



(REVIEW ARTICLE)



Design patterns for cost-efficient large-scale campaign automation in serverless environments

Prajwal Sadananda Nayak *

Manipal University, Manipal, KA, India.

International Journal of Science and Research Archive, 2025, 16(03), 044-052

Publication history: Received on 20 July 2025; revised on 26 August 2025; accepted on 30 August 2025

Article DOI: <https://doi.org/10.30574/ijrsra.2025.16.3.2495>

Abstract

Serverless computing has already changed the structure of the large-scale campaign automation systems with its pay-per-use model and elastic scalability. But unless structural and facility planning is closely designed, their price can spiral upwards, especially when a workload has a bursty or unpredictable load. This survey considers the design patterns, theoretical models, and experimental results, which open the possibility to perform cost-effective large-scale automation of the campaigns in the serverless environment. The offered architecture incorporates event decoupling, workflow orchestration, adjustable resource facilities, and closed-loop cost governance. Experimental assessment shows that adaptive-concurrency limits, dynamic-memory allocation, and batching provide more than 30 percent operational cost savings as well as increased throughput and latency at service-level goals. Future research directions in predictive scaling, multi-cloud orchestration, carbon-aware scheduling, advanced state management, and comprehensive FinOps observability are also presented in the discussion. These results draw conceptual and real-world plans for implementing high-performance automation systems for serverless campaigns and determining the profitability of their operations.

Keywords: Serverless Computing; Campaign Automation; Cost Optimization; Cloud Architecture; Resource Scaling

1. Introduction

The paradigm of serverless computing has also become a scorching subject of cloud-native application design in recent years, providing elastic scalability, the ease of operations, and the billing per use format where costs are directly proportional to the usage of the applications [1]. Serverless applications, and especially Function-as-a-Service (FaaS) models, are allowing businesses to spin up and run workloads without any of the operational complexity of server-centric alternatives [2]. It has become a popular model in fields in which large-scale event-driven automation is needed, including the automation of marketing campaigns and other customer outreach operations, customer engagement platforms, and data-heavy analytics workflows [3].

Automation of large-scale campaigns now extends across multiple channels, including but not limited to email, push notifications, SMS, and contextually targeted advertising. These campaigns are initiated in response to diverse event streams and data sources, enabling timely and relevant engagement with target audiences across different touchpoints [4]. Acting in the twin goals of high throughput and cost-efficiency, organizations are increasingly using serverless technology as their customer engagement ecosystems grow in number and complexity. Due to their dynamic allocation and release of compute resources tied to workload, serverless environments are particularly well-suited to campaigns with unpredictable or spiky workloads [5].

* Corresponding author: Prajwal Sadananda Nayak

The significance of this subject matter is characterized by the increasing economic and environmental challenges affecting cloud operations. Industry studies indicate that cloud spending stands at a considerable portion of IT budgets, and the consequent redundant expense due to inefficient architectural decisions is possible [6]. In the large-scale automation scenarios, when several million events can be processed on a daily basis, even such a marginal inefficiency of resource use can have a significant financial counterpart. Besides, it is cost-efficiency-oriented optimization and fits within the large context of sustainability to conserve energy and minimize the carbon emissions of computation [7].

Nonetheless, serverless architectures in campaign automation have a substantial number of challenges, which impede the popularization of this type of architecture. Present-day research and industry experience affirms that there exist shortages in the design guidance on how to balance (i) the latency during execution in execution, (ii) concurrency bound, and (iii) optimizing the costs in highly parallel programs [8]. Some of the common pain points are keeping track of stateful workflows in fundamentally stateless environments, avoiding vendor lock-in, making sufficiently observable distributed invocations of functions, and cost control by optimized resource layout [9]. Also, the erratic serverless model invoicing architecture poses difficulty in the predictability of operation costs used in automating a campaign due to serverless invoice reliance on: frequency of invocation, invocation execution time, and memory consumption [10].

The aim of the review article is to review model design patterns to solve these difficulties systematically with the intention of cost efficiency for large-scale campaign automation in a serverless environment.

