

ARe-SQL: An Online, Automatic and Non-Intrusive Approach for Rewriting SQL Queries

Arlino H. M. de Araújo^{1,2}, José Maria Monteiro², José Antônio F. de Macêdo²,
Júlio A. Tavares³, Angelo Brayner³, Sérgio Lifschitz⁴

¹ Federal University of Piauí, Brazil

² Federal University of Ceará, Brazil

{arlino, monteiro, jose.macedo}@lia.ufc.br

³ University of Fortaleza, Brazil

{julio, brayner}@unifor.br

⁴ Pontifical Catholic University of Rio de Janeiro

sergio@inf.puc-rio.br

Abstract. Database applications have become very complex, dealing with a huge volume of data and database objects. Concurrently, low query response time and high transaction throughput have emerged as mandatory requirements. Among other possible interventions regarding database performance, SQL query rewriting has proved an efficient alternative. The idea is to write a new SQL statement equivalent to the initially formulated statement, where the new SQL statement provides performance gains w.r.t. query response time. In this article, we propose an online, automatic and non-intrusive approach for rewriting SQL queries, denoted ARe-SQL. ARe-SQL rewrites SQL statements, using a set of eleven heuristics. Based on ARe-SQL, two different strategies for SQL query rewriting, denoted assisted and automatic, were implemented. These strategies were evaluated in three different DBMS considering three distinct scenarios. The results indicate that both strategies can provide performance gains, which, in some cases, reach 92.5%.

Categories and Subject Descriptors: H.3.2 [Information Storage and Retrieval]: Information Storage—*File organization*; H.3.3 [Information Storage and Retrieval]: Information Search and Retrieval—*Retrieval models*

Keywords: database tuning, query processing, query rewriting

1. INTRODUCTION

Database queries are expressed by means of high-level declarative languages, such as SQL (Structured Query Language), for instance. Such queries are submitted to the query engine, which is responsible for processing queries in database management systems (DBMS). Thus, query engines should implement four main activities: query parsing, logical execution plan (LEP) generation, physical execution plan (PEP) generation and PEP execution. LEP and PEP generation and PEP execution are often called the query optimization. The main functionality of the query optimizer is to find an efficient execution plan for a given query. For that purpose, cost-based optimizers search, in a large space of alternatives, called search space, the PEP, which is expected to be executed in the less amount of time.

An efficient strategy to support optimizers in the task of building efficient execution plans is query rewriting. The rewriting technique is based on the principle that two queries are equivalent if and only if their execution produces the same result [Ramakrishnan and Gehrke 2002]. Thus, rewriting a query Q_a consists in specifying a new SQL statement Q_b equivalent to Q_a , where Q_a and Q_b return the same result and the execution of Q_b provides performance gains w.r.t. Q_a 's execution. In order to illustrate performance gains achieved by means of query rewriting, consider the queries Q_a and Q_b

Copyright©2014 Permission to copy without fee all or part of the material printed in JIDM is granted provided that the copies are not made or distributed for commercial advantage, and that notice is given that copying is by permission of the Sociedade Brasileira de Computação.

```

SELECT * FROM orders
WHERE o_totalprice < ALL (SELECT o_totalprice
FROM orders
WHERE o_orderpriority = '2-HIGH')
Execution time = 13,685 ms
    
```

Fig. 1. SQL expression for Q_a

```

SELECT * FROM orders
WHERE o_totalprice < (SELECT MIN(o_totalprice)
FROM orders
WHERE o_orderpriority = '2-HIGH')
Execution time = 1,825ms
    
```

Fig. 2. SQL expression for Q_b

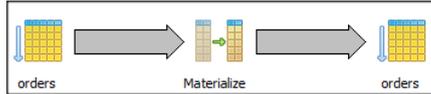


Fig. 3. PEP P_{Q_a} yielded by PostgreSQL for Q_a

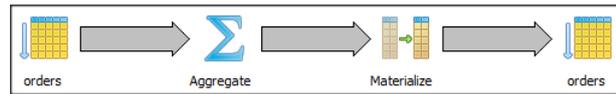


Fig. 4. PEP P_{Q_b} yielded by PostgreSQL for Q_b

depicted in Figures 1 and 2, respectively. In order to execute those queries, the TPC-H benchmark’s database has been created in a PostgreSQL database server.

Looking more closely to Figures 1 and 2, one can observe that Q_a and Q_b are equivalent. Nonetheless, Q_a execution lasts $13.685ms$, while Q_b is executed in $1.825ms$. In fact, Q_b was obtained by removing the *ALL* operator from Q_a and including the *MIN* operator in subquery of Q_a .

Therefore, we may formulate the following hypothesis, the way a query Q is formulated in terms of SQL commands induces the query optimizer to produce a given PEP for Q . To show the veracity of our hypothesis, consider the execution plans for queries Q_a and Q_b depicted in Figures 3 and 4. Those plans have been yielded by PostgreSQL’s query engine. The PEP showed in Figure 4 is more efficient ($1,825$ ms) than PEP depicted in Figure 3 ($13,685$ ms). This fact stems from the least amount of data materialized in the “Materialize” operation, which is performed after the “Aggregate” operation (in the PEP illustrated in Figure 4).

