

Lecture 24: Query Execution

Monday, November 27, 2006

Outline

- Query optimization: algebraic laws 16.2

Example

Product(pname, maker), **Company**(cname, city)

```
Select Product.pname  
From Product, Company  
Where Product.maker=Company.cname  
and Company.city = "Seattle"
```

- How do we execute this query ?

Example

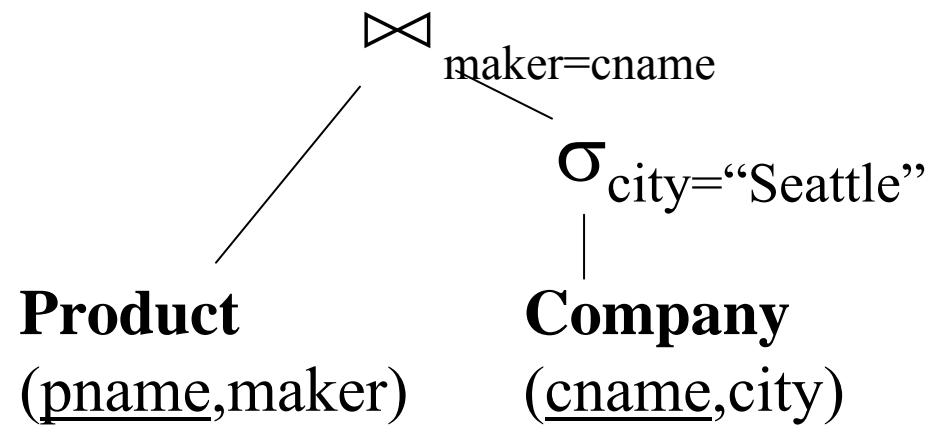
Product(pname, maker), **Company**(cname, city)

Assume:

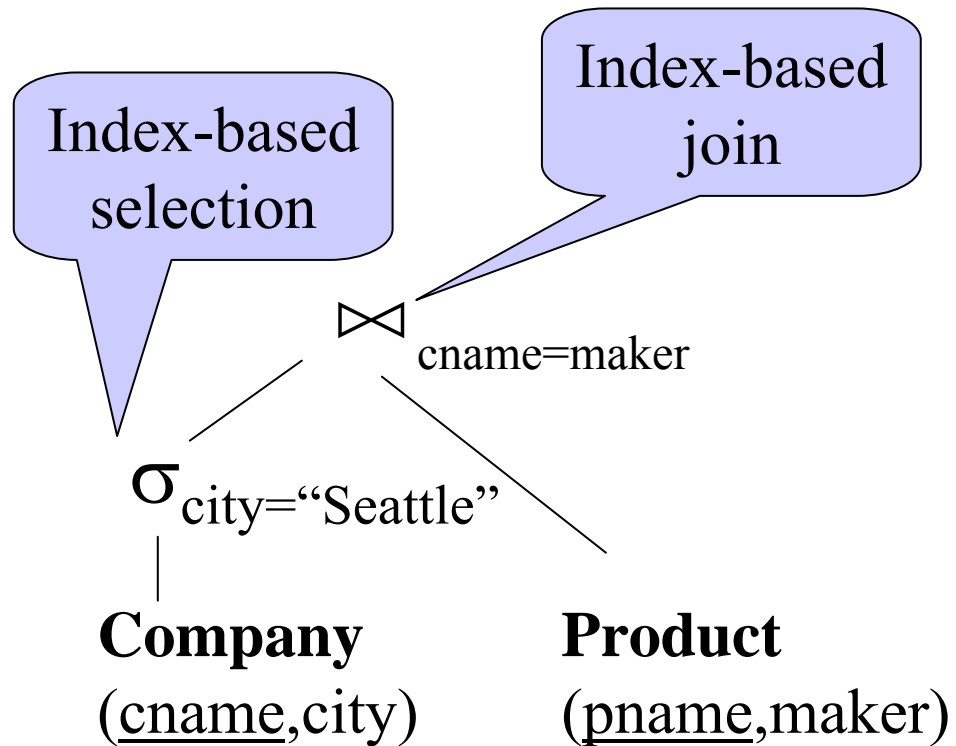
Clustered index: **Product**.pname, **Company**.cname

