# CSE524 Parallel Computation

Lawrence Snyder
www.cs.washington.edu/CSEp524

24 April 2007

# Announcements

o New homework assigned at end of class
o Discuss this week's assignment in 2 parts
  n Red/Blue, now
  n Essay, later this class
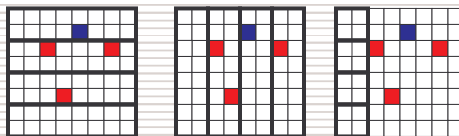o Project description circulated at break

Project thinking begins this week

# Red/Blue Problem

o Write Peril-L code to implement Red/Blue
o Goal is to get a scalable program
o Inputs:
  n $n$ x $n$ board, initialized and allocated
  n $t$, a small region used to sample
  n Termination condition: some $t$ x $t$ on left > 90%
o Rules of the game
  n Red moves 1 right on 1st half step if cell free
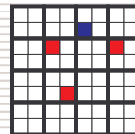  n Blue moves 1 down on 2nd half step if cell free

3

# View Possible Solutions Globally
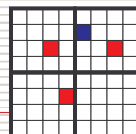
o Shared allows threads to process rows/cols



**All references to global memory**

o Allocate $t$ x $t$ threads



**Must be multiplexed**

o Scalable blocks



**Maximize local work; multiplexing possible**

4

# Red-Blue … One Solution
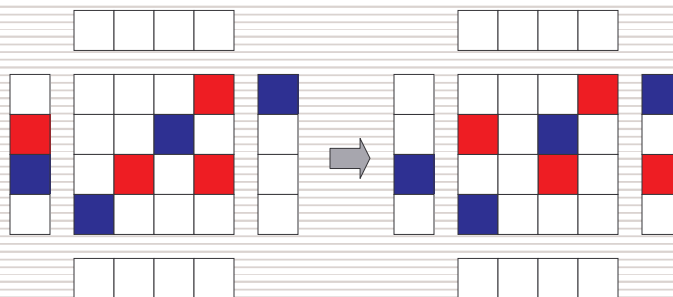
```
int RB_G[n][n];                      t/n, colors 0=w,1=r,2=b
/* Initialize RB with data */        Get external data
int size=sqrt(P); my=n/size;         Procs along row/col
int bin_GO[size][size][my];          MM for passing cols
int bou_GO[size][size][my];          MM for passing rows
int thresh=.9;                       90% threshold
forall thr,thc in (0:size-1,0..size-1) {   2D threads
  int Lrb[my][my]=is_local(RB_G);  Work on local assigned
  int lft[my],rgt[my],top[my],bot[my];Local neighbr storage
  int xl = RB_G_myLo_1; yl = RB_G_myHi_1;  Global index
  int xh = RB_G_myLo_2; yh = RB_G_myHi_2;  Global index
  int done=0,countR, countB;       Termination variables
```

5

# Cartoon of Red Move

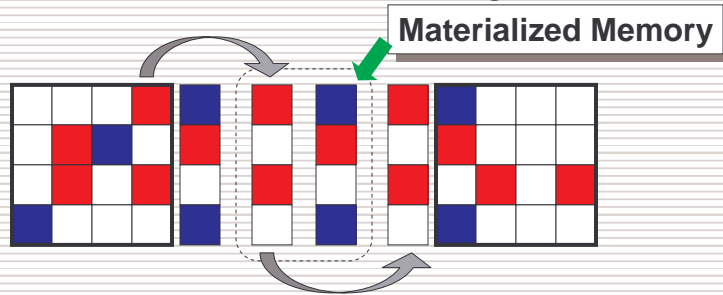o **Get adjacent context from neighbors**



o **Compute local state**

Similarly for blue moves

6

3

# Cartoon of State Transfer

o Save values in materialized global memory



**Materialized Memory**

---

# Solution (Continued)

```
while (done==0) {                              Go till threshold
  bou_GO[thr][thc][0:my-1]=Lrb[xl][yl:yh];  pass left col
  bin_GO[trh][thc][0:my-1]=Lrb[xh][yl:yh];  pass right col
  rgt[0:my-1]=bou_GO[thr][(thc+1)%my][0:my-1];last col+1
  lft[0:my-1]=bin_GO[thr][(thc-1)%my][0:my-1];first col-1
  moveRed(Lrb[][],rgt[],lft[]); barrier;
  bou_GO[thr][thc][0:my-1]=Lrb[xl:xh][yl];  pass top row
  bin_GO[trh][thc][0:my-1]=Lrb[xl:xh][yh];  pass bot row
  top[0:my-1]=bin_GO[(thr-1)%my][thc][0:my-1];  top row-1
  bot[0:my-1]=bou_GO[(thr+1)%my][thc][0:my-1];  last row+1
  moveBlue(Lrb[][],bot[],top[]);
```