Database query tuning tools explore these aspects in order to help query optimizers to produce better PEPs. These tools may be categorized as (i) online or offline; (ii) automatic or non-automatic; and (iii) intrusive or non-intrusive. Offline approaches rely on human intervention for specifying the time window for capturing the workload. In addition, they do not allow the system to react to workload dynamics [Weikum et al. 2002; Chaudhuri and Weikum 2006]. On the other hand, an online tuning tool is able to dynamically capture database workload. Automatic query tuning tools have the ability of automatically triggering adjustments to improve query performance. Conversely, non-automatic tools transfer to the DBA (Database Administrator) the responsibility for making such adjustments. This includes the definition of the most appropriate moment to apply the recommended query changes and to choose, among the suggested actions, those that in fact should be executed and in which order. Intrusive solutions are those that require changes in the DBMS code and, therefore, are tightly coupled to a particular DBMS. Non-intrusive solutions are those that are disconnected from a particular DBMS implementation and that are usually not affected by any new DBMS release.

In fact, there are some tools that aim at tuning database query performance by providing rewriting recommendations. However, they adopt an offline approach, since they provide recommendations only when a human (e.g., DBA) requires such recommendations. Most of them are intrusive, which means that they are DBMS-specific. Additionally, they present a strong issue, for a given query they may provide several rewriting recommendations. In this case, the DBA is in charge to choose one of them without any indication (from the tool) that the chosen recommendation will provide the best performance gains regarding response time.

In this article, we propose an online, automatic and non-intrusive approach for rewriting SQL queries, denoted ARe-SQL. The proposed approach provides the necessary support to implement two different strategies for SQL-query rewriting. In the first strategy, denoted automatic, ARe-SQL

autonomously captures SQL queries submitted by users or applications and identifies which queries should be rewritten. Thereafter, queries are sent to the database query engine. The second strategy, denoted assisted, works as an autonomous advisor, which analyze previously executed SQL statements in order to recommend (through alerts, reports and wizards) SQL tuning opportunities. In both strategies, the rewritten queries profit from performance gains as we show in Section 4.

The key goal of ARe-SQL is to assist optimizers in building efficient query physical execution plan. In a proactive manner, ARe-SQL rewrites SQL statements which will induce query engine to produce better plans than those that would be built by executing the original SQL statements. We advocate that the main features presented by ARe-SQL are the following.

- Implementation of a non-intrusive strategy.** Due to this feature ARe-SQL may be applied to different database systems. In other words, it is not DBMS-specific.
- Autonomy.** ARe-SQL has the ability of automatically triggering query rewrites or recommendations. This includes the definition of the most appropriate moment to apply the recommended rewrites and to choose, among the suggested actions, those that in fact should be executed and in which order.
- Low-overhead tuning activity.** By using our approach, the database system keeps running while our automatic tool runs continuously capturing the workload, in a lightweight manner.

We have implemented the proposed approach, including automatic and assisted strategies. These strategies were evaluated in three different DBMS (Oracle, SQL Server and PostgreSQL) considering three different scenarios, comprising the TPC-H benchmark. The experimental results indicate that both strategies can provide performance gains. In some cases, such performance gains reach 92.5%.

The rest of the article is structured as follows. Section 2 discusses most relevant related work. The proposed approach is presented and analyzed in Section 3. In turn, Section 4 analyzes experimental results. Finally, Section 5 concludes this article.

2. RELATED WORK

Bruno et al. [2009] propose a framework, called Power Hints, which enables the creation and use of hints. In the proposed framework, hints are created by means of regular expressions, making it easy and flexible to create restrictions for query execution plans, allowing more precise tuning.

Herodotou and Babu [2009] present a tool called zTuned. Its purpose is to facilitate experiments related to tune SQL queries. The tool produces sets of plans with operators that have the same cardinality (called neighborhood plans) and chooses the optimal plan among the best plans for each neighborhood, using the mechanism of cost estimative of the DBMS. The tool works decoupled from the query optimizer and can potentially be used with any DBMS that has a cost-based optimizer.

Belknap et al. [2009] describe a new Oracle 11g database feature, called Automatic SQL Tuning. This tool extends the SQL Tuning Advisor (present on Oracle 10g) fully automating the SQL tuning workflow and solving some SQL performance problems without any DBA intervention. Quest SQL Optimizer for Oracle [Quest 2010] provides semantically equivalent SQL statements for a given query. In this case, the DBA should pass to the tool which query should be analyzed. Moreover, it is up to the DBA to choose one of the several recommended SQL statements.

Markl et al. [2003] discuss an autonomic query optimizer, called LEO (LEarning Optimizer for DB2), that automatically self-validates its model without requiring any user interaction to repair incorrect statistics or cardinality estimates. By monitoring queries as they execute, the autonomic optimizer compares the optimizer's estimates with actual cardinalities at each step in a PEP, and computes adjustments to its estimates that may be used during future optimizations of similar queries.

The tool IBM Optim Development Studio [IBM 2010] collects query performance metrics in DB2. The metrics are query execution frequency, cost and time. Thus, with such metrics, it is possible to identify queries, for which query rewriting recommendation are worthwhile. It is important to mention that this tool collects metrics for a given period of time (defined by the DBA).

Embarcadero DB Optimizer XE [Embarcadero 2010] identifies hints which should be encoded in a given SQL statement, i.e., it does not propose query rewritings. Its goal is mainly to eliminate unnecessary outer joins and Cartesian products. Additionally, Embarcadero may provide recommendations for improving index configuration.

