

# **MapReduce algorithms for processing relational data**

# Design Pattern: Secondary Sorting

- MapReduce sorts input to reducers by key
  - Values are arbitrarily ordered
- What if want to sort value also?
  - E.g.,  $k \rightarrow (v_1, r), (v_3, r), (v_4, r), (v_8, r) \dots$

# Secondary Sorting: Solutions

- Solution 1:
  - Buffer values in memory, then sort
  - Why is this a bad idea?
- Solution 2:
  - “Value-to-key conversion” design pattern:  
form composite intermediate key,  $(k, v_1)$
  - Let execution framework do the sorting
  - Preserve state across multiple key-value pairs  
to handle processing
  - Anything else we need to do?

# Value-to-Key Conversion

## Before

$k \rightarrow (v_1, r), (v_4, r), (v_8, r), (v_3, r) \dots$

Values arrive in arbitrary order...

## After

$(k, v_1) \rightarrow (v_1, r)$

$(k, v_3) \rightarrow (v_3, r)$

$(k, v_4) \rightarrow (v_4, r)$

$(k, v_8) \rightarrow (v_8, r)$

...

Values arrive in sorted order...

Process by preserving state across multiple keys

Remember to partition correctly!

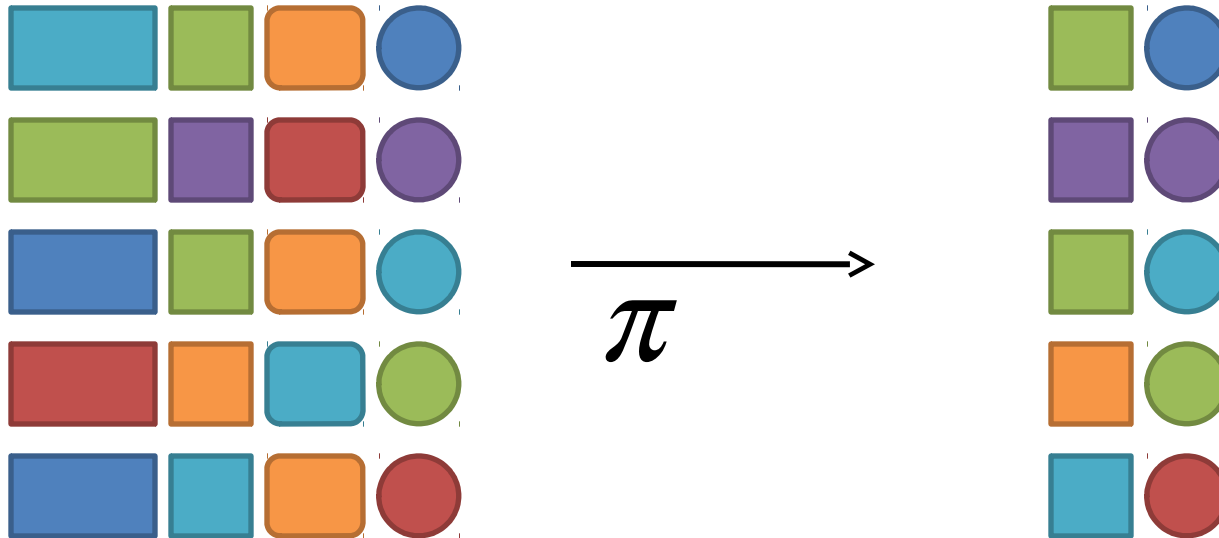
# Working Scenario

- Two tables:
  - User demographics (gender, age, income, etc.)
  - User page visits (URL, time spent, etc.)
- Analyses we might want to perform:
  - Statistics on demographic characteristics
  - Statistics on page visits
  - Statistics on page visits by URL
  - Statistics on page visits by demographic characteristic
  - ...

# Relational Algebra

- Primitives
  - Projection ( $\pi$ )
  - Selection ( $\sigma$ )
  - Cartesian product ( $\times$ )
  - Set union ( $\cup$ )
  - Set difference ( $-$ )
  - Rename ( $\rho$ )
- Other operations
  - Join ( $\bowtie$ )
  - Group by... aggregation
  - ...

