

Models In Parallel Computation

It is difficult to write programs without a good idea of how the target computer will execute the code. The most important information is knowing how expensive the operations are in terms of time, space, and communication costs

First ... The Quick Sort Essay

- Did Quick Sort seem like a good parallel algorithm initially?
- Is it clear why it might not be?
- Thoughts?

Last Week

- Matrix Multiplication was used to illustrate different parallel solutions
 - **Maximum parallelism**, $O(\log n)$ time, $O(n^3)$ processors, PRAM model
 - **Basic** (strips x panels), $O(n)$ time, $O(n^2)$ processors
 - **Pipelined** (systolic), $O(n)$ time, $O(n^2)$ processors, VLSI model
 - **SUMMA algorithm**, used many techniques, $O(n)$ time, $O(n^2)$ processors, scalable, Distributed Memory Model

Last Week (continued)

- Different techniques illustrated --
 - Decompose into independent tasks
 - Pipelining
 - Overlapping computation and communication
- Optimizations
 - Enlarge task size, e.g. several rows/columns at once
 - Improve caching by blocking
 - Reorder computation to “use data once”
 - Exploit broadcast communication

The SUMMA algorithm used all of these ideas

Goal For Today ...

Understand how to think of a parallel computer independently of any hardware, but specifically enough to program effectively

Equivalently, *Find a Parallel Machine Model*

- It's tricky because unlike sequential computers parallel architectures are very different from each other
 - Being too close to a physical HW (low level) means embedding features that may not be on all platforms
 - Being too far from physical HW (high level) means writing code taking too many software layers of build, and so, slow

Plan for Today

- Importance of von Neumann model & C programming language
- Recall PRAM model
 - Valiant's Maximum Algorithm
 - Analyze result to evaluate model
- Introduce CTA Machine Model
 - Analyze result to evaluate model
- Alternative Models
 - LogP is too specific
 - Functional is too vague

Successful Programming

When we write programs in C they are ...

- **Efficient** -- programs run fast, especially if we use performance as a goal
 - traverse arrays in row major order to improve caching
- **Economical** -- use resources well
 - represent data by packing memory
- **Portable** -- run well on any computer with C compiler
 - all computers are universal, but with C fast programs are fast everywhere
- **Easy to write** -- we know many ‘good’ techniques
 - reference data, don’t copy

These qualities all derive from von Neumann model

Von Neumann (RAM) Model

- Call the ‘standard’ model of a random access machine (RAM) the von Neumann model
 - A processor interpreting 3-address instructions
 - PC pointing to the next instruction of program in memory
 - “Flat,” randomly accessed memory requires 1 time unit
 - Memory is composed of fixed-size addressable units
 - One instruction executes at a time, and is completed before the next instruction executes
- The model is not literally true, e.g., memory is hierarchical but made to “look flat”

C directly implements this model in a HLL

Why Use Model That's No Literally True?

- Simple is better, and many things--GPRs, floating point format--don't matter at all
- Avoid embedding assumptions where things could change ...
 - Flat memory, tho originally true, is no longer right, but we don't retrofit the model; we don't want people “programming to the cache”
 - Yes, exploit spatial locality
 - No, avoid blocking to fit in cache line, or tricking cache into prefetch, etc.
 - Compilers bind late, particularize and are better than you are!

vN Model Contributes To Success

- The cost of C statements on the vN machine is “understood” by C programmers ...
- How much time does **$A[r][s] += B[r][s];$** require?
 - Load row_size_A, row_size_B, r, s, A_base, B_base (6)
 - tempa = (row_size_A * r + s) * data_size (3)
 - tempb = (row_size_B * r + s) * data_size (3)
 - A_base + tempa; B_base + tempb; load both values (4)
 - Add values and return to memory (2)
 - Same for many operations, any data size
- Result is measured in “instructions” not time

Widely known and effectively used

Portability

- Most important property of the C-vN coupling:
It is approximately right everywhere
- Why so little variation in sequential computers?

