

Chapter 4

Efficiency and Optimization

PHYSICAL DATA ORGANIZATION

- the **conceptual schema** defines which data is described and its semantics.
- the **physical schema** defines the **physical database** where the data is actually stored.
⇒ efficiency
- system: the data is actually stored in **files**: data that semantically belongs together (a relation, a part of a relation (**hashing**), some relations (**cluster**)).
- additionally, there are files that contain auxiliary information (**indexes**).
- data is accessed **pagewise** or **blockwise** (typically, 4KB – 8KB).
- each page contains some **records** (tuples). Records consist of **fields** that are of an elementary type, e.g., bit, integer, real, string, or pointer.

127

DB SERVER ARCHITECTURE: SECONDARY STORAGE AND CACHING

runtime server system: accessed by user queries/updates

- parser: translates into algebra, determines the required relations + indexes
- file manager: determines the file/page where the requested data is stored
- buffer/cache manager: provides relevant data in the cache
- query/update processing: uses only the cache

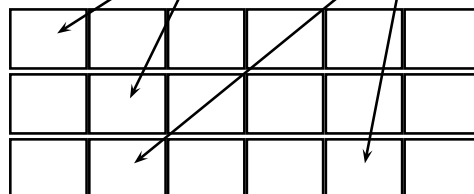
Cache (main memory): pagewise organized

- Accessed pages are fetched into the cache
- pages are also changed in the cache
- and written to the database later ...



Secondary Storage (Harddisk): pagewise organized

- data pages with tuples
- index pages with tree indexes (see later)
- database log etc. (see later)



128

DATABASE ACCESS MECHANISM

Records must be loaded from (and written to) the secondary memory for processing:

- the **file manager** determines the page where the record is stored.
- the **buffer/cache manager** is responsible to provide the page in the buffer (**buffer management**):
 - maintains a **pool** of pages (organized as frames).
for every page, it is stored if the page has been changed, how often/frequently it has been used, and if it is currently used by transactions
 - if the required page is not in the cache, some stored page is replaced (if it has been changed, it must be written to the secondary memory)
- complex **prefetching** strategies, based on knowledge about transactions.
[see lecture on Operating Systems]
- for now, it is sufficient to note that we have to deal with pagewise access.

129

Storage of Files, Pages, and Records

- Inside a file, every tuple/record has a **tuple identifier** of the form (p, n) where p is the page number and n is its index inside the page.
Each page then contains a **directory** that assigns a physical address to each n .
- memory management for deleted records
- different strategies for fixed-length and variable-length records

... so far to the physical facts ...

130

4.1 Efficient Data Access

- efficiency depends on the detailed organization and additional algorithms and data structures
- support **generic** operations:
 - **Scan**: all pages that contain records are read.
 - **Equality Search**: all records that satisfy some equality predicate are read.
`SELECT * FROM City WHERE Country = "D";`
 - **Range Search**: all records that satisfy some comparison predicate are read.
`SELECT * FROM City WHERE Population > 100.000;`
 - **Modify, Delete**: analogously
 - **Insert**: analogously: search for an appropriate place where to put the record.
- linear search (scan) ??
- **Need for efficient searching (equality and/or range)**

131

INDEXING

Indexes (for a file) are auxiliary structures that support special (non-linear) **access paths**

- Based on **search keys**
- not necessarily the relational “keys”, but *any* combination of attributes
- a relation may have several search keys
- an *index* is a set of data entries on some pages together with efficient access mechanisms for locating an entry according to its search key value.
- different types of indexes, depending on the operations to be supported:
 - equality search
 - range search (ordered values)
 - search on small domains
- in general, joins by key/foreign-key references are supported by indexes.

132

TREE INDEXES

This topic brings data structures and databases (= applications of data structures) together.

- introductory lecture “Computer Science I”: store numbers in trees.
- databases: tree *index* over the values of a column of a relation
 - search tree based on the values (numbers, strings)
 - the tuples themselves are not stored in the tree
 - entries (or leaf entries only) hold the values *and* point to the respective tuples in the database
- special trees with higher degrees:
 - each node (of the size of a storage page) has multiple entries and multiple children.

133

B- AND B*-TREES

A **B-tree** (R. Bayer & E. McCreight, 1970) of order (m, ℓ) , $m \geq 3$, $\ell \geq 1$, is a search tree [see lecture on Algorithms and Data Structures]:

- the root is either a leaf or it has at least 2 children
- every inner node has at least $\lceil m/2 \rceil$ and at most m children
- all leaves are on the same level (balanced tree), and hold at most ℓ entries
- inner nodes have the form $(p_0, k_1, p_1, k_2, p_2, \dots, k_n, p_n)$ where $\lceil m/2 \rceil - 1 \leq n \leq m - 1$ and
 - k_i are search key values, ordered by $k_i < k_j$ for $i < j$
 - p_i points to the $i + 1$ th child
 - all search key values in the left (i.e., p_{i-1}) subtree are less than the value of k_i (and all values in the right subtree are greater or equal)

B-trees are used for “simply” organizing items of an ordered set (e.g. for sorting) as an extension of binary search trees.

