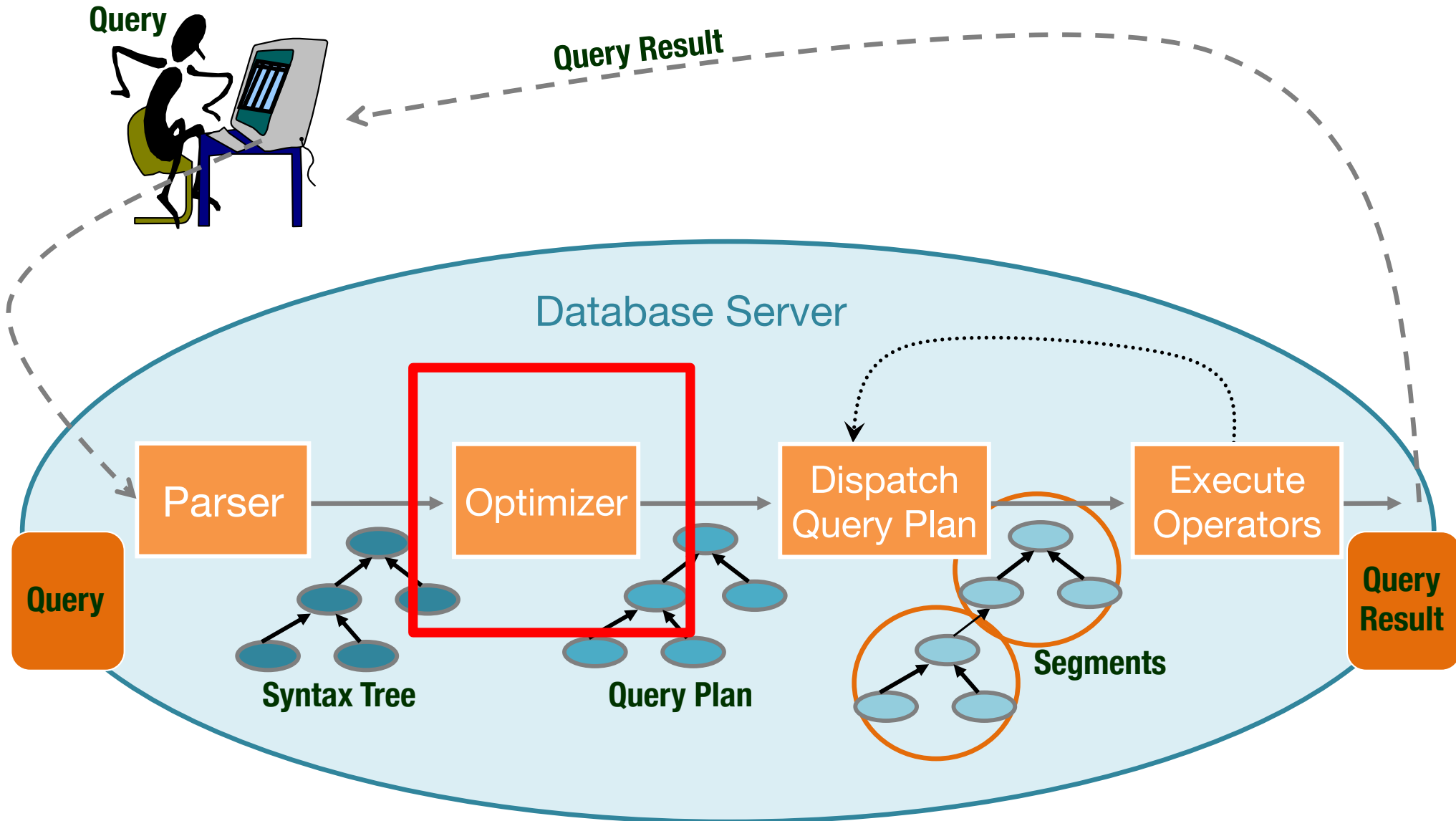


# Query Optimization

## Chapter 15

# Query Execution Life-Cycle



# Query Optimization

Query optimizer selects an evaluation plan with the least cost in two steps:

