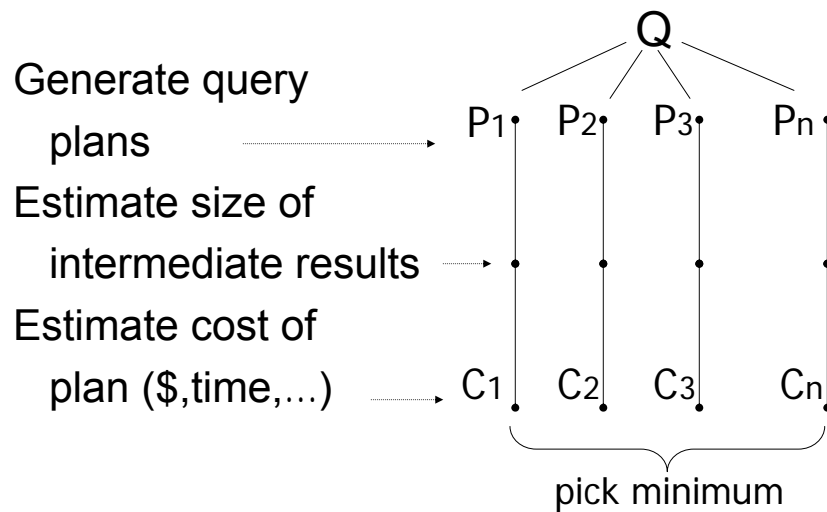


# Query Optimization

It is safer to accept any chance that offers itself, and extemporize a procedure to fit it, than to get a good plan matured, and wait for a chance of using it.

Thomas Hardy (1874)  
in *Far from the Madding Crowd*

## Query optimization



# Query Optimization

- Query:  $R1 \bowtie R2 \bowtie R3$
- $R1@Site1, R2@Site2, R3@Site3$
- $Result@Site1$
- Possible plans
  - Plan 1: Send R2 and R3 to Site1. Perform query at Site 1
  - Plan 2: Send R3 to Site2; Evaluate  $I = R2 \bowtie R3$ ; send I to Site1; Evaluate result =  $I \bowtie R1$
  - Many other plans ... (including types of joins, number of sites, semijoins, etc)

CS5225

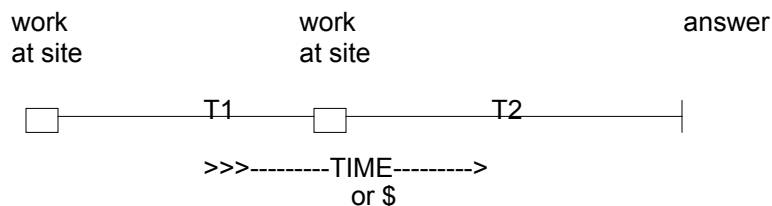
Query Optimization

3

# Cost estimation

- ☛ As in centralized system: estimate result sizes
- ☛ But: # IOs may not be best metric

e.g., Transmission time may dominate



CS5225

Query Optimization

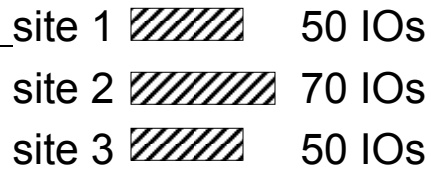
4

## Another reason why plain IOs not enough: Parallelism

Plan A



Plan B

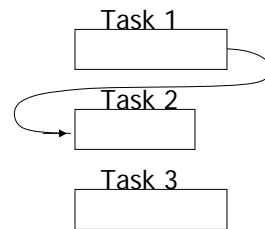


CS5225

Query Optimization

5

- Cost metrics
  - IOs, Bytes transmitted, \$, ...
  - Can add together
- Response time metric
  - cannot add
  - need scheduling and dependency info
  - skew important



CS5225

Query Optimization

6

## ⇔ Take into account: (in parallel/distributed system)

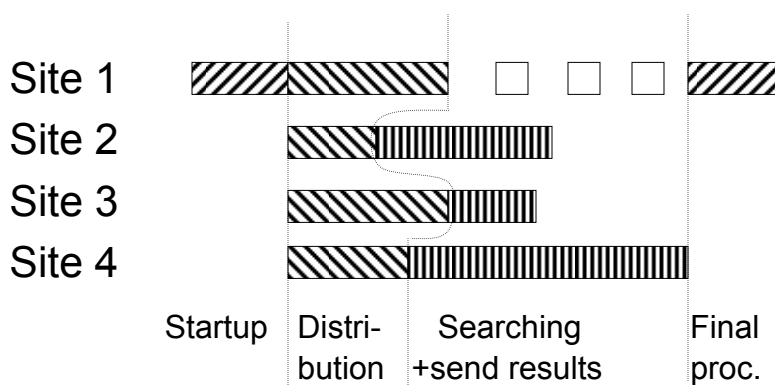
- Start up costs (for parallel operation)
- Data distribution costs/time
- Contention
  - memory, disk, network,...
- Assembling result

