

**CSCI317 Database Performance Tuning
Singapore 2023-3**

Assignment 1

Published on 14 July 2023

Scope

This assignment includes the tasks related to database administration, implementation of queries in relational algebra, interpretation of query processing plans, discovery of `SELECT` statement from a query processing plan, and optimal distribution of relational table over the persistent storage devices.

This assignment is due by **Saturday, 5 August 2023, 9.00 pm (sharp) Singaporean Time.**

Please read very carefully information listed below.

This assignment contributes to 15% of the total evaluation in the subject.

A submission procedure is explained at the end of specification.

This assignment consists of 5 tasks and specification of each task starts from a new page.

It is recommended to solve the problems before attending a laboratory class in order to efficiently use supervised laboratory time.

A submission marked by Moodle as "late" is treated as a late submission no matter how many seconds it is late.

A policy regarding late submissions is included in the subject outline.

A submission of compressed files (zipped, gzipped, rared, tared, 7-zipped, lhzed, ... etc) is not allowed. The compressed files will not be evaluated.

All files left on Moodle in a state "`Draft (not submitted)`" will not be evaluated.

It is expected that all tasks included within **Assignment 1** will be solved **individually without any cooperation** with the other students. If you have any doubts, questions, etc. please consult your lecturer or tutor during lab classes or office hours. Plagiarism will result in a **FAIL** grade being recorded for the assessment task.

Please read very carefully information included in Prologue section below about software environment to be used in the subject.

Prologue

In this subject we shall use Oracle 19c database server running under Oracle Linux 7.4 operating system on a virtual machine hosted by VirtualBox. To start Oracle database server you have to start VirtualBox first. If you have not installed VirtualBox on your system yet then it is explained in Cookbook for CSIT115 Recipe 1.1, Step 1 "How to use VirtualBox ?" (<https://www.uow.edu.au/~jrg/115/cookbook/e1-1-frame.html>) how to install and how to start VirtualBox.

When VirtualBox is started, import an appliance included in a file `OracleLinux7.4-64bits-Oracle19c-22-JAN-2020.ova`. You can download ova image of the appliance using one of the links published on Moodle in a section OTHER RESOURCES.

When ready, power on a virtual machine `OracleLinux7.4-64bits-Oracle19c-22-JAN-2020`.

A password to a Linux user `ORACLE` is `oracle` and a password to Oracle users `SYSTEM` and `SYS` (database administrators) is also `oracle`. Generally, whenever you are asked about a password then it is always `oracle`, unless you change it.

When logged as a Linux user, you can access Oracle database server either through a command line interface (CLI) `SQLcl` or through Graphical User Interface (GUI) `SQL Developer`.

You can find in Cookbook for CSCI317, Recipe 1, information on How to access Oracle 19c database server, how to use SQL Developer, how to use basic SQL and `SQLcl`, and how to create a sample database ? (<https://documents.uow.edu.au/~jrg/317sim/cookbook/e1-2-frame.html>) more information on how to use `SQLcl` and `SQL Developer`.

Tasks

Task 1 (3 marks)

The objectives of this task are to learn how to create a database user, how to create a tablespace, how to create relational table in a given tablespace, how to use loader to load data into a relational table and how to find information persistent storage structures of a relational database. An additional hidden objective is to refresh your SQL skills and the ways how to use software installed on a virtual machine.

To install a sample database, perform 14 steps listed below. Note, that no report is expected from the installation of a sample database. However, please note that you will not be able to implement the assignments without a sample database.

- (1) After starting a virtual machine and logging as Linux ORACLE user start Terminal command line shell program.
- (2) Use `cd` command to move to a folder TPCHR. If you do not use a provided virtual machine with Oracle 19c and instead you use your own installation of Oracle DBMS then you can download a zipped file with TPCHR folder from Moodle. A zipped file is available in a folder OTHER RESOURCES.
- (3) A script `tbscreate.sql` can be used to create a tablespace for the relational tables included in TPC-HR benchmark database. A conceptual schema of the database is available in a file `tpchr.pdf`.

