

Hardware/Software Performance Tradeoffs (plus Msg Passing Finish)

CS 258, Spring 99
David E. Culler
Computer Science Division
U.C. Berkeley

Message Passing Grid Solver

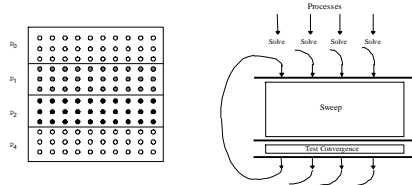
- Cannot declare A to be global shared array
 - compose it logically from per-process private arrays
 - usually allocated in accordance with the assignment of work
 - process assigned a set of rows allocates them locally
- Transfers of entire rows between traversals
- Structurally similar to SPMD SAS
- Orchestration different
 - data structures and data access/naming
 - communication
 - synchronization
- Ghost rows

```

10. procedure Solver()
11. begin
12.   int i, j, pid, n = n/nprows, done = 0;
13.   float temp, tempdiff, mydiff = 0; /*private variables*/
14.   float temp, tempdiff, mydiff = 0; /*private variables*/
15.   myA = malloc(a 2-d array of size [n/nprows + 2] by n+2); /*allocate my rows of A, in an unspecified way*/
16.   while (!done) do
17.     mydiff = 0; /*set local diff to 0*/
18.     /* Exchange border rows of neighbor into myA[i] and myA[i+1] */
19.     if (pid != 0) then SEND myA[i], n*sizeof(float), pid-1, ROW;
20.     if (pid == nprows-1) then
21.       SEND myA[n], n*sizeof(float), pid+1, ROW;
22.     if (pid != 0) then RECEIVE myA[i], n*sizeof(float), pid-1, ROW;
23.     if (pid == nprows-1) then
24.       RECEIVE myA[n], n*sizeof(float), pid+1, ROW;
25.     for i = 1 to n do
26.       for j = 1 to n do
27.         /*for each of my (nonghost) rows*/
28.         /*for all nonborder elements in that row*/
29.         temp = myA[i, j];
30.         myA[i, j] = 0.5 * ( myA[i, j] + myA[i-1, j] + myA[i+1, j] +
31.           myA[i, j-1] + myA[i, j+1] );
32.         mydiff += abs( myA[i, j] - temp );
33.       endfor
34.     endfor
35.     /*communicate local diff values and determine if
36.     done; can be replaced by reduction and broadcast*/
37.     if (pid == 0) then
38.       SEND mydiff, sizeof(float), 0, DIFF; /*process 0 holds global total diff*/
39.     else
40.       RECEIVE mydiff, sizeof(float), 0, DIFF;
41.     endif
42.     /*pid 0 does this*/
43.     for i = 1 to nprows-1 do
44.       for each other process*
45.         RECEIVE mydiff, sizeof(float), i, DIFF;
46.     endfor
47.     mydiff += tempdiff; /*accumulate into total*/
48.     /*pid 0 does this*/
49.     if (tempdiff/(n*n) < TOL) then done = 1;
50.     for i = 1 to nprows-1 do
51.       for each other process*
52.         SEND mydiff, sizeof(float), i, DONE;
53.     endfor
54.   end
55. endwhile
56. end procedure

```

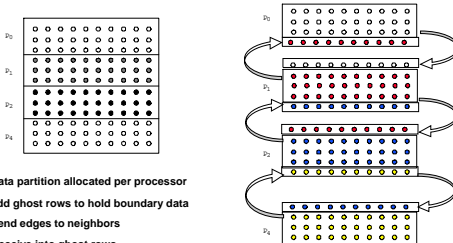
SAS Recap



- Partitioning = Decomposition + Assignment
- Orchestration = coordination and communication
 - SPMD, Static Assignment
 - Implicit communication
 - Explicit Synchronization: barriers, mutex, events

2/3/99 CS258 S99.5 2

Data Layout and Orchestration



Data partition allocated per processor
Add ghost rows to hold boundary data
Send edges to neighbors
Receive into ghost rows
Compute as in sequential program

