

DATA516/CSED516  
Scalable Data Systems and  
Algorithms

Lecture 4

Distributed Query Evaluation

# Announcements

- Project proposals due on Friday
- Reviews due every week
- HW2 due next Monday

Coming soon: guest lecturers!

- Cloud Databases: Shan Shan Huang, RelationalAI
- Graph Databases: Mingxi Wu, Tigergraph

# Distributed Query Processing Algorithms

# Horizontal Data Partitioning

- **Block Partition, a.k.a. Round Robin:**
  - Partition tuples arbitrarily s.t.  $\text{size}(R_1) \approx \dots \approx \text{size}(R_P)$
- **Hash partitioned on attribute A:**
  - Tuple  $t$  goes to chunk  $i$ , where  $i = h(t.A) \bmod P + 1$
- **Range partitioned on attribute A:**
  - Partition the range of  $A$  into  $-\infty = v_0 < v_1 < \dots < v_P = \infty$
  - Tuple  $t$  goes to chunk  $i$ , if  $v_{i-1} < t.A < v_i$

# Notation

When a relation  $R$  is distributed to  $p$  servers, we draw the picture like this:



Here  $R_1$  is the fragment of  $R$  stored on server 1, etc

$$R = R_1 \cup R_2 \cup \dots \cup R_p$$

# Uniform Load and Skew

- $|R| = N$  tuples, then  $|R_1| + |R_2| + \dots + |R_p| = N$
- We say the load is uniform when:  
$$|R_1| \approx |R_2| \approx \dots \approx |R_p| \approx N/p$$
- Skew means that some load is much larger:  
$$\max_i |R_i| \gg N/p$$

We design algorithms for uniform load, discuss skew later

# Parallel Algorithm

- Selection  $\sigma$
- Join  $\bowtie$
- Group by  $\gamma$

# Parallel Selection

Data:  $R(\underline{K}, A, B, C)$

Query:  $\sigma_{A=v}(R)$ , or  $\sigma_{v1 < A < v2}(R)$

- Block partitioned:
  - All servers must scan and filter the data
- Hash partitioned:
  - Can have all servers scan and filter the data
  - Or can optimize and only have some servers do work
- Range partitioned
  - Also only some servers need to do the work



# Parallel GroupBy

Data:  $R(\underline{K}, A, B, C)$

Query:  $\gamma_{A, \text{sum}(C)}(R)$

- Discuss in class how to compute in each case:
- R is hash-partitioned on A
- R is block-partitioned or hash-partitioned on K

# Parallel GroupBy

Data:  $R(\underline{K}, A, B, C)$

Query:  $\gamma_{A, \text{sum}(C)}(R)$

- Discuss in class how to compute in each case:
- R is hash-partitioned on A
  - Each server  $i$  computes locally  $\gamma_{A, \text{sum}(C)}(R_i)$
- R is block-partitioned or hash-partitioned on K
  - Need to reshuffle data on A first (next slide)
  - Then compute locally  $\gamma_{A, \text{sum}(C)}(R_i)$

# Basic Parallel GroupBy

Data:  $R(\underline{K}, A, B, C)$

Query:  $\gamma_{A, \text{sum}(C)}(R)$

- R is block-partitioned or hash-partitioned on K



# Basic Parallel GroupBy

Data:  $R(\underline{K}, A, B, C)$

Query:  $\gamma_{A, \text{sum}(C)}(R)$

- R is block-partitioned or hash-partitioned on K

Reshuffle R  
on attribute A

$R_1$

$R_2$

$R_p$

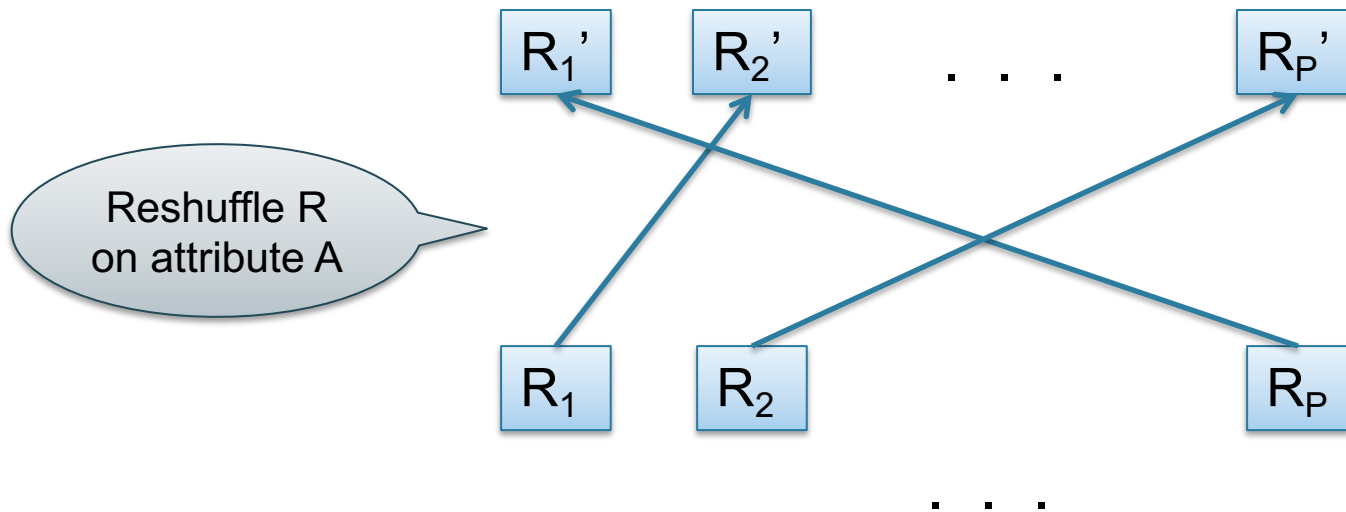
...

