# Database Resource Allocation Based on Resilient Intermediates

Martin Kersten, Ying Zhang, Pavlos Katsogridakis,
Panagiotis Koutsourakis and Joeri van Ruth
*MonetDB Solutions*
Amsterdam, The Netherlands
<*lastname*>@monetdbsolutions.com

*Abstract*—Scale-out of big data analytics applications often does not pay off due to the poor performance in response time and the increasing bill due to a longer execution time on a resource limited machine. To enable a stable DBMS workload environment it helps to maintain several virtual machines with difference resource configurations (CPU, memory, disk, etc) hosting part of the database, so that users can send their tasks to those machines that have the best price/performance characteristics. This, however, requires a method to decide which VM should be used for a given query.

When choosing the VM, the memory usage of a query is a particularly important factor, especially for the main-memory (optimised) DBMSs which are generally used for analytical queries today. In this paper, we introduce MALCOM, a memory footprint predictor for queries based on resilient intermediates in MonetDB. Unlike traditional cost-based approaches, MALCOM uses an empirical approach (i.e. using the memory usage information of queries executed in the past) to incrementally update its model to improve its predictions. Our preliminary experiment results show that this approach is robust against varying data distributions.

## I. INTRODUCTION

Since the start of database research, database designers have keenly looked at the opportunities to use large, distributed processing platforms. Cluster-based products are readily available, such as in appliance products from Oracle Exadata [1], SQL Parallel Data Warehouse [2], IBM Blu [3] and Teradata [4], but they are often limited to a few tens of compute nodes. A plethora of research activities [5] has shown that in all but the simplest cases achieving a good performance is at least hard, especially when a query involves joins spread over multiple compute nodes and thus requires expensive data exchange.

The predominant way out nowadays, taken by NoSQL systems such as Cassandra [6] and Impala [7], is to address part of the problem space by focusing on select-aggregate queries. This choice has proven to be pivotal to support big data analytics in many real-world circumstances, as shown by the widespread use of Apache Spark [8]. The basic abstraction in Spark is a Resilient Distributed Dataset (RDD), which represents an immutable and partitioned collection of elements that can be operated on in parallel using operators, such as map, filter, persist and aggregates.

Although in many cases it is easy to scale-up for improved response time, partitioning a database to benefit from a low cloud service price tag and to overcome resource limita-tions of smaller machines is still a much sought-after skill. This product space is addressed by Snowflake [9] and AWS Redshift [10]. Snowflake has been designed from a cloud perspective, taking resource management as its key driving factor. It conceptually provides every user with a complete copy of the database and relies on multi-level caching. AWS Redshift is an improved version of PostgreSQL, which has been further tuned towards better I/O bandwidth use.

In this paper we take a fresh look at resource allocation for query processing in the context where intermediates in a query plan are fully materialised before passed on towards the next operator. This model fits not only the Apache Spark programming model, but also the query execution model of our database system MonetDB [11]. Resilient intermediates provide new avenues for query optimisation and scheduling as its underlying computation model is based on materialisa-tion of all intermediate steps. Furthermore, in most practical business analytic cases the past is a reasonable predictor of the future.

The main contributions of this paper are

- we develop a simulator, called MALCOM to predict the memory footprint for queries based on resilient interme-diates in MonetDB.
- We demonstrate that the approach is robust against vary-ing data distributions.
- We demonstrate the opportunities using an extensive evaluation against TPC-H and a real-world data set.

The approach taken here differs from traditional cost-based query optimisers [12] deployed in distributed database sys-tems by learning about actual resource claims over time, i.e. after each query execution we have precise knowledge of the resources used. This information can be harvested and used to predict future operations of a similar nature. The rationale stems from the common knowledge that any database application environment has a limited number of "business transactions" or "business intelligence templates" where only some parameters are changed with each call. This knowledge has been used in the past to drive development of DBA wizards [13] for index selection by humans and self-tuning optimisers [14] to avoid expensive join paths in individual queries.

