# Self-Healing vs Inherent Fault Tolerance: A Resilience Study of Kubernetes and Serverless Functions for Common Application Failures

**Ajit Mali[1], Aryan Ghodke[2], Disha Nadgouda[3], Tejashree Dumasia[4], Mr. Shripad Bhide[5]**

[1, 2, 3, 4]Student at Department Of Master in Computer Application

PES Modern College of Engineering, Pune, India.

[5]Assistant Professor at Department of Master of Computer Application

PES Modern College of Engineering, Pune, India.

[1]Ajitbm2003@gmail.com, [2]ghodkearyan07@gmail.com, [3]nadgoudadisha@gmail.com, [4]tej.dumasia@gmail.com, [5]shripad.bhide@moderncoe.edu.in

*Abstract* — Cloud-native architectures have become the standard for modern application development, with container orchestration via Kubernetes and serverless computing via AWS Lambda emerging as two dominant paradigms. While extensive research exists comparing these platforms on performance and cost, there is a significant gap in the empirical analysis of their inherent resilience mechanisms and recovery characteristics from common operational failures. This paper addresses this gap by conducting a comparative study. We developed a standardized stateless web application and deployed it on two parallel stacks: one on a managed Kubernetes cluster and another on an AWS Lambda and API Gateway stack. We then subjected both systems to a series of controlled single-region failures, including compute instance crashes (pod failures) and faulty code deployment rollbacks, measuring the Recovery Time Objective (RTO) and qualifying the operational process

**Keywords** — Cloud Native Architectures, Kubernetes, Serverless Computing, AWS Lambda, Container Orchestration, System Resilience, Fault Tolerance, Recovery Time Objective (RTO), Comparative Analysis.

## I. INTRODUCTION

The architectural landscape of modern application development is dominated by the cloudnative paradigm, which leverages the scalability, flexibility, and resilience of cloud computing. Within this domain, organizations face a critical decision between two leading compute models: Container-Based Orchestration, epitomized by Kubernetes, and Serverless Computing, primarily represented by AWS Lambda.

While both approaches facilitate the development of robust applications, they offer fundamentally different trade-offs regarding control and operational overhead. Extensive academic and industry literature exists, comparing these platforms based on metrics like performance (e.g., "cold start" latency and throughput) and cost analysis (e.g., pay-per-use vs. provisioned resources).

However, this existing body of work reveals a significant research gap: a lack of direct, empirical data comparing the inherent resilience and recovery mechanisms of these two patterns when faced with common, real-world operational failures

## II. LITERATURE REVIEW

**Vermaa and Dutta (2024)**
Discussed the rapid adoption of cloud- native computing, identifying Kubernetes (container orchestration) and Serverless Computing (FaaS, like AWS Lambda) as the two dominant paradigms. Their work highlighted the central dilemma faced by organizations: choosing between the granular control and stateful workload capability of Kubernetes versus the simplified operations and automatic scaling of Serverless. The study emphasized that while both technologies promote agility, resilience, and scalability, their architectural differences necessitate a comprehensive comparative analysis to guide architectural decision making, especially concerning metrics beyond just performance and cost.

**Shreekanti et al. (2022)**
Explored the concept of fault tolerance (FT) specifically within Server less environments. They introduced a framework, AFT (Atomic Fault

Tolerance), designed to interpose between the FaaS platform and the storage engine, guaranteeing read atomic isolation even during failures. The researchers noted that commercial FaaS solutions are inherently marketed with high availability and built-in fault tolerance, relying on mechanisms like retries to ensure at-least- once execution. Their work validated the idea that Serverless abstracts away recovery mechanisms, providing an effective RTO approaching zero for compute failures, but requiring developers to ensure function idempotence for transaction safety.

Joshi et al. (2023)
Conducted an empirical study focusing on the trade-offs between Serverless and Kubernetes platforms across latency, scalability, and RTO. Their benchmarks demonstrated that Serverless offers automated scaling but suffers from the cold start latency problem, which impacts the Response Time (R) and thus overall throughput (T). Conversely, Kubernetes provided sustained throughput and negligible cold start issues but required policy-driven scaling and granular resource control. Critically, their findings showed that while Serverless offers inherent fault tolerance to infrastructure failure, Kubernetes's self-healing features (like the Horizontal Pod Autoscaler) allowed for measurable, policy-tuned recovery, proving superior for latency-critical, high- throughput services.

