

RECAP

Computation →
Communication within the processor
(or cores)

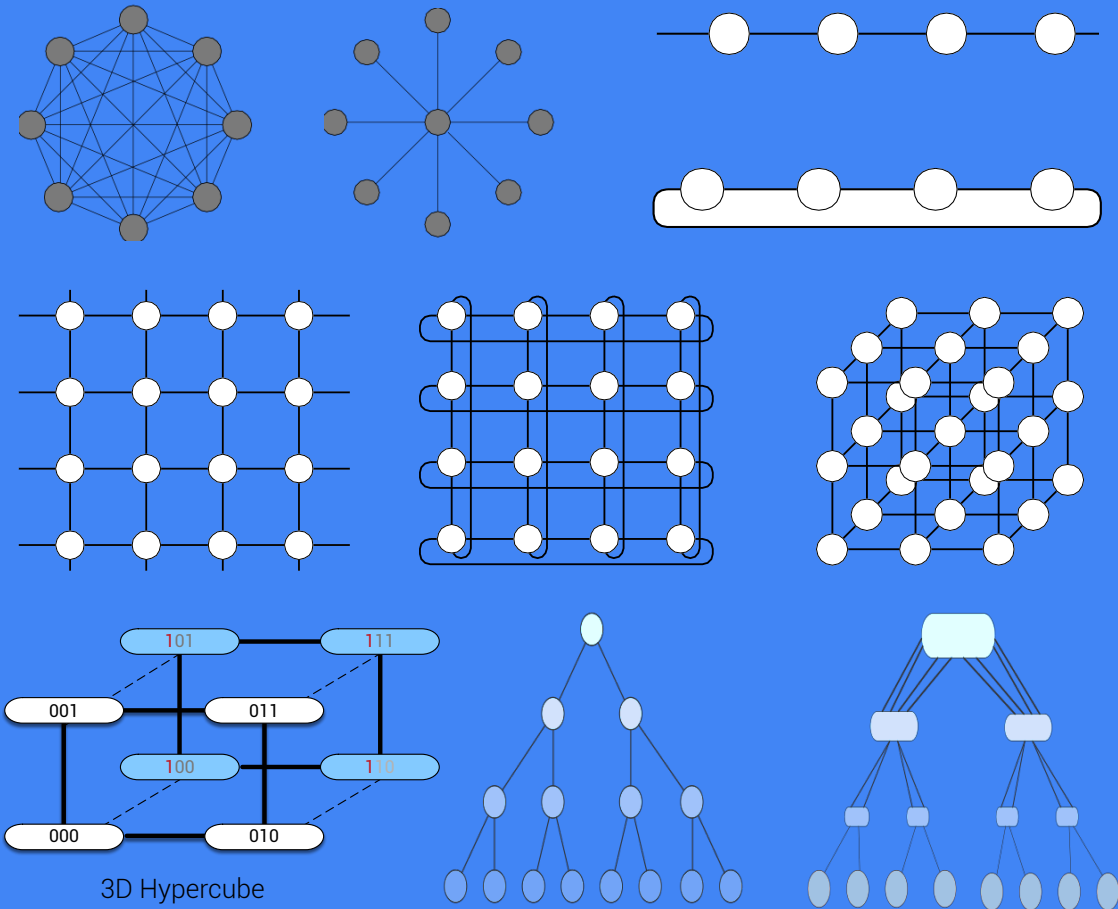
SIMD vs. MIMD (SPMD)

Buses, Crossbars

NUMA vs. UMA

RECAP

Communication (across machines) →
Network Topologies →



Degree.

Distance vs. Diameter

Connectivity vs. Bisection Width

Latency vs. Bandwidth

This Week

Measuring Performance

Basic Communication Operations

Communication Costs

Merging Networks →

Best of all involved

Student Presentation (next week)

