

# Introduction to Data Management

## CSE 344

### Lecture 25: MapReduce and Spark

# Today

- MapReduce Review
- Parallel Join Algorithms
- Spark



# Parallel Data Processing @ 2000



# Map Reduce Data Model

**Instance:** Files!

- where a file = a bag of (key, value) pairs

**Schema:** None!

- just like other key-value data models

**Query language:** a MapReduce program:

- Input: a bag of (inputkey, value) pairs
- Output: a bag of (outputkey, value) pairs

# Step 1: the **MAP** Phase

User provides the **map** function:

- Input: **(input key, value)**
- Output: bag of **(intermediate key, value)**

System applies map **in parallel** to all **(input key, value)** pairs in the input file

## Step 2: the REDUCE Phase

User provides the `reduce` function:

- Input: (intermediate key, bag of values w/ same key)
- Output: bag of output (values)

System groups all pairs with the same intermediate key, and passes the bag of values to reduce

# Example

- Counting the number of occurrences of each word in a large collection of documents
- Each Document
  - The **key** = document id (**did**)
  - The **value** = set of words (**word**)

```
map(String key, String value):  
  // key: document name  
  // value: document contents  
  for each word w in value:  
    EmitIntermediate(w, "1");
```

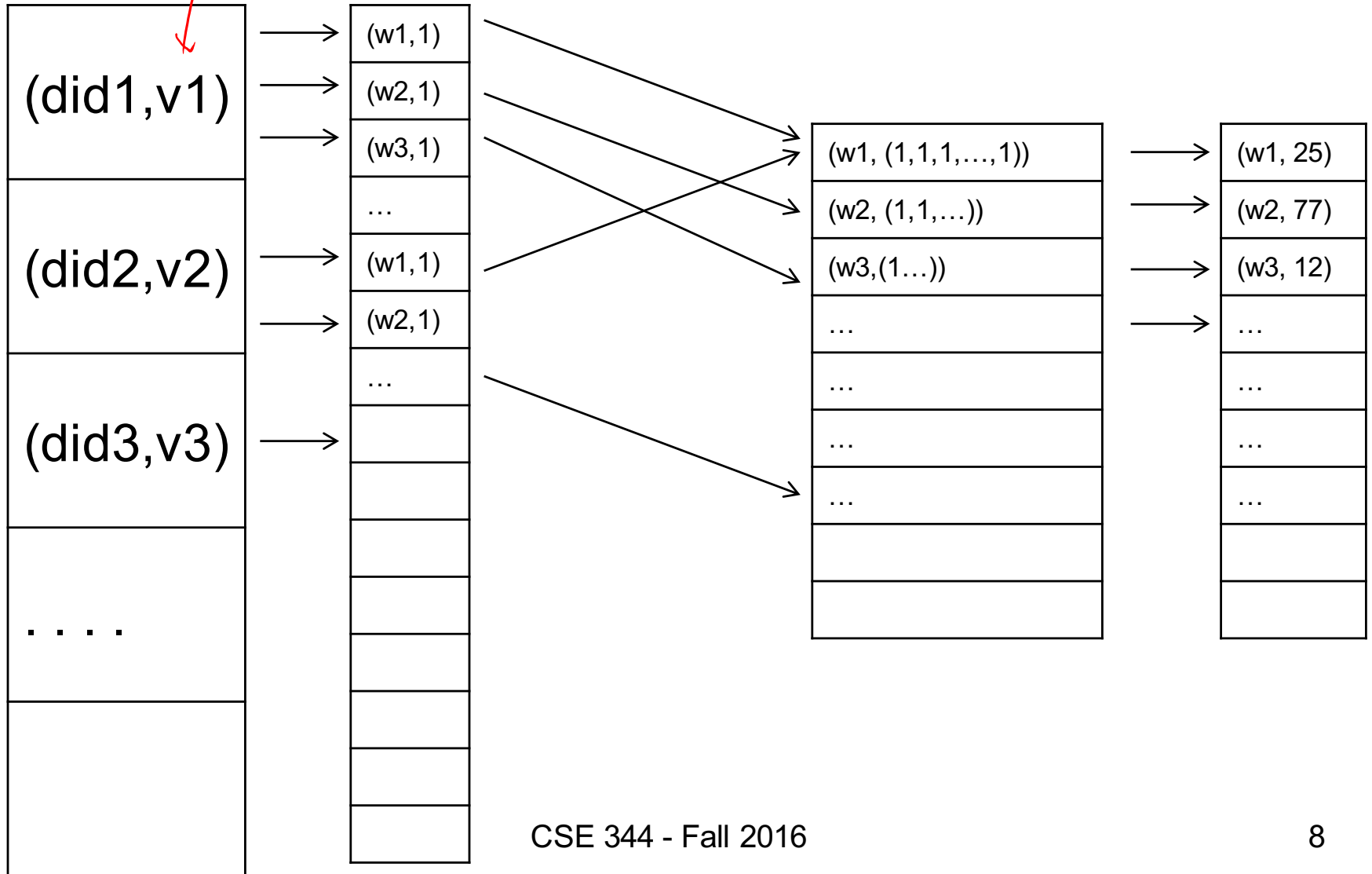
```
reduce(String key, Iterator values):  
  // key: a word  
  // values: a list of counts  
  int result = 0;  
  for each v in values:  
    result += ParseInt(v);  
  Emit(AsString(result));
```

