

Evolving High-Volume Systems: Reactive Execution Models for Resilient Operations

Kishore Subramanya Hebbar*

International Business Machines, Atlanta, USA

*Corresponding author

Email address: hebbar.kishore@gmail.com (Kishore Subramanya Hebbar)

Abstract

Conventional software architectures often struggle to scale effectively because of the variability in task demands and rising traffic volumes, which challenge synchronous processing methods. These load scenarios can involve sudden spikes in traffic followed by downstream slowdowns. Typically, these platforms do not have a straightforward or quick response. Most of the earlier literature focused on reactivity at an abstract and overarching level. Unfortunately, there is a gap in practical systems when it comes to documenting the transformations suitable for synchronous pipelines. The main goal of this project is to address this gap by proposing a methodology that covers everything from the design phase to implementation, utilizing both synchronous and reactive architectures within a unified application framework. The approach typically starts with a synchronous service that uses blocking input and output, executing a straightforward high-throughput workflow with tightly coupled steps. The same workflow could be redesigned into non-blocking, event-driven processing through which each stage is under explicit flow control and all stages are sufficiently isolated from one another. Both implementations share identical interfaces, data models, and persistence layers to enable a fair comparison. Multiple controlled load and stress tests are performed to observe the behavior of the respective paradigms under both continuous and burst traffic conditions. The results show that a reactive design cuts the tail latency by 45% under peak loads, ensures higher throughput with fewer threads, and degrades more gracefully under downstream slowness. The study highlights that using a reactive platform can be an effective strategy for moving large transaction-heavy and latency-critical systems, resulting in improved scalability, durability, and operational stability through fundamental architectural changes rather than incremental tuning.

Keywords: Application modernization, Reactive architecture, High-volume systems, Non-blocking processing, Event-driven design, System scalability, Latency optimization

1. Introduction

Over the past decade, software systems have evolved into more extensive networks that link a greater number of devices and accommodate increased operational demands. The applications which used to process predictable request traffic now need to deliver constant service while handling unexpected traffic increases from users who access the system across the world [1]. Organizations update their outdated systems to enable real-time decision-making and event-driven system operations and cloud-native deployment methods. The present system architecture design establishes performance and stability features, which developers use to create their systems. The request processing system together with resource allocation methods and execution control mechanisms determine whether systems achieve smooth scaling or experience operational failures during high-demand periods [2, 3]. As the amount of traffic rises, numerous production systems start to reveal their shortcomings in managing operations that depend on synchronous blocking processes. The traditional pipeline system uses thread-per-request handling together with tightly coupled service calls, which creates three problems that include thread exhaustion, unpredictable latency, and cascading failures that occur when downstream components experience performance issues [4]. The problems become most evident in environments which process large amounts of data, including transaction processing platforms and data ingestion services and real-time analytics systems. The temporary solution of infrastructure scaling generates the illusion of problem resolution, yet it fails to solve the core architectural limitations which restrict system throughput and system resilience. Instead of seeking new execution flow techniques, the teams continue to focus on tweaking timeouts, enlarging thread pools, and adding hardware components. The existing literature provides minimal practical guidance on converting complete synchronous pipelines into reactive systems because it either maintains a conceptual framework or focuses on specific frameworks [5]. The lack of comparative data in performance assertions makes it hard for practitioners to judge the right moment for making architectural changes instead of pursuing small-scale optimizations. There is a distinct gap between theoretical understanding and practical implementation. What is missing is a concrete, system-level examination of how synchronous architectures behave under sustained and burst traffic, how reactive redesign alters

execution characteristics, and what measurable benefits result from this shift. The existing papers provide insufficient direction about establishing architectural redesign processes which separate architectural effects from changes in business logic. Lacking this clarification causes modernization projects to become unstructured code tweaks rather than carefully planned architectural upgrades supported by confirmed information. This research paper examines the current gap in knowledge by providing a comprehensive comparison between synchronous and reactive system designs that function within the same high-capacity processing environment. This paper presents a redesign method which uses design patterns to maintain existing system functions while switching from blocking execution to non-blocking event-driven execution. The dual architecture implementation with matching interfaces and data models and system workloads enables us to test system performance and scalability and failure behavior under controlled conditions. The study demonstrates through concrete cases how to identify scenarios where reactive platforms lead to better stability and throughput, confirming that reactive architecture is an effective method for system modernization in high-volume operational contexts [5, 6].

