

Query Rewriting via LLMs

Sriram Dharwada¹, Himanshu Devrani¹, Jayant Haritsa¹, and Harish Doraiswamy²

¹Indian Institute of Science

{sriramd,himanshud,haritsa}@iisc.ac.in

²Microsoft Research India

harish.doraiswamy@microsoft.com

Abstract

When complex SQL queries suffer slow executions despite query optimization, DBAs typically invoke automated query rewriting tools to recommend “lean” equivalents that are conducive to faster execution. The rewritings are usually achieved via transformation rules, but these rules are limited in scope and difficult to update in a production system. Recently, LLM-based techniques have also been suggested, but they are prone to semantic and syntactic errors.

We investigate here how the remarkable cognitive capabilities of LLMs can be leveraged for performant query rewriting while incorporating safeguards and optimizations to ensure correctness and efficiency. Our study shows that these goals can be progressively achieved through incorporation of (a) an ensemble suite of basic prompts, (b) database-sensitive prompts via redundancy removal and selectivity-based rewriting rules, and (c) LLM token probability-guided rewrite paths. Further, a suite of logic-based and statistical tools can be used to check for semantic violations in the rewrites prior to DBA consideration.

We have implemented the above LLM-infused techniques in the `LITHE` system, and evaluated complex analytic queries from standard benchmarks on contemporary database platforms. The results show significant performance improvements for slow queries, over both SOTA rewriters and the native optimizer. For instance, with TPC-DS on PostgreSQL, the GM of runtime speedups was a high **13.2** over the native optimizer, whereas SOTA only gave **4.9**.

Overall, `LITHE` is a promising step toward viable LLM-based advisory tools for ameliorating enterprise query performance.

1 Introduction

The SQL queries embedded in enterprise applications are often riddled with inefficiencies and redundancies [51, 2, 6], especially when machine-generated by modeling software such as ORM tools (e.g., Entity Framework [37], Hibernate [27]). Query optimizers are expected to, in principle, automatically remove such wasteful bloat while constructing efficient execution plans. However, in practice, they are often led astray by complex representations, resulting in poor response times.

As a case in point, consider the blog-processing query [6] shown in Figure 1, which computes a daily summary of rating metrics for highly-rated blogs. When this seemingly complex query is given to the current PostgreSQL engine (v16) [3], its execution plan essentially mimics the hierarchical structure of the query. This strategy leads to multiple scans and joins of the `Blogs` and `Posts` tables, making the

```

SELECT t.Key, SUM(t.Rating) AS PostRating,
  (SELECT SUM(b0.Rating) FROM (
    SELECT p0.PostId, p0.BlogId, p0.Content, p0.CreatedDate,
      p0.Rating, p0.Title, b1.BlogId AS BlogId0,
      b1.Rating AS Rating0, b1.Url, p0.day AS Key
    FROM Posts AS p0 INNER JOIN Blogs AS b1 ON p0.BlogId = b1.BlogId
    WHERE b1.Rating > 5) AS t0
  INNER JOIN Blogs AS b0 ON t0.BlogId = b0.BlogId
  WHERE t.Key = t0.Key) AS BlogRating
FROM (SELECT p.Rating, p.day AS Key
  FROM Posts AS p INNER JOIN Blogs AS b ON p.BlogId = b.BlogId
  WHERE b.Rating > 5) AS t
GROUP BY t.Key;

```

Figure 1: Complex SQL Representation

```

SELECT p.day AS Key, SUM(p.Rating) AS PostRating,
  SUM(b.Rating) AS BlogRating
FROM Posts AS p INNER JOIN Blogs AS b ON p.BlogId = b.BlogId
WHERE b.Rating > 5
GROUP BY p.day;

```

Figure 2: Lean Equivalent Query

query take several *minutes* to complete. However, the query can be equivalently rewritten (assuming NOT NULL column constraints and key-joins) in the “lean” flat version shown in Figure 2 – this alternative requires only a single join, and runs more than an *order-of-magnitude* faster, completing within seconds!

For such optimizer-failure scenarios, an alternative option to rectify slow-running queries is to carry out *external tuning* – in particular, DBAs typically resort to SQL-to-SQL (**S2S**) *query rewriting tools* for recommending performant alternatives. A viable S2S query transformer should cover a wide range of queries and ensure the rewrite is: (1) semantically equivalent to the original; (2) ideally performance-beneficial, but at least not causing regression; and (3) incurring overheads that are practical for deployment.

Prior Work

A range of innovative *Rule-based* (e.g. Learned Rewrite [57]) and *Model-based* (e.g. Gen-Rewrite [34]) approaches have been proposed for S2S transformations. These state-of-the-art (SOTA) techniques have foregrounded the potential benefits of query rewriting. However, as explained later in Section 9, their realization of these benefits is curtailed by: (a) restrictions in rewrite scope, (b) susceptibility to semantic and syntactic errors, and (c) transformations via the plan space (i.e optimization on the nodes of the execution plan tree) rather than directly in query space (i.e. transforming the query structure itself). In this paper, we address these lacuna and substantively amplify the effectiveness of query rewriting.

The LITHE Rewriter

We propose LITHE (LLM Infused Transformations of HEfty queries), an LLM-based query rewriting assistant to aid DBAs in tuning slow-running queries. As illustrated in the architectural diagram of Figure 3, LITHE takes the user query Q^U as input and outputs a transformed query Q^T , together with (a) the *expected* performance improvement of Q^T , in terms of optimizer estimated cost; (b) a *verification label* indicating the semantic equivalence mechanism (provable or statistical) between Q^T and Q^U ; and (c) a *reasoning* for why the LLM expects Q^T to be helpful wrt performance. Armed

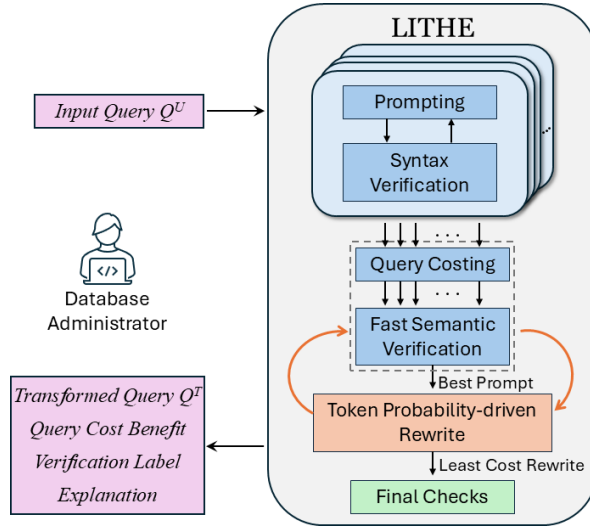


Figure 3: High-level architecture of LITHE

with this information, DBAs can leverage their expertise to decide whether or not to use Q^T . Note that having the DBA in the loop is a common practice in commercial query advisory systems [17].

To help the LLM adapt to different query patterns and structures, we introduce database-aware rules into the prompts, which make use of the rich metadata (e.g., logical schema, selectivities) available in database environments (Section 4). Our rules are invoked directly in *query space*, providing the LLM with the latitude to generalize the rule usage to a wide range of queries. As shown later in our experiments, this is a powerful mechanism for ensuring performant rewrites across database environments. Furthermore, we leverage the rich telemetry provided by LLMs – particularly, the *token probabilities* output at each step in the prediction sequence. Whenever the LLM lacks high confidence in the next token, we follow multiple alternative paths in the decision process. To ensure practical overheads on this enumerative approach, a Monte Carlo tree search (MCTS) technique is incorporated in our implementation (Section 5).

LITHE also incorporates a suite of logic-based and result equivalence-based semantic equivalence tests to verify the correctness of the rewrites (Section 6.2). Finally, given that there are often material discrepancies between optimizer cost predictions and actual execution times [35], LITHE includes heuristic mechanisms, described in Section 6.3, to identify “brittle” rewrites and discard them in favor of the original.

Results. Our experiments on benchmark environments demonstrate that LITHE achieves, for many slow queries, semantically correct transformations that significantly reduce the abstract (optimizer-estimated) costs. In particular, for TPC-DS, LITHE constructed “highly productive” ($> 1.5x$ estimated speedup) rewrites for as many as 26 queries, whereas SOTA promised such rewrites for only about half the number. Further, the GM (Geometric Mean) of LITHE’s cost reductions reached **11.5**, almost double the **6.1** offered by SOTA.

We also evaluated whether the projected cost reductions translated into real execution speedups. Here, we find that LITHE is indeed often substantively faster at run-time as well. Specifically, the geometric mean of the runtime speedups for slow queries was as high as **13.2** over the native optimizer, whereas SOTA delivered **4.9** in comparison. Finally, our regression identification techniques were successful in guarding against “brittle” rewrites that are promising cost-wise, but eventually poor at run-time.

Contributions

In summary, our study makes the following contributions:

1. Assesses LLM suitability for S2S transformation.
2. Transforms directly in query space instead of plan space intermediates, leading to performant rewrites.
3. Incorporates database-sensitive rules in LLM prompts, covering both schematic and statistical dimensions.
4. Leverages LLM token probabilities to guide navigation of the rewrite search space and minimize LLM errors.
5. Evaluates rewrite quality over a range of database environments, demonstrating substantial performance benefits.
6. Identifies learnings that could help guide research directions for industrial-strength query rewriting.

2 LITHE Overview

The LITHE architecture (Figure 3) comprises a five-stage pipeline:

1. Prompt-based rewriting. This component consists of two modules: LLM Prompting and Syntax Verification. The user query and an initial prompt are fed to the LLM prompting module, requesting a rewrite. The LLM output’s syntax is then verified with the database engine’s parser – if invalid, the error is returned to the LLM via an updated prompt asking for a correction. This goes on iteratively (within a threshold) until a valid SQL query is obtained.

This process is repeated for each of the different prompts described in Sections 3 and 4, and the resulting syntactically-valid rewrites are fed to the query costing component. We evaluate each prompt within a fresh LLM context, thereby making the prompt processing to be order-independent.

2. Query Costing. The costs of candidate rewrites are evaluated via the database engine’s optimizer. Rewrites whose costs are greater than that of the original query are immediately discarded. Whereas, the potentially beneficial rewrites (if any) are checked for semantic equivalence to the original query.

3. Fast Semantic Equivalence. Statistical (result-based) techniques are employed to quickly and cheaply assess the semantic equivalence of a recommended rewrite. If the rewrite is deemed valid by this module, it is returned along with the prompt that generated it; otherwise an invalid label is returned.

4. Token probability-driven rewrite. The prompt producing the most beneficial (and valid) rewrite is used as input to a Monte-Carlo tree search (MCTS)-based procedure to further refine the rewrite quality. Specifically, multiple exploration paths are followed whenever the LLM lacks high confidence in the predicted token (details in Section 5). The query costing and semantic equivalence modules are also used internally within this procedure.

5. Final Checks and Output. The least expensive valid rewrite (as identified by the MCTS module) is tested using a suite of provable (logic-based) techniques to generate a final equivalence label (provable or statistical), and its cost benefit over the original query is computed. Further, the execution “brittleness” is assessed using the robustness heuristics. If no valid rewrite is identified, or if the rewrite is expected to be a regression, the original query itself is returned to the DBA. Whereas, if a beneficial rewrite is recommended, an LLM-generated reasoning for the expected performance improvement is also extracted.

2.1 Performance Framework

We consider a query that takes more than T seconds to complete on the native database engine as a “slow query”, potentially triggering intervention by the DBA. Based on common industry perception (e.g. [7]), T is set at 10 seconds in our study. For this context, we define a **Cost Productive Rewrite (CPR)** as a rewrite that improves a slow query’s performance by at least **1.5** times wrt the optimizer-estimated cost – this aggressive choice of threshold is so that: (a) there is enough headroom in the optimizer prediction that a runtime regression is unlikely; and (b) the benefits of the rewrite substantively outweigh the transformation overheads.

The overall benefit of a rewriting tool is quantified by the number of CPR obtained on the slow query workload. Additionally, we measure **CSGM**, the *geometric mean* (GM) of the cost speedups obtained by these rewrites. Finally, to assess the actual run-time benefits, we evaluate **TSGM**, the geometric mean (GM) of the response-time speedups obtained by these rewrites. These speedups are measured relative to the original query on the native optimizer.

