

## Options for Solving the PPP

o Leave the problem to the compiler …

## Options for Solving the PPP

o Leave the problem to the compiler …
  - n Very low level parallelism (ILP) is already being exploited
  - n Sequential languages cause us to introduce unintentional sequentiality
  - n Parallel solutions often require a paradigm shift
  - n Compiler writers' track record over past 3 decades not promising … recall HPF
  - n Bottom Line: Compilers will get more helpful, but they probably won't solve the PPP

# Options for Solving the PPP

o Adopt a very abstract language that can target to any platform …

31

# Options for Solving the PPP

o Adopt a very abstract language that can target to any platform …
  - n No one wants to learn a new language, no matter how cool
  - n How does a programmer know how efficient or effective his/her code is? Interpreted code?
  - n What are the "right" abstractions and statement forms for such a language?
    - o Emphasize programmer convenience?
    - o Emphasize compiler translation effectiveness?

32

## Options for Solving the PPP

o Agree on a set of parallel primitives (spawn process, lock location, etc.) and create libraries that work w/ sequential code …

33

## Options for Solving the PPP

o Agree on a set of parallel primitives (spawn process, lock location, etc.) and create libraries that work w/ sequential code …
  n Libraries are a mature technology
  n To work with multiple languages, limit base language assumptions … L.C.D. facilities
  n Libraries use a stylized interface (fcn call) limiting parallelism-specific abstractions possible
  n Achieving consistent semantics is difficult

34

## Options for Solving the PPP

o Create an abstract machine model that accurately describes common capabilities and let the language facilities catch up …

## Options for Solving the PPP

o Create an abstract machine model that accurately describes common capabilities and let the language facilities catch up …
  - n Not a full solution until languages are available
  - n The solution works in sequential world (RAM)
  - n Requires discovering (and predicting) what the common capabilities are
  - n Solution needs to be (continually) validated against actual experience

## Options for Solving the PPP

- Leave the problem to the compiler …
- Adopt a very abstract language that can target to any platform …
- Agree on a set of parallel primitives (spawn process, lock location, etc.) and create libraries that work w/ sequential code …
- Create an abstract machine model that accurately describes common capabilities and let the language facilities catch up …

37

## Break

38

# Reason by Analogy: RAM Model

o The Random Access Machine is our friend
  - n Control, ALU, (Unlimited) Memory, [Input, Output]
  - n Fetch/execute cycle runs 1 inst. pointed at by PC
  - n Memory references are "unit time" independent of location
    - o Gives RAM it's name in preference to von Neumann
    - o "Unit time" is not literally true, but caches fake it
  - n Executes "3-address" instructions

**It's so intuitive, it seems like there's no other way to compute!**

39

# How To Use the RAM

o When reasoning about performance …
  - n Worry about how many instructions executed because execution time proportional to cycles
  - n Treat memory references (operand fetch) as a negligible part of the instruction execution
  - n Estimate time and space needs based on increasing problem size, O($n$)
    - o Linear search vs Binary search
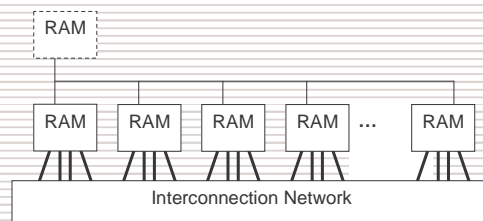  - n Crucial to effectively using C, etc.

40

# Generalization of RAM: PRAM

o Parallel Random Access Machine (PRAM)
  n Unlimited number of processors
  n Processors are standard RAM machines, executing synchronously
  n Memory reference is "unit time"
  n Outcome of collisions at memory specified
    o EREW, CREW, CRCW …

o Model fails bc synchronous execution w/ unit cost memory reference does not scale

41

# CTA Model

o Candidate Type Architecture: A model with $P$ standard processors, $d$ degree, $\lambda$ latency



o Node == processor + memory + NIC

Key Property: Local memory ref is 1, global memory is $\lambda$

42

# What CTA Doesn't Describe

o CTA has no global memory … but memory could be globally *addressed*

o Mechanism for referencing memory not specified: shared, message passing, 1-side

o Interconnection network not specified

o $\lambda$ is not specified beyond $\lambda \gg 1$ -- cannot be because every machine is different

o Controller, combining network "optional"

43

# Typical Values for $\lambda$

o Lambda can be estimated for any machine (given numbers include no contention or congestion)

| | | |
|---|---|---|
| CMP | AMD | 100 |
| SMP | Sun Fire E25K | 400-660 |
| Cluster | Itanium + Myrinet | 4100-5100 |
| Super | BlueGene/L | 5000 |

As with merchandizing: It's location, location, location!

44

# PRAM Mispredicts Preferred Alg

o Consider finding maximum of $n$ numbers
o Best algorithm
- n CRCW PRAM: Valiant's algorithm O(log log $n$)
- n CTA Model: Tournament algorithm O(log $n$)
o Observed performance real implementation
- n PRAM: O(log $n$ log log $n$)
- n CTA: O(log $n$)
o We do not want a model that directs us to an impractical solution

45

# Apply CTA to Count 3s

o How does CTA guide us for Count 3s pgm
- n Array segments will be allocated to local mem
- n Each processor should count 3s in its segment
- n Global total should be formed using reduction
- n Performance is
  - o Full parallelism for local processing
  - o λ log n for combining (and broadcast)
  - o Base of log should be large, i.e high degree nodes
o Same solution as before, but by different rt

46

# Communication Mechanisms

o Shared addressing

- n One consistent memory image; primitives are load and store
- n Must protect locations from races
- n Widely considered most convenient, though it is often tough to get a program to perform
- n CTA implies that best practice is to keep as much of the problem private; use sharing only to communicate

A common pitfall: Logic is too fine grain

47

# Communication Mechanisms

o Message Passing

- n No global memory image; primitives are send() and recv()
- n Required for most large machines
- n User writes in sequential language with message passing library:
  - o Message Passing Interface (MPI)
  - o Parallel Virtual Machine (PVM)
- n CTA implies that best practice is to build and use own abstractions

Lack of abstractions makes message passing brutal

48

## Communication Mechanisms

o One Sided Communication

  n One global address space; primitives are get() and put()

  n Consistency is the programmer's responsibility

  n Elevating mem copy to a comm mechanism

  n Programmer writes in sequential language with library calls -- not widely available unfortunately

  n CTA implies that best practice is to build and use own abstractions

  > One-sided is lighter weight than message passing

49

## Assignment for next week

o Read Chapter 3

o Homework Problem: Analyze the complexity of the Odd/Even Interchange Sort: Given array A[n], exchange o/e pairs if not ordered, then exchange e/o pairs if not ordered, then repeat until sorted

o Analyze in CTA model (i.e. for $P$, $\lambda$, $d$), and charge the o/e-e/o pair $c$ time if operands are local; ignore all other local computation

50