Unclustered index: **Product**.maker, **Company**.city

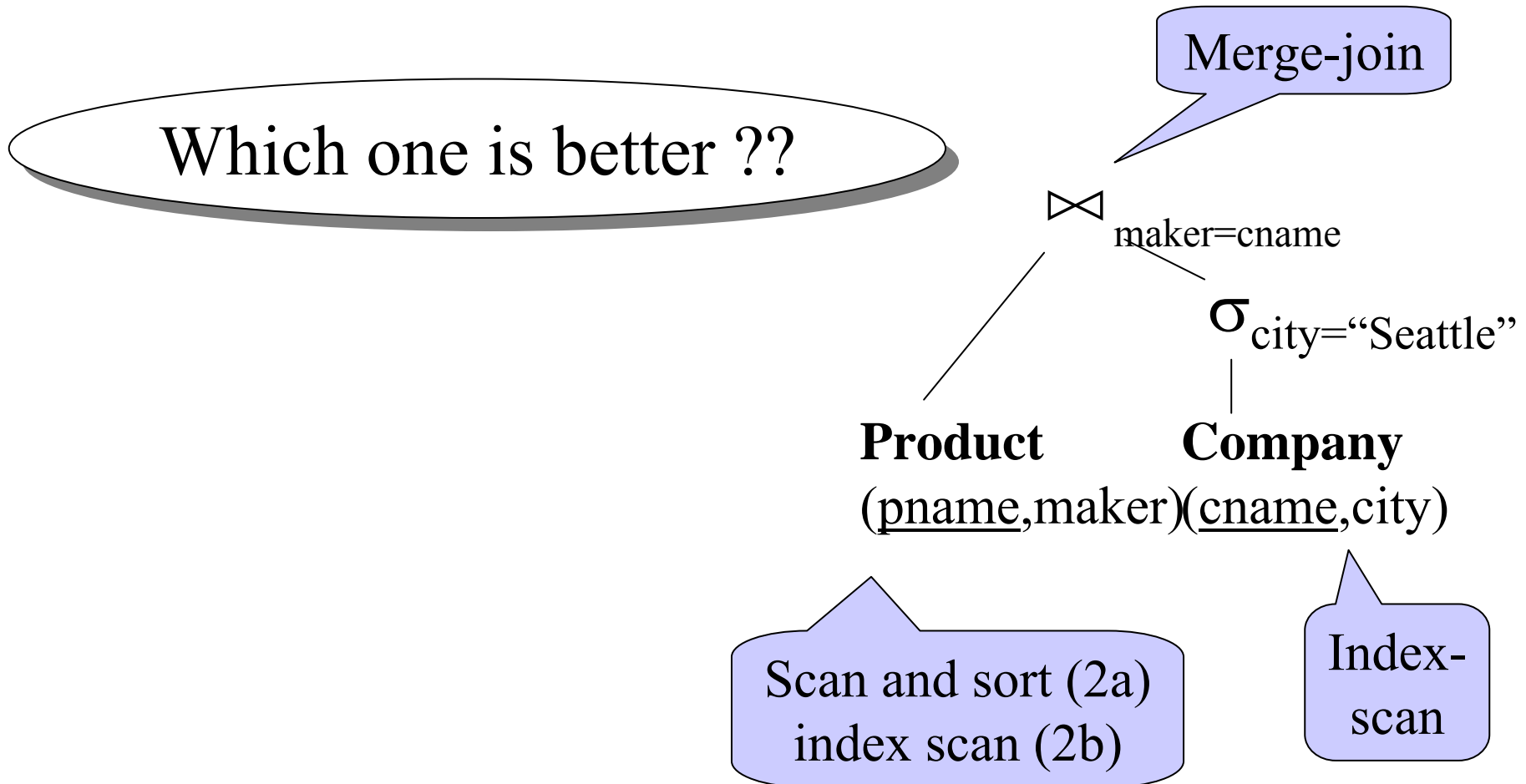
Logical Plan:

