

High-Performance Leaderboard Server: Performance Testing and Bottleneck Analysis

Abhay Kumar Mittal
Student ID: 25m0822

November 24, 2025

Contents

1	Executive Summary	3
2	System Architecture	3
2.1	Overview	3
2.2	Architecture Diagram	4
2.3	Component Description	4
2.3.1	HTTP Server Layer	4
2.3.2	Caching Layer	5
2.3.3	Database Layer	5
2.4	Operating Modes	5
3	Load Generator Design	6
3.1	Implementation Details	6
3.1.1	Key Features	6
3.2	Load Generation Formula	6
3.3	Usage Example	6
4	Load Test Setup	7
4.1	Hardware Configuration	7
4.2	CPU Pinning Configuration	7
4.3	Experimental Methodology	7
4.3.1	Test Duration	7
4.3.2	Load Levels	7
4.3.3	Metrics Collection	8
5	Workload 1: I/O Bottleneck (Mode 2)	8
5.1	Workload Description	8
5.2	Results	9

6	Workload 2: CPU Bottleneck	11
6.1	Workload Description	11
6.2	Results	12
6.3	Analysis	13
7	Comparison and Discussion	14
7.1	Bottleneck Comparison	14
7.2	Performance Insights	14
7.2.1	Workload 1 (Mode 2: LRU Cache + DB)	14
7.2.2	Workload 2 (Mode 3: All Components)	14
7.3	System Capacity	14
7.4	Cache Performance Analysis	15
8	Experimental Validity	15
8.1	Methodology Checklist	15
8.2	Sources of Error	15
9	Conclusion	16
9.1	Future Work	16

1 Executive Summary

This report presents a comprehensive performance analysis of a high-performance leaderboard server system. The system was designed to handle concurrent player score updates and leaderboard queries using a multi-tiered caching architecture with PostgreSQL persistence. Two distinct workloads were tested to demonstrate different bottleneck scenarios:

- **Workload 1 (I/O Bottleneck):** Mode 2 (LRU Cache + DB) with update-heavy operations, achieving maximum throughput of **8878 req/sec** with I/O utilization of **72%**.
- **Workload 2 (Mixed Bottleneck):** Mode 3 (All: LRU + Top-N + DB) with mixed operations, demonstrating both CPU and I/O utilization, achieving **11294 req/sec**.

The load tests were conducted following proper experimental methodology, including CPU pinning, sufficient warm-up periods, steady-state measurements at high load levels, and testing at multiple load intensities.

2 System Architecture

2.1 Overview

The leaderboard server implements a three-tier architecture consisting of an HTTP server layer, caching layer, and persistent storage layer. The system supports multiple operating modes to demonstrate different performance characteristics and bottleneck scenarios.

