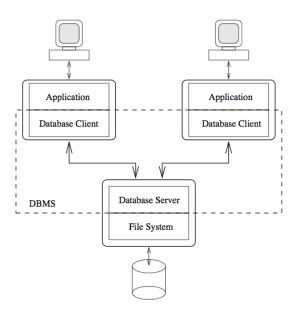
# Database Tuning and Physical Design: Basics of Query Execution Fall, 2018

School of Computer Science University of Waterloo

Databases CS348

## The Client/Server Architecture



# The Client/Server Architecture (cont'd)

## **Applications**

- User interaction: query input, presentation of results
- Application-specific tasks

#### **Database Server**

- DDL evaluation
- DML compilation: selection of a query plan for a query
- DML execution
- Concurrency control
- Buffer management: rollback and failure recovery

## File System

Storage and retrieval of unstructured data

# **Basics of Query Execution**

#### Goal

Develop a simple *relational calculator* that answers queries.

#### Considerations:

- How is data physically represented?
- 2 How to compute answers to complex queries?
- How are intermediate results managed?

4 / 43

## How do we Execute Queries?

- 1 Parsing, typechecking, etc.
- 2 Relational Calculus (SQL) translated to Relational Algebra
- Optimization:
  - ⇒ generates an efficient *query plan*
  - ⇒ uses statistics collected about the stored data
- Plan execution:
  - ⇒ access methods to access stored relations
  - ⇒ physical relational operators to combine relations

# Relational Algebra

#### Idea

Define a *set of operations* on the universe of finite relations...
...called a *RELATIONAL ALGEBRA*.

$$(\mathcal{U}; R_0, \ldots, R_k, \times, \sigma_{\varphi}, \pi_{V}, \cup, -)$$

#### Constants:

 $R_i$ : one for each relational scheme

## Unary operators:

 $\sigma_{\varphi}$ : selection (keeps only tuples satisfying  $\varphi$ )

 $\pi_V$ : projection (keeps only attributes in V)

## Binary operators:

x: Cartesian product

∪: union

-: set difference

# Examples

Account				
Acnt#	Type	Balance	Bank	Branch
1234	CHK	\$1000	TD	1
1235	SAV	\$20000	TD	2
1236	CHK	\$2500	CIBC	1
1237	CHK	\$2500	Royal	5
2000	BUS	\$10000	Royal	5
2001	BUS	\$10000	TD	3

Bank	
Name	Address
TD	TD Centre
CIBC	CIBC Tower

## **Projection**

#### Definition:

$$\pi_{V}(R) = \{(x_{i_1}, \dots, x_{i_k}) : (x_1, \dots, x_n) \in R, i_j \in V\}$$

where *V* is an *ordered list* of column *numbers*.

## Example:

$$\pi_{\#1,\#2}(Account) =$$

CHK
SAV
CHK
CHK
BUS
BUS

8 / 43

## Selection

Definition:

$$\sigma_{\varphi}(R) = \{(x_1, \dots, x_n) : (x_1, \dots, x_n) \in R, \\ \land \varphi(x_1, \dots, x_n)\}$$

where  $\varphi$  is a *built-in* selection condition.

## Example:

$$\sigma_{\#3>5000}(Account) =$$

1235	SAV	\$20000	TD	2
2000	BUS	\$10000	Royal	5
2001	BUS	\$10000	TD	3

## **Product**

#### Definition:

$$R \times S = \{((x_1, \dots, x_n, y_1, \dots, y_m) : (x_1, \dots, x_n) \in R, (y_1, \dots, y_n) \in S\}$$