**HW vendors must
run installed SW
so follow vN rules**

**SW vendors must
run on installed HW
so follow vN rules**

**Everyone wins ... no
motive to change**

Von Neumann Summary

- The von Neumann model “explains” the costs of C because C expresses the facilities of the von Neumann machines in a set of useful programming facilities
- Knowing the relationship between C and the von Neumann machine is essential for writing efficient programs
- Following the rules produces good results everywhere because everyone benefits
- These ideas are “in our bones” ... it’s how we think

What is the parallel version of vN?

PRAM Often Proposed As A Candidate

PRAM (Parallel RAM) ignores memory organization, collisions, latency, conflicts, etc.

Ignoring these are *claimed* to have benefits ...

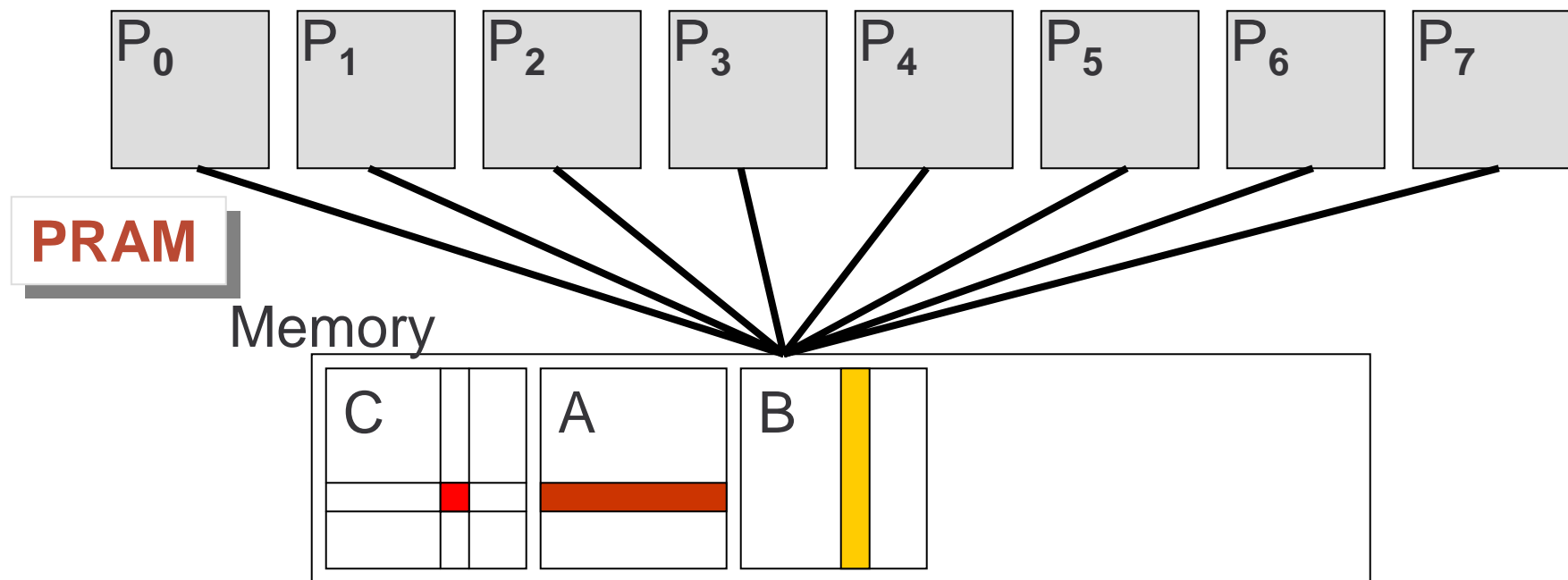
- Portable everywhere since it is very general
- It is a simple programming model ignoring only insignificant details -- off by “only log P”
- Ignoring memory difficulties is OK because hardware can “fake” a shared memory
- Good for getting started: Begin with PRAM then refine the program to a practical solution if needed