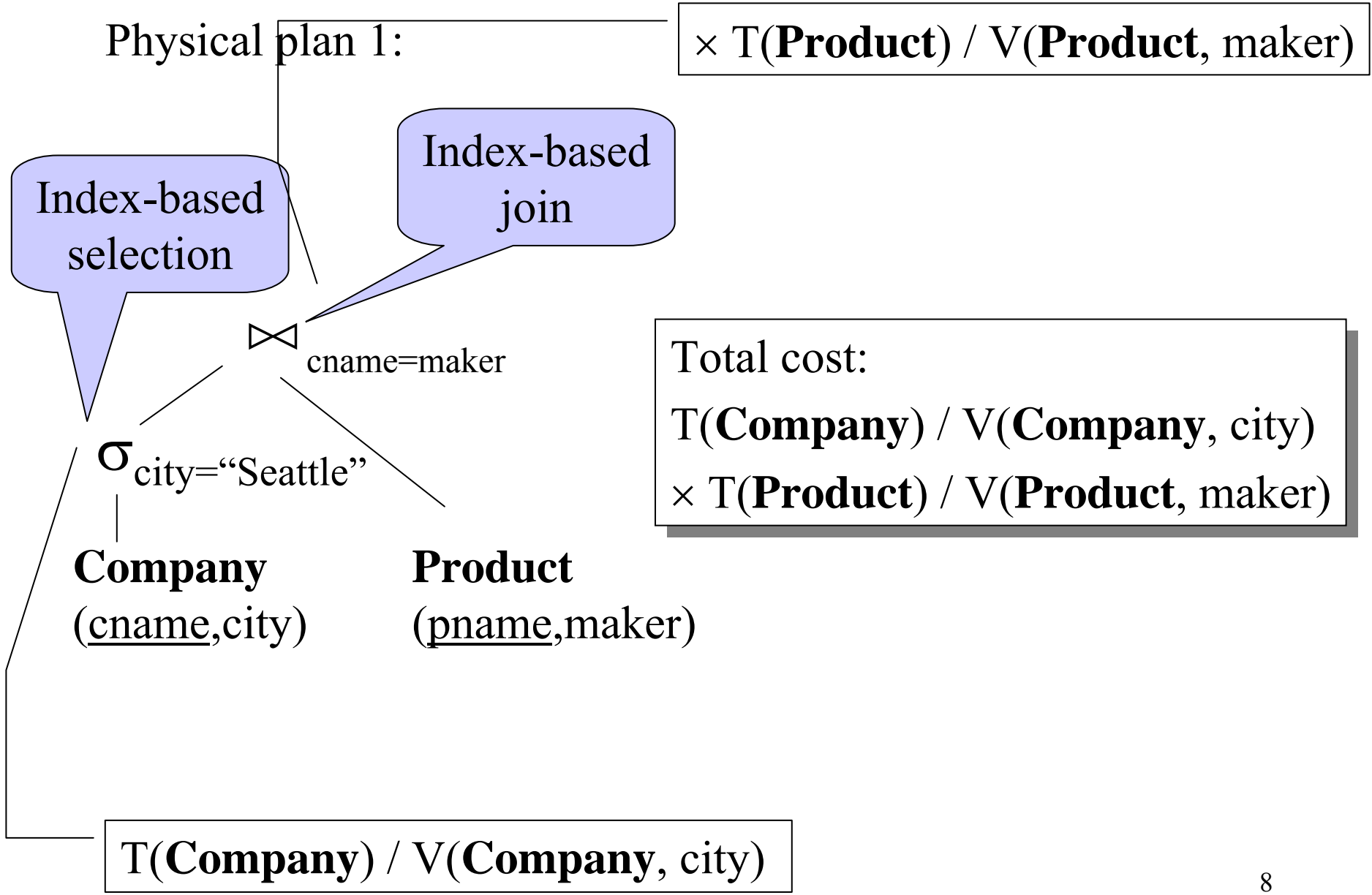


Physical plan 1:



