

# Distributed Databases

CS347  
Lecture 14  
May 30, 2001

1

## Topics for the Day

- Query processing in distributed databases
  - Localization
  - Distributed query operators
  - Cost-based optimization

2

## Query Processing Steps

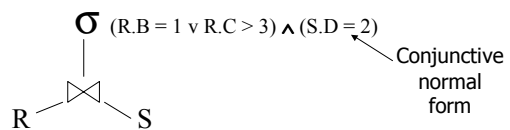
- Decomposition
  - Given SQL query, generate one or more algebraic query trees
- Localization
  - Rewrite query trees, replacing relations by fragments
- Optimization
  - Given cost model + one or more localized query trees
  - Produce minimum cost query execution plan

3

## Decomposition

- Same as in a centralized DBMS
- Normalization (usually into relational algebra)

Select A,C  
From R Natural Join S  
Where (R.B = 1 and S.D = 2) or (R.C > 3 and S.D = 2)



4

## Decomposition

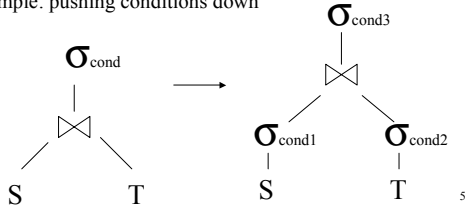
- Redundancy elimination

$$(S.A = 1) \wedge (S.A > 5) \Rightarrow \text{False}$$

$$(S.A < 10) \wedge (S.A < 5) \Rightarrow S.A < 5$$

- Algebraic Rewriting

- Example: pushing conditions down

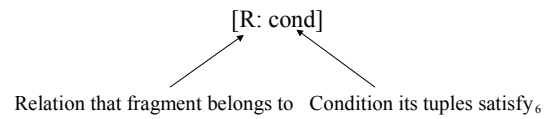


5

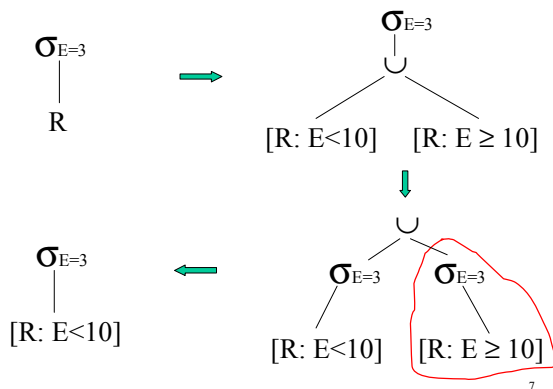
## Localization Steps

1. Start with query tree
2. Replace relations by fragments
3. Push  $\cup$  up &  $\pi, \sigma$  down (CS245 rules)
4. Simplify – eliminating unnecessary operations