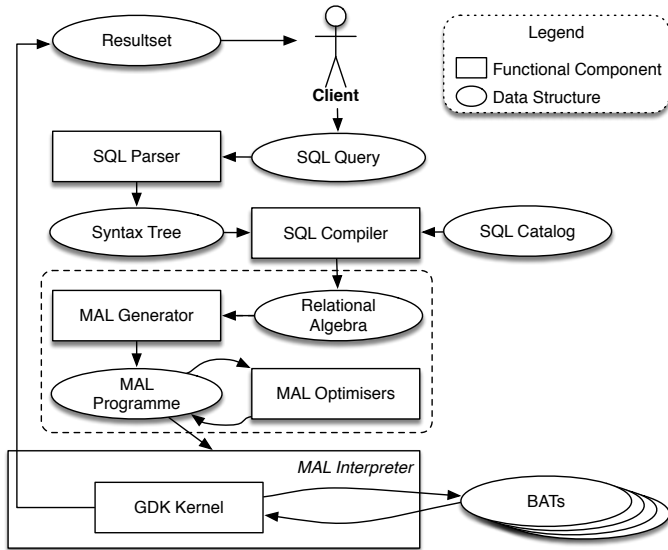**Paper Outline.** Section II provides a short overview of the

Fig. 1. MonetDB query execution architecture.

MonetDB architecture. Section III introduces the components and algorithms for our resource estimator MALCOM. Section IV illustrates the effectiveness of our approach in two use cases: TPC-H and air traffic.

## II. BACKGROUND

In this section we introduce the query execution engine of MonetDB and the profiling information available for our task.

### A. MonetDB Architecture

MonetDB is a widely used columnar DBMS that internally uses resilient intermediates to break up query processing in well identifying steps. A query plan is broken up into independent steps, glued together into a dataflow dependency graph. The dataflow graph is greedily consumed by the database kernel assigning a dedicated core to each operation. The resource pressure is kept at a minimum by trimming down the degree of parallel processing when the main memory resource is heavily used. The system can be instructed to produce an event record for each completed instruction. This provides a.o. insights into the input/output sizes and timing.

Figure 1 illustrates the components of MonetDB to execute an SQL query. The MonetDB Assembly Language (MAL)[1] is the MonetDB internal language into which SQL queries are compiled and executed. MAL is purely designed as intermediate language to express the relational operations. The SQL Parser and MAL Optimiser deploy well-known rewriting rules (e.g. parallelisation, and dead code/common expression/constant elimination) to reduce the intermediate sizes and processing time. They do not rely on any cost-model or pre-computed statistics.

The middle layer (in the dashed box) is a sequence of specialised optimisers that morph a logical plan produced by

---

[1]https://www.monetdb.org/Documentation/Manuals/MonetDB/MALreference

the SQL compiler into a physical execution plan containing relational operators expressed in MAL statements. The bottom layer (under the dashed box) contains the implementation of the MAL statements. Each operator takes as input the resilient intermediates produced by operators executed before or the persistent data on disk.

As an example, consider this simple SQL query: `SELECT COUNT(*) FROM sys._tables`, which is eventually translated into a physical execution plan in MAL statements to be executed by the MonetDB kernel. Below is an excerpt of the trace we have captured when executing our example query.

```
C_5=<tmp_2141>[100]:bat[:oid] := sql.tid("sys":str, "_tables":str);
X_8=<tmp_154>[100]:bat[:int] := sql.bind
        ("sys":str, "_tables":str, "id":str, 0:int);
X_17=<tmp_2022>[100]:bat[:int] := algebra.projection
        (C_5=<tmp_2141>[100]:bat[:oid], X_8=<tmp_154>[100]:bat[:int]);
X_18=100:lng := aggr.count(X_17=<tmp_2022>[100]:bat[:int]);
sql.resultSet("sys.L3":str, "L3":str, "bigint":str,
              64:int, 0:int, 7:int, X_18=100:lng);
```

This MAL program is straightforward. It first loads (`sql.bind`) and projects (`algebra.projection`) the data of one column (`sys._tables.id`) from the disk. Then the data is passed to `aggr.count` to compute the COUNT. Finally, `sql.resultSet` emits the query result.

