

# Theoretical Basics

## Week 1 – Part 2

**Amittai Aviram – 2020-09-08 – Boston University**

# Speedup

Let  $p$  be the number of processing units,  $T(p)$  be the amount of time taken by  $p$  processors to accomplish a constant unit of work, and  $S(p)$  be the *speedup* provided by parallel processing on  $p$  processors as compared with use of a single processor. Then the speedup is given by

$$S = \frac{T(1)}{T(p)}$$

The maximal speedup is usually *linear*.

# Efficiency and Cost

The *efficiency* of a parallel computing system is how much speedup you get per processing unit—given as a percent.

$$E = \frac{S}{p} = \frac{T(1)}{T(p) \times p}$$

The ideal is 100%, which is linear speedup.

The *cost* of the system is measured in *time*—the total time taken by all processors.

$$C = T(p) \times p$$

# Computation and Communication

- The cost of communicating data can significantly affect the efficiency of the system.
- The higher the ratio of communication to computation time, the less efficient the system.

# Scalability

## The Maintenance of Efficiency as $p$ Increases

- Strong Scalability  
The algorithm or system remains efficient as we add processors while keeping the data input size constant.
- Weak Scalability  
The algorithm or system remains efficient as both the number of processors *and* the size of the input grow.
- An algorithm or system is *not scalable* if its performance degrades significantly as the size of the system (and input) grows.