Note: To denote fragments in query trees

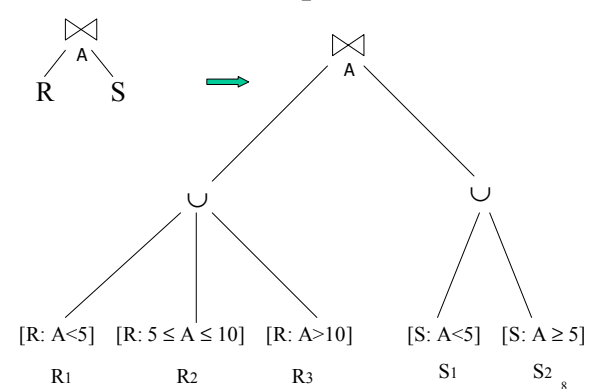


### Example 1

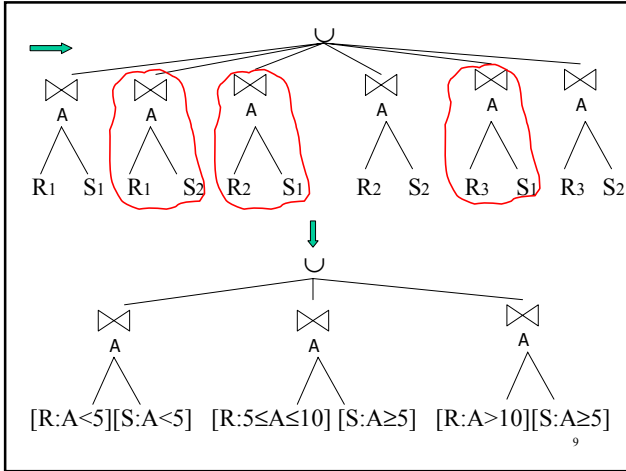


7

### Example 2



8



### Rules for Horiz. Fragmentation

- $\sigma_{C_1}[R: C_2] \Rightarrow [R: C_1 \wedge C_2]$
- $[R: \text{False}] \Rightarrow \emptyset$
- $[R: C_1] \bowtie_A [S: C_2] \Rightarrow [R \bowtie_A S: C_1 \wedge C_2 \wedge R.A = S.A]$
- In Example 1:  
 $\sigma_{E=3}[R_2: E \geq 10] \Rightarrow [R_2: E=3 \wedge E \geq 10]$   
 $\Rightarrow [R_2: \text{False}] \Rightarrow \emptyset$
- In Example 2:  
 $[R: A < 5] \bowtie_A [S: A \geq 5]$   
 $\Rightarrow [R \bowtie_A S: R.A < 5 \wedge S.A \geq 5 \wedge R.A = S.A]$   
 $\Rightarrow [R \bowtie_A S: \text{False}] \Rightarrow \emptyset$

10

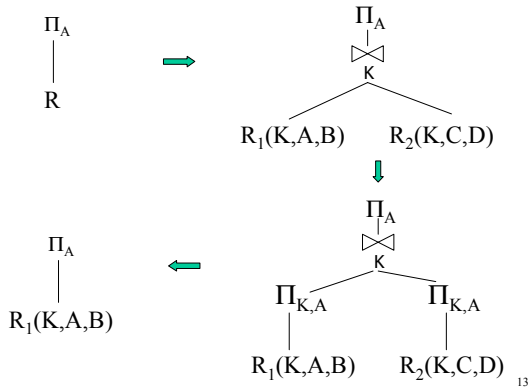
### Example 3 – Derived Fragmentation

S's fragmentation is derived from that of R.

11

12

### Example 4 – Vertical Fragmentation



### Rule for Vertical Fragmentation

- Given vertical fragmentation of  $R(A)$ :  
 $R_i = \Pi_{A_i}(R), A_i \subseteq A$
- For any  $B \subseteq A$ :  
 $\Pi_B(R) = \Pi_B \left[ \bowtie_i R_i \mid B \cap A_i \neq \emptyset \right]$

14

### Parallel/Distributed Query Operations

- Sort
  - Basic sort
  - Range-partitioning sort
  - Parallel external sort-merge
- Join
  - Partitioned join
  - Asymmetric fragment and replicate join
  - General fragment and replicate join
  - Semi-join programs
- Aggregation and duplicate removal

15

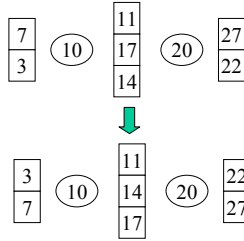
### Parallel/distributed sort

- Input: relation  $R$  on
  - single site/disk
  - fragmented/partitioned by sort attribute
  - fragmented/partitioned by some other attribute
- Output: sorted relation  $R$ 
  - single site/disk
  - individual sorted fragments/partitions

16

## Basic sort

- Given  $R(A, \dots)$  range partitioned on attribute A, sort R on A

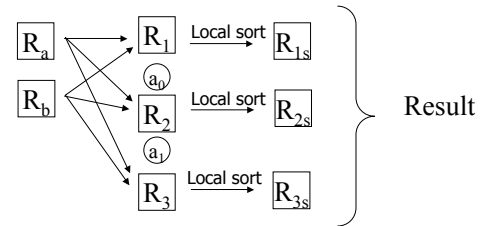


- Each fragment is sorted independently
- Results shipped elsewhere if necessary

17

## Range partitioning sort

- Given  $R(A, \dots)$  located at one or more sites, not fragmented on A, sort R on A
- Algorithm: range partition on A and then do basic sort



18

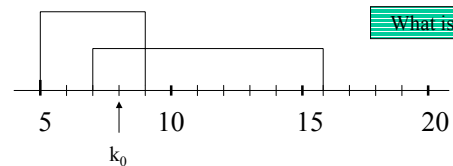
## Selecting a partitioning vector

- Possible centralized approach using a “coordinator”
  - Each site sends *statistics* about its fragment to coordinator
  - Coordinator decides # of sites to use for local sort
  - Coordinator computes and distributes partitioning vector
- For example,
  - Statistics could be (min sort key, max sort key, # of tuples)
  - Coordinator tries to choose vector that equally partitions relation

