

## Performance and Scalability of Parallel Systems

20/11/03

Parimala Thulasiraman

1

## Parallel Systems

- Sequential Algorithm: Execution time expressed as a function of the size of its input.
- 

20/11/03

Parimala Thulasiraman

2

## Run Time

- Serial Runtime: time elapsed between beginning and end of its execution on a sequential computer ( $T_s$ )
- Parallel Runtime: the moment that a parallel computation starts to the moment that the last processor finishes execution ( $T_p$ )

20/11/03

Parimala Thulasiraman

5

## Speedup

- In evaluating a parallel system, interested in knowing how much performance gain is achieved by parallelizing a given application over a sequential implementation.
- Speedup (S): A measure that captures the relative benefit of solving a problem in parallel.

20/11/03

Parimala Thulasiraman

6

## Speedup

- Serial Algorithm:
  - More than one to sequential algorithm
  - Not all suitable for parallelization
  - Serial computer: use the algorithm that solves the problem in the least amount of time.
  - Given parallel algorithm, fair to judge its performance w.r.t the fastest sequential algorithm on a single proc.
  - Best sequential algorithm = fastest sequential algorithm
- Definition: Ratio of the serial time taken of the best sequential algorithm to solve a problem on a single processor to the time required to solve the same problem on a parallel computer with p processors.

20/11/03

Parimala Thulasiraman

7

## Speedup

- Absolute Speedup :  $S = \frac{T_s}{T_p}$
- Execution time of best sequential algorithm/Execution time on p processors
- Relative Speedup:  $S = \frac{T_1}{T_p}$
- Execution time on 1 processor/Execution on p processors

20/11/03

Parimala Thulasiraman

8

## Theoretical analysis

- Speedup = number of computational steps on 1 processor/number of computational steps with p processors.

20/11/03

Parimala Thulasiraman

9

## Example

- Adding  $n$  numbers on an  $n$  processor hypercube
  - Assume  $n$  numbers,  $n$  processors
  - At the end of the computation one of the processor holds the sum
  - $n$  is a power of 2.

20/11/03

Parimala Thulasiraman

10

## Example

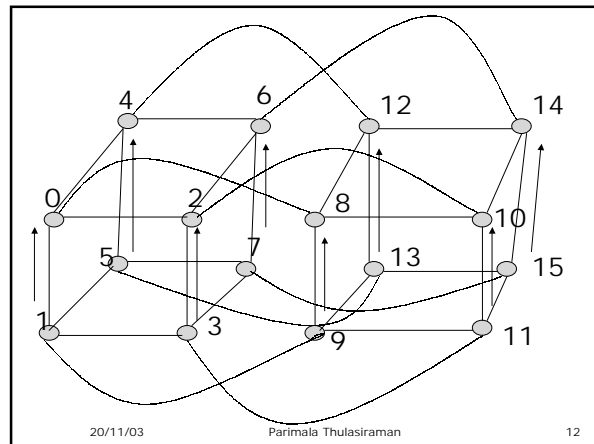
- Consists of 1 addition & 1 communication
- Processors that communicate are directly connected with each other in a hypercube; labels differ in one bit position
- Addition and communication constant amount of time