# Projection

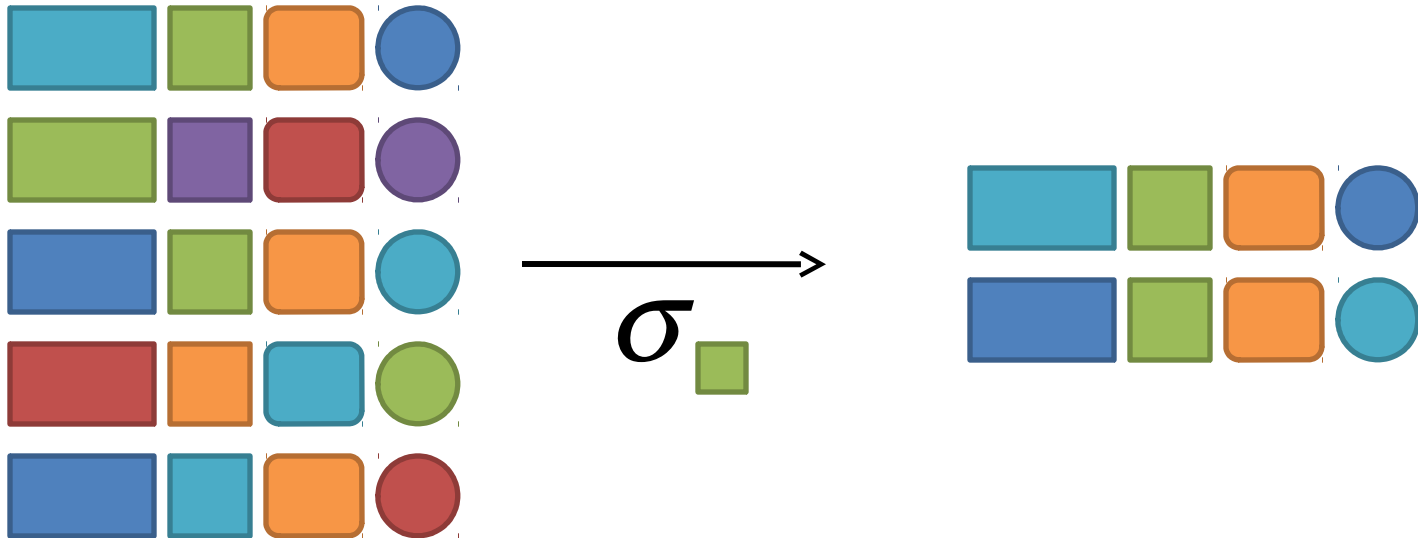


# Projection in MapReduce

- Easy!
  - Map over tuples, emit new tuples with appropriate attributes
  - No reducers, unless for regrouping or resorting tuples
  - Alternatively: perform in reducer, after some other processing
- Basically limited by HDFS streaming speeds
  - Speed of encoding/decoding tuples becomes important
  - Relational databases take advantage of compression
  - Semistructured data? No problem!



# Selection



# Selection in MapReduce

- Easy!
  - Map over tuples, emit only tuples that meet criteria
  - No reducers, unless for regrouping or resorting tuples
  - Alternatively: perform in reducer, after some other processing
- Basically limited by HDFS streaming speeds
  - Speed of encoding/decoding tuples becomes important
  - Relational databases take advantage of compression
  - Semistructured data? No problem!

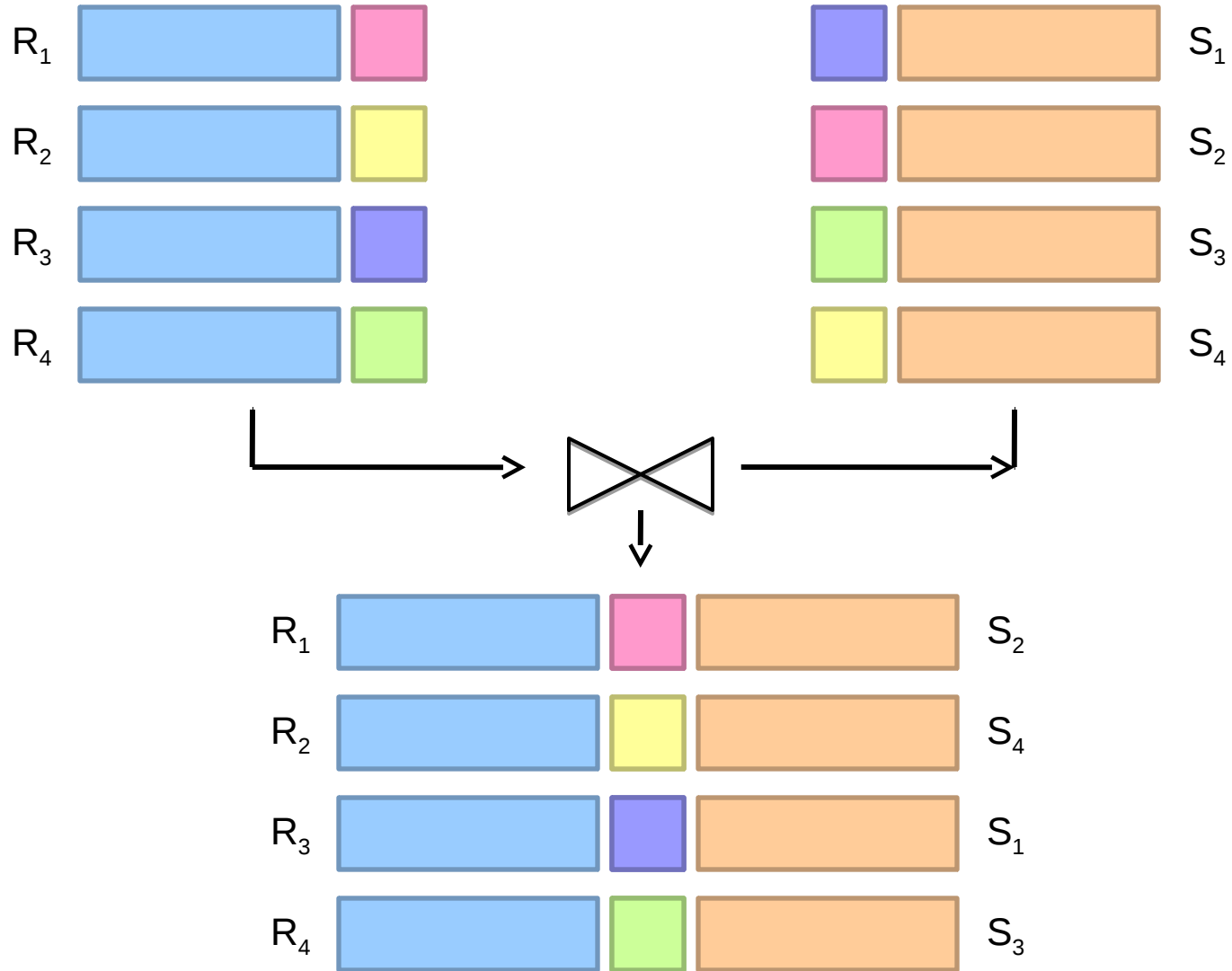
# Group by... Aggregation

- Example: What is the average time spent per URL?
- In SQL:
  - `SELECT url, AVG(time) FROM visits GROUP BY url`
- In MapReduce:
  - Map over tuples, emit time, keyed by url
  - Framework automatically groups values by keys
  - Compute average in reducer
  - Optimize with combiners

# Relational Joins



# Relational Joins



# Natural Join Operation – Example

- Relations  $r$ ,  $s$ :

$A$	$B$	$C$	$D$
$\alpha$	1	$\alpha$	a
$\beta$	2	$\gamma$	a
$\gamma$	4	$\beta$	b
$\alpha$	1	$\gamma$	a
$\delta$	2	$\beta$	b