Physical plans 2a and 2b:



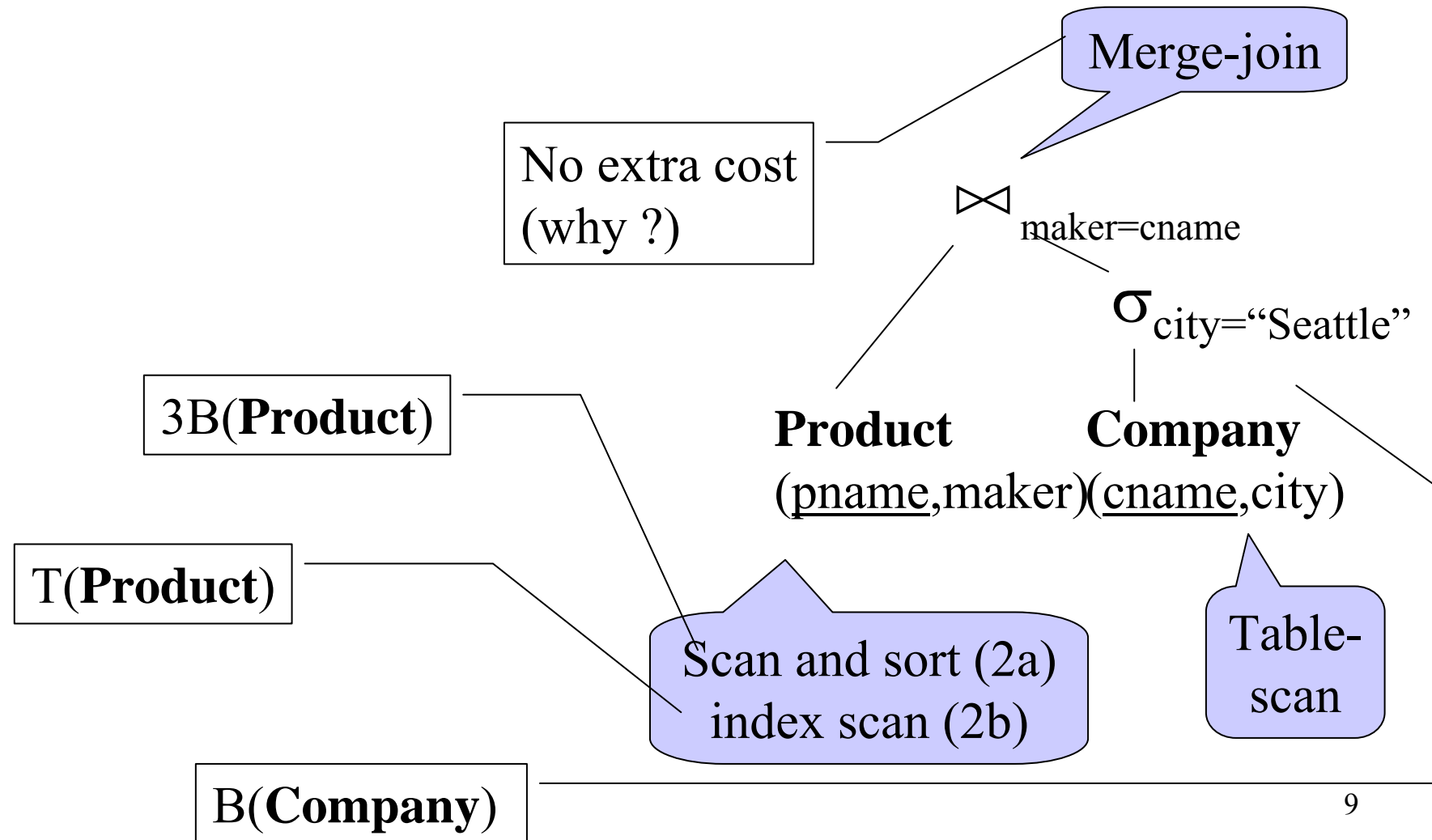


Total cost:

(2a): $3B(\text{Product}) + B(\text{Company})$

(2b): $T(\text{Product}) + B(\text{Company})$

Physical plans 2a and 2b:



Plan 1: $T(\mathbf{Company})/V(\mathbf{Company},\text{city}) \times$
 $T(\mathbf{Product})/V(\mathbf{Product},\text{maker})$

Plan 2a: $B(\mathbf{Company}) + 3B(\mathbf{Product})$

Plan 2b: $B(\mathbf{Company}) + T(\mathbf{Product})$

Which one is better ??