We will make these more precise next week

Recall Parallel Random-Access Machine

PRAM has any number of processors

- Every processor references any memory in “time 1”
- Memory read and write collisions must be resolved



SMPs implement PRAMs for small P ... not scalable

Variations on PRAM

Resolving the memory conflicts considers read and write conflicts separately

- Exclusive read/exclusive write (EREW)
 - The most limited model
- Concurrent read/exclusive write (CREW)
 - Multiple readers are OK
- Concurrent read/concurrent write (CRCW)
 - Various write-conflict resolutions used
- There are at least dozen other variations

All theoretical -- not used in practice

Find Maximum with PRAM (Valiant)

Task: Find largest of n integers w/ n processors

Model: CRCW PRAM (writes OK if same value)

How would YOU do it?

L.G.Valiant, “Parallelism in comparison problems,”
SIAM J. Computing 4(3):348-355, 1975

L.G. Valiant, “A Bridging Model for Parallel
Computation,” CACM 33(8):103-111, 1990

R.J. Anderson & L. Snyder, “A Comparison of
Shared and Nonshared Memory Models for
Parallel Computation,” *Proc. IEEE* 79(4):480-487

Algorithm Sketch

Algorithm: T rounds of $O(1)$ time each

In round, process groups of m vals, v_1, v_2, \dots, v_m

- Fill m memory locations x_1, x_2, \dots, x_m with 1
- For each $1 \leq i, j \leq m$ a processor tests ...
`if $v_i < v_j$ then $x_i = 0$ else $x_j = 0$`
- If $x_k = 1$ it's max of group; pass v_k to next round

The 'trick' is to pick m right to minimize T

Finding Max (continued)

Round 1: $m = 3$

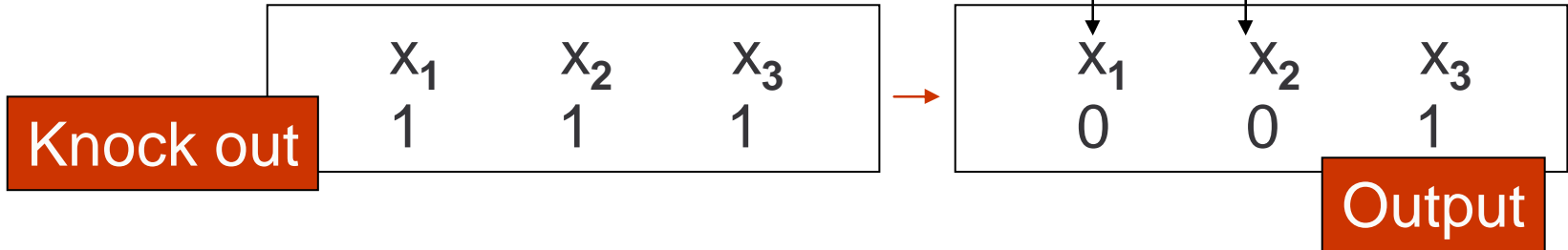
Input

V_1	V_2	V_3
20	3	34

Schedule

	V_1	V_2	V_3
V_1	-	$V_1 : V_2$	$V_1 : V_3$
V_2	-	-	$V_2 : V_3$
V_3	-	-	-

For groups of size 3, three tests can find max, i.e. 3 processors



Solving Whole Problem

- Round 1 uses P processors to find the max in groups of $m=3$... producing $P/3$ group maxes
- Round 2 uses P processors to find the max in groups of $m=7$... producing $P/21$ group maxes
- Generally to find the max of a group requires $m(m-1)/2$ comparisons
- Picking m when there are P processors, r maxes ... largest m s.t. $(r/m)(m(m-1)/2) \leq P$ i.e. $r(m-1) \leq 2P$

Finding Max (continued)

- Initially, $r = P$, so
 $r(m-1) \leq 2P$
implies $m = 3$, producing $r = P/3$
- For $(P/3)(m-1) \leq 2P$ implies next group = 7
- Etc.
- Group size increases quadratically implying the maximum is found in $O(\log \log n)$ steps on CRCW PRAM

It's very clever, but is it of any practical use?