$r$

$B$	$D$	$E$
1	a	$\alpha$
3	a	$\beta$
1	a	$\gamma$
2	b	$\delta$
3	b	$\epsilon$

$s$

$r \bowtie s$

$A$	$B$	$C$	$D$	$E$
$\alpha$	1	$\alpha$	a	$\alpha$
$\alpha$	1	$\alpha$	a	$\gamma$
$\alpha$	1	$\gamma$	a	$\alpha$
$\alpha$	1	$\gamma$	a	$\gamma$
$\delta$	2	$\beta$	b	$\delta$

# Natural Join Example

<u>sid</u>	<u>bid</u>	<u>day</u>
22	101	10/10/96
58	103	11/12/96

**R1**

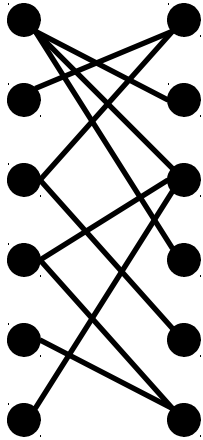
<u>sid</u>	sname	rating	age
22	dustin	7	45.0
31	lubber	8	55.5
58	rusty	10	35.0

**S1**

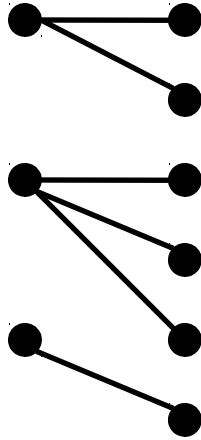
**R1** ⋈ **S1** =

sid	sname	rating	age	bid	day
22	dustin	7	45.0	101	10/10/96
58	rusty	10	35.0	103	11/12/96

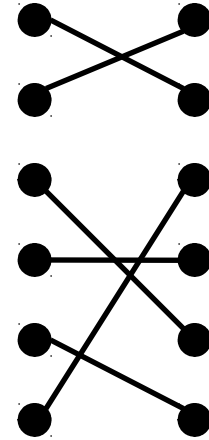
# Types of Relationships



**Many-to-Many**



**One-to-Many**



**One-to-One**



# Join Algorithms in MapReduce

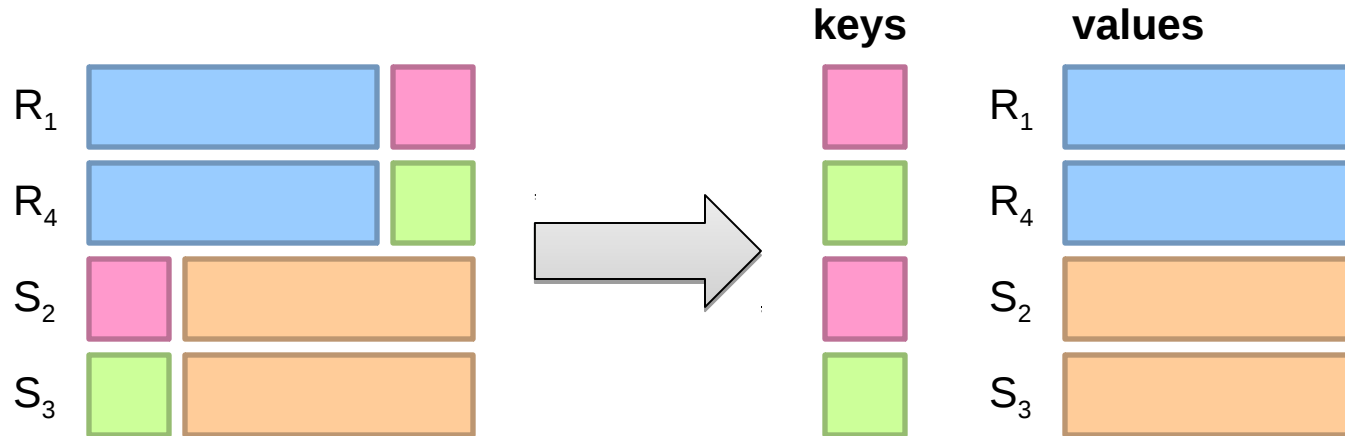
- Reduce-side join
- Map-side join
- In-memory join
  - Striped variant
  - Memcached variant

# Reduce-side Join

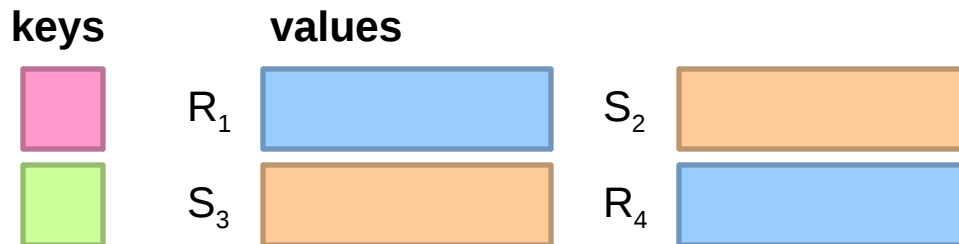
- Basic idea: group by join key
  - Map over both sets of tuples
  - Emit tuple as value with join key as the intermediate key
  - Execution framework brings together tuples sharing the same key
  - Perform actual join in reducer
  - Similar to a “sort-merge join” in database terminology
- Two variants
  - 1-to-1 joins
  - 1-to-many and many-to-many joins

