Parallel Computations & Applications

National Tsing Hua University 2018, Fall Semester



Outline

- Embarrassingly Computations
- Load Balancing & Termination
- Divide-And-Conquer Computations
- Pipelined Computations
- Synchronous Computations



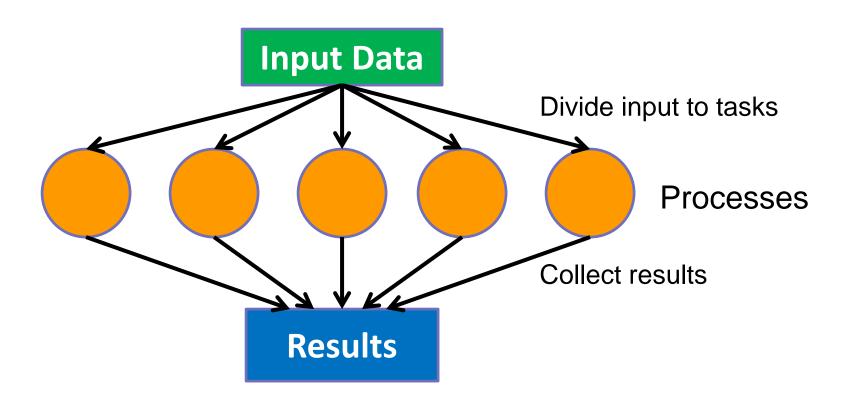
Outline

- Embarrassingly Computations
 - Image Transformations
 - Monte Carlo Methods
 - Mandelbrot Set
- Load Balancing & Termination
- Divide-And-Conquer Computations
- Pipelined Computations
- Synchronous Computations



What is Embarrassingly Parallel

A computation that can be divided into a number of completely independent tasks



Example 1: Image Transformations

- Low-level image operations:
 - Shifting: object shifted by Δx in the x-dimension and Δy in the y-dimension:

$$x' = x + \Delta x$$
, $y' = y + \Delta y$



 \triangleright Scaling: object scaled by a factor of Sx in the x-direction and Sy in the y-direction;

$$x' = xS_x$$
, $y' = yS_y$





 \triangleright Rotation: object rotated through the angle θ about the origin of the coordinate system:

$$x' = x \cos \theta + y \sin \theta$$

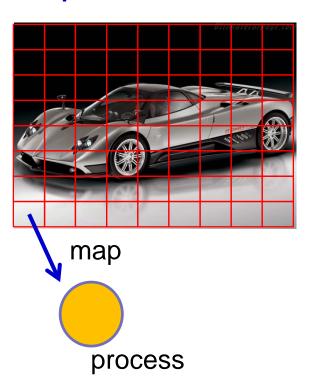
$$y' = -x \sin \theta + y \cos \theta$$

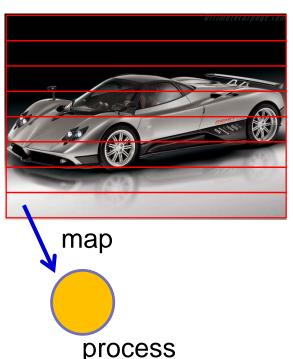


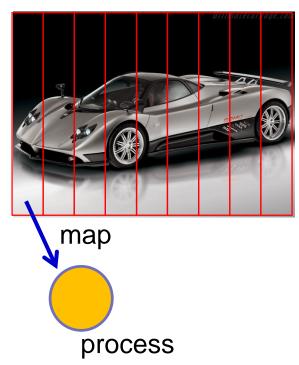




- partition
- partition
- Square region Row region Column region partition







Pseudo-code for Image Shift

Partition region by ROW with width 10 480

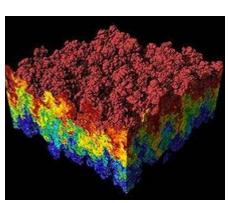
```
//master process
for(i=0, row=0; i<48; i++, row+=10) // for each of 48 processes
    send(row, P<sub>i</sub>);
                                    // send row no.
for(i=0; i<480; i++) for(j=0; j<640; j++) temp_map[i][j] = 0; // initialize temp
for(i=0; i<(480*640); i++) {
                                                          // for each pixel
                                                          // accept new coordinates
    recv(oldrow, oldcol, newrow, newcol, P<sub>ANY</sub>);
    if !((newrow<0)||((newrow>=480)||(newcol<0)||((newcol>=640))
         temp map[newrow][newcol] = map[oldrow][oldcol];
for(i=0; i<480; i++) for(j=0; j<640; j++) map[i][j] = temp_map[i][j]; // update map
// slave process
recv (row, Pmaster);
for (oldrow = row; oldrow < (row+10); oldrow++) // for each row in the partition
    for (oldcol = 0; oldcol < 640; oldcol++) {
                                                        // for each column in the row
         newrow = oldrow + delta x;
                                                        // shift along x-dimension
                                                        // shift along y-dimension
         newcol = oldcol + delta_y;
         send(oldrow, oldcol, newrow, newcol, Pmaster); // send out new coordinates
```

640

10 x 640

Example 2: Monte Carlo Methods

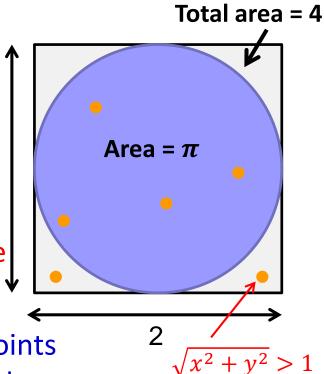
- Monte Carlo methods: a class of computational algorithms that rely on repeated random sampling to compute their results
 - Invented in 1940s by John von Neumann, Stanislaw Ulam and Nicholas Metropolis, while they were working on nuclear weapon (Manhattan Project)
 - Especially useful for simulating systems with many coupled degrees of freedom, such as fluids, disordered material



HISTORY

Monte Carlo Methods --- π calculation

- How to compute π ???
 - \triangleright Definition of π : the area of a circle with unit radius
 - ightharpoonup We know: $\frac{\text{Area of circle}}{\text{Area of square}} = \frac{\pi}{4}$
 - Randomly choose points from the square
 - \triangleright Giving sufficient number of samples, the fraction of points **within** the circle will be $\pi/4!!!$
 - ➤ E.g.: With 10,000 randomly sample points we expect 7854 points within the circle
 - \rightarrow 7854/10000 = π /4 \rightarrow π = 7854/10000*4 = 3.1416



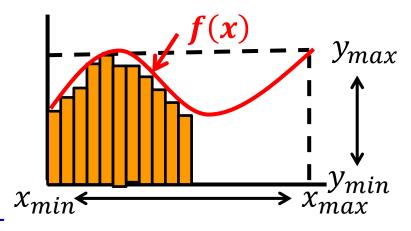
Monte Carlo Methods --- Integral

- Monte Carlo Method can compute ANY **definite** integral!
 - max and min values of the integral must be known
 - Very inefficient....

■ Method:

- \triangleright Randomly choose point (x, y):
 - $\star x_{max} \le x \le x_{min}$
 - $y_{max} \le y \le y_{min}$
- Compute the area (integral) according to the ratio of points inside and outside the area
 - \rightarrow just like the computation of π

$$Area = \int_{x_{min}}^{x_{max}} f(x) dx$$



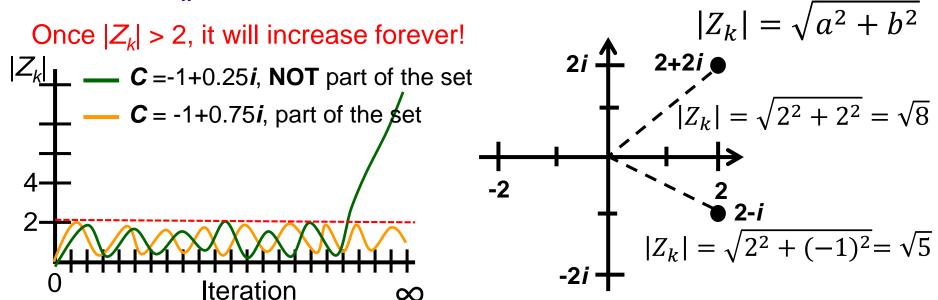
 \triangleright Given any point (x, y), outside means : y > f(x)



■ The Mandelbrot Set is a set of **complex numbers** that are quasi-stable when computed by iterating the function:

$$Z_0 = C$$
, $Z_{k+1} = Z_k^2 + C$

- \triangleright **C** is some complex number: $\mathbf{C} = a + b\mathbf{i}$
- $ightharpoonup Z_{k+1}$ is the (k+1)th iteration of the complex number
- ightharpoonup If $|Z_k| \le 2$ for ANY $k \to C$ belongs to Mandelbrot Set

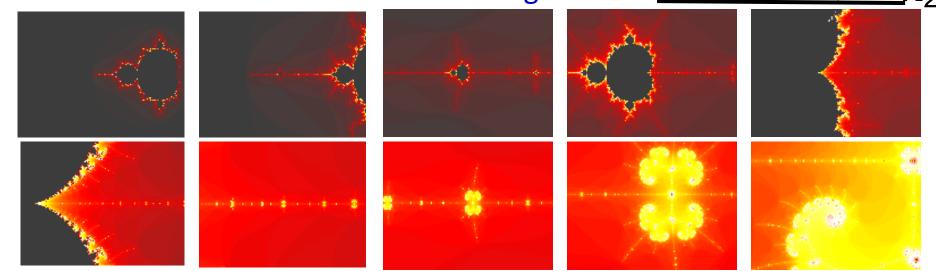


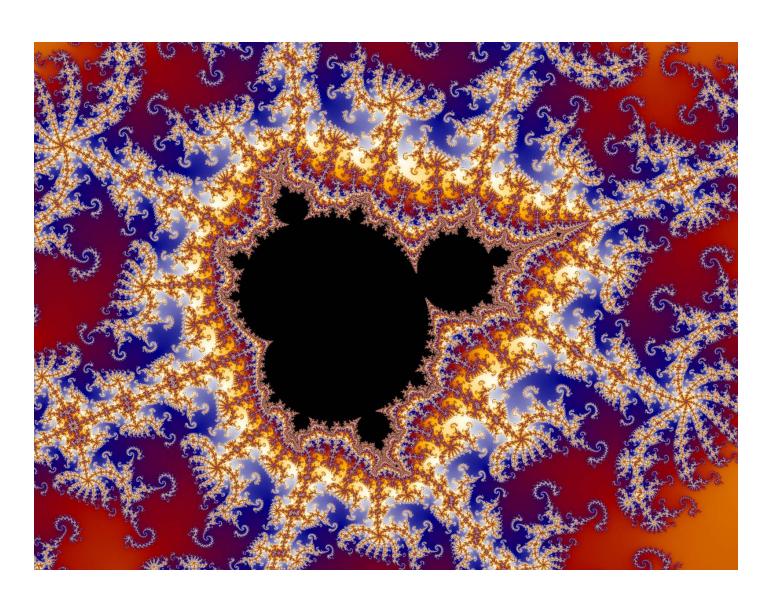
Parallel Programming - NTHU LSA Lab

11

Fractal

- What exact is Mandelbrot Set?
 - ➤ It is a **fractal**: An object that display self-similarity at various scale; Magnifying a fractal reveals small-scale details similar to the large-scale characteristics
 - ➤ After plotting the Mandelbrot Set determined by thousands of iteration:
 - Add color to the points outside the set &
 zoom in at the center of the image: -2<u>i</u>





.

