

## Perspective on Parallel Programming

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

## Outline for Today

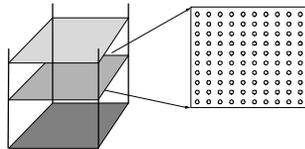
- Motivating Problems (application case studies)
- Process of creating a parallel program
- What a simple parallel program looks like
  - three major programming models
  - What primitives must a system support?
- Later: Performance issues and architectural interactions

1/29/99

CS258 S99

2

## Simulating Ocean Currents



(a) Cross sections

(b) Spatial discretization of a cross section

- Model as two-dimensional grids
  - Discretize in space and time
  - finer spatial and temporal resolution => greater accuracy
- Many different computations per time step
  - » set up and solve equations
  - Concurrency across and within grid computations
- Static and regular

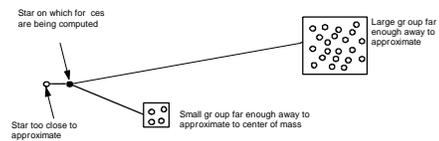
1/29/99

CS258 S99

3

## Simulating Galaxy Evolution

- Simulate the interactions of many stars evolving over time
- Computing forces is expensive
  - $O(n^2)$  brute force approach
  - Hierarchical Methods take advantage of force law:  $G \frac{m_1 m_2}{r^2}$



- Many time-steps, plenty of concurrency across stars within one

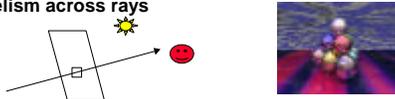
1/29/99

CS258 S99

4

## Rendering Scenes by Ray Tracing

- Shoot rays into scene through pixels in image plane
- Follow their paths
  - they bounce around as they strike objects
  - they generate new rays: ray tree per input ray
- Result is color and opacity for that pixel
- Parallelism across rays



- How much concurrency in these examples?

1/29/99

CS258 S99

5

## Creating a Parallel Program

- Pieces of the job:
  - Identify work that can be done in parallel
    - » work includes computation, data access and I/O
  - Partition work and perhaps data among processes
  - Manage data access, communication and synchronization

1/29/99

CS258 S99

6

## Definitions

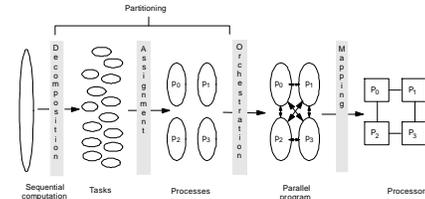
- **Task:**
  - Arbitrary **piece of work** in parallel computation
  - Executed sequentially; concurrency is only across tasks
  - E.g. a particle/cell in Barnes-Hut, a ray or ray group in Raytrace
  - Fine-grained versus coarse-grained tasks
- **Process (thread):**
  - Abstract entity that performs the tasks assigned to processes
  - Processes communicate and synchronize to perform their tasks
- **Processor:**
  - Physical engine on which process executes
  - Processes virtualize machine to programmer
    - » write program in terms of processes, then map to processors

1/29/99

CS258 S99

7

## 4 Steps in Creating a Parallel Program



- **Decomposition** of computation in tasks
- **Assignment** of tasks to processes
- **Orchestration** of data access, comm, synch.
- **Mapping** processes to processors

1/29/99

CS258 S99

8

## Decomposition

- Identify concurrency and decide level at which to exploit it
- Break up computation into tasks to be divided among processes
  - Tasks may become available dynamically
  - No. of available tasks may vary with time
- Goal: Enough tasks to keep processes busy, but not too many
  - Number of tasks available at a time is upper bound on achievable speedup