## 1. Plan Enumeration

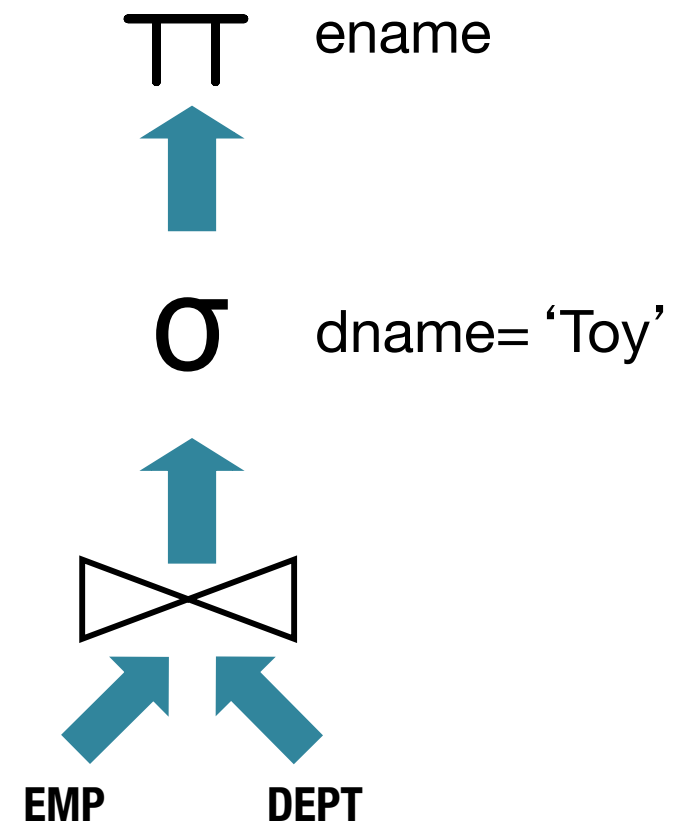
- Today →
- Different query plans
  - Different implementations (i.e., evaluation algorithms) for each operator
2. Cost Estimation
- Cost of each operator
  - Overall cost of the plan
- Previous lectures
- 
- ```
graph LR; Today[Today] --> P1[1. Plan Enumeration]; P1 --> P2[2. Cost Estimation]; PrevLec[Previous lectures] --> P1_2[• Different implementations (i.e., evaluation algorithms) for each operator]; PrevLec --> P2_1[• Cost of each operator];
```

# Example: RA Tree

EMP (ssn, ename, addr, sal, did)

DEPT (did, dname, floor, mgr)

```
SELECT DISTINCT E.ename
FROM Emp E, Dept D
WHERE D.dname = 'Toy'
AND D.did = E.did
```

