

---

---

# Parallelization Principles

---

---

R. Govindarajan  
SERC, IISc  
govind@iisc.ac.in

# Overview

---

- Introduction
- Parallelization Steps
- Example
  - Shared Address Space
  - Distributed Address Space

## Acknowledgments:

Slides for this tutorial are taken from presentation materials available with the book “Parallel Computing Architecture: A Hardware/Software Approach” (Culler, Singh and Gupta, Morgan Kaufmann Pub.) and the associated course material. They have been suitably adapted.

# Parallel Programming

---

- Shared, global, address space, hence called *Shared Address Space*
  - Any processor can *directly* reference any memory location
  - Communication occurs implicitly as result of loads and stores
- *Message Passing Architecture*
  - Memory is private to each node
  - Processes communicate by messages

# Definitions

---

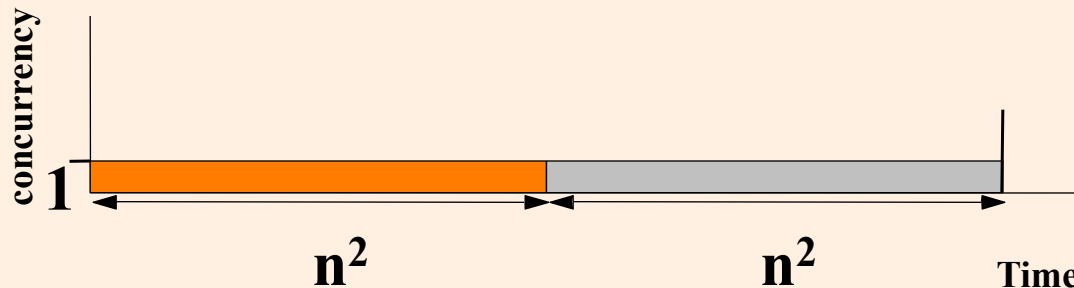
- $\text{Speedup} = \frac{\text{Exec. Time in UniProcessor}}{\text{Exec. Time in } n \text{ processors}}$
- $\text{Efficiency} = \frac{\text{Speedup}}{n}$
- Amdahl's Law:
  - For a program with  $s$  part sequential execution, speedup is limited by  $1/s$ .

# Understanding Amdahl's Law

---

## Example: 2-phase calculation

- sweep over  $n \times n$  grid and do some independent computation
- sweep again and add each value to global sum



(a) Serial

- Serial Execution Time =  $n^2 + n^2 = 2n^2$

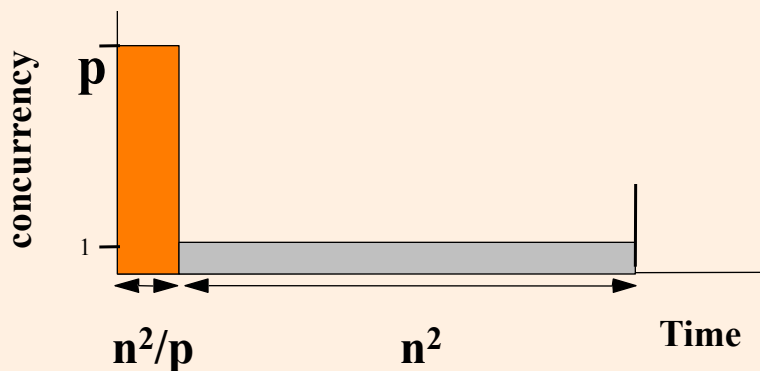
# Understanding Amdahl's Law

## Parallel Execution time:

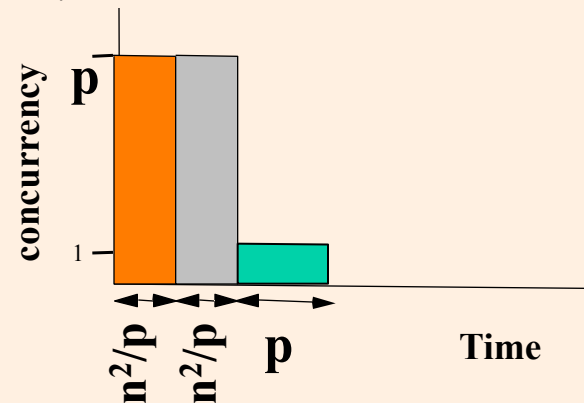
- Time for first phase =  $n^2/p$
- Second phase serialized at global variable =  $n^2$ ;
- Speedup =  $(2n^2/(n^2 + n^2/p))$  or at most 2

## Improved Parallel Execution

- Localize the sum in  $p$  procs and then do serial sum.
- Speedup =  $(2n^2/(2n^2/p + p)) \approx p$



(b) Naïve Parallel



(c) Parallel

# Definitions

---

---

## ■ Task

- Arbitrary *piece of work* in parallel computation
- Executed sequentially; concurrency is only across tasks
- Fine-grained vs. coarse-grained tasks

## ■ Process (thread)

- Abstract entity that performs the tasks
- Communicate and synchronize to perform the tasks