2. Methodology

The study presents a novel approach that treats reactive architecture as a measurable system-level redesign instead of a programming abstraction. The methodology establishes performance benchmarks through testing complete synchronous systems and complete reactive systems which operate at high-volume capacity while using controlled design elements and matching interfaces and standardized operating conditions to assess architectural effects [7].

2.1. Baseline Synchronous Architecture Design

The first methodology step establishes a system baseline through implementation of a synchronous processing model which matches current operational standards used by most existing production systems. The system architecture operates through a request-driven pipeline system which processes incoming requests with dedicated threads while performing downstream tasks through blocking operations. The system executes service-to-service communication together with database access and external calls through a synchronous execution process which creates tightly integrated stages that control the request lifecycle. The baseline architecture exists as a complete representation of real-world systems instead of an oversimplified model. The system processes multiple stages which include request validation business logic execution persistence operations and downstream service interaction. The organization establishes thread pools according to standard operational methods while it sets timeouts at safe values which match actual system reliability needs. The performance characteristics which we observe during tests show actual architectural behavior because no artificial constraints were present. The implementation of this system through synchronous mode provides two distinct benefits. The first purpose establishes a working system which enables comparison of the reactive redesign process. The second purpose of the system shows common failure patterns which occur when the system faces high demand through thread saturation and increased queueing time and decreased throughput during downstream slowdowns. The current system does not use any optimizations for reactive execution because the goal is to record how standard blocking pipelines function when they handle large workloads. Figure 1 illustrates the baseline synchronous architecture, where each request is processed through a blocking, thread-per-request pipeline with tightly coupled execution stages. [13]

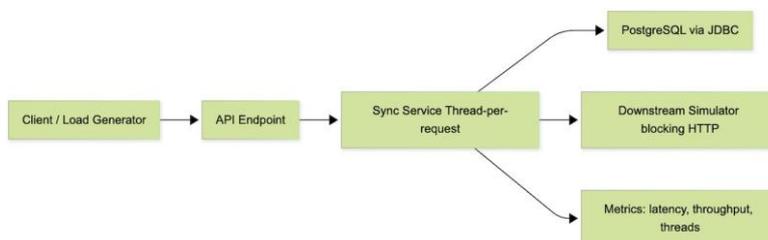


Figure 1: Synchronous Architecture Diagram.

2.2. Reactive Platform Redesign and Pattern Application

The second phase of the methodology requires the creation of a new workflow design which uses reactive principles to maintain the same functional outcomes. The goal is to create an execution model which enables better system performance through increased system capacity and enhanced system reliability without introducing new features or modifying existing business rules. The study maintains all external interfaces alongside request payloads and response semantics because these elements need to remain intact for accurate evaluation. The reactive redesign transforms processing operations by replacing all blocking functions with systems that use non-blocking event-driven processes.

The system processes incoming requests through asynchronous pipelines which execute tasks based on incoming data instead of relying on specific threads to control the process. The system uses reactive streams to separate processing stages which enables it to control data flow while applying backpressure when its downstream components experience delays or complete breakdowns. Architectural patterns serve as the main design framework for this redesign process [8]. The system creates event-driven boundaries which separate execution stages to stop delays from spreading throughout the entire system. The system implements flow control mechanisms which enable resource protection while establishing asynchronous error handling paths to prevent system-wide failures. The system establishes consistent pattern application throughout the pipeline which creates predictable system performance during times of uneven system demand. The system redesign prohibits the use of framework-specific shortcuts and proprietary optimizations. By concentrating on the structural aspects of architecture and excluding tooling, the research aims to produce universally applicable insights for reactive system implementations. The platform demonstrates how a production system can undergo modernization through incremental architectural changes without requiring a total system rewrite. Figure 2 presents the reactive system architecture, highlighting non-blocking execution, decoupled processing stages and explicit flow control between components.