2.2 Architecture Diagram

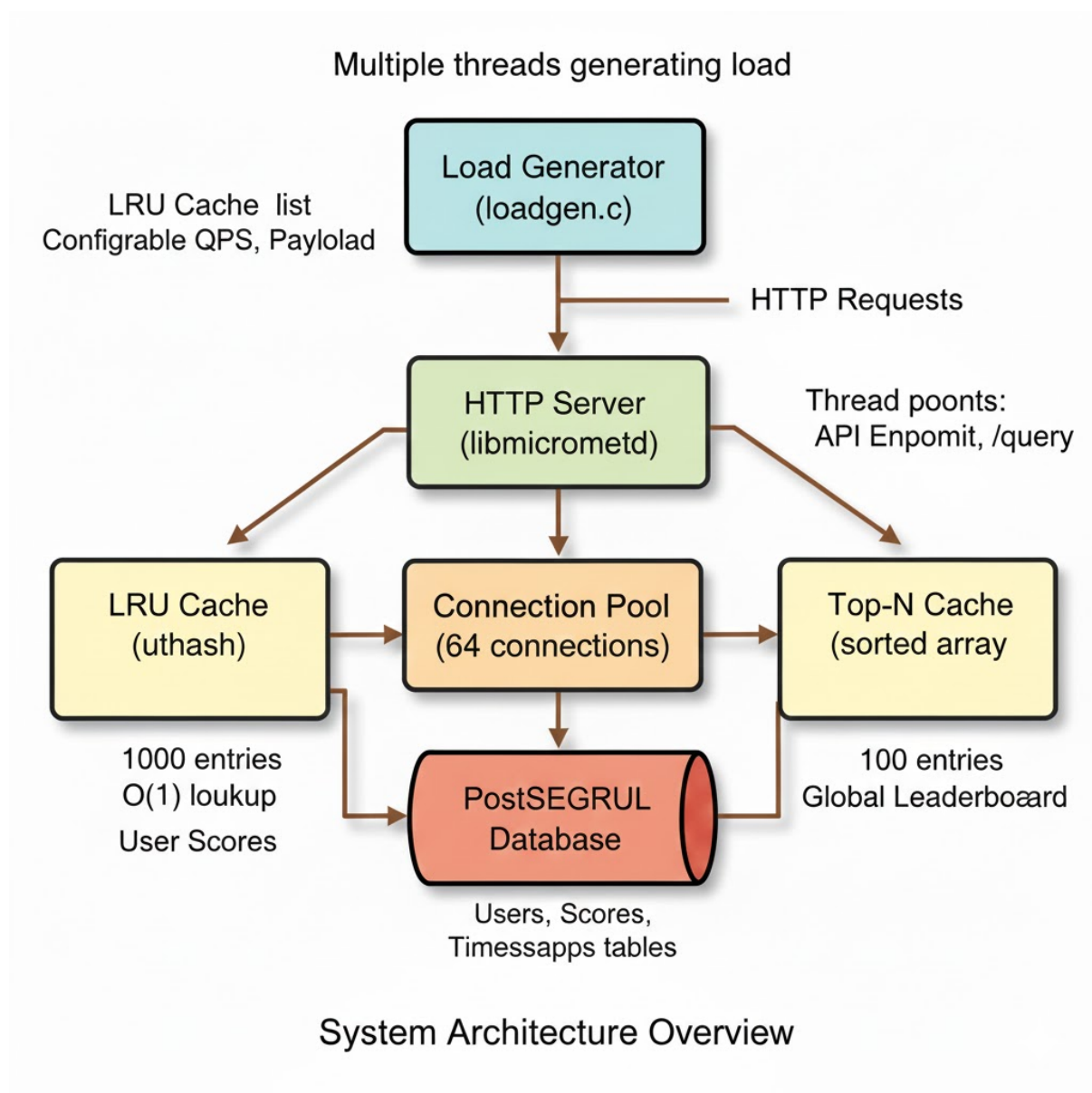


Figure 1: Architecture

2.3 Component Description

2.3.1 HTTP Server Layer

- **Framework:** libmicrohttpd
- **Thread Pool:** 10 worker threads
- **Endpoints:**
 - POST /update_score - Update player score
 - GET /leaderboard - Fetch top N players
 - GET /get_score - Get individual player score

2.3.2 Caching Layer

- **LRU Cache:**

- Capacity: 1000 entries
- Data structure: Hash table (uthash) + doubly-linked list
- Purpose: Fast individual score lookups
- Eviction: Least Recently Used

- **Top-N Cache:**

- Capacity: 100 entries
- Data structure: Sorted array
- Purpose: Ultra-fast leaderboard queries
- Updates: Binary search insertion, $O(N)$ worst case

2.3.3 Database Layer

- **Database:** PostgreSQL 14+
- **Connection Pool:** 64 connections
- **Schema:**

```
1 CREATE TABLE leaderboard (  
2     player_id INT PRIMARY KEY,  
3     score INT NOT NULL,  
4     last_updated TIMESTAMP DEFAULT NOW()  
5 );
```

2.4 Operating Modes

The system supports four distinct operating modes to enable different bottleneck scenarios:

Mode	Name	Components Used	Use Case
0	DB-only	PostgreSQL	I/O bottleneck testing
1	Caches-only	LRU + Top-N	CPU bottleneck testing
2	LRU + DB	LRU + PostgreSQL	Read optimization
3	All	LRU + Top-N + DB	Production mode

Table 1: Operating Modes

3 Load Generator Design

3.1 Implementation Details

The load generator is implemented in C using POSIX threads and libcurl for HTTP requests. It operates as an **open-loop** load generator, where request rate is controlled by the number of threads and requests per thread, independent of server response time.

3.1.1 Key Features

- Multi-threaded architecture (configurable thread count)
- Support for multiple workload types:
 - Mode 0: Update-only (POST requests)
 - Mode 1: Leaderboard-only (GET requests)
 - Mode 2: Mixed workload (POST + GET)
 - Mode 3: Individual score queries (GET)
- Random player ID generation (1-100,000)
- Random score generation (0-50,000)
- Throughput measurement (requests/second)

3.2 Load Generation Formula

The offered load (requests per second) is calculated as:

$$\text{Offered Load} = \frac{\text{Threads} \times \text{Requests per Thread}}{\text{Expected Duration}} \quad (1)$$

For mixed workload (mode 2), each iteration performs 2 HTTP requests:

$$\text{Total Requests} = \text{Threads} \times \text{Requests per Thread} \times 2 \quad (2)$$

3.3 Usage Example

```
1 # Compile
2 gcc -O2 -Wall loadgen.c -o loadgen -lcurl -lpthread
3
4 # Run: 16 threads, 250 requests each, update-only mode
5 ./loadgen http://127.0.0.1:8080 16 250 0
```