134

B*-TREES

(sometimes also called B⁺-Trees)

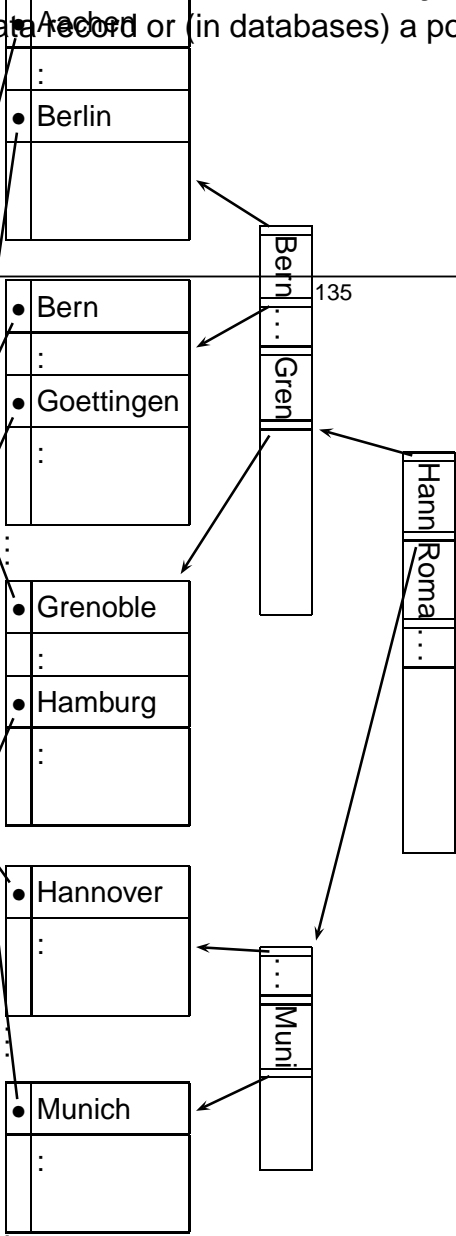
A **B*-tree** of order (m, ℓ) , $m \geq 3$, $\ell \geq 1$, is closely related, except:

- they are intended to associate *data* with search key values:
- the inner nodes do not hold additional data, but are still intended to guide the search.

• The leaves are of the form $([k_1, s_1], [k_2, s_2], \dots, [k_g, s_g])$ where $g \leq \ell$, k_i is a search key value, and s_i is a data record or (in databases) a pointer to the corresponding record.

Example B*-Tree over city.name

| | | |
|-----------|----|---------|
| : | : | : |
| Munich | D | 1244676 |
| : | : | : |
| Grenoble | F | 150758 |
| : | : | : |
| Aachen | D | 247113 |
| : | : | : |
| Bern | CH | 134393 |
| : | : | : |
| Göttingen | D | 127519 |
| : | : | : |
| Hannover | D | 525763 |
| : | : | : |
| Berlin | D | 3472009 |
| : | : | : |
| Hamburg | D | 1705872 |
| : | : | : |



Properties

- Let N the number of entries. Then, for the height h of the tree,
 $h \leq \lceil \log_{m/2}(2N/\ell) \rceil$ ($\ell/2$ entries per leaf, inner nodes half filled) and
 $h \geq \lceil \log_m(N/\ell) \rceil$ (ℓ entries/leaf, inner nodes completely filled)
- equality search needs h steps
inside each of the inner nodes, binary search is applied
- if the leaves are connected by pointers, ordered sequential access (range search) is also supported
- insertions and modifications may be expensive (tree reorganization)

Use of B*-Trees as Access Paths in Databases

- databases: cities stored in data files, index trees hold *pointers* to city records in their tree entries.
- separation between index files/pages and data files/pages.
- multiple search trees for each relation possible.

137

Example

Example 4.1

Consider the MONDIAL database with 3000 cities, with an index over the name. Assume the following sizes:

- every leaf (tuple) page contains 10 cities,
- every inner node contains 20 pointers

Then

- every inner node on the lowest level covers 200 cities
- every inner node on the second lowest level covers 4000 cities
- minimal: only one level of inner nodes
- maximal: two levels of inner nodes (nodes about only 2/3 filled)
- access every city with *WHERE Name = "..."* in 3 or 4 steps
- index on population, e.g., for *WHERE Population > 1.000.000 ORDER BY Population*
- realistic numbers: block size 4K: lowest level (keys+pointers to DB): 100 cities; inner nodes: 100 references.

□

138

HASH-INDEX (DICTIONARY)

Hash index over the value of a column:

The values are distributed over k **tiles**.

- A **hash function** h is a function that maps each value to a tile number.
operation: lookup(value)
- each tile holds pointers to all tuples whose column value is mapped to this tile; each tile consists of one or more pages.
- common technique: convert value to an integer i . $i \bmod k$ gives the tile number.

Example 4.2

Consider a hash index on $City.(name, province, country)$. h computes the sum of the ASCII numbers of the letters and takes the remainder modulo 111. □

Properties:

- equality search and insert in constant time (+ time for searching in the tile)
- does not support range queries or ordered output

139

BIT LISTS

If the number of possible values of a search key value is small wrt. the number of tuples, **bit lists** are useful as indexes

Let a an attribute of the records stored in a file f , where k values of a exist. Then, a bit list index to f consists of k bit vectors $B(v_i)$, $i \leq 1 \leq k$.

- $B(v_i)(j) = 1$ if the j th record in f has the value v_i for the attribute a .