# Basic Parallel GroupBy

Data:  $R(\underline{K}, A, B, C)$

Query:  $\gamma_{A, \text{sum}(C)}(R)$

- $R$  is block-partitioned or hash-partitioned on  $K$

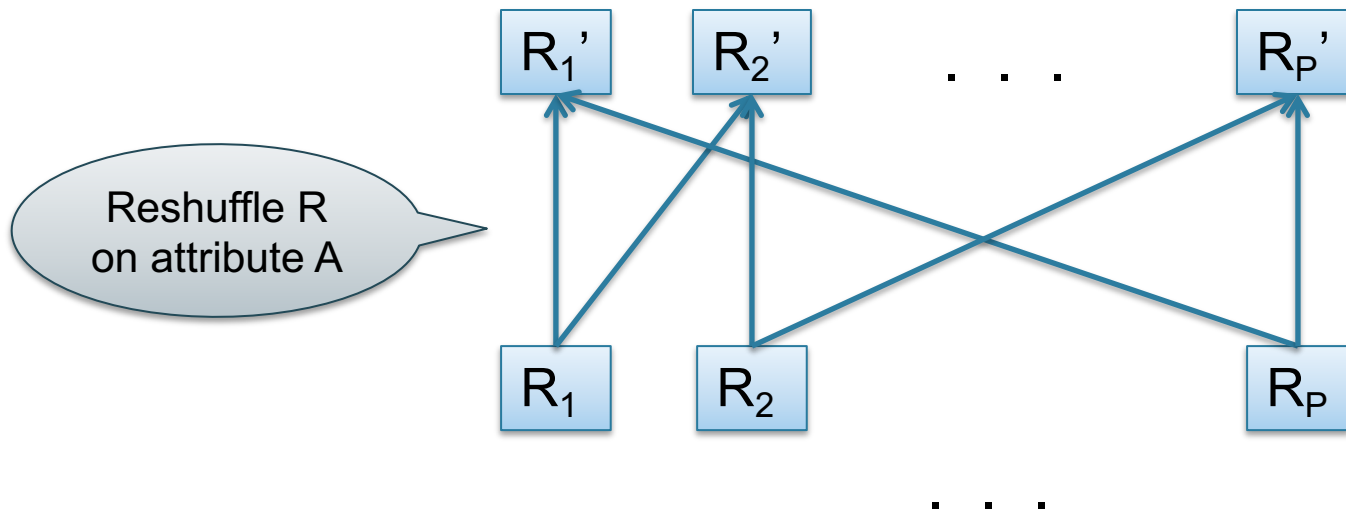


# Basic Parallel GroupBy

Data:  $R(\underline{K}, A, B, C)$

Query:  $\gamma_{A, \text{sum}(C)}(R)$

- $R$  is block-partitioned or hash-partitioned on  $K$

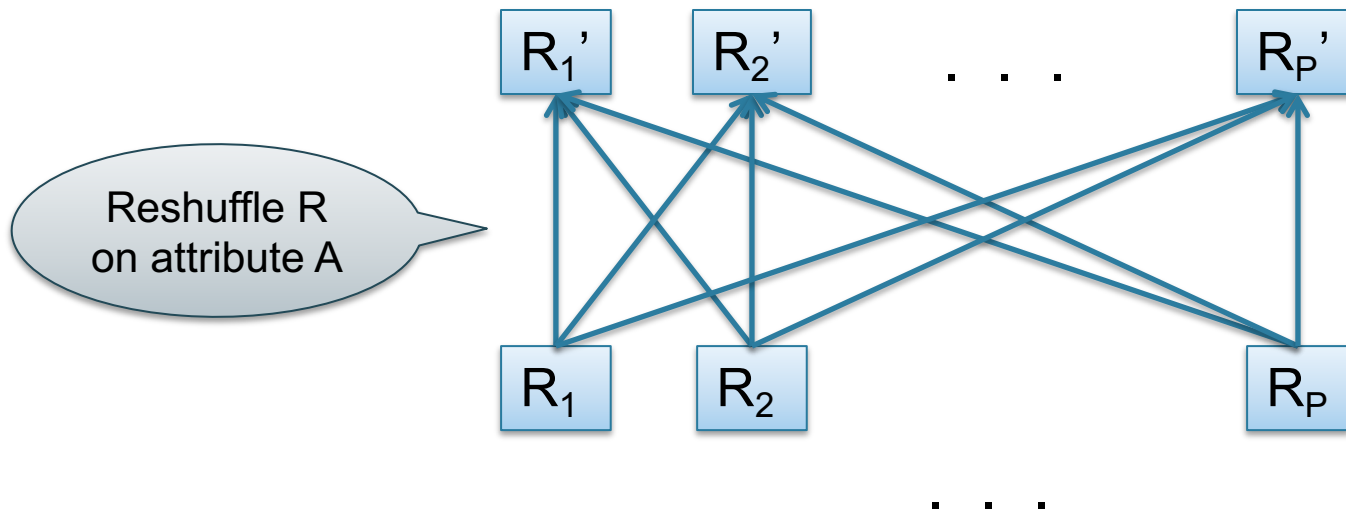


# Basic Parallel GroupBy

Data:  $R(\underline{K}, A, B, C)$

Query:  $\gamma_{A, \text{sum}(C)}(R)$

- $R$  is block-partitioned or hash-partitioned on  $K$

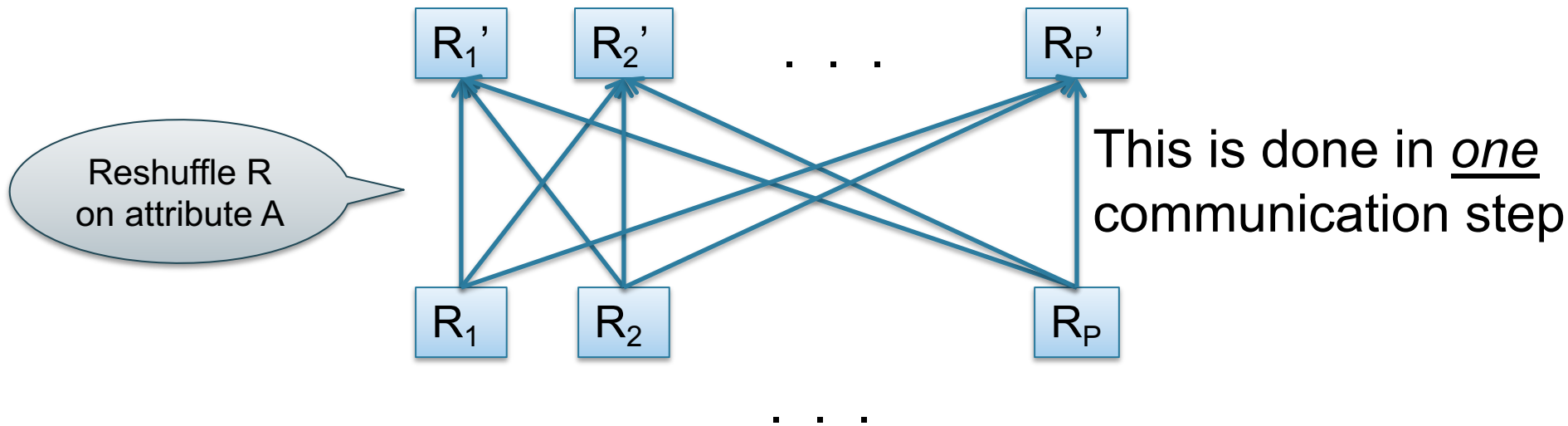


# Basic Parallel GroupBy

Data:  $R(\underline{K}, A, B, C)$

Query:  $\gamma_{A, \text{sum}(C)}(R)$

- $R$  is block-partitioned or hash-partitioned on  $K$



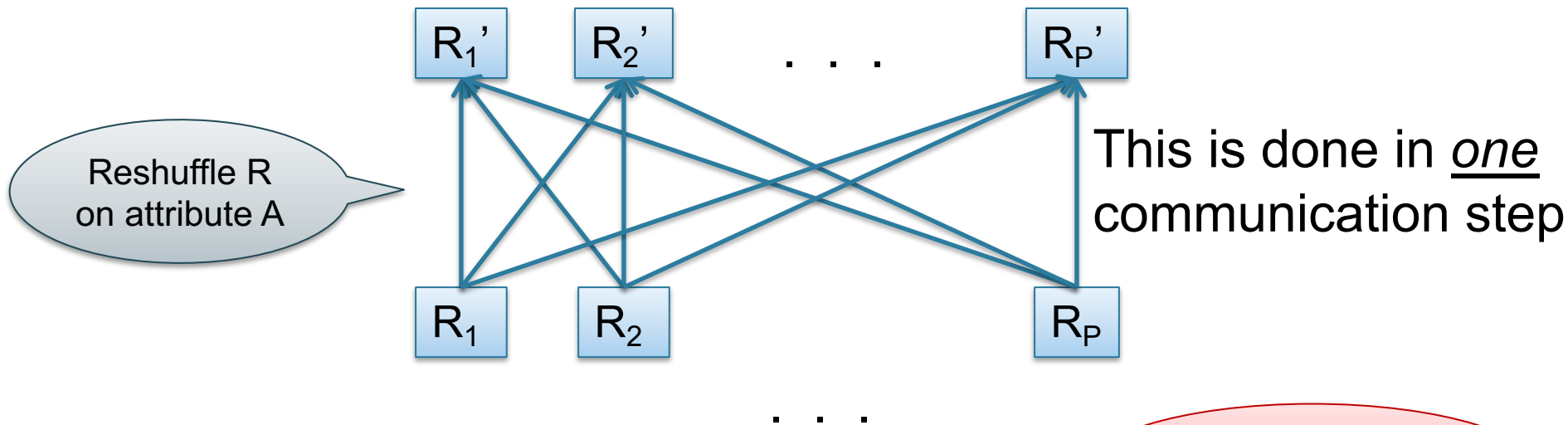


# Basic Parallel GroupBy

Data:  $R(\underline{K}, A, B, C)$

Query:  $\gamma_{A, \text{sum}(C)}(R)$

- $R$  is block-partitioned or hash-partitioned on  $K$



# Basic Parallel GroupBy

Data:  $R(\underline{K}, A, B, C)$

Query:  $\gamma_{A, \text{sum}(C)}(R)$

# Basic Parallel GroupBy

Data:  $R(\underline{K}, A, B, C)$

Query:  $Y_{A, \text{sum}(C)}(R)$

**Step 0:** [**Optimization**] each server  $i$  computes local group-by:

$$T_i = Y_{A, \text{sum}(C)}(R_i)$$

# Basic Parallel GroupBy

Data:  $R(\underline{K}, A, B, C)$

Query:  $\gamma_{A, \text{sum}(C)}(R)$

**Step 0:** [**Optimization**] each server  $i$  computes local group-by:

$$T_i = \gamma_{A, \text{sum}(C)}(R_i)$$

**Step 1:** partitions tuples in  $T_i$  using hash function  $h(A)$ :

$T_{i,1}, T_{i,2}, \dots, T_{i,p}$   
then send fragment  $T_{i,j}$  to server  $j$

# Basic Parallel GroupBy

**Data:**  $R(\underline{K}, A, B, C)$

**Query:**  $\gamma_{A, \text{sum}(C)}(R)$

**Step 0:** [**Optimization**] each server  $i$  computes local group-by:

$$T_i = \gamma_{A, \text{sum}(C)}(R_i)$$

**Step 1:** partitions tuples in  $T_i$  using hash function  $h(A)$ :

$T_{i,1}, T_{i,2}, \dots, T_{i,p}$   
then send fragment  $T_{i,j}$  to server  $j$

**Step 2:** receive fragments, union them, then group-by

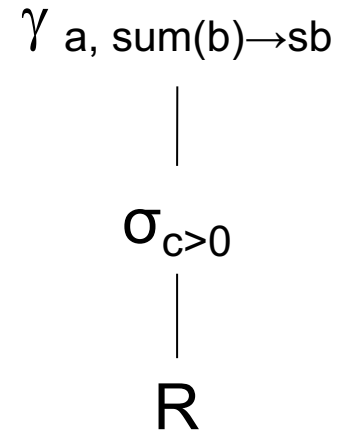
$$R'_j = T_{1,j} \cup \dots \cup T_{p,j}$$
$$\text{Answer}_j = \gamma_{A, \text{sum}(C)}(R'_j)$$

# Example Query with Group By

```
SELECT a, sum(b) as sb  
FROM R WHERE c > 0  
GROUP BY a
```

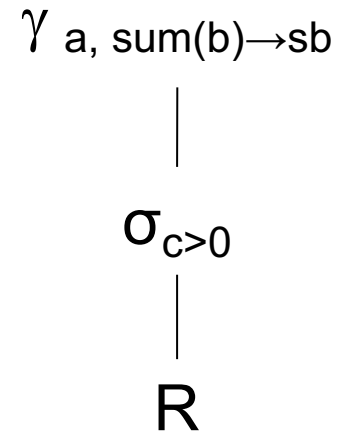
# Example Query with Group By

```
SELECT a, sum(b) as sb  
FROM R WHERE c > 0  
GROUP BY a
```



# Example Query with Group By

```
SELECT a, sum(b) as sb  
FROM R WHERE c > 0  
GROUP BY a
```



Machine 1

1/3 of R

Machine 2

1/3 of R

Machine 3

1/3 of R



```
SELECT a, sum(b) as sb FROM R WHERE c > 0 GROUP BY a
```

Machine 1

1/3 of R

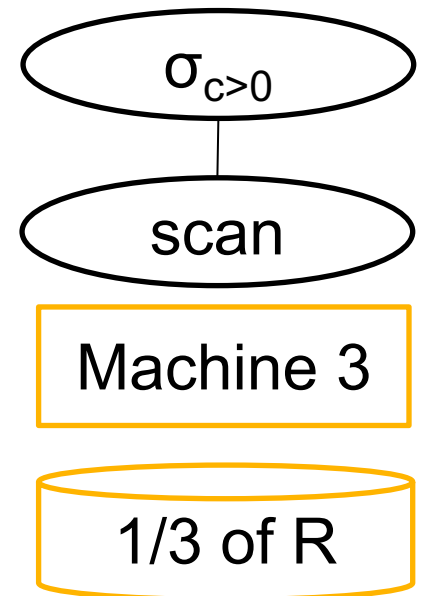
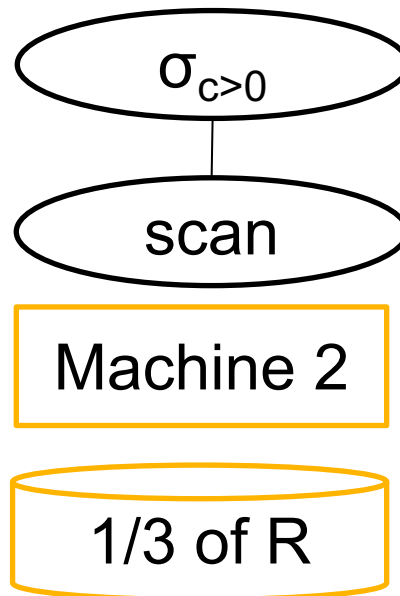
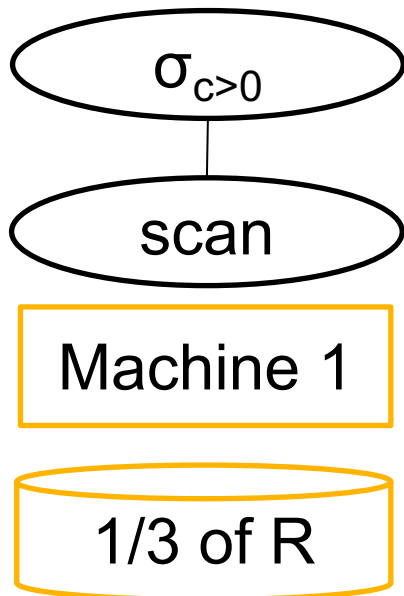
Machine 2

1/3 of R

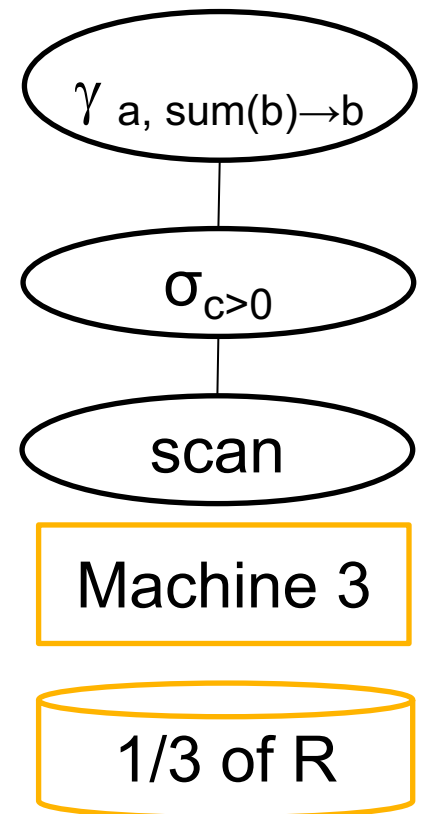
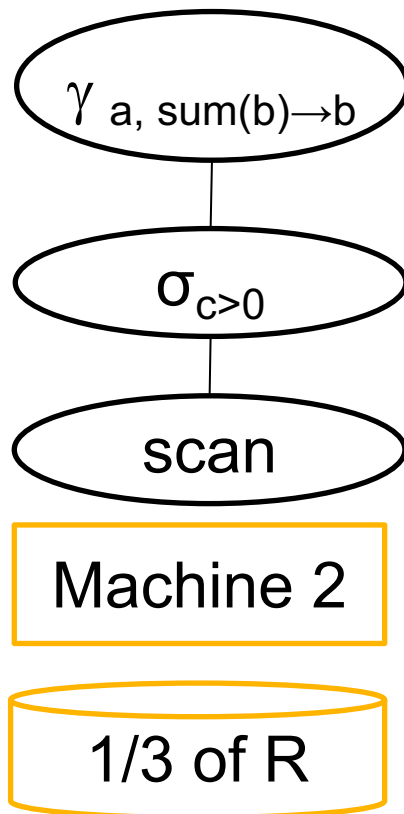
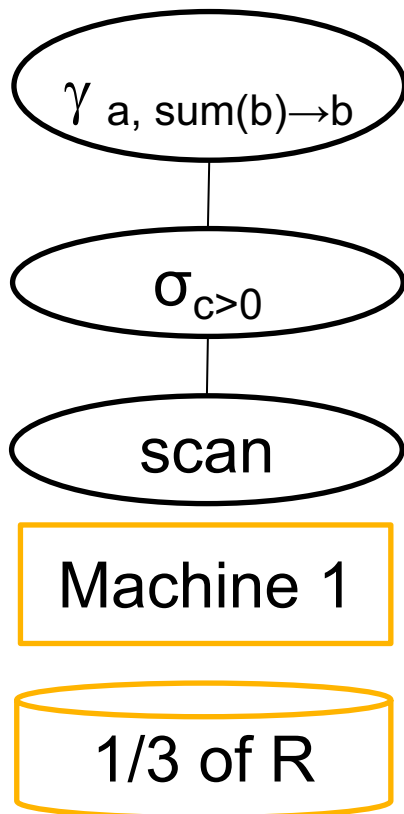
Machine 3

