# Query Processing Notes
## CSE232

## Query Processing

- The query processor turns user queries and data modification commands into a query plan - a sequence of operations (or algorithm) on the database
  - from high level queries to low level commands
- Decisions taken by the query processor
  - Which of the algebraically equivalent forms of a query will lead to the most efficient algorithm?
  - For each algebraic operator what algorithm should we use to run the operator?
  - How should the operators pass data from one to the other? (eg, main memory buffers, disk buffers)

Example

Select B,D
From R,S
Where R.A = "c" $\wedge$ S.E = 2 $\wedge$ R.C=S.C

| R | A | B | C | | S | C | D | E |
|---|---|---|---|---|---|---|---|---|
| | a | 1 | 10 | | | 10 | x | 2 |
| | b | 1 | 20 | | | 20 | y | 2 |
| | c | 2 | 10 | | | 30 | z | 2 |
| | d | 2 | 35 | | | 40 | x | 1 |
| | e | 3 | 45 | | | 50 | y | 3 |

Answer

| B | D |
|---|---|
| 2 | x |

• How do we execute query eventually?

One idea

- Scan relations
- Do Cartesian product
- Select tuples
- Do projection

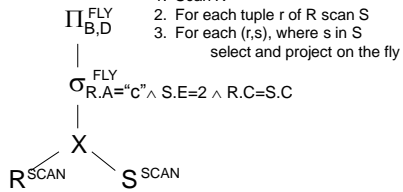| RxS | R.A | R.B | R.C | S.C | S.D | S.E |
|-----|-----|-----|-----|-----|-----|-----|
| | a | 1 | 10 | 10 | x | 2 |
| | a | 1 | 10 | 20 | y | 2 |
| | . | | | | | |
| | . | | | | | |
| Bingo! → Got one... | C | 2 | 10 | 10 | x | 2 |
| | . | | | | | |
| | . | | | | | |

<u>Relational Algebra</u> - can be
enhanced to describe plans...

<u>Ex:</u> Plan I

$$\Pi_{B,D}^{FLY}$$

1. Scan R
2. For each tuple r of R scan S
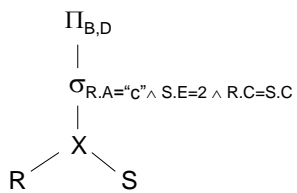3. For each (r,s), where s in S
   select and project on the fly

$$\sigma_{R.A="c" \wedge S.E=2 \wedge R.C=S.C}^{FLY}$$

$$X$$

$$R^{SCAN} \quad S^{SCAN}$$

<u>OR:</u> $\Pi_{B,D}^{FLY} [ \sigma_{R.A="c" \wedge S.E=2 \wedge R.C=S.C}^{FLY} (R^{SCAN} X \ S^{SCAN})]$

"FLY" and "SCAN" are the defaults

<u>Ex:</u> Plan I

$$\Pi_{B,D}$$

$$\sigma_{R.A="c" \wedge S.E=2 \wedge R.C=S.C}$$

$$X$$

$$R \quad S$$

Another idea:

Plan II

$$\Pi_{B,D}$$

$$\bowtie^{HASH}$$

$\bowtie$ natural join

$$\sigma_{R.A = "c"} \quad \sigma_{S.E = 2}$$

$$R \quad S$$

Scan R and S, perform on the fly selections, do hash join, project

R

| A | B | C |
|---|---|---|
| a | 1 | 10 |
| b | 1 | 20 |
| c | 2 | 10 |
| d | 2 | 35 |
| e | 3 | 45 |

$\sigma(R)$

| A | B | C |
|---|---|---|
| c | 2 | 10 |

$\sigma(S)$

| C | D | E |
|---|---|---|
| 10 | x | 2 |
| 20 | y | 2 |
| 30 | z | 2 |

S

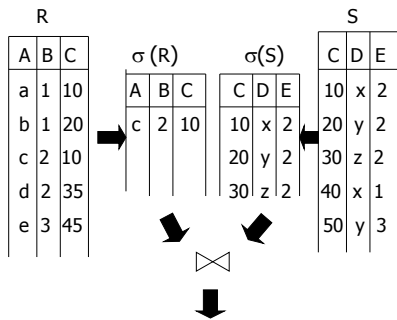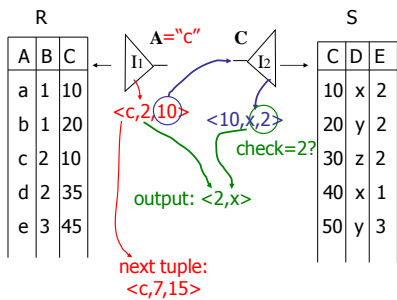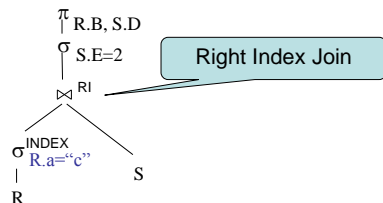| C | D | E |
|---|---|---|
| 10 | x | 2 |
| 20 | y | 2 |
| 30 | z | 2 |
| 40 | x | 1 |
| 50 | y | 3 |

## Plan III

Use R.A and S.C Indexes

(1) Use R.A index to select R tuples with R.A = "c"

(2) For each R.C value found, use S.C index to find matching join tuples

(3) Eliminate join tuples S.E ≠ 2

(4) Project B,D attributes



R

| A | B | C |
|---|---|---|
| a | 1 | 10 |
| b | 1 | 20 |
| c | 2 | 10 |
| d | 2 | 35 |
| e | 3 | 45 |

**A="c"**  **C**
I1 → I2

<c,2,10>  <10,x,2>  check=2?

output: <2,x>

next tuple: <c,7,15>

S

| C | D | E |
|---|---|---|
| 10 | x | 2 |
| 20 | y | 2 |
| 30 | z | 2 |
| 40 | x | 1 |
| 50 | y | 3 |

**Algebraic Form of Plan**

$\pi$ R.B, S.D

$\sigma$ S.E=2

$\bowtie$ RI          Right Index Join

$\sigma^{INDEX}_{R.a="c"}$          S

R

# From Query To Optimal Plan

- Complex process
- Algebra-based logical and physical plans
- Transformations
- Evaluation of multiple alternatives