Both an execution plan and an execution trace contain a sequence of MAL statements in the following general form:

```
VAR=<FILENAME>[COUNT]:VAR_TYPE :=
    MOD.FUNC(PARAM1=<FILENAME>[COUNT]:PARAM1_TYPE, ...);
```

Every function (`FUNC`) belongs to a module (`MOD`). The arguments are either typed scalar values (`:type`) or a reference to a column (`:bat[:type]`). If a variable (`VAR`) or parameter (e.g. `PARAM1`) refers to a column, which can be memory mapped, it is also tagged with its base `FILENAME` on disk and the number of values in this column (`COUNT`). The `=<FILENAME>[COUNT]` is an optional part of a MAL statement. It is added to the corresponding execution trace when a MAL execution plan is being carried out.

The MAL statements play an important role in MALCOM. Before executing a query, one can pass the query's MAL execution plan to MALCOM to get an estimation of the memory footprint of this query. After a query execution, one can pass the query's execution trace to MALCOM so that MALCOM can update its model with the statistics of the actual query execution, e.g. the number of values contained in the input and output variables of each MAL statement.

### B. MonetDB Profiling Information

The MonetDB kernel can be instructed to emit profiling events for the execution of MAL statements, e.g. by establishing a connection using MonetDB's profiling tool[2]. The computational and memory overhead introduced by this is negligible, since the profiling information is mostly already available as a by-product of query executions and the MonetDB server merely write it to a socket. It is up to the external tools to capture those events, and process and/or store them.

---

[2]https://www.monetdb.org/Documentation/Manuals/MonetDB/Profiler

| JSON object at "start" time | JSON object at "done" time |
|---|---|
| {"source":"trace",<br>"ctime":1528314717302449,<br>"module":"algebra",<br>"instruction":"projection",<br>"state":"start",<br>"usec":0,<br>"rss":87,<br>"size":0,<br>"nvcsw":1,<br>"stmt":"X_17?:= **algebra.projection**(?);",<br>**"ret"**:[{<br>  "index":"0",<br>  "name":"X_17",<br>  "alias":"sys._tables.id",<br>  "type":"bat[:int]",<br><br>  "count":"0",<br>  "size":0,<br>  **"eol":0**}],<br>**"arg"**:[{<br>  "index":"1",<br>  "name":"C_5",<br>  "kind":"transient",<br>  "bid":"863",<br>  "count":"92",<br>  "size":736,<br>  **"eol":1**},<br>{"index":"2",<br>  "name":"X_8",<br>  ...<br>  **"eol":1**}]} | {"source":"trace",<br>"ctime":1528314717302680,<br>"module":"algebra",<br>"instruction":"projection",<br>"state":"done",<br>"usec":230,<br>"rss":87,<br>"size":0,<br><br>"stmt":"X_17?:= **algebra.projection**(?);",<br>"ret":[{<br>  "index":"0",<br>  "name":"X_17",<br>  "alias":"sys._tables.id",<br>  "type":"bat[:int]",<br>  "kind":"transient",<br>  "count":"92",<br>  "size":368,<br>  **"eol":0**}],<br>**"arg"**:[{<br>  "index":"1",<br>  "name":"C_5",<br>  "kind":"transient",<br>  "bid":"863",<br>  "count":"92",<br>  "size":736,<br>  **"eol":1**},<br>{"index":"2",<br>  "name":"X_8",<br>  ...<br>  **"eol":1**}]} |

Fig. 2. JSON profiling objects produced for an `algebra.projection` operation.

Figure 2 shows an excerpt of the before/after profiling events produced for the `algebra.projection` operation when executing the MAL program above. The left column shows the event at the `"start"` and the right column the event when the operation is `"done"`. Differences between the two objects are marked in red. Of most interest to our estimation are the properties shown for the arguments (`"arg"`) and return variables (`"ret"`). For instance, in a `"ret"` object, the field `"size"` is a good estimation of how much memory the result set of this function consumes; while in an `"arg"` object, the field `"eol":1` indicates this argument has reached its end-of-life. This information together with the `"size"` allows us to estimate how much memory is freed after this operation.