CS5225

Query Optimization

7

## Example: Response time

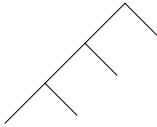
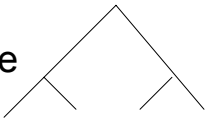


CS5225

Query Optimization

8

# Query Optimizer

- Cost model
- Plan space
  - Deep tree  vs bushy tree 
- Enumeration/Search strategy
  - Exhaustive (with pruning)
  - Hill climbing (greedy)
  - Query separation
  - SDD-1 (semi-join based)

CS5225

Query Optimization

9

## (1) Exhaustive

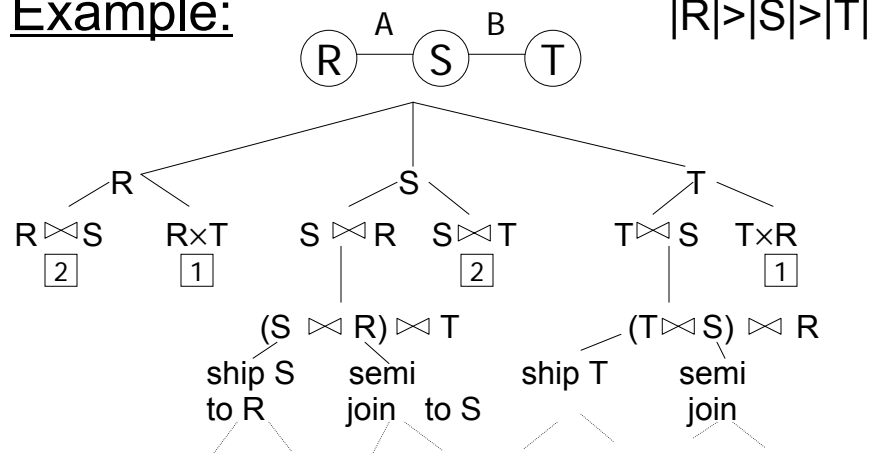
- consider "all" query plans  
with a set of techniques
- prune some plans (heuristics)

CS5225

Query Optimization

10

**Example:**



Heuristics:

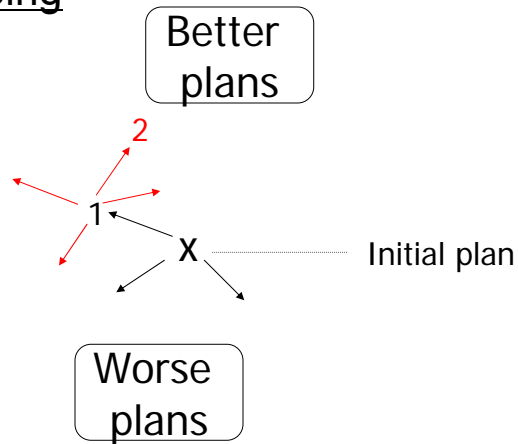
- 1 Prune because cross-product not necessary
- 2 Prune because larger relation first

**⊗** In generating plans, keep goal in mind:

e.g.: Goal is parallelism in system with fast net, consider partitioning relation(s) first

e.g.: Goal is reduction of net traffic, consider semi-joins

## (2) Hill climbing



CS5225

Query Optimization

13

## Hill Climbing Algorithm

**Step 1:** Do initial processing

**Step 2:** Select initial feasible solution ( $P_0$ )

- 2.1 Determine the candidate result sites - sites where a relation referenced in the query exist
- 2.2 Compute the cost of transferring all the other referenced relations to each candidate site
- 2.3  $P_0$  = candidate site with minimum cost

**Step 3:** Determine candidate splits of  $P_0$  into

$$P_1 = \{P_{1a}, P_{1b}\}$$

- 3.1  $P_{1a}$  consists of sending one of the relations to the other relation's site
- 3.2  $P_{1b}$  consists of sending the join of the relations to the final result site