Mandelbrot Set Program

■ Compute $Z_{k+1} = Z_k^2 + C$ > Let $C = C_{real} + C_{imag}i$, $Z_k = Z_{real} + Z_{imag}i$ > $Z_{k+1} = (Z_{real}^2 - Z_{imag}^2 + 2Z_{real}Z_{imag}i) + (C_{real} + Z_{real}^2)$

```
Struct complex {
    float real;
    float imag;
};
```

Sequential Mandelbrot Set Program

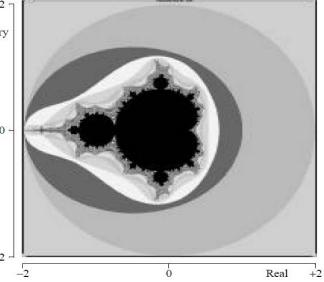
- Testing program:
 - Giving a complex number
 - \triangleright Return the iteration number when $|Z_k| > 2$
 - Let the maximum iteration is 256

```
int cal_pixel (complex c) {
  int count = 0;
                                   // number of iterations
                                   // maximum iteration is 256
  int max= 256;
  float temp, lengthsq;
                                   // initialize complex number z
  complex z;
  z.real = 0; z.imag = 0;
  do {
    temp = (z.real * z.real) - (z.imag * z.imag) + c.real; // compute next z.real
    z.imag = (2 * z.real * z.imag) + c.imag;
                                                         // compute next z.imag
    z.real = temp;
    lengthsq = (z.real * z.real) + (z.imag * z.imag);
    count++;
                                                     // update iteration counter
  } while ((lengthsq < 4.0) && (count < max));</pre>
  return count;
                                                                              15
```

Sequential Mandelbrot Set Program

Scaling Coordinate Display Program:

- Plot the Mandelbrot Set from the coordinate system
- Color indicate the iteration number black=256, white=0
- Points are apart with a fixed distance read_disk, imag_dist



```
for (x=real_min; x < real_max; x += real_dist) {
    for (y=imag_min; y < imag_max; x += imag_dist) {
        c.real = x; c.img = y;
        color = cal_pixel (c);
        display(x, y, color);
    }
}</pre>
```

Parallelizing Mandelbrot Set Program

■ Partition screen 640*480 by row using 48 processes

```
//master process
for(i=0, row=0; i<48; i++, row+=10)
    send(row, P<sub>i</sub>);

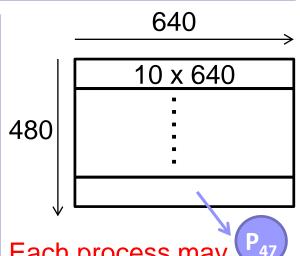
for(i=0; i<(480*640); i++) {
    recv(&x, &y, &color, P<sub>ANY</sub>);
    display(x, y, color);

}

// for each process
// send row no.

// for each pixel point
// receive coordinate/colors
// display pixel
}
```

```
//slave process
recv (&row, P<sub>master</sub>);
for (x=0; x < 640; x++) {
    for (y=row; y < (row+10); y++) {
        c.real = min_real + (x * scale_real);
        c.imag = min_imag + (y * scale_image);
        color = cal_pixel (c);
        send(x, y, &color, P<sub>master</sub>);
    }
}
```



Each process may have different load!



Outline

- Embarrassingly Computations
- Load Balancing & Termination
- Divide-And-Conquer Computations
- Pipelined Computations
- Synchronous Computations



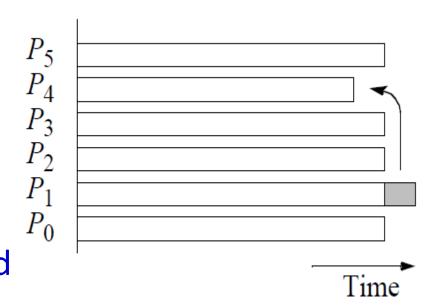
Load-Balancing & Termination

Load-balancing

Used to distribute computations fairly across processors in order to obtain the highest possible execution speed

■ Termination detection

- Detecting when a computation has been completed
- More difficult when the computation is distributed





Static Load-Balancing

- Static means ...
 - Pre-determine assignment between tasks and processes

- Round robin algorithm selects processes in turn
- Randomized algorithms selects processes randomly
- Recursive bisection recursively divides the problem into sub-problems of equal computational effort



Static Load-Balancing

- Several fundamental flaws with static load balancing even if a mathematical solution exists:
 - Very difficult to estimate accurately the execution times of various parts of a program without actually executing the parts
 - Communication delays that vary under different circumstances
 - > Some problems have an indeterminate number of steps to reach their solution (e.g. Manderbrot set)



Dynamic Load Balancing

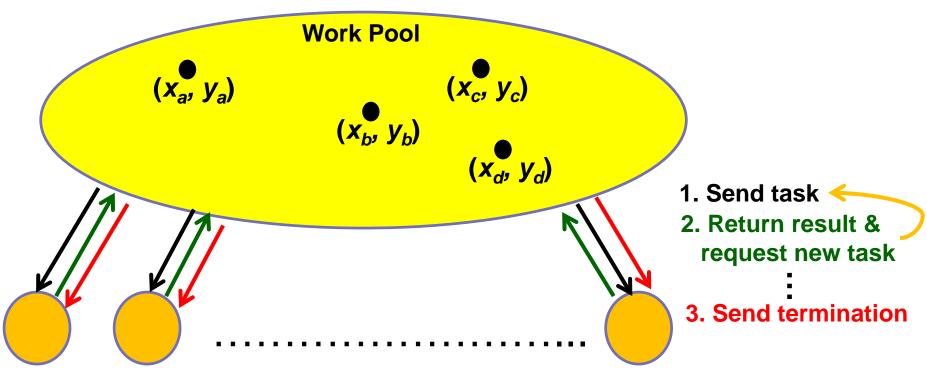
- **Dynamic** means ...
 - > Assign tasks to processes during the execution
 - Does incur an additional overhead during execution, but it is much more effective than static load balancing

Approaches

- Centralized Work Pool
- Decentralized Work Pool
- Fully Distributed Work Pool

Centralized Work Pool

- Work pool / Processor Farm
 - > Useful when tasks require different execution time
 - Dynamic load balancing



Coding for Work Pool Approach

```
//master process
                                  // # of active processes
count = 0;
                                // row being sent
row = 0;
for (i=0; i<num_proc; i++) { // send initial row to each processes
    send(row, P<sub>i</sub> , data_tag);
    count++;
    row++;
do {
    recv(&slave, &r, color, P<sub>ANY</sub>, result_tag);
    count--;
    if (row < num_row) {</pre>
                              // keep sending until no new task
        send(row, P<sub>slave</sub>, data_tag); // send next row
         count++;
                                          Tag is needed to distinguish
         row++;
                                          between data and termination msg
    } else {
        send(row, P<sub>slave</sub>, terminate_tag); // terminate
    display(r, color);
                                         // display row
} while(count > 0);
```



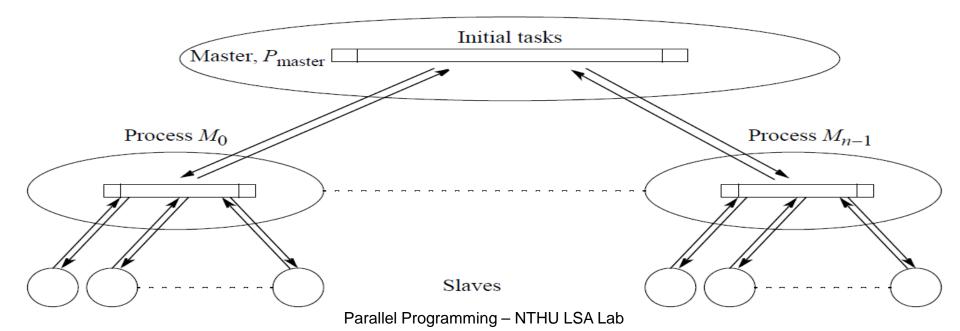
Coding for Work Pool Approach

```
//slave process P ( i )
recv(&row, P<sub>master</sub> , source_tag);
while (source_tag == data_tag) { // keep receiving new task
    c.imag = min_imag + (row * scale_image);
    for (x=0; x<640; x++) {
        c.real = min_real + (x * scale_real);
        color[x] = cal_pixel (c); // compute color of a single row
    send(i, row, color, P<sub>master</sub>, result_tag); // send process id and results
    recv(&row, P<sub>master</sub> , source_tag);
```

м.

Decentralized Work Pool

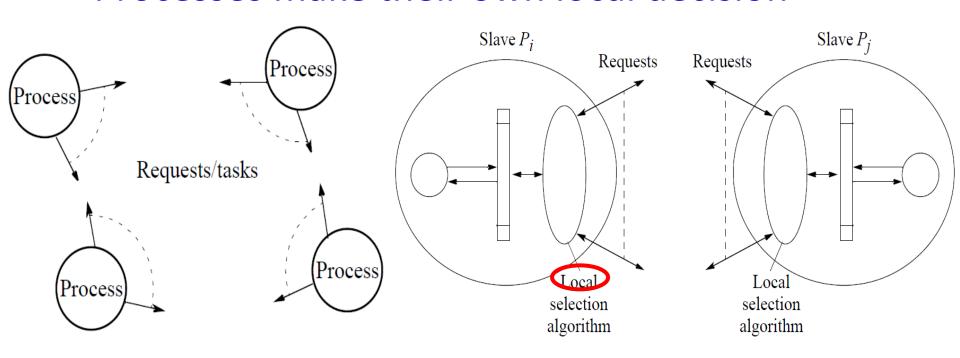
- Hierarchical structure
- The master divides the initial work pool into parts to each of the "mini-masters"
- Each mini-master controls one group of slaves





Fully Distributed Work Pool

- Processes have or generate their own tasks
- Tasks may be transferred among processes
- Processes make their own local decision





Fully Distributed Work Pool

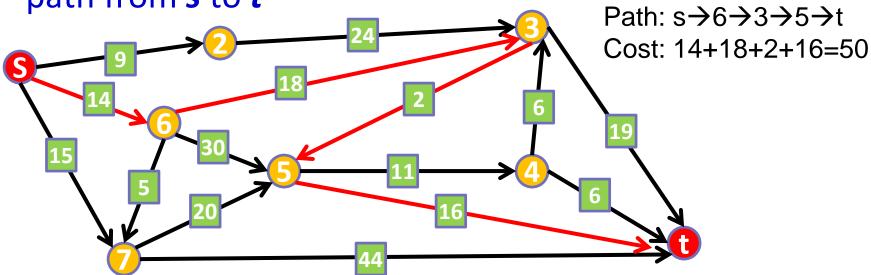
- Receiver-initiated method (Receiver has lighter load)
 - Request tasks from heavy loaded processes
 - More suitable (less request) for heavy loaded system
- Sender-initiated method (sender has higher load)
 - Send tasks to light loaded processes
 - More suitable (less request) for light loaded system
- Challenges:
 - > it could be difficult or expensive to determine process loads
 - > may NOT provide immediate load balance



Shortest Path in a Weighted Diagraph

Given a weighted graph, find the shortest directed

path from s to t

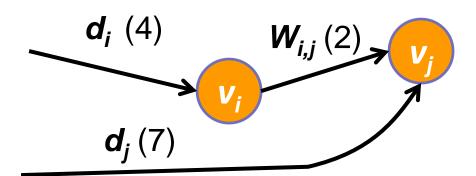


- Note: weight are arbitrary numbers
 - Not necessary distance
 - Need not satisfy the triangle inequality

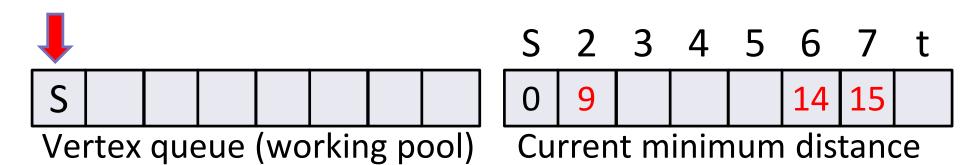
1

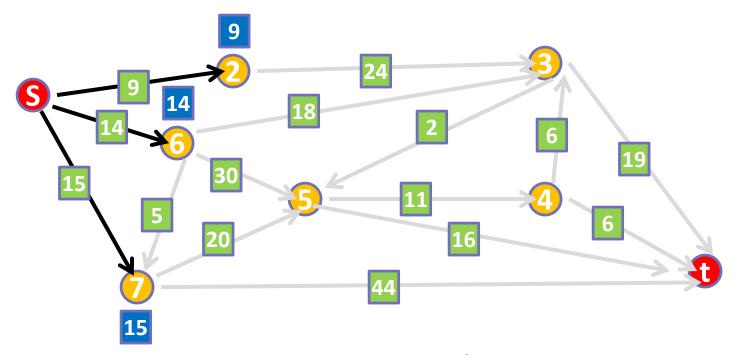
Moore's Shortest Path Algorithm

- Let d_i be the shortest distance from source to v_i
- Let $\mathbf{w}_{i,j}$ be the weight between \mathbf{v}_i and \mathbf{v}_j
- Iteratively search over the graph
 - \triangleright If d_i be the new shortest distance from source to v_i
 - For each \mathbf{v}_{j} adjacent to \mathbf{v}_{i} , update the shortest distance to \mathbf{v}_{j} $new_{-}d_{j} = \min(d_{j}, d_{i} + w_{i,j})$