1/29/99

CS258 S99

9

## Limited Concurrency: Amdahl's Law

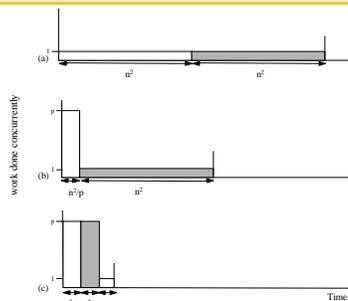
- Most fundamental limitation on parallel speedup
- If fraction  $s$  of seq execution is inherently serial, speedup  $\leq 1/s$
- Example: 2-phase calculation
  - sweep over  $n$ -by- $n$  grid and do some independent computation
  - sweep again and add each value to global sum
- Time for first phase =  $n^2/p$
- Second phase serialized at global variable, so time =  $n^2$
- Speedup  $\leq \frac{2n^2}{\frac{n^2}{p} + n^2}$  or at most 2
- Trick: divide second phase into two
  - accumulate into private sum during sweep
  - add per-process private sum into global sum
- Parallel time is  $n^2/p + n^2/p + p$ , and speedup at best  $\frac{2n^2}{2n^2 + p^2}$

1/29/99

CS258 S99

10

## Understanding Amdahl's Law

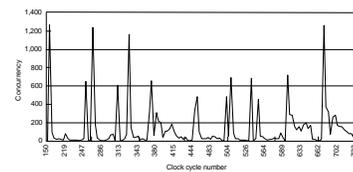


1/29/99

CS258 S99

11

## Concurrency Profiles



- Area under curve is total work done, or time with 1 processor
- Horizontal extent is lower bound on time (infinite processors)
- Speedup is the ratio:  $\frac{\sum_{k=1}^n f_k k}{\sum_{k=1}^n \frac{k}{p}}$ , base case:  $\frac{1}{s + \frac{1-s}{p}}$
- Amdahl's law applies to any overhead, not just limited concurrency

1/29/99

CS258 S99

12

## Assignment

- Specify **mechanism** to divide work up among processes
  - E.g. which process computes forces on which stars, or which rays
  - Balance workload, reduce communication and management cost
- Structured approaches usually work well
  - Code inspection (parallel loops) or understanding of application
  - Well-known heuristics
  - *Static* versus *dynamic* assignment
- As programmers, we worry about partitioning first
  - Usually independent of architecture or prog model
  - But cost and complexity of using primitives may affect decisions

1/29/99

CS258 S99

13

## Orchestration

- Naming data
- Structuring communication
- Synchronization
- Organizing data structures and scheduling tasks temporally
- **Goals**
  - Reduce cost of communication and synch.
  - Preserve locality of data reference
  - Schedule tasks to satisfy dependences early
  - Reduce overhead of parallelism management
- Choices depend on Prog. Model., comm. abstraction, efficiency of primitives
- Architects should provide appropriate primitives efficiently

1/29/99

CS258 S99

14

## Mapping

- **Two aspects:**
  - Which process runs on which particular processor?
    - » mapping to a network topology
  - Will multiple processes run on same processor?
- **space-sharing**
  - Machine divided into subsets, only one app at a time in a subset
  - Processes can be pinned to processors, or left to OS
- **System allocation**
- **Real world**
  - User specifies desires in some aspects, system handles some
- **Usually adopt the view: process <-> processor**

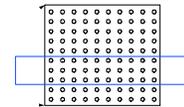
1/29/99

CS258 S99

15

## Parallelizing Computation vs. Data

- Computation is decomposed and assigned (partitioned)
- Partitioning Data is often a natural view too
  - Computation follows data: *owner computes*
  - Grid example; data mining;
- Distinction between comp. and data stronger in many applications
  - Barnes-Hut
  - Raytrace



1/29/99

CS258 S99

16

## Architect's Perspective

- What can be addressed by better hardware design?
- What is fundamentally a programming issue?

1/29/99

CS258 S99

17

## High-level Goals

Table 2.1 Steps in the Parallelization Process and Their Goals

| Step          | Architecture-Dependent? | Major Performance Goals   |
|---------------|-------------------------|---|
| Decomposition | Mostly no               | Expose enough concurrency but not too much  |
| Assignment    | Mostly no               | Balance workload<br>Reduce communication volume   |
| Orchestration | Yes                     | Reduce noninherent communication via data locality<br>Reduce communication and synchronization cost as seen by the processor<br>Reduce serialization at shared resources<br>Schedule tasks to satisfy dependences early |
| Mapping       | Yes                     | Put related processes on the same processor if necessary<br>Exploit locality in network topology  |

- High performance (speedup over sequential program)
- But low resource usage and development effort
- Implications for algorithm designers and architects?