CS5225

Query Optimization

14

## Hill Climbing Algorithm

**Step 4:** Replace  $P_0$  with  $P_1$  that gives

$$\text{cost}(P_{1a}) + \text{cost}(\text{local join}) + \text{cost}(P_{1b}) < \text{cost}(P_0)$$

**Step 5:** Recursively apply steps 3–4 on  $P_1$  until no such plans can be found

**Step 6:** Check for redundant transmissions in the final plan and eliminate them.

CS5225

Query Optimization

15

Example  $R \bowtie S \bowtie \sigma(U) \bowtie V$

Rel	Site	Size	tuple size = 1
R	1	10	
S	2	20	
U	3	90	
V	4	40	

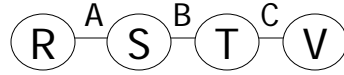
Goal: minimize data transmission

CS5225

Query Optimization

16

Step 1:  $R \bowtie S \bowtie T \bowtie V$



Rel Site Size

tuple size = 1

R 1 10

S 2 20

T 3 30

V 4 40

$T = \sigma(U)$

Selectivity is 1/3

Step 2: Initial plan – send relations to one site

What site do we send all relations to?

To site 1: cost=20+30+40=90

To site 2: cost=10+30+40=80

To site 3: cost=10+20+40=70

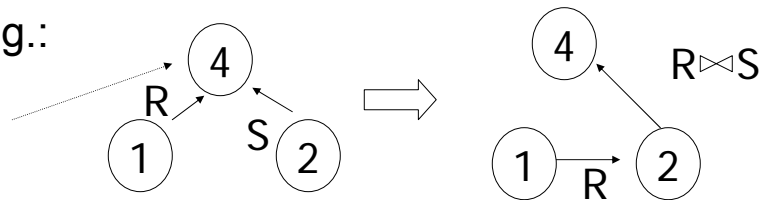
To site 4: cost=10+20+30=60

$P_0$ :  $R (1 \rightarrow 4)$   
 $S (2 \rightarrow 4)$   
 $T (3 \rightarrow 4)$   
 Compute  $R \bowtie S \bowtie T \bowtie V$  at site 4

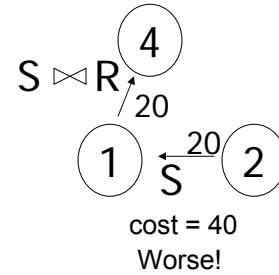
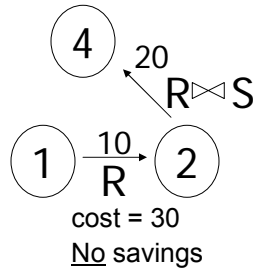
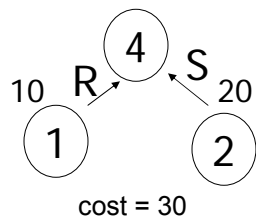
## Steps 3 & 4

- Consider sending each relation to neighbor:

e.g.:



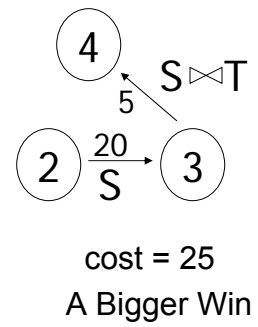
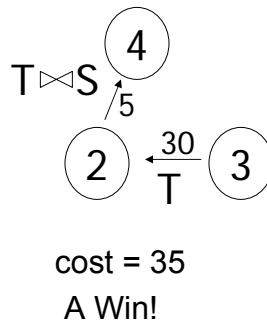
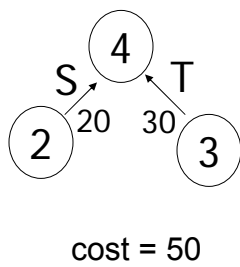
Assume: Size     $R \bowtie S = 20$   
                           $S \bowtie T = 5$   
                           $T \bowtie V = 1$