## III. MALCOM MICRO MODELS

With an abundance of profiling events we can derive a micro-model for each MAL instruction to estimate their footprint. The goal of MALCOM is, given a MAL execution plan of an SQL query, to estimate the resource needs by *only using information from our memory footprint estimation model*.

The algorithm to estimate an upper bound of the memory needed to execute a MAL plan, is shown in pseudo code below. When a MAL plan is received, MALCOM first annotates each MAL statement with an estimation of how much memory it will consume (`i.mem_fprint`) and release (`i.free_size`). Then the algorithm iterates over the MAL plan (i.e. the `mal_statements` below). At each iteration, it adds the memory footprint of this MAL statement (`i.mem_fprint`) to the current total memory consumption (`curr_mem`) and updates `max_mem` if necessary. After that, it adjusts `curr_mem` with the amount of memory that will be freed by this statement (`i.free_size`).

```
max_mem  = 0
curr_mem = 0
for i in mal_statements:
  curr_mem += i.mem_fprint
  max_mem = max(max_mem, curr_mem)
  curr_mem -= i.free_size
```

After the execution of the MAL plan, we update our memory footprint estimation model with the observed execution information. We initialise the estimation model with basic column statistics (`min`, `max`, `count`, etc.) that can be gathered using MonetDB's `ANALYZE` command.

The estimation model is built by dividing MAL instructions with similar functionality (most of them represent a relational operator each) into several groups and abstracting away their specific signatures. Currently, the model includes ~10 groups. We briefly analyse each of them below. Note that we only consider bulk operators here (i.e. taking columns as operands), which are the default ones in MonetDB.

### A. Load instructions

A `bind` instruction loads (or memory maps) a column into memory, thus the return size is the size of the column. This is a worst case assumption, because in practice not all of the column needs to be loaded.

### B. Arithmetic Operators

These operators always return the same number of values as their operands (MonetDB requires both operands to have equal size). However, the data type of the output can be a larger-sized data type than both operands to capture possible overflow. Hence, their result size is computed as:

```
arith.rsize = sizeOf(arg1) * sizeof(ret.datatype)
```

### C. Aggregate Operators

This category includes operations such as `sum`, `avg`, `min`, `max`, `count`, `single` and `dec_round`. The number of values returned by these operators equals the number of groups in which the input data column is divided (by earlier `GROUP BY` statements, or 1 if there is no `GROUP BY`). The most general signature of these operators takes two operands: `arg1` is a column containing the actual values to work on; `arg2` is a column containing the group IDs, one for each value in `arg1`. The output size of an aggregate operator is computed by multiplying the number of unique values in `arg2` with the size of the return data type.

```
aggr.rsize = COUNT(DISTINCT arg2) * sizeof(ret.datatype)
```

### D. Limit Operators

This group includes `firstn` and `sample`. They return at most N values from its input column `arg` as specified by the limit. Hence, their output size is computed as:

```
limit.rsize = MIN(COUNT(arg), N) * sizeof(arg.datatype)
```

## E. Grouping Operators

MonetDB currently has 24 grouping operators for different situations. For instance, the position of the input data column (`arg1`) in a `GROUP BY` SQL clause determines the use of a `GROUP` operator or a `SUBGROUP` operator in MAL. More variations of `GROUP` or `SUBGROUP` operators are used depending on the availability of auxiliary information (e.g. some statistics of the input column). However, all grouping operators generally return three columns of results: (i) a `groups` column containing the group IDs, one for each value in `arg1`; (ii) an `extents` column containing the `OID` (MonetDB internal type for Object Identifiers, denoting positions of data values in a column) of a representative of each group; and (iii) a `histo` column containing the number of values in each group corresponding the values in `extents`. The data type of `groups` and `extents` are both `OID`, and the data type of `histo` is `LNG` (MonetDB internal type for Long integers). The number of values in `extends` and `histo` is the same, and is estimated using a simple kNN algorithm based on the statistics of previous queries or basic statistics of the involved columns. Putting everything together, the total output size of a grouping operator is estimated as:

```
group.rsize = COUNT(arg1) * sizeof(OID) +
              estimate_nr_groups(arg1) * (sizeof(OID) + sizeof(LNG))
```

