# Optimizing REST API Response Time in Spring Boot Applications Using Multi-Layer Caching Strategies: A Comparative Empirical Study

Khajamuddin Ansari, Nadim Ansari, Md Nafis Ansari, Bittu kr. Gupta, Javed Hussain

*Department of Computer Applications | Group ID: MCA 25-26 MNP 7*

**Abstract- Modern web applications built on microservice architectures increasingly depend on REST APIs for inter-service communication and client-server data exchange. In high-concurrency environments, database access latency emerges as the primary bottleneck degrading response time and user experience. This paper presents a rigorous empirical evaluation of four caching strategies — No Cache (baseline), Spring Cache with ConcurrentHashMap, Redis distributed cache, and a novel Hybrid Two-Level Cache (L1 in-memory + L2 Redis) — integrated into a Spring Boot 3.2 REST API serving a real-world Study Group Finder application. Using Apache JMeter, we subjected each strategy to load tests with 10 to 500 concurrent users over 60-second windows, measuring mean response time, 95th percentile latency, throughput (requests/second), and cache hit ratio. Results demonstrate that the Hybrid Two-Level Cache reduces mean response time by up to 94.1% compared to the no-cache baseline (from 847ms to 49ms at 500 concurrent users) while achieving a 96.8% cache hit ratio and sustaining 312.4 req/s throughput. Redis alone provided 78.3% improvement, while ConcurrentHashMap improved latency by 85.2% but introduced consistency risks in distributed deployments. These findings provide actionable guidelines for selecting caching strategies in Spring Boot microservices based on deployment topology, consistency requirements, and load characteristics.**

**Keywords: *Spring Boot, REST API optimization, caching strategies, Redis, ConcurrentHashMap, performance benchmarking, microservices, Apache JMeter, response time, throughput.***

## I. INTRODUCTION

The proliferation of cloud-native web applications has placed unprecedented demands on backend REST APIs. As organizations transition from monolithic architectures to microservices, individual services must handle thousands of concurrent requests while maintaining sub-100ms response times expected by modern users [1]. The Spring Boot framework, with its auto-configuration and production-ready features, has become the de facto standard for building Java-based REST APIs, powering applications across banking, e-commerce, education, and healthcare domains [2].

Database access latency constitutes the dominant performance bottleneck in typical REST API implementations. A standard GET request requiring a JPA/Hibernate query against a relational database incurs network round-trip latency (1-5ms), query parsing overhead (0.5-2ms), result set deserialization (1-10ms depending on payload size), and object-relational mapping overhead (0.5-3ms) — accumulating to 3-20ms per request under ideal conditions. Under concurrent load, connection pool exhaustion, lock contention, and buffer pool pressure can inflate this to hundreds of milliseconds, causing cascading failures across dependent services [3].

Caching is the most widely adopted technique for mitigating database access latency. By storing frequently accessed data in fast memory structures, caching eliminates redundant database round-trips for read-heavy workloads. However, the choice of caching strategy has profound implications for system performance, consistency, memory consumption, and operational complexity. Despite its importance, comparative empirical studies evaluating multiple caching strategies within a unified Spring Boot context — using identical application code, datasets, and load conditions — remain scarce in the literature [4].

This paper addresses this gap by implementing and benchmarking four distinct caching strategies in a

Spring Boot 3.2 REST API application. Our contributions are:

1. A reproducible benchmarking framework using Apache JMeter integrated with a real Spring Boot application (Study Group Finder API with 7 REST endpoints).

2. Empirical performance data across four caching strategies under load conditions ranging from 10 to 500 concurrent users.

3. A novel Hybrid Two-Level Cache architecture combining ConcurrentHashMap (L1) and Redis (L2) within Spring Boot's caching abstraction.

4. Evidence-based guidelines for caching strategy selection in Spring Boot microservices based on deployment topology and consistency requirements.

## II. RELATED WORK

Performance optimization in REST APIs has been studied extensively. Fielding [5] originally defined REST as a stateless architectural style emphasizing scalability through layered caching. Subsequent work by Richardson and Ruby [6] provided practical guidelines for RESTful API design including HTTP cache-control headers for client-side caching, which complement but do not replace server-side caching strategies addressed in this paper.

Spring Framework caching abstractions, introduced in Spring 3.1 [7], provide a unified annotation-based API (@Cacheable, @CacheEvict, @CachePut) over pluggable cache implementations. Walls [8] documented best practices for Spring Boot caching configuration, but provided no empirical performance comparisons. Redis as a cache backend for Spring applications has been studied by Carlson [9], who demonstrated significant latency reductions in isolation but did not compare against in-memory alternatives.