1/29/99

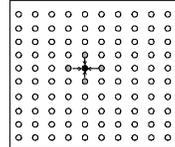
CS258 S99

18

## What Parallel Programs Look Like

## Example: iterative equation solver

- Simplified version of a piece of Ocean simulation
- Illustrate program in low-level parallel language
  - C-like pseudocode with simple extensions for parallelism
  - Expose basic comm. and synch. primitives
  - State of most real parallel programming today



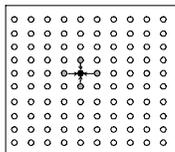
Expression for updating each interior point:  
 $A[i,j] = 0.2 * (A[i,j] + A[i,j-1] + A[i-1,j] + A[i,j+1] + A[i+1,j])$

1/29/99

CS258 S99

20

## Grid Solver



Expression for updating each interior point:  
 $A[i,j] = 0.2 * (A[i,j] + A[i,j-1] + A[i-1,j] + A[i,j+1] + A[i+1,j])$

- Gauss-Seidel (near-neighbor) sweeps to convergence
  - interior  $n$ -by- $n$  points of  $(n+2)$ -by- $(n+2)$  updated in each sweep
  - updates done in-place in grid
  - difference from previous value computed
  - accumulate partial diffs into global diff at end of every sweep
  - check if has converged
    - » to within a tolerance parameter

1/29/99

CS258 S99

21

## Sequential Version

```

1. int n;
2. float **A, diff = 0; //size of matrix: (n+2)-by-(n+2) elements*/
3. main()
4. begin
5.   read(n); //read input parameter: matrix size*/
6.   A ← malloc (a 2-d array of size n + 2 by n + 2 doubles);
7.   initialize(A); //initialize the matrix A somehow*/
8.   Solve (A); //call the routine to solve equations*/
9. end main

10.procedure Solve (A) //solve the equation system*/
11. float **A; //A is an (n+2)-by-(n+2)array*/
12.begin
13. int i, j, done = 0;
14. float diff = 0, temp;
15. while (!done) do //outermost loop over sweeps*/
16.   diff = 0; //initialize maximum difference to 0*/
17.   for i ← 1 to n do //sweep over nonborder points of grid*/
18.     for j ← 1 to n do
19.       temp = A[i,j]; //save old value of element*/
20.       A[i,j] ← 0.2 * (A[i,j] + A[i,j-1] + A[i-1,j] +
21.         A[i,j+1] + A[i+1,j]); //compute average*/
22.       diff ← abs(A[i,j] - temp);
23.     end for
24.   end for
25.   if (diff/(n*n) < TOL) then done = 1;
26. end while
27.end procedure
    
```

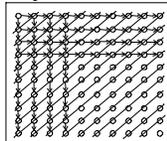
1/29/99

CS258 S99

22

## Decomposition

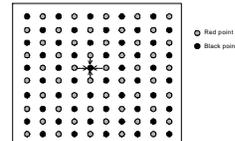
- Simple way to identify concurrency is to look at loop iterations
  - *dependence analysis*; if not enough concurrency, then look further
- Not much concurrency here at this level (all loops *sequential*)
- Examine fundamental dependences



- Concurrency  $O(n)$  along anti-diagonals, serialization  $O(n)$  along diag.
- Retain loop structure, use pt-to-pt synch; Problem: too many synch ops.
- Restructure loops, use global synch; Imbalance and too much synch<sup>3</sup>

## Exploit Application Knowledge

- Reorder grid traversal: red-black ordering



- Different ordering of updates: may converge quicker or slower
- Red sweep and black sweep are each fully parallel:
- Global synch between them (conservative but convenient)
- Ocean uses red-black
- We use simpler, asynchronous one to illustrate
  - » no red-black, simply ignore dependences within sweep
  - » parallel program *nondeterministic*