2/3/99 CS258 S99.5 4

Notes on Message Passing Program

- Use of ghost rows
- Receive does not transfer data, send does
 - unlike SAS which is usually receiver-initiated (load fetches data)
- Communication done at beginning of iteration, so no asynchrony
- Communication in whole rows, not element at a time
- Core similar, but indices/bounds in local rather than global space
- Synchronization through sends and receives
 - Update of global diff and event synch for done condition
 - CANNOT implement locks and barriers with messages
- REDUCE and BROADCAST simplify code

```

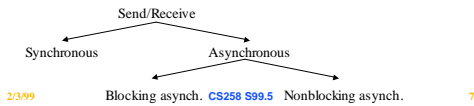
/*communicate local diff values and determine if done, using reduction and broadcast*/
25b. REDUCE 0, mydiff, sizeof(float), ADD;
25c. if (pid == 0) then
25d.   if (mydiff/(n*n) < TOL) then done = 1;
25e. endif
25f. BROADCAST 0, done, sizeof(int), DONE;
25g.
25h.

```

2/3/99 CS258 S99.5 6

Send and Receive Alternatives

- **extended functionality:** stride, scatter-gather, groups
- **Synchronization semantics**
 - Affect when data structures or buffers can be reused at either end
 - Affect event synch (mutual excl. by fiat: only one process touches data)
 - Affect ease of programming and performance
- **Synchronous messages provide built-in synch. through match**
 - Separate event synchronization may be needed with asynch. messages
- **With synch. messages, our code may hang. Fix?**



Orchestration: Summary

- **Shared address space**
 - Shared and private data (explicitly separate ??)
 - Communication implicit in access patterns
 - Data distribution not a correctness issue
 - Synchronization via atomic operations on shared data
 - Synchronization explicit and distinct from data communication
- **Message passing**
 - Data distribution among local address spaces needed
 - No explicit shared structures
 - » implicit in comm. patterns
 - Communication is explicit
 - Synchronization implicit in communication
 - » mutual exclusion by fiat

2/3/99 CS258 S99.5 8

Correctness in Grid Solver Program

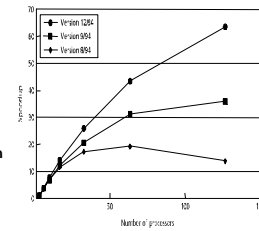
	SAS	Msg-Passing
Explicit global data structure?	Yes	No
Assignment indept of data layout?	Yes	No
Communication	Implicit	Explicit
Synchronization	Explicit	Implicit
Explicit replication of border rows?	No	Yes

- **Decomposition and Assignment similar in SAS and message-passing**
- **Orchestration is different**
 - Data structures, data access/naming, communication, synchronization
 - Performance?

2/3/99 CS258 S99.5 9

Performance Goal => Speedup

- **Architect Goal**
 - observe how program uses machine and improve the design to enhance performance
- **Programmer Goal**
 - observe how the program uses the machine and improve the implementation to enhance performance
- **What do you observe?**
- **Who fixes what?**



2/3/99 CS258 S99.5 10

Analysis Framework

$$\text{Speedup} \leq \frac{\text{Sequential Work}}{\text{Max (Work + Synch Wait Time + Comm Cost + Extra Work)}}$$

- **Solving communication and load balance NP-hard in general case**
 - But simple heuristic solutions work well in practice
- **Fundamental Tension among:**
 - balanced load
 - minimal synchronization
 - minimal communication
 - minimal extra work
- **Good machine design mitigates the trade-offs**

2/3/99 CS258 S99.5 11

Load Balance and Synchronization

$$\text{Speedup}_{\text{problem}(p)} \leq \frac{\text{Sequential Work}}{\text{Max Work on any Processor}}$$

- **Instantaneous load imbalance revealed as wait time**
 - at completion
 - at barriers
 - at receive
 - at flags, even at mutex

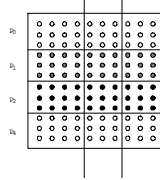
$$\frac{\text{Sequential Work}}{\text{Max (Work + Synch Wait Time)}}$$

