

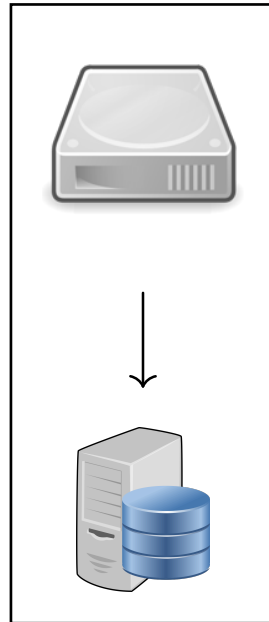
PARALLELIZING QUERIES

CSE 4/562: Database Systems | Lecture 14

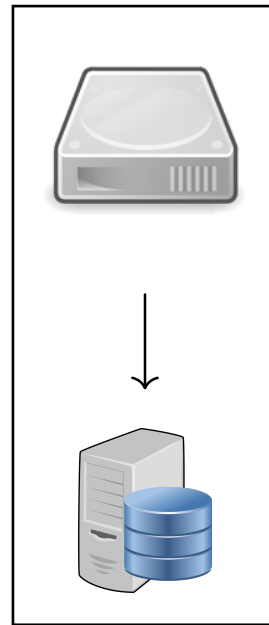
DB. Sys.: T.C.B.: Ch. 20.1-20.4

Why Scale?

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

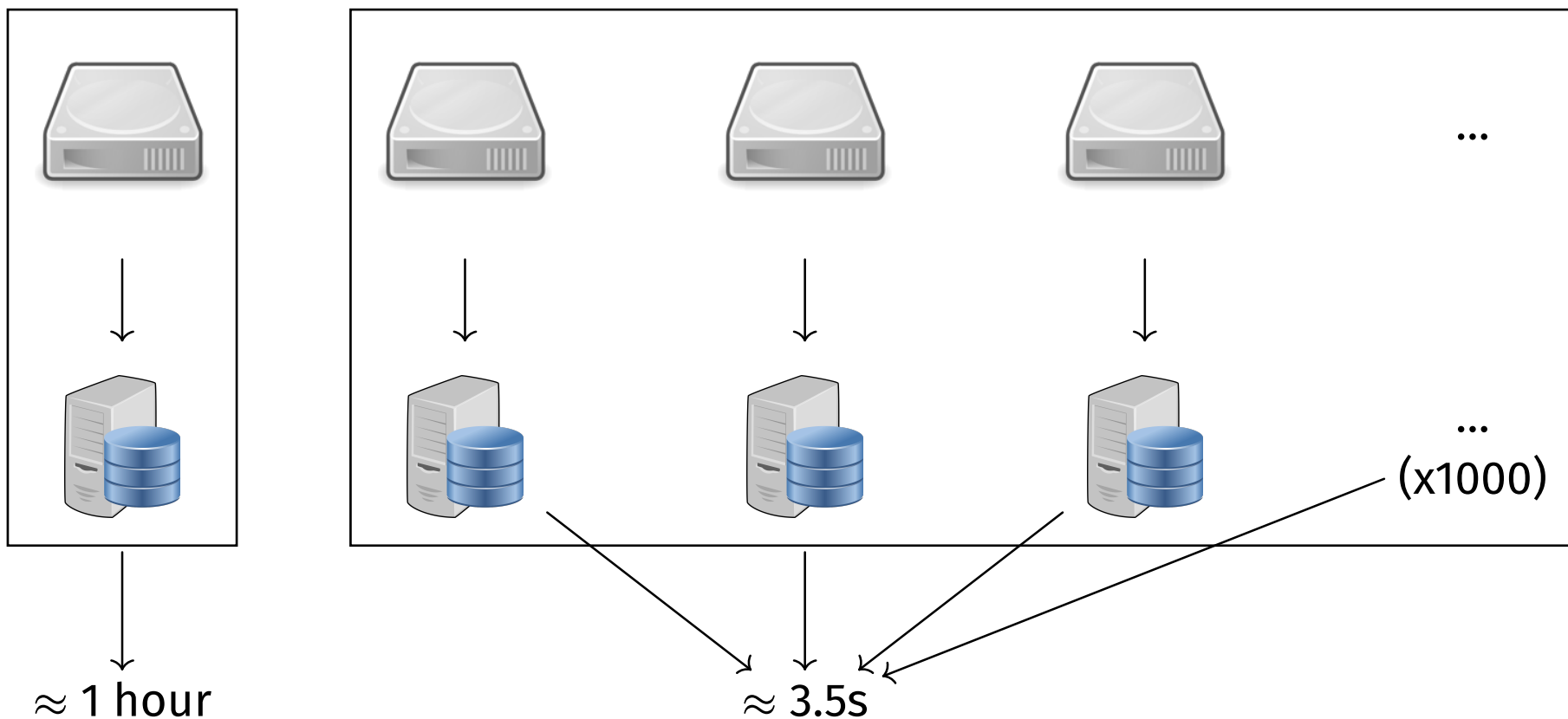


Scan 1 PB at 300MB/s (SATA r2)