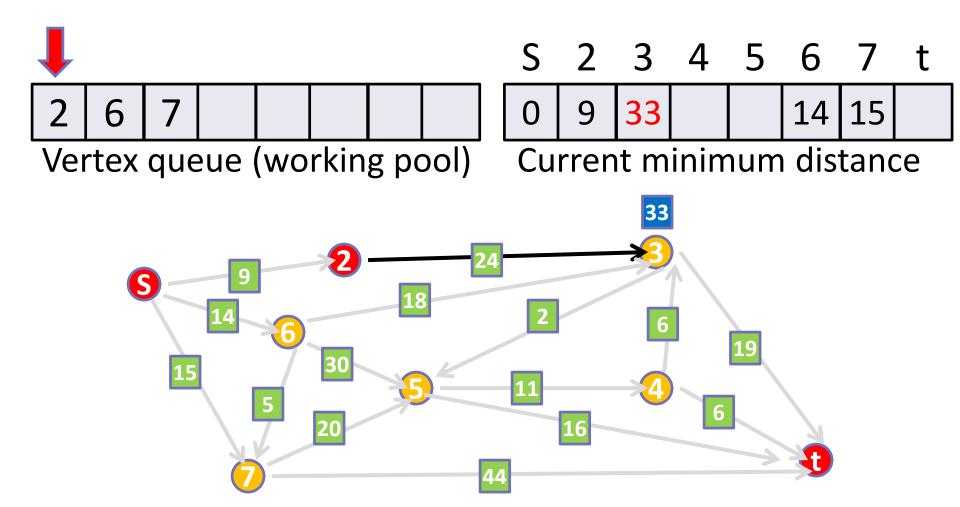




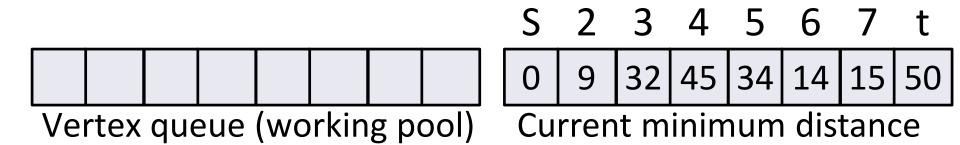


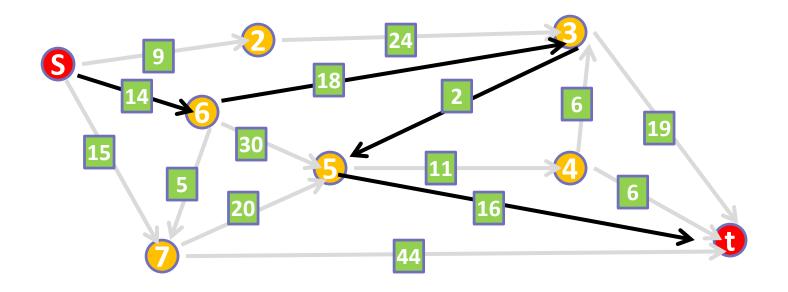
Parallel Programming – NTHU LSA Lab





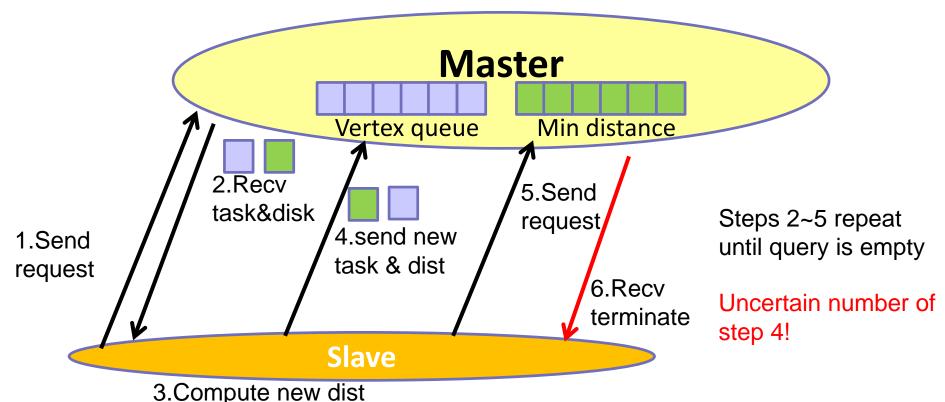






Centralized Pool Parallel Code

- Centralized work pool holds the vertex queue and distance array
- Each slave takes a vertex and returns new vertices and distances
- Since graph weights is fixed, it could be copied into each slave.



Parallel Programming – NTHU LSA Lab

M

Distributed Pool Parallel Code

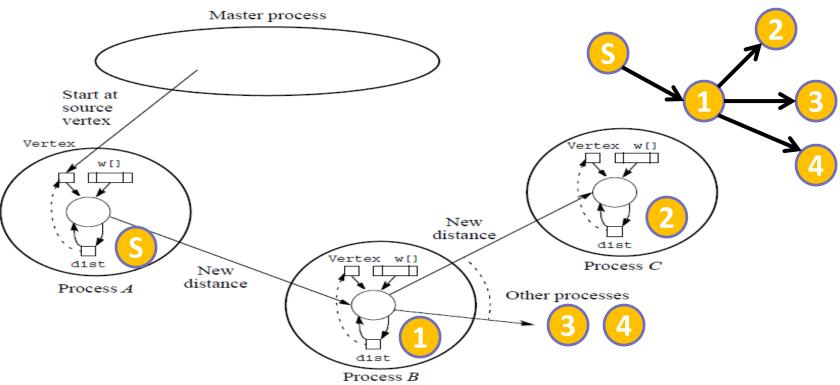
- Distribute vertex queue and distance array
 - > Each process holds all the information of a vertex
 - > Including queue, distance and adjacent weight

Algorithm

- > For each new distance to a vertex V_i is computed
 - → send the distance to P_i and re-activate P_i
 - if the new distance is shorter:
 - compute new distance for each adjacent vertex V_j
 - send the new distance to P_i



Distributed Pool Parallel Code



```
recv(newdist, PANY);
if (newdist < dist)
    dist = newdist; /* start searching around vertex */
    for (j = 1; j < n; j++) /* get next edge */
        if (w[j] != infinity) {
            d = dist + w[j];
            send(&d, Pj); /* send distance to proc j */
        }</pre>
```

Q: How do we know when to terminate?

A: Use dual-pass ring algorithm...



Distributed Termination Detection

- Termination at time t, must satisfy...
 - 1. Application-specific local termination conditions exist throughout the collection of processes
 - There are no messages in transit between processes

- ➤ 2nd condition is necessary because a message in transit might **reactivate** a terminated process
- ➤ However, it is difficult to recognize because the message transfer time is not known in advance

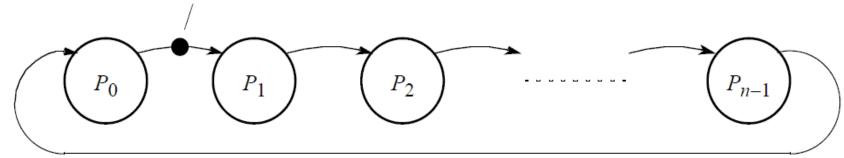


Single-Pass Ring Termination Algorithm

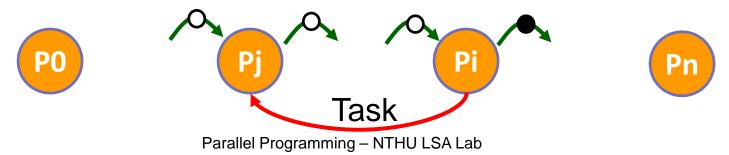
Steps:

- When P0 terminates, it generates a token that is passed to the next process P1
- When Pi (1 ≤ i ≤n) receives the token, it waits for its location termination condition and then passes the token to the next process Pi+1
- When P0 receives back the token, a global termination message can be sent to all processes
- Must assume processes can NOT be re-activated

Token passed to next processor when reached local termination condition

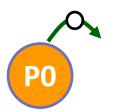


- It can handle processes being reactivated but requires two passes around the ring
- Reactivation occurs when ...
 - \triangleright P_i , to pass a task to P_j where j < i and after a token has passed P_j
- To handle reactivation, the token must recirculate through the ring a second time
 - > The token is colored BLACK when it goes through Pi
 - > The token becomes WHITE again once it goes through PO
 - ➤ If P0 receive WHITE tokens, all processes are terminated





 P0 becomes white when it has terminated and generates a white token to P1.





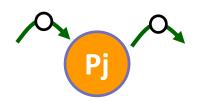






- 2. The token is passed through the ring from one process to the next when each process has terminated, but the color of the token may be changed.
 - A white process will pass on the token in its original color (either black or white). After Pi has passed on a token, it becomes a white process. Pn-1 passes the token to P0.





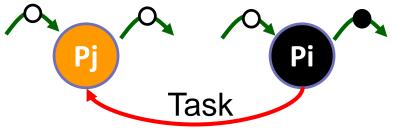






- 2. The token is passed through the ring from one process to the next when each process has terminated, but the color of the token may be changed.
 - ▶ If Pi passes a task to Pj where j < i (that is, before this process in the ring), it becomes a black process; otherwise it is a white process.
 - A black process will color a token black and pass it on.

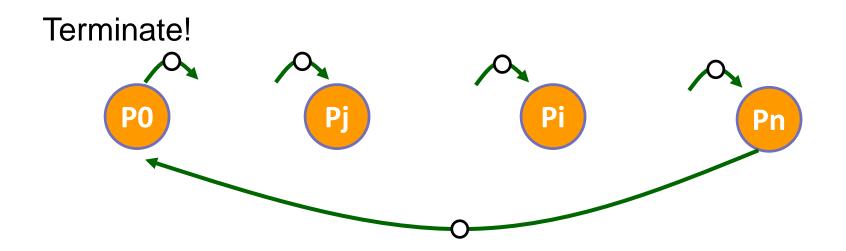




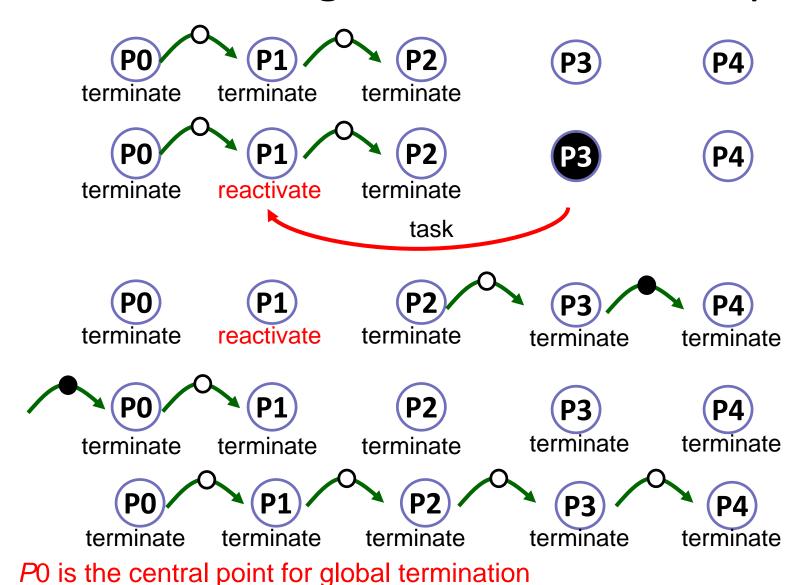




3. When P0 receives a black token, it passes on a white token; if it receives a white token, all processes have terminated.



Dual-Pass Ring Termination Example





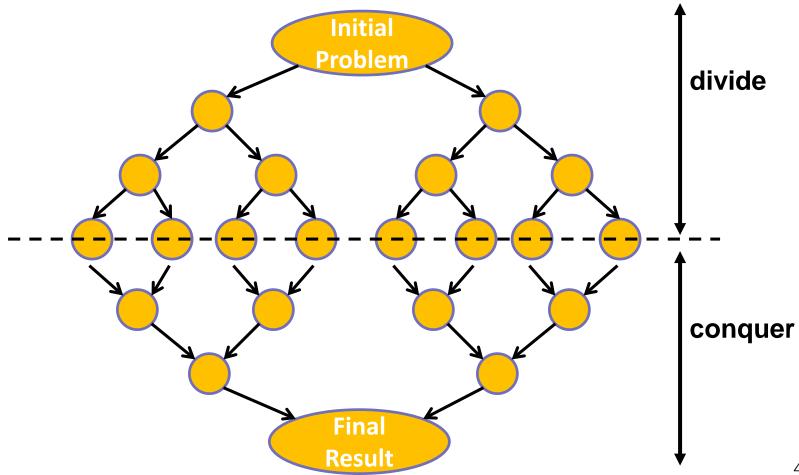
Outline

- Embarrassingly Computations
- Load Balancing & Termination
- Divide-And-Conquer Computations
 - Sorting Algorithms
 - N-Body Simulation
- Pipelined Computations
- Synchronous Computations

M

What is Divide & Conquer

■ Recursively divide a problem into sub-problems that are of the same form as the larger problem



46



Sorting

Re-arranging a list of numbers into increasing (strictly non-decreasing) order

> One of the most common and critical operation in

LSA Lab

large data processing

Who gets admission to the school?