Properties:

- access all tuples with a given value:
without bit list: linear search over all pages
with bit list: access bit list, and access those pages where a "1"-tuple is located.
- modifications: no problem
- deletions: depends (if gaps are allowed on the pages)

140

Example 4.3

Consider the relation $is_member(organization, country, type)$ where $type$ has only the values “member”, “applicant”, and “observer”:

| <i>is_member</i> | | |
|------------------|----------------|------------------|
| <i>Org.</i> | <i>Country</i> | <i>Type</i> |
| <i>EU</i> | <i>D</i> | <i>member</i> |
| <i>UN</i> | <i>D</i> | <i>member</i> |
| <i>UN</i> | <i>CH</i> | <i>observer</i> |
| <i>UN</i> | <i>R</i> | <i>member</i> |
| <i>EU</i> | <i>PL</i> | <i>applicant</i> |
| <i>EU</i> | <i>CZ</i> | <i>applicant</i> |
| <i>:</i> | <i>:</i> | <i>:</i> |

Bit list for type=member: 1 1 0 1 0 0 ...

Bit list for type=observer: 0 0 1 0 0 0 ...

Bit list for type=applicant: 0 0 0 0 1 1 ...

Bit list for org=EU: 1 0 0 0 1 1 ...

Bit list for org=UN: 0 1 1 1 0 0 ...

- Search for members of the UN:
without bit list: linear search over all pages
with bit list: access bit list, and access those pages where a 1-tuple is located.
- 2nd bit list on organization column: “members of UN” by logical “and” of bit lists. □

141

CLUSTERING

Data from the secondary storage is fetched by **pages**.

Pages are cached in the main memory.

- keep data that semantically belongs together on the same pages
- obviously done for relations
- data of several relations $R_1(\bar{X}_1), \dots, R_n(\bar{X}_n)$ can also be grouped to **clusters**:
 - choose a set $\bar{Y} \subseteq \bar{X}_1 \cap \dots \cap \bar{X}_n$ as **cluster key**.
 - combine relations by their \bar{Y} -values.
 - provides obvious advantages when evaluating joins over the cluster key.

142

Example 4.4

Consider the relation schemata

organization(name,abbrev,established,...)

and

is_member(organization, country, type).

The foreign key

is_member.organization → *organization.abbrev*

is used as cluster key.

| Organization | | |
|--------------|----------------|--------------------|
| Abbrev | | |
| EU | Name | established |
| | European Union | 07.02.1992 |
| | Country | Type |
| | D | member |
| | A | member |
| | : | : |
| UN | Name | established |
| | United Nations | 26.06.1945 |
| | Country | Type |
| | D | member |
| | CH | observer |
| | : | : |

4.2 Efficient Query Evaluation

Queries are formulated *declaratively* (e.g., SQL or algebra trees), actually built over a small set of basic operations (cf. the definition of the relational algebra).

Semantical optimization: consider integrity constraints in the database.

Example: *population* > 0, thus, a query that asks for negative values can be answered without explicit computation.

- not always obvious
- general case: first-order theorem proving.
- special cases: [see lecture on Database Theory]

Logical/algebraic optimization: search for an equivalent algebra expression that performs better:

- size of intermediate results,
- implementation of operators as algorithms,
- presence of indexes and order.

BASIC TECHNIQUES

Iteration: all tuples in the input relations are processed iteratively. If possible, instead of the tuples themselves, an index can be used.

Indexes: selections and joins can be based on processing of an index for determining the tuples that satisfy the selection condition (a join condition is also a selection condition).

Indexes on single attributes: obvious.

Indexes on multiple attributes:

- a **hash index** to a conjunction of **equality** predicates of the form **field = value** can be used if every field of the search key occurs exactly in one predicate together with a constant,
- a **tree index** to a conjunction of **comparison** predicates of the form **field θ value** can be used if a prefix of the search key exists such that each field of the prefix occurs in exactly one predicate together with a constant.

Example: if (Country,Province,CityName) is the search key of an index, it can also be used as an index for the key's prefix (Country,Province).

145

4.2.1 Algebra Operations

SELECTION

Selection: $\sigma[R.field \theta value]R$.

- no index, unordered tuples: linear scan of the file
 - no index, tuples ordered wrt. *field*: find the tuple(s) that satisfy "*field θ value*" (binary search) and process the these tuples.
 - tree index: find the tuple(s) that satisfy "*field θ value*" using the tree index and process these tuples.
 - hash: suitable only if θ is equality.
- + strategies for compound selection conditions.

146

PROJECTION

Projection: $\pi[\text{field}_1, \text{field}_2, \dots, \text{field}_m]R$.

Main problem: remove duplicates (relational algebra does not allow duplicates, SQL does).

- by sorting:
 - scan the file, create a new file with projected tuples.
 - sort the new file (over all fields, $n \log n$).
 - scan the result, remove duplicates.
 - by hash: scan the file, put the projected tuples into a hash table.
 - duplicates all end up in the same tile (still in files).
 - remove duplicates separately for each tile (iterating the hashing process with different functions until tile fits in main memory)
 - collect the tiles in the output file.
- ⇒ hashing not only as an index for search *structures*, but also as an *algorithm* for partitioning.

147

JOIN

Consider an Equijoin: $R \bowtie_{R.A=S.B} S$.

Nested-loop-join

```
foreach tuple  $r \in R$  do
  foreach tuple  $s \in S$  do
    if  $r_A = s_B$  then add  $[r, s]$  to result
```

Very inefficient:

- assume: only one page of each relation fits into main memory
- m tuples of R per page, n tuples of S per page:
 - $|R|/m + |R| \cdot |S|/n$ filesystem accesses (for each tuple of R , loop over all pages of S).
- $|R| \cdot |S|$ comparisons.

obvious: if possible, process the smaller relation in inner loop to keep it in main memory ($|R|/m + |S|/n$ filesystem accesses).