## ■ Processor

- Physical engine on which process executes

# Tasks involved in Parallelization

---

- 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



# Parallelizing Computation vs. Data

---

- Computation is decomposed and assigned (partitioned) - *task decomposition*
  - Task graphs, synchronization among tasks
- Partitioning Data is often a natural view too - *data or domain decomposition*
  - Computation follows data: *owner computes*
  - Grid example; data mining;

# Domain Decomposition: Example

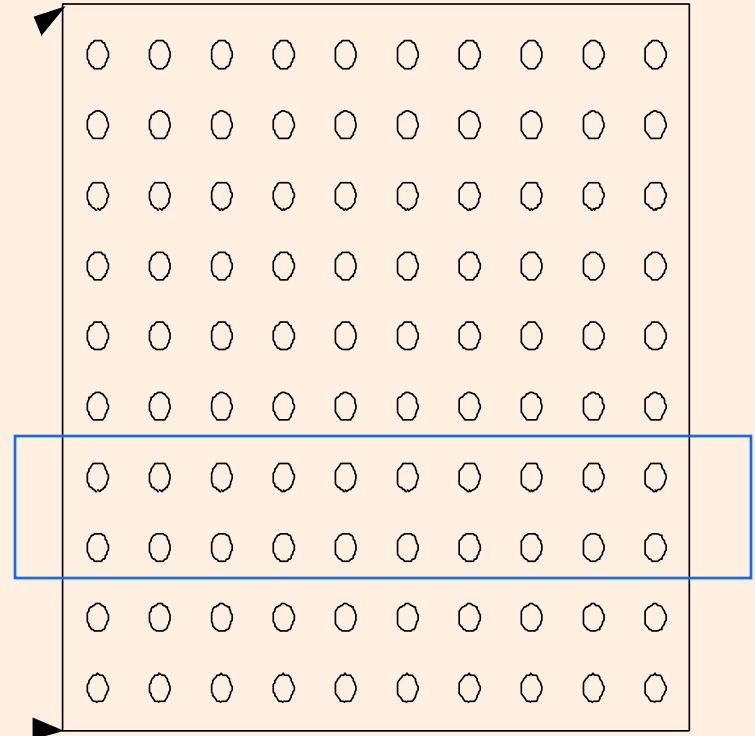
---

- Some computation performed on all elts. of the array

for  $i=1$  to  $m$

  for  $j= 1$  to  $n$

$$a[i,j] = a[i,j] + v[i]$$

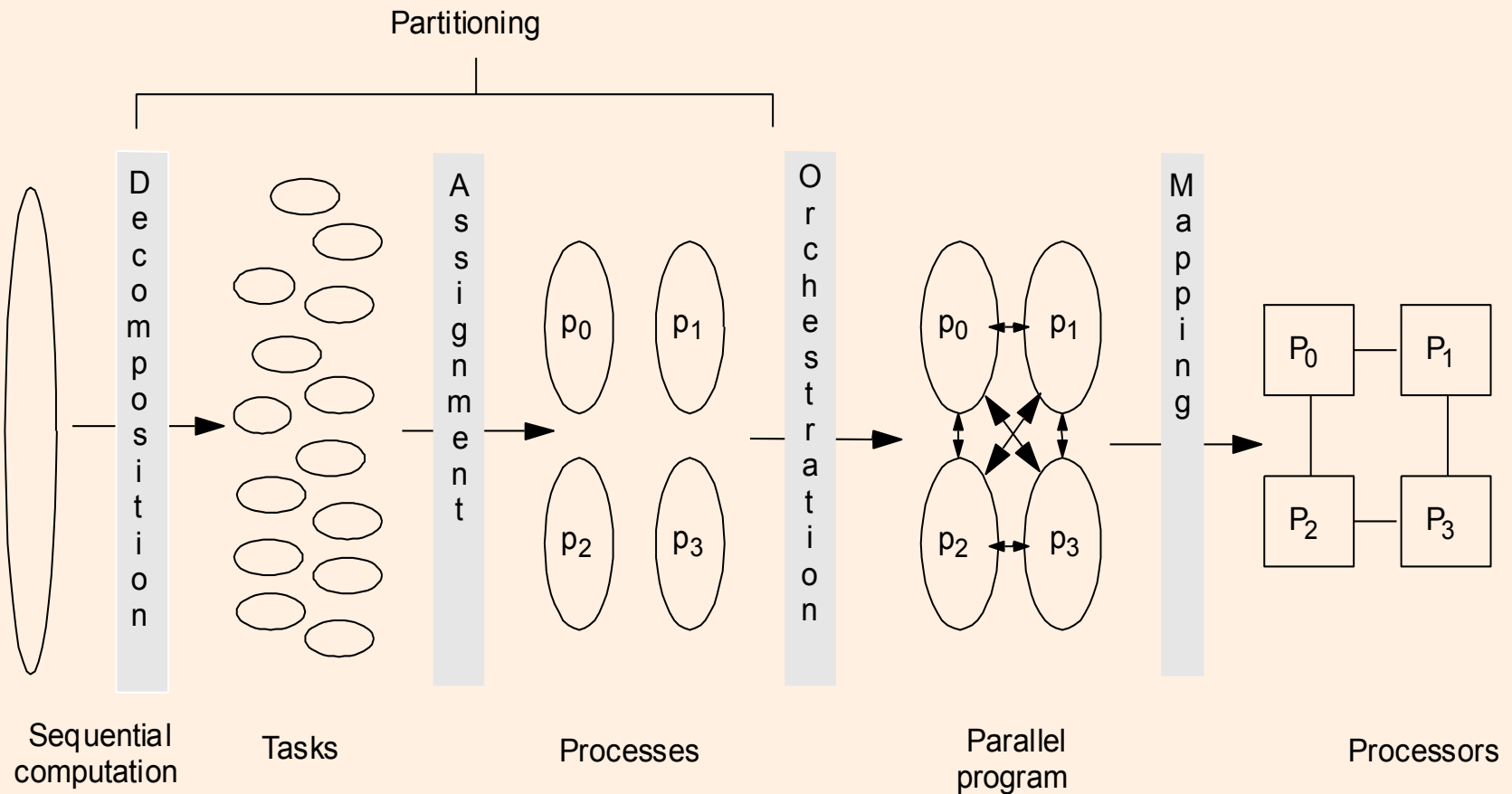


# Steps in Creating a Parallel Program

---

- Decomposition of computation into tasks
- Assignment of tasks to processes
- Orchestration of data access, communication, and synchronization.
- Mapping processes to processors

# Steps in Creating a Parallel Program



# Decomposition

---

- Identify concurrency
- Break up computation into tasks to be divided among processes
  - Tasks may become available dynamically
  - No. of available tasks may vary with time
- Goal: Expose available parallelism → enough tasks to keep all processors busy

# Assignment

---

- Specifies how to group tasks together for a process
  - Balance workload, reduce communication and management cost
- Structured approaches usually work well