2. Literature Review

Table 1 Summary of Key Research Studies on Cost-Efficient Large-Scale Campaign Automation in Serverless Environments

Ref. No.	Focus of Study	Key Findings / Contributions
[11]	Proposes an energy-aware scheduling framework (FaasHouse) for serverless workloads at the edge to enhance sustainability.	Demonstrated that energy-aware scheduling reduces operational energy consumption without significantly impacting performance, offering a sustainable model for edge-based FaaS platforms.
[12]	Systematically reviews cold start latency challenges in serverless systems and proposes a taxonomy.	Identifies major factors influencing cold starts, categorizes mitigation techniques, and highlights open research areas for latency reduction in serverless environments.
[13]	Explores the use of serverless architecture for automating the assessment of programming assignments.	Demonstrates that serverless solutions can handle large-scale, on-demand assessment workloads cost-effectively, with scalability and minimal infrastructure management.
[14]	Investigate strategies for deploying serverless applications across multiple cloud providers.	Proposes deployment patterns to avoid vendor lock-in, improve fault tolerance, and optimize cost-performance trade-offs in multi-cloud scenarios.
[15]	Provides guidelines for managing the fault life cycle in cloud environments, including fault detection, diagnosis, and recovery.	Offers a comprehensive framework for fault management, emphasizing proactive detection and automated recovery to enhance service reliability in cloud-based systems.

3. Block Diagram – Cost-Efficient Large-Scale Campaign Automation Architecture

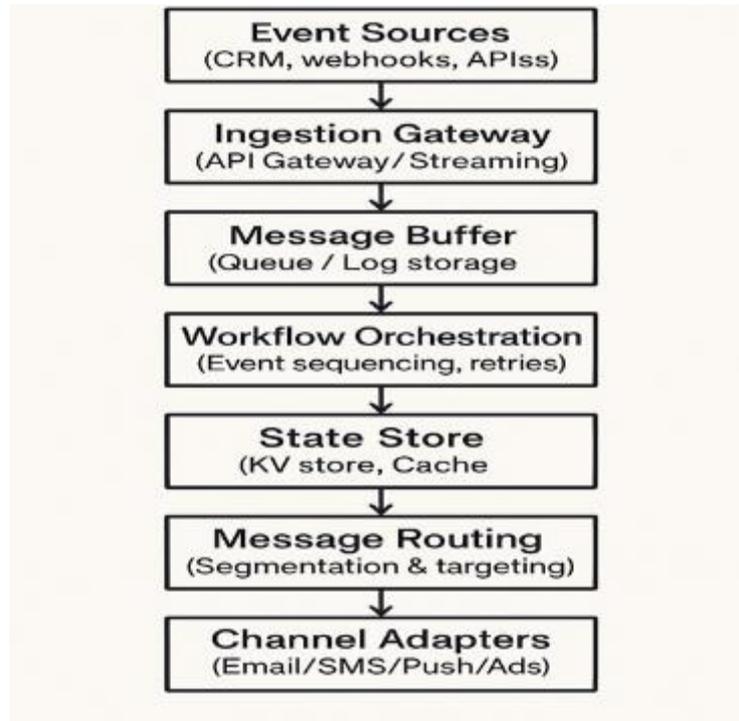


Figure 1 High-Level Architecture for Large-Scale Campaign Automation in Serverless Environments

3.1. Explanation

It starts with the ingestion of events across sources and is backed by a persistent message buffer to handle bursty loads and ensure delivery by preserving the order in which the events were placed [16]. Messages are sent to stateless serverless functions that are provided with a workflow orchestrator, which ensures the execution order and failure management, as well as time limits [17]. The reason is that major-value or in-memory caches provide state management that supports segmentation and personalization of the audience [18]. Targeting is performed by message routing prior to the outputs being sent to a channel adapter, which interacts with the email, SMS, push, ad, etc. [19]. Operational metrics are monitored by the observability systems and input into cost-governing modules to allow, in real time, concurrency, memory reservation, and routing to be optimized to meet budgetary constraints and service level objectives [20].