1/3 of R

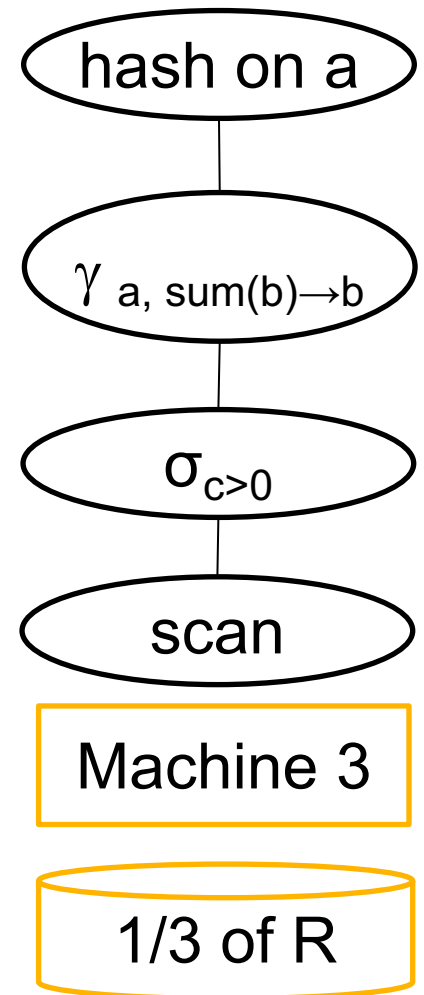
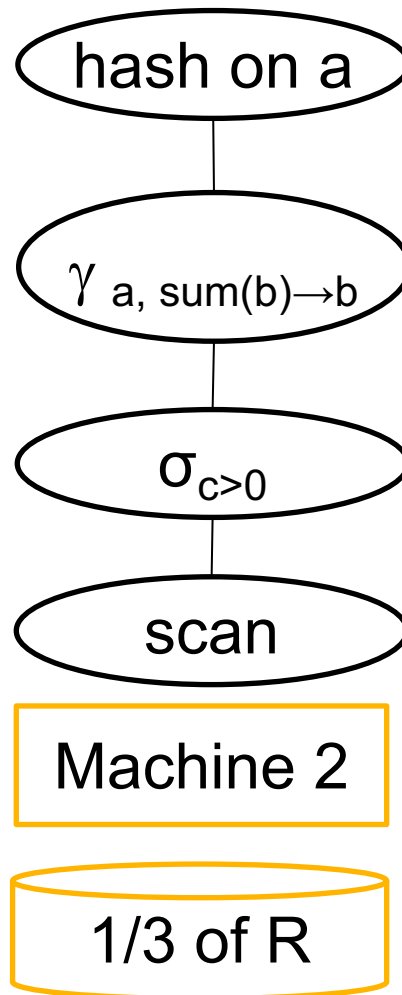
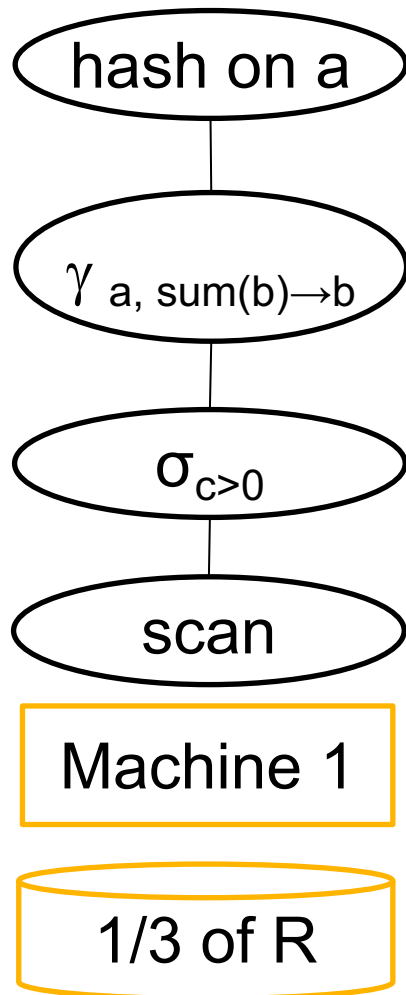
```
SELECT a, sum(b) as sb FROM R WHERE c > 0 GROUP BY a
```



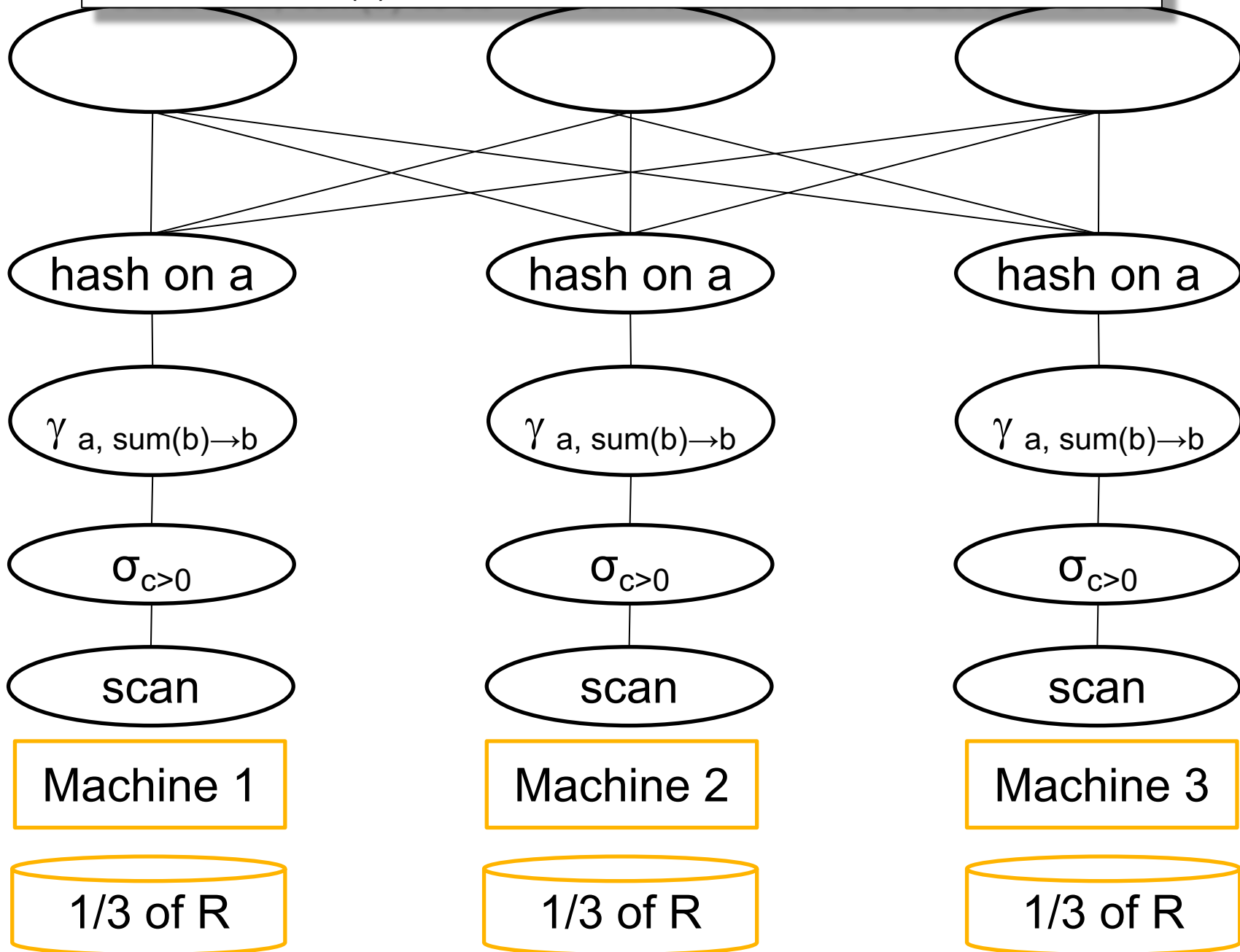
SELECT a, sum(b) as sb FROM R WHERE c > 0 GROUP BY a



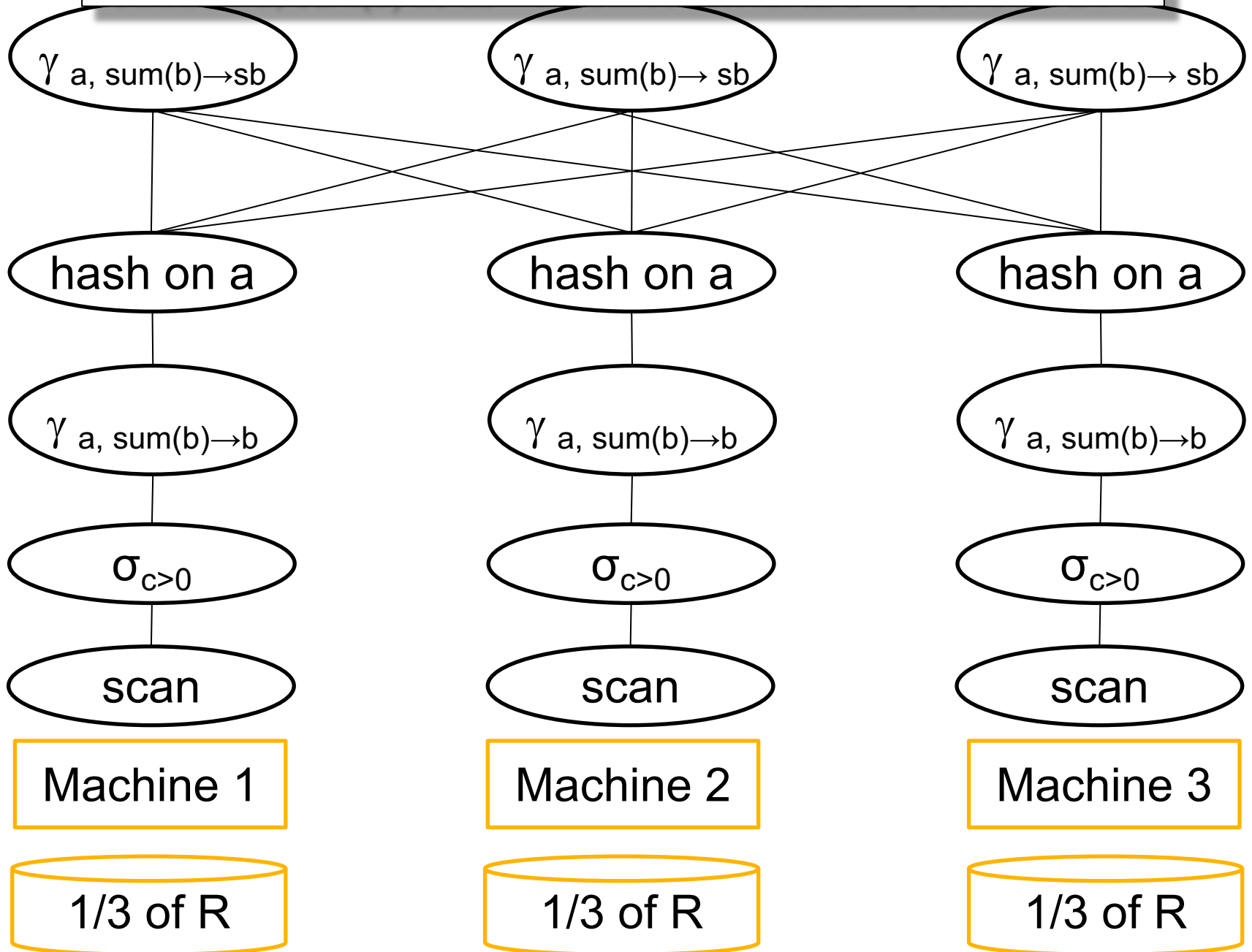
```
SELECT a, sum(b) as sb FROM R WHERE c > 0 GROUP BY a
```



SELECT a, sum(b) as sb FROM R WHERE c > 0 GROUP BY a



SELECT a, sum(b) as sb FROM R WHERE c > 0 GROUP BY a



# Pushing Aggregates Past Union

The rule that allowed us to do early summation is:

$$\begin{aligned}\gamma_{A, \text{sum}(B) \rightarrow C}(R_1 \cup R_2) &= \\ &= \gamma_{A, \text{sum}(D) \rightarrow B}(\gamma_{A, \text{sum}(B) \rightarrow D}(R_1) \cup \gamma_{A, \text{sum}(B) \rightarrow D}(R_2))\end{aligned}$$

For example:

- $R_1$  has  $B = x, y, z$ ;  $R_2$  has  $B = u, w$
- Then:  $x + y + z + u + w = (x + y + z) + (u + w)$

# Pushing Aggregates Past Union

Which other rules can we push past union?