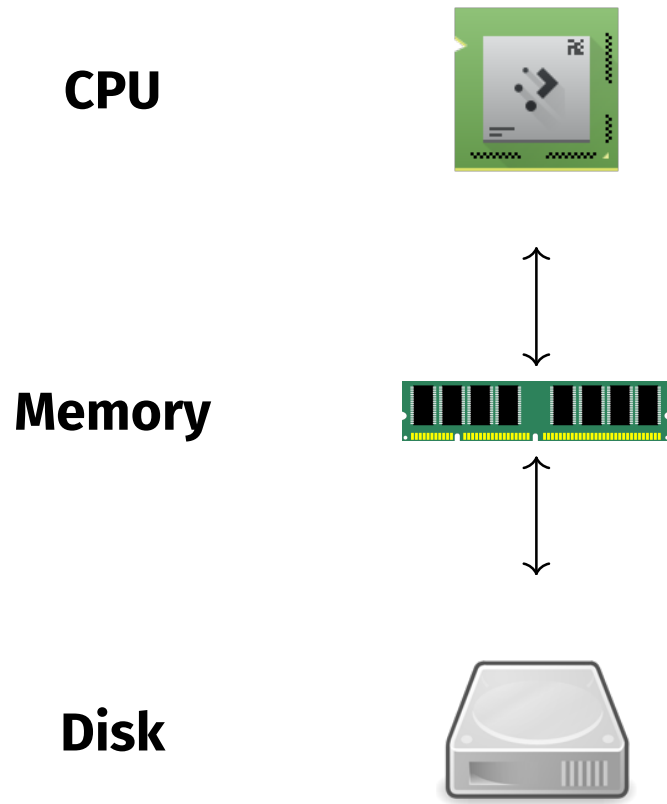


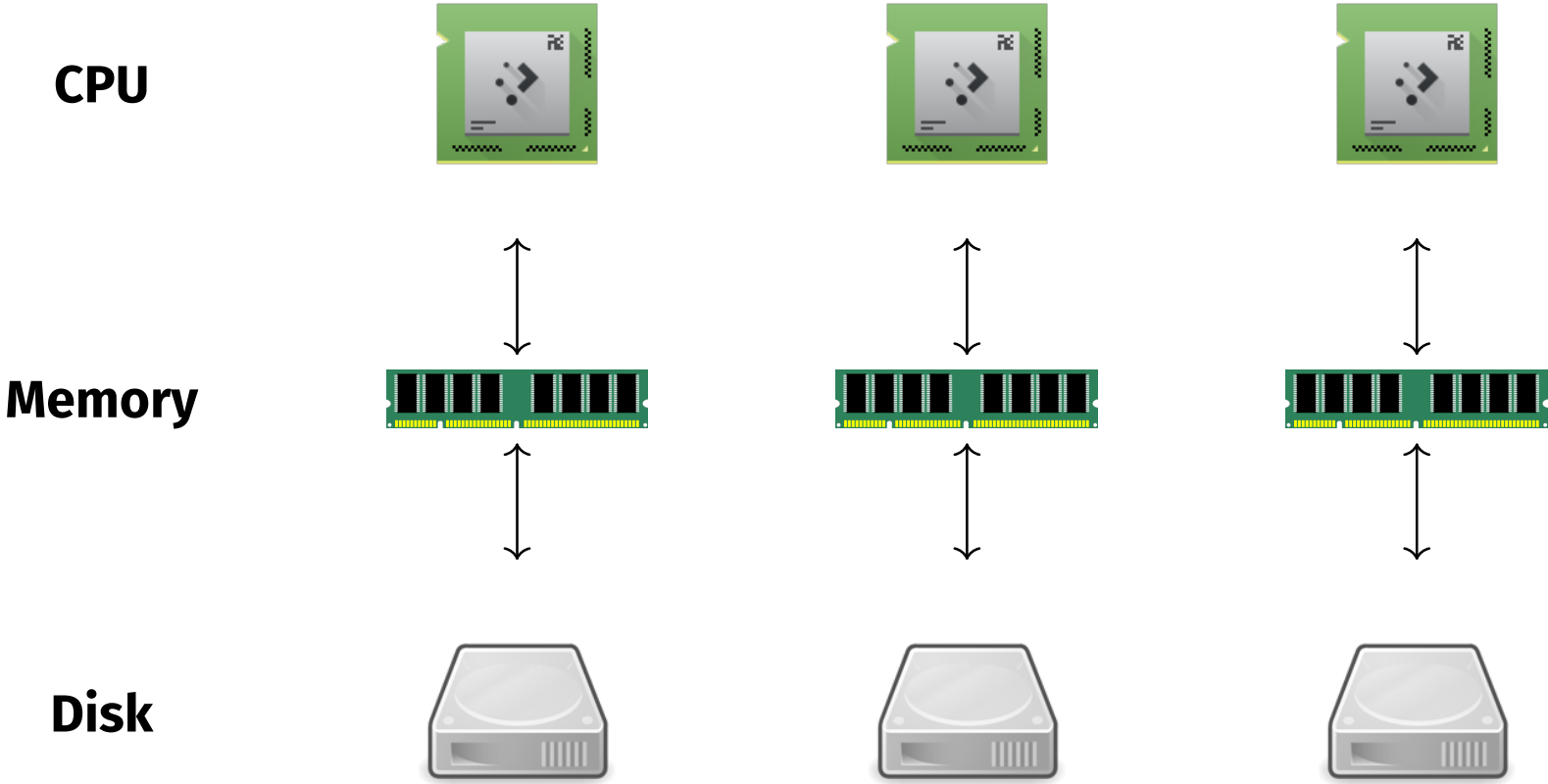
≈ 1 hour

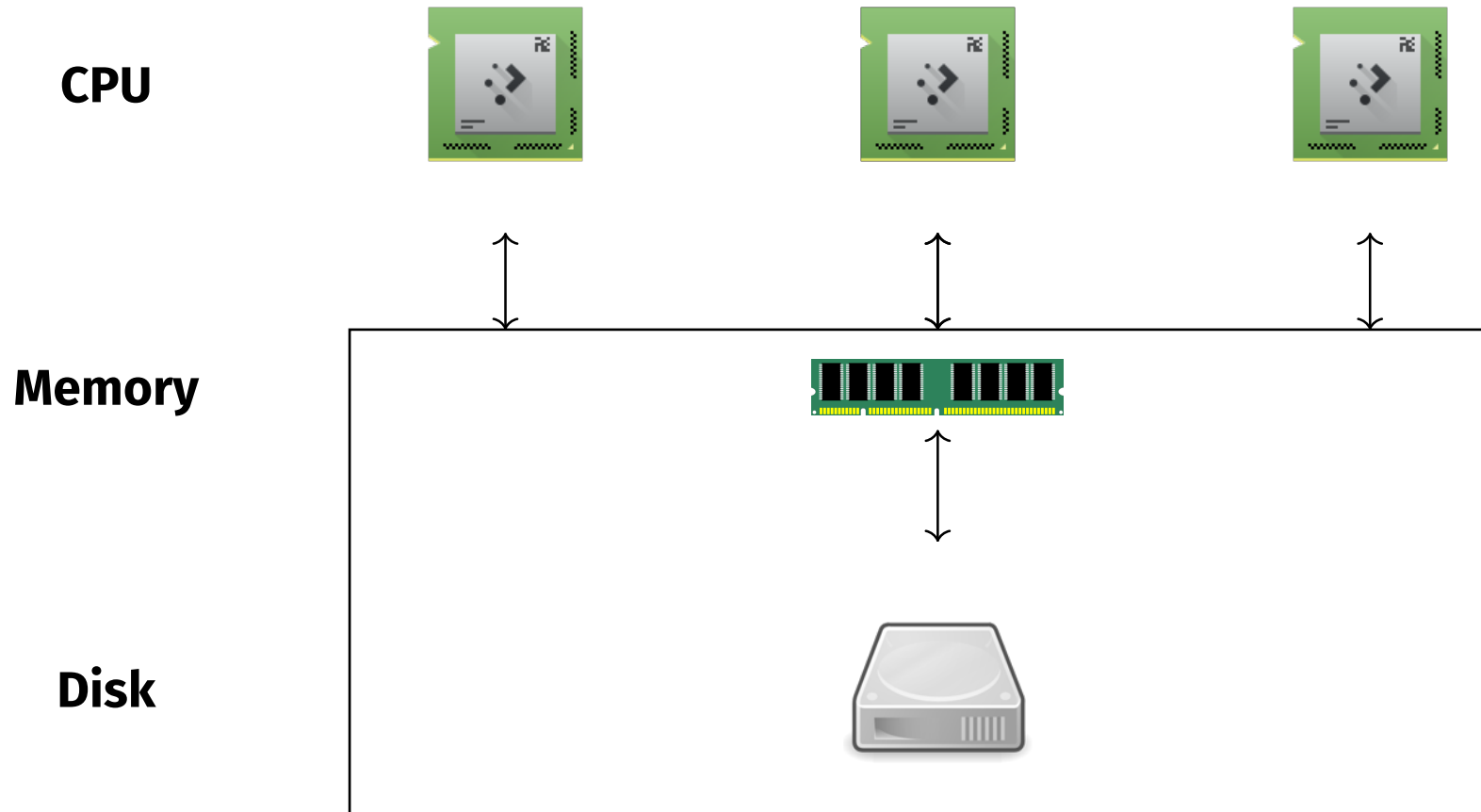
Scan 1 PB at 300MB/s (SATA r2)



