

Parallel DBs

April 23, 2018

Why Scale?

Scan of 1 PB at 300MB/s (SATA r2 Limit)

Why Scale Up?

Scan of 1 PB at 300MB/s (SATA r2 Limit)



~ 1 Hour

Why Scale Up?

Scan of 1 PB at 300MB/s (SATA r2 Limit)



Data Parallelism

Replication



Partitioning

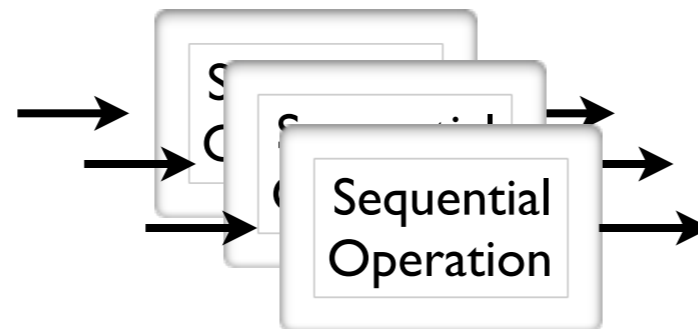


Operator Parallelism

- Pipeline Parallelism: A task breaks down into stages; each machine processes one stage.



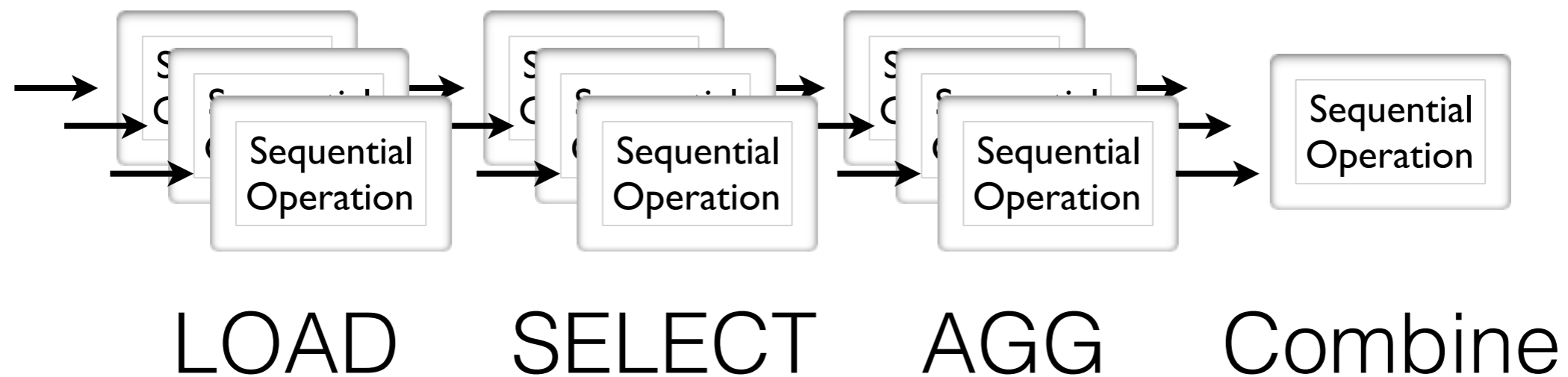
- Partition Parallelism: Many machines doing the same thing to different pieces of data.



Types of Parallelism

- Both types of parallelism are natural in a database management system.

SELECT SUM(...) FROM Table WHERE ...

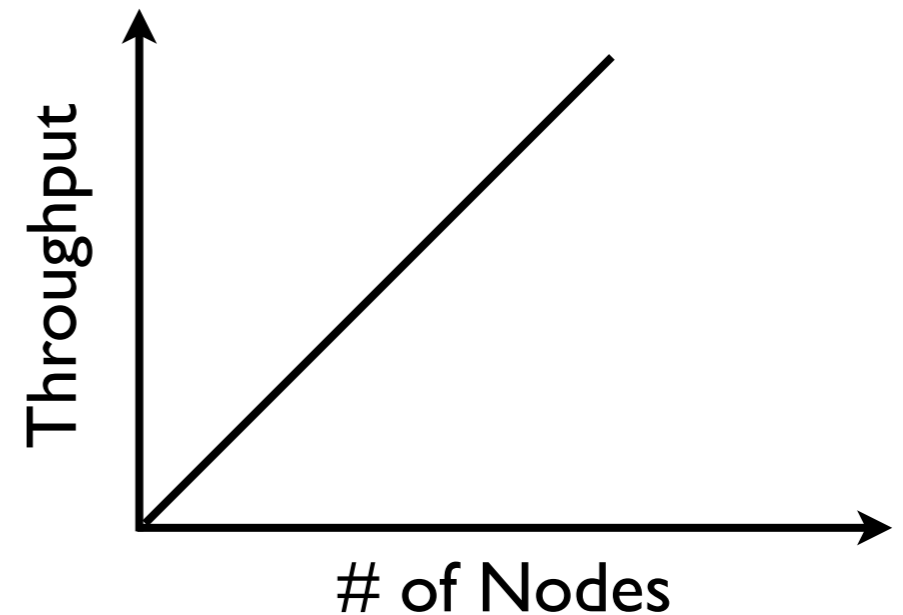
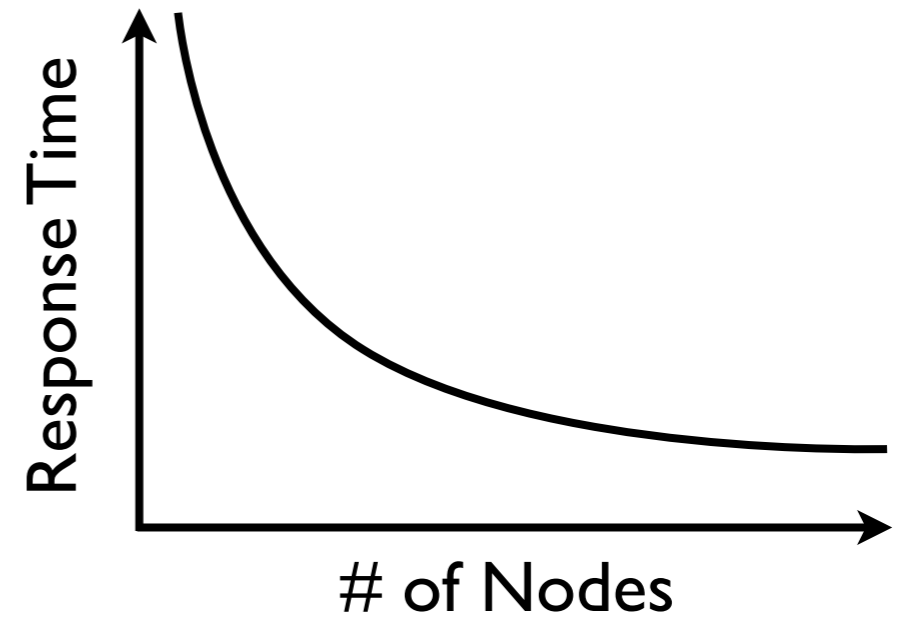


DBMSes: The First || Success Story

- Every major DBMS vendor has a || version.
- Reasons for success:
 - Bulk Processing (Partition ||-ism).
 - Natural Pipelining in RA plan.
 - Users don't need to think in ||.

Types of Speedup

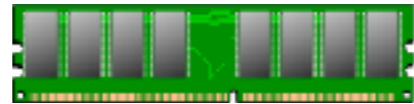
- Speed-up ||-ism
 - More resources = proportionally less time spent.
- Scale-up ||-ism
 - More resources = proportionally more data processed.