# Issues in Query Processing and Optimization

- Generate Plans
  - employ efficient execution primitives for computing relational algebra operations
  - systematically transform expressions to achieve more efficient combinations of operators
- Estimate Cost of Generated Plans
  - Statistics
- "Smart" Search of the Space of Possible Plans
  - always do the "good" transformations (relational algebra optimization)
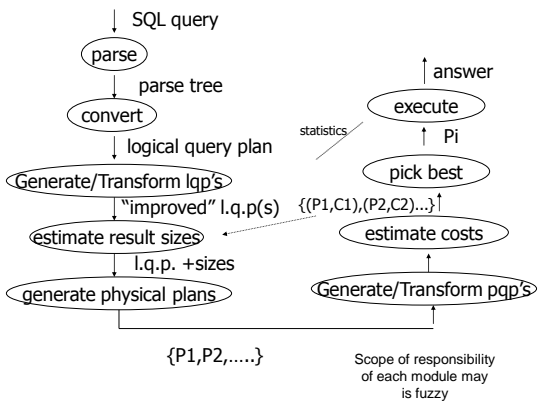  - prune the space (e.g., System R)
- Often the above steps are mixed

The flow diagram (top):
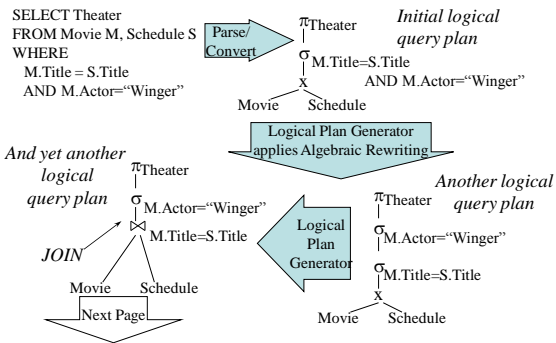
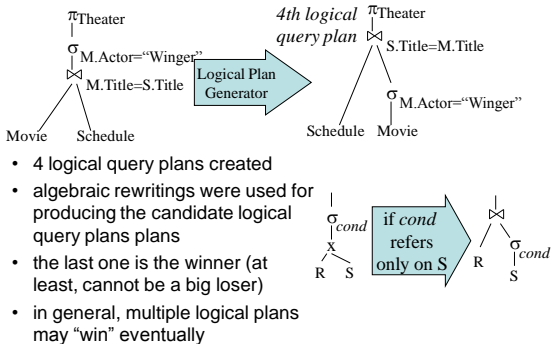SQL query → parse → parse tree → convert → logical query plan → Generate/Transform lqp's → "improved" l.q.p(s) → estimate result sizes → l.q.p. +sizes → generate physical plans → {P1,P2,.....}

estimate costs ← {(P1,C1),(P2,C2)...} ← pick best ← Pi ← execute → answer

Generate/Transform pqp's → estimate costs

statistics

Scope of responsibility of each module may is fuzzy

## Example: The Journey of a Query

SELECT Theater
FROM Movie M, Schedule S
WHERE
    M.Title = S.Title
    AND M.Actor="Winger"

Parse/Convert →

*Initial logical query plan*

$\pi$Theater
$\sigma$M.Title=S.Title AND M.Actor="Winger"
X
Movie    Schedule

Logical Plan Generator applies Algebraic Rewriting ↓

*And yet another logical query plan*

$\pi$Theater
$\sigma$M.Actor="Winger"
⋈ M.Title=S.Title
*JOIN*
Movie    Schedule
Next Page ↓

← Logical Plan Generator

*Another logical query plan*

$\pi$Theater
$\sigma$M.Actor="Winger"
$\sigma$M.Title=S.Title
X
Movie    Schedule

## The Journey of a Query cont'd: Summary of Logical Plan Generator

$\pi$Theater
$\sigma$M.Actor="Winger"
⋈ M.Title=S.Title
Movie    Schedule

Logical Plan Generator →

*4th logical query plan*

$\pi$Theater
⋈ S.Title=M.Title
Schedule
$\sigma$M.Actor="Winger"
Movie

- 4 logical query plans created
- algebraic rewritings were used for producing the candidate logical query plans plans
- the last one is the winner (at least, cannot be a big loser)
- in general, multiple logical plans may "win" eventually

$\sigma_{cond}$
X
R    S

if *cond* refers only on S →

⋈
R    $\sigma_{cond}$
      S

6

## The Journey of a Query Continues at the Physical Plan Generator

$\pi$Theater
$\bowtie$ S.Title=M.Title

$\sigma$ M.Actor="Winger"

Schedule   Movie

Physical Plan Generator

index on Actor, tables Schedule sorted on Title,

$\pi$Theater
$\bowtie$ SORT-MERGE S.Title=M.Title

$\sigma$ INDEX M.Actor="Winger

Schedule   Movie

*Physical Plan 2*

Physical Plan Generators chooses execution primitives and data passing

index on Actor and Title, unsorted tables, tables>>memory

$\pi$Theater
$\bowtie$ LEFT INDEX S.Title=M.Title

$\sigma$ INDEX Actor="Winger"

*Physical Plan 1*

Schedule   Movie

More than one plans may be generated by choosing different primitives

---

## Example:   Nested SQL query

```
SELECT title
FROM StarsIn
WHERE starName IN (
        SELECT name
        FROM MovieStar
        WHERE birthdate LIKE '%1960'
);
```

(Find the movies with stars born in 1960)

---

## Example:   Parse Tree

```
                        <Query>
                         <SFW>
SELECT  <SelList>   FROM   <FromList>   WHERE   <Condition>
          <Attribute>       <RelName>          <Tuple> IN <Query>
            title            StarsIn         <Attribute>  ( <Query> )
                                             starName      <SFW>

SELECT   <SelList>   FROM   <FromList>   WHERE   <Condition>
           <Attribute>       <RelName>        <Attribute>  LIKE  <Pattern>
             name             MovieStar        birthDate         '%1960'
```

<u>Example:</u>   Generating Relational Algebra

$$\Pi_{title}$$

$$\sigma$$

StarsIn          &lt;condition&gt;

&lt;tuple&gt;     IN  $\Pi_{name}$

&lt;attribute&gt;    $\sigma_{birthdate\ LIKE\ '\%1960'}$

starName         MovieStar

An expression using a two-argument $\sigma$, midway
between a parse tree and relational algebra


<u>Example:</u>   Logical Query Plan (Relational
Algebra)

$$\Pi_{title}$$

$$\sigma_{starName=name}$$

$$\times$$

StarsIn     $\Pi_{name}$

$$\sigma_{birthdate\ LIKE\ '\%1960'}$$

MovieStar

May consider "IN" elimination as a rewriting
in the logical plan generator or may consider it
a task of the converter


<u>Example:</u>   Improved Logical Query Plan

$$\Pi_{title}$$

$$\bowtie$$
starName=name

StarsIn     $\Pi_{name}$

$$\sigma_{birthdate\ LIKE\ '\%1960'}$$

MovieStar

Question:
Push project to
StarsIn?