- Sum?
- Count?
- Avg?
- Max?
- Median?



# Pushing Aggregates Past Union

Which other rules can we push past union?

- Sum?
- Count?
- Avg?
- Max?
- Median?

Distributive	Algebraic	Holistic
$\text{sum}(a_1+a_2+\dots+a_9)=$ $\text{sum}(\text{sum}(a_1+a_2+a_3)+$ $\text{sum}(a_4+a_5+a_6)+$ $\text{sum}(a_7+a_8+a_9))$	$\text{avg}(B) =$ $\text{sum}(B)/\text{count}(B)$	$\text{median}(B)$

# Speedup and Scaleup

Consider the query  $\gamma_{A, \text{sum}(C)}(R)$

Assume the local runtime for group-by is linear  $O(|R|)$

If we double number of nodes  $P$ , what is the runtime?

If we double both  $P$  and size of  $R$ , what is the runtime?

# Speedup and Scaleup

Consider the query  $\gamma_{A, \text{sum}(C)}(R)$

Assume the local runtime for group-by is linear  $O(|R|)$

If we double number of nodes  $P$ , what is the runtime?

- Half (chunk sizes become  $\frac{1}{2}$ )

If we double both  $P$  and size of  $R$ , what is the runtime?

- Same (chunk sizes remain the same)

# Speedup and Scaleup

Consider the query  $Y_{A, \text{sum}(C)}(R)$

Assume the local runtime for group-by is linear  $O(|R|)$

If we double number of nodes  $P$ , what is the runtime?

- Half (chunk sizes become  $\frac{1}{2}$ )

If we double both  $P$  and size of  $R$ , what is the runtime?

- Same (chunk sizes remain the same)

But only if the data is without skew!

# Parallel Join: $R \bowtie_{A=B} S$

Data:  $R(\underline{K1}, A, C), S(\underline{K2}, B, D)$

Query:  $R \bowtie_{A=B} S$



Initially, R and S are block partitioned.

Notice: they may be stored in DFS (recall MapReduce)

Some servers hold R-chunks, some hold S-chunks, some hold both

# Parallel Join: $R \bowtie_{A=B} S$

Data:  $R(\underline{K1}, A, C), S(\underline{K2}, B, D)$

Query:  $R \bowtie_{A=B} S$

Reshuffle R on R.A  
and S on S.B

$R_1, S_1$

$R_2, S_2$

...

$R_P, S_P$

Initially, R and S are block partitioned.

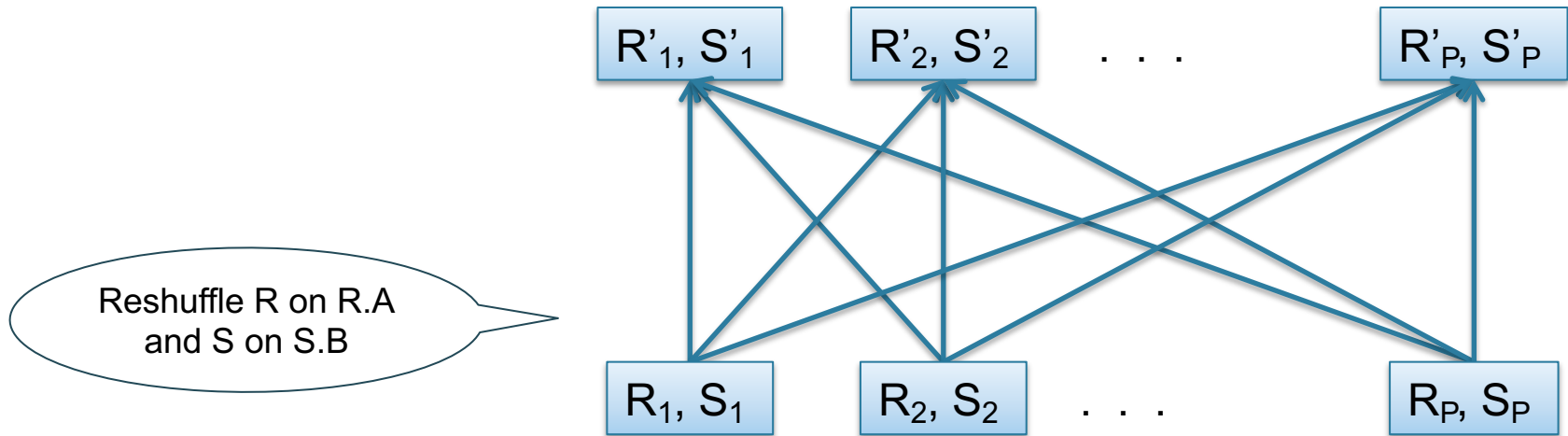
Notice: they may be stored in DFS (recall MapReduce)

Some servers hold R-chunks, some hold S-chunks, some hold both

# Parallel Join: $R \bowtie_{A=B} S$

Data:  $R(\underline{K1}, A, C), S(\underline{K2}, B, D)$

Query:  $R \bowtie_{A=B} S$



Initially, R and S are block partitioned.

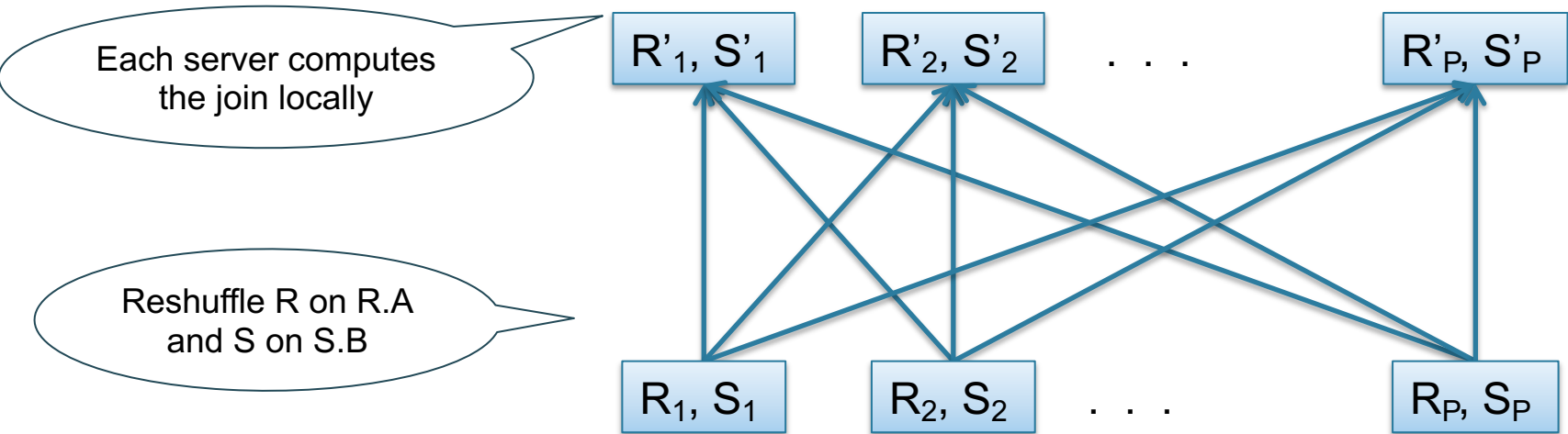
Notice: they may be stored in DFS (recall MapReduce)

Some servers hold R-chunks, some hold S-chunks, some hold both

# Parallel Join: $R \bowtie_{A=B} S$

Data:  $R(\underline{K1}, A, C), S(\underline{K2}, B, D)$

Query:  $R \bowtie_{A=B} S$



Initially, R and S are block partitioned.

Notice: they may be stored in DFS (recall MapReduce)

Some servers hold R-chunks, some hold S-chunks, some hold both



# Parallel Join: $R \bowtie_{A=B} S$

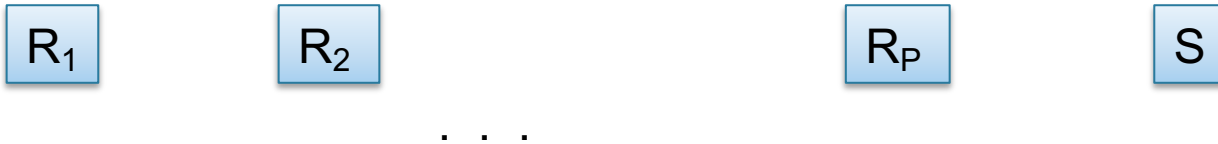
- Step 1
  - Every server holding any chunk of R partitions its chunk using a hash function  $h(t.A)$
  - Every server holding any chunk of S partitions its chunk using a hash function  $h(t.B)$
- Step 2:
  - Each server computes the join of its local fragment of R with its local fragment of S

# Optimization for Small Relations

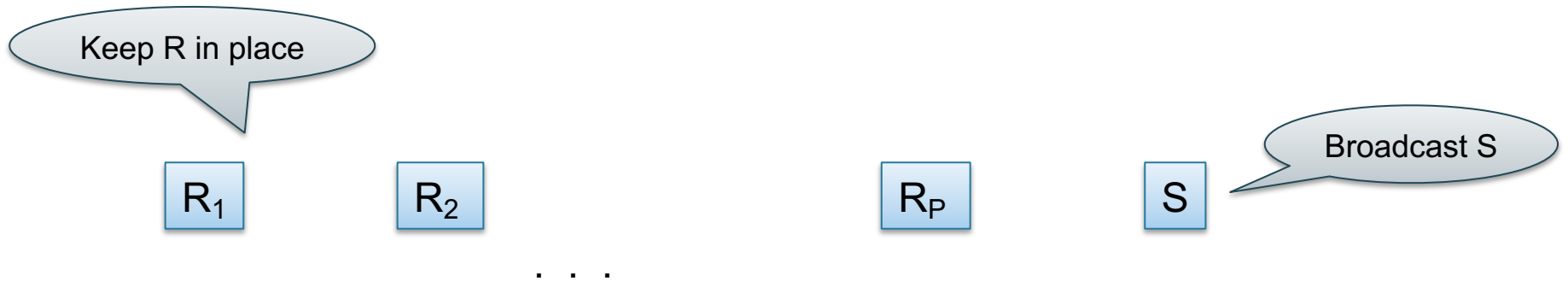
- When joining R and S
- If  $|R| \gg |S|$ 
  - Leave R where it is
  - Replicate entire S relation across nodes
- Also called a **small join** or a **broadcast join**