Parallelism Models



CPU

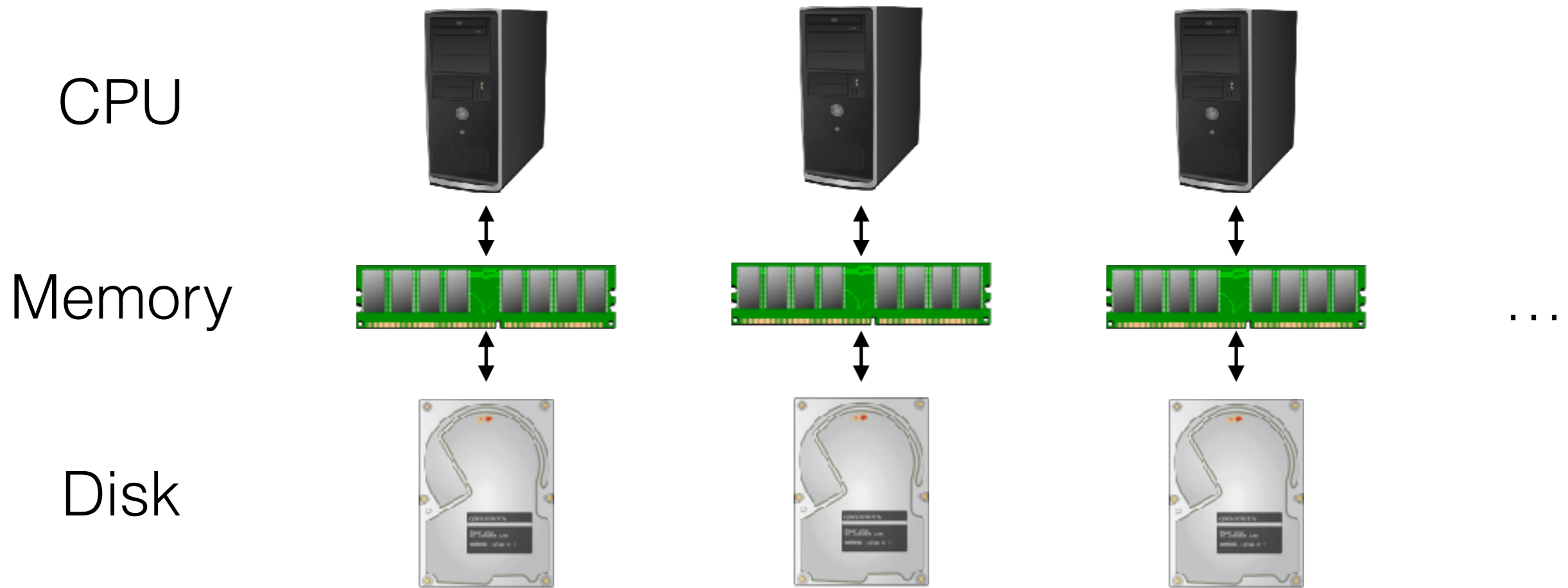


Memory



Disk

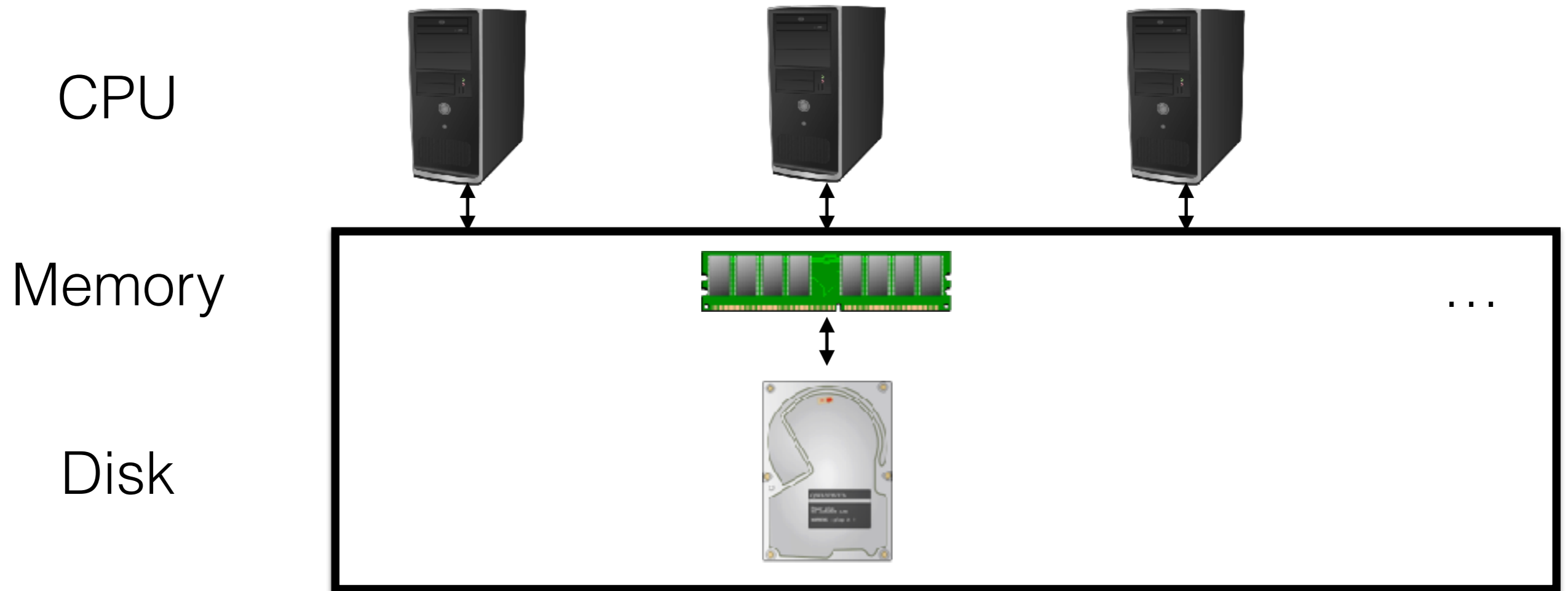
Parallelism Models



How do the nodes communicate?

Parallelism Models

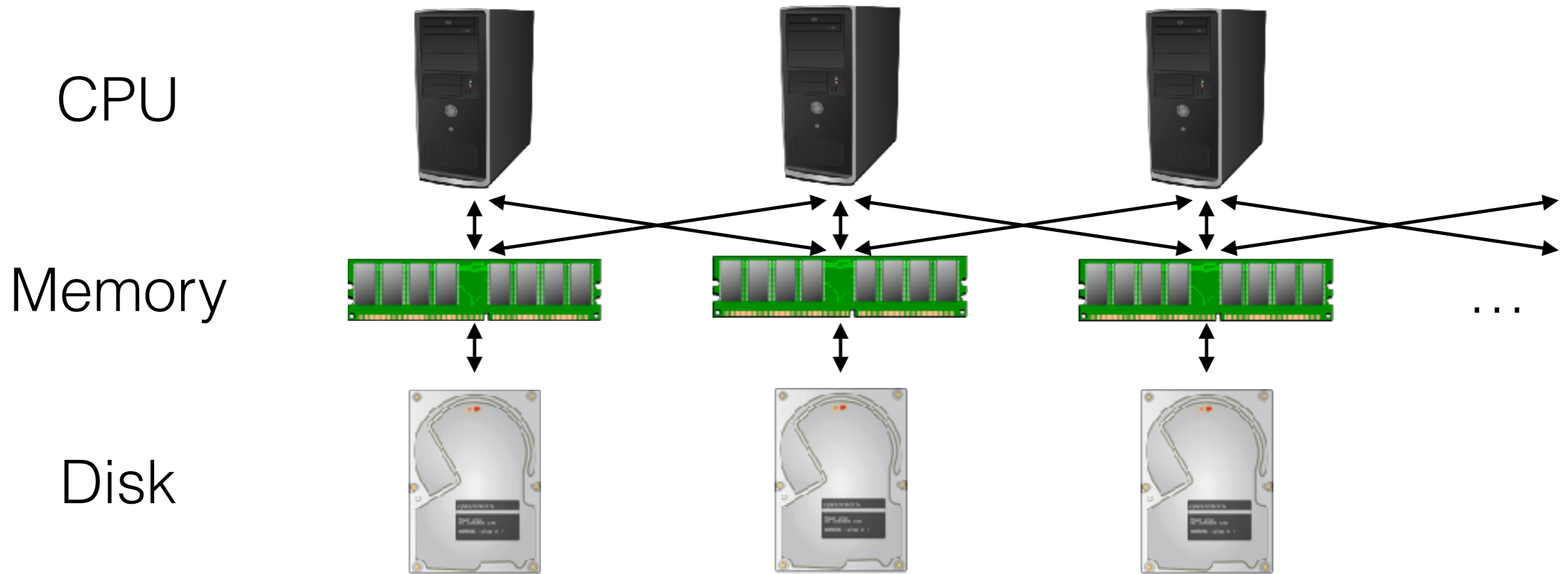
Option 1: “Shared Memory” available to all CPUs



e.g., a Multi-Core/Multi-CPU System

Parallelism Models

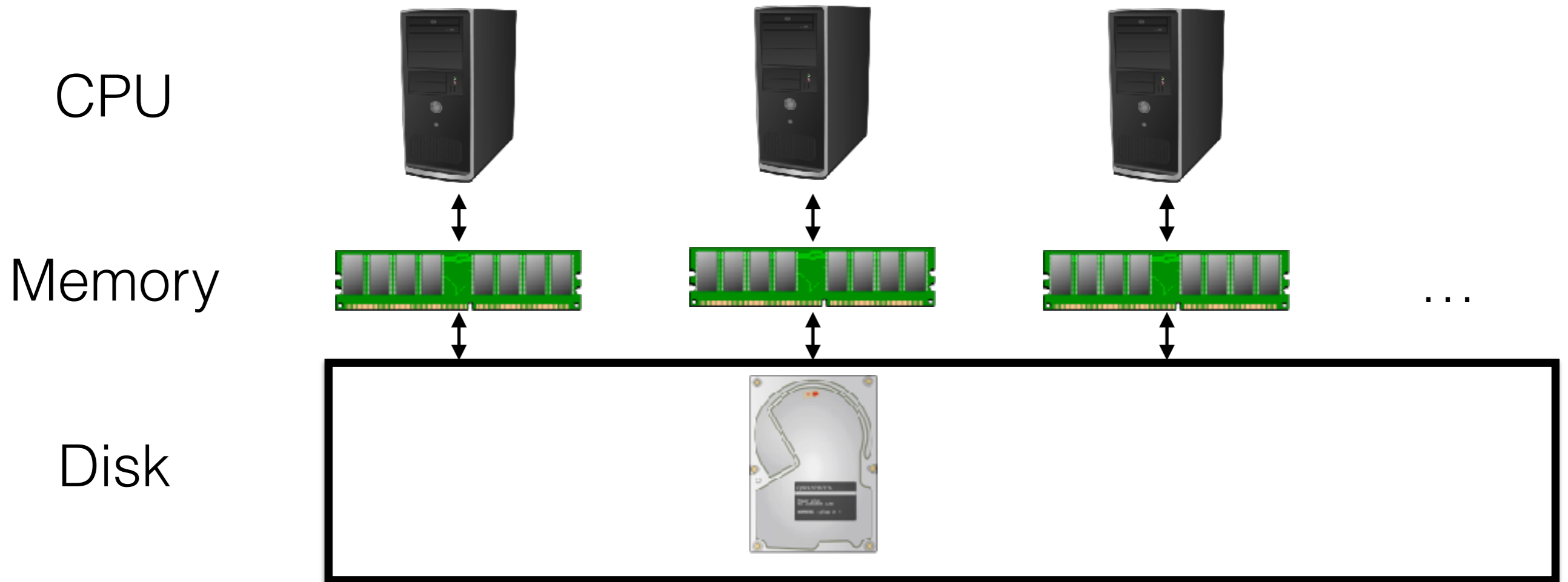
Option 2: Non-Uniform Memory Access.



