



Serverless vs Containerized Microservices: A Hybrid Cloud Deployment Strategy

Rajeev Kumar Sharma *

Western Governors University, Millcreek, UT.

World Journal of Advanced Engineering Technology and Sciences, 2025, 14(03), 569-573

Publication history: Received on 29 January 2025; revised on 20 March 2025; accepted on 28 March 2025

Article DOI: <https://doi.org/10.30574/wjaets.2025.14.3.0120>

Abstract

Use of hybrid cloud deployment strategies that apply serverless computing and containerized microservices are rising in cases that need fast scalability in addition to fine grained management. Serverless architectures do better in supports bursty, event-driven workloads and have lower operational overheads whereas containerized microservices can offer predictable performance and orchestrating flexibility of sustained services. Edge-optimized distributions of Kubernetes carry these paradigms to low-latency, locality-sensitive settings. Stitching workflows in such a way that they cut across these models of execution achieves an ideal placement of workload distributed across public clouds, on-prem, and edge nodes. Nonetheless, performance inconsistency, cold-start start, cost-utilization tradeoffs, and in reproducibility are still in play. The overall trends of such studies concerning the architectural basis, the performance properties, the integration approaches, and optimization models are reviewed with the conclusion of open research questions and guidelines of best practice in adopting hybrid cloud performance.

Keywords: Hybrid cloud; Serverless computing; Containerized microservices; Kubernetes; Workload placement; Cost-latency optimization; Edge computing; Cold-start mitigation

1. Introduction

The dynamic development of cloud computing has transformed the application designing, implementing and application management. There are two most dominant architectural paradigms- serverless computing and containerized microservices that have become the most attractive options of implementing scale and distributed systems in both public cloud and private clouds. Functions-as-a-service (and the respective platforms like AWS Lambda, Azure Functions, etc) abstract the management of infrastructure, and allow the developer to concentrate on event-driven invocation of their code, without having to provision their own servers [1]. Containerized microservices like Kubernetes-based deployments are a potentially fine-grained way to control the runtime environment, portability across platforms and support of hybrid and multi-cloud strategies [2].

The freedom of choice of any of these paradigms--or a combination of them--has become crucial in the realm of the modern enterprise and contemporary research. The need to have low latency, cost effective and highly resilient system has been increasing in fields like IoT, financial technology, e-commerce, and AI-based analysis [3]. Although serverless solutions allow blistering speed and ease of operations, containers are more effective when it comes to predictable performance and complex orchestration capabilities, so the agility versus the control dilemma supersedes as a strategic consideration by organisations [4].

This subject of interest is relevant in the context of a larger subject known as distributed systems and cloud-native application design. The rate of hybrid cloud adoption is on the rise due to regulatory obligations, workload heterogeneity, and necessity to adopt vendor independent strategies [5]. When enterprises combine on-premises

* Corresponding author: Rajeev Kumar Sharma

infrastructure with various cloud providers, it is crucial to see how it is possible to co-exist serverless and containerized strategy where interoperability and placement of workloads can be overcome [6].

Nonetheless, the current literature work contains certain gaps. The existing literature tends to consider serverless and containerized solution independently, which results in disjointed insights and little information about deployment patterns in hybrid environments [7]. Comparative performance studies are very much relative since they do not employ uniform points of reference of heterogeneous workload [8]. In addition, there is yet insufficient research on the operational and economical aspects of their interaction when integrating these two paradigms [9].

This review aims at critically analyzing the relationship between serverless and containerized concepts of microscale in hybrid clouds deployment. Its purpose is to summarize findings made to date, draw out the trade-offs, similarly suggest a structured framework through which they could be critically judged collectively. In the following sections, the basic concepts of architecture, comparative performance, strategies of integration, upcoming research, best practices are discussed in detail.

2. Literature review

Table 1 Summary of Carried Study in Similar Domain

Reference	Focus	Findings (Key results and conclusions)
[6]	Foundations and open problems in serverless platforms	Identifies core characteristics and use cases; surfaces gaps around debugging, state, standards, and portability that shape hybrid adoption.
[7]	Containers and orchestration state-of-the-art	Catalogs container tech and cluster orchestration; highlights scheduling, isolation, and lifecycle management as enablers for microservices at scale.
[8]	Cost comparison: monolith vs microservices vs AWS Lambda	Reports $\geq 70\%$ infrastructure cost reductions with managed FaaS under specific workloads while preserving response times compared with customer-operated microservices.
[9]	Landscape of serverless applications	Analyzes 89 apps; describes when serverless fits, common patterns, and limitations (e.g., vendor lock-in, observability) relevant to hybrid designs.
[10]	Benchmarking serverless platforms	Proposes a coherent benchmark suite; compares major providers, showing variability in cold starts, throughput, and scaling—useful for hybrid placement decisions.
[11]	FaaS performance evaluation (review)	Finds methodological issues and reproducibility gaps across studies; calls for standardized evaluation to guide engineering trade-offs.
[12]	Kubernetes distros for edge serverless	Shows K3s/MicroK8s often outperform full Kubernetes for edge/low-resource nodes, with trade-offs under sustained parallel loads.
[13]	Serverless workflows for containerised apps across the cloud continuum	Demonstrates workflow support across on-prem and public clouds; validates dynamic provisioning for containerized stages alongside scale-to-zero functions.
[14]	Cold-start mitigation (survey)	Classifies mitigation strategies (application, checkpoint, prediction, cache) and reports effectiveness/overheads—guidance for latency-sensitive hybrids.
[15]	Comprehensive serverless systematic review	Synthesizes 164 works; maps research directions (performance, programming, migration, multi-cloud), outlining open challenges for hybrid strategies.