# Example

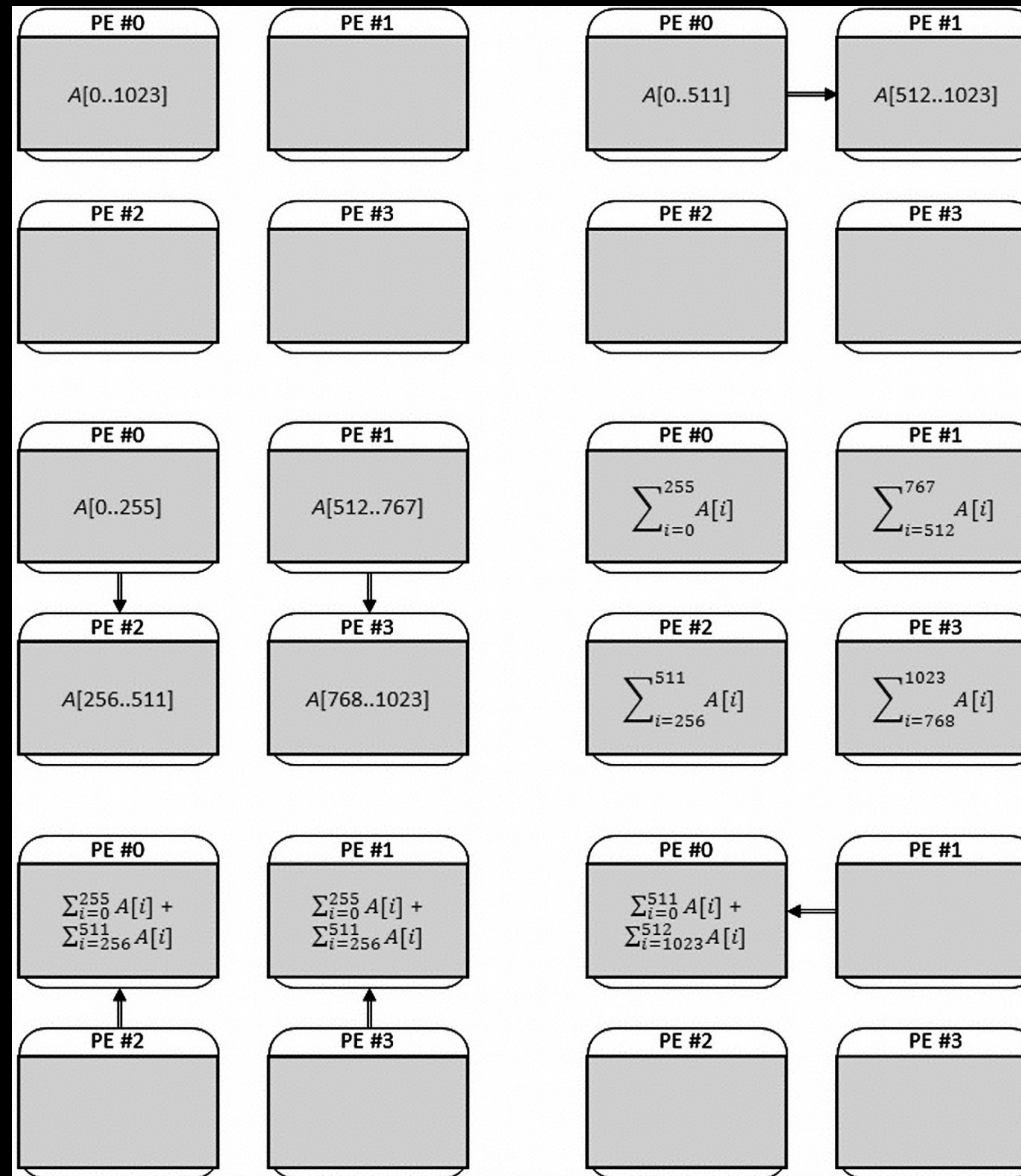
## Summing 1024 Numbers

- Given: a distributed system of interconnected processors.
- We hold the input size constant so as to measure *strong scalability*.
- Assumptions:
  - Each processor takes 1 time unit to add two numbers.
  - Each communication between processors takes 3 time units.
  - We have  $p = 2^k$  processors (for some small integer  $k$ ).