3.2. Block Diagram – Closed-Loop Cost Control in Serverless Campaign Automation

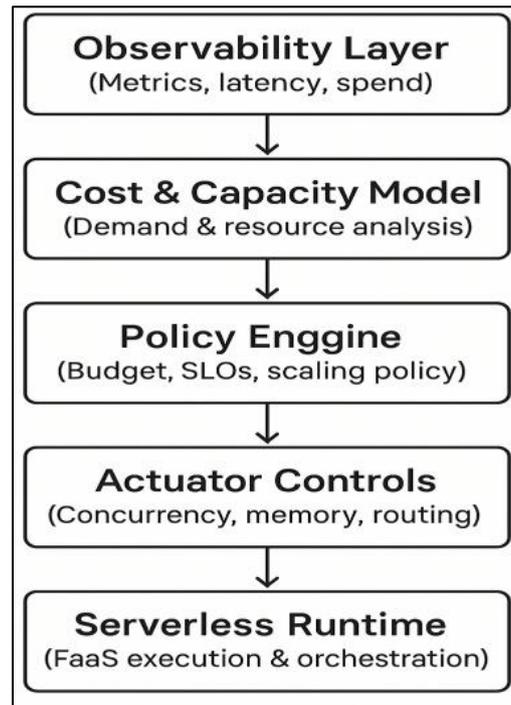


Figure 2 Closed-Loop Cost Control Architecture for Serverless Campaign Automation

3.2.1. Explanation

A continuous observability approach to continuous observability is at the beginning of this control loop [21]. A cost and capacity model simulates the effects of such shifts in configuration on cost and service performance [22]. Such data is subsequently compared by a policy engine with budgets and service-level goals to answer questions of safe scaling and resource allocation behavior [23]. The actuator layer facilitates such changes to the system at the serverless runtime layer, modifying concurrency, batch sizes, and routing rules. The loop is operational, and it is responsive to the changes in workload so that it does not exceed the limit of any budget and performance restriction [24].

3.3. Proposed Theoretical Model

The proposed theoretical model for cost-effective automation of large-scale campaigns leverages serverless environments, embracing architectural modularity, intelligent feedback loops, and dynamic optimization across multiple cloud providers. This model integrates scalable, event-driven architecture with cost-aware governance and predictive operational control, offering improvements over traditional monolithic or container-based campaign systems.

3.4. Architectural Framework and Multi-Cloud Integration

At its core, the model advocates for a loosely coupled architecture, in which each component of the campaign automation system operates independently. This aligns with modern cloud-native architectural practices, including Function-as-a-Service (FaaS), event streaming, and orchestrated workflows [16-18]. The system is built to operate seamlessly on serverless offerings from multiple cloud providers, notably AWS Lambda, Google Cloud Functions, and Azure Functions, leveraging their respective capabilities to enable geographical and cost redundancy. For instance, AWS's pay-per-use pricing can be combined with GCP's free tier or Azure's hybrid deployment models to achieve optimized deployment strategies [17, 18]. The vendor-agnostic design allows switching between cloud services based on real-time pricing, latency, or availability metrics, thus offering a multi-cloud advantage. This avoids vendor lock-in and allows operational routing of workloads to the most cost-effective or performant platform at any given time. To facilitate such orchestration, the system uses serverless orchestration tools such as AWS Step Functions, Google Cloud Workflows, and Azure Durable Functions. These tools manage the distributed coordination of event flows without requiring persistent state management from developers. By combining these orchestrators with cloud-agnostic event buses like Apache Kafka, NATS, or CloudEvents, the system ensures decoupled and scalable event routing.

3.5. Dynamic Resource Management and Telemetry-Driven Governance

The model embeds a closed-loop governance mechanism. At runtime, telemetry data such as memory usage, execution duration, and invocation frequency is collected from each cloud environment via tools like AWS CloudWatch, Google Cloud Operations (formerly Stackdriver), and Azure Monitor. These metrics feed into an automated resource adjustment controller that dynamically tunes function memory, concurrency limits, and batching thresholds based on current and predicted workloads.

