



ARCHITECTURAL TRADE-OFFS BETWEEN STATELESS AND STATEFUL MICROSERVICES IN LARGE-SCALE CLOUD SYSTEMS

Anup Rao

Software Engineer 2, Microsoft, Redmond, WA, USA, ANUP.RAO@microsoft.com

ORCID : 0009-0008-7306-1046

Abstract

Through a controlled experimental review, this study examined the architectural trade-offs of stateless and stateful microservices in large-scale cloud systems. To enable a fair comparison, Kubernetes was used to implement parallel installations of both designs on AWS with similar hardware specs. Three scenarios—normal load, peak load, and failure conditions—were used to evaluate performance, scalability, fault tolerance, and operational overhead. The findings showed that, especially in situations with fault injection and high traffic, stateless microservices produced higher throughput, quicker recovery times, and reduced operational complexity. On the other hand, stateful microservices performed exceptionally well in terms of transactional consistency and session retention, albeit at the expense of slower scalability and more maintenance work. According to the results, a hybrid architectural strategy that combines externalized state management with stateless service layers may offer large-scale cloud deployments the best possible balance between data integrity, flexibility, and resilience.

Keywords: *Stateless microservices, Stateful microservices, Cloud computing, Scalability, Fault tolerance, Kubernetes, AWS, Distributed systems, Architectural trade-offs, Hybrid architecture.*

1. INTRODUCTION

Microservices architecture, which offers modularity, scalability, and operational flexibility, has emerged as a major trend in the design of large-scale cloud systems in recent years. The choice to implement services as stateless or stateful components was a crucial architectural factor in this paradigm. The benefits of stateless microservices, which did not save client-specific information between requests, were evident in their flexibility, fault recovery, and deployment simplicity. These services are ideal for workloads with varying demand because they may be horizontally scaled with no coordination overhead. However, they had to rely on external storage or caching solutions due to their inability to retain persistent state internally, which could have added complexity to the system and caused network latency.

Conversely, stateful microservices had maintained internal data across requests, enabling rich session-based interactions, transactional consistency, and complex workflow management. While this state retention had been essential for applications such as financial processing, collaborative platforms, and real-time data analytics, it had also introduced challenges in scaling, recovery, and data synchronization. Stateful services had required more sophisticated orchestration, persistent storage integration, and replication strategies, making their management inherently more complex in distributed environments.

The trade-off between stateless and stateful architectures had therefore represented a critical decision point for system architects, influencing not only performance and scalability but also

resilience, operational overhead, and long-term maintainability. While prior research had explored microservices performance optimization and cloud-native resilience, comparative empirical studies focusing specifically on the scalability, fault tolerance, and operational implications of these two architectural patterns in large-scale cloud deployments had been limited. This study had aimed to address that gap by systematically evaluating the performance, recovery characteristics, and operational demands of stateless and stateful microservices within a controlled cloud environment, offering insights to guide future architectural decision-making.

2. LITERATURE REVIEW

Andell (2020) had examined the architectural ramifications of Function-as-a-Service (FaaS) and serverless computing, highlighting their fundamentally stateless characteristics. According to the study, serverless designs made it possible for highly scalable execution models and eliminated the strain of managing infrastructure. However, these functions' lack of persistent state forced them to rely on external storage systems, which could present latency and consistency issues. Externalizing state was a crucial design choice for stateless microservices, where these insights were especially pertinent.

Chowdhury, Salahuddin, Limam, and Boutaba (2019) has offered a thorough examination of the use of microservices to re-architect Network Function Virtualization (NFV) ecosystems. The advantages of breaking down monolithic VNFs into smaller, loosely linked components to improve flexibility and scalability were described in their study. Crucially, when moving to microservices-based NFV, they found that state management was still a major obstacle, particularly when it came to preserving data integrity and service continuity across dispersed network functions. The trade-offs between stateful and stateless microservices in large-scale systems were directly analogous to this.

Furda, Fidge, Zimmermann, Kelly, and Barros (2017) had looked into moving enterprise traditional systems to microservices, with a particular emphasis on data consistency, statefulness, and multitenancy. According to the authors, stateful microservices frequently added complications in synchronization, scaling, and fault recovery, even while they were useful for preserving user context and session data. Stateless systems, on the other hand, were easier to scale and deploy, but they needed strong externalized state methods to manage ongoing data needs. Their results supported the idea that state management techniques were essential to the maintainability and performance of microservices.

