



University of Trieste
Data Management for Big Data Course
Academic Year 2022–2023

Data Warehouse case study

Davide Capone

Sandro Junior Della Rovere

Enrico Stefanel

Contents

1	Introduction	2
1.1	TPC-H benchmark database	2
1.1.1	Database statistics	2
1.1.2	Database SQL definition	3
1.2	Organization of work	3
2	Set of queries	3
2.1	Export/import revenue value	3
2.2	Late delivery	5
2.3	Returned item loss	7
2.4	Execution times	8
3	Materialization	8
3.1	Lineitem - Orders View	8
3.2	Customer - Location View	9
3.3	Supplier - Location View	9
3.4	Queries with materialized views	9
3.5	Execution times	12
4	Indexes design	12
4.1	Indexes on relations	12
4.1.1	Execution times	14
4.2	Indexes on Materialized Views	14
4.2.1	Execution times	16
5	Fragmentation	16
5.0.1	Execution times	22
5.1	Indexes on fragmented tables	22
5.1.1	Execution times	22
6	Conclusions	22

1 Introduction

The aim of this project is to study an efficient implementation of a suite of business oriented ad-hoc queries over the public TPC-H benchmark, which can be considered as a Big Data datawarehouse. The chosen DMBS for the implementation is PostgreSQL.

The way in which this efficient implementation is going to be realized is by exploiting *materialized views*, *indexes* and *(vertical) fragmentation*. The total weight of the database after the implementation of these solutions is requested to stay below 1.5 times the original database weight.

1.1 TPC-H benchmark database

The TPC-H benchmark is a decision support benchmark that can be downloaded from the TPC official website. The data generator lets the user specify a *scale factor* in order to control the size of the resulting database. Our choices was to use a *scale factor* of 10, meaning that the overall database size is approximately 14 GB. The sizes of the individual tables that compose the database are also shown below.

1.1.1 Database statistics

The benchmark is composed by eight tables:

- CUSTOMER, with 16 columns and 1 500 000 tuples (312 MB);
- LINEITEM, with 32 columns and 59 986 052 tuples (11 GB); the main attributes that are going to be used are:
 - l_extendedprice (1 351 462 distinct values, i.e. there is an average of 44 tuples with the same value, that range from 900.91 to 104 949.50),
 - l_discount (11 distinct values, i.e. there is an average of 5 453 277 tuples with the same value, that range from 0.00 to 0.10),
 - l_returnflag (which can assume values A→accepted, R→returned, N→not yet delivered; the percentage of tuples for A and R are almost 25 %, while the percentage of tuples where l_returnflag is N is about 50 %),
 - l_commitdate (2466 distinct values, i.e. there is an average of 24 325 tuples with the same value, that range from 1992-01-31 to 1998-10-31),
 - l_receiptdate (2555 distinct values, i.e. there is an average of 23 478 tuples with the same value, that range from 1992-01-03 to 1998-12-31);
- NATION, with 8 columns and 25 tuples (24 kB);
- ORDERS, with 18 columns and 1 500 000 tuples (2481 kB); the main attribute that is going to be used is:
 - o_orderdate (2406 distinct values, i.e. there is an average of 6234 tuples with the same value, that range from 1992-01-01 to 1998-08-02);

- PART, with 18 columns and 2 000 000 tuples (363 MB); the main attribute that is going to be used is:
 - `p_type` (150 distinct values, i.e. there is an average of 13 333 tuples with the same value);
- PARTSUPP, 10 columns and with 8 000 000 tuples (1535 MB);
- REGION, 6 columns and with 5 tuples (24 kB);
- SUPPLIER, 14 columns and with 100 000 tuples (20 MB).

Other attributes have been used, but statistics about them have been omitted for lack of usefulness (e.g., keys of the tables, for which the cardinality is exactly the cardinality of the corresponding table).

1.1.2 Database SQL definition

The SQL definition of the tables can be found on the official benchmark documentation[1].

1.2 Organization of work

Each component of the group implemented a query (specifically, Della Rovere designed the first query, Capone the second one, and Stefanel the last one). After regular meetings between the group members comparing the solo works, common substructures between the queries have been detected, and the following work proceeded in a coral way.

All the execution times reported in the present report are the results of independent runs on the same machine, with roughly the same external factors. In particular, a computer with an Apple M1 processor, 8 GB of RAM, and macOS operating system has been used. Furthermore, no consecutive runs have been performed on the same query, in order to reduce external biases (following a Round-Robin schema).

2 Set of queries

2.1 Export/import revenue value

It is asked to return the *export/import revenue* between two different nations (E, I) where E is the nation of the lineitem supplier and I the nation of the lineitem customer, and where the *revenue* is defined as

$$\text{SUM}(l_extendedprice * (1 - l_discount))$$

The aggregation should be performed with the Month \rightarrow Quarter \rightarrow Year, (Part) Type and Nation \rightarrow Region roll-ups.

