# International Journal of Research Publication and Reviews

# Serverless Computing: Challenges and Insights

*Kannaji Jyothi Venkata Nagendra [1], A. Sudhakar [2]*

**GMR Institute of Technology**

**ABSTRACT:**

Serverless computing has fast revolutionized cloud infrastructure and its lightweight, cost-effective model allows developers to write code based on execution rather than taking care of the server in which it will be executed. It even utilizes Function as a Service and Backend as a Service to provide automatic scaling plus dynamic resource allocation. Even with rising popularity, some of the key issues that serverless computing poses come in the form of reducing start-up latency, enhancing scheduling precision, and handling complex networking communications. Furthermore, it maximizes resource efficiency by not allowing any resources to go idle, thus saving considerable costs both for the developers and the providers.

Today, serverless computing is applied in a wide area, including machine learning, real-time data analytics, and distributed computing. As a core component of future cloud architectures, it will be most pertinent. Emerging fields, including serverless edge computing, will gain more attention to address the demands by developing latency- and fault-tolerance-oriented solutions. As it expands, it presents a fine opportunity for innovation in cloud computing and opens up the way to next-generation scalable and efficient platforms.

**Keywords:** Serverless computing, cloud infrastructure, Function as a Service (FaaS), Backend as a Service (BaaS), scalability.

## 1. Introduction:

One of the latest and most innovative trends in relation to cloud-based application deployment and management is serverless computing, advancing at a tremendous rate in the clouds. In contrast with old paradigms that required developers to provision, configure, and manage servers, serverless computing makes it possible for applications to run on managed infrastructure and eliminates these issues. This makes it possible for cloud providers to perform resource provisioning, scaling, or maintenance in order for the developers to shift their focus solely on programming. This in the end allows for higher efficiency and lower costs.

The advancement of serverless computing is most useful in the domain of highly dynamic and event-driven applications where the workload can be highly unstable. Serverless architectures, such as AWS Lambda or Google Cloud Functions, solve the problem posed by unstable workloads by charging them for the time of a computer that has been utilized. However, the model shows several unique disadvantages such as cold-start delays, execution time limits, and the problem of active resource utilization.

Recent innovations both in research and industry have pushed serverless computing approaches even further in recent years. For instance, its potentiality of changing the landscape in NFV has been discussed in [1], where greater flexibility and performance can be offered by serverless approaches. Lightweight virtualization technologies such as Firecracker [2] have made it possible to enable the efficiency of isolating serverless workloads for security at minimal overhead. Some statistical learning methods, like COSE [3], have been proposed to optimize serverless configurations as a symptom of the increased interest in intelligent adaptive solutions.

High-performance serverless frameworks, like SAND [4], alleviate some of the significant performance bottlenecks, allowing applications to execute at low latency and high throughput. The Cherry Pick [5] shows dynamic identification of the best cloud configuration for resource-intensive tasks. A new role of adaptive systems in relation to increasing the efficiency of serverless has then emerged.

Companies like Amazon Web Services introduced highly robust and scalable platforms that have inspired the adoption of serverless. AWS Wavelength [6] supports serverless computing with 5G edge infrastructure to leverage ultra-low-latency applications, services such as DynamoDB [7] and Amazon S3 [8] are foundational storage and database solutions tailored to the requirements of serverless technologies.

The present paper explores the current advancement, challenges, and future directions in serverless computing. Composed by synthesizing insights from both academic research and industrial implementations to attempt at delivering a holistic view of the transformative potential of serverless computing, the key areas of address ability in this regard include its applications, performance optimization strategies, and innovations that shape the future of cloud-based computing.

## 2. Analysis:

### 2.1 Methodology

### 2.1.1 Brief Introduction to Serverless computing

Serverless In serverless computing, the new direction in cloud architecture allows developers to build and deploy applications without such needs to manage servers. Developers in a classical cloud environment have to configure and maintain virtual servers; often, they have to over-provision to accommodate peak loads, increasing operational costs. Function abstracts away all these activities. Instead, developers have come up with functions which the cloud provider executes swiftly concerning any activity-taking place-wisely enough, nothing more or less than FaaS, Functions as a Service. It charges the developer only for the usage time of the computation while provisioning and scaling the resources automatically at the instance of their requirement. It leads to a cost-effective model mainly for applications characterized by having intermittent or erratic workloads where it is inefficient to pay for constant server up-time.