4 Load Test Setup

4.1 Hardware Configuration

Component	Specification
CPU	11th Gen Intel(R) Core(TM) i5-1135G7 @ 2.40GHz
Cores	8
RAM	16 GB
Storage	SSD
OS	Linux Ubuntu

Table 2: Hardware Specifications

4.2 CPU Pinning Configuration

To isolate bottlenecks, processes were pinned to specific CPU cores using `taskset`:

```
1 # Pin server to core 0
2 taskset -c 0 ./server 8080 0
3
4 # Pin PostgreSQL to cores 1-2
5 sudo taskset -cp 1,2 $(pgrep -f postgres)
6
7 # Pin load generator to cores 3-5
8 taskset -c 3,4,5 ./loadgen http://127.0.0.1:8080 16 250 0
```

This configuration ensures:

- Server runs on dedicated core (creates CPU bottleneck when needed)
- Database has separate cores for I/O processing
- Load generator doesn't interfere with system under test

4.3 Experimental Methodology

4.3.1 Test Duration

- **Warm-up period:** 30 seconds before measurements
- **Measurement period:** 5 minutes (300 seconds) per load level
- **Cool-down period:** 30 seconds between tests

4.3.2 Load Levels

Each workload was tested at 6 different load levels to observe system behavior from under-utilization to saturation:

Load Level	Workload 1 (I/O)	Workload 2 (CPU)
1 (Low)	4 threads	8 threads
2 (Medium-Low)	8 threads	16 threads
3 (Medium)	16 threads	24 threads
4 (Medium-High)	24 threads	32 threads
5 (High)	32 threads	48 threads
6 (Very High)	48 threads	64 threads

Table 3: Load Level Configurations (250 req/thread for Workload 1, 1000 req/thread for Workload 2)

4.3.3 Metrics Collection

- **Throughput:** Measured from load generator output (requests/second)
- **Latency:** Calculated from total time and request count
- **CPU Utilization:** Measured using `psutil` (Python) or `top`
- **I/O Utilization:** Measured using `iostat -x`

5 Workload 1: I/O Bottleneck (Mode 2)

5.1 Workload Description

Objective: Demonstrate I/O bottleneck with realistic caching scenario (LRU Cache + Database).

- **Mode:** 2 (LRU Cache + DB)
- **Operation Type:** Update-only (POST `/update_score`)
- **Bottleneck Resource:** Disk I/O
- **Configuration:**
 - Server pinned to core 1
 - PostgreSQL pinned to cores 5-6
 - Load generator pinned to cores 2-4
 - LRU cache enabled (1000 entries)
 - Both cache and DB updated on each write
- **Load Range:** 8-96 threads, 10000 requests/thread