$$T_p = \Theta(\log n)$$

20/11/03

Parimala Thulasiraman

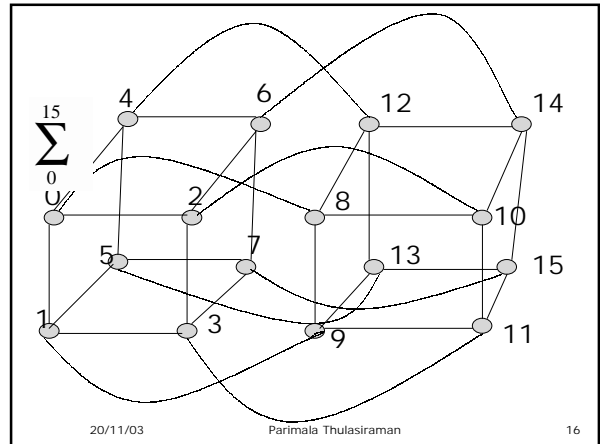
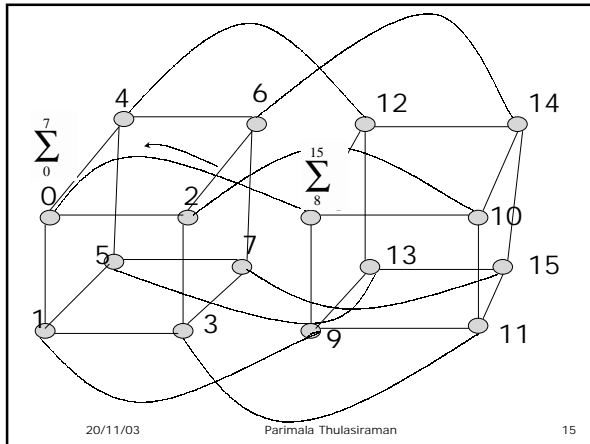
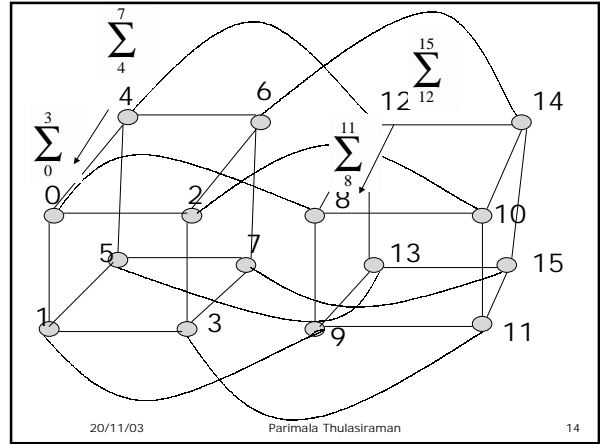
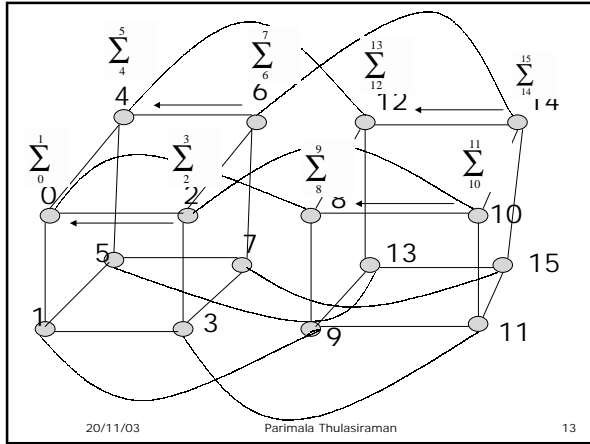
11



20/11/03

Parimala Thulasiraman

12



## Example

- Addition and communication constant amount of time

$$T_p = \Theta(\log n)$$

20/11/03

Parimala Thulasiraman

17

## Example

- Single processor :  $T_s = \Theta(n)$

- $S = \frac{T_s}{T_p} = \Theta\left(\frac{n}{\log n}\right)$

Theoretically, speedup cannot exceed the number of processors  $P$ .  $T_s$  Units of time to solve a sequential algorithm. A speedup of  $P$  can be obtained on

$P$  processors if none of the processors spends more than  $T_s / P$  time.

