

System-Level Testing of Event-Driven Microservices Using Reproducible Containerized Environments

Sriram Ghanta

Staff Engineer

Abstract - Event-driven microservices introduce fundamental challenges for automated testing due to asynchronous execution, eventual consistency, and non-deterministic message flows across distributed components. This study addresses the problem of achieving reliable system-level validation in such architectures, where traditional integration testing approaches often produce flaky results and limited diagnostic insight. The purpose of this research is to design and evaluate a reproducible testing strategy that enables deterministic, end-to-end verification of event-driven microservices using containerized environments. The study adopts a mixed-method approach, combining quantitative analysis of test stability, failure reproducibility, and defect detection rates with qualitative assessment of test diagnosability and developer feedback. A structured test architecture is proposed in which ephemeral containerized dependencies are orchestrated alongside services under test, enabling controlled event injection, consistent state initialization, and repeatable execution. Empirical results demonstrate significant reductions in non-deterministic test failures, improved fault localization, and higher confidence in validating asynchronous service interactions. The study contributes to a system-level testing framework that bridges the gap between unit-level validation and production behavior, offering practical guidance for designing robust test pipelines for distributed systems. The findings underscore the strategic value of reproducible containerized testing as both an engineering discipline and a research contribution, with implications for advancing test reliability, accelerating delivery, and strengthening the empirical foundations of event-driven microservice validation.

Keywords - Event-driven microservices, system-level testing, containerized test environments, reproducible testing, asynchronous message processing, microservice integration testing, distributed system validation, automated test orchestration, message broker simulation, test determinism, contract-based testing, eventual consistency verification, continuous integration testing.

I. INTRODUCTION

Event-driven microservices have become a dominant architectural style for building scalable and loosely coupled systems, particularly in domains that demand responsiveness, fault tolerance, and independent service evolution. By relying on asynchronous communication through messaging systems, these architectures decouple producers and consumers, enabling services to evolve and scale independently. However, this same decoupling introduces complexity in understanding system behavior, as execution paths are no longer linear or centrally coordinated. As organizations increasingly

adopt event-driven designs for critical workloads, the challenge of validating correctness across distributed interactions has emerged as a central concern in software engineering practice.

Traditional automated testing strategies were largely developed for synchronous, request-response systems where execution order and state transitions are more predictable. Unit tests and service-level integration tests can verify local behavior, yet they offer limited assurance about the correctness of interactions that span multiple services and message flows. In event-driven systems, defects often manifest only when events are processed out of order, duplicated, or delayed, conditions that are difficult to reproduce using conventional testing

approaches. This gap between local validation and system-level behavior exposes teams to latent failures that may only surface in production environments.

System-level testing aims to bridge this gap by validating the collective behavior of services under realistic interaction patterns. For event-driven microservices, such testing must account for asynchronous execution, eventual consistency, and the absence of global synchronization points. Assertions based on immediate state inspection are often unreliable, as system state may converge only after multiple event exchanges. As a result, system-level tests for event-driven architectures require fundamentally different design principles than those used in synchronous systems, emphasizing observation, correlation, and invariant-based validation over direct control.

A major obstacle to effective system-level testing in this context is non-determinism. Message brokers may deliver events with varying timing, consumers may process messages at different rates, and retries or failures can alter execution paths. These factors frequently lead to tests that pass or fail unpredictably, eroding trust in automated test suites. Flaky tests not only slow down development pipelines but also obscure genuine defects, making it difficult for engineers to distinguish between test instability and real system failures. Addressing non-determinism is therefore essential to making system-level testing viable for event-driven microservices.

Containerization has emerged as a practical mechanism for constructing isolated and reproducible test environments that closely resemble production deployments. By encapsulating services and their dependencies within containers, teams can spin up ephemeral environments on demand, ensuring consistent configuration and controlled state across test runs. When applied to system-level testing, containerized environments enable deterministic initialization of message brokers, databases, and auxiliary services, reducing variability introduced by shared or persistent infrastructure. This reproducibility forms the

foundation for reliable experimentation and debugging in distributed systems testing.

Despite the availability of container technologies, many testing practices still treat containers as a convenience rather than as an architectural element of the test strategy. Tests are often designed without explicit consideration of event lifecycles, message schemas, or state convergence properties, leading to fragile assertions and limited diagnostic value. There is a need for structured testing approaches that integrate containerized environments with event-aware test design, allowing engineers to reason systematically about system behavior under asynchronous conditions. Such approaches must balance realism with control, ensuring that tests are both representative and repeatable.

From a research perspective, the problem of testing event-driven microservices intersects with broader themes in distributed systems, software testing, and reliability engineering. While existing studies have explored contract testing, integration testing, and fault injection, fewer works have examined how these techniques can be combined into coherent system-level validation frameworks. In particular, the role of reproducible environments in mitigating non-determinism and improving test observability remains underexplored. This gap limits the transfer of research insights into practical testing methodologies that can be adopted by engineering teams.