148

Block-nested-loop-Join

Optimization by page-oriented strategy:

Divide S into blocks that fit in main memory and process each of the blocks separately.

- assume: only one page of each relation fits into main memory
- m tuples of R per page, n tuples of S per page:
 - take first page (R_1) of R and first page (S_1) of S ;
 - combine *each* tuple of R_1 with each tuple with S_1 .
 - continue with R_1 and second page (S_2) of S ... up to R_1 with S_{last}
 - continue with R_2 and S_1 etc.
- $|R|/m + |R|/m \cdot |S|/n$ filesystem accesses (for each page of R , loop over all pages of S).
- still $|R| \cdot |S|$ comparisons.

Further optimization:

when a page of S is loaded, generate a temporary index over the join attribute(s) that is used in the comparison (then, only matching pairs are considered).

149

Index-nested-loop-Join

If for one of the input relations there is an index over the join attribute, choose it as the inner relation.

- given index on S 's join attribute
- loop over R , for each value access matching tuple(s) in S .
- $|R|/m + |Result|$ filesystem accesses,
- up to $|R|/m + |R| \cdot |S|$ filesystem accesses,
- maybe less efficient than block-nested-loop!

Parallelization

Divide one of the relations into partitions and process each of them on a separate processor.

150

Sort-merge-join

Algorithm type: "Scan line"

- compute sorted indexes of both relations on the join attributes
(if tree indexes on join attributes are available, simply use them)
- search for matches as follows:
 - proceed through the ordered indexes R and S stepwise, always doing a step in the index where the "smaller" value is.
- if a match is found:
 - access tuples and generate a result tuple.
 - check for tuples which have the same values of the join attributes (must immediately follow this match in both relations).
- $|R|/m + |S|/n + |Result|$ filesystem accesses.
- much less, if indexes are already available or relations are *stored* in this order.

151

Sort-merge-join: Algorithm

```
if R not sorted on attribute A, sort it;
if S not sorted on attribute B, sort it;
Tr := first tuple in R;
Ts := first tuple in S;
Gs := first tuple in S;
while Tr ≠ eof and Ts ≠ eof do {
    while Tr ≠ eof and Tr.A < Ts.B do Tr = next tuple in R after Tr;
    while Ts ≠ eof and Tr.A > Ts.B do Ts = next tuple in S after Ts;
    // now, Tr.A = Ts.B: match found
    while Tr ≠ eof and Tr.A = Ts.B do {
        Gs := Ts;
        while Gs ≠ eof and Gs.B = Tr.A do {
            add [Tr, Gs] to result;
            Gs = next tuple in S after Gs;}
        Tr = next tuple in R after Tr;
    }
    Ts := Gs;
}
```

152

Hash-join

Algorithm type: “Divide and Conquer”

Partitioning (building) phase:

- Partition the smaller relation R by a hash function h_1 applied to $R.A$.
- Partition the larger relation S by the same hash function applied to $S.B$.
(into different hash tables)

Matching (probing) phase:

- potential(!) matches have been mapped to “face-to-face” partitions.
- thus, consider each pair of corresponding partitions:
- if partition of R does not fit into main memory, proceed recursively with another h .
- otherwise – i.e., R -partition fits in main memory:
 - if corresponding S -partition also fits into main memory, compute the join.
 - otherwise proceed recursively with another (finer) h_2 inside main memory and process S -partition pagewise.

153

EXERCISE

Exercise 4.1

Consider the join of two relations R and S wrt. the join condition $R_i = S_j$. R uses M pages with p_R tuples each, and S uses N pages with p_S tuples each. Let $M = 1000$, $p_R = 100$, $N = 500$, $p_S = 80$; the cache can keep 100 pages. Compute the number of required I/O-operations for computing the join (without writing the result relation) for:

- simple nested-loop-join,
- pagewise nested-loop-join,
- index-nested-loop-join (Index on S_j),
- sort-merge-join,
- hash-join.

□

154

GROUPING AND AGGREGATION

General Structure:

SELECT A_1, \dots, A_n list of attributes
FROM R_1, \dots, R_m list of relations
WHERE F condition(s)
GROUP BY B_1, \dots, B_k list of grouping attributes
HAVING G condition on groups
ORDER BY H sort order

- SUM, AVG, COUNT: linear scan necessary (need to consider all tuples)
- MIN, MAX: index can be used
- support for grouping (SQL: GROUP BY) using an index or a hash table

UNION, DIFFERENCE, INTERSECTION \Leftrightarrow CARTESIAN PRODUCT

- intersection and cartesian product are special cases of join
- union and difference: analogous to sort-merge-join or hash-join

155

4.2.2 Algebraic Optimization

The operator tree of an algebra expression provides a base for several optimization strategies:

- reusing intermediate results
- equivalent restructuring of the operator tree
- “shortcuts” by melting several operators into one (e.g., join + equality predicate \rightarrow equijoin)
- combination with actual situation: indexes, properties of data

