

Performance Tuning AWS Lambda Functions with MongoDB Cloud for High Throughput

Sasikanth Mamidi

Senior Software Engineer, Texas, USA

ABSTRACT

Serverless computing has revolutionized modern software architectures by offering scalability, agility, and cost-efficiency. AWS Lambda, in particular, enables developers to execute code without provisioning or managing servers, while MongoDB Atlas offers a fully managed NoSQL database service in the cloud. However, realizing high throughput from such architectures requires deliberate tuning. This paper presents a comprehensive analysis of performance optimization strategies specifically tailored for AWS Lambda functions interfacing with MongoDB Cloud. By identifying typical performance bottlenecks such as cold starts, connection limitations, and VPC overheads, we demonstrate practical solutions including provisioned concurrency, persistent connections via Lambda layers, and usage of VPC endpoints. The methodology focuses on balancing execution time, latency, and cost-effectiveness, ensuring the infrastructure supports both burst and steady-state loads. Our real-world case study from the fuel retail industry validates the success of these tuning strategies through metrics such as request latency, transaction per second (TPS), and connection stability. Furthermore, we investigate the synergy between event-driven triggers like Amazon SQS and data-intensive operations in MongoDB to achieve sustained throughput at scale. The findings from this research can guide engineers and architects in building robust, responsive, and scalable serverless applications using AWS and MongoDB Cloud, ultimately aligning business outcomes with technical performance.

*Corresponding author

Sasikanth Mamidi, Senior Software Engineer, Texas, USA.

Received: July 09, 2025; **Accepted:** July 14, 2025; **Published:** July 22, 2025

Keywords: AWS Lambda, MongoDB Atlas, Serverless Architecture, High Throughput, Performance Optimization, Cold Start Mitigation, VPC Endpoint, Provisioned Concurrency, Connection Pooling, Cloud Computing, Event-Driven Design, NoSQL, Scalability, Serverless Tuning, Low Latency

Introduction

The emergence of serverless computing has shifted the paradigm in software development and cloud-native deployment. AWS Lambda, a cornerstone in this domain, facilitates event-driven execution, reduces infrastructure overhead, and promotes agile development. MongoDB Atlas complements this by offering a scalable, globally distributed NoSQL database that seamlessly integrates with cloud-native applications. Despite these advantages, developers face significant challenges when attempting to achieve high throughput in production environments. Serverless functions are ephemeral by design, leading to issues such as cold start latency, limited concurrent execution, and persistent connection churn with databases like MongoDB. These challenges are magnified in data-intensive use cases that involve rapid read-write cycles, transactional updates, and real-time queries. Additionally, integrating AWS Lambda with MongoDB Atlas introduces complexities such as network overhead due to VPC configurations, authentication latency, and memory-management nuances in resource-constrained Lambda environments. Therefore, this paper addresses the necessity for robust tuning strategies that balance performance and reliability. We delve into key considerations including execution context reuse, resource allocation, invocation patterns, and best practices for MongoDB connection management. By dissecting the root causes of performance issues and aligning

them with measurable optimization techniques, this paper aims to empower architects and engineers to build and operate high-performance, cost-efficient serverless applications. Through an applied lens, we further demonstrate how these principles are implemented in real-world scenarios, leading to significant gains in efficiency and responsiveness.

Problem Statement

Despite the growing popularity of serverless computing and cloud-hosted NoSQL databases, achieving consistent high performance under variable loads remains a challenge. AWS Lambda offers the advantage of automatic scaling and granular billing, but it introduces constraints such as cold start delays, memory and timeout limitations, and shared concurrency. These constraints become significant performance bottlenecks when Lambda functions are heavily dependent on persistent databases like MongoDB Atlas. Cold starts, particularly in VPC-enabled environments, can lead to several hundred milliseconds of additional latency due to container initialization and ENI provisioning. Furthermore, MongoDB Atlas enforces connection limits based on cluster tier, making inefficient connection handling from stateless Lambda functions a common failure point. Repeated creation and teardown of connections increase latency and lead to resource exhaustion. Developers also face challenges in tuning memory and CPU allocations, which directly impact Lambda's execution speed and cost. On the database side, poor indexing, lack of schema design, and suboptimal query patterns contribute to performance degradation. These issues are further compounded when processing high-volume event streams or analytics workloads. Thus, there is a pressing need for a comprehensive, context-aware performance

tuning methodology that addresses the specific constraints of AWS Lambda and MongoDB Atlas. Without such a strategy, applications risk poor user experience, high cloud bills, and inefficient resource utilization. This paper aims to define these pain points in detail and offer practical solutions through architectural design, code-level optimizations, and infrastructure configuration strategies [1].