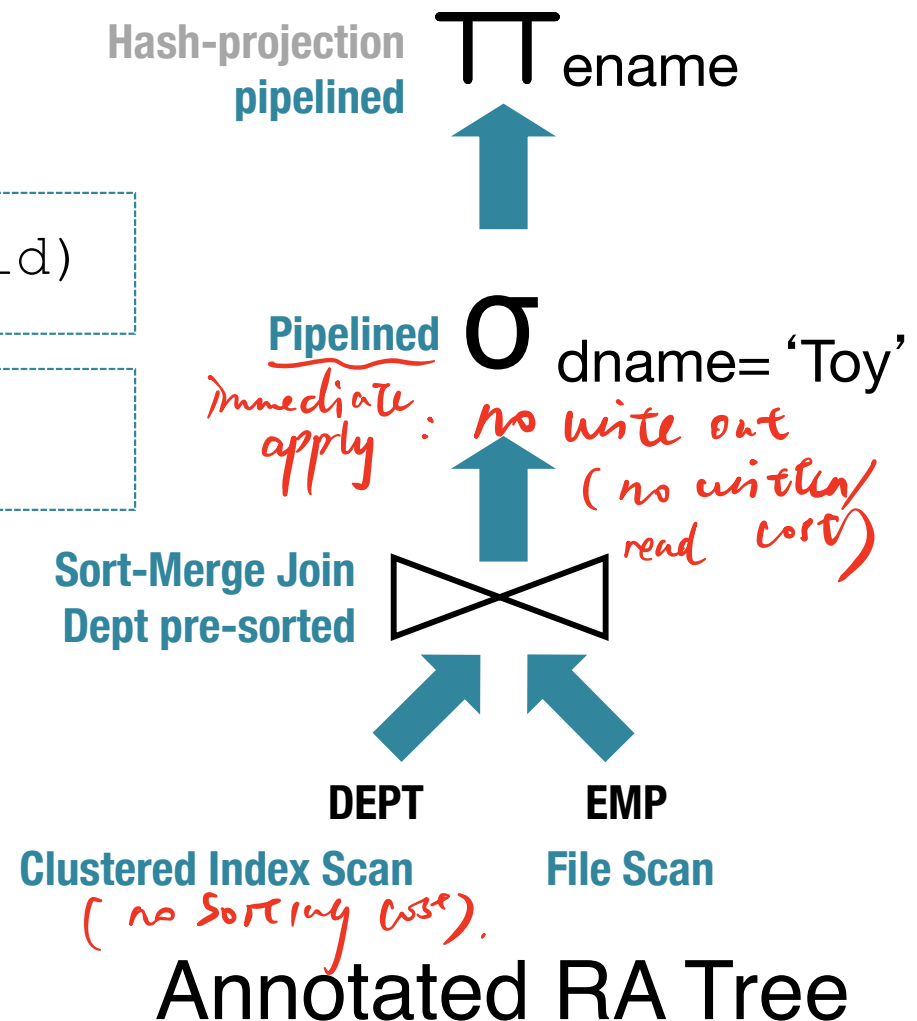


# Example: Annotated RA Tree

EMP (ssn, ename, addr, sal, did)

DEPT (did, dname, floor, mgr)

```
SELECT DISTINCT E.ename
FROM Emp E, Dept D
WHERE D.dname = 'Toy'
AND D.did = E.did
```



# Annotated RA Expressions

- Which algorithm to use for each operator
- Where Intermediate Results are:
  - ① **Pipelined:** Tuples resulting from one operator fed directly into the next
  - ② **Materialized:** Create a temporary table to store intermediate results

# RA Equivalence - Selections

$$\sigma_{P_1} (\sigma_{P_2} (R)) \equiv \sigma_{P_2} (\sigma_{P_1} (R)) \quad (\sigma \text{ commutativity})$$

$$\sigma_{P_1 \wedge P_2 \dots \wedge P_n} (R) \equiv \sigma_{P_1} (\sigma_{P_2} (\dots \sigma_{P_n} (R))) \quad (\text{cascading } \sigma)$$

Selection operation is commutative and multiple selections on the same relation can be combined into a single selection

# RA Equivalence – Projections

$$\Pi_{a_1}(R) \equiv \Pi_{a_1}(\Pi_{a_2}(\dots \Pi_{a_k}(R) \dots)), \quad a_i \subseteq a_{i+1} \text{ (cascading } \Pi)$$

*previous proj. must be superset of following projections*

- In the above, each  $a_i$  is a set of attributes.
- For example: Let's say  $R$  has attributes  $c1, c2, c3, c4$

$$\Pi_{c1, c3}(R) \equiv (\Pi_{c1, c3}(\Pi_{c1, c2, c3}(R)))$$

- Basically, we are eliminating one attribute at a time
- The above is useful when optimizing expressions that involve both projections and joins



# RA Equivalence: Cross-Products & Joins

$$R \bowtie S \equiv S \bowtie R \text{ (commutativity)}$$

$$R \bowtie (S \bowtie T) \equiv (R \bowtie S) \bowtie T \text{ (associativity)}$$

$$R \times S \equiv S \times R \text{ (commutativity)}$$

$$R \times (S \times T) \equiv (R \times S) \times T \text{ (associativity)}$$

- $\Rightarrow$  joins and cross products can be performed in any order  
(SQL can reorder the columns at final presentation time)

# Question?

For particular tables R and S in a database instance, a query optimizer has determined that the join method chosen to compute  $R \bowtie S$  should use R as outer and S as inner.

If asked to compute  $S \bowtie R$ , this query optimizer will likely choose:

- A. S as outer and R as inner
- B. S as inner and R as outer ✓
- C. The choice does not matter.

