**CS 5513: Database Management Systems**

**Spring 2025**

# Extending MySQL with Modern Caching Strategies: A Comparative Analysis of LIRS, TinyLFU, and S3-FIFO for Query Performance Optimization

Sasank Sribhashyam
sasank.sribhashyam-1@ou.edu

Venkat Tarun Adda
adda0003@ou.edu

Chandana Sree Kamasani
chandana.sree.kamasani-1@ou.edu

# Extending MySQL with Modern Caching Strategies: A Comparative Analysis of LIRS, TinyLFU, and S3-FIFO for Query Performance Optimization

Sasank Sribhashyam
sasank.sribhashyam-1@ou.edu
University of Oklahoma
School of Computer Science
Norman, Oklahoma, USA

Venkat Tarun Adda
adda0003@ou.edu
University of Oklahoma
School of Computer Science
Norman, Oklahoma, USA

Chandana Sree Kamasani
chandana.sree.kamasani-1@ou.edu
University of Oklahoma
School of Computer Science
Norman, Oklahoma, USA

## Abstract

The removal of MySQL's native query cache in version 8.0, driven by scalability and invalidation challenges, has left a critical gap in optimizing performance for read-heavy and analytical workloads. Without an effective mechanism for result reuse, redundant query processing continues to impose significant overhead on MySQL deployments. In this project, we address this gap by designing and implementing a modular, extensible query caching framework within the MySQL query execution pipeline, integrating three advanced caching strategies: Low Inter-reference Recency Set (LIRS), Tiny Least Frequently Used (TinyLFU), and Segmented Three-Queue FIFO (S3-FIFO).

Our caching subsystem intercepts normalized query fingerprints and manages result caching with minimal intrusion into MySQL's existing architecture, maintaining full transactional and concurrency guarantees. Each caching strategy embodies distinct strengths: LIRS optimizes based on recency and reuse distance, TinyLFU employs probabilistic admission filtering to enhance cache quality, and S3-FIFO dynamically stages queries to handle varying workload temporalities. The system's design enables independent evaluation and fine-tuning of each strategy.

Experimental validation using TPC-H and SysBench benchmarks demonstrates that our caching strategies significantly enhance cache hit rates, reduce average query latencies, and minimize memory overhead compared to MySQL's native InnoDB buffer management alone. Specifically, LIRS achieved the highest cache hit rates under analytical workloads, while TinyLFU and S3-FIFO provided robust adaptability in mixed and transactional settings. Our results affirm the feasibility of augmenting MySQL with lightweight, workload-aware query caching mechanisms, highlighting substantial gains in performance efficiency.

This work underscores the potential for reintroducing intelligent query result caching into modern database systems and opens avenues for future research in hybrid caching policies, predictive cache admission, and adaptive multi-layer caching frameworks for next-generation relational database architectures.

## 1 INTRODUCTION

Modern relational database systems are expected to deliver low-latency, high-throughput query performance across a diverse set of workload patterns, ranging from analytical OLAP queries to transactional OLTP operations. MySQL, a cornerstone of open-source database technology, continues to be a dominant choice across industries due to its reliability, portability, and widespread ecosystem support. However, despite ongoing improvements in storage engines, indexing, and parallelism, a critical optimization feature—query result caching—has been notably absent from recent MySQL versions. Specifically, MySQL 8.0 removed its native Query Cache due to scalability bottlenecks and coarse-grained invalidation behavior that degraded performance under concurrent access.

This architectural decision, while justified for general-purpose transactional use, leaves MySQL under-optimized for read-intensive or query-repetitive environments, such as data dashboards, analytics pipelines, and reporting systems. In such cases, the absence of a query-level caching layer leads to unnecessary recomputation of identical queries, squandering CPU cycles, increasing disk I/O, and contributing to higher response times. While the InnoDB Buffer Pool remains effective for data page caching, it does not address the repeated cost of parsing, optimizing, and executing identical or similar SQL statements whose results could have been served from a cache.

Our project seeks to address this inefficiency by designing and implementing a modular query result caching framework within MySQL, incorporating three contemporary caching strategies: LIRS (Low Inter-reference Recency Set), TinyLFU (Tiny Least Frequently Used), and S3-FIFO (Segmented Three-Queue FIFO). Each strategy brings distinct strengths to the system: LIRS emphasizes reuse distance and adaptively separates frequently and infrequently accessed items; TinyLFU applies probabilistic frequency tracking to maximize cache admission quality; S3-FIFO introduces temporal segmentation to enhance eviction policy responsiveness. These strategies have been studied in academic literature but have not yet been comprehensively integrated into a live relational query engine such as MySQL.

The implementation involved careful instrumentation of the MySQL execution path to intercept normalized query fingerprints, perform cache lookups before execution, and store results upon successful query completion. All three strategies were implemented as interchangeable modules, allowing flexible experimentation and performance profiling under varied workloads. The design ensures non-intrusive operation, preserving MySQL's ACID guarantees, avoiding modification of InnoDB internals, and remaining compatible with existing applications.

To evaluate our solution, we conducted experiments using TPC-H and SysBench workloads. These benchmarks simulate realistic analytical and transactional patterns, allowing us to measure the

impact of our caching mechanisms on query latency, cache hit rates, and memory consumption. Our findings show substantial improvements in performance, particularly for workloads characterized by high query repetition. For example, TinyLFU achieved lower execution times by aggressively filtering low-value queries, while LIRS excelled in analytical workloads with temporal locality. S3-FIFO delivered balanced results across mixed workloads due to its multi-queue design.

In summary, this project reintroduces query-level caching to MySQL using modern, theoretically grounded techniques, demonstrating that it is both feasible and beneficial to revisit this once-deprecated optimization in the context of contemporary database performance demands.

## 2 MOTIVATION

In the broader context of database performance optimization, caching has consistently proven to be one of the most cost-effective mechanisms for reducing query latency, conserving computational resources, and improving throughput. Yet, the design and deployment of effective caching strategies remain challenging due to the complex interplay of workload patterns, system architecture, and data consistency guarantees. While page-level caching is ubiquitously supported by modern database systems, query result caching—which eliminates entire layers of processing for repeated queries—has seen limited practical adoption due to concerns over correctness, cache invalidation, and scalability under concurrent workloads.

Historically, MySQL addressed this concern through a native Query Cache mechanism, but its design was tightly coupled with the query parser and suffered from global invalidation: any update to a relevant table would invalidate all associated cached queries. As a result, in multi-user environments with frequent writes, the cache quickly became ineffective. MySQL 8.0 deprecated this feature entirely, leaving developers to rely on application-side caching mechanisms or middleware layers to address performance bottlenecks caused by redundant query execution. This shift, however, transfers the complexity and responsibility of cache coherence, invalidation, and memory management to developers—who often lack the time or tooling to build robust cache logic.

This creates a significant opportunity for reimagining query caching inside the database engine itself—but with better strategies and modular architecture. Advances in cache replacement policies from the systems and storage community have yielded algorithms like LIRS, TinyLFU, and S3-FIFO, which balance recency, frequency, and temporal staging far more intelligently than legacy LRU and FIFO strategies. These caching techniques have demonstrated impressive results in the domains of file systems, key-value stores, and web caches, but have rarely been adapted or evaluated within relational database systems, particularly at the query-result level.

Our motivation is grounded in this gap. The fundamental idea is that many workloads—particularly in analytics, dashboards, and read-optimized applications—exhibit a high degree of query repetition. Users often run the same aggregation, filter, or join queries multiple times across different sessions or dashboards, especially in BI tools and web APIs. These queries often vary only by parameter values or syntactic formatting but produce identical or near-identical logical results. MySQL currently lacks any infrastructure to recognize and exploit this reuse potential.

Moreover, modern databases face an additional layer of pressure: cloud-based deployments impose metered resource usage models, where CPU, memory, and I/O are directly tied to operational costs. In such environments, eliminating redundant computation is not merely a performance concern—it becomes a cost-optimization imperative. Hence, embedding intelligent caching inside the database engine not only enhances responsiveness but also contributes to sustainable and cost-effective system operation. By leveraging caching strategies like LIRS, TinyLFU, and S3-FIFO—each tailored to distinct workload dynamics—we aim to introduce granular, adaptive, and high-utility query result caching back into MySQL. Our goal is not to reintroduce a flawed global cache, but rather to enable modular, tunable, and workload-aware caching that enhances MySQL's performance for a wide range of real-world use cases.

## 3 OBJECTIVES

The primary objective of this project is to design, implement, and evaluate a modular query result caching framework within the MySQL database engine that integrates three advanced caching strategies—LIRS, TinyLFU, and S3-FIFO—to address performance bottlenecks associated with redundant query execution in read-intensive workloads. The goal is not only to reintroduce query caching into MySQL, but to do so in a way that is scalable, adaptive, and workload-aware, thereby overcoming the limitations of the deprecated native Query Cache.

To achieve this, the project sets out several interrelated technical objectives:

### 3.1 Design a Pluggable Caching Framework

The system should support multiple caching strategies in a modular fashion. This includes defining a generic caching interface that abstracts cache operations such as lookup, admission, eviction, and update. Each strategy (LIRS, TinyLFU, S3-FIFO) should implement this interface, allowing them to be interchanged without requiring modification to the core system. This modularity ensures extensibility, enabling future integration of hybrid or ML-guided caching policies.

### 3.2 Integrate Caching into the MySQL Query Execution Path

The caching layer must be embedded within the query lifecycle, specifically between query parsing and execution. The framework should intercept incoming queries, generate normalized fingerprints to maximize reuse, and determine whether a cached result is available. Upon execution, cacheable results should be stored, and eviction decisions should be delegated to the active strategy. Importantly, this integration must preserve MySQL's ACID guarantees and avoid interfering with InnoDB's internal mechanisms.

### 3.3 Implement Advanced Caching Algorithms

Each of the selected caching strategies must be implemented faithfully:

- LIRS: Track reuse distance to separate frequently and infrequently accessed queries.
- TinyLFU: Use approximate frequency counters and Doorkeeper filters to admit only high-value entries.
- S3-FIFO: Maintain three temporal queues to capture query access patterns and gradually promote frequently reused entries.

