

Data Management - Final Project

Sara Carpenè, Alessio Valentinis, Marco Zampar

July 21, 2024

1 Introduction

This project has the aim of optimize a set of four queries from the TPC benchmark H. The database consists of eight tables: customer, lineitem, nation, orders, part, partsupp, region, supplier. The relations between tables can be seen in the schema in figure 1.

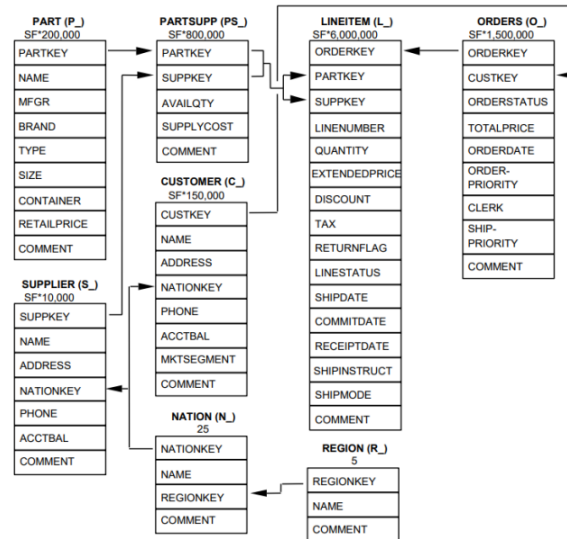


Figure 1: Schema of the relation between tables in the TPC-H benchmark

1.1 Creation and population of the database

The size of the database is scalable, and depends on a scale factor (SF), and we were given data generated from $SF = 10$. Each table, from the documentation, has its own primary key and one or more foreign keys. Complete description of the table can be found here, while the complete SQL tables implementation and creation can be found here.

For the purpose of space economy, we opted for an initial "vanilla" creation of the tables, without any kind of keys: this choice is furthermore supported by the fact that we are dealing with a Data Warehousing context, so we expect the data to be already prepared and cleaned from duplicates.

After having populated the tables, we inserted all the keys reported in the documentation in order to work properly with the relations. However, always for the purpose of gaining some space, we decided not to implement the *Primary key* to the Main table `lineitem`, as (for the same purposes described above,) we don't need to check for uniqueness constraints during the population of the Database, and for the sake of optimization, a simple index should help us to spare some pretty useful space.

In the table 1, we provide some information about the dimension of the various tables, in increasing order.

Table	Number of rows	Dim without keys (MB)	Dim with keys (MB)
region	5	0.01	0.02
nation	25	0.01	0.02
supplier	100000	17.35	19.51
customer	1500000	290.17	322.32
part	2000000	320.14	363.00
partsupp	8000000	1362.80	1535.12
orders	15000000	2038.97	2360.30
lineitem	59986052	8787.95	10073.67

Table 1: General statistics

1.2 Measurement techniques

In order to assess the performance of the various queries, we used the `EXPLAIN ANALYZE` feature of *PostgreSQL* (per gli amici Postgress). It provided a complete summary of the execution plan of the queries, and provides also the execution time in milliseconds. In order to have some significance on the result, we executed every query 5 times, and then we computed the mean and standard deviation of the measurements.

In order to keep track of the dimension of the database, we used the `SELECT pg_database_size();` command, and convert the dimension in MB or GB.

1.3 Hardware specifications

All the test were conducted on a MacBook Air laptop with the following characteristics:

- CPU: Chip Apple M2, 8 core (4 performance cores and 4 efficiency cores);
- RAM: 8GB;
- SSD: SATA 256GB;
- GPU: built-in 10 cores GPU;
- OS: macOS Sonoma 14.2.1.

2 Statistics of the DB

In this section, we will present some useful statistics of the database.

The original database has a total dimension of 14681.63 MB (?? 13.075GB ??), this encompasses both the physical size of the tables and the dimensions of primary and secondary indexes across various attributes.

In section A it can be found the complete statistics of the tables detailing its attributes along with the count of unique values, minimum, and maximum values for each attribute.