MAP

$\{w_1, w_2, \dots, w_1, \dots\}$

Shuffle

REDUCE





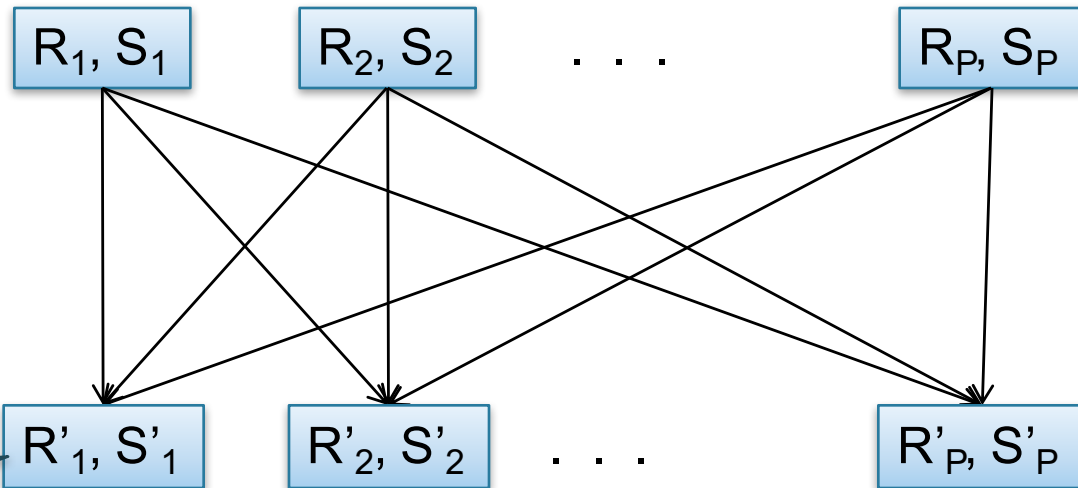
# Fault Tolerance

- If one server fails once every year...  
... then a job with 10,000 servers will fail in less than one hour
- MapReduce handles fault tolerance by writing intermediate files to disk:
  - Mappers write file to local disk
  - Reducers read the files (=reshuffling); if the server fails, the reduce task is restarted on another server

$$R(A,B) \bowtie_{B=C} S(C,D)$$

# Partitioned Hash-Join

Initially, both R and S are horizontally partitioned



Reshuffle R on R.B  
and S on S.C

Each server computes  
the join locally

$R(A,B) \bowtie_{B=C} S(C,D)$

## Partitioned Hash-Join

```
map(String relationName, String s):  
  Tuple t = parse(s);  
  switch (relationName):  
    case 'R': EmitIntermediate(t.B, KV('R', t));  
    case 'S': EmitIntermediate(t.C, KV('S', t));
```