  - Code inspection (parallel loops) or understanding of application
  - *Static versus dynamic assignment*
- Both decomposition and assignment are *usually* independent of architecture or programming model
  - But cost and complexity of using primitives may affect decisions

# Orchestration

---

- **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 Programming Model, Communication abstraction, and efficiency of primitives
- Architecture should provide appropriate primitives efficiently

# Mapping

---

- Two aspects:
  - Which process runs on which particular processor?
  - 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 some aspects, system handles some



# 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 Reduce communication volume
Orchestration	Yes	Reduce noninherent communication via data locality Reduce communication and synchronization cost as seen by the processor Reduce serialization at shared resources Schedule tasks to satisfy dependences early
Mapping	Yes	Put related processes on the same processor if necessary Exploit locality in network topology

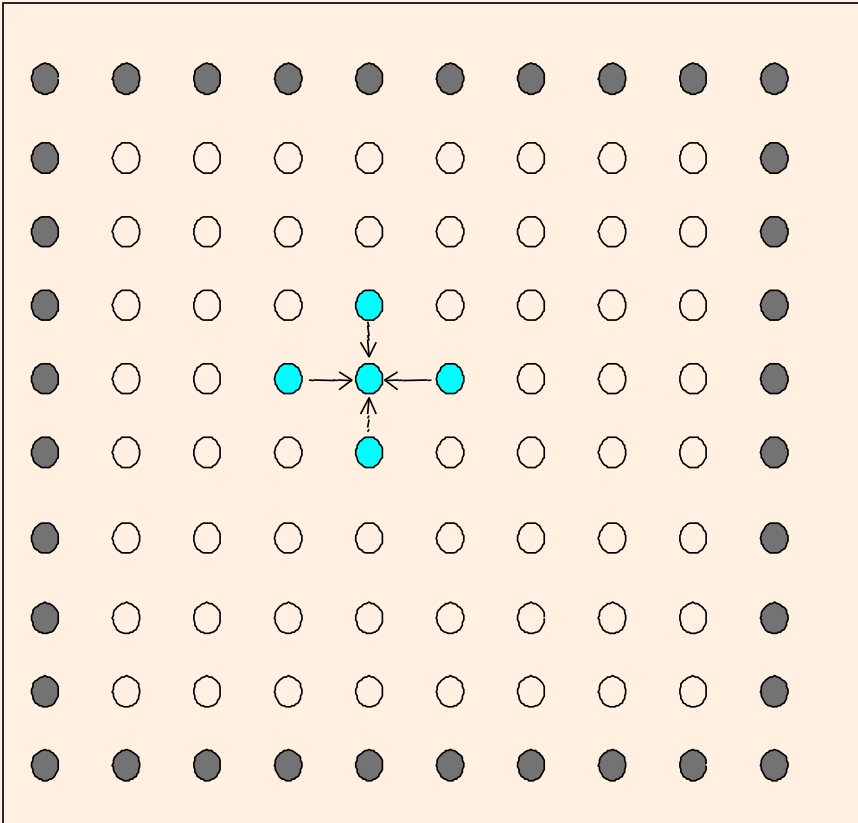
# Example: Grid Solver

---

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

# Grid solver (Simple Version)

---



```
for i = 1 to n
  for j = 1 to n
  {
    B[i,j] = 0.2 * (A[i,j] +
      A[i-1,j] + A[i+1,j] +
      A[i,j-1] + A[i,j+1]);
    diff += abs(B[i,j] - A[i,j]);
  }
for i = 1 to n
  for j = 1 to n
    A[i,j] = B[i,j];
```