2.2 Query Micro-benchmark

To motivate the progression of strategies incorporated in **LITHE**, we create an initial micro-benchmark comprising 10 diverse TPC-DS queries for which we were able to *hand-craft* high-quality CPRs. These queries are processed on GPT-4o, the popular OpenAI LLM, and the rewrites are evaluated on the PostgreSQL v16 database engine. The human rewrites deliver a CSGM of **11.84**, serving as an aspirational target to attain computationally. Later, in Section 7, we extend the evaluation to complete benchmarks.

Further, for ease of presentation, we focus on the CSGM metric in the **LITHE** design sections (Sections 3, 4, 5). The TSGM performance is subsequently discussed in Section 7.

3 Basic Prompts

In this section, we explore the simplest interface to LLMs, namely *prompting*, for query rewriting. We evaluate four basic prompts, enumerated in Figure 4, which cover a progressive range of instructional detail and test the effectiveness of the LLM’s base knowledge.

Prompt 1: This is the baseline prompt used in [34], which simply asks the LLM to rewrite a given query to improve performance.

Prompt 2: Explicit instructions are included to maintain semantic and functional equivalence while rewriting.

Prompt 3: Verbose instructions are given to rewrite the query, providing step-by-step guidance to the LLM to think rationally. It is first asked to pick out potential inefficiencies in the input query, and then tasked to identify approaches to address these inefficiencies. Finally, it is instructed to apply the identified solution. Essentially, the prompt tries to make the LLM reason akin to human experts.

Prompt 4: The sequence of instructions in Prompt 3 is split into sub-prompts, and provided to the LLM in an *iterative* manner instead of all at once. This breaks down the complex instructions of Prompt 3 into digestible steps that help the LLM focus on individual tasks.

Performance

Table 1 shows the performance of the four prompt templates on the micro-benchmark. We find that less than half the rewrites are productive with individual prompts. However, a drill-down shows that the best prompt in the ensemble is *query-specific* – this opens up the possibility of using all four prompts in

| | |
|---------------------|------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| Prompt 1 | <p>**Enter Query here**</p> <p>Rewrite this query to improve performance.</p> |
| Prompt 2 | <p>You are a database expert and SQL optimizer. Your role involves identifying inefficient SQL queries and transforming them into optimized, functionally equivalent <Database Engine> versions.</p> <p>**Enter Query here**</p> <p>Rewrite this query to improve performance of query while maintaining semantic equivalence.</p> |
| Prompt 3 / Prompt 4 | <p>You are a database expert and SQL optimizer. Your role involves identifying inefficient SQL queries and transforming them into optimized, functionally equivalent <Database Engine> versions. Your tone is analytical and instructional.</p> <p>A user has provided the following <Database Engine> query that is potentially inefficient:</p> <p>**Enter Query here**</p> <p>The task is to first identify whether the query is inefficient or not. If it is inefficient, you must rewrite the query to make it more efficient while maintaining semantic equivalence. Here are a few steps that can help you complete the task:</p> <ol style="list-style-type: none"> 1. Start by identifying specific inefficiencies in the provided query. If you feel the original query is already efficient, skip the next two steps, and simply return the query as-is. 2. Next, provide guidelines for optimizations, and explain the rationale behind the recommended optimizations. Correspond how these changes would map onto the original query to maintain syntactic and semantic equivalence. 3. Finally craft the new, optimized <Database Engine> query which includes all the enhancements discussed. Give complete rewritten query with no manual involvement by user. |

Figure 4: Templates used for Basic Prompts

parallel, and then choosing the best among them. This ensemble approach raises the CPRs to 6 (Row 5 in Table 1) – however, there remain four queries that are not productively rewritten by these prompts.

Table 1: Performance of Basic Prompts.

| Prompt | # CPR | CSGM |
|------------------------|-------|------|
| Prompt 1 | 3 | 1.99 |
| Prompt 2 | 2 | 1.85 |
| Prompt 3 | 4 | 2.83 |
| Prompt 4 | 4 | 3.00 |
| Prompt Ensemble | 6 | 3.23 |
| SOTA Ensemble | 3 | 2.49 |

The CSGM, shown in the last column of Table 1, is at most 3 for the individual prompts, while the ensemble reaches 3.23. But these speedups, although productive, are all lower than those delivered by the human rewrites.

Finally, an ensemble of SOTA techniques (described in Section 7) was also processed on the same platform. They delivered 3 CPRs with a CSGM of 2.49 (last row in Table 1), indicating the wide gap between the current reality and what is humanly possible.

4 Database-Sensitive Prompts

As discussed above, basic prompting needs to be improved on two fronts: (1) Ensuring productive rewrites where feasible; and (2) Maximizing the impact of these productive rewrites. Our discussions with industry experts revealed that query inefficiencies are most commonly attributed to redundant operations within the queries themselves. In fact, contemporary (non-LLM) optimizers often try to reduce redundancies during the initial logical optimization phase, prior to physical plan enumeration (e.g., redundant join elimination [9]). Furthermore, our analysis of execution plans for slow queries revealed that the optimizer often fails to substitute filter-related operators, even when beneficial to do so.

To address these issues, we incorporate database domain knowledge into the prompts. Specifically, we design a *one shot*-based prompting template, augmented with a set of database-aware rewrite rules. As a proof of concept, we explore two categories of rewrites here: (a) Rules that eliminate redundancy in the input queries; and (b) Predicate selectivity-based rules that implicitly guide, via query space reformulations, the query optimizer towards efficient query execution plans. Of course, this basic set of rules can be expanded further, but as shown by our experiments, even this minimal set is capable of delivering substantive improvements over a broad set of database environments.

The templates used for these prompts (Figure 5) apply only a *single* rule. This was a conscious design choice because LLMs can be overwhelmed by excessive information given in monolithic form [46]. A related question is whether, while retaining the one-rule-per-prompt design, the rules could be *progressively* applied, thereby benefiting queries with multiple types of redundancies. However, the combinatorially large number of rule permutation sequences makes the selection process appear infeasibly expensive. Therefore, we apply each rule using a separate prompt, finally returning the rewrite providing the best improvement. Despite this restrictive policy, highly performant rewrites are obtained in our evaluation (Section 7).

4.1 Redundancy Removal

There are different types of redundancy that can occur in a SQL query – repeated computations, superfluous filter predicates, unnecessary joins, etc. Rules **R1 through R4** in Table 2 are designed to tackle such redundancies. The relevant schematic information (e.g. table names, column names, constraints) required by these rules is also provided in the prompt.

Table 2: Rules for Database-sensitive prompts.

| | Redundancy Removal Rules |
|----|-------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| R1 | Use CTEs (Common Table Expressions) to avoid repeated computation of a given expression. |
| R2 | When multiple subqueries use the same base table, rewrite to scan the base table only once. |
| R3 | Remove redundant conjunctive filter predicates. |
| R4 | Remove redundant key (PK-FK) joins. |
| | Statistics-based Rules |
| R5 | Choose EXIST or IN from subquery selectivity (high/low). |
| R6 | Pre-filter tables involved in self-joins and with low selectivities on their filter and/or join predicates. Remove any redundant filters from the main query. Do not create explicit join statements. |

The template for such rule-based prompts is shown in Figure 5(a) and includes an *example rewrite* to demonstrate the rule application to the LLM. Our empirical observation is that the precise wording of

| | | |
|----------------------------|-----|------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| Redundancy-Removal Prompt | (a) | <p>You are a database expert and SQL optimizer. Your role involves identifying inefficient SQL queries and transforming them into optimized, functionally equivalent <Database Engine> versions. Your tone is analytical and instructional.</p> <p>This is a task to rewrite queries to improve performance using the following rewrite rule:</p> <p>**Enter Rewrite Rule here**</p> <p>[ORIGINAL] **Enter Example Original Query here**</p> <p>[REWRITTEN] **Enter Example Rewritten Query here**</p> <p>Now consider the query below and try to rewrite it to improve the performance.</p> <p>[ORIGINAL] **Enter User Input Query here**</p> <p>[REWRITTEN]</p> |
| | (b) | <p>You are a database optimizer and your task is to analyze the given <Database Engine> query using database schema and statistical information available to you. Find inefficiency in the query. If query is already efficient then return the original query. Otherwise, rewrite the query in a more efficient way to improve its performance, and explain your rewrites.</p> <p>This is a task to improve query performance using the following rewrite rule:</p> <p>**Enter Rewrite Rule here**</p> <p>Here is an example to help you</p> <p>[ORIGINAL] **Enter Example Original Query here**</p> <p>[STATISTICS] **Enter Statistics for example Query here**</p> <p>[REWRITTEN] **Enter Example Rewritten Query here**</p> <p>Now consider the query and statistics given below and try to rewrite the query to improve performance.</p> <p>[ORIGINAL] **Enter User Input Query here**</p> <p>[Statistics] **Enter Input Query statistics here**</p> <p>[REWRITTEN]</p> |
| Statistics Guidance Prompt | | |

Figure 5: Templates for Database-sensitive Prompts

the rule instructions is not significant. Instead, what matters are the example rewrites, which ensure the LLM applies the rules effectively, and does not go off on tangential trajectories. The specific examples used with our rules are listed in Section D.

Performance

The performance improvement achieved on the micro-benchmark by an ensemble that adds the redundancy-removing prompts to the basic set (Section 3) is shown in Table 3. We observe that the CPR increases to **7**, and CSGM grows to **6.85**.

Table 3: Performance with Redundancy Removal Rules.

| Prompt | # CPR | CSGM |
|----------------------------------------|-------|------|
| Basic Prompts $\cup \{R1, \dots, R4\}$ | 7 | 6.85 |


```

SELECT c_birth_country, count(*) cnt1
FROM customer c
WHERE (EXISTS (
    SELECT * FROM web_sales, date_dim
    WHERE c.c_customer_sk = ws_bill_customer_sk
    and ws_sold_date_sk = d_date_sk
    and d_year >= 1999 )
or EXISTS (
    SELECT * FROM catalog_sales, date_dim
    WHERE c.c_customer_sk = cs_ship_customer_sk
    and cs_sold_date_sk = d_date_sk
    and d_year >= 1999 )
)
GROUP BY c_birth_country;

```

Use EXIST clause when subquery selectivity is high

```

SELECT c_birth_country, count(*) cnt1
FROM customer c
WHERE ( c_customer_sk IN (
    SELECT ws_bill_customer_sk
    FROM web_sales, date_dim
    WHERE ws_sold_date_sk = d_date_sk
    and d_year = 1999)
OR c_customer_sk IN (
    SELECT cs_ship_customer_sk
    FROM catalog_sales, date_dim
    WHERE cs_sold_date_sk = d_date_sk
    and d_year = 1999 )
)
GROUP BY c_birth_country;

```

Use IN clause when subquery selectivity is low

Figure 6: Example Queries illustrating Rule 5

4.2 Selectivity-based Guidance

We now consider rules whose applicability to a query is conditional on the specific database environment, specifically its *statistical* aspects. For example, consider the alternative rewrites shown in Figure 6 using the EXIST and IN clauses (highlighted in red), respectively. Here the appropriate choice is dictated by the *selectivity* of the inner subquery – EXISTS for high selectivity values and IN for low values. This observation is encoded in Rule **R5** of Table 2. Similarly rule **R6**, which pre-filters tables that are involved in self-joins and with low selectivity filter and/or join predicates. Note that a specific instruction to not create explicit joins is added in rule R6. This is because in the presence of CTEs, the LLM is prone to schematic confusion regarding which attribute belongs to which table, leading it to construct *invalid joins*.

The input prompt for these rules, as shown in Figure 5(b), is modified to include the following:

1. Estimated selectivities of columns appearing in the WHERE and JOIN clauses – these values are obtained via calls to the cardinality estimation modules of the query optimizer.
2. Clause rewrite rules and instructions based on statistics.
3. Examples relevant to the chosen rewrite rules.

Performance

The performance improvements following addition of selectivity-guided prompts are shown in Table 4. We observe that CPRs are now obtained for **9** of the ten micro-benchmark queries. Moreover, the resulting CSGM increases to **10.57**, quite close to the human target of 11.84.

Table 4: Performance of Metadata-infused Prompts.

| Prompt | # CPR | CSGM |
|----------------------------------------|-------|-------|
| Basic Prompts $\cup \{R1, \dots, R6\}$ | 9 | 10.57 |

We note in closing that rules R1 through R6 not only add queries to the productive category, but also deliver greater improvement for those already deemed to be productive via the standard prompts of Section 3. Further, to minimize selection overheads, a classifier could be designed to choose the appropriate rule – this option is examined in Section B.