Example:  $Account \times Bank =$ 

1234	CHK	\$1000	TD	1	TD	TD Centre
1235	SAV	\$20000	TD	2	TD	TD Centre
1236	CHK	\$2500	CIBC	1	TD	TD Centre
1237	CHK	\$2500	Royal	5	TD	TD Centre
2000	BUS	\$10000	Royal	5	TD	TD Centre
2001	BUS	\$10000	TD	3	TD	TD Centre
1234	CHK	\$1000	TD	1	CIBC	CIBC Tower
1235	SAV	\$20000	TD	2	CIBC	CIBC Tower
1236	CHK	\$2500	CIBC	1	CIBC	CIBC Tower
1237	CHK	\$2500	Royal	5	CIBC	CIBC Tower
2000	BUS	\$10000	Royal	5	CIBC	CIBC Tower
2001	BUS	\$10000	TD	3	CIBC	CIBC Tower

## Union

#### Definition:

$$\label{eq:sum} \begin{split} \textit{R} \cup \textit{S} &= \{ (x_1, \dots, x_n) : \; (x_1, \dots, x_n) \in \textit{R} \\ &\quad \lor (x_1, \dots, x_n) \in \textit{S} \} \end{split}$$

## Example:

$$\pi_{\#1,\#2}(\sigma_{\#2='CHK'}(Account)) \cup \pi_{\#1,\#2}(\sigma_{\#2='SAV'}(Account)) =$$

1234	CHK
1236	CHK
1237	CHK
1235	SAV

## Difference

Definition:

$$\{ (x_1, \dots, x_n) : \ (x_1, \dots, x_n) \in R, \\ \wedge (x_1, \dots, x_n) \not \in S \}$$

## Example:

Is there an account without a bank?

$$\pi_{\#1,\#4}(Account) - \pi_{\#1,\#4}(\sigma_{\#4=\#6}(Account \times Bank)) =$$

1237 Royal 2000 Royal

12 / 43

# Relational Calculus/SQL to Algebra

How do we know that these operators are sufficient to evaluate *all* Relational Calculus queries?

## Theorem (Codd)

For every domain independent Relational Calculus query there is an equivalent Relational Algebra expression.

```
RCtoRA(R_i(x_1, ..., x_k)) = R_i
RCtoRA(Q \land x_i = x_j) = \sigma_{\#i = \#j}(RCtoRA(Q))
RCtoRA(\exists x_i.Q) = \pi_{FV(Q) - \{\#i\}}(RCtoRA(Q))
RCtoRA(Q_1 \land Q_2)) = RCtoRA(Q_1) \lor RCtoRA(Q_2)
RCtoRA(Q_1 \lor Q_2)) = RCtoRA(Q_1) \cup RCtoRA(Q_2)
RCtoRA(Q_1 \land \neg Q_2)) = RCtoRA(Q_1) - RCtoRA(Q_2)
```

- $\dots$  queries in  $\wedge$  must have disjoint sets of free variables
- ... we must *invent* consistent way of referring to attributes

## Iterator Model for RA

How do we avoid (mostly) storing intermediate results?

#### Idea

We use the cursor OPEN/FETCH/CLOSE interface.

Every implementation of an Relational Algebra operator:

- implements the cursor interface to produce answers
- uses the same interface to get answers from its children

... we make (at least) one *physical implementation* per operator.

# Physical Operators (example: selection)

```
// select {#i=#j}(Child)
  OPERATOR child:
  int i, j;
public:
  OPERATOR selection (OPERATOR c, int i0, int j0)
                { child = c; i = i0; j = j0; };
  void open() { child.open(); };
  tuple fetch() { tuple t = child.fetch();
                  if (t==NULL || t.attr(i)=t.attr(j))
                       return t;
                  return this.fetch();
  void close() { child.close(); }
```

# Physical Operators (cont.)

The rest of the lot:

product:

simple nested loops algorithm

projection:

eliminate unwanted attributes from each tuple

union:

simple concatenation

set difference:

nested loops algorithm that checks for tuples on r.h.s.