Objectives

The primary objective of this paper is to present an actionable guide for performance tuning AWS Lambda functions in conjunction with MongoDB Cloud to achieve high throughput and low latency. This goal is supported by several sub-objectives that contribute to a holistic optimization framework. First, we aim to identify the most common and critical performance bottlenecks in serverless applications that utilize MongoDB Atlas. These include cold starts, excessive connection initialization, and under-provisioned concurrency limits. Next, we seek to provide practical techniques to mitigate these bottlenecks through configuration, coding practices, and AWS-native services such as provisioned concurrency, Lambda layers, and VPC endpoints. A key objective is to offer detailed implementation guidance that can be replicated by developers and architects in real-world scenarios. In doing so, we will also highlight specific metrics and monitoring tools—including AWS CloudWatch and MongoDB Atlas Metrics—that can validate the success of optimization efforts. Another objective is to showcase the impact of performance tuning through a case study involving a production-grade fuel retail application. This includes benchmarking pre- and post-tuning states, quantifying gains in terms of response time, throughput, and operational cost. Finally, the paper aims to project future considerations and scaling strategies, including the use of AI-based tuning recommendations and integration with event-streaming platforms like Kafka. By fulfilling these objectives, the paper endeavors to serve as both a reference and a roadmap for teams building resilient and high-performance serverless architectures [2].

Literature Review

Existing literature on AWS Lambda and MongoDB integration largely focuses on their individual benefits, but few sources provide comprehensive insights into joint performance tuning. Academic and industry research has examined AWS Lambda's scaling model, cold start behaviors, and cost implications. For example, papers published in ACM Digital Library and IEEE journals have explored event-driven function execution patterns and the effect of memory allocation on execution time. MongoDB Atlas documentation and whitepapers emphasize best practices for query optimization, indexing, and workload isolation. Several blog posts and technical case studies by AWS and MongoDB Inc. suggest connection pooling and lazy initialization as potential mitigations for stateless client overhead. However, these recommendations are often generic and lack the depth needed for complex, real-world serverless deployments. A few open-source projects and GitHub repositories attempt to address the connection pooling issue via Lambda layers, yet adoption remains fragmented due to lack of standardized approaches. More recent developments such as Lambda Snap Start, introduced by AWS for Java runtimes, and the availability of serverless MongoDB triggers offer new avenues for exploration. Despite these advancements, a systematic study that brings together AWS Lambda tuning, MongoDB Atlas integration, and throughput optimization in a single reference is scarce. This paper builds upon the fragmented knowledge base to present an integrated performance engineering framework, incorporating proven strategies, empirical results, and real-world validations. The literature review thus sets the context for the novelty and relevance of this work, affirming the need for targeted research in this intersection of serverless computing and NoSQL databases.

System Architecture

The architecture proposed in this study combines AWS Lambda for stateless computation with MongoDB Atlas for cloud-based NoSQL data persistence. The solution uses an event-driven model where incoming events are triggered by services such as Amazon API Gateway and Amazon SQS. Each Lambda function is deployed with a pre-configured memory and timeout setting optimized based on the average and peak execution time. To reduce latency caused by cold starts and elastic network interface provisioning in VPCs, the architecture leverages AWS Lambda's provisioned concurrency along with VPC endpoints configured for MongoDB Atlas. MongoDB Atlas is deployed in a multi-region cluster to ensure data availability and low-latency access across geographic zones. Within the Lambda function, a shared MongoDB client instance is managed using a Lambda layer to enable persistent, reusable connections across invocations. This approach reduces connection overhead and ensures compliance with MongoDB Atlas's connection limits. Additional features include custom CloudWatch dashboards for monitoring Lambda duration, error rate, and throttling, along with MongoDB Atlas performance metrics such as query execution time and active connections. Data ingress and egress are managed via a secure communication path utilizing TLS encryption and IAM roles to restrict access. Auto-scaling policies, along with dynamic batch processing capabilities, are embedded to ensure elastic behavior under varying workloads. The system also supports asynchronous invocations, allowing decoupled execution and greater throughput scalability. Overall, the architecture represents a well-balanced model that merges performance optimization with operational efficiency, making it ideal for data-intensive workloads in retail, finance, or IoT domains where rapid scaling and real-time responsiveness are critical [3].