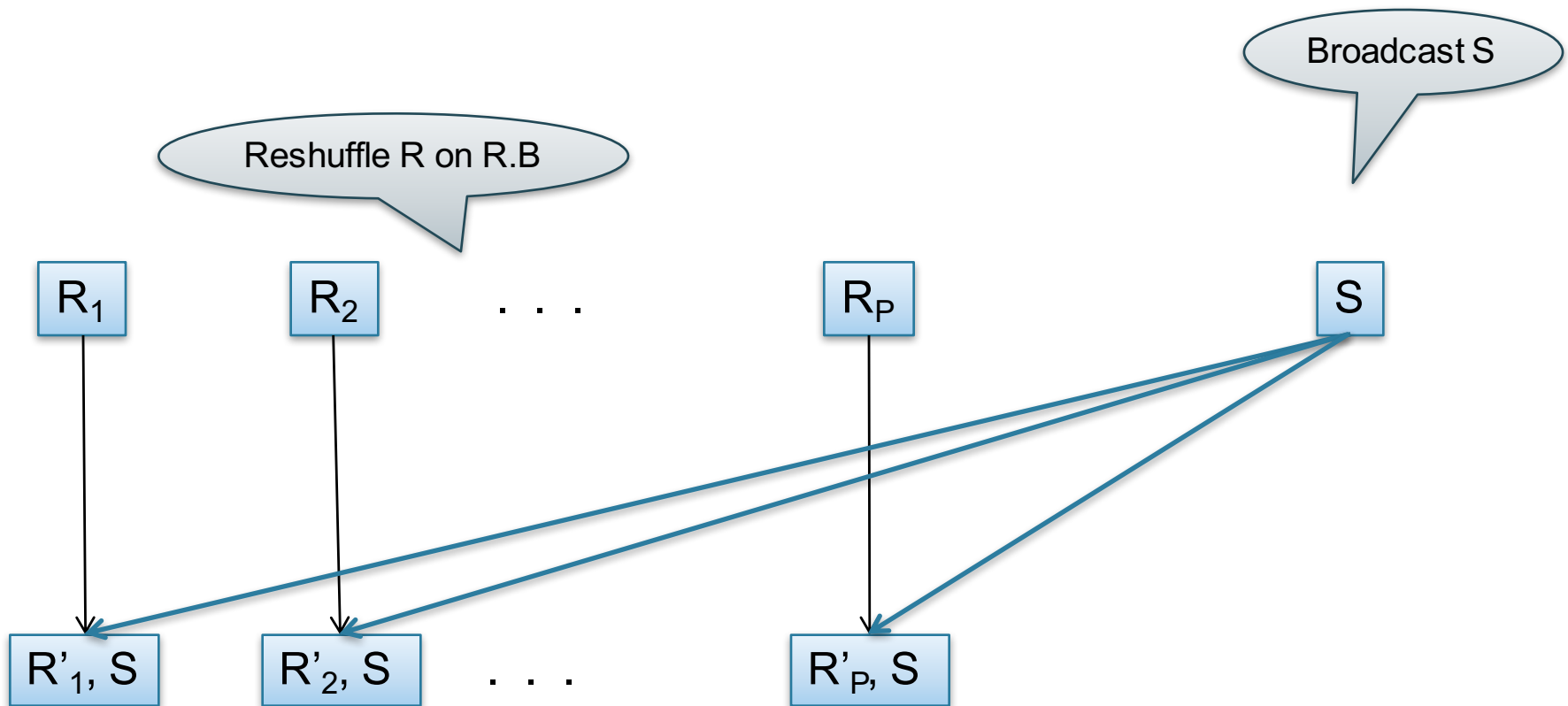
key  
Or call hash(t.B)

```
reduce(String k, Iterator values):  
  R = []; S = [];  
  for each v in values:  
    switch (v.relationName):  
      'R': R.insert(v);  
      'S': S.insert(v);  
  for r in R, for s in S  
    Emit(Tuple(r,s));
```

All tuples here must join

$$R(A,B) \bowtie_{B=C} S(C,D)$$

# Broadcast Join



$R(A,B) \bowtie_{B=C} S(C,D)$

## Broadcast Join

```
map(String key, String s):  
  Tuple [] rs = parse(s)  
  S = readFromNetwork();  
  ht = new Hashtable()  
  for each w in S:  
    ht.insert(w.C, w)  
  
  for each r in rs:  
    for each s in ht.find(r.B):  
      Emit(Tuple(r,s));
```

`map` should read  
several records of R:  
`value` = some group  
of records

Read entire table S,  
build a Hash Table

```
reduce(...):  
  /* empty: map-side only */
```

# Let's consider a 3-way join

- How do we compute this query?

$Q(x,y,z) :- R(x,y), S(y,z), T(z,x)$

- This computes all “triangles”
- e.g., let  $Follows(x,y)$  be all pairs of Twitter users s.t. x follows y.
- If  $R=S=T=Follows$ , then  $Q$  computes all triples of people that follow each other.