Performance Evaluation

Basic Messaging

Serial vs. Parallel

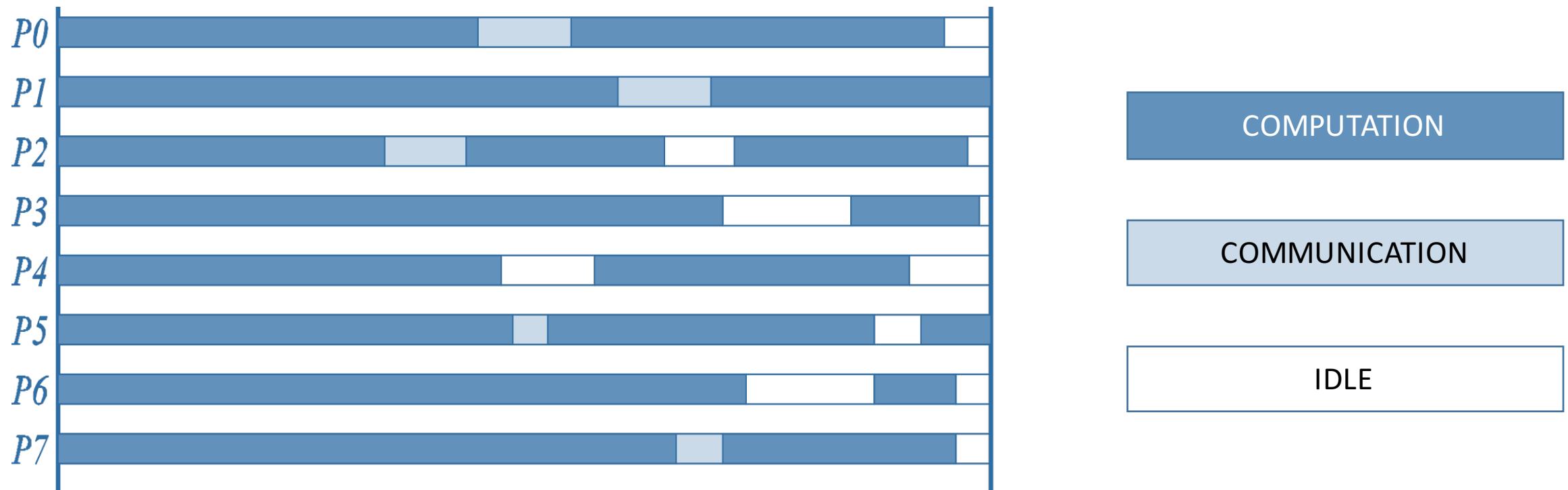
5

- A sequential algorithm is evaluated by its runtime (in general, asymptotic runtime as a function of input size).
- The **asymptotic runtime** of a sequential program is **identical on any serial platform**.
- The **parallel runtime** of a program depends on the **input size**, the number of **processors**, and the **communication parameters** of the machine.
- A parallel system → **algorithm + architecture**.

Serial vs. Parallel (How do we quantify benefit)

6

- Wall clock time - the time from the start of the first processor to the stopping time of the last processor in a parallel ensemble
- If time take by a serial code is T , and we use p processors, would a parallelized version take T/p ?