1/29/99

CS258 S99

24

## Decomposition

```

15. while (!done) do /*a sequential loop*/
16.   diff = 0;
17.   for_all i ← 1 to n do /*a parallel loop nest*/
18.     for_all j ← 1 to n do
19.       temp = A[i,j];
20.       A[i,j] ← 0.2 * (A[i,j] + A[i,j-1] + A[i-1,j] +
21.         A[i,j+1] + A[i+1,j]);
22.       diff += abs(A[i,j] - temp);
23.     end for_all
24.   end for_all
25.   if (diff/(n*n) < TOL) then done = 1;
26. end while

```

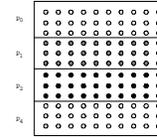
- Decomposition into elements: degree of concurrency  $n^2$
- Decompose into rows? Degree ?
- for\_all assignment ??

1/29/99

CS258 S99

25

## Assignment



- Static assignment: decomposition into rows
  - block assignment of rows: Row  $i$  is assigned to process  $p$
  - cyclic assignment of rows: process  $i$  is assigned rows  $i, i+p, \dots$

- Dynamic assignment

× get a row index, work on the row, get a new row, ...

- What is the mechanism?
- Concurrency? Volume of Communication?

1/29/99

CS258 S99

26

## Data Parallel Solver

```

1. int n, nprocs; /*grid size (n+2-by-n+2) and number of processes*/
2. float **A, diff = 0;
3. main()
4. begin
5.   read(n); read(nprocs); /*read input grid size and number of processes*/
6.   A ← G_MALLOC (a 2-d array of size n+2 by n+2 doubles);
7.   initialize(A); /*initialize the matrix A somehow*/
8.   Solve(A); /*call the routine to solve equation*/
9. end main
10. procedure Solve(A) /*solve the equation system*/
11.   float **A; /*A is an (n+2-by-n+2) array*/
12. begin
13.   int i, j, done = 0;
14.   float mydiff = 0, temp;
15.   DECOMP A[BLOCK*, nprocs];
16.   while (!done) do /*outermost loop over sweeps*/
17.     mydiff = 0; /*initialize maximum difference to 0*/
18.     for_all i ← 1 to n do /*sweep over non-border points of grid*/
19.       for_all j ← 1 to n do /*save old value of element*/
20.         temp = A[i,j];
21.         A[i,j] ← 0.2 * (A[i,j] + A[i,j-1] + A[i-1,j] +
22.           A[i,j+1] + A[i+1,j]); /*compute average*/
23.         mydiff += abs(A[i,j] - temp);
24.       end for_all
25.     end for_all
26.     mydiff += abs(mydiff, diff, ADD);
27.     if (diff/(n*n) < TOL) then done = 1;
28.   end while
29. end procedure

```

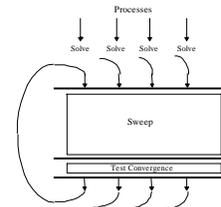
1/29/99

CS258 S99

27

## Shared Address Space Solver

Single Program Multiple Data (SPMD)



- Assignment controlled by values of variables used as loop bounds

1/29/99

CS258 S99

28

## Generating Threads

```

1. int n, nprocs; /*matrix dimension and number of processors to be used*/
2a. float **A, diff; /*A is global (shared) array representing the grid*/
2b. LOCKDEC(diff_lock); /*diff is global (shared) maximum difference in current sweep*/
2c. BARDEC (bar1); /*barrier declaration for global synchronization between sweeps*/
3. main()
4. begin
5.   read(n); read(nprocs); /*read input matrix size and number of processes*/
6.   A ← G_MALLOC (a two-dimensional array of size n+2 by n+2
7.     doubles); /*initialize A in an unspecified way*/
8a. CREATE (nprocs-1, Solve, A); /*main process becomes a worker too*/
8b. WAIT_FOR_END (nprocs-1); /*wait for all child processes created to terminate*/
9. end main
10. procedure Solve(A) /*A is entire n+2-by-n+2 shared array,
11.   float **A; as in the sequential program*/
12. begin
13. ---
14. end procedure