Kang, Le, and Tao (2016) has looked at cloud infrastructure DevOps procedures in relation to container-based microservice designs. They have shown that containers allowed for scalability, isolation, and quick deployment when paired with microservices. Although operational agility had been their main concern, they also saw that container orchestration tactics were greatly impacted by the decision between stateful and stateless designs. While stateful services needed careful orchestration to maintain data integrity during scaling or failure recovery, stateless services were simpler to duplicate and replace.

Thumala (2020) We discussed how to create extremely resilient cloud infrastructures, with a focus on methods to guarantee service continuation in the face of adversity. The study made clear that the architectural decision between stateful and stateless systems frequently determined resilience. While stateful services are more difficult to recover from, they were required for some transactional and session-dependent workloads. In contrast, stateless services

have been demonstrated to recover from failures more quickly because they rely less on permanent state.

3. RESEARCH METHODOLOGY

Large-scale cloud systems' quick adoption of microservices architecture has revolutionized how businesses develop, implement, and grow their applications. In such systems, deciding whether to construct services as stateless or stateful components was a crucial architectural choice. Stateless microservices provided advantages including easier deployment, elastic scalability, and enhanced fault tolerance because they did not store any client-specific data between requests. Stateful microservices, on the other hand, preserved session or transactional data between requests, allowing for more functionality but posing problems with data consistency, fault recovery, and scaling. Performance, dependability, resource usage, and operational complexity were all significantly impacted by this architectural trade-off. Prior studies had looked at resilience and performance optimization in microservices, but there wasn't much practical data comparing stateful versus stateless methods in the same large-scale cloud environment. By empirically assessing these trade-offs in a regulated cloud-based environment, this work sought to close this gap.

3.1. Research Design

Two parallel microservices architectures, one stateless and one stateful, were deployed in a controlled cloud infrastructure under identical settings as part of the study's comparative experimental research approach. By removing external variable factors, this method had made it possible to assess and compare performance, scalability, and resilience qualities directly.

3.2. Experimental Environment

Using Kubernetes as the orchestration framework, all experiments were carried out on Amazon Web Services (AWS). Docker was used to containerize application services in order to guarantee deployment consistency. A PostgreSQL database was used for service state maintenance and AWS Elastic Block Store (EBS) for permanent storage in the stateful configuration. The stateless setup relied on transitory caching rather than any long-term storage techniques, processing queries completely in-memory. To guarantee fairness in performance evaluation, the virtual machine characteristics in all contexts were the same.

3.3. Data Collection Methods

Grafana for real-time visualization and Prometheus for system metrics were used to gather performance statistics. Throughout the testing process, throughput, response times, CPU utilization, and memory use were all continuously tracked. Failures had been purposefully introduced using Chaos Monkey to simulate node crashes, network outages, and forced restarts in order to evaluate fault tolerance. Using Apache JMeter, scalability testing was carried out by gradually increasing workloads from 500 to 50,000 concurrent users. Additionally, throughout a 30-day evaluation period, deployment durations, configuration modifications, and debugging activities were tracked in order to quantify operational overhead.

3.4. Test Scenarios

For both architectures, three different test scenarios had been run. A Normal Load Scenario with 1,000 concurrent users and no fault injections was the first. In order to assess scaling behavior, the second was a Peak Load Scenario that replicated traffic spikes up to 50,000 concurrent users. The third was a Failure Scenario, which evaluated recovery time and resilience by randomly terminating nodes and injecting lag during moderate workloads.

3.5.Data Analysis

To find important distinctions between the two systems, collected data was processed and statistically examined. Paired t-tests were used to assess quantitative performance outcomes and ascertain whether differences in latency, throughput, and resource use were statistically significant. To find reoccurring trends and problems, theme analysis was used to examine qualitative data from operational overhead observations. For clarity, the results were displayed using tables and comparative charts.

3.6.Ethical Considerations

Since there were no human subjects in the study, ethical concerns were reduced. However, in order to guarantee adherence to organizational regulations and financial limitations, all cloud resources had been made available through institutional accounts. To avoid inadvertently affecting production systems, testing has been limited to isolated areas.