# Sequential Version

---

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

# Sequential Version (contd.)

---

```
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 non-border points of grid*/
18.             for j ← 1 to n do
19.                 B[i,j] ← 0.2 * (A[i,j] + A[i,j-1] + A[i-1,j] +
20.                    A[i,j+1] + A[i+1,j]); /*compute average*/
21.                 diff += abs(B[i,j] - A[i,j]);
22.             end for
23.         end for
24.         if (diff/(n*n) < TOL) then done = 1;
25.         else Copy_Array (A ← B)
26.     end while
27. end procedure
```

# Decomposition & Assignment

---

```
for i = 1 to n
  for j = 1 to n
  {
    B[i,j] = 0.2 * (A[i,j] +
      A[i-1,j] + A[i+1,j] +
      A[i,j-1] + A[i,j+1]);
    diff += abs(B[i,j] - A[i,j]);
  }
```

```
for i = 1 to n
  for j = 1 to n
    A[i,j] = B[i,j];
```

## ■ Decomposition

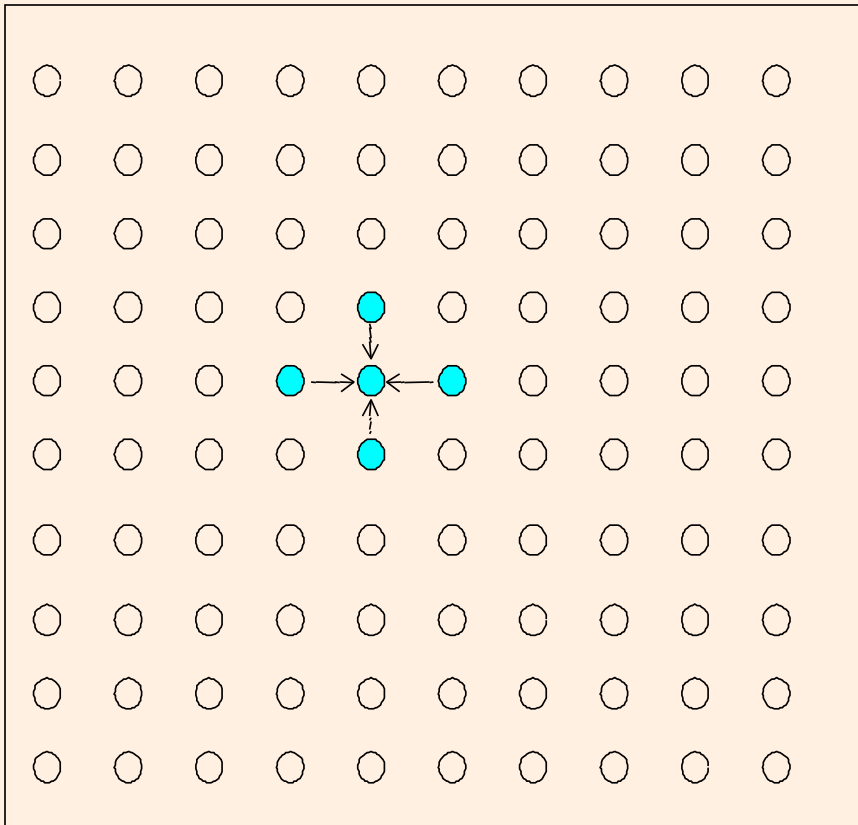
- Both i and j loops can be parallelized - no data dependences
- Each grid point can be a task
- To compute **diff**, some coordination would be required!