# Doing projections and selections first

$$\Pi_A(\sigma_c(R)) \equiv \sigma_c(\Pi_A(R)) \quad (\text{if } c \text{ only uses attr in } A)$$

$$\sigma_P(R \times S) \equiv (R \bowtie_P S)$$

$$\sigma_P(R \times S) \equiv \sigma_P(R) \times S \quad (\text{if } p \text{ is only on } R)$$

*P only applies on R*

$$\sigma_P(R \bowtie S) \equiv \sigma_P(R) \bowtie S \quad (\text{if } p \text{ is only on } R)$$

**Key ideas:** When possible, consider doing:

- Joins rather than cross-products
- Selections before joins
- (Partial) Projections before all operators.

# Doing selections before join

$$\sigma_P (R \times S) \equiv \sigma_{P_4} (\sigma_{P_1} (R) \bowtie_{p_3} \sigma_{P_2} (S))$$

Note:  $P = p_1 \wedge p_2 \wedge p_3 \wedge p_4$

**Idea:** Replace  $P$  by a cascade of conditions  $p_1$ ,  $p_2$ ,  $p_3$ , and  $p_4$  such that:

$P_1$ : conditions only on  $R$  (e.g.,  $R.a = 5$ )

$P_2$ : conditions only on  $S$  (e.g.,  $S.d > 9$ )

$P_3$ : join conditions with equality on  $R$  and  $S$  ( $R.a = S.a$ )

$P_4$ : other conditions involving both  $R$  and  $S$  columns  
(e.g.,  $R.a > S.b$ )

# Example

$$\sigma_P (R \times S) \equiv \sigma_{p_4} (\sigma_{p_1} (R) \bowtie_{p_3} \sigma_{p_2} (S))$$

Note:  $P = p_1 \wedge p_2 \wedge p_3 \wedge p_4$

$$\sigma_{\text{major}=\text{CS} \wedge R.\text{umid}=S.\text{umid} \wedge \text{grade}>\text{C}} (R \times S) \equiv$$

$$\sigma_{\text{true}} (\sigma_{\text{major}=\text{CS}} (R) \bowtie_{\text{umid}} \sigma_{\text{grade}>\text{C}} (S))$$

| umid | sname | major |
|------|-------|-------|
| 23   | Alice | CS    |
| 71   | Bob   | CE    |
| 11   | Mary  | CS    |
| 13   | John  | CS    |

| umid | course | semes | grade |
|------|--------|-------|-------|
| 13   | 484    | W17   | C+    |
| 23   | 484    | W17   | C+    |
| 71   | 484    | W17   | C+    |
| 23   | 482    | W17   | C     |
| 71   | 482    | W17   | D     |
| 11   | 482    | W17   | D-    |

# RA Equivalence – Multiple Ops

$$\Pi_P (R \times S) \equiv \Pi_{P_1}(R) \times \Pi_{P_2}(S)$$

Columns  $p$  in the cross-product consist of columns  $p_1$  from  $R$  and columns  $p_2$  from  $S$

$$\Pi_P (R \bowtie_C S) \equiv \Pi_P (\Pi_{P_1}(R) \bowtie_C \Pi_{P_2}(S))$$

*\* for join & projection*  
 $p_1$  are attrs of  $R$  that appear in  $p$  or  $c$   
 $p_2$  are attrs of  $S$  that appear in  $p$  or  $c$

# Example

$$\Pi_P (R \bowtie_C S) \equiv \Pi_P (\Pi_{P_1}(R) \bowtie_C \Pi_{P_2}(S))$$

$$\begin{aligned} & \Pi_{\text{sname, course}} (R \bowtie_{\text{umid}} S) \\ & \equiv \Pi_{\text{sname, course}} (\Pi_{\text{sname, umid}}(R) \bowtie_{\text{umid}} \Pi_{\text{course, umid}}(S)) \end{aligned}$$

| umid | sname | major |
|------|-------|-------|
| 23   | Alice | CS    |
| 71   | Bob   | CE    |
| 11   | Mary  | CS    |
| 13   | John  | DS    |

| umid | course | semes | grade |
|------|--------|-------|-------|
| 13   | 484    | W17   | C+    |
| 23   | 484    | W17   | C+    |
| 71   | 484    | W17   | C+    |
| 23   | 482    | W17   | C     |
| 71   | 482    | W17   | D     |
| 11   | 482    | W17   | D-    |

# RA Equivalence – More Rules

$$\Pi_{A1, A2, \dots, An} (\sigma_P (R)) \equiv \Pi_{A1, A2, \dots, An} (\sigma_P (\Pi_{A1, \dots, An, B1, \dots, BM} R))$$