Here, to keep the focus on the four queries that we want to optimize, we will limit the presentation to some of that statistics. Specifically we decided to mention only the tables and the attributes that will be involved in at least one of the chosen query.

Attribute	Distinct values	Min value	Max value
c_custkey	1500000	1	1500000
c_name	1500000	'Customer#000000001'	'Customer#001500000'
c_address	1500000	-	-
c_nationkey	25	0	24
c_phone	1499963	-	-
c_acctbal	818834	'-999.99'	'9999.99'
c_comment	1496636	-	-

Table 2: Costumer statistics

Attribute	Distinct values	Min value	Max value
p_partkey	2000000	1	2000000
p_type	150	-	-
p_container	40	-	-

Table 3: Part statistics

Attribute	Distinct values	Min value	Max value
l_orderkey	15000000	1	60000000
l_partkey	2000000	1	2000000
l_quantity	50	1	50
l_extendedprice	1351462	900.91	104949.5
l_discount	11	'0.0'	'0.1'
l_tax	9	'0.0'	'0.08'
l_returnflag	3	-	-
l_linestatus	2	-	-
l_shipdate	2526	'1992-01-02'	'1998-12-01'
l_linenum	7	1	7

Table 4: Lineitem statistics

Attribute	Distinct values	Min value	Max value
o_orderkey	15000000	1	60000000
o_custkey	999982	1	1499999
o_orderdate	2406	'1992-01-01'	'1998-08-02'

Table 5: Orders statistics

Attribute	Distinct values	Min value	Max value
n_nationkey	25	0	24
n_name	25	-	-

Table 6: Nation statistics

3 Query schemas

The assignment consists in using TPC-Benchmark H to test and optimize a set of four queries, using indexes, materialized views, a mixed approach of the two and fragmentation.

The set of queries selected for the assignment are Q1, Q10, Q14, Q17 of the Official Documentation. An overall view on the description and SQL implementation is given below.

3.1 Query 1

Brief description The Pricing Summary Report Query provides a summary pricing report for all lineitems shipped as of a given date. The date is within 60 - 120 days of the greatest ship date contained in the database. The query lists totals for extended price, discounted extended price, discounted extended price plus tax, average quantity, average extended price, and average discount. These aggregates are grouped by RETURNFLAG and LINESTATUS, and listed in ascending order of RETURNFLAG and LINESTATUS. A count of the number of lineitems in each group is included.

Functional definition

Listing:

```
SELECT
  l_returnflag,
  l_linestatus,
  SUM(l_quantity) AS sum_qty,
  SUM(l_extendedprice) AS sum_base_price,
  SUM(l_extendedprice * (1 - l_discount)) AS sum_disc_price,
  SUM(l_extendedprice * (1 - l_discount) * (1 + l_tax)) AS sum_charge,
  AVG(l_quantity) AS avg_qty,
  AVG(l_extendedprice) AS avg_price,
  AVG(l_discount) AS avg_disc,
  COUNT(*) AS count_order
FROM
  lineitem
WHERE
```

```

l_shipdate <= DATE '1998-12-01' - INTERVAL '[DELTA]' DAY
GROUP BY
  l_returnflag,
  l_linestatus
ORDER BY
  l_returnflag,
  l_linestatus;

```

For a matter of simplicity we decided to take as slicing values the ones proposed by the validation paragraph in the official documentation. (So
 '[DELTA]' = '90')

3.2 Query 10

Brief description The Returned Item Reporting Query finds the top 20 customers, in terms of their effect on lost revenue for a given quarter, who have returned parts. The query considers only parts that were ordered in the specified quarter. The query lists the customer's name, address, nation, phone number, account balance, comment information and revenue lost. The customers are listed in descending order of lost revenue. Revenue lost is defined as $\text{sum}(l_extendedprice * (1 - l_discount))$ for all qualifying lineitems.