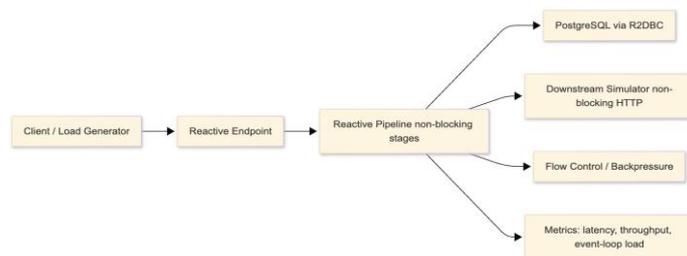


Figure 2: Reactive System Architecture Diagram.

2.3. Implementation Strategy and Code Structure

Both the synchronous and reactive systems are implemented using the same data models, persistence layer and deployment environment. This alignment is critical to ensuring that differences in behavior can be attributed to architectural choices rather than infrastructure or configuration variance. Shared components, such as data access logic and validation rules, are reused wherever possible to reduce divergence. Code structure is intentionally organized to highlight execution flow rather than framework mechanics. In the synchronous version, control flow follows a linear, call-driven path, making thread ownership implicit and tightly coupled to request handling. In the reactive version, execution is expressed as a sequence of asynchronous stages, with explicit transitions between processing steps and clear separation between computation and input or output operations. To make these execution differences explicit, Algorithms 1 and 2 present algorithmic representations of the synchronous and reactive processing models. Rather than focusing on framework-specific syntax, the algorithms emphasize control flow, thread ownership, and decision points that influence scalability and failure behavior under load. Minimal but representative code snippets are extracted to illustrate these differences. Rather than presenting full implementations, the methodology focuses on execution boundaries, error propagation paths, and resource usage patterns. This approach keeps the discussion accessible while still grounding the analysis in concrete implementation details. Special attention is given to observability. Both systems are instrumented to capture metrics related to latency distribution, throughput, thread usage, and error rates. Logging and tracing are configured consistently to allow comparable visibility into execution behavior. This ensures that performance observations are supported by measurable signals rather than inferred outcomes.

Algorithm 1 Synchronous Pipeline Execution

- 1: Receive incoming request
- 2: Assign a dedicated worker thread
- 3: Record request start time
- 4: **if** request payload is invalid **then**
- 5: Release worker thread 6: Return error response 7: **end if**
- 8: Perform blocking database read
- 9: **if** database call times out or fails **then**
- 10: Release worker thread 11: Return failure response 12: **end if**
- 13: Invoke downstream service using blocking call

14: **if** downstream service is slow or unavailable **then**
15: Worker thread remains blocked
16: Queue incoming requests
17: **end if**
18: Process combined response
19: Perform blocking database write 20: Record request completion time 21: Release worker thread
22: Return response to client

2.4. *Load Testing, Measurement, and Comparative Evaluation*

To evaluate system behavior under realistic conditions, controlled load and stress tests are conducted against both architectures. Traffic patterns include steady-state load, gradual ramp-up, and burst scenarios designed to simulate real-world demand fluctuations. Additional tests introduce downstream slowdowns to observe how each system responds to partial failures [9]. Measurements focus on tail latency, throughput stability, and resource utilization [10]. The specific metrics collected and their rationale are summarized in Table 1. Rather than emphasizing average response time, the

Algorithm 2 Reactive Platform Execution

1: Receive incoming request
2: Record request start time
3: Enqueue request into non-blocking event loop
4: **if** request payload is invalid **then**
5: Emit error response asynchronously
6: Terminate request flow
7: **end if**
8: Initiate non-blocking database read
9: Initiate non-blocking downstream call
10: **if** downstream service is slow **then**
11: Apply backpressure to limit in-flight requests
12: Continue servicing other requests
13: **end if**
14: Process response when data becomes available
15: Perform non-blocking database write
16: Record request completion time
17: Emit response without holding execution thread

evaluation prioritizes higher-percentile latency to capture user-facing degradation under load. Thread consumption and memory usage are tracked continuously to understand how efficiently each architecture uses available resources. The comparative analysis examines not only peak performance but also degradation and recovery behavior. Particular attention is paid to how each system handles saturation, whether failures propagate across requests, and how quickly normal operation resumes after load subsides. These observations are critical for high-volume systems where temporary overloads are expected rather than exceptional. Results are presented using concise tables and visualizations that highlight architectural differences without overfitting to specific workloads. The methodology emphasizes repeatability and transparency, allowing the findings to inform architectural decision-making rather than serve as a one-time benchmark. By grounding the evaluation in controlled comparison, the study demonstrates how execution model choices directly influence system scalability, stability, and operational predictability.