#### WARNING

This doesn't quite work: projection and union may produce *duplicates*... need to be followed by a *duplicate elimination operator* 

# How to make it FAST(er)?

#### Observation

Naive implementation for each operator will work

...very (very very very) slowly

#### What to do?

- Use (disk-based) data structures for efficient searching INDEXING (used, e.g., in selections)
- Use better algorithms to implement the operators commonly based on SORTING or HASHING
- Rewrite the RA expression to an equivalent, but more efficient one remove unnecessary operations (e.g., duplicate elimination) enable the use of better algorithms/data structures

# Atomic Relations and Indexing

■ When an index  $R_{index}(x)$  (where x is the search attribute) is available we replace a subquery of the form

$$\sigma_{\mathsf{X}=\mathsf{C}}(\mathsf{R})$$

with accessing  $R_{index}(x)$  directly,

Otherwise: check all file blocks holding tuples for *R*.

Even if an index is available, scanning the entire relation may be faster in certain circumstances:

- The relation is very small
- The relation is large, but we expect most of the tuples in the relation to satisfy the selection criteria

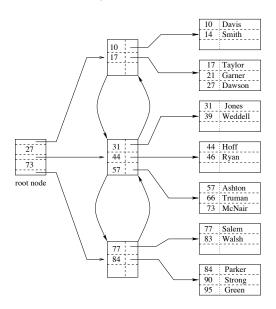
# Clustering vs. Non-Clustering Indexes

- An index on attribute *A* of a relation is a **clustering** index if tuples in the relation with similar values for *A* are stored together in the same block.
- Other indices are **non-clustering** (or secondary) indices.

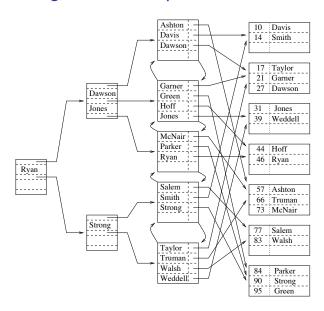
#### Note

A relation may have at most one clustering index, and any number of non-clustering indices.

# Clustering Index Example



# Non-Clustering Index Example

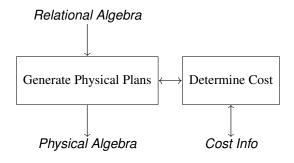


# **Query Optimization**

- Many possible query plans for a single query:
  - 1 equivalences in Relational Algebra
  - 2 choice of Operator Implementation
  - ⇒ performance differs greatly
- How do we choose the best plan?
  - 1 "always good" transformations
  - 2 cost-based model
  - $\Rightarrow$  finding an **optimal plan** is computationally not feasible: we look for a *reasonable* one

# General Approach

- Generate all physical plans equivalent to the query
- Pick the one with the lowest cost



## ... All Equivalent Plans?!

- Cannot be done in general:
  - ⇒ it is **undecidable** if a query is (un-)satisfiable (equivalent to an *empty* plan)
- Very expensive even for conjunctive queries
  - ⇒ the *Join-ordering* problem
- In practice:
  - $\Rightarrow$  only plans of certain form are considered (restrictions on the search space.)
  - ⇒ the goal is to eliminate the really bad ones

## ... and Pick the Best one?!

- How do we determine which plan is the best one?
  - ⇒ cannot simply run the plan to find out
  - $\Rightarrow$  instead, estimate the cost based on stats collected by the DBMS for all relations
- A Simple Cost Model for disk I/O; Assumptions:

Uniformity: all possible values of an attribute are equally likely to appear in a relation.

Independence: the likelihood that an attribute has a particular value (in a tuple) does not depend on values of other attributes.

# A Simple Cost Model (cont.)

- For a stored relation *R* with an attribute *A* we keep:
  - |R|: the cardinality of R (the number of tuples in R)
  - b(R): the blocking factor for R
  - 3 min(R, A): the minimum value for A in R
  - $\blacksquare$  max(R, A): the maximum value for A in R
  - 5 distinct(R, A): the number of distinct values of A
- Based on these values we try to estimate the cost of physical plans.