Query:  $R \bowtie S$

# Broadcast Join

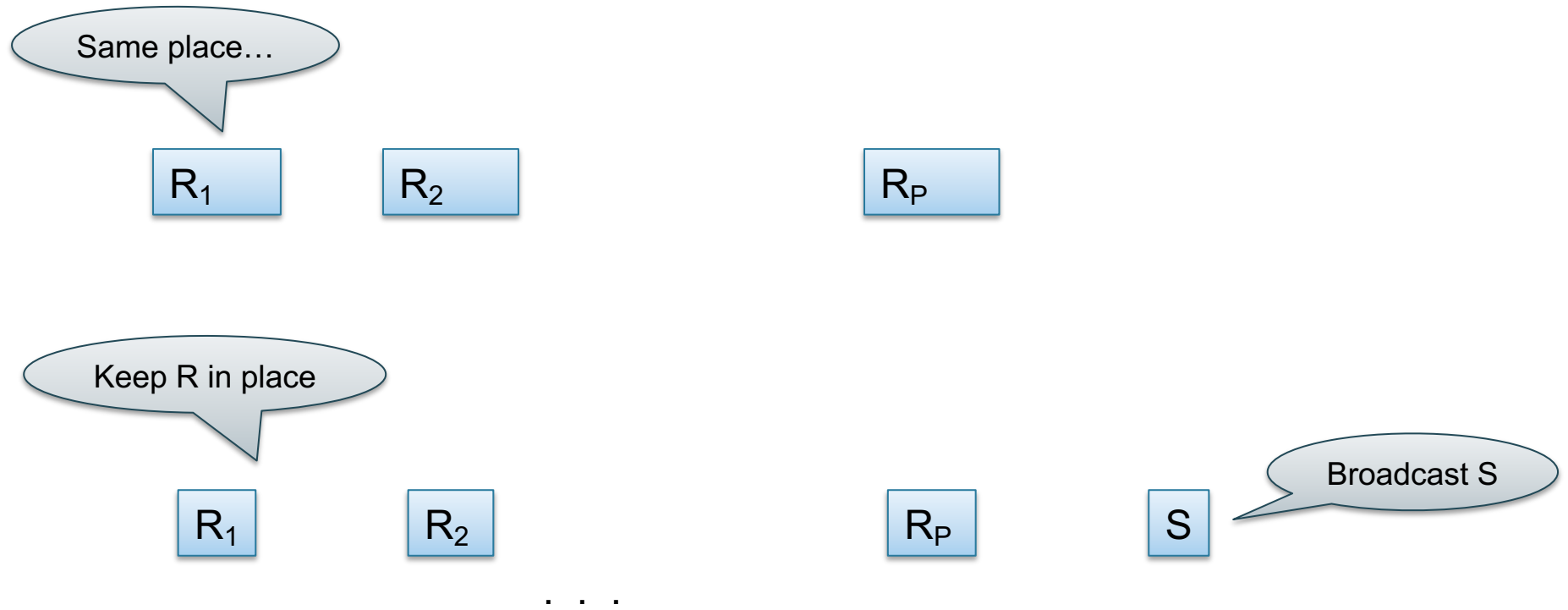


# Broadcast Join



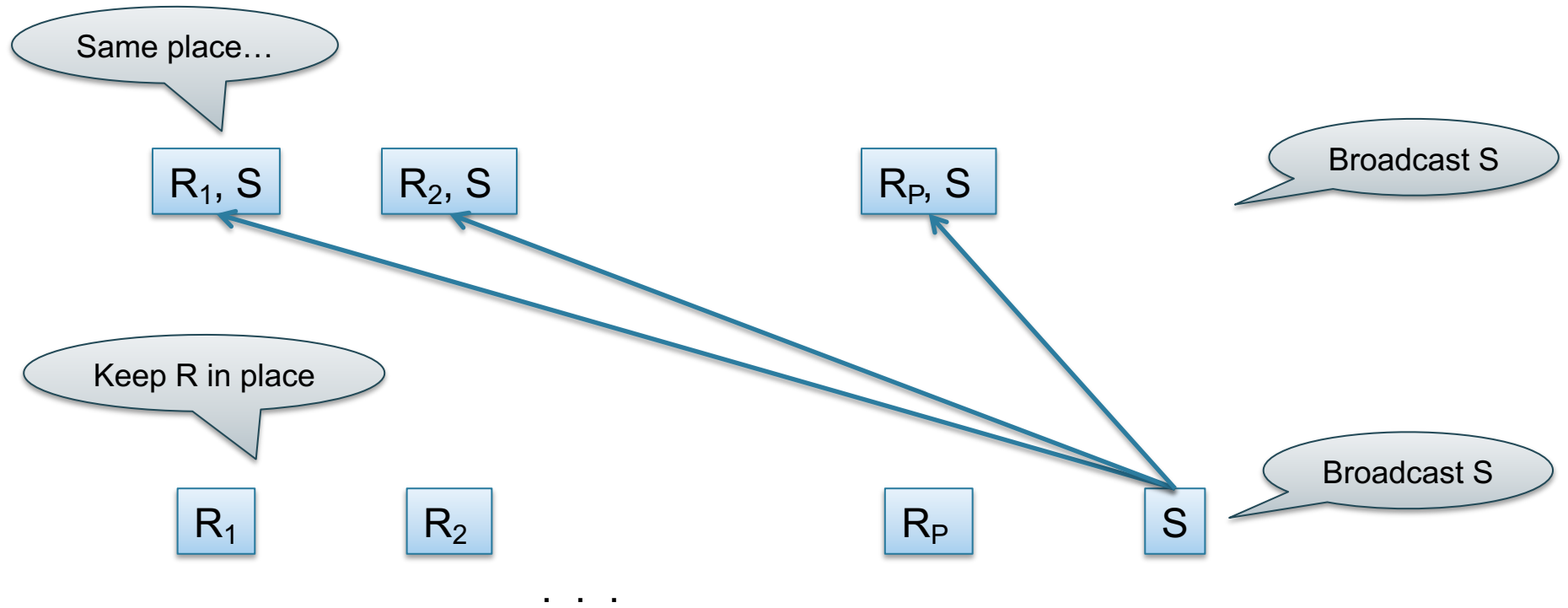
Query:  $R \bowtie S$

# Broadcast Join



Query:  $R \bowtie S$

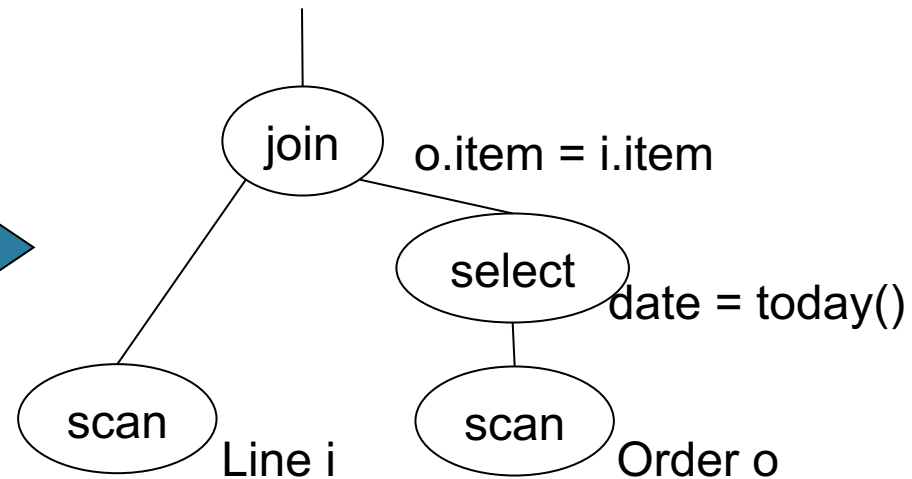
# Broadcast Join

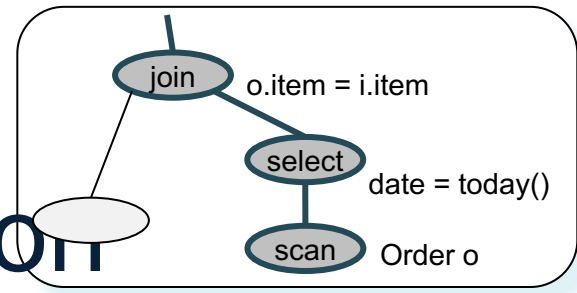


# Example Query Execution

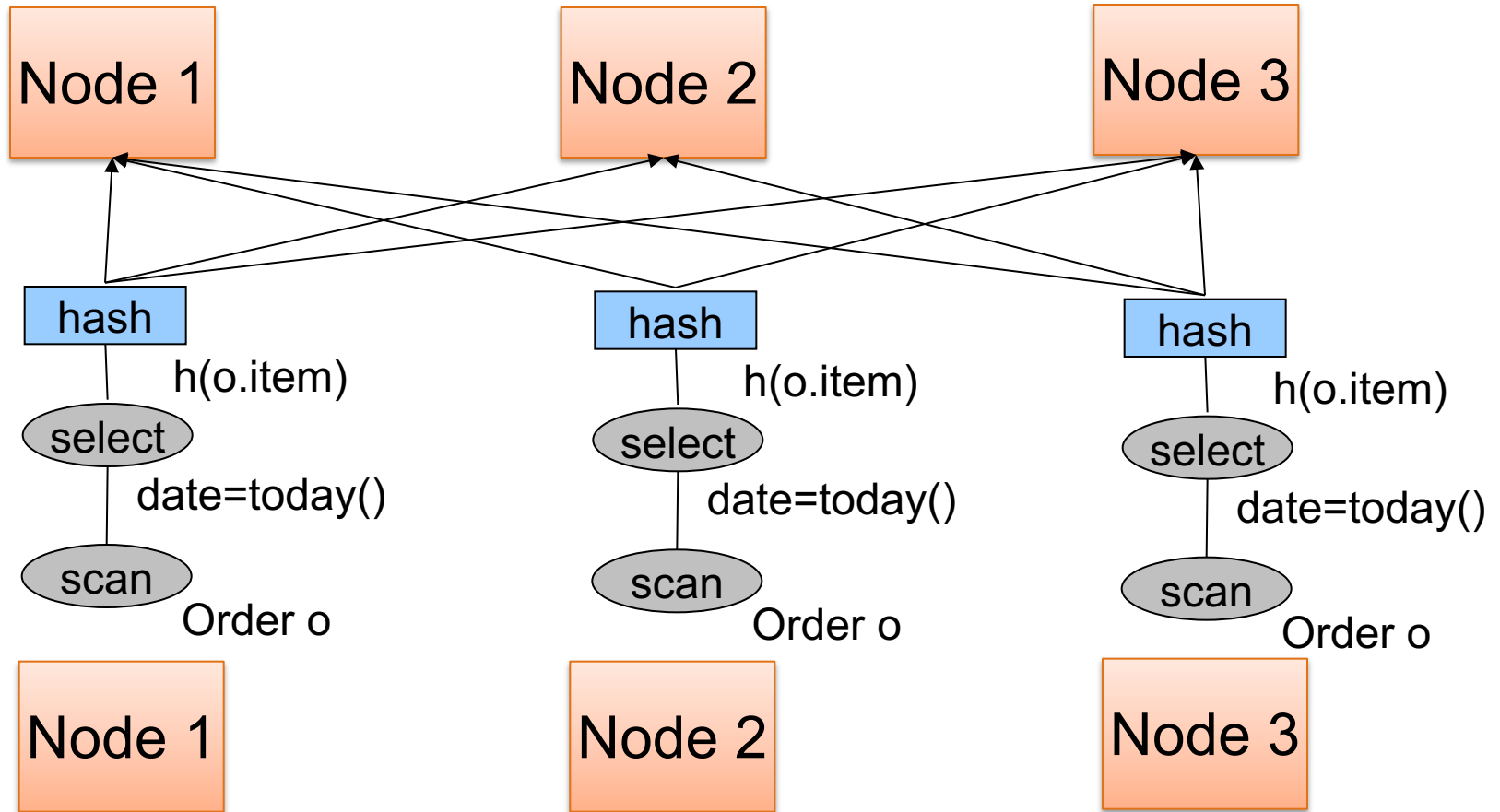
*Find all orders from today, along with the items ordered*

```
SELECT *  
FROM Order o, Line i  
WHERE o.item = i.item  
      AND o.date = today()
```



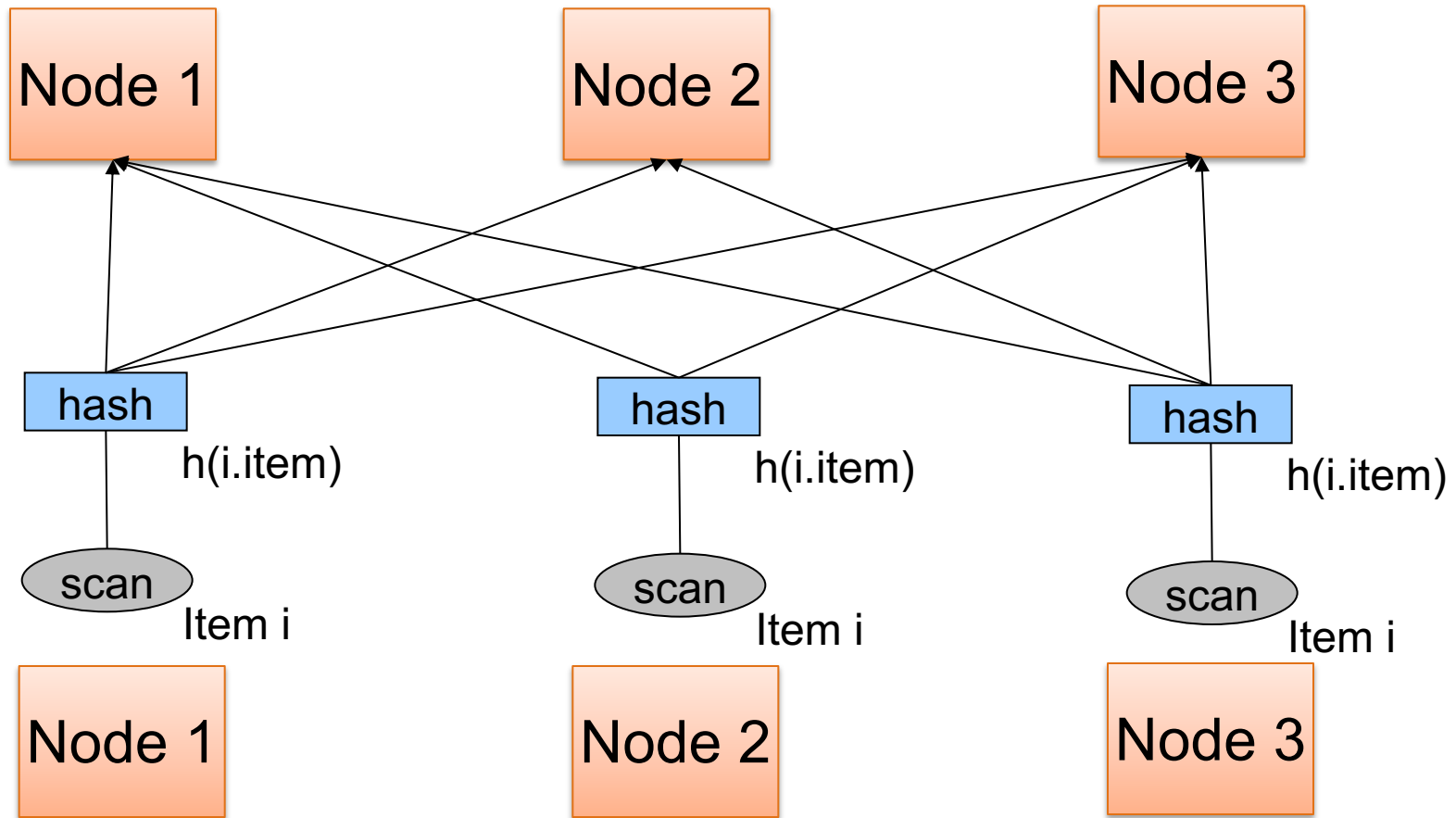
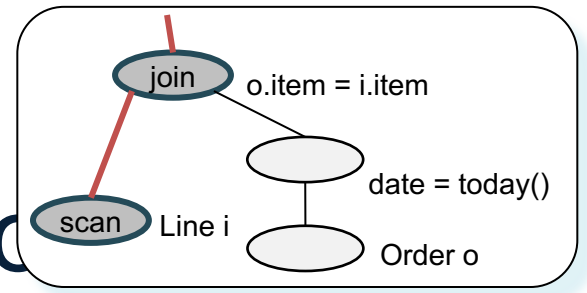