The query is implemented in such a way that it allows the slicing over (Part) Type and Exporting Nation.

```

1  WITH lineitem_orders AS (
2      SELECT
3          l_partkey,
4          l_suppkey,
5          o_orderdate,
6          o_custkey,
7          l_extendedprice,
8          l_discount
9      FROM lineitem JOIN orders ON (l_orderkey = o_orderkey)
10 ), customer_location AS (
11     SELECT
12         c_custkey,
13         c_name,
14         n_nationkey AS c_nationkey,
15         n_name AS c_nationname,
16         r_regionkey AS c_regionkey,
17         r_name AS c_regionname
18     FROM customer
19         JOIN nation ON (c_nationkey = n_nationkey)
20         JOIN region ON (n_regionkey = r_regionkey)
21 ), supplier_location AS (
22     SELECT
23         s_suppkey,
24         s_name,
25         n_nationkey AS s_nationkey,
26         n_name AS s_nationname,
27         r_regionkey AS s_regionkey,
28         r_name AS s_regionname
29     FROM supplier
30         JOIN nation ON (s_nationkey = n_nationkey)
31         JOIN region ON (n_regionkey = r_regionkey)
32 ), query1 AS (
33     SELECT
34         EXTRACT (YEAR FROM o_orderdate) AS _year,
35         EXTRACT (QUARTER FROM o_orderdate) AS _quarter,
36         EXTRACT (MONTH FROM o_orderdate) AS _month,
37         c_regionname,
38         c_nationname,
39         c_name,
40         s_regionname,
41         s_nationname,
42         s_name,

```

```

43     p_type,
44     SUM(l_extendedprice * (1 - l_discount)) AS revenue
45 FROM lineitem_orders
46     JOIN part ON l_partkey = p_partkey
47     JOIN supplier_location ON (s_suppkey = l_suppkey)
48     JOIN customer_location ON (c_custkey = o_custkey)
49 WHERE
50     s_nationkey <> c_nationkey
51     AND p_type = 'PROMO BURNISHED COPPER'
52     AND s_nationname = 'UNITED STATES'
53 GROUP BY
54     _year,
55     _quarter,
56     _month,
57     c_regionkey,
58     c_regionname,
59     c_nationkey,
60     c_nationname,
61     c_custkey,
62     c_name,
63     s_regionkey,
64     s_regionname,
65     s_nationkey,
66     s_nationname,
67     s_suppkey,
68     s_name,
69     p_type
70 )
71 SELECT * FROM query1;

```

2.2 Late delivery

It is asked to retrieve the number of orders where at least one “lineitem” has been received later than the committed date.

The aggregation should be performed with the Month \rightarrow Year roll-up, and the (Customer’s) Nation \rightarrow Region roll-up.

The query has been implemented in such a way that it allows the slicing over a specific Month, and a specific (Part) Type.

```

1 WITH lineitem_orders AS (
2     SELECT
3         o_orderkey,
4         l_partkey,

```

```

5      l_suppkey,
6      o_orderdate,
7      o_custkey,
8      l_commitdate,
9      l_receiptdate
10     FROM lineitem JOIN orders ON (l_orderkey = o_orderkey)
11 ), customer_location AS (
12     SELECT
13         c_custkey,
14         n_nationkey AS c_nationkey,
15         n_name AS c_nationname,
16         r_regionkey AS c_regionkey,
17         r_name AS c_regionname
18     FROM customer
19         JOIN nation ON (c_nationkey = n_nationkey)
20         JOIN region ON (n_regionkey = r_regionkey)
21 ), query2 AS (
22     SELECT
23         EXTRACT(YEAR FROM o_orderdate) AS _year,
24         EXTRACT(MONTH FROM o_orderdate) AS _month,
25         c_regionname,
26         c_nationname,
27         COUNT(DISTINCT(o_orderkey)) AS orders_no
28     FROM lineitem_orders
29         JOIN part ON l_partkey = p_partkey
30         JOIN customer_location ON (c_custkey = o_custkey)
31     WHERE
32         l_receiptdate > l_commitdate
33         AND _month = 1
34         AND p_type = 'PROMO BURNISHED COPPER'
35     GROUP BY
36         _year,
37         _month,
38         c_regionkey,
39         c_regionname,
40         c_nationkey,
41         c_nationname
42 )
43 SELECT * FROM query2;

```

2.3 Returned item loss

It is asked to retrieve the *revenue loss* for customers who might be having problems with the parts that are shipped to them, where a *revenue loss* is defined as

$$\text{SUM}(l_extendedprice * (1 - l_discount))$$

for all qualifying *lineitems*.

The aggregations should be performed with the Month → Quarter → Year and Customer roll-ups.