Functional definition

Listing:

```

SELECT
  c_custkey,
  c_name,
  SUM(l_extendedprice * (1 - l_discount)) AS revenue,
  c_acctbal,
  n_name,
  c_address,
  c_phone,
  c_comment
FROM
  customer,
  orders,
  lineitem,
  nation
WHERE
  c_custkey = o_custkey
  AND l_orderkey = o_orderkey
  AND o_orderdate >= DATE '[DATE]'
  AND o_orderdate < DATE '[DATE]' + INTERVAL '3' MONTH
  AND l_returnflag = 'R'
  AND c_nationkey = n_nationkey
GROUP BY
  c_custkey,
  c_name,
  c_acctbal,
  c_phone,
  n_name,
  c_address,

```

```

c_comment
ORDER BY
revenue DESC;

```

For a matter of simplicity we decided to take as slicing values the ones proposed by the validation paragraph in the official documentation. (So ' [DATE]' = '1993-10-01')

3.3 Query 14

Brief description The Promotion Effect Query determines what percentage of the revenue in a given year and month was derived from promotional parts. The query considers only parts actually shipped in that month and gives the percentage. Revenue is defined as $(l_extendedprice * (1 - l_discount))$.

Functional definition

Listing:

```

SELECT
    100.00 * SUM (CASE WHEN p_type like 'PROMO%'
                      THEN l_extendedprice*(1-l_discount)
                      ELSE 0 END) / SUM(l_extendedprice * (1 - l_discount))
    AS promo_revenue
FROM
    lineitem,
    part
WHERE
    l_partkey = p_partkey
    AND l_shipdate >= date '[DATE]'
    AND l_shipdate < date '[DATE]' + interval '1' month;

```

For a matter of simplicity we decided to take as slicing values the ones proposed by the validation paragraph in the official documentation. (So ' [DATE]' = '1995-09-01')

3.4 Query 17

Brief description The Small-Quantity-Order Revenue Query considers parts of a given brand and with a given container type and determines the average lineitem quantity of such parts ordered for all orders (past and pending) in the 7-year database. What would be the average yearly gross (undiscounted) loss in revenue if orders for these parts with a quantity of less than 20 % of this average were no longer taken?

Functional definition

Listing:

```

SELECT
    SUM(l_extendedprice) / 7.0 AS avg_yearly
FROM
    lineitem,
    part
WHERE
    p_partkey = l_partkey

```

```

AND p_brand = '[BRAND]'
AND p_container = '[CONTAINER]'
AND l_quantity < (
    SELECT
        0.2 * AVG(l_quantity)
    FROM
        lineitem
    WHERE
        l_partkey = p_partkey
);

```

For a matter of simplicity we decided to take as slicing values the ones proposed by the validation paragraph in the official documentation. (So
 '[BRAND]' = 'Brand#23' and '[CONTAINER]' = 'MED BOX')

4 Baseline

We decided to test the execution time of queries without additional indexes or views, other than the keys suggested by the documentation, and to use it as a baseline for further improvement in the management of the queries.

In order to keep the results as similar as possible to the theoretical case, we decided to disable any kind of hash-related operations, that in PostgreSQL are very much used and optimized, and can alter the improvement of the execution time when introducing queries.

The results of the tests are reported in the plot at table 7.

The execution times for Query 17 were not reported in the table because it was not feasible to execute it within reasonable time frames to obtain mean and standard deviation values. Since after two hours of computation it had still not produced a result, we decided to estimate empirically the possible execution times. We ran the query on a limited dataset, attempting to understand the relationship between the number of rows in the table and execution times. Given that executing the query on data related to one month, or 103.000 rows, requires a cost of 144.035 operations and takes 2 seconds, and running the query on data related to three months, or 8.814 rows, requires a cost of 6.091.570 and takes 173 seconds, we reasonably assumed that executing the query on the entire table composed of 60 million rows, with a cost of 117.789.668.871 operations, would take significantly longer than what we had available. Additionally, since the purpose of data warehousing is to support the decision-making, this query without optimization becomes useless and thus reporting its cost would be irrelevant if it exceeds a couple of hours.