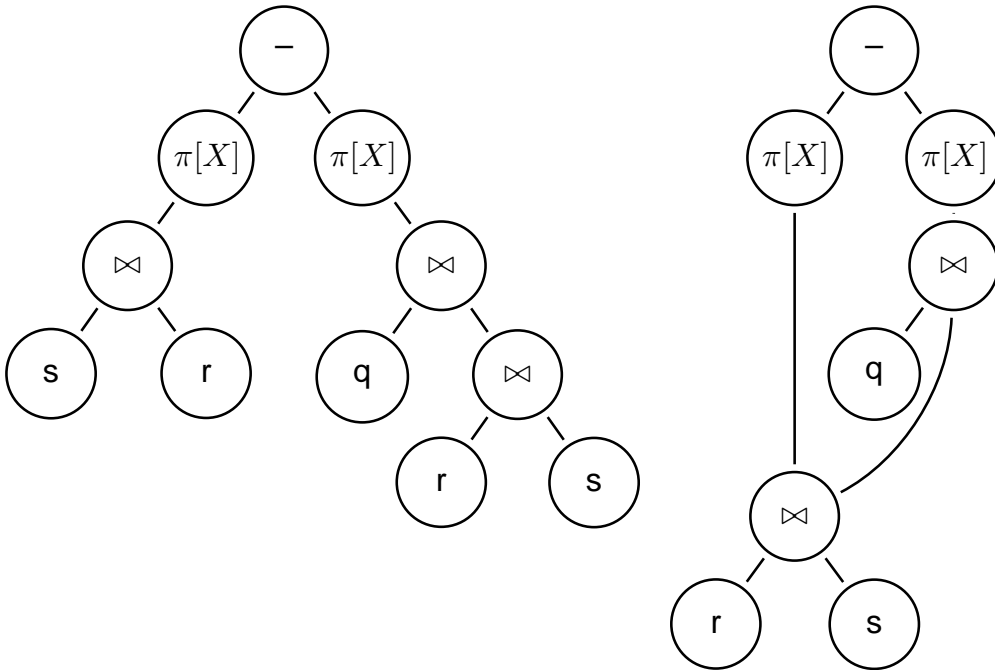
Real-life databases implement this functionality.

- SQL: **declarative** specification of a query
- internal: algebra tree + optimizations

156

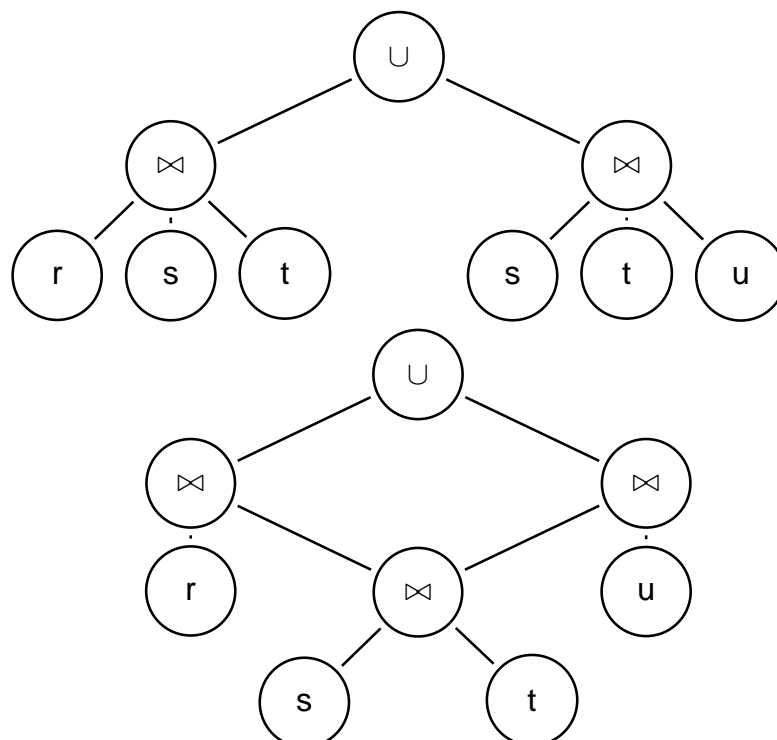
REUSING INTERMEDIATE RESULTS

- Multiply occurring subtrees can be reused
(directed acyclic graph (DAG) instead of algebra tree)



157

Reusing intermediate results



158

OPTIMIZATION BY TREE RESTRUCTURING

- Equivalent transformation of the operator tree that represents an expression
- Based on the equivalences shown on Slide 101.
- minimize the size of intermediate results
(reject tuples/columns as early as possible during the computation)
- selections reduce the number of tuples
- projections reduce the size of tuples
- apply both as early as possible (i.e., before joins)
- different application order of joins

159

Push Selections Down

Assume $r, s \in \text{Rel}(\bar{X})$, $\bar{Y} \subseteq \bar{X}$.

$$\sigma[\text{cond}](\pi[\bar{Y}](r)) \equiv \pi[\bar{Y}](\sigma[\text{cond}](r))$$