Communication Models



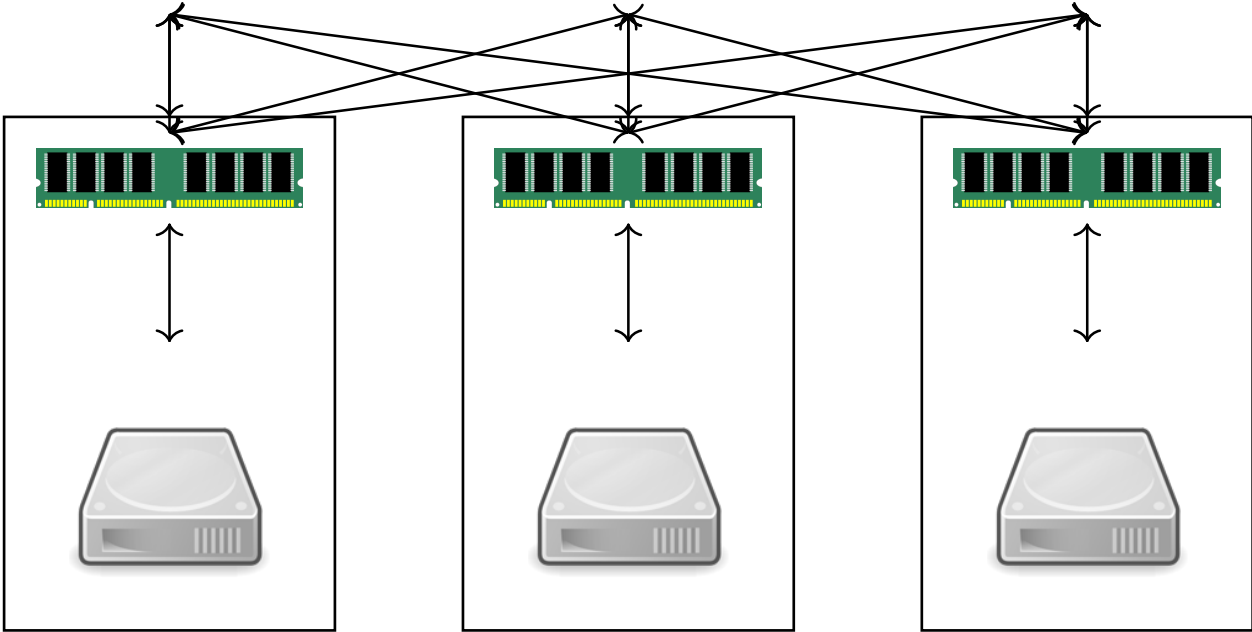




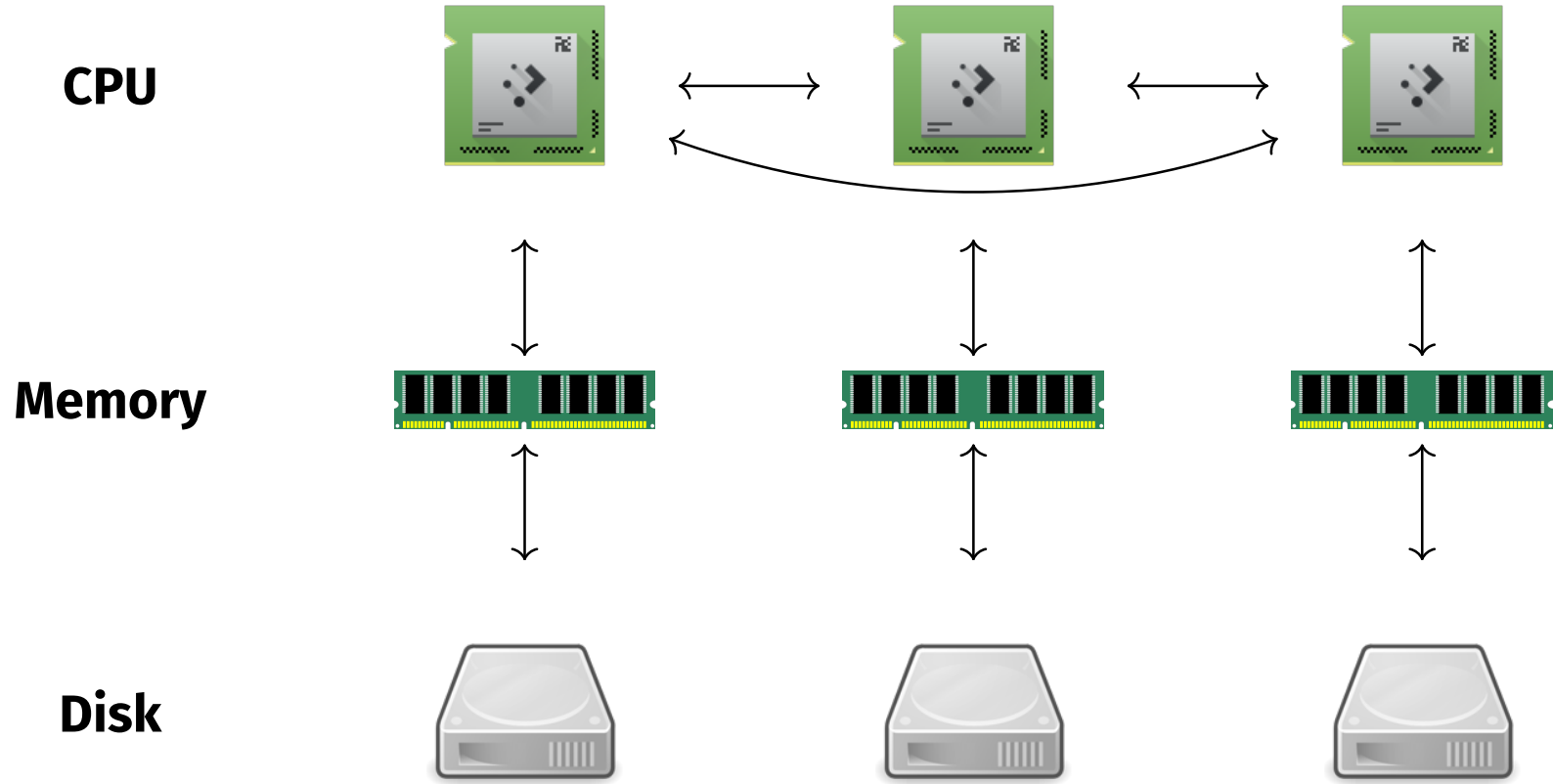
CPU



Memory



Disk



Shared Memory

- S3
- RAM/Modern OSes

NUMA

- AMD CPUs
- Hadoop

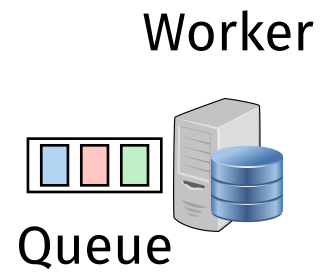
Shared Nothing

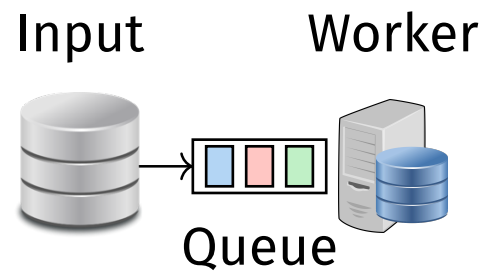
- MPP
- Apache Spark

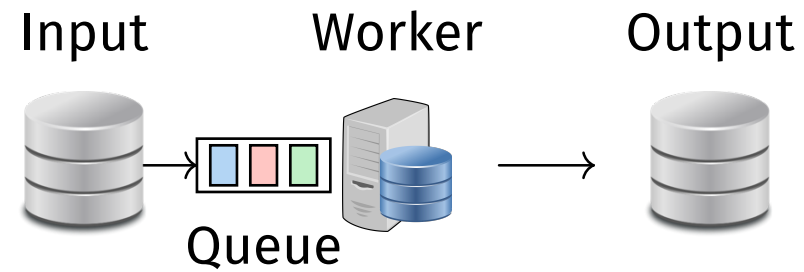
Today: Shared Nothing

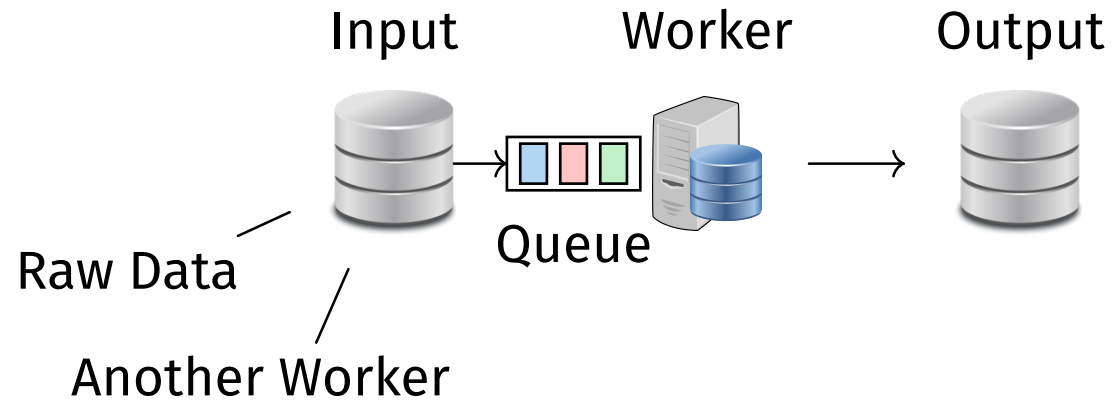
Worker

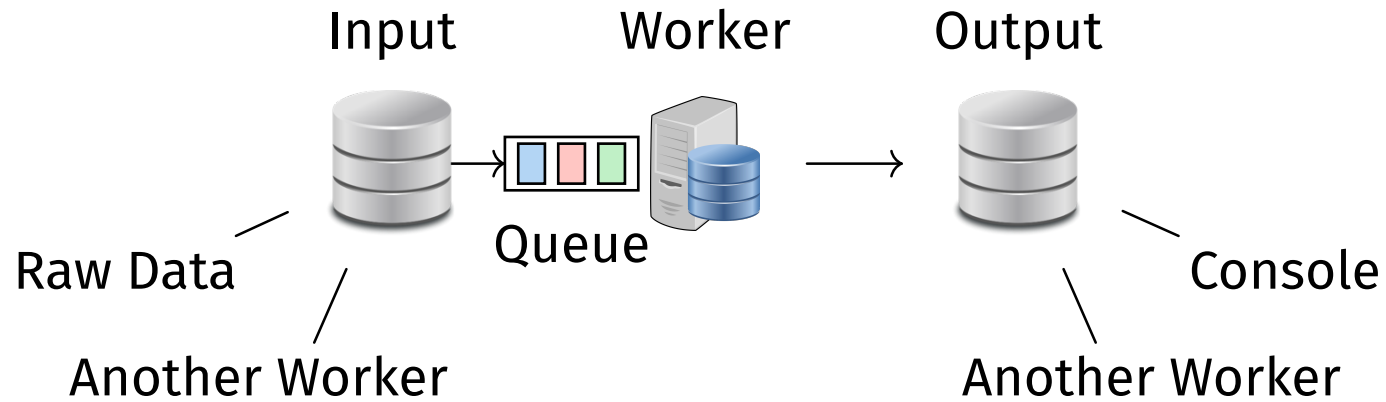


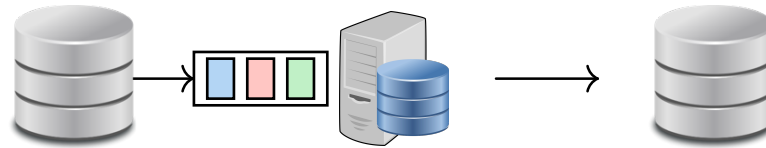




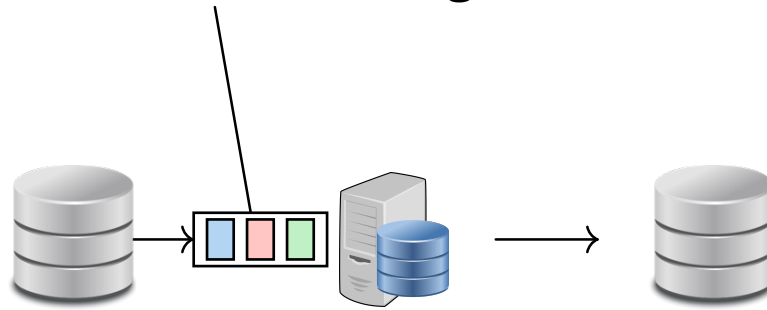




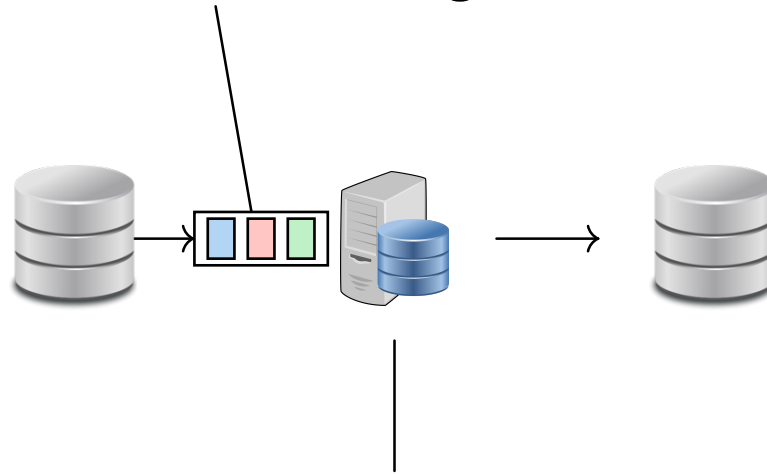




Data doesn't arrive fast enough.