## Solution Continued

```
  if (thc == 0)                    Am I on left edge?
    {for (k=0; k<my/t; k++) {        Do all my txt blocks
      countR = 0; countB = 0;        Count up colors
      for (i=0; i<t; i++) {
       for (j=0; j<t; j++) {
         if (Lrb[i][j]== 1)
            countR++;
         if (Lrb[i][j]== 2)
            countB++;
       }
      }
```

9

## Solution Continued

```
/* After a t x t block ...  */
     if (countR/(countR+countB)>=thresh ||
        countB/(countR+countB)>=thresh)
           done = 1;              1 color is over thresh
    }                             End of t x t tiles
   }                              End of if for term test
   done=+/done;                   Find out how others did;synchs
  }                               End of while
}                                 End of thread
```

10

5

## Recap Solution

- o Scaled by available processors; $t$ x $t$ is min
- o Materialized memory implements a bulk exchange, benefiting from fast transfer
- o Constructing context allows local compute
- o Termination performed by reduction
- o Why was the barrier needed?
- o Is less data motion possible?

A good exercise: Analyze R/B complexity for different P

## Schwartz's Algorithm

- o Jack Schwartz (NYU) asked: What is optimal number of processors to combine $n$ values?
  - n Reasonable Answer: binary tree w/ values at leaves has O($log\ n$) complexity
  - n To this solution add $log\ n$ more values to leaves
  - n Same complexity (O(log n)), but $nlog\ n$ values!
- o Generally $P$ is not a variable, and $P << n$
- o Use Schwartz as heuristic: Prefer to work at leaves rather than enlarge (make a deeper) tree, implying tree will have $log\ P$ height
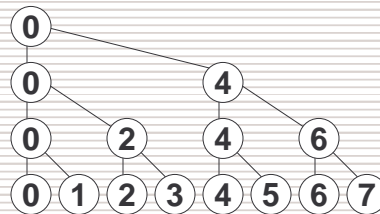
12

# Expressing A Tree

o Consider writing tree-based sum in Peril-L
  - n Processors will have several roles: leaf, intr node
  - n Solve synchronization with materialized memory
o Stages:
  - n Compute result at leaves
  - n Threads have IDs in 0 - $2^p$-1, pass if 0 != ID%2
  - n Even threads add, pass value if 0 != ID%4
  - n Etc.

13

# Peril-L Code for Collective Ops

o Thread logic is easy ... initialize `stride = 1`
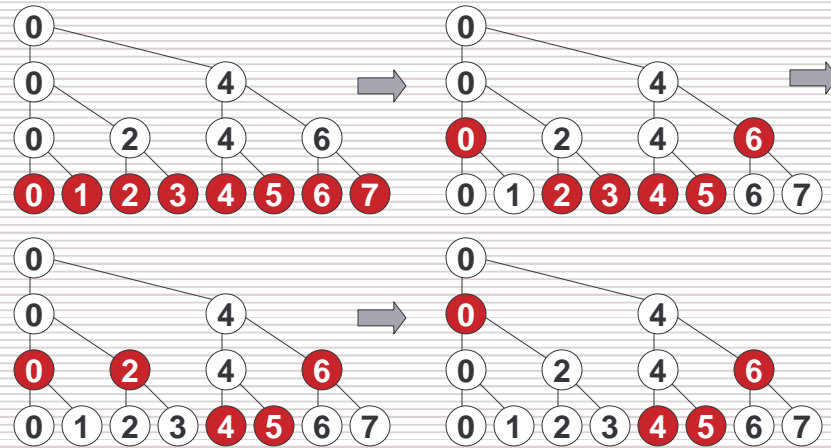
```
nodeval_GO[tID] =  tally;          Send local val to tree node
while (stride < P) {               Begin logic for tree
  if (tID %(2*stride) == 0) {
      nodeval_GO[tID]= nodeval_GO[tID] +
      nodeval_GO[tID + stride]);
      stride = 2*stride;
  }
    else
      break;
}
```

Reverse to broadcast



14

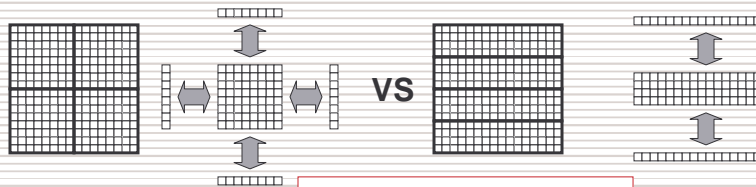# Materialized Memory Self-Synchs

---

# Asynchronous Trees

- o MM gives asynchronous tree, preferred over a "leveled" tree using barriers

- o Notice that once a thread's role is complete it can continue computing

- o For collective op **+** broadcast, a tree implemented with materialized memory allows half of the threads to continue

- o **The Principle**: Relax synchronization rules wherever possible!