## ■ Assignment

- Each grid point
- Each row or column
- A group of rows or columns

# Grid solver (Update-in-place Version)

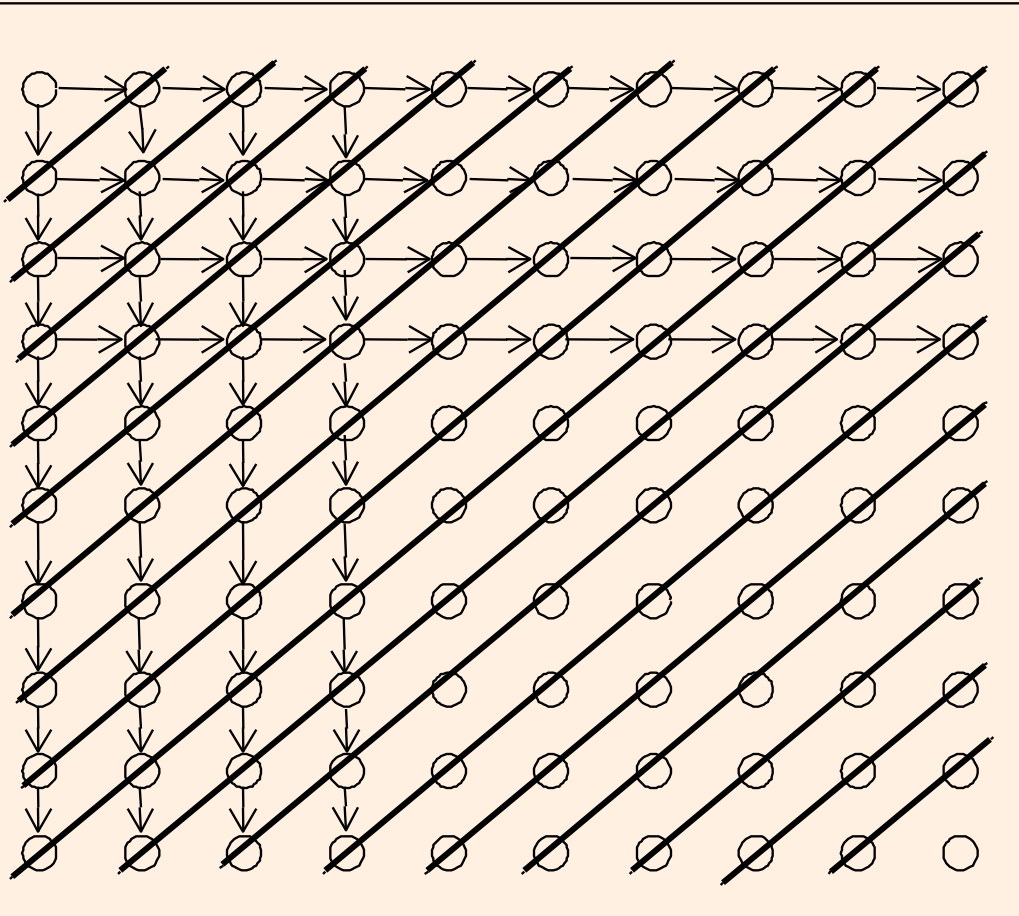
---



```
for i = 1 to n
  for j = 1 to n
  {
    temp = A[i,j];
    A[i,j] = 0.2 * (A[i,j] +
      A[i-1,j] + A[i+1,j] +
      A[i,j-1] + A[i,j+1]);
    diff += abs(temp - A[i,j]);
  }
```

# Decomposition & Assignment

---



## ■ Decomposition

- Dependence on both  $i$  and  $j$  loops
- Each grid point can be a task
- Need point-to-point synchronization -- Very expensive

## ■ Assignment

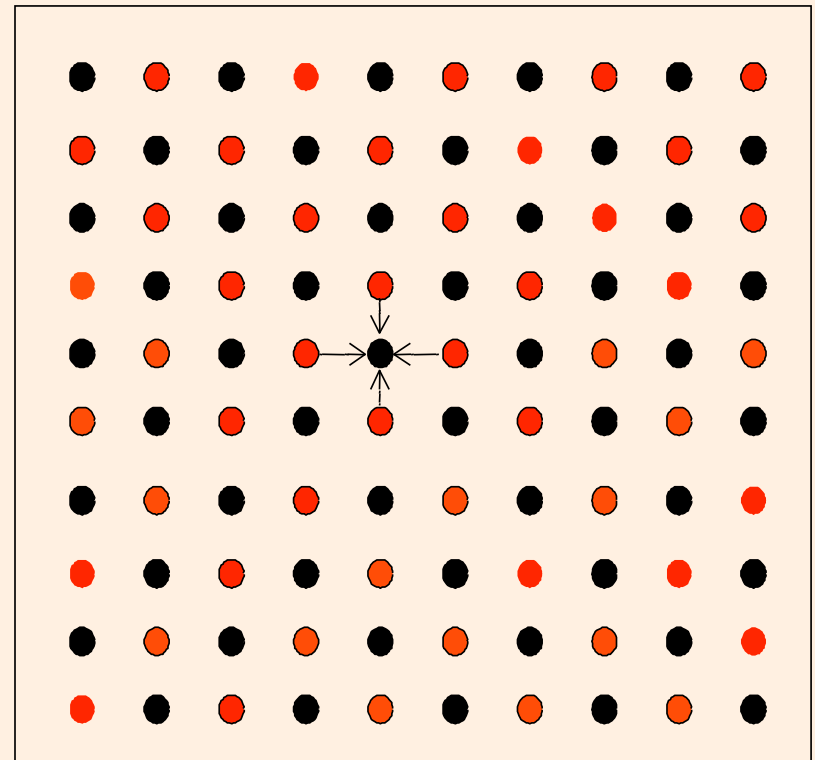
- Grid points along diagonal form a task
- Restructure loop and global synchronization
- Load imbalance



# Exploiting Application Knowledge

---

- Reorder grid traversal: red-black ordering
- Red sweep and black sweep are each fully parallel:
- Global synch between them (conservative but convenient)
- Different ordering of updates: may converge slower



# Red-Black Parallel Version

---

```
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.         forall i ← 1 to n step 2 do/*sweep black points of grid*/
18.             forall j ← 2 to n+1 step 2 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 forall
24.         end forall
24a        /* similarly forall loop for red points */
25.         if (diff/(n*n) < TOL) then done = 1;
26.     end while
27. end procedure
```

Ensure  
computation for  
all black points  
are complete!