20/11/03

Parimala Thulasiraman

18

## Superlinear Speedup

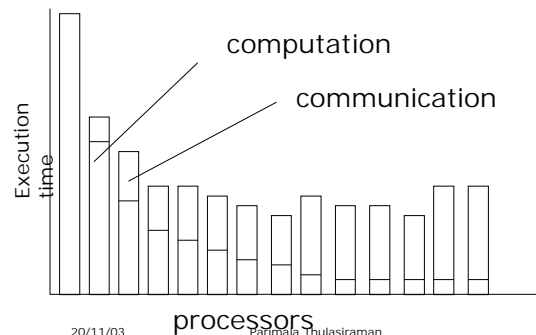
- In some cases, a speedup greater than  $P$  is observed.
  - Due to non-optimal sequential algorithm
  - H/W characteristics that put the sequential algorithm at a disadvantage.
    - Data too big to fit in memory-use of secondary storage degrades performance
    - Partition data among processor, much easier to fit into local memory

20/11/03

Parimala Thulasiraman

19

For a fixed problem size there is an optimum Number of processors that minimizes Overall execution time



20/11/03

Parimala Thulasiraman

20

## Efficiency

- Only an ideal parallel system with P processors can deliver speedup = P.
- Processors cannot devote their entire time to computation (communication, synchronization, idling)
- Efficiency: A measure of the fraction of time for which the processor is usefully employed; Ratio of speedup to the number of processors

20/11/03

Parimala Thulasiraman

21

## Efficiency

- Ideal: speedup = P, efficiency = 1
- Practice: speedup < P; 0 < E < 1

20/11/03

Parimala Thulasiraman

22

## Efficiency

- Example:

$$\frac{S}{P} = \frac{n / \log n}{n} = \frac{1}{\log n}$$

$$\therefore E = \Theta\left(\frac{1}{\log n}\right)$$

20/11/03

Parimala Thulasiraman

23

## Cost

- Cost = Parallel Runtime \* # of processors

$$T_p \times P$$

- Reflects the sum of the time each processor spends solving the problem

$$E = S / P = \frac{T_s / T_p}{P} = \frac{T_s}{T_p \times P} = \frac{\text{time seq. alg.}}{\text{cost}}$$

$$\text{Cost} = T_s / E$$

20/11/03

Parimala Thulasiraman

24

## Cost Optimal

- Cost of solving a problem on a  $||1$  computer is proportional to the execution time of the fastest sequential algorithm on a single processor– Cost Optimal
- Cost optimal system has an efficiency of  $\Theta(1)$
- Cost is referred to as Work or processor time product
- Example not cost optimal  $\therefore E = \Theta\left(\frac{1}{\log n}\right)$   
 $cost = PT_p = n \log n$   
 $T_s = n$

20/11/03

Parimala Thulasiraman

25

## Granularity

- In the previous example we used as many processors as there are inputs.
- This is too excessive.
- In practice, we assign larger pieces of input data to processors.

20/11/03

Parimala Thulasiraman

26

## Scaling Down

- Using fewer than the maximum possible number of processors to execute a parallel algorithm is called scaling down
- $P < n$
- Mapping data appropriately to maintain cost optimality

20/11/03

Parimala Thulasiraman

27

## Scalability

- Adding  $n$  numbers on  $p$  processors hypercube
- 1 unit time to add and 1 unit time to communicate
- $n/P$  elements per processor
- Addition :  $n/P-1$  time locally
- $P$  partial sums added in  $\log P$  steps, one addition, one communication

20/11/03

Parimala Thulasiraman

28

Input = 16, P = 4

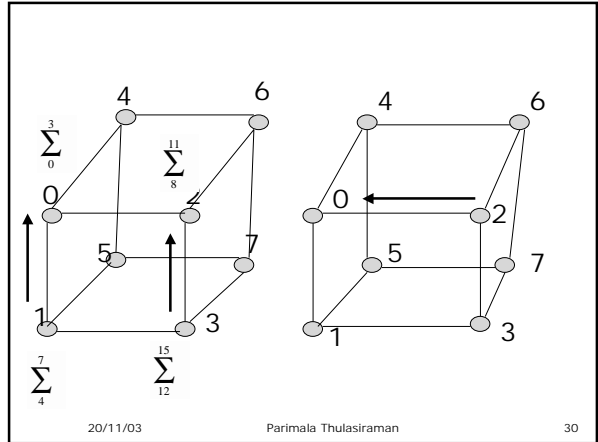
3	7	11	15
2	6	10	14
1	5	9	13
0	4	8	12
P0	P1	P2	P3