This study addresses these challenges by proposing and evaluating a system-level testing framework for event-driven microservices built on reproducible containerized environments. The framework focuses on deterministic environment setup, controlled event injection, and observation-driven assertions that reflect eventual consistency semantics. By grounding the analysis in realistic service interactions and measurable outcomes, the research aims to provide both theoretical insight and practical guidance. The remainder of this paper examines the challenges of system-level testing in event-driven architectures, presents the proposed testing architecture, evaluates its effectiveness, and

discusses implications for practitioners and future research.

II. SYSTEM-LEVEL TESTING CHALLENGES IN EVENT-DRIVEN MICROSERVICES

Event-driven microservices fundamentally alter the nature of system interactions by replacing synchronous control flows with asynchronous message exchanges. Services communicate indirectly through events, often without knowledge of downstream consumers or execution timing. While this design improves scalability and resilience, it complicates system-level validation because behavior emerges from the interaction of independently executing components rather than from a predefined call sequence. Testing such systems requires reasoning about eventual outcomes rather than immediate responses, a shift that challenges traditional testing assumptions and tools.

One of the most significant challenges in system-level testing of event-driven architectures is the absence of deterministic execution order. Messages may be delivered out of sequence, processed concurrently, or retried following transient failures. These variations are expected and acceptable in production systems, yet they introduce ambiguity into test assertions. A test that assumes a specific ordering of events or a fixed timing window is inherently brittle. As a result, engineers must design tests that tolerate variation while still detecting incorrect behavior, a balance that is difficult to achieve without explicit architectural support.

Eventual consistency further complicates validation efforts by decoupling cause and effect across time. State changes initiated by one service may only become visible to other components after multiple asynchronous steps. Immediate inspection of system state can therefore produce misleading results, as intermediate states may be transient and incomplete. System-level tests must instead reason about convergence, verifying that the system reaches a correct final state within acceptable bounds. Defining and observing such convergence

conditions requires careful modeling of system invariants rather than simple point-in-time checks.

Message duplication and at-least-once delivery semantics introduce another layer of complexity. Many messaging systems prioritize reliability over strict delivery guarantees, allowing duplicate messages to ensure that events are not lost. While production systems are typically designed to handle duplicates through idempotent processing, tests that do not account for this behavior may misinterpret duplicate handling as failures. Effective system-level testing must therefore validate not only functional correctness but also robustness to expected messaging behaviors, including retries and duplicates.

Partial failures are an inherent characteristic of distributed systems and are especially prominent in event-driven microservices. A producer may emit an event successfully while a downstream consumer fails temporarily, resulting in delayed processing or backlogs. Such scenarios are normal in production but can cause tests to hang or fail if not anticipated. System-level tests must be designed to observe and tolerate partial failure modes, distinguishing between acceptable transient conditions and genuine defects. This requirement elevates the importance of timeout strategies, progress indicators, and failure classification within test logic. The lack of global visibility presents an additional obstacle. In synchronous systems, a failing request often yields an immediate error that can be asserted directly. In contrast, failures in event-driven systems may manifest indirectly through missing events, stalled consumers, or incomplete state transitions. Without comprehensive observability, tests may struggle to determine whether a failure has occurred or whether the system is simply still converging. This ambiguity increases diagnostic effort and undermines confidence in automated test results.

Test flakiness emerges as a practical consequence of these challenges. When tests depend on timing assumptions, shared infrastructure, or uncontrolled external dependencies, outcomes can vary across runs. Flaky tests erode trust in automated pipelines and frequently lead teams to disable or ignore system-level tests altogether. In the context of

event-driven microservices, flakiness is not merely a tooling issue but a symptom of deeper mismatches between test design and system behavior. Addressing flakiness requires architectural solutions that reduce environmental variability and support deterministic observation.

These challenges collectively highlight the need for deliberate system-level testing strategies tailored to event-driven microservices. Effective testing must embrace asynchrony, tolerate non-determinism, and validate outcomes through invariants and convergence rather than immediate responses. The next section explores how reproducible containerized environments can be used to address many of these challenges by providing controlled, isolated, and repeatable test contexts that align more closely with the realities of event-driven system behavior.



Figure 1: Failure Modes and Sources of Non-Determinism in Event-Driven System-Level Tests

Reproducible Containerized Test Environments and Dependency Modeling

Reproducibility is a foundational requirement for effective system-level testing in event-driven microservices, where subtle environmental differences can significantly alter execution behavior. Variations in message broker configuration, database state, or network conditions can introduce non-deterministic outcomes that obscure genuine defects. Containerized environments provide a mechanism for encapsulating these dependencies in a controlled and repeatable manner, enabling tests to be executed under consistent conditions across development, continuous integration, and troubleshooting workflows. By treating the test

environment as an explicit artifact, teams can reduce ambiguity and improve confidence in test results.