Assessing Valiant's Max Algorithm

The PRAM model caused us to ...

- Exploit the “free use” of read and write collisions, which are not possible in practice
- Ignore the costs of data motion, so we adopt an algorithm that runs *faster* than the time required to bring all data values together, which is $\Omega(\log n)$
- So what?

Running Valiant's Algorithm

- PRAM's don't exist and can't be built
- To run the algorithm we need a simulator for the CRCWPRAM
- In order to simulate the concurrent reads and the concurrent writes, a parallel computer will need $\Omega(\log P)$ time per step, though there are bandwidth requirements and serious engineering problems to attain that goal [details in future lecture]
- *Observed* performance of Valiant's Max:
 $O(\log n \log \log n)$





Alternative Solution

- What is the best (practical) way of computing max?
 - A tree algorithm a variation on global sum
 - $O(\log P)$ time on P processors (as with sum, the problem size can be $P \log P$)
 - The tree algorithm doesn't need to be simulated ... it runs in the stated time directly on all existing parallel processors
- Since $O(\log n) < O(\log n \log \log n)$ the PRAM model mispredicted the best practical algorithm

The PRAM didn't help, it hurt our effort

Is The PRAM A Good Abstraction?

Different Opinions ...

- OK for finding theoretical limits to parallelism 
- It is a simple programming model ignoring only insignificant details -- off only by $\log P$ 
- Ignoring memory difficulties is OK because hardware can “fake” a shared memory 
- Start with PRAM then evolve to more realistic solution -- good for getting started 