19

## Example

- Coordinator receives:
  - From site 1: Min 5, Max 9, 10 tuples
  - From site 2: Min 7, Max 16, 10 tuples
- Assume sort keys distributed uniformly within  $[\min, \max]$  in each fragment
- Partition R into two fragments



20

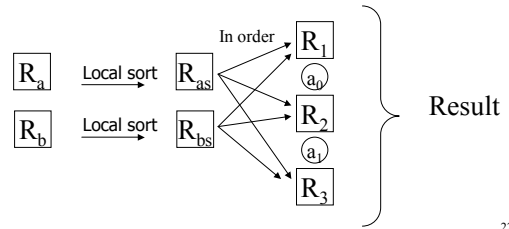
## Variations

- Different kinds of statistics
    - Local partitioning vector
    - Histogram
- $$\begin{array}{ccccccc} & & & \text{Site } l & & & \\ & & & \leftarrow \# \text{ of tuples} & & & \\ \hline & 3 & 4 & 3 & & & \\ & 5 & 6 & 8 & 10 & \leftarrow \text{local vector} & \end{array}$$
- Multiple rounds between coordinator and sites
    - Sites send statistics
    - Coordinator computes and distributes initial vector  $V$
    - Sites tell coordinator the number of tuples that fall in each range of  $V$
    - Coordinator computes final partitioning vector  $V_f$

21

## Parallel external sort-merge

- Local sort
- Compute partition vector
- Merge sorted streams at final sites



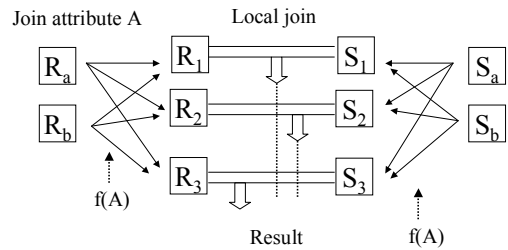
22

## Parallel/distributed join

- Input: Relations  $R, S$   
 May or may not be partitioned
- Output:  $R \bowtie S$   
 Result at one or more sites

23

## Partitioned Join



Note: Works only for equi-joins

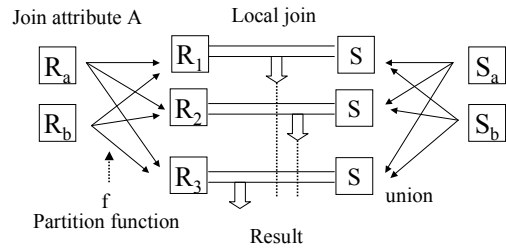
24

## Partitioned Join

- Same partition function (f) for both relations
- f can be range or hash partitioning
- Any type of local join (nested-loop, hash, merge, etc.) can be used
- Several possible scheduling options. Example:
  - partition R; partition S; join
  - partition R; build local hash table for R; partition S and join
- Good partition function important
  - Distribute join load evenly among sites

25

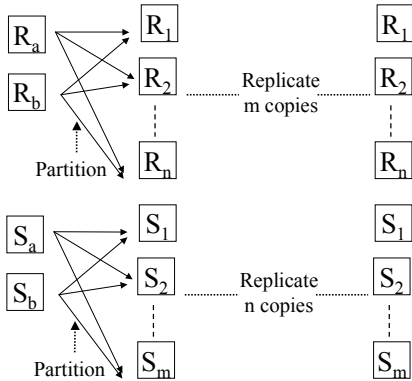
## Asymmetric fragment + replicate join



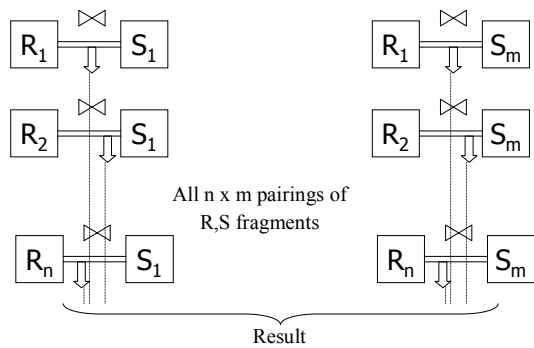
- Any partition function f can be used (even round-robin)
- Can be used for any kind of join, not just equi-joins

26

## General fragment + replicate join



27



- Asymmetric F+R join is a special case of general F+R.
- Asymmetric F+R is useful when S is small.