Example: Result sizes are important for selecting physical plans



Need expected size

StarsIn

$\Pi$
$\sigma$
MovieStar

Example: One Physical Plan



$\Pi_{title}$

Additional parameters:
memory size, result sizes...

starName=name

$\Pi_{name}$

StarsIn$^{SCAN}$

$\sigma^{INDEX}_{birthdate \ LIKE \ '\%1960'}$

MovieStar

## Topics

- Bag Algebra, List Algebra and other extensions
  - name & value conversions, functions, aggregation

## Algebraic Operators: A Bag version

- *Union of R and S*: a tuple t is in the result as many times as the sum of the number of times it is in R plus the times it is in S
- *Intersection of R and S*: a tuple t is in the result the minimum of the number of times it is in R and S
- *Difference of R and S*: a tuple t is in the result the number of times it is in R minus the number of times it is in S
- $\delta(R)$ converts the bag $R$ into a set
  - SQL's *R UNION S* is really $\delta(R \cup S)$
- **Example:** Let $R=\{A,B,B\}$ and $S=\{C,A,B,C\}$. Describe the union, intersection and difference...

## Extended Projection

- We extend the relational project $\pi_A$ as follows:
  - The attribute list may include $x \rightarrow y$ in the list $A$ to indicate that the attribute $x$ is renamed to $y$
  - Arithmetic, string operators and scalar functions on attributes are allowed. For example,
    - $a+b \rightarrow x$ means that the sum of $a$ and $b$ is renamed into $x$.
    - $c||d \rightarrow y$ concatenates the result of $c$ and $d$ into a new attribute named $y$
- The result is computed by considering each tuple in turn and constructing a new tuple by picking the attributes names in $A$ and applying renamings and arithmetic and string operators
- **Example:**

## An Alternative Approach to Arithmetic and Other 1-1 Computations

- Special purpose operators that for every input tuple they produce one output tuple
  - *MULT* $_{A,B \rightarrow C}R$: for each tuple of $R$, multiply attribute $A$ with attribute $B$ and put the result in a new attribute named $C$.
  - *PLUS* $_{A,B \rightarrow C}R$
  - *CONCAT* $_{A,B \rightarrow C}R$
- **Exercise:** Write the above operators using extended projection. Assume the schema of $R$ is $R(A,B,D,E)$.

## Products and Joins

- *Product of R and S (R×S)*:
  - If an attribute named *a* is found in both schemas then rename one column into *R.a* and the other into *S.a*
  - If a tuple *r* is found *n* times in *R* and a tuple *s* is found *m* times in *S* then the product contains *nm* instances of the tuple *rs*
- Joins
  - *Natural Join $R \bowtie S = \pi_A \sigma_C (R×S)$* where
    - *C* is a condition that equates all common attributes
    - *A* is the concatenated list of attributes of *R* and *S* with no duplicates
    - you may view tha above as a rewriting rule
  - *Theta Join*
    - arbitrary condition involving multiple attributes

## Grouping and Aggregation

- $\gamma_{GroupByList;\, aggrFn1 \to attr1, \ldots, aggrFnN \to attrN}$
- Conceptually, grouping leads to nested tables and is immediately followed by functions that aggregate the nested table
- Example: $\gamma_{Dept;\, AVG(Salary) \to AvgSal, \ldots, SUM(Salary) \to SalaryExp}$

Find the average salary for each department
SELECT Dept, AVG(Salary) AS AvgSal,
 SUM(Salary) AS SalaryExp
FROM Employee
GROUP-BY Dept

**Employee**

| Name | Dept | Salary |
|------|------|--------|
| Joe  | Toys | 45     |
| Nick | PCs  | 50     |
| Jim  | Toys | 35     |
| Jack | PCs  | 40     |

| Dept | Nested Table Name  Salary | |
|------|------|------|
| Toys | Joe | 45 |
|      | Jim | 35 |
| PCs  | Nick | 50 |
|      | Jack | 40 |

| Dept | AvgSal | SalaryExp |
|------|--------|-----------|
| Toys | 40     | 80        |
| PCs  | 45     | 90        |

## Grouping and Aggregation: An Alternate approach

- Operators that combine the *GROUP-BY* clause with the aggregation operator (*AVG,SUM,MIN,MAX,…*)
- $SUM_{GroupbyList;\, GroupedAttribute \to ResultAttribute}\, R$ corresponds to   SELECT *GroupbyList,*
         SUM(*GroupedAttribute*) AS *ResultAttribute*
      FROM *R*
      GROUP BY *GroupbyList*
- Similar for *AVG,MIN,MAX,COUNT…*
- Note that $\delta(R)$ could be seen as a special case of grouping and aggregation
- **Example**

## Sorting and Lists

- SQL and algebra results are ordered
- Could be non-deterministic or dictated by SQL ORDER BY, algebra τ
- $τ_{OrderByList}$
- A result of an algebraic expression o(exp) is ordered if
  - If o is a τ
  - If o retains ordering of exp and exp is ordered
    - Unfortunately this depends on implementation of o
  - If o creates ordering
  - Consider that leaf of tree may be SCAN(R)

## Relational algebra optimization

- Transformation rules (preserve equivalence)
- What are good transformations?

## Algebraic Rewritings: Commutativity and Associativity



**Question 1**: Do the above hold for both sets and bags?
**Question 2**: Do commutativity and associativity hold for arbitrary Theta Joins?

## Algebraic Rewritings: Commutativity and Associativity (2)

Commutativity    Associativity

Union

Intersection

**Question 1**: Do the above hold for both sets and bags?
**Question 2**: Is difference commutative and associative?

## Algebraic Rewritings for Selection: Decomposition of Logical Connectives

$\sigma_{cond1}$
$\sigma_{cond2}$
R

$\sigma_{cond1 \text{ AND } cond2}$
R

$\sigma_{cond2}$
$\sigma_{cond1}$
R

$\sigma_{cond1 \text{ OR } cond2}$
R

$\sigma_{cond1}$    $\sigma_{cond2}$
R

Does it apply to bags?