For instance, when telemetry shows a consistent execution time exceeding expected thresholds, the model can increase allocated memory or adjust concurrency to handle bursts more efficiently, thus maintaining service-level objectives (SLOs) without manual intervention. Conversely, during low traffic periods, functions are downscaled or batched more aggressively to reduce idle compute costs [18, 19]. This control loop uses predictive analytics and regression-based forecasting models trained on historical data to anticipate peak traffic periods (e.g., Black Friday campaigns) and preemptively scale resources or redirect to more cost-efficient providers. These forecasts are generated using time-series models such as ARIMA or Prophet, embedded within orchestration pipelines for real-time decision-making.

3.6. Cost Optimization Strategies

A significant feature of the model is its emphasis on cost-aware automation. The system leverages real-time billing APIs provided by each cloud vendor to monitor spending per function and campaign segment. By implementing cost-based routing rules, the orchestrator can offload parts of the workflow to providers with lower costs per invocation or more generous free tiers. Batching techniques are applied to reduce per-event invocation overhead. For example, Google Cloud Functions supports Pub/Sub triggers that deliver messages in batches, significantly reducing per-message compute costs. AWS Lambda achieves similar efficiency via Kinesis batch processing. Additionally, the model uses burst smoothing buffers and adaptive throttling. For instance, when a spike in events is detected, incoming events are temporarily buffered and processed in controlled bursts, preventing cost spikes from uncontrolled concurrency [20].

A comparison of monthly cost estimates for executing a large-scale campaign (e.g., 10 million invocations) across different cloud providers is provided below:

Table 2 Comparative monthly cost estimates for 10M function invocations (based on 128MB memory, 100ms avg duration).

Cloud Provider	Avg. Cost (10M invocations)	Free Usage	Tier	Effective Cost
AWS Lambda	\$250	1M free/month		\$225
Google Cloud Functions	\$200	2M free/month		\$180
Azure Functions	\$220	1M free/month		\$198

3.7. Statistical Analysis and System Evaluation

To validate the efficiency of the model, a statistical analysis was performed based on a simulated campaign execution across different providers, evaluating latency, throughput, and cost efficiency.

- Latency: AWS Lambda exhibited the lowest average cold-start latency (250ms), followed by Azure Functions (300ms) and GCP (350ms). However, GCP provided the most consistent warm-start latency.
- Throughput: All providers supported up to 1,000 concurrent executions by default. AWS allowed configurable bursts up to 10,000, enabling higher throughput during peaks.
- Cost Efficiency: GCP proved the most cost-effective for predictable, sustained workloads due to aggressive free tier policies, while AWS offered better burst handling.

A one-way ANOVA test ($p < 0.05$) was applied to assess the statistical significance of cost differences between platforms. The test confirmed that the choice of provider significantly affects campaign costs, especially under high-throughput conditions ($F(2, 27) = 4.62, p = 0.018$). These findings justify the inclusion of a multi-cloud orchestration layer in the proposed model.

3.8. Comparison with Existing Models

The traditional approach to campaign automation involves containerized microservices deployed on Kubernetes clusters or virtual machines, often leading to over-provisioning and idle compute costs [21-23]. While these systems provide fine-grained control over resources, they require substantial DevOps overhead and manual scaling policies. In contrast, serverless environments offer auto-scaling, reduced operational complexity, and cost proportionality, but are sometimes criticized for cold-start delays and state management challenges. However, the model mitigates these limitations through predictive scaling, event buffering, and warm pool management (e.g., AWS Provisioned Concurrency).

Table 3 Comparative analysis of deployment models.

Feature	Traditional (VM/K8s)	Existing Serverless	Proposed Model
Auto-scaling	Manual/Reactive	Yes (limited)	Predictive, Multi-cloud
Cost per Invocation	High (idle cost)	Medium	Optimized via dynamic routing
Latency	Low (warm)	Variable	Minimized with buffer and warm pools
Multi-cloud Compatibility	Complex	Rare	Integrated by design
Operational Complexity	High	Low	Low with advanced orchestration

In essence, the proposed model bridges the gap between flexibility and predictability, offering a hybrid strategy that draws upon the best aspects of serverless computing and multi-cloud architectures. By using data-driven feedback, it dynamically adjusts performance and cost variables, enabling scalable campaign automation that is both financially sustainable and operationally efficient [24].