Which is the most relevant page?



Google 搜索 手气不错



Which

friend to be



Sorting in Parallel

- In sequential
 - \triangleright We all know it is $n \log n$
 - > i.e. n is the number of elements
- In parallel with *n* processors
 - Optimal parallel time complexity

$$\frac{O(n\log n)}{n} = O(\log n)$$

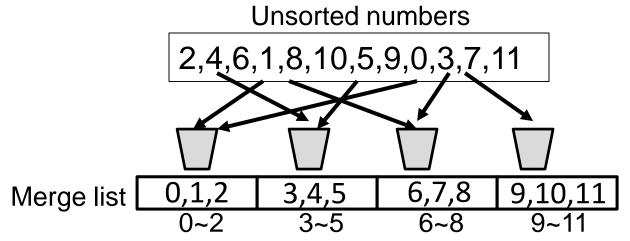
Has been obtained but the constant hidden in the order notation extremely large

LSA Lab 48



Bucket Sort

- Algorithm
 - Range of numbers is divided into m equal regions
 - 2. One bucket is assigned for each region
 - 3. Place numbers to buckets based on the region
 - 4. Use sequential sort for each bucket



- Only effective if number of items per bucket is similar!!
 - Numbers should have a known interval ([max, min])
 - Numbers better to be uniformly distributed

10

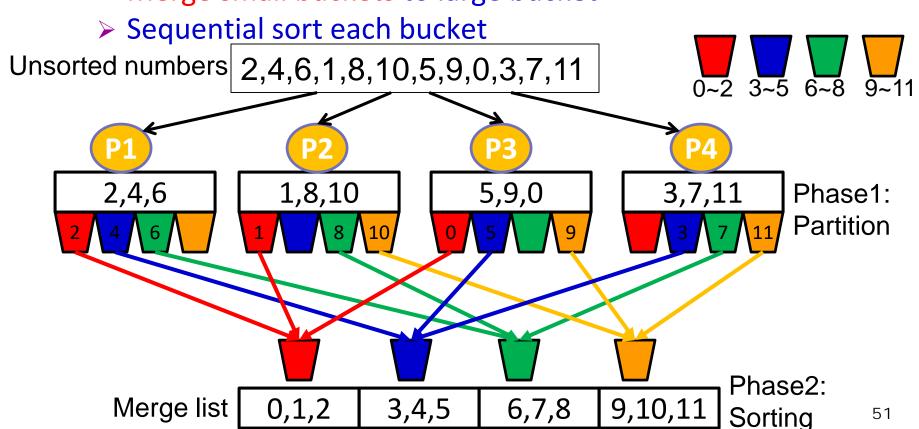
Complexity Analysis

- Sequential:
 - 1. Distribute numbers to bucket: O(n)
 - 2. Sequential sort each bucket: (n/m)log(n/m) x m
 - Overall: O(n log(n/m))
- Parallelize sorting: one process per bucket
 - 1. Distribute numbers to bucket: O(n)
 - Sequential sort each bucket: (n/m)log(n/m)
 - Overall: O(n + n/m log(n/m))
 - > A single process must scan through all numbers in step1



Further Parallelized Bucket Sort

- Parallelize partitioning and sorting:
 - Partition numbers to m parts/processes
 - > Each process divides its numbers to small buckets
 - Merge small buckets to large bucket





Merge Sort

■ Divide & Conquer

 \triangleright Sequential: $O(n \log n)$

 \triangleright Parallel: O(n)

 $T_{comm} = O\left(2\left(\frac{n}{2} + \frac{n}{4} + \dots + 1\right)\right) = O(n)$

#elements

Unsorted list
$$T_{comp} = 0 \left(n + \frac{n}{2} + \frac{n}{4} + \cdots 2 \right) = O(n)$$

Pivide list $T_{comp} = 0 \left(n + \frac{n}{2} + \frac{n}{4} + \cdots 2 \right) = O(n)$

A 2 7 8 5 1 3 6 P₀ P₂ P₄ P₅ P₆ P₇

A 2 7 8 5 1 3 6 P₀ P₁ P₂ P₃ P₄ P₅ P₆ P₇

Merge

1 2 3 4 5 6 7 8

Sorted list

3

5

LSA Lab

Process allocation



Quick Sort

- Most popular sequential sorting algorithm
 - Parallel: Iteratively pick pivot and partition numbers
- Complexity:
 - \triangleright Sequential: $O(n \log n)$
 - \triangleright Parallel: O(n)

Unsorted list

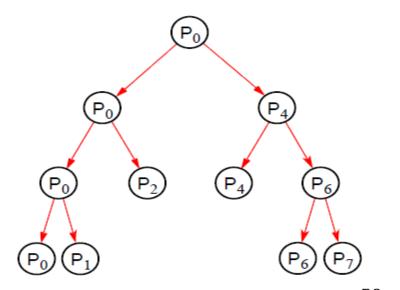
4 2 7 8 5 1 3 6

3 2 1 4 5 7 8 6

2 1 3 4 5 7 8 6

1 2 3 6 7 8

Still best choose in parallel? Not really & load might not be balanced



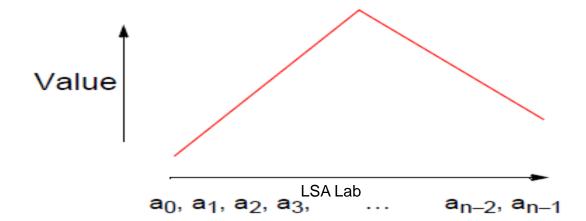
53



Bitonic Mergesort

- A parallel sorting algorithm based on bitonic sequence
 - Remove the O(N) bottleneck for merging
- Monotonic sequence:
 - A sequence of increasing numbers

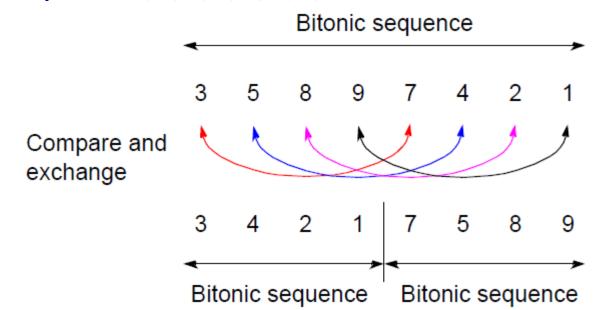
- Bitonic sequence:
 - Two sequences, one increasing and one decreasing



M.

Characteristic of Bitonic Sequence

- If we perform a compare-and-exchange operation on a_i with $a_{(i+n/2)}$ for all i
 - Get TWO bitonic sequences
 - The numbers in one sequence are all less than the numbers in the other sequence
- Example: <3,5,8,9,7,4,2,1>

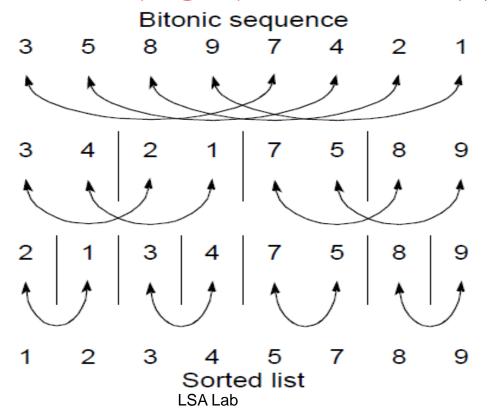




Bitonic Sorting Operation

- Given a bitonic sequence, recursively performing compare-and-exchange will sort the list
- Complexity is only O(log n) instead of O(n)

Compare and exchange

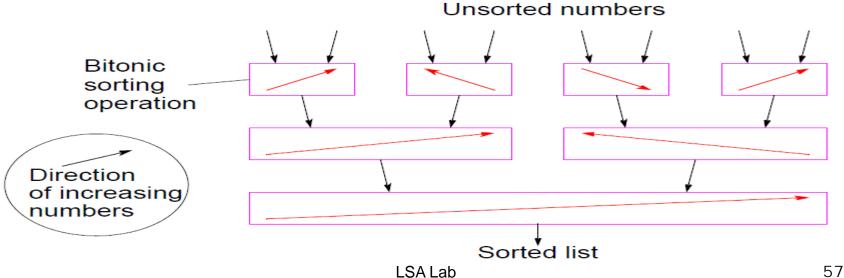




Bitonic Mergesort

- Combine bitonic sorting into merge sorting
- Complexity:
 - \triangleright Let $n=2^k$: there will be k phases
 - \triangleright kth phase has size $2^k \rightarrow$ need only k steps with bitonic sorting

$$T_p = O\left(\sum_{i=1}^k i\right) = O\left(\frac{k(k+1)}{2}\right) = O\left(\frac{\log n(\log n)}{2}\right) = O(\log^2 n)$$





Rank Sort

- Count the number of numbers that are smaller than each of the selected number
- Sequential: $O(n^2)$

```
for(i=0; i<n; i++) {
    rank=0;
    for(j=0; j<n; j++) {
        if (a[i] < a[j]) rank++;
    }
    b[rank] = a[i];
}</pre>
```

```
rank: 1 0 2 3
a[i]: 4 1 6 8
b[i]: 1 4 6 8
```

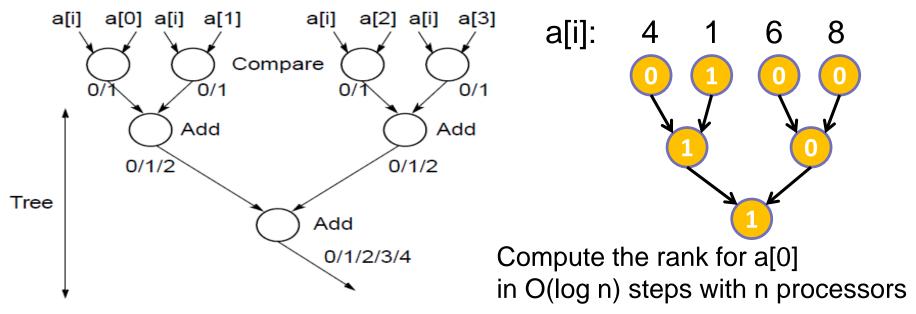
- Parallel: O(n)
 - use one processor for each number

LSA Lab 58



Rank Sort with n² Processors

- We can use n processor to compute the rank of a number in parallel (summation parallel algo.)
- Time complexity? $O(\log n)$
- How about parallel efficiency? Actually decreased!

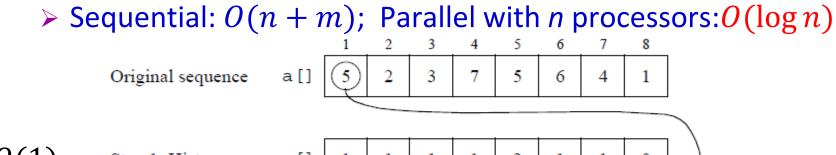


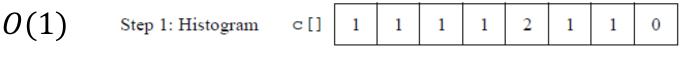
LSA Lab 59



Counting Sort

- If the numbers are integers and the range is known
 - \triangleright Coding the rank sort algorithm to reduce the sequential time from $O(n^2)$ to O(n)
- Use array c[] to count the histogram of values
- Complexity

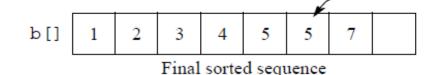




Then decrement c [5]

Move 5 to position 6.

$$O(\log n)$$
 Step 2: Prefix sum



Summary

■ With #processors = n

	Sequential	Parallel	
Bucket Sort	$O(n \log n)$	O(n)	
Merge Sort	$O(n \log n)$	O(n)	
Quick Sort	$O(n \log n)$	O(n)	
Bitonic	$O(n \log n)$	$O(\log^2 n)$	
mergesort	2		With n ²
Rank Sort	$O(n^2)$	$O(n)$ $O(\log n)$	processors
Counting Sort*	O(n)	$O(\log n)$	

^{*}Special case with known value range

LSA Lab 61

N-Body Problem

- Newtonian laws of physics
 - > The gravitational force between two bodies of masses $m_a \& m_b$:

$$F = \frac{Gm_a m_b}{r^2}$$

Subject to the force, acceleration occurs $F = m \times a$

- Let the time interval be Δt & current velocity v^t , position x^t
 - \triangleright New velocity v^{t+1} :

$$F = m \frac{v^{t+1} - v^t}{\Delta t} \Rightarrow v^{t+1} = v^t + \frac{F\Delta t}{m}$$