## Algebraic Rewritings for Selection: Decomposition of Negation

**Question**                    **Complete**

$\sigma_{cond1 \text{ AND NOT } cond2}$
R

$\sigma_{\text{NOT } cond2}$
R

$\sigma_{cond1 \text{ OR NOT } cond2}$
R

## Pushing the Selection Thru Binary Operators: Union and Difference

$\sigma_{cond}$

$\cup$

R   S

⟺

$\cup$

$\sigma_{cond}$ R   $\sigma_{cond}$ S

Union

$\sigma_{cond}$

−

R   S

⟺

−

$\sigma_{cond}$ R   $\sigma_{cond}$ S

Difference

⟺

−

$\sigma_{cond}$   S

R

**Exercise**: Do the rule for intersection

## Pushing Selection thru Cartesian Product and Join

$\sigma_{cond}$

×

R   S

The right direction requires that *cond* refers to *S* attributes only

×

R   $\sigma_{cond}$ S

$\sigma_{cond}$

⋈

R   S

The right direction requires that *cond* refers to *S* attributes only

⋈

R   $\sigma_{cond}$ S

The right direction requires that all *the attributes used by cond appear in both R and S*

⋈

$\sigma_{cond}$ R   $\sigma_{cond}$ S

**Exercise**: Do the rule for theta join

Rules:  $\pi, \sigma$  combined

Let x = subset of R attributes
   z = attributes in predicate P
        (subset of R attributes)

$$\pi_x[\sigma_{P(R)}] = \pi_x\{\sigma_P[\underset{\pi_{xz}}{\cancel{\pi_x}}(R)]\}$$

## Pushing Simple Projections Thru Binary Operators

A projection is simple if it only consists of an attribute list

$$\pi_A(R \cup S) \Longleftrightarrow \pi_A(R) \cup \pi_A(S) \qquad \text{Union}$$

**Question 1**: Does the above hold for both bags and sets?
**Question 2:** Can projection be pushed below
   intersection and difference?
   Answer for both bags and sets.

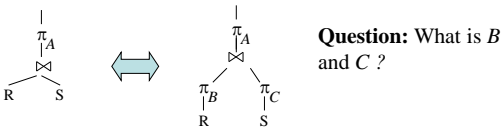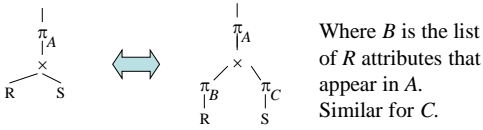## Pushing Simple Projections Thru Binary Operators: Join and Cartesian Product

$$\pi_A(R \times S) \Longleftrightarrow \pi_A(\pi_B(R) \times \pi_C(S))$$

Where *B* is the list of *R* attributes that appear in *A*. Similar for *C*.

$$\pi_A(R \bowtie S) \Longleftrightarrow \pi_A(\pi_B(R) \bowtie \pi_C(S))$$

**Question:** What is *B* and *C ?*

**Exercise:** Write the rewriting rule that pushes projection below theta join.

## Projection Decomposition

$$\pi_X(R) \Longrightarrow \pi_X(\pi_{XY}(R))$$

## Some Rewriting Rules Related to Aggregation: SUM

- $\sigma_{cond} SUM_{GroupbyList;GroupedAttribute \rightarrow ResultAttribute} R$
  $\Leftrightarrow SUM_{GroupbyList;GroupedAttribute \rightarrow ResultAttribute} \sigma_{cond} R,$
  if *cond* involves only the *GroupbyList*

- $SUM_{GL;GA \rightarrow RA}(R \cup S)$
  $\Leftrightarrow PLUS_{RA1, RA2:RA}$
  $((SUM_{GL;GA \rightarrow RA11} R) \triangleright \triangleleft (SUM_{GL;GA \rightarrow RA2} S)))$

- $SUM_{GL2;RA1 \rightarrow RA2} SUM_{GL1;GA \rightarrow RA1} R \Leftrightarrow$
  $SUM_{GL2:GA \rightarrow RA2} R$
  - **Question:** does the above hold for both bags and sets?

---

Derived Rules:  $\sigma + \bowtie$ combined

More Rules can be Derived:

$\sigma_{p \wedge q} (R \bowtie S) =$

$\sigma_{p \wedge q \wedge m} (R \bowtie S) =$

$\sigma_{p \vee q} (R \bowtie S) =$

**p** only at **R**, **q** only at **S**, **m** at both **R** and **S**

---

--> Derivation for first one:

$\sigma_{p \wedge q} (R \bowtie S) =$

$\sigma_p [\sigma_q (R \bowtie S)] =$

$\sigma_p [R \bowtie \sigma_q (S)] =$

$[\sigma_p (R)] \bowtie [\sigma_q (S)]$

Which are always "good"
transformations?

- $\sigma_{p1 \wedge p2} (R) \rightarrow \sigma_{p1} [\sigma_{p2} (R)]$

- $\sigma_p (R \bowtie S) \rightarrow [\sigma_p (R)] \bowtie S$
- $R \bowtie S \rightarrow S \bowtie R$

- $\pi_x [\sigma_p (R)] \rightarrow \pi_x \{\sigma_p [\pi_{xz} (R)]\}$

_____

_____

_____

_____

_____

_____

_____

In textbook: more transformations

- Eliminate common sub-expressions
- Other operations: duplicate elimination

_____

_____

_____

_____

_____

_____

_____

Bottom line:

- No transformation is <u>always</u> good at the l.q.p level
- Usually good:
  – early selections
  – elimination of cartesian products
  – elimination of redundant subexpressions
- Many transformations lead to "promising" plans
  – Commuting/rearranging joins
  – In practice too "combinatorially explosive" to be handled as rewriting of l.q.p.