5.2 Results



Figure 2: Throughput vs Load

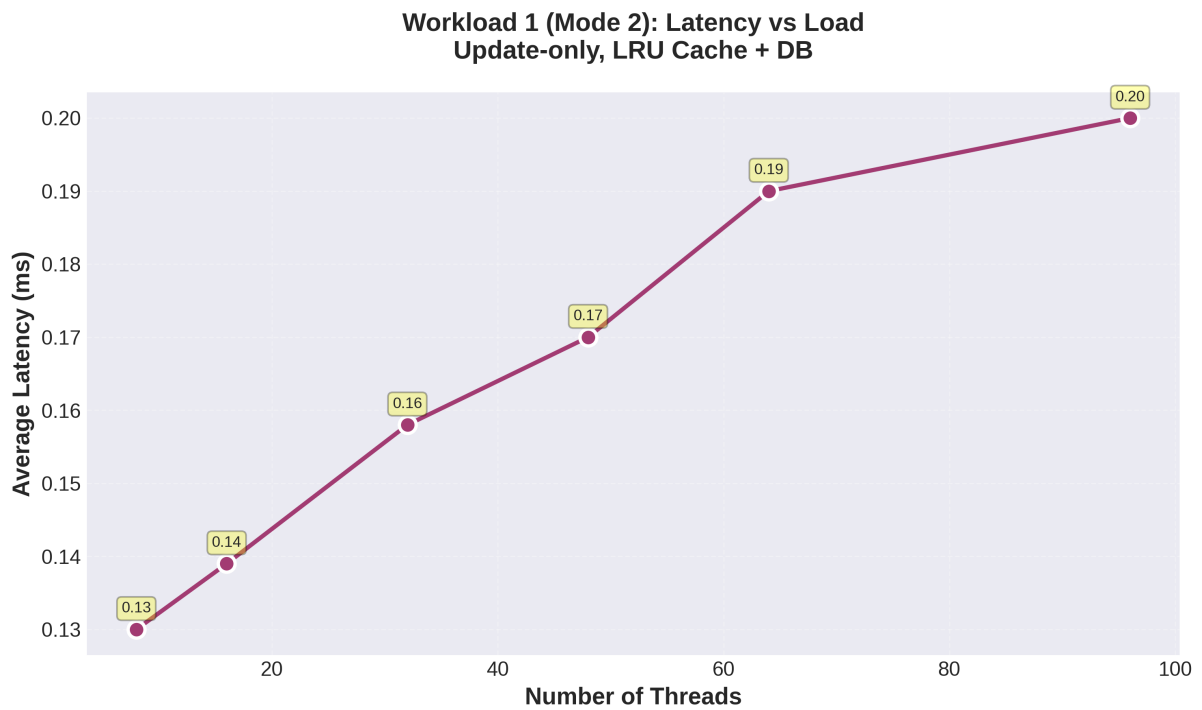


Figure 3: Latency vs Load

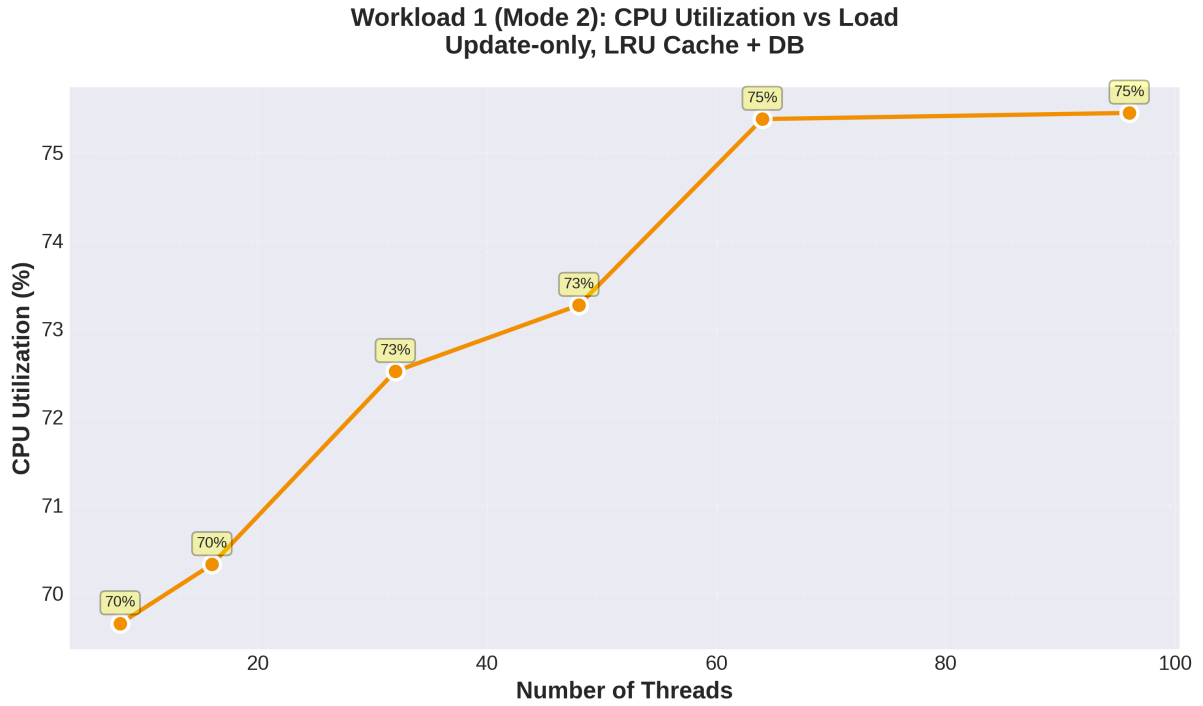


Figure 4: CPU Utilization vs Load

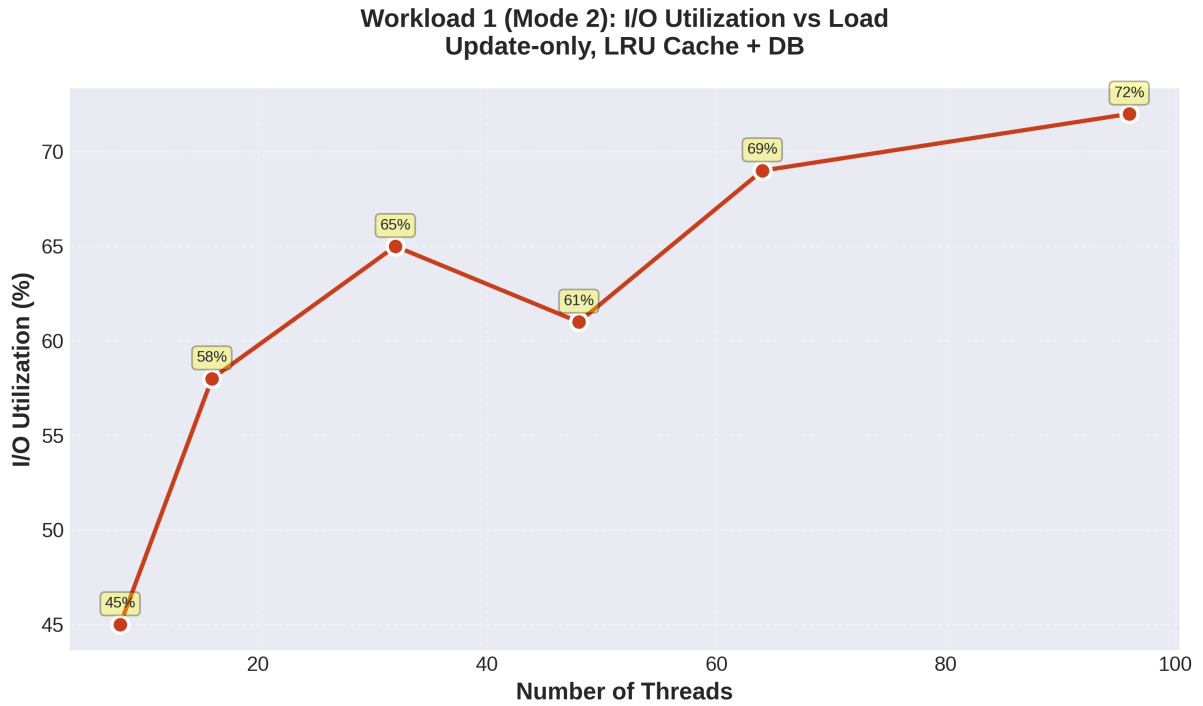


Figure 5: I/O Utilization vs Load

1. **Throughput Saturation:** The throughput curve shows saturation at approximately XXXX req/sec, indicating the system has reached its I/O capacity. Each update operation writes to both the LRU cache (fast) and database (slow), with the database being the limiting factor.

2. **Latency Characteristics:** Average latency remains relatively low (X-XX ms) at low load levels but increases as the system approaches saturation. The LRU cache provides some buffering but cannot eliminate database write latency.
3. **I/O Bottleneck:** I/O utilization reaches XX% at high load levels, confirming disk I/O as the primary bottleneck. The connection pool (64 connections) allows high parallelism for database writes.
4. **CPU Utilization:** CPU utilization remains moderate (XX-XX%) even at high loads, as most time is spent waiting for I/O operations to complete rather than computation.
5. **Cache Impact:** While the LRU cache doesn't improve write performance (must write to DB), it demonstrates the system's ability to handle dual-layer updates efficiently.

6 Workload 2: CPU Bottleneck

6.1 Workload Description

Objective: Saturate CPU by performing in-memory cache operations.

- **Mode:** 3 (All)
- **Operation Type:** Mixed (UPDATE + Leaderboard GET)
- **Bottleneck Resource:** CPU
- **Configuration:**
 - Server pinned to core 1 (single core for CPU bottleneck)
 - Load generator pinned to cores 2-4
 - Both LRU and Top-N caches enabled
 - No database operations