# Red-Black Parallel Version (contd.)

---

- Decomposition into elements: degree of concurrency  $n^2/2$ ; 2 global synchronizations per  $n^2$  computation
- **forall** loop to express the parallelism.
- Too fine-grain parallelism  $\Rightarrow$  group tasks to form a process.
- Decompose into rows? Computation vs. communication overhead?

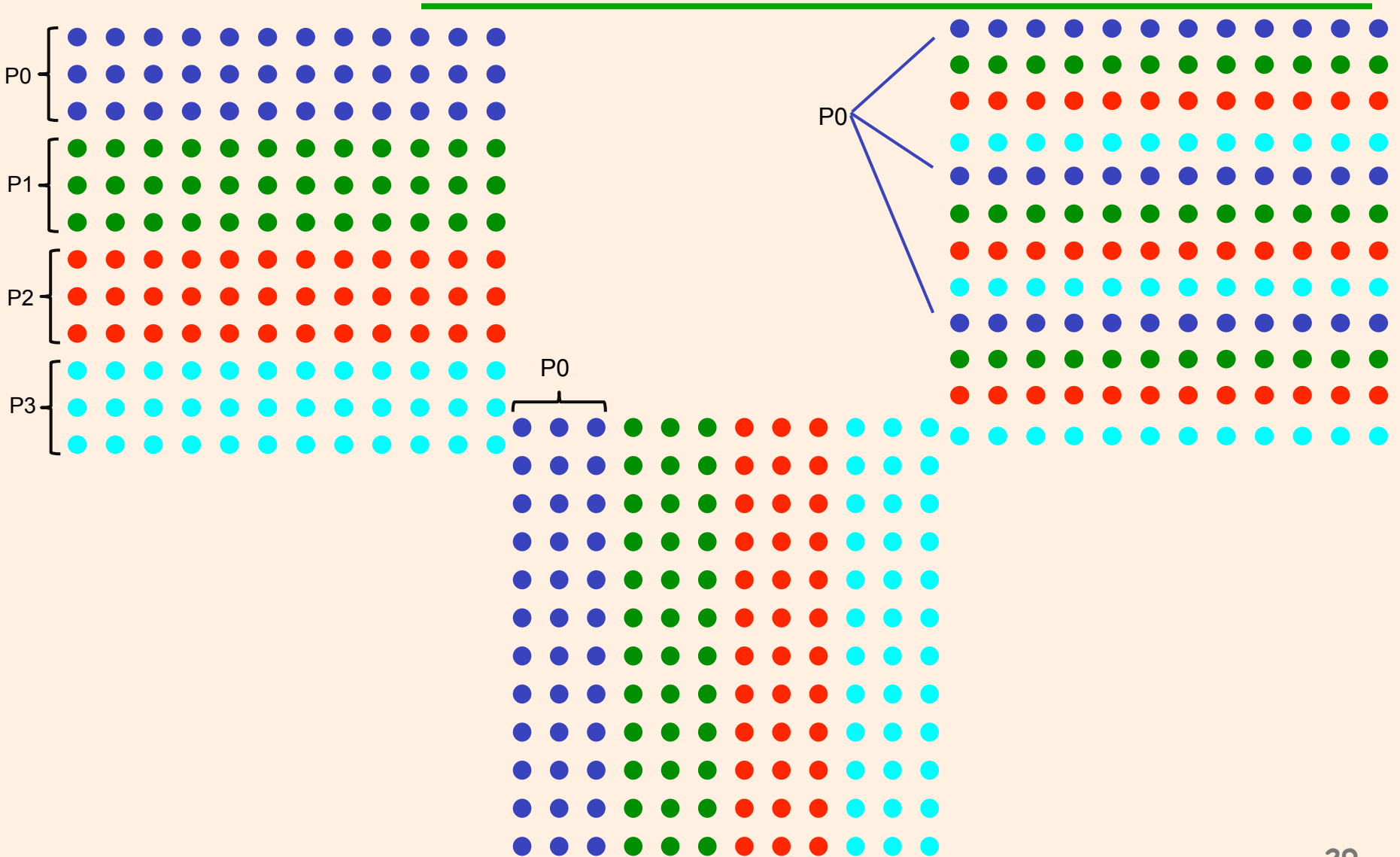
# Assignment

---

- **Static assignment: decomposition into rows**
  - **Block** assignment of rows: Rows  $i*(n/p), \dots, (i+1)*(n/p) - 1$  are assigned to process  $i$
  - **Cyclic** assignment of rows: process  $i$  is assigned rows  $i, i+p, i+2p \dots$
- **Dynamic assignment**
  - get a row index, work on the row, get a new row, ...
- **Concurrency? Volume of Communication?**

# Assignment (contd.)

---



# Orchestration

---

- Different for different programming models/architectures
  - Shared address space
    - Naming: global addr. Space
    - Synch. through barriers and locks
  - Distributed Memory / Message passing
    - Non-shared address space
    - Send-receive messages + barrier for synch.