# Reduce-side Join: 1-to-1

## Map



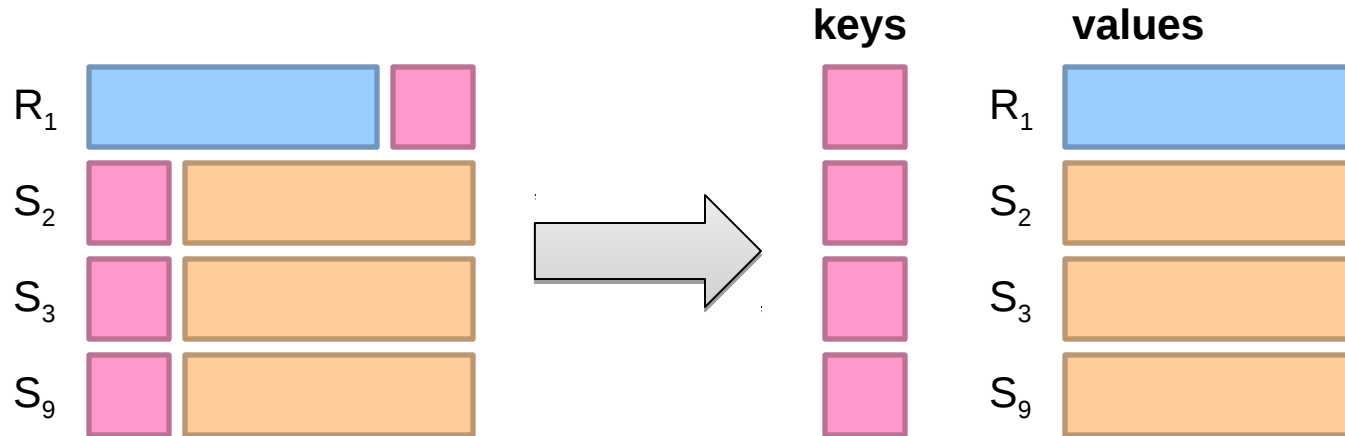
## Reduce



**Note: no guarantee if R is going to come first or S**

# Reduce-side Join: 1-to-many

## Map



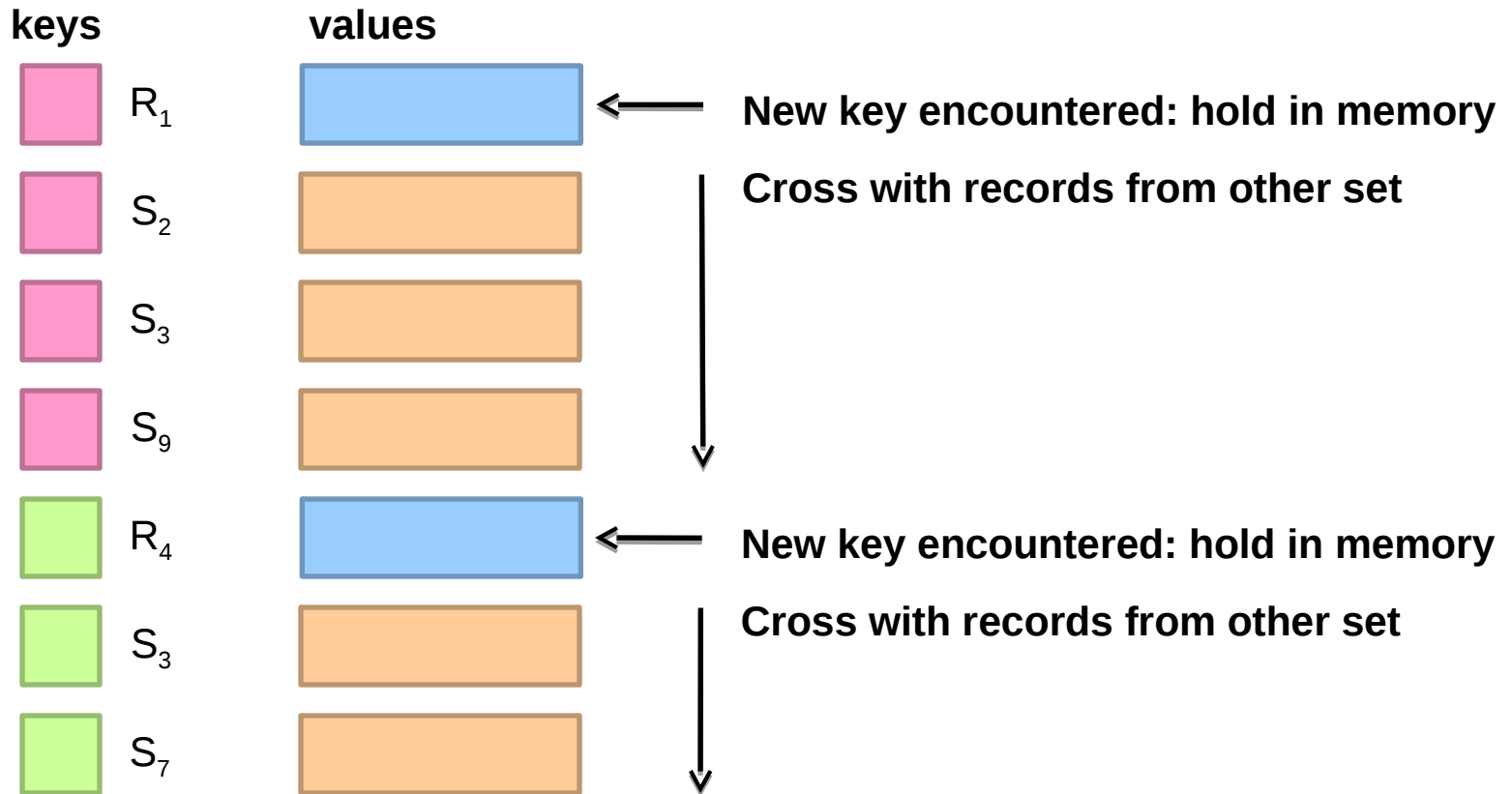