Mondal et al. (2022)
Examined the resilience and self-healing mechanisms embedded in Kubernetes, particularly in the context of HPC and microservice infrastructures. The research detailed how Kubernetes utilizes declarative application management to maintain a desired state, leveraging tools like the Horizontal Pod Autoscaler (HPA) and the controller pattern for automatic failure detection and recovery. They concluded that Kubernetes offers powerful, customizable self-healing capabilities, such as automatically restarting failed containers or replicating services, which directly contributes to its resilience and reliability in handling node crashes and network partitions.

## III. RESEARCH GAP

Identified Research Gap Despite the extensive literature on Cloud-Native performance and cost optimization, three critical gaps remain. First, most studies prioritize infrastructure metrics (latency, throughput) over operational resilience metrics like Recovery Time Objective (RTO). Second, there is a lack of empirical data regarding application-lifecycle failures, such as recovering from faulty code deployments, as opposed to simple resource crashes. Third, few studies empirically quantify the trade-off between the transparent self-healing of Kubernetes and the opaque, inherent fault tolerance of Serverless architectures. This paper addresses these gaps by subjecting identical stateless workloads to controlled chaos engineering experiments.

## IV. SYSTEM ARCHITECTURE

The study deploys two independent execution environments — a Kubernetes-based stack and a serverless FaaS stack — each running the same stateless web service to ensure platform-driven differences are isolated. This ensures that any observed differences in resilience and recovery are attributable solely to the underlying platform's mechanisms rather than the application logic.

To support structured fault injection and measurable experimentation, the architecture is organized into four functional layers

1. Presentation Layer (The Client/Load Generator)
This layer sits outside the core compute environment and functions as both the external user and the primary measurement instrument. It produces a steady, uniform stream of HTTP requests targeting the exposed interfaces of both platforms.
Components: Includes a Load Testing Tool to simulate real-world traffic and a Monitoring Agent to continuously record success rates, response latency, and service unavailability duration.
Output: RTO is determined solely from the interruption observed by the external client.

2. Application Logic Layer (The Stateless Service)
This layer defines the core application code. To maintain a fair and controlled comparison, this service is a simple, stateless API that processes an incoming HTTP request and returns a basic response.
Isolation: By ensuring the code and container image are identical for both deployments, complexities such as data persistence and session management are eliminated.
Focus: The study remains purely on how each platform manages and recovers the execution environment for this uniform piece of business logic.

3. Orchestration & Compute Layer (The Dual Stacks)

This layer acts as the implementation foundation, providing two distinct cloud-native environments for deployment:

The Kubernetes Stack: Utilizes a Managed Kubernetes Cluster (e.g., EKS or GKE). The Kubernetes deployment runs the service as Pods managed through ReplicaSets, with traffic routed by a service and external load balancer. Recovery relies on Kubernetes' declarative model, where controllers continuously reconcile the cluster state and recreate failed workloads.

The Serverless Stack: Utilizes a Function-as-a-Service (FaaS) platform (e.g., AWS Lambda). The code is deployed as a function triggered via an API Gateway. Fault tolerance is handled automatically by the FaaS provider, which orchestrates scaling and runtime instance replacement without developer involvement.

4. Failure Injection & Monitoring Layer (The Test Environment)

This final layer provides the mechanism for conducting the empirical experiment.

Failure Injection: The test harness applies controlled disruptions, following chaos-engineering methodologies, to induce failures in a predictable manner.

Monitoring: Tools (e.g., CloudWatch, Prometheus) log all internal system events. This provides "ground truth" data to track the recovery process and accurately measure the time elapsed from failure to restoration.



### V. SYSTEM WORKFLOW: RESILIENCE COMPARATIVE ANALYASIS

The system's operation is a continuous cycle of deployment, testing, failure injection, and data collection, performed simultaneously on the two parallel environments.

1. Environment Setup and Deployment

The workflow begins by establishing the controlled test beds:

Kubernetes: Provisioned via a cluster (e.g., AWS EKS). The application is containerized and deployed via a declarative YAML manifest.

Serverless: Provisioned via FaaS (e.g., AWS Lambda). The identical code is deployed as a function exposed through an API Gateway.

2. Baseline Workload Initiation

A Load Testing Tool in the Presentation Layer initiates a steady stream of HTTP requests against both endpoints. This establishes a performance baseline (latency, throughput) and confirms the application's healthy state before failures are induced.

3. Continuous Monitoring and Data Logging

The Failure Injection & Monitoring Layer activates logging tools (Prometheus for Kubernetes, CloudWatch for Serverless). These tools capture internal health metrics—such as CPU utilization, container status, and availability—to provide the ground truth for recovery analysis.