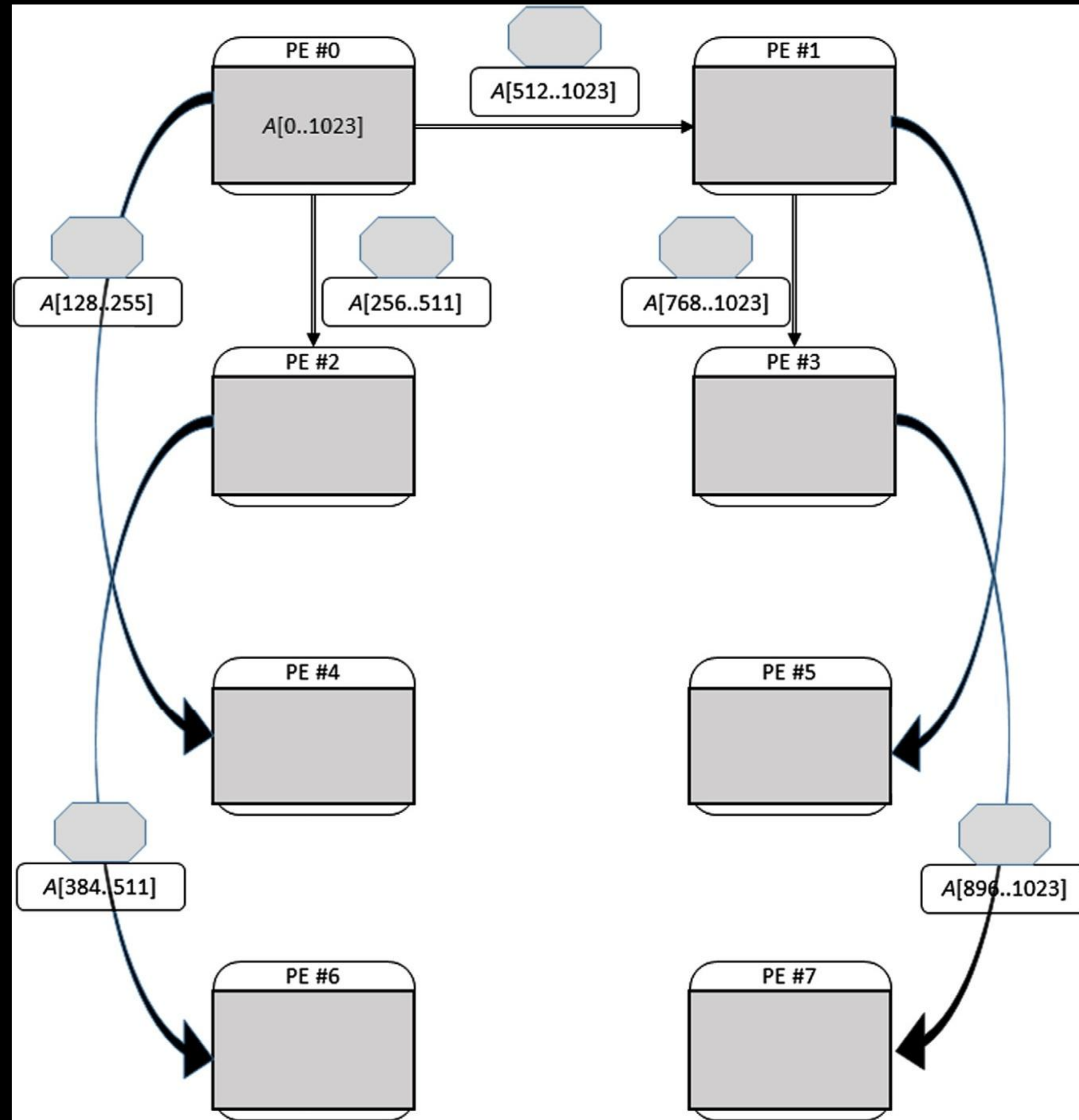
# Algorithm

## Distributed Divide-and-Conquer

- The first node divides the input list in half and sends one half to its neighbor.
- This division is repeated recursively until each of the processors receives an even share of the input.
  - E.g., for  $p = 8 = 2^3$  processors, node 0 divides its input in half three times, node 1 divides its input (from node 0) twice, etc.
- After dividing up the input, the nodes must gather the results recursively back from their neighbors until the total winds up in node 0.