5 Token Probability Driven Rewrite

A key challenge with LLMs is “hallucinations” – responses that range from being mildly incorrect to completely made up. This is often due to the output tokens having low confidence, which “confuse” the LLM and generate suboptimal outputs [23]. To have a robust approach for such cases, we take inspiration from the code generation literature [55]. Specifically, we propose a *Monte Carlo Tree Search* (MCTS) based decoding approach to search for a sequence of LLM-generated tokens that results in both a valid query rewrite and performance improvements. While MCTS has been previously used for *ordering* rule applications in query rewrites [57], our goal is to guide the LLM in *exploring* these rewrites.

This approach models the problem of query rewriting as a decision tree denoting a *Markov Decision Process* (MDP) [41]. The root node of the tree corresponds to the initial prompt. Figure 7 illustrates one such example tree. An edge from a parent node to a child represents a possible token generated by the LLM and is associated with a value denoting the probability of generating this token given the path taken thus far. This is illustrated by the text (token) and number (probability), on the left and right side respectively, of an edge in Figure 7. Here, each edge can be considered as an *action* of the MDP. A node is considered *terminal* if the incoming edge corresponds to the “;” token, signaling end of the textual query.

The *state* of a node n is represented by the partial rewrite created by following the path from the root to n – it is obtained by concatenating the tokens on this path. The root’s state is an empty rewrite, and each terminal node’s state is a complete rewritten query. For example, the state of Node 4 in the example decision tree shown in Figure 7 is the partial rewrite ```select * from```.

Given the hundreds of thousands of tokens in LLM vocabularies, it may be very expensive (both financially and computationally) to construct the entire tree. It is therefore essential to significantly reduce the token search space while exploring the tree for valid rewrites. This is precisely the purpose of MCTS which applies an *Upper Confidence Bound* (UCB) heuristic [44] to identify the best paths in a tree without computing the entire tree.

5.1 MCTS Search Process

The search procedure is codified in Algorithm 1, and comprises four stages – *Selection*, *Expansion*, *Simulation* and *Back Propagation* – described below, which are repeated for $iter_{max}$ iterations.

1. Selection: The first stage identifies the most suitable node of the decision tree that is yet to be expanded (i.e., the tokens corresponding to this node have not yet been processed by the search procedure). Starting from the root, this process greedily selects successive edges (actions) till an unprocessed non-terminal node is reached (Lines 6–8 in Algorithm 1). Actions are picked using a UCB that balances exploration and exploitation. The goal is to pick those actions that have either (1) a higher potential to

Algorithm 1 Token-augmented Rewrite

```
root      # Start State
k         # Maximum number of child node expansions
 $\theta$       # Probability threshold for node expansion
 $iter_{max}$  # Maximum number of iterations

1: Potential, visits, V  $\leftarrow$  empty Map
2: for  $i \leftarrow 1, 2, \dots, iter_{max}$  do
3:   visits[root]  $\leftarrow$  visits[root] + 1
4:    $n_{cur} \leftarrow root$ 
5:   # Stage 1: Selection
6:   while  $len(n_{cur}.children) > 0$  do
7:      $n_{cur} \leftarrow \arg \max_{a \in Actions(n_{cur}.children)} UCB(n_{cur}, a)$ 
8:     visits[ $n_{cur}$ ]  $\leftarrow$  visits[ $n_{cur}$ ] + 1
9:   # Stage 2: Expansion
10:  expand  $\leftarrow$  True
11:  while expand and ‘,’  $\notin n_{cur}.state$  do
12:    tokensnext, Pnext  $\leftarrow$  Model( $n_{cur}, k$ )
13:    if  $P_{next}[0] \leq \theta$  then
14:      for  $token \in tokens_{next}$  do
15:         $n_{new} \leftarrow$  new Node with State  $n_{cur}.state \cdot token$ 
16:        Append  $n_{new}$  to  $n_{cur}.children$ 
17:      expand  $\leftarrow$  False
18:    else
19:       $n_{cur}.state \leftarrow n_{cur}.state \cdot tokens_{next}[0]$ 
20:  # Stage 3: Simulation - Expand from  $n_{cur}$  to full rewrite
21:  query  $\leftarrow$  GreedyExpand( $n_{cur}$ )
22:   $v \leftarrow$  ComputePotential(query)
23:  Potential[query]  $\leftarrow$   $v$ 
24:  # Stage 4: Backpropagation
25:  while  $n_{cur} \neq \text{Null}$  do
26:     $V[n_{cur}] \leftarrow \max(V[n_{cur}], v)$ 
27:     $n_{cur} \leftarrow \text{Parent}(n_{cur})$ 
28: # Return valid rewrite with maximum Potential > 1
29: if  $\exists q \in Potential \mid Potential[q] > 1$  then
30:   return  $q$  having maximum value of Potential[ $q$ ]
31: else
32:   return the original query
```

produce correct and faster rewrites (exploitation); or (2) been selected fewer times in the past (exploration).

Specifically, given a node n and a set of possible actions $a \in A$, the next node in this traversal is chosen as:

$$n_{next} = \arg \max_{a \in A} UCB(n, a) \quad (1)$$

where UCB is a heuristic sourced from [44], modified to reflect our formulation where values are

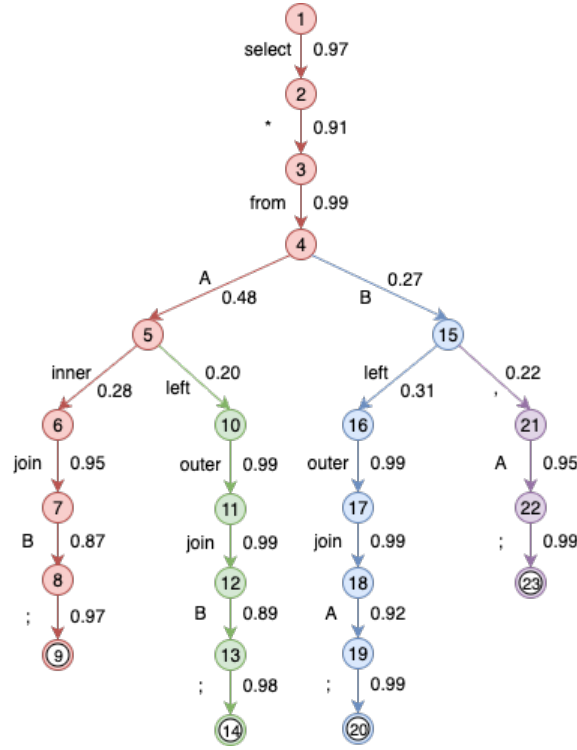


Figure 7: Sample decision tree traversal (Algorithm 1)

associated with nodes rather than edges in the tree. It is defined as follows:

$$UCB(n, a) = V(n') + \beta(n) * P_{LLM}(a|n.state) * \frac{\sqrt{\log(visits[n])}}{1 + visits[n']} \quad (2)$$

Here, n' is the node reached from n by taking action a , and the first component $V(n')$ represents the exploitation potential of n' to produce correct and faster queries (this notion is formalized below in Stage 3). The second component in the equation represents exploration – it is higher for those child nodes of n that are visited less often. Here, P_{LLM} represents the next token probability, which can be queried from the LLM by enabling the `logprobs` option in the API, and $visits[n]$ is the number of times n has been visited during the search process. β is a function that controls the balance between exploration and exploitation. It depends on two hyperparameters c_{base} and c – a higher value of c_{base} makes the algorithm favor exploitation, whereas a higher value of c increases the incentive to explore. β is defined as:

$$\beta(n) = \log\left(\frac{visits[n] + c_{base} + 1}{c_{base}}\right) + c \quad (3)$$

For example, consider the state of the tree in Figure 7 when Nodes 5 and 15 are the current unexpanded nodes. At this juncture, the selection procedure will use the UCB values of these two nodes to choose which node to expand next.

2. Expansion: The second stage is used to expand the unprocessed node n_{cur} chosen by the Selection stage. Specifically, it retrieves from the LLM the top k probable next tokens from n_{cur} 's state (Line 12), and expands the decision tree by adding k new child nodes corresponding to these tokens. To make the expansion tractable, multiple child nodes are added only if the probability of the highest token falls below a threshold θ (Line 13). Otherwise, the tokens are generated in a greedy fashion from the current node until a point where the LLM is again unsure of the next token, or it reaches a terminal node (i.e. completes a query rewrite).

In our implementation, k and θ are set to 2 and 0.7, respectively. Since branching occurs *only* when this relatively high θ threshold is breached, most nodes have only *one* child, resulting in long linear sequences. This design, combined with at most two nodes being processed during a branch, significantly reduces the effective fan-out of the tree. Moreover, it ensures that MCTS focuses computational resources on the brittle segments of the rewrite, and does not waste them on robust modifications. For example, when the root node in Figure 7 is first processed, then Nodes 1, 2, 3 and 4 are expanded and created sequentially since the token probabilities are higher than θ (Lines 24–26 of the while loop). Only when Node 4 is reached, two new nodes (5 and 15) are created as its children.

3. Simulation: The node n_{cur} is expanded greedily based on the highest-probability tokens until a terminal node is reached (Line 21). The rewritten query obtained from this greedy expansion is then used to compute the potential $V(n_{cur})$ (Line 22) as follows: for a valid rewrite, $V(n_{cur})$ represents the *speedup* it provides w.r.t. the original query. However, if invalid (i.e. syntactically or semantically incorrect), $V(n_{cur})$ is assigned a zero value. After every simulation, the rewritten query obtained after the greedy expansion of n_{cur} is cached along with $V(n_{cur})$ in a map, *Potential* (Line 23).

For instance, to compute $V(n_{15})$ in Figure 7, the blue path in the tree is greedily expanded to identify its rewrite potential. That is, the partial rewrite “select * from B left outer join A;” is used to evaluate $V(n_{15})$.

4. Back Propagation: The V value of the simulation for n_{cur} is back-propagated to all its ancestor nodes, and their values are updated iff the new value is higher than the existing value.

Rewritten Query. At the end of all iterations, $q \in Potential$ with highest value $Potential[q]$ that is greater than 1 is returned as the rewritten query (Lines 29–30). In case no such rewrite exists, implying that all the valid rewrites are slower than the original query, the original query itself is returned (Line 32).

5.2 Input Prompt to MCTS

The root state in Algorithm 1 corresponds to the state just after the prompt is fed to the LLM. One way to use this algorithm is to execute it for all the various prompts discussed in the previous sections, and choose the rewrite that provides the best performance. This, however, is expensive both from the aspect of query rewrite time, as well as the number of LLM tokens used. To minimize these costs, given a user query, the LITHE workflow first selects the prompt yielding the most effective rewrite from among the techniques of Sections 3 and 4. It then employs this prompt to initiate the MCTS-based rewrite. In case no prompt provides a lower-cost rewrite, baseline Prompt 1 of Section 3 is used as the fallback option.

5.3 Performance

LITHE’s performance with MCTS-based rewrites fully matches the human target, yielding **10** CPRs and the maximum CSGM of **11.84**. While these gains over just prompting may seem marginal, we wish make two observations: (1) MCTS extracts maximum improvement wherever possible, and (2) As shown in Section 7.5.2, its impact becomes significant for smaller LLMs like LLaMA.

Discussion. Our approach implicitly assumes that LLMs *do* possess some degree of knowledge about query optimization – for example, when presented with a poorly written query versus an optimized one, the LLM is usually capable of identifying why the optimized version would execute better. This expertise arises from the extensive training corpus of LLMs, which includes a large repository of books, papers, and articles on query optimization and execution.

However, this core knowledge is not fully sufficient to provide good rewrites, and there is a clear need for augmentation with explicit query-specific guidance. This is particularly so when there are multiple

plausible optimization paths, or when the LLM may be subtly misled – this is the reason we leverage MCTS to systematically explore alternatives based on token probabilities. Further, there is explicit evaluation via the database optimizer of the expected improvement in quantitative performance. This separation ensures that our method leverages the LLM’s generative capabilities to posit rewrites while relying on an external ground-truth signal (i.e. the database optimizer) to confirm rewrite quality.

6 Implementation Choices

In this section, we briefly discuss the design choices made in our implementation of LITHE.

6.1 LITHE Parameter Settings

The “*temperature*” parameter of GPT-4o, which ranges over $[0,1]$, controls randomness of the model’s response. While higher values suit creative tasks that require diverse and exploratory outputs, we seek focused, deterministic answers. Thus, we set the temperature to 0, forcing the model to greedily sample the next token.