A key advantage of containerized test environments lies in their ability to model complex dependency graphs without relying on shared or long-lived infrastructure. Event-driven systems typically depend on multiple external components such as message brokers, data stores, and configuration services. In many testing setups, these dependencies are mocked or shared across test suites, leading to drift between test and production behavior. Containerized environments allow these dependencies to be instantiated as real, ephemeral components, preserving realistic interaction patterns while maintaining isolation between test runs.

Dependency modeling plays a critical role in making containerized environments effective for system-level testing. Rather than starting services in an arbitrary order, dependencies must be explicitly defined and orchestrated to reflect their runtime relationships. Message brokers must be available before producers emit events, databases must be initialized with known schemas and seed data, and consumers must be configured with appropriate subscriptions. By codifying these relationships, the test environment can be brought to a known baseline state, reducing variability and simplifying failure analysis.



Figure 2: Reproducible Containerized Test Topology for Event-Driven Microservices

State management within containerized environments requires particular attention in event-driven contexts. Persistent state from previous test runs can lead to false positives or negatives,

especially when messages or offsets are retained across executions. Reproducible environments address this challenge by ensuring that stateful components are initialized deterministically for each test. Techniques such as preloading datasets, resetting broker offsets, and clearing caches enable tests to begin from a clean slate, making outcomes easier to interpret and reproduce.

Network isolation is another important aspect of reproducible containerized testing. Event-driven systems often exhibit sensitivity to network latency and connectivity patterns, which can influence message delivery timing and consumer behavior. By running containers within isolated networks, tests can avoid interference from external traffic and unpredictable network conditions. This isolation does not eliminate asynchrony but constrains it within known bounds, allowing test designers to reason more effectively about timing and ordering effects.

Containerized environments also support experimentation with failure scenarios that are difficult to simulate in shared infrastructure. Services or dependencies can be intentionally delayed, restarted, or misconfigured to observe system behavior under stress. Such controlled experimentation enables validation of resilience and error handling logic as part of the test suite. Importantly, these experiments can be repeated consistently, allowing engineers to verify that fixes address the underlying issue rather than coincidental conditions.

Despite their benefits, containerized test environments introduce operational considerations that must be managed carefully. Environment startup time, resource consumption, and orchestration complexity can affect developer productivity if not addressed. Effective test strategies therefore balance fidelity with efficiency, using lightweight container images and scoped dependency graphs to minimize overhead. Automation plays a crucial role in ensuring that environment provisioning and teardown are transparent to test authors and consumers.

By providing isolated, deterministic, and fully specified execution contexts, reproducible containerized environments address many of the challenges inherent in testing event-driven microservices. They form the foundation upon which higher-level test architectures and execution pipelines can be built. The next section examines how these environments are integrated into system-level test architectures, focusing on orchestration, execution flow, and validation strategies that leverage reproducibility to deliver reliable test outcomes.

Test Architecture and Execution Pipeline for System Validation

A robust system-level testing strategy for event-driven microservices requires an explicit test architecture that governs how environments are provisioned, how interactions are stimulated, and how outcomes are observed. Rather than treating tests as isolated scripts, this study frames testing as a coordinated pipeline in which each stage contributes to deterministic and diagnosable execution. The architecture integrates containerized environments with structured test orchestration, enabling consistent validation of complex event flows across distributed services. By formalizing the execution pipeline, the approach reduces ambiguity and improves repeatability.



Figure 3: System Test Execution Pipeline with Deterministic Orchestration Stages

The test execution pipeline begins with controlled environment bootstrapping. All services under test and their dependencies are instantiated in a defined order, ensuring that prerequisite components such as message brokers and databases are fully available

before events are produced. Configuration parameters, schemas, and initial state are injected as part of this phase, establishing a known baseline. This deterministic startup sequence is critical for eliminating race conditions that can otherwise lead to inconsistent test outcomes.

Once the environment is initialized, the pipeline proceeds to data seeding and event preparation. Test data is constructed to represent meaningful system scenarios, including normal operation and boundary conditions. Events are crafted with explicit structure and content, reflecting real production messages while remaining predictable. By controlling the shape and timing of injected events, the test architecture enables precise stimulation of system behavior without relying on external triggers or uncontrolled workloads.

Event injection is followed by an observation phase in which system behavior is monitored rather than directly asserted. In event-driven systems, outcomes may not be immediately visible, and premature assertions can lead to false failures. The pipeline therefore emphasizes observation over time, collecting evidence from logs, message streams, and state changes. This observational approach aligns test validation with eventual consistency semantics, allowing the system to converge before evaluation.