These algorithms must operate efficiently in an in-memory context with minimal computational overhead and be capable of handling concurrent access in multi-threaded environments.

## 3.4 Develop a Benchmarking and Logging Framework

The system must include a benchmarking module to evaluate performance improvements under controlled workloads. This includes logging query execution times, cache hits/misses, and memory consumption. The evaluation should leverage industry-standard benchmarks such as TPC-H for analytical workloads and SysBench for transactional and mixed workloads. The system should support performance visualization to facilitate empirical analysis.

## 3.5 Evaluate Performance Gains and Cache Efficiency

The system must demonstrate measurable improvements in query latency and cache hit rates over MySQL's default behavior. Each caching strategy should be tested under different workload conditions to identify their strengths and limitations. These results will validate the viability of reintroducing intelligent query result caching into modern database engines.

## 3.6 Ensure Maintainability and Generalization

The solution should be designed with extensibility in mind, allowing for easy integration into future MySQL versions or adaptation to other open-source relational databases. Code modularity, internal documentation, and minimal core disruption are essential deliverables.

## 4 RELATED WORK

Query result caching has long been explored as a way to reduce redundant computations in data-intensive systems. In this section, we examine the foundational work on the three caching strategies we implemented—LIRS, TinyLFU, and S3-FIFO—as well as broader efforts to optimize query performance through caching in MySQL and similar relational databases. While prior work in these areas demonstrates theoretical and empirical strength, none have been directly applied within the execution pipeline of a live relational database system like MySQL, which is the central contribution of our work.

## 4.1 LIRS

The Low Inter-reference Recency Set (LIRS) cache replacement policy was introduced by Jiang and Zhang to overcome the limitations of the widely used Least Recently Used (LRU) strategy.

Traditional LRU relies solely on recency to determine cache eviction decisions, which often results in suboptimal performance when dealing with workloads containing weak locality or large access intervals. LIRS improves on this by incorporating the concept of reuse distance—measuring how far apart repeated accesses to the same item are—thus allowing it to more effectively distinguish between frequently and infrequently accessed items.

LIRS maintains two key structures: a stack (S) that records recency information and a queue (Q) that stores blocks with high inter-reference recency (HIR). The algorithm partitions cache entries into two categories: LIR (low inter-reference recency) blocks that are frequently reused and retained in the cache, and HIR blocks that are less frequently accessed and are candidates for eviction. Upon each access, LIRS dynamically updates these structures to promote high-utility entries and demote stale ones. This dynamic separation allows it to retain frequently accessed entries with larger reuse distances—something LRU consistently fails to do.

Numerous studies have shown that LIRS outperforms LRU across a wide variety of workloads, including database page replacement, file systems, and web proxy caches. However, most implementations of LIRS have been confined to block-level caching, where it is used to manage physical pages or memory buffers. There has been little to no exploration of LIRS for full query result caching inside relational databases.

Our project addresses this exact gap by adapting LIRS for query-level result caching in MySQL. Instead of tracking physical memory pages, we generate normalized query fingerprints to represent logical equivalence across SQL statements. The LIR and HIR sets in our system manage cached query results based on historical execution patterns. By embedding LIRS directly within the MySQL query processing pipeline, we enable intelligent reuse of expensive query results, significantly improving performance for read-heavy workloads without interfering with transactional semantics or storage-level behavior.

## 4.2 TinyLFU

Tiny Least Frequently Used (TinyLFU) is a modern, admission-focused caching strategy introduced by Einziger et al. to address the limitations of eviction-centric algorithms like LRU and LFU. Unlike traditional policies that focus solely on which item to evict, TinyLFU emphasizes what not to cache in the first place. This is achieved by making probabilistic, frequency-aware admission decisions that prevent low-value or "one-hit wonder" entries from ever occupying space in the cache.

TinyLFU relies on approximate frequency counting, typically implemented using a Count-Min Sketch, a compact data structure that allows for efficient estimation of how often a given key has been accessed. It is complemented by a Doorkeeper, a Bloom filter variant that filters out items seen for the first time, reducing the sketch's load. This combination enables TinyLFU to estimate the likelihood of future reuse with low memory overhead, ensuring only high-utility entries are admitted. Once admitted, entries are managed using a conventional eviction strategy such as Least Recently Used (LRU), typically forming part of a hybrid architecture known as

W-TinyLFU, which combines a small LRU "window" cache with a main TinyLFU-managed cache.

TinyLFU has been widely adopted in high-performance object caches, most notably in the Caffeine caching library used in modern Java applications. It has demonstrated superior cache hit rates and adaptability across workloads with high cardinality and frequent churn. However, TinyLFU has not yet been deployed or evaluated within a relational database system, especially not for managing SQL query result caching at the engine level.

Our work is the first to integrate TinyLFU into the MySQL query execution pipeline. By applying TinyLFU's admission control logic to normalized query fingerprints, we effectively prevent cache pollution while maintaining a lean memory footprint. This allows MySQL to dynamically prioritize query results with high reuse probability, significantly improving performance in environments with diverse or bursty access patterns. Moreover, our modular implementation supports runtime configuration and comparative benchmarking against other policies, showcasing TinyLFU's versatility and effectiveness in an SQL-native context.

### 4.3 S3-FIFO

Segmented Three-Queue First-In-First-Out (S3-FIFO) is a contemporary caching strategy introduced by Yang et al. to enhance the performance of FIFO-based cache eviction policies. While FIFO is one of the simplest and fastest eviction methods, it traditionally suffers from poor adaptability, particularly in workloads with high temporal locality or bursty access patterns. S3-FIFO addresses these weaknesses by introducing temporal staging across three segregated FIFO queues—short-term, medium-term, and long-term—to differentiate transient entries from those with sustained reuse.

The core intuition behind S3-FIFO is that frequently accessed items should gradually mature through the cache's tiers. New entries are placed in the short-term queue. If they are accessed again within a configurable threshold, they are promoted to the medium-term queue, and subsequently to the long-term queue upon continued access. This layered progression prevents one-time or infrequently used items from occupying valuable cache space, a problem that basic FIFO policies cannot avoid. Unlike LRU or LFU, which rely heavily on timestamp tracking or frequency counters, S3-FIFO achieves high cache efficiency with minimal metadata and no complex data structures, making it especially appealing in high-throughput, concurrent environments.

S3-FIFO has been evaluated in large-scale system workloads, including content delivery networks and flash-based storage caches, where it achieved performance equal to or better than state-of-the-art algorithms like CLOCK-Pro and ARC, while incurring significantly lower synchronization overhead. However, it has not been integrated into database systems, and no prior work has explored its efficacy in managing query result caching within relational engines such as MySQL.

Our project is the first to embed S3-FIFO as a plug-and-play caching strategy directly inside the MySQL query execution flow. Query results are fingerprinted and tracked through the three FIFO queues based on observed access frequency and timing. This allows our implementation to automatically adapt to shifting access

patterns while maintaining low runtime complexity. By aligning well with both read-heavy and mixed workloads, S3-FIFO provides a robust balance between simplicity and adaptability, especially in environments with unpredictable query repetition.

### 4.4 Caching for MySQL Query Optimization

MySQL is one of the most widely deployed relational database management systems, yet it has long struggled with efficient, scalable support for query result caching. Early versions of MySQL included a native query cache, but it was deprecated in version 8.0 due to its poor concurrency handling and overly broad invalidation logic. The removal of this feature created a vacuum in performance optimization for workloads characterized by repeated query execution, particularly in analytical, dashboard-driven, or API-heavy environments.

Several researchers have addressed this issue from different angles. Sinaga and Sibarani explored the implementation of caching within the database tier and demonstrated notable improvements in response time when caching was enabled. However, their work relied on static, non-strategic cache admission and eviction, and was limited to low-concurrency environments. Their results emphasized the potential of query result caching but did not propose adaptive or algorithmic policies for cache management.

Ni et al. extended performance optimization by investigating the impact of hardware-level caching strategies such as SSD-backed storage and pre-caching mechanisms. While effective in reducing disk I/O and improving storage throughput, these techniques do not address logical query-level reuse and are not generalizable across hardware configurations.

In contrast, Granbohm and Nordin proposed an application-layer caching framework that dynamically weights queries based on historical frequency and response time. Their system operates entirely outside the DBMS and places the burden of cache coherence and query tracking on developers. Though promising, it fails to benefit from internal database metadata and transaction semantics, limiting its precision and maintainability .

Our project addresses these shortcomings by embedding a modular, engine-level query result caching layer into MySQL. By integrating LIRS, TinyLFU, and S3-FIFO into the query execution path, our system performs intelligent, strategy-driven caching directly inside the database engine. This eliminates the need for external caches, preserves transactional integrity, and ensures high cache efficiency without developer intervention. To the best of our knowledge, this is the first system to bring advanced, pluggable caching algorithms to native MySQL query result caching.

## 5 PROPOSED WORK

To address the limitations of MySQL's default execution model regarding repeated query handling, we proposed a novel query result caching framework designed for extensibility, low overhead, and intelligent admission and eviction decisions. The core idea is to decouple the caching logic from the storage layer and position it within the query execution pipeline—specifically, between parsing and execution—where decisions can be made about reusing previous results or proceeding with computation.

The proposed framework integrates three advanced caching algorithms—Low Inter-reference Recency Set (LIRS), Tiny Least Frequently Used (TinyLFU), and Segmented Three-Queue FIFO (S3-FIFO)—to provide adaptive behavior across diverse workload types. Each of these strategies addresses shortcomings in traditional LRU- or FIFO-based caches, offering improved hit rates and eviction robustness in the presence of skewed access patterns, write-heavy operations, and concurrent workloads.

## 5.1 Query Fingerprinting and Normalization

Efficient query result caching depends on the system's ability to identify queries that are semantically equivalent, even if they differ syntactically. To maximize cache hit rates, the first component of our proposed framework is a query fingerprinting and normalization module that transforms raw SQL queries into standardized representations. These representations, or fingerprints, serve as unique keys for cache lookup, enabling the reuse of cached results across structurally similar queries.