## F. Set Operators

For the set operators we mostly compute an upper bound of the result size using heuristics:

```
unionall.rsize = sizeof(arg1.datatype) * (COUNT(arg1) + COUNT(arg2))
union.rsize    = sizeof(arg1.datatype) *
                 (COUNT(DISTINCT arg1) + COUNT(DISTINCT arg2))
ntsct.rsize    = sizeof(arg1.datatype) *
                 MIN(COUNT(DISTINCT arg1) + COUNT(DISTINCT arg2))
xcpt_all.size  = sizeof(arg1.datatype) * COUNT(arg1)
xcpt.size      = sizeof(arg1.datatype) * COUNT(DISTINCT arg1)
```

Both `UNION` and `UNION ALL` return a concatenation of their two input columns `arg1` and `arg2`, except that `UNION` eliminates the duplicates in its result. Hence, `unionall.rsize` is the precise result size of a `UNION ALL`, while `union.rsize` is an upper bound of the result size of a `UNION`, because its computation does not exclude unique values that exist in both `arg1` and `arg2`.

`INTERSECT` returns values that exist in both its input columns `arg1` and `arg2` with duplicates eliminated. Hence, `ntsct.rsize` is an upper bound of the result size, as its formula does not exclude unique values that are only in `arg1` or only in `arg2`.

Both `EXCEPT` and `EXCEPT ALL` return all values that are in its first input column `arg1` but not in the second column `arg2`. In addition, `EXCEPT` eliminates duplicates. Therefore, both `xcpt_all.size` and `xcpt.size` are upper bounds of their respective result sizes, since their computations does not exclude (unique) values that also exist in `arg2`.

## G. Projection Operators

Projection operators extract a small part of a column. The arguments are a candidate list `cand` containing the OIDs of the to-be-projected values and a reference to the (persistent) column `col`. The number of elements in the output equals the number of elements of the candidate list. Hence, their exact output size is computed as:

```
proj.rsize = COUNT(cand) * sizeof(col.datatype)
```

## H. Selection Operators

This operator group includes the filter operations `theta-select` and `select`. For these operators, we know that the output is always smaller than or equal to the candidate tuples considered. To estimate the result size, traditional cost-based models assume a uniform distribution of the data and calculate the fraction of the domain, i.e. the selectivity factor. In practice, however, the uniform distribution of the data assumption does not always hold, so these models have limited accuracy. In our model, we keep the results of a series of actual filter operations, so as to use them to find a "historical nearest-neighbour" for any filter operation in a MAL plan whose cost we need to estimate.

The select operators are abstracted into a single template with three operands `sel(col, range, op)`, where `col` is a reference to the (persistent) column, `range` is the selection range (low, high), and `op` is the comparison operator ($<$, $>$, $<=$, $>=$, etc). In our model, we keep a dictionary of all the selections executed so far in the format of this signature, with an extra attribute `cnt` to denote the number of values selected.

The estimation for a selection operator `sel(col, range, op)` works as follows. First, we find in the `dictionary` records of all previous selections on the same `col` with the same `op`. Then, we use a $k$ nearest neighbour (kNN) procedure to find the 5 nearest records based on the selection `range`. Next, for each of the 5 records, we extrapolate the number of selected values based on the selectivity and input column size. Finally, we compute the estimated memory footprint of this selection as the average of the 5 extrapolations, multiplied by the data size of the input column. This estimation procedure is shown in the pseudo code below:

```
extrap = 0
for dict in kNN(dictionary, sel, 5)
  extrap += dict.cnt * (COUNT(sel.col) / COUNT(dict)) *
            (sel.range / dict.range)
sel.rsize = extrap/5 * sizeof(sel.col.datatype)
```

## I. Join Operators

For a cross product of two columns (`col1`, `col2`) we know it will return `COUNT(col1) * COUNT(col2)` number of values. The `cross product` operator of MonetDB takes two data columns as its inputs, and returns two columns where each column contains OIDs referring to data values in an input column. The two OID columns together denote how the values from the input columns are aligned in the cross product result. So the exact output size of a cross product is computed as:

```
cp.rsize = COUNT(col1) * COUNT(col2) * sizeof(OID) * 2
```

The estimation model for the other join operators is similar to that of selection operators. Again, the signatures of all join operators can be abstracted into a single one with three operands `join(col1, col2, op)`, where `col1` and `col2` are the input column, and `op` the join operator (eq, left, outer, etc).