Query	Mean [s]	Std [s]
Q1	41.778	1.412
Q10	33.077	1.634
Q14	28.253	1.239
Q17	N/A	N/A

Table 7: Execution times of query

[MAYBE ADD MEDIAN]

5 Indexes

Our first attempt was to add indexes on foreign keys, but in almost all cases, they weren't used, or didn't bring too many advantages, compared with their size.

Our second attempt was to add indexes on the attributes used for slicing, so involved in the **WHERE** condition.

Table	Attribute	Used in Query	Creation time [s]	Index size [MB]
lineitem	l_shipdate	Q1	32.43	397.54
lineitem	l_returnflag	Q1	61.83	396.46
lineitem	l_partkey	Q17	46.99	429.50
orders	o_orderdate	Q10	8.10	100.18

Table 8: Indexes dimensions

With these indexes, which are ensured to be used in the execution of the queries, resulted in a total database size of *20.55GB*.

The execution time of the queries is summarized in the table below.

Query	Mean [s]	Std [s]
Q1	30.009	1.224
Q10	25.677	1.799
Q14	23.607	0.327
Q17	11.232	0.967

Table 9: Execution times of query with indexes

[Comment on time gains in percentage]. See that doing only this we made query 17 feasible

6 Materialized views

In this section we will propose some materialized views that aim to improve the execution time of the chosen queries. For the creation of this views we enabled all the hash related operations, since the main purpose of the section is to evaluate performances related to materialization of some tables, rather than delving into the specifics of this process.

6.1 Lineitem-part

In order to improve performances in executing query 14 we decided to create a materialized view as follow:

Listing:

```
CREATE MATERIALIZED VIEW part_lineitem AS
SELECT
  l_returnflag,
```



```

l_linestatus,
l_quantity,
l_extendedprice,
l_discount,
l_tax,
l_shipdate,
l_partkey,
p_partkey,
p_brand,
p_container,
SUBSTRING(p_type FROM 1 FOR 5) AS p_type_prefix,
0.2 * AVG(l_quantity) OVER (PARTITION BY l_partkey) AS avg_quantity
FROM
lineitem l
JOIN
part p ON l.l_partkey = p.p_partkey;

```

Statistics of this view:

- Required time to create the view 310.952 seconds.
- Size of the view: 6.43 GB.

We rewrote the query 14 and 17 in order to exploit the materialized views.

Materializing the join operation of query 14 we expected to have a lower execution time with respect to the baseline, but this did not happen. By observing the output of the *EXPLAIN ANALYZE* function we understood that the problem with the lack of gain in performance is that the optimizer performs a sequential scan of the lineitem table, since there is no index on shipdate.

The greatest improvement in using this materialization was expected to be in query 17. Indeed the execution time dropped to 15.265 seconds (on a single run), which is a sensible gain in terms of performances.

6.2 Costumer-orders-lineitem-nation

In order to spare the most time of the joins in query 10 we created a bigger materialized view.

Listing:

```

CREATE MATERIALIZED VIEW customer_order_lineitem_nation AS
SELECT
c.c_custkey,
c.c_name,
c.c_acctbal,
n.n_name,
c.c_address,
c.c_phone,
c.c_comment,
l.l_returnflag,
l.l_discount,
l.l_extendedprice,
o.o_orderdate

```

```

FROM
  customer c
JOIN
  orders o ON c.c_custkey = o.o_custkey
JOIN
  lineitem l ON l.l_orderkey = o.o_orderkey
JOIN
  nation n ON c.c_nationkey = n.n_nationkey;

```

Statistics of this view:

- Required time to create the view 437.107 seconds.
- Size of the view: 12.63 GB.

We tested the performances of query 10 rewritten using this materialization, but there was no gain in efficiency, other than the fact that this materialization was made ad-hoc for this single query.

6.3 Lineitem-part-orders

As a last option we opted for a mixed approach creating a view which is neither as big nor as specific as the previous ones.