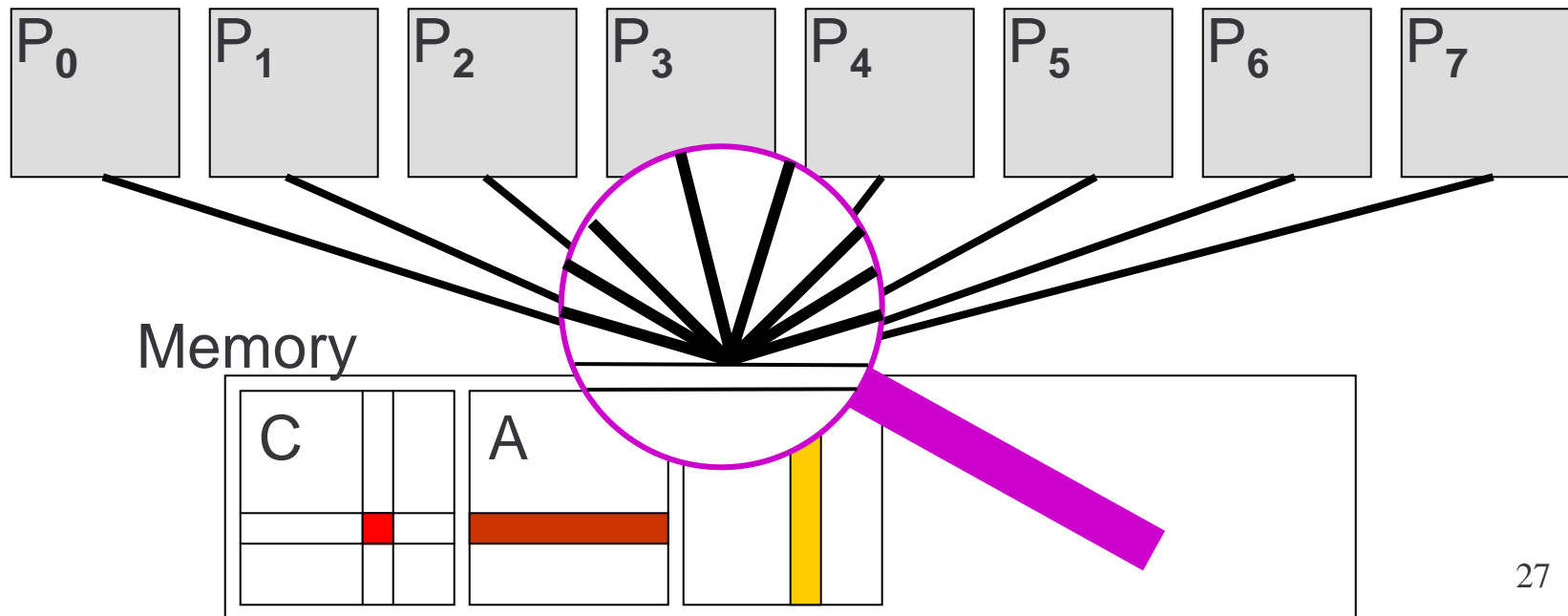
Break

Requirements Of A Practical Model

- Like von Neumann ...
 - Capture key features of HW implementations
 - Reflect costs realistically
 - Portable across physical machines
 - KISS: Ignore everything that “doesn’t count”
- Problems for parallel programmers ...
 - Parallel computers have widely different architectures -- steady-state cost not established
 - Existing SW for Cray vector machines, poor guide
 - Sequential computers keep getting faster

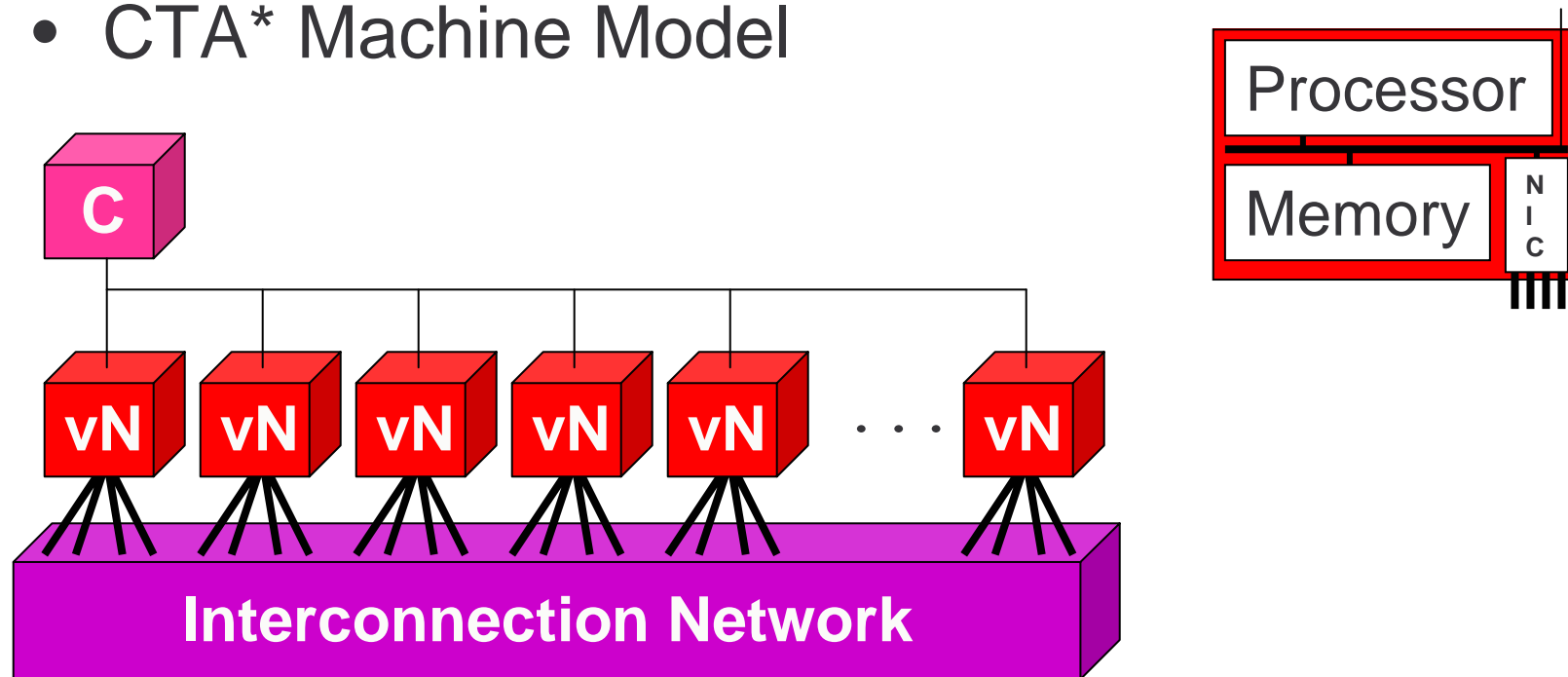
Model Must Be “Realizable”