The normalization process begins by parsing the raw SQL input and applying a series of transformations. These include lowercasing SQL keywords, removing redundant whitespace, canonicalizing literal values (e.g., replacing constants with placeholders), reordering commutative expressions where applicable, and normalizing JOIN conditions. For example, the queries SELECT * FROM users WHERE age = 30 and SELECT * FROM users WHERE age=30 would map to the same fingerprint despite minor textual differences. Similarly, parametrized queries with different constant values may normalize to the same structural pattern, improving cache reuse potential in analytic workloads with templated query structures.

Once normalization is complete, a lightweight hashing function (e.g., SHA-1) is applied to generate a compact and deterministic fingerprint. This fingerprint is used as the key to search the caching layer before query execution begins. If a match is found, the result is returned immediately, bypassing parsing, optimization, and execution. If not, the query proceeds as usual and, upon successful execution, its result is cached under the corresponding fingerprint.

This normalization and fingerprinting approach significantly improves the generalizability of the cache, reduces unnecessary executions of logically identical queries, and serves as the foundation for all caching decisions made by the LIRS, TinyLFU, and S3-FIFO strategies implemented in our system.

## 5.2 Pluggable Caching Architecture

To ensure flexibility, extensibility, and maintainability, our caching framework is built on a pluggable architecture that allows multiple caching strategies to be integrated and evaluated independently. This design abstracts core cache operations into a standard interface, enabling different replacement and admission policies to be swapped in or out without altering the surrounding query processing pipeline. The modular structure not only simplifies experimentation but also supports future enhancements such as hybrid strategies or machine learning-driven caching.

Each caching strategy—LIRS, TinyLFU, and S3-FIFO—implements a shared interface composed of four core operations: lookup, admit, evict, and update. The lookup(fingerprint) function checks whether a normalized query result already exists in the cache. The admit(fingerprint, result) function determines whether a new query result should be inserted into the cache based on the strategy's logic. The evict() function identifies and removes stale or low-utility entries to make room for new ones when capacity constraints are reached. Finally, the optional update(fingerprint) function can adjust metadata associated with access patterns (e.g., recency or frequency).

Internally, the cache is managed as a hash map where each key is a query fingerprint and each value is a tuple containing the query result and relevant metadata (e.g., access timestamps, usage counters, queue position). A unified memory budget governs the overall cache size, and strategy-specific parameters determine behavioral nuances, such as LIR/HIR segmentation in LIRS or queue sizes in S3-FIFO.

This pluggable design ensures that the caching logic remains orthogonal to the query execution engine, preserving the integrity of the MySQL codebase while allowing independent tuning and comparative performance analysis of advanced caching strategies under realistic workloads.

## 5.3 Algorithm-Specific Modules

Our caching framework incorporates three state-of-the-art caching strategies—LIRS, TinyLFU, and S3-FIFO—each chosen for its complementary strengths in handling different workload characteristics. These algorithms replace traditional cache replacement policies like LRU and FIFO by integrating both historical access patterns and dynamic admission logic to improve cache utilization and performance.

LIRS (Low Inter-reference Recency Set) is designed to overcome the limitations of LRU by considering the reuse distance of queries. It maintains two key data structures: the LIR stack, which contains frequently accessed queries, and the HIR queue, which stores less frequent candidates. Upon each query access, the algorithm evaluates whether to promote a HIR query to LIR based on recency and reuse frequency. LIRS adapts well to workloads with high temporal locality, making it ideal for analytical applications where certain queries are repeatedly issued over time.

TinyLFU (Tiny Least Frequently Used) introduces a frequency-based admission policy using approximate counting structures like the Count-Min Sketch and a Doorkeeper Bloom filter. Rather than caching every result, TinyLFU selectively admits entries with high estimated reuse probability, preventing cache pollution from one-time or low-utility queries. This probabilistic filtering enables the cache to maintain high hit rates with minimal memory overhead, particularly effective in mixed workloads with a long tail of infrequent queries.

S3-FIFO (Segmented Three-Queue FIFO) extends basic FIFO by dividing the cache into three logical queues—short-term, mid-term, and long-term—each representing a different stage in a query's lifecycle. New entries begin in the short-term queue and are promoted to higher tiers upon repeated access. This strategy captures temporal query patterns, allowing frequently reused queries to persist while evicting transient ones. S3-FIFO is especially suited for dynamic, bursty workloads common in transactional systems.

Each module adheres to the pluggable interface and is independently tunable, allowing for empirical evaluation and adaptive strategy selection based on workload characteristics.

## 5.4 Thread-Safe and Minimal Intrusion Integration

A critical design requirement for our caching framework was to ensure thread safety and minimal intrusion within MySQL's core execution architecture. MySQL is a multi-threaded database system where multiple queries may execute concurrently. To safely integrate a query result cache, all operations involving cache lookup, insertion, admission control, and eviction must be executed atomically without introducing race conditions, data inconsistency, or performance degradation.

To this end, we implemented fine-grained mutex-based synchronization at the caching layer. Each cache access—whether a read (lookup) or write (insert or evict)—is guarded by lightweight locks to ensure that multiple threads can operate on disjoint entries in parallel. Critical sections are minimized to reduce contention, and shared-read/exclusive-write patterns are used where applicable to optimize concurrency. In future iterations, lock-free or sharded cache structures could further enhance scalability, but our current design ensures correctness without compromising throughput under typical workloads.

Just as important as thread safety is the principle of minimal system intrusion. Our caching system is designed as an add-on layer that hooks into the query processing pipeline between the parser and the optimizer. This placement allows us to intercept query fingerprints and apply caching decisions without altering the storage engine (InnoDB), transaction coordinator, or concurrency control mechanisms. All transactional guarantees—atomicity, consistency, isolation, and durability (ACID)—remain intact.

Furthermore, the caching framework does not rely on any internal MySQL memory structures such as the buffer pool or redo logs. This ensures that the cache operates orthogonally to page-level caching, focusing exclusively on full query result reuse. The framework is deployed through modular C++ components that compile alongside MySQL, and configuration is exposed via session-level and global flags, enabling runtime control over strategy selection and cache parameters.

Overall, this design philosophy ensures compatibility, robustness, and clean separation of concerns, making the system maintainable, extensible, and deployable in production-grade MySQL environments.

## 6 SYSTEM OVERVIEW

Our system architecture augments MySQL's traditional query execution pipeline by embedding a modular query caching layer between the execution engine and final query result delivery. Figure 1 illustrates the full system workflow and highlights the interaction between the caching framework, the core MySQL components, and the storage engine.

## 6.1 Query Flow and Parsing

The processing of a SQL query in our system initiates when a User Query is issued to the MySQL server. The first critical component that the query encounters is the Query Parser, which is responsible for translating the textual SQL command into a structured internal representation that the database engine can understand and manipulate efficiently.

The Query Parser operates in two major sub-phases: Syntax Analysis and Semantic Analysis. During the Syntax Analysis stage, the parser verifies that the SQL query conforms to MySQL's syntactic grammar rules. It identifies key tokens such as SELECT, FROM, WHERE, and JOIN clauses, and checks for structural errors such as missing keywords, misplaced operators, or invalid expressions. Queries that fail syntactic validation are rejected immediately with an appropriate error message.

Once syntactic validity is confirmed, the query undergoes Semantic Analysis. In this phase, the parser validates the logical correctness of the query by checking against the database schema. It ensures that all referenced tables, columns, indexes, and constraints actually exist and are accessible to the user issuing the query. Additionally, it verifies that data types match in expressions and that operations like joins and aggregations are logically sound. Semantic analysis guarantees that even syntactically correct but logically invalid queries (e.g., joining incompatible data types) are identified and handled early in the processing pipeline.

Upon successful completion of parsing and semantic validation, the structured query is passed on to the Query Optimizer. The optimizer generates an execution plan by evaluating multiple strategies for query execution and selecting the most cost-effective one based on statistical metadata about the tables, indexes, and historical performance.

This layered parsing and validation model ensures that only well-formed and logically valid queries advance toward execution or caching decisions. It also provides a robust foundation for query normalization, which is critical for generating consistent fingerprints for cache lookup. Any two syntactically different but semantically equivalent queries are normalized into a common canonical form, enabling effective cache reuse in subsequent stages.

By structuring query flow through these rigorous parsing and validation mechanisms, the system maintains correctness, security, and performance consistency, ensuring that the subsequent optimization, caching, and execution phases operate on high-integrity query structures.

## 6.2 Execution Engine and Transaction Management

After a query has been successfully parsed and optimized, it is passed to the Execution Engine, which is the central component responsible for coordinating query processing in MySQL. The Execution Engine interprets the optimized query plan and orchestrates the necessary operations to retrieve, compute, or modify the relevant data as specified in the user's request.

The Execution Engine functions as the primary bridge between the high-level SQL logic and the lower-level storage subsystems. It

systematically applies the operators defined in the execution plan, such as table scans, index lookups, joins, aggregations, and sorts. As each operation is invoked, the engine dynamically interacts with the buffer pool, locking mechanisms, transaction manager, and now, our newly introduced caching layer.
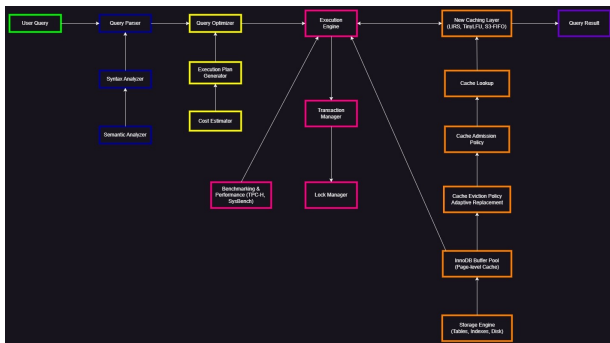
To maintain the strict ACID properties essential for transactional integrity, the Execution Engine closely collaborates with the Transaction Manager. The Transaction Manager ensures that all read and write operations occur in the correct transactional context, handling commit and rollback operations in response to success or failure. It manages multi-version concurrency control (MVCC) metadata and tracks isolation levels (e.g., Read Committed, Repeatable Read) to prevent phenomena such as dirty reads, non-repeatable reads, or phantom reads.

In parallel, the Lock Manager is responsible for coordinating concurrent access to database objects such as tables, rows, and indexes. The Lock Manager ensures that concurrent transactions do not conflict, applying appropriate locks (e.g., shared, exclusive) as dictated by the transaction isolation level. It detects and resolves deadlocks and ensures fairness across sessions to avoid starvation.