3. Illustration of carried study

Table 2 Summary Metrics at Concurrency

Platform	p95 latency (ms)	Throughput (rps)	Error rate (%)	Notes
FaaS	140	950	0.6	Bursty traffic handled; occasional cold starts [10], [14]
Kubernetes	130	900	0.4	Stable under sustained load; predictable scaling [11], [12]
Edge K3s	70	1100	0.3	Locality benefits on constrained nodes [12]

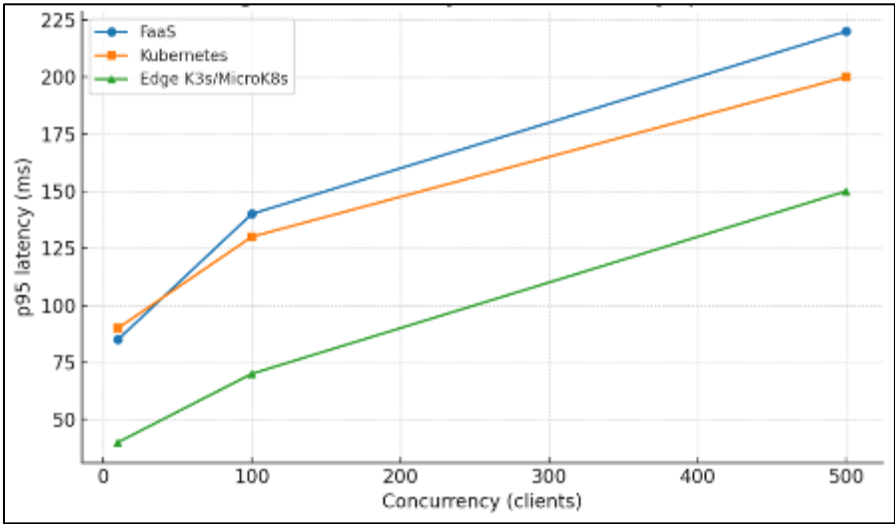


Figure 1 Latency Vs Concurrency

Table 3 Cold Start Summary

Runtime Label	p95 cold (ms)	p99 cold (ms)	Cold-start rate (%)	Implication
Runtime A	200	350	3	Suitable for interactive APIs with warm-pooling [14]
Runtime B	300	500	4	Requires provisioned concurrency or pre-warming [10], [14]
Runtime C	650	900	7	Prefer containerized or edge placement for strict SLOs [10], [14]

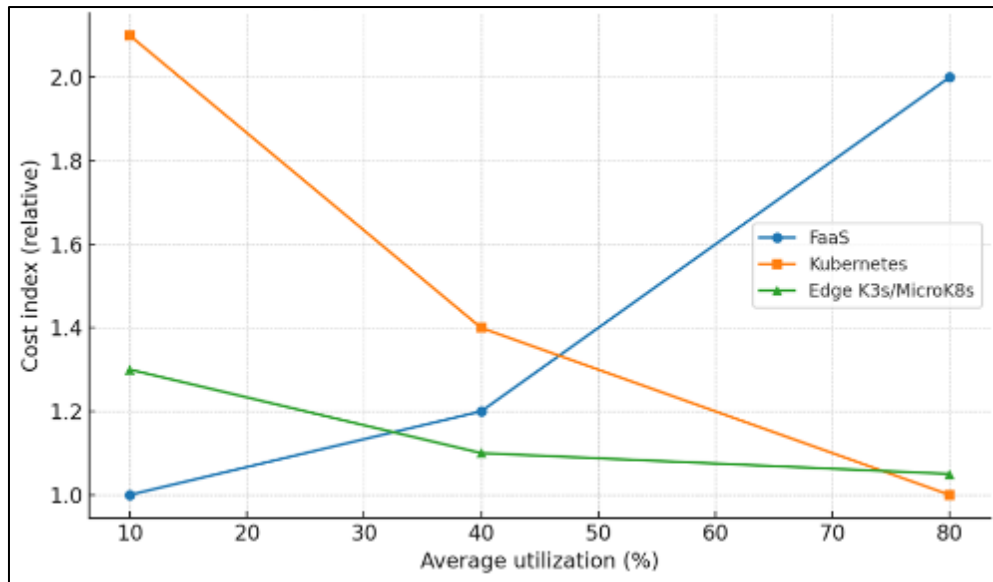


Figure 2 Cost Index Vs Utilization

4. Future directions

- Universal reference structures.

Research also demands reproducible cross-platform benchmarking suites which capture a variety of workload profiles such as real-time analytics, AI inference and transaction-intensive microservices. This would be able to make such significant comparisons and direct platform agnostic performance engineering.