2/3/99 CS258 S99.5 12

Improving Load Balance

- Decompose into more smaller tasks ($\gg P$)
- Distribute uniformly
 - variable sized task
 - randomize
 - bin packing
 - dynamic assignment
- Schedule more carefully
 - avoid serialization
 - estimate work
 - use history info.

```
for_all i = 1 to n do
  for_all j = 1 to n do
    A[i, j] = A[i-1, j] + A[i, j-1] + ...
```

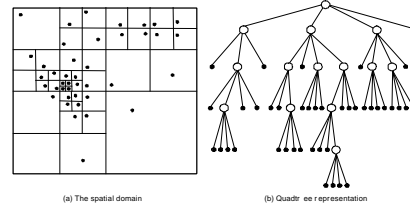


2/3/99

CS258 S99.5

13

Example: Barnes-Hut



- Divide space into roughly equal # particles
- Particles close together in space should be on same processor
- Nonuniform, dynamically changing

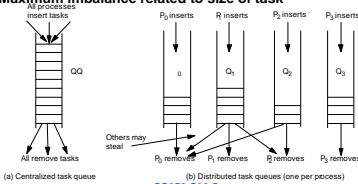
2/3/99

CS258 S99.5

14

Dynamic Scheduling with Task Queues

- Centralized versus distributed queues
- Task stealing with distributed queues
 - Can compromise comm and locality, and increase synchronization
 - Whom to steal from, how many tasks to steal, ...
 - Termination detection
 - Maximum imbalance related to size of task



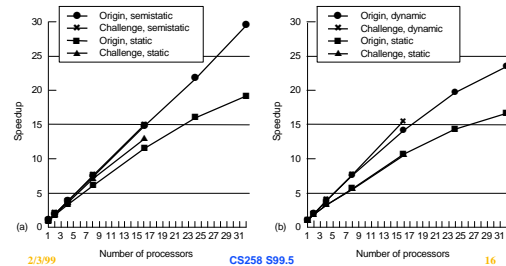
2/3/99

CS258 S99.5

15

Impact of Dynamic Assignment

- Barnes-Hut on SGI Origin 2000 (cache-coherent shared memory):



2/3/99

CS258 S99.5

16

Self-Scheduling

```
volatile int row_index = 0; /* shared index variable */

while (not done) {
  initialize row_index; barrier;
  while ((i = fetch_and_inc(&row_index) < n) {
    for (j = 1; j < n; j++) {
      A[i, j] = A[i-1, j] + A[i, j-1] + ...
    }
  }
}
```

2/3/99

CS258 S99.5

17

Reducing Serialization

- Careful about assignment and orchestration
 - including scheduling
- Event synchronization
 - Reduce use of conservative synchronization
 - » e.g. point-to-point instead of barriers, or granularity of pt-to-pt
 - But fine-grained synch more difficult to program, more synch ops.
- Mutual exclusion
 - Separate locks for separate data
 - » e.g. locking records in a database: lock per process, record, or field
 - » lock per task in task queue, not per queue
 - » finer grain \Rightarrow less contention/serialization, more space, less reuse
 - Smaller, less frequent critical sections
 - » don't do reading/testing in critical section, only modification
 - Stagger critical sections in time

2/3/99

CS258 S99.5

18

Impact of Efforts to Balance Load

- **Parallelism Management overhead?**
- **Communication?**
 - amount, size, frequency?
- **Synchronization?**
 - type? frequency?
- **Opportunities for replication?**
- **What can architecture do?**

2/3/99

CS258 S99.5

19

Arch. Implications of Load Balance

- **Naming**
 - global position independent naming separates decomposition from layout
 - allows diverse, even dynamic assignments
- **Efficient Fine-grained communication & synch**
 - more, smaller
 - » msgs
 - » locks
 - point-to-point
- **Automatic replication**