Listing:

```

CREATE MATERIALIZED VIEW part_lineitem_order AS
SELECT
  l_returnflag,
  l_linestatus,
  l_quantity,
  l_extendedprice,
  l_discount,
  l_tax,
  l_shipdate,
  l_partkey,
  p_partkey,
  p_brand,
  p_container,
  SUBSTRING(p_type FROM 1 FOR 5) AS p_type_prefix,
  0.2 * AVG(l_quantity) OVER (PARTITION BY l_partkey) AS avg_quantity,
  o_orderkey,
  o.o_custkey,
  o.o_orderdate
FROM
  lineitem l
JOIN
  part p ON l.l_partkey = p.p_partkey
JOIN
  orders o ON l.l_orderkey = o.o_orderkey;

```

Statistics of this view:

- Required time to create the view 366.508 seconds.
- Size of the view: 6.93 GB.

As this is the expected best materialization, we report the modified queries that use it.

Query 10

```
SELECT
  c_custkey,
  c_name,
  SUM(l_extendedprice * (1 - l_discount)) AS revenue,
  c_acctbal,
  n_name,
  c_address,
  c_phone,
  c_comment
FROM
  part_lineitem_order
  JOIN customer c ON c.c_custkey = o_custkey
  JOIN nation n ON c.c_nationkey = n.n_nationkey
WHERE
  o_orderdate >= DATE '1993-10-01'
  AND o_orderdate < DATE '1993-10-01' + INTERVAL '3' MONTH
  AND l_returnflag = 'R'
GROUP BY
  c_custkey,
  c_name,
  c_acctbal,
  c_phone,
  n_name,
  c_address,
  c_comment
ORDER BY
  revenue DESC;
```

Query 14

```
SELECT
  100.00 * SUM(CASE
    WHEN p_type_prefix LIKE 'PROMO'
    THEN l_extendedprice * (1 - l_discount)
    ELSE 0
  END) / SUM(l_extendedprice * (1 - l_discount)) AS promo_revenue
FROM
  part_lineitem_order
WHERE
  l_shipdate >= DATE '1995-09-01'
  AND l_shipdate < DATE '1995-09-01' + INTERVAL '1' MONTH;
```

Query 17

```
SELECT
```

```

SUM(l_extendedprice) / 7.0 AS avg_yearly
FROM
part_lineitem_order
WHERE
  p_brand = 'Brand#23'
  AND p_container = 'MED BOX'
  AND l_quantity < avg_quantity;

```

In this case we performed a complete benchmark, and the results are reported in the table below 10.

Query	Mean	Std
Q1	24.770	2.034
Q10	50.135	59.566
Q14	20.677	0.353
Q17	20.465	0.655

Table 10: Execution times of query with materialized views

For the sake of completeness, we will briefly present an idea to optimize with materialization also query one, which is the only query in the chosen set that isn't directly influenced by the already introduced materialized views. We proposed a materialization that returns a "slimmer" version of the lineitem table with less attributes than the original one. Specifically we decided to keep only l_returnflag, l_linestatus, l_extendedprice, l_discount, l_tax, l_quantity.

This approach led to a slight decrease in the execution time of query 1, but this improvement was not so significant to balance the size of the materialized view and its specificity. Indeed this materialization is "tailored" to query 1 and it could hardly be used in the optimization of other queries.

Since, in completing this project, we have always considered our set of four queries as a small subset of the possible query that a company may see as interesting in the decision making process, we didn't want to burden the database with too specific materialized views. For this reason we decided to not consider also this materialization, and to continue the project with only the materialized view that involves lineitem, part and orders tables.

7 Mixed approach

In this section we explored a mixed approach with both materialization and indexing. As we already sapred the join conditions, we will just consider indexes that are useful for slicing. Remember that the chosen materialized view results to be *Lineitem-part-orders*, but for the sake of completeness, we reported also indexing on all the materialization attempts, in order to record possible time improvement.

For every materialization we will report a table with the introduced indexes and their creation time and size.

7.1 Lineitem-part

Attribute	Creation time [s]	Index size [MB]
l_shipdate	30.16	397.55
p_brand	58.95	403.14
p_container	60.34	403.15

Table 11: Indexes dimensions

Using this indexes the total size of tables is: 7.60 GB.