# Query Execution

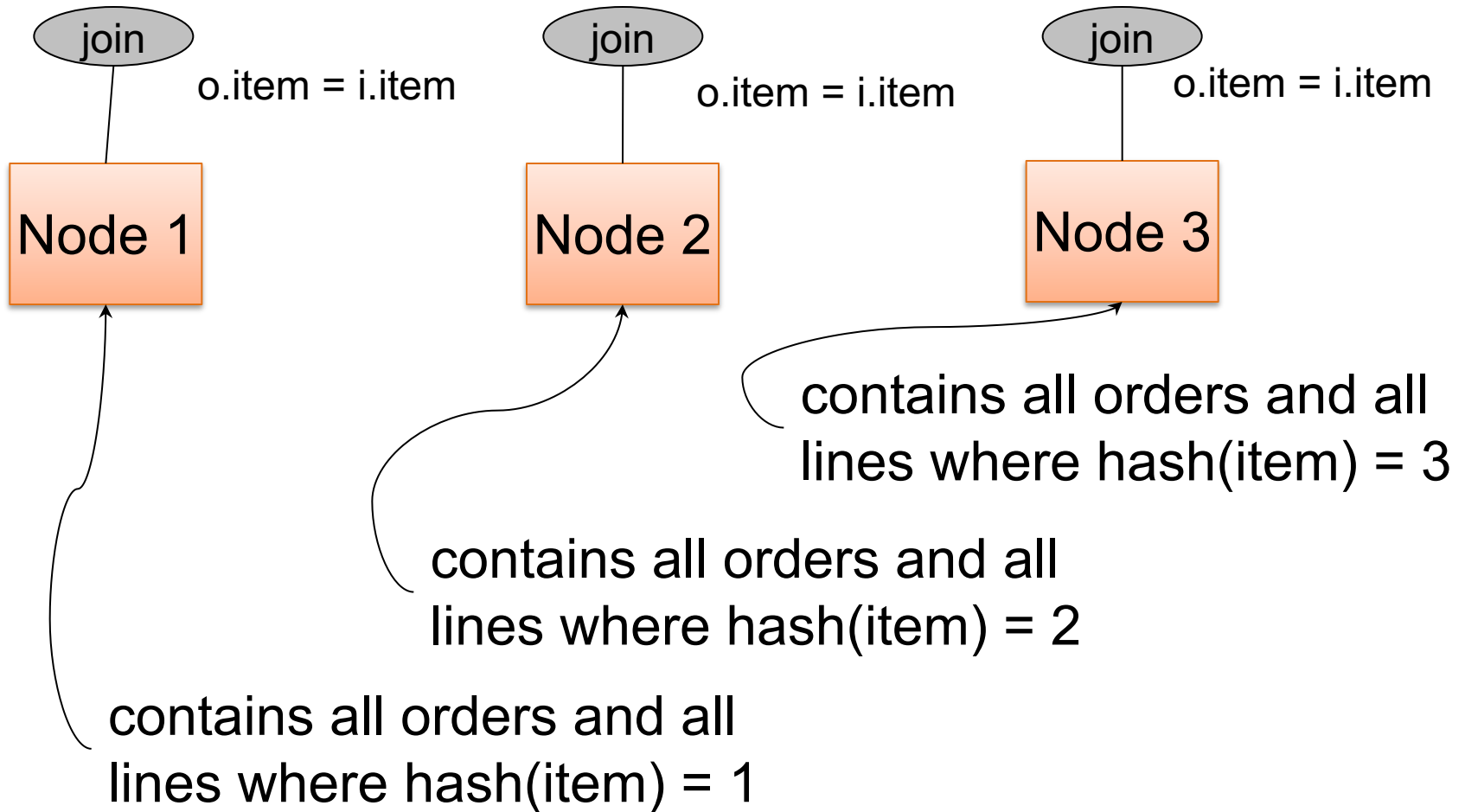




# Query Execution



# Query Execution



## Example 2

```
SELECT *  
FROM R, S, T  
WHERE R.b = S.c AND S.d = T.e AND (R.a - T.f) > 100
```

Machine 1

1/3 of R, S, T

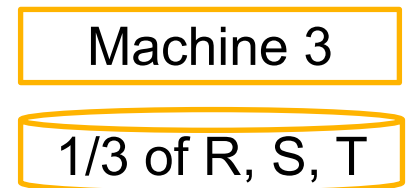
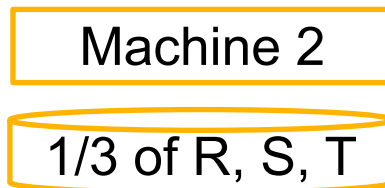
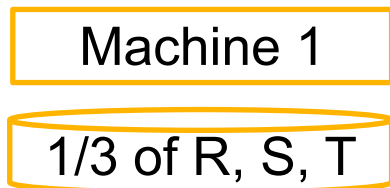
Machine 2

1/3 of R, S, T

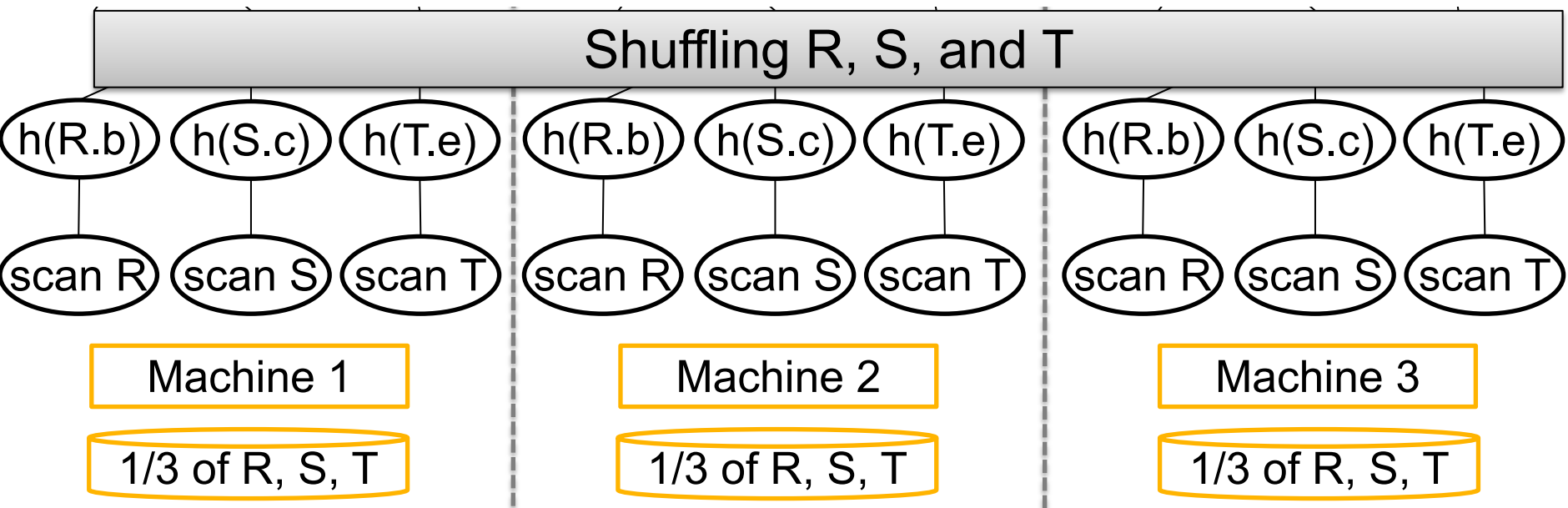
Machine 3

1/3 of R, <sup>51</sup>S, T

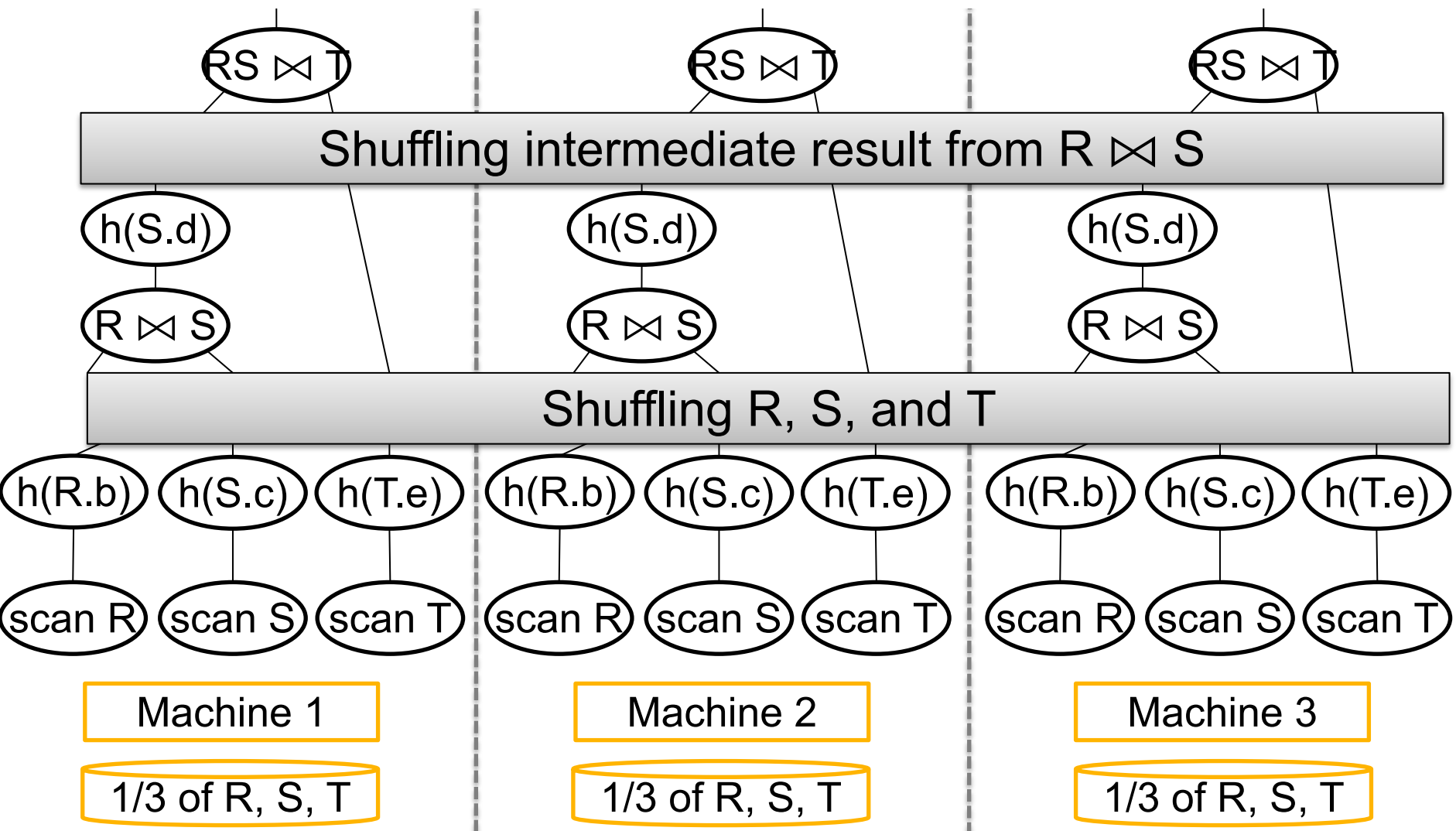
... WHERE  $R.b = S.c$  AND  $S.d = T.e$  AND  $(R.a - T.f) > 100$