# Block Allocations

o The Red/Blue computation illustrated a 2D-block data parallel allocation of the problem

o Generally block allocations are better for data transmission: surface to volume advantage … since only edges are x-mitted
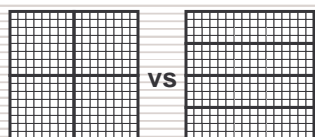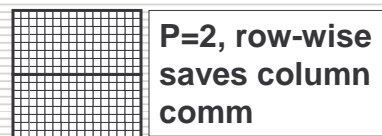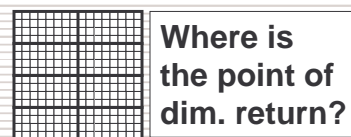
**VS**

Now scale problem 4x

17

# Different Regimens

o Though block is generally a good allocation it's not absolute:

**P=1, all comm wasted**

**P=2, row-wise saves column comm**

vs **P=4, rows and blocks are a wash**

**Where is the point of dim. return?**
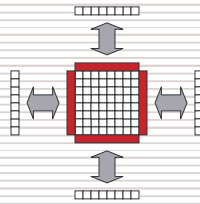
18

## Shadow Buffers/Fluff

o To simplify local computation in cases where nearest neighbors values x-mitted, allocate in-place memory to store values:
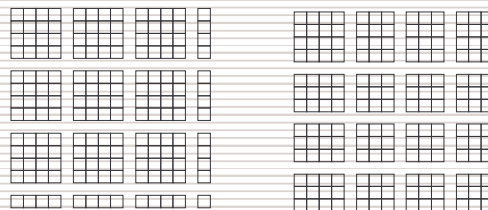
o Array can be referenced as if it's all local

Edge storage (rgt,…) illustrated in R/B

19

## Aspect Ratio

o Generally *P* and *n* do not allow for a perfectly balanced allocation …
o Several ways to assign arrays to processors

**Generally a small effect**

**Quotient + remainder**

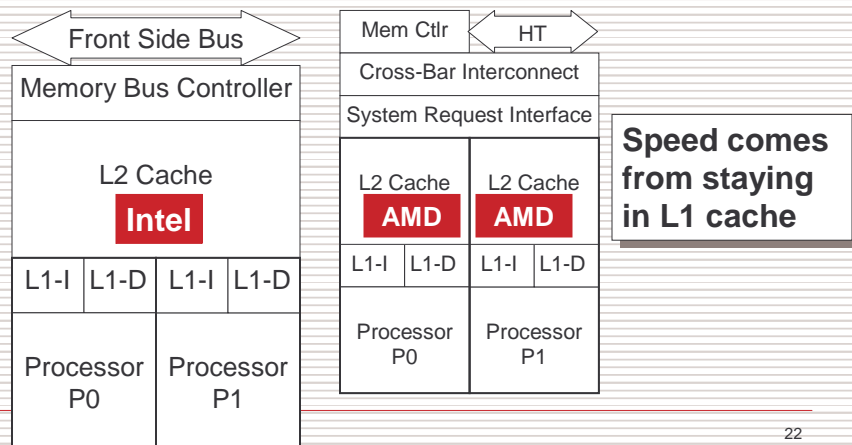**Ceiling + floor**

20

10

# Assigning Processor 0 Work

- $P_0$ is often assigned "other duties", such as
  - Orchestrate I/O
  - Root node for combining trees
  - Work Queue Manager …
- Assigning $P_0$ the smallest quantum of work helps it avoid becoming a bottleneck
  - For either quotient + remainder or ceiling/floor $P_0$ should be the last processor

This is a late-stage tuning matter
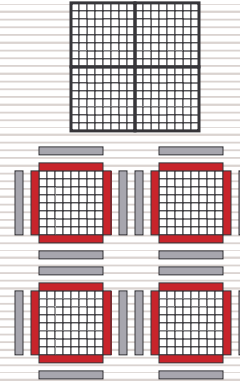
21

---

# Solving Array Problems on CMPs

- How important are these topics for CMPs?

| 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 | L1-I | L1-D | L1-I | L1-D |
| Processor P0 | Processor P1 | Processor P0 | Processor P1 |