$Q(x,y,z) :- R(x,y),S(y,z),T(z,x)$

Shuffle join

Partition

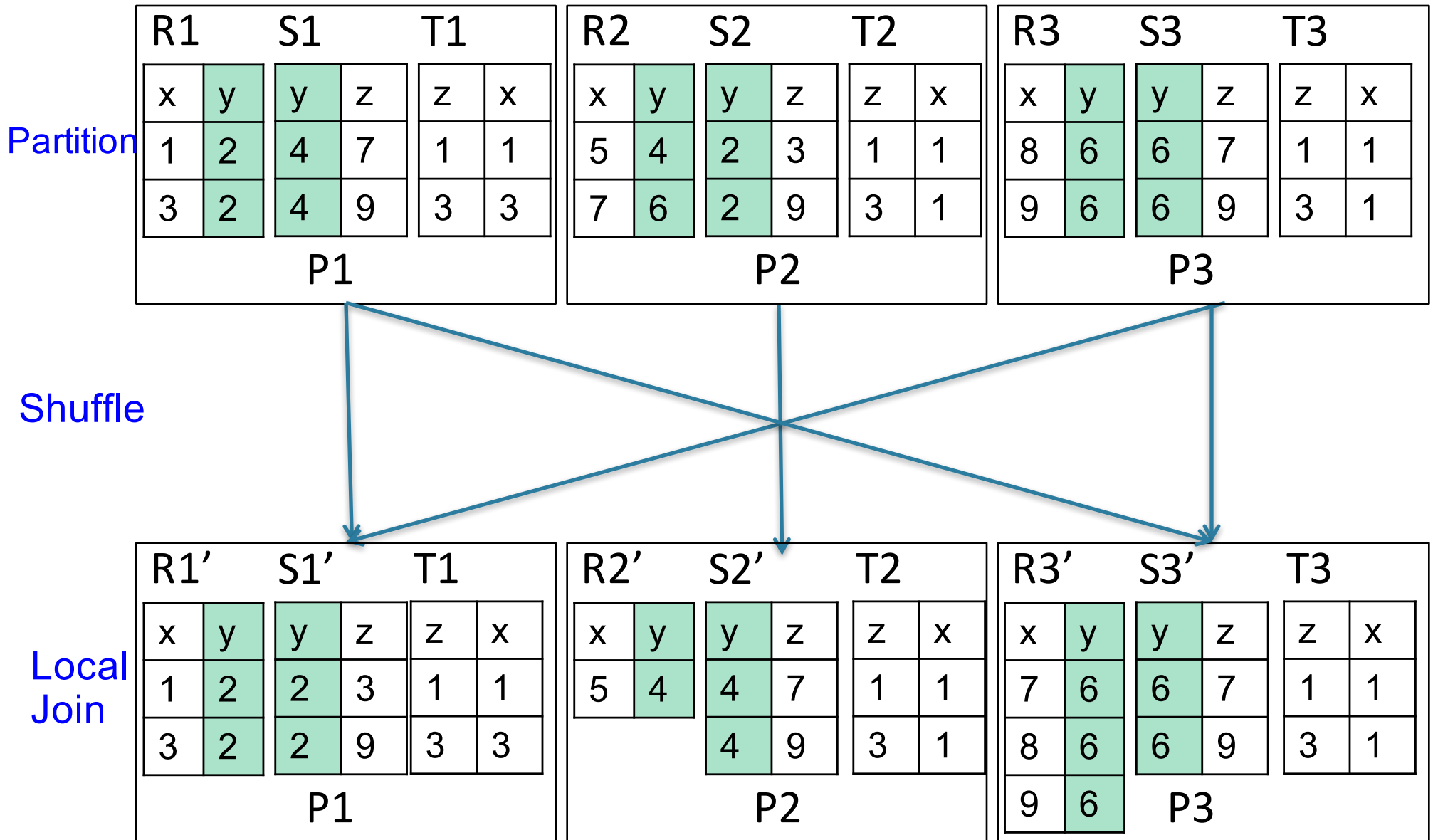
R1		S1		T1	
x	y	y	z	z	x
1	2	4	7	1	1
3	2	4	9	3	3
P1					