The query has been implemented in such a way that it allows the slicing over the Name of a Customer combined with a specific Quarter.

```
1  WITH lineitem_orders AS (  
2      SELECT  
3          o_orderkey,  
4          o_orderdate,  
5          o_custkey,  
6          l_extendedprice,  
7          l_discount,  
8          l_returnflag  
9      FROM lineitem JOIN orders ON (l_orderkey=o_orderkey)  
10 ),  
11 query3 AS (  
12     SELECT  
13         EXTRACT(YEAR FROM o_orderdate) AS _year,  
14         EXTRACT(QUARTER FROM o_orderdate) AS _quarter,  
15         EXTRACT(MONTH FROM o_orderdate) AS _month,  
16         c_name,  
17         SUM(l_extendedprice*(1-l_discount)) AS returnloss  
18     FROM  
19         lineitem_orders  
20         JOIN customer ON (o_custkey=c_custkey)  
21     WHERE  
22         l_returnflag='R'  
23         AND c_name='Customer#000129976'  
24         AND EXTRACT(QUARTER FROM o_orderdate) = 1  
25     GROUP BY  
26         _year,  
27         _quarter,  
28         _month,  
29         c_custkey,  
30         c_name  
31 )  
32 SELECT * FROM query3;
```

2.4 Execution times

The query timings have been measured as previously discussed in subsection 1.2.

Query	Run 1	Run 2	Run 3	Run 4	Run 5	μ	σ
1	40 944	39 423	41 053	39 928	38 800	40 029	971
2	45 150	51 300	46 270	46 677	50 235	47 626	2679
3	8245	6886	8604	8790	6943	7893	915

Table 1: Naïve query timings, in milliseconds.

3 Materialization

After a study on the given set of queries, some common intermediate results have been detected between the three, and they have been chosen as candidate views. In order to try lowering the average execution cost, materialized views have been defined starting from the said results.

Furthermore, since the problem considers dealing with a data warehouse (OLAP), there would not be frequent updates, so materialized views are a smart choice.

3.1 Lineitem - Orders View

Since all the queries perform a join between relations *lineitem* and *orders*, an idea is to pre-process this join by creating a primary view.

```
1 CREATE MATERIALIZED VIEW lineitem_orders_mv AS
2   SELECT
3     o_orderkey,
4     l_partkey,
5     l_suppkey,
6     o_orderdate,
7     o_custkey,
8     l_extendedprice,
9     l_discount,
10    l_returnflag,
11    l_commitdate,
12    l_receiptdate
13  FROM lineitem JOIN orders ON (l_orderkey = o_orderkey);
```

This materialized view is composed by 59 986 052 rows, and weighs 4378 MB.

3.2 Customer - Location View

Queries 1 and 2 execute a join operation between *customer*, *nation* and *region*, also this step has been pre-processed by creating a primary view.

```
1 CREATE MATERIALIZED VIEW customer_location_mv AS
2     SELECT
3         c_custkey,
4         c_name,
5         n_nationkey AS c_nationkey,
6         n_name AS c_nationname,
7         r_regionkey AS c_regionkey,
8         r_name AS c_regionname
9     FROM customer
10    JOIN nation ON (c_nationkey = n_nationkey)
11    JOIN region ON (n_regionkey = r_regionkey);
```

This materialized view is composed by 1 500 000 rows, and weighs 167 MB.

3.3 Supplier - Location View

Finally, queries 1 and 2 execute a join operation between *supplier*, *nation* and *region*, also this step has been pre-processed by creating another primary view.

```
1 CREATE MATERIALIZED VIEW supplier_location_mv AS
2     SELECT
3         s_suppkey,
4         s_name,
5         n_nationkey AS s_nationkey,
6         n_name AS s_nationname,
7         r_regionkey AS s_regionkey,
8         r_name AS s_regionname
9     FROM supplier
10    JOIN nation ON (s_nationkey = n_nationkey)
11    JOIN region ON (n_regionkey = r_regionkey);
```

This materialized view is composed by 100 000 rows, and weighs 12 MB.

3.4 Queries with materialized views

The original queries have been re-written, in order to use materialized views instead of using JOIN operations on the related tables.

```
1 WITH query1 AS (
2     SELECT
3         EXTRACT(YEAR FROM o_orderdate) AS _year,
```

```

4      EXTRACT(QUARTER FROM o_orderdate) AS _quarter,
5      EXTRACT(MONTH FROM o_orderdate) AS _month,
6      c_regionname,
7      c_nationname,
8      c_name,
9      s_regionname,
10     s_nationname,
11     s_name,
12     p_type,
13     SUM(l_extendedprice * (1 - l_discount)) AS revenue
14 FROM lineitem_orders_mv
15     JOIN part ON l_partkey = p_partkey
16     JOIN supplier_location_mv ON (s_suppkey = l_suppkey)
17     JOIN customer_location_mv ON (c_custkey = o_custkey)
18 WHERE
19     s_nationkey <> c_nationkey
20     AND p_type = 'PROMO BURNISHED COPPER'
21     AND s_nationname = 'UNITED STATES'
22 GROUP BY
23     _year,
24     _quarter,
25     _month,
26     c_regionkey,
27     c_regionname,
28     c_nationkey,
29     c_nationname,
30     c_custkey,
31     c_name,
32     s_regionkey,
33     s_regionname,
34     s_nationkey,
35     s_nationname,
36     s_suppkey,
37     s_name,
38     p_type
39 )
40 SELECT * FROM query1;