_____

_____

_____

_____

_____

_____

_____

## Algorithms for Relational Algebra Operators

- Three primary techniques
  - Sorting
  - Hashing
  - Indexing
- Three degrees of difficulty
  - data small enough to fit in memory
  - too large to fit in main memory but small enough to be handled by a "two-pass" algorithm
  - so large that "two-pass" methods have to be generalized to "multi-pass" methods (quite unlikely nowadays)

The dominant cost of operators running on disk:

- Count # of disk blocks that must be read (or written) to execute query plan

To estimate costs, we use additional parameters:

$B(R)$ = # of blocks containing R tuples

$f(R)$ = max # of tuples of R per block

$M$ = # memory blocks available

Sorting information

$HT(i)$ = # levels in index I

Caching information (eg, first levels of index always cached)

$LB(i)$ = # of leaf blocks in index i

## Clustering index

Index that allows tuples to be read in an
order that corresponds to a sort order

A

| 10 | |
|----|--|
| 15 | |
| 17 | |

| 19 | |
|----|--|
| 35 | |
| 37 | |

A index

## Clustering can radically change cost

- Clustered relation

  | R1 R2 R3 R4 |   | R5 R5 R7 R8 |   …..

- Clustering index

## Pipelining can radically change cost

- Interleaving of operations across multiple operators
- Smaller memory footprint, fewer object allocations
- Operators support:
  - open()
  - getNext()
  - close()
- Simple for unary
- Pipelined operation for binary discussed along with physical operators

parent

open()
getNext()
close()

class project
open()
{ return child.open() }

getNext()
{ return child.getNext() }

π

child

Example   R1 ⋈ R2 over common attribute C

First we will see main memory-based
implementations

_____
_____
_____
_____
_____
_____
_____

- Iteration join (conceptually – without
  taking into account disk block issues)
      for each r ∈ R1 do
        for each s ∈ R2 do
            if r.C = s.C then output r,s pair

_____
_____
_____
_____
_____
_____
_____

- Merge join (conceptually)
    (1) if R1 and R2 not sorted, sort them
    (2) i ← 1; j ← 1;
        While (i ≤ T(R1)) ∧ (j ≤ T(R2)) do
          if R1{ i }.C = R2{ j }.C then outputTuples
          else if R1{ i }.C > R2{ j }.C then j ← j+1
          else if R1{ i }.C < R2{ j }.C then i ← i+1

_____
_____
_____
_____
_____
_____

Procedure Output-Tuples
   While (R1{ i }.C = R2{ j }.C) $\wedge$ (i $\le$ T(R1)) do
      [jj $\leftarrow$ j;
      while (R1{ i }.C = R2{ jj }.C) $\wedge$ (jj $\le$ T(R2)) do
                [output pair R1{ i }, R2{ jj };
                    jj $\leftarrow$ jj+1  ]
      i $\leftarrow$ i+1  ]

## Example

| i | R1{i}.C | R2{j}.C | j |
|---|---------|---------|---|
| 1 | 10 | 5 | 1 |
| 2 | 20 | 20 | 2 |
| 3 | 20 | 20 | 3 |
| 4 | 30 | 30 | 4 |
| 5 | 40 | 30 | 5 |
|   |    | 50 | 6 |
|   |    | 52 | 7 |

• Join with index  (Conceptually)

For each r $\in$ R1 do               | Assume R2.C index |
   [ X $\leftarrow$ index (R2, C, r.C)
      for each s $\in$ X do
            output r,s pair]

Note:  X $\leftarrow$ index(rel, attr, value)
       then X = set of rel tuples with attr = value

- Hash join (conceptual)
  - Hash function h, range $0 \to k$
  - Buckets for R1: G0, G1, ... Gk
  - Buckets for R2: H0, H1, ... Hk

  Algorithm
  (1) Hash R1 tuples into G buckets
  (2) Hash R2 tuples into H buckets
  (3) For i = 0 to k do
      match tuples in Gi, Hi buckets

Simple example    hash: even/odd

| R1 | R2 | | | Buckets | |
|----|----|----|------|---------|---------|
| 2 | 5 | Even | 2 4 8 | | 4 12 8 14 |
| 4 | 4 | | R1 | | R2 |
| 3 | 12 | Odd: | 3 5 9 | | 5 3 13 11 |
| 5 | 3 | | | | |
| 8 | 13 | | | | |
| 9 | 8 | | | | |
| | 11 | | | | |
| | 14 | | | | |

Factors that affect performance

(1)   Tuples of relation stored
          physically together?

(2)   Relations sorted by join attribute?

(3)   Indexes exist?

## Disk-oriented Computation Model

- There are $M$ main memory buffers.
  - Each buffer has the size of a disk block
- The input relation is read one block at a time.
- The cost is the number of blocks read.
- If $B$ consecutive blocks are read the cost is $B/d$.
- The output buffers are not part of the $M$ buffers mentioned above.
  - *Pipelining* allows the output buffers of an operator to be the input of the next one.
  - We do not count the cost of writing the output.

## Notation

- $B(R)$ = number of blocks that $R$ occupies
- $T(R)$ = number of tuples of $R$
- $V(R,[a_1, a_2, \ldots, a_n])$ = number of distinct tuples in the projection of $R$ on $a_1, a_2, \ldots, a_n$

## One-Pass Main Memory Algorithms for Unary Operators

- Assumption: Enough memory to keep the relation
- Projection and selection:
  - Scan the input relation $R$ and apply operator one tuple at a time
  - Incremental cost of "on the fly" operators is 0

- Duplicate elimination and aggregation
  - create one entry for each group and compute the aggregated value of the group
  - it becomes hard to assume that CPU cost is negligible
    - main memory data structures are needed

## One-Pass Nested Loop Join

- Assume $B(R)$ is less than $M$
- Tuples of $R$ should be stored in an efficient lookup structure
- **Exercise:** Find the cost of the algorithm below

```
for each block Br of R do
  store tuples of Br in main memory
for each each block Bs of S do
  for each tuple s of Bs
    join tuples of s with matching tuples of R
```

## Generalization of Nested-Loops

```
for each chunk of M-1 blocks Br of R do
  store tuples of Br in main memory
  for each each block Bs of S do
    for each tuple s of Bs
      join tuples of s with matching tuples of R
```