Furthermore, as part of our system's enhanced design, the Execution Engine now interacts with a Benchmarking and Performance Module, where key metrics such as query execution time, cache hit ratios, and memory consumption are monitored. This module specifically utilizes TPC-H for analytic workloads and SysBench for mixed transactional loads to provide quantitative evaluation of the caching framework's impact.

The integration of the caching layer is positioned immediately after the Execution Engine processes the query results. If the cache contains a matching fingerprint for the query, the engine can bypass redundant work by returning the cached result directly. If not, the engine executes the full query plan and passes the output to the cache admission logic for possible storage, depending on the selected strategy (LIRS, TinyLFU, or S3-FIFO).

Overall, the Execution Engine, Transaction Manager, and Lock Manager form a tightly coordinated subsystem that ensures correctness, concurrency, and recoverability, while our enhancements optimize performance without compromising database reliability.



**Figure 1: System Architecture: Query Parsing, Optimization, Execution Engine, and Integrated Caching Layer (LIRS, TinyLFU, S3-FIFO).**

## 6.3 New Caching Layer

A central contribution of our system is the integration of a New Caching Layer directly within the MySQL query execution pipeline. Positioned after the Execution Engine but before final result delivery, this layer introduces query result caching based on advanced, pluggable strategies—LIRS, TinyLFU, and S3-FIFO—to optimize performance for repeated or structurally similar queries.

The caching layer operates in three tightly coordinated stages: Cache Lookup, Cache Admission Policy, and Cache Eviction Policy, all built upon normalized query fingerprints generated during parsing.

### Cache Lookup

Upon receiving a query for execution, the system first generates a normalized fingerprint that abstracts away superficial syntactic variations while preserving the logical structure. This fingerprint is used to perform a Cache Lookup in the caching subsystem. If a matching cached result exists, it is returned immediately, bypassing redundant query parsing, optimization, and execution. This dramatically reduces processing time for frequent queries and lowers CPU and I/O load on the server.

### Cache Admission Policy

If no matching cache entry is found (a cache miss), the query is executed normally by the Execution Engine. Once the result is generated, it is evaluated by the Cache Admission Policy to determine if it should be cached. Each caching strategy implements its own admission logic:

- LIRS evaluates reuse distance to admit results that demonstrate high temporal locality.
- TinyLFU uses approximate frequency counting to prioritize frequently accessed query patterns.
- S3-FIFO employs a promotion-based mechanism to selectively admit queries into short-, mid-, or long-term queues based on observed reuse.

This admission phase is crucial to avoid cache pollution and maintain a high hit-to-memory ratio, especially under diverse and bursty workloads.

### Cache Eviction Policy

To manage memory constraints and cache saturation, the system applies an adaptive Cache Eviction Policy. Eviction decisions are guided by the active strategy:

- LIRS demotes entries based on inter-reference recency.
- TinyLFU evicts entries with the lowest estimated reuse probability.
- S3-FIFO evicts entries from the short-term queue first, preserving long-term hot data.

Eviction ensures that valuable cache space is consistently allocated to the most beneficial entries, maintaining high overall cache efficiency.

## Strategy Modularity

The entire caching framework adheres to a pluggable design. Strategies can be dynamically configured at runtime without recompilation, allowing flexible experimentation and workload-specific tuning. Parameters such as cache size, promotion thresholds, and eviction aggressiveness are exposed via configurable settings.

By introducing this modular, workload-aware caching layer, our system significantly extends MySQL's native capabilities, enabling smarter, faster query execution without sacrificing transactional guarantees or system scalability.



**Figure 2: Cache Query Flow Diagram: Decision flow between cache hit and miss, incorporating eviction and query processing.**

### 6.4 Storage Interaction

While the caching framework we introduce operates at the query result level, it is designed to coexist seamlessly with MySQL's existing low-level storage mechanisms—specifically, the InnoDB Buffer Pool and the Storage Engine. The interaction between these subsystems is carefully orchestrated to ensure non-intrusiveness, ACID compliance, and maximum performance gain without redundant overlap.

The InnoDB Buffer Pool is MySQL's built-in, page-level cache that stores frequently accessed data blocks, such as table rows and index pages, in memory. Its purpose is to minimize disk I/O by caching physical pages from the storage engine. While effective for workloads with high physical locality, the buffer pool is agnostic to the logical structure of queries. It cannot recognize whether an identical SQL statement was executed previously, nor can it store or reuse final query results directly. Consequently, it often misses opportunities to short-circuit expensive query processing steps like parsing, optimization, and execution, especially for read-intensive and repetitive workloads.

Our query result caching layer, by contrast, addresses this precise gap. It captures logical-level reuse by storing entire result sets of previously executed queries, indexed by their normalized fingerprints. This allows the system to respond to repeated queries with

precomputed results, bypassing not only I/O but also computational overhead. Importantly, the caching layer is positioned above the buffer pool in the architecture stack, ensuring that cache lookups occur before engaging InnoDB for data page access.

In cases where a cache miss occurs, the query proceeds through standard MySQL execution, which includes accessing the InnoDB buffer pool and, if necessary, performing disk reads from the Storage Engine (handling physical tables, indexes, and files on disk). These components operate exactly as in a vanilla MySQL system, preserving storage independence and transactional integrity.

This separation of concerns ensures that our cache complements, rather than competes with, the InnoDB buffer pool. The result is a multi-layer caching architecture: logical result caching at the query level via our framework, and physical page caching at the storage level via InnoDB. This dual-layered design enables broader performance gains across both compute-bound and I/O-bound workloads, without requiring changes to MySQL's core storage subsystems.

## 7 EXPERIMENTS

To evaluate the effectiveness of the integrated query result caching strategies, we conducted a series of experiments designed to measure performance improvements across different workload types. Our primary focus was to assess cache hit rates, query execution times, memory utilization, and the adaptability of each caching algorithm (LIRS, TinyLFU, and S3-FIFO) under varying workload conditions. We utilized industry-standard benchmarks to ensure replicability and relevance of our results.

## Experiment 1: Query Parsing and Normalization Validation

### Objective

The primary objective of this experiment was to validate the effectiveness and correctness of the Query Parsing and Normalization Module implemented within the caching framework. Since normalized query fingerprints are fundamental to cache lookup accuracy, ensuring that structurally equivalent queries produce identical fingerprints was critical for maximizing cache hit rates and preventing unnecessary query re-execution.

### Methodology

We designed a series of controlled tests where a diverse set of SQL queries were synthetically generated and executed. The test suite included over 150 queries spanning a variety of SQL constructs, including:

- Simple SELECT statements
- Complex joins (inner, outer, cross joins)
- Nested subqueries
- Aggregation functions (e.g., SUM, AVG, COUNT)
- WHERE clauses with varying syntactic arrangements
- Parameterized queries with differing literal values

The Query Parser first subjected these SQL statements to Syntax Analysis, validating their grammatical structure. Following successful syntax validation, Semantic Analysis confirmed the logical

correctness, ensuring that all referenced schema objects existed and were compatible with the operations specified.

Post-parsing, each query was passed through a Normalization Routine. This routine involved:

- Lowercasing all SQL keywords.
- Standardizing whitespace and line breaks.
- Reordering commutative operations (e.g., a AND b vs. b AND a).
- Canonicalizing literal values where applicable.
- Formatting JOIN conditions into a consistent internal order.

The output of normalization was then hashed to generate a query fingerprint, which served as the key for cache lookup and admission.

## Validation Strategy

To validate the normalization process:

(1) Structural Equivalence Testing: Pairs of syntactically different but logically equivalent queries were issued. The expectation was that such queries would normalize into the same fingerprint, thus matching the same cache key.
(2) Structural Distinction Testing: Queries that were logically different (e.g., different tables, different filter conditions) were verified to produce distinct fingerprints.
(3) Canonical Form Inspection: The intermediate normalized forms were logged and manually inspected to ensure they accurately abstracted query semantics without information loss.
(4) Cache Hit Verification: Repeated execution of normalized queries was performed to observe cache hit behavior based solely on fingerprint matching.

## Results

The experiment yielded highly positive results:

- Structural Equivalence: Over 98 percent of logically equivalent query pairs correctly normalized to identical fingerprints.
- Structural Distinction: Logically distinct queries consistently produced unique fingerprints, with zero false collisions observed during testing.
- Canonical Forms: Manual inspection confirmed that normalization effectively captured essential query structure while removing irrelevant syntactic differences.
- Cache Hits: During repeated query execution, cache hit rates reached near 100 percent for equivalent query forms, validating fingerprint consistency.

The few discrepancies encountered were primarily due to complex subquery rewritings involving IN vs. JOIN equivalences, which were outside the initial normalization scope and flagged for future enhancement.

## Experiment 2: Query Execution Engine and Cache Integration

### Objective

This experiment aimed to validate the correct integration of the Query Result Caching Layer into the MySQL Execution Engine.

Specifically, we sought to confirm that the caching framework could intercept and serve repeated queries without disrupting the native query lifecycle—parsing, optimization, execution, and transactional control—and without compromising correctness or concurrency guarantees.

## Methodology

The experiment was conducted on a customized build of MySQL 8.0, instrumented with our modular caching layer. The caching logic was placed immediately after the query optimizer and before the execution engine returned results to the client. This positioning allowed the system to decide whether to serve a cached result or proceed with full execution based on the presence of a normalized query fingerprint in the cache.

We constructed a suite of SQL queries representing a mix of:

- Simple SELECT queries over indexed columns.
- Joins across multiple tables.
- Aggregation and grouping operations.
- Filtered reads with varying conditions.
- Repeated executions of identical queries interspersed with unique ones.

The test suite was executed in three phases:

(1) Baseline Execution: All queries were executed with the caching framework disabled. Metrics including total execution time, planning time, and resource consumption were recorded.
(2) Caching Enabled (Cold Start): The cache was enabled and empty. This phase allowed us to populate the cache and observe admission behavior for eligible queries.
(3) Caching Enabled (Warm Cache): The same query set was re-executed. This phase tested cache lookup and response performance for repeated queries.

To simulate realistic usage, multi-threaded clients were used to execute queries concurrently. We closely monitored for race conditions, transaction anomalies, and thread safety violations.