```

```

1 WITH query2 AS (
2 SELECT
3     EXTRACT(YEAR FROM o_orderdate) AS _year,
4     EXTRACT(MONTH FROM o_orderdate) AS _month,
5     c_regionname,

```

```

6      c_nationname,
7      COUNT(DISTINCT(o_orderkey)) AS orders_no
8 FROM lineitem_orders_mv
9      JOIN part ON l_partkey = p_partkey
10     JOIN customer_location_mv ON (c_custkey = o_custkey)
11 WHERE
12     l_receiptdate > l_commitdate
13     AND EXTRACT(MONTH FROM o_orderdate) = 1
14     AND p_type = 'PROMO BURNISHED COPPER'
15 GROUP BY
16     _year,
17     _month,
18     c_regionkey,
19     c_regionname,
20     c_nationkey,
21     c_nationname
22 )
23 SELECT * FROM query2;

```

```

1 WITH query3 AS (
2 SELECT
3     EXTRACT(YEAR FROM o_orderdate) AS _year,
4     EXTRACT(QUARTER FROM o_orderdate) AS _quarter,
5     EXTRACT(MONTH FROM o_orderdate) AS _month,
6     c_name,
7     SUM(l_extendedprice * (1 - l_discount)) AS returnloss
8 FROM
9     lineitem_orders_mv
10    JOIN customer ON o_custkey = c_custkey
11 WHERE
12     l_returnflag = 'R'
13     AND c_name = 'Customer#000129976'
14     AND EXTRACT(QUARTER FROM o_orderdate) = 1
15 GROUP BY
16     _year,
17     _quarter,
18     _month,
19     c_custkey,
20     c_name
21 )
22 SELECT * FROM query3;

```

The total weight of the database with the materialized views is 19 GB (the constraint

on the data warehouse size is respected).

3.5 Execution times

As for the queries with materialized views, we collected statistics on execution timings and results can be observed on Table 2. Further considerations are reported in the section 6.

Query	Run 1	Run 2	Run 3	Run 4	Run 5	μ	σ
1	16 207	15 048	15 257	15 421	15 028	15 392	483
2	14 957	17 235	15 519	16 654	16 400	16 153	910
3	13 742	15 604	14 050	14 822	14 822	14 608	732

Table 2: Query timings with materialized views, in milliseconds.

4 Indexes design

Indexes may help to further reduce the queries execution times. In order to design the indexes, the queries have been executed with the `EXPLAIN ANALYSE` tool, which helps to detect the most expensive operations that are involved.

It turns out that the `JOIN` and `GROUP BY` operations are the most costly, so the indexes were built on the dimensions involved in the aforementioned operations and also on the ones involved in `WHERE` clauses. Only attributes that are used by at least two queries have been considered.

4.1 Indexes on relations

```

1 CREATE INDEX IF NOT EXISTS lineitem_l_orderkey_idx
2   ON lineitem USING btree
3   (l_orderkey ASC NULLS LAST)
4   TABLESPACE pg_default;
5
6 CREATE INDEX IF NOT EXISTS lineitem_l_suppkey_idx
7   ON lineitem USING btree
8   (l_suppkey ASC NULLS LAST)
9   TABLESPACE pg_default;
10
11 CREATE INDEX IF NOT EXISTS lineitem_l_partkey_idx
12   ON lineitem USING btree
13   (l_partkey ASC NULLS LAST)
14   TABLESPACE pg_default;
15
16 CREATE INDEX IF NOT EXISTS order_o_orderdate_idx

```

```

17     ON orders USING btree
18     (o_orderdate ASC NULLS LAST)
19     TABLESPACE pg_default;
20
21 CREATE INDEX IF NOT EXISTS order_o_custkey_idx
22     ON orders USING btree
23     (o_custkey ASC NULLS LAST)
24     TABLESPACE pg_default;
25
26 CREATE INDEX IF NOT EXISTS part_p_type_idx
27     ON part USING btree
28     (p_type ASC NULLS LAST)
29     TABLESPACE pg_default;
30
31 CREATE INDEX IF NOT EXISTS nation_n_name_idx
32     ON nation USING btree
33     (n_name ASC NULLS LAST)
34     TABLESPACE pg_default;
35
36 CREATE INDEX IF NOT EXISTS region_r_name_idx
37     ON region USING btree
38     (r_name ASC NULLS LAST)
39     TABLESPACE pg_default;
40
41 CREATE INDEX IF NOT EXISTS supplier_s_nationkey_idx
42     ON supplier USING btree
43     (s_nationkey ASC NULLS LAST)
44     TABLESPACE pg_default;
45
46 CREATE INDEX IF NOT EXISTS customer_c_nationkey_idx
47     ON customer USING btree
48     (c_nationkey ASC NULLS LAST)
49     TABLESPACE pg_default;
50
51 CREATE INDEX IF NOT EXISTS customer_c_name_idx
52     ON customer USING btree
53     (c_name ASC NULLS LAST)
54     TABLESPACE pg_default;

```

The total weight of the database with the indexes (without the materialized views) is 16 GB, in this way the bound given in section 1 is satisfied for the mandatory part of the project. What follows (i.e., indexes on materialized views) concerns optimizations that have been defined as *optional* in the project statement, and so it is not guaranteed

that the limit remains fulfilled.