 \triangleright New position x^{t+1} :

$$x^{t+1} = x^t + v^{t+1} \Delta t$$





Three-Dimensional Space

- Considering 2 bodies at $(x_a, y_a, z_a) \& (x_b, y_b, z_b)$ $r = \sqrt{(x_a - x_b)^2 + (y_a - y_b)^2 + (z_a - z_b)^2}$
- The forces, velocities and positions can be resolved in the three direction independently

$$F_{x} = \frac{Gm_{a}m_{b}}{r^{2}} \left(\frac{x_{b} - x_{a}}{r} \right)$$

$$F_{y} = \frac{Gm_{a}m_{b}}{r^{2}}\left(\frac{y_{b} - y_{a}}{r}\right)$$

$$F_z = \frac{Gm_am_b}{r^2} \left(\frac{z_b - z_a}{r}\right)$$

N-Body Sequential Code

Assume all bodies have the same mass m

```
for (t=0; t<T; t++) {
   for (i=0; i<N; i++) {
       F = Compute_Force(i); // compute force in O(N^2)
       v_new[i] = v[i] + F *dt / m; // compute new velocity
       x_new[i] = x[i] + v_new[i] * dt; // compute new position
   for(i=0; i<N; i++){
       x[i] = x_new[i];
                                    // update position
                                    // update velocity
       v[i] = v_new[i];
```

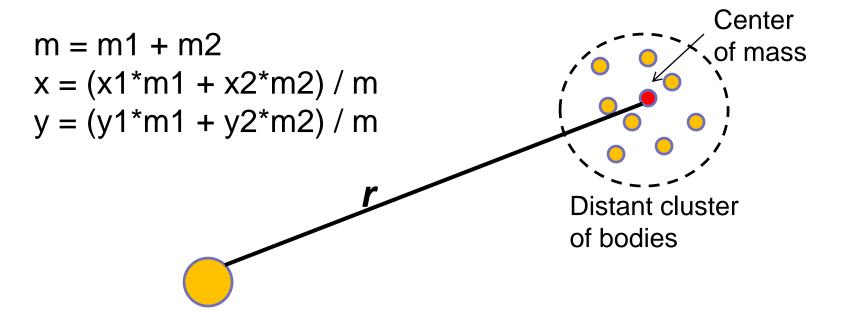
■ Non-feasible as N increases due to $O(N^2)$ complexity



Approximate Algorithms

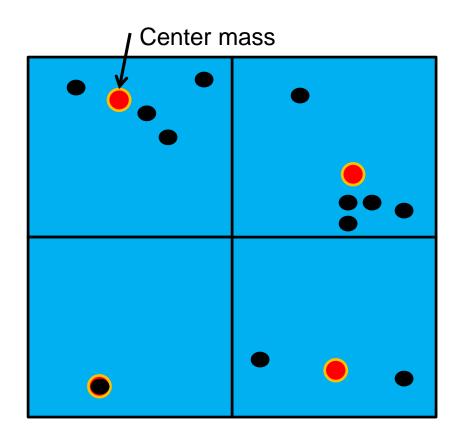
 Reduce time complexity by approximating a cluster of bodies as a single distant body

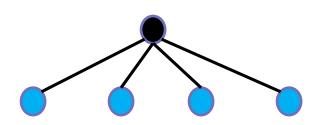
How to find those clusters of bodies?





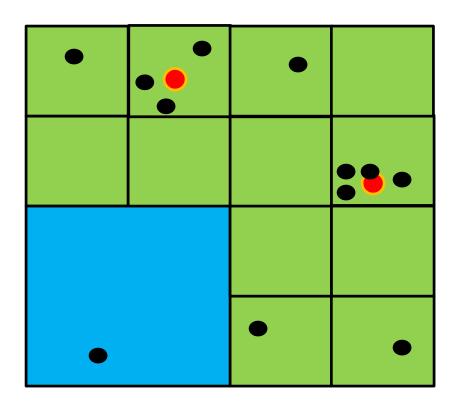
- Step1: Recursively divide space by two in each dimensions
 - > Record the center mass and position of each internal node

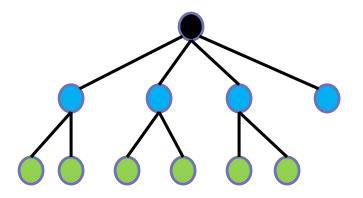






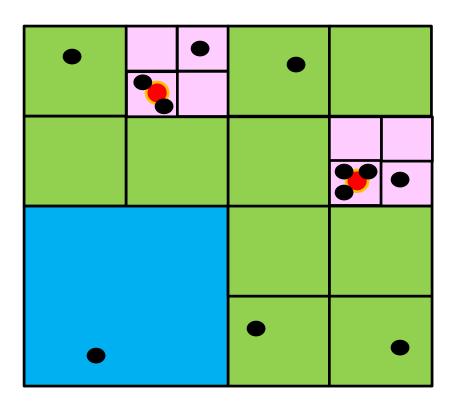
- Step1: Recursively divide space by two in each dimensions
 - > Record the center mass and position of each internal node

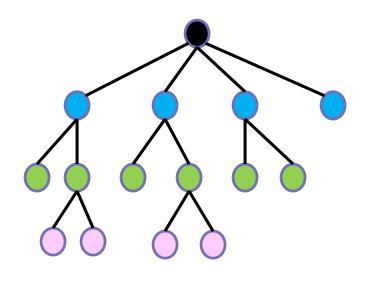






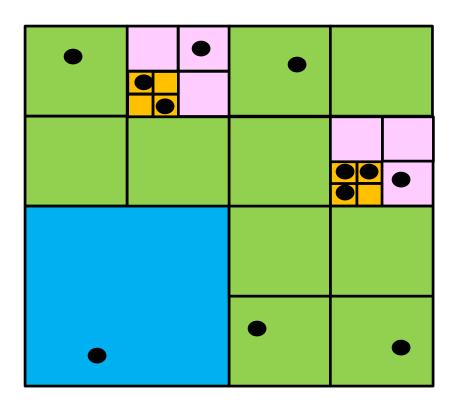
- Step1: Recursively divide space by two in each dimensions
 - > Record the center mass and position of each internal node

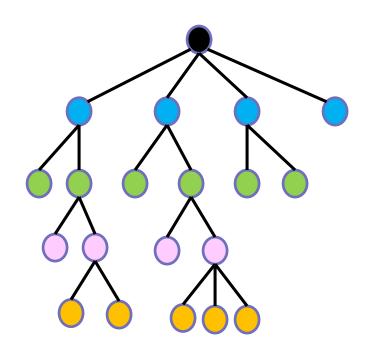






- Step1: Recursively divide space by two in each dimensions
 - > Record the center mass and position of each internal node

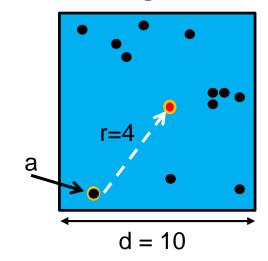


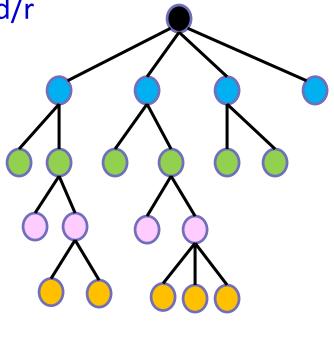


- Step2: Compute approximate forces on each object
 - 1. traverse the nodes of the tree, starting from the root.
 - 2. If the center-of-mass of an **internal node** is **sufficiently far** from the body, approximate the internal node as a single body
 - \triangleright Far is determined by a parameter: $\theta=d/r$
 - r: the distance between the body and the node's center-of-mass
 - ♦ d: the width of the region

Example: θ =0.5

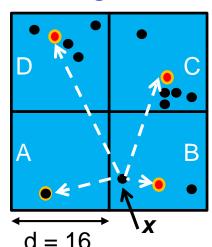
 $d/r=2.5 > \theta$

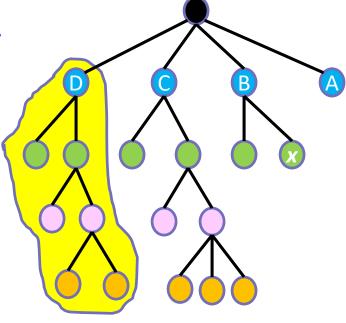




- Step2: Compute approximate forces on each object
 - 1. Traverse the nodes of the tree, starting from the root.
 - 2. If the center-of-mass of an **internal node** is **sufficiently far** from the body, approximate the internal node as a **single body**
 - \triangleright Far means $d/r < \theta$ (e.t. $0 < \theta < 1$)
 - r: the distance between the body and the node's center-of-mass
 - ♦ d: the width of the region

Example: θ =1 d/r_A =16/10 > θ d/r_B =16/2 > θ d/r_C =16/15 > θ d/r_D =16/20 < θ

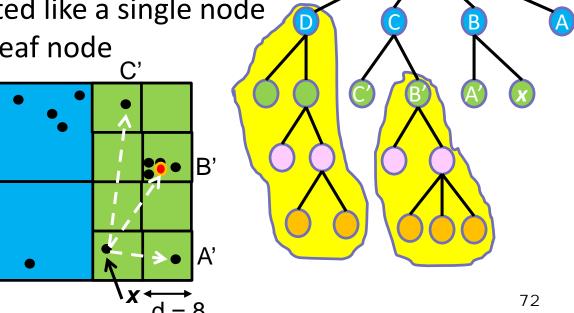




- Step2: Compute approximate forces on each object
 - 3. If it is a leaf node, calculate the force and add to the object.
 - 4. Otherwise, recursively compute the force from children of the internal node.

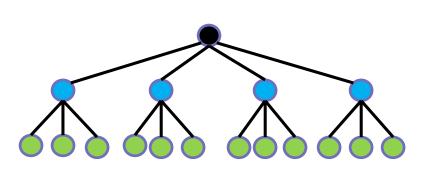
Example: $\theta=1$ $d/r_{A'}=8/7 > \theta \rightarrow A'$ is a leaf node $d/r_{B'}=8/15 < \theta \rightarrow B'$ treated like a single node

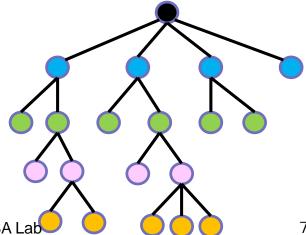
 $d/r_{C'}=8/20 < \theta \rightarrow C'$ is a leaf node



Barnes-Hut Algorithm

- lacksquare 0 controls the accuracy and approximation error of the algorithm
 - $\Rightarrow \theta = 0 \Rightarrow d/r$ ALWAYS larger than $\theta \Rightarrow$ same as brute force
 - \Rightarrow 0 = 1 \Rightarrow most likely only need to consider the object within the same cluster/region
- If the tree is balanced, the complexity is $O(n \log n)$
 - > But in general, the tree could be very unbalanced
- The tree must be re-built for each time interval

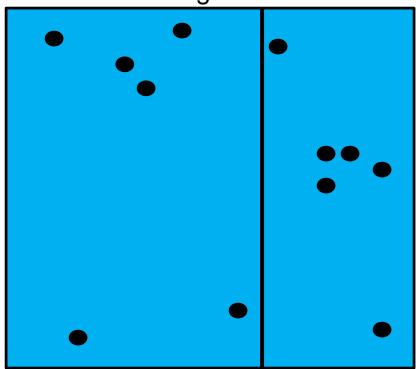


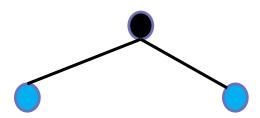




Recursively evenly divide space with the same number of bodies in each of the dimensions

Divide along x dimension

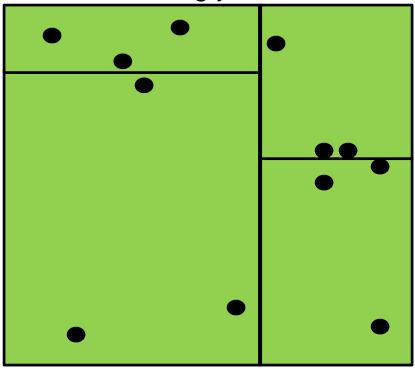


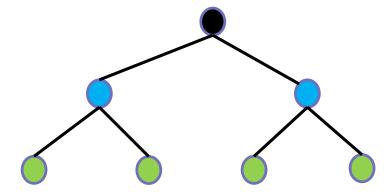




Recursively evenly divide space with the same number of bodies in each of the dimensions