R2		S2		T2	
x	y	y	z	z	x
5	4	2	3	1	1
7	6	2	9	3	1
P2					

R3		S3		T3	
x	y	y	z	z	x
8	6	6	7	1	1
9	6	6	9	3	1
P3					

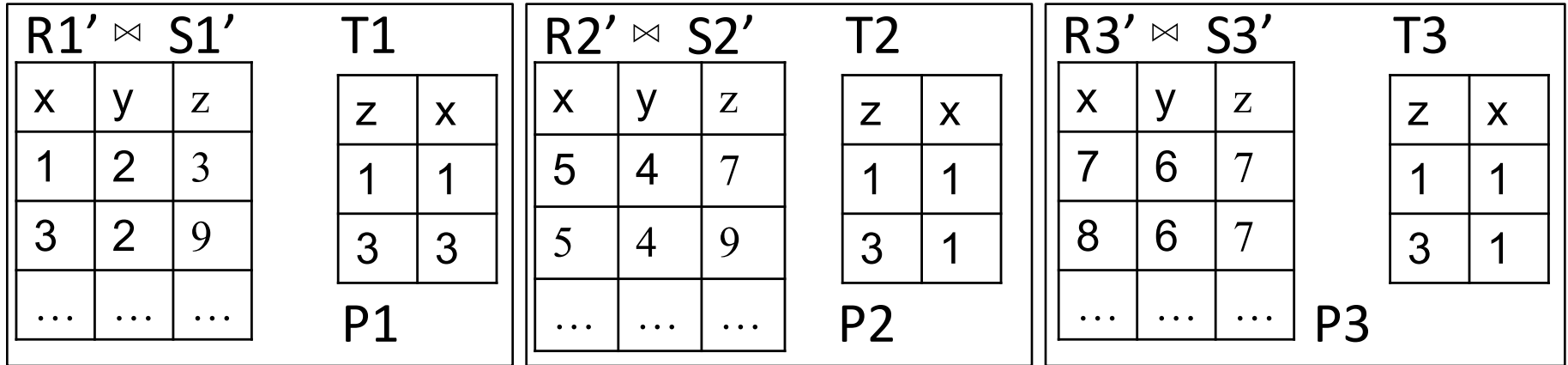
$Q(x,y,z) :- R(x,y),S(y,z),T(z,x)$

Shuffle join

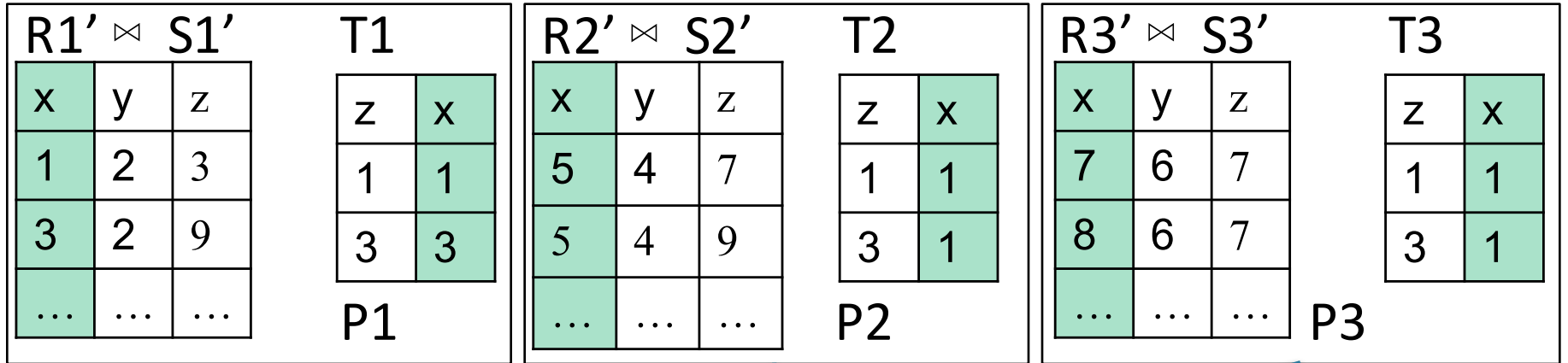




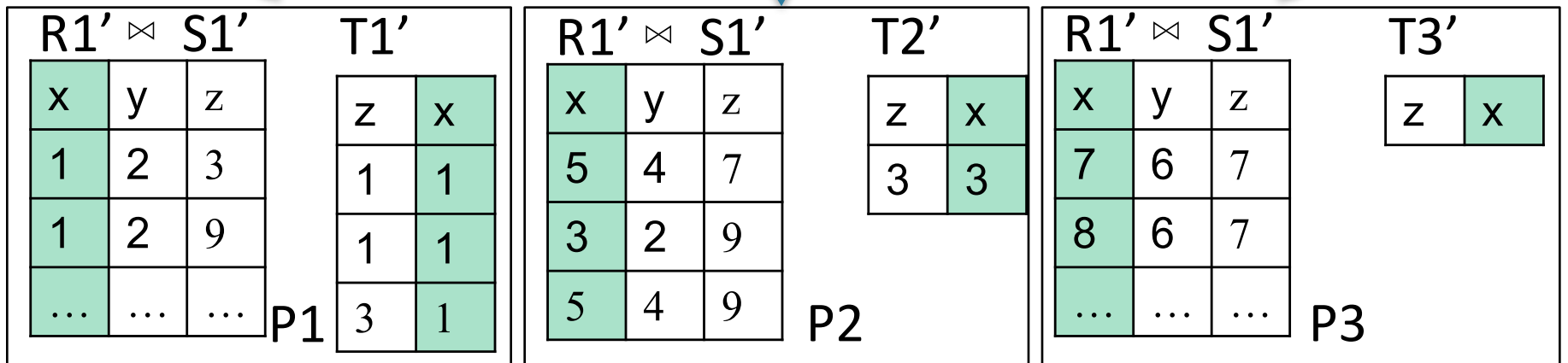
Local  
Join



Local Join



Shuffle



$Q(x,y,z) :- R(x,y), S(y,z), T(z,x)$

Broadcast join