CS5225

Query Optimization

21



CS5225

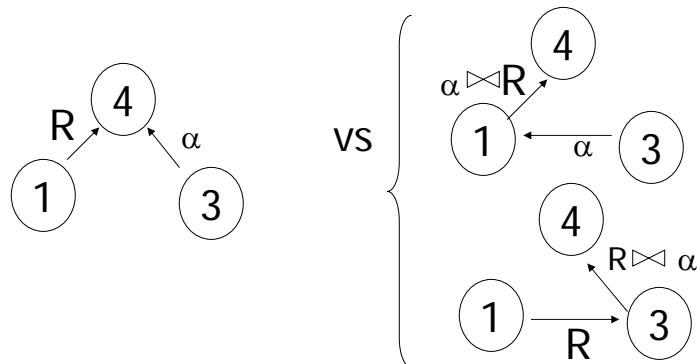
Query Optimization

22

P1: P1a: S (2 → 3)  
 $\alpha = S \bowtie T$   
 P1b: R (1 → 4)  
 $\alpha$  (3 → 4)  
 compute answer at site 4

## Step 5: Repeat Steps 3 & 4

- Treat  $\alpha = S \bowtie T$  as relation



# Hill climbing may miss best plan!

Example: best plan could be:

$P_B$ :  $T (3 \rightarrow 4)$   $\boxtimes 30$

$\beta = T \bowtie V$

$\beta (4 \rightarrow 2)$   $\boxtimes 1$

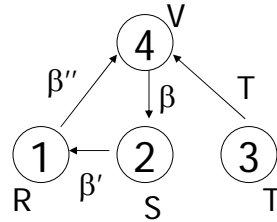
$\beta' = \beta \bowtie S$

$\beta' (2 \rightarrow 1)$   $\boxtimes 1$

$\beta'' = \beta' \bowtie R$

[optional]  $\beta'' (1 \rightarrow 4)$   $\boxtimes 1$

Compute answer 33 = total



Costs could be low because is very selective

## (3) Query separation

- separate query into 2 or more steps
- optimize each step independently

## Example: simple queries

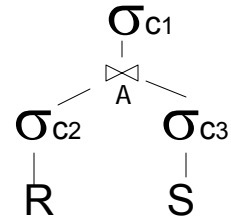
1. Compute  $R' = \Pi_A[\sigma_{c_2} R]$

$S' = \Pi_A[\sigma_{c_3} S]$

2. Compute  $J = R' \bowtie S'$

3. Compute

$$\text{Ans} = \sigma_{c_1} \{ [J \bowtie \sigma_{c_2} R] \bowtie [J \bowtie \sigma_{c_3} S] \}$$



CS5225

Query Optimization

27

## In other words:

(a) Compute  $\underline{A}$  values in answer  
(steps 1,2)

(b) Get tuples from sites with matching  
 $\underline{A}$  values and compute answer  
(step 3)

CS5225

Query Optimization

28

## Simple query

- Relations have a single attribute
  - Output has a single attribute
- e.g.,  $J \leftarrow R' \bowtie S'$

## Idea

- Decompose query into
  - Local processing
  - Simple query (or queries)
  - Final processing
- Optimize simple query
- Philosophy
  - Hard part is distributed join
  - Do this part with only keys; get rest of data later
  - Simpler to optimize simple queries

## SDD-1 Algorithm

**Step 1:** In the execution strategy (call it *ES*), include all the local processing

**Step 2:** Reflect the effects of local processing on the database profile

**Step 3:** Construct a set of beneficial semijoin operations (*BS*) as follows :

$$BS = \emptyset$$

For each semijoin  $SJ_i$

$$BS \leftarrow BS \cup SJ_i \text{ if } cost(SJ_i) < benefit(SJ_i)$$

CS5225

Query Optimization

31

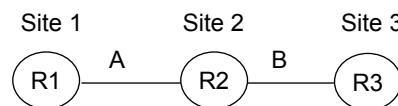
## SDD-1 Algorithm – Example

Consider the following query

```

SELECT R3.C
FROM   R1, R2, R3
WHERE  R1.A = R2.A
AND    R2.B = R3.B
    
```

which has the following query graph and statistics:



SF: selectivity factor  
Result = SF \* R

CS5225

Query Optim...

relation	card	tuple size	relation size
R1	30	50	1500
R2	100	30	3000
R3	50	40	2000

attribute	SF <sub>X</sub>	size( $\Pi_{\text{attribute}}$ )
R1.A	0.3	36
R2.A	0.8	320
R2.B	1.0	400
R3.B	0.4	80