**B1 ... BM attributes in P**

*need for  
selection &  
projection*

$$\Pi_{umid} (\sigma_{major=CE} (R)) \equiv \Pi_{umid} (\sigma_{major=CE} (\Pi_{umid, major} R))$$

| umid | sname | major |
|------|-------|-------|
| 23   | Alice | CS    |
| 71   | Bob   | CE    |
| 11   | Carl  | CS    |
| 13   | Debra | DS    |



# Optimization Strategies

- Rule-Based
  - Choose more likely cheaper plan based on RA tree
  - E.g. Push selections in.
  - E.g. Push projections in.
- Cost-Based
  - Estimate cost and choose lower-cost plan
  - E.g. choosing between join methods.

# Cost Based Optimization

- Traditional techniques work well for  $< 10$  joins
- **Cost estimation:** Approximate at best
  - Use statistics from systems catalogs
  - Combination of CPU and I/O costs
    - EECS 484 focuses on I/O costs only
    - We also use **System R** approach
      - very inexact but ok in practice (better techniques are known now)

# Estimating the Cost of a Plan

## 1. Estimate *cost* of each operation in plan tree

- Depends on input cardinalities (# of rows)
- Algorithm cost (see previous lectures)

## 2. Estimate *size of result*

- Use information about the input relations
- For selections and joins, **assume independence of predicates**

# Collecting & Maintaining Statistics

- Statistics stored in the catalogs (pg\_stats & pg\_class in Postgres)
  - Relation
    - Cardinality (# rows)
    - Size in pages
  - Index or attribute
    - ① Cardinality (# distinct key values)
    - ② For index: Size in pages and Height (of b-tree)
    - ③ Range of values; Possibly histogram as well e.g. employee age 18-75
- Catalogs **update periodically**
  - Can be slightly inconsistent/inaccurate

# Estimating Output Size

```
SELECT attribute list  
FROM relation list  
WHERE term1 AND ... AND termk
```

Question: What is the cardinality of the result set?

- Max # tuples: product of input relation cardinalities
- Each term “filters” out some tuples: Reduction factor
- Result cardinality = Max # tuples \* product of all RF's
- Assumption: terms are independent!
- Term  $col=value$       RF:  $1/Range(I)$
- Term  $col>value$       RF:  $(High(I)-value)/(High(I)-Low(I))$
- Term  $col1=col2$       RF:  $1/MAX(Range(I1), Range(I2))$

# Cost of different plans



EMP (ssn, ename, addr, sal, did)

```
SELECT E.ename
FROM   Emp E
WHERE  E.did=8
AND    E.sal > 40K
```

Estimate the cost of different access methods:

- Index on did:
  - Clustered index:
  - Unclustered index:
- Index on sal:
  - Clustered index:
  - Unclustered index:
- File scan: 1,000 pages

1,000 data pages, 10K tuples  
100 pages in B+-tree  
# depts: 10  
Salary Range: 0K – 200K

# Cost of different plans

EMP (ssn, ename, addr, sal, did)

```
SELECT E.ename
FROM   Emp E
WHERE  E.did=8
AND    E.sal > 40K
```

- Index on did:

- Tuples Retrieved:  $(1/10) * 10,000$

- Clustered index:  
*leaf pages* *data pages*  
 $(1/10) * (100 + 1,000)$  pages of I/O  
= 110 I/Os

- Unclustered index:

$$(1/10) * (100 + 10,000) \text{ pages I/O} \\ = 1010 \text{ I/Os}$$

*leaf* *each tuple cost an 10 for the unclustered index*

Given:

- 1,000 data pages, 10K tuples
- 100 leaf pages in B+-tree
- All tree levels except leaves cached in memory
- # depts: 10
- Salary Range: 0K – 200K

# Question?

EMP (ssn, ename, addr, sal, did)

- Index on did:
  - Index on sal:
    - Tuples Retrieved = ??
- A. 1000
- B. 4000
- C. 8000 ✓
- D. 10,000

```
SELECT E.ename
FROM   Emp E
WHERE  E.did=8
AND    E.sal > 40K
```

80%

Given:

- 1,000 data pages, 10K tuples
- 100 leaf pages in B+-tree
- All tree levels except leaves cached in memory
- # depts: 10
- Salary Range: 0K – 200K