**Speed comes from staying in L1 cache**

22

11

## Locality Always Matters

o Array computations on CMPs
  - n Dense Allocation vs Fluff
  - n Issue is cache invalidation
  - n Keeping MM managed intermediate buffers keeps array and fluff local (L1)
  - n Sharing causes elements at edge to repeatedly invalidate killing locality

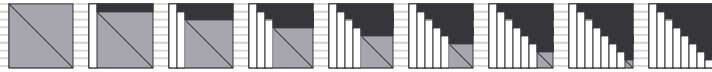  False sharing an issue, too

## Reduce For Global Sum

o Reduce abstraction vs `sum_G += 1`
  - n Reduce raises level of abstraction
  - n Reduce provides private var implementation
o Tree implementation of reduce
  - n Irrelevant for tiny numbers of processors

**Bottom Line: CTA machine model not strictly required for CMPs. Using it yields scalable programs that benefit CMPs in most ways without causing any harm.**
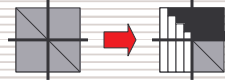
## Load Balancing

o Certain computations are inherently imbalanced … LU Decomposition is one



gray is balanced work, white & black are finished

o Standard block decomposition quickly becomes very biased



  n Cyclic and block cyclic allocation are one fix

Details in the appendix

25

## Work Queue, A Balancing Act(ivity)

o When work is generated dynamically, a work queue leads to reasonable balance

o Solution components:
  n Queue -- often a cyclically managed global array
  n Work or task descriptors -- succinct statements of the task to perform, often only index or bound

o Process:
  n As work is created, *append* to queue
  n On task completion, *remove* work from Q-head

26

# Work Queue Considerations

o Synchronization

    n Threads modifying the queue do so exclusively

    n Independent modification of head (*remove*) & tail (*append*) is dangerous, tough to do right



**tail**↑    ↑**head**

o Work quantum (grain) must exceed the overhead of queue manipulation and comm

    n Large quantum may lead to poor "end game"

---

# Considerations (Continued)

o Ordering, though strict at queue *remove,* may not be strict when viewed at completion

o Work assignment interacts with locality

    n Want strategies where descriptor embodies all data needed for task: chess move search

    n Avoid global memory reference by caching data: floor planning, place and route, etc.

    n When global reference inevitable, grain size should must cover comm costs too

Work quantum must be variable, not scalable

# Break

# Essay On 'Quicksort'

o Discuss

"Quicksort is an example of a sequential algorithm that is a good candidate to be incrementally transformed into a parallel algorithm"

o Where does the belief come from?

o What aspects are compatible/incompatible with the CTA?
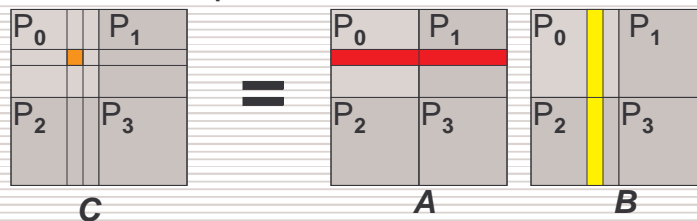
# Reconceptualizing a Computation

o Good parallel solutions result from rethinking a computation …

  n Sometimes that amounts to ordering scalar operations

  n Sometimes it requires starting from scratch

o The SUMMA matrix multiplication algorithm is the poster computation for rethinking!
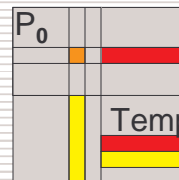
This computation is part of homework assignment

31

# Return To A Lecture 1 Computation

o Matrix Multiplication on Processor Grid



**C**          **A**          **B**

  n Matrices **A** and **B** producing $n$ x $n$ result **C** where

  $$C_{rs} = \sum_{1 \le k \le n} A_{rk} * B_{ks}$$

Temp

32

16

## Applying Scalable Techniques

o Assume each processor stores block of **C**, **A**, **B**; assume "can't" store all of any matrix

o To compute $c_{rs}$ a processor needs all of row $r$ of **A** and column $s$ of **B**.

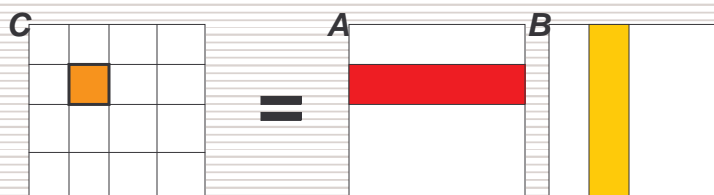o Consider strategies for minimizing data movement, because that is the greatest cost -- what are they?