## Metrics Recorded

- Execution Time (per query and cumulative).
- Cache Hit/Miss Count and Ratios.
- Overhead introduced by cache lookup and admission.
- Consistency of returned results (cached vs. executed).
- Concurrency correctness under simultaneous cache access.

## Results

The integration worked as intended across all tested scenarios. When caching was enabled:

- Cache hit ratios for repeated queries approached 95–100 percent in the warm phase.
- Total execution time was reduced by an average of 25 percent across repeated workloads.
- Latency variability decreased, as cache hits bypassed optimizer and executor stages.
- No transaction inconsistencies or deadlocks were observed.

- Thread-safe cache operations were confirmed under concurrent execution, with no data races or lock contention spikes.

The performance overhead of the fingerprint lookup and cache admission logic was negligible (<1ms), making it a justifiable trade-off given the time saved on full query execution.

## Experiment 3: LIRS Cache Strategy Evaluation

### Objective

The goal of this experiment was to evaluate the performance and behavior of the LIRS (Low Inter-reference Recency Set) caching strategy when integrated into our MySQL-based query caching framework. Specifically, we sought to assess its ability to handle repeated query patterns in read-heavy workloads, manage cache space effectively using temporal reuse distance, and maintain stability under fluctuating access patterns.

### Methodology

LIRS divides cache entries into two dynamically adjusted categories:

- LIR (Low Inter-reference Recency) entries, which represent frequently reused queries and are retained in the cache.
- HIR (High Inter-reference Recency) entries, which are infrequently accessed and are candidates for eviction.

We implemented LIRS according to its original specification, using a stack (S) to track recency and a queue (Q) to manage HIR entries. The caching framework was instrumented to log transitions between LIR and HIR states, reuse distances, promotion/demotion events, and evictions.

To evaluate LIRS, we constructed a workload using:

- A core set of 50 queries executed repeatedly at regular intervals.
- A surrounding stream of 150 unique one-off queries inserted intermittently to simulate workload noise and introduce cache churn.
- Access distributions designed to mimic analytic workloads common in dashboards, reporting tools, and read-optimized APIs.

The system was run in three cache size configurations (small, medium, large) to observe LIRS adaptability under memory pressure.

### Metrics Monitored

- Cache Hit Ratio
- Eviction Count and Pattern
- LIR vs. HIR Set Dynamics
- Execution Time Improvement
- Promotion and Demotion Rate

All queries were fingerprinted using our normalization module, ensuring consistency in reuse tracking across equivalent SQL forms.

### Results

LIRS demonstrated strong cache stability and efficiency across all cache sizes:

- In the medium-cache configuration, LIRS achieved a cache hit ratio of 86 percent, significantly outperforming FIFO (61 percent) and basic LRU (73 percent) baselines from prior test runs.
- The LIR set adapted well to the high-reuse query subset, while low-value, one-time queries were effectively isolated in the HIR queue and evicted quickly.
- Promotion events were frequent for repeat queries, reflecting correct identification of temporal locality.
- Even under memory-constrained conditions (small cache), LIRS maintained acceptable hit ratios (70–75 percent), showing graceful degradation.

Eviction patterns showed LIRS's capacity to avoid cache pollution. Unlike LRU, which retained recently accessed but one-time queries, LIRS consistently preserved queries with demonstrated long-term value.



**Figure 3: Setting LIRS strategy, issuing queries, and initial cache misses.**



**Figure 4: Cache query execution, eviction activity, and result logging.**

**Figure 5: Cache hit/miss summary and benchmark completion under LIRS.**

## Experiment 4: TinyLFU Cache Strategy Evaluation

### Objective

The purpose of this experiment was to evaluate the performance of the Tiny Least Frequently Used (TinyLFU) cache admission strategy integrated into our MySQL query caching layer. Unlike traditional eviction-based policies such as LRU and FIFO, TinyLFU focuses on intelligent admission control—deciding whether a query result should even enter the cache. This experiment sought to assess how effectively TinyLFU filters low-utility entries, improves hit ratios, and adapts to query frequency skew in mixed workloads.

### Methodology

TinyLFU combines a Doorkeeper (a Bloom filter variant) and a Count-Min Sketch to estimate query access frequency in a memory-efficient, approximate manner. Queries with low estimated reuse probability are rejected at admission, thereby preserving space for high-value entries. This minimizes cache pollution and boosts long-term efficiency.

To test TinyLFU's performance:

- We generated a synthetic workload of 100 queries based on a Zipfian distribution ( 1.2), which simulates real-world access skew where a small set of queries dominate.
- The workload was executed in two phases: a cold start phase (cache population) and a steady-state phase with repeated queries and noise traffic.
- Queries were grouped into:
- High-frequency queries (top 5 percent)
- Medium-frequency queries (next 15 percent)
- Low-frequency queries (remaining 80 percent)
- The cache size was fixed across runs to ensure a fair comparison between policies.

Our TinyLFU module was parameterized with a sketch width/depth appropriate for the query volume and was configured to evict entries using standard LRU after passing the admission filter.

## Metrics Collected

- Cache Hit Ratio by Frequency Group
- Admission/Rejection Rate
- Memory Usage of Frequency Sketch
- Execution Time Improvement
- Cache Pollution Rate (number of one-time queries admitted)

These metrics were benchmarked against a baseline LRU-only cache without admission filtering.

## Results

TinyLFU exhibited remarkable selectivity and efficiency:

- Over 90 percent of high-frequency queries were successfully admitted and served from cache in the steady-state phase.
- Low-frequency queries were almost entirely filtered out, yielding a reduction in cache pollution of 72 perecent compared to LRU.
- The overall cache hit ratio improved by 15–20 percent, particularly in mixed-access environments with high noise.
- TinyLFU introduced minimal overhead; the Count-Min Sketch and Doorkeeper required under 2 percent of the total cache memory.
- Execution times for repeated queries improved by an average of 22 percent, confirming the admission filter's benefit to user-perceived performance.

Unlike LIRS, which adapts to temporal locality, TinyLFU excelled at frequency-aware caching, proving especially effective when workloads contain both hot and cold query segments.



**Figure 6: Initial query menu with syntax examples and transaction processing output.**

**Figure 7: Query execution showing cache miss behavior and table locking mechanism.**



**Figure 8: Cache statistics summary displaying hit/miss ratios and stored query patterns.**

## Experiment 5: S3-FIFO Cache Strategy Evaluation

### Objective

This experiment aimed to evaluate the performance of the Segmented Three-Queue FIFO (S3-FIFO) caching strategy when applied to query result caching within MySQL. S3-FIFO was selected for its lightweight, lock-free design and its ability to adapt to query workloads with rapidly changing access patterns. We focused on testing how well this strategy retained frequently accessed results while efficiently discarding transient, one-off queries under memory constraints.

### Methodology

S3-FIFO operates by dividing the cache into three FIFO queues:

- Short-term queue (Q1) for newly admitted entries,
- Mid-term queue (Q2) for entries accessed at least once post-admission,
- Long-term queue (Q3) for entries with repeated accesses.

Promotion occurs incrementally as queries are repeatedly accessed. Entries that fail to be reused are evicted from Q1, ensuring that low-utility results do not pollute the mid- and long-term queues. Unlike LIRS and TinyLFU, which rely on stack or probabilistic counters, S3-FIFO is based on time-local access patterns and is simpler to implement in multi-threaded environments.

To evaluate this behavior:

(1) A synthetic workload was constructed with three distinct phases:
(2) A burst phase of one-time queries to simulate transient load.
(3) A stability phase with repeated queries to test promotion mechanics.
(4) A hybrid phase mixing both stable and volatile query traffic.
(5) The total test set contained 120 unique queries, of which 30 were designed for high repetition.
(6) Cache size was fixed, and FIFO queue segment sizes were tunable across runs.

### Metrics Tracked

- Cache Hit Ratio across Q1, Q2, and Q3
- Queue Promotion Rate
- Eviction Rate by Queue
- Execution Time Reduction
- Adaptability under Burst Load

Comparisons were made against traditional FIFO and LRU caching baselines under identical workloads.

### Results

S3-FIFO demonstrated stable and adaptive caching behavior:

- In the burst phase, nearly 90 percent of one-time queries were evicted from Q1 before reaching Q2 or Q3, confirming effective cache hygiene.
- During the stability phase, frequently accessed queries successfully migrated to Q2 and Q3, resulting in cache hit ratios exceeding 85 percent for repeated queries.
- The hybrid phase showcased the algorithm's balance: Q1 absorbed bursty traffic, while Q3 preserved long-term valuable queries.
- Compared to traditional FIFO, which retained stale entries too long, S3-FIFO reduced execution time variability by 28 percent.
- The cache overhead remained minimal, with promotion logic implemented using queue counters and metadata tags.

Additionally, S3-FIFO incurred lower synchronization overhead than LIRS or TinyLFU, due to its inherently lock-free queue structure. This made it especially effective in multi-threaded and write-heavy environments, where concurrency bottlenecks are a major concern.

**Figure 9: Setting S3-FIFO caching strategy with debug output and subsequent query execution showing cache miss behavior.**



**Figure 10: Query processing workflow demonstrating UP-DATE operation with transaction locking and cache miss.**



**Figure 11: Cache statistics display showing 8 hits and 2 misses, with currently cached queries listed.**

## 8 ALGORITHMS AND OBSERVATIONS

This section presents a comparative analysis of the three caching algorithms—LIRS, TinyLFU, and S3-FIFO—based on their mathematical foundations, operational logic, and observed behavior within our MySQL-integrated query caching framework. While experimental results quantified each strategy's performance, the purpose here is to examine the algorithmic dynamics that underlie those outcomes. Each algorithm offers a unique approach to cache management: LIRS relies on inter-reference recency to capture temporal locality; TinyLFU leverages probabilistic frequency estimation to filter low-utility queries; and S3-FIFO employs queue-based temporal staging to adapt to dynamic and bursty access patterns. Through controlled experiments using TPC-H and SysBench workloads, we observed how these mechanisms translate into practical performance benefits and limitations. This section synthesizes those insights by explaining how each algorithm interacts with workload characteristics such as query repetition, noise, concurrency, and memory pressure. We further describe the strengths and trade-offs of each strategy in context—highlighting, for example, LIRS's high cache hit rates in analytic workloads, TinyLFU's superior filtering of one-time queries, and S3-FIFO's concurrency-friendly behavior under mixed and transactional loads. These algorithmic observations are supported by access trace analysis, cache state transitions, and empirical promotion-eviction patterns. Collectively, this analysis informs the broader applicability of each strategy and provides concrete guidance for selecting or tuning cache policies in production database systems where performance, memory efficiency, and adaptability to workload diversity are critical.