Therefore, most of the investigated SQL tuning tools (Automatic SQL Tuning Advisor [Dageville and Dias 2006], IBM Optim Development Studio [IBM 2010], Embarcadero DB Optimizer XE [Embarcadero 2010] and Quest SQL Optimizer for Oracle [Quest 2010]) adopt an offline approach. In this sense, they transfer to the DBA the responsibility for defining the set of queries to be evaluated for choosing one of the several alternatives provided by them. Observe that DBAs work in a reactive way, i.e., they only trigger a tool or an advisor when the problem already exists. Furthermore, after identifying the problem (a time consuming query) and a possible solution, a DBA should rewrite the SQL statement, test and send it to a programmer to change the application code which is using the SQL statement. This whole process is rather time consuming. On the other hand, Oracle Automatic SQL Tuning is an online tool. However, it's an intrusive solution and runs only with Oracle 11g. The approach proposed in this article differs from previous work since ARe-SQ is an online, automatic and non-intrusive approach. Furthermore, most of these tools are proprietary, involving high acquisition costs. These aspects hinder the performance comparison between ARe-SQL and previous solutions.

Finally, this article extends our previous work [Araújo et al. 2013] in some aspects. First, we present a more detailed related work. Next, we include a discussion about how ARe-SQL was implemented in both assisted and automatic approaches. We also present new experiments to evaluate overhead.

3. ARE-SQL

ARe-SQL's main functionality is to influence query optimizer to choose an effectual query execution plan. ARe-SQL proactively rewrites SQL statements which will induce query engine to produce better plans. In order to achieve its goal, ARe-SQL works directed by a set of heuristics. The heuristics consist of rules to identify potential opportunities for tuning SQL statements and the ways to rewrite the statements. Table I brings the eleven heuristics applied by ARe-SQL. Furthermore, it indicates whether or not each heuristic is currently implemented by three major DBMS: PostgreSQL 8.3, Oracle 11g and SQL Server 2008.

3.1 ARe-SQL Architecture

In order to implement ARe-SQL, two different approaches were employed: assisted and automatic. For each approach, a different tool was implemented. However, both approaches present the following features. First, it is non-intrusive, since it is completely decoupled from the source code of the DBMS. This allows that the conceived solution can be used with any DBMS. Second, it is independent of location. it can run on a machine different than that used to host the DBMS, not consuming server resources where the DBMS is hosted.

3.1.1 Assisted Approach. A query-rewriting tool based on assisted approach consists of an advisor which has the ability: *(i)* to capture the previously SQL statements executed by DBMS; *(ii)* to analyze these statements, and; *(iii)* to recommend (through alerts, reports or wizards) SQL tuning opportunities. Thus, the advisor can identify SQL statements that could be rewritten. Additionally, such a type of tool allows the DBA to interact with the tuning process. For instance, a DBA may select a subset of available heuristics to be applied for the SQL-query tuning process. In other words,

Table I. Heuristics for SQL Tuning

	Heuristics for SQL Tuning	PostgreSQL	Oracle	SQL Server
H1	Transform queries which create and use temporary table into an equivalent sub-query.	No	No	No
H2	Eliminate unnecessary GROUP BY.	No	No	No
H3	Remove having clause whose predicates do not have any aggregate function. The predicates should be moved to a WHERE clause.	No	No	No
H4	Change query with disjunction in the WHERE to a union of query results.	No	Yes	No
H5	Remove ALL operation with greater/less-than comparison operators by including a MAX or MIN aggregate function in the subquery.	No	Yes	No
H6	Remove SOME operation with greater/less-than comparison operators by including a MAX or MIN aggregate function in the subquery.	No	Yes	No
H7	Remove ANY operation with greater/less-than comparison operators by including a MAX or MIN aggregate function in the subquery.	No	Yes	No
H8	Replace IN set operation by a join operation.	No	Yes	No
H9	Eliminate unnecessary DISTINCT.	No	No	No
H10	Move function applied to a column index to another position in the expression.	No	No	No
H11	Move arithmetic expression applied to a column index to another position in the expression.	No	No	No

the DBA can specify that some heuristics are unnecessary or inappropriate for a given database system or a given statement.

We have implemented a tool, called ARE-SQL Advisor, which implements an assisted approach for ARE-SQL. Figure 5 depicts an abstract model of ARE-SQL Advisor architecture. To utilize ARE-SQL Advisor in a given DBMS, it is necessary to instantiate two specific drivers: a driver for workload access (DWA) and a driver for statistics access (DSA). These drivers are the only components of ARE-SQL Advisor which are DBMS-specific. Nevertheless, these drivers are implemented in a non-intrusive fashion. We have so far implemented drivers for the following DBMS: PostgreSQL 8.4, Oracle 11g and SQL Server 2008 (see Section 4). The main components of the architecture illustrated in Figure 5 are the following.