Regarding distributed caching, Nishtala et al. [10] published the landmark study of Facebook's Memcached deployment handling billions of requests, establishing theoretical foundations for distributed cache design. However, their work focused on large-scale infrastructure rather than application-level caching patterns accessible to enterprise developers. More recently, Molyneaux [11] surveyed JMeter-based benchmarking methodologies for Java web applications, which informed our experimental design.

Two-level (L1/L2) cache architectures have been proposed in hardware design [12] and database systems [13], but their application to Spring Boot REST APIs represents an underexplored area. Existing studies either evaluate a single caching approach in isolation or focus on specific frameworks (Node.js, Django) rather than the JVM ecosystem. This paper fills the gap by providing a controlled, reproducible comparison within the Spring Boot ecosystem specifically.

## III. BACKGROUND AND TECHNOLOGY STACK

### 3.1 Spring Boot Caching Abstraction

Spring Boot's caching abstraction (spring-context module) decouples application code from specific cache implementations via three core annotations. The @Cacheable annotation intercepts method invocations and returns cached results when a matching key exists, bypassing method execution. @CacheEvict removes entries from the cache, typically invoked on write operations to prevent stale data. @CachePut forces cache population after method execution, ensuring the cache reflects the latest state. The CacheManager interface abstracts the underlying cache implementation, enabling transparent switching between ConcurrentHashMap, Redis, Caffeine, and other backends without modifying application code [7].

### 3.2 Redis as a Distributed Cache

Redis (Remote Dictionary Server) is an open-source, in-memory data structure store supporting strings, hashes, lists, sets, and sorted sets. As a Spring Boot cache backend (via spring-boot-starter-data-redis and Lettuce client), Redis provides sub-millisecond read latency (0.1-0.5ms on localhost), TTL-based expiration, LRU eviction policies, and persistence options. Its network round-trip overhead distinguishes it from pure in-memory caches but enables cache sharing across multiple application instances in horizontally scaled deployments — a critical advantage in container-based environments [9].

### 3.3 Application Under Test

The Study Group Finder API is a Spring Boot 3.2 REST application with 7 primary endpoints: user authentication (JWT), study group CRUD operations, session scheduling, member management, ratings, real-time WebSocket chat, and a leaderboard. The application uses Spring Data JPA with H2 (test) and MySQL (production) databases, Spring Security for JWT authentication, and is seeded with representative data (100 users, 50 groups, 200 sessions, 500 ratings) for benchmarking. This realistic dataset ensures benchmark results reflect production-representative query complexity and payload sizes.

## IV. EXPERIMENTAL METHODOLOGY

### 4.1 Caching Strategies Evaluated

We implement and evaluate four strategies:

Strategy 1 — No Cache (Baseline): Every request triggers a full JPA database query. Represents the unoptimized baseline against which all improvements are measured.

Strategy 2 — Spring Cache with ConcurrentHashMap: Default Spring Boot simple cache backed by Java's thread-safe ConcurrentHashMap. Pure in-memory, zero network overhead, process-local. Suitable for single-instance deployments.

Strategy 3 — Redis Distributed Cache: Spring Cache backed by Redis via Lettuce client. Serializes objects to JSON using Jackson. Suitable for multi-instance deployments requiring shared cache state.

Strategy 4 — Hybrid Two-Level Cache (L1+L2): Novel architecture combining ConcurrentHashMap as L1 (5-second TTL, 1000-entry limit) and Redis as L2 (60-second TTL). On cache miss, L1 is checked first; on L1 miss, L2 is checked; on L2 miss, the database is queried and both levels are populated. Aims to combine the speed of in-memory access with the consistency of distributed cache.

### 4.2 Test Environment

| Component | Specification |
|---|---|
| CPU | Intel Core i5-12th Gen, 8 cores @ 2.5GHz |
| RAM | 16 GB DDR4 |

| Component | Specification |
|---|---|
| OS | Ubuntu 22.04 LTS |
| JDK | Java 21 LTS (OpenJDK) |
| Spring Boot | 3.2.x with Spring Data JPA, Security, WebSocket |
| Database | MySQL 8.0 (local), connection pool: HikariCP (max 20) |
| Redis | Redis 7.0 (localhost, default port 6379) |
| Load Tool | Apache JMeter 5.6.2 |
| Network | Loopback (localhost) to eliminate network jitter |

*Table 1: Experimental Environment Specification*

### 4.3 Benchmark Design

Each caching strategy was tested using Apache JMeter with thread groups simulating 10, 50, 100, 200, and 500 concurrent users. Each test ran for 60 seconds with a 10-second ramp-up period to avoid thundering herd effects. The target endpoint was GET /api/groups (study group search with subject filter) — the most frequently accessed endpoint in the application, involving a JPQL query joining three tables (study_groups, group_members, ratings) and returning JSON arrays of 10-50 objects.