Used by most AMD servers

Parallelism Models

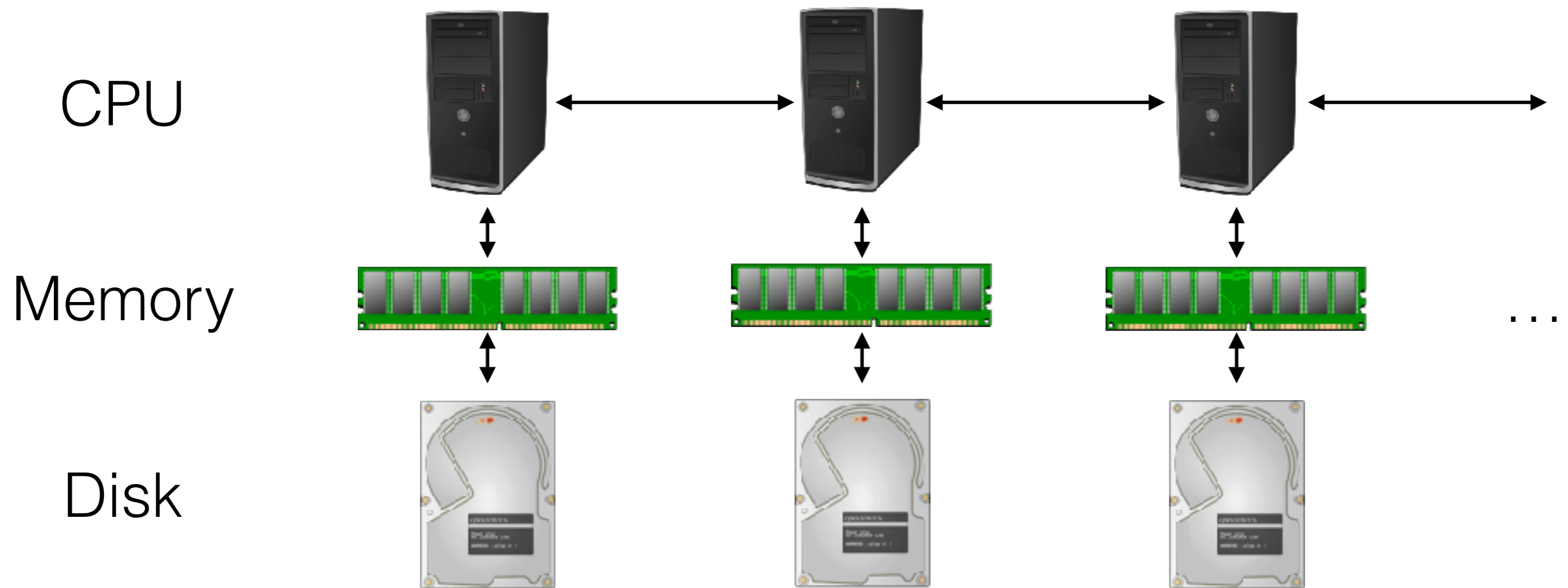
Option 3: “Shared Disk” available to all CPUs



Each node interacts with a “disk” on the network.

Parallelism Models

Option 4: “Shared Nothing” in which all communication is explicit.



Examples include MPP, Map/Reduce. Often used as basis for other abstractions.

Parallelizing

OLAP - Parallel Queries

OLTP - Parallel Updates

Parallelizing

OLAP - Parallel Queries

OLTP - Parallel Updates

Parallelism & Distribution

- Distribute the Data
 - Redundancy
 - Faster access
- Parallelize the Computation
 - Scale up (compute faster)
 - Scale out (bigger data)

Operator Parallelism

- **General Concept:** Break task into individual units of computation.
- **Challenge:** How much data does each unit of computation need?
- **Challenge:** How much data *transfer* is needed to allow the unit of computation?

Same challenges arise in Multicore, CUDA programming.