Other ways to Quantify

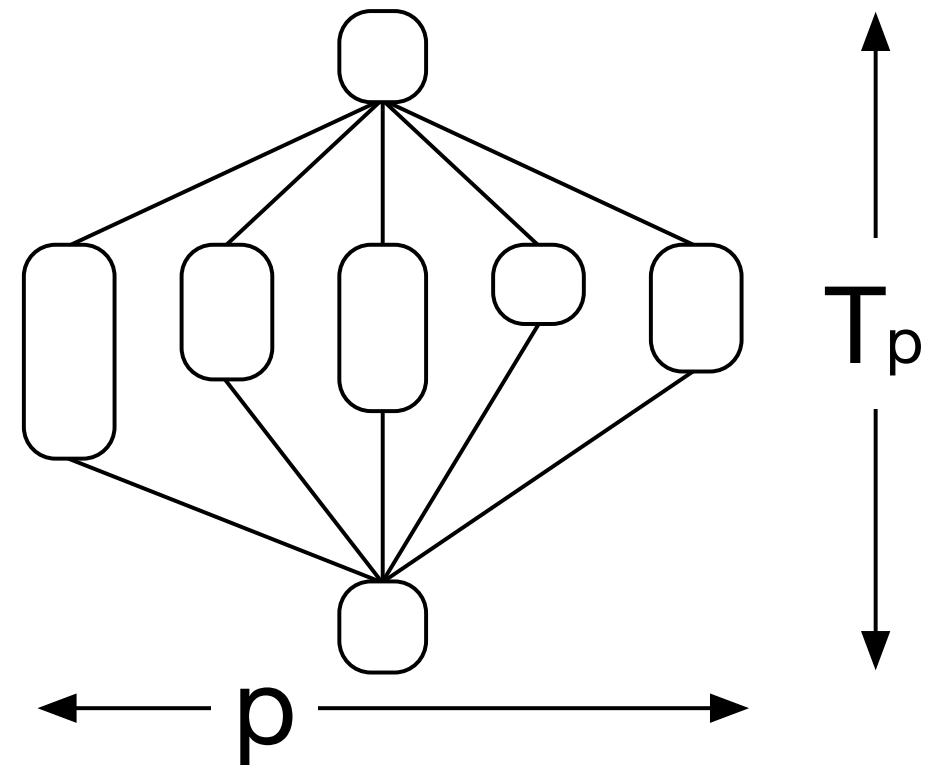
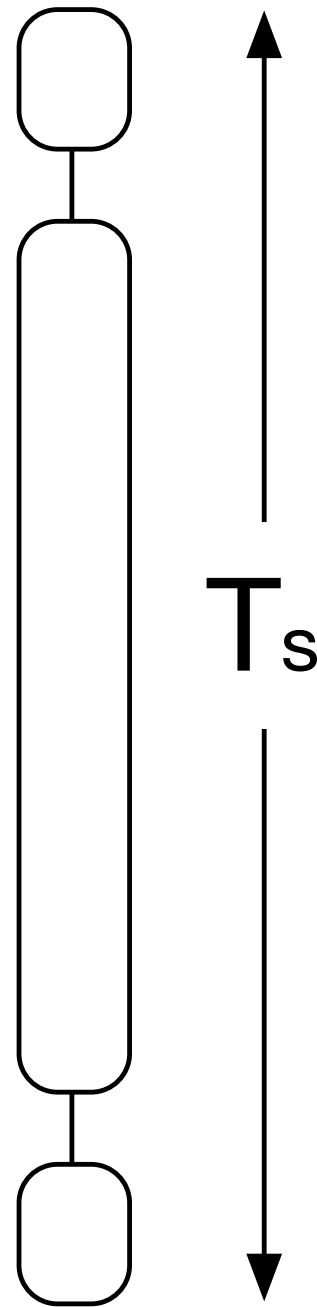
Serial $\rightarrow T_s$