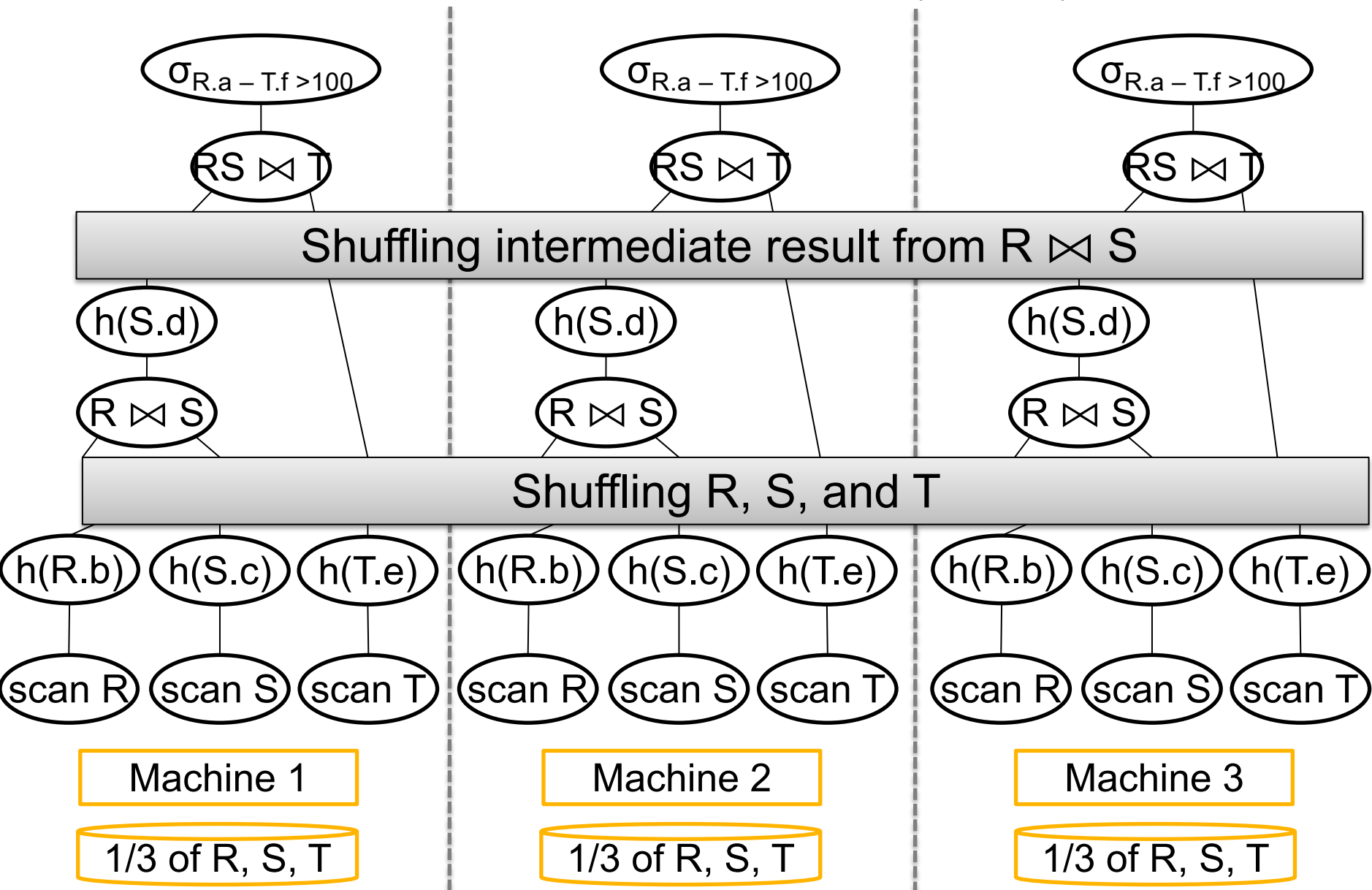
... WHERE  $R.b = S.c$  AND  $S.d = T.e$  AND  $(R.a - T.f) > 100$



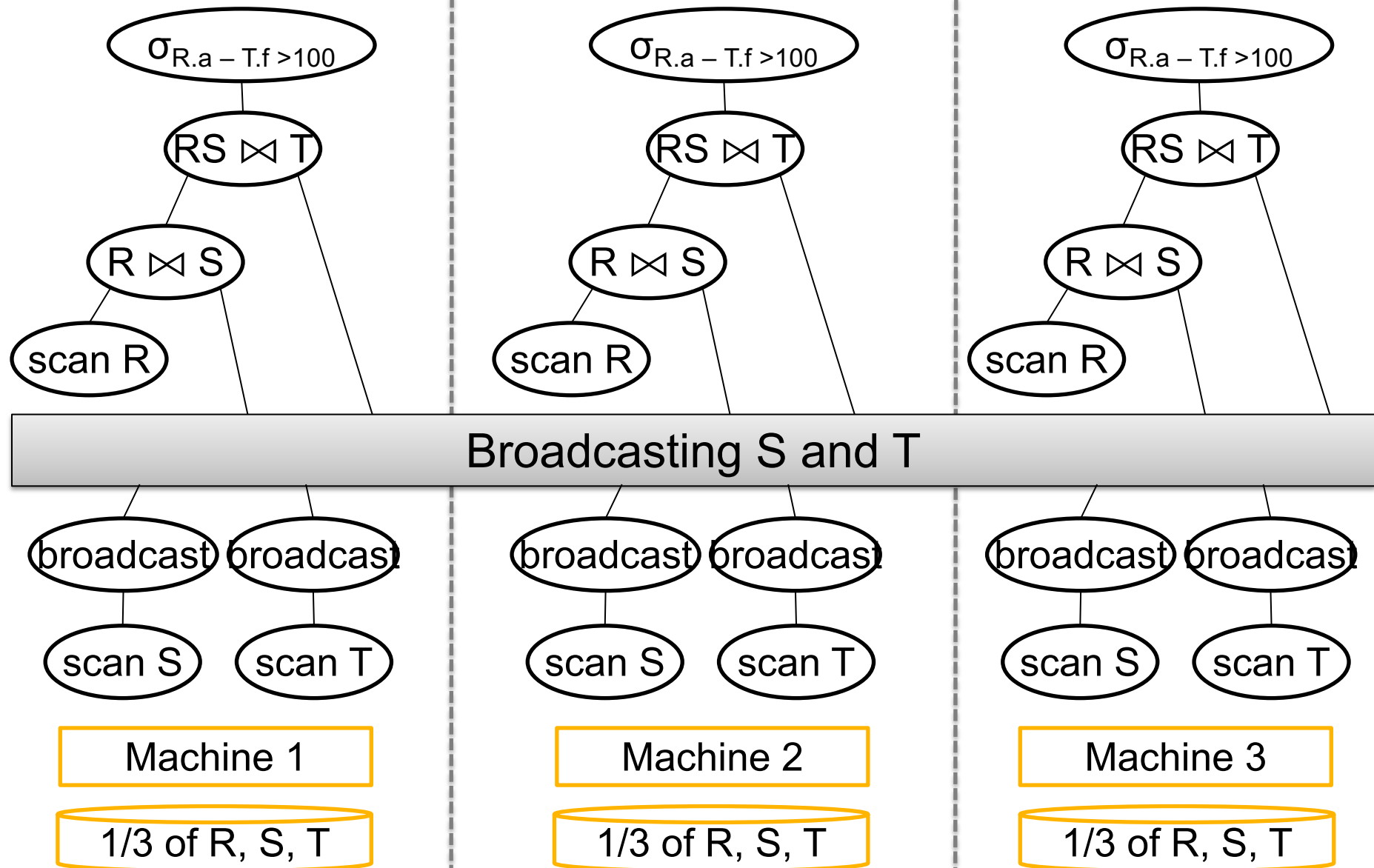
... WHERE R.b = S.c AND S.d = T.e AND (R.a - T.f) > 100



... WHERE  $R.b = S.c$  AND  $S.d = T.e$  AND  $(R.a - T.f) > 100$



... WHERE  $R.b = S.c$  AND  $S.d = T.e$  AND  $(R.a - T.f) > 100$





# Skew

# Skew

- Skew in the input: a data value has much higher frequency than others
- Skew in the output: a server generates many more values than others, e.g. join
- Skew in the computation

# Simple Skew Handling Techniques

For range partition:

- Ensure each range gets same number of tuples
- E.g.:  $\{1, 1, 1, 2, 3, 4, 5, 6\} \rightarrow [1,2]$  and  $[3,6]$
- Eq-depth v.s. eq-width histograms

# Simple Skew Handling Techniques

Skew in the computation:

- Create more partitions than nodes
  - “virtual servers”
- And be smart about scheduling the partitions
- Note: MapReduce uses this technique

# Skew for Hash Partition

Relation  $R(A,B,C,\dots)$ , we hash-partition on  $A$   
If  $A$  is a key: we expect a uniform partition

# Skew for Hash Partition

Relation  $R(A,B,C,\dots)$ , we hash-partition on  $A$

If  $A$  is a key: we expect a uniform partition

If  $A$  is not a key:

- Some value  $A=v$  may occur very many times
  - The “Justin Bieber” effect 😊
  - $v$  is called a “heavy hitter”

# Skew for Hash Partition

Relation  $R(A,B,C,\dots)$ , we hash-partition on  $A$

If  $A$  is a key: we expect a uniform partition

If  $A$  is not a key:

- Some value  $A=v$  may occur very many times
  - The “Justin Bieber” effect 😊
  - $v$  is called a “heavy hitter”
- All records with same value  $v$  are hashed to the same server  $i$
- Partition  $R_i$  is much larger than  $|R|/p$ ; skew!!

# Discussion

Distributed joins: usually hash- or broadcast-join

Heavy hitter values will significantly degrade performance of a hash-join

- **Observation 1:** there are “few” heavy hitter values (why?)
- **Observation 2:** we can compute the heavy hitter values rather easily (how?)

Rest of the lecture: How many times can  $v$  occur before it is a heavy hitter?



# Analyzing Heavy Hitters

- We will discuss how to choose the threshold such that a value that occurs more times than the threshold becomes a “heavy hitters”
- This analysis is based on Chernoff bounds, which is a general technique that is useful in statistics and randomized algorithm

# Problem Statement

Given:  $N$  data items  $v_1, \dots, v_N$

- We hash-partition them to  $P$  nodes
- When is the partitioning uniform?

# Problem Statement

Given:  $N$  data items  $v_1, \dots, v_N$

- We hash-partition them to  $P$  nodes
- When is the partitioning uniform?

**Uniform:** each node has  $O(N/P)$  items

# Problem Statement

Given:  $N$  data items  $v_1, \dots, v_N$

- We hash-partition them to  $P$  nodes
- When is the partitioning uniform?

**Uniform:** each node has  $O(N/P)$  items

**Skew:** some node has  $\gg N/P$  items

# Problem Statement

Given:  $N$  data items  $v_1, \dots, v_N$

- We hash-partition them to  $P$  nodes
- When is the partitioning uniform?