Parallel Data Flow



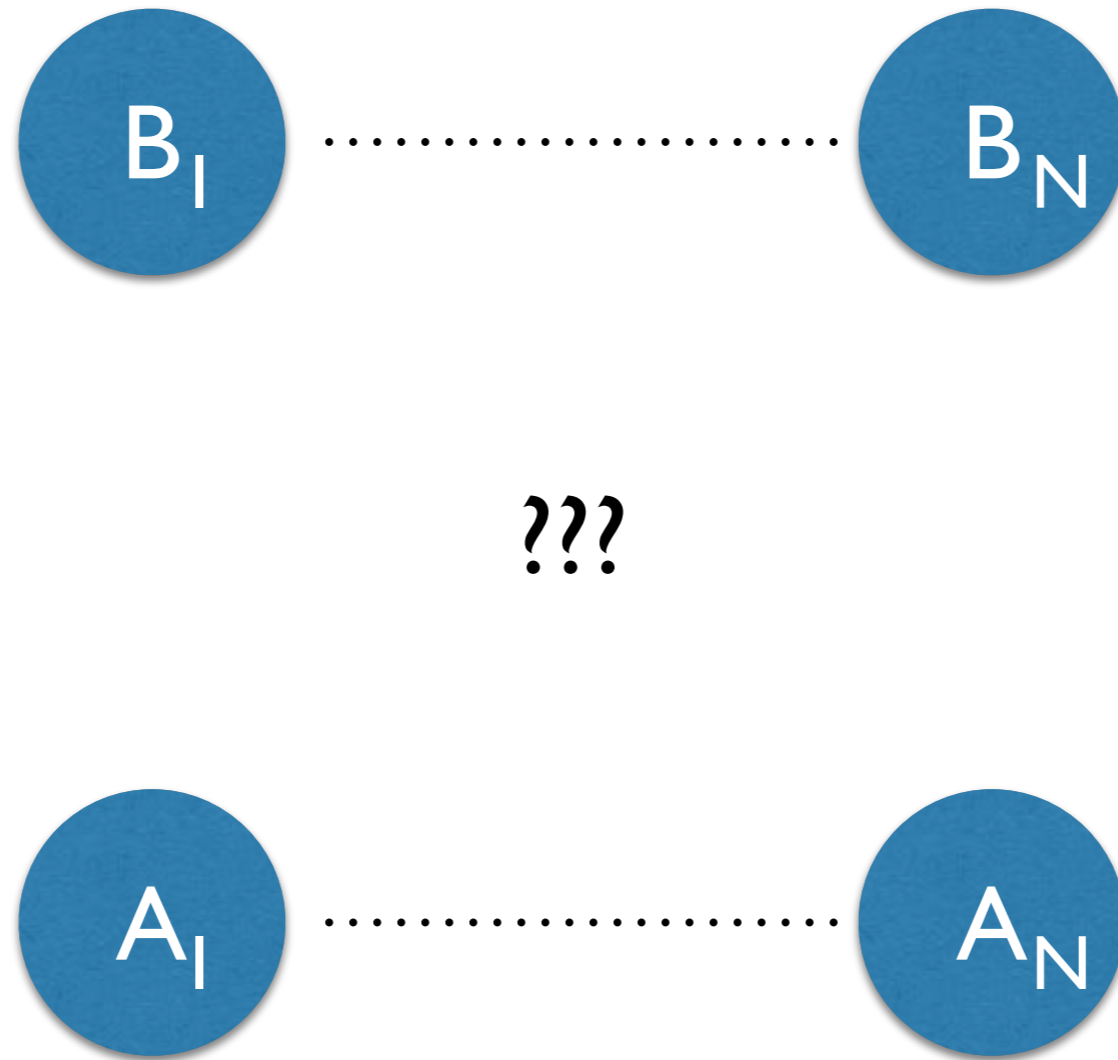
No Parallelism

Parallel Data Flow



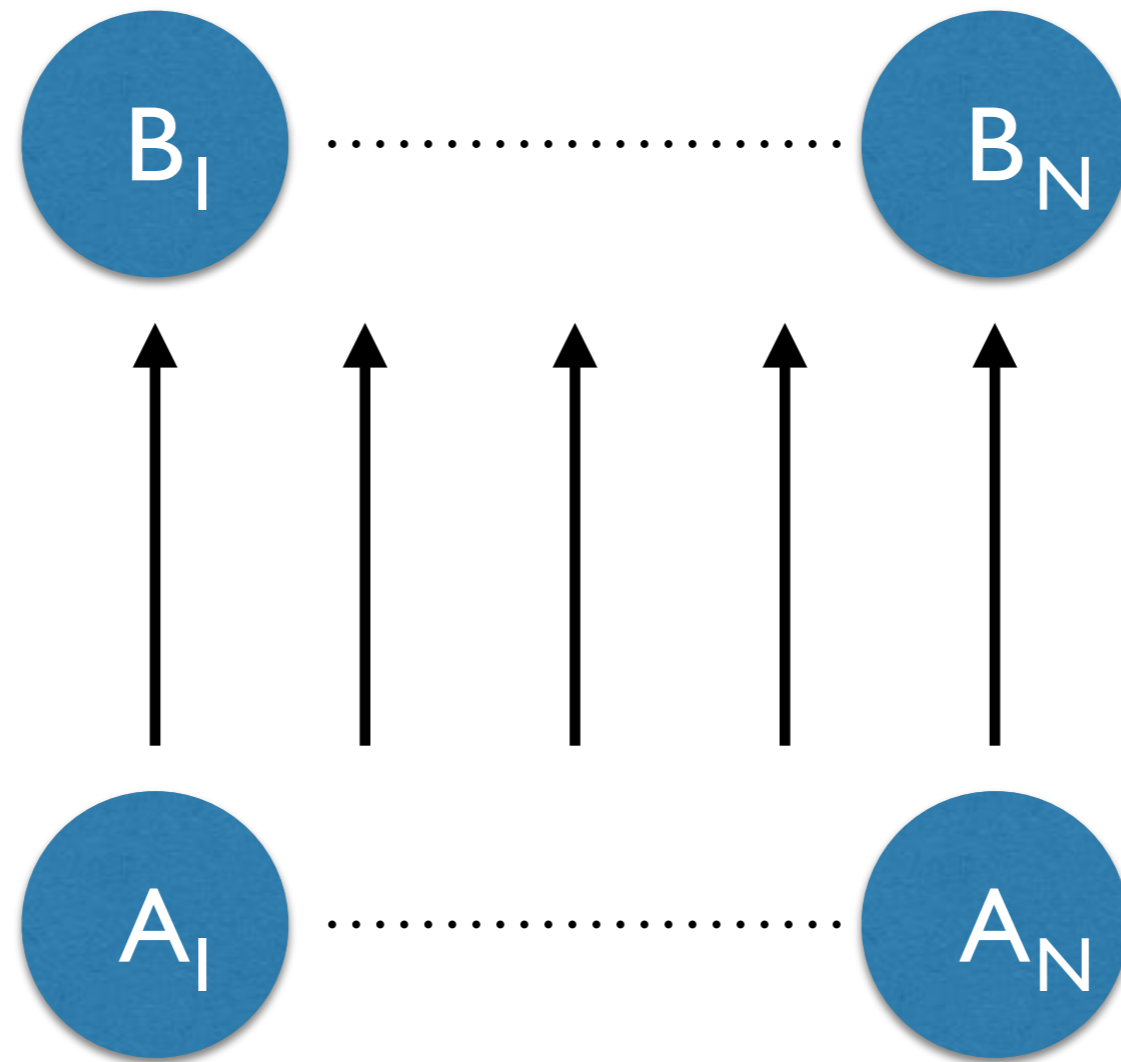
N-Way Parallelism

Parallel Data Flow



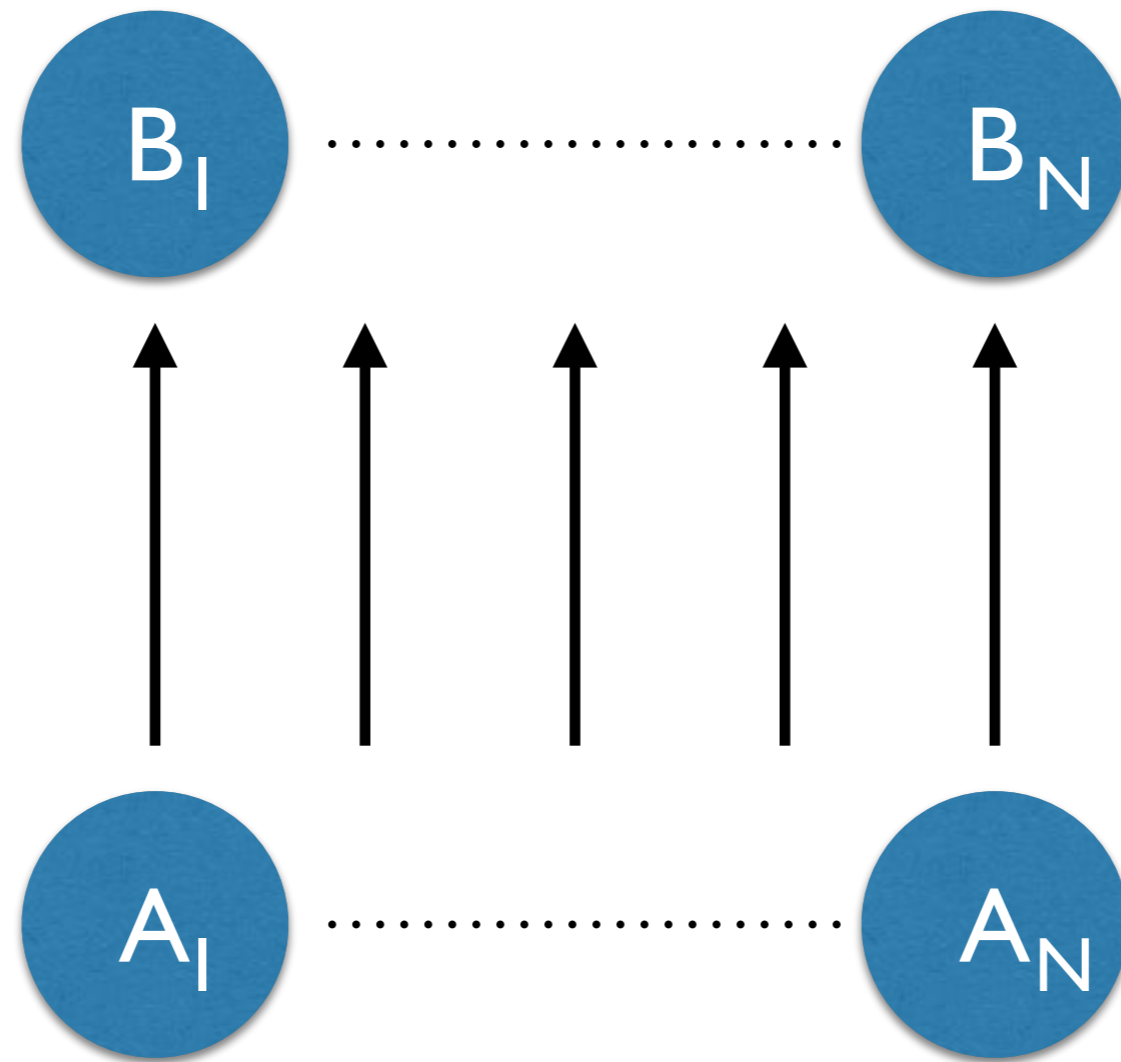
Chaining Parallel Operators

Parallel Data Flow



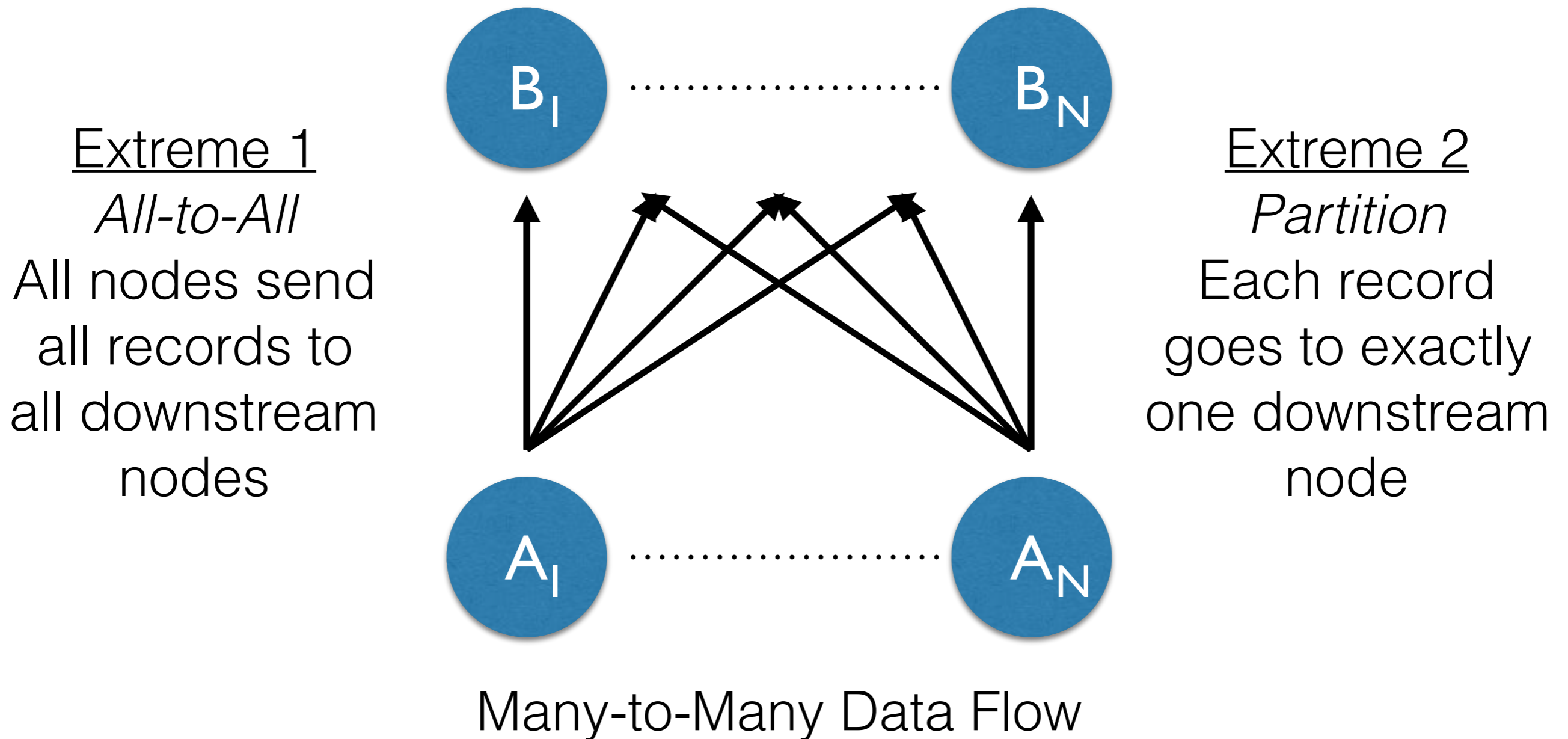
One-to-One Data Flow (“Map”)

Parallel Data Flow

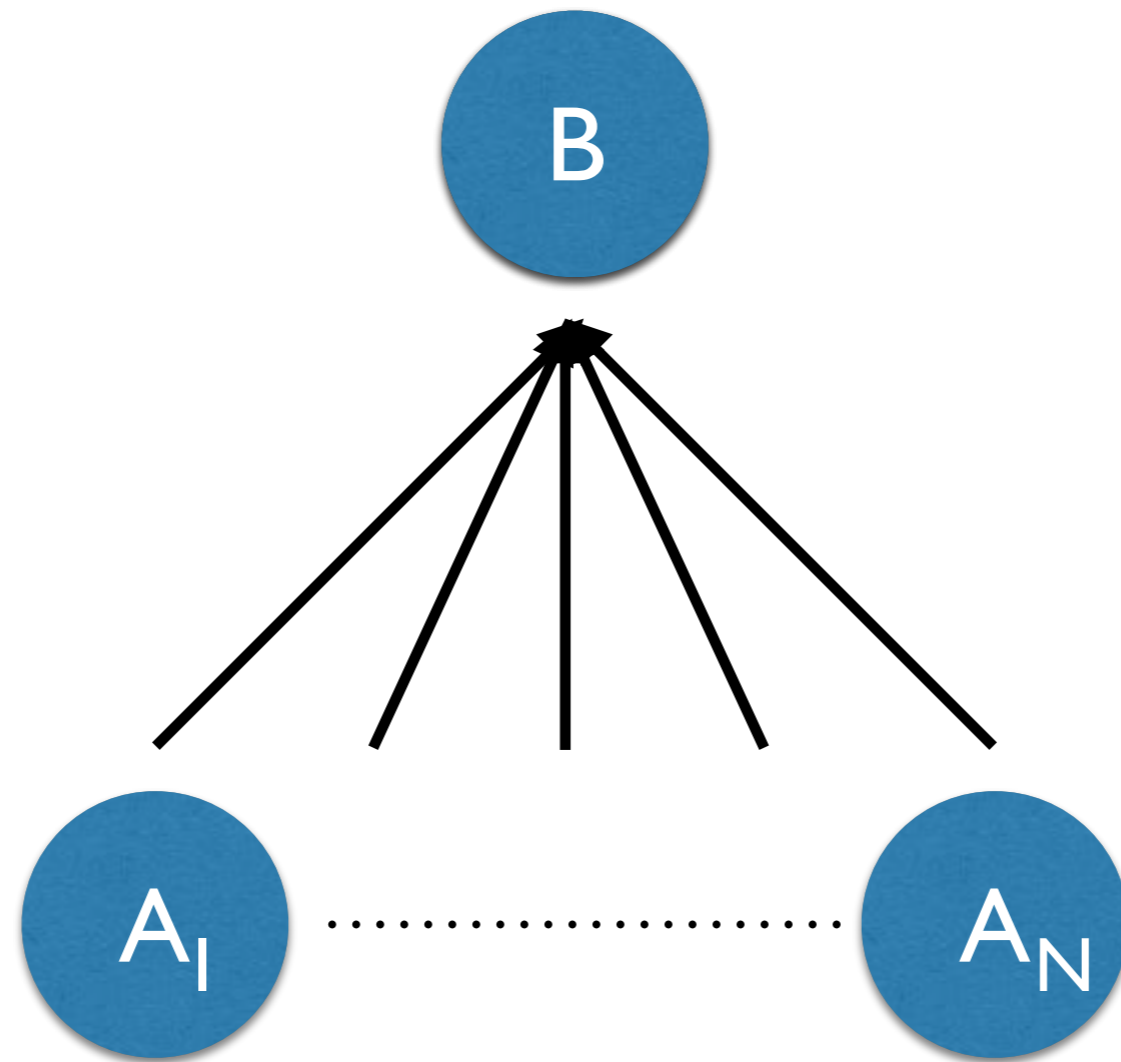


One-to-One Data Flow

Parallel Data Flow



Parallel Data Flow



Many-to-One Data Flow (“Reduce/Fold”)

Parallel Operators

Select

Project

Union (bag)

What is a logical “unit of computation”?