4.1.1 Execution times

The tests have been performed without the materialized views. In subsection 4.2 also materialized views are going to be considered. Results are reported in Table 3.

Query	Run 1	Run 2	Run 3	Run 4	Run 5	μ	σ
1	35740	34379	32891	32285	32067	33472	1556
2	55374	53992	53613	53575	55976	54506	1100
3	61	101	45	86	46	67	25

Table 3: Query timings with indexes, in milliseconds.

4.2 Indexes on Materialized Views

Trying to additionally cut the query costs, also indexes on the materialized views have been designed, following the same strategy as before.

```

1 CREATE INDEX IF NOT EXISTS lineitem_orders_o_orderkey_idx
2   ON lineitem_orders_mv USING btree
3   (o_orderkey ASC NULLS LAST)
4   TABLESPACE pg_default;
5
6 CREATE INDEX IF NOT EXISTS lineitem_orders_l_suppkey_idx
7   ON lineitem_orders_mv USING btree
8   (l_suppkey ASC NULLS LAST)
9   TABLESPACE pg_default;
10
11 CREATE INDEX IF NOT EXISTS lineitem_orders_l_partkey_idx
12   ON lineitem_orders_mv USING btree
13   (l_partkey ASC NULLS LAST)
14   TABLESPACE pg_default;
15
16 CREATE INDEX IF NOT EXISTS lineitem_orders_o_orderdate_idx
17   ON lineitem_orders_mv USING btree
18   (o_orderdate ASC NULLS LAST)
19   TABLESPACE pg_default;
20
21 CREATE INDEX IF NOT EXISTS lineitem_orders_o_custkey_idx
22   ON lineitem_orders_mv USING btree
23   (o_custkey ASC NULLS LAST)
24   TABLESPACE pg_default;
25

```

```

26 CREATE INDEX IF NOT EXISTS supplier_location_s_nationkey_idx
27     ON supplier_location_mv USING btree
28     (s_nationkey ASC NULLS LAST)
29     TABLESPACE pg_default;
30
31 CREATE INDEX IF NOT EXISTS supplier_location_s_nationname_idx
32     ON supplier_location_mv USING btree
33     (s_nationname ASC NULLS LAST)
34     TABLESPACE pg_default;
35
36 CREATE INDEX IF NOT EXISTS supplier_location_s_regionkey_idx
37     ON supplier_location_mv USING btree
38     (s_regionkey ASC NULLS LAST)
39     TABLESPACE pg_default;
40
41 CREATE INDEX IF NOT EXISTS supplier_location_s_regionname_idx
42     ON supplier_location_mv USING btree
43     (s_regionname ASC NULLS LAST)
44     TABLESPACE pg_default;
45
46 CREATE INDEX IF NOT EXISTS customer_location_c_nationkey_idx
47     ON customer_location_mv USING btree
48     (c_nationkey ASC NULLS LAST)
49     TABLESPACE pg_default;
50
51 CREATE INDEX IF NOT EXISTS customer_location_c_nationname_idx
52     ON customer_location_mv USING btree
53     (c_nationname ASC NULLS LAST)
54     TABLESPACE pg_default;
55
56 CREATE INDEX IF NOT EXISTS customer_location_c_regionkey_idx
57     ON customer_location_mv USING btree
58     (c_regionkey ASC NULLS LAST)
59     TABLESPACE pg_default;
60
61 CREATE INDEX IF NOT EXISTS customer_location_c_regionname_idx
62     ON customer_location_mv USING btree
63     (c_regionname ASC NULLS LAST)
64     TABLESPACE pg_default;
65
66 CREATE INDEX IF NOT EXISTS customer_location_c_name_idx
67     ON customer_location_mv USING btree
68     (c_name ASC NULLS LAST)

```

The total weight of the database with the indexes on materialized views (which have been defined in section 3) is 23 GB (1.6 times the size of the initial database).

4.2.1 Execution times

The results of execution times for five independent runs on each query with all the aforementioned indexes (on tables and on materialized views) can be seen in Table 4.

Query	Run 1	Run 2	Run 3	Run 4	Run 5	μ	σ
1	22314	21998	21013	21356	20749	21486	658
2	24344	25531	22948	23063	24373	24052	1069
3	83	44	66	44	42	56	18

Table 4: Query timings with indexes on tables and materialized views, in milliseconds.

5 Fragmentation

To conclude the testings, it has been decided to try *vertical fragmentation*.