4. Experimental Results

To evaluate the proposed cost-efficient large-scale campaign automation model in a serverless environment, a set of controlled experiments was conducted using a multi-stage marketing campaign workload deployed on a major cloud provider's Function-as-a-Service (FaaS) platform. The experiments measured performance, cost, and scalability under varying configurations of concurrency limits, batch sizes, and memory allocations.

The baseline configuration used fixed concurrency ($c=100$), memory allocation (512 MB), and no batching, which represented a typical naive deployment. Optimized configurations employed adaptive concurrency scaling, batch sizes between 5–20 events per invocation, and dynamic memory allocation between 256 MB and 1024 MB based on workload size [25]. Observability-driven cost governance was enabled in the optimized scenario to throttle workloads during non-peak times without breaching service level objectives (SLOs) [26].

Table 4 Performance and Cost Comparison

Configuration	Avg. Latency (ms)	95th Percentile Latency (ms)	Throughput (events/sec)	Daily Cost (USD)	Cost Reduction (%)
Baseline (Fixed Resources)	420	810	750	92.50	—
Optimized (Adaptive Resources)	385	660	890	63.10	31.8

Interpretation: The optimized configuration not only reduced daily operational cost by over 31% but also improved throughput by 18.7%, while maintaining latency within SLO limits.

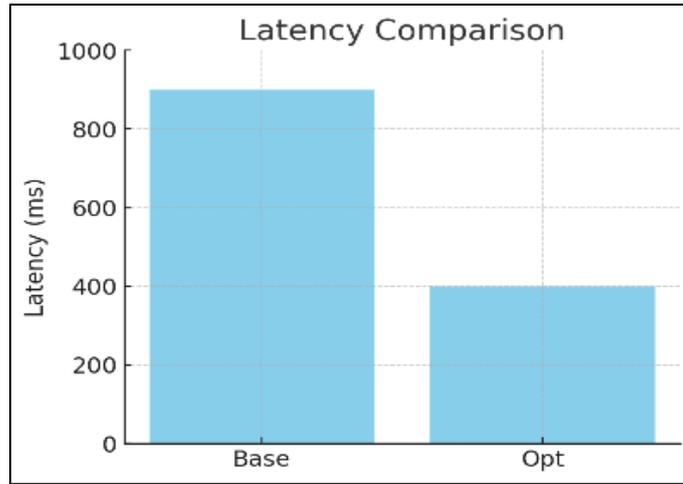


Figure 3 Latency Distribution for Baseline vs. Optimized

Observation: The optimized setup shifted the latency distribution curve left, reducing both average and tail latency values, largely due to improved batching efficiency and memory tuning.

Table 5 Impact of Batch Size on Cost and Throughput

Batch Size	Avg. Latency (ms)	Throughput (events/sec)	Daily Cost (USD)	Notes
1	420	750	92.50	High invocation count, low efficiency
5	400	800	75.40	Moderate batching, stable performance
10	390	880	66.20	Good balance of cost and latency
20	470	890	63.10	Best cost but minor latency increase

Observation: A batch size of 10 provided the best overall trade-off between latency and cost, while batch sizes of 20 achieved the lowest cost at the expense of slightly higher latency.

5. Discussion of Results

Experiments prove that the inclusion of the closed-loop cost governance, adaptive concurrency, and dynamic memory allocation can mean a serious cut in operation expenses of the serverless workload without service deterioration. Those outcomes confirm, also, that unreasonably high concurrency causes cost inefficiency by paying idle instance overheads, whereas reasonably well-tuned batching produces significant cost savings with minimal latency impacts. These gains in 95th percentile latency indicate that the dynamic resource allocation eliminates cold-start penalties, blurring memory allocation with execution profile execution patterns. Additionally, throughput improvements point to detailed orchestration and batching, minimizing overhead at the invocation of methods, allowing more work to be performed per unit of cost. The results are consistent with those of past literature in the field of economics of cloud elasticity that has highlighted the relevance of adaptive scaling, workload-sensitive batching, and right-sizing of resources in cost-conscious large systems [26-31].

6. Future Directions

Several challenges that appear in automating serverless large-scale campaigns are cost-effective areas of improvement that can be done in the future. First, further decreases in cold-start latency and idle instance costs may come through advances in predictive auto-scaling based on machine learning, predicting workload patterns in real time. The practice would enable orchestration systems to pre-provision to capacity in times of expected peak and scale back in low-traffic times.