# Shared Memory Version

---

```
1.  int n, nprocs;      /* matrix: (n + 2-by-n + 2) elts.*/
2.  float **A, diff = 0;
2a. LockDec (diff_lock);
2b. BarrierDec (barrier1);
3.  main()
4.  begin
5.      read(n) ;      /*read input parameter: matrix size*/
5a.      Read (nprocs);
6.      A ← g_malloc (a 2-d array of (n+2) × (n+2) doubles);
6a.      Create (nprocs -1, Solve, A);
7.      initialize(A);    /*initialize the matrix A somehow*/
8.      Solve (A);        /*call the routine to solve equation*/
8a.      Wait_for_End (nprocs-1);
9.  end main
```

# Shared Memory Version

```
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, pid, done = 0;
14.     float mydiff, temp;
14a.         mybegin = 1 + (n/nprocs)*pid;
14b.         myend = mybegin + (n/nprocs);
15.     while (!done) do /*outermost loop over sweeps*/
16.         mydiff = diff = 0; /*initialize local difference to 0*/
16a.         Barrier (barrier1, nprocs);
17.         for i ← mybeg to myend do/*sweep
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 new value
22.                 mydiff += abs(A[i,j] - temp);
23.             end for
24.         end for
24a.         Reduce (mydif, diff);
24b.
24c.
24d.         barrier (barrier1, nprocs);
25.         if (diff/(n*n) < TOL) then done = 1;
26.     end while
27. end procedure
```

- No red-black, simply ignore dependences within sweep
- Simpler asynchronous version, may take longer to converge!

Why do we need this barrier?

Reduce (mydif, diff):

Why do we need this barrier?



# Shared Memory Program : Remarks

---

- done condition evaluated redundantly by all
- Each process has private mydiff variable
- Most interesting special operations are for synchronization provided by LOCK-UNLOCK around *criticalsection*
  - Set of operations we want to execute atomically
  - accumulations into shared **diff** have to be mutually exclusive
- Good global reduction?