## Reduce



**What's the problem?**

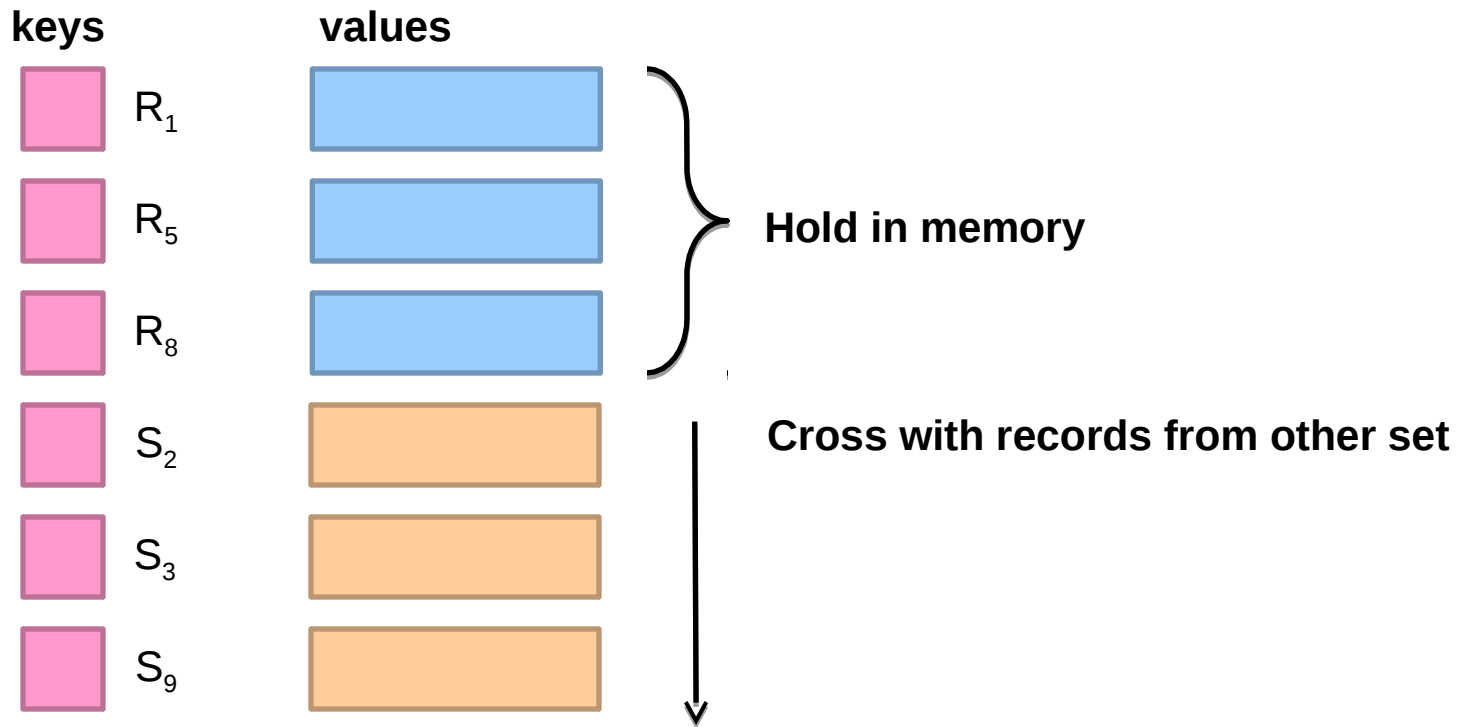
# Reduce-side Join: V-to-K Conversion

In reducer...



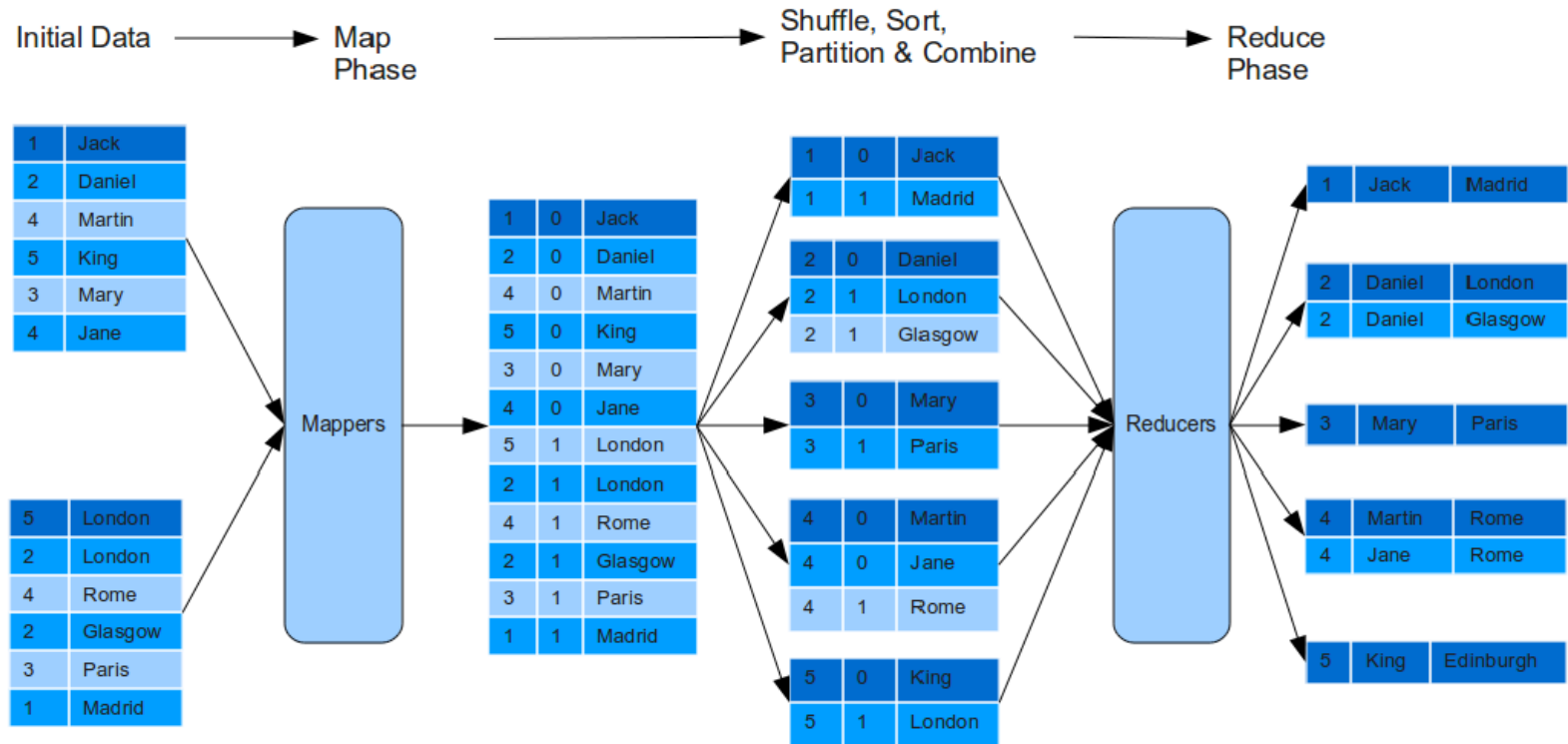
# Reduce-side Join: many-to-many

In reducer...