(1 tuple)

Is there a data dependency between units?

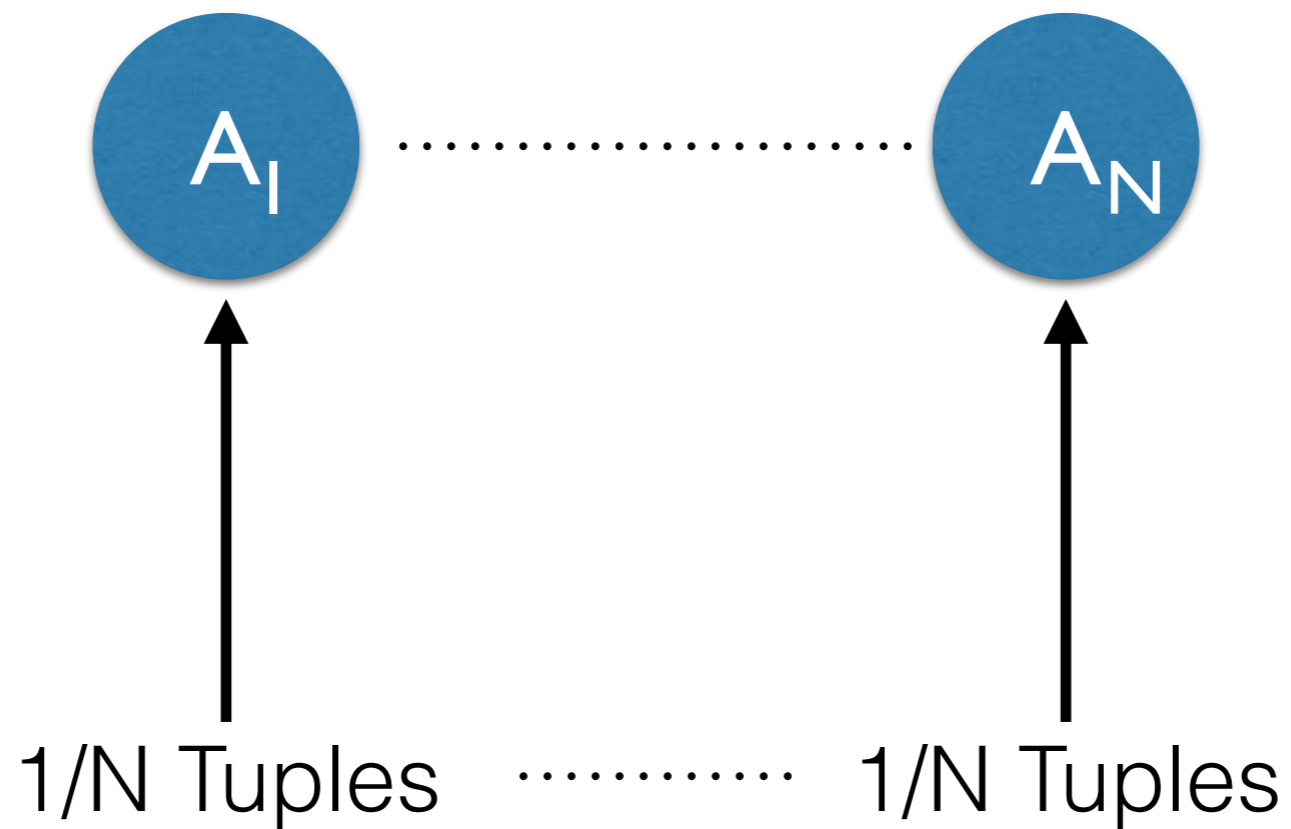
(no)

Parallel Operators

Select

Project

Union (bag)



Parallel Joins

```
FOR i IN 1 to N  
  FOR j IN 1 to K  
    JOIN(Block i of R,  
         Block j of S)
```

One Unit of Computation

Parallel Joins

K Partitions of S

N Partitions of R

Block I of R
⋈

Block I of S

N
⋮
↓

Block N of R
⋈

Block I of S

K

⋯→

⋮
↘

K

⋯→

Block I of R
⋈

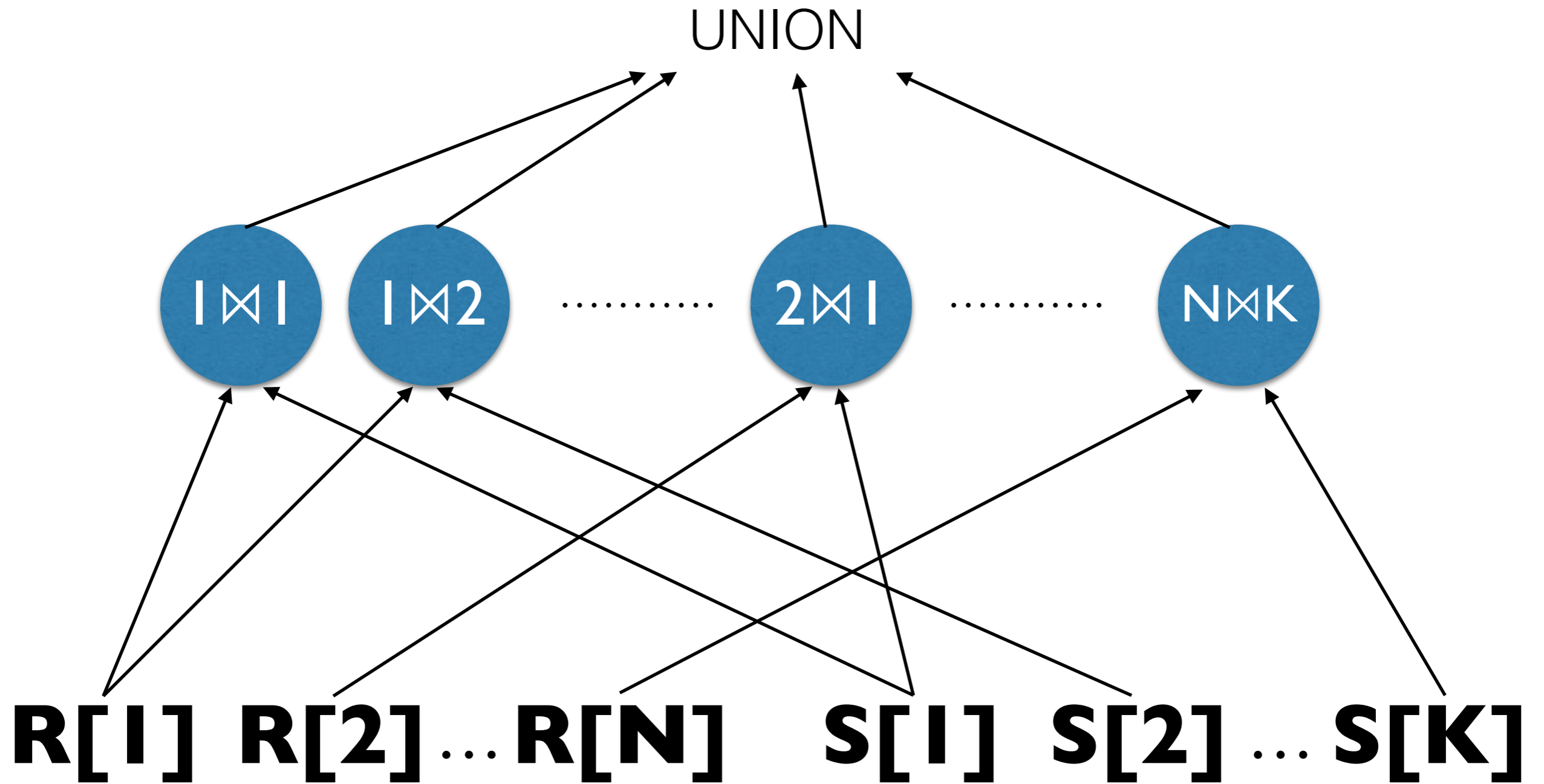
Block K of S

N
⋮
↓

Block N of R
⋈

Block K of S

Parallel Joins



Parallel Joins

How much data needs to be transferred?

How many “units of computation” do we create?

Parallel Joins

What if we partitioned “intelligently”?

Parallel Joins

	0	1	2	3
0	✓	✗	✗	✗
1		✓		
2			✓	
3				✓

R \bowtie_B **S**: Which Partitions of S Join w/ Bucket 0 of R?

Parallel Joins

$\frac{B}{S}$	$B < 25$	$25 \leq B < 50$	$50 \leq B < 75$	$75 \leq B$
<u>R.B</u> $B < 25$	✓	✓	✓	✓
$25 \leq B < 50$	✗	✓	✓	✓
$50 \leq B < 75$	✗	✗	✓	✓
$75 \leq B$	✗	✗	✗	✓

R $\bowtie_{R.B < S.B}$ **S**: Which Partitions of S Can Produce Output?

Parallel Joins

Use partitioning to eliminate
units of computation

Exactly the same idea as External Hash Join
(Called Theta Join for Inequalities)

Bloom Join

No Specific Partitioning

No Specific Partitioning

✓	✓	✓	✗
✗	✓	✓	✗
✓	✗	✗	✗
✗	✓	✓	✗

What if the join is highly selective... Can we detect which tuples are useful?