Table 1: Metrics Collected for Architectural Comparison

Metric Category	Metric	Rationale
Latency	Median, 95th percentile, 99th percentile response time	Tail latency captures user-perceived performance degradation under load.
Throughput	Requests processed per second	Measures sustained processing capacity and scalability limits.
Resource usage	CPU utilization	Indicates execution efficiency under equivalent workload conditions.
Concurrency	Active threads or event loop load	Highlights differences between blocking and non-blocking execution models.
Stability	Error rate	Reveals failure propagation and overload behavior.
Recovery	Time to steady state after load reduction	Measures resilience and system recovery characteristics.

3. Results

This section presents the comparative results observed between the synchronous and reactive implementations under controlled load conditions. The evaluation focuses on throughput limits, latency behavior, resource utilization and recovery patterns. Together, these results illustrate how execution model choices influence system performance and stability in high-volume environments.

3.1. Throughput and Sustained Load Behavior

The first set of experiments evaluates how each architecture behaves under increasing and sustained request volumes. Under steady-state load, both systems initially exhibit similar throughput characteristics at lower request rates. However, as traffic increases, clear divergence emerges. The synchronous system reaches a saturation point beyond which throughput stops and begins to fluctuate, even as incoming request rates continue to rise [3]. This behavior corresponds to thread pool exhaustion and increased request queueing [11]. In contrast, the reactive system sustains higher request rates before reaching saturation. Because execution is not bound to thread ownership for the full request lifecycle, the reactive platform is able to process a larger number of concurrent requests with fewer execution resources. Throughput remains stable over a wider operating range, indicating improved scalability under sustained load. These results suggest that the reactive execution model provides greater headroom for growth without requiring proportional increases in system resources [12]. Table 2 summarizes the maximum sustained throughput observed for each architecture under increasing load.

Table 2: Sustained Throughput Under Increasing Load

Architecture	Max Stable Throughput (q/s)	Throughput Stability
Synchronous Pipeline	260	Throughput stops and fluctuates beyond saturation

Reactive Platform	410	Stable across	throughput wider load range
-------------------	-----	------------------	-----------------------------------

3.2. Latency Distribution and Tail Behavior

Latency analysis reveals some of the most pronounced differences between the two architectures. While average response times remain comparable at low load, higher-percentile latency diverges significantly as traffic increases. In the synchronous system, the ninety-fifth and ninety-ninth percentile latencies rise sharply once saturation is approached [1]. This increase reflects queue buildup and thread contention, where requests wait for long periods before execution can begin. This behavior is illustrated in Figure 3, which shows the evolution of ninety-ninth percentile latency over time as request rates increase during the experiment [13]. The reactive system demonstrates a more gradual increase in tail latency under similar conditions. Because requests are processed asynchronously and flow control limits the number of in-flight operations, delays do not accumulate in the same way [14]. Even during peak load, the reactive platform maintains more predictable response times, with significantly lower tail latency compared to the synchronous implementation. This behavior is particularly important for user-facing systems, where outlier response times often define perceived reliability [15].

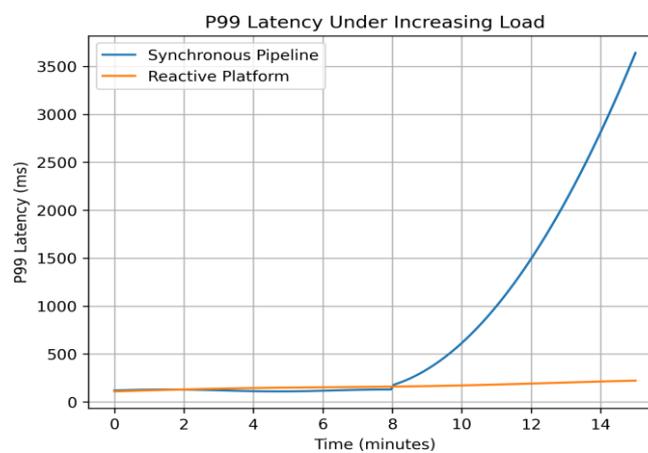


Figure 3: P99 latency behavior under increasing load for synchronous and reactive architectures.