# Cost of different plans



EMP (ssn, ename, addr, sal, did)

- Index on did:
  - Tuples Retrieved:  $(1/10) * 10,000$
  - Clustered index:  $(1/10) * (100+1,000) = 110$  pages
  - Unclustered index:  $(1/10) * (100+10,000) = 1,010$  pages
- Index on sal:
  - Clustered index:  $(200-40)/(200-0) * (100+1,000) = 880$  pages
  - Unclustered index:  $(200-40)/(200-0) * (100+10,000) = 8,080$  pages
- File scan: 1,000 pages

```
SELECT E.ename
FROM   Emp E
WHERE  E.did=8
AND    E.sal > 40K
```

1,000 data pages, 10K tuples  
100 pages in B+-tree  
# depts: 10  
Salary Range: 0K – 200K

# Estimating Output Size

```
SELECT attribute list  
FROM relation list  
WHERE term1 AND ... AND termk
```

Question: What is the cardinality of the result set?

- Max # tuples: product of input relation cardinalities
- Each term “filters” out some tuples: Reduction factor
- **Result cardinality = Max # tuples \* product of all RF's**
- Assumption: terms are independent!
- Term  $col=value$       RF:  $1/Range(I)$
- Term  $col>value$       RF:  $(High(I)-value)/(High(I)-Low(I))$
- Term  $col1=col2$       RF:  $1/MAX(Range(I1), Range(I2))$

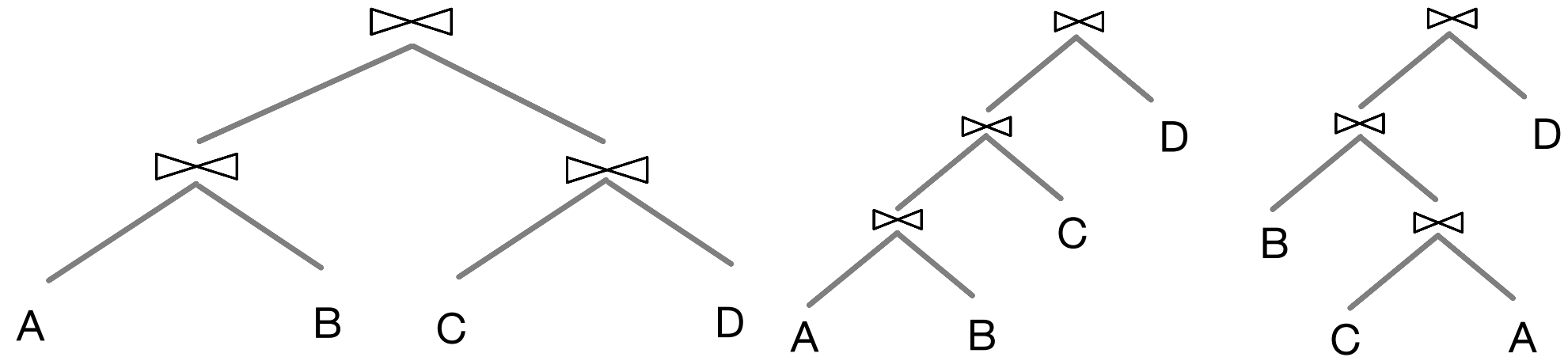
# Example: Estimate size of join

- Given:
  - Parts Relation containing 1 million distinct values of (the primary key) attribute PartID *unique parts in Parts  $\Rightarrow$  1 million*
  - Orders relation comprising 2000 different orders of 1000 distinct Parts. (Some parts have more than one order). *unique parts in Orders  $\Rightarrow$  1000*
- Q: What is estimated # of tuples in natural join between Parts and Orders on the attribute PartID?
- Solution:**  *$\max(1000, 1m) = 1m$* 
  - Size of cross-product: 2000 x 1 million tuples. (Call this A)
  - RF of [Orders.PartID = Parts.PartID selection predicate] =  $\max(1 \text{ million}, 1000) = 1 \text{ million}$  (call this B)
  - Estimated size of natural join =  $A/B = \mathbf{2000 \text{ tuples.}}$  (Answer)