Each processor locally  
Adds its  $n/P$  elements in  
 $\Theta(n/P)$  time.  
Add P partial sums on P  
processors

20/11/03

Parimala Thulasiraman

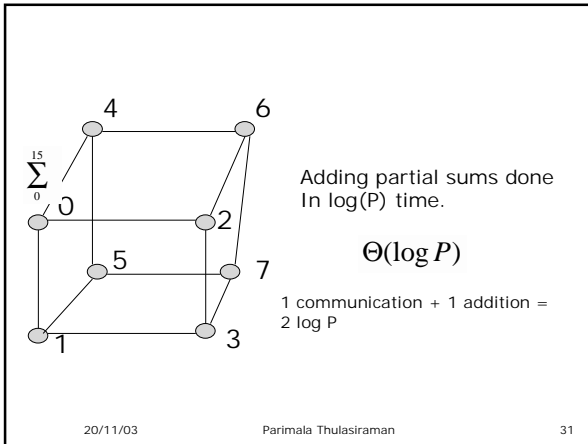
29



20/11/03

Parimala Thulasiraman

30



20/11/03

Parimala Thulasiraman

31

Compute time

$$T_p = (n/P + \log P)$$

$$\text{Cost} = T_p \times P = n + P \log P$$

Serial time =  $n-1$   
 $S = n/(n/P + \log P) = nP/(n + P \log P)$

$$E = S/P = n/(n + P \log P)$$

20/11/03

Parimala Thulasiraman

32

(See paper)

The expressions can be used to calculate the speedup and efficiency of any pair of  $n$  and  $P$ .

Observation:

1. Speedup does not increase linearly as the number of processors increases  
Speedup becomes saturated, speedup curve flattens (Amdahl's law)
  2. Efficiency drops with an increasing number of processors
  3. Efficiency increases with  $n$  and same number of processors
- Lots of Parallel systems exhibit this phenomena

20/11/03

Parimala Thulasiraman

33

## Scalability

- Increasing number of processors reduces efficiency
- Increasing size of input, increases efficiency
- Increasing size of input and number of processors simultaneously keeps efficiency fixed.
- $N=64, P=4, E = 0.80$
- $N=192, P=8, E = 0.80$
- $N=512, p=16, E=0.80$ 
  - The ability to maintain efficiency at a fixed value by simultaneously increasing the number of processors and size of the problem is exhibited by many parallel systems. These are **scalable** parallel systems

20/11/03

Parimala Thulasiraman

34

## Scalability

- Scalability: A measure of the capacity to increase speedup in proportion to number of processors
- It reflects a parallel systems ability to utilize increasing processing resources effectively.
- Isoefficiency Metric: Relate problem size to number of processors to maintain fixed efficiency.
- Question: At what rate should the problem size be increased with respect to number of processors to keep efficiency fixed? (page 5)
  - Degree of scalability

20/11/03

Parimala Thulasiraman

35

## Overhead Function

- Causes for performance loss are termed **overhead (communication, idling, contention, synchronization)**
- Total overhead or overhead function: cost that is not incurred by the fastest known serial algorithm on a sequential computer  $T_o(W, P)$
- Total time spent by all processors =

$$PT_p$$

$$T_o = PT_p - T_1$$

Time for a seq. alg. To execute on 1 processor :  $T_1$

20/11/03

Parimala Thulasiraman

36

- Problem size :
  - Total number of basic operations inherent in the problem.
  - is a function of the input size.
- W : problem size