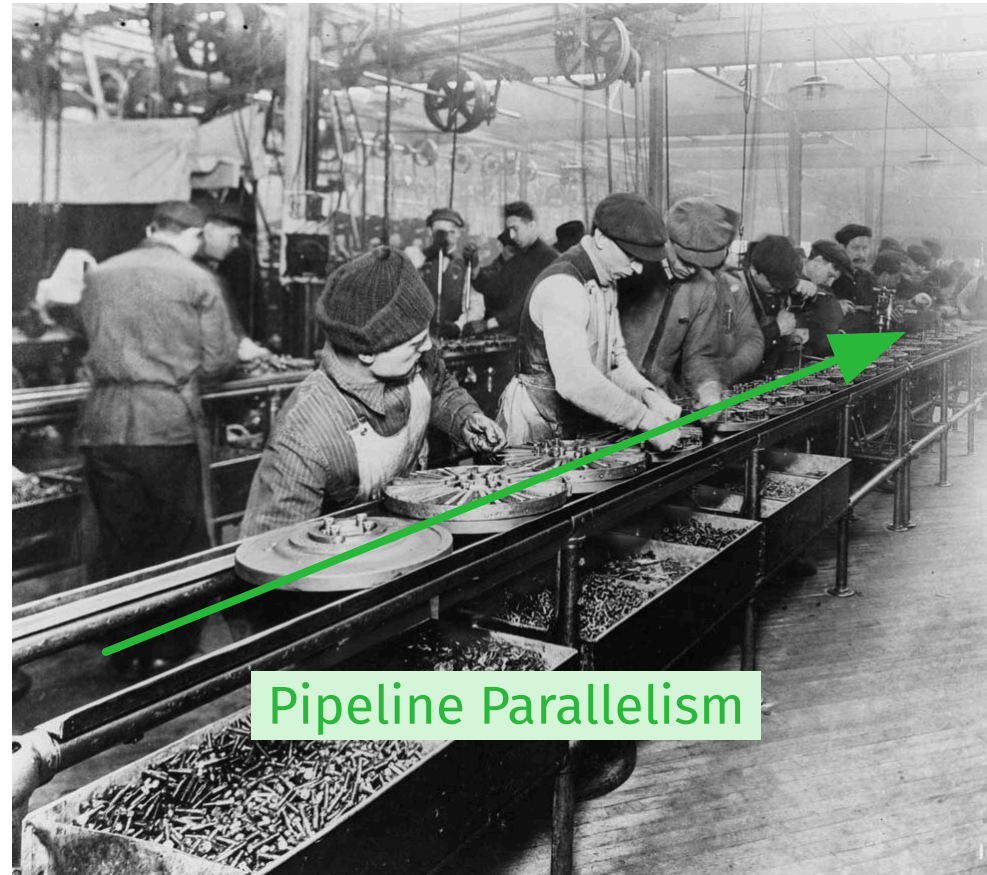
Data doesn't arrive fast enough.

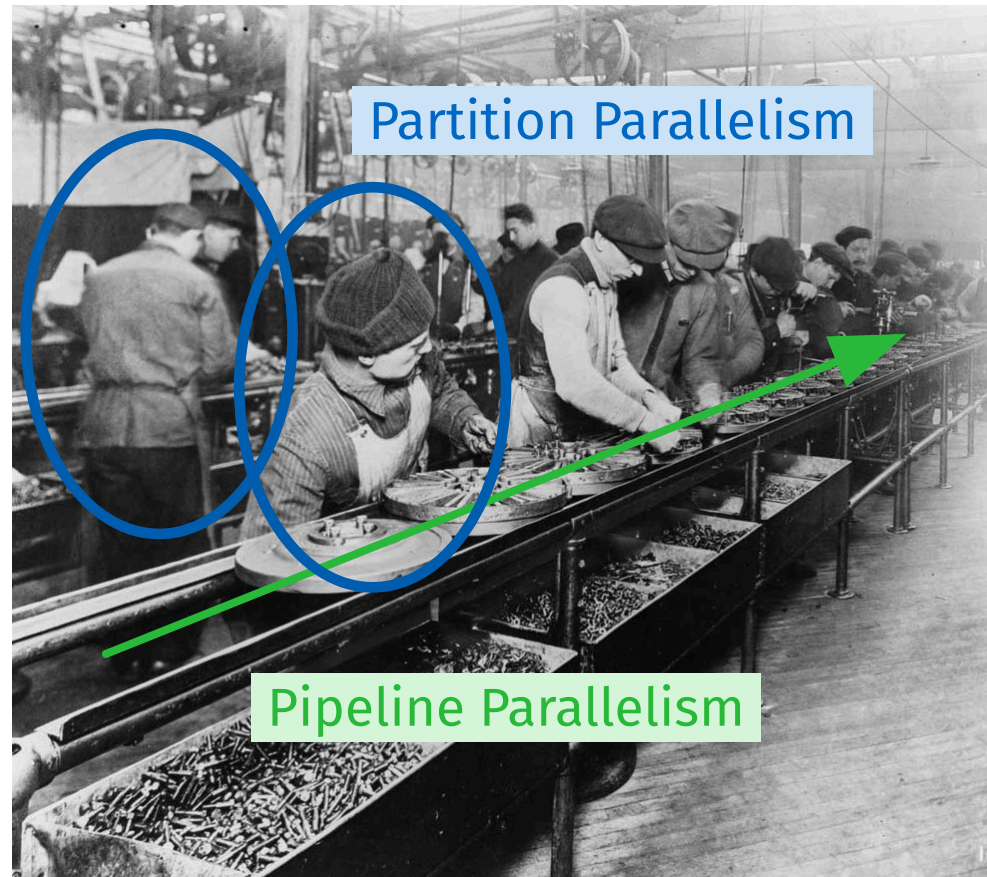


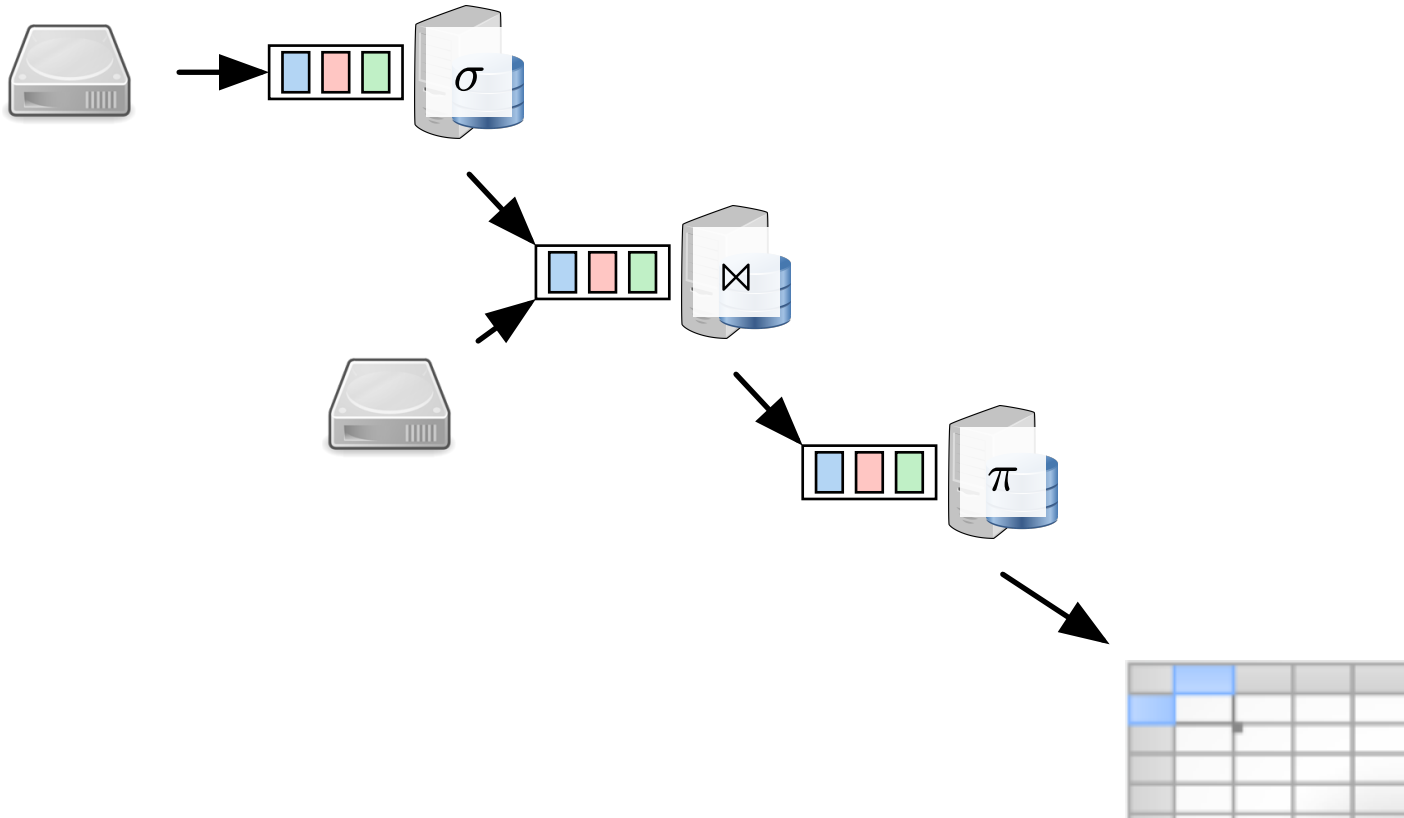
Worker processes data slower than it arrives.