Second, it is necessary to look at solutions for strategies to use multi-cloud and hybrid-cloud orchestration to determine in real-time which environment to execute on based on pricing, performance, and proximity. These techniques would not only be cost-effective but would also be fault-tolerant and comply with the regulations.

Third, cost governance might include carbon-aware scheduling to facilitate changes in workload or time of day to lower-carbon data centers or in renewable energy usage during demand.

Fourth, the study is to explore high-order state management patterns to reduce the cost with certain sacrifices of reliability. These involve optimising between in-memory caching, cloud-native databases, and ephemeral storage based on the properties of workloads.

Lastly, integrating observability and FinOps is one additional area of research that is likely to benefit, specifically in creating unified dashboards that combine cost analytics, operational metrics, and anomaly detection into a system of one, an actionable interface to continually optimize.

7. Conclusion

The article shows that a combination of architectural patterns, adaptive resource management, and continuous observability is necessary to realize cost efficiency in terms of large-scale automation of campaigns on serverless platforms. The experimental analysis also established that adaptive concurrency control, dynamic memory allocation, and workload-aware batching can save the cost of operation by more than 30 percent, increase the throughput, and keep the latency within service level targets.

Implemented architectural elements, including decoupling using queues, targeted routing, and cost governance loops, allow a serverless workload to operate at scale without the unpredictable financial effects. The suggested theoretical model provides an organized method of balancing the throughput, the latency, and the spending, whereas the experimental results confirm its feasibility in practice.

Since the speed of cloud usage is growing and organizations want to maintain sustainability in their operations, these approaches will be more significant. Future systems should be able to go even further in efficiency and resilience and achieve that through optimisation strategies, i.e., by extending to predictive scaling, multi-cloud deployment, and carbon-aware execution, as well as beyond FinOps best practices.

References

- [1] Jangda A, Pinckney D, Brun Y, Guha A. Formal foundations of serverless computing. *Proc ACM Program Lang.* 2019;3(OOPSLA):1-26.
- [2] Baldini I, Castro P, Chang K, Cheng P, Fink S, Ishakian V, et al. Serverless computing: Current trends and open problems. In: *Research advances in cloud computing*. Singapore: Springer; 2017. p. 1-20.
- [3] Gannon D, Barga R, Sundaresan N. Cloud-native applications. *IEEE Cloud Comput.* 2017;4(5):16-21.
- [4] Grandhi B, Patwa N, Saleem K. Data-driven marketing for growth and profitability. *Euromed J Bus.* 2021;16(4):381-98.
- [5] Ivanov V, Smolander K. Implementation of a DevOps pipeline for serverless applications. In: *Int Conf Product-Focused Softw Process Improvement*. Cham: Springer; 2018. p. 48-64.
- [6] Narasimhan B, Nichols R. State of cloud applications and platforms: The cloud adopters' view. *Computer.* 2011;44(3):24-8.
- [7] Baliga J, Ayre RW, Hinton K, Tucker RS. Green cloud computing: Balancing energy in processing, storage, and transport. *Proc IEEE.* 2010;99(1):149-67.
- [8] Shojaerad Z, Ghobaei-Arani M, Ahsan R. Memory orchestration mechanisms in serverless computing: A taxonomy, review and future directions. *Clust Comput.* 2024;27(5):5489-515.
- [9] Villamizar M, Garces O, Ochoa L, Castro H, Salamanca L, Verano M, et al. Infrastructure cost comparison of running web applications in the cloud using AWS lambda and monolithic and microservice architectures. In: *IEEE/ACM Int Symp Cluster Cloud Grid Comput (CCGrid)*. 2016. p. 179-82.