We also keep a dictionary of all the previous joins in the format of this signature, annotated with an extra value `cnt` to denote the number of values returned by that particular join operation. To estimate the result size for a join operator `join(col1, col2, op)`, we first find in the dictionary records of all previous selections on the same columns with the same `op`. Then, we run a kNN to find the 5 nearest records based on the sizes of the input columns. Finally, we extrapolate the result count based on the input column sizes, and compute the result size (like cross product, two OID columns are returned). The pseudo code is shown below:

```
extrap = 0
for dict in kNN(dictionary, join, 5)
  extrap += dict.cnt * (COUNT(join.col1) / COUNT(dict.col1)) *
                       (COUNT(join.col2) / COUNT(dict.col2))
sel.rsize = extrap/5 * sizeof(OID) * 2
```

**Optimiser Simulator Architecture.** With the micro-models in place, we can use a simple MAL-simulator to obtain a fairly accurate indication of the memory footprint. In its basic form it simulates a sequential execution of the query. The full-blown version uses the same scheduling method as within the MonetDB code base to approximate the parallel behaviour.

## IV. EVALUATION

In this section we discuss the results obtained with the MALCOM prototype. For our experiments we used both TPC-H and the air traffic[3] benchmarks. The former is the baseline against which most database systems are evaluated. Its major weakness is the uniform data distribution. Although TPC-H simplifies the analysis of the quality of a space predictor, it is not representative of real-life workloads. Therefore, we also use the air traffic benchmark, which consists of a single table with >120M rows and 100 columns of flight information. The data in this benchmark is skewed.

### A. TPC-H

To test MALCOM against TPC-H we created a query generator, which changes the parameters in the official benchmark queries. As a training set for each query, we randomised every selection point and range to produce 200 random versions of the query. As a test set we used the original queries. A key factor is to understand how long it takes before MALCOM's prediction converges to the actual memory footprint. Therefore, we update MALCOM's model with the traces of the training queries one by one. After each iteration, we evaluate the accuracy of MALCOM's prediction using the test queries. We have run all 22 TPC-H queries. Most of their predictions rapidly converge to the actual memory footprint due to the uniform data distribution of TPC-H. Hence, for reasons of space, we only show the results of some more remarkable queries here.

Initial experiments led to the results such as shown in Figure 3, where the x-axis shows the number of training queries used and y-axis shows the percentage of deviation of the estimation from the actual footprint. Q1 is a simple large

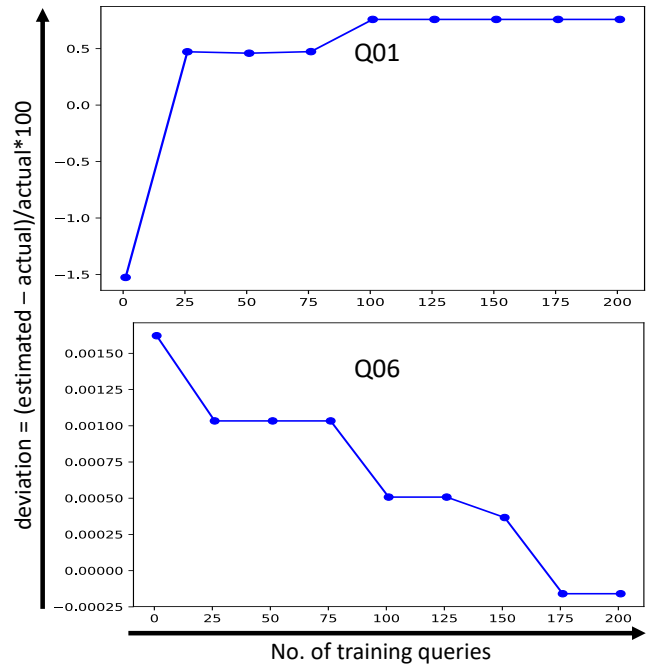[3]https://github.com/MonetDBSolutions/airtraffic-benchmark



Fig. 3. TPC-H queries Q1 and Q6 prediction deviations

scan followed by an aggregation. The sole parameter is the value range. The experiment shows a steep learning curve, which can be attributed to the uniform distribution. However, it also results in a little overfitting. Q6 is also mostly a simple scan and aggregate query, but here the number of tuples are less. This experiment shows that MALCOM takes much longer to reach an almost perfect prediction of the memory footprint.