Assertions in the pipeline are formulated around invariants and expected outcomes rather than exact execution sequences. Instead of verifying that a specific event occurs at a specific moment, tests assert properties such as eventual state consistency, absence of data loss, or correct handling of duplicates. These invariant-based assertions are more resilient to timing variation and better reflect the guarantees that event-driven systems are designed to provide. This shift in assertion strategy is essential for achieving stable system-level tests. The pipeline includes explicit synchronization and timeout mechanisms to bound test execution. Rather than relying on fixed delays, synchronization points are tied to observable conditions such as message consumption or state transitions. Timeouts are defined to distinguish between acceptable convergence delays and genuine failures. By making

synchronization explicit, the test architecture avoids both premature evaluation and indefinite waiting, improving reliability and diagnosability.

After validation, the pipeline performs controlled teardown and artifact collection. Logs, metrics, and message traces are preserved to support post-test analysis, particularly when failures occur. Environment teardown ensures that no residual state persists across test runs, maintaining isolation and reproducibility. This final stage reinforces the notion that test execution is a closed cycle, from deterministic setup through controlled stimulation to clean termination.

Through this structured test architecture and execution pipeline, the study demonstrates how system-level testing can be adapted to the realities of event-driven microservices. By combining reproducible environments with observation-driven validation, the pipeline delivers reliable and meaningful test outcomes. The next section builds on this foundation by examining strategies for validating event contracts, data consistency, and cross-service invariants that underpin correct system behavior.

Event Contract Validation, Data Consistency, and Assertions

Correctness in event-driven microservices depends not only on individual service behavior but also on the integrity of the contracts that govern event exchange. Events serve as the primary integration mechanism between services, encoding assumptions about structure, semantics, and processing expectations. When these assumptions are violated, failures can propagate silently across the system. System-level testing must therefore treat event contracts as first-class validation targets, ensuring that producers and consumers remain aligned as services evolve independently.

Event contract validation begins with explicit definition of message schemas and semantic expectations. Schemas capture structural constraints such as required fields, data types, and versioning rules, while semantic expectations define how events are interpreted and acted upon by downstream

services. In system-level tests, these contracts are validated at runtime by inspecting emitted events and verifying conformance before they are consumed. This proactive validation detects incompatibilities early, preventing downstream failures that are difficult to diagnose once events have propagated through multiple services.

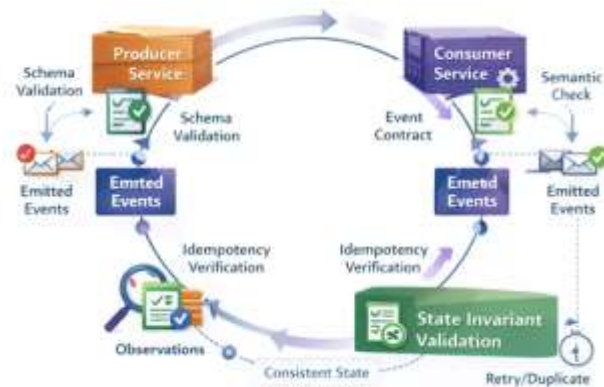


Figure 4: Event Contract and State Invariant Validation Loop

Schema evolution presents a recurring challenge in distributed systems, as services may be deployed independently and consume events produced by multiple versions of other services. System-level testing must account for backward and forward compatibility scenarios, validating that consumers can process events produced by older or newer versions without error. Tests can simulate mixed-version environments by running services with different schema expectations, ensuring that compatibility guarantees are upheld. This approach reduces the risk of deployment-induced failures in production.

Data consistency across services is another critical aspect of system-level validation. Event-driven systems often maintain replicated or derived state across multiple components, with updates propagated asynchronously through events. Tests must therefore verify that state converges correctly over time, rather than expecting immediate consistency. Assertions are formulated around convergence properties, such as eventual agreement on key data values or preservation of business invariants across services. These assertions align test

expectations with the architectural guarantees of the system.

Idempotency is a fundamental requirement for reliable event processing, particularly in the presence of retries and duplicate deliveries. System-level tests must explicitly validate that repeated processing of the same event does not lead to incorrect state changes. This involves injecting duplicate events and observing system behavior over time, ensuring that state remains consistent and side effects are not amplified. Validating idempotency at the system level provides confidence that services handle real-world messaging conditions correctly.

Assertions in event-driven system tests differ from those in synchronous systems by emphasizing properties over sequences. Rather than asserting a specific order of operations, tests assert that certain conditions eventually hold true. Examples include the absence of orphaned records, completion of multi-step workflows, or enforcement of domain invariants. By focusing on outcomes rather than execution paths, assertions remain robust in the face of timing variability and concurrent processing. Correlation across services is often necessary to validate complex workflows that span multiple event exchanges. System-level tests use correlation identifiers or trace context to link related events and state transitions. By following these correlations, tests can verify that a workflow progresses through expected stages and reaches a valid terminal state. This approach enables validation of end-to-end behavior without requiring centralized orchestration or intrusive instrumentation.