Parallel $T_{ALL} \rightarrow pT_p$

Overhead

$$T_o \rightarrow T_{ALL} - T_s$$

$$T_o \rightarrow pT_p - T_s$$



Example 1

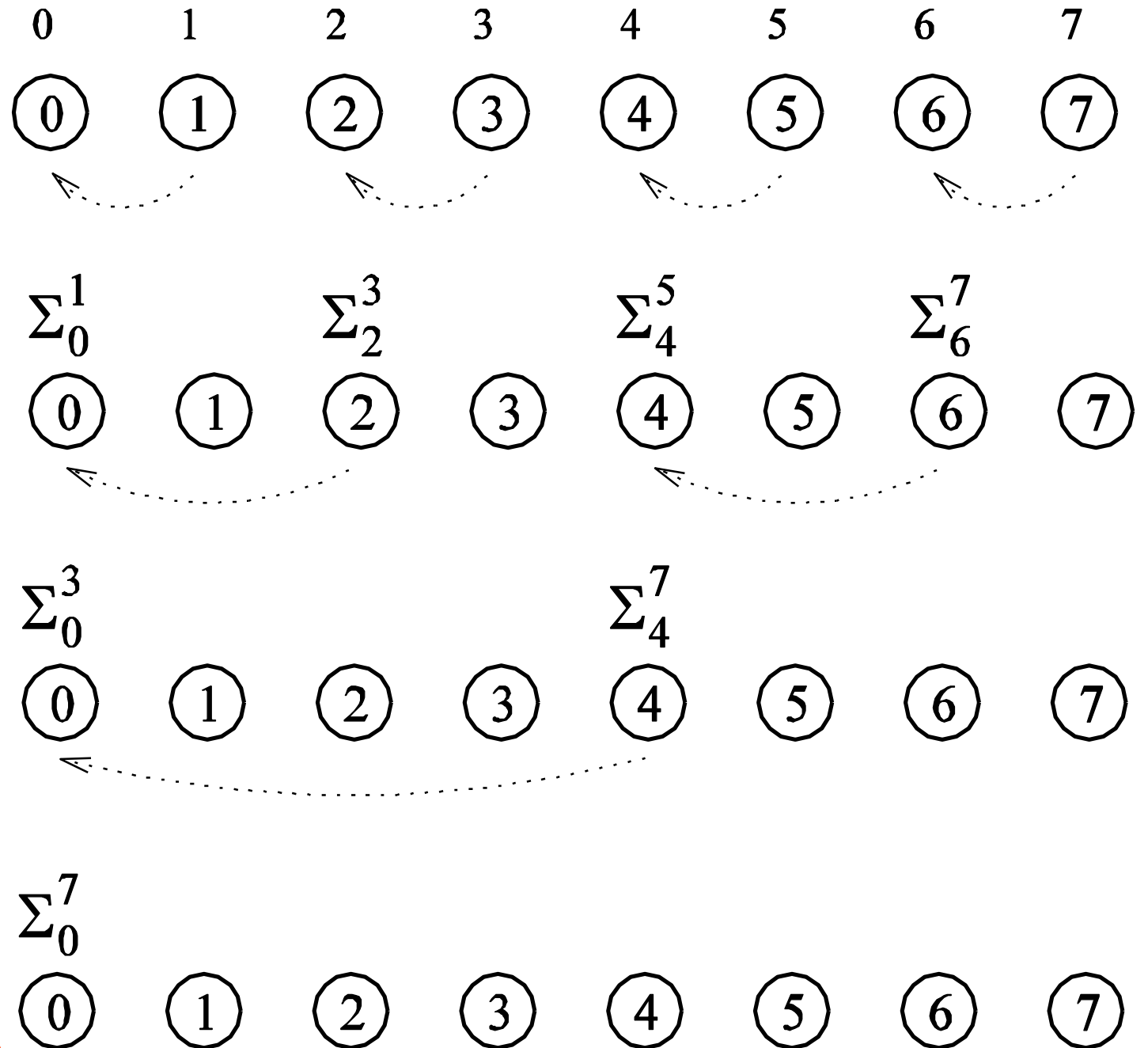
Add n numbers

$T_s \rightarrow O(n)$

How would one do
this in **log n** steps

$T_P = O(\log n)$

$S \rightarrow O(n / \log(n))$



What's the benefit: SPEEDUP

Speedup (S) is the ratio of the time taken to solve a problem on a single processor to the time required to solve the same **problem** on a parallel computer with p identical processing elements.

Speedup (Example)

10

- Consider the problem of parallel bubble sort.
- The serial time for bubble-sort is 150 seconds.
- The parallel time for odd-even sort (efficient parallelization of bubble sort) is 40 seconds.
- The speedup would appear to be $150/40 = 3.75$.
- But is this really a fair assessment of the system?
- What if serial quicksort only took 30 seconds? In this case, the speedup is $30/40 = 0.75$. This is a more realistic assessment of the system.

Speedup (Max and Min)

11

- Speedup can be as low as 0 (the parallel program never terminates).
- Max Speedup?
 - Theoretically $\rightarrow p$ (number of processors)
 - Can it be greater than p ?
 - SUPERLINEAR SPEEDUP
 - Parallel does less work than serial $\rightarrow S > p$
 - Parallel can use more resources than serial $\rightarrow S > p$

12

The diagram shows a hierarchical tree structure. The root node is labeled "Processing element 0" and "Processing element 1". It branches into two children. The left child branches into two children, one of which branches into two children. The right child branches into two children, one of which branches into two children. The diagram illustrates a hierarchical structure with 16 nodes in total, representing a 4x4 grid of processing elements.

Searching an unstructured tree for a node with a given label, 'S', on two processing elements using depth-first traversal. The two-processor version with processor 0 searching the left subtree and processor 1 searching the right subtree expands only the shaded nodes before the solution is found. The corresponding serial formulation expands the entire tree. It is clear that the serial algorithm does more work than the parallel algorithm.

Efficiency

13

- Speedup \rightarrow Serial / Parallel
- Efficiency \rightarrow Are you keeping your processors busy?

$$E = \left(\frac{\textit{Speedup}}{\textit{processors}} \right) \rightarrow \left(\frac{T_s / T_p}{p} \right) \rightarrow \left(\frac{T_s}{p T_p} \right)$$

What is the Efficiency of the number adder?

$$S \rightarrow \frac{Ts}{Tp}$$

Efficiency = S / p
In adder $n = p$

$$\rightarrow \Theta\left(\frac{n/\log n}{n}\right)$$

$$\rightarrow \Theta\left(\frac{1}{\log n}\right)$$

Cost

15

- Speedup \rightarrow Serial / Parallel
- Efficiency \rightarrow Speedup / Number of Processors

$$\text{COST} \rightarrow p \times T_p$$

- Cost reflects the sum of the time that each processing element spends solving the problem.
- A parallel system is said to be *cost-optimal* if the cost of solving a problem on a parallel computer is asymptotically identical to serial cost.