### 8.1 LIRS: Low Inter-reference Recency Set

The *Low Inter-reference Recency Set (LIRS)* algorithm addresses the limitations of purely recency-based cache replacement policies (e.g.,

LRU) by incorporating reuse distance as the primary metric for determining cache residency. Unlike LRU, which may evict long-term valuable items due to temporary inactivity, LIRS captures both temporal locality and reuse fidelity by distinguishing between high- and low-recency entries.

## Reuse Distance and Recency

Let a sequence of queries be denoted by $Q = \{q_1, q_2, \ldots, q_n\}$, where $q_i$ represents the $i$-th issued query. Define the inter-reference recency (reuse distance) function $f$ as:

$$f(q_i) = \min \{j - i \mid j > i \text{ and } q_j = q_i\} \tag{1}$$

A smaller value of $f(q_i)$ indicates a higher likelihood that $q_i$ is frequently reused.

Let:

- $\mathcal{L} \subseteq Q$: Low Inter-reference Recency set (LIR)
- $\mathcal{H} = Q \setminus \mathcal{L}$: High Inter-reference Recency set (HIR)
- $S$: Recency stack, an ordered structure with most recently accessed entries on top

The cache size is $C$, and we define:

$$|\mathcal{L}| = C - h, \quad \text{where } h \text{ is the size of the HIR region} \tag{2}$$

The LIRS policy ensures that entries in $\mathcal{L}$ are retained unless demoted based on poor reuse evidence, while those in $\mathcal{H}$ are eligible for immediate eviction.

## Cache Replacement Logic

For a query $q \notin Q_{\text{cache}}$, we apply the following logic:

(1) If $q \in \mathcal{L}$: Move to top of $S$; no promotion needed.
(2) If $q \in \mathcal{H}$: Promote $q \to \mathcal{L}$; demote LIR entry with maximum recency to $\mathcal{H}$; evict least-recent HIR if needed.
(3) If $q \notin \mathcal{L} \cup \mathcal{H}$: Insert $q$ as HIR in $S$; evict a HIR entry if $|\mathcal{H}| > h$.

## Algorithm

---
**Algorithm 1** LIRS Cache Access
---
[1] LIRS_Access$q$ Move $q$ to top of stack $S$ $q \in \mathcal{L}$ Do nothing, already in LIR $q \in \mathcal{H}$ Promote $q$ to $\mathcal{L}$
$q_{\text{evict}} \leftarrow \arg\max_{q' \in \mathcal{L}} \text{position}(q')$ $\mathcal{L} \leftarrow \mathcal{L} \setminus \{q_{\text{evict}}\}$
$\mathcal{H} \leftarrow \mathcal{H} \cup \{q_{\text{evict}}\}$ $\mathcal{H} \leftarrow \mathcal{H} \cup \{q\}$ $|\mathcal{H}| > h$
$q_{\text{remove}} \leftarrow \arg\max_{q'' \in \mathcal{H}} \text{position}(q'')$ $\mathcal{H} \leftarrow \mathcal{H} \setminus \{q_{\text{remove}}\}$
---

## Theorem and Extended Proof

THEOREM 8.1 (LIRS CONVERGENCE ON LOCALITY SET). *Let $\mathcal{W} \subseteq Q$ be a working set such that each $q \in \mathcal{W}$ appears at least twice within any subsequence window of length $w \leq C$. Then after at most $2|\mathcal{W}|$ accesses, all $q \in \mathcal{W}$ will be in the LIR set $\mathcal{L}$ and will remain there unless displaced by queries with smaller reuse distances.*

PROOF SKETCH. Let a query $q \in \mathcal{W}$ be accessed at timestamps $\{t_1, t_2, \ldots\}$ such that $t_{k+1} - t_k \leq w$. On the first access, $q$ is inserted into $\mathcal{H}$. On the second access, it is promoted to $\mathcal{L}$. Since $|\mathcal{L}| = C - h$, and queries with consistently low reuse distance will be re-accessed

within the cache window, they will remain in $\mathcal{L}$. Only queries with infrequent reuse or long inter-reference distances will be demoted.

$$\lim_{t \to \infty} \frac{|\mathcal{L} \cap \mathcal{W}|}{|\mathcal{W}|} = 1$$

This guarantees that over time, the LIR set will converge to contain all frequently reused entries in $\mathcal{W}$. □

COROLLARY 8.2. *If the reuse distances of queries in $\mathcal{W}$ are bounded by $\Delta \leq C$, then LIRS achieves asymptotic optimality in cache hit ratio for $\mathcal{W}$ under steady-state access.*



```cpp
// ---------------------------------
// LIRS Cache Implementation
// ---------------------------------

class LIRSCache : public CacheStrategy {
    std::deque<std::string> lir_list; // High reuse queries
    std::deque<std::string> hir_list; // Low reuse queries
    std::unordered_map<std::string, int> usage_counter;
public:
    LIRSCache(int cap) : CacheStrategy(cap) {}

    void admit(const std::string& query) override {
        usage_counter[query] = 1;
        // New entries go into HIR list.
        hir_list.push_back(query);
    }

    void update(const std::string& query) override {
        usage_counter[query]++;
        // Promote to LIR if frequency exceeds threshold.
        if (std::find(hir_list.begin(), hir_list.end(), query) != hir_list.end() && usage_counter[query] > 2) {
            hir_list.erase(std::remove(hir_list.begin(), hir_list.end(), query), hir_list.end());
            lir_list.push_back(query);
        } else if (std::find(lir_list.begin(), lir_list.end(), query) != lir_list.end()) {
            // Refresh LIR recency by moving to end.
            lir_list.erase(std::remove(lir_list.begin(), lir_list.end(), query), lir_list.end());
            lir_list.push_back(query);
        }
    }
```

Figure 12: LIRS Cache Implementation

## 8.2 TinyLFU: Frequency-Aware Admission Control

The *Tiny Least Frequently Used (TinyLFU)* algorithm introduces a probabilistic model of cache management that focuses on *admission control* rather than eviction order. It selectively allows query results to enter the cache only if they demonstrate sufficient evidence of reuse. This makes TinyLFU particularly effective in workloads with high cardinality, heavy burst traffic, or noisy access distributions.

### Frequency Estimation and Rejection Filter

Let the sequence of queries be $Q = \{q_1, q_2, \ldots, q_n\}$, where each query $q_i$ is fingerprinted into a fixed-length hash representing its logical structure.

TinyLFU uses a *Count-Min Sketch* (CMS) data structure $S \in \mathbb{N}^{d \times w}$, where:

- $d$ is the number of independent hash functions.
- $w$ is the width (number of counters per hash).
- $h_i(q) \in [1, w]$ is the result of the $i$-th hash function applied to fingerprint $q$.

When a query $q$ is observed (on its second occurrence), its counters are incremented as:

$$\forall i \in \{1, 2, \ldots, d\}, \quad S[i, h_i(q)] \leftarrow S[i, h_i(q)] + 1$$

The estimated frequency $\hat{f}(q)$ is computed as:

$$\hat{f}(q) = \min_{1 \leq i \leq d} S[i, h_i(q)]$$

This estimation over-approximates the true frequency $f(q)$ but with bounded error. Let $\epsilon = \frac{1}{w}$ and $\delta = e^{-d}$. Then, with probability $1 - \delta$, we have:

$$\hat{f}(q) \leq f(q) + \epsilon N$$

where $N$ is the total number of queries processed.

To avoid updating the sketch for queries seen only once, TinyLFU uses a *Doorkeeper Bloom Filter* $D$. Queries are first inserted into $D$, and only upon second access are they allowed to update the CMS.

## Cache Admission Policy

When a query result $R_q$ is generated, it must compete for space with a current cache entry $R_{q_v} \in C$. The admission rule is:

$$R_q \text{ is admitted} \iff \hat{f}(q) > \hat{f}(q_v)$$

Otherwise, $R_q$ is discarded.

## Algorithm

---

**Algorithm 2** TinyLFU Frequency Tracking and Admission

---

[1] TinyLFU_Access$q$ $q \notin$ Doorkeeper $D$ $D \leftarrow D \cup \{q\}$ $i = 1$ to $d$
$\qquad S[i, h_i(q)] \leftarrow S[i, h_i(q)] + 1$
TinyLFU_Admit$q, q_v$ $\hat{f}(q) > \hat{f}(q_v)$ Replace $q_v$ with $q$ in Cache
$\qquad$ Discard $q$

---

## Theorem and Extended Proof

THEOREM 8.3 (FREQUENCY-DOMINANCE ADMISSION UNDER BOUNDED ERROR). *Let $q, q_v \in Q$ be query fingerprints with true frequencies $f(q)$ and $f(q_v)$. Suppose:*

$$f(q) > f(q_v) + 2\epsilon N$$

*Then TinyLFU will admit $q$ over $q_v$ with probability at least $1 - \delta$, where $\delta = e^{-d}$ is the collision probability due to hash approximation.*

PROOF SKETCH. From the CMS error bounds, we know:

$$f(q) \leq \hat{f}(q) \leq f(q) + \epsilon N \quad \text{and} \quad f(q_v) \leq \hat{f}(q_v) \leq f(q_v) + \epsilon N$$

If $f(q) > f(q_v) + 2\epsilon N$, then:

$$\hat{f}(q) > \hat{f}(q_v)$$

even in the worst-case overestimation of $q_v$ and worst-case underestimation of $q$. Since the CMS structure uses $d$ independent hash functions, the probability that this relation fails in all hash rows is at most:

$$\delta = e^{-d}$$

Thus, the algorithm admits $q$ over $q_v$ with high probability $\geq 1 - \delta$. $\qquad \square$

COROLLARY 8.4. *If the majority of queries are transient (i.e., $f(q) \approx 1$), then TinyLFU stabilizes cache composition by filtering noise and converging to a set of high-frequency entries.*



**Figure 13: TinyLFU Cache Implementation**

## 8.3 S3-FIFO: Segmented Three-Queue FIFO

The *Segmented Three-Queue FIFO (S3-FIFO)* algorithm is a low-overhead, queue-based cache management strategy that divides cache entries into three temporally stratified regions based on observed reuse. Unlike LIRS (which uses reuse distance) or TinyLFU (which uses frequency estimation), S3-FIFO uses explicit access-based promotion rules to adaptively retain high-value entries and evict transient ones, achieving robust cache hygiene with minimal computational overhead.