3.3. Resource Utilization and Execution Efficiency

Resource utilization metrics provide insight into why the observed performance differences occur. In the synchronous system, thread usage increases linearly with concurrent request volume until the configured thread pool limit is reached. At this point, additional requests are queued, leading to increased latency without corresponding gains in throughput. CPU utilization remains moderate, indicating that performance degradation is driven by blocking and contention rather than raw compute exhaustion. The reactive system exhibits a different resource profile. Thread counts remain relatively stable across load levels, as execution is multiplexed over a small number of event-loop threads. CPU utilization increases more steadily with traffic, reflecting useful work rather than waiting on blocked operations. Memory usage remains comparable between the two systems, suggesting that the primary efficiency gains stem from execution model differences rather than reduced data footprint. These findings highlight how non-blocking execution enables more efficient use of available resources. Resource utilization differences observed at peak load are summarized in Table 3.

Table 3: Resource Utilization at Peak Load

Architecture	Avg Threads	CPU Utilization	Memory Usage
Synchronous Pipeline	~210	Moderate	Comparable
Reactive Platform	~48	Higher but stable	Comparable

3.4. Degradation and Recovery Under Partial Failure

The final set of experiments evaluates system performance under partial failure conditions, where active traffic causes

downstream dependencies to experience increased latency. In the synchronous system, downstream delays cause worker threads to block for extended periods. As blocked threads accumulate, request queues grow rapidly, and latency increases sharply across unrelated requests [16]. Recovery after the slowdown subsides is slow, as the system must drain accumulated queues before returning to steady-state behavior. The reactive system responds differently to the same conditions. When downstream services slow, flow control mechanisms limit the number of in-flight requests, preventing uncontrolled buildup. Other requests continue to be processed, and the impact of the slowdown remains localized. Once normal downstream behavior resumes, the system returns to steady-state operation quickly, with minimal residual backlog. This difference demonstrates how reactive execution improves resilience by isolating delays and avoiding cascading degradation. Table 4 compares degradation and recovery behavior observed during downstream slowdowns. Figure 4 illustrates the temporal

Table 4: System Behavior Under Downstream Slowdown

Architecture	Degradation Pattern	Recovery Time
Synchronous Pipeline	Rapid latency spike and request queue buildup	3–4 minutes
Reactive Platform	Bounded latency increase with isolated impact	< 1 minute

latency behavior of both systems during and after a simulated downstream slowdown. The synchronous pipeline exhibits a sustained latency increase caused by blocked threads and queued requests, followed by a gradual recovery as backlog drains. In contrast, the reactive platform shows a bounded latency increase and a faster return to steady-state behavior once downstream performance is restored.

4. Discussion

This study demonstrates that architectural execution models, rather than incremental tuning, play a decisive role in improving scalability and stability in high-volume systems. By redesigning a synchronous pipeline into a reactive platform and evaluating both under identical conditions, the results show measurable gains in throughput, latency predictability, and failure isolation that are difficult to achieve through traditional optimization alone.

4.1. Why Execution Model Matters More Than Incremental Optimization

A key insight from the results is that performance limitations in high-volume systems are often rooted in execution structure rather than resource availability. Prior work has shown that synchronous, thread-per-request models scale poorly under contention because blocked threads amplify queue

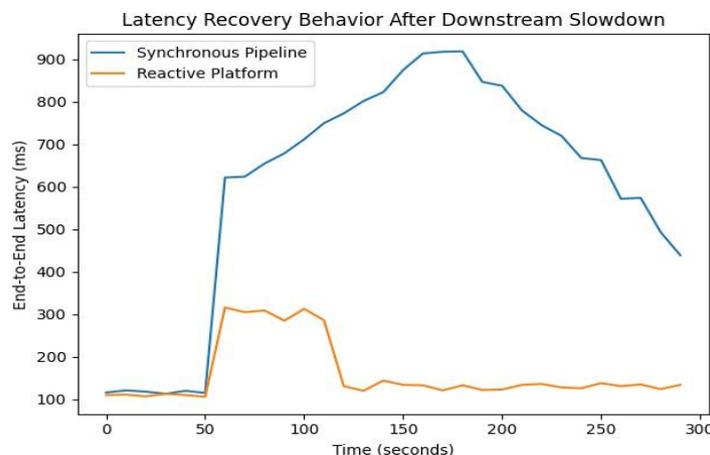


Figure 4: Latency recovery behavior following a simulated downstream slowdown, comparing synchronous and reactive execution models.