### B. Air Traffic

The air traffic benchmark is a real-world example of business intelligence application. The data is represented by a single table, and it is highly skewed and sparse. The benchmark contains 19 queries range from simple select-group-aggregate to complex self joins that require careful processing strategies. We have created training queries and test queries, and evaluated MALCOM's prediction in the same way as what we did with the TPC-H evaluation. We tested MALCOM with all 19 queries. In general, the prediction gradually converges to the actual memory footprint. Here we high light several remarkable cases.

Figure 4 show how MALCOM's predictions converge for queries 4, 10, 15 and 19. Again we can observe both some under- and over-estimations, which all improve over time. Albeit they take more queries to stabilise than in the TPC-H experiments. This is because the air traffic data is highly skewed. In addition, the graphs of Q04, Q15 and Q19 illustrate that MALCOM does not always improve its predictions monotonically, which is also caused by the skewness of data.

## V. RELATED WORK

The approach taken in this project can best be compared with the long tradition in database query optimisers and
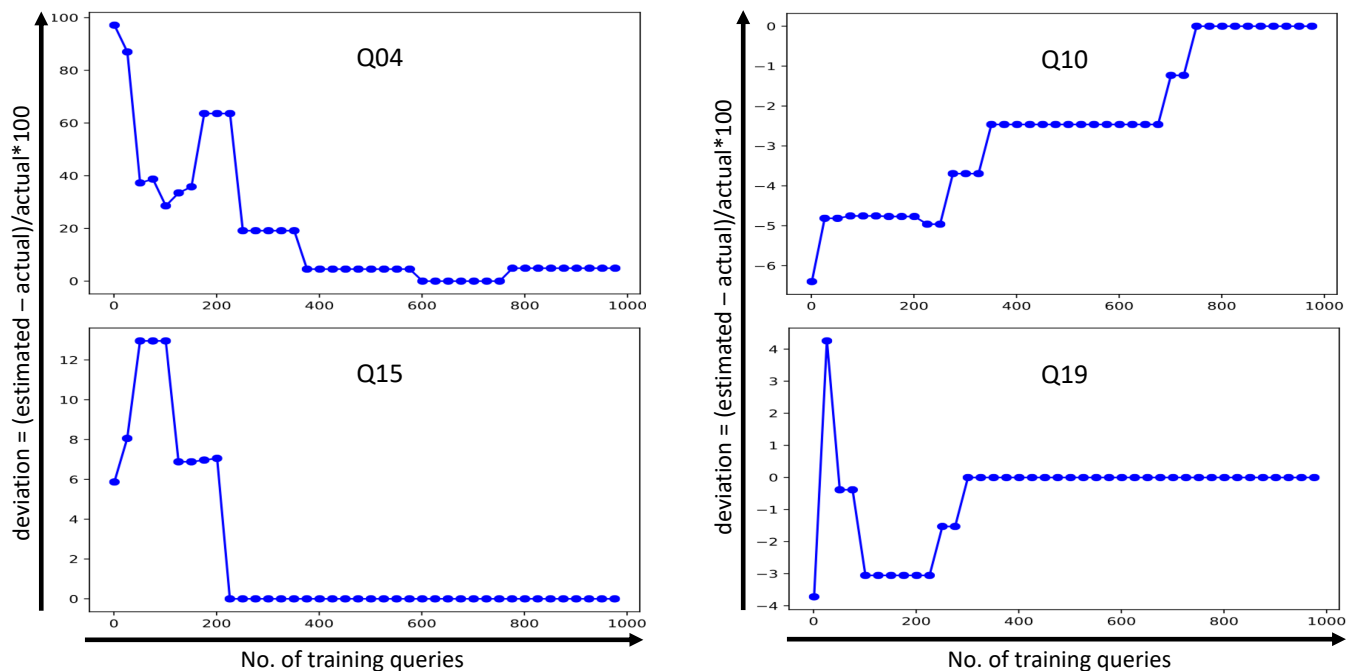
Fig. 4. Air traffic queries prediction deviations