2/3/99

CS258 S99.5

20

Reducing Extra Work

- **Common sources of extra work:**
 - **Computing a good partition**
 - » e.g. partitioning in Barnes-Hut or sparse matrix
 - **Using redundant computation to avoid communication**
 - **Task, data and process management overhead**
 - » applications, languages, runtime systems, OS
 - **Imposing structure on communication**
 - » coalescing messages, allowing effective naming
- **Architectural Implications:**
 - Reduce need by making communication and orchestration efficient

$$\text{Speedup} \leq \frac{\text{Sequential Work}}{\text{Max (Work + Synch Wait Time + Comm Cost + Extra Work)}}$$

2/3/99

CS258 S99.5

21

Reducing Inherent Communication

$$\text{Speedup} \leq \frac{\text{Sequential Work}}{\text{Max (Work + Synch Wait Time + Comm Cost)}}$$

- **Communication is expensive!**
- **Measure: communication to computation ratio**
- **Inherent communication**
 - Determined by assignment of tasks to processes
 - One produces data consumed by others
- ⇒ **Use algorithms that communicate less**
- ⇒ **Assign tasks that access same data to same process**
 - same row or block to same process in each iteration

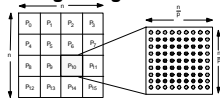
2/3/99

CS258 S99.5

22

Domain Decomposition

- **Works well for scientific, engineering, graphics, ... applications**
- **Exploits local-biased nature of physical problems**
 - Information requirements often short-range
 - Or long-range but fall off with distance
- **Simple example: nearest-neighbor grid computation**



Perimeter to Area comm-to-comp ratio (area to volume in 3-d)
 • Depends on n, p : decreases with n , increases with p

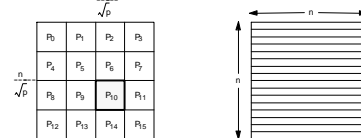
2/3/99

CS258 S99.5

23

Domain Decomposition (contd)

Best domain decomposition depends on information requirements
 Nearest neighbor example: block versus strip decomposition:



- **Comm to comp:** $\frac{d^{*n} p}{n}$ for block, $2^{*n} p$ for strip
- **Application dependent: strip may be better in other cases**
 - E.g. particle flow in tunnel

2/3/99

CS258 S99.5

24

Relation to load balance

- Scatter Decomposition, e.g. initial partition in Raytrace

12	
3	4

Domain decomposition

12		12		12		12
3	4	3	4	3	4	3
12		12		12		12
3	4	3	4	3	4	3
12		12		12		12
3	4	3	4	3	4	3
12		12		12		12
3	4	3	4	3	4	3

Scatter decomposition

Preserve locality in task stealing

- Steal large tasks for locality, steal from same queues, ...

2/3/99

CS258 S99.5

25

Implications of Comm-to-Comp Ratio

- Architects examine application needs to see where to spend effort
 - bandwidth requirements (operations / sec)
 - latency requirements (sec/operation)
 - time spent waiting
- Actual impact of comm. depends on structure and cost as well
- Need to keep communication balanced across processors as well

$$\text{Speedup} \leq \frac{\text{Sequential Work}}{\text{Max}(\text{Work} + \text{Synch Wait Time} + \text{Comm Cost})}$$

2/3/99

CS258 S99.5

26

Structuring Communication

- Given amount of comm, goal is to reduce cost
- Cost of communication as seen by process:

$$C = f * (o + l + \frac{n_c/m}{B} + t_c - \text{overlap})$$

- f = frequency of messages
 - o = overhead per message (at both ends)
 - l = network delay per message
 - n_c = total data sent
 - m = number of messages
 - B = bandwidth along path (determined by network, NI, assist)
 - t_c = cost induced by contention per message
 - overlap = amount of latency hidden by overlap with comp. or comm.
- Portion in parentheses is cost of a message (as seen by processor)
- ignoring overlap, is *latency* of a message

2/3/99

CS258 S99.5

27

– Goal: reduce terms in latency and increase overlap

Reducing Overhead

- Can reduce no. of messages m or overhead per message o
- o is usually determined by hardware or system software
 - Program should try to reduce m by coalescing messages
 - More control when communication is explicit
- Coalescing data into larger messages:
 - Easy for regular, coarse-grained communication
 - Can be difficult for irregular, naturally fine-grained communication
 - may require changes to algorithm and extra work
 - coalescing data and determining what and to whom to send
 - will discuss more in implications for programming models later