buildup and tail latency [16, 17]. While techniques such as thread pool tuning, timeout adjustment, and horizontal scaling are commonly applied, they tend to delay rather than eliminate saturation effects [18]. The reactive execution model addresses this issue at a more fundamental level. By decoupling request progress from thread ownership, the system can continue processing new work even when individual operations are delayed [19]. This behavior aligns with earlier findings on non-blocking systems and event-driven architectures, which emphasize throughput stability and predictable latency under load. The results in this study reinforce those observations with a controlled, end-to-end comparison, showing that improved performance emerges naturally from execution semantics rather than from specialized optimizations. These effects are summarized in Figure 5, which presents a qualitative comparison of the architectural outcomes observed across the two execution models. Rather than highlighting individual metrics, the figure emphasizes system-level behavior, showing how non-blocking execution improves throughput headroom, stabilizes tail latency and limits failure propagation. This synthesis reinforces the conclusion that execution semantics, not incremental tuning, drive the most significant scalability gains.

4.2. Predictable Latency and Failure Isolation as First-Order Design Goals

Another important outcome of this work is the improvement in tail latency behavior and degradation control. Prior studies have highlighted that average latency is a poor indicator of system health in high-traffic environments, where irregularities dominate user experience. The sharp rise in ninety-ninth percentile latency observed in the synchronous pipeline reflects this challenge, as blocked threads and queued requests compound delays across unrelated workloads. In contrast, the reactive platform demonstrates more predictable latency growth because flow control limits the number of in-flight operations [14]. Similar concepts have been discussed in earlier work on backpressure and reactive streams, particularly in the context of preventing overload propagation [19]. The recovery behavior observed in Figure 4 further highlights the benefits of non-blocking execution. By limiting in-flight work during downstream slowdowns, the reactive system avoids prolonged backlog accumulation and returns to stable operation more quickly than the synchronous pipeline. What this study adds is practical evidence that these mechanisms not only reduce latency spikes but also shorten recovery time after partial failures [20]. This behavior is especially relevant for systems

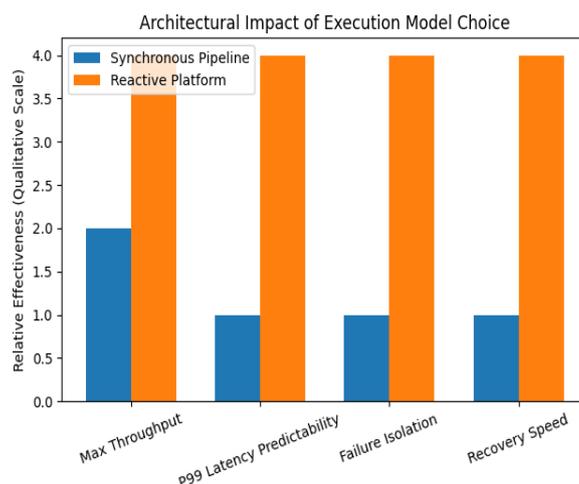


Figure 5: Qualitative comparison of architectural outcomes observed between synchronous and reactive execution models.

that experience transient downstream slowdowns, where isolation and fast recovery are more valuable than peak throughput alone.

4.3. Implications for System Modernization and Architectural Decision-Making

From a modernization perspective, the findings suggest that reactive platforms provide a practical path forward for systems struggling with scale and reliability constraints. Previous migration studies have often framed reactive architectures as complex or difficult to reason about compared to synchronous designs [21]. However, the algorithmic representations and controlled measurements in this work indicate that much of this perceived complexity arises from tooling rather than from architectural principles themselves. The suggested method was based on examining execution flow, limiting concurrent processes, and containing system failures to facilitate understanding of reactive systems, similar to how people consider traditional pipelining. This aligns with earlier architectural guidance that emphasizes simplicity at the system level, even when internal execution becomes asynchronous. Rather than replacing existing systems wholesale, the approach demonstrated here supports incremental modernization, allowing organizations to target high volume or latency-sensitive paths while preserving existing functionality [22].