The hyperparameters used by LITHE for MCTS are as follows: The maximum number of iterations $iter_{max}$ is set to 8, expansion threshold θ is 0.7, and number of expansions k is 2. The values of c_{base} and c were set to 10 and 4, respectively. These settings were determined after an empirical evaluation of the various trade-offs, providing a robust balance between efficiency and quality.

Finally, we try a maximum of 5 times to fix, via prompt corrections, any rewrite that exhibits syntax errors (Section 2) – if this threshold is crossed, we drop the rewrite.

6.2 Query Equivalence Testing

We use a multi-stage approach, described below, to help the DBA test semantic equivalence between the original query and a recommended rewrite. As mentioned earlier, having the DBA in the loop is a common practice in commercial query advisory systems [17].

1. Result Equivalence via Sampling. We use a sampling-based approach to quickly test equivalence in the rewrite generation stages of the pipeline. The idea here is to execute the queries on several small samples of the database and verify equivalence based on the sample results. However, while this test is a necessary condition for query equivalence, it is not sufficient. That is, false positives may be present because the sampled database may not cover all the predicates featured in the query. To minimize this likelihood, we use a combination of (1) *correlated sampling* [54] for maintaining join integrity in the sample, (2) adjusting constants in the filter predicates to produce populated results, and (3) injecting query-specific synthetic tuples in the sample, similar to the XData mutant-killing tool [42], to cover predicate boundary conditions – the complete details are in Section C.

2. Logic-based Equivalence. Although verifying the equivalence of a general pair of SQL queries is NP-complete [10], a variety of logic-based tools (e.g. Cosette[15], SQL-Solver [20], VeriEQL [26], QED [49]) are available for proving equivalence over restricted classes of queries. Once the least-cost rewrite is identified, LITHE uses QED [49] and SQLSolver [20] since these two approaches covered a larger set of queries compared to the alternatives. The advantage of such a logic-based approach is that it is definitive in outcome. Note that this rewrite has already passed the sampling-based tests described above.

3. Result Equivalence on the Entire Database. If the logic-based test is inconclusive, result equivalence is evaluated on the entire database itself. The DBA may choose to prematurely terminate this test in case the checking time is found to be excessive.

6.3 Regression Identification

We use heuristics to predict whether the actual execution time of a given promising rewrite may turn out to be slower than the original query. The first heuristic is based on the SEER robust plan identification algorithm [22], which is predicated on the database engine providing a “force-plan” feature (e.g. [1]). This approach is based on the notion of *replacement safety*, wherein the replacement plan is always either better or within a small sub-optimality factor of the original optimizer-recommended choice. It was theoretically shown that a safe plan replacement at the *perimeter* of the selectivity space is sufficient to infer safe replacement in the interior of the space as well. The proof used a generic characterization of plan cost functions and the first and second derivatives of these functions at the perimeter. As a matter of efficiency, it was empirically found that instead of the entire perimeter, it was usually sufficient to consider safety just at the *corners*.

We use the above result to design a regression test as follows: let P^U and P^T be the plans corresponding to the original query, Q^U , and the recommended rewrite, Q^T , respectively. We construct *parametrized* versions of these queries, where the constants in the filter predicates are replaced by variables. Then, by assigning appropriate values to these parameters, we construct queries that are located at the *corners* of the selectivity space. The plans P^U and P^T are forced at each of these corner locations, and if the rewrite’s cost is lower than the original at *all* of them, the rewrite is deemed to be robust.

At this time, plan forcing is not available in all database engines. For such limited environments, we fall back to an alternative heuristic – we compare the query runtimes obtained on the series of executions carried out on sampled databases during semantic equivalence testing. If the rewrite’s runtimes are lower in *all* of these intermediate executions, then it is estimated to be robust. We found this extrapolation heuristic to be effective in our experiments due to two reasons: First, the sampled databases, while relatively small compared to the original, are still large in absolute terms (1 GB or more). Second, since most rewrites are related to redundancy removals, a relative performance improvement obtained due to such a rule holds true irrespective of the data sizes.

7 Experimental Evaluation

In this section, we report on LITHE’s performance profile. We first describe the experimental setup, including comparative baselines, query suites and evaluation platforms. Then we present the speedup results for both aggregate benchmark and individual queries, followed by characterization of the rewrite overheads in computational and financial terms. We finally discuss the impact of alternative platforms wrt database engine, database schema and LLMs.

Rewrite Baselines. We compare LITHE with a collection of contemporary rewrite techniques, collectively referred to as SOTA – the details of these techniques are provided in Section 9. Specifically, the SOTA collection consists of the following approaches:

1. Baseline LLM prompt [34]: This is Prompt 1 from Section 3.
2. Learned Rewrite [57], a purely rule-based rewriter.
3. LLM-R² [33], an LLM-guided rule-based rewriter.
4. GenRewrite [34], a purely LLM-based rewriter.

Given an input query, each of the SOTA approaches is independently invoked to perform a rewrite, and the rewrite with the *best* performance is used as the baseline for comparison. Note that these approaches

may occasionally generate rewrites that are expected by the optimizer’s costing module to regress the performance. For safety, we immediately discard such rewrites, similar to LITHE.

Query Set. Our evaluation in this paper primarily focuses on complex analytical queries from the standard TPC-DS decision-support benchmark [16], which models a retail supplier environment. The benchmark is used at its default size of 100 GB. As mentioned earlier, we focus on “slow queries” that take over 10 seconds to execute in their original form, operating in a cold-cache environment.

LITHE has also been evaluated on other benchmarks, including DSB [18], ARCHER [56], JOB [29] and StackOverflow [36]. The overall performance characteristics were found to be similar to those presented here – see Section A for details.

Testbed. The experiments were carried out on the following data processing platform: Sandbox server with Intel(R) Xeon(R) CPU E5-1660, 32 GB RAM, and 12 TB HDD, running Ubuntu 22.04 LTS. A majority of the experiments used PostgreSQL v16 database engine and GPT-4o LLM for both LITHE and SOTA. Commercial engines and other LLMs are considered in Sections 7.4 and 7.5.

Metrics. For each rewrite technique, we identified the number of queries for which a CPR (cost productive rewrite with ≥ 1.5 speedup) could be constructed. Subsequently, we computed the CSGM (Cost Speedup Geometric Mean) and TSGM (Time Speedup Geometric Mean) performance obtained by each technique over the set of all CPRs (i.e. CPRs arising from either LITHE or SOTA). Cost speedups are computed relative to the native optimizer, and runtime speedups are measured as the ratio of original query runtime to rewritten query runtime in a cold-cache environment.

From the investment perspective, we measured the average rewrite time per query, and additionally for the LLM-based techniques, the number of tokens used in the rewrite process.

7.1 Rewrite Quality (Cost and Time)

7.1.1 Estimated Cost

LITHE produces a rewrite with a positive cost speedup ($\geq 1x$) for 46 of the 88 TPC-DS queries deemed to be slow by our threshold. Of these 46, there were **26** CPRs resulting in a *highly productive* CSGM of **11.5**. On the other hand, SOTA delivers only **13** CPRs (out of 42 positive rewrites) with a CSGM of **6.1**. All but one of the SOTA CPRs also feature in the LITHE CPRs, making the total number of CPRs considered to be 27. Of these 27, we were able to formally verify 11 using the logic-based tools, whereas the remaining 16 passed our statistical tests. Furthermore, we also manually verified the correctness of these rewritten queries.

A drill-down into the cost speedup performance at the granularity of individual queries is shown in Figure 8, which compares LITHE (orange bars) and SOTA (blue bars) on each of the 27 CPR queries – note that the cost speedups on the x -axis are tabulated on a \log_{10} scale, and the queries are sequenced in decreasing order of LITHE speedup. The vertical dotted line at 1 represents the normalized baseline cost of the original query with the native optimizer, while the vertical line at 1.5 is the CPR threshold.

We first observe, gratifyingly, that rewrites are indeed capable of promising dramatic cost speedups over the native engine – take, for instance, Q41, which improves by a whopping *five orders-of-magnitude* for both SOTA and LITHE. This improvement in query performance is due to replacing the “WHERE (SELECT COUNT(*) from ...) > 0” clause with “WHERE EXIST (SELECT 1 from ...)” – the latter efficiently checks result existence in an inner subquery since it removes the costly aggregation function.

Second, in as many as 15 queries, LITHE’s cost speedup *substantively exceeds* SOTA, while in the remainder it largely matches SOTA. In fact, for several queries (e.g. Q45, Q25, Q4) LITHE produces a highly beneficial CPR but SOTA returns the original query. The only case where SOTA is appreciably

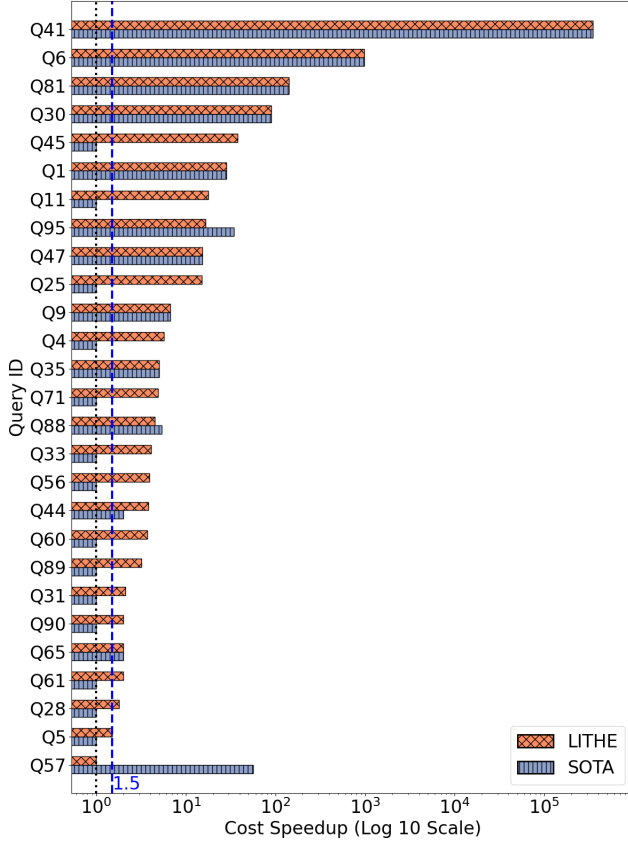


Figure 8: Plan Cost Speedups

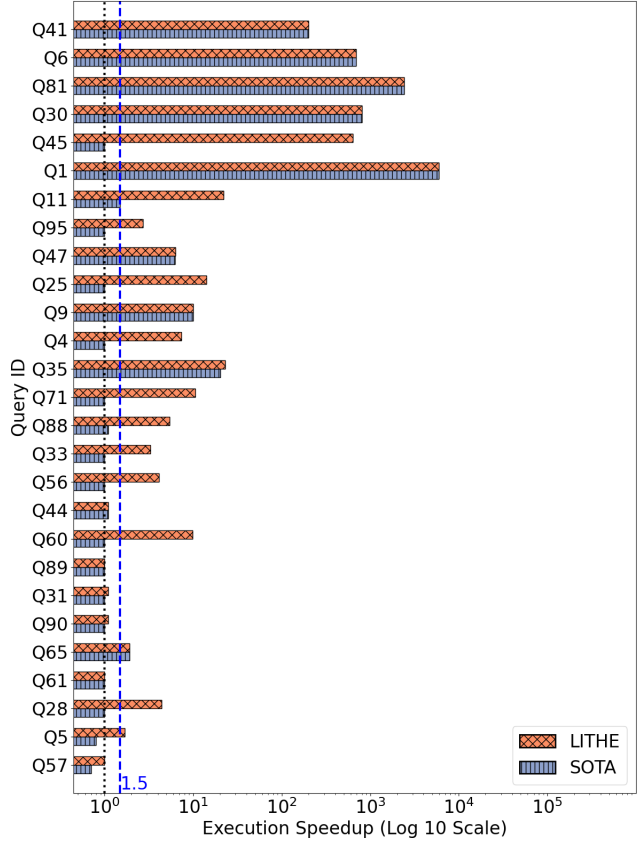


Figure 9: Execution Time Speedups

better is Q57 where it projects a large speedup but LITHE settles for the original query. And in Q88 and Q95, SOTA performs marginally better.