**Exercise**: Compute cost

## Simple Sort-Merge Join

- Assume natural join on *C*
- Sort *R* on *C* using the two-phase multiway merge sort
  - if not already sorted
- Sort *S* on *C*
- Merge (opposite side)
  - assume two pointers **Pr,Ps** to tuples on disk, initially pointing at the start
  - sets **R'**, **S'** in memory
- Remarks:
  - Very low average memory requirement during merging (but no guarantee on how much is needed)
  - **Cost:**

```
while Pr!=EOF and Ps!=EOF
  if *Pr[C] == *Ps[C]
    do_cart_prod(Pr,Ps)
  else if *Pr[C] > *Ps[C]
    Ps++
  else if *Ps[C] > *Pr[C]
    Pr++

function do_cart_prod(Pr,Ps)
  val=*Pr[C]
  while *Pr[C]==val
    store tuple *Pr in set R'
  while *Ps[C]==val
    store tuple *Ps in set S';
  output cartesian product
      of R' and S'
```

## Efficient Sort-Merge Join

- Idea: Save two disk I/O's per block by combining the second pass of sorting with the ``merge''.
- Step 1: Create sorted sublists of size $M$ for $R$ and $S$
- Step 2: Bring the first block of each sublist to a buffer
  - assume no more than $M$ sublists in all
- Step 3: Repeatedly find the least $C$ value $c$ among the first tuples of each sublist. Identify all tuples with join value $c$ and join them.
  - When a buffer has no more tuple that has not already been considered load another block into this buffer.

## Efficient Sort-Merge Join Example

$R$

| $C$ | $RA$ |
|-----|------|
| 1 | $r_1$ |
| 2 | $r_2$ |
| 3 | $r_3$ |
| ... | |
| 20 | $r_{20}$ |

$S$

| $C$ | $SA$ |
|-----|------|
| 1 | $s_1$ |
| ... | |
| 5 | $s_5$ |
| 16 | $s_{16}$ |
| ... | |
| 20 | $s_{20}$ |

Assume that after first phase of multiway sort we get 4 sublists, 2 for $R$ and 2 for $S$.
Also assume that each block contains two tuples.

$R$

| 3  7 | 8 10 | 11 13 | 14 16 | 17 18 |
|------|------|-------|-------|-------|
| 1  2 | 4  5 | 6  9  | 12 15 | 19 20 |

$S$

| 1  3 | 5 17 | |
|------|------|--|
| 2  4 | 16 18 | 19 20 |

## Two-Pass Hash-Based Algorithms

- General Idea: Hash the tuples of the input arguments in such a way that all tuples that must be considered together will have hashed to the same hash value.
  - If there are $M$ buffers pick $M-1$ as the number of hash buckets
- Example: Duplicate Elimination
  - Phase 1: Hash each tuple of each input block into one of the $M-1$ bucket/buffers. When a buffer fills save to disk.
  - Phase 2: For each bucket:
    - load the bucket in main memory,
    - treat the bucket as a small relation and eliminate duplicates
    - save the bucket back to disk.
  - **Catch:** Each bucket has to be less than $M$.
  - **Cost:**

## Hash-Join Algorithms

- Assuming natural join, use a hash function that
  – is the same for both input arguments $R$ and $S$
  – uses only the join attributes
- Phase 1: Hash each tuple of $R$ into one of the $M-1$ buckets $R_i$ and similar each tuple of $S$ into one of $S_i$
- Phase 2: For $i=1...M-1$
  – load $R_i$ and $S_i$ in memory
  – join them and save result to disk
- **Question:** What is the maximum size of buckets?
- **Question:** Does hashing maintain sorting?

## Index-Based Join: The Simplest Version

Assume that we do natural join of $R(A,B)$ and $S(B,C)$ and there's an index on $S$

```
for each Br in R do
  for each tuple r of Br with B value b
    use index of S to find
          tuples {s₁ ,s₂ ,...,sₙ} of S with B=b
    output {rs₁ ,rs₂ ,...,rsₙ}
```

**Cost:** Assuming $R$ is clustered and non-sorted and the index on $S$ is clustered on $B$ then
$B(R)+T(R)B(S)/V(S,B)$ + some more for reading index
**Question:** What is the cost if $R$ is sorted?

## Opportunities in Joins Using Sorted Indexes

- Do a conventional Sort-Join avoiding the sorting of one or both of the input operands

• Estimating cost of query plan

(1) Estimating <u>size</u> of results
(2) Estimating # of IOs

<u>Estimating result size</u>

• Keep statistics for relation R
  – $T(R)$ : # tuples in R
  – $S(R)$ : # of bytes in each R tuple
  – $B(R)$: # of blocks to hold all R tuples
  – $V(R, A)$ : # distinct values in R
            for attribute A

<u>Example</u>

R

| A | B | C | D |
|-----|---|----|---|
| cat | 1 | 10 | a |
| cat | 1 | 20 | b |
| dog | 1 | 30 | a |
| dog | 1 | 40 | c |
| bat | 1 | 50 | d |

A: 20 byte string
B: 4 byte integer
C: 8 byte date
D: 5 byte string

$T(R) = 5$    $S(R) = 37$
$V(R,A) = 3$          $V(R,C) = 5$
$V(R,B) = 1$          $V(R,D) = 4$

Size estimates  for W = R1 x R2

T(W) =

$$T(R1) \times T(R2)$$

S(W) =

$$S(R1) + S(R2)$$

_____

_____

_____

_____

_____

_____

_____

Size estimate  for W = $\sigma_{Z=val}$ (R)

S(W) = S(R)

T(W) = ?

_____

_____

_____

_____

_____

_____

_____

Example

R

| A | B | C | D |
|-----|---|----|---|
| cat | 1 | 10 | a |
| cat | 1 | 20 | b |
| dog | 1 | 30 | a |
| dog | 1 | 40 | c |
| bat | 1 | 50 | d |

V(R,A)=3
V(R,B)=1
V(R,C)=5
V(R,D)=4

$$W = \sigma_{Z=val}(R) \quad T(W) = \frac{T(R)}{V(R,Z)}$$

_____

_____

_____

_____

_____

_____

_____

What about $W = \sigma_{z \geq val}(R)$ ?