- Agent for Workload Obtainment (AWO): This agent observes the operations submitted to the DBMS and retrieves the SQL statements and then stores them in the local metabase (Local MetaData - LM). This agent can be configured to run continuously (On-the-fly).
- Local Meta Data (LM): Database that stores workloads captured by AWO.
- Driver for Workload Access (DWA): This component enables ARE-SQL Advisor to access the metabase (catalog) of a given DBMS.
- Agent for Statistics Obtainment (ASO): This component is in charge of accessing statistics information of the target DBMS, such as table cardinality, the amount of disk pages required to store a database table, the height of (B^+ tree) index structures and so forth.
- Driver for Statistics Access (DSA): driver that allows the ASO to retrieve statistics for a specific DBMS.
- Heuristic Set (HS): set of heuristics used by the agents to identify SQL statements with oportunities of tuning and in order to rewrite them. HS is composed of 11 heuristics, which are depicted in Table I. However, new heuristics can be defined and inserted into HS.
- Agent for SQL Tuning (AST): Component responsible for tuning a particular SQL statement using the set of heuristics (HS).

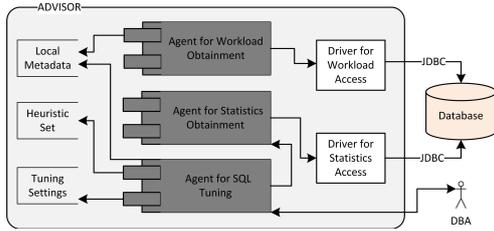


Fig. 5. ARE-SQL Advisor’s Architecture

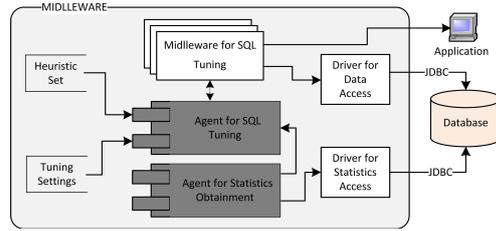


Fig. 6. Architecture of ARE-SQL’s Automatic Approach

—Tuning Settings (TS): This component is a preference file containing pairs $\langle \text{SQL statement } Q, \text{subset of heuristics } H \rangle$ defined by the DBA. Each pair indicates a heuristic subset chosen by the DBA from HS that he/she wants to be applied for tuning Q . If there isn’t an entry in this file for a given statement Q , all 11 heuristics are applied for tuning Q (in both, Assisted and Automatic Approaches).

3.1.2 *Automatic Approach.* The automatic approach consists of a middleware, ARE-SQL Mid, that works between the application and the DBMS. This middleware is responsible for: (i) automatically receiving SQL statements sent by applications; (ii) analyzing and rewriting them whenever necessary, and; (iii) submitting the statements (rewritten or not) to the DBMS. The architecture designed for the automatic approach is illustrated in Figure 6. The main components of this architecture, in addition to components that are also present in Figure 5, are the following.

- Middleware for SQL Tuning (MST): is responsible for receiving SQL statements from the applications; send them to the AST agent; receive from the AST agent the rewritten (or not) statements; send these SQL commands to the DBMS; receive the SQL results and send them to the applications.
- Driver for Data Access (DDA): driver that allows the engine of the middleware (Middleware for SQL Tuning) to send the SQL statements to the target DBMS.

Additionally, in order to reduce the overhead in the automatic approach, a query rewriting cache (QRC) can be used. This component is a cache file containing pairs $\langle \text{SQL statement } Q_a, \text{rewritten SQL } Q_b \rangle$, where Q_b was produced by ARE-SQL using Q_a as input. It is important to note that QRC is a rewriting cache. Then, it does not store PEPs, just SQL statements. The DBA can enable or not using QRC. If QRC is disabled or there isn’t an entry in this file for a given statement Q_a , ARE-SQL will need to apply the heuristics to produce Q_b .

3.2 Implementing ARE-SQL

Given a particular SQL statement Q , the rewriting process, by means of heuristics presented in Table I, needs to sweep the text of Q searching for tuning opportunities (eg, the existence of an unnecessary *DISTINCT* operator) and then changes the syntax of Q , generating a new command Q' . For this, we use a SQL parser called ZQL [Thuraisingham et al. 2010]. It is important to note that other SQL parsers have also been investigated and analyzed, such as *General SQL Parser* [Merlo et al. 2007] and *SQL Query Parser* [Elliott et al. 2009]. Nonetheless, ZQL has been selected to be an easy tool to use, free (open-source) and additionally it allows the inclusion of new rules in its SQL grammar.

The ZQL parser was developed using JavaCC (Java Compiler Compiler) [Kodaganallur 2004] with a SQL grammar. In order to recognize SQL statements, ZQL stores elements which may compose SQL statements in Java data structures. Thus, for a given SQL statement, ZQL generates objects that store the column list in SELECT clause, the table list in FROM clause, arithmetic or logical expressions, among others. These objects were used to implement the heuristic set presented in Table I.

```

SELECT c_custkey, c_name, c_address, c_phone
FROM (SELECT * FROM customer), nation, region
WHERE c_nationkey = n_nationkey
AND n_regionkey = r_regionkey
GROUP BY c_custkey, c_name, c_address, c_phone

```

Fig. 7. SQL expression for Q_c

```

Function heuristicH5( $Q_j$ )
   $oq_j$  = a ZQuery instance representing  $Q_j$ 
  Vector operators =  $oq_j.getWhere().getOperators()$ 
  For each operator  $\in$  operators do
    If ((operator == ">=") or (operator == ">") or (operator == "<=") or (operator == "<"))
      leftText = text on the left side of the operator
      rightText = text on the right side of the operator
      If rightText startsWith "ALL"
        sq = the sub-query on rightText
        remove "ALL" from sq
        If ((operator == ">=") or (operator == ">"))
          add "MAX" function in sq
        End If
        If ((operator == "<=") or (operator == "<"))
          add "MIN" function in sq
        End If
      End If
    End If
  End For
  Return  $Q'_j$  //The SQL Statement produced by applying H5
End Function

```

Fig. 8. Algorithm for Heuristic H5

The ZQL receives as input a string containing a given SQL statement Q (select, insert, update, delete, commit, rollback or set transaction) and creates a Java data structure to represent Q . For this, ZQL uses the Java Class *ZStatement* or one of its subclasses *ZQuery* or *ZUpdate*.