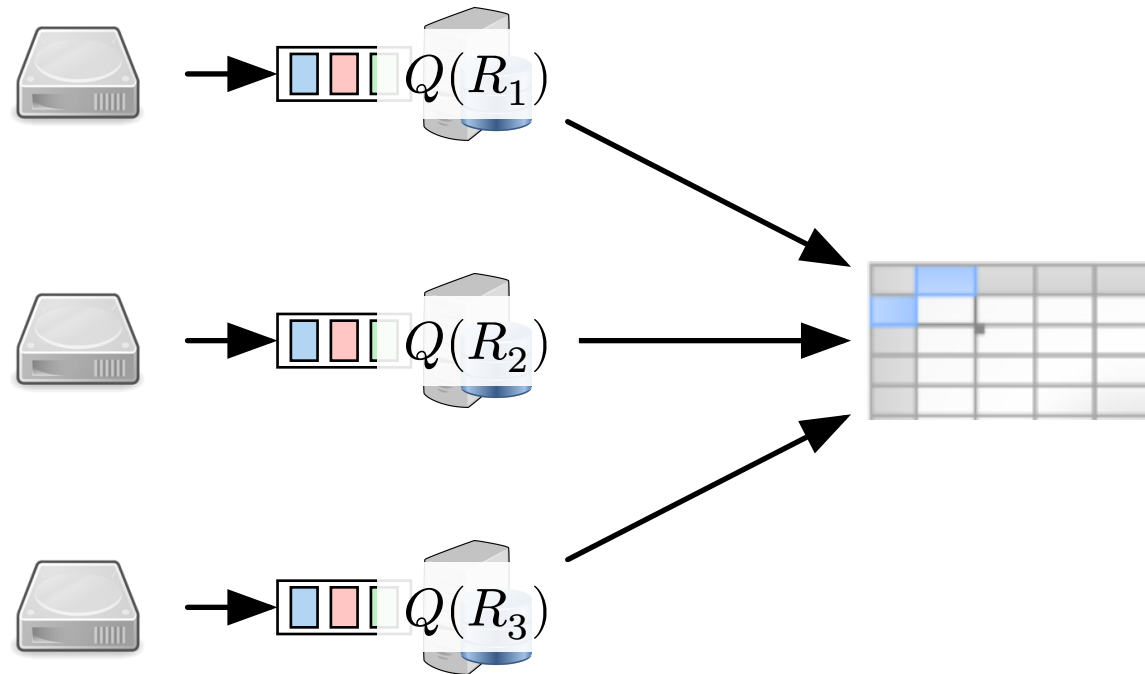


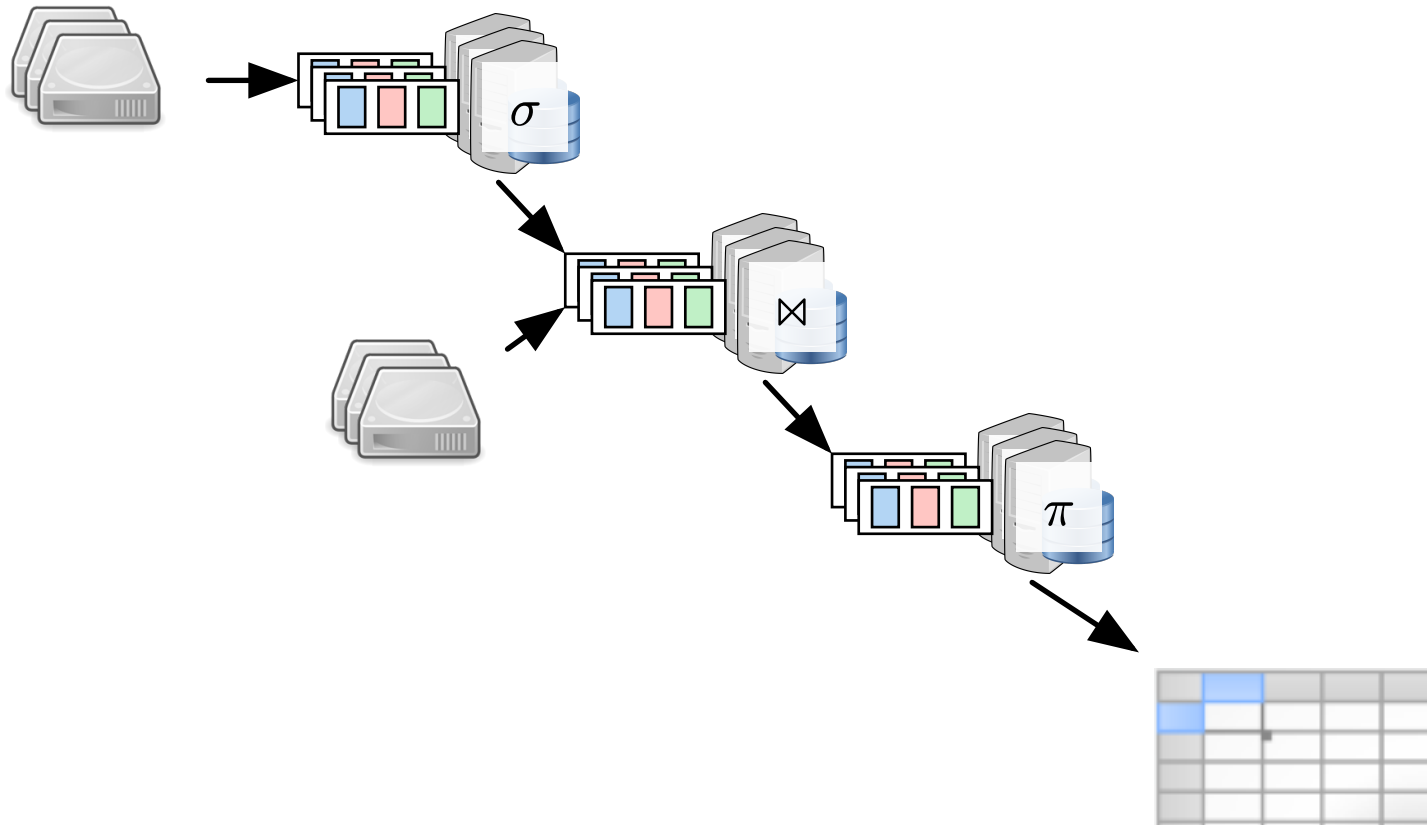
Image credit: [Wikipedia](#) - Public Domain Image - Photographer Unknown











Replication



Partitioning (Sharding)



How do we pick what goes into each partition?

- Arbitrary
- Range
- Hash

How do we pick what goes into each partition?

- Arbitrary
- Range
- Hash

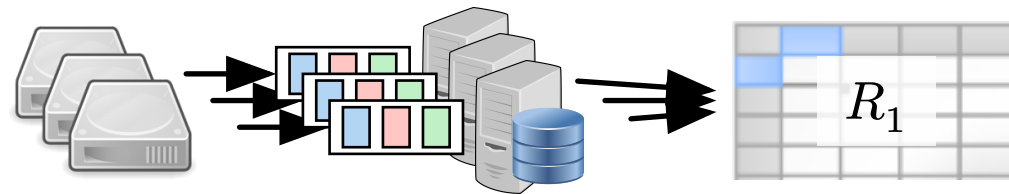
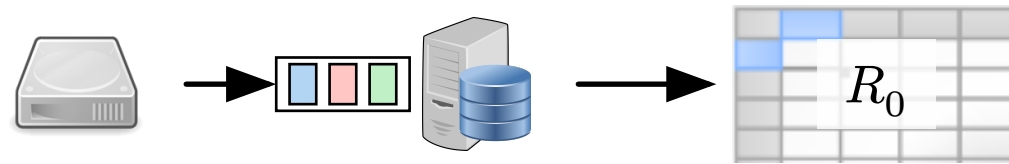
look familiar?

How do we pick what goes into each partition?

- Arbitrary
- Range
- Hash

look familiar?

How do we parallelize?



$$R_0 \stackrel{?}{=} R_1$$

Can we run each worker independently on one partition?

$$Q(R_0 \cup \dots \cup R_N) \stackrel{?}{=} Q(R_0) \cup \dots \cup Q(R_N)$$

Selection

- $\sigma(R_0 \cup \dots \cup R_N) =$

Projection

- $\pi(R_0 \cup \dots \cup R_N) =$

Union

- $(R_0 \cup \dots \cup R_N) =$

Selection

- $\sigma(R_0 \cup \dots \cup R_N) = \sigma(R_0) \cup \dots \cup \sigma(R_N)$

Projection

- $\pi(R_0 \cup \dots \cup R_N) =$

Union

- $(R_0 \cup \dots \cup R_N) =$

Selection

- $\sigma(R_0 \cup \dots \cup R_N) = \sigma(R_0) \cup \dots \cup \sigma(R_N)$

Projection

- $\pi(R_0 \cup \dots \cup R_N) = \pi(R_0) \cup \dots \cup \pi(R_N)$

Union

- $(R_0 \cup \dots \cup R_N) =$

Selection

- $\sigma(R_0 \cup \dots \cup R_N) = \sigma(R_0) \cup \dots \cup \sigma(R_N)$

Projection

- $\pi(R_0 \cup \dots \cup R_N) = \pi(R_0) \cup \dots \cup \pi(R_N)$

Union

- $(R_0 \cup \dots \cup R_N) = (R_0) \cup \dots \cup (R_N)$

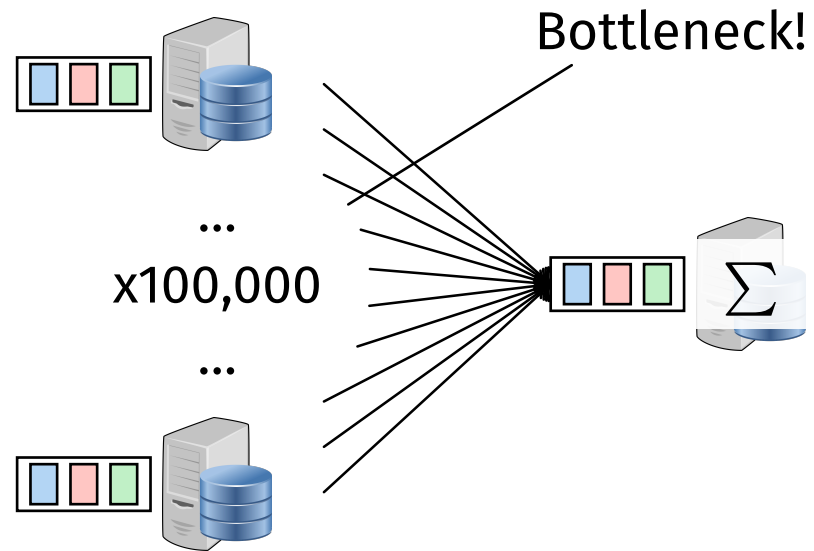
Aggregation

Single Row Aggregation

- $\sum(R_0 \cup \dots \cup R_N) =$

Single Row Aggregation

- $\sum(R_0 \cup \dots \cup R_N) = ???$



Algebraic Aggregates

- `Init()`: The “default” value of the aggregate
- `Accumulate(agg, value)`: Incorporate a new value
- `Finalize(agg) -> value`: Post-process the aggregate