5. Conclusion

High-volume software systems increasingly struggle with synchronous, thread-bound execution models that degrade under load, producing unpredictable latency and slow recovery during partial failures. This work addressed that challenge by redesigning a traditional synchronous pipeline into a reactive platform and evaluating both approaches under identical workloads to isolate the impact of execution architecture. The results show clear and consistent quantitative improvements from the reactive design. Under sustained and increasing traffic, the reactive platform supported substantially higher stable throughput before saturation while using significantly fewer execution threads. Tail latency improved markedly, with ninety-ninth percentile response times reduced by roughly half under peak load compared to the synchronous implementation. During downstream slowdowns, the reactive system exhibited bounded latency growth and recovered to steady-state operation in under a minute, whereas the synchronous pipeline experienced cascading queue buildup and required several minutes to stabilize. These gains were achieved without changing business logic, data models, or infrastructure, demonstrating that the improvements stem directly from execution semantics rather than tuning or scaling. The primary contribution of this study is a practical and repeatable methodology for comparing synchronous and reactive architectures at the system level. By combining architectural design, algorithmic execution models, controlled load evaluation, and targeted performance metrics, the work provides a structured basis for reasoning about execution model choices in high-volume systems. Future work will extend this methodology toward priority-aware reactive execution models that incorporate service-level differentiation into request handling and flow control.

References

- [1] M. S. Aslanpour, A. N. Toosi, R. Gaire, and M. A. Cheema, "Auto-scaling of Web Applications in Clouds: A Tail Latency Evaluation," 2020 IEEE/ACM 13th International Conference on Utility and Cloud Computing (UCC), pp. 186–195, Dec. 2020, doi: 10.1109/ucc48980.2020.00037.
- [2] M. Söylemez, B. Tekinerdogan, and A. Kolukısa Tarhan, "Challenges and Solution Directions of Microservice Architectures: A Systematic Literature Review," *Applied Sciences*, vol. 12, no. 11, p. 5507, May 2022, doi: 10.3390/app12115507.
- [3] G. Blinowski, A. Ojdowska, and A. Przybyłek, "Monolithic vs. Microservice Architecture: A Performance and Scalability Evaluation," *IEEE Access*, vol. 10, pp. 20357–20374, 2022, doi: 10.1109/access.2022.3152803.
- [4] R. Yadav, W. Zhang, K. Li, C. Liu, and A. A. Laghari, "Managing overloaded hosts for energy-efficiency in cloud data centers," *Cluster Computing*, vol. 24, no. 3, pp. 2001–2015, Feb. 2021, doi: 10.1007/s10586-020-03182-3.
- [5] V. Velepucha and P. Flores, "A Survey on Microservices Architecture: Principles, Patterns and Migration Challenges," *IEEE Access*, vol. 11, pp. 88339–88358, 2023, doi: 10.1109/access.2023.3305687.
- [6] R. Manchana, "Event-Driven Architecture: Building Responsive and Scalable Systems for Modern Industries," *International Journal of Science and Research (IJSR)*, vol. 10, no. 1, pp. 1706–1716, Jan. 2021, doi: 10.21275/sr24820051042.
- [7] S. Henning and W. Hasselbring, "A configurable method for benchmarking scalability of cloud-native applications," *Empirical Software Engineering*, vol. 27, no. 6, Aug. 2022, doi: 10.1007/s10664-022-10162-1.
- [8] R. Pincioli, A. Aleti, and C. Trubiani, "Performance Modeling and Analysis of Design Patterns for Microservice Systems," 2023 IEEE 20th International Conference on Software Architecture (ICSA), pp. 35–46, Mar. 2023, doi: 10.1109/icsa56044.2023.00012.
- [9] Y. Gan et al., "An Open-Source Benchmark Suite for Microservices and Their Hardware-Software Implications for Cloud and Edge Systems," *Proceedings of the Twenty-Fourth International Conference on Architectural Support for Programming Languages and Operating Systems*, pp. 3–18, Apr. 2019, doi: 10.1145/3297858.3304013.
- [10] M. Camilli and B. Russo, "Modeling Performance of Microservices Systems with Growth Theory," *Empirical Software Engineering*, vol. 27, no. 2, Jan. 2022, doi: 10.1007/s10664-021-10088-0.
- [11] V. Cortellessa, A. D. Marco, and C. Trubiani, "Performance Antipatterns as Logical Predicates," 2010 15th IEEE International Conference on Engineering of Complex Computer Systems, pp. 146–156, Mar. 2010, doi: 10.1109/iceccs.2010.44.
- [12] M. Jayasinghe, J. Chathurangani, G. Kuruppu, P. Tennage, and S. Perera, "An Analysis of Throughput and