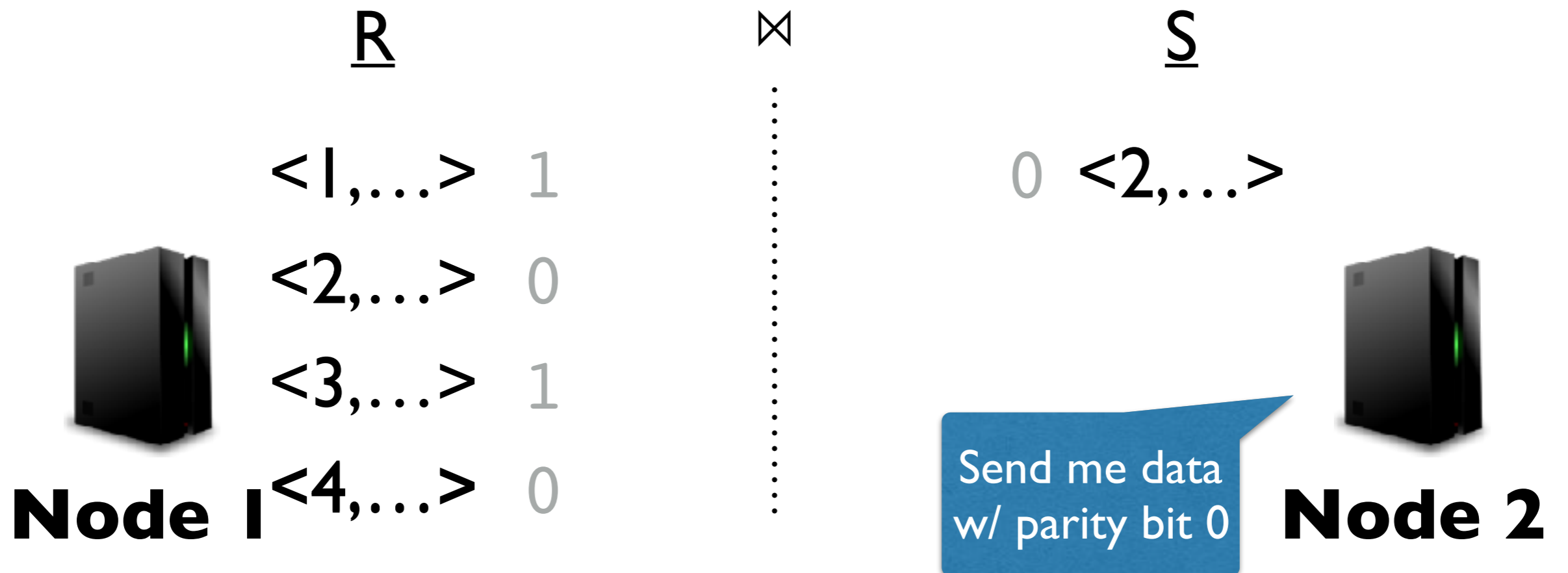
Bloom Join

Goal: Summarize which tuples are useful for the join?

False positives: OK
False negatives: NOT OK

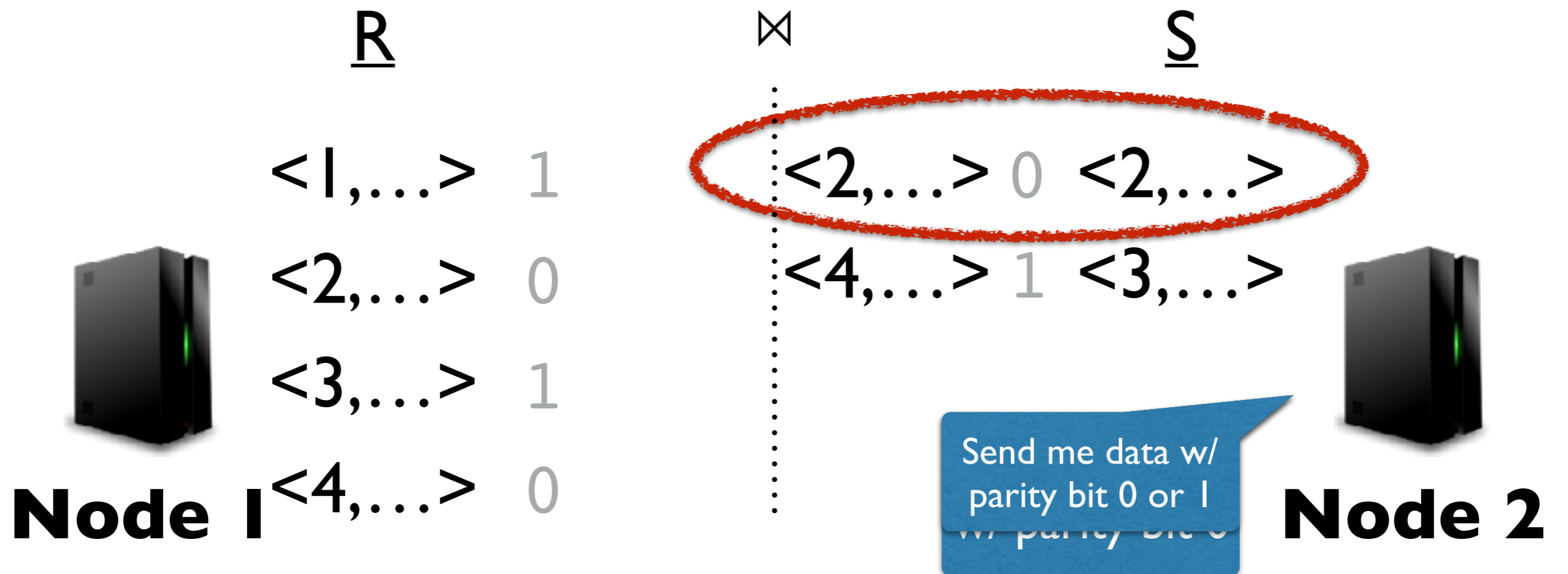
Bloom Join

Strategy 1: Parity Bit



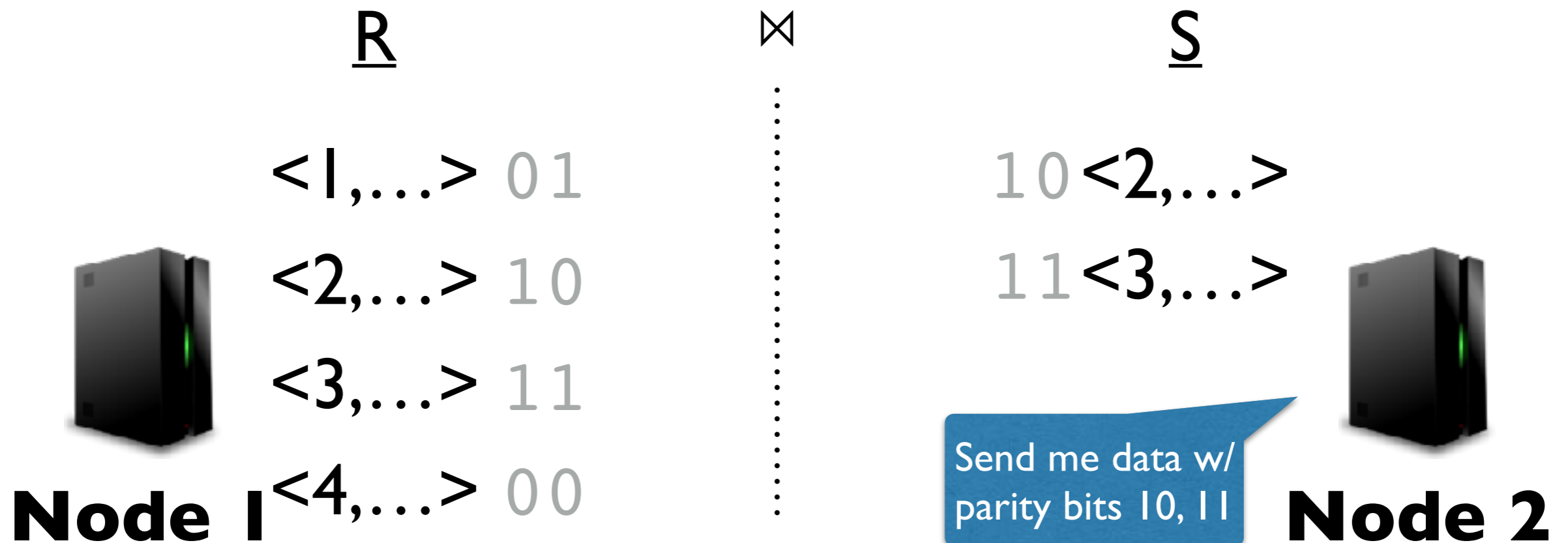
Bloom Join

Strategy 1: Parity Bit



Bloom Join

Strategy 2: Multiple Parity Bits



What's the problem with this?

A Simplified Bloom Join

Key 1 00101010

Key 2 01010110

Key 3 10000110

Key 4 01001100

How do we summarize?

Bitwise OR

e.g. (Key 1 | Key 2)

= 01111110

How do we test for inclusion?

(Key & Summary) == Key?

(Key 1 & S) = 00101010 ✓

(Key 3 & S) = 00000110 ✗

(Key 4 & S) = 01001100 ✓

False Positive

Bloom Filters

Generating a bit vector for a key:

M - # of bits in the bit vector

K - # of hash functions

For ONE key/record:

For i between 0 and K:

$\text{bitvector}[\text{hash}_i(\text{key}) \% M] = 1$

Each bit vector has $\sim K$ bits set

Bloom Filters

Probability that 1 bit is set by 1 hash fn

$$1/m$$

Bloom Filters

Probability that 1 bit is not set by 1 hash fn

$$1 - 1/m$$

Bloom Filters

Probability that 1 bit is not set by k hash fns

$$(1 - 1/m)^k$$

Bloom Filters

Probability that 1 bit is not set by k hash fns
for n records

$$(1 - 1/m)^{kn}$$

So for an arbitrary record, what is the probability
that all of its bits will be set?