32

## SDD-1 Algorithm – Example

- Beneficial semijoins: Assume:  $T_{MSG} = 0$ 
  - $SJ_1 = R2 \bowtie R1$ , whose benefit is
    - $2100 = (1 - 0.3) * 3000$  and cost is 36
  - $SJ_2 = R2 \bowtie R3$ , whose benefit is
    - $1800 = (1 - 0.4) * 3000$  and cost is 80
- Nonbeneficial semijoins:
  - $SJ_3 = R1 \bowtie R2$ , whose benefit is
    - $300 = (1 - 0.8) * 1500$  and cost is 320
  - $SJ_4 = R3 \bowtie R2$ , whose benefit is 0 and cost is 400

$$\begin{aligned} \text{Cost} &= \text{transfer semijoin attribute} \\ &= T_{MSG} + T_{TR} * (\text{size}) \end{aligned}$$

$$\begin{aligned} \text{Benefit} &= \text{cost of transferring irrelevant tuples} \\ &\quad \text{(avoided by the semijoin)} \\ &= (1 - SF) * \text{size} * T_{TR} \end{aligned}$$

CS5225

Query Optimization

33

## SDD-1 Algorithm

### Iterative Process

- Step 4:** Remove the most beneficial  $SJ_i$  from  $BS$  and append it to  $ES$
- Step 5:** Modify the database profile accordingly
- Step 6:** Modify  $BS$  appropriately
  - compute new benefit/cost values
  - check if any new semijoin need to be included in  $BS$
- Step 7:** If  $BS \neq \emptyset$ , go back to Step 4.

CS5225

Query Optimization

34

## SDD-1 Algorithm – Example

- Iteration 1:
  - Remove  $SJ_1$  from  $BS$  and add it to  $ES$ .
  - Update statistics

$$R2' = size(R2) = 900 (= 3000 \times 0.3)$$

$$size(R2'.A) = 320 \times 0.3 = 96$$

$$SF_{\times}(R2'.A) = \sim 0.8 \times 0.3 = \sim 0.24$$

- Iteration 2:
  - Two beneficial semijoins:

$$SJ_2 = R2' \bowtie R3, \text{ whose benefit is } 540 = (1 - 0.4) \times 900 \text{ and cost is } 80$$

$$SJ_3 = R1 \bowtie R2', \text{ whose benefit is } \mathbf{300 = (1 - 0.8) \times 1500} \text{ and cost is } 96$$

- Add  $SJ_2$  to  $ES$  NOT  $1140 = (1 - 0.24) \times 1500!!$

- Update statistics
  - $size(R2') = 360 (= 900 \times 0.4)$

Note: selectivity of  $R2'$  may also change, but not important in this example.

relation	card	tuple size	relation size
R1	30	50	1500
<b>R2'</b>	<b>30</b>	<b>30</b>	<b>900</b>
R3	50	40	2000

attribute	$SF_{\times}$	$size(\Pi_{\text{attribute}})$
R1.A	0.3	36
<b>R2'.A</b>	<b>0.24</b>	<b>96</b>
<b>R2'.B</b>	<b>1.0</b>	<b>120</b>
R3.B	0.4	80

CS5225

Query Optimization

35

## SDD-1 Algorithm – Example

- Iteration 3:
  - No new beneficial semijoins.
  - Remove remaining beneficial semijoin  $SJ_3$  from  $BS$  and add it to  $ES$ .
  - Update statistics
    - $size(R1) = 1200 (= 1500 \times 0.8)$
    - $SF_{\times}(R1.A) = \sim 0.3 \times 0.8 = 0.24$

CS5225

Query Optimization

36

## SDD-1 Algorithm

### Assembly Site Selection

**Step 8:** Find the site where the largest amount of data resides and select it as the assembly site

Example:

Amount of data stored at sites:

Site 1: 1200

Site 2: 360

Site 3: 2000

Therefore, Site 3 will be chosen as the assembly site.

## SDD-1 Algorithm

### Postprocessing

**Step 9:** For each  $R_i$  at the assembly site, find the semijoins of the type  $R_i \bowtie R_j$  where the total cost of  $ES$  without this semijoin is smaller than the cost with it and remove the semijoin from  $ES$ .

**Step 10:** Permute the order of semijoins if doing so would improve the total cost of  $ES$ .

- Example: Final strategy:
  - Send  $(R2 \bowtie R1) \bowtie R3$  to Site 3
  - Send  $(R1 \bowtie R2)$  to Site 3

## Summary: Query Optimization

- Cost/result estimation
- Three key components in optimizer
  - Cost model, search space, enumeration algorithm
- In practice, avoid bad plans (rather than find optimal)
- Strategies
  - Exhaustive
  - Hill climbing
  - Separation
  - SDD-1 (semi-join based)