database design wizards. They all collect query traces from an actual production system and use them to derive, e.g., an optimal set of search accelerators [15]. This process seeks a balance between index creation and maintenance, but primarily deals with performance optimisation. The memory footprint is of less concern. Alternatively, it extends the work on gathering query traces to (semi-)automatically improve the cost model for query optimisation [16]. A better statistics improves both performance and resource use. All these systems focus on a relative small and fixed compute cluster or database appliance.

In a more recent project [17] the authors gather sub-plans from the query trace log and use it as the building block for new queries. They show that reuse of good plans, i.e. based on past behaviour, leads to both a faster optimisation step and overall better performance. MALCOM does not address the optimiser itself, but assumes that a physical plan has already been produced. It merely determines which virtual machine can handle the plan comfortably.

## VI. SUMMARY AND OUTLOOK

In this paper we proposed a novel resource allocation technique for a database engine using resilient intermediates. The MALCOM prototype is a crucial tool in the design of cloud-based database management solution. The initial results are promising, i.e. a relatively low number of queries are sufficient to get a good memory footprint estimation.

In the near future, we plan to extend MALCOM to also look at the memory footprint in relationship of the parallel execution. This should lead to a lower bound on the memory footprint at the cost of running slower. Both memory footprint and degree of parallelism enables users to optimise their systems based on costs or response times.

## REFERENCES

[1] M. Bach *et al.*, *Expert Oracle Exadata*, 2nd ed. Berkely, CA, USA: Apress, 2015.

[2] T. Stöhr, H. Märtens, and E. Rahm, "Multi-dimensional database allocation for parallel data warehouses," in *Proceedings of the 26th International Conference on Very Large Data Bases*, 2000.

[3] V. Raman *et al.*, "Db2 with blu acceleration: So much more than just a column store," *Proc. VLDB Endow.*, vol. 6, no. 11, Aug. 2013.

[4] W. O'Connell *et al.*, "A teradata content-based multimedia object manager for massively parallel architectures," in *Proceedings of the ACM SIGMOD International Conference on Management of Data*, 1996.

[5] M. T. Ozsu, *Principles of Distributed Database Systems*, 3rd ed. Upper Saddle River, NJ, USA: Prentice Hall Press, 2007.

[6] "Apache cassandra," http://cassandra.apache.org.

[7] "Apache impala," https://impala.apache.org.

[8] "Apache spark," https://spark.apache.org.

[9] B. Dageville *et al.*, "The snowflake elastic data warehouse," in *Proceedings of the International Conference on Management of Data*, 2016.

[10] A. Gupta *et al.*, "Amazon redshift and the case for simpler data warehouses," in *SIGMOD*, 2015, pp. 1917–1923.

[11] "MonetDB," http://www.monetdb.org.

[12] S. Chaudhuri, "An overview of query optimization in relational systems," in *Proceedings of the Seventeenth ACM SIGACT-SIGMOD-SIGART Symposium on Principles of Database Systems*, ser. PODS '98, 1998, pp. 34–43.

[13] https://www.microsoft.com/en-us/research/project/autoadmin/.

[14] https://www.ibm.com/us-en/marketplace/datarcs-optimizer.

[15] S. Chaudhuri and V. R. Narasayya, "Self-tuning database systems: A decade of progress," in *VLDB*, September 2007, pp. 3–14.

[16] V. Markl, G. M. Lohman, and V. Raman, "LEO: an autonomic query optimizer for DB2," *IBM Systems Journal*, vol. 42, no. 1, 2003.

[17] B. Ding *et al.*, "Plan stitch: Harnessing the best of many plans," *PVLDB*, vol. 11, no. 10, pp. 1123–1136, 2018.