Bloom Filters

Probability that 1 bit is set by k hash fns
for n records

$$1 - \left(1 - \frac{1}{m}\right)^{kn}$$

Bloom Filters

Probability that all k bits are set by k hash fns
for n records

$$\approx \left(1 - \left(1 - 1/m \right)^{kn} \right)^k$$

$$\approx \left(1 - e^{-kn/m} \right)^k$$

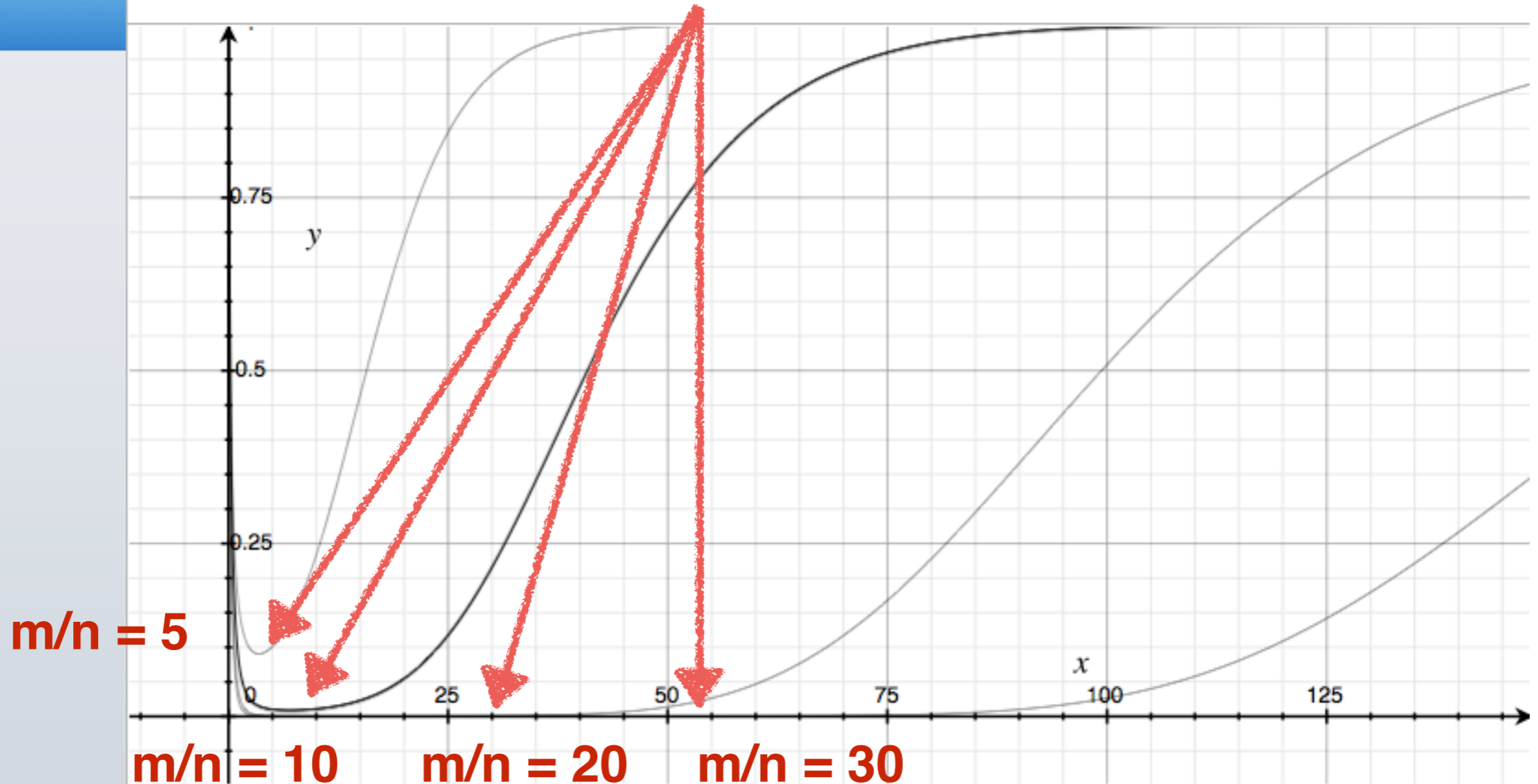
Bloom Filters

- $y = \left(1 - e^{-\frac{1}{5}x}\right)^x$
- $y = \left(1 - e^{-\frac{1}{10}x}\right)^x$
- $y = \left(1 - e^{-\frac{1}{20}x}\right)^x$
- $y = \left(1 - e^{-\frac{1}{30}x}\right)^x$

$$y = \left(1 - e^{-\frac{1}{10}x}\right)^x$$

Minimal P[collision]

$\rightarrow \Sigma^2$



Minimal P[collision] is at $k \approx c \cdot m/n$

Bloom Filters

$$k \approx c \cdot m/n$$

$$\frac{m}{k} \approx cn$$

m is linearly related to n (for a fixed k)

Bloom Join

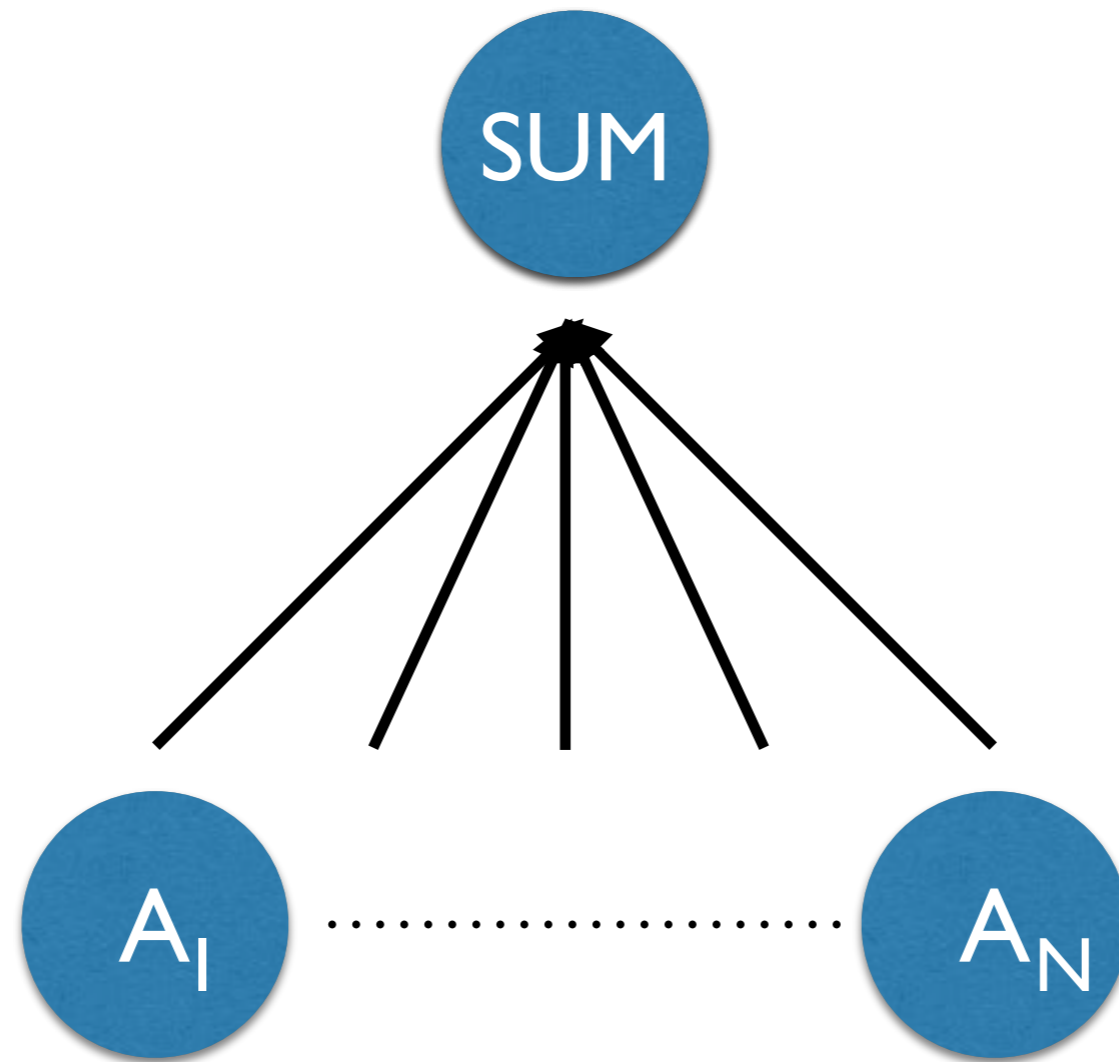
- Node 2 Computes Bloom Filter for Local Records
- Node 2 Sends Bloom Filter to Node 1
- Node 1 Matches Local Records Against Bloom Filter
- Node 1 Sends Matched Records to Node 2
 - Superset of “useful” records
- Node 2 Performs Join Locally

Parallel Aggregates

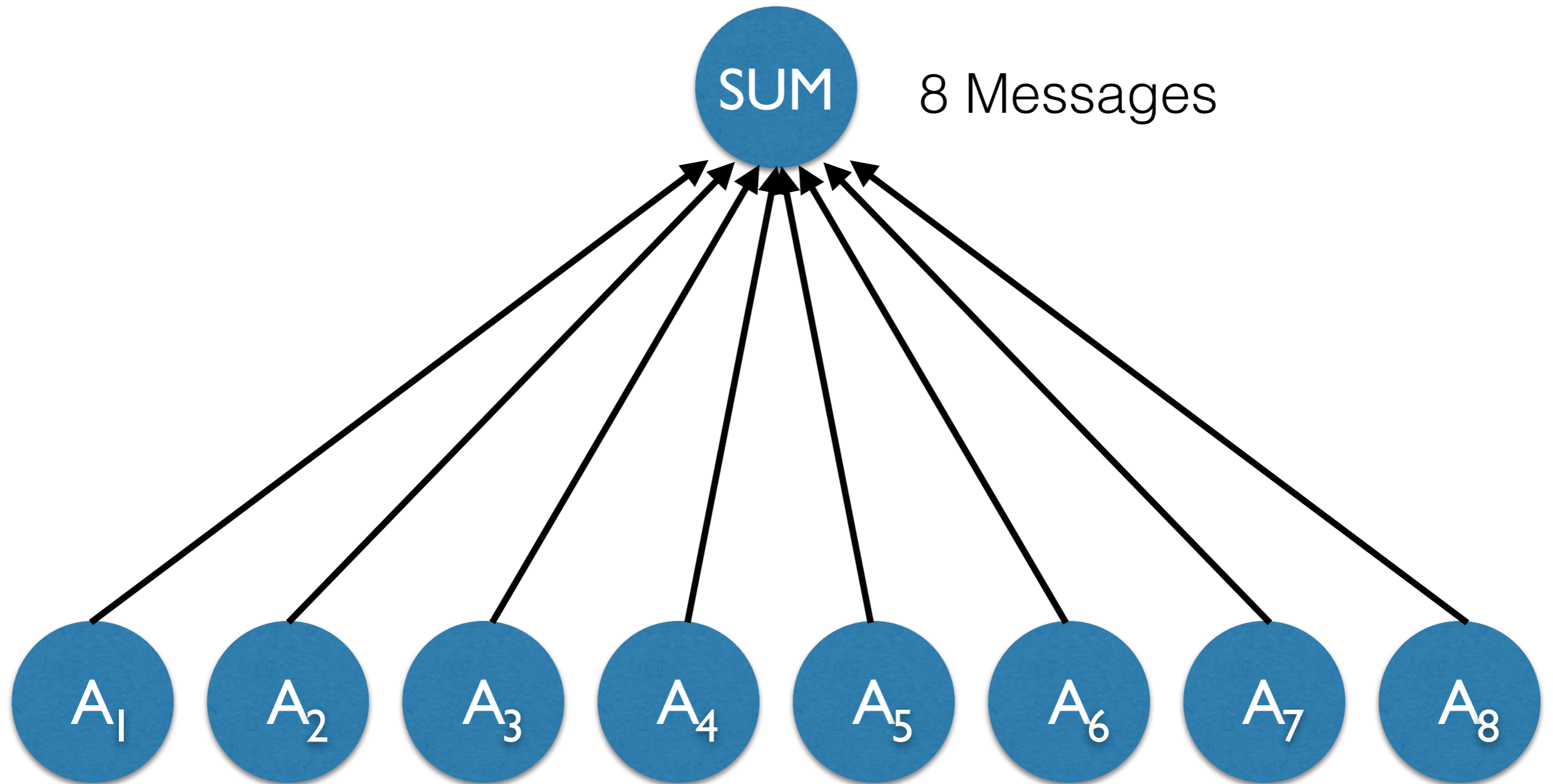
Algebraic: Bounded-size intermediate state
(Sum, Count, Avg, Min, Max)