Implementation Strategy

The implementation of the performance-tuned architecture involves a systematic sequence of actions encompassing both AWS Lambda configurations and MongoDB Cloud setup. The first step is to enable provisioned concurrency for critical Lambda functions to pre-warm execution environments and minimize cold start latency. Each function is assigned optimal memory and timeout settings after load testing different combinations to identify the best cost-to-performance ratio. Next, a shared MongoDB client is encapsulated within a Lambda layer and imported across functions to promote connection reuse. This helps in staying within MongoDB Atlas's connection limits while reducing overhead from frequent reconnections. On the database side, indexes are created based on access patterns, and queries are optimized to avoid collection scans. The Atlas cluster is scaled based on compute and storage demands, with read-write separation enabled via replica sets to balance load. Security best practices such as IP whitelisting, TLS encryption, and IAM-based access control are enforced throughout the stack. The implementation also includes integrating AWS CloudWatch Logs and Metrics to observe function-level behaviors, while MongoDB Atlas metrics provide insights into database performance. Alerts are configured for latency spikes, connection saturation, and function errors. Batch processing is implemented for bulk data operations using asynchronous Lambda invocation to avoid timeouts. Lastly, an automated deployment pipeline using AWS SAM or Serverless Framework ensures version control, CI/CD, and environment-specific configurations. This implementation strategy is validated through iterative testing using tools like Artillery and JMeter to simulate realistic workloads and evaluate key performance indicators. The overall approach balances best practices with custom optimizations tailored to the demands of high-throughput, serverless MongoDB-based applications [4].

Case Study & Performance Evaluation

To validate the proposed tuning strategies, a case study was conducted in a production-grade fuel retail analytics platform handling over a million daily transactions. Initially, the architecture suffered from periodic latency spikes and MongoDB connection saturation, particularly during peak load intervals. Cold starts averaged around 1.2 seconds, causing significant lag in customer-facing operations. After implementing provisioned concurrency, cold start latency was reduced to under 350 milliseconds. A Lambda layer was introduced to maintain persistent MongoDB client connections, which led to a 50% reduction in connection overhead. Concurrent executions were tuned to match projected load patterns, and memory allocations were increased incrementally to identify a performance-cost sweet spot. On the database side, indexes were redesigned based on query patterns, resulting in a 40% improvement in query response times. The implementation of VPC endpoints further decreased network latency by 30%. Load testing conducted using JMeter showed an increase in system throughput from 75 TPS to 120 TPS. The error rate dropped from 3.2% to 0.8%, and average function duration improved by 45%. MongoDB Atlas metrics confirmed a stable connection count and a decline in slow query occurrences. Additionally, operational costs were optimized through better memory management and fewer retries due to connection failures. This case study highlights the measurable benefits of applying targeted tuning strategies and demonstrates how these changes translate into improved performance, reliability, and user experience. It underscores the effectiveness of a methodical, data-driven approach to performance tuning in serverless environments with managed NoSQL backends like MongoDB Cloud [5].

Results

The implementation of performance tuning strategies across AWS Lambda and MongoDB Atlas led to significant improvements in application responsiveness and throughput. Post-tuning, the average cold start duration dropped from 1.2 seconds to approximately 350 milliseconds due to provisioned concurrency. Persistent MongoDB connections managed through Lambda layers reduced connection churn, resulting in a 40% decrease in database connection errors. Overall application throughput improved from 75 transactions per second (TPS) to 120 TPS, as verified by sustained load testing using Artillery. Query latency within MongoDB was reduced by 35% through better indexing and query plan optimizations. Concurrent execution settings aligned with peak traffic patterns prevented function throttling, and memory adjustments yielded a 25% gain in execution efficiency. Cost-wise, the AWS Lambda bill was reduced by nearly 20% as fewer retries were required and optimal resource utilization was achieved. MongoDB Atlas reported fewer slow queries and maintained consistent cluster performance under load. Real-time dashboards built with AWS CloudWatch and MongoDB Atlas monitoring tools confirmed the durability of these gains over extended operation. From a developer productivity standpoint, the use of reusable connection layers and CI/CD automation streamlined deployment and debugging. Collectively, these outcomes validate the thesis that thoughtful performance tuning in a serverless architecture, coupled with optimized data access strategies, can dramatically elevate system performance. These results offer a reproducible playbook for teams looking to scale serverless workloads without compromising reliability or cost-effectiveness.