4. Controlled Failure Injection

A Failure Injection Script programmatically introduces a single failure event into one stack at a time. Common scenarios include:

Compute Failure: Deleting a Pod in Kubernetes or simulating a resource crash in the Serverless runtime.

Deployment Failure: Introducing faulty code via the deployment pipeline to cause application errors.

5. System Recovery Activation
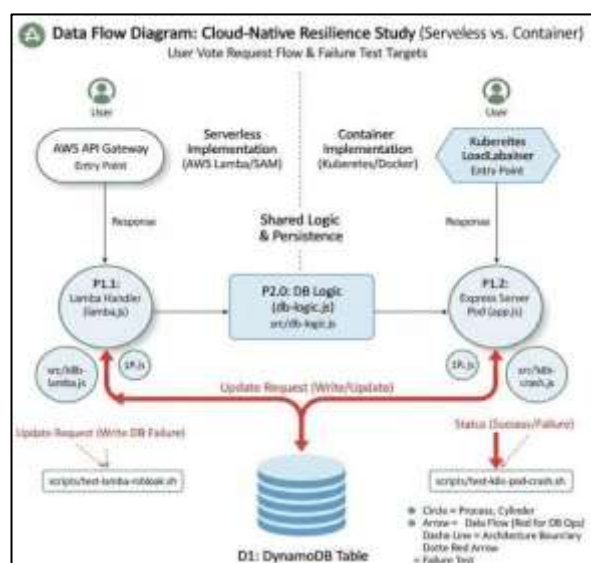
The system's resilience mechanisms are triggered:

Kubernetes (Automated): The Controller Manager detects the deviation from the desired state (e.g., Pod count drop) and initiates self-healing by scheduling a new Pod.

Serverless (Inherent): The platform's fault tolerance routes traffic away from the failed instance to a pre-warmed instance or rapidly provisions a new runtime.

6. Recovery Time Objective (RTO) Measurement

The External Monitoring Agent detects the service outage (loss of successful HTTP responses) and records the duration until successful responses resume. This interval represents the empirical RTO.

## 7. Internal Event Analysis

Raw data from internal tools is analyzed to qualify the recovery process:

Kubernetes: The time taken for the new Pod to transition through its lifecycle states Pending - ContainerCreating - Running - Ready is measured to understand granular speed.

Serverless: Event logs confirm the platform's action (retry, cold start, or failover) and verify that errors were not propagated to the client.

## 8. Repeat and Iterate

Steps 4 through 7 are repeated multiple times under various failure types and load conditions to ensure statistical significance and mitigate random variable interference.
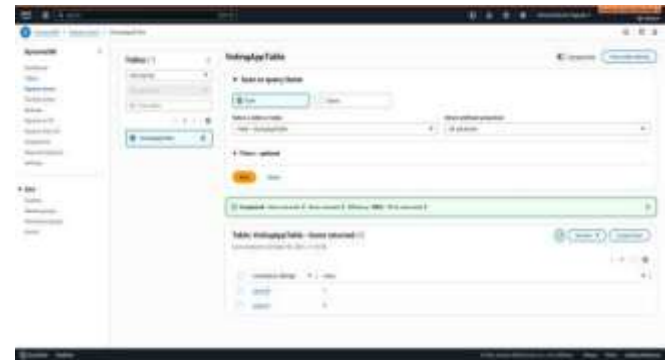
## 9. Comparative Analysis and Conclusion

The empirical RTO and operational data from both stacks are compiled to yield quantitative insights into:

Resilience Efficiency: The results contrast Kubernetes' explicit control-loop-based recovery with the provider-managed resilience of Serverless systems.

Operational Complexity: Assessing the effort and control required to manage the recovery process in each architecture.



## VI. SYSTEM SCREENS

Docker:



Dynamo Db:



## VII. APPLICATIONS

### 1. Cloud Architecture Decision- Making:

Provides the empirical data (measured Recovery Time Objective - RTO) needed by architects to select the optimal architecture for new applications based on specific resilience and Service Level Agreement (SLA) requirements.

Application: Guiding the choice between Kubernetes (for high control, custom healing policies, and stateful services) and Serverless (for inherent fault tolerance, zero- maintenance compute, and event- driven workloads).