Since $\mathbf{E} = \mathbf{T_s} / \mathbf{p T_p}$, for cost optimal systems, $\mathbf{E} = O(1)$.

Cost of a Parallel System: Example

16

Consider the problem of adding numbers

We have, $T_p = \log n$ (for $p = n$).

The cost of this system is given by $p T_p = n \log n$.

Since the serial runtime of this operation is $\Theta(n)$, the algorithm is not cost optimal.

Consider a parallel sorting algorithm that uses n processing elements to sort the list in time $(\log n)^2$.

- Since the serial runtime of a (comparison-based) sort is $T_s = n \log n$, the speedup and efficiency of this algorithm are given by $n / \log n$ and $1 / \log n$, respectively.
- COST of this algorithm is $n (\log n)^2$.
- This algorithm is not cost optimal but only by a factor of $\log n$.
- If $p < n$, assigning n tasks to p processors gives $T_p = n (\log n)^2 / p$.
- The corresponding speedup of this formulation is $p / \log n$.
- This speedup goes down as the problem size n is increased for a given p !

Non-Optimal Cost

18

- Assume $n = 1024$ and $p = 32$ processors
- Speedup $\rightarrow p / \log(n) \rightarrow 32/10 = 3.2$
- If $n = 10^6 \rightarrow S = 32/\sim 20 = 1.6$
- If $n = 10^9 \rightarrow S = 32/\sim 30 = \sim 1.0$