4. RESULTS AND DISCUSSION

In large-scale cloud environments, the experimental evaluation's findings demonstrated clear distinctions between stateless and stateful microservices with regard to scalability, fault tolerance, operational overhead, and resource consumption. Stateful microservices demonstrated advantages in managing session-based workloads and preserving transactional consistency, while stateless microservices consistently achieved higher throughput and faster recovery from failures, according to data gathered from the controlled AWS Kubernetes deployments. The three predetermined test scenarios' complete results, backed by statistical analysis and visual data visualization, are presented in the ensuing subsections.

4.1.Performance Under Normal Load

The stateless design demonstrated improved average throughput, handling about 8.2% more requests per second than the stateful setup under typical load conditions (1,000 concurrent users, no failure injections). Due to the lack of persistent storage operations, stateless services have continuously had reduced average response times.

Table 1: Performance Metrics – Normal Load Scenario

Metric	Stateless Microservices	Stateful Microservices	Difference (%)
Throughput (req/sec)	12,450	11,500	+8.2
Average Response Time (ms)	95	112	-15.2
CPU Utilization (%)	62	66	-6.0
Memory Utilization (%)	48	55	-12.7

Since stateless services avoided the latency caused by database interactions, these outcomes were in line with predictions. Nonetheless, under continuous, moderate load, the performance difference was not significant, indicating that both topologies were feasible.

4.2. Scalability Under Peak Load

The stateless version showed noticeably improved elasticity during the peak load scenario (50,000 concurrent users), scalable horizontally with little increase in latency. During traffic spikes, the stateful design, which was limited by persistent storage I/O, had a 37% higher average response time and occasionally encountered queue delays.

Table 2: Scalability Metrics – Peak Load Scenario

Metric	Stateless Microservices	Stateful Microservices	Difference (%)
Throughput (req/sec)	48,200	42,700	+12.9
Average Response Time (ms)	155	212	-37.0
Horizontal Scaling Time (s)	35	51	-31.3
Error Rate (%)	0.4	1.2	-66.7

Stateless services were better suited for workloads with erratic traffic patterns, according to the findings. Because there was no data migration or state synchronization during container replication, the scaling times were faster.

4.3. Resilience and Fault Recovery

Because the stateless design may restart services without state restoration, it recovered 42% faster on average in failure scenarios with injected node terminations and network latency spikes. Stateful microservices reduced data loss in transactional operations by preserving ongoing sessions, despite their delayed recovery time.

Table 3: Fault Tolerance Metrics – Failure Scenario

Metric	Stateless Microservices	Stateful Microservices	Difference (%)
Mean Time to Recovery (MTTR, s)	28	48	-41.7
Data Loss Incidents	0	0	N/A
Session Preservation (%)	0	100	N/A
Post-Recovery Error Rate (%)	0.6	0.9	-33.3

These findings highlighted the trade-off: stateless services excelled in recovery speed, whereas stateful services offered continuity for long-lived sessions.

4.4. Operational Overhead

During the 30-day study period, operational monitoring revealed that stateless installations needed 27% fewer maintenance work. The lack of state management made tasks like scalability, upgrades, and container restarts easier. The stateful architecture, on the other hand, required careful control of volume, replication consistency, and database migrations.

Table 4: Operational Overhead Summary

Parameter	Stateless Microservices	Stateful Microservices	Difference (%)
Avg. Deployment Time (min)	4.8	7.1	-32.4

Avg. Debugging Time per Issue (min)	18.2	24.7	-26.3
Maintenance Incidents	6	9	-33.3

These results reinforced the perception that stateless architectures reduced day-to-day operational complexity, making them favorable for rapidly evolving environments.

Discussion

According to the study's findings, stateless microservices are perfect for highly elastic and fault-tolerant systems because they offer better scalability, quicker recovery, and fewer operational overhead. They did not, however, have built-in support for complicated transactional workflows or client sessions without the need of additional state management tools. Stateful microservices, on the other hand, were better at managing persistent, session-dependent tasks, which made them useful for use cases involving long-running processes, financial transactions, or real-time collaboration apps.