**Uniform:** each node has  $O(N/P)$  items

**Skew:** some node has  $\gg N/P$  items

1. Due to the hash function  $h$ , or
2. Due to skew in the data

# Role of the Hash Function

Assume  $v_1, \dots, v_N$  are distinct

Hash function computes  $h(v_i) \in \{1, \dots, P\}$

- If  $h$  is fixed then we can find bad items that will overload one server; **how?**
- If  $h$  is random: balls-in-bins problem; we analyze it using the Chernoff bound

Note:  
very many  
variants

# The Chernoff Bound

Bernoulli r.v.:  $X_1, \dots, X_N \in \{0,1\}$

For all  $i$ ,  $\Pr(X_i = 1) = \mu \in (0,1)$

We are interested in  $Y = X_1 + X_2 + \dots + X_N$

**Fact:**  $E[Y] = N\mu$

**Theorem** (Chernoff bound)

$$\Pr(Y > (1 + \delta)E[Y]) \leq \exp\left(-\frac{\delta^2}{3}E[Y]\right)$$

# Role of the Hash Function

Fix one server  $j$ ;

Define indicator variables:

$$X_1 = [h(v_1) = j], \dots, X_N = [h(v_N) = j]$$

$$\Pr(X_1 = 1) = \dots = \Pr(X_N = 1) = 1/P$$

**Load of server  $j$ :**  $\text{Load}(j) = X_1 + X_2 + \dots + X_N$

**Expected load:**  $E[\text{Load}(j)] = N/P$



# Role of the Hash Function

Load of server  $j$ :  $\text{Load}(j) = X_1 + X_2 + \cdots + X_N$

Expected load:  $E[\text{Load}(j)] = \frac{N}{P}$

# Role of the Hash Function

Load of server  $j$ :  $\text{Load}(j) = X_1 + X_2 + \dots + X_N$

Expected load:  $E[\text{Load}(j)] = \frac{N}{P}$



Why?

**Case 1:**  $v_1, \dots, v_N$  distinct; then  $X_1, \dots, X_N$  are iid.

# Role of the Hash Function

Load of server  $j$ :  $\text{Load}(j) = X_1 + X_2 + \dots + X_N$

Expected load:  $E[\text{Load}(j)] = \frac{N}{P}$

Why?

**Case 1:**  $v_1, \dots, v_N$  distinct; then  $X_1, \dots, X_N$  are iid.

Skew at  $j$

Cernoff:  $\Pr\left(\text{Load}(j) > (1 + \delta) \frac{N}{P}\right) \leq \exp\left(-\frac{\delta^2 N}{3P}\right)$

# Role of the Hash Function

Load of server  $j$ :  $\text{Load}(j) = X_1 + X_2 + \dots + X_N$

Expected load:  $E[\text{Load}(j)] = \frac{N}{P}$

Why?

**Case 1:**  $v_1, \dots, v_N$  distinct; then  $X_1, \dots, X_N$  are iid.

Skew at  $j$

Cernoff:  $\Pr\left(\text{Load}(j) > (1 + \delta) \frac{N}{P}\right) \leq \exp\left(-\frac{\delta^2 N}{3P}\right)$

Union bound:  $\Pr(\text{Skew}) \leq P \cdot \exp\left(-\frac{\delta^2 N}{3P}\right)$

Skew at 1 or at 2 ... or at  $P$

# Role of the Hash Function

**Case 1:**  $v_1, \dots, v_N$  distinct:

$$\Pr(\text{Skew}) \leq P \cdot \exp\left(-\frac{\delta^2 N}{3 P}\right)$$

Discussion: usually  $N \gg P$

# Role of the Hash Function

**Case 1:**  $v_1, \dots, v_N$  distinct:

$$\Pr(\text{Skew}) \leq P \cdot \exp\left(-\frac{\delta^2 N}{3P}\right)$$

Discussion: usually  $N \gg P$

- E.g. want load/server  $< 30\%$  above expected, then  $\delta = 0.3$  Assume  $N=10^9$  and  $P=10^3$

# Role of the Hash Function

**Case 1:**  $v_1, \dots, v_N$  distinct:

$$\Pr(\text{Skew}) \leq P \cdot \exp\left(-\frac{\delta^2 N}{3P}\right)$$

Discussion: usually  $N \gg P$

- E.g. want load/server  $< 30\%$  above expected, then  $\delta = 0.3$  Assume  $N=10^9$  and  $P=10^3$

$$\Pr(\text{Skew}) \leq 1000 \cdot e^{-\frac{0.09}{3}10^6} = 1000 \cdot e^{-3 \cdot 10^4} \approx 0$$

# Role of the Hash Function

**Case 1:**  $v_1, \dots, v_N$  distinct:

$$\Pr(\text{Skew}) \leq P \cdot \exp\left(-\frac{\delta^2 N}{3P}\right)$$

Discussion: usually  $N \gg P$

- Start worrying only when  $N \approx P \ln P$  (why?)



# Role of the Hash Function

- Don't write your own has function!
- Randomize it (how?)
- Make sure  $N \gg P$  (if not, why parallelize?)
- Then Load =  $O(N/P)$

Take away: a good hash function shall not cause skew!

# Role of the Data Skew

**Case 2:**  $v_1, \dots, v_N$  have duplicates

Call  $v_i$  a heavy hitter if it occurs  $\gg N/P$  times

# Role of the Data Skew

**Case 2:**  $v_1, \dots, v_N$  have duplicates

Call  $v_i$  a heavy hitter if it occurs  $\gg N/P$  times

**Fact** if there exists a heavy hitter, then there exists a server  $j$  s.t.  $\text{Load}(j) \gg \frac{N}{P}$

# Role of the Data Skew

**Case 2:**  $v_1, \dots, v_N$  have duplicates

Call  $v_i$  a heavy hitter if it occurs  $\gg N/P$  times

**Fact** if there exists a heavy hitter, then there exists a server  $j$  s.t.  $\text{Load}(j) \gg \frac{N}{P}$

Therefore:  $\Pr(\text{Skew})=1$

# Role of the Data Skew

**Case 2:**  $v_1, \dots, v_N$  have duplicates

Call  $v_i$  a heavy hitter if it occurs  $\gg N/P$  times

**Fact** if there exists a heavy hitter, then there exists a server  $j$  s.t.  $\text{Load}(j) \gg \frac{N}{P}$

Therefore:  $\Pr(\text{Skew})=1$

No hash function can handle heavy hitters

# Role of the Data Skew

**Case 3:**  $v_1, \dots, v_N$  have duplicates, no heavy hitters

Assume each value occurs  $\frac{N}{cP}$  times, for  $c > 1$

$v_1, v_1, \dots, v_1, v_2, v_2, \dots, v_2, \dots$

$\underbrace{\hspace{10em}}_{\frac{N}{cP}} \quad \underbrace{\hspace{10em}}_{\frac{N}{cP}}$

$cP$  distinct values

# Role of the Data Skew

**Case 3:**  $v_1, \dots, v_N$  have duplicates, no heavy hitters

Assume each value occurs  $\frac{N}{cP}$  times, for  $c > 1$

$$\underbrace{v_1, v_1, \dots, v_1}_{\frac{N}{cP}}, \underbrace{v_2, v_2, \dots, v_2}_{\frac{N}{cP}}, \dots$$

$$X_1 = [h(v_1) = j], X_2 = [h(v_2) = j], \dots$$

$cP$  distinct values

# Role of the Data Skew

**Case 3:**  $v_1, \dots, v_N$  have duplicates, no heavy hitters

Assume each value occurs  $\frac{N}{cP}$  times, for  $c > 1$

$$\underbrace{v_1, v_1, \dots, v_1}_{\frac{N}{cP}}, \underbrace{v_2, v_2, \dots, v_2}_{\frac{N}{cP}}, \dots$$

$cP$  distinct values

$$X_1 = [h(v_1) = j], X_2 = [h(v_2) = j], \dots$$

$$Y = \sum_i X_i \quad E[Y] = c \quad Load(j) = Y \frac{N}{cP}$$

$$\Pr(\text{Skew}) \leq P \cdot \Pr(Y > (1 + \delta)E[Y])$$



# Role of the Data Skew

**Case 3:**  $v_1, \dots, v_N$  have duplicates, no heavy hitters

Assume each value occurs  $\frac{N}{cP}$  times, for  $c > 1$

$$\underbrace{v_1, v_1, \dots, v_1}_{\frac{N}{cP}}, \underbrace{v_2, v_2, \dots, v_2}_{\frac{N}{cP}}, \dots$$

$cP$  distinct values

$$X_1 = [h(v_1) = j], X_2 = [h(v_2) = j], \dots$$

$$Y = \sum_i X_i \quad E[Y] = c \quad \text{Load}(j) = Y \frac{N}{cP}$$

$$\Pr(\text{Skew}) \leq P \cdot \Pr(Y > (1 + \delta)E[Y]) \leq P \cdot \exp\left(-\frac{\delta^2 c}{3}\right)$$

# Role of the Data Skew

**Case 3:**  $v_1, \dots, v_N$  have duplicates, no heavy hitters

Assume each value occurs  $\frac{N}{cP}$  times, for  $c > 1$

$$\underbrace{v_1, v_1, \dots, v_1}_{\frac{N}{cP}}, \underbrace{v_2, v_2, \dots, v_2}_{\frac{N}{cP}}, \dots$$

$cP$  distinct values

$$X_1 = [h(v_1) = j], X_2 = [h(v_2) = j], \dots$$

$$Y = \sum_i X_i \quad E[Y] = c \quad Load(j) = Y \frac{N}{cP}$$

$$Pr(\text{Skew}) \leq P \cdot Pr(Y > (1 + \delta)E[Y]) \leq P \cdot \exp\left(-\frac{\delta^2 c}{3}\right)$$

Need  $c \gtrsim \ln P$

# Discussion

Use library hash function! Randomize!

- When each value occurs  $\leq \frac{N}{P \cdot \ln P}$  times, then  $Load \leq (1 + \delta) \frac{N}{P}$  with high probability
- When some value occurs  $\gg \frac{N}{P}$  times, the load will be skewed
- Gray area: when values occur  $\approx \frac{N}{P}$  times: it can be shown that  $Load \approx \frac{N \cdot \ln(P)}{P}$

# SkewJoin

Main idea: separate the heavy hitters from the light hitters

- Hash join the light hitters: the partition is uniform because they are light
- Broadcast join the heavy hitters: works because there are very few heavy hitters

# SkewJoin: Details

Query:  $R \bowtie_{A=B} S$ ,  $R.A = \text{foreign key}$ ,  $S.A = \text{key}$

# SkewJoin: Details

Query:  $R \bowtie_{A=B} S$ ,  $R.A$  = foreign key,  $S.A$ =key

- Step 1: find the heavy hitters in  $R.A$ 
  - I.e. find the values  $v=R.A$  that occur  $\geq \frac{N}{P}$  times
  - There are  $\leq P$  heavy hitters! Broadcast them

# SkewJoin: Details

Query:  $R \bowtie_{A=B} S$ ,  $R.A$  = foreign key,  $S.A$ =key

- Step 1: find the heavy hitters in  $R.A$ 
  - I.e. find the values  $v=R.A$  that occur  $\geq \frac{N}{P}$  times
  - There are  $\leq P$  heavy hitters! Broadcast them

- Step 2: each sever partitions locally:

$$R = R_{light} \cup R_{heavy}, S = S_{light} \cup S_{heavy}$$

Notice:  $|S_{heavy}| \leq P$  (i.e. it is small)

# SkewJoin: Details

Query:  $R \bowtie_{A=B} S$ ,  $R.A = \text{foreign key}$ ,  $S.A = \text{key}$

- Step 1: find the heavy hitters in  $R.A$ 
  - I.e. find the values  $v=R.A$  that occur  $\geq \frac{N}{P}$  times
  - There are  $\leq P$  heavy hitters! Broadcast them
- Step 2: each server partitions locally:  
$$R = R_{light} \cup R_{heavy}, S = S_{light} \cup S_{heavy}$$

Notice:  $|S_{heavy}| \leq P$  (i.e. it is small)
- Step 3: hash-join  $R_{light} \bowtie S_{light}$



# SkewJoin: Details

Query:  $R \bowtie_{A=B} S$ ,  $R.A$  = foreign key,  $S.A$ =key

- Step 1: find the heavy hitters in  $R.A$ 
  - I.e. find the values  $v=R.A$  that occur  $\geq \frac{N}{P}$  times
  - There are  $\leq P$  heavy hitters! Broadcast them
- Step 2: each server partitions locally:  
 $R = R_{light} \cup R_{heavy}$ ,  $S = S_{light} \cup S_{heavy}$   
Notice:  $|S_{heavy}| \leq P$  (i.e. it is small)
- Step 3: hash-join  $R_{light} \bowtie S_{light}$
- Step 4: broadcast join  $R_{heavy} \bowtie S_{heavy}$

# Discussion

- Many distributed query processors do not handle skew well
- (Project idea: how does your favorite engine handle skewed data?)
- In practice, you may need to partition skewed data manually