2/3/99

CS258 S99.5

28

Reducing Network Delay

- Network delay component = $f * h * t_h$
 - h = number of hops traversed in network
 - t_h = link+switch latency per hop
- Reducing f : communicate less, or make messages larger
- Reducing h :
 - Map communication patterns to network topology
 - e.g. nearest-neighbor on mesh and ring; all-to-all
 - How important is this?
 - used to be major focus of parallel algorithms
 - depends on no. of processors, how t_{pr} compares with other components
 - less important on modern machines
 - overheads, processor count, multiprogramming

2/3/99

CS258 S99.5

29

Reducing Contention

- All resources have nonzero occupancy
 - Memory, communication controller, network link, etc.
 - Can only handle so many transactions per unit time
- Effects of contention:
 - Increased end-to-end cost for messages
 - Reduced available bandwidth for individual messages
 - Causes imbalances across processors
- Particularly insidious performance problem
 - Easy to ignore when programming
 - Slow down messages that don't even need that resource
 - by causing other dependent resources to also congest
 - Effect can be devastating: *Don't flood a resource!*

2/3/99

CS258 S99.5

30

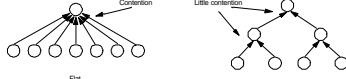
Types of Contention

- Network contention and end-point contention (hot-spots)

- Location and Module Hot-spots

- Location: e.g. accumulating into global variable, barrier

» solution: tree-structured communication



- Module: all-to-all personalized comm. in matrix transpose

- solution: stagger access by different processors to same node temporally

- In general, reduce burstiness; may conflict with making messages larger

2/3/99

CS258 S99.5

31

Overlapping Communication

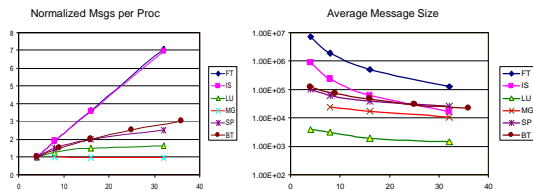
- Cannot afford to stall for high latencies
 - even on uniprocessors!
- Overlap with computation or communication to hide latency
- Requires extra concurrency (*slackness*), higher bandwidth
- Techniques:
 - Prefetching
 - Block data transfer
 - Proceeding past communication
 - Multithreading

2/3/99

CS258 S99.5

32

Communication Scaling (NPB2)

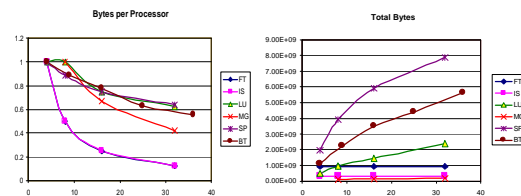


2/3/99

CS258 S99.5

33

Communication Scaling: Volume



2/3/99

CS258 S99.5

34

What is a Multiprocessor?

- A collection of communicating processors
 - View taken so far
 - Goals: balance load, reduce inherent communication and extra work
- A multi-cache, multi-memory system
 - Role of these components essential regardless of programming model
 - Prog. model and comm. abstr. affect specific performance tradeoffs

2/3/99

CS258 S99.5

35

Memory-oriented View

- Multiprocessor as Extended Memory Hierarchy
 - » as seen by a given processor
- Levels in extended hierarchy:
 - Registers, caches, local memory, remote memory (topology)
 - Glued together by communication architecture
 - Levels communicate at a certain granularity of data transfer
- Need to exploit spatial and temporal locality in hierarchy
 - Otherwise extra communication may also be caused
 - Especially important since communication is expensive

2/3/99

CS258 S99.5

36

Uniprocessor

- Performance depends heavily on memory hierarchy
- Time spent by a program
 - $Time_{prog}(1) = Busy(1) + Data\ Access(1)$
 - Divide by cycles to get CPI equation
- Data access time can be reduced by:
 - Optimizing machine: bigger caches, lower latency...
 - Optimizing program: temporal and spatial locality