You can use `gedit` editor (you can also use `emacs` editor) to view the contents of SQL script `tbscreate.sql`. If you like, you can also change a name of the tablespace to any valid name (for example `csci317`) and location of a file that implements the tablespace, however, it is not necessary.

- (4) To connect SQLcl client as a user SYSTEM process the following command at a shell prompt.

```
sql system
```

You can also use SQL Developer client.

- (5) Next, at SQL> prompt process a command

```
cd ~/TPCHR
```

to move a client to TPCHR folder.

- (6) Next, process a script `tbscreate.sql` to create a new tablespace.
- (7) While connected as SYSTEM user process the scripts `sfs.sql` and `sf.sql` to list free space per tablespace and to list the names of files implementing the tablespaces.

- (8) A script `usrcreate.sql` can be used to create a new database user `tpchr` and to grant appropriate privileges and resources to the new user. Use `gedit` editor (you can also use `emacs` editor) to view and, if you like, to update the contents of SQL script `usrcreate.sql`. To do so you have to use a command

```
!gedit usrcreate.sql
```

in front of `SQL>` prompt.

- (9) Return to `SQLcl` client and process a script `usrcreate.sql`.
- (10) Connect to Oracle server as a user `SYS` and process the scripts `setprivs.sql` and `setplustrace.sql` to grant appropriate privileges to a user `tpchr`.
- (11) Connect to Oracle server as a new user `tpchr` and process a script `dbcreate.sql`. The script creates the relational tables of TPC-HR benchmark database.
- (12) A shell script `dbload.sh` loads data into a sample database. Use an editor to view the contents of shell script `dbload.sh`. If you changed Oracle user name and/or password in step (6) then you have to update the script.
- (13) Disconnect from `SQLcl` and process a shell script `dbload.sh` at a command line prompt in the following way.

```
./dbload.sh
```

and be patient, ... it will take some time to load data ...

- (14) Connect to Oracle server as a user `tpchr` and process SQL script `dbcount.sql` to find the total number of rows in each table.

There is no need to create any reports from the steps listed above.

Implement SQL script `solution1.sql` that performs the following actions.

- (1) First the script connects to a database server as a user `system` and lists the following columns from a dynamic performance view `V$INSTANCE`:

```
INSTANCE_NAME,  
HOST_NAME,  
STARTUP_TIME,  
DATABASE_STATUS,  
DATABASE_TYPE,  
LOGINS and  
INSTANCE_MODE.
```

To connect to a database server within SQL script as a user `system` with a password `oracle` insert the following line into the script.

```
connect system/oracle
```

- (2) Next, the script connects as a user `tpchr` and processes `ANALYZE TABLE` statement to load into a data dictionary statistical information related to the relational tables and indexes implementing a sample database `tpchr` created earlier.

To connect to a database server within SQL script as a user `tpchr` with a password `oracle` insert the following line into the script.

```
connect tpchr/oracle
```

You can find a sample `ANALYZE TABLE` statement in Cookbook, Recipe 5.1, Step 1 How to use `ANALYZE` statement ?

- (3) Next, the script creates a relational table `LINEITEM1994` and loads into the table all rows from a relational table `LINEITEM` and such that the rows contain information about the orders submitted in 1994. A relational table `LINEITEM1994` must be created in the same tablespace as a relational table `LINEITEM`. There is no need to enforce consistency constraints on a relational table `LINEITEM1994`. Do not forget to analyse a relational table `LINEITEM1994`.

- (3) Next, the script connects as a user `system`. To connect to a database server within SQL script as a user `system` with a password `oracle` insert the following line into the script.

```
connect system/oracle
```

- (4) Next, while still connected as a user `system` the script retrieves and lists the following information from a data dictionary.

- (i) The current timestamp obtained from an application of a function `systimestamp`.