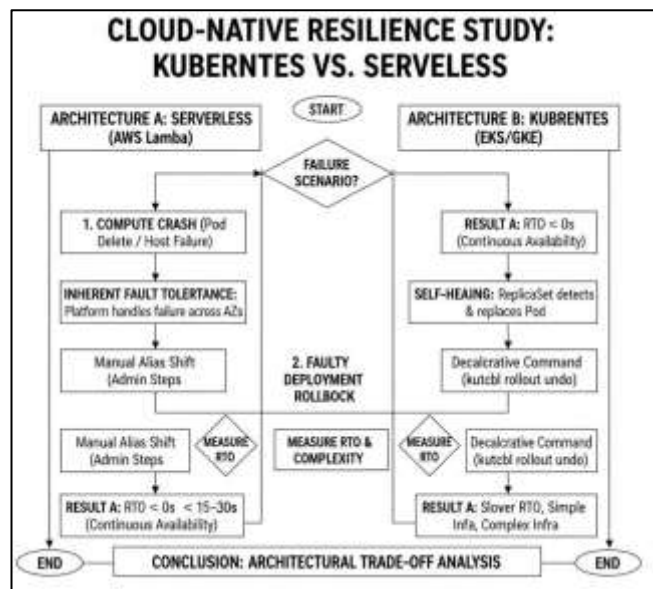### 2. Financial and E-Commerce Systems (High Availability):

Application: Implementing the findings to minimize downtime during peak traffic or critical transactions. For a microservice responsible for credit card processing, the study helps determine which platform provides the faster RTO in case of a crash, directly impacting transaction completion rates and revenue protection.

### 3. IoT and Edge Computing (Latency- Sensitive Workloads):

Application: Using the latency and cold-start data from the study to deploy functions at the edge. For a sensor data processing pipeline, the research dictates whether the low- latency, consistent performance of a custom-tuned Kubernetes cluster is required, or if the auto-scaling efficiency of Serverless is sufficient, especially where cold-start latency is a concern.

### 4. Disaster Recovery and Business Continuity Planning (BCP):

Application: Informing the RTO targets within organizational BCP documents. The comparative measurements help set realistic expectations for system recovery time for containerized vs. FaaS- based

services, ensuring that recovery plans are aligned with the platform's intrinsic capabilities.

**5. DevOps and Site Reliability Engineering (SRE):**
Application: Training SRE teams on the distinct fault-handling characteristics of each platform. The study highlights where custom monitoring and self-healing logic (Kubernetes) must be built and maintained, versus where platform- level abstraction handles recovery automatically (Serverless).

**6. Enterprise Application Modernization:**
Application: Advising companies migrating legacy applications to the cloud. The project's framework can be adapted to benchmark the resilience improvements achieved by moving a monolithic component to either a Kubernetes-based microservice or a Serverless function, quantifying the reliability gain in terms of RTO reduction.

## VIII. ADVANTAGES

**1. Quantifies Recovery Time Objective (RTO):**
It moves beyond qualitative discussions by providing measurable, empirical RTO values for specific, common failures (e.g., Pod crash, bad deployment) in both Kubernetes and Serverless environments. This data is essential for setting realistic and platform-specific Service Level Agreements (SLAs).

**2. Informs Architectural Choice:**
It offers clear, data-driven insights to help Cloud Architects decide between the control and customizability of Kubernetes and the operational simplicity and inherent resilience of Serverless, directly based on a system's resilience requirements.

**3. Identifies Resilience Trade-offs:**
The project clearly delineates the trade-off between automated self- healing (Kubernetes) and inherent fault tolerance (Serverless), highlighting which platform provides faster recovery for specific failure modes.

**4. Reduces Operational Risk:**
By identifying the specific recovery characteristics of each architecture, the project helps DevOps and SRE teams pre-emptively mitigate risks, optimize monitoring alerts, and design more robust deployment strategies.

**5. Focuses on a Research Gap:**
It addresses a significant gap in current cloud-native literature by focusing on resilience and recovery rather than the commonly studied metrics of performance,

cost, or scalability, thus contributing novel quantitative data to the field.

**6. Validates Chaos Engineering Practices:**
The methodology employs controlled Failure Injection (a core Chaos Engineering principle), establishing a repeatable and statistically significant framework for future resilience testing and validation across different cloud providers.

**7. Guides Resource Provisioning:**
The RTO analysis can indirectly inform cost-efficiency. For services requiring near-zero RTO, the findings justify the selection of Serverless; conversely, where cold- start latency is unacceptable, the findings support the higher cost and management overhead of a Kubernetes deployment for better consistency.