To illustrate the process of parsing a SQL statement, consider the SQL query Q_c depicted in Figure 7. Now, assume that Q_c is sent to the ZQL. Thus, to represent Q_c , ZQL creates an instance of the *ZQuery* class, denoted by o_c . The method *getSelect()* of o_c builds a Vector of objects V containing one object (instance of the class *ZColumn*) for each column in the *SELECT* clause of Q_c : $c_custkey$, c_name , $c_address$ and c_phone . Similarly, the method *getFrom()* returns a Vector of objects containing one object for each table in the *FROM* clause of Q_c . It is noteworthy that, in this case, the vector will include a *ZQuery* instance to represent the sub-query *SELECT * FROM customer* and two objects (instances of the class *ZTable*) to depict the tables *nation* and *region*. The method *getWhere()* returns an object of the *ZExpression* class. A *ZExpression* object can contain instances of the *ZQuery* class. The method *groupBy()* returns an object of the *ZGroupby* class. A *ZGroupby* object can contain instances of the *ZQuery* class.

For each heuristic in *HS* (Table I) we have implemented an algorithm. Figure 8 describes the process of applying the heuristic H5. This algorithm receives as input a SQL statement Q_j and produces as output a SQL statement Q'_j , a new SQL statement equivalent to Q_j produced by heuristic H5.

The algorithm illustrated in Figure 9 describes the process of rewriting a SQL query Q_j , using the heuristic set presented in Table I. This algorithm receives as input a SQL statement Q_j and produces as output a SQL statement Q'_j , a new SQL statement equivalent to Q_j .

```

Function runHeuristics( $Q_j$ )
  For each heuristic  $h_i \in HS$  do
    run  $h_i(Q_j)$ 
  End For
  For each sub-query  $Q_s$  in the WHERE clause of  $Q_j$  do
    runHeuristics( $Q_s$ )
  End For
  For each sub-query  $Q_s$  in the FROM clause of  $Q_j$  do
    runHeuristics( $Q_s$ )
  End For
  For each sub-query  $Q_s$  in the SELECT clause of  $Q_j$  do
    runHeuristics( $Q_s$ )
  End For
  For each sub-query  $Q_s$  in the HAVING clause of  $Q_j$  do
    runHeuristics( $Q_s$ )
  End For
  Return  $Q'_j$  //The SQL Statement produced by the rewriting process
End Function

```

Fig. 9. Recursive algorithm for SQL rewriting

Note that the DBA has the ability to interact with the tuning suggestions provided by the *advisor*, in the assisted strategy. For example, the DBA can decide to apply only a subset of heuristics belonging heuristic set (HS), the DBA can force the application of heuristics that have not been effectively used by the *advisor* or the DBA can choose the order in which the heuristics are performed, among other options. Besides, the DBA can choose and set a subset of the available heuristics so that they are always executed in the rewriting process for queries in general or for SQL statement in particular. In this case, this information is stored in the Tuning Settings (TS) repository. For that reason, before rewriting a SQL statement Q , the AST agent checks, at the Tuning Settings (TS) component, if there is a subset of heuristics selected (defined) for rewriting Q . If not, all available heuristics are used to rewrite Q .

4. EXPERIMENTAL RESULTS

In order to show the potentials of ARE-SQL, several experiments have been conducted. The main results achieved so far are presented and discussed in this section. Thus, we first provide information on how the experimentation environment was set up. Thereafter, empirical results are quantitatively presented and qualitatively discussed.

4.1 Simulation Setup

In order to validate both approaches implemented by ARE-SQL, assisted and automatic, we have evaluated ARE-SQL in three different scenarios. In the first scenario, TPC-H benchmark has been used, including its database and workload, which is composed of 23 queries. In the second scenario, a workload of 30 synthetic queries has been executed on TPC-H database. In both scenarios, we have used TPC-H benchmark with scale factor of 2 GB with the indexes suggested by OSDL DBT-3 [OSDL 2014]. Finally, in the third scenario, we exploit the database of the Integrated Management System (IMS), a system used in several Brazilian universities, and another synthetic workload.