Under high load and failure scenarios, the performance disparity increased, and stateless services took advantage of their more straightforward architecture to grow and recover more quickly. However, these advantages came at the cost of shifting state management to outside systems, which might have increased architectural complexity in other places.

Practically speaking, the findings indicated that hybrid techniques, in which application services maintain their statelessness but important state is externalized to specialized data stores, would provide the optimal balance between scalability and consistency in large-scale cloud deployments.

5. CONCLUSION

According to the study's findings, stateless architectures worked better for large-scale cloud systems that needed high scalability, quick fault recovery, and lower operating overhead, even if both stateless and stateful microservices had clear benefits. They continuously produced easier maintenance procedures, quicker scaling times, and higher throughput under peak loads. Stateful microservices, however, continued to play a vital role in situations requiring real-time data continuity, transactional integrity, or session persistence. System architects could attain the best balanced performance by implementing a hybrid strategy—using stateless services for elasticity and resilience while assigning crucial state management to specialized, fault-tolerant data stores—according to the trade-off analysis, which showed that no single approach was consistently optimal.

REFERENCES

1. O. Andell, *Architectural Implications of Serverless and Function-as-a-Service*. 2020.
2. S. R. Chowdhury, M. A. Salahuddin, N. Limam, and R. Boutaba, "Re-architecting NFV ecosystem with microservices: State of the art and research challenges," *IEEE Network*, vol. 33, no. 3, pp. 168–176, 2019.
3. A. Furda, C. Fidge, O. Zimmermann, W. Kelly, and A. Barros, "Migrating enterprise legacy source code to microservices: On multitenancy, statefulness, and data consistency," *IEEE Software*, vol. 35, no. 3, pp. 63–72, 2017.

4. H. Kang, M. Le, and S. Tao, "Container and microservice driven design for cloud infrastructure devops," in *Proc. 2016 IEEE Int. Conf. on Cloud Engineering (IC2E)*, Apr. 2016, pp. 202–211.
5. S. Thumala, "Building highly resilient architectures in the cloud," *Nanotechnology Perceptions*, vol. 16, no. 2, 2020.
6. A. Pekkala, *Migrating a Web Application to Serverless Architecture*. 2019.
7. A. Singhvi, S. Banerjee, Y. Harchol, A. Akella, M. Peek, and P. Rydin, "Granular computing and network intensive applications: Friends or foes?," in *Proc. 16th ACM Workshop on Hot Topics in Networks*, Nov. 2017, pp. 157–163.
8. P. García-López, M. Sánchez-Artigas, S. Shillaker, P. Pietzuch, D. Breitgand, G. Vernik, ... and A. J. Ferrer, "Servermix: Tradeoffs and challenges of serverless data analytics," *arXiv preprint arXiv:1907.11465*, 2019.
9. G. Sayfan, *Mastering Kubernetes: Level up your container orchestration skills with Kubernetes to build, run, secure, and observe large-scale distributed apps*. Packt Publishing Ltd., 2020.
10. I. S. Ayebo, "Comparative analysis of serverless and traditional cloud architectures," *ResearchGate*, Sep. 2017.
11. G. Sayfan, *Mastering Kubernetes: Master the art of container management by using the power of Kubernetes*. Packt Publishing Ltd., 2018.
12. R. T. J. Bolscher, *Leveraging Serverless Cloud Computing Architectures: Developing a Serverless Architecture Design Framework Based on Best Practices Utilizing the Potential Benefits of Serverless Computing*, M.S. thesis, Univ. of Twente, 2019.
13. D. Fireman, J. Brunet, R. Lopes, D. Quaresma, and T. E. Pereira, "Improving tail latency of stateful cloud services via GC control and load shedding," in *Proc. 2018 IEEE Int. Conf. on Cloud Computing Technology and Science (CloudCom)*, Dec. 2018, pp. 121–128.
14. Y. Al-Dhuraibi, F. Paraiso, N. Djarallah, and P. Merle, "Elasticity in cloud computing: State of the art and research challenges," *IEEE Transactions on Services Computing*, vol. 11, no. 2, pp. 430–447, 2017.
15. A. Sarkar and A. Shah, *Learning AWS: Design, Build, and Deploy Responsive Applications Using AWS Cloud Components*. Packt Publishing Ltd., 2018.