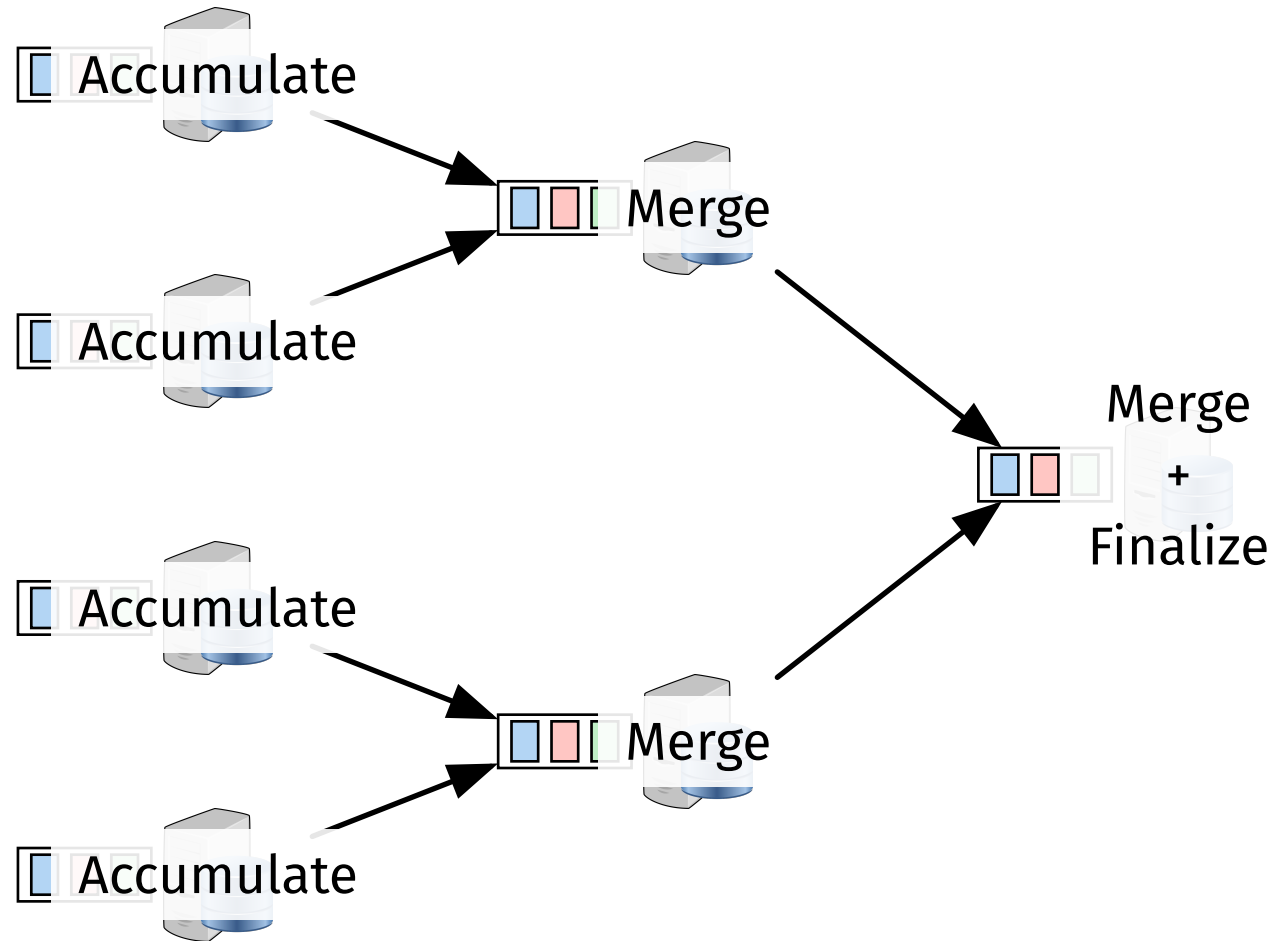
Algebraic Aggregates

- `Init()`: The “default” value of the aggregate
- `Accumulate(agg, value)`: Incorporate a new value
- `Finalize(agg) -> value`: Post-process the aggregate
- `Merge(agg, agg)`: Combine independent intermediate aggregates

Single Row Aggregation

- $\sum(R_0 \cup \dots \cup R_N) = \sum(R_0) + \dots + \sum(R_N)$

Final “Merge” step only processes N tuples



Multi-Row Aggregation

- **Arbitrary Partitioning:** As single-row
- **Hash/Range Partitioning:** Aggregate groups do not cross partitions

Joins

$$(R_1 \cup R_2) \bowtie (S_1 \cup S_2)$$

$$(R_1 \cup R_2) \bowtie (S_1 \cup S_2)$$

$$(R_1 \bowtie (S_1 \cup S_2)) \cup (R_2 \bowtie (S_1 \cup S_2))$$

$$(R_1 \cup R_2) \bowtie (S_1 \cup S_2)$$

$$(R_1 \bowtie (S_1 \cup S_2)) \cup (R_2 \bowtie (S_1 \cup S_2))$$

$$(R_1 \bowtie S_1) \cup (R_1 \bowtie S_2) \cup (R_2 \bowtie (S_1 \cup S_2))$$

$$(R_1 \bowtie S_1) \cup (R_1 \bowtie S_2) \cup (R_2 \bowtie S_1) \cup (R_2 \bowtie S_2)$$

$$\begin{aligned} & (R_1 \cup \dots \cup R_N) \bowtie (S_1 \cup \dots \cup S_K) \\ & = \\ & (R_1 \bowtie S_1) \cup \dots \cup (R_1 \bowtie S_K) \\ & \quad \dots \cup \dots \cup \dots \\ & (R_N \bowtie S_1) \cup \dots \cup (R_N \bowtie S_K) \end{aligned}$$

	S_0	S_1	S_2	S_3	S_4
R_0	$R_0 \bowtie S_0$	$R_0 \bowtie S_1$	$R_0 \bowtie S_2$	$R_0 \bowtie S_3$	$R_0 \bowtie S_4$
R_1	$R_1 \bowtie S_0$	$R_1 \bowtie S_1$	$R_1 \bowtie S_2$	$R_1 \bowtie S_3$	$R_1 \bowtie S_4$
R_2	$R_2 \bowtie S_0$	$R_2 \bowtie S_1$	$R_2 \bowtie S_2$	$R_2 \bowtie S_3$	$R_2 \bowtie S_4$
R_3	$R_3 \bowtie S_0$	$R_3 \bowtie S_1$	$R_3 \bowtie S_2$	$R_3 \bowtie S_3$	$R_3 \bowtie S_4$
R_4	$R_4 \bowtie S_0$	$R_4 \bowtie S_1$	$R_4 \bowtie S_2$	$R_4 \bowtie S_3$	$R_4 \bowtie S_4$

	S_0	S_1	S_2	S_3	S_4
R_0	$R_0 \bowtie S_0$	$R_0 \bowtie S_1$	$R_0 \bowtie S_2$	$R_0 \bowtie S_3$	$R_0 \bowtie S_4$
R_1	$R_1 \bowtie S_0$	$R_1 \bowtie S_1$	$R_1 \bowtie S_2$	$R_1 \bowtie S_3$	$R_1 \bowtie S_4$
R_2	$R_2 \bowtie S_0$	$R_2 \bowtie S_1$	$R_2 \bowtie S_2$	$R_2 \bowtie S_3$	$R_2 \bowtie S_4$
R_3	$R_3 \bowtie S_0$	$R_3 \bowtie S_1$	$R_3 \bowtie S_2$	$R_3 \bowtie S_3$	$R_3 \bowtie S_4$
R_4	$R_4 \bowtie S_0$	$R_4 \bowtie S_1$	$R_4 \bowtie S_2$	$R_4 \bowtie S_3$	$R_4 \bowtie S_4$

N workers gets us \sqrt{N} scaling

How do we pick what goes into each partition?