At this stage, one might expect that adding more rules to LITHE could bring it on par with SOTA for queries like Q57. However, we deliberately include only broad-brush rules in LITHE to ensure generalizability and efficiency. Moreover, as the following timing section shows, due to this conservative approach LITHE actually *outperforms* SOTA on these queries (Q57, Q88, Q95) at runtime!

We also analyzed queries where LITHE was unable to provide CPRs to determine whether there were common patterns that led to unsuccessful outcomes. Our audit indicated that the primary reasons were: (a) Structural Simplicity – for instance, flat SPJ queries without nesting that are already conducive to current optimizers (e.g. Q3, Q15, Q52), or (b) Structural Tightness – for instance, queries without repeated computations where the scope for improvement is inherently limited (e.g. Q7, Q19, Q22). Further, some failures may be due to the *specific* LLM being used, and not the queries themselves (see Section 7.5.2). Therefore, one could consider implementing an *ensemble* of diverse LLMs to gain enhanced beneficial coverage.

7.1.2 Execution Time

Thus far, we had considered optimizer-estimated costs. We now move to wall-clock runtimes for query executions – Figure 9 shows the runtime speedups (on a \log_{10} scale) obtained by LITHE and SOTA. We again observe that there are indeed several queries where substantial time benefits are achieved by the rewrites over the native engine, even exceeding *order-of-magnitude* benefits in some cases – for instance, LITHE improves Q45 by a huge factor of 700! Further, for 12 of the 26 queries, LITHE

significantly outperforms SOTA – as a case in point, by **14** times in Q11. In the remainder, it matches SOTA, including as mentioned above, queries where SOTA’s optimizer costs were better.

From a modeling perspective, we see that the well-documented gap between optimizer predictions and actual run-times is prevalent in the rewrite space as well. On the one hand, there is Q45 where the projected speedup of 40 increases to 700 at runtime, whereas on the other, the 10^5 speedup for Q41 decreases to 200. But for SOTA, the reductions can be severe – a striking case in point is Q57, where SOTA actually causes *regression* despite a speedup projection of close to 100. However, the good news is that with LITHE, although the runtime speedups did not always match the projections, regressions were not encountered among the CPR rewrites thanks to the sampling-based checks. Overall, LITHE produces more robust rewrites resulting in a high TSGM of **13.2**, whereas SOTA only delivers TSGM of **4.9**.

7.1.3 Reasoning

As a confirmatory exercise, we compared the LLM-generated explanation output by LITHE with our own manual analysis of the query plans generated for the original and rewritten queries. For example, consider Q45 – the provided explanation correctly lists the main reasons for the 700x speedup as “*the use of CTEs*”, “*Reduction in data volume early on*” due to pre-filtering based on join predicates, etc. Overall, we found that the explanations provided by LITHE matched our manual analysis of the plans, indicating that model-based reasoning is well aligned with human-backed reasoning in these scenarios.

7.2 Ablation Analysis of Rewrite Quality

7.2.1 Components of LITHE

A natural question at this stage is the role of the various techniques in LITHE towards achieving its large performance benefits. This analysis is captured in Table 5 which lists the CPRs for each technique when invoked in *isolation* as well as the *cumulative* number of CPRs when the different techniques are combined in the order of listing.

Table 5: CPRs contributed by LITHE components.

| | # CPR | |
|-----------------------------|----------|------------|
| | Isolated | Cumulative |
| Basic Prompts (Section 3) | 15 | 15 |
| Rules R1 — R4 (Section 4.1) | 10 | 18 |
| Rules R5, R6 (Section 4.2) | 11 | 25 |
| MCTS (Section 5) | 26 | 26 |

We observe that the Basic Prompt ensemble and Redundancy Rules (R1 – R4) contribute to around two-thirds of the CPRs (18/26). When the statistics-infused Rules (R5, R6) are added, this number jumps to 25. Finally, MCTS in isolation (with Prompt 1 as the seed prompt) produces only 11 CPRs. However, by using the best prompt as seed, an additional CPR is obtained and the cost speedup of a pre-existing CPR is also improved.

7.2.2 Database-Sensitive Prompts.

Since MCTS explores the search space at a fine granularity, one could ask whether just the Basic Prompts in conjunction with MCTS would suffice to provide good performance. The motivation is that

it would relieve us from using the database-sensitive rules R1–R6 which incur significant computational and financial overheads. When this experiment was conducted, the CSGM dropped precipitously to a paltry **5**, a far cry from the 11.5 obtained with the database-sensitive rules. These results highlight the need to reflect database awareness for effective query rewriting, and not rely solely on prior LLM knowledge.

7.3 Rewrite Overheads (Time/Money)

Having established the performance benefits of rewrites, we now turn our attention to their time and financial overheads.

7.3.1 Transformation Time

We first look at the end-to-end time required to perform a rewrite (i.e. to run the entire pipeline shown in Figure 3). Table 6 shows the average processing time per CPR query, where we see that the LITHE rewrite process takes a few minutes. However, note that this investment may be acceptable in deployment given that the execution benefits typically far outweigh the compilation overheads. For instance, with Q11, the original query took nearly *an hour* to complete, whereas the LITHE rewrite executed in under 3 *minutes*. Further, many applications tend to use a set of canned queries which are run thousands of times. Thus, even a large one-time investment can be easily recovered over repeat executions of such queries.

Table 6: Rewrite Overheads of LITHE and SOTA.

| | Avg. Time (min) | Avg. Tokens | Avg. Cost (USD) |
|-------|-----------------|-------------|-----------------|
| LITHE | 5 | 18427 | 0.045 |
| SOTA | 1.7 | 20076 | 0.050 |

Notwithstanding the above, we also observe that LITHE is considerably slower than SOTA in producing rewrites. A drill-down showed the lion’s share of the time is taken by the initial prompt ensemble and the final MCTS module. We explore techniques to improve their efficiency in Section B.

7.3.2 Monetary Outlay

The average number of LLM tokens required by LITHE and SOTA, and their associated financial costs¹, are also shown in Table 6. The good news is that the inference charges per query are just a few cents, making rewriting practical from a deployment perspective.

7.4 Commercial DBMS

A legitimate question could be whether the rewrites made amends for the PostgreSQL optimizer but may fail to be useful in highly-engineered database engines. To evaluate this issue, we performed TPC-DS rewrites on a pair of popular commercial DBMS, which we refer to as **OptA** and **OptB**.

The performance of LITHE and SOTA on these two systems is shown in Table 7, with LITHE continuing to do better than SOTA. Despite the apparent lack of optimization headroom, LITHE still produces 12 and 9 CPRs resulting in a healthy CSGM of **3.6** and **4.1**, respectively. Further, the TSGM provided by these rewrites are a useful **2.1** and **1.9**, respectively.

¹At the time of writing, GPT-4o costs USD 2.5 per million tokens.

Table 7: Performance on Commercial Database Engines.

| | # CPR | | CSGM | | TSGM | |
|-------|-------|------|------|------|------|------|
| | OptA | OptB | OptA | OptB | OptA | OptB |
| LITHE | 12 | 9 | 3.6 | 4.1 | 2.1 | 1.9 |
| SOTA | 3 | 5 | 1.5 | 1.3 | 1.4 | 1.2 |

Interestingly, although we did not observe any regressions with PostgreSQL, a few did surface in the commercial systems. Nevertheless, our regression identification mechanisms effectively caught these brittle rewrites. As a case in point, a promising rewrite estimated by one of the commercial optimizers for Q23, actually takes 37 minutes to complete as compared to 18 minutes for the original query – this doubling slowdown was successfully flagged by the SEER heuristic (Section 6.3), and the rewrite was abandoned.

The above results suggest LITHE has a useful role to play in industrial environments. From a different perspective, a company building a new database engine could use LITHE to non-invasively overcome the limitations of early versions of its optimizer.

7.5 Dependence on LLM

7.5.1 Impact of Training

An interesting question to ask now is whether the performance benefits seen thus far could be an artifact of GPT-4o having already been trained well on the TPC-DS benchmark, which is prominent in the public domain.

To investigate this issue, we ran LITHE on two datasets that were confirmed to be *unknown* to the GPT-4o version used in our evaluation. The first is the recently released Football benchmark [24] for Text-to-SQL evaluation. From the query suite, which largely consists of simple SPJ and Union queries, we selectively picked the few complex queries. Further, we created our own complex queries to include in the test suite. On this enriched workload, LITHE on PostgreSQL produced 11 CPRs with CSGM of 12 and TSGM of 1.5, while SOTA gave 6 CPRs with CSGM of 3 and TSGM of 1.1.

The second dataset is a proprietary customer benchmark used extensively by an industrial development team. For this environment, LITHE on a commercial engine produced 8 CPRs with a CSGM of 2 and a TSGM of 3.3. Note that these improvements are achieved *despite* this engine having been fine-tuned on this workload over an extended period. Moreover, the numbers are consistent with those reported for TPC-DS in Section 7.4.

7.5.2 Impact of Model

In our concluding experiment, we evaluated LITHE’s performance on the LLaMA 3.1 70 billion parameter instruct model, substantially smaller compared to GPT-4o, which may have several hundred billion parameters [8]. To attain practical inference times, the model was loaded using a low 4-bit quantization. Further, to ensure reproducibility and deterministic answers, the *do_sample* parameter was set to *False*, which forces the LLM to perform greedy decoding. To make up for the huge reduction in model parameters as compared to GPT-4o, we include up to two example demonstrations for each rule-based prompt.

For this environment, Table 8 shows LITHE’s performance on the TPC-DS queries with and without MCTS. Although certainly lower than the corresponding numbers with GPT-4o (Section 7.1), it is encouraging to see that, in absolute terms, significant performance benefits can be obtained for most

queries, especially with MCTS support. So, the message is that smaller models can also be fruitfully used in real-world environments.

Table 8: LITHE Rewrite Performance with LLaMA

| | # CPR | CSGM | TSGM |
|--------------------|-------|------|------|
| LLaMA without MCTS | 18 | 5.6 | 6.5 |
| LLaMA with MCTS | 22 | 8.5 | 10.9 |

Finally, we also evaluated LITHE on **Gemini 2.5 Flash**. Here, LITHE produced 17 CPR with CSGM of 7.2 and TSGM of **5.8**. Although not as strong as GPT-4o, these results highlight the potency of LITHE’s recommendations across a range of LLM platforms.

8 The Road Ahead

Based on our study, we now present a few observations with implications for the future design and deployment of rewriting tools.

8.1 Rewrite Space Coverage by LLMs

Given the decades-long research on database query optimization, we expected the potential for performance improvement via rewriting to be limited. What came as a surprise was the substantial scope for improvement still available, as showcased by the large CSGM and TSGM values, even on commercial platforms. These results suggest that LLMs explore optimization spaces that are well outside the purview of contemporary database engines. Specifically, rewrites in the query space appear to guide the optimizer to explore fresh regions of the plan space not evaluated with the original query formulation. Further, this enhanced space could be augmented, in a two-stage process, with the recent proposals for LLM-based “plan hints” that steer the optimizer in fruitful directions within a plan space [11]. We plan to explore these space relationships and options in our future work.

8.2 Rewrite Migration to Optimizer

The above demonstrated the potent exploratory power of LLMs. But from an overheads perspective, such rewrites should ideally be within the optimizer’s native search space rather than recommended from outside. Therefore, it would be a useful exercise to try and distill fresh optimization rules from these instances, leveraging the extensibility features of contemporary optimizers [19] to facilitate their incorporation in existing systems.

On the flip side, there appears to be an “impedance mismatch” against such integration for certain classes of rewrites. For example, consider the TPC-DS Q90 rewrite in Figure 10. The original query individually computed AM (morning) sales and PM (evening) sales, which were then used to compute the AM to PM ratio. The rewrite, however, extracted all relevant rows in one shot and computed the ratio using CASE statements – encoding such transformations as generic rules in the optimizer appears challenging, given the combinatorial ways in which such transformations can occur.