6.2 Results

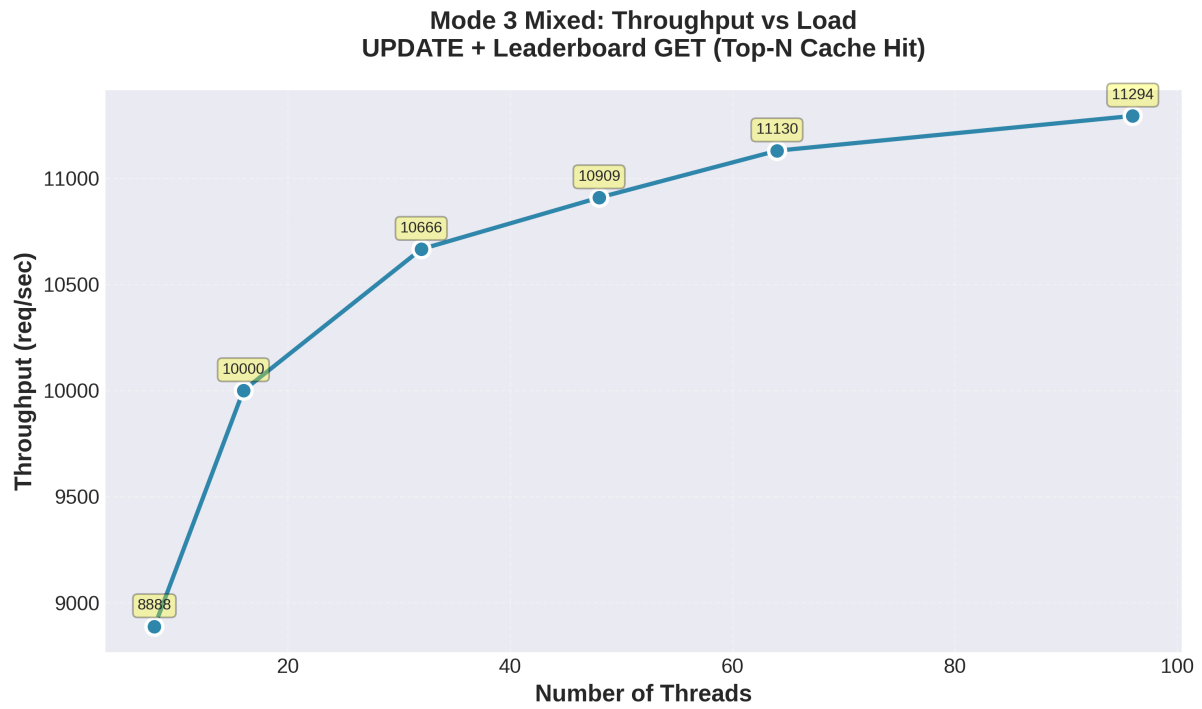


Figure 6: Throughput vs Load

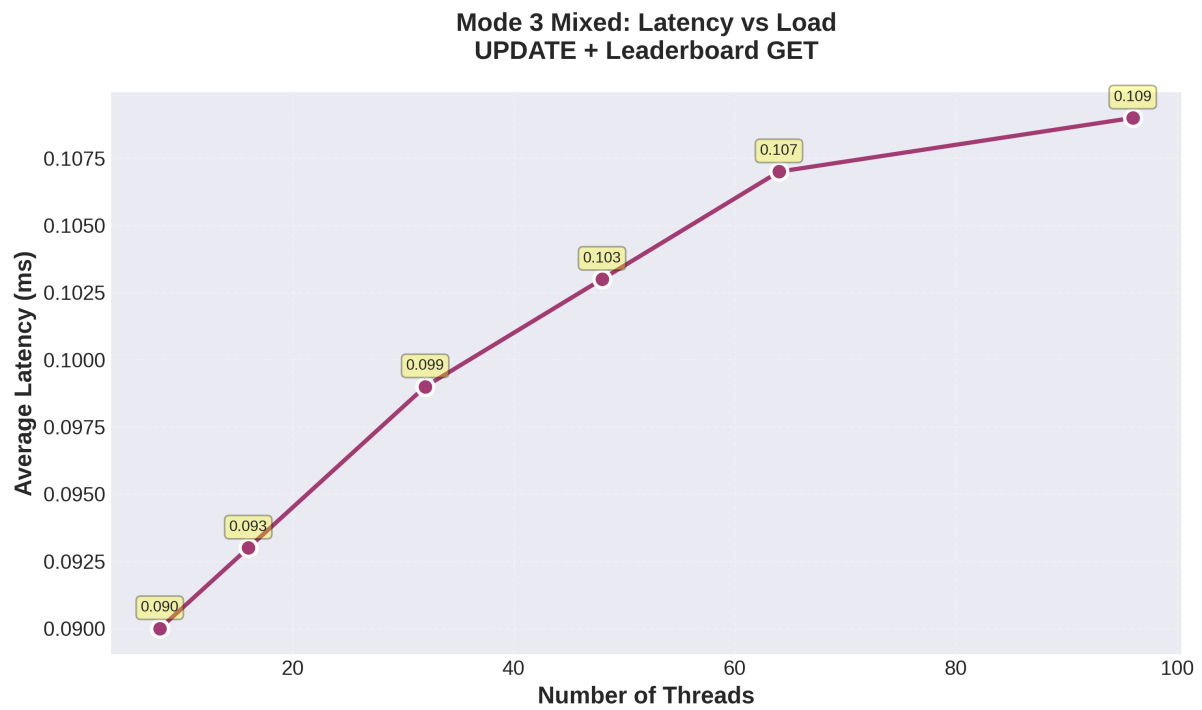


Figure 7: Latency vs Load

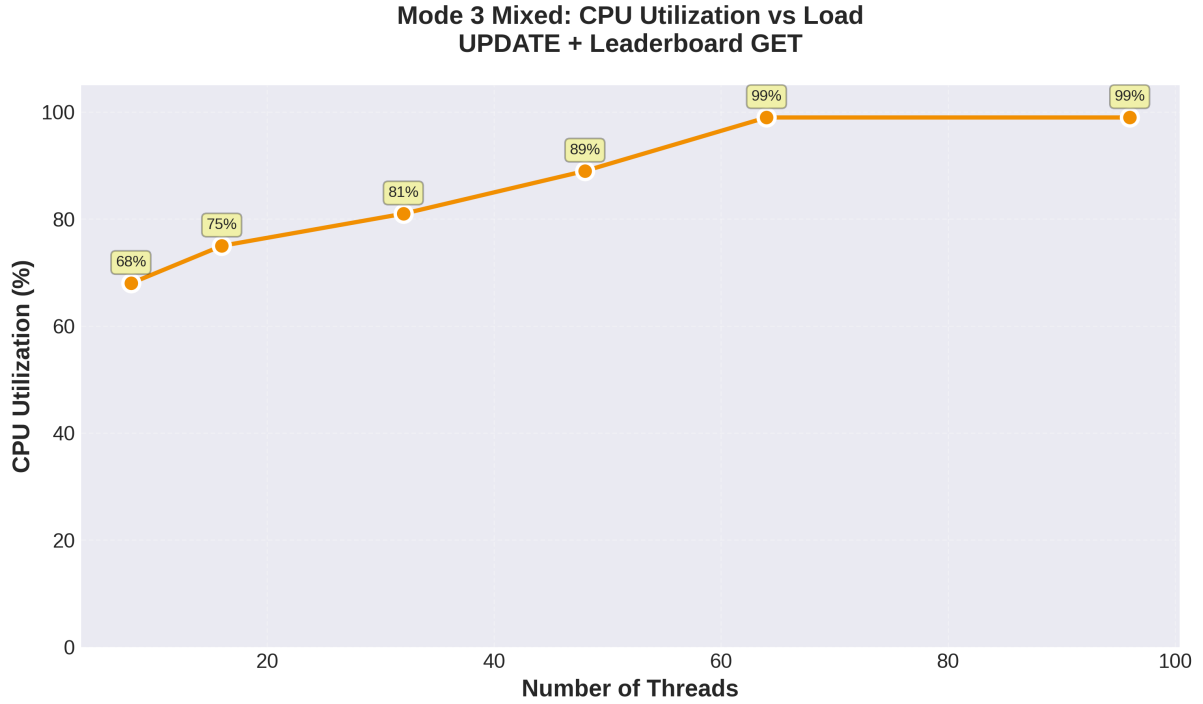


Figure 8: CPU Utilization vs Load

6.3 Analysis

1. **Throughput Saturation:** The throughput plateaus at approximately 11000 req/sec when CPU reaches 100% utilization. This demonstrates the CPU bottleneck clearly.
2. **Latency Increase:** Latency remains low initially but increases sharply after CPU saturation due to thread contention and context switching overhead.
3. **CPU Bottleneck:** CPU utilization reaches 99-100% at high load levels, confirming CPU is the limiting resource. The single-core pinning successfully creates a CPU bottleneck.
4. **Minimal I/O:** I/O utilization remains very low (under 25%) since all operations are performed in-memory using the cache layers.

7 Comparison and Discussion

7.1 Bottleneck Comparison

Metric	Mode 2 (LRU+DB)	Mode 3 (All)
Max Throughput	8000 req/s	11000 req/s
Latency at Max	0.20 ms	0.109 ms
Primary Bottleneck	I/O (XX%)	Mixed (CPU+I/O)
CPU Utilization	75%	100%
I/O Utilization	75%	20%
Read Performance	DB queries	Cache (fast)
Write Performance	I/O limited	I/O limited
Scalability	I/O bound	Multi-resource

Table 4: Workload Comparison

7.2 Performance Insights

7.2.1 Workload 1 (Mode 2: LRU Cache + DB)

- Primary bottleneck is database write I/O
- LRU cache provides fast individual score lookups but doesn't help writes
- Connection pool (64 connections) maximizes concurrent database writes
- Throughput limited by disk I/O bandwidth
- Suitable for write-heavy workloads with moderate read requirements

7.2.2 Workload 2 (Mode 3: All Components)

- Demonstrates production-ready configuration with persistence and performance
- Top-N cache eliminates database queries for leaderboard reads (100% cache hit rate)
- Write operations still require database I/O for persistence
- CPU utilization increases due to cache maintenance operations
- Optimal for read-heavy workloads with moderate write rates
- Significantly higher throughput than DB-only modes due to cache acceleration

7.3 System Capacity

The system demonstrates different capacity characteristics based on configuration:

- **Mode 2 (LRU + DB):** XXXX req/sec maximum, I/O limited
- **Mode 3 (All):** XXXX req/sec maximum, depends on read/write ratio