$P_0$

Temp

$\blacksquare = \dfrac{\blacksquare_1}{\square_1} + \dfrac{\blacksquare_2}{\square_2} + \ldots + \dfrac{\blacksquare_n}{\square_n}$

33

---

## Grab All Rows/Columns At Once

o If all rows/columns are present, it's local

**C**   **A**   **B**

=

•**Each element requires O(n) operations**
•**Modern pipelined processors benefit**
  **from large blocks of work**
•**But memory space and BW are issues**

34

---

17

# Process *t* x *t* Blocks
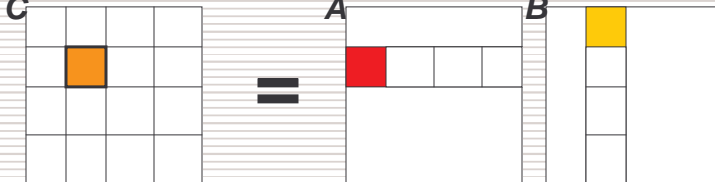
○ Referring to local storage

```
for (r=0; r < t; r++){
    for (s=0; s < t; s++){
        c[r][s] = 0.0;
        for (k=0; k < n; k++){
            c[r][s] += a[r][k]*b[k][s];
        }
    }
}
```
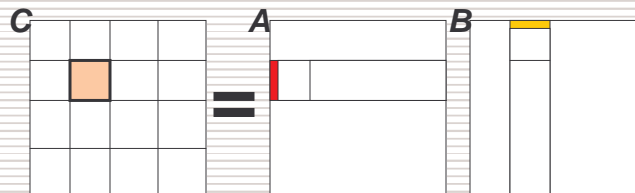
*Sweeter caching*



**C**     **A**     **B**

=

Only move a *t* x *t* block at a time

35

---

# Change Of View Point

○ Don't think of row-times-column



**C**     **A**     **B**

=

$$a_{11} \begin{array}{|c|c|} \hline & b_{11} \quad b_{12} \\ \hline a_{11}b_{11} & a_{11}b_{12} \\ \hline a_{21}b_{11} & a_{21}b_{12} \\ \hline \end{array} a_{21}$$

$$\blacksquare = \blacksquare_1 + \blacksquare_2 + \ldots + \blacksquare_n$$
$$* \quad * \quad \ldots \quad *$$
$$\square_1 \quad \square_2 \quad \square_n$$

**Switch orientation -- by using a *column* of *A* and a *row* of *B* compute all 1st terms of the dot products**
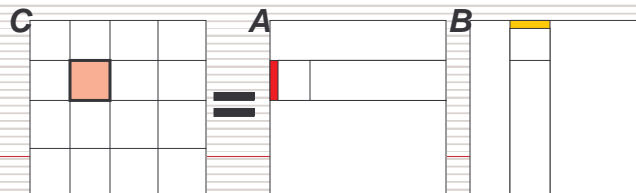
36

18

# SUMMA

o Scalable Universal Matrix Multiplication Alg
  - n Invented by van de Geijn & Watts of UT Austin
  - n Claimed best machine independent MM
o Whereas MM is usually A row x B column, SUMMA is A column x B row because computation switches from sense
  - n Normal: Compute all terms of dot product
  - n SUMMA: Computer first term of all dot products
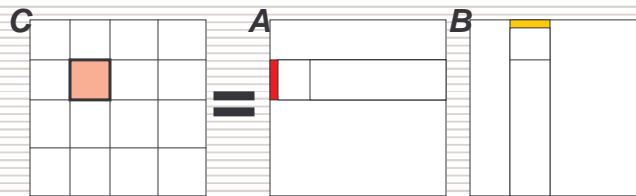
Strange. But fast!

37

# SUMMA Assumptions

o Threads have two indices, handle t x t block
o Let $p = P^{1/2}$, then thread u,v
  - n reads all columns of A for indices u*t:(u+1)*t-1,j
  - n reads all rows of B for indices i,v*t:(v+1)*t-1
  - n The arrays will be in global memory and referenced as needed

C        A        B

38

## Solve SUMMA In Peril-L

o Arrays in Global Memory

o Consider how processors reference arrays

    n Where do operands come from?

    n Where are results saved?

*C*        *A*        *B*



www.cs.utexas.edu/users/rvdg/abstracts/SUMMA.html

39

## Homework

o Read the remainder of Chapter 5

o Write a version of SUMMA in Peril-*L*

o Decide on Project Topic -- write a paragraph description of the computation that you expect to solve for the term project

o We will use a cluster making MPI and ZPL the most amenable languages we'll discuss

40