```

1/29/99

CS258 S99

29

## Assignment Mechanism

```

10. procedure Solve(A)
11.   float **A; /*A is entire n+2-by-n+2 shared array,
12.   as in the sequential program*/
13. begin
14.   int i, j, pid, done = 0;
15.   float temp, mydiff = 0; /*private variables*/
16a. int mymin = 1 + (pid * n/nprocs); /*assume that n is exactly divisible by*/
16b. int mymax = mymin + n/nprocs - 1 /*nprocs for simplicity here*/
17. while (!done) do /*outer loop sweeps*/
18.   mydiff = diff = 0; /*set global diff to 0 (okay for all to do it)*/
19.   BARRIER(bar1, nprocs); /*ensure all reach here before anyone modifies diff*/
20.   for i ← mymin to mymax do /*for each of my rows*/
21.     temp = A[i,j]; /*for all nonborder elements in that row*/
22.     A[i,j] = 0.2 * (A[i,j] + A[i,j-1] + A[i-1,j] +
23.       A[i,j+1] + A[i+1,j]);
24.     mydiff += abs(A[i,j] - temp);
25.   endfor
26.   LOCK(diff_lock); /*update global diff if necessary*/
27.   diff = mydiff;
28.   UNLOCK(diff_lock);
29.   BARRIER(bar1, nprocs); /*ensure all reach here before checking if done*/
30.   if (diff/(n*n) < TOL) then done = 1; /*check convergence; all get
31.     same answer*/
32. endwhile
33. end procedure

```

1/29/99

CS258 S99

30

## SAS Program

- SPMD: not lockstep. Not necessarily same instructions
- Assignment controlled by values of variables used as loop bounds
  - unique pid per process, used to control assignment
- done condition evaluated redundantly by all
- Code that does the update identical to sequential program
  - each process has private mydiff variable
- Most interesting special operations are for synchronization
  - accumulations into shared diff have to be mutually exclusive
  - why the need for all the barriers?
- Good global reduction?
  - Utility of this parallel accumulate???

1/29/99

CS258 S99

31

## Mutual Exclusion

- Why is it needed?
- Provided by LOCK-UNLOCK around *critical section*
  - Set of operations we want to execute atomically
  - Implementation of LOCK/UNLOCK must guarantee mutual excl.
- Serialization?
  - Contention?
  - Non-local accesses in critical section?
  - use private mydiff for partial accumulation!

1/29/99

CS258 S99

32

## Global Event Synchronization

- BARRIER(nprocs): wait here till nprocs processes get here
  - Built using lower level primitives
  - Global sum example: wait for all to accumulate before using sum
  - Often used to separate phases of computation
- Process P 1      Process P 2      Process P nprocs
- set up eqn system      set up eqn system      set up eqn system
- Barrier (name, nprocs)      Barrier (name, nprocs)      Barrier (name, nprocs)
- solve eqn system      solve eqn system      solve eqn system
- Barrier (name, nprocs)      Barrier (name, nprocs)      Barrier (name, nprocs)
- apply results      apply results      apply results
- Barrier (name, nprocs)      Barrier (name, nprocs)      Barrier (name, nprocs)
- Conservative form of preserving dependences, but easy to use
- WAIT\_FOR\_END (nprocs-1)

1/29/99

CS258 S99

33

## Pt-to-pt Event Synch (Not Used Here)

- One process notifies another of an event so it can proceed
  - Common example: producer-consumer (bounded buffer)
  - Concurrent programming on uniprocessor: semaphores
  - Shared address space parallel programs: semaphores, or use ordinary variables as flags

```

                P1                P2
a: while (flag is 0) do nothing;   A = 1;
print A;                          b: flag = 1;
```

- Busy-waiting or spinning

1/29/99

CS258 S99

34

## Group Event Synchronization

- Subset of processes involved
  - Can use flags or barriers (involving only the subset)
  - Concept of producers and consumers
- Major types:
  - Single-producer, multiple-consumer
  - Multiple-producer, single-consumer
  - Multiple-producer, single-consumer