- (ii) The total number of rows, total number of data blocks, total number of extents and the total number of bytes occupied by a relational table `LINEITEM` and the total number of rows, total number of data blocks, total number of extents and the total number of bytes occupied by a relational table `LINEITEM1994`. The storage parameters must be listed in the following format.

```
table-name total-rows total-blocks total-extents total-bytes
```

You can find more information about the segment related contents of data dictionary in Cookbook, Recipe 10.2, Step 1 How to find the size of a segment ?

When ready start SQLcl client, connect to Oracle database server, and process SQL script solution1.sql. Save a report from processing of the script in a file solution1.lst. It is explained in Cookbook, Recipe 1.5, Step 9, "How to create and to save a report" how to save a report from processing of SQL script in a text file.

The script must be processed with SQLcl options ECHO and FEEDBACK set to ON such that all SQL statements processed are included in the report !

A good habit is to put SQLcl statements

```
SPOOL solution1
SET ECHO ON
SET FEEDBACK ON
SET LINESIZE 300
SET PAGESIZE 300
```

at the beginning of each SQL script implemented and the following statement at the end of the script

```
SPOOL OFF
```

A report from processing of the script must have NO errors ! A report with the errors scores no marks !

Deliverables

A file solution1.lst that contains a report from the processing of a script solution1.sql.

Task 2 (3 marks)

An objective of this task is to express the queries in a notation of the relational algebra expressions

Consider the following queries related to the relational tables included in TPC-HR sample database.

- (1) Find the names of parts (`P_NAME`) ordered by the customers in 1994 (`year(O_ORDERDATE)`).
- (2) Find the names of customers (`C_NAME`) living in Singapore (`city(C_ADDRESS)`) and the names of suppliers (`S_NAME`) living in Sydney (`city(S_ADDRESS)`) ordered in the ascending order of names. Use an operation `sort` to order the names.
- (3) Find the keys and the names of customers (`C_CUSTKEY`, `C_NAME`) who submitted no orders in 1994.

Write the implementations of the queries listed above as expressions of the relational algebra.

Please note that you MUST use the relational algebra operations explained to you during the lecture classes in the subject AND NOT the operations used in the query processing plans created by the Oracle query optimizer.

Save the relational algebra expressions implementing the queries listed above in a file `solution2.pdf`.

Deliverables

A file `solution2.pdf` that contains implementation of the queries listed above as the expressions of the relational algebra. The handwritten and scanned/photographed implementations of the queries are acceptable. Please note that you MUST use the relational algebra operations explained to you during the lecture classes in the subject AND NOT the operations used in the query processing plans created by the Oracle query optimizer.

Task 3 (3 marks)

An objective of this task is to interpret a query processing plan created by a query optimizer and to draw a syntax tree of a query processing plan

Consider the following fragment of query processing plan.

Id	Operation	Name	Rows	Bytes	TempSpc	Cost (%CPU)	Time
0	SELECT STATEMENT		134K	63M		33989 (1)	00:00:02
* 1	HASH JOIN		134K	63M	51M	33989 (1)	00:00:02
* 2	HASH JOIN		135K	50M	1232K	26145 (1)	00:00:02
* 3	HASH JOIN		4565	1172K		2258 (1)	00:00:01
* 4	TABLE ACCESS FULL	PARTSUPP	4565	633K		1857 (1)	00:00:01
* 5	TABLE ACCESS FULL	PART	60000	7089K		401 (1)	00:00:01
* 6	TABLE ACCESS FULL	LINEITEM	1800K	214M		12160 (1)	00:00:01
7	TABLE ACCESS FULL	ORDERS	450K	46M		2697 (1)	00:00:01

Predicate Information (identified by operation id):

- 1 - access("L_ORDERKEY"="O_ORDERKEY")
- 2 - access("LINEITEM"."L_PARTKEY"="P_PARTKEY")
- 3 - access("PART"."P_PARTKEY"="PS_PARTKEY")
- 4 - filter("PARTSUPP"."PS_SUPPLYCOST"<20)
- 5 - filter("PART"."P_RETAILPRICE">100)
- 6 - filter("LINEITEM"."L_PARTKEY">=0)

Find and draw a syntax tree of the query processing plan listed above. To draw a syntax tree, use the relational algebra operations explained during the lecture classes. Assume that the operations HASH JOIN used in a query processing plan is the same as the operations of join in the relational algebra. Please remember, that you must create a syntax tree with the relational algebra operations explained to you during the lecture classes and NOT with the implementations of such operations by Oracle database system. Save a drawing of a syntax tree in a file `solution3.pdf`.