- Arbitrary
- Range
- Hash

For $R \bowtie_A S$:

- $R_i \leftarrow \sigma_{\text{hash}(A)=i}(R)$
- $S_i \leftarrow \sigma_{\text{hash}(A)=i}(S)$

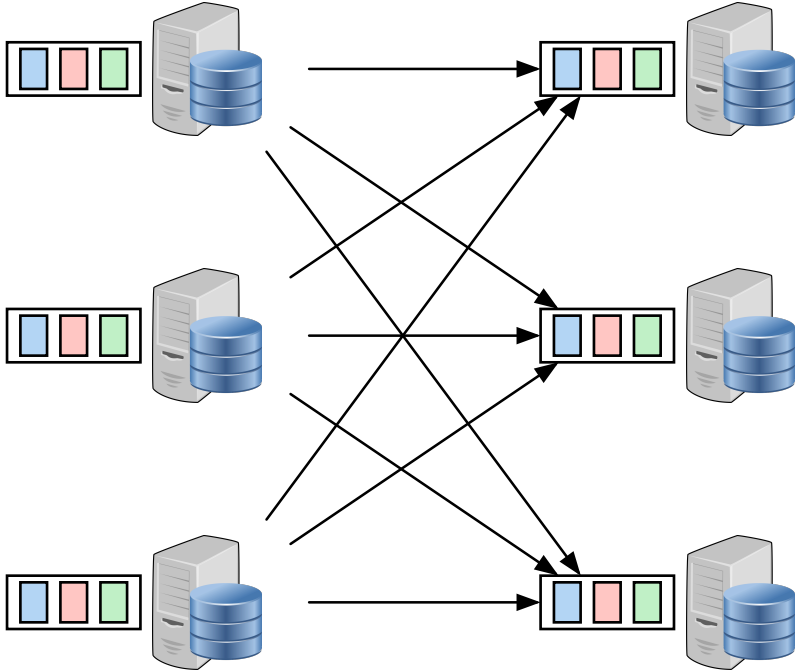
	S_0	S_1	S_2	S_3	S_4
R_0	$R_0 \bowtie S_0$	$R_0 \bowtie S_1$	$R_0 \bowtie S_2$	$R_0 \bowtie S_3$	$R_0 \bowtie S_4$
R_1	$R_1 \bowtie S_0$	$R_1 \bowtie S_1$	$R_1 \bowtie S_2$	$R_1 \bowtie S_3$	$R_1 \bowtie S_4$
R_2	$R_2 \bowtie S_0$	$R_2 \bowtie S_1$	$R_2 \bowtie S_2$	$R_2 \bowtie S_3$	$R_2 \bowtie S_4$
R_3	$R_3 \bowtie S_0$	$R_3 \bowtie S_1$	$R_3 \bowtie S_2$	$R_3 \bowtie S_3$	$R_3 \bowtie S_4$
R_4	$R_4 \bowtie S_0$	$R_4 \bowtie S_1$	$R_4 \bowtie S_2$	$R_4 \bowtie S_3$	$R_4 \bowtie S_4$

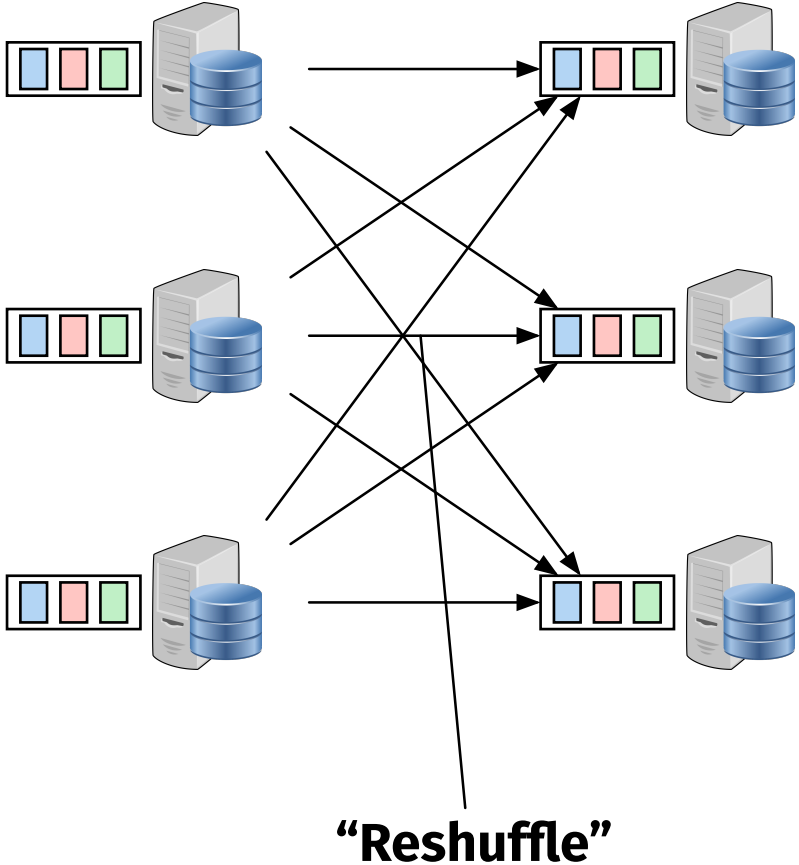
	S_0	S_1	S_2	S_3	S_4
R_0	$R_0 \bowtie S_0$	$R_0 \bowtie S_1$	$R_0 \bowtie S_2$	$R_0 \bowtie S_3$	$R_0 \bowtie S_4$
R_1	$R_1 \bowtie S_0$	$R_1 \bowtie S_1$	$R_1 \bowtie S_2$	$R_1 \bowtie S_3$	$R_1 \bowtie S_4$
R_2	$R_2 \bowtie S_0$	$R_2 \bowtie S_1$	$R_2 \bowtie S_2$	$R_2 \bowtie S_3$	$R_2 \bowtie S_4$
R_3	$R_3 \bowtie S_0$	$R_3 \bowtie S_1$	$R_3 \bowtie S_2$	$R_3 \bowtie S_3$	$R_3 \bowtie S_4$
R_4	$R_4 \bowtie S_0$	$R_4 \bowtie S_1$	$R_4 \bowtie S_2$	$R_4 \bowtie S_3$	$R_4 \bowtie S_4$

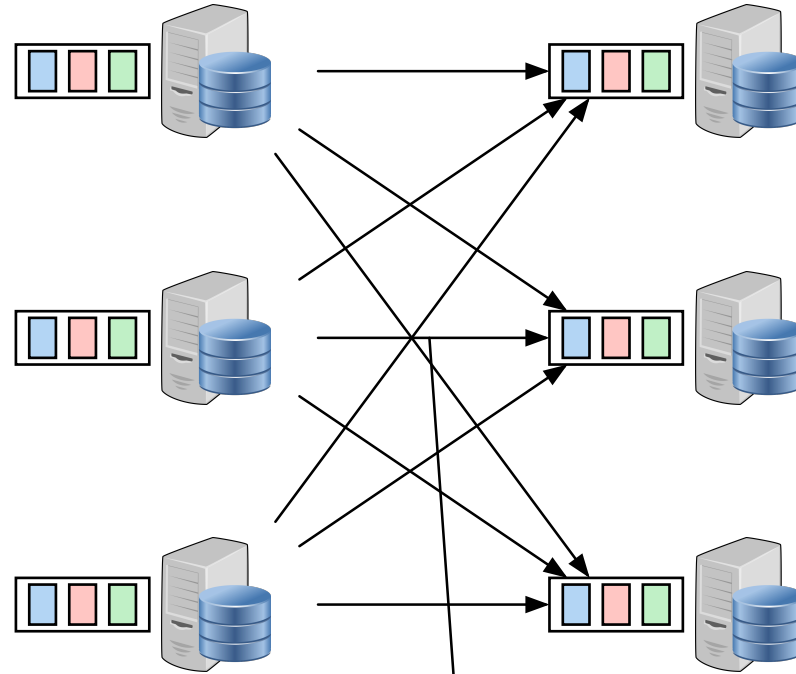
	S_0	S_1	S_2	S_3	S_4
R_0	$R_0 \bowtie S_0$	$R_0 \bowtie S_1$	$R_0 \bowtie S_2$	$R_0 \bowtie S_3$	$R_0 \bowtie S_4$
R_1	$R_1 \bowtie S_0$	$R_1 \bowtie S_1$	$R_1 \bowtie S_2$	$R_1 \bowtie S_3$	$R_1 \bowtie S_4$
R_2	$R_2 \bowtie S_0$	$R_2 \bowtie S_1$	$R_2 \bowtie S_2$	$R_2 \bowtie S_3$	$R_2 \bowtie S_4$
R_3	$R_3 \bowtie S_0$	$R_3 \bowtie S_1$	$R_3 \bowtie S_2$	$R_3 \bowtie S_3$	$R_3 \bowtie S_4$
R_4	$R_4 \bowtie S_0$	$R_4 \bowtie S_1$	$R_4 \bowtie S_2$	$R_4 \bowtie S_3$	$R_4 \bowtie S_4$

Hash gets us N scaling for N workers

What if the partitions aren't aligned so nicely?







“Reshuffle”

(Each tuple goes to exactly one place)

Reshuffling

- R_i is initially located at **X**
- S_i is initially located at **Y**

Then...