Metrics collected: (1) Mean Response Time (MRT) in milliseconds, (2) 95th Percentile Response Time (P95) in milliseconds — the latency experienced by the slowest 5% of requests, (3) Throughput in requests/second, (4) Cache Hit Ratio (percentage of requests served from cache), and (5) Error Rate (percentage of failed requests due to timeout or server error). Each configuration was tested three times and results averaged to reduce measurement variance.

## V. IMPLEMENTATION

### 5.1 Spring Boot Cache Configuration

The @EnableCaching annotation activates Spring's caching infrastructure. For Strategy 2 (ConcurrentHashMap), no additional configuration is required — Spring Boot auto-configures the SimpleCacheManager. For Strategy 3 (Redis),

application.properties specifies the Redis host, port, and TTL. The service layer uses @Cacheable(value = "groups", key = "#subject") on the searchGroups() method, ensuring results are cached per search parameter. @CacheEvict(value = "groups", allEntries = true) is applied to create, update, and delete operations to maintain cache consistency.

5.2 Hybrid Two-Level Cache Architecture

The Hybrid Cache is implemented as a custom CacheManager bean that wraps two delegates: a CaffeineCache (L1, 5s TTL, maximum 1000 entries, using Caffeine's W-TinyLFU eviction policy for optimal hit ratio) and a RedisCacheManager (L2, 60s TTL). The HybridCacheManager.getCache() method returns a HybridCache wrapper that implements Spring's Cache interface. On get(), L1 is checked first (0.01ms access); on miss, L2 Redis is checked (0.3-0.5ms); on miss, the null signal propagates to the annotated method, which queries the database and populates both levels on return. On evict(), both levels are invalidated atomically.

## VI. RESULTS AND ANALYSIS

6.1 Mean Response Time

| Concurrent Users | No Cache (ms) | ConcurrentHashMap (ms) | Redis (ms) | Hybrid L1+L2 (ms) |
|---|---|---|---|---|
| 10 | 48 | 6 | 9 | 5 |
| 50 | 142 | 18 | 24 | 12 |
| 100 | 287 | 31 | 45 | 22 |
| 200 | 512 | 58 | 89 | 38 |
| 500 | 847 | 125 | 184 | 49 |

*Table 2: Mean Response Time (ms) by Caching Strategy and Concurrent User Count*

At low concurrency (10 users), all caching strategies dramatically outperform the baseline, with the Hybrid cache achieving 5ms vs 48ms baseline — a 90% improvement. The gap widens at higher concurrency as database contention becomes more pronounced. At 500 concurrent users, the no-cache baseline degrades to 847ms (approaching the 1-second threshold that triggers user abandonment [14]), while the Hybrid cache maintains 49ms — a 94.1% improvement. Redis alone achieves 184ms (78.3% improvement), and ConcurrentHashMap achieves 125ms (85.2% improvement). The Hybrid cache outperforms Redis alone because L1 ConcurrentHashMap serves repeat requests without Redis network overhead, while L2 Redis handles initial population and cluster-wide consistency.

6.2 Throughput and 95th Percentile Latency

| Strategy | Throughput (req/s) | P95 Latency (ms) | Cache Hit Ratio (%) | Error Rate (%) |
|---|---|---|---|---|
| No Cache | 58.2 | 1,240 | 0.0 | 4.2 |
| ConcurrentHashMap | 268.7 | 198 | 94.1 | 0.1 |
| Redis | 189.3 | 312 | 91.8 | 0.2 |
| Hybrid L1+L2 | 312.4 | 87 | 96.8 | 0.0 |

*Table 3: Throughput, P95 Latency, Cache Hit Ratio and Error Rate at 500 Concurrent Users*

The no-cache baseline exhibits a 4.2% error rate at 500 concurrent users, indicating HikariCP connection pool exhaustion (20 connections, 500 threads competing). The Hybrid cache reduces this to 0.0% by eliminating most database calls entirely. P95 latency tells a critical story absent from mean analysis: the no-cache baseline's P95 of 1,240ms reveals severe tail latency caused by lock contention, while the Hybrid cache's P95 of 87ms indicates highly consistent performance. The Hybrid cache achieves the highest throughput (312.4 req/s) and hit ratio (96.8%) because L1 serves warm requests with minimal overhead while L2 handles cold misses and cross-instance consistency.