- There should be no “magic”
- Explain how the processors connect to the memory



Proposed Practical Parallel Model

- CTA* Machine Model



Benefit from everything the vN model does right

CTA = Candidate Type Architecture

CTA Citation

- Lawrence Snyder, “Type Architecture, Shared Memory and the Corollary of Modest Potential,” *Annual Reviews of Computer Science*, 1986

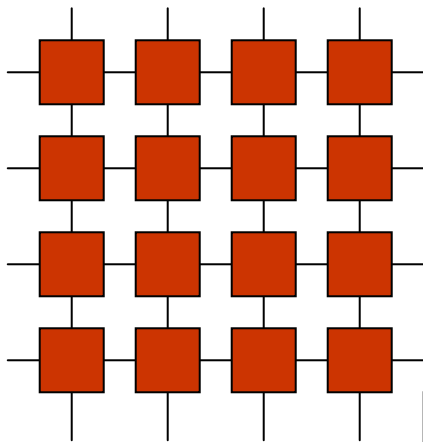
Properties of the CTA ...

- Processor Elements are von Neumann machines -- they could be SMPs or other parallel designs, but they execute vN ISA
- PEs are connected by a sparse ($O(n)$ edges) point-to-point network, not specified
- PE's connection to network is “narrow”
- Latency $\lambda \gg 1$ for non-local memory reference
- Controller has “thin” broadcast connection to all processors

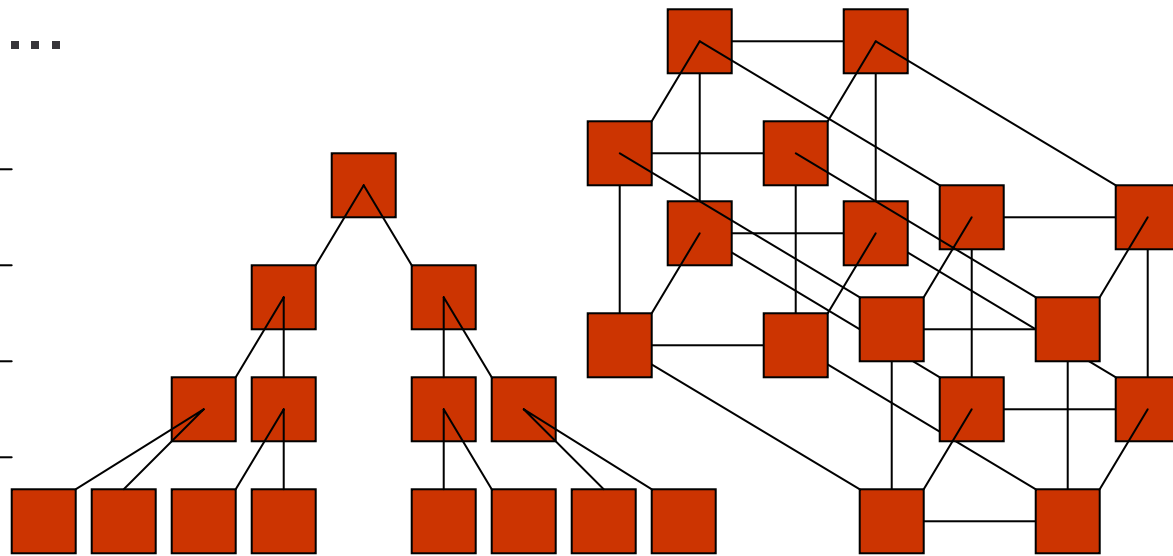
The distributed memory model of last week