The synthetic workloads, used in the second and third scenarios, were designed using the following strategy. Each synthetic workload was formed by 30 SQL queries. Each SQL query was written in order to enable the application of one of the heuristics presented in the Table I. It is important to note that the heuristics provide different gains. Then, a particular heuristic can provide a much greater benefit than another. However, in all queries that formed the synthetic workloads (used in scenarios

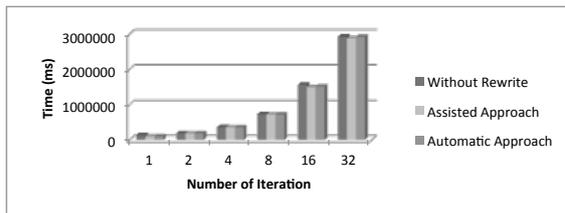


Fig. 10. Benchmark TPC-H on PostgreSQL

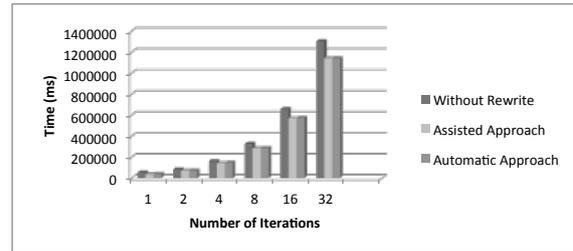


Fig. 11. Benchmark TPC-H on SQL Server

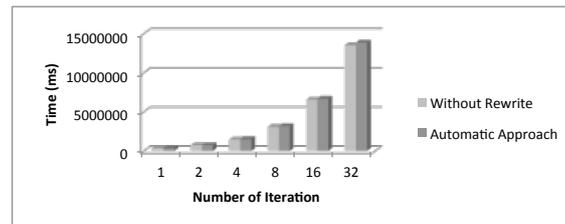


Fig. 12. Benchmark TPC-H on Oracle

2 and 3), the used heuristics provided some gain (a positive profit). That is, any query written by ARe-SQL had a response time smaller than the original query execution time in synthetic workloads.

All experiments were run on a Core i3-2100 (3.10GHz) server, with 4GB RAM and 500 GB HD. PostgreSQL 8.1, Oracle 11g and SQL Server 2008 have been used as database systems. For each experimentation scenario, three different experiments have been performed: i) The workload containing the original queries has been submitted to a database system. The results of this experiment have been used as baseline; ii) The workload submitted to ARe-SQL advisor (assisted approach). Thereafter, the workload with tuned SQL queries was manually submitted to a database system, and; iii) Each query from original workload has been sent to the automatic approach of ARe-SQL. In this case, after tuning SQL queries, ARe-SQL automatically submits those queries to a database system. Then, in the experimental results, the execution time for the assisted strategy includes the query rewriting overhead. For each test, we have executed the set of queries belonging to the used workload once, twice, 4, 8, 16, and 32 times. For each different number of executions (iterations), the sum of execution time of the whole set of queries was computed. For instance, in scenario 1 (Full TPC-H Benchmark), to get the execution time for 32 executions, we have execute the set of TPC-H queries 32 times, measured the time interval for executing all queries and computed the sum of the 32 measured execution times.

4.2 Scenario 1: Full TPC-H Benchmark

Figures 10 and 11 bring the result of using the assisted and automatic approach of ARe-SQL running on PostgreSQL and SQL Server, respectively. It is important to observe that from the 23 queries that comprise the TPC-H benchmark, two (18 and 20) have been rewritten by ARe-SQL advisor (see Section 3.1.1). Regarding the results presented in Figures 10 and 11, the assisted approach had a small decrease in workload runtime. This can be explained by the fact that only two TPC-H queries have presented opportunities for tuning. Table II shows the execution time of TPC-H queries 18 and 20 in their original formats and after being rewritten. Is important to emphasize that this is not a bad result, since the TPC-H queries were written by experts in SQL.

On Oracle, however, TPC-H queries 18 and 20 have not been rewritten by ARe-SQL advisor, because Oracle implements heuristic H8. For that reason, the ARe-SQL automatic approach had, for 32 iterations, a slight performance decrease w.r.t. the baseline (2.45%), see Figure 12 due to the fact that the automatic approach had the overhead of trying to rewrite both queries.

Table II. Execution times of the TPC-H queries 18 and 20 in their original formats and after being rewritten

	PostgreSQL		SQL Server	
	Original	Rewritten	Original	Rewritten
Query TPC-H 18	134807ms	111933ms	19000ms	14000ms
Query TPC-H 20	912ms	712ms	15000ms	11000ms