Remember that this materialization is useful mainly for queries 14 and 17, so the indexes were placed on the suitable slicing attributes.

7.2 Costumer-orders-lineitem-nation

Attribute	Creation time [s]	Index size [MB]
o_orderdate	58.30	397.51
l_returnflag	50.62	396.46

Table 12: Indexes dimensions

Using this indexes the total size of tables is: 13.41 GB.

Recalling that this materialization is pretty useful only for one query, we exploited indexing only for the proper attributes, but given the resulting size of the table, we furthermore opted for dropping this big materialization in favor of smaller ones.

7.3 Lineitem-part-orders

Attribute	Creation time [s]	Index size [MB]
l_shipdate	47.12	397.55
o_orderdate	44.72	397.51
l_returnflag	47.30	396.46
p_brand	49.68	403.14
p_container	59.57	403.15

Table 13: Indexes dimensions

Using this indexes the total size of tables is: 8.11 GB.

As this materialization resulted to be useful for three of the four proposed queries, we added indexes for all the useful slicing attributes involved in queries 10, 14, 17.

Recalling the fact that the chosen final materialization is the *lineitem-part-orders*, with all the indexes the final database size is of 32.23 GB.

We reported in table ?? the results of testing the queries with the described structure.

Query	Mean	Std
Q1	23.184	1.360
Q10	34.227	4.548
Q14	24.086	1.504
Q17	1.587	0.103

Table 14: Execution times of query with both materialized views and indexes

8 Fragmentation

We considered to implement the fragmentation only for the tables of `lineitem` and `orders` since they are the most computationally expansive to scan entirely and since they are strictly involved in the set of our chosen queries.

While designing the fragmentation, we have always considered the broader aspect of the database as a decision-making tool, avoiding introducing too specific partitions that could have improved the specific set of chosen queries, but would have unnecessarily burdened the database as they were not generalizable to other queries. Therefore, we decided to consider a temporal fragmentation, as the temporal dimension is often involved in slicing conditions, even outside of the chosen queries.

In particular we fragmented the `orders` table with respect to the `o_orderdate` attribute. Each partitioned table contains a time-span of three months to allow a significant improvement in executing query 10. Furthermore we introduced in the partitioned tables a primary key in `o_orderkey` and a foreign key on `o_custkey` referencing `c_custkey`.

For the `lineitem` table we decided to use a partition on `l_shipdate` and a sub-partition on `l_returnflag`, to allow exploiting this partitioning in query 1, 10, 14.

The results are reported in the table 15

Query	Mean	Std
Q1	60.109	4.050
Q10	19.369	4.091
Q14	3.834	0.160
Q17	311.702	85.586

Table 15: Execution times of query with fragmented db

9 Conclusions

The queries selected for this project were pretty different from each other, and each of them asked for a different optimization technique.

This said, we can have a comprehensive look at what happened to execution times of all the queries given all the techniques analyzed.

If we compare this behavior with the comprehensive size of the dataset, we can have a more critical opinion on what can be the best technique to adopt in our case.

Regarding query one we can see that the most significant improvement is reached with the usage of indexing. We can notice also that the fragmentation of the `lineitem` table worsen

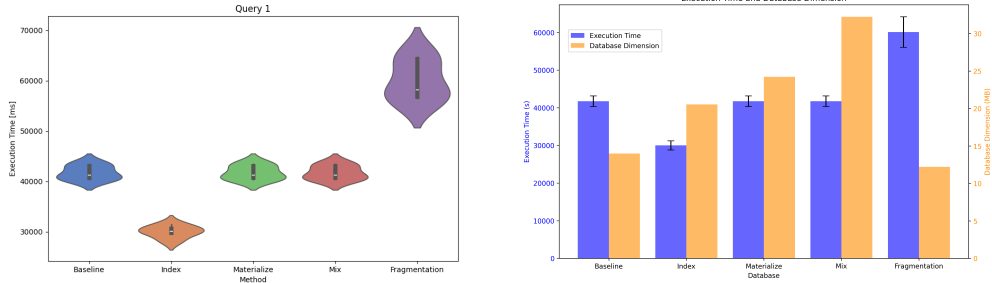


Figure 2: Query 1

the performances, leading to an execution time higher than the one of the baseline. This may be expected since the slicing condition is highly non selective, and scanning (almost) all the tables in the partition have an higher cost than scanning the same quantity of rows in a single table. The performances with the strategy of materialization didn't change since the materialized views are not involved in this query.