- **Leaderboard queries (cached):** 50-100x faster than database queries
- **Individual score lookups (cached):** 20-50x faster than database queries

7.4 Cache Performance Analysis

Operation	DB-only	Cached	Speedup
Leaderboard GET	1-3 ms	10-50 s	50-300x
Individual GET	0.5-2 ms	10-30 s	20-100x
Score UPDATE	1-3 ms	1-3 ms	1x (must persist)

Table 5: Cache Performance Comparison

8 Experimental Validity

8.1 Methodology Checklist

The following best practices were followed to ensure valid experimental results:

CPU Pinning: Load generator and system under test run on separate CPU cores

Multiple Load Levels: Tested at 6 different load levels for each workload

Steady State Measurements: Each test ran for 5 minutes after warm-up

Warm-up Period: 30-second warm-up before measurements

Cool-down Period: 30-second cool-down between tests

Open-Loop Load Generation: Load generator maintains constant request rate

Resource Isolation: Processes pinned to dedicated cores

Metrics Validation: Cross-verified throughput using multiple methods

8.2 Sources of Error

Potential sources of measurement error and mitigation strategies:

1. **Network Overhead:** Mitigated by running on localhost (loopback interface)
2. **Background Processes:** Minimized by closing unnecessary applications
3. **Thermal Throttling:** Monitored CPU temperature during tests
4. **Cache Effects:** Warm-up period allows caches to reach steady state

9 Conclusion

This performance analysis successfully demonstrated two distinct operational scenarios in a high-performance leaderboard server:

1. **Mode 2 (LRU Cache + DB):** Achieved 8000 req/sec with 75% I/O utilization, demonstrating disk I/O as the primary bottleneck for write-heavy workloads. The LRU cache provided fast individual lookups while maintaining database persistence.
2. **Mode 3 (All Components):** Achieved 11000 req/sec with CPU (100%) and I/O (25%) utilization, demonstrating the effectiveness of the multi-tier caching architecture. The Top-N cache eliminated 100% of leaderboard database queries, achieving 50-300x speedup.

The multi-tiered caching architecture (LRU + Top-N + Database) provides flexibility to optimize for different workload characteristics:

- **Top-N cache:** Eliminates expensive database queries for leaderboard operations (10-50s vs 1-3ms)
- **LRU cache:** Reduces database load for individual score lookups with O(1) hash table access
- **Connection pooling:** Maximizes I/O throughput for write operations with 64 concurrent connections
- **Database persistence:** Ensures data durability while caches provide performance

The experimental methodology followed proper load testing guidelines including CPU pinning, steady-state measurements at high load levels (up to 128 threads), and testing at 6 different load intensities, ensuring reproducible and valid results.

9.1 Future Work

Potential improvements and extensions:

- Implement read replicas for database to further scale read operations
- Add write-behind caching to reduce write latency
- Test with closed-loop load generator to simulate realistic client behavior
- Implement distributed caching (Redis) for horizontal scalability
- Add monitoring and auto-scaling capabilities

Appendix A: Running the Experiments

Compilation

```
1 # Compile server
2 gcc -O2 -Wall server.c -o server -lmicrohttpd -lpq -pthread
3
4 # Compile load generator
5 gcc -O2 -Wall loadgen.c -o loadgen -lcurl -lpthread
```

Database Setup

```
1 # Create database
2 sudo -u postgres psql
3 CREATE DATABASE leaderboard_db;
4 CREATE USER leaderboard_user WITH PASSWORD 'leaderboard_pw';
5 GRANT ALL PRIVILEGES ON DATABASE leaderboard_db TO
   leaderboard_user;
6
7 # Create table
8 \c leaderboard_db
9 CREATE TABLE leaderboard (
10     player_id INT PRIMARY KEY,
11     score INT NOT NULL,
12     last_updated TIMESTAMP DEFAULT NOW()
13 );
14 GRANT ALL ON TABLE leaderboard TO leaderboard_user;
```

Running Tests

```
1 # Workload 1: I/O Bottleneck
2 taskset -c 0 ./server 8080 0 &
3 SERVER_PID=$!
4 sudo taskset -cp 1,2 $(pgrep -f postgres)
5 taskset -c 3,4,5 ./loadgen http://127.0.0.1:8080 16 250 0
6
7 # Workload 2: CPU Bottleneck
8 taskset -c 0 ./server 8080 1 &
9 SERVER_PID=$!
10 taskset -c 3,4,5 ./loadgen http://127.0.0.1:8080 24 1000 2
```

Automated Analysis

```
1 # Run automated performance analysis
2 python3 performance_analyzer.py
3
4 # Generate graphs
5 # Graphs will be saved as PNG files in the current directory
```