2/3/99

CS258 S99.5

37

Extended Hierarchy

- Idealized view: local cache hierarchy + single main memory
- But reality is more complex
 - Centralized Memory: caches of other processors
 - Distributed Memory: some local, some remote; + network topology
 - Management of levels
 - » caches managed by hardware
 - » main memory depends on programming model
 - SAS: data movement between local and remote transparent
 - message passing: explicit
 - Levels closer to processor are lower latency and higher bandwidth
 - Improve performance through architecture or program locality
 - Tradeoff with parallelism; need good node performance and parallelism

2/3/99

CS258 S99.5

38

Artifactual Communication

- Accesses not satisfied in local portion of memory hierarchy cause communication
 - Inherent communication, implicit or explicit, causes transfers
 - » determined by program
 - Artifactual communication
 - » determined by program implementation and arch. interactions
 - » poor allocation of data across distributed memories
 - » unnecessary data in a transfer
 - » unnecessary transfers due to system granularities
 - » redundant communication of data
 - » finite replication capacity (in cache or main memory)
 - Inherent communication assumes unlimited capacity, small transfers, perfect knowledge of what is needed.

2/3/99

CS258 S99.5

39

Communication and Replication

- Comm induced by finite capacity is most fundamental artifact
 - Like cache size and miss rate or memory traffic in uniprocessors
 - Extended memory hierarchy view useful for this relationship
- View as three level hierarchy for simplicity
 - Local cache, local memory, remote memory (ignore network topology)
- Classify “misses” in “cache” at any level as for uniprocessors
 - » compulsory or cold misses (no size effect)
 - » capacity misses (yes)
 - » conflict or collision misses (yes)
 - » communication or coherence misses (no)
- Each may be helped/hurt by large transfer granularity (spatial locality)

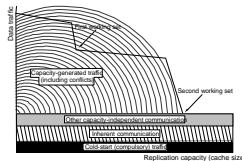
2/3/99

CS258 S99.5

40

Working Set Perspective

- At a given level of the hierarchy (to the next further one)



- Hierarchy of working sets
- At first level cache (fully assoc, one-word block), inherent to algorithm
 - » working set curve for program
- Traffic from any type of miss can be local or nonlocal (communication)

2/3/99

CS258 S99.5

41

Orchestration for Performance

- Reducing amount of communication:
 - Inherent: change logical data sharing patterns in algorithm
 - Artifactual: exploit spatial, temporal locality in extended hierarchy
 - » Techniques often similar to those on uniprocessors
- Structuring communication to reduce cost

2/3/99

CS258 S99.5

42

Reducing Artfactual Communication

- **Message passing model**
 - Communication and replication are both explicit
 - Even artifactual communication is in explicit messages
 - » send data that is not used
- **Shared address space model**
 - More interesting from an architectural perspective
 - Occurs transparently due to interactions of program and system
 - » sizes and granularities in extended memory hierarchy
- **Use shared address space to illustrate issues**

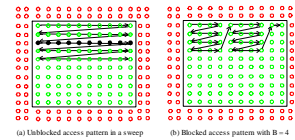
2/3/99

CS258 S99.5

43

Exploiting Temporal Locality

- Structure algorithm so working sets map well to hierarchy
 - » often techniques to reduce inherent communication do well here
 - » schedule tasks for data reuse once assigned
- Multiple data structures in same phase
 - » e.g. database records: local versus remote
- Solver example: blocking



- More useful when $O(n^{k+1})$ computation on $O(n^k)$ data
 - many linear algebra computations (factorization, matrix multiply)

2/3/99

CS258 S99.5

44

Exploiting Spatial Locality

- Besides capacity, granularities are important:
 - Granularity of allocation
 - Granularity of communication or data transfer
 - Granularity of coherence
- Major spatial-related causes of artifactual communication:
 - Conflict misses
 - Data distribution/layout (allocation granularity)
 - Fragmentation (communication granularity)
 - False sharing of data (coherence granularity)
- All depend on how spatial access patterns interact with data structures
 - Fix problems by modifying data structures, or layout/alignment
- Examine later in context of architectures
 - one simple example here: data distribution in SAS solver