Divide along y dimension

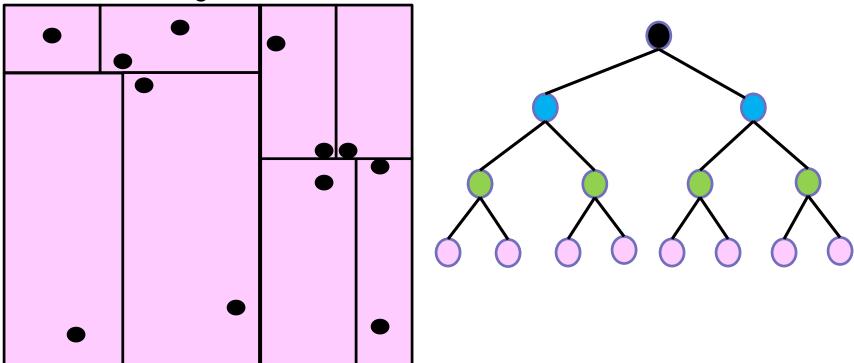






Recursively evenly divide space with the same number of bodies in each of the dimensions

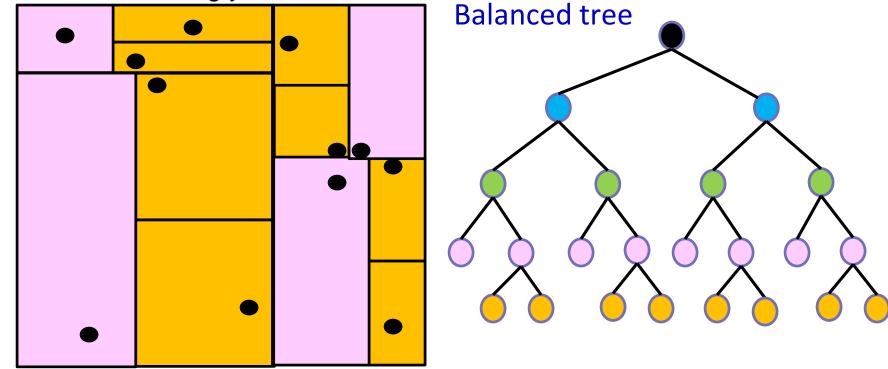
Divide along x dimension





Recursively evenly divide space with the same number of bodies in each of the dimensions

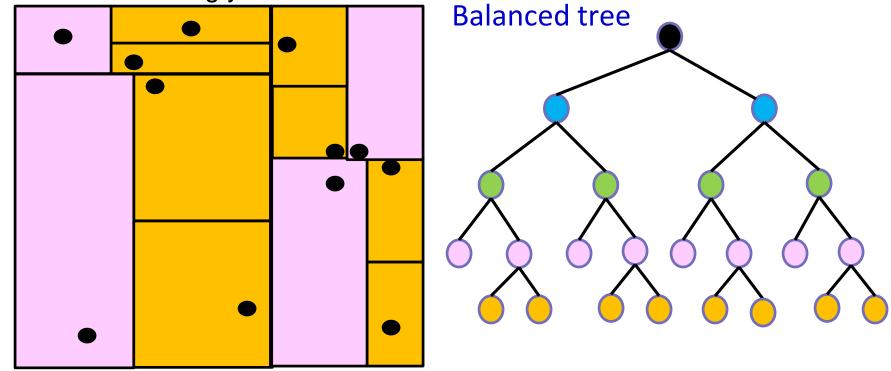
Divide along y dimension





- It is more balanced, but less accurate
 - Objects close to each other may not in the same cluster

Divide along y dimension





Outline

- Embarrassingly Computations
- Load Balancing & Termination
- Divide-And-Conquer Computations
- Pipelined Computations
 - Adding Numbers
 - Sorting Numbers
 - Linear Equation Solver
- Synchronous Computations

What is Pipelined Computations

- A problem is divided into a series of tasks
- Tasks have to be completed one after the other
- Each task will be executed by a separate process or processor



$$\rightarrow$$
 P0 \rightarrow P1 \rightarrow P2 \rightarrow P3 \rightarrow



Types of Pipelined Computations

- Pipelined approach can provide increased speed under three types of computations:
- If more than one instance of the complete problem is to be executed
- If a single instance has a series of data items must be processed, each requiring multiple operations
- If information to start the next process can be passed forward before the process has completed all its internal operations

Type 1 Pipelined Computations

1. If more than one **instance** of the complete problem is to be executed

	_		p-1			_		m			
						Instance	Instance	Instance	Instance	Instance	
P_5						1	2	3	4	5	
ъ					Instance	Instance	Instance	Instance	Instance	Instance	
P_4					1	2	3	4	5	6	
P_3				Instance							
				1	2	3	4	5	6	7	
P_2			Instance								
			1	2	3	4	5	6	7		
P_1		Instance									
		1	2	3	4	5	6	7			
P_0	Instance										
	1	2	3	4	5	6	7				
	- 1 -	- •		- •	1.0		_				

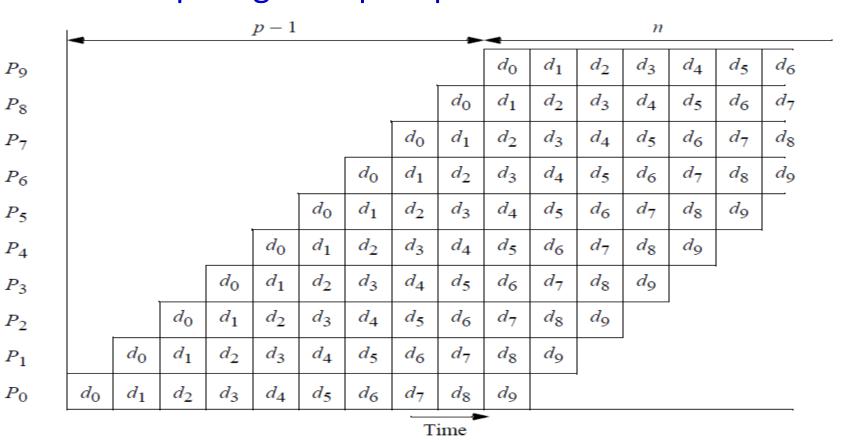
(Alternative space-time diagram) Time

- After the first (p-1) cycles, one problem instance is completed in each pipeline cycle
- The number of instance should be >> the number of



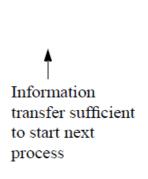
Type 2 Pipelined Computations

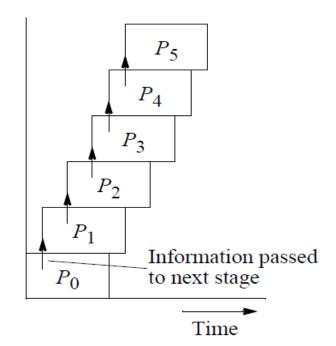
2. If a series of **data** items must be processed, each requiring multiple operations

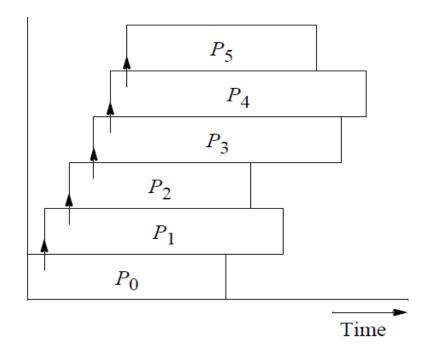


Types 3 Pipelined Computations

 Only one problem instance, but each process can pass on information to the next process, before it has completed





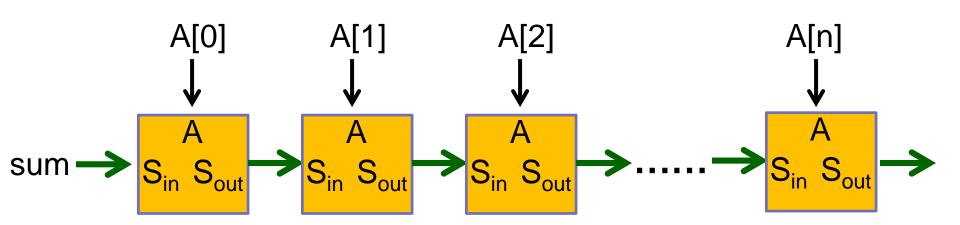


(a) Processes with the same execution time

(b) Processes not with the same execution time

Example1: Adding Numbers

- Compute sum of an array:
 - for(i=0; i<n; i++) sum += A[i]</pre>
- Pipeline for an unfolded loop:
 - \triangleright sum += A[0], sum += A[1], sum += A[2],





Example1: Adding Numbers

■ The basic code for Pi:

```
recv(&sum, P<sub>i-1</sub>);
sum += number;
send(&sum, P<sub>i+1</sub>);
```

■ For the first process, P0:

```
send(&sum, P<sub>i+1</sub>);
```

■ For the last process, Pn-1:

```
recv(&sum, P<sub>i-1</sub>);
sum += number;
```

■ SPMD Program:

```
// code for process Pi
if (Pi != P0) {
  recv(&sum, P<sub>i-1</sub>);
  sum += number;
}
if (Pi != Pn) {
  send(&sum, P<sub>i+1</sub>);
}
```



Example2: Sorting Numbers

■ Insertion Sort:



- 5 2

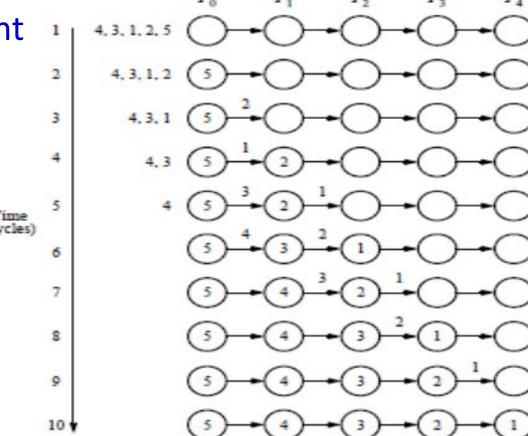


Example2: Sorting Numbers

- Insertion Sort:
 - > Each process holds one number
 - Compare & move the smaller

```
number to the right
```

```
recv(&number, P<sub>i-1</sub>);
if (number > x) {
    send(&x, P<sub>i+1</sub>);
    x= number;
} else {
    send(&number, P<sub>i+1</sub>);
}
```



10

Example 3: Linear Equation Solver

- Special linear equations of "upper-triangular" form
 - > a's and b's are constants, x's are unknown to be found

$$\begin{pmatrix}
a_{n-1,0} & a_{n-1,1} & a_{n-1,2} & \dots & a_{n-1,n-1} \\
a_{n-2,0} & a_{n-2,1} & \dots & a_{n-2,n-2} & 0 \\
\vdots & \vdots & \vdots & 0 & 0 \\
a_{1,0} & a_{1,1} & 0 & 0 & 0 \\
a_{0,0} & 0 & 0 & 0
\end{pmatrix}
\begin{pmatrix}
x_0 \\
x_1 \\
\vdots \\
x_{n-2} \\
x_{n-1}
\end{pmatrix} = \begin{pmatrix}
b_{n-1} \\
b_{n-2} \\
\vdots \\
b_1 \\
b_0
\end{pmatrix}$$

$$a_{n-1,0}x_0 + a_{n-1,1}x_1 + a_{n-1,2}x_2 + \dots + a_{n-1,n-1}x_{n-1} = \mathbf{b}_{n-1}$$

$$\vdots$$

$$a_{2,0}x_0 + a_{2,1}x_1 + a_{2,2}x_2 = \mathbf{b}_2$$

$$a_{1,0}x_0 + a_{1,1}x_1 = \mathbf{b}_1$$

$$a_{0,0}x_0 = \mathbf{b}_0$$



Example 3: Linear Equation Solver

Back Substitution

 $\triangleright x_0$ is found from the last equation

$$x_0 = \frac{b_0}{a_{0,0}}$$

 \triangleright Value for x_0 is substituted into the next equation

$$x_1 = \frac{b_1 - a_{1,0} x_0}{a_{1,1}}$$

 \triangleright Values for x_0 , x_1 are substituted into the next equation

$$x_2 = \frac{b_2 - a_{2,0}x_0 - a_{2,1}x_1}{a_{2,2}}$$

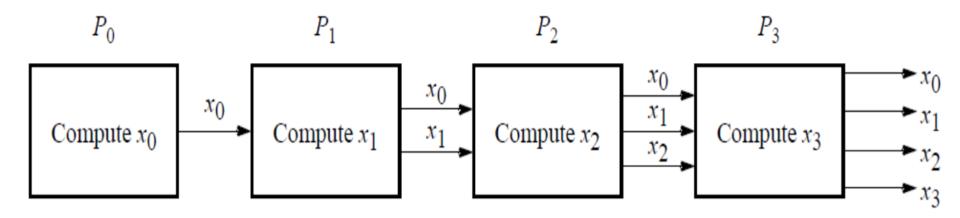
> So on until all unknowns are found ...

$$x_i = \frac{b_i - \sum_{j=0}^{i-1} a_{i,j} x_j}{a_{i,i}}$$

т.