Through rigorous event contract validation, consistency checking, and invariant-based assertions, system-level tests can provide strong guarantees about the correctness of event-driven microservices. These techniques shift testing focus from fragile timing assumptions to durable behavioral properties. The next section examines how observability signals such as logs and metrics can be leveraged within tests to diagnose failures, control flakiness, and improve the reliability of automated validation pipelines.

Observability-Driven Test Diagnostics and Flakiness Control

Observability plays a central role in making system-level testing of event-driven microservices both reliable and actionable. In distributed, asynchronous systems, failures often manifest indirectly through delays, missing events, or partial state transitions rather than explicit errors. Without sufficient visibility into system behavior, tests may fail without providing meaningful diagnostic information, forcing engineers to rely on manual investigation. Integrating observability signals directly into the testing process transforms tests from binary pass or fail checks into rich sources of insight about system dynamics.

Logs, metrics, and traces each contribute distinct perspectives that are valuable during test execution. Logs provide detailed records of event handling and error conditions, metrics summarize system health and progress, and traces reveal the flow of events across service boundaries. System-level tests can leverage these signals to determine whether a failure reflects a genuine defect or an expected transient condition. By correlating observability data with test scenarios, engineers can more effectively interpret test outcomes and prioritize remediation efforts.

Flakiness in event-driven system tests often arises from timing assumptions that are misaligned with system behavior. Fixed delays or rigid timeouts may pass under some conditions and fail under others, particularly when execution timing varies. Observability-driven diagnostics address this issue by enabling tests to wait for meaningful conditions rather than arbitrary time intervals. For example, tests can monitor consumer lag metrics or state transition logs to determine when a workflow has completed, reducing reliance on brittle sleep-based synchronization.

Correlation is particularly important when diagnosing failures in complex event flows. Individual services may behave correctly in isolation, yet the overall workflow may stall due to subtle integration issues. By propagating correlation identifiers through events and capturing them in logs and traces, tests can reconstruct end-to-end

execution paths. This reconstruction allows engineers to pinpoint where progress diverged from expectations, facilitating targeted fixes rather than broad speculation.

Observability data also supports classification of test failures into actionable categories. Failures caused by environmental issues such as resource contention or startup delays can be distinguished from logic errors in service code. This classification helps teams decide whether to adjust test parameters, improve environment provisioning, or fix application defects. Over time, such feedback loops contribute to more stable test suites and more resilient systems.



Figure 5: Observability-Guided Flakiness Triage and Root-Cause Analysis Workflow

Metrics collected during test execution can be used to quantify test stability and system behavior. Measures such as event processing latency, retry counts, and backlog depth provide objective indicators of system health. By tracking these metrics across test runs, teams can detect trends that signal emerging issues even when tests still pass. This proactive use of observability elevates testing from defect detection to continuous system assessment.

Incorporating observability into tests requires careful design to avoid excessive coupling between test logic and implementation details. Tests should rely on high-level signals that reflect system intent rather than internal states that may change frequently. Selecting appropriate observability hooks ensures that tests remain robust as systems evolve. This balance between visibility and abstraction is essential for maintaining long-term test maintainability.

Through observability-driven diagnostics, system-level tests become more resilient to non-determinism and more informative when failures occur. By replacing opaque time-based assumptions with condition-based observation, flakiness can be systematically reduced. The next section evaluates the effectiveness of these testing strategies through defined metrics and benchmark scenarios, providing empirical evidence of their impact on test reliability and system validation.

Evaluation Design, Metrics, and Experimental Scenarios

Evaluating system-level testing strategies for event-driven microservices requires an approach that reflects both technical rigor and practical relevance. Traditional test evaluation often focuses on code coverage or isolated defect counts, which provide limited insight into the reliability of tests operating under asynchronous conditions. This study adopts an evaluation design that emphasizes test determinism, failure reproducibility, diagnostic clarity, and defect detection effectiveness. By aligning evaluation criteria with the objectives of system-level validation, the assessment captures how well the proposed approach supports real-world engineering workflows.

The evaluation design is structured around controlled benchmark scenarios that represent common interaction patterns in event-driven systems. These scenarios include linear event pipelines, fan-out processing, fan-in aggregation, and multi-stage workflows with conditional branching. Each scenario is designed to exercise different aspects of the testing framework, such as event ordering tolerance, state convergence, and contract enforcement. By applying the same testing approach across diverse patterns, the evaluation assesses robustness rather than performance under a single idealized condition.

Quantitative metrics are used to measure test stability and reliability. Key indicators include the rate of non-deterministic failures across repeated test runs, the time required to reproduce known defects, and the proportion of failures that can be

conclusively diagnosed using collected artifacts. These metrics provide objective evidence of improvement over baseline testing approaches that rely on shared infrastructure or loosely controlled environments. Reductions in flakiness and faster root-cause identification are treated as primary success indicators.