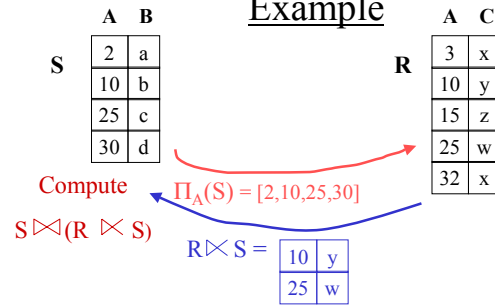
28

## Semi-join programs

- Used to reduce communication traffic during join processing
- $R \bowtie S = (R \bowtie S) \bowtie S$   
 $= R \bowtie (S \bowtie R)$   
 $= (R \bowtie S) \bowtie (S \bowtie R)$

29

## Example



- Using semi-join, communication cost = 4 A + 2 (A + C) + result
- Directly joining R and S, communication cost = 4 (A + B) + result

30

## Comparing communication costs

- Say R is the smaller of the two relations R and S
- $(R \bowtie S) \bowtie S$  is cheaper than  $R \bowtie S$  if  
 $\text{size}(\Pi_A S) + \text{size}(R \bowtie S) < \text{size}(R)$
- Similar comparisons for other types of semi-joins
- Common implementation trick:
  - Encode  $\Pi_A S$  (or  $\Pi_A R$ ) as a bit vector
  - 1 bit per domain of attribute A

001101000010100

31

## n-way joins

- To compute  $R \bowtie S \bowtie T$ 
  - Semi-join program 1:  $R' \bowtie S' \bowtie T$   
 where  $R' = R \bowtie S$  &  $S' = S \bowtie T$
  - Semi-join program 2:  $R'' \bowtie S' \bowtie T$   
 where  $R'' = R \bowtie S'$  &  $S' = S \bowtie T$
  - Several other options
- In general, number of options is exponential in the number of relations

32

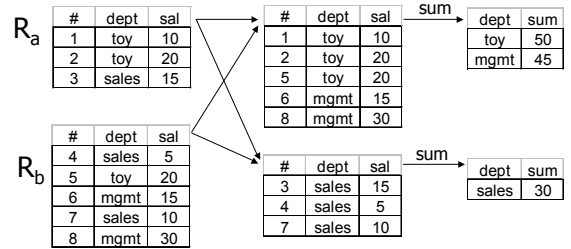


## Other operations

- Duplicate elimination
  - Sort first (in parallel), then eliminate duplicates in the result
  - Partition tuples (range or hash) and eliminate duplicates locally
- Aggregates
  - Partition by grouping attributes; compute aggregates locally at each site

33

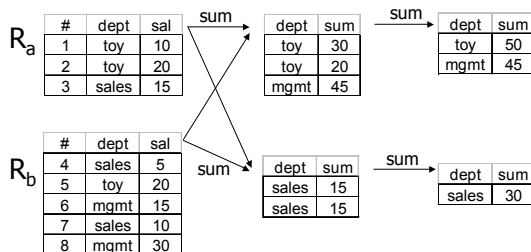
## Example



sum(sal) group by dept

34

## Example



Does this work for all kinds of aggregates?

Aggregate during partitioning to reduce communication cost

35

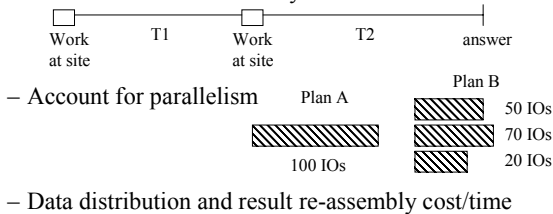
## Query Optimization

- Generate query execution plans (QEPs)
- Estimate cost of each QEP (\$, time, ...)
- Choose minimum cost QEP
- What's different for distributed DB?
  - New strategies for some operations (semi-join, range-partitioning sort, ...)
  - Many ways to assign and schedule processors
  - Some factors besides number of IO's in the cost model

36

## Cost estimation

- In centralized systems - estimate sizes of intermediate relations
- For distributed systems
  - Transmission cost/time may dominate



37

## Optimization in distributed DBs

- Two levels of optimization
- Global optimization
  - Given localized query and cost function
  - Output optimized (min. cost) QEP that includes relational and communication operations on fragments
- Local optimization
  - At each site involved in query execution
  - Portion of the QEP at a given site optimized using techniques from centralized DB systems

38

## Search strategies

1. Exhaustive (with pruning)
2. Hill climbing (greedy)
3. Query separation

39

## Exhaustive with Pruning

- A fixed set of techniques for each relational operator
- Search space = “all” possible QEPs with this set of techniques
- Prune search space using heuristics
- Choose minimum cost QEP from rest of search space

40