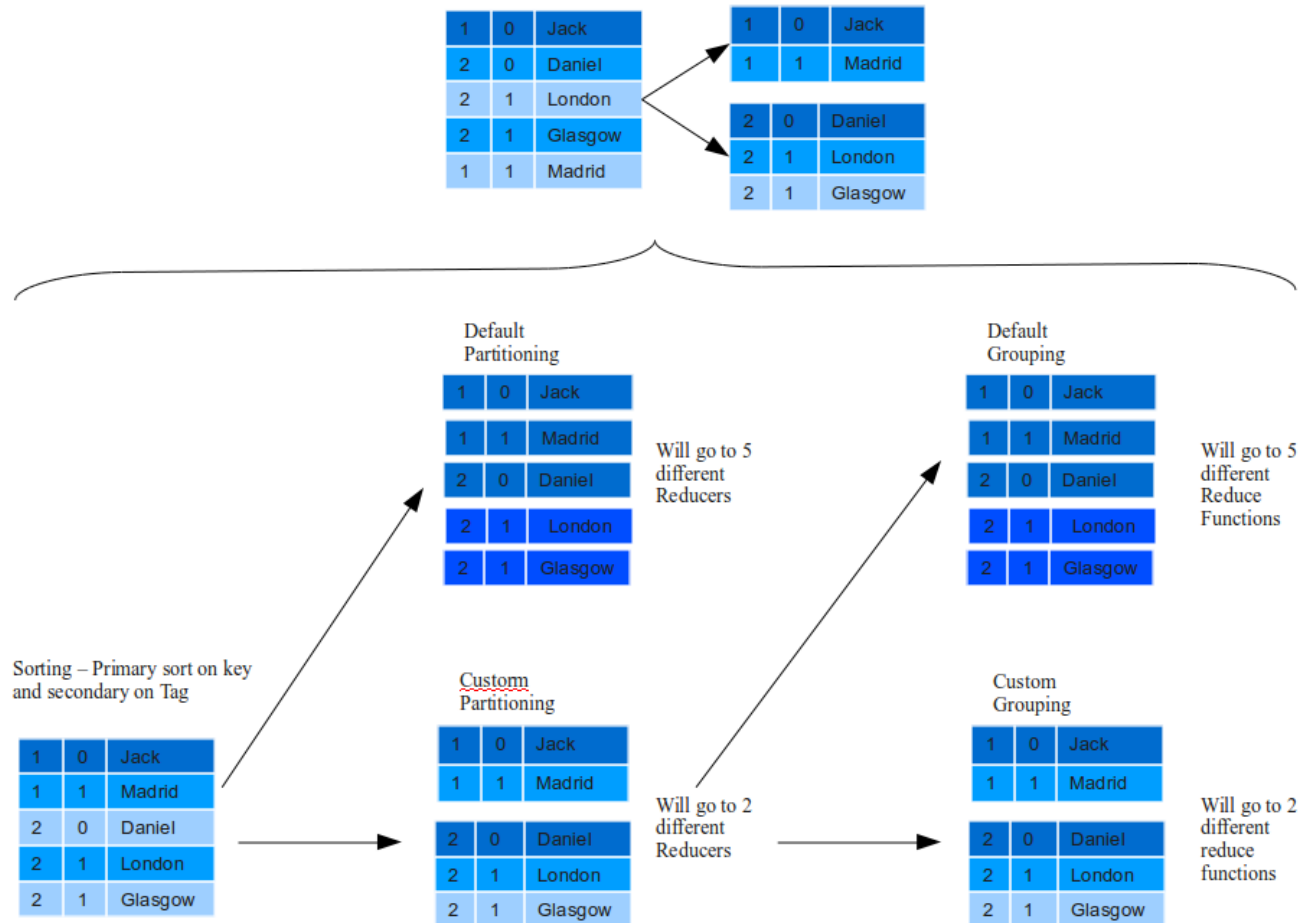
**What's the problem?**

# Reduce-side Join: many-to-many



Produce mapper output with composite key that includes foreign key and table name

# Reduce-side Join: many-to-many

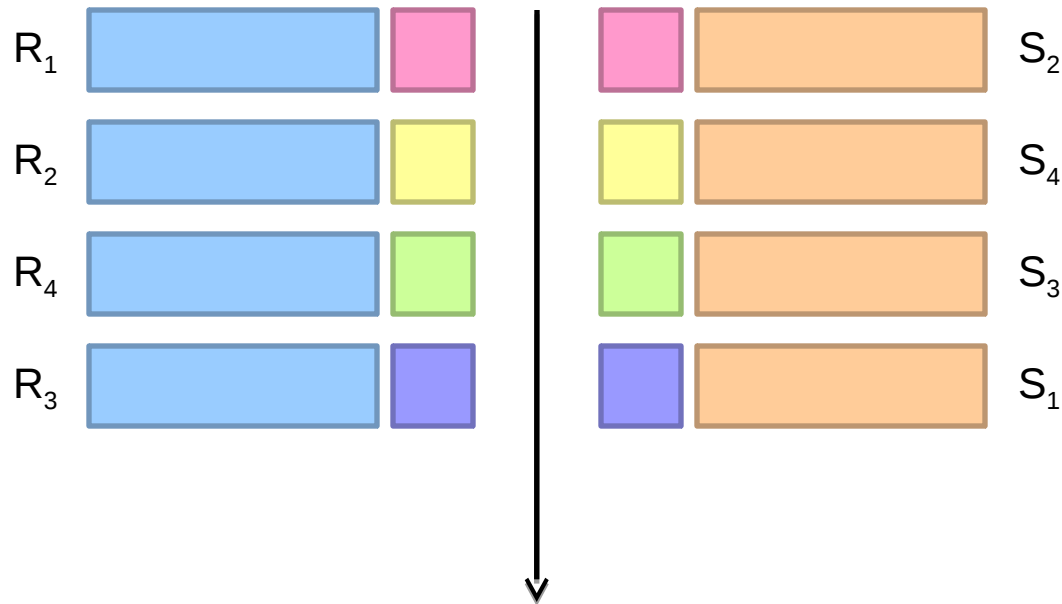


Use custom partitioning and grouping to send data with same key to a single reducer