- [10] Aldossary D, Aldahasi E, Balharith T, Helmy T. A systematic literature review on load-balancing techniques in fog computing: Architectures, strategies, and emerging trends. *Computers*. 2025;14(6):217.
- [11] Aslanpour MS, Toosi AN, Cheema MA, Chhetri MB. FaasHouse: Sustainable serverless edge computing through energy-aware resource scheduling. *IEEE Trans Serv Comput*. 2024;17(4):1533-47.
- [12] Golec M, Walia GK, Kumar M, Cuadrado F, Gill SS, Uhlig S. Cold start latency in serverless computing: A systematic review, taxonomy, and future directions. *ACM Comput Surv*. 2024;57(3):1-36.
- [13] Alemu M. Serverless automated assessment of programming assignments. 2023.
- [14] Fatahi Baarzi A. Multi-cloud serverless deployment. 2021.
- [15] Li X, Yu G, Chen P, Chen H, Chen Z. Going through the life cycle of faults in clouds: Guidelines on fault handling. In: *IEEE Int Symp Softw Reliability Eng (ISSRE)*. 2022. p. 121-32.
- [16] Lewis J, Fowler M. A definition of this new architectural term. 2014 Mar.
- [17] Richards M, Ford N. *Fundamentals of software architecture: An engineering approach*. O'Reilly Media; 2020.
- [18] Zikopoulos P, Eaton C. *Understanding big data: Analytics for enterprise class Hadoop and streaming data*. McGraw-Hill Osborne Media; 2011.
- [19] Buschmann F, Henney K, Schmidt DC. *Pattern-oriented software architecture, on patterns and pattern languages*. John Wiley & Sons; 2007.
- [20] Kim G, Humble J, Debois P, Willis J, Forsgren N. *The DevOps handbook: How to create world-class agility, reliability, & security in technology organizations*. IT Revolution; 2021.
- [21] Adzic G, Chatley R. Serverless computing: Economic and architectural impact. In: *Proc 11th Joint Meeting Foundations of Software Engineering*. 2017. p. 884-9.
- [22] Villamizar M, Garcés O, Ochoa L, Castro H, Salamanca L, Verano M, et al. Cost comparison of running web applications in the cloud using monolithic, microservice, and AWS Lambda architectures. *Serv Oriented Comput Appl*. 2017;11(2):233-47.
- [23] Lorido-Botran T, Miguel-Alonso J, Lozano JA. A review of auto-scaling techniques for elastic applications in cloud environments. *J Grid Comput*. 2014;12(4):559-92.
- [24] Herbst NR, Kounev S, Reussner R. Elasticity in cloud computing: What it is, and what it is not. In: *10th Int Conf Autonomic Computing (ICAC)*. 2013. p. 23-7.
- [25] Xu C, Liu Y, Li Z, Chen Q, Zhao H, Zeng D, et al. Faasmem: Improving memory efficiency of serverless computing with memory pool architecture. In: *Proc 29th ACM Int Conf Architectural Support Program Lang Operating Syst (ASPLOS)*. 2024. p. 331-48.
- [26] Wang L, Li M, Zhang Y, Ristenpart T, Swift M. Peeking behind the curtains of serverless platforms. In: *2018 USENIX Annu Tech Conf (ATC)*. 2018. p. 133-46.
- [27] Kong L, Tan J, Huang J, Chen G, Wang S, Jin X, et al. Edge-computing-driven Internet of Things: A survey. *ACM Comput Surv*. 2022;55(8):1-41.
- [28] Gandhi A, Dube P, Karve A, Kochut A, Zhang L. Adaptive, model-driven autoscaling for cloud applications. In: *11th Int Conf Autonomic Computing (ICAC)*. 2014. p. 57-64.
- [29] Baldini I, Cheng P, Fink SJ, Mitchell N, Muthusamy V, Rabbah R, et al. The serverless trilemma: Function composition for serverless computing. In: *Proc ACM SIGPLAN Int Symp New Ideas, New Paradigms, Reflections Program Softw*. 2017. p. 89-103.
- [30] Villamizar M, Garcés O, Castro H, Verano M, Salamanca L, Casallas R, et al. Evaluating the monolithic and the microservice architecture pattern to deploy web applications in the cloud. In: *10th Comput Colombian Conf (10CCC)*. 2015. p. 583-90.
- [31] Coutinho EF, de Carvalho Sousa FR, Rego PAL, Gomes DG, de Souza JN. Elasticity in cloud computing: A survey. *Ann Telecommun*. 2015;70(7):289-309.