SIGNIFICANT COSTS WITH NOT BEING COST OPTIMAL

Cost of solving in parallel has the same asymptotic growth as the serial version

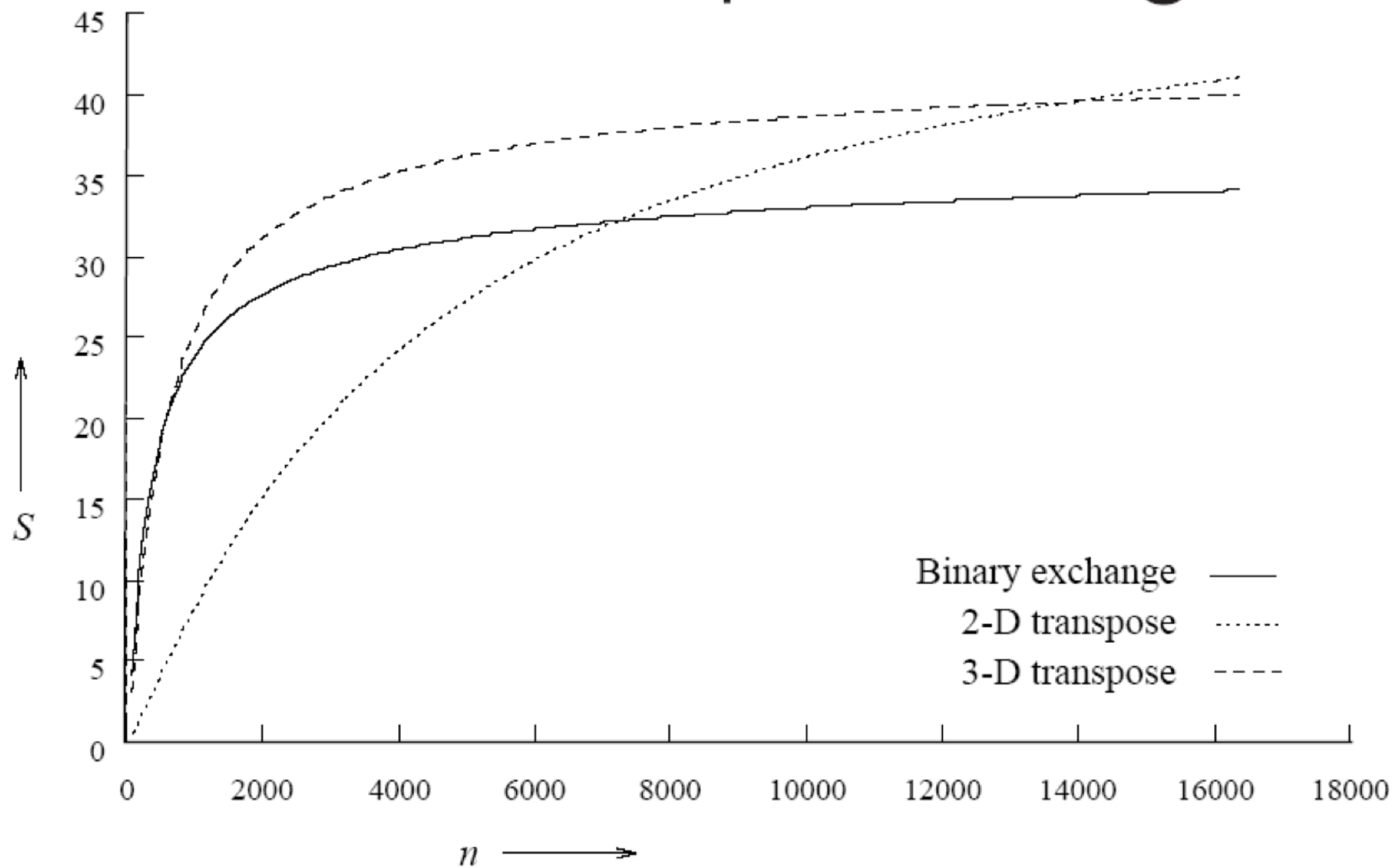
Granularity

- The adder example was not optimal
 - One processor per element
- In practice: assign large pieces to each element → Increase GRANULARITY

Scaling Characteristics of Parallel Programs

20

- For a given problem size (i.e., the value of T_s remains constant), as we increase the number of processing elements, T_o increases.
- The overall efficiency of the parallel program goes down.
This is the case for all parallel programs.



Scaling Characteristics of Parallel Programs

22

- The efficiency of a parallel program can be written as:

$$E = \frac{S}{p} = \frac{T_S}{pT_P}$$

or

$$E = \frac{1}{1 + \frac{T_o}{T_S}} \quad (4)$$

- The total overhead function T_o is an increasing function of p .

Adder Example (n elements, p processors)

23

- Assume a unit time to add two numbers
- First Local summation $\rightarrow n / p$ time
- Communication $\rightarrow \log p \rightarrow \rightarrow$ subsequent summation $\rightarrow \log p$

- $T_p = n / p + 2 \log p$

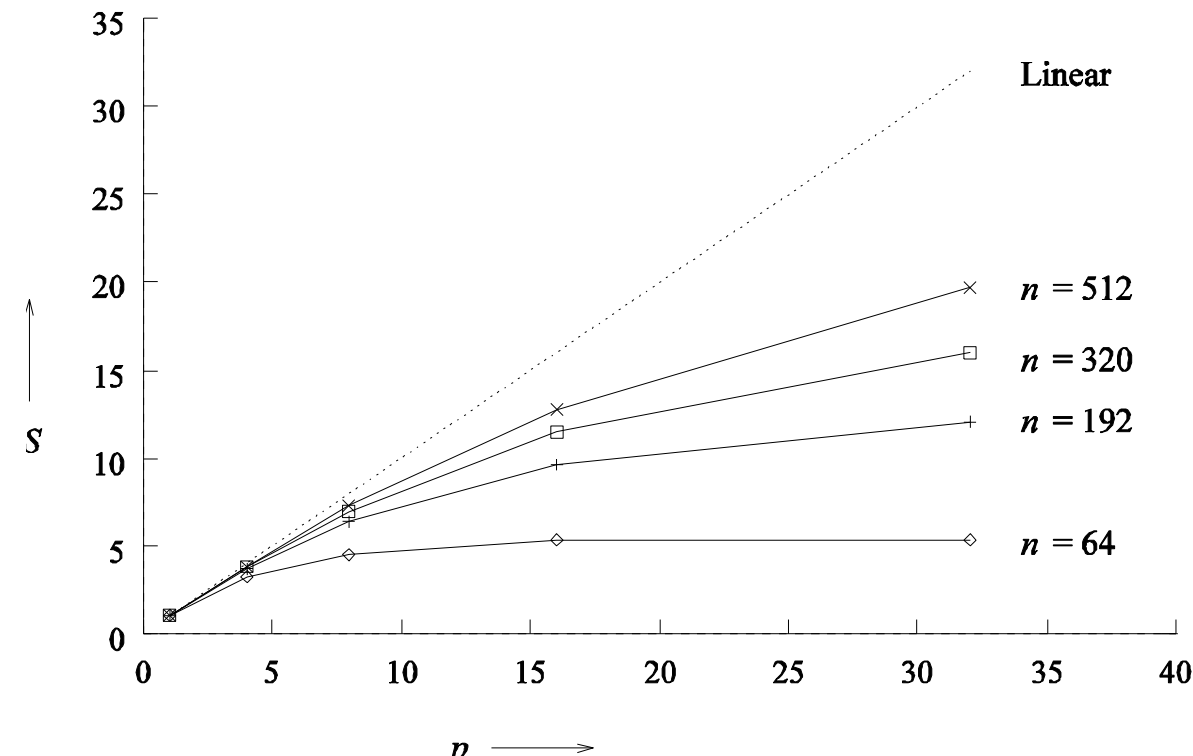
- $S = T_s / T_p \rightarrow \frac{n}{\frac{n}{p} + 2 \log p}$