# Cost Comparisons

Example with  $n = 1024$

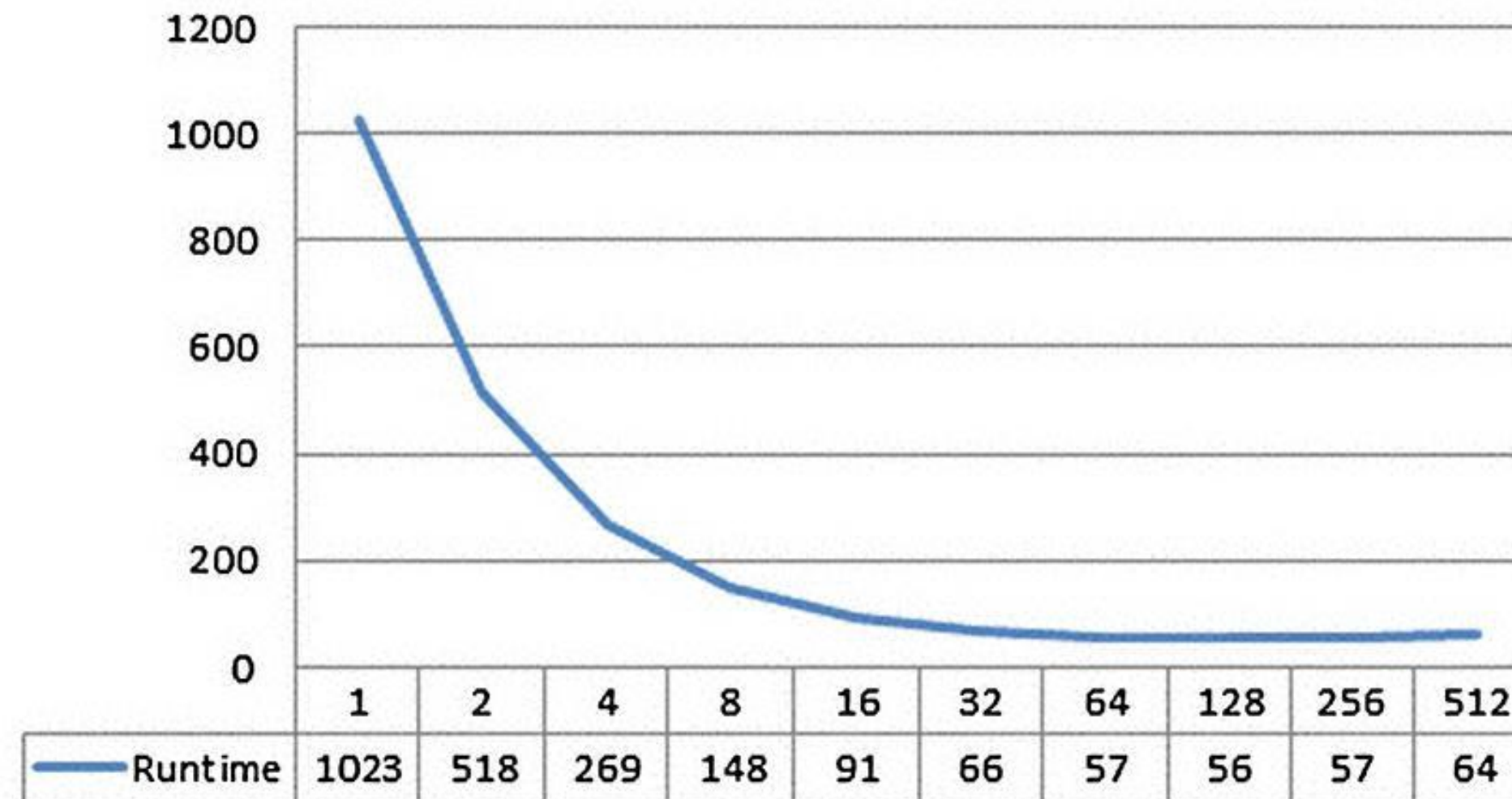
- $p = 1: T(1,n) = n - 1$   
 $T = 1023$
- $p = 2: T(2,n) = 3 + n/2 - 1 + 3 + 1$   
 $T = 517 \quad S = 1023/517 = 1.98 \quad E = 1023/(517 \times 2) = 99 \%$
- $p = 4: T(4,n) = 3 + 3 + n/4 - 1 + 3 + 1 + 3 + 1$   
 $T = 269 \quad S = 1023/269 = 3.80 \quad E = 1023/(269 \times 4) = 95 \%$
- $p = 8: T(8,n) = 3 + 3 + 3 + n/8 - 1 + 3 + 1 + 3 + 1 + 3 + 1$   
 $T = 148 \quad S = 1023/148 = 6.91 \quad E = 1023/(148 \times 8) = 86 \%$