It depends on the data !!

Example

$$T(\mathbf{Company}) = 5,000 \quad B(\mathbf{Company}) = 500 \quad M = 100$$

$$T(\mathbf{Product}) = 100,000 \quad B(\mathbf{Product}) = 1,000$$

We may assume $V(\mathbf{Product}, \text{maker}) \approx T(\mathbf{Company})$ (why ?)

- Case 1: $V(\mathbf{Company}, \text{city}) \approx T(\mathbf{Company})$

$$V(\mathbf{Company}, \text{city}) = 2,000$$

- Case 2: $V(\mathbf{Company}, \text{city}) \ll T(\mathbf{Company})$

$$V(\mathbf{Company}, \text{city}) = 20$$

Which Plan is Best ?

Plan 1: $T(\mathbf{Company})/V(\mathbf{Company},\text{city}) \times T(\mathbf{Product})/V(\mathbf{Product},\text{maker})$

Plan 2a: $B(\mathbf{Company}) + 3B(\mathbf{Product})$

Plan 2b: $B(\mathbf{Company}) + T(\mathbf{Product})$

Case 1:

Case 2:

Lessons

- Need to consider several physical plan
 - even for one, simple logical plan
- No magic “best” plan: depends on the data
- In order to make the right choice
 - need to have *statistics* over the data
 - the B’s, the T’s, the V’s

Query Optimization

- Have a SQL query Q
- Create a plan P
- Find equivalent plans $P = P' = P'' = \dots$
- Choose the “cheapest”.

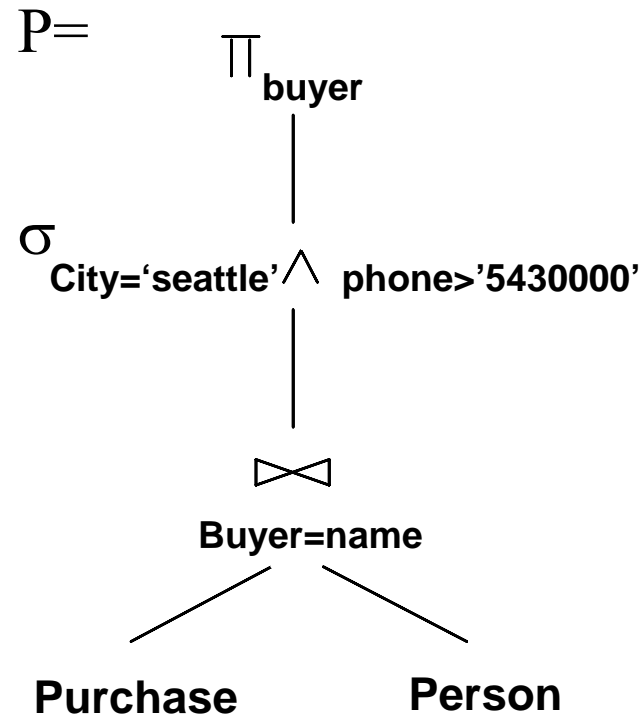


HOW ??