Example 3: Linear Equation Solver

■ First pipeline stage computes x_0 and passes x_0 onto the second stage, which computes x_1 from x_0 and passes both x_0 and x_1 onto the next stage, which computes x_2 from x_0 and x_1 , and so on



Example 3: Linear Equation Solver

■ Parallel Code

```
// code for Pi
sum = 0;
for (j=0; j<i; j++) { // compute partial result recv(&x[j], P<sub>i-1</sub>); // once data is available send(&x[j], P<sub>i+1</sub>);
sum += a[i][j]*x[j];
}
x[i] = (b[i] - sum) / a[i][j]; // send out final result to send(&x[j], P<sub>i+1</sub>); // next process

Final computed value

First value passed onward
```

■ Time complexity:

- $> O(n^2)$ without passing forward
- $\triangleright O(n)$ with passing forward

Time



Outline

- Embarrassingly Computations
- Load Balancing & Termination
- Divide-And-Conquer Computations
- Pipelined Computations
- Synchronous Computations
 - Prefix Sum
 - System of Linear Equations
 - > Heat Distribution

.

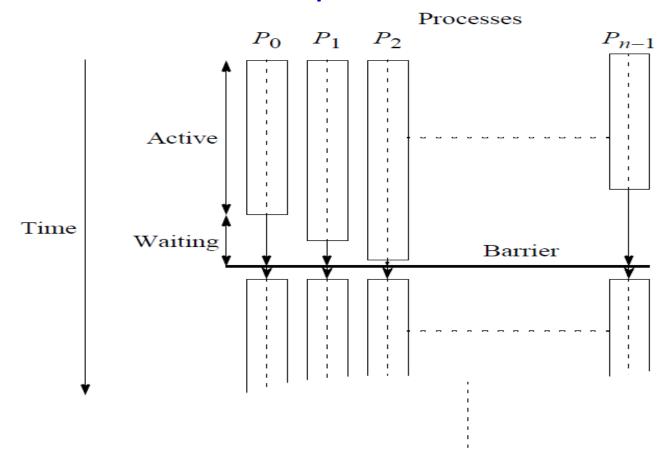
Synchronous Computations

- **Definition:** all the processes *synchronized* at regular points
- Barrier: Basic mechanism for synchronizing processes
 - Inserted at the point in each process where it must wait
 - Message (token) is passed among processes for synchronization
- Deadlock: Common problem occurs from synchronization
 - Two or multiple processes waiting for each other



Barrier

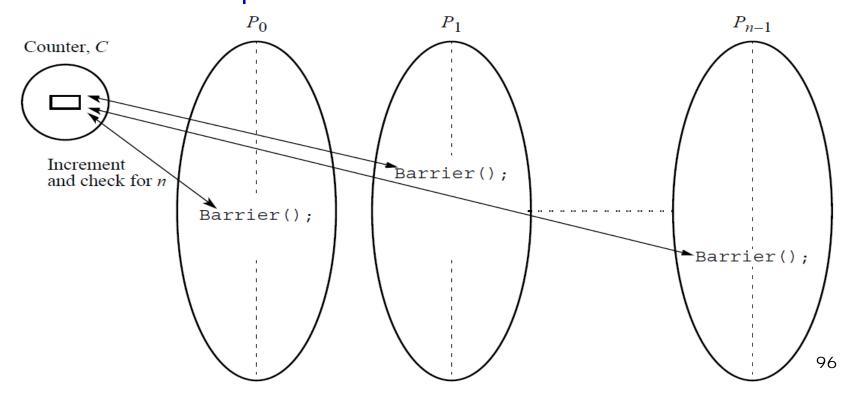
■ All processes can only continue from this POINT when all the processes have reached it



M

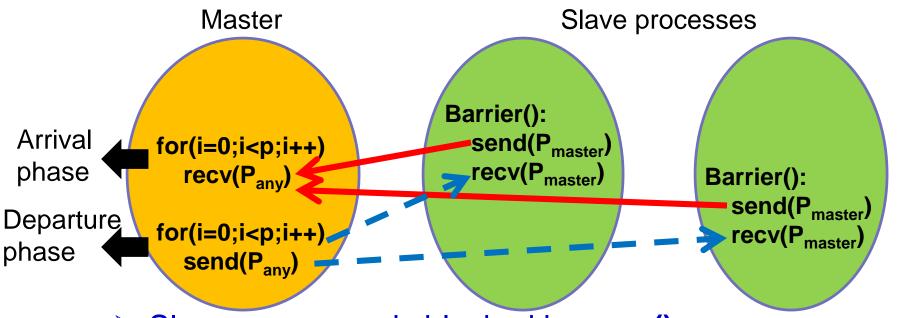
Counter Barrier Implementation

- A.k.a: Linear Barrier
 - Centralized counter: count # of processes reaching the barrier
 - Increase & check the counter for each barrier call
 - Processes is locked by the barrier call until counter == # processes



Counter Barrier Implementation

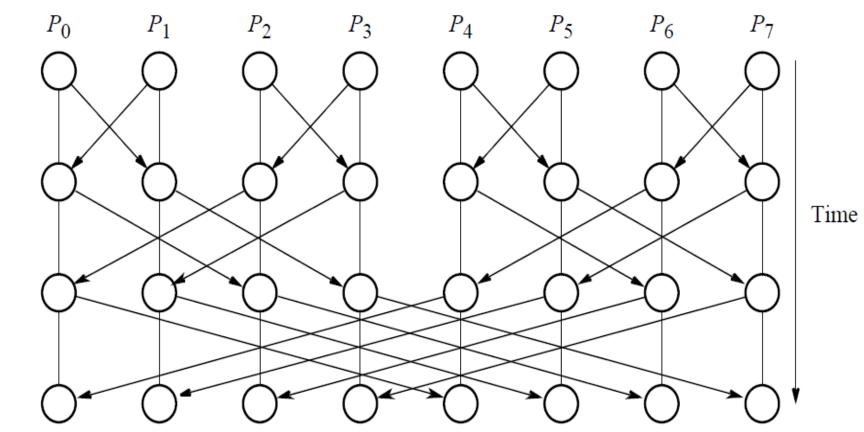
- Counter-based barrier often have two phases
 - > Arrival phase: a process enters arrival phase and does not leave this phase until all processes have arrived in this phase
 - Departure phase: Processes are released after moving to the departure phase



- Slave processes is blocked by recv()
- Master could be a bottleneck

Butterfly Barrier Implementation

At stage i, each process passes a token to the process with 2ⁱ distance away



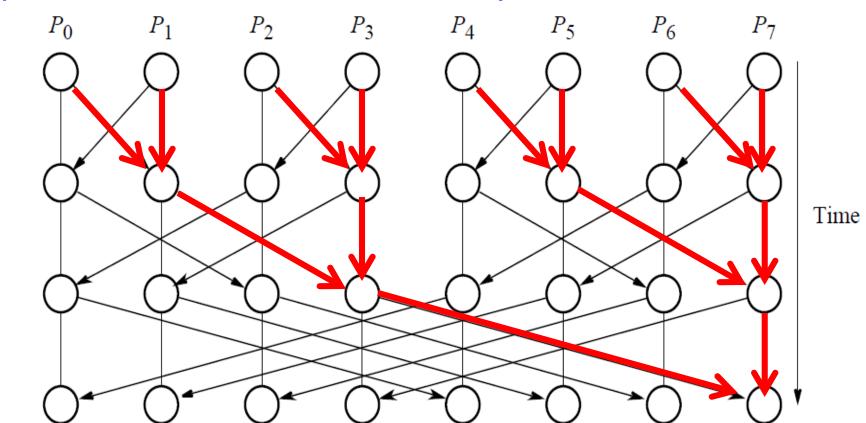
1st stage

2nd stage

3rd stage

Butterfly Barrier Implementation

At stage i, each process passes a token to the process with 2i distance away



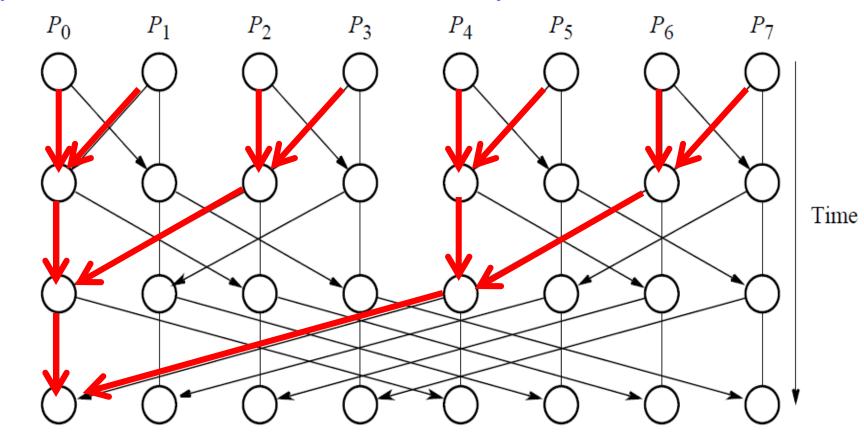
1st stage

2nd stage

3rd stage

Butterfly Barrier Implementation

At stage i, each process passes a token to the process with 2i distance away



1st stage

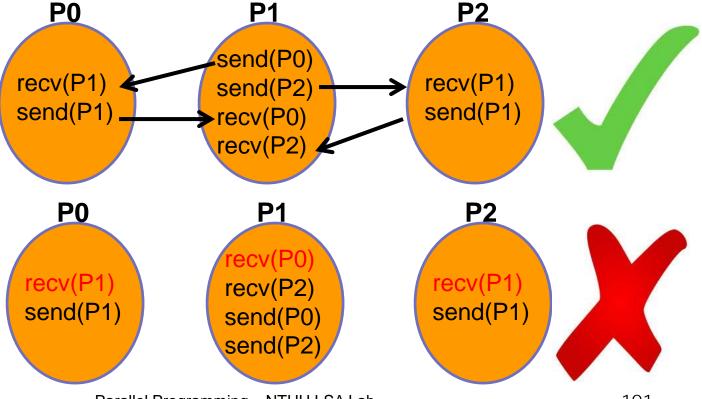
2nd stage

3rd stage

Deadlock Problem

A set of blocked processes each holding some resources and waiting to acquire a resource held by another process in the set

■ Example:



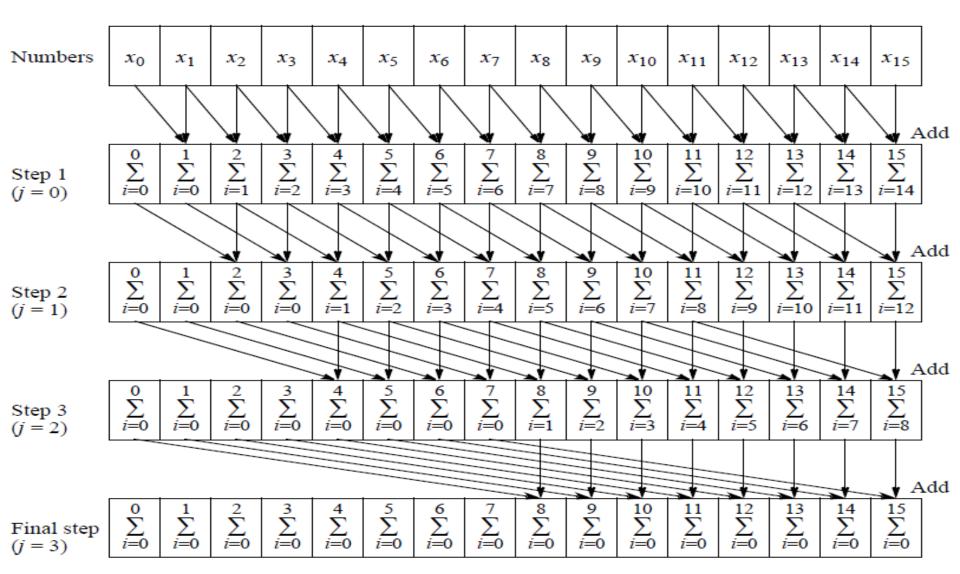


Example 1: Prefix Sum

- Given a list of numbers $x_0, x_1, ..., x_{n-1}$, compute all *partial* summations
 - $> x_0; x_0 + x_1; x_0 + x_1 + x_2; \dots$
 - ➤ Could also replace operator + with AND, OR, *, etc.
- Example:
 - x = 1,2,3,4,5
 - > Sum = 1,3,6,10,15
- Sequential code: O(n²)

```
//sequential code
for(i = 0; i < n; i++) {
    sum[i] = 0;
    for (j = 0; j <= i; j++)
        sum[i] = sum[i] + x[j];
}
```

Data Parallelism Solution



ne.