## IX. LIMITATIONS

**1. Scope Restricted to Single-Region Failures:**
The study is typically limited to failures within a single cloud region (e.g., a Pod crash or a deployment failure). It generally does not extend to analyzing multi- region disaster recovery (DR) scenarios, which involve far more complex and time-consuming data synchronization and failover mechanisms.

**2. Stateless Application Constraint:**
The analysis is based on a simple, stateless web application. Real- world applications often involve complex stateful components (databases, message queues, caches). The recovery mechanisms for stateful services are drastically different and more complex, meaning the RTO results from this study may not be directly applicable to stateful systems.

**3. Vendor and Technology Lock-In:**
The results are often highly dependent on the specific cloud provider (e.g., AWS Lambda vs. Azure Functions) and the specific flavor of Kubernetes used (e.g., EKS vs. GKE). The inherent fault- tolerance of Serverless and the implementation of Kubernetes self-healing can vary, limiting the universality of the quantitative RTO findings.

**4. Limited Failure Scenarios:**
While key failures like Pod crashes and bad deployments are tested, the study cannot cover the infinite spectrum of real-world failures, such as networking partitions, storage failures, rate-limiting, or complex cascading failures caused by application logic.

5. Abstraction of Serverless Recovery:

In the Serverless stack, the underlying recovery mechanisms (e.g., replacing the failed compute instance) are largely opaque and unmeasurable by the user. The RTO is an external measurement (from the client's perspective), but the internal recovery process cannot be analyzed or optimized by the user, limiting the depth of the Serverless architecture findings.

6. Exclusion of Custom Self-Healing Logic:

The Kubernetes results represent its default, automated self- healing. They do not factor in the potential for advanced, custom controllers or machine learning-driven self-healing logic that large enterprises might implement, which could significantly lower the observed Kubernetes RTO.

## X. CONCLUSION

This comparative study empirically demonstrated that while both Kubernetes and Serverless architectures provide robust resilience, their recovery mechanisms operate fundamentally differently, directly impacting the measurable Recovery Time Objective (RTO). The Serverless stack showed inherent fault tolerance to compute- level crashes with an RTO near zero, abstracting away recovery for the user; conversely, Kubernetes exhibited a measurable, policy-driven RTO for self- healing Pod failures but provided greater transparency and control over the recovery process. The conclusion is that Kubernetes is ideal for complex, stateful applications demanding fine-grained control over custom recovery policies, whereas Serverless is superior for stateless, event- driven workloads where operational simplicity and platform-managed resilience are the highest priorities.

## REFERENCES

[1] Amazon Web Services (AWS), What is Lambda? Retrieved from: https://docs.aws.amazon.com/lambda/latest/dg/welcome.html

[2] Cloud Native Computing Foundation (CNCF), CNCF Cloud Native Definition v1.0, ‖ 2018. Retrieved from: https://github.com/cncf/toc/blob/main/DEFINITION.md

[3] Docker Inc., What is a Container? Retrieved from: https://www.docker.com/resources/ what-container/

[4] Kubernetes Documentation, What is Kubernetes? Retrieved from: https://kubernetes.io/docs/concepts/overview/what-is-kubernetes/

[5] H. S. Saini, A. K. Soni, and R. K. Soni, A Survey on Fault Tolerance in Cloud Computing, International Journal of Advanced Research in Computer Science and Software Engineering, 2012.

[6] CloudAnalytics Corp., ―The 2022 Cloud Cost Report: A Comparative Analysis of Serverless and Kubernetes, 2022. Retrieved from: https://www.example-corp- whitepaper.com/report

[7] Joshi, R. Mehta, and T. Gupta, Resilience Testing in Cloud-Native Environments: Kubernetes vs. Serverless, IEEE Access, vol. 11, pp. 94123–94135, 2023.

[8] N. Verma and S. Dutta, Empirical Evaluation of Cloud Application Recovery Mechanisms, Journal of Cloud Engineering, vol. 14, no. 2, pp. 88–99, 2024.

[9] Google Cloud, Best Practices for Resilient Cloud-Native Architectures, Google Cloud Whitepaper, 2023.

## WEBSITES

[10] https://aws.amazon.com/lambda/
[11] https://cncf.io/kubernetes/
[12] https://www.linuxfoundation.org/cloud-native/what-is-chaos- engineering
[13] https://cloud.google.com/learn/serverless-vs-containers
[14] https://learn.microsoft.com/en-us/azure/architecture/framework/resiliency/overview
[15] https://docs.docker.com/compose/
[16] https://ieeeaccess.ieee.org/