# Generalizing

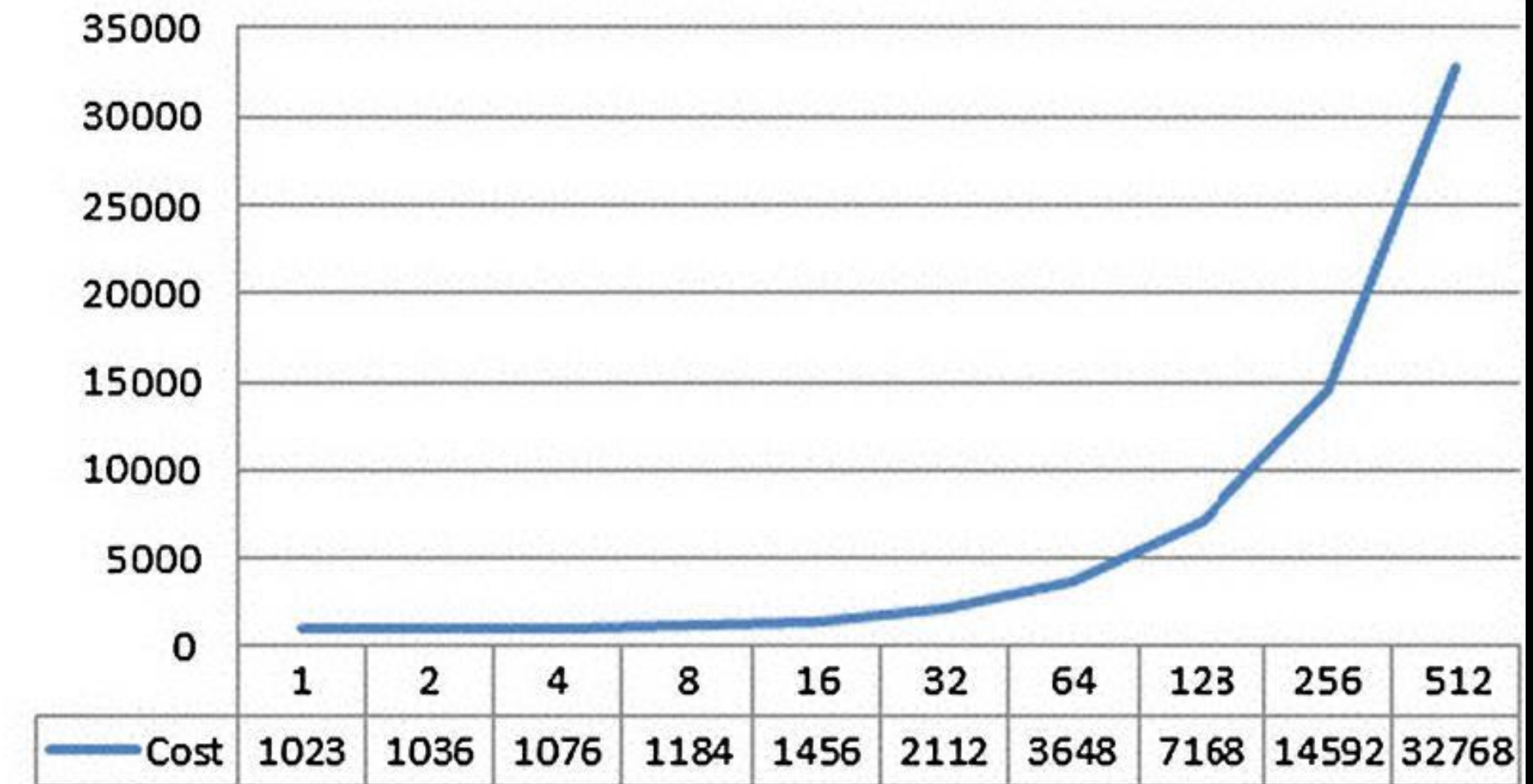
- For  $p = 2^q$  processors and  $n = 2^k$  input numbers
- Components
  - Data distribution time:  $3 \times q$
  - Local sum computation time:  $n/p - 1 = 2^{k-q} - 1$
  - Time to collect partial results:  $3 \times q$
  - Time to add partial results locally:  $q$
- $T(p, n) = T(2^q, 2^k) = 3q + 2^{k-q} - 1 + 3q + q = 2^{q-k} - 1 + 7q$