(condition: *cond* does not use attributes from $\bar{X} - \bar{Y}$,
otherwise left term is undefined)

$$\sigma_{\text{pop} > 1E6}(\pi[\text{name, pop}](\text{country})) \equiv \pi[\text{name, pop}](\sigma_{\text{pop} > 1E6}(\text{country}))$$

$$\sigma[\text{cond}](r \cup s) \equiv \sigma[\text{cond}](r) \cup \sigma[\text{cond}](s)$$

$$\begin{aligned} \sigma_{\text{pop} > 1E6}(\pi[\text{name, pop}](\text{country}) \cup \pi[\text{name, pop}](\text{city})) \\ \equiv \sigma_{\text{pop} > 1E6}(\pi[\text{name, pop}](\text{country})) \cup \sigma_{\text{pop} > 1E6}(\pi[\text{name, pop}](\text{city})) \end{aligned}$$

$$\sigma[\text{cond}](\rho[N](r)) \equiv \rho[N](\sigma[\text{cond}'](r))$$

(where *cond'* is obtained from *cond* by renaming according to *N*)

$$\sigma[\text{cond}](r \cap s) \equiv \sigma[\text{cond}](r) \cap \sigma[\text{cond}](s)$$

$$\sigma[\text{cond}](r - s) \equiv \sigma[\text{cond}](r) - \sigma[\text{cond}](s)$$

π : see comment above. Optimization uses only left-to-right.

160

Push Selections Down (Cont'd)

Assume $r \in \text{Rel}(\bar{X})$, $s \in \text{Rel}(\bar{Y})$. Consider $\sigma[\text{cond}](r \bowtie s)$.

Let $\text{cond} = \text{cond}_{\bar{X}} \wedge \text{cond}_{\bar{Y}} \wedge \text{cond}_{\overline{\bar{X}\bar{Y}}}$ such that

- $\text{cond}_{\bar{X}}$ is concerned only with attributes in \bar{X}
- $\text{cond}_{\bar{Y}}$ is concerned only with attributes in \bar{Y}
- $\text{cond}_{\overline{\bar{X}\bar{Y}}}$ is concerned both with attributes in \bar{X} and in \bar{Y} .

Then,

$$\sigma[\text{cond}](r \bowtie s) \equiv \sigma[\text{cond}_{\overline{\bar{X}\bar{Y}}}] (\sigma[\text{cond}_{\bar{X}}](r) \bowtie \sigma[\text{cond}_{\bar{Y}}](s))$$

Example 4.5

Names of all countries that have been founded earlier than 1970, their capital has more than 1.000.000 inhabitants, and more than half of the inhabitants live in the capital. □

161

Example 4.5 (Continued)

(Solution)

$$\begin{aligned} & \pi[\text{Name}] (\sigma[\text{establ} < \text{"01 01 1970"} \wedge \text{city.pop} > 1.000.000 \wedge \text{country.pop} < 2 \cdot \text{city.pop}] \\ & \quad (\text{country} \times_{\text{country}.\text{(capital,prov,code)}=\text{city}(\text{name,prov,country})} \text{city}) \\ & \equiv \pi[\text{Name}] (\sigma[\text{country.pop} < 2 \cdot \text{city.pop}] \\ & \quad (\sigma[\text{establ} < \text{"01 01 1970"}](\text{country}) \\ & \quad \quad \times_{\text{country}.\text{(capital,prov,code)}=\text{city}(\text{name,prov,country})} \\ & \quad \quad \quad \sigma[\text{city.pop} > 1.000.000](\text{city}))) \end{aligned}$$

□

- Nevertheless, if cond is e.g. a complex mathematical calculation, it can be cheaper first to reduce the number of tuples by \cap , $-$, or \bowtie

⇒ data-dependent strategies (see later)

162

Push Projections Down

Assume $r, s \in \text{Rel}(\bar{X})$, $\bar{Y} \subseteq \bar{X}$.

Let $cond = cond_{\bar{X}} \wedge cond_{\bar{Y}}$ such that

- $cond_{\bar{Y}}$ is concerned only with attributes in \bar{Y}
- $cond_{\bar{X}}$ is the remaining part of $cond$ that is also concerned with attributes $\bar{X} \setminus \bar{Y}$.

$$\pi[\bar{Y}](\sigma[cond](r)) \equiv \sigma[cond_{\bar{Y}}](\pi[\bar{Y}](\sigma[cond_{\bar{X}}](r)))$$

$$\pi[\bar{Y}](\rho[N](r)) \equiv \rho[N](\pi[\bar{Y}'](r))$$

(where \bar{Y}' is obtained from \bar{Y} by renaming according to N)

$$\pi[\bar{Y}](r \cup s) \equiv \pi[\bar{Y}](r) \cup \pi[\bar{Y}](s)$$

- Note that this does *not* hold for “ \cap ” and “ $-$ ”!
- advantages of pushing “ σ ” vs. “ π ” are data-dependent
Default: push σ lower.

Assume $r \in \text{Rel}(\bar{X})$, $s \in \text{Rel}(\bar{Y})$.

$$\pi[\bar{Z}](r \bowtie s) \equiv \pi[Z](\pi[\bar{X} \cap \bar{Z}\bar{Y}](r) \bowtie \pi[\bar{Y} \cap \bar{Z}\bar{X}](s))$$

- complex interactions between reusing subexpressions and pushing selection/projection

163

Application Order of Joins

Minimize intermediate results:

```
SELECT organization.name, country.name
FROM organization, country, is_member
WHERE organization.abbrev = is_member.organization
      AND country.code = is_member.country
```

Exploit selectivity of join:

- $\underbrace{(\text{org} \times \text{country})}_{200 \cdot 200 = 40000} \bowtie \text{is_member}$
7000

- $\underbrace{(\text{org} \bowtie \text{is_member})}_{200, 7000 \rightsquigarrow 7000} \bowtie \text{country}$
7000

If indexes on `country.code` and `organization.abbrev` are available:

- loop over `is_member`
- extend each tuple with matching `country` and `organization` by using the indexes.