- $E = S / p \rightarrow \frac{1}{1 + \frac{2p \log p}{n}}$

Scaling Characteristics of Parallel Programs: Example (continued)

24

Plotting the speedup for various input sizes gives us:



Speedup versus the number of processing elements for adding a list of numbers.
Speedup tends to saturate and efficiency drops as a consequence of Amdahl's law

Weak vs. Strong Scaling

25

- Say you have a task X (matrix multiplication)
- **WEAK** → Increase the size of the problem but keep the resource fixed.
Size/Resource is fixed.
Bigger Matrix → Same Processor
Gives you an idea on how your algorithm behaves with fixed size/core
- **STRONG** → Do the job faster but keep the problem fixed.
Size/Resource varies.
Faster Multiplication
Gives you an idea on how your problem behaves with a fixed total size

Basic Communication Operations: Introduction

- Many interactions in practical parallel programs occur in well-defined patterns involving groups of processors.
- Efficient implementations of these operations can improve performance, reduce development effort and cost, and improve software quality.
- Efficient implementations must leverage underlying architecture. For this reason, we refer to specific architectures here.
- We select a descriptive set of architectures to illustrate the process of algorithm design.

Basic Communication Operations: Introduction

27

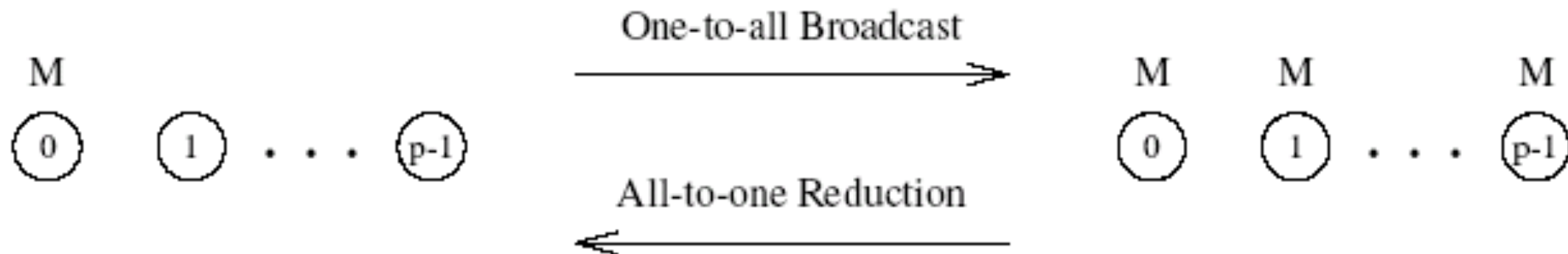
- Group communication operations are built using point-to-point messaging primitives.
- We assume that the network is bidirectional and that communication is single-ported.

One-to-All Broadcast and All-to-One Reduction

- One processor has a piece of data (of size m) it needs to send to everyone.
- The dual of one-to-all broadcast is *all-to-one reduction*.
- In all-to-one reduction, each processor has m units of data. These data items must be combined piece-wise (using some associative operator, such as addition or min), and the result made available at a target processor.

One-to-All Broadcast and All-to-One Reduction on Rings

- Simplest way is to send $p-1$ messages from the source to the other $p-1$ processors - this is not very efficient.
- Use recursive doubling: source sends a message to a selected processor. We now have two independent problems defined over halves of machines.
- Reduction can be performed in an identical fashion by inverting the process.



Next

Messaging
Algorithms and
their analysis