Therefore, a fruitful area of future research could be achieving a middle-ground between the disparate world-views of LLMs and traditional optimizers.

```

SELECT CAST(
  SUM( CASE WHEN time_dim.t_hour BETWEEN 8 AND 9
            THEN 1 ELSE 0 END) AS DECIMAL(15, 4)
) / CAST(
  SUM( CASE WHEN time_dim.t_hour BETWEEN 19 AND 20
            THEN 1 ELSE 0 END) AS DECIMAL(15, 4)
) AS am_pm_ratio
FROM web_sales
JOIN household_demographics
  ON ws_ship_hdemo_sk = household_demographics.hd_demo_sk
JOIN time_dim ON ws_sold_time_sk = time_dim.t_time_sk
JOIN web_page ON ws_web_page_sk = web_page.wp_web_page_sk
WHERE household_demographics.hd_dep_count = 6
  AND web_page.wp_char_count BETWEEN 5000 AND 5200
  AND (time_dim.t_hour BETWEEN 8 AND 9
       OR time_dim.t_hour BETWEEN 19 AND 20)
ORDER BY am_pm_ratio
limit 100;

```

Figure 10: Rewritten TPC-DS Q90

8.3 Agentic LLMs for Query Rewriting

An effective way to extend LITHE is to use an agentic LLM that actively interacts with the database environment. This would allow a “LITHE agent” to leverage database-related “tools” such as the query optimizer’s cardinality estimation and cost estimation modules, as well as rewrite-related tools such as rule generation, semantic validation, and regression checking. A memory store could be attached to the agent to log prior interactions with users and rewrite rules learnt over time, thus supporting continuous learning. The agent can decide the order in which these tools are invoked and make appropriate choices during the rewrite process. Finally, combining our MCTS formulation with tool-augmented agentic reasoning would also be a promising avenue for future work.

8.4 Scope of Semantic Equivalence Tools

As seen in the experiments section, logic-based query equivalence testing covers industrial-strength queries only to a limited extent. On the other hand, while it is highly likely that the statistics-verified rewrites are valid, it still requires the DBA to make a final call on the correctness. This limitation restricts the use of LITHE in a fully automated scenario, i.e., as a direct preprocessor to the query engine. Therefore, a key challenge is to improve logic-based coverage.

8.5 Enhancing NL2SQL approaches with Rewrites

LITHE could potentially be used to enhance NL2SQL approaches to produce performant SQL queries. One way would be to incorporate some rewrite rules directly into NL2SQL prompts. Based on our empirical assessment (Figure 11b), we would recommend the simplest Prompt 1 from the Basic Prompts in conjunction with a subset of Rules R1–R4. To avoid multiple prompts, one could use the LLM itself to identify the appropriate rule based on the NL query. While these prompts might not always provide the best rewrites, our experiments show that they are sufficient for many scenarios while incurring only small overheads as compared to using the entire prompt suite and MCTS search.

An alternate non-invasive approach would be to simply use LITHE as a post-processor after the initial translation, but again restricted to using only a subset of the prompts/rules to keep the overheads within permissible limits.

9 Related Work

Rule-based SQL rewriting. Most of the recent work on SQL query rewriting is rule-based [57, 51, 12, 52, 14, 39]. For instance, WeTune [51] uses a rule generator to enumerate a set (up to a maximum size) of logically valid plans for a given query to create new rewrite rules, and uses an SMT solver to prove the correctness of the generated rules. While this approach can generate a large set of new rewrite rules, it often fails in coming up with transformation rules for complex queries due to the computational overheads of verifying rule correctness. As such, it is unable to rewrite any of the TPC-DS queries [21, 34].

Learned Rewrite [57] uses existing Calcite [13] rules and aims to learn the optimal subset of rules along with the order in which they must be applied. Since the rewrite search space grows exponentially with the number of rules, it uses MCTS scheme to efficiently navigate this space and find the rewritten query with maximum cost reduction. This is in contrast to LITHE where MCTS is used to search over the output token space of an LLM in order to prevent LLM hallucinations. Specifically, our use of MCTS ensures that the LLM correctly follows the prompt instructions.

LLM-R² [33] is also rule-based but takes a different approach to identify the order for rewrite rule applications: it uses an LLM to find the best Calcite rules and the order in which to apply them to improve the query performance. R-Bot [45] also leverages an LLM to optimize the order of Calcite rules, but employs advanced contemporary techniques such as retrieval-augmented generation (RAG) and step-by-step self-reflection to improve the outcomes.

Query Booster [12] implements human-centered rewriting – it provides an interface to specify rules using an expressive rule language, which it generalizes to create rewrite rules to be applied on the query. There are also rule-based rewrite approaches designed for specific types of rewrites such as optimizing correlated window aggregations [52] and common expression elimination [14].

All of the above approaches operate via the query *plan space*, which can restrict the kind of rewrites that can be accomplished. Whereas, LITHE uses a small set of general rewrite rules that work directly in the *query space*.

LLM-based rewriting. GenRewrite [34] is the first LLM-based approach to use the LLM for end-to-end query rewriting. Instead of using predefined rules from Calcite [13], they employ the LLM to create Natural Language Rewrite Rules (NLR2s) to be used as hints, and perform several iterations of prompting to get the rewritten query. They show that LLMs can outperform rule-based approaches due to their ability to understand contexts, and demonstrate a significant improvement in query rewriting compared to prior methods.

A limitation, however, is that LLM-generated rewrite rules often fail to generalize beyond specific query pairs. Even when generalized rules are present, LLMs can struggle to apply rules correctly if not provided with accompanying examples. Finally, it must be noted that LLMs are unaware of the underlying database which restricts their ability to produce efficient metadata-aware rules.

Query rewriting in SQL engines. Query rewrite prior to plan enumeration has long been a feature of contemporary (non-LLM) database optimizers [9, 25]. However, the difference in our context is that the LLM gives us the flexibility to choose from a wide variety of approaches, not restricted to the hardwired options within individual systems.

LLMs for Database Modules. LLM technologies have been advocated for a variety of database modules. For instance, they have been extensively used for Text-to-SQL transformations [32, 53, 4, 5, 46, 58]. Very recently, statistical metadata was leveraged to also improve the Text-to-SQL generation process [43]. The main focus of these techniques is to correctly ascertain the information necessary to formulate the SQL query [48, 40, 46].

On the other hand, the goal of SQL-to-SQL rewriting is on improving the performance of an existing

SQL query. Thus, unlike Text-to-SQL transformations where the input text is inherently ambiguous, SQL queries are precisely defined, and therefore equivalence to a precise ground-truth has to be provably maintained.

In recent times, LLMs have also been considered for a variety of database modules, including plan-hinting [11], join-order optimization [47], index selection [58], data pipelines [28], data management [31], and multi-modal data optimization [50]. These approaches can be used in conjunction with LITHE since they address orthogonal segments of the query processing pipeline.

10 Conclusions and Future Work

We investigated how the latent power of LLM technologies can be productively materialized in the context of SQL-to-SQL rewriting. Our study progressively infused database domain knowledge, such as redundancy removal rules and schematic+statistical metadata, into the LLM prompts. Further, the output telemetry of LLMs, in the form of token probabilities, was used to signal situations where the LLM lacked confidence, triggering exploration of a larger search space. Finally, a combination of logic-based and statistical tests was employed to verify the equivalence of the rewrites.

An empirical evaluation over common database benchmarks showed that rewriting is a potent mechanism to improve query performance. In fact, even order-of-magnitude speedups were routinely achieved with regard to both abstract costing and execution times. However, our results also showed a significant semantic distance between foundation models and query optimizers, with regard to both scope and precision, which would have to be bridged to fully leverage the latent power of LLMs. Further, our focus here was primarily on prompting-based strategies – in our future research, we plan to investigate how domain-specific *fine-tuning* could be leveraged to provide GPT-4o-like rewrites on small open models.

References

- [1] Sqlserver xml plan forcing, 2012.
- [2] View SQL Queries From Your Code With Prefix, 2019.
- [3] Postgresql release 16, 2023.
- [4] Defog sql-eval, 2024.
- [5] Defog sqlcoder, 2024.
- [6] Groupby with joined table produces sub optimal query, april 2024.
- [7] MySQL 8.4 Reference Query Manual – The Slow QUery Log, 2024.
- [8] OpenAI’s GPT 4o vs GPT 4o mini: Which AI model to use and why., 2024.
- [9] R. Ahmed, A. Lee, A. Witkowski, D. Das, H. Su, M. Zait, and T. Cruanes. Cost-based query transformation in oracle. In *Proceedings of the 32nd International Conference on Very Large Data Bases*, VLDB ’06, page 1026–1036. VLDB Endowment, 2006.
- [10] A. V. Aho, Y. Sagiv, and J. D. Ullman. Equivalences among relational expressions. *SIAM J. Comput.*, 8(2):218–246, 1979.

- [11] P. Akioyamen, Z. Yi, and R. Marcus. The unreasonable effectiveness of llms for query optimization. *CoRR*, abs/2411.02862, 2024.
- [12] Q. Bai, S. Alsudais, and C. Li. Querybooster: Improving SQL performance using middleware services for human-centered query rewriting. *Proc. VLDB Endow.*, 16(11):2911–2924, 2023.
- [13] E. Begoli, J. Camacho-Rodríguez, J. Hyde, M. J. Mior, and D. Lemire. Apache calcite: A foundational framework for optimized query processing over heterogeneous data sources. In G. Das, C. M. Jermaine, and P. A. Bernstein, editors, *Proceedings of the 2018 International Conference on Management of Data, SIGMOD Conference 2018, Houston, TX, USA, June 10-15, 2018*, pages 221–230. ACM, 2018.
- [14] N. Bruno, J. Debrodt, C. Song, and W. Zheng. Computation reuse via fusion in amazon athena. In *38th IEEE International Conference on Data Engineering, ICDE 2022, Kuala Lumpur, Malaysia, May 9-12, 2022*, pages 1610–1620. IEEE, 2022.
- [15] S. Chu, C. Wang, K. Weitz, and A. Cheung. Cosette: An automated prover for SQL. In *8th Biennial Conference on Innovative Data Systems Research, CIDR 2017, Chaminade, CA, USA, January 8-11, 2017, Online Proceedings*. www.cidrdb.org, 2017.
- [16] T. T. P. P. Council. TPC Benchmark™ DS (tpc-ds). In *TPC Benchmark DS (Decision Support)*, 2006.
- [17] B. Dageville, D. Das, K. Dias, K. Yagoub, M. Zaït, and M. Ziauddin. Automatic SQL tuning in oracle 10g. In M. A. Nascimento, M. T. Özsu, D. Kossmann, R. J. Miller, J. A. Blakeley, and K. B. Schiefer, editors, *(e)Proceedings of the Thirtieth International Conference on Very Large Data Bases, VLDB 2004, Toronto, Canada, August 31 - September 3 2004*, pages 1098–1109. Morgan Kaufmann, 2004.
- [18] B. Ding, S. Chaudhuri, J. Gehrke, and V. R. Narasayya. DSB: A decision support benchmark for workload-driven and traditional database systems. *Proc. VLDB Endow.*, 14(13):3376–3388, 2021.
- [19] B. Ding, V. Narasayya, and S. Chaudhuri. Extensible query optimizers in practice. *Foundations and Trends® in Databases*, 14(3-4):186–402, 2024.
- [20] H. Ding, Z. Wang, Y. Yang, D. Zhang, Z. Xu, H. Chen, R. Piskac, and J. Li. Proving query equivalence using linear integer arithmetic. *Proc. ACM Manag. Data*, 1(4):227:1–227:26, 2023.
- [21] R. Dong, J. Liu, Y. Zhu, C. Yan, B. Mozafari, and X. Wang. Slabcity: Whole-query optimization using program synthesis. *Proc. VLDB Endow.*, 16(11):3151–3164, 2023.
- [22] H. Doraiswamy, P. N. Darera, and J. R. Haritsa. Identifying robust plans through plan diagram reduction. *Proc. VLDB Endow.*, 1(1):1124–1140, 2008.
- [23] E. Fadeeva, A. Rubashevskii, A. Shelmanov, S. Petrakov, H. Li, H. Mubarak, E. Tsymbalov, G. Kuzmin, A. Panchenko, T. Baldwin, P. Nakov, and M. Panov. Fact-checking the output of large language models via token-level uncertainty quantification. In L. Ku, A. Martins, and V. Srikumar, editors, *Findings of the Association for Computational Linguistics, ACL 2024, Bangkok, Thailand and virtual meeting, August 11-16, 2024*, pages 9367–9385. Association for Computational Linguistics, 2024.