Partition

R1		S1		T1	
x	y	y	z	z	x
1	2	4	7	1	1
3	2	4	9	3	3

P1

R2		S2		T2	
x	y	y	z	z	x
5	4	2	3	9	5
7	6	2	9	3	1

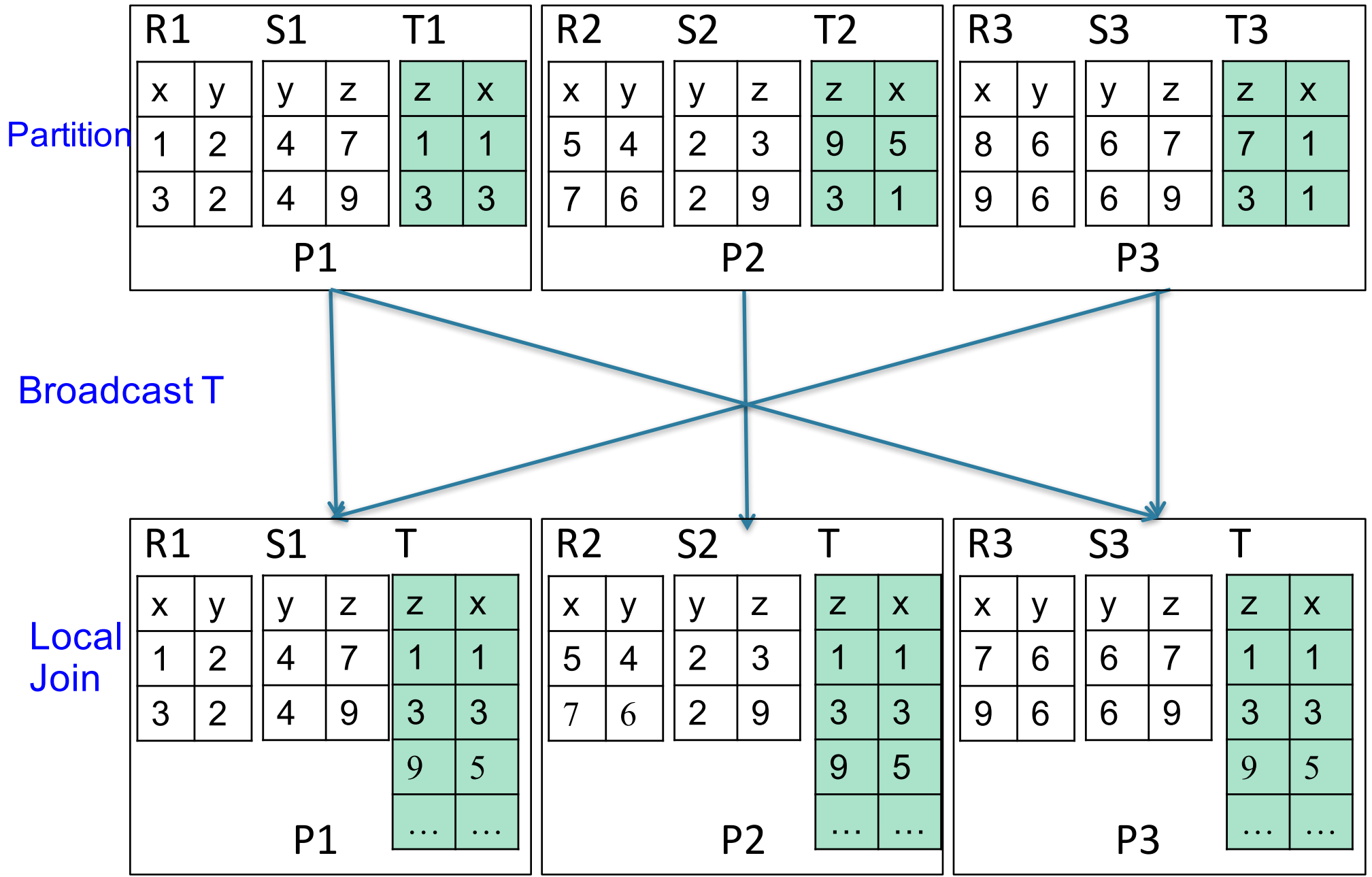
P2

R3		S3		T3	
x	y	y	z	z	x
8	6	6	7	7	1
9	6	6	9	3	1

P3

$Q(x,y,z) :- R(x,y),S(y,z),T(z,x)$

Broadcast join



$Q(x,y,z) :- R(x,y), S(y,z), T(z,x)$

Broadcast + shuffle join

Partition

R1		S1		T1	
x	y	y	z	z	x
1	2	4	7	1	1
3	2	4	9	3	3