$$A. \frac{1}{1m} = 2000 \text{ tuples}$$

# Optimizing Join Order

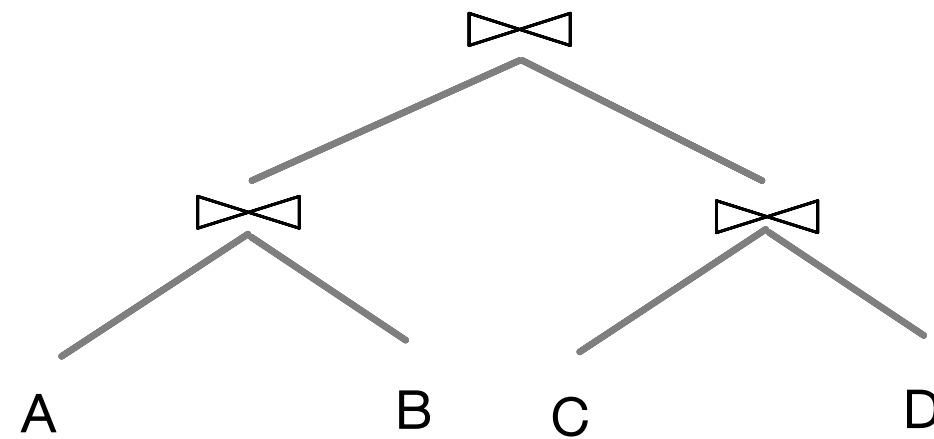
- Consider (A join B join C join D).
- The query planner has many ways it can accomplish the join.
- Given n joins in the expression, way more than n! choices!



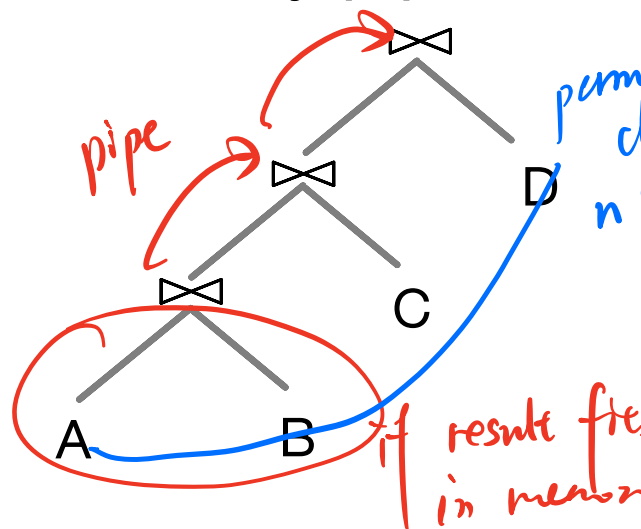
(Here, C and A would be joined if they share an attribute. Else we would do a cross product)

# Optimizing Join Order

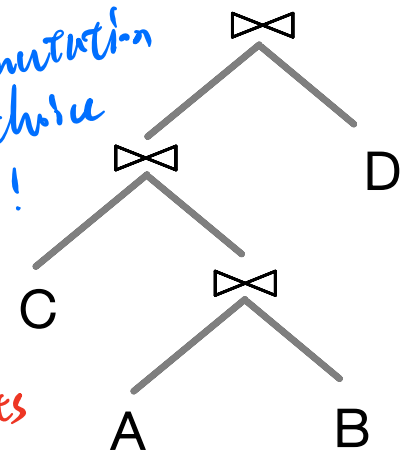
- System R : Only consider **left-deep** join trees
    - Used to restrict the search space
    - Left-deep plans can be **fully pipelined** (usually)
      - Intermediate results not written to temporary files
      - Not all left-deep trees are fully pipelined (e.g., SM join)
- only left child is allowed to have children*