$$T_o = PT_p - T_1$$

$$T_o = PT_p - W_t$$

20/11/03

Parimala Thulasiraman

37

## Isoefficiency Function

$$T_p = \frac{Wt_c + T_o}{P}$$

Function of problem size,  
Overhead function, # of  
processors

$$S = \frac{Wt_c}{T_p} = \frac{Wt_c P}{Wt_c + T_o}$$

$$E = \frac{S}{P} = \frac{Wt_c}{Wt_c + T_o} = \frac{1}{1 + \frac{T_o}{Wt_c}}$$

20/11/03

Parimala Thulasiraman

38

## Observation

$$E = \frac{S}{P} = \frac{1}{1 + \frac{PT_p - Wt_c}{Wt_c}}$$

$$= \frac{1}{1 + \frac{PT_p}{Wt_c} - 1}$$

$$= \frac{Wt_c}{PT_p}$$

- W is kept constant and P is increased, E decreases because total overhead increases with P.
- If W is increases, fixed P, E increases

20/11/03

Parimala Thulasiraman

39

## Isoefficiency Function

- For different parallel systems, W must be increased at different rates w.r.t P in order to maintain fixed efficiency
- If W increases linearly w.r.t P, then scalable
- If W increases exponentially as P increases, not scalable: to get good speedup for large number of processors, problem size has to be enormous

20/11/03

Parimala Thulasiraman

40

## Isoefficiency Function

- For scalable parallel systems, efficiency can be maintained at a fixed rate (between 0 and 1) if the ratio  $T_o/W$  is maintained at const. value

20/11/03

Parimala Thulasiraman

41

$$E = \frac{1}{1 + \frac{T_o}{Wt_c}}$$

$$\frac{T_o}{Wt_c} = \frac{1 - E}{E}$$

$$W = \frac{E}{1 - E} T_o t_c$$

20/11/03

Parimala Thulasiraman

42

## Isoefficiency Function

- $K = E/(1-E)$ ,  $W = KT_o$
- Growth rate of W required to keep efficiency fixed as P increases : **isoefficiency function**
- Determines the ease with which a parallel system can maintain a constant efficiency and hence achieve speedups increasing in proportion to the number of processors
- Small isoefficiency means that small increments in the problem size are sufficient for efficient utilization of an increasing number of processors. Parallel system is highly scalable.
- High isoefficiency, poorly scalable

20/11/03

Parimala Thulasiraman

43

## Example

$$T_p = \frac{n}{P} - 1 + 2 \log P$$

$$T_p = T_1 + T_o$$

$$T_o = 2 \log P$$

$$W = 2K P \log P$$

The asymptotic isoefficiency function is  $\Theta(P \log P)$

If P is increased to P', the problem size, n, must be increased by a factor of  $P' \log P' / P \log P$  to increase speedup by a factor of P'/P (to get the same efficiency as on P processors).

20/11/03

Parimala Thulasiraman

44

## Complicated Expression

$$T_o = P^{3/2} + P^{3/4}W^{3/4}$$

$$W = KP^{3/2} = \Theta(P^{3/2})$$

$$W = KP^{3/4}W^{3/4}$$

$$= K^4 P^3 = \Theta(P^3)$$

Take the higher of the 2 rates.

20/11/03

Parimala Thulasiraman

45

## Isoefficiency Funtion

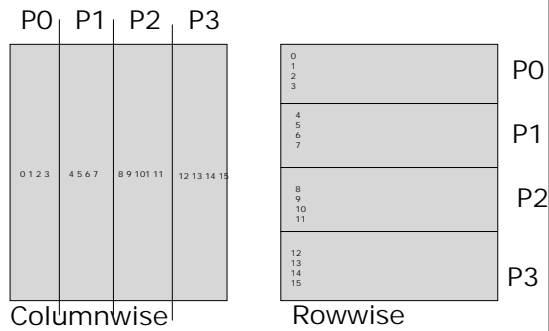
- Captures the characteristics of the parallel algorithm as well as the parallel architecture.
- Characterizes the amount of parallelism inherent in the algorithm
- Studies in h/w parameters: speed of processors and communication channels