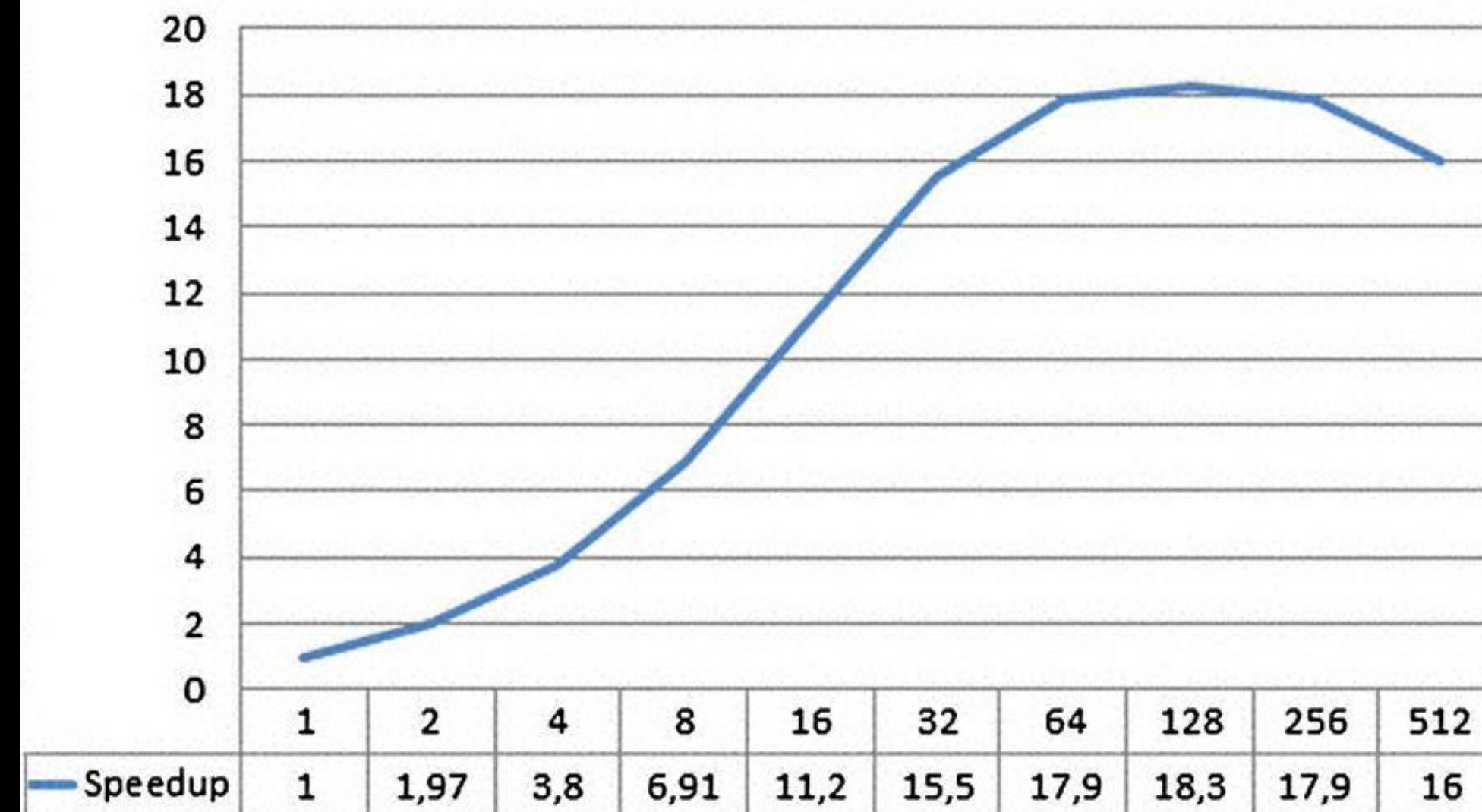
### Runtime



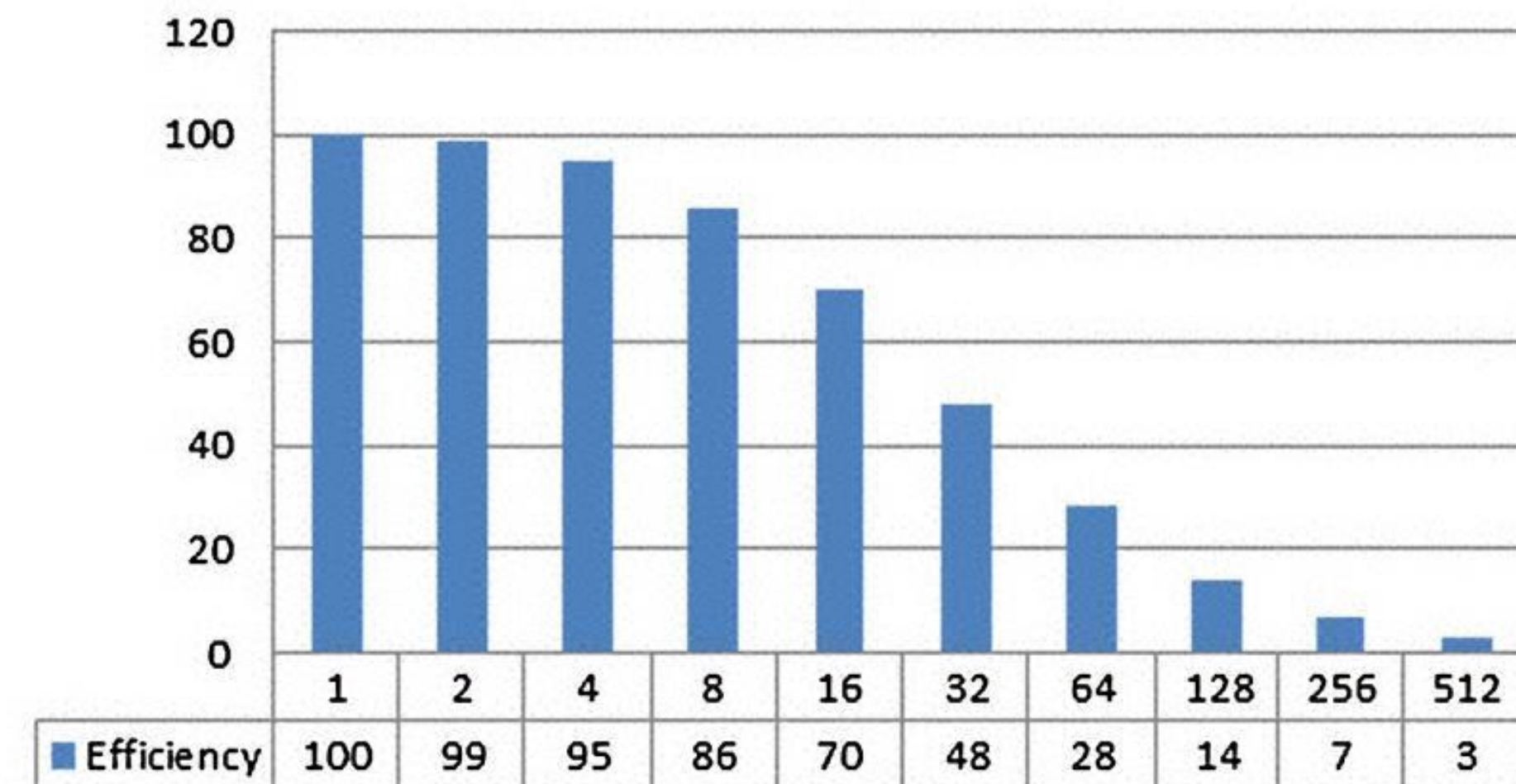
### Cost



### Speedup



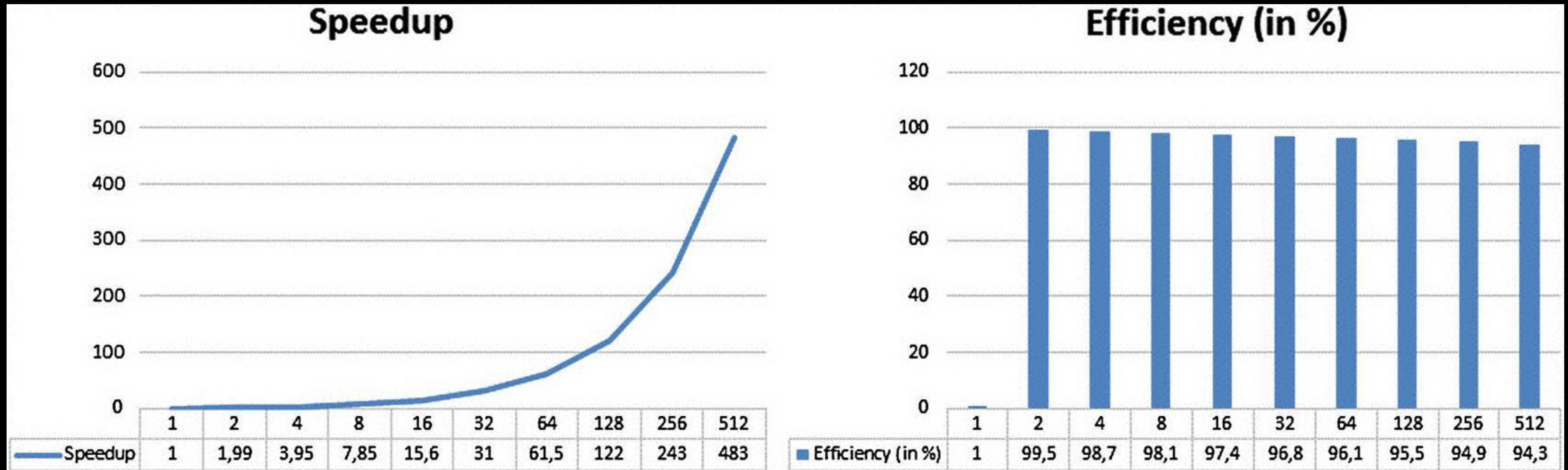
### Efficiency (in %)





# Weak Scalability Analysis

For  $n = 1024 \times p$  and  $1 \leq p \leq 512$



# Conclusions

- Our algorithm *is not* strongly scalable.
- BUT our algorithm *is* weakly scalable.

# Looking Ahead

## Options in Parallel Programming and Systems

- Data distribution
  - Distributed memory
  - Shared memory
  - Partitioned shared memory
- Partitioning the data
  - Managing data dependencies
- Balancing workloads