P1

R2		S2		T2	
x	y	y	z	z	x
5	4	2	3	9	5
7	6	2	9	3	1

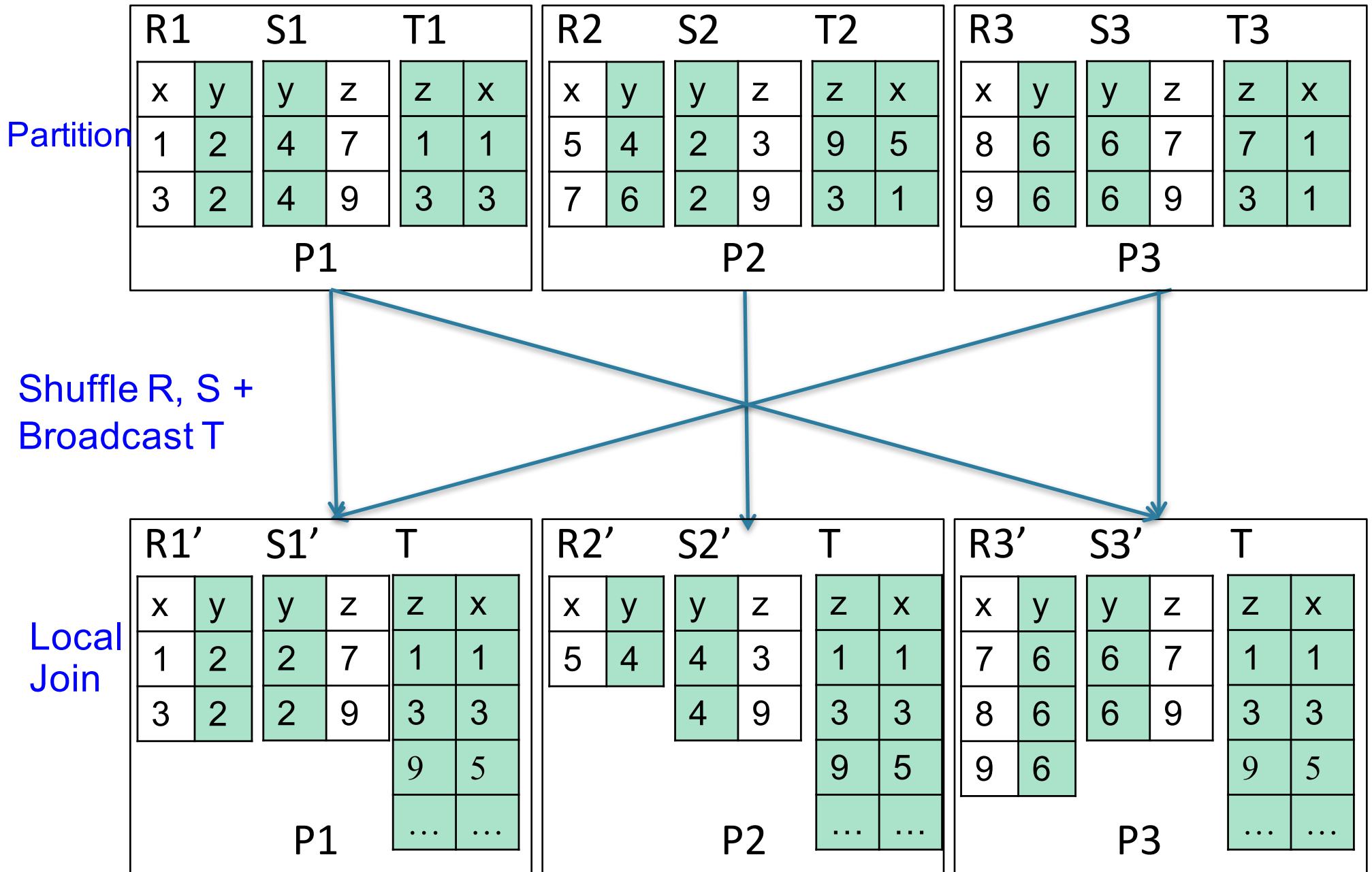
P2

R3		S3		T3	
x	y	y	z	z	x
8	6	6	7	7	1
9	6	6	9	3	1

P3

$Q(x,y,z) :- R(x,y), S(y,z), T(z,x)$

Broadcast + shuffle join



# Issues

- Shuffle join
  - Requires two shuffles
- Broadcast join
  - Requires two broadcasts
- Shuffle + broadcast join
  - Redundant broadcast
- Can we do better?