## Cost of Retrieval

```
Mark(Studnum, Course, Assignnum, Mark)
SELECT Studnum, Mark
FROM Mark
WHERE Course = 'PHYS'
AND Studnum = 100 AND Mark > 90
```

#### Indices:

- clustering index CourseInd on Course
- non-clustering index StudnumInd on Studnum

#### Assume:

- |Mark| = 10000
- b(Mark) = 50
- 500 different students
- 100 different courses
- 100 different marks

# Strategy 1: Use CourseInd

Assuming *uniform distribution* of tuples over the courses, there will be about |Mark|/100 = 100 tuples with Course = PHYS.

Searching the CourseInd index has a cost of 2. Retrieval of the 100 matching tuples adds a cost of 100/b(Mark) data blocks. The total cost of 4.

Selection of N tuples from relation R using a clustered index has a cost of 2 + N/b(R).

28 / 43

# Strategy 2: Use StudnumInd

Assuming *uniform distribution* of tuples over student numbers, there will be about |Mark|/500 = 20 tuples for each student.

Searching the StudnumInd has a cost of 2. Since this is not a clustered index, we will make the pessimistic assumption that each matching record is on a separate data block, i.e., 20 blocks will need to be read. The total cost is 22.

Selection of N tuples from relation R using a clustered index has a cost of 2 + N.

# Strategy 3: Scan the Relation

The relation occupies 10,000/50 = 200 blocks, so 200 block I/O operations will be required.

Selection of N tuples from relation R by scanning the entire relation has a cost of |R|/b(R).

# Cost of other Relational Operations

Costs of **physical** operations (in I/O's):

- Selection:  $cost(\sigma_c(E)) = (1 + \epsilon_c) cost(E)$ .
- Nested-Loop Join (*R* is the **outer** relation):

$$cost(R \bowtie S) = cost(R) + (|R|/b) cost(S)$$

■ Index Join (R is the outer relation, and S is the inner relation: B-tree with depth  $d_S$ ):

$$cost(R \bowtie S) = cost(R) + d_S|R|$$

Sort-Merge Join:

$$cost(R \bowtie S) = cost(sort(R)) + cost(sort(S))$$
where  $cost(sort(E)) = cost(E) + (|E|/b) \log(|E|/b)$ .

...

Why don't we always use the Merge-Sort Join?

31 / 43

## Size Estimation

In the cost estimation we need to know sizes of results of operations: we use the **selectivity**, defined, for a condition  $\sigma_{\text{condition}}(R)$ , as:

$$\operatorname{sel}(\sigma_{\operatorname{condition}}(R)) = \frac{|\sigma_{\operatorname{condition}}(R)|}{|R|}$$

Again, the optimizer will *estimate* selectivity using simple rules based on its statistics:

$$\operatorname{sel}(\sigma_{A=c}(R)) pprox rac{1}{\operatorname{distinct}(R,A)}$$
 $\operatorname{sel}(\sigma_{A \leq c}(R)) pprox rac{c - \min(R,A)}{\max(R,A) - \min(R,A)}$ 
 $\operatorname{sel}(\sigma_{A \geq c}(R)) pprox rac{\max(R,A) - c}{\max(R,A) - \min(R,A)}$ 

# Size Estimation (cont.)

#### For Joins:

■ General Join (on attribute *A*):

$$|R \bowtie S| \approx |R| \frac{|S|}{\operatorname{distinct}(S, A)}$$

or as

$$|R \bowtie S| \approx |S| \frac{|R|}{\operatorname{distinct}(R, A)}$$

Foreign key Join (Student and Enrolled joined on Sid):

$$|R \bowtie S| = |S| \frac{|R|}{|S|} = |R|$$

May joins are foreign key joins, like this one.

## More Advanced Statistics

So far, have only a very primitive cost estimation approach

In practice: more complex approaches

- histograms to approximate non-uniform distributions
- correlations between attributes
- uniqueness (keys) and containment (inclusions)
- sampling methods
- etc, etc

## Plan Generation

- Apply "always good" transformations
  - ⇒ heuristics that work in the majority of cases
- Cost-based join-order selection
  - ⇒ applied on conjunctive subqueries (the "select blocks")
  - $\Rightarrow$  still computationaly not tractable.