# Map-side Join: Basic Idea

Assume two datasets are sorted by the join key:

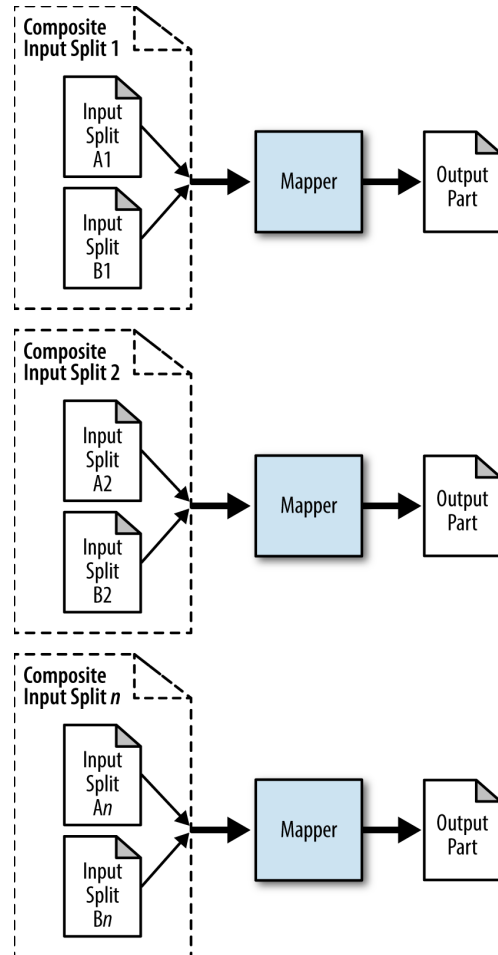


A sequential scan through both datasets to join  
(called a “merge join” in database terminology)

# Map-side Join: Parallel Scans

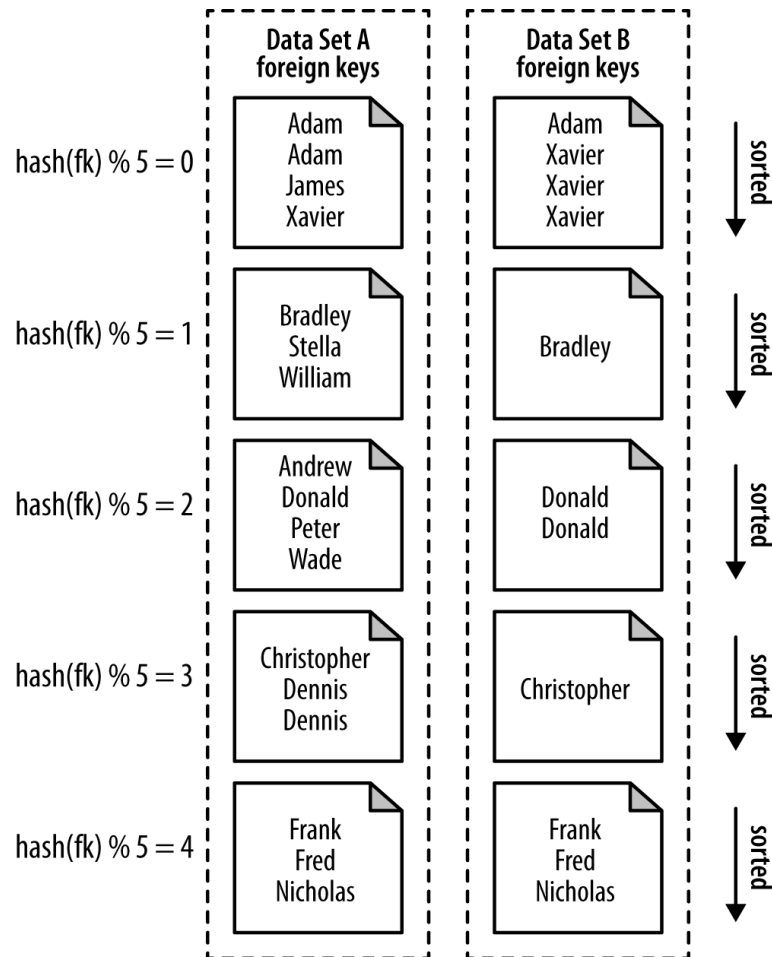
- If datasets are sorted by join key, join can be accomplished by a scan over both datasets
- How can we accomplish this in parallel?
  - Partition and sort both datasets in the same manner
- In MapReduce:
  - Map over one dataset, read from other corresponding partition
  - No reducers necessary (unless to repartition or resort)
- Consistently partitioned datasets: realistic to expect?

# Map-side Join: Parallel Scans



and split both A and B before sending to mapper. Mapper will produce output, reducer needed.