20/11/03

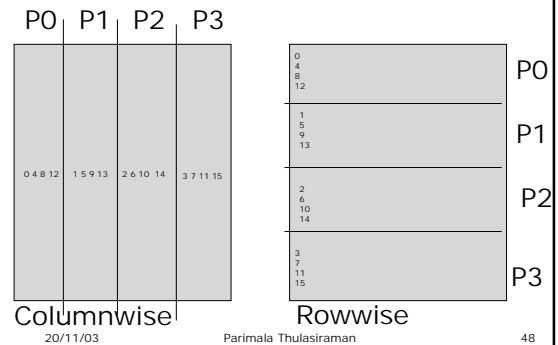
Parimala Thulasiraman

46

## Matrix-Vector Multiplication (Block-Stripping)



## Matrix-Vector Multiplication (Cyclic-Stripping)



## Matrix-Vector Multiplication

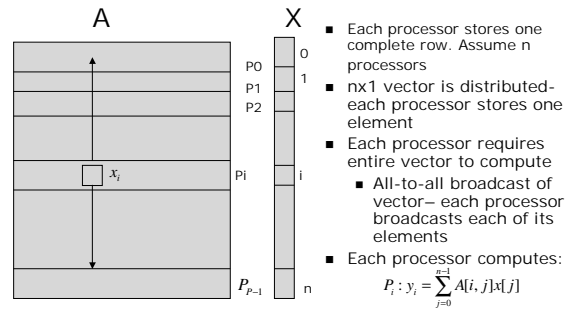
- Sequential Time :
    - Requires  $n^2$  Additions and multiplications
- $$W = n^2$$
- $$T_1 = n^2 t_c = \Theta(n^2)$$

20/11/03

Parimala Thulasiraman

49

## One row per processor –Simple case (Rowwise stripping)



20/11/03

Parimala Thulasiraman

50

## Parallel Run time

- Each processor broadcasts  $x : \Theta(n)$
- Multiplication of a single row of  $A$  with  $X$  by each processor is  $\Theta(n)$
- Addition of a single row of  $A$  with  $X$  by each processor is  $\Theta(n)$

$$T_p = \Theta(n)$$

$$\text{cost} = P \times T_p = n \times n = n^2$$

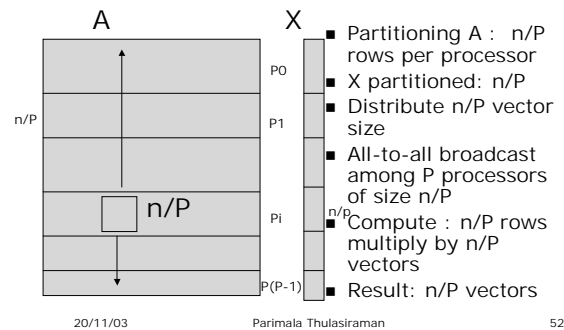
Cost optimal:  $T = \Theta(n^2)$

20/11/03

Parimala Thulasiraman

51

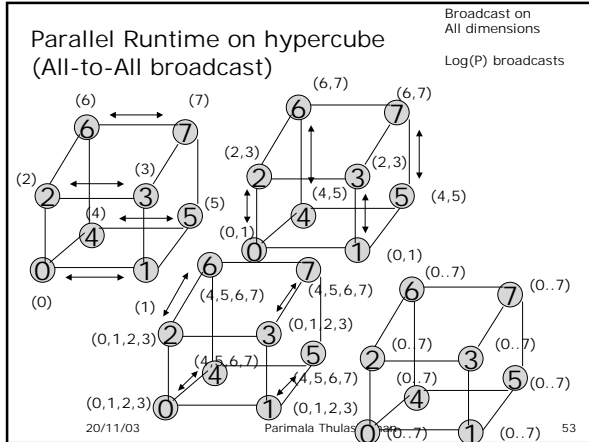
## $P < n$ processors



20/11/03

Parimala Thulasiraman

52



## All-to-All broadcast

- Log P steps.
- Communication along every dimension
- Each step data size doubles and is transmitted
- Size of message exchanged at  $i$ th step is  $2^{i-1}m$
- Time to send and receive messages  $t_s + t_w 2^{i-1}m$
- To complete the entire iteration:  $t_s$  : startup