$T(W) = ?$

- Solution # 1:
    $T(W) = T(R)/2$

- Solution # 2:
    $T(W) = T(R)/3$

- Solution # 3:   Estimate values in range

<u>Example</u>  R



Min=1    $V(R,Z)=10$

$W = \sigma_{z \geq 15}(R)$

Max=20

$f = \dfrac{20-15+1}{20-1+1} = \dfrac{6}{20}$    (fraction of range)

$T(W) = f \times T(R)$

Equivalently:
    $f \times V(R,Z)$ = fraction of distinct values
$T(W) = [f \times V(Z,R)] \times \dfrac{T(R)}{V(Z,R)} = f \times T(R)$

Size estimate for $W = R1 \bowtie R2$

Let x = attributes of R1
y = attributes of R2

Case 1     $X \cap Y = \varnothing$

Same as R1 x R2

Case 2   $W = R1 \bowtie R2$     $X \cap Y = A$

| R1 | A | B | C |   | R2 | A | D |
|----|---|---|---|---|----|---|---|
|    |   |   |   |   |    |   |   |

Assumption:

$\Pi_A R1 \subseteq \Pi_A R2 \Rightarrow$ Every A value in R1 is in R2
(typically A of R1 is foreign key
of the primary key of A of R2)
$\Pi_A R2 \subseteq \Pi_A R1 \Rightarrow$ Every A value in R2 is in R1
"containment of value sets"   (justified by primary
key – foreign key relationship)

Computing T(W)  when A of R1 is the

foreign key $\Pi_A R1 \subseteq \Pi_A R2$



| R1 | A | B | C |   | R2 | A | D |
|----|---|---|---|---|----|---|---|

1 tuple of R1 matches with exactly 1 tuple
of R2

so    T(W)  =    T(R1)

Another way to approach when

$\Pi_A R1 \subseteq \Pi_A R2$

R1 | A | B | C    R2 | A | D

Take
1 tuple       Match

1 tuple matches with $\dfrac{T(R2)}{V(R2,A)}$ tuples...

so $\quad T(W) = \dfrac{T(R2)}{V(R2, A)} \times T(R1)$


• $V(R1,A) \leq V(R2,A) \quad T(W) = \dfrac{T(R2)\, T(R1)}{V(R2,A)}$

• $V(R2,A) \leq V(R1,A) \quad T(W) = \dfrac{T(R2)\, T(R1)}{V(R1,A)}$

[A is common attribute]


In general $\quad W = R1 \bowtie R2$

$T(W) = \dfrac{T(R2)\, T(R1)}{\max\{ V(R1,A), V(R2,A) \}}$

## Combining estimates on subexpressions: Value preservation

Result = |
$\sigma_{C=1}$
$\bowtie$
R    S

Result = |
$\bowtie$
    S
$\sigma_{C=1}$
|
R

R(A, C)          S(A, B)
$T(R) = 10^3$     $T(S) = 10^2$
$V(A, R) = 10^3$  $V(A, S) = 50$
$V(C, R) = 10^2$

$T(R \bowtie S) =$
$T(R) \times T(S) / max(V(A,R), V(A, S)) = 10^2$

$V(C, R \bowtie S) = 10^2$ ← (Big) assumption:
                           Value preservation of C

$T(Result) = T(R \bowtie S) / V(C, R \bowtie S) = 1$

---

## Value preservation may have to be pushed to a weird assumption (but there's logic behind it!)

Result = |
$\sigma_{C=1}$
$\bowtie$
R    S

Result = |
$\bowtie$
    S
$\sigma_{C=1}$
|
R

Ideally, the size estimation should not depend on which of the two equivalent formulas for Result one uses. However, to achieve this we may need to push the value preservation assumption to artificial intermediate estimates…

R(A, C)          S(A, B)
$T(R) = 10^3$     $T(S) = 10^2$
$V(A, R) = 10^3$  $V(A, S) = 50$
$V(C, R) = 10^2$

$T(R \bowtie S) = 10^2$

$V(C, R \bowtie S) = 10^2$

$T(Result) = 1$

$T(\sigma_{c=1}R) = T(R) / V(C, R) = 10$
$V(A, \sigma_{c=1}R) = \mathbf{10^3}$
$T(Result) =$
$T(\sigma_{c=1}R) \times T(S) / max(V(A, \sigma_{c=1}R), V(A, S)) = 1$

We had to extend value preservation to the weird assumption that attribute A has more values than the number of tuples in R. In this way the number of S tuples matching an R tuple stays steady

---

## Value preservation of join attribute

Foreign-to-primary
CSEenroll(EID, SID, …)    Students(SID, …)    Honors (HID, SID, …)

T(CSEenroll) = 10,000     T(Students) = 20,000    T(Honors) = 5,000
V(SID, CSEenroll) = 1,000 V(SID, Students) = 20,000 V(SID, Honors) = 500

T(CSEenroll(EID, SID, …) $\bowtie$ Students(SID, …) $\bowtie$ Honors (HID, SID, …)) = ?

$\bowtie$          T(.) = **10,000 x 5,000 / max(500, <u>20,000</u>) = 2,500  CORRECT**
                   10,000 x 5,000 / max(500, <u>1,000</u>) = 50,000  WRONG
        Honors

$\bowtie$  T(.) = 10,000
           V(SID, .) ?= 1,000 (preservation of SIDs in CSEenroll)
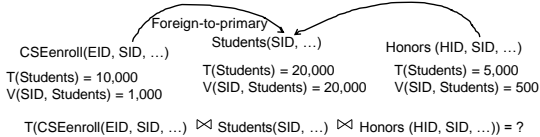CSEenroll  Students   or 20,000 (preservation of SIDs in Students) ?

### If in doubt, think in terms of probabilities and matching records

- A SID of Student appears in CSEEnroll with probability 1000/20000
  - i.e., 5% of students are enrolled in CSE
- A SID of Student appears in Honors with probability 500/20000
  - i.e., 2.5% of students are honors students
=> An SID of Student appears in the join result with probability 5% x 2.5%
- On the average, each SID of CSEEnroll appears in 10,000/1,000 tuples
  - i.e., each CSE-enrolled student has 10 enrollments
- On the average, each SID of Honors appears in 5,000/500 tuples
  - i.e., each honors' student has 10 honors
⇒ Each Student SID that is in both Honors and CSEEnroll is in 10x10 result tuples
⇒ T(result) = 20,000 x 5% x 2.5% x 10 x 10 = 2,500 tuples

Foreign-to-primary

CSEenroll(EID, SID, …)     Students(SID, …)      Honors (HID, SID, …)

T(Students) = 10,000      T(Students) = 20,000      T(Students) = 5,000
V(SID, Students) = 1,000   V(SID, Students) = 20,000   V(SID, Students) = 500

T(CSEenroll(EID, SID, …) ⋈ Students(SID, …) ⋈ Honors (HID, SID, …)) = ?