Logical Query Plan

```
SELECT P.buyer
FROM Purchase P, Person Q
WHERE P.buyer=Q.name AND
      P.city='seattle' AND
      Q.phone > '5430000'
```

Purchase(buyer, city)
Person(name, phone)



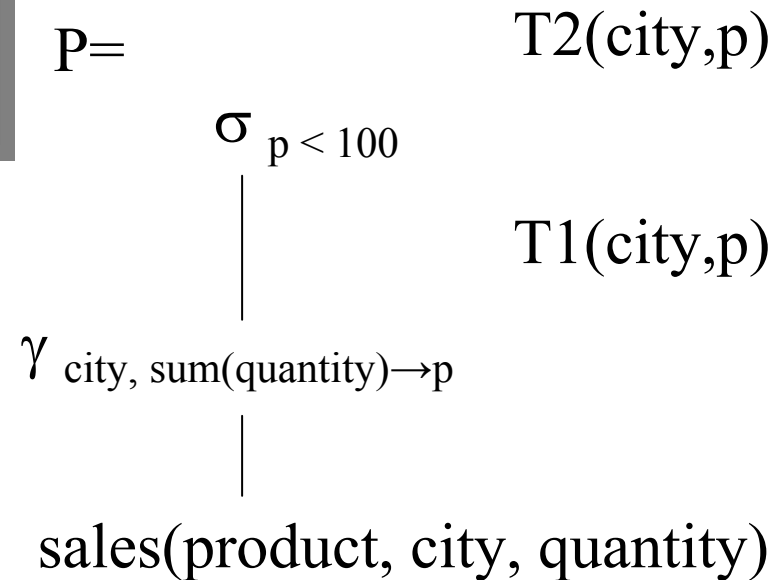
In class:
find a “better” plan P’

Logical Query Plan

Q=

```
SELECT city, sum(quantity)
FROM sales
GROUP BY city
HAVING sum(quantity) < 100
```

P=



In class:
find a “better” plan P’

The three components of an optimizer

We need three things in an optimizer:

- Algebraic laws
- An optimization algorithm
- A cost estimator

Algebraic Laws

- Commutative and Associative Laws

$$R \cup S = S \cup R, \quad R \cup (S \cup T) = (R \cup S) \cup T$$

$$R \times S = S \times R, \quad R \times (S \times T) = (R \times S) \times T$$

$$R \times S = S \times R, \quad R \times (S \times T) = (R \times S) \times T$$

- Distributive Laws

$$R \times (S \cup T) = (R \times S) \cup (R \times T)$$

Algebraic Laws

- Laws involving selection:

$$\sigma_{C \text{ AND } C'}(R) = \sigma_C(\sigma_{C'}(R)) = \sigma_C(R) \cap \sigma_{C'}(R)$$

$$\sigma_{C \text{ OR } C'}(R) = \sigma_C(R) \cup \sigma_{C'}(R)$$

$$\sigma_C(R \times S) = \sigma_C(R) \times S$$

- When C involves only attributes of R

$$\sigma_C(R - S) = \sigma_C(R) - S$$

$$\sigma_C(R \cup S) = \sigma_C(R) \cup \sigma_C(S)$$

$$\sigma_C(R \times S) = \sigma_C(R) \times S$$

Algebraic Laws

- Example: $R(A, B, C, D), S(E, F, G)$

$$\sigma_{F=3} (R \bowtie_{D=E} S) = \quad ?$$

$$\sigma_{A=5 \text{ AND } G=9} (R \bowtie_{D=E} S) = \quad ?$$

Algebraic Laws

- Laws involving projections

$$\Pi_M(R \mid \times \mid S) = \Pi_M(\Pi_P(R) \mid \times \mid \Pi_Q(S))$$

$$\Pi_M(\Pi_N(R)) = \Pi_{M,N}(R)$$

- Example $R(A,B,C,D)$, $S(E, F, G)$

$$\Pi_{A,B,G}(R \mid \times \mid_{D=E} S) = \Pi_{?}(\Pi_{?}(R) \mid \times \mid_{D=E} \Pi_{?}(S))$$

Algebraic Laws

- Laws involving grouping and aggregation:

$$\delta(\gamma_{A, \text{agg}(B)}(R)) = \gamma_{A, \text{agg}(B)}(R)$$

$$\gamma_{A, \text{agg}(B)}(\delta(R)) = \gamma_{A, \text{agg}(B)}(R) \text{ if agg is "duplicate insensitive"}$$