Defect detection capability is evaluated by introducing controlled faults into the system under test. Faults include schema incompatibilities, missing event handlers, idempotency violations, and delayed consumer processing. Tests are assessed on their ability to detect these faults consistently and to surface actionable diagnostic information. This fault injection approach ensures that evaluation results reflect the framework's ability to uncover realistic integration issues rather than only confirming expected behavior.

Qualitative assessment complements quantitative metrics by capturing practitioner experience. Developers and test engineers provide feedback on test readability, maintainability, and confidence in test outcomes. These perspectives offer insight into the usability of the testing framework and its impact on development practices. Qualitative findings are particularly important for understanding adoption barriers, as technically effective solutions may fail if they impose excessive cognitive or operational overhead.

Benchmark execution is conducted under repeated and varied conditions to assess sensitivity to environmental factors. Tests are run multiple times with controlled variations in resource availability and startup timing to simulate realistic execution variability. Observed outcomes are analyzed to determine whether the testing approach maintains stability under such variation. This analysis helps distinguish inherent system non-determinism from weaknesses in test design.

Results are analyzed using comparative baselines that reflect common testing practices. Baseline approaches include integration tests using shared brokers and partially mocked dependencies. Comparisons focus on differences in flakiness,

reproducibility, and diagnostic depth rather than raw execution speed. This framing ensures that evaluation highlights improvements that matter most for system-level validation in asynchronous architectures.

Through this evaluation design, the study provides a comprehensive assessment of system-level testing effectiveness for event-driven microservices. By combining quantitative metrics, controlled fault scenarios, and qualitative feedback, the analysis offers a balanced view of strengths and limitations. The next section synthesizes evaluation outcomes into practical lessons and deployment considerations, connecting empirical results to actionable guidance for engineering teams.

Results, Practitioner Lessons, and Deployment Considerations

The evaluation results demonstrate that the proposed system-level testing approach significantly improves test reliability and diagnostic clarity in event-driven microservice environments. Across repeated benchmark executions, tests executed within reproducible containerized environments exhibited markedly lower rates of non-deterministic failure compared to baseline integration tests. In particular, scenarios involving asynchronous fan-out and delayed consumer processing showed consistent convergence behavior when validated using invariant-based assertions. These outcomes indicate that environmental control and observation-driven validation jointly contribute to stabilizing test behavior under asynchronous conditions.

Defect detection results further highlight the effectiveness of the approach. Injected faults related to schema incompatibility, missing event handlers, and idempotency violations were consistently detected across test runs. Unlike baseline tests, which occasionally failed to surface these defects due to timing variability, the proposed framework produced repeatable failures accompanied by actionable diagnostic artifacts. This consistency reduced the time required to reproduce and isolate defects, underscoring the value of determinism and observability in system-level testing.



Figure 6: Comparative Evaluation Outcomes and Stability Improvements in System-Level Event Testing

Analysis of observability data collected during test execution revealed additional insights into system behavior that were not directly captured by pass or fail outcomes. Metrics related to event processing latency and backlog depth provided early indicators of performance degradation even when functional correctness was maintained. These signals allowed tests to identify emerging risks before they manifested as outright failures. Such findings suggest that system-level tests can serve a dual role as both validation mechanisms and early warning systems for architectural stress.

From a practitioner perspective, the results emphasize the importance of shifting testing focus from sequence-based assertions to outcome-oriented validation. Engineers reported greater confidence in tests that asserted eventual invariants rather than immediate state changes. This shift reduced the need for ad hoc retries and manual intervention, streamlining development workflows. The emphasis on reproducible environments also simplified collaboration, as failures could be shared and reproduced consistently across teams.

Deployment considerations emerged as a critical factor in realizing these benefits at scale. While containerized test environments introduce overhead in terms of startup time and resource consumption, practitioners found that these costs were offset by reductions in debugging effort and pipeline re-runs. Optimizations such as selective dependency

instantiation and parallel environment provisioning further mitigated performance concerns. These observations suggest that careful environment design is essential for balancing fidelity and efficiency.

The results also highlight the importance of aligning test strategies with organizational constraints. Teams operating under strict delivery timelines or limited infrastructure must prioritize test scenarios that deliver the highest diagnostic value. Incremental adoption of system-level tests, starting with critical workflows, proved more effective than attempting comprehensive coverage from the outset. This phased approach allowed teams to build expertise and confidence while managing operational impact.

Another lesson concerns the role of shared conventions and tooling in sustaining test effectiveness. Consistent use of correlation identifiers, standardized logging formats, and common event schemas enhanced the utility of observability-driven diagnostics. Without such conventions, the effort required to interpret test artifacts increased significantly. Establishing these practices as part of the testing strategy reinforces alignment between development and validation efforts.