---

## Plan Enumeration

- A smart exhaustive algorithm
  - According to textbook's Section 16.6
  - no ppt notes
- The INGRES heuristic for plan enumeration

---

## Arranging the Join Order: the Wong-Youssefi algorithm (INGRES)

**Sample TPC-H Schema**
```
Nation(NationKey, NName)
Customer(CustKey, CName, NationKey)
Order(OrderKey, CustKey, Status)
Lineitem(OrderKey, PartKey, Quantity)
Product(SuppKey, PartKey, PName)
Supplier(SuppKey, SName)
```

Find the names of suppliers that sell a product that appears in a line item of an order made by a customer who is in Canada

```
SELECT SName
FROM Nation, Customer, Order, LineItem, Product, Supplier
WHERE Nation.NationKey = Cuctomer.NationKey
    AND Customer.CustKey = Order.CustKey
    AND Order.OrderKey=LineItem.OrderKey
    AND LineItem.PartKey= Product.Partkey
    AND Product.Suppkey = Supplier.SuppKey
    AND NName = "Canada"
```
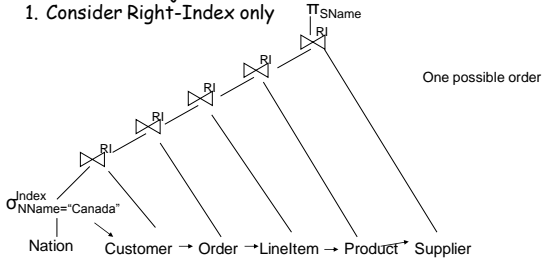
## Challenges with Large Natural Join Expressions

For simplicity, assume that in the query
1. All joins are natural
2. whenever two tables of the FROM clause have common attributes we join on them
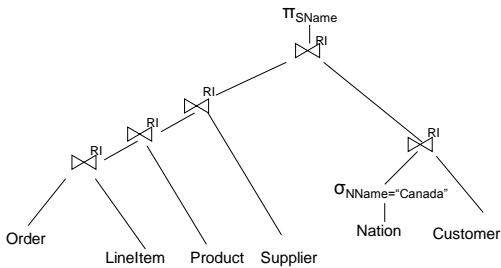1. Consider Right-Index only

$\pi_{SName}$

One possible order

$\sigma_{NName="Canada"}^{Index}$

Nation → Customer → Order → LineItem → Product → Supplier

## Multiple Possible Orders

$\pi_{SName}$

Order LineItem Product Supplier $\sigma_{NName="Canada"}$ Nation Customer

## Wong-Yussefi algorithm assumptions and objectives

- Assumption 1 (weak): Indexes on all join attributes (keys and foreign keys)
- Assumption 2 (strong): At least one selection creates a *small* relation
  – A join with a small relation results in a small relation
- Objective: Create sequence of index-based joins such that all intermediate results are small

# Hypergraphs



- relation hyperedges
  - two hyperedges for same relation are possible
- each node is an attribute
- can extend for non-natural equality joins by merging nodes

# Small Relations/Hypergraph Reduction
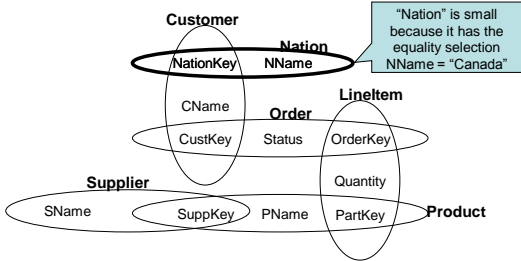


"Nation" is small because it has the equality selection NName = "Canada"

$\sigma^{Index}_{NName="Canada"}$
|
Nation

Pick a small relation (and its conditions) to start the plan



Remove small relation (hypergraph reduction) and color as "small" any relation that joins with the removed "small" relation

$$\bowtie^{RI}$$

$\sigma^{Index}_{NName="Canada"}$
|
Nation          Customer

Pick a small relation (and its conditions if any) and join it with the small relation that has been reduced

35

## After a bunch of steps…



$$\pi_{SName}$$

$$\sigma^{Index}_{NName="Canada"}$$

Nation    Customer → Order → LineItem → Product → Supplier

## Multiple Instances of Each Relation

SELECT S.SName
FROM Nation, Customer, Order, LineItem L, Product P, Supplier S,
         LineItem LE, Product PE, Supplier Enron
WHERE Nation.NationKey = Cuctomer.NationKey
     AND Customer.CustKey = Order.CustKey
     AND Order.OrderKey=L.OrderKey
     AND L.PartKey= P.Partkey
     AND P.Suppkey = S.SuppKey
     AND Order.OrderKey=LE.OrderKey
     AND LE.PartKey= PE.Partkey
     AND PE.Suppkey = Enron.SuppKey
     AND Enron.Sname = "Enron"
     AND NName = "Cayman"

> Find the names of suppliers whose products appear in an order made by a customer who is in Cayman Islands and an Enron product appears in the same order

## Multiple Instances of Each Relation

# Multiple choices are possible

## The basic dynamic programming approach to enumerating plans

for each sub-expression

$op(e_1 e_2 \ldots e_n)$ of a logical plan

– (recursively) compute the best plan and cost for each subexpression $e_i$
– for each physical operator $op^p$ implementing $op$
  • evaluate the cost of computing $op$ using $op^p$ and the best plan for each subexpression $e_i$
  • (for faster search) memo the best $op^p$

## Local suboptimality of basic approach and the Selinger improvement

• Basic dynamic programming may lead to (globally) suboptimal solutions
• Reason: A suboptimal plan for $e_1$ may lead to the optimal plan for $op(e_1 e_2 \ldots e_n)$
  – Eg, consider $e_1 \bowtie e_2$ and
  – assume that the optimal computation of $e_1$ produces unsorted result
  – Optimal $\bowtie$ is via sort-merge join on A
  – It could have paid off to consider the suboptimal computation of $e_1$ that produces result sorted on A
• Selinger improvement: memo also any plan (that computes a subexpression) and produces an order that may be of use to ancestor operators

## Using dynamic programming to optimize a join expression

- Goal: Decide the join order and join methods
- Initiate with n-ary join $\bowtie_C (e_1 \, e_2 \, \ldots \, e_n)$, where $c$ involves only join conditions
- Bottom up: consider 2-way non-trivial joins, then 3-way non-trivial joins etc
  - "non trivial" -> no cartesian product