Holistic: Unbounded-size intermediate state
(Median, Mode/Top-K Count, Count-Distinct;
Not Distribution-Friendly)

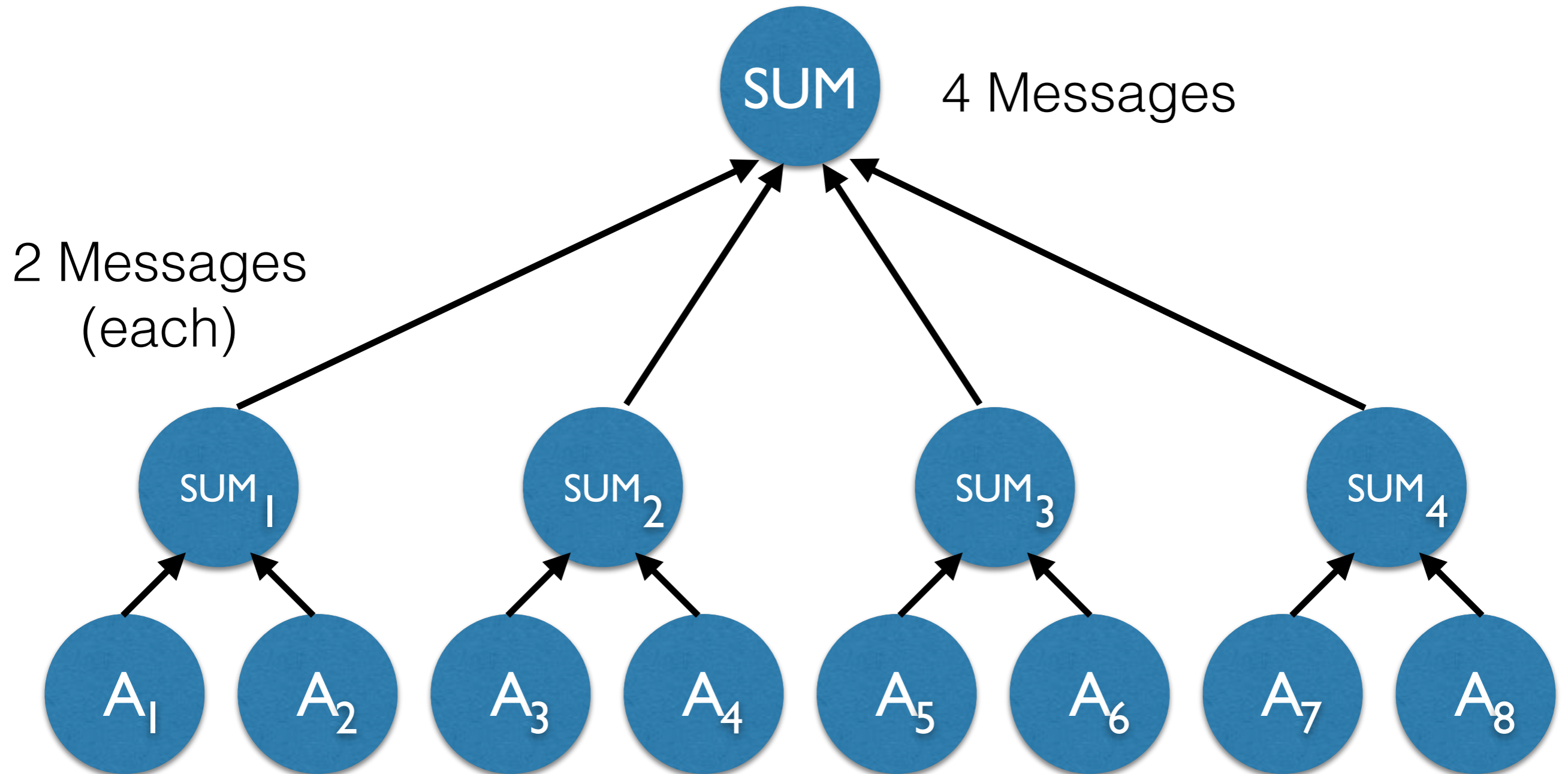
Fan-In Aggregation



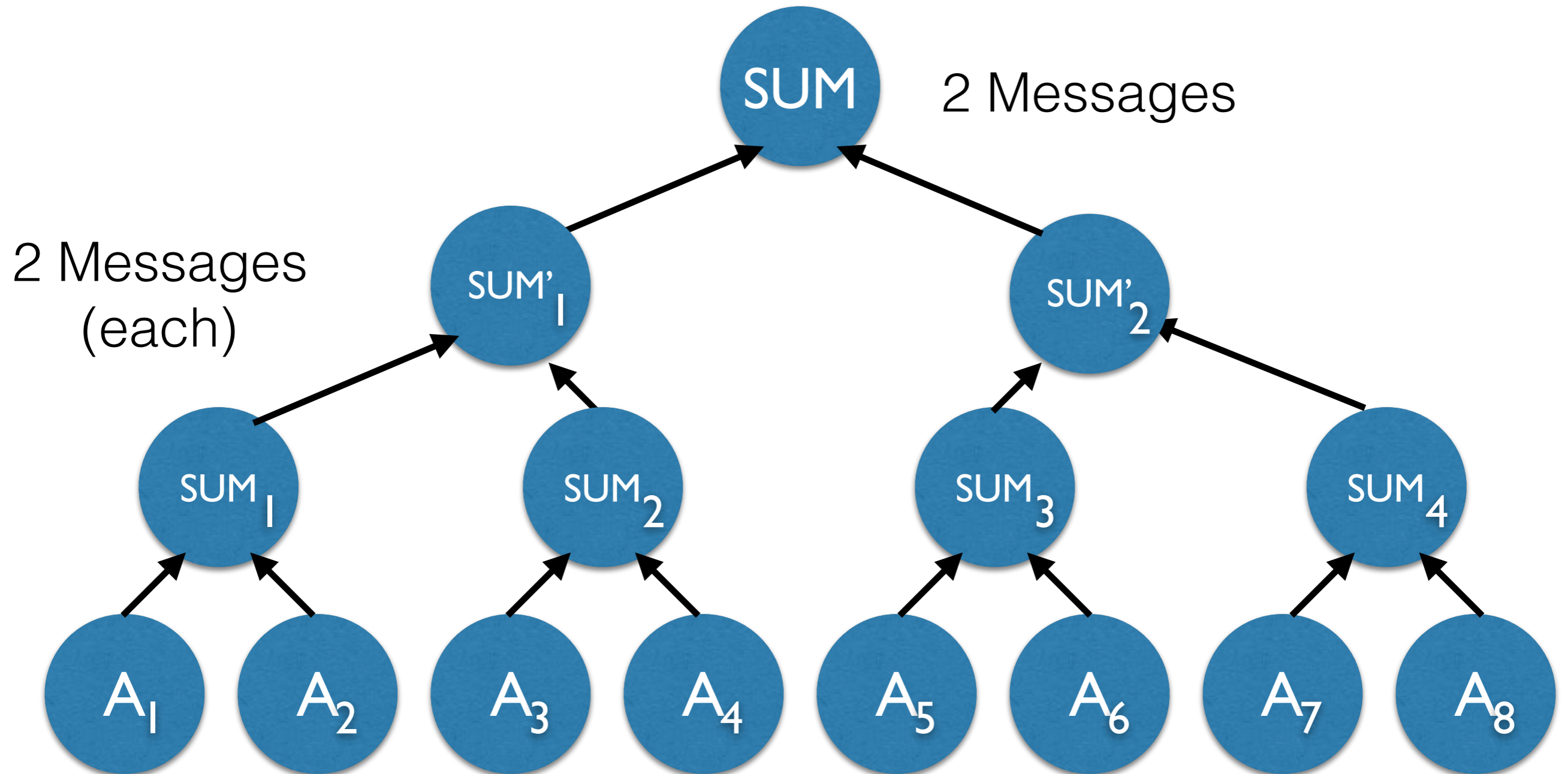
Fan-In Aggregation



Fan-In Aggregation



Fan-In Aggregation



Fan-In Aggregation

If Each Node Performs K Units of Work...
(K Messages)
How Many Rounds of Computation Are Needed?

$$\log_k(N)$$

Fan-In Aggregation Components

Combine(Intermediate₁, ..., Intermediate_N)
= Intermediate

$$\langle \text{SUM}_1, \text{COUNT}_1 \rangle \otimes \dots \otimes \langle \text{SUM}_N, \text{COUNT}_N \rangle \\ = \langle \text{SUM}_1 + \dots + \text{SUM}_N, \text{COUNT}_1 + \dots + \text{COUNT}_N \rangle$$

Compute(Intermediate) = Aggregate

$$\text{Compute}(\langle \text{SUM}, \text{COUNT} \rangle) = \text{SUM} / \text{COUNT}$$