

## Assignment 4

### 100 points

Professor Yanlei Diao

#### Question 1: OLTP versus OLAP Applications

Please indicate for each of the following applications, whether it is more suitable to use an OLTP database or an OLAP database. Briefly explain your answer if needed.

- (a) Banking
- (b) Ticketing
- (c) Telecommunications: e.g., phone location updates and lookups
- (d) E-commerce: searching and buying books
- (e) Retailing: user profiling, promotion strategy analysis
- (f) Financial services: claims analysis, risk analysis, credit card fraud detection
- (g) Telecommunications: call pattern analysis, fraud detection
- (h) Utilities: power usage analysis

#### Question 2: OLAP Operations

Consider the instance of Sales relation shown below.

<i>pid</i>	<i>timeid</i>	<i>locid</i>	<i>sales</i>
11	1	1	25
11	2	1	8
11	3	1	15
12	1	1	30
12	2	1	20
12	3	1	50
13	1	1	8
13	2	1	10
13	3	1	10
11	1	2	35
11	2	2	22
11	3	2	10
12	1	2	26
12	2	2	45
12	3	2	20
13	1	2	20
13	2	2	40
13	3	2	5

- (a) Show the result of pivoting the relation on *pid* and *timeid*.
- (b) Write a collection of SQL queries to contain the same result as in Part (a).
- (c) Show the result of pivoting the relation on *pid* and *locid*.

### Question 3: Implementation of the Cube Operator

Consider a relation SALES consisting of 19 dimension attributes,  $A_1, A_2, A_3, \dots$  and 1 numerical measure, Sales.

```
SALES (A1, A2, A3, ..., A19, sales)
```

Assume that all attributes are of the equal length.

We would like to run the following cube operator:

```
Select A1, A2, A3, sales  
From SALES  
Group by Cube (A1, A2, A3) ;
```

Assume that the relation contains  $N$  pages and we have  $B$  buffer pages to run this query, where  $\sqrt{N} < B < N$ . Please describe the most efficient algorithm to run this query and analyze its I/O cost.

**Hint:** please first think how to run: Group SALES By  $A_1, A_2, A_3$ . Then add other Group By operations covered by the cube.

### Question 4: Query Rewriting

(1) Define a schema  $S$ , a database instance  $D$ , a set of views  $V$ , and a query  $Q$  such that  $Q(D)$  under the closed-world assumption given the extensions of the views in  $V$  is different from  $Q(D)$  under the open-world assumption given the extensions of the views in  $V$ .

(2) Consider the following two questions:

(a) Find a query  $Q$ , a query language  $L$  and a set of views  $V$  such that there is no maximally contained rewriting of  $Q$  using  $V$  with respect to  $L$ .

(b) For the same  $Q$  and  $V$ , find another language  $L'$  for which there exists a maximally contained rewriting of  $Q$  using  $V$  with respect to  $L'$ .

### Question 5: Windows and Materialized Views in PostgreSQL

In this problem set, we continue to use the DBLP publication database containing information of over 3 million papers published in computer science conferences and journals (this data was derived from the DBLP system, maintained by Michael Ley at <http://www.informatik.uni-trier.de/~ley/db/>).

\* **Schema.** This database consists of five tables: (1) an authors table, containing the names of authors, (2) the venue table, containing information about conferences or journals where papers are published, (3) the papers table, describing the papers themselves, (4) the paperauths table which indicates which authors wrote which papers, and (5) the papertypes. The schema is the following:

```
authors (id: INTEGER NOT NULL, name: VARCHAR(200))
```

```
venue (id: INTEGER NOT NULL, name: VARCHAR(200) NOT NULL, year: INTEGER NOT NULL, school:  
VARCHAR (200), volume: VARCHAR(50), number: VARCHAR(50), type: INTEGER NOT NULL)
```

```
papers (id: INTEGER, name: VARCHAR NOT NULL, venue: INTEGER REFERENCES VENUE(id),  
pages: VARCHAR(50), url: VARCHAR);
```

```
paperauths (paperid: INTEGER REFERENCES PAPERS(id), authid: INTEGER REFERENCES AUTHORS(id))
```

```
papertypes (id: INTEGER NOT NULL, type: VARCHAR(20) NOT NULL)
```

\* **Indexes.** Your database has **at least** the following indexes:

- All primary key indexes,
- A B+ tree on the <name> attribute of the **authors** table, and
- A B+ tree on the <paperid, authid> attributes of the **paperauths** table.