It makes sense to try such an approach, since queries only use a subset of attributes of the tables. *Vertical fragmentation* is usually implemented in distributed systems but it can still be useful in this case to reduce the workload induced by the queries.

```

1  ---- NATION:
2  ----- no fragmentation, the fragment used by the queries will have 3
   columns and the other 1 column only.
3
4  ---- REGION:
5  ----- no fragmentation, the fragment used by the queries will have 2
   columns and the other 1 column only.
6
7  ---- CUSTOMER:
8  CREATE TABLE IF NOT EXISTS customer_frag_1
9  (
10     c_custkey integer NOT NULL,
11     c_name character varying(25) COLLATE pg_catalog."default" NOT
        NULL,
12     c_nationkey integer NOT NULL,
13     CONSTRAINT customer_frag_1_pkey PRIMARY KEY (c_custkey),
14     CONSTRAINT customer_frag_1_fk1 FOREIGN KEY (c_nationkey)
15         REFERENCES nation (n_nationkey) MATCH SIMPLE
16         ON UPDATE NO ACTION

```



```

17         ON DELETE NO ACTION
18     ) AS
19     SELECT
20         c_custkey,
21         c_name,
22         c_nationkey
23     FROM customer;

```

```

1  CREATE TABLE IF NOT EXISTS customer_frag_2
2  (
3      c_custkey integer NOT NULL,
4      c_address character varying(40) COLLATE pg_catalog."default" NOT
        NULL,
5      c_phone character(15) COLLATE pg_catalog."default" NOT NULL,
6      c_acctbal numeric(15,2) NOT NULL,
7      c_mktsegment character(10) COLLATE pg_catalog."default" NOT NULL,
8      c_comment character varying(117) COLLATE pg_catalog."default" NOT
        NULL,
9      CONSTRAINT customer_frag_2_pkey PRIMARY KEY (c_custkey)
10 ) AS
11     SELECT
12         c_custkey,
13         c_address,
14         c_phone,
15         c_acctbal,
16         c_mktsegment,
17         c_comment
18     FROM customer;

```

```

1  ---- SUPPLIER:
2  CREATE TABLE IF NOT EXISTS supplier_frag_1
3  (
4      s_suppkey integer NOT NULL,
5      s_name character(25) COLLATE pg_catalog."default" NOT NULL,
6      s_nationkey integer NOT NULL,
7      CONSTRAINT supplier_frag_1_pkey PRIMARY KEY (s_suppkey),
8      CONSTRAINT supplier_frag_1_fk1 FOREIGN KEY (s_nationkey)
9          REFERENCES nation (n_nationkey) MATCH SIMPLE
10         ON UPDATE NO ACTION
11         ON DELETE NO ACTION
12 ) AS
13     SELECT

```

```

14     s_suppkey,
15     s_name,
16     s_nationkey
17 FROM supplier;

```

```

1 CREATE TABLE IF NOT EXISTS supplier_frag_2
2 (
3     s_suppkey integer NOT NULL,
4     s_address character varying(40) COLLATE pg_catalog."default" NOT
      NULL,
5     s_phone character(15) COLLATE pg_catalog."default" NOT NULL,
6     s_acctbal numeric(15,2) NOT NULL,
7     s_comment character varying(101) COLLATE pg_catalog."default" NOT
      NULL,
8     CONSTRAINT supplier_frag_2_pkey PRIMARY KEY (s_suppkey)
9 ) AS
10 SELECT
11     s_suppkey,
12     s_address,
13     s_phone,
14     s_acctbal,
15     s_comment
16 FROM supplier;
17
18 ---- PARTSUPP: no fragmentation (table not used in queries)

```

```

1 ---- PART:
2 CREATE TABLE IF NOT EXISTS part_frag_1
3 (
4     p_partkey integer NOT NULL,
5     p_type character varying(25) COLLATE pg_catalog."default" NOT
      NULL,
6     CONSTRAINT part_frag_1_pkey PRIMARY KEY (p_partkey)
7 ) AS
8 SELECT
9     p_partkey,
10    p_type
11 FROM part;

```

```

1 CREATE TABLE IF NOT EXISTS part_frag_2
2 (
3     p_partkey integer NOT NULL,

```

```

4      p_name character varying(55) COLLATE pg_catalog."default" NOT
      NULL,
5      p_mfgr character(25) COLLATE pg_catalog."default" NOT NULL,
6      p_brand character(10) COLLATE pg_catalog."default" NOT NULL,
7      p_size integer NOT NULL,
8      p_container character(10) COLLATE pg_catalog."default" NOT NULL,
9      p_retailprice numeric(15,2) NOT NULL,
10     p_comment character varying(23) COLLATE pg_catalog."default" NOT
      NULL,
11     CONSTRAINT part_frag_2_pkey PRIMARY KEY (p_partkey)
12 ) AS
13 SELECT
14     p_partkey,
15     p_name,
16     p_mfgr,
17     p_brand,
18     p_size,
19     p_container,
20     p_retailprice,
21     p_comment
22 FROM part;

```

```

1  ---- ORDERS:
2  CREATE TABLE IF NOT EXISTS orders_frag_1
3  (
4      o_orderkey integer NOT NULL,
5      o_custkey integer NOT NULL,
6      o_orderdate date NOT NULL,
7      CONSTRAINT orders_frag_1_pkey PRIMARY KEY (o_orderkey),
8      CONSTRAINT orders_frag_1_fk1 FOREIGN KEY (o_custkey)
9          REFERENCES customer (c_custkey) MATCH SIMPLE
10         ON UPDATE NO ACTION
11         ON DELETE NO ACTION
12 ) AS
13 SELECT
14     o_orderkey,
15     o_custkey,
16     o_orderdate
17 FROM orders;