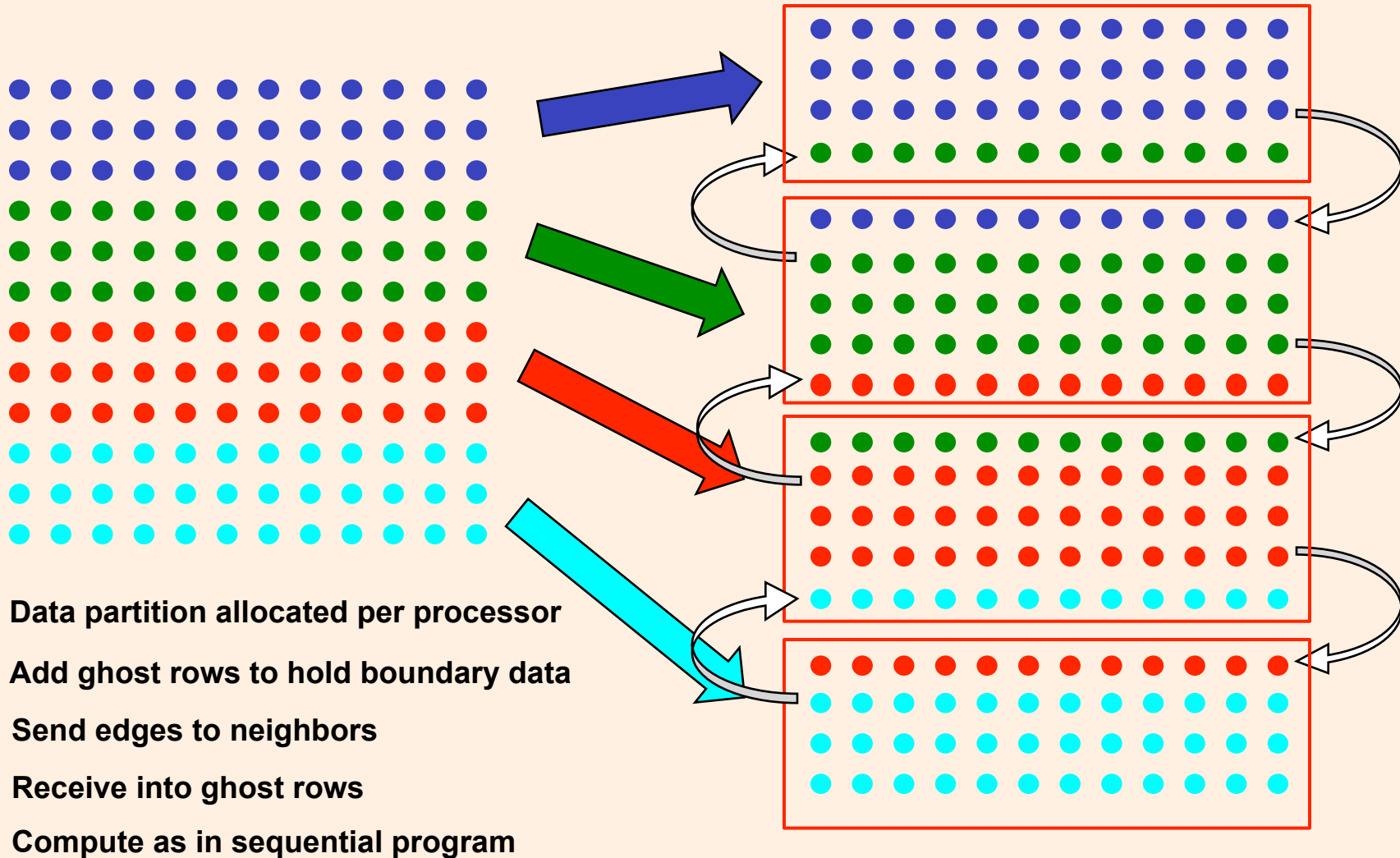
# Message Passing Version

---

- Cannot declare  $A$  to be global shared array
  - compose it from per-process private arrays
  - usually allocated in accordance with the assignment of work -- owner-compute rule
    - process assigned a set of rows allocates them locally
- Structurally similar to SPMD Shared Memory Version
- Orchestration different
  - data structures and data access/naming
  - communication
  - synchronization
- Ghost rows

# Data Layout and Orchestration

---



# Message Passing Version

---

```
1.  int n, nprocs;      /* matrix: (n + 2-by-n + 2) elts.*/
2.  float **myA;
3.  main()
4.  begin
5.      read(n) ;      /*read input parameter: matrix size*/
5a.     read (nprocs);
/* 6.  A ← g_malloc (a 2-d array of (n+2) × (n+2) doubles); */
6a.     Create (nprocs -1, Solve, A);
/* 7.  initialize(A);      */ /*initialize the matrix A somehow*/
8.      Solve (A);      /*call the routine to solve equation*/
8a.     Wait_for_End (nprocs-1);
9.  end main
```

# Message Passing Version

---

```
10. procedure Solve (A) /*solve the equation system*/
11.     float A[n+2][n+2]; /*A is an (n + 2)-by-(n + 2) array*/
12. begin
13.     int i, j, pid, done = 0;
14.     float mydiff, temp;
14a.     myend = (n/nprocs) ;
14b.     myA = malloc (array of ((n/nprocs)+2) x (n+2) floats );
14c.     If (pid == 0)
           Initialize (A)
14d.     GetMyArray (A, myA); /* get n x (n+2) elts. from process 0 */
15.     while (!done) { /*outermost loop over sweeps*/
16.         mydiff = 0; /*initialize local difference to 0*/
16a.         if (pid != 0) then
           SEND (&myA[1,0] , n*sizeof(float), (pid-1), row);
16b.         if (pid != nprocs-1) then
           SEND (&myA[myend,0], n*sizeof(float), (pid+1), row);
16c.         if (pid != 0) then
           RECEIVE (&myA[0,0], n*sizeof(float), (pid - 1), row);
16d.         if (pid != nprocs-1) then
           RECEIVE (&myA[myend+1,0], n*sizeof(float), (pid - 1), row);
16e.         ...     ...     ...
```

# Message Passing Version - Solver

---

12.begin

```
15.     while (!done) do /*outermost loop over sweeps*/
17.         for i ← 1 to myend do/*sweep for all points of grid*/
18.             for j ← 1 to n do
19.                 temp = myA[i,j]; /*save old value of element*/
20.                 myA[i,j] ← 0.2 * (myA[i,j] + myA[i,j-1] +myA[i-1,j] +
21.                     myA[i,j+1] + myA[i+1,j]); /*compute average*/
22.                 mydiff += abs(myA[i,j] - temp);
23.             end for
24.         end for
24a        if (pid != 0) then
24b.            SEND (mydif, sizeof(float), 0, DIFF);
24c.            RECEIVE (done, sizeof(int), 0, DONE);
24d.        else
24e.            for k ← 1 to nprocs-1 do
24f.                RECEIVE (tempdiff, sizeof(float), k , DIFF);
24g.                mydiff += tempdiff;
24h.            Endfor
24i.            if (diff/(n*n) < TOL) then done = 1;
24j.            for k ← 1 to nprocs-1 do
24k.                SEND (done, sizeof(float), k , DONE);
26.     end while
27.end procedure
```

# Message Passing Version : Remarks

---

- Communication in whole rows, not element at a time
- Code 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
- Communication done at beginning of iteration, synchronization only between neighboring processes

# Orchestration: Summary

---

## ■ Shared address space

- Shared and private data explicitly separate
- Communication implicit in access patterns
- 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)