- Which of the following are “duplicate insensitive” ?
sum, count, avg, min, max

$$\gamma_{A, \text{agg}(D)}(R(A,B) \mid \times \mid_{B=C} S(C,D)) = \\ \gamma_{A, \text{agg}(D)}(R(A,B) \mid \times \mid_{B=C} (\gamma_{C, \text{agg}(D)} S(C,D)))$$

Optimizations Based on Semijoins

THIS IS ADVANCED STUFF; NOT ON
THE FINAL

- $R \bowtie S = \Pi_{A_1, \dots, A_n} (R \bowtie S)$
- Where the schemas are:
 - Input: $R(A_1, \dots, A_n), S(B_1, \dots, B_m)$
 - Output: $T(A_1, \dots, A_n)$

Optimizations Based on Semijoins

Semijoins: a bit of theory (see [AHV])

- Given a query:

$$Q :- \Pi (\sigma (R_1 \bowtie R_2 \bowtie \dots \bowtie R_n))$$

- A full reducer for Q is a program:

$$\begin{array}{l} R_{i1} := R_{i1} \bowtie R_{j1} \\ R_{i2} := R_{i2} \bowtie R_{j2} \\ \dots \\ R_{ip} := R_{ip} \bowtie R_{jp} \end{array}$$

- Such that no dangling tuples remain in any relation

Optimizations Based on Semijoins

- Example: $Q(A,E) :- R1(A,B) \bowtie R2(B,C) \bowtie R3(C,D,E)$

- A full reducer is:
 $R2(B,C) := R2(B,C) \bowtie R1(A,B)$
 $R3(C,D,E) := R3(C,D,E) \bowtie R2(B,C)$
 $R2(B,C) := R2(B,C) \bowtie R3(C,D,E)$
 $R1(A,B) := R1(A,B) \bowtie R2(B,C)$

The new tables have only the tuples necessary to compute $Q(E)$

Optimizations Based on Semijoins

- Example:

$$Q(E) :- R1(A,B) \bowtie R2(B,C) \bowtie R3(A,C, E)$$

- Doesn't have a full reducer (we can reduce forever)

Theorem a query has a full reducer iff it is “acyclic”

Optimizations Based on Semijoins

- Semijoins in [Chaudhuri'98]

```
CREATE VIEW DepAvgSal As (  
    SELECT E.did, Avg(E.Sal) AS avgsal  
    FROM Emp E  
    GROUP BY E.did)  
  
SELECT E.eid, E.sal  
FROM Emp E, Dept D, DepAvgSal V  
WHERE E.did = D.did AND E.did = V.did  
    AND E.age < 30 AND D.budget > 100k  
    AND E.sal > V.avgsal
```

Optimizations Based on Semijoins

- First idea:

```
CREATE VIEW LimitedAvgSal As (  
    SELECT E.did, Avg(E.Sal) AS avgsal  
    FROM Emp E, Dept D  
    WHERE E.did = D.did AND D.budget > 100k  
    GROUP BY E.did)
```

```
SELECT E.eid, E.sal  
FROM Emp E, Dept D, LimitedAvgSal V  
WHERE E.did = D.did AND E.did = V.did  
    AND E.age < 30 AND D.budget > 100k  
    AND E.sal > V.avgsal
```

Optimizations Based on Semijoins

- Better: full reducer

```
CREATE VIEW PartialResult AS
  (SELECT E.id, E.sal, E.did
   FROM Emp E, Dept D
   WHERE E.did=D.did AND E.age < 30
   AND D.budget > 100k)
```

```
CREATE VIEW Filter AS
  (SELECT DISTINCT P.did FROM PartialResult P)
```

```
CREATE VIEW LimitedAvgSal AS
  (SELECT E.did, Avg(E.Sal) AS avgsal
   FROM Emp E, Filter F
   WHERE E.did = F.did GROUP BY E.did)
```

Optimizations Based on Semijoins

```
SELECT P.eid, P.sal  
FROM PartialResult P, LimitedDepAvgSal V  
WHERE P.did = V.did AND P.sal > V.avgsal
```