# Map-side Join: Parallel Scans



# Parallel Scan & Join

$p \in P; q \in Q; gq \in Q$

**while** more tuples in inputs **do**

**while**  $p.a < gq.b$  **do**

        advance  $p$

**end while**

**while**  $p.a > gq.b$  **do**

        advance  $gq$  {a group might begin here}

**end while**

**while**  $p.a == gq.b$  **do**

$q = gq$  {mark group beginning}

**while**  $p.a == q.b$  **do**

            Add  $\langle p, q \rangle$  to the result

            Advance  $q$

**end while**

        Advance  $p$  {move forward}

**end while**

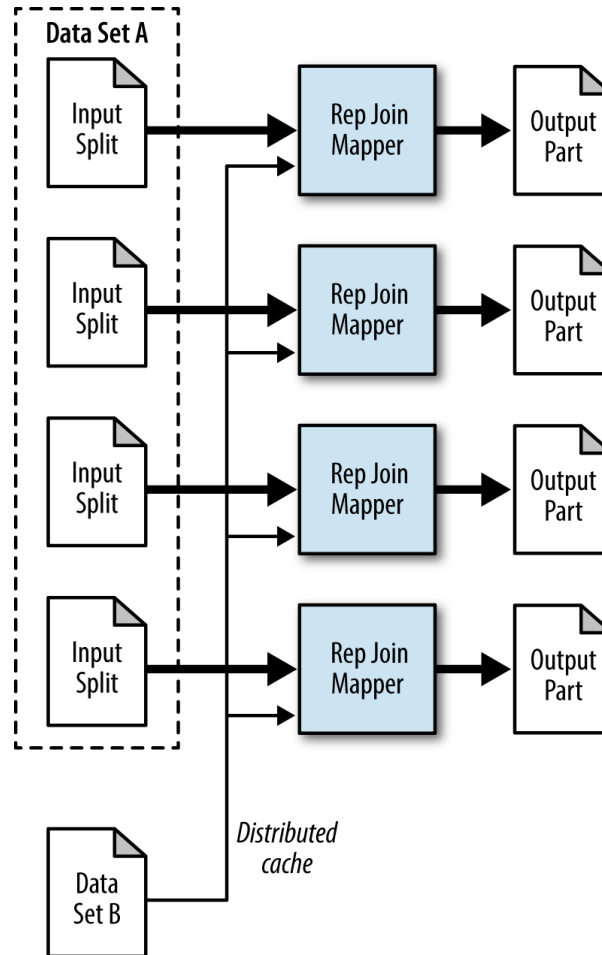
$gq = q$  {candidate to begin next group}

**end while**

# In-Memory Join

- Basic idea: load one dataset into memory, stream over other dataset
  - Works if  $R \ll S$  and  $R$  fits into memory
  - Called a “hash join” in database terminology
- MapReduce implementation
  - Distribute  $R$  to all nodes
  - Map over  $S$ , each mapper loads  $R$  in memory, hashed by join key
  - For every tuple in  $S$ , look up join key in  $R$
  - No reducers, unless for regrouping or resorting tuples

# In-Memory Join



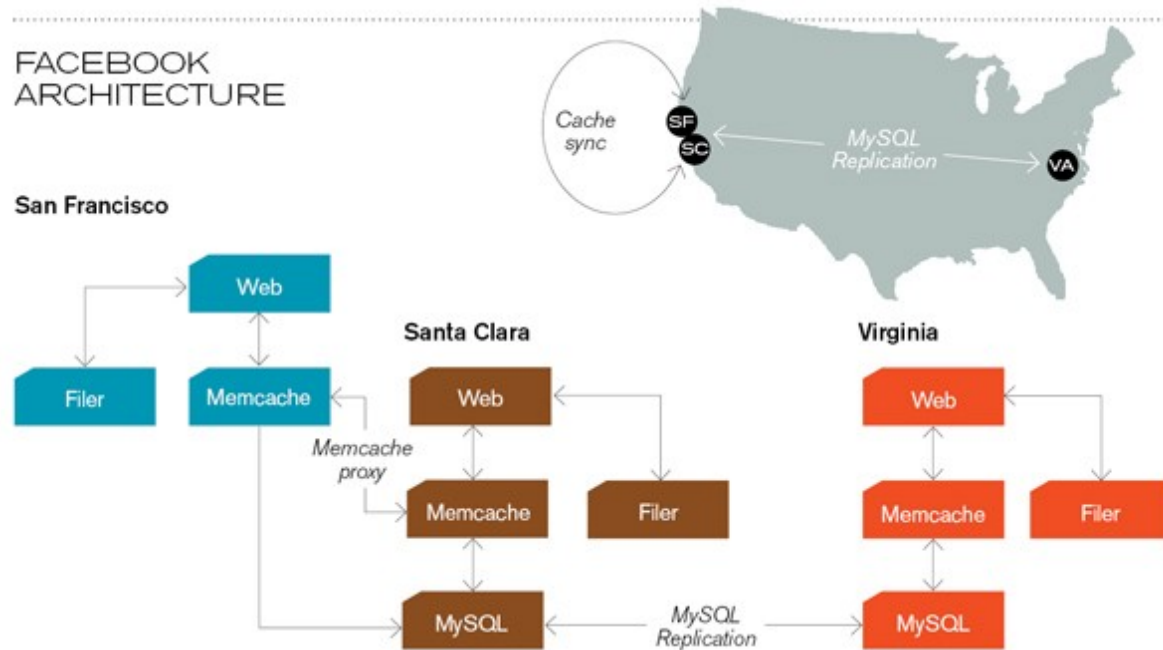
and distribute dataset B to all the mappers. For each key in B, iterate over all the  
r joining

# In-Memory Join: Variants

- Striped variant:
  - R too big to fit into memory?
  - Divide R into  $R_1, R_2, R_3, \dots$  s.t. each  $R_n$  fits into memory
  - Perform in-memory join:  $\forall n, R_n \bowtie S$
  - Take the union of all join results
- Memcached join:
  - Load R into memcached
  - Replace in-memory hash lookup with memcached lookup



# Memcached



# Memcached Join

- Memcached join:
  - Load R into memcached
  - Replace in-memory hash lookup with memcached lookup
- Capacity and scalability?
  - Memcached capacity  $\gg$  RAM of individual node
  - Memcached scales out with cluster
- Latency?
  - Memcached is fast (basically, speed of network)
  - Batch requests to amortize latency costs

# Which join to use?

- In-memory join > map-side join > reduce-side join
  - Why?
- Limitations of each?
  - In-memory join: memory
  - Map-side join: sort order and partitioning
  - Reduce-side join: general purpose

# Processing Relational Data: Summary

- MapReduce algorithms for processing relational data:
  - Group by, sorting, partitioning are handled automatically by shuffle/sort in MapReduce
  - Selection, projection, and other computations (e.g., aggregation), are performed either in mapper or reducer
  - Multiple strategies for relational joins
- Complex operations require multiple MapReduce jobs
  - Example: top ten URLs in terms of average time spent
  - Opportunities for automatic optimization