Deliverables

A file `solution3.pdf` with a drawing of syntax tree of the given query processing plan. A syntax tree must use the relational algebra operations explained to you during the lecture classes. You are allowed to use any line drawing tool to draw a syntax tree. A scanned/photographed copy of a neat hand drawing is also acceptable.

Task 4 (3 marks)

An objective of this task is to discover **SELECT** statement from a given query processing plan.

Consider the following fragment of query processing plan.

Id	Operation	Name	Rows	Bytes	Cost (%CPU)	Time
0	SELECT STATEMENT		1666K	219M	40597	(1) 00:00:02
* 1	HASH JOIN RIGHT ANTI		1666K	219M	40597	(1) 00:00:02
2	VIEW	VW_NSO_1	33743	428K	8217	(1) 00:00:01
* 3	HASH JOIN		33743	1186K	8217	(1) 00:00:01
* 4	TABLE ACCESS FULL	CUSTOMER	2250	49500	1035	(1) 00:00:01
* 5	TABLE ACCESS FULL	ORDERS	449K	6150K	7181	(1) 00:00:01
* 6	TABLE ACCESS FULL	LINEITEM	1800K	214M	32375	(1) 00:00:02

Predicate Information (identified by operation id):

```
-----  
1 - access("L_ORDERKEY"="O_ORDERKEY")  
3 - access("O_CUSTKEY"="C_CUSTKEY")  
4 - filter("CUSTOMER"."C_NAME" LIKE '%James%')  
5 - filter("ORDERS"."O_TOTALPRICE">1000 AND "O_CUSTKEY">=0)  
6 -  
filter("L_TAX"<1000)
```

Discover **SELECT** statement that may have a query processing plan listed above. If you are not be able to get exactly the same query processing plan then, try to find **SELECT** statement with a query processing plan the most similar to the given one.

Use **EXPLAIN PLAN** statement of SQL and SQL script `showplan.sql` to find a query processing plan of a given **SELECT** statement. Save **EXPLAIN PLAN** statement and invocation of a script `showplan.sql` in SQL script `solution4.sql`.

When ready start `SQLcl` client, connect to Oracle database server, and process SQL script `solution4.sql`. Save a report from processing of the script in a file `solution4.lst`. It is explained in Cookbook, Recipe 1.5, Step 9, "How to create and to save a report" how to save a report from processing of SQL script in a text file.

The script must be processed with `SQLcl` options `ECHO` and `FEEDBACK` set to `ON` such that all SQL statements processed are included in the report !

A good habit is to put `SQLcl` statements

```
SPOOL solution4  
SET ECHO ON  
SET FEEDBACK ON  
SET LINESIZE 300  
SET PAGESIZE 300
```

at the beginning of each SQL script implemented and the following statement at the end of the script

```
SPOOL OFF
```

A report from processing of the script must have NO syntax errors !

Deliverables

A file `solution4.lst` with the outcomes from processing of `EXPLAIN PLAN` statement and a listing of query processing plan.

Task 5 (3 marks)

An objective of this task is to find the best distribution of relational tables over the persistent storage devices.

Assume, that to avoid the conflicts when accessing the relational tables of TPC-HR sample database we would like to distribute the relational tables over three different persistent storage devices. Then, the relational tables that are joined together can be simultaneously read from two or more persistent storage devices. Do not worry if your system does not have the persistent storage devices. We shall simulate the persistent storage devices as three different tablespaces `DRIVE_C`, `DRIVE_D` and `DRIVE_E`. You do not have to create the tablespaces. The tablespaces are already created such that each one is located on a different persistent storage device.