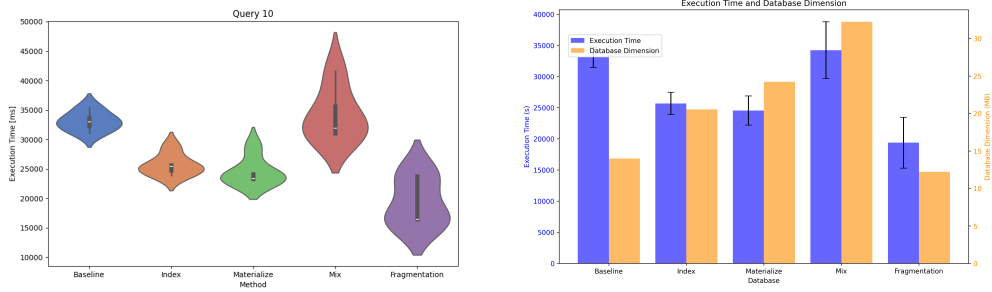


Figure 3: Query 10

Regarding query 10 we can see that the performances of all the approaches are comparable, with a slight improvement given by the fragmentation. Indeed the fragmentation on orderdate is fully exploited by the slicing condition that involved a three month linespan. Given this consideration it may be expected an even bigger improvement, but the scan of the lineitem partition to retrieve the subpartition with "l_returnflag=R" moderated the gains.

Query 14 is the one that obtained the most evident improvement with the fragmentation. This is explained by the fact that the lineitem table is the most computationally expensive to be scanned, by introducing the fragmentation on l_shipdate, the DBMS needs to consider only one subtable, furthermore the group by condition is facilitated by the subpartition on l_returnflag.

Regarding query 17 all the proposed form of optimization significantly improved the execution time with respect to the baseline (in principle with an execution time greater than 2 hours). A mixed approach of materialization and indexed allows to perform the slicing in a really efficient way, without having to perform the join at every execution since it's included in the pre-computed view.

As a general overview of this project we observe that a combination of different types