- [24] J. Fürst, C. Kosten, F. Nooralahzadeh, Y. Zhang, and K. Stockinger. Evaluating the data model robustness of text-to-sql systems based on real user queries. In A. Simitsis, B. Kemme, A. Queralt, O. Romero, and P. Jovanovic, editors, *Proceedings 28th International Conference on Extending Database Technology, EDBT 2025, Barcelona, Spain, March 25-28, 2025*, pages 158–170. Open-Proceedings.org, 2025.
- [25] C. Galindo-Legaria and M. Joshi. Orthogonal optimization of subqueries and aggregation. In *Proceedings of the 2001 ACM SIGMOD International Conference on Management of Data, SIGMOD ’01*, page 571–581, New York, NY, USA, 2001. Association for Computing Machinery.
- [26] Y. He, P. Zhao, X. Wang, and Y. Wang. Verieql: Bounded equivalence verification for complex SQL queries with integrity constraints. *CoRR*, abs/2403.03193, 2024.
- [27] Hibernate ORM, 2025.
- [28] S. B. Junior, P. Ceravolo, S. Groppe, M. Jarrar, S. Maghool, F. Sèdes, S. Sahri, and M. van Keulen. Are large language models the new interface for data pipelines? In P. Cudré-Mauroux, A. Kö, and R. Wrembel, editors, *Proceedings of the International Workshop on Big Data in Emergent Distributed Environments, BiDEDE 2024, Santiago, Chile, June 9-15, 2024*, pages 6:1–6:6. ACM, 2024.
- [29] V. Leis, A. Gubichev, A. Mirchev, P. A. Boncz, A. Kemper, and T. Neumann. How good are query optimizers, really? *Proc. VLDB Endow.*, 9(3):204–215, 2015.
- [30] V. Leis, A. Gubichev, A. Mirchev, P. A. Boncz, A. Kemper, and T. Neumann. How good are query optimizers, really? *Proc. VLDB Endow.*, 9(3):204–215, 2015.
- [31] G. Li, X. Zhou, and X. Zhao. LLM for data management. *Proc. VLDB Endow.*, 17(12):4213–4216, 2024.
- [32] J. Li, B. Hui, G. Qu, J. Yang, B. Li, B. Li, B. Wang, B. Qin, R. Geng, N. Huo, X. Zhou, C. Ma, G. Li, K. C. Chang, F. Huang, R. Cheng, and Y. Li. Can LLM already serve as A database interface? A big bench for large-scale database grounded text-to-sqls. In A. Oh, T. Naumann, A. Globerson, K. Saenko, M. Hardt, and S. Levine, editors, *Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023*, 2023.
- [33] Z. Li, H. Yuan, H. Wang, G. Cong, and L. Bing. LLM-R2: A large language model enhanced rule-based rewrite system for boosting query efficiency. *Proc. VLDB Endow.*, 18(1):53–65, 2024.
- [34] J. Liu and B. Mozafari. Query rewriting via large language models. *CoRR*, abs/2403.09060, 2024.
- [35] G. Lohman. Is Query Optimization a Solved Problem?, 2014.
- [36] R. Marcus. Stack dataset, 2021.
- [37] Entity Framework, 2025.
- [38] K. Opsahl-Ong, M. J. Ryan, J. Purtell, D. Broman, C. Potts, M. Zaharia, and O. Khattab. Optimizing instructions and demonstrations for multi-stage language model programs. In Y. Al-Onaizan, M. Bansal, and Y. Chen, editors, *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing, EMNLP 2024, Miami, FL, USA, November 12-16, 2024*, pages 9340–9366. Association for Computational Linguistics, 2024.

- [39] H. Pirahesh, J. M. Hellerstein, and W. Hasan. Extensible/rule based query rewrite optimization in starburst. In M. Stonebraker, editor, *Proceedings of the 1992 ACM SIGMOD International Conference on Management of Data, San Diego, California, USA, June 2-5, 1992*, pages 39–48. ACM Press, 1992.
- [40] M. Pourreza and D. Rafiei. DIN-SQL: decomposed in-context learning of text-to-sql with self-correction. In A. Oh, T. Naumann, A. Globerson, K. Saenko, M. Hardt, and S. Levine, editors, *Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023*, 2023.
- [41] M. L. Puterman. *Markov Decision Processes: Discrete Stochastic Dynamic Programming*. Wiley Series in Probability and Statistics. Wiley, 1994.
- [42] S. Shah, S. Sudarshan, S. Kajbaje, S. Patidar, B. P. Gupta, and D. Vira. Generating test data for killing SQL mutants: A constraint-based approach. In S. Abiteboul, K. Böhm, C. Koch, and K. Tan, editors, *Proceedings of the 27th International Conference on Data Engineering, ICDE 2011, April 11-16, 2011, Hannover, Germany*, pages 1175–1186. IEEE Computer Society, 2011.
- [43] V. Shkapenyuk, D. Srivastava, T. Johnson, and P. Ghane. Automatic metadata extraction for text-to-sql, 2025.
- [44] D. Silver, T. Hubert, J. Schrittwieser, I. Antonoglou, M. Lai, A. Guez, M. Lanctot, L. Sifre, D. Kumaran, T. Graepel, T. P. Lillicrap, K. Simonyan, and D. Hassabis. Mastering chess and shogi by self-play with a general reinforcement learning algorithm. *CoRR*, abs/1712.01815, 2017.
- [45] Z. Sun, X. Zhou, and G. Li. R-bot: An llm-based query rewrite system. *CoRR*, abs/2412.01661, 2024.
- [46] S. Talaei, M. Pourreza, Y. Chang, A. Mirhoseini, and A. Saberi. CHESS: contextual harnessing for efficient SQL synthesis. *CoRR*, abs/2405.16755, 2024.
- [47] J. Tan, K. Zhao, R. Li, J. X. Yu, C. Piao, H. Cheng, H. Meng, D. Zhao, and Y. Rong. Can large language models be query optimizer for relational databases? *CoRR*, abs/2502.05562, 2025.
- [48] B. Wang, C. Ren, J. Yang, X. Liang, J. Bai, Q. Zhang, Z. Yan, and Z. Li. MAC-SQL: A multi-agent collaborative framework for text-to-sql. *CoRR*, abs/2312.11242, 2023.
- [49] S. Wang, S. Pan, and A. Cheung. QED: A powerful query equivalence decider for SQL. *Proc. VLDB Endow.*, 17(11):3602–3614, 2024.
- [50] Y. Wang, H. Ma, and D. Z. Wang. No more optimization rules: Llm-enabled policy-based multi-modal query optimizer. *CoRR*, abs/2403.13597, 2024.
- [51] Z. Wang, Z. Zhou, Y. Yang, H. Ding, G. Hu, D. Ding, C. Tang, H. Chen, and J. Li. Wetune: Automatic discovery and verification of query rewrite rules. In Z. G. Ives, A. Bonifati, and A. E. Abbadi, editors, *Proceedings of the 2018 International Conference on Management of Data, SIGMOD Conference 2022, Philadelphia, PA, USA, June 12 - 17, 2022*, pages 94–107. ACM, 2022.

- [52] W. Wu, P. A. Bernstein, A. Raizman, and C. Pavlopoulou. Factor windows: Cost-based query rewriting for optimizing correlated window aggregates. In *38th IEEE International Conference on Data Engineering, ICDE 2022, Kuala Lumpur, Malaysia, May 9-12, 2022*, pages 2722–2734. IEEE, 2022.
- [53] S. Xue, C. Jiang, W. Shi, F. Cheng, K. Chen, H. Yang, Z. Zhang, J. He, H. Zhang, G. Wei, W. Zhao, F. Zhou, D. Qi, H. Yi, S. Liu, and F. Chen. DB-GPT: empowering database interactions with private large language models. *CoRR*, abs/2312.17449, 2023.
- [54] F. Yu, W. Hou, C. Luo, D. Che, and M. Zhu. CS2: a new database synopsis for query estimation. In K. A. Ross, D. Srivastava, and D. Papadias, editors, *Proceedings of the ACM SIGMOD International Conference on Management of Data, SIGMOD 2013, New York, NY, USA, June 22-27, 2013*, pages 469–480. ACM, 2013.
- [55] S. Zhang, Z. Chen, Y. Shen, M. Ding, J. B. Tenenbaum, and C. Gan. Planning with large language models for code generation. In *The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023*. OpenReview.net, 2023.
- [56] D. Zheng, M. Lapata, and J. Z. Pan. Archer: A human-labeled text-to-sql dataset with arithmetic, commonsense and hypothetical reasoning. In Y. Graham and M. Purver, editors, *Proceedings of the 18th Conference of the European Chapter of the Association for Computational Linguistics, EACL 2024 - Volume 1: Long Papers, St. Julian’s, Malta, March 17-22, 2024*, pages 94–111. Association for Computational Linguistics, 2024.
- [57] X. Zhou, G. Li, J. Wu, J. Liu, Z. Sun, and X. Zhang. A learned query rewrite system. *Proc. VLDB Endow.*, 16(12):4110–4113, 2023.
- [58] X. Zhou, Z. Sun, and G. Li. DB-GPT: large language model meets database. *Data Sci. Eng.*, 9(1):102–111, 2024.

Appendix

A Additional Experiments

A.1 Rewrite Quality (Other Benchmarks)

We also evaluated `LITHE` on the following set of benchmarks.

1. DSB [18]: A TPC-DS variant with complex data distributions and additional query templates featuring many-to-many joins and non-equi-joins.
2. ARCHER [56]: A Text-to-SQL benchmark spanning 10 databases with availability of ground-truth SQL queries.
3. JOB [30]: An optimizer stress-test benchmark featuring queries with large and complex join graphs.
4. StackOverflow [36]: A real-world benchmark with query templates modeling questions and answers from experts. A random instance of each template is taken.

Note that the queries in these workloads (with the exception of DSB) are mostly *fast running* (i.e. completing in less than 10 seconds). Therefore, with the primary goal to test the coverage of `LITHE` over a diverse variety of workloads (and not the deployment angle), we include all the queries in our analysis. The number of CPRs produced and the cost speedup delivered by `LITHE` and `SOTA` over these benchmarks are shown in Table 9.

Table 9: Comparing `LITHE` with `SOTA` on different benchmarks.

| Benchmark | # CPR | | | CSGM | |
|---------------|-------|------|-------|-------|------|
| | LITHE | SOTA | Union | LITHE | SOTA |
| DSB | 9 | 3 | 9 | 7.7 | 1.7 |
| ARCHER | 22 | 19 | 22 | 2.1 | 1.95 |
| JOB | 4 | 2 | 4 | 1.9 | 1.4 |
| StackOverflow | 2 | 1 | 2 | 8.7 | 7.5 |

In case of DSB, `LITHE` produces CPR for 9 queries resulting in a highly productive CSGM of **7.7**. Similar to the case with TPC-DS, `LITHE` performs better than `SOTA` both with respect to its coverage, as well as the cost speedups.

Turning our attention to the other benchmarks (ARCHER, JOB, StackOverflow), the number of CPR queries is smaller due to the predominance of flat SPJ formulations in these benchmarks, which limits the scope for productive rewriting. Nevertheless, `LITHE` continues to achieve better CPR coverage, whereas `SOTA` misses quite a few opportunities. Further, the CSGM of `LITHE` is visibly better than `SOTA`.

A.2 Rewrite Quality (Rule-based vs LLM-based approaches)

Our `SOTA` baseline comprises two theoretical rule-based approaches—Learned Rewrite and LLM-R2—and two LLM-based approaches—GenRewrite and Baseline Prompt. The detailed breakdown of CPR achieved by these methods is presented in Table 10.

Table 10: # CPR of Different SOTA Approaches in Comparison with LITHE

| Tools/Techniques | Approach | Individual # CPR | $\geq 1.1x$ |
|------------------|------------|------------------|-------------|
| LITHE | LLM Based | 26 | 30 |
| Genrewrite | LLM Based | 11 | 12 |
| Baseline | LLM Based | 10 | 10 |
| Learned Rewrite | Rule Based | 1 | 2 |
| LLMR2 | Rule Based | 0 | 0 |

While the theoretical rule-based approaches provide formal correctness guarantees, our TPC-DS experiments revealed that the LLM-based approaches consistently outperformed them. Among the 13 CPRs (i.e., query rewrites with cost speedup $\geq 1.5X$), only one was produced by a theoretical rule-based approach. Moreover, even when we decreased this threshold to as low as $\geq 1.1X$, only one more rewrite became a CPR!

LLM-based approaches, on the other hand, seem more suited for query rewrites, resulting in 12 (GenRewrite) and 30 (LITHE) rewrites with a cost speedup $\geq 1.1X$. This is likely a consequence of the LLM having a “global” vision of the entire query, giving it more leeway in considering rewrites in the query-space. Whereas, theoretical rule-based approaches operate mainly at the operator level within the plan tree, thus being confined to only local optimizations.