Consider the following queries.

- (i) Find the names and addresses (`C_NAME`, `C_ADDRESS`) of customers located in a region with a given name (`R_NAME`).
- (ii) Find the names of parts (`P_NAME`) included in the orders that have a given shipment date (attribute `L_SHIPDATE`) and a given supply cost (`PS_SUPPLYCOST`).
- (iii) Find the region names (`R_NAME`) of suppliers that have an account balance (`S_ACCTBAL`) lower than a given value.
- (iv) Find the names of customers (`C_NAME`) who ordered at least one part shipped in a given year (`L_SHIPDATE`).
- (v) Find the total number of customers per each nation (`N_NAME`).

Note, that the prefixes of the column names indicate the relational tables the columns are located at. For example, `R_NAME` denotes a column in a relational table `REGION`.

An objective of this task is to analyze the queries listed above to find the relational tables used by each query. Then, to distribute the relational tables over the persistent storage devices simulated by the tablespaces `DRIVE_C`, `DRIVE_D` and `DRIVE_E` such, that each relational table used by the same query is located on a different persistent storage device.

Such approach reduces the total number of conflicts when accessing the persistent storage devices and it speeds up the query processing. If it is impossible to distribute the relational tables used by the same application on the different persistent storage devices then try to minimize the total number of conflicts. You do not need to worry about the distribution of indexes used for processing of the queries.

Create a document `solution5.pdf` that contains the following information.

- (1) For each one of the queries listed above find what relational tables are used by a query and draw an undirected hypergraph such that each one of its hyperedges contains the names of tables used by one query. The names of tables are the nodes of the hypergraph.

- (2) Use the hypergraph created in the previous step to find distribution of the relational tables over the persistent storage devices `DRIVE_C`, `DRIVE_D` and `DRIVE_E` such, that each relational table used by the same query is located on the different persistent storage device. If it is impossible to do it, locate smaller relational tables on the same device and larger relational tables on the different devices.

Include a drawing of a hypergraph obtained from a step (1) and information which relational table is assigned to which device obtained from a step (2) in a document `solution5.pdf`.

Hint

You can find a definition and visualization of an undirected hypergraph at:
<https://en.wikipedia.org/wiki/Hypergraph>

Deliverables

A file `solution5.pdf` that contains a drawing of a hypergraph obtained from a step (1) and information which relational table is assigned to which device obtained from a step (2). You are allowed to use any line drawing tool to draw a hypergraph. A scanned/photographed copy of a neat hand drawing is also acceptable.

Submission

Note, that you have only one submission. So, make it absolutely sure that you submit the correct files with the correct contents. No other submission is possible!

Submit the files **solution1.lst**, **solution2.pdf**, **solution3.pdf**, **solution4.lst**, and **solution5.pdf** through Moodle in the following way:

- (1) Access Moodle at **http://moodle.uowplatform.edu.au/**
- (2) To login use a **Login** link located in the right upper corner the Web page or in the middle of the bottom of the Web page
- (3) When logged select a site **CSCI317 (SP323) Database Performance Tuning**
- (4) Scroll down to a section **Submissions**
- (5) Click at a link **In this place you can submit the outcomes of Assignment 1**
- (6) Click at a button **Add Submission**
- (7) Move a file **solution1.lst** into an area **You can drag and drop files here to add them**. You can also use a link **Add...**
- (8) Repeat step (7) for the files **solution2.pdf**, **solution3.pdf**, **solution4.lst**, and **solution5.pdf**.
- (9) Click at a button **Submit assignment** for the bottom of the current web page.
- (10) Click at the checkbox with a text attached: **By checking this box, I confirm that this submission is my own work, ...** in order to confirm the authorship of your submission.
- (11) Click at a button **Continue**
- (12) Check if **Submission status** is **Submitted for grading**.

End of specification