$$T_{total} = \sum_{i=1}^{\log P} (t_s + 2^{i-1}t_w m)$$

$$= \log P t_s + m(P-1)t_w$$

$$= \log P t_s + \left(\frac{n}{P}\right)(P-1)t_w$$

$$= \log P t_s + nt_w, n \gg P$$

$t_w$  : per word transfer time

20/11/03 Parimala Thulasiraman 54

## Computation time

- An addition and multiplication takes times  $t_c$
- Each processor: dot product of  $(n/P)$  elements and vector  $n$ .

$$\therefore \frac{n^2}{P} t_c$$

20/11/03 Parimala Thulasiraman 55

$$T_p = t_c \frac{n^2}{P} + (t_s \log P + nt_w)$$

$$S = \frac{T_1}{T_p} = \frac{t_c n^2}{t_c \frac{n^2}{P} + (t_s \log P + nt_w)}$$

$$= \frac{P}{1 + P(t_s \log P + t_w n) / t_c n^2}$$

$$E = \frac{S}{P}$$

$$E = \frac{1}{1 + P(t_s \log P + t_w n) / t_c n^2}$$

20/11/03 Parimala Thulasiraman 56

## Scalability Analysis (isoefficiency function)

- Rate at which problem size needs to increase with number of processors to maintain fixed E :  $\Theta(P^2)$

$$T_o = PT_p - T_1 = n^2 t_c + t_s P \log P + t_w n P - n^2 t_c$$

$$= t_s P \log P + t_w n P$$

$$W = KT_o = K t_s P \log P$$

$$W = K t_w n P$$

$$= n^2 = K t_w n P$$

$$n = K t_w P \quad (\text{substitute into } W)$$

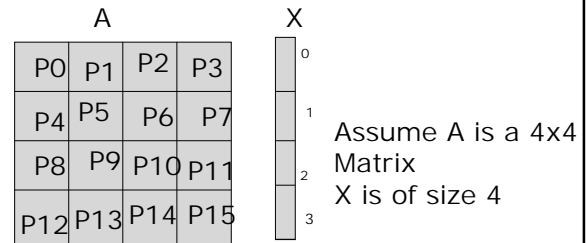
$$W = K t_w^2 P^2$$

20/11/03

Parimala Thulasiraman

57

## Block CheckerBoard Partitioning



20/11/03

Parimala Thulasiraman

58

## One element per processor

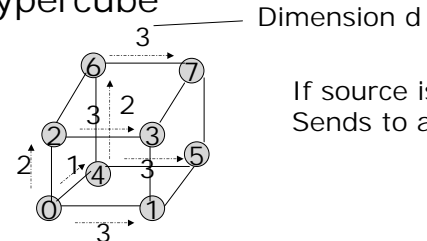
- Assume  $n^2$  processors
- Each processor stores a single element
- $i$ th element of  $x$  is available at  $i$ th element of each row in matrix
- Distribute  $x$  to the last column.
- Align the  $x$  in the diagonal
- Copy the vector from each diagonal to the other processors— $n$  simultaneous (one to all broadcast)
- Each processor multiplies each element with  $x$  element.
- Result : add each of the computed product—leave the sum in the last column of processors—**single node accumulation**

20/11/03

Parimala Thulasiraman

59

## Parallel run time on a hypercube



$$\log \sqrt{n^2} \text{ communication steps } P = n^2$$

Each  $\log n$  step takes  $t_s + t_w m$  in each dimension—for all dimensions  $\log n (t_s + t_w m)$

20/11/03

Parimala Thulasiraman

60

### 3 communication steps

- Align x along main diagonal :  $\Theta(\log n)$
- Broadcast columnwise to selected processors:  $\Theta(\log n)$
- Single node accumulation:  $\Theta(\log n)$
- Single multiplication in constant time