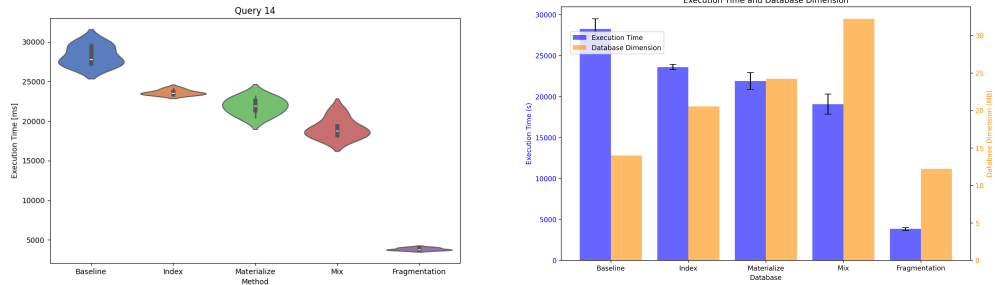


Figure 4: Query 14

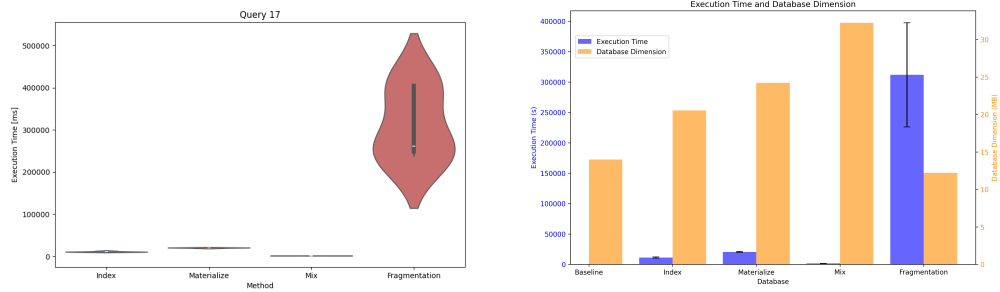


Figure 5: Query 17

of optimization is the best way to improve the performances of different queries. Also we can consider that sometimes it is important to find a balance between the benefits that a technique may carry in speeding up the computation, and the space (in term of memory) that this technique may require. Also it is important to average the effect of an optimization technique over different types of queries to avoid introducing too specific tools/procedures(?) [Supercazzola finale una volta che abbiamo i risultati]

A Appendix A

Complete statistics of the database.

Attribute	Distinct values	Min value	Max value
c_custkey	1500000	1	1500000
c_name	1500000	'Customer#000000001'	'Customer#001500000'
c_address	1500000	-	-
c_nationkey	25	0	24
c_phone	1499963	-	-
c_acctbal	818834	'-999.99'	'9999.99'
c_mktsegment	5	-	-
c_comment	1496636	-	-

Table 16: Costumer statistics

Attribute	Distinct values	Min value	Max value
s_suppkey	100000	1	100000
s_name	100000	'Supplier#000000001'	'Supplier#000100000'
s_address	100000	-	-
s_nationkey	25	0	24
s_phone	100000	-	-
s_acctbal	95588	'-999.92'	'9999.93'
s_comment	99983	-	-

Table 17: Supplier statistics

Attribute	Distinct values	Min value	Max value
p_partkey	2000000	1	2000000
p_name	1999828	-	-
p_mfgr	5	Manufacturer#1	Manufacturer#5
p_brand	25	Brand#11	Brand#55
p_type	150	-	-
p_size	50	1	50
p_container	40	-	-
p_retailprice	31681	900.91	2098.99
p_comment	806046	-	-

Table 18: Part statistics

Attribute	Distinct values	Min value	Max value
ps_partkey	2000000	1	2000000
ps_supplierkey	100000	1	100000
ps_availqty	9999	1	9999
ps_supplycost	99901	1.0	1000.0
ps_comment	7914164	-	-

Table 19: Partsupp statistics

Attribute	Distinct values	Min value	Max value
l_orderkey	15000000	1	60000000
l_partkey	2000000	1	2000000
l_supplierkey	100000	1	100000
l_linenumber	7	1	7
l_quantity	50	1	50
l_extendedprice	1351462	900.91	104949.5
l_discount	11	0.0	0.1
l_tax	9	0.0	0.08
l_returnflag	3	-	-
l_linestatus	2	-	-
l_shipdate	2526	'1992-01-02'	'1998-12-01'
l_commitdate	2466	'1992-01-31'	'1998-10-31'
l_receiptdate	2555	'1992-01-03'	'1998-12-31'
l_shipinstruction	4	-	-
l_shipmode	7	-	-
l_comment	34378943	-	-

Table 20: Lineitem statistics

Attribute	Distinct values	Min value	Max value
o_orderkey	15000000	1	60000000
o_custkey	999982	1	1499999
o_orderstatus	3	-	-
o_totalprice	11944103	838.05	558822.56
o_orderdate	2406	'1992-01-01'	'1998-08-02'
o_orderpriority	5	-	-

Table 21: Orders statistics

Attribute	Distinct values	Min value	Max value
r_regionkey	5	0	4
r_name	5	-	-
r_comment	5	-	-

Table 22: Region statistics

Attribute	Distinct values	Min value	Max value
n_nationkey	25	0	24
n_name	25	-	-
n_regionkey	25	0	4
n_comment	25	-	-

Table 23: Nation statistics