# Hypercube join

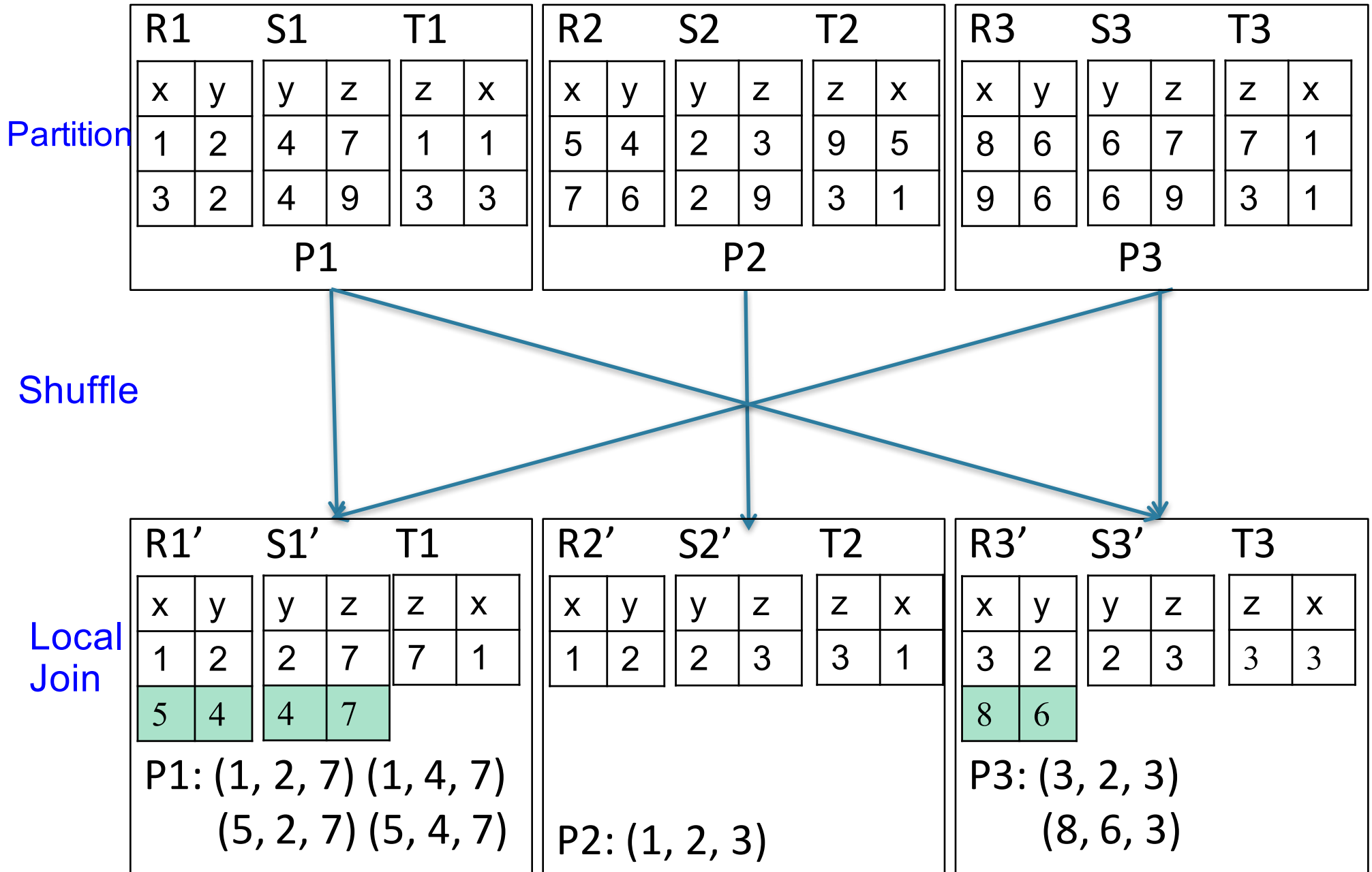
- Assign each server a specific value of  $(x,y,z)$ 
  - e.g., P1:  $(1,2,3)$  means P1 gets all tuples with  $x=1$ ,  $y=2$ , and  $z=3$
  - Join can now proceed in one step!
- What if we have more values of  $x$  than # of servers?
  - Use modulo!





$Q(x,y,z) :- R(x,y),S(y,z),T(z,x)$

Hypercube join





# Parallel Data Processing @ 2010



# Issues with MapReduce

- Difficult to write more complex queries
- Need multiple MapReduce jobs: dramatically slows down because it writes all results to disk

# Implementing Relational Operators in MapReduce

Given relations  $R(A,B)$  and  $S(B, C)$  compute:

- Selection:  $\sigma_{A=123}(R)$
- Group-by:  $\gamma_{A, \text{sum}(B)}(R)$
- Join:  $R \bowtie S$

# Selection $\sigma_{A=123}(R)$

```
map(String key, Tuple t):  
  if t.A = 123:  
    EmitIntermediate(key, t);
```

```
reduce(String k, Iterator values):  
  for each v in values:  
    Emit(v);
```

# Selection $\sigma_{A=123}(R)$

```
map(String key, Tuple t):  
  if t.A = 123:  
    EmitIntermediate(key, t);
```

```
reduce(String k, Iterator values):  
  for each v in values:  
    Emit(v);
```

No need for reduce.  
But need system hacking  
to remove reduce from MapReduce

# Group By $\gamma_{A, \text{sum}(B)}(R)$

```
map(String key, Tuple t):  
    EmitIntermediate(t.A, t.B);
```

Can't use hashtable to map  
A's to B's

```
reduce(String k, Iterator values):  
    s = 0  
    for each v in values:  
        s = s + v  
    Emit(k, v);
```



# Spark

- Open source system from Berkeley
- Distributed processing over HDFS
- Differences from MapReduce:
  - Multiple steps, including iterations
  - Stores intermediate results in main memory
  - Closer to relational algebra (familiar to you)
- Details: <http://spark.apache.org/examples.html>

# Spark Interface

- Spark supports a Scala interface
- Scala = extension of Java with functions/closures
- We will illustrate Scala/Spark in the lectures
  
- Spark also supports a SQL interface, and compiles SQL to its Scala interface
- For HW8: you only need the SQL interface!

# RDD

- RDD = Resilient Distributed Datasets
  - A distributed relation, together with its *lineage*
  - Lineage = expression that says how that relation was computed = a relational algebra plan
- Spark stores intermediate results as RDD
- If a server crashes, its RDD in main memory is lost. However, the driver (=master node) knows the lineage, and will simply recompute the lost partition of the RDD