$\therefore T_p = \Theta(\log n)$   
 $\text{cost} = PT_p = n^2 \log n$  Not cost optimal  
 $T_1 = n^2$

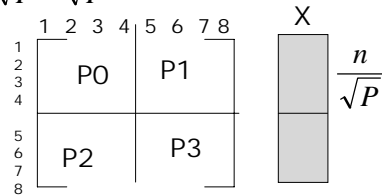
20/11/03

Parimala Thulasiraman

61

$P \ll n$

$\frac{n}{\sqrt{P}} \times \frac{n}{\sqrt{P}}$  blocks per processor



20/11/03

Parimala Thulasiraman

62

### Parallel run time

- Vector is aligned along diagonal
- One-to-all broadcast of the vector
- Each performs computation on its

$$\frac{n}{\sqrt{P}} \times \frac{n}{\sqrt{P}} = \frac{n^2}{P}$$

Data size                      vector  
 Node accumulation of the partial sums

20/11/03

Parimala Thulasiraman

63

### Parallel run time

- Vector to main diagonal from last processor in each column:

$$t_s + t_w \frac{n}{\sqrt{P}} \log \sqrt{P}$$

- One-to-all broadcast:

$$t_s \log \sqrt{P} + t_w \frac{n}{\sqrt{P}} \log \sqrt{P}$$

- Single node accumulation:

$$t_s \log \sqrt{P} + t_w \frac{n}{\sqrt{P}} \log \sqrt{P}$$

20/11/03

Parimala Thulasiraman

64

## Parallel run time

- computation:

$$T_p = t_c \frac{n^2}{P} + 2t_s \log\sqrt{P} + 3t_w \frac{n}{\sqrt{P}} \log\sqrt{P}$$

$$= t_c \frac{n^2}{P} + t_s \log P + \frac{3}{2} t_w \frac{n}{\sqrt{P}} \log P$$

20/11/03

Parimala Thulasiraman

65

## Scalability analysis

$$T_p = t_c \frac{n^2}{P} + 2t_s \log\sqrt{P} + 3t_w \frac{n}{\sqrt{P}} \log\sqrt{P}$$

$$\text{cost} = T_p P = t_c n^2 + t_s P \log P + \frac{3}{2} t_w n \sqrt{P} \log P$$

$T_o$

20/11/03

Parimala Thulasiraman

66

## Isoefficiency function

- Determine the isoefficiency as we did last time for both startup and per word time parts. (see paper for details)

$$\Theta(P \log^2 P)$$

20/11/03

Parimala Thulasiraman

67

## Parallel Overhead

- Major sources of overhead in a parallel system:
  - Interprocessor communication
  - Load imbalance
  - Extra computation
- Interprocessor communication: if P processor spend  $t_{comm}$  time performing communication, then interprocessor communication contributes  $t_{comm} \times P$  to the overhead function

20/11/03

Parimala Thulasiraman

68

## Load Imbalance

- If different processors have different work load, some processors may be idle during part of the time that others are working on the problem.
  - Statically partitioning the problem is difficult
  - Having equal workload is difficult
- Synchronization: If not all ready to synchronize, time wasted
- Idle time is an overhead
- In some applications there is a sequential part that only one processor can work on.
- $W_s$ , work due to sequential component,  $W_p$  is the work due to parallelizable component. While one processor works on  $W_s$  the remaining (p-1) are idle.
- A serial component contributes  $(P-1) W_s$  to the overhead function of a P-processor parallel system.

20/11/03

Parimala Thulasiraman

69

## Extra Computation

- Fastest sequential algorithm maybe difficult or impossible to parallelize, forcing us to use a parallel algorithm based on a poorer but easily parallelizable sequential one (with high degree of concurrency)
- Let X be the execution time of the fastest sequential alg. Let X' be the execution time of the poorer algorithm (easily parallelizable) for that problem
- $X-X'$  = overhead, extra amount of work performed to solve the problem

20/11/03

Parimala Thulasiraman

70