A.3 Rewrite quality (Coverage of Basic Prompts)

Figure 11a shows the CPR distribution across the Basic Prompts on the Microbenchmark. Among the 6 rewritten queries, 3 are unique to a single prompt, while the other 3 are produced by multiple prompts. At first glance, one might think that Prompt 3 could be removed from our prompt suite since it does not produce any unique rewrite. However, this is an artifact of the micro-benchmark – on the corresponding picture for the full TPC-DS benchmark (Figure 11b), Prompt 3 does contribute unique rewrites.

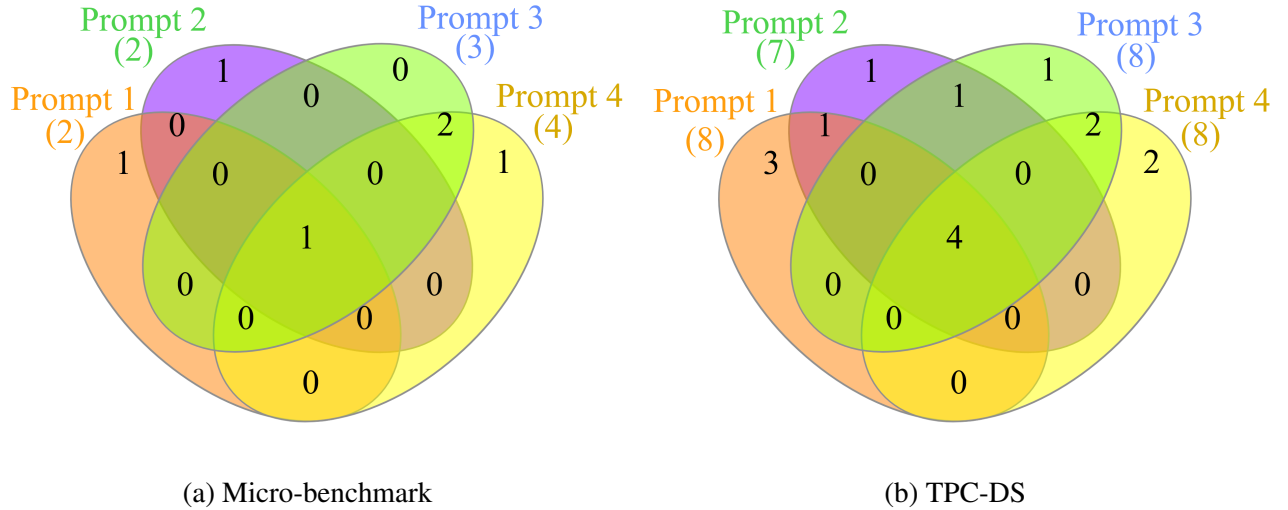


Figure 11: CPR distribution of Basic Prompts

B Overheads Reduction for LITHE

Rule Selection Classifier

Reducing the query rewrite time while still obtaining the same performance would be possible if we could directly use the MCTS-driven rewrite with the appropriate prompt. Towards this end, we build a classifier to pick the most appropriate rewrite rule for a given input query. Specifically, the classifier identifies which, if any, of Rules R1–R6 is appropriate for a given query. If none are appropriate, then it falls back to just the set of basic prompts to identify the best prompt to be given as input to the MCTS module. In addition to reducing the rewrite times, using a classifier can also reduce the financial costs of the rewrite.

We design an LLM based classifier to accomplish this as follows: The LLM is given the rewrite rules discussed so far, and additionally, for each rule, an example demonstrating when the rule can be applied and a counter-example demonstrating when the rule cannot be applied. For the database schema and statistics-based rules, the relevant information is also fed to the classifier so that it can make an informed decision. Then, given an input query, the classifier is tasked with selecting the most appropriate rewrite rule.

Table 11: Impact of Classifier (TPC-DS)

| Metrics | Without Classifier | With Classifier |
|-------------|--------------------|-----------------|
| # CPR | 26 | 23 |
| CSGM | 11.5 | 8.5 |
| Avg. Tokens | 18427 | 16003 |
| Avg. Time | 5 min | 2.4 min |

Table 11 compares the performance of LITHE with and without the classifier. The time overheads do visibly go down by about 52 percent, and the tokens by about 13 percent. However, there is a price to be paid – the CPR is reduced to 23 and the CSGM comes down to 8.5. In our future work, we plan to look into whether a better tradeoff could be achieved between quality and overheads.

As an alternative to the LLM based classifier, we also experimented with a DSPy MIPRO [38] based classifier to identify which, if any, of Rules R1–R6 is appropriate for a given query. However, our initial results indicated that it was weaker than our LLM-based classifier. While the average time taken by this approach reduced to just 10 seconds per query, the CPRs are reduced to 14, and the CSGM and TSGM fall down to 3.3 and 4.5 respectively.

But we hasten to add that this outcome is likely an artifact of our small rule set – looking to the future, where there may be a larger corpus, a data-driven selection method like MIPRO may come to the fore.

Pruning in MCTS

A bottleneck in the MCTS-based exploration is the need to greedily expand a node (during the simulation stage) until an entire rewrite is output. In principle, if we could quickly check for semantic and syntactic correctness at intermediate stages, then unproductive paths could be terminated early. We are currently working on the design and implementation of such checks.

C Query Equivalence Testing via Sampling

We use a sampling-based approach to quickly test equivalence in the rewrite generation stages of the pipeline. The idea here is to execute the queries on several small samples of the database and verify equivalence based on the sample results. However, while this test is a necessary condition for query equivalence, it is not a sufficient condition. That is, there are no false negatives, but there can be false positives. This is because the sampled database may not cover all the predicates present in the query. This can cause two types of problems: First, it is possible for two different queries to return the same result. This can happen when, for example, the entirety of the sampled data satisfies a predicate of one query, while the same predicate is not present in the other. Second, if the underlying sample does not satisfy any of the predicates in either query, then an empty result will be returned by both queries. This again does not imply that the queries are equivalent.

To statistically address the first problem (false positives), we create multiple samples of the database with different seeds, and run the test on all these samples. The goal is to reduce the likelihood of non-equivalent queries returning the same results.

To minimize the occurrence of the second problem (empty results), the following approach is taken:

1. We use *correlated sampling* [54] to sample the database. This technique leverages the join graph of the schema to produce a sample that maintains join integrity between the tables participating in the query.
2. Given a pair of queries to test for equivalence, we adjust the constants in the filter predicates to reduce the chances of an empty result. We make use of the rows in the sampled data for this purpose. For example, say an equality predicate is present in the query and the associated constant is absent in the sampled database. We then replace the query constant with a value already present in the sample. Similarly, the constants for other comparison operators are adjusted based on the ranges of the corresponding columns in the sampled database.

Note that these modifications only change the selectivity of the query, but not its semantics.

3. To cover predicate boundary conditions, we inject query-specific synthetic tuples into the sample—similar to the approach used by the XData mutant-killing tool [42]—in order to reduce the likelihood of non-equivalent queries producing the same results.

D Examples used in Prompts for Rules 1–6

R1: Use CTEs (Common Table Expressions) to avoid repeated computation.

Original Query

```
SELECT emp.employee_name,  
       mgr.manager_name  
FROM   employees emp,  
       managers mgr  
WHERE  emp.manager_id = mgr.manager_id  
       AND emp.employee_id IN (SELECT manager_id  
                                FROM   (SELECT manager_id,  
                                              manager_name  
                                FROM   managers  
                                WHERE  job_id = 'IT_PROG'  
                                AND manager_id > 100))  
       AND mgr.manager_name IN (SELECT manager_name  
                                FROM   (SELECT manager_id,  
                                              manager_name  
                                FROM   managers  
                                WHERE  job_id = 'IT_PROG'  
                                AND manager_id > 100));
```

Rewritten Query

```
WITH cte  
    AS (SELECT manager_id,  
              manager_name  
        FROM   managers  
        WHERE  job_id = 'IT_PROG'  
              AND manager_id > 100)  
SELECT emp.employee_name,  
       mgr.manager_name  
FROM   employees emp,  
       managers mgr  
WHERE  emp.manager_id = mgr.manager_id  
       AND emp.employee_id IN (SELECT manager_id  
                                FROM   it_prog_managers)  
       AND mgr.manager_name IN (SELECT manager_name  
                                FROM   it_prog_managers);
```

R2: When multiple subqueries use the same base table, rewrite to scan the base table only once.

Original Query

```
SELECT (SELECT Avg(salary)
        FROM employees
        WHERE department = 'Sales'
              AND experience_years BETWEEN 1 AND 5
              AND salary BETWEEN 50000 AND 60000) AS Sales_Avg,
       (SELECT Avg(salary)
        FROM employees
        WHERE department = 'HR'
              AND experience_years BETWEEN 5 AND 10
              AND salary BETWEEN 80000 AND 90000) AS HR_Avg;
```

Rewritten Query

```
SELECT avg(
        CASE
            WHEN department = 'Sales' THEN salary) AS sales_avg,
       avg(
        CASE
            WHEN department = 'HR' THEN salary) AS hr_avg
FROM employees
WHERE (
        department = 'Sales'
        AND experience_years BETWEEN 1 AND 5
        AND salary BETWEEN 50000 AND 60000)
OR (
        department = 'HR'
        AND experience_years BETWEEN 5 AND 10
        AND salary BETWEEN 80000 AND 90000);
```

R3: Eliminate overlapping subqueries.

Original Query

```
SELECT c.*
FROM   customer c
WHERE  c.address_id IN (SELECT a.address_id
                        FROM   address)
      AND c.address_id IN (SELECT a.address_id
                        FROM   address
                        WHERE  a.pin_code = '560012');
```

Rewritten Query

```
SELECT c.*
FROM   customer c
WHERE  c.address_id IN (SELECT a.address_id
                        FROM   address
                        WHERE  a.pin_code = '560012');
```

R4: Remove unnecessary joins between a primary key and a foreign key.

Schema

```
CREATE TABLE products
(
    p_product_id INTEGER NOT NULL,
    PRIMARY KEY (p_product_id)
);

CREATE TABLE fact_sales
(
    f_sales_id    INTEGER NOT NULL,
    f_units_sold  INTEGER NOT NULL,
    f_product_id  INTEGER NOT NULL,
    PRIMARY KEY (f_sales_id),
    FOREIGN KEY (f_product_id) REFERENCES products(p_product_id)
);
```

Original Query

```
SELECT p_product_id,
       f_units_sold
FROM   fact_sales,
       products
WHERE  f_product_id = p_product_id;
```

Rewritten Query

```
SELECT f_product_id,
       f_units_sold
FROM   fact_sales;
```

R5: Choose EXIST or IN based on subquery selectivity.

Original Query

```
Select item.id, item.code, item.price
from item
where item.sourceid in (
    Select element.sourceid
    from element
    where element.zip > 1100
)
order by item.id;
```

Statistics

Selectivity of different predicates is given below :
(1) source_id > 1100 on table element :: 0.7385

Rewritten Query

```
Select item.id, item.code, item.price
from item
where exists(select 1
    from element
    where item.sourceid = element.sourceid
    and element.sourceid > 1100
)
order by item.id;
```

R6: Pre-filter tables that are involved in self-joins and have low selectivities on their filter and/or join predicates. Remove any redundant filters from the main query. Do not create explicit join statements.

Original Query

```
with total_price_cte as (  
    select item.id, colour.colorcode, sum(item.price) total_price  
    from item, color  
    where item.colorcode = colour.colorcode  
    group by item.id, colour.colorcode  
)  
select t_sec.id, t_first.colorcode  
from total_price_cte t_first, total_price_cte t_sec  
where t_sec.id = t_first.id  
    and t_first.colorcode = 'R'  
    and t_sec.colorcode = 'R'  
    and t_first.total_price > 0  
order by t_sec.id  
limit 100;
```

Statistics

Selectivity of different predicates is given below :
(1) colour.colorcode = 'R' :: 0.01

Rewritten Query

```
with total_price_cte as (  
    select item.id, colour.colorcode, sum(item.price) total_price  
    from item, color  
    where item.colorcode = colour.colorcode  
    group by item.id, colour.colorcode  
) ,  
filtered_total_price_cte as (  
    select * from total_price_cte  
    where colorcode = 'R'  
)  
select t_sec.id, t_first.colorcode  
from filtered_total_price_cte t_first, filtered_total_price_cte t_sec  
where t_sec.id = t_first.id  
    and t_first.total_price > 0  
order by t_sec.id  
limit 100;
```