Network Not Specified

- Historically, programmers have “programmed to the network” harming their computations
- Examples ...



OK



Bottlenecks

Not Sparse

Network is like the system bus -- forget it exists

Using CTA What Is Best Max Algorithm?

CTA charges for operations, memory reference, communication ... PRAM solution of little use

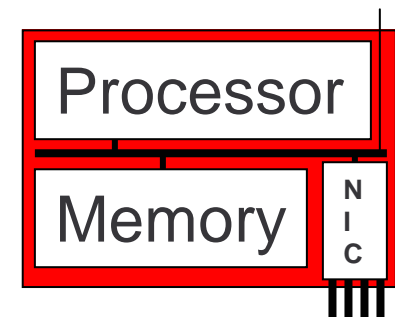
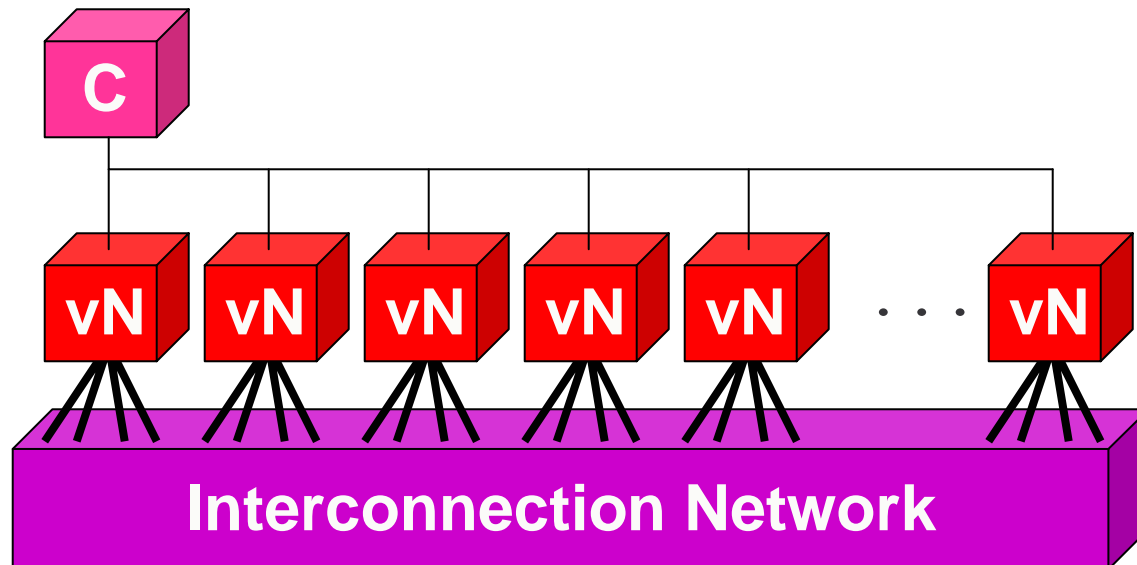
- Tournament algorithm finds maximum of n elements in $O(\log n)$ time like the global sum
 - Odd PEs send value to next lower processor
 - Even PEs recv, compare and send larger value to “parent”
 - IDivide PE index by 2, and repeat until one processor left
 - Time is $O(\log n)$ communication + parallel compare steps