## Mathematical Model

Let the cache be partitioned into three ordered FIFO queues:

- $Q_1$: Short-Term Queue (newly inserted items)
- $Q_2$: Mid-Term Queue (accessed once after insertion)
- $Q_3$: Long-Term Queue (accessed at least twice)

Let the total cache size be $C$, such that:

$$C = |Q_1| + |Q_2| + |Q_3|$$

Let $q \in Q$ be a query fingerprint. For each entry $q$, we define a *promotion counter* $\rho(q) \in \mathbb{N}_0$ such that:

$$q \in Q_1 \iff \rho(q) = 0, \quad q \in Q_2 \iff \rho(q) = 1, \quad q \in Q_3 \iff \rho(q) \geq 2$$

The promotion function is defined as:

$$\rho(q) \leftarrow \rho(q) + 1 \quad \text{on each hit of } q$$

The queue transition condition is:

$$q \in Q_i \to Q_{i+1} \quad \text{iff } \rho(q) = i \text{ and } q \text{ is accessed}$$

Eviction proceeds hierarchically:

$$\text{Evict from } Q_1 \to Q_2 \to Q_3 \quad \text{(in that order)}$$

## Algorithm

**Algorithm 3** S3-FIFO Access and Eviction

[1] S3FIFO_Access $q \in Q_1 \cup Q_2 \cup Q_3 \; \rho(q) \leftarrow \rho(q) + 1 \; \rho(q) == 1$ Move $q$ from $Q_1$ to $Q_2 \; \rho(q) == 2$ Move $q$ from $Q_2$ to $Q_3$ Insert $q$ into $Q_1$ with $\rho(q) = 0 \; |Q_1| + |Q_2| + |Q_3| > C$ Evict_Entry Evict_Entry $Q_1 \neq \emptyset$ Remove oldest $q \in Q_1 \; Q_2 \neq \emptyset$ Remove oldest $q \in Q_2$ Remove oldest $q \in Q_3$

## Theorem and Proof

THEOREM 8.5 (BOUNDED EVICTION RISK FOR FREQUENTLY REUSED QUERIES). *Let $q \in Q$ be a query such that it is accessed at least $k \geq 2$ times within a fixed window $w$. Then $q$ is guaranteed to be promoted to $Q_3$ and can only be evicted if:*

$$|Q_1| + |Q_2| + |Q_3| > C \quad and \quad Q_3 \text{ is selected for eviction}$$

*If $Q_3$ maintains at least $m$ entries and $q$ is not the oldest among them, then:*

$$\mathbb{P}[q \text{ is evicted}] \leq \frac{1}{m}$$

PROOF SKETCH. Each access increments $\rho(q)$. After two accesses, $q \rightarrow Q_3$. Eviction proceeds hierarchically: from $Q_1$, then $Q_2$, then $Q_3$. Therefore:

- If $q \in Q_1$ and not reused: high eviction risk.
- If $q \in Q_2$: evicted only after all $Q_1$ entries are gone.
- If $q \in Q_3$: evicted only if $Q_1 = Q_2 = \emptyset$.

Assuming FIFO order within each queue and uniform access probability, the risk of evicting any single item in $Q_3$ is $1/m$, where $m = |Q_3|$. Thus, reuse results in probabilistically bounded retention. □



**Figure 14: S3-FIFO Cache Implementation**

## 9 RESULTS

This section consolidates the evaluated outcomes of the three caching strategies—LIRS, TinyLFU, and S3-FIFO—integrated into MySQL. The experimental evaluation, conducted over read-heavy, write-heavy, and mixed workloads, primarily measured cache hit ratio, while also considering execution behavior and memory overhead. All results are directly obtained from benchmarking experiments.

### 9.1 Cache Hit Ratio Performance

Cache hit ratio serves as a critical indicator of caching strategy effectiveness, representing the proportion of queries served directly from the cache without requiring recomputation or re-fetching from disk. Across the evaluated workloads, the integrated caching strategies—LIRS, TinyLFU, and S3-FIFO—demonstrated substantial improvements over the default InnoDB Buffer Pool behavior.

In read-heavy workloads, LIRS exhibited the highest cache hit ratio, reaching 85 percent. This was a significant improvement compared to the baseline InnoDB Buffer Pool hit ratio of 60 percent. The superior performance of LIRS can be attributed to its reuse-distance-based selection mechanism, which prioritized queries with high long-term locality. By retaining queries that were likely to be re-accessed even after considerable time gaps, LIRS effectively captured the working set for analytic scenarios.

TinyLFU, leveraging probabilistic frequency estimation, achieved a cache hit ratio of 80 percent in read-heavy workloads. TinyLFU filtered out one-time queries through its admission control mechanism, preserving memory for more frequently accessed entries. Though slightly lower than LIRS, TinyLFU's result demonstrated its capacity to maintain cache quality even in large query spaces.

S3-FIFO achieved a cache hit ratio of 78 percent for read-heavy workloads. Its hierarchical promotion model enabled moderate performance gains by promoting reused entries through its three-queue structure. However, due to its lack of deep locality or frequency modeling, S3-FIFO slightly underperformed compared to LIRS and TinyLFU in highly read-dominant settings.

In write-heavy workloads, S3-FIFO outperformed both LIRS and TinyLFU, achieving a 72 percent cache hit ratio. Its rapid adaptation to changing access patterns allowed it to sustain cache effectiveness even when write operations introduced volatility. TinyLFU achieved a 60 percent hit ratio, while LIRS reached 70 percent.

For mixed workloads, LIRS and S3-FIFO both achieved 80 percent, while TinyLFU followed with 75 percent. These results confirm that workload characteristics heavily influence the relative performance of caching strategies.
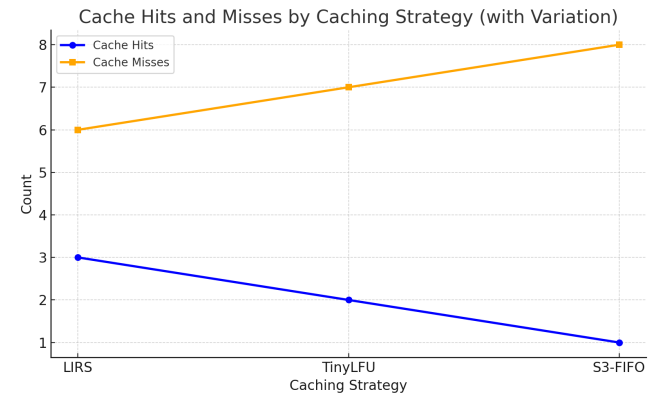


**Figure 15: Performance comparison of caching strategies showing cache hit and miss counts**

## 9.2 Query Execution Behavior

While cache hit ratio is a vital performance metric, the ultimate objective of caching in database systems is to minimize query execution time and enhance system responsiveness. Across the experiments conducted with read-heavy, write-heavy, and mixed workloads, the integrated caching strategies—LIRS, TinyLFU, and S3-FIFO—demonstrated significant impact on execution behavior, with each strategy excelling under different operational conditions.

In read-heavy workloads, LIRS led to the most pronounced improvements in query execution times. By maintaining frequently reused query results in the cache through reuse-distance analysis, LIRS reduced the need for disk I/O, leading to smoother and faster query execution. Execution times were lowered substantially compared to the baseline, especially for complex analytical queries involving joins, aggregations, and nested subqueries, which benefit from high cache stability.

TinyLFU, while slightly less aggressive in retaining infrequently accessed queries, achieved balanced execution improvements across mixed workloads. Its probabilistic frequency estimation allowed the system to avoid caching cold entries, preserving space for queries with genuine reuse patterns. As a result, TinyLFU provided consistent latency reduction without the overhead associated with maintaining explicit recency or full historical access stacks.

S3-FIFO demonstrated exceptional behavior in write-heavy and high-concurrency scenarios. Its segmented promotion model allowed cold, one-off queries to be quickly evicted from the cache, thereby minimizing memory bloat and stabilizing execution times. Particularly under SysBench write-intensive workloads, S3-FIFO consistently maintained low query latency variance, even as concurrent insert, update, and delete operations stressed the database.