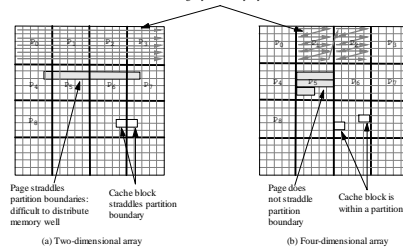
2/3/99

CS258 S99.5

45

Spatial Locality Example

- Repeated sweeps over 2-d grid, each time adding 1 to elements
- Natural 2-d versus higher-dimensional array representation



2/3/99

CS258 S99.5

46

Architectural Implications of Locality

- Communication abstraction that makes exploiting it easy
- For cache-coherent SAS, e.g.:
 - Size and organization of levels of memory hierarchy
 - » cost-effectiveness: caches are expensive
 - » caveats: flexibility for different and time-shared workloads
 - Replication in main memory useful? If so, how to manage?
 - » hardware, OS/runtime, program?
 - Granularities of allocation, communication, coherence (?)
 - » small granularities => high overheads, but easier to program
- Machine granularity (resource division among processors, memory...)

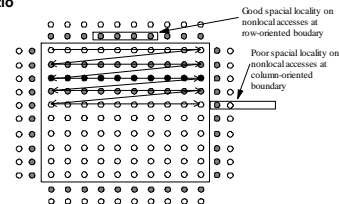
2/3/99

CS258 S99.5

47

Tradeoffs with Inherent Communication

- Partitioning grid solver: blocks versus rows
 - Blocks still have a spatial locality problem on remote data
 - Rowwise can perform better despite worse inherent c-to-c ratio



- Result depends on n and p

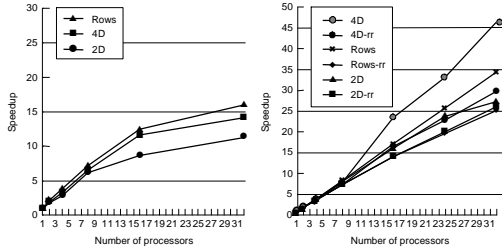
2/3/99

CS258 S99.5

48

Example Performance Impact

• Equation solver on SGI Origin2000

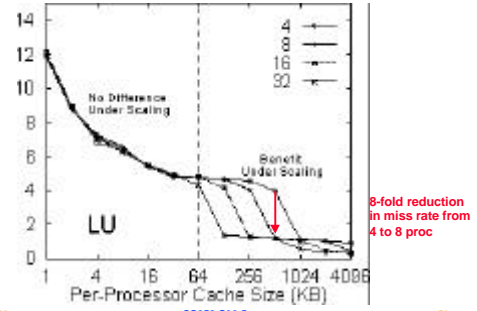


2/3/99

CS258 S99.5

49

Working Sets Change with P

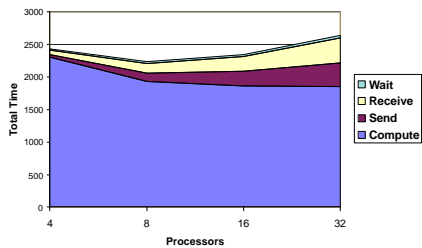


2/3/99

CS258 S99.5

50

Where the Time Goes: LU-a



2/3/99

CS258 S99.5

51

Summary of Tradeoffs

- Different goals often have conflicting demands
 - Load Balance
 - » fine-grain tasks
 - » random or dynamic assignment
 - Communication
 - » usually coarse grain tasks
 - » decompose to obtain locality: not random/dynamic
 - Extra Work
 - » coarse grain tasks
 - » simple assignment
 - Communication Cost:
 - » big transfers: amortize overhead and latency
 - » small transfers: reduce contention

2/3/99

CS258 S99.5

52