Data Parallelism Code

■ Sequential Code: O(n²), optimal: O(n)

```
for (j = 0; j < log(n); j++) /* at each step */
for (i = 2<sup>j</sup>; i < n; i++) /* add to accumulating sum */
x[i] = x[i] + x[i - 2<sup>j</sup>]
```

Parallel Code: O(log n)

```
for (j = 0; j < log(n); j++) /* at each step */
forall (i = 0; i < n; i++) /* add to accumulating sum */
if (i >= 2^{j}) x[i] = x[i] + x[i - 2^{j}];
```



Each iteration composed of several processes that start together at beginning of iteration and next iteration cannot begin until all processes have finished previous iteration

openMP

MPI

```
for (j=0; j<n; j++) { // each iteration
  i = myrank;
  body(i);
  barrier(mygroup);
}</pre>
```

Example 2: System of Linear Equations

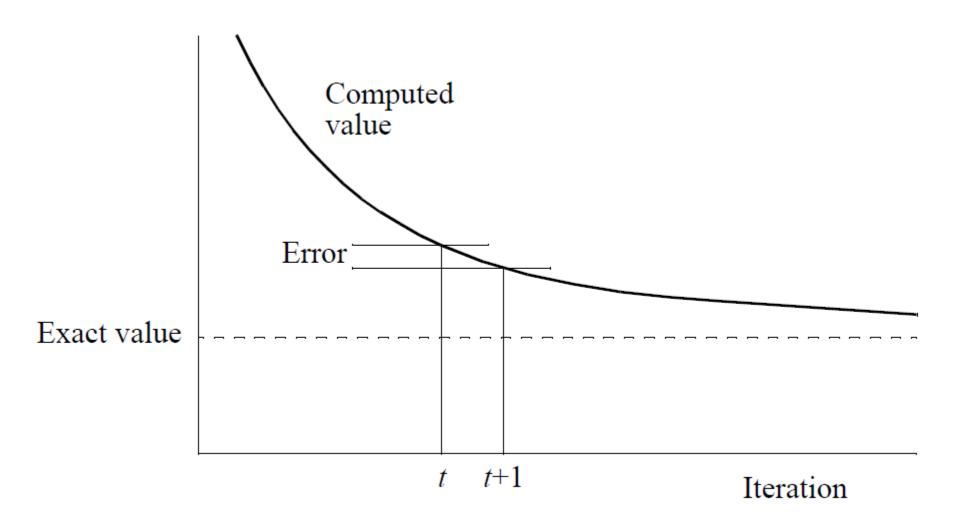
System of linear equations

$$\begin{pmatrix}
a_{0,0} & a_{0,1} & a_{0,2} & \dots & a_{0,n-1} \\
a_{1,0} & a_{1,1} & a_{1,2} & \dots & a_{1,n-1} \\
\vdots & \vdots & \vdots & \vdots & \vdots \\
a_{n-2,0} & a_{n-2,1} & a_{n-2,2} & \dots & a_{n-2,n-1} \\
a_{n-1,0} & a_{n-1,1} & a_{n-1,2} & \dots & a_{n-1,n-1}
\end{pmatrix}
\begin{pmatrix}
x_0 \\
x_1 \\
\vdots \\
x_{n-2} \\
x_{n-1}
\end{pmatrix} =
\begin{pmatrix}
b_0 \\
b_1 \\
\vdots \\
b_{n-2} \\
b_{n-1}
\end{pmatrix}$$

- Jacobi iteration algorithm:
 - > Convert *i*th iteration to $x_i = \frac{1}{a_{i,i}} [b_i \sum_{i \neq j} a_{i,j} x_j]$
 - \triangleright Initial guess with $x_i = b_i$, and calculate new x_i values
 - ightharpoonupRepeat until $\begin{vmatrix} x_i & t-1 \\ x_i & -x_i \end{vmatrix}$ < error tolerance Parallel Programming NTHU LSA Lab



Jacobi iteration algorithm



Jacobi iteration algorithm example

$$\begin{cases} -x_0 + 2x_1 - x_2 = 2\\ 2x_0 + x_1 - 2x_2 = 2\\ 2x_0 - x_1 + 2x_2 = 2 \end{cases}$$
$$x_i = \frac{1}{a_{i,i}} [b_i - \sum_{i \neq j} a_{i,j} x_j]$$

$$\begin{cases} x_0 = 2 - \frac{(2x_1 - x_2)}{-1} \\ x_1 = 2 - \frac{(2x_0 - 2x_2)}{1} \\ x_2 = 2 - \frac{(2x_0 - x_1)}{2} \end{cases}$$

■ Iter1:
$$x_0^1 = 2$$
, $x_1^1 = 2$, $x_2^1 = 2$

⇒
$$x_0^2 = 2 - \frac{2x_1^1 - x_2^1}{-1} = 4$$
, $x_1^2 = 2$, $x_2^2 = 1$
⇒ $e_0 = |2 - 4| = 2$, $e_1 = 0$, $e_2 = 1$

$$rac{1}{2}e_0 = |2 - 4| = 2$$
, $e_1 = 0$, $e_2 = 1$

■ ilter2:
$$x_0^2 = 2 - \frac{2x_1^2 - x_2^2}{-1} = 5$$
, $x_1^3 = -2$, $x_2^3 = -1$

$$ightharpoonup e_0 = |4 - 5| = 1, \qquad e_1 = 4, \qquad e_2 = 2$$

Jacobi iteration algorithm

- Sequential Code
 - > a[][] and b[] holding constants in the equations
 - > x[] holding unknowns
 - fixed number of iterations

```
x_{i} = \frac{1}{a_{i,i}} [b_{i} - \sum_{i \neq j} a_{i,j} x_{j}]
```

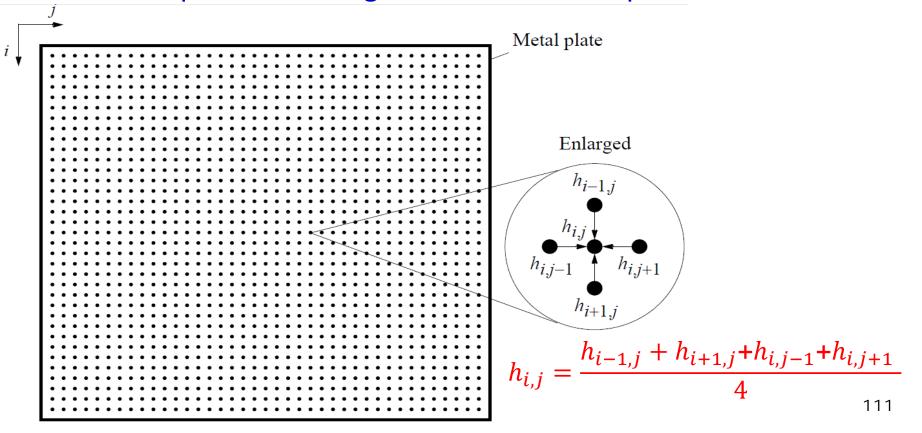


Jacobi iteration algorithm

- Parallel Code
 - Process i handles unknown x[i]

```
x[i] = b[i]; /*initialize unknown*/
for (iteration = 0; iteration < limit; iteration++) {
   sum = -a[i][i] * x[i];
   for (j = 0; j < n; j++) /* compute summation */
       sum = sum + a[i][j] * x[j];
   new_x[i] = (b[i] - sum) / a[i][i]; /* compute unknown */
   allGather(&new_x[i]); /* gather & broadcast new value */
                                    /* wait for all processes */
   barrier();
```

- lacksquare Consider an area as a mesh of points, $h_{i,j}$
 - Known temperatures along its edges
 - Heat is propagated until a fixed # of iterations or until temperature change is less than small pre-described amount





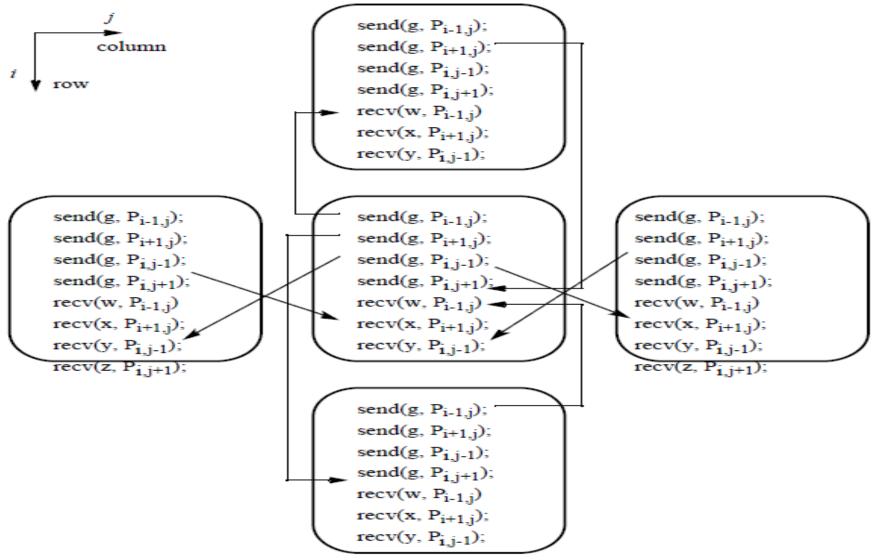
Sequential Code

```
for (iteration = 0; iteration < limit; iteration++) {
    for (i = 1; i < n; i++)
        for (j = 1; j < n; j++)
            g[i][j] = 0.25*(h[i-1][j]+h[i+1][j]+h[i][j-1]+h[i][j+1]);
    for (i = 1; i < n; i++) /* update points */
        for (j = 1; j < n; j++)
            h[i][j] = g[i][j];
}</pre>
```



Parallel Code

```
for (iteration = 0; iteration < limit; iteration++) {
         g = 0.25 * (w + x + y + z);
         Isend(&g, P<sub>i-1,i</sub>); /* non-blocking sends */
         Isend(&g, P_{i+1,i});
         Isend(&g, P<sub>i,i-1</sub>);
         Isend(&g, P_{i,i+1});
         recv(&w, P<sub>i-1,i</sub>); /* synchronous receives */
         recv(&x, P_{i+1,i});
         recv(&y, P<sub>i,i-1</sub>);
         recv(&z, P<sub>i,i+1</sub>);
```





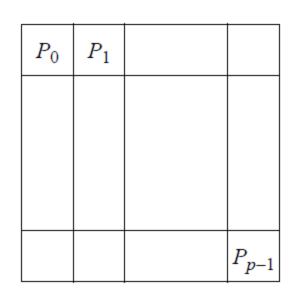
Partitioning

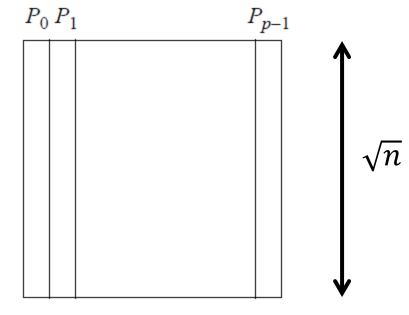
■ Blocks: Block partitioning requires lower bandwidth but higher latency

$$t_{comm} = 8(t_{startup} + \sqrt{n/p} \ t_{data})$$

Strips:

$$t_{comm} = 4(t_{startup} + \sqrt{n} \ t_{data})$$





 $\sqrt{n/p}$ \uparrow

Blocks

Strips (columns)

Partitioning

Blocks partition has a larger communication time (worse) than strip partition if:

$$\begin{split} &8 \Big(t_{startup} + \sqrt{n/p} \ t_{data} \Big) > 4 \Big(t_{startup} + \sqrt{n} \ t_{data} \Big) \\ &\Rightarrow t_{startup} > \sqrt{n} (1 - 2/\sqrt{p}) t_{data} \end{split}$$