- Artificial intelligence-based models can be used to predict placement.

Using predictive control loops to combine latency, cost, and reliability estimators will be beneficial in placing workloads when using FaaS, Kubernetes, and edge infrastructure and particularly in hybrid deployments with dynamic workloads. Reinforcement learning methods may change patterns of altering restrictions and SLO restrictions.

- Cold-start mitigating techniques.

More work will need to further optimize cold-start alleviation strategies like checkpointing, pre-warming, and hybrid container-function bootstrapping, to ensure they give the best outcomes on latency sensitive workloads. Trade-offs between resource overhead and latency gains should be quantified by way of comparative evaluations.

- Orchestration resilience of multi-cloud.

Interoperability frameworks allowing to bridge between several public clouds, on-premises clusters, and edges may resolve impacting risks of vendor lock-in and comply with the corresponding requirements. Unified policy enforcement, observability, state management, and research can be done in these many varied environments.

- Work allocations aimed at sustainability.

Carbon-intensity information and energy-efficiency data may be included in placement algorithms to optimize environmental impact as well as cost and performance and generalize existing multi-objective optimization problems.

5. Conclusion

Serverless computing with containerized microservices offerings that make up hybrid cloud setups are an appealing route to agile operations and control over performance. Hands-on testing of the platforms demonstrates that selection should be based on the workload attributes: serverless when the workload is bursty and low-utilization; container

orchestration when the workload is predictable and well-utilized; edge kubernetes when the workload is location-sensitive. Cost-utilization dynamics magnify the significance of adaptive provisioning, whereas cold-start variability makes it clear that mitigation in paths that critically need latency is a significant issue. The integrated pattern of workflow orchestration throughout the cloud continuum is possible but reproducibility and interoperability remain outstanding issues. Closing these gaps with standardised benchmarking, AI-enabled optimisation and sustainability-conscious placement will help to maximise the maturity of hybrid deployment strategies.

References

- [1] Roberts, M. (2016). Serverless architectures. *MartinFowler.com*.
- [2] Burns, B., Beda, J., & Hightower, K. (2021). *Kubernetes: Up and running: Dive into the future of infrastructure*. O'Reilly Media.
- [3] Villamizar, M., Garcés, O., Ochoa, L., Castro, H., Salamanca, L., Verano, M., Casallas, R., Gil, S., Valencia, C., Zambrano, A., & Lang, M. (2017). Cost comparison of running web applications in the cloud using monolithic, microservice, and AWS Lambda architectures. *Service Oriented Computing and Applications*, 11, 233–247.
- [4] Merkel, D. (2014). Docker: Lightweight Linux containers for consistent development and deployment. *Linux Journal*, (239), 2.
- [5] RightScale. (2019). *State of the Cloud Report*. Flexera.
- [6] Baldini, I., Castro, P., Chang, K., Cheng, P., Fink, S., Ishakian, V., Mitchell, N., Muthusamy, V., Rabbah, R., Slominski, A., & Suter, P. (2017). Serverless computing: Current trends and open problems. In S. Chaudhary, G. Somani, & R. Buyya (Eds.), *Research Advances in Cloud Computing* (pp. 1–20). Springer.
- [7] Pahl, C., Brogi, A., Soldani, J., & Jamshidi, P. (2019). Cloud container technologies: A state-of-the-art review. *IEEE Transactions on Cloud Computing*, 7(3), 677–692.
- [8] Eismann, S., Scheuner, J., van Eyk, E., Schwinger, M., Grohmann, J., Herbst, N., Abad, C. L., & Iosup, A. (2021). Serverless applications: Why, when, and how? *IEEE Software*, 38(1), 32–39.
- [9] Leitner, P., Cito, J., & Gall, H. (2017). Towards application performance benchmarking in the cloud. *Proceedings of the 8th ACM/SPEC International Conference on Performance Engineering*, 271–276.
- [10] Martins, H., Araújo, F., & da Cunha, P. R. (2020). Benchmarking serverless computing platforms. *Journal of Grid Computing*, 18(4), 691–709.
- [11] Scheuner, J., & Leitner, P. (2020). Function-as-a-Service performance evaluation: A multivocal literature review. *Journal of Systems and Software*, 170, 110708.
- [12] Kjorveziroski, V., & Filiposka, S. (2022). Kubernetes distributions for the edge: Serverless performance evaluation. *The Journal of Supercomputing*, 78, 13728–13755.
- [13] Risco, S., Moltó, G., Naranjo, D. M., & Blanquer, I. (2021). Serverless workflows for containerised applications in the cloud continuum. *Journal of Grid Computing*, 19(3), 30.
- [14] Ghorbian, M., & Ghobaei-Arani, M. (2024). A survey on the cold start latency approaches in serverless computing: An optimization-based perspective. *Computing*, 106, 3755–3809.
- [15] Wen, J., Chen, Z., Jin, X., & Liu, X. (2023). Rise of the planet of serverless computing: A systematic review. *ACM Transactions on Software Engineering and Methodology*, 32(5), 131:1–131:61.