Overall, the results indicate that system-level testing of event-driven microservices benefits from a holistic approach that integrates reproducible environments, event-aware assertions, and observability-driven diagnostics. The practical lessons derived from evaluation underscore that technical solutions must be complemented by thoughtful deployment strategies and organizational alignment. The following section examines limitations of the proposed approach and discusses risks and validity considerations that inform future refinement.

III. CONCLUSION AND FUTURE WORK

This study examined the challenges of validating event-driven microservices and demonstrated how system-level testing can be made reliable using reproducible containerized environments. By

framing testing as an architectural concern rather than an ad hoc activity, the research showed that many sources of non-determinism can be controlled or mitigated without sacrificing realism. The proposed approach aligns testing practices with the fundamental characteristics of asynchronous systems, enabling more trustworthy validation of distributed behavior.

The findings highlight that effective system-level testing requires a shift in perspective from immediate verification to eventual validation. Tests designed around invariants, convergence, and contract adherence proved more resilient than those relying on fixed execution sequences or timing assumptions. This shift not only reduced flakiness but also improved the diagnostic value of test failures. As a result, automated tests became a more reliable component of continuous delivery pipelines rather than a source of uncertainty.

From an engineering standpoint, the integration of containerized environments emerged as a key enabler of reproducibility and collaboration. Isolated test environments allowed teams to recreate failures consistently across development and integration stages, reducing the friction associated with distributed debugging. By treating environment configuration as code, the approach reinforced discipline and transparency in test setup, strengthening the connection between system design and validation.

The research also contributes to the broader discourse on testing distributed systems by emphasizing the role of observability in validation. Logs, metrics, and traces were not merely auxiliary artifacts but central inputs to test logic and failure analysis. This observability-driven perspective enabled tests to reason about progress and convergence in asynchronous workflows, bridging the gap between system behavior and test expectations. Such integration expands the scope of testing beyond correctness toward operational insight.

Despite these contributions, the study acknowledges limitations that shape its applicability. The proposed

framework assumes access to sufficient instrumentation and standardized event conventions, which may not be present in all systems. Additionally, the resource overhead associated with containerized environments may constrain adoption in highly resource-limited settings. These factors suggest that the approach is most effective when applied selectively to critical workflows rather than indiscriminately across all components.

Future research can extend this work by exploring more adaptive testing strategies that respond dynamically to observed system behavior. Techniques that adjust assertion windows or observation criteria based on runtime signals may further improve resilience to variability. There is also potential to investigate how fault injection and resilience testing can be more tightly integrated into system-level validation frameworks, expanding coverage of failure modes without increasing flakiness.

Another promising direction involves studying the human aspects of system-level testing adoption. Understanding how developers interpret and act on observability-driven test feedback can inform improvements in tooling and documentation. Empirical studies of team practices may reveal patterns that influence test effectiveness and long-term maintenance. Such insights would complement technical advances with organizational learning.

In conclusion, this research demonstrates that reproducible containerized environments provide a practical foundation for reliable system-level testing of event-driven microservices. By combining architectural discipline, event-aware assertions, and observability-driven diagnostics, the proposed approach advances both engineering practice and research understanding. Continued exploration of adaptive techniques and human-centred factors offers a pathway for further strengthening validation strategies in increasingly complex distributed systems.

REFERENCES

1. Brewer, E. A. (2012). CAP twelve years later: How the rules have changed. *Computer*, 45(2), 23–29. <https://doi.org/10.1109/MC.2012.37>
2. Vogels, W. (2009). Eventually consistent. *Communications of the ACM*, 52(1), 40–44. <https://doi.org/10.1145/1435417.1435432>
3. Dean, J., & Barroso, L. A. (2013). The tail at scale. *Communications of the ACM*, 56(2), 74–80. <https://doi.org/10.1145/2408776.2408794>
4. Parasa, M. (2020). Designing future ready compensation systems with data driven fairness and performance alignment in SAP SuccessFactors. *International Journal of Scientific Research and Engineering Trends*, 6(4). <https://doi.org/10.5281/zenodo.17698304>.
5. Dragoni, N., Giallorenzo, S., Lluch-Lafuente, A., Mazzara, M., Montesi, F., Mustafin, R., & Safina, L. (2017). Microservices: Yesterday, today, and tomorrow. In *Present and Ulterior Software Engineering* (pp. 195–216). Springer. https://doi.org/10.1007/978-3-319-67425-4_12
6. Lehvä, J., Mäkitalo, N., & Mikkonen, T. (2019). Consumer-driven contract tests for microservices: A case study. In *Product-Focused Software Process Improvement (PROFES)* (pp. 497–512). Springer. https://doi.org/10.1007/978-3-030-35333-9_35
7. Luo, Q., Hariri, F., Eloussi, L., & Marinov, D. (2014). An empirical analysis of flaky tests. *Proceedings of the ACM SIGSOFT Symposium on the Foundations of Software Engineering*. <https://doi.org/10.1145/2635868.2635920>
8. Nithin Nanchari. (2020). The Role of Internet of Things (IoT) in Healthcare. *European Journal of Advances in Engineering and Technology*, 7(4), 67–69. Zenodo. <https://doi.org/10.5281/zenodo.15968914>
9. Ghaleb, T. A., da Costa, D. A., & Zou, Y. (2019). An empirical study of the long duration of continuous integration builds. *Empirical Software Engineering*, 24, 2102–2139. <https://doi.org/10.1007/s10664-019-09695-9>
10. Rahman, A., Rigby, P. C., & Adams, B. (2018). Characterizing the influence of continuous integration: Empirical results from 250+ open source and proprietary projects. *Proceedings of the International Conference on Automated*