1/29/99

CS258 S99

35

## Message Passing Grid Solver

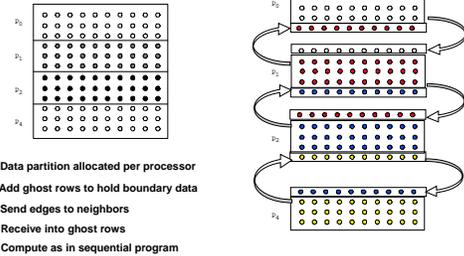
- Cannot declare A to be global shared array
  - compose it logically from per-process private arrays
  - usually allocated in accordance with the assignment of work
    - » process assigned a set of rows allocates them locally
- Transfers of entire rows between traversals
- Structurally similar to SPMD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

1/29/99

CS258 S99

36

## Data Layout and Orchestration



1/29/99 CS258 S99 37

```

10. procedure Solve()
11. begin
13.   int i, j, pid, n' = n/nprocs, done = 0;
14.   float temp, mydiff, mydiff = 0; /*private variables*/
15.   myA = malloc(a 2-d array of size [nprocs + 2] by n+2); /*initialize my rows of A, in an unspecified way*/
16.   mydiff = 0; /*local diff to 0*/
17.   while (!done) do
18.     /*Exchange border rows of neighbors into myA[N] and myA[n'+1,N]*/
19.     if (pid != 0) then SEND(myA[1,0], n*sizeof(float), pid-1, 0);
20.     if (pid == nprocs-1) then SEND(myA[n',0], n*sizeof(float), pid-1, 0);
21.     if (pid != 0) then RECEIVE(myA[0,0], n*sizeof(float), pid-1, 0);
22.     if (pid == nprocs-1) then RECEIVE(myA[n'+1,0], n*sizeof(float), pid-1, 0);
23.     for i = 1 to n' do /*for each of my (nonghost) rows*/
24.       for j = 1 to n do /*for all nonborder elements in that row*/
25.         myA[i,j] = 0.2 * myA[i,j] + myA[i,j-1] + myA[i-1,j] +
26.           myA[i,j+1] + myA[i+1,j];
27.         mydiff += abs(myA[i,j] - temp);
28.       endfor
29.     endfor
30.     /*communicate local diff values and determine if
31.     done can be replaced by reduction and broadcast*/
32.     if (pid == 0) then SEND(mydiff, sizeof(float), 0, 0);
33.     RECEIVE(done, sizeof(int), 0, DONE);
34.     /*pid does this*/
35.     for i = 1 to nprocs-1 do /*for each other process*/
36.       RECEIVE(tempdiff, sizeof(float), i, 0);
37.       mydiff += tempdiff; /*accumulate into total*/
38.     endfor
39.     if (mydiff/n/n' < TOL) then done = 1;
40.     for i = 1 to nprocs-1 do /*for each other process*/
41.       SEND(done, sizeof(int), i, DONE);
42.     endfor
43.   endwhile
44. end procedure
1/29/99 CS258 S99 38

```

## 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
  - Could implement locks and barriers with messages
- Can use REDUCE and BROADCAST library calls to simplify code
 

```

25k. REDUCE(0, mydiff, sizeof(float), ADD);
25l. if (pid == 0) then
25m.   if (mydiff/(n*n) < TOL) then done = 1;
25n.   endif
25o.   BROADCAST(0, done, sizeof(int), DONE)

```

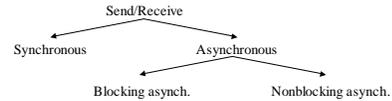
1/29/99 CS258 S99 39

## Send and Receive Alternatives

Can extend functionality: stride, scatter-gather, groups

Semantic flavors: based on when control is returned

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 needed with asynch. messages
- With synch. messages, our code is deadlocked. Fix?

## Orchestration: Summary

- Shared address space
  - Shared and private data explicitly separate
  - Communication implicit in access patterns
  - No correctness need for data distribution
  - 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 (at least in synch. case)
    - » mutual exclusion by fiat

1/29/99 CS258 S99 41

## 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?

1/29/99 CS258 S99 42