```

```

1  CREATE TABLE IF NOT EXISTS orders_frag_2
2  (

```

```

3      o_orderkey integer NOT NULL,
4      o_orderstatus character(1) COLLATE pg_catalog."default" NOT NULL,
5      o_totalprice numeric(15,2) NOT NULL,
6      o_orderpriority character(15) COLLATE pg_catalog."default" NOT
      NULL,
7      o_clerk character(15) COLLATE pg_catalog."default" NOT NULL,
8      o_shippriority integer NOT NULL,
9      o_comment character varying(79) COLLATE pg_catalog."default" NOT
      NULL,
10     CONSTRAINT orders_frag_2_pkey PRIMARY KEY (o_orderkey)
11 ) AS
12 SELECT
13     o_orderkey,
14     o_orderstatus,
15     o_totalprice,
16     o_orderpriority,
17     o_clerk,
18     o_shippriority,
19     o_comment
20 FROM orders;

```

```

1  ---- LINEITEM:
2  CREATE TABLE IF NOT EXISTS lineitem_frag_1
3  (
4      l_orderkey integer NOT NULL,
5      l_partkey integer NOT NULL,
6      l_suppkey integer NOT NULL,
7      l_linenummer integer NOT NULL,
8      l_extendedprice numeric(15,2) NOT NULL,
9      l_discount numeric(15,2) NOT NULL,
10     l_returnflag character(1) COLLATE pg_catalog."default" NOT NULL,
11     l_commitdate date NOT NULL,
12     l_receiptdate date NOT NULL,
13     CONSTRAINT lineitem_frag_1_pkey PRIMARY KEY (l_orderkey,
        l_linenummer),
14     CONSTRAINT lineitem_frag_1_fk1 FOREIGN KEY (l_orderkey)
        REFERENCES orders (o_orderkey) MATCH SIMPLE
15         ON UPDATE NO ACTION
16         ON DELETE NO ACTION,
17     CONSTRAINT lineitem_frag_1_fk2 FOREIGN KEY (l_partkey, l_suppkey)
        REFERENCES partsupp (ps_partkey, ps_suppkey) MATCH SIMPLE
18         ON UPDATE NO ACTION
19         ON DELETE NO ACTION
20
21

```

```

22 ) AS
23 SELECT
24     l_orderkey,
25     l_linenumber,
26     l_partkey,
27     l_suppkey,
28     l_extendedprice,
29     l_discount,
30     l_returnflag,
31     l_commitdate,
32     l_receiptdate
33 FROM lineitem;

```

```

1 CREATE TABLE IF NOT EXISTS lineitem_frag_2
2 (
3     l_orderkey integer NOT NULL,
4     l_linenumber integer NOT NULL,
5     l_quantity numeric(15,2) NOT NULL,
6     l_tax numeric(15,2) NOT NULL,
7     l_linestatus character(1) COLLATE pg_catalog."default" NOT NULL,
8     l_shipdate date NOT NULL,
9     l_shipinstruct character(25) COLLATE pg_catalog."default" NOT
    NULL,
10    l_shipmode character(10) COLLATE pg_catalog."default" NOT NULL,
11    l_comment character varying(44) COLLATE pg_catalog."default" NOT
    NULL,
12    CONSTRAINT lineitem_frag_2_pkey PRIMARY KEY (l_orderkey,
        l_linenumber),
13    CONSTRAINT lineitem_frag_2_fk FOREIGN KEY (l_orderkey)
14        REFERENCES orders (o_orderkey) MATCH SIMPLE
15        ON UPDATE NO ACTION
16        ON DELETE NO ACTION
17 ) AS
18 SELECT
19     l_orderkey,
20     l_linenumber,
21     l_quantity,
22     l_tax,
23     l_linestatus,
24     l_shipdate,
25     l_shipinstruct,
26     l_shipmode,
27     l_comment

```

28 `FROM lineitem;`

The weight of the data warehouse with fragmented tables is roughly the same as the original one, since no additional data structures have been defined and the only action performed is a physical *split* of relations.

5.0.1 Execution times

Timings have been calculated using the queries defined in section 2 by only changing tables names.

Query	Run 1	Run 2	Run 3	Run 4	Run 5	μ	σ
1	14977	15803	15851	16517	16047	15839	558
2	20368	20474	19809	20284	20866	20360	380
3	2478	2150	2145	2350	2346	2294	147

Table 5: Query timings using fragmentation, in milliseconds.

5.1 Indexes on fragmented tables

Since the results shown in Table 5 are promising, it has been decided to implement the *indexes* (the ones defined in subsection 4.1) on the corresponding fragments.

The total size of the data warehouse at this point is 15 GB (again, the size constraint defined in section 1 is respected).