1. Copy S_i from **Y** to **X**
2. Compute $R_i \bowtie S_i$ at **X**

**Can we reduce network
use more?**

Worker 1

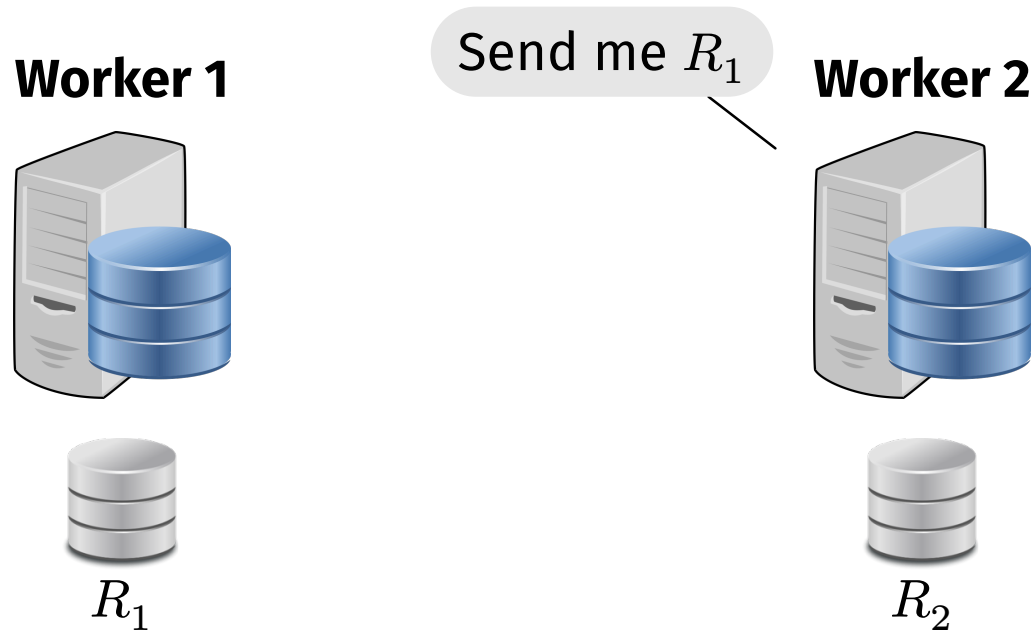


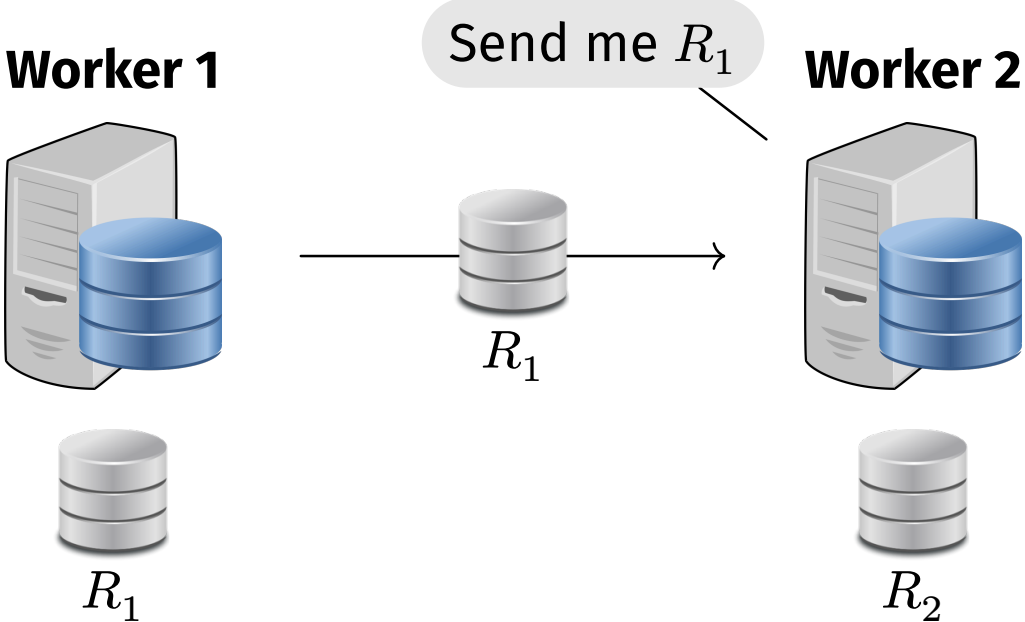
R_1

Worker 2



R_2





Worker 1

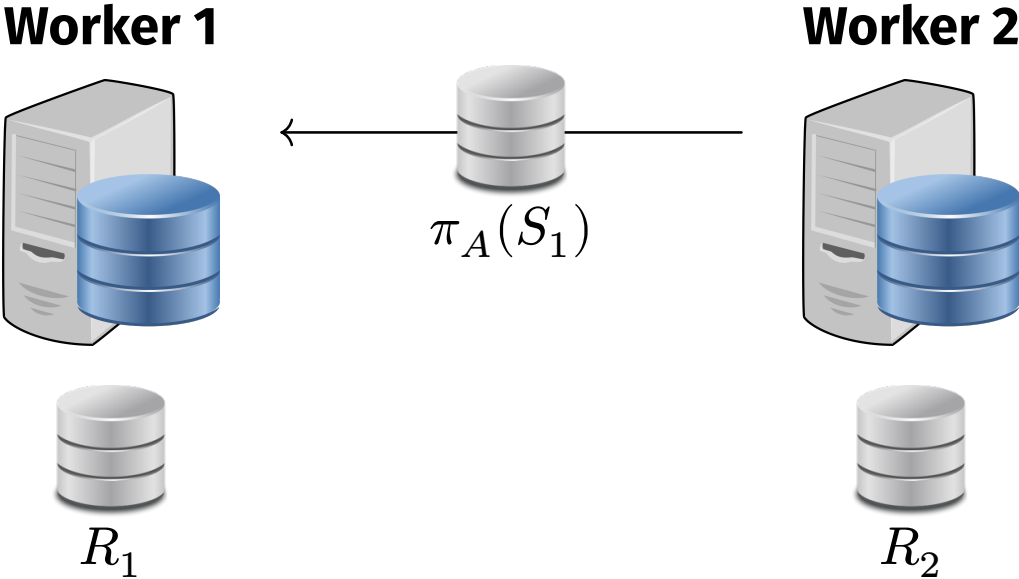


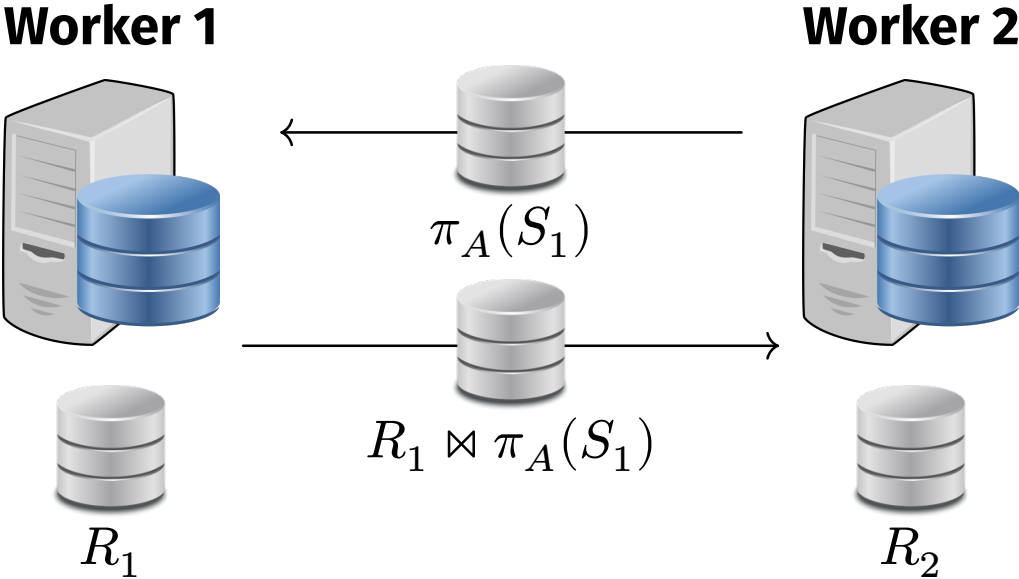
R_1

Worker 2



R_2





Worker 1



<M, 1>
<N, 2>
<O, 2>
<P, 3>
<Q, 4>

Worker 2



<2, X>
<3, Y>
<6, Y>

Worker 1



<M, 1>

<N, 2>

<O, 2>

<P, 3>

<Q, 4>

Send me rows
where $A \in \{2, 3, 6\}$

Worker 2



<2, X>

<3, Y>

<6, Y>

Worker 1



<M, 1>
<N, 2>
<O, 2>
<P, 3>
<Q, 4>

<N, 2>
<O, 2>
<P, 3>

Send me rows
where $A \in \{2, 3, 6\}$

Worker 2



<2, X>
<3, Y>
<6, Y>

Worker 1



<M, 1>
<N, 2>
<O, 2>
<P, 3>
<Q, 4>

Send me rows
where $A \in \{2, 3, 6\}$

<N, 2>
<O, 2>
<P, 3>

Worker 2



<2, X>
<3, Y>
<6, Y>

This is called a **semi-join**

Problem: Worker 2 is sending a lot of data

Problem: Worker 2 is sending a lot of data

Idea: Compress $\pi_A(S_1)$ using a bloom filter

Worker 1



<M, 1>

<N, 2>

<O, 2>

<P, 3>

<Q, 4>

Worker 2



<2, X>

<3, Y>

<6, Y>

Worker 1



- <M, 1>
- <N, 2>
- <O, 2>
- <P, 3>
- <Q, 4>

Send me rows
where $A \in \text{Bloom}(\{2, 3, 6\})$

Worker 2



- <2, X>
- <3, Y>
- <6, Y>

Worker 1



- <M, 1>
- <N, 2>
- <O, 2>
- <P, 3>
- <Q, 4>

Send me rows
where $A \in \text{Bloom}(\{2, 3, 6\})$

- <N, 2>
- <O, 2>
- <P, 3>
- <Q, 4>

Worker 2



- <2, X>
- <3, Y>
- <6, Y>

Worker 1



- <M, 1>
- <N, 2>
- <O, 2>
- <P, 3>
- <Q, 4>

Send me rows
where $A \in \text{Bloom}(\{2, 3, 6\})$

- <N, 2>
- <O, 2>
- <P, 3>
- <Q, 4>

Worker 2



- <2, X>
- <3, Y>
- <6, Y>

This is called a **Bloom Join**

- Deadline for checkpoint 3 extended by one week
- No class on Thursday (video lecture to be released soon)