164

Example/Exercise

Consider the equivalent (to the previous example) query:

```
SELECT organization.name, country.name
FROM organization, country,
WHERE EXISTS
  (SELECT *
   FROM is_member
   WHERE organization.abbrev = is_member.organization
        AND country.code = is_member.country)
```

- suggests the non-optimal evaluation!
- transform the above query into algebra
- ... yields the same “broad” join as before ...
- ... and leads to the optimized join ordering.

165

OPERATOR EVALUATION BY PIPELINING

- above, each algebra operator has been considered separately
- if a query consists of several operators, the materialization of intermediate results should be avoided
- **Pipelining** denotes the immediate propagation of tuples to subsequent operators

Example 4.6

- $\sigma[A = 5 \wedge B > 6]R$:

Assume an index that supports the condition $A = 5$.

- without pipelining: compute $\sigma[A = 5]R$ using the index, obtain R' . Then, compute $\sigma[B > 6]R'$.
- with pipelining: compute $\sigma[A = 5]R$ using the index, and check **on-the fly** each qualifying tuple against $\sigma[B > 6]$. □

- **Unary** (i.e., selection and projection) operations can always be pipelined with the next lower binary operation (e.g., join)

166

Example 4.6 (Continued)

- $\sigma[\text{cond}](R \bowtie S)$:
 - *without pipelining: compute $R \bowtie S$, obtain RS , then compute $\sigma[\text{cond}](RS)$.*
 - *with pipelining: during computing $(R \bowtie S)$, each tuple is immediately checked whether it satisfies cond .*
- $(R \bowtie S) \bowtie T$:
 - *without pipelining: compute $R \bowtie S$, obtain RS , then compute $RS \bowtie T$.*
 - *with pipelining: during computing $(R \bowtie S)$, each tuple is immediately propagated to one of the described join algorithms for computing $RS \bowtie T$.* □

Most database systems use iterator-based implementation of algebra operators with pipelining.

- no materialization of intermediate results
 - ⇒ no indexes on them, no sorting, no hashing
 - index-based joins only if a base relation is joined
 - look for order-preserving algorithms, then “merge”-based joins are still possible on higher levels.

167

SITUATION-DEPENDENT OPTIMIZATION

Determining the optimal **execution plan** depends on **cost models**, **heuristics**, and the actual physical schema and the actual database contents:

- index structures
- statistics on data (i.e., to some extent state dependent)
 - cardinality of relations
 - distribution of values

168

Notions

The **selectivity** sel of operations can be used for estimating the size of the result relation:

- Selection with condition p :

$$sel_p = \frac{|\sigma[p](R)|}{|R|}$$

Proportion of tuples that satisfy the selection condition.

- for p an equality test $R.A = c$ where A is a key of R , $sel_p = 1/|R|$.
 - if $R.A$ distributes evenly on i different values, $sel_p = 1/i$.
- Join of R and S :

$$sel_{RS} = \frac{|R \bowtie S|}{|R \times S|} = \frac{|R \bowtie S|}{|R| \cdot |S|}$$

is the proportion of result tuples wrt. the cartesian product.

for $R \bowtie_{R.A=S.B} S$ with A a key attribute, $|R \bowtie_{R.A=S.B} S| \leq |S|$, thus $sel_{RS} \leq 1/|R|$.

- optimal order of applying joins
- optimal order of applying selections

169

APPLICATION-LEVEL ALGORITHMIC OPTIMIZATION

And **never** forget about using efficient algorithms for querying the database!

- analyze the problem from the **algorithmic** point of view
- before hacking

EXAMPLES

Use the MONDIAL database for the following examples.

Example 4.7

Compute all pairs of european countries that are adjacent to the same set of seas. □

Example 4.8

Compute all political organizations that have at least one member country on every continent (this operation is called relational division). □

170

4.2.3 Exercises and Examples

Solution to Exercise 4.1

R: 1000 pages, each of them with 100 tuples

S: 500 pages, each of them with 80 tuples

Cache: 100 pages

a) simple nested loop join

– outer loop: *R*

Load each page of *R*. For all tuples, iterate over *S*'s pages.

$$1000 \cdot (1 \text{ (load)} + 100 \text{ (tuples per page)} \cdot 500 \text{ (pages of } S)) = 50,001,000$$

– outer loop: *S*

Load each page of *S*. For all tuples, iterate over *R*'s pages.

$$500 \cdot (1 \text{ (load)} + 80 \text{ (tuples per page)} \cdot 1000 \text{ (pages of } R)) = 40,000,500$$

171

Solution to Exercise 4.1 (cont'd)

R: 1000 pages, each of them with 100 tuples

S: 500 pages, each of them with 80 tuples

Cache: 100 pages

b) pagewise nested loop:

– outer loop: *R*. Load each page of *R*. Combine all tuples of that page with all tuples from each page of *S*.

$$1000 \cdot (1 + 500 \text{ (pages of } S)) = 501,000$$

– outer loop: *S*. $500 \cdot (1 + 1000 \text{ (pages of } R)) = 505,000$

b2) maximum-pages nested loop: load as many (first 99) pages of *R* as possible in the cache and join with one page of *S* after the other. Then continue with 2nd 99 pages.