5.1.1 Execution times

Query	Run 1	Run 2	Run 3	Run 4	Run 5	μ	σ
1	24459	22156	21858	21829	21471	22355	1201
2	45378	40845	41757	40048	40449	41695	2154
3	143	164	61	67	93	106	46

Table 6: Query timings using fragmentation and indexes, in milliseconds.

6 Conclusions

Considering the overall results shown in Figure 1, the optimal approach for optimising the efficiency of queries 1 (Export/import revenue value) and 2 (Late delivery) might be to use the materialized views proposed in section 3 but this happens to be the worst-case scenario for query 3 (Returned item loss). The opposite situation occurs when using the indexes defined in section 4: the query 3 is optimised, but queries 1 and 2 show roughly the same run times as the naïve solution.

Assuming that all the three queries have the same importance (i.e., none of them is being executed a lot more frequently than the others), a good trade-off may seem to use *materialized views* and *indexes* (subsection 4.2). This solution, by the way, does not meet the size constraint of the project.

The final solution that is being proposed for the given problem concerns using the (*vertical*) *fragmentation* reported in section 5 without any additional index. Query 1 is executed in 15.84s, i.e. 2.5 times faster than the naïve solution; query 2 runs in 20.36s, i.e. it is 2.3 times faster than the related naïve solution; query 3 is executed in 2.29 s, i.e. 3.5 times faster than the corresponding naïve solution.

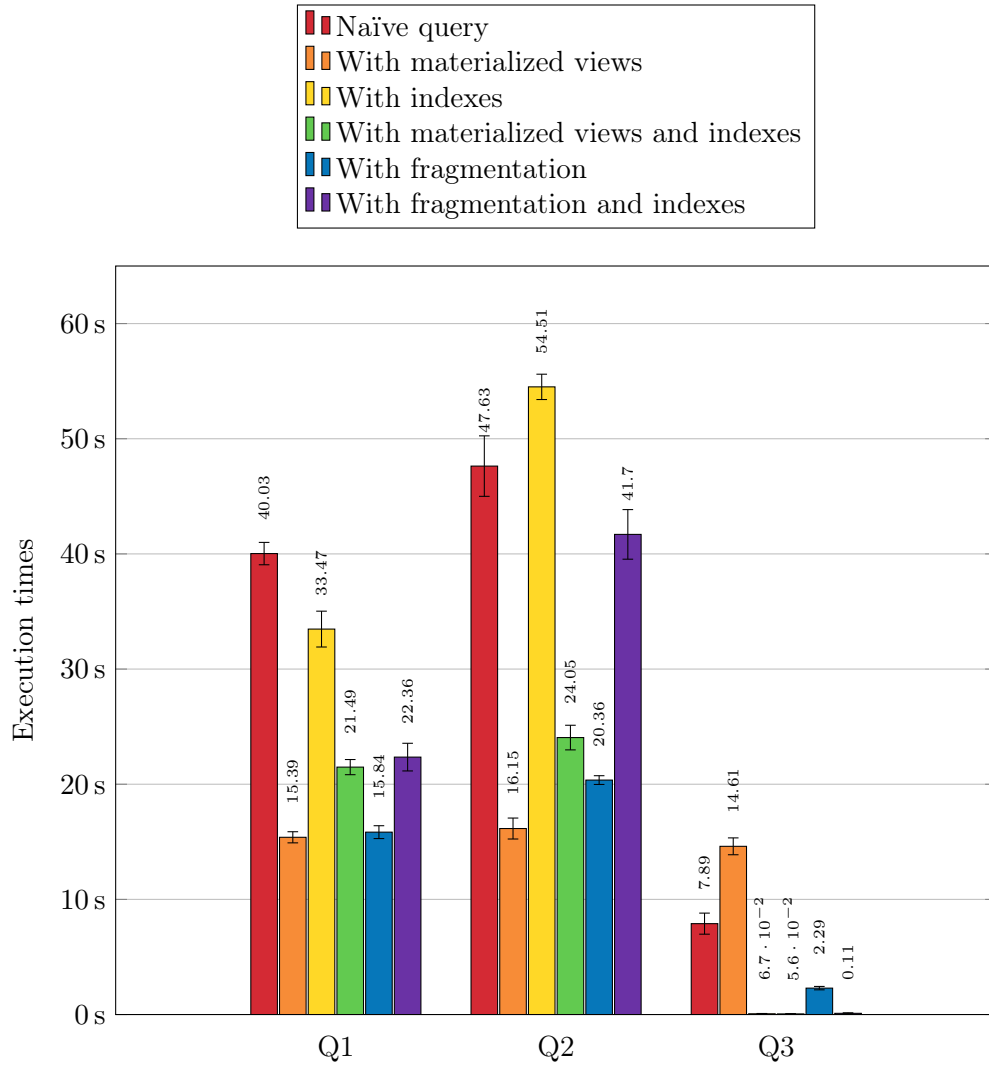


Figure 1: Query timings

References

- [1] Transaction Processing Performance Council (TPC). *TPC BENCHMARK™ H (Decision Support) Standard Specification*. 2022. URL: https://www.tpc.org/TPC_Documents_Current_Versions/pdf/TPC-H_v3.0.1.pdf.