This is a practical algorithm that can be tested

CTA Emphasizes Locality

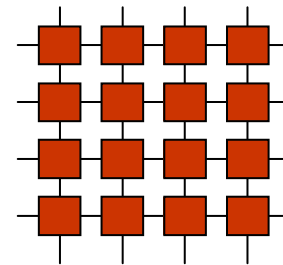
CTA models ... where the data is & where it goes

- The λ parameter captures the benefit of exploiting locally stored values
 - λ says moving data to other PEs has significant cost



The λ value

- λ is a purposely vague quantity
 - The cost is not known and cannot be known because every machine is different
 - There is no agreement yet on *the* parallel architecture
 - λ depends on bandwidth, engineering, interfaces, mem
 - The cost cannot be a single number because of network topology, conflicts, etc.
 - Optimizing for communication -- if possible -- should be done by the compiler for the machine not the programmer for eternity



CTA Abstracts Existing Architectures

- The MIMD parallel machines -- all machines since mid 1980s -- are modeled by CTA
- iPSC, HyperCube, UltraComputer, Cedar, SP-1, iPSC2, KSR-1, CM-5, Sequent, Touchstone, T3D, Dash, J-Machine, SP-2, KSR-2, T3E, Tera, -2000, etc., etc., etc.
- A few minor adjustments may be needed
 - Some machines do not have a controller (front end machine), but PE_0 can be both a PE and Controller

Clusters and SMPs

- (Beowulf) Clusters are popular because they are easy to build, simple to operate and most of the software is Open Source
 - Clusters built using a point-to-point network (Myrinet) are CTA machines
 - Clusters with bus-architectures fail to meet the point-to-point network feature, but a decent bus is an OK approximation (and there is no other model!)
- SMPs are technically not CTA machines, but since they are “better” ($\lambda \approx 1$), using CTA guidelines is OK

Programming

- Since the CTA contains the von Neumann architecture C is good for programming PEs but more global abstractions are needed
- The machine model that message passing programmers use matches the CTA, so message passing solutions --but not the actual code-- will work
 - Message passing libraries are “portable” in only a limited way
 - Lack of abstractions limits key compiler help and key optimizations

ZPL is built on the CTA and exploits its benefits

Assessing the CTA

The considerations were

- Capture key features of HW implementations
- Reflect costs realistically
- Portable across physical machines
 - The above have been established experimentally
- Be as simple as possible relative to the above
 - Difficult to measure, but programmers can use it

We'll use the CTA in this class when programming

Improving CTA may be possible, but not in some ways

Alternatives ...

Many models have been proposed, but they have had weaknesses

- LogP [UC Berkeley] is CTA+parameters
- The model measures latency and other variables related to communication -- too many free variables
 - The model was effective at predicting behavior on the Thinking Machines CM-5, which motivated it
 - With many parameters, it was difficult to decide which were significant, especially when their values are free
 - It remains a model of theoretical interest

LogP's params try being exact re: λ . Not possible

Alternatives (continued)

- Functional programming has motivated many models
 - *Concurrency is natural part of parameter evaluation*
 - Processors not specifically identified -- they simply work on next task
 - Write-once memory avoids races and indeterminacy
 - No notion of locality since it is not possible to associate memory with any logical processor
 - Functional models presume shared memory computers and have considerable overhead; functional has not become popular

Summary

- The von Neumann model allows us to be successful C programmers
- A parallel alternative is needed
 - The PRAM is not a good alternative because it leads us to impractical algorithms like the $O(\log \log n)$ max
 - The CTA does work because it charges for data motion and communication using the parameter λ
 - Capturing locality is a central benefit of the CTA
 - CTA abstracts all contemporary MIMD machines

Finding a model is a balancing act in which the “right” things are modeled and everything else is ignored

Assignment

- Read the Type Architecture paper ... it's a classic paper, so will seem quaint in places, but try to see the machine model ideas as they emerged ...