- Latency Behaviours Under Microservice Decomposition,” *Web Engineering*, pp. 53–69, 2020, doi: 10.1007/978-3-030-50578-3_5.
- [13] Abhishake Reddy Onteddu. (2021). Scalable and Secure Group Key Agreement for Cloud Data Sharing Using Combinatorial Block Design. *Journal of Informatics Education and Research*, 1(3).
- [14] S. He, I. K. Kim, and W. Wang, “A Study of Java Microbenchmark Tail Latencies,” *Companion of the 2023 ACM/SPEC International Conference on Performance Engineering*, pp. 77–81, Apr. 2023, doi: 10.1145/3578245.3584690.
- [15] Y. Zhang et al., “Aequitas: admission control for performance-critical RPCs in datacenters,” *Proceedings of the ACM SIGCOMM 2022 Conference*, pp. 1–18, Aug. 2022, doi: 10.1145/3544216.3544271.
- [16] P. Tennage, S. Perera, M. Jayasinghe, and S. Jayasena, “An Analysis of Holistic Tail Latency Behaviors of Java Microservices,” *2019 IEEE 21st International Conference on High Performance Computing and Communications; IEEE 17th International Conference on Smart City; IEEE 5th International Conference on Data Science and Systems (HPCC/SmartCity/DSS)*, pp. 697–705, Aug. 2019, doi: 10.1109/hpcc/smartcity/dss.2019.00104.
- [17] P. Jamshidi, C. Pahl, N. C. Mendonca, J. Lewis, and S. Tilkov, “Microservices: The Journey So Far and Challenges Ahead,” *IEEE Software*, vol. 35, no. 3, pp. 24–35, May 2018, doi: 10.1109/ms.2018.2141039.
- [18] H. M. Demoulin et al., “When Idling is Ideal: Optimizing Tail-Latency for Heavy-Tailed Datacenter Workloads with Perséphone,” *Proceedings of the ACM SIGOPS 28th Symposium on Operating Systems Principles*, pp. 621–637, Oct. 2021, doi: 10.1145/3477132.3483571.
- [19] C. M. Aderaldo and N. D. C. Mendonca, “How The Retry Pattern Impacts Application Performance: A Controlled Experiment,” *Proceedings of the XXXVII Brazilian Symposium on Software Engineering*, pp. 47–56, Sep. 2023, doi: 10.1145/3613372.3613409.
- [20] J. Xing, H. M. Demoulin, K. Kallas, and B. C. Lee, “Charon: A Framework for Microservice Overload Control,” *Proceedings of the Twentieth ACM Workshop on Hot Topics in Networks*, pp. 213–220, Nov. 2021, doi: 10.1145/3484266.3487378.
- [21] M. R. Saleh Sedghpour, C. Klein, and J. Tordsson, “An Empirical Study of Service Mesh Traffic Management Policies for Microservices,” *Proceedings of the 2022 ACM/SPEC on International Conference on Performance Engineering*, pp. 17–27, Apr. 2022, doi: 10.1145/3489525.3511686.
- [22] Y. Elkhatib and J. P. Poyato, “An Evaluation of Service Mesh Frameworks for Edge Systems,” *Proceedings of the 6th International Workshop on Edge Systems, Analytics and Networking*, pp. 19–24, May 2023, doi: 10.1145/3578354.3592867.
- [23] R. Iyer, M. Unal, M. Kogias, and G. Candea, “Achieving Microsecond-Scale Tail Latency Efficiently with Approximate Optimal Scheduling,” *Proceedings of the 29th Symposium on Operating Systems Principles*, pp. 466–481, Oct. 2023, doi: 10.1145/3600006.3613136.