## "Always good" transformations

Push selections:

$$\sigma_{\varphi}(E_1 \bowtie_{\theta} E_2) = \sigma_{\varphi}(E_1) \bowtie_{\theta} E_2$$

for  $\varphi$  involving columns of  $E_1$  only (and vice versa).

Push projections:

$$\pi_V(R \bowtie_{\theta} S) = \pi_V(\pi_{V_1}(R) \bowtie_{\theta} \pi_{V_2}(S))$$

where  $V_1$  is the set of all attributes of R involved in  $\theta$  and V (similarly for  $V_2$ ).

Replace products by joins:

$$\sigma_{\varphi}(R \times S) = R \bowtie_{\varphi} S$$

36 / 43

⇒ also reduces the space of plans we need to search

## Example

- Assume that
  - there are |S| = 1000 students,
  - enrolled in |C| = 500 classes.
  - the enrollment table is |E| = 5000,
  - and, on average, each student is registered for five courses.
- Then:

$$\operatorname{cost}(\sigma_{\operatorname{name}='\operatorname{Smith}'}(S \bowtie (E \bowtie C))) >> \\ \operatorname{cost}(\sigma_{\operatorname{name}='\operatorname{Smith}'}(S) \bowtie (E \bowtie C))$$

## Join Order Selection

- Joins are associative  $R \bowtie S \bowtie T \bowtie U$  can be equivalently expressed as

  - $(R \bowtie S) \bowtie (T \bowtie U)$
  - $R \bowtie (S \bowtie (T \bowtie U))$
  - $\Rightarrow$  try to minimize the intermediate result(s).
- Moreover, we need to decide which of the subexpressions is evaluated first
  - ⇒ e.g., Nested Loop join's cost is **not** symmetric!

## Example

We have the following two join orders to pick from:

- oname='Smith'(S)  $\bowtie$  ( $E \bowtie C$ ) we produce  $E \bowtie C$ , which has one tuple for each course registration (by any student)  $\sim$  5000 tuples.
- 2  $(\sigma_{\text{name}='\text{Smith}'}(S) \bowtie E) \bowtie C$  we produce an intermediate relation which has one tuple for each course registration by a student named Smith. If there are only a few Smith's among the 1,000 students (say there are 10), this relation will contain about 50 tuples.

## **Pipelined Plans**

- All operators (except sorting) operate without storing intermediate results
  - ⇒ iterator protocols in constant storage
  - ⇒ no recomputation for left-deep plans

# **Temporary Store**

- General pipelined plans lead to recomputation
- We introduce an additional **store** operator
  - ⇒ allows us to store intermediate results in a relation
  - ⇒ we can also built a (hash) index on top of the result
- Semantically, the operator represents the identity
- The costs of plans:
  - cumulative cost—to compute the value of the expression and store then in a relation (once):

$$cost_c(store(E)) = cost_c(E) + cost_s(E) + |E|/b$$

2 scanning cost—to "read" all the tuples in the stored result of the expression:

$$cost_s(store(E)) = |E|/b$$

41 / 43

# Parallelism in Query Execution

Another approach to improving performance:

take advantage of parallelism in hardware

- Mass storage usually reads/writes data in blocks
- Multiple mass storage units can be accessed in parallel
- Relational operators amenable to parallel execution

# Summary

## Relational Algebra is the basis for efficient implementation of SQL

- Provides a connection between conceptual and physical level
- Expresses query execution in (easily) manageable pieces
- Allows the use of efficient algorithms/data structures
- Provides a mechanism for query optimization based on logical transformations (including simplifications based on integrity constraints, etc.)

Performance of database operations depends on the way queries (and updates) are executed against a particular *physical schema/design*.

... understanding *basics* of query processing is necessary
to making *physical design decisions*... performance also depends on *transaction management* (later)

(University of Waterloo) Query Execution 43 / 43