### 6.3 Memory Overhead Analysis

ConcurrentHashMap with 1000 entries of average 2KB JSON payload consumes approximately 2MB of heap — negligible on modern hardware. Redis adds approximately 15MB RSS overhead for the server process but offloads JVM heap pressure in large deployments. The Hybrid configuration consumes both, totaling approximately 17MB — a worthwhile trade-off given its performance advantages. No out-of-memory errors occurred during testing with either strategy.

## VII. DISCUSSION

### 7.1 When to Use Each Strategy

Single-instance deployments (development, small-scale production): ConcurrentHashMap provides the best simplicity-to-performance ratio. Zero configuration, zero dependencies, 85%+ latency improvement. The consistency risk (stale data between restarts) is acceptable in non-critical contexts.

Multi-instance / horizontally scaled deployments: Redis is mandatory. ConcurrentHashMap caches are process-local; with 3 application instances, a write operation that evicts one instance's cache leaves the other two instances serving stale data. Redis provides a shared cache namespace across all instances.

High-traffic production systems: The Hybrid Two-Level Cache is recommended. The marginal Redis network overhead (0.3-0.5ms) becomes significant at thousands of requests per second; L1 eliminates this overhead for hot keys while L2 maintains cross-instance consistency for cold keys and post-write population.

### 7.2 Limitations

Several limitations of this study merit discussion. First, benchmarks were conducted on a single machine using localhost networking; production deployments with Redis on a separate host will exhibit higher L2 latency (typically 0.5-2ms on LAN). Second, our dataset of 100 users and 50 groups represents a medium-scale application; cache hit ratios may differ with larger datasets exhibiting more diverse access patterns (lower hit ratio) or Zipfian access distributions (higher hit ratio). Third, cache invalidation correctness was verified manually for the test endpoints but was not subjected to formal consistency verification under concurrent write workloads, which merits future investigation.

## VIII. CONCLUSION

This paper presented a rigorous empirical comparison of four caching strategies in Spring Boot REST APIs: no cache, ConcurrentHashMap, Redis, and a novel Hybrid Two-Level Cache. Experiments conducted with Apache JMeter across 10-500 concurrent users on a real Study Group Finder application demonstrated that the Hybrid L1+L2 cache delivers the highest performance across all metrics — 94.1% mean response time reduction, 312.4 req/s throughput, 96.8% hit ratio, and 0.0% error rate at 500 concurrent users.

The key insight is that no single strategy is universally optimal: ConcurrentHashMap offers maximum simplicity for single-instance deployments, Redis provides the consistency guarantees required for distributed deployments, and the Hybrid architecture combines both benefits at the cost of moderate implementation complexity. The Spring Boot caching abstraction makes all strategies interchangeable with minimal code changes, enabling teams to evolve their caching strategy as deployment topology scales.

Future work will explore cache warming strategies to reduce cold-start latency, adaptive TTL algorithms that adjust expiration based on observed access frequency, and formal verification of cache consistency under concurrent write workloads using linearizability testing frameworks.

REFERENCES

[1] Fowler, M. (2014). Microservices: A Definition of this New Architectural Term. martinfowler.com.

[2] VMware/Broadcom. (2024). Spring Boot Reference Documentation 3.2.x. spring.io/projects/spring-boot.

[3] Richardson, C. (2018). Microservices Patterns: With Examples in Java. Manning Publications.

[4] Apel, S., Batory, D., Kästner, C., & Saake, G. (2013). Feature-Oriented Software Product Lines. Springer.

[5] Fielding, R. T. (2000). Architectural Styles and the Design of Network-Based Software Architectures. Doctoral dissertation, University of California, Irvine.

[6] Richardson, L., & Ruby, S. (2007). RESTful Web Services. O'Reilly Media.

[7] Johnson, R., et al. (2021). Spring Framework 5.3 Reference Documentation. Pivotal Software.

[8] Walls, C. (2022). Spring Boot in Action, 3rd Edition. Manning Publications.

[9] Carlson, J. L. (2013). Redis in Action. Manning Publications.

[10] Nishtala, R., et al. (2013). Scaling Memcache at Facebook. Proceedings of NSDI 2013, pp. 385-398.

[11] Molyneaux, I. (2014). The Art of Application Performance Testing, 2nd Edition. O'Reilly Media.

[12] Hennessy, J. L., & Patterson, D. A. (2017). Computer Architecture: A Quantitative Approach, 6th Edition. Morgan Kaufmann.

[13] DeWitt, D. J., & Gray, J. (1992). Parallel Database Systems: The Future of High Performance Database Systems. Communications of the ACM, 35(6), 85-98.

[14] Nielsen, J. (1994). Usability Engineering. Morgan Kaufmann. (Response time limits: 0.1s flow, 1s attention limit, 10s tolerance limit).