If you have created other indexes, please drop them using the `drop index 'index-name'` command.

\* **Useful commands:**

mysql: **SHOW TABLES**

postgresql: `\d`

mysql: **SHOW INDEXES**

postgresql: `\di`

mysql: **SHOW COLUMNS**

postgresql: `\d table`

mysql: **DESCRIBE TABLE**

postgresql: `\d+ table`

\* **Materialized Views.** PostgreSQL supports materialized views using the following command:

<http://www.postgresql.org/docs/current/static/sql-creatematerializedview.html>

\* **Window Functions.** PostgreSQL supports the window functions as defined below:

<http://www.postgresql.org/docs/current/static/tutorial-window.html>

### (1) Query writing using Views:

Define a single view that allows to rewrite the following type queries:

- Find the name of the n authors having the maximum number of publications (n integer) ;
- Find the n years when a maximum number of author have published ;
- Find the average on the number of year on how many authors published exactly once in that year.

Try to minimize the size (number of rows and of columns) of the defined view.

### (2) Query writing using Views:

Let us consider the following query:

Find the venue where the author with id 17060 has published most.

Can you answer that query using only the above defined view? If yes, express it. Otherwise, prove this is not possible.

### (3) Window Functions

Please write the following queries using SQL99 with the window function. Please submit both your SQL queries and answers of running these queries over our DBLP database.

**TIP:** The number of pages of a paper is stored in a VARCHAR field. We need to perform some data conversion to be able to get an INT value for it. In addition, this field includes different formats of pages (range of pages: *page1 – page2*, punctual noncontiguous pages: *page1, page3*, etc). To keep things simple, let's just consider the case of ranges of pages (which is the majority). You can run this query to successfully convert VARCHAR ranges of pages to the total number of pages for all papers.

```
SELECT
CASE
  WHEN p.pages ~ E'^\\d+\\d+$' THEN
    NULLIF(SPLIT_PART(p.pages, '-', 2), ' '::int - NULLIF(SPLIT_PART(p.pages, '-', 1), ' '::int)
  ELSE 0
END
FROM papers p;
```

(a) Compare the number of pages of each paper of each author with the average number of pages of all the papers of that author. Your resulting table should include 5 columns (names of these are in parentheses): ID of the author (id), name of the author (author), name of the paper (title), number of pages of the paper (pages), average number of pages of all papers of that author (avg). Finally, order your result set by author id (ascending order) and limit results to only the first 10 rows.

(b) Find each author's longest paper (paper with the most pages). In case of ties (if an author's longest paper has the same number of pages as other of his or her papers), order papers by alphabetical order. Your resulting table should include 4 columns (names of these are in parentheses): ID of the author (id), name of the author (author), name of the paper (title), and number of pages of the paper (pages). Finally, order your result set by author id (ascending order) and limit results to only the first 10 rows.

(c) First, find the number of papers each author wrote or co-wrote each year. Now, for these results, store them in a temporary table R, and for each year an author has published a paper, find the minimum and maximum number of papers the author has published in a span of 3 years (note that the years are not necessarily contiguous: for each year when an author has published a paper, look at the previous row and also at the next row in table R).

Your resulting table should include 5 columns (names of these are in parentheses): ID of the author (id), name of the author (author), year of publication (year), minimum number of papers published in a span of 3 publications (min), maximum number of papers published in a span of 3 publications (max). Finally, order your result set by author id (ascending order) and year (ascending order) and limit results to only the first 10 rows.

Example:

Let's suppose author X published 4 papers in 2000, 5 papers in 2001, 3 papers in 2003, and 1 paper in 2005. Your first result set (table R) would look something like this (depending on the names you give to your columns, of course).

id	author	year	count
1	X	2000	4
1	X	2001	5
1	X	2003	3
1	X	2005	1

And this would be your final result set:

id	author	year	min	max
1	X	2000	4	5
1	X	2001	3	5
1	X	2003	1	5
1	X	2005	1	3