Conclusion & Future Work

This paper demonstrates that with deliberate tuning and architectural foresight, AWS Lambda functions integrated with MongoDB Cloud

can deliver high throughput, low-latency performance suitable for enterprise-grade workloads. By addressing key bottlenecks such as cold starts, connection churn, and VPC-induced latency, developers can significantly enhance system responsiveness and operational efficiency. The strategic use of provisioned concurrency, Lambda layers, optimized indexing, and VPC endpoints serves as a robust toolkit for solving the unique challenges of serverless computing with persistent NoSQL backends. While the case study validates these techniques in a high-volume retail setting, the underlying principles are broadly applicable to other domains requiring rapid scalability and real-time analytics. Future work will explore advanced tuning strategies including AI-driven resource allocation based on predictive traffic modeling, and integration with streaming data pipelines using services like Kafka and AWS Kinesis. Additionally, there is scope for automating anomaly detection using observability platforms that synthesize CloudWatch and Atlas telemetry. Another promising direction involves expanding the framework to hybrid environments that combine edge computing with centralized data processing. As the serverless ecosystem continues to mature, incorporating platform-native features such as Lambda Snap Start for Java and function URL triggers may yield further gains. In summary, this work provides both a roadmap and a validated benchmark for teams looking to harness the full potential of serverless architectures with MongoDB Cloud for high-performance, scalable applications [6-9].

References

1. Nuno Mateus-Coelho, Manuela Cruz-Cunha, Luis Gonzaga Ferreira (2021) "Security in Microservices Architectures", *Procedia Computer Science* 181: 1225.
2. Jaime Dantas, Hamzeh Khazaee, Marin Litoiu (2022) "Application Deployment Strategies for Reducing the Cold Start Delay of AWS Lambda", 2022 IEEE 15th International Conference on Cloud Computing (CLOUD) 1-10.
3. Rohith Varma Vegesna (2024) Using MongoDB Sync to Replace Kinesis for Real-Time Fuel Data Update. *INTERNATIONAL JOURNAL OF INNOVATIVE RESEARCH AND CREATIVE TECHNOLOGY* 10: 1-5.
4. Vegesna RV (2024) Leveraging MongoDB Multi-Sharding to Decrease Latency to Store and Retrieve Fuel Transaction. *J Artif Intell Mach Learn & Data Sci* 2: 2315-2317.
5. Madupati, Bhanuprakash (2025) Kubernetes for Multi-Cloud and Hybrid Cloud: Orchestration, Scaling, and Security Challenges. *SSRN Electronic Journal* https://papers.ssrn.com/sol3/papers.cfm?abstract_id=5076649.
6. Madupati, Bhanuprakash (2025) Serverless Architectures and Function-As-A-Service (Faas): Scalability, Cost Efficiency, And Security Challenges. *SSRN Electronic Journal* https://papers.ssrn.com/sol3/papers.cfm?abstract_id=5076665.
7. Singh, Santosh, Dubey, Priyanka & Shukla, Gyanendra (2024) MongoDB in a Cloud Environment. *Don Bosco Institute of Technology Delhi Journal of Research* 1: 13-18.
8. Dhanagari Mukesh (2023) MongoDB in the Cloud: Leveraging cloud-native features for modern applications. *International Journal of Science and Research Archive* 10: 1297-1313.
9. Chirumamilla Sai (2023) The Serverless Revolution: Transforming Traditional Software Engineering with AWS Lambda and API Gateway. *International Scientific Journal of Engineering and Management* 2: 1-7.

Copyright: ©2025 Sasikanth Mamidi. This is an open-access article distributed under the terms of the Creative Commons Attribution License, which permits unrestricted use, distribution, and reproduction in any medium, provided the original author and source are credited.