– 11 times (10 · 99 *R* pages + 1 · 10 *R* pages) through all *S* pages.

$$\text{Overall: } 1000 + 11 \cdot 500 = 6500.$$

– symmetric with (5 · 99 + 1 · 5) pages of *S*:

$$\text{Overall: } 500 + 6 \cdot 1000 = 6500.$$

– can be algorithmically optimized by an on-the-fly index during the main-memory join: index the 99 pages, loop over the single page.

172

Solution to Exercise 4.1 (cont'd)

R : 1000 pages, each of them with 100 tuples

S : 500 pages, each of them with 80 tuples; Cache: 100 pages

c) Index-nested-loop (index on S_j): Assumptions:

- every R_i matches – average – $k = 4$ S tuples times.
- index over S already exists and can be kept in main memory (creating an index over S makes 500 accesses)

Iterate over R , for each tuple search for the S_j key and access the tuples

$$1000 \cdot (1 \text{ (load)} + 100 \cdot 4 \text{ (access tuples)}) = 401,000.$$

Note: the number of page accesses depends on the number of results since for every actual result there is one page access.

$$\text{In case that } S \text{ is ordered wrt. } S_j: 1000 \cdot (1 \text{ (load)} + 100 \cdot 1 \text{ (access tuples)}) = 101,000.$$

d) sort-merge: sort R according to $R.i$:

sort each 100 pages in place (Quicksort), store them:

$$10 \text{ (times)} \cdot 100 \text{ (pages)} \cdot 2 \text{ (read and write)} = 2000$$

mergesort the sorted pages (linear) – read and write: 2000

sort S according to S_j analogously: 2000

merge-join of sorted relations: linear scan: $1000 + 500 = 1500$

Overall: 7500

173

Solution to Exercise 4.1 (cont'd)

R : 1000 pages, each of them with 100 tuples

S : 500 pages, each of them with 80 tuples

Cache: 100 pages

e) Hash-Join:

– Hash R :

100 pages fit in memory - read 34 pages and distribute over 60 partitions. Do this 30 times.

$$\text{(about } 30 \cdot (34 \text{ (read)} + 60 \text{ (write)}) < 3000 \text{ accesses).}$$

In the average, each partition them contains about 17 pages.

– Same for S . Get 60 partitions of average 9 pages. (about 1500 accesses).

– now we have 60 corresponding pairs of partitions (one from R , one from S).

Join each pair: read $17 + 9$ (overestimate) pages per partition and process them.

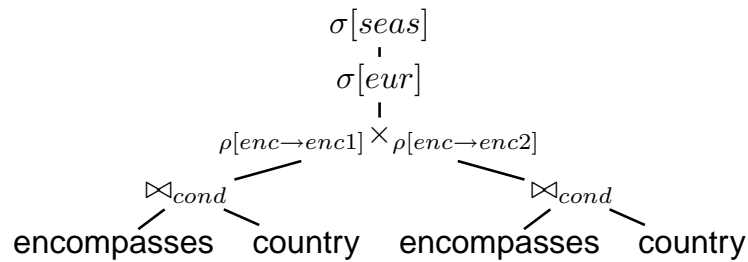
$$\text{(about } 60 \cdot (17 + 9) = 1560 \text{ accesses).}$$

Overall: 6060 accesses.

174

Solution to Exercise 4.7

First, compute all pairs of european countries. Note that we are interested in their names, thus we need also the *Country* relation.



- $cond = \text{"enc.country=country.code"}$.
- $eur = \text{"enc1.continent='europe' and enc2.continent='europe'"}$.
- $seas$ expresses that both countries border the same set of seas. It is a correlated subquery.
- add suitable projections.
- move $\sigma[eur]$ downwards both sides directly to *encompasses*
- obviously, both subtrees of \times are identical.

175

Solution to Exercise 4.7 (cont'd)

- $\sigma[seas(C_1, C_2)]$ is a correlated subquery that takes two country codes as input:

$$\begin{aligned} \sigma[seas(C_1, C_2)] &= seas(C_1) - seas(C_2) = \emptyset \wedge seas(C_2) - seas(C_1) = \emptyset \\ &= (seas(C_1) - seas(C_2)) \cup (seas(C_2) - seas(C_1)) = \emptyset \end{aligned}$$
- $seas(C) = \pi[sea](\sigma[country = C](geo_sea))$
- for each country, $seas(C)$ is computed only once and then reused.

Resulting SQL skeleton (using subqueries in the FROM clause):

```
SELECT ...
FROM (SELECT european countries) as C1,
     (SELECT european countries) as C2
WHERE  $\sigma[seas(C_1, C_2)]$ 
```

176

Solution to Exercise 4.8

- $\pi[abbrev](\sigma[all_continents(org)])(organization)$ where $\sigma[all_continents(org)]$ is true for
 $\{org \mid \forall cont : cont \text{ is a continent} \rightarrow org \text{ has a member on } cont\}$
- convert \forall into $\neg exists$:
 $\{org \mid \neg \exists cont : cont \text{ is a continent and } org \text{ has no member on } cont\}$
- Thus, $\sigma[all_continents(org)]$ checks if
 $\pi[name](continent) - \pi[enc.continent](\sigma[organization = org]is_member) \bowtie encompasses$
is empty.

Resulting SQL skeleton (uses a correlated subquery):

```
SELECT ...
FROM organization
WHERE NOT EXISTS
  ((SELECT continents)
   MINUS
   (SELECT continents where org has a member))
```