- Software Engineering.
<https://doi.org/10.1145/3278142.3278149>
11. Ståhl, D., & Bosch, J. (2014). Modeling continuous integration practice differences in industry software development. *Journal of Systems and Software*, 87, 48–59. <https://doi.org/10.1016/j.jss.2013.08.032>
 12. Di Tommaso, P., Palumbo, E., Chatzou, M., Prieto, P., Heuer, M. L., & Notredame, C. (2015). The impact of Docker containers on the performance of genomic pipelines. *PeerJ*, 3, e1273. <https://doi.org/10.7717/peerj.1273>
 13. Kozhimbayev, Z., & Sinnott, R. O. (2017). A performance comparison of container-based technologies for the cloud. *Future Generation Computer Systems*, 68, 175–182. <https://doi.org/10.1016/j.future.2016.08.025>
 14. Medel, V., Rana, O. F., Bañares, J. A., & Arronategui, U. (2016). Modelling performance and resource management in Kubernetes. *Proceedings of the International Conference on Utility and Cloud Computing*. <https://doi.org/10.1145/2996890.3007869>
 15. Schwarzkopf, M., Konwinski, A., Abd-El-Malek, M., & Wilkes, J. (2013). Omega: Flexible, scalable schedulers for large compute clusters. *Proceedings of the ACM European Conference on Computer Systems (EuroSys)*. <https://doi.org/10.1145/2465351.2465386>
 16. Verma, A., Pedrosa, L., Korupolu, M., Oppenheimer, D., Tune, E., & Wilkes, J. (2015). Large-scale cluster management at Google with Borg. *Proceedings of the ACM European Conference on Computer Systems (EuroSys)*. <https://doi.org/10.1145/2741948.2741964>
 17. Chandola, V., Banerjee, A., & Kumar, V. (2009). Anomaly detection: A survey. *ACM Computing Surveys*, 41(3), Article 15. <https://doi.org/10.1145/1541880.1541882>
 18. Padur, S. K. R. (2020). From centralized control to democratized insights: Migrating enterprise reporting from IBM Cognos to Microsoft Power BI. *International Journal of Scientific Research in Computer Science, Engineering and Information Technology*, 6(1), 218–225. <https://doi.org/10.32628/CSEIT2390625>
 19. He, S., Zhu, J., He, P., & Lyu, M. R. (2016). Experience report: System log analysis for anomaly detection. *IEEE International Symposium on Software Reliability Engineering (ISSRE)*. <https://doi.org/10.1109/ISSRE.2016.21>
 20. Routhu, K. K. (2018). Seamless HR finance interoperability: A unified framework through Oracle Integration Cloud. *International Journal of Science, Engineering and Technology*, 6(1). <https://doi.org/10.5281/zenodo.17292100>
 21. Bauer, A., Leucker, M., & Schallhart, C. (2011). Runtime verification for LTL and TLTL. *ACM Transactions on Software Engineering and Methodology*, 20(4), Article 14. <https://doi.org/10.1145/2000799.2000800>
 22. Sudhir Vishnubhatla. (2018). From Risk Principles to Runtime Defenses: Security and Governance Frameworks for Big Data in Finance. In *International Journal of Science, Engineering and Technology (Vol. 6, Number 1)*. Zenodo. <https://doi.org/10.5281/zenodo.17452405>
 23. Uргаonkar, B., Shenoy, P., Chandra, A., Goyal, P., & Wood, T. (2008). Agile dynamic provisioning of multi-tier internet applications. *ACM Transactions on Autonomous and Adaptive Systems*, 3(1), Article 1. <https://doi.org/10.1145/1342171.1342172>
 24. Barroso, L. A., & Hölzle, U. (2007). The case for energy-proportional computing. *Computer*, 40(12), 33–37. <https://doi.org/10.1109/MC.2007.443>