**Not left-deep join tree**



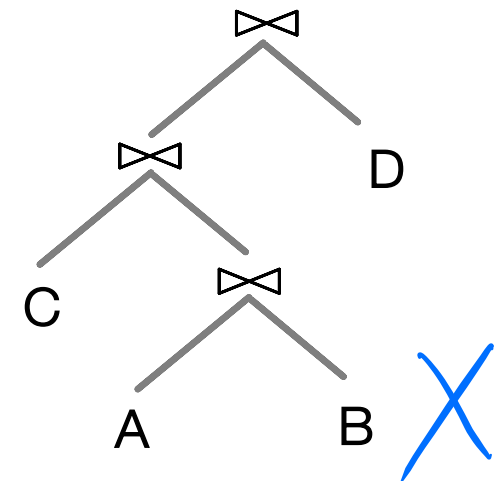
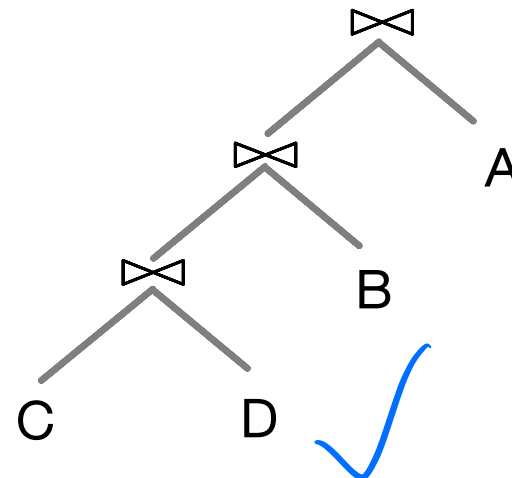
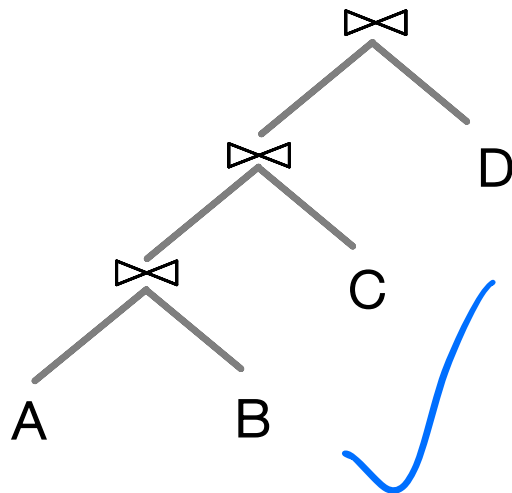
**Left-deep join tree**



**Not left-deep join tree**

# Example


Which one (if any) is a potential plan that System R will consider for joining A, B, C, and D?



# Finding Best Left-Deep Plan

- Decide:
  - Join order
  - Join method for each join
- Enumerated using N passes (if N relations joined):
- **Pass 1:** Find best 1-relation plan for each relation (apply selections and projections first and consider using indexing)
- For each relation, retain only:
  - Cheapest plan overall (e.g. File scan), plus
  - Cheapest plans that produce **ordered** tuples. Order may be useful for later stages (e.g., sort-merge-join)

# Enumeration of Left-Deep Plans

Pass 2: Find best way to join result of each 1-relation plan (as outer) to another relation (All 2-relation plans) 

- For each pair (subset) of relations, retain only:
  - Cheapest plan overall, plus
  - Cheapest plan for sorted **order** of tuples, if order is possibly useful for later operations.
- Also, apply selections and projections aggressively, when possible, to reduce result size
- Assume pipelining of results to avoid additional I/O



# Enumeration of Left-Deep Plans

**Pass N:** Find best way to join result of retained (N-1)-relation plans (as outer) to the N<sup>th</sup> relation

e.g. Sailor Reserves Boats  
aim:

~~S~~ ~~R~~ ~~B~~

For each subset of relations, retain only:

- Cheapest plan overall, plus *find these to do* ① S, R, B *(no connection)*
- Cheapest plan for each interesting order of the tuples ② ~~S~~ ~~R~~, ~~R~~ ~~B~~ *no S ~~B~~*
- Only consider joining relations if there is a connecting join condition, i.e., **avoid Cartesian products if possible**
- ORDER BY, GROUP BY handled as a final step

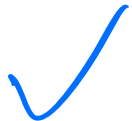
# Question?

Given  $N$  relations, the number of possible left-deep plans is close to:

A.  $N$

B.  $N^2$

C.  $N!$



D.  $\binom{N}{2} = N \text{ choose } 2$

# Summary

- Query optimization critical to the DBMS performance
  - Helps understand performance impact of database design
- Two parts to optimizing a query:
  - Enumerate alternative plans. (Typically only consider left-deep plans)
  - Estimate cost of each plan: size of result and cost of algorithm

# Optional Exercises

12.1 (all parts), 15.1, 15.5, 15.7, 15.9