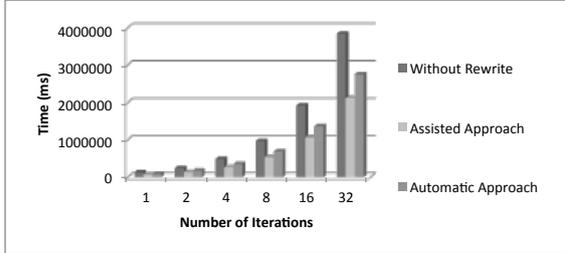


Fig. 13. Synthetic queries in PostgreSQL TPC-H database

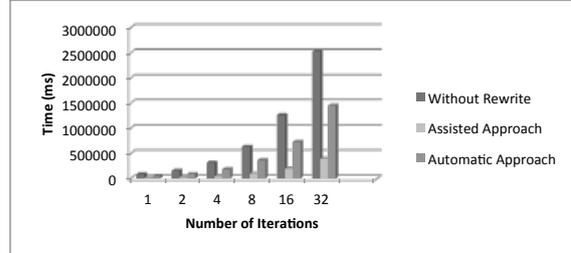


Fig. 14. Synthetic queries in SQL Server TPC-H database

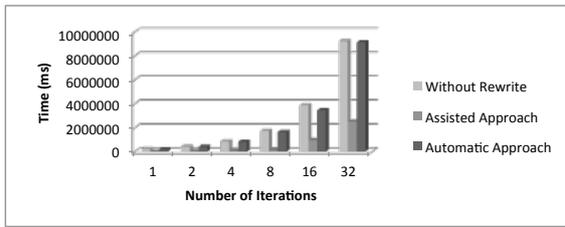


Fig. 15. Synthetic queries in TPC-H database on Oracle

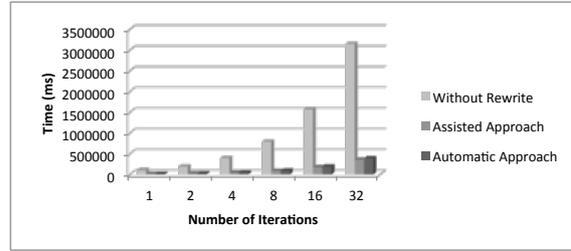


Fig. 16. Synthetic queries in IMS database on PostgreSQL

4.3 Scenario 2: TPC-H Database with Synthetic Workload

Figures 13, 14 and 15 show that using ARe-SQL advisor ensures a significant reduction in query response time. On the other hand, ARe-SQL’s automatic approach presents smaller benefits than ARe-SQL advisor, since it involves the overhead of tuning the SQL statements received in runtime.

It is important to note that in the test on Oracle (Figure 15), the automatic approach presents a small decrease in the execution time of the workload. This is explained by the fact that from the total of eleven heuristics five are already implemented by Oracle (H4, H5, H6, H7 and H8). Besides, from the six remaining heuristics, three of them (H9, H10 and H11) make use of statistical information, and, therefore, require additional access to DBMS, which significantly increased the overhead to rewrite SQL queries.

In the experiment with TPC-H Database and Synthetic Workload on Oracle, the runtime for the 32 iterations of the original workload (without rewriting) took 9226863ms, while the runtime for the 32 iterations of the workload generated after applying the assisted approach lasted 2517264ms (Figure 16). In this case, the assisted approach yielded a gain of 6709599ms (72.7%).

4.4 Scenario 3: IMS Database with Synthetic Workload

For this scenario, only PostgreSQL has been used, as the IMS database is only available for that DBMS. Again, the proposed approaches have provided a high reduction in time execution of the submitted workload (Figure 16). In the experiment with IMS Database and Synthetic Workload on PostgreSQL, the runtime for the 32 iterations of the original workload (without rewriting) took 3145344ms, while the runtime for the 32 iterations of the workload generated after applying the assisted approach lasted

Table III. SQL Optimizer for Oracle using TPC-H benchmark

Query	# generated SQL commands	Rewriting Overhead	# SQL clauses with cost estimate less than the original SQL	# SQL clauses with response time estimate less than the original SQL	Original SQL response time	Execution time for the best rewritten SQL
2	115	116000ms	3	3	310ms	128ms
3	53	58000ms	3	1	6393ms	4900ms
4	14	15000ms	1	0	4204ms	60000ms
5	228	229000ms	5	0	4658ms	4666ms
7	110	111000ms	1	0	4658ms	7808ms
11	139	140000ms	8	4	600ms	259ms
12	20	21000ms	1	1	1995ms	1993ms
14	5	6000ms	1	0	512ms	538ms
16	35	36000ms	2	0	278244ms	280190 ms
18	160	161000ms	1	0	17384ms	22306 ms
20	92	93000ms	8	3	3261ms	3212ms
21	120	121000ms	5	0	8660ms	8787ms
22	1	1000ms	1	0	407ms	418ms

Table IV. Overhead Analysis using TPC-H benchmark

Query	2	3	4	5	7	11	12	14	16	18	20	21	22
SQL Optimizer	116s	58s	15s	229s	111s	140s	21s	6s	36s	161s	93s	121s	1s
ARe-SQL	.01s	.31s	.12s	.08s	.05s	.13s	.16s	.002s	.03s	.11s	.24s	.28s	.08s

352560ms. In this case, assisted approach yielded a gain of 2792784ms (88.79%).

4.5 Overhead Analysis

To analyze the overhead in SQL rewriting process we made an experiment using TPC-H benchmark. The idea was to measure the time needed to rewrite each one of the 23 SQL queries in the TPC-H benchmark workload. For this, we have used as a baseline the rewritten time produced by *Quest SQL Optimizer for Oracle*. This tool was chosen because it is widely used in real scenarios.

Table III shows a piece of the report generated by *Quest SQL Optimizer for Oracle* using TPC-H queries as input. This report indicated that from the total of 23 SQL queries 16 could be rewritten successfully. Only 7 SQL queries couldn't be tuned (1, 6, 10, 13, 17, 19 and 23). However, after performing (with empty cache), for each one of these 16 SQL queries, the best SQL command (w.r.t cost estimate) generated by the *Quest SQL Optimizer for Oracle* tool, we noted that only 5 SQL queries had reduced their runtimes (queries 2, 3, 11, 12 and 20). In practice, queries 4, 5, 7, 8, 9, 14, 16, 18, 21 and 22 had a worsening in their execution times. It is important to note that in some cases, like queries 4 and 7, the rewritten SQL response time was much greater than the original query runtime. Additionally, *Quest SQL Optimizer for Oracle* tool is DBMS-specific (running only with Oracle). Besides, in some cases, this tool may generate tens or hundreds of rewritten SQL instructions for tune the same SQL statement, such as, for example, for query 5, 228 different rewritten SQL instructions were generated. This fact could increase so much the rewritten time, producing a high overhead, and making this tool unfeasible for online scenarios, that is, those where the SQL statement needs to be adjusted automatically (at runtime).