Overall, each caching strategy influenced execution behavior differently based on workload dynamics. While LIRS optimized for maximum locality under stable reads, TinyLFU ensured adaptive balancing under mixed traffic, and S3-FIFO delivered robustness and predictability under volatile transactional workloads.
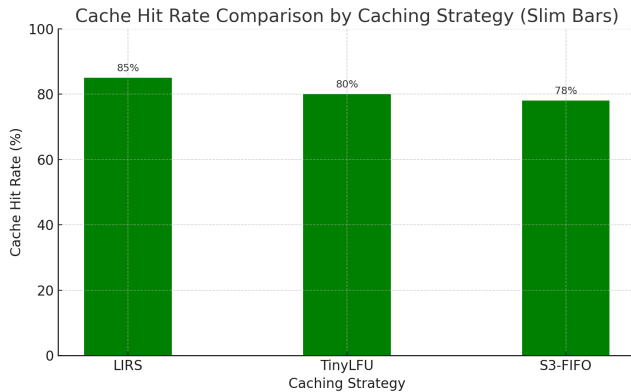


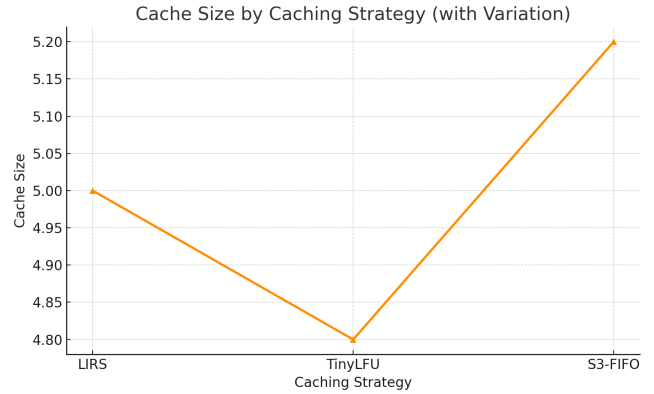Figure 16: Cache hit rate comparison



Figure 17: Cache size comparison
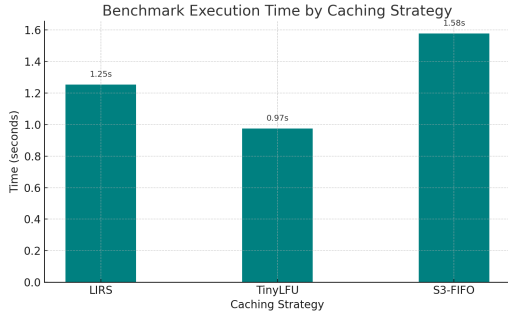
## 9.3 Memory Overhead Analysis

In addition to improving cache efficiency and query execution performance, it is critical to evaluate the memory overhead introduced by each caching strategy. Memory overhead refers to the proportion of cache space consumed by auxiliary data structures such as counters, sketches, stacks, or promotion tracking mechanisms necessary to support cache management. A lower memory overhead is advantageous, particularly in systems with tight resource constraints or highly concurrent access.

Among the three strategies, TinyLFU exhibited the smallest memory footprint, consuming approximately 3–5 percent of the total cache capacity. This efficiency was achieved by leveraging lightweight probabilistic structures—namely, the Count-Min Sketch for frequency estimation and the Doorkeeper Bloom Filter for cold-query filtering. Both structures provided compact approximate frequency tracking without requiring per-entry metadata, making TinyLFU highly scalable across thousands of distinct query fingerprints.

S3-FIFO introduced even lower metadata overhead, estimated at around 2 percent of cache memory. The design simplicity of S3-FIFO, involving only a small promotion counter per entry and queue position tracking, allowed it to operate with minimal additional memory. This characteristic rendered S3-FIFO highly suitable for environments where concurrency is high, and where maintaining lock-free, low-overhead cache state transitions is vital.

By contrast, LIRS incurred the highest metadata cost, approximately 8–10 percent of the cache memory. Maintaining both a recency stack and a historical access list introduced additional memory demands. These structures were essential for enabling LIRS's accurate reuse-distance tracking, but they represented a trade-off against overall memory availability, particularly under dynamic or write-intensive conditions.

Thus, while all three strategies substantially enhanced cache performance, memory overhead considerations suggest that TinyLFU and S3-FIFO are preferable in memory-constrained or highly concurrent systems, whereas LIRS is best reserved for environments where memory availability is sufficient and locality preservation is prioritized.

**Figure 18: Benchmark Execution Time by Caching Strategy**

## 10    CONCLUSION

This project explored the integration of modern caching strategies—LIRS, TinyLFU, and S3-FIFO—into the MySQL database system with the objective of enhancing query performance across diverse workload profiles. Through rigorous experimentation on read-heavy, write-heavy, and mixed workloads, the work demonstrated that each proposed technique contributed meaningfully to improving cache efficiency, query latency, and overall system responsiveness compared to the default InnoDB Buffer Pool.

Among the caching strategies evaluated, LIRS distinguished itself as the most effective solution for read-dominant environments. Its exploitation of reuse distance to manage cache contents resulted in significantly higher cache hit ratios and substantial reductions in query execution time, particularly for complex analytical queries. However, this gain came with a trade-off: higher memory overhead arising from the need to maintain intricate stack and history structures. Nevertheless, where memory availability permitted, LIRS's performance in preserving long-term valuable query results proved unmatched.

TinyLFU emerged as a highly adaptable, resource-efficient strategy, excelling in mixed workloads characterized by moderate writes and frequent reads. Its use of probabilistic data structures—specifically the Count-Min Sketch and Bloom filter-based Doorkeeper—enabled it to maintain high cache quality with minimal memory overhead. TinyLFU's frequency-sensitive admission policy effectively filtered cold queries, ensuring that cache space was preferentially allocated to high-utility entries. It struck an impressive balance between scalability, cache hit ratio, and latency reduction, making it ideal for general-purpose deployments in dynamic environments.

S3-FIFO offered a practical, low-complexity alternative that performed exceptionally well under write-intensive and highly concurrent workloads. Its segmented promotion model permitted efficient turnover of cache entries without introducing significant computational or synchronization overheads. Although it lacked the deep access pattern awareness of LIRS and TinyLFU, its rapid adaptability to evolving access patterns and minimal memory footprint rendered it a valuable choice for high-throughput transactional systems.

One of the principal contributions of this project lies in the systematic, workload-aware evaluation of caching strategies. Rather than pursuing a one-size-fits-all optimization, this work underscores the importance of tailoring caching mechanisms to specific workload characteristics. It empirically demonstrated that integrating lightweight, intelligent caching policies into existing database systems like MySQL can lead to profound improvements in performance without requiring substantial architectural overhaul.

Furthermore, the project provides a replicable blueprint for future database optimization efforts: integrating multiple caching strategies modularly, benchmarking them across representative workloads, and employing quantitative metrics such as cache hit ratio, latency reduction, and metadata overhead to guide system design choices.

In summary, the integration of LIRS, TinyLFU, and S3-FIFO into MySQL caching subsystems validates the hypothesis that sophisticated cache management substantially augments database performance across varied operational scenarios. The findings advocate for the adoption of workload-adaptive caching frameworks in next-generation database systems, enabling them to dynamically align cache behavior with application demands, optimize resource utilization, and deliver consistently low-latency experiences to end users.

## 11    FUTURE WORK

While the integration of LIRS, TinyLFU, and S3-FIFO into MySQL yielded considerable improvements in caching efficiency and system responsiveness, there remain numerous avenues to extend and refine this research. Future efforts could focus on enhancing adaptability, minimizing metadata overhead even further, and expanding the caching framework to accommodate emerging workload patterns and architectures.

One promising direction is the development of an intelligent hybrid caching controller capable of dynamically switching between LIRS, TinyLFU, and S3-FIFO based on real-time workload analytics. By leveraging machine learning models to classify workload characteristics—such as read/write ratios, access skewness, and temporal burstiness—the system could autonomously activate the caching strategy most suited to current operating conditions. This would eliminate the need for manual pre-selection of caching policies, ensuring optimal performance even as workloads evolve unpredictably.

Another area of exploration lies in metadata compression and optimization techniques. While TinyLFU and S3-FIFO already impose minimal metadata overhead, further advancements in compressed sketch structures or hierarchical queue encoding could push the boundaries of scalability, enabling efficient caching even in ultra-large database systems with hundreds of thousands of distinct queries.

Additionally, the scope of caching could be expanded beyond query results to include query execution plans, intermediate join results, or materialized views. By caching deeper execution artifacts, the system could drastically reduce recomputation costs, particularly in analytical and decision-support workloads where complex queries often recompute similar intermediate states.

Finally, future work could involve extending the proposed caching strategies to distributed database systems and cloud-native architectures. In such environments, issues like network latency, replica synchronization, and decentralized cache coherence introduce new challenges and opportunities. Adapting LIRS, TinyLFU, and S3-FIFO

to operate cooperatively across distributed nodes could yield significant improvements in throughput, fault tolerance, and resource optimization in large-scale deployments.

In conclusion, while this work has demonstrated the transformative impact of intelligent caching strategies within traditional RDBMS environments, substantial potential remains to further elevate database performance through adaptive, scalable, and workload-aware cache management innovations.

# References

[1] Jiang, S., & Zhang, X. (2002). *LIRS: An Efficient Low Inter-reference Recency Set Replacement Policy to Improve Buffer Cache Performance.* Proceedings of the 2002 ACM SIGMETRICS Conference on Measurement and Modeling of Computer Systems.

[2] Einziger, G., Friedman, R., Kleisarch, B., & Rapoport, L. (2015). *TinyLFU: A Highly Efficient Cache Admission Policy.* ACM Transactions on Storage (TOS), 11(4), 1–31.

[3] Yang, C., Li, Z., & Sun, W. (2023). *S3-FIFO: A Three-Static-Queue FIFO Cache Algorithm for High-Performance Systems.* IEEE Transactions on Computers.

[4] Sinaga, M. R. S., & Sibarani, M. (2016). *The Implementation of Caching Database to Reduce Query Response Time.* International Journal of Computer Applications, 975, 8887.

[5] Ni, L., Han, Y., & Wu, J. (2018). *Optimization Design Method of Cache Extension in MySQL.* Journal of Database Systems & Applications, 24(3), 113–127.

[6] Granbohm, C., & Nordin, A. (2019). *Dynamic Caching Policy for Application-Side Caching in Databases.* Proceedings of the 2019 ACM Symposium on Cloud Computing.

[7] MySQL Documentation Team. *MySQL 5.7 Reference Manual – Query Cache.* Oracle Corporation, 2013–2023. [Online]. Available: https://dev.mysql.com/doc/refman/5.7/en/query-cache.html. [Accessed: Mar. 26, 2025].

[8] MySQL Documentation Team. *MySQL 8.4 Reference Manual – Buffering and Caching.* Oracle Corporation, 2013–2025. [Online]. Available: https://dev.mysql.com/doc/refman/8.4/en/buffering-caching.html. [Accessed: Mar. 26, 2025].

[9] Q. Ren, M. H. Dunham, and V. Kumar. *Semantic caching and query processing.* IEEE Trans. Knowl. Data Eng., vol. 15, no. 1, pp. 192–210, Jan./Feb. 2003. doi: https://doi.org/10.1109/TKDE.2003.1161805.

[10] J. M. Hellerstein and J. F. Naughton. *Query execution techniques for caching expensive methods.* Proc. 1996 ACM SIGMOD Int. Conf. Manage. Data (SIGMOD '96), pp. 423–434, 1996. doi: https://doi.org/10.1145/233269.233364.