Table IV analyzes and compares the rewriting overhead in both tools: *Quest SQL Optimizer for Oracle* and ARe-SQL. It is important to note that ARe-SQL showed a lower rewriting overhead in all analyzed queries. For example, the *Quest SQL Optimizer for Oracle* spent 161s to rewrite Query 18 while ARe-SQL took only 0.11s in this task.

5. CONCLUSION

In this article, we presented an online, automatic and non-intrusive approach for rewriting SQL queries, denoted ARe-SQL. The key goal of ARe-SQL is to assist optimizers in building efficient query physical execution plans. For this, ARe-SQL rewrites SQL statements, using a set of 11 heuristics. Based on ARe-SQL, two different strategies for SQL-query rewriting, denoted assisted and automatic, were implemented. These strategies were evaluated in three different DBMS considering three different scenarios. The experimental results indicate that both strategies can provide performance gains. In some cases, such performance gains reach 92.5%. The proposed solutions are applicable to situations where the query optimizer cannot produce optimal plans, even using access methods and assessment strategies supported by the DBMS. ARe-SQL is even more important in contexts in which SQL statements are written by non-experts. As future work we intend to investigate and implement heuristics based on query hints to tune SQL queries.

REFERENCES

- ARAÚJO, A. H. M., MONTEIRO, J. M., TAVARES, J. A., BRAYNER, A., AND LIFSCHITZ, S. On Using an Automatic, Autonomous and Non-Intrusive Approach for Rewriting SQL Queries. In *Proceedings of the Brazilian Symposium on Databases (Short Paper)*. Recife, PE, Brazil, pp. 115–120, 2013.
- BELKNAP, P., DAGEVILLE, B., DIAS, K., AND YAGOUB, K. Self-Tuning for SQL Performance in Oracle Database 11G. In *Proceedings of the IEEE International Conference on Data Engineering*. Washington, DC, USA, pp. 1694–1700, 2009.
- BRUNO, N., CHAUDHURI, S., AND RAMAMURTHY, R. Power Hints for Query Optimization. In *Proceedings of the IEEE International Conference on Data Engineering*. IEEE Computer Society, Washington, DC, USA, pp. 469–480, 2009.
- CHAUDHURI, S. AND WEIKUM, G. Foundations of Automated Database Tuning. In *Proceedings of the International Conference on Very Large Data Bases*. pp. 1265–1265, 2006.
- DAGEVILLE, B. AND DIAS, K. Oracle’s Self-Tuning Architecture and Solutions. *IEEE Data Engineering Bulletin* 29 (3): 24–31, 2006.
- ELLIOTT, B., CHENG, E., THOMAS-OGBUJI, C., AND OZSOYOGLU, Z. M. A Complete Translation from SPARQL into Efficient SQL. In *Proceedings of the International Database Engineering & Applications Symposium*. Cetraro, Calabria, Italy, pp. 31–42, 2009.
- EMBARCADERO. Embarcadero DB Optimizer XE, 2010. Available at <http://www.embarcadero.com/products/db-optimizer-xe>. Access in June 2011.
- HERODOTOU, H. AND BABU, S. Automated SQL Tuning through Trial and (sometimes) Error. In *Proceedings of the Second International Workshop on Testing Database Systems*. Rhode Island, USA, 2009.
- IBM. IBM Optim Development Studio, 2010. Available at <http://www-01.ibm.com/software/data/optim/development-studio>. Access in June 2011.
- KODAGANALLUR, V. Incorporating Language Processing into Java Applications: a JavaCC tutorial. *IEEE Software* 21 (4): 70–77, 2004.
- MARKL, V., LOHMAN, G. M., AND RAMAN, V. LEO: an autonomic query optimizer for DB2. *IBM Syst. J.* 42 (1): 98–106, Jan., 2003.
- MERLO, E., LETARTE, D., AND ANTONIOL, G. Automated Protection of PHP Applications Against SQL-injection Attacks. In *Proceedings of the European Conference on Software Maintenance and Reengineering*. pp. 191–202, 2007.
- OSDL. OSDL, 2014. Available at <http://osldb.sourceforge.net/>.
- QUEST. Quest SQL Optimizer for Oracle, 2010. Available at <http://www.quest.com/SQL-Optimizer-for-Oracle>. Access in June 2011.
- RAMAKRISHNAN, R. AND GEHRKE, J. *Database Management Systems*. McGraw Hill, 2002.
- THURASINGHAM, B., KHADILKAR, V., GUPTA, A., KANTARCIOGLU, M., AND KHAN, L. Secure Data Storage and Retrieval in the Cloud. In *Proceedings of the International Conference on Networking, Applications and Worksharing*. Texas, USA, 2010.
- WEIKUM, G., MOENKEBERG, A., HASSE, C., AND ZABBACK, P. Self-tuning Database Technology and Information Services: from wishful thinking to viable engineering. In *Proceedings of the International Conference on Very Large Data Bases*. pp. 20–31, 2002.