The two main building blocks of a serverless architecture are FaaS, an acronym for Function-as-a-Service, and BaaS, or Backend as a Service. FaaS enables one to execute various parts of code in a stateless fashion, making it highly adaptable to specific event-driven architectures, for example, such as an API response, file uploads, and triggers from a database. BaaS keeps the core backend services up and running with databases, authentication, storage, and real-time data processing to your application but does not require direct control over underlying infrastructure. This design simplifies development: The output is rapid building, testing, and deployment of applications that include less complexity. Serverless models support microservices architecture whereby the applications are split into different modules and each module performs one role specifically. Modularity facilitates quick cycle time, better scalability, and isolation in testing, which helps a more agile software development process.

### 2.2 Serverless Computing: A Comprehensive Survey

Serverless computing is, therefore, a disruptive cloud model that not only lets developers code but allows them to forget that there is infrastructure to manage in the first place. A comprehensive survey of the subject undertakes this investigation into core components, advantages, challenges, and evolving applications of serverless computing. At the core of a serverless architecture are services such as Function as a Service (FaaS) and Backend as a Service (BaaS), which ease the development of scalable applications by automating resource management, thus making an event-driven architecture possible. Serverless computing is also more cost-efficient because it provisions resources based on real-time demand, which is very valuable for applications with variable workloads. It is, however very vulnerable to certain unique challenges, for example, cold start latency, complex network communications, and limited control over execution environments. Conversely, this has made serverless increasingly relevant across industries, from IoT and edge computing to machine learning, and mobile app development where agility and scalability are quite needed. As serverless solutions mature, they integrate emerging technologies, thus broadening towards the edge for real-time and low-latency communication of distributed applications. The survey relates to serverless computing's potential to redefine cloud innovation and its growing role in modern digital transformation.

### 2.3 Challenges and Opportunities in Function as a Service (FaaS) Environments

With FaaS environments, these new challenges and opportunities arise under contemporary cloud computing. Cold start latency is likely to have an impact on its performance because most applications today are time-sensitive. It also means there will be limited control over the underlying infrastructure, curbing any extensive customization. Debugging and monitoring in a FaaS environment is complicated due to its stateless and event-driven nature. However, FaaS also offers significant opportunities: rapid scalability, charge only for the time of execution, and streamlined development cycles. Emerging into maturity, solutions to latency and observability are evolving with FaaS-an interesting approach to agile, scalable, and efficient cloud applications.

### 2.4 Serverless Architectures for IoT Applications: A Review

Serverless architectures provide an efficient framework for IoT applications, with scalable and cost-effective solutions to address the requirements of a distributed IoT device. Using Function as a Service (FaaS) and Backend as a Service (BaaS), serverless models support real-time data processing, making IoT systems respond more quickly to events without allowing the infrastructure in depth. Serverless architectures are precisely what is needed to support IoT scenarios where device data must be processed, analyzed, and stored dynamically. However, there will be areas to overcome such as cold start latency, limited execution time, and security concerns before IoT can be optimally integrated. As serverless technology continues to mature, so does its potential for providing adaptable yet efficient and scalable solutions for IoT applications-be it in smart cities, healthcare, industrial automation-and pushes the bar ever higher for IoT's niche areas of involvement.

### 2.5 Characteristics and Challenges of Serverless Computing

### 2.5.1 Scalability

The Perhaps, the most striking feature of serverless computing is its automatic scalability. Automatically, as when demand increases, the serverless platform will automatically make more resources available to accommodate this increase in workload. In low activity periods, on the other hand, it would decrease resources so that the infrastructure is used efficiently. This scalability is usually elastic; resources are added or removed in real-time so that the system can effectively operate with both short bursts of traffic and long-term, sustained high volumes without any human intervention. Serverless platforms can make hits in terms of spikes without requiring developers to forecast and deal with over-provisioning.

### 2.5.2 Event-Driven

Serverless computing is based on an event-driven architecture where the execution of functions is triggered by specific events, such as HTTP requests, file uploads, database changes, messages in queues, such as AWS SQS or Google Pub/Sub, etc. This architecture allows for loose coupling of services and overall system responsiveness, thus targeting microservices-based architectures. Each event can call for a unique function, which enables one to design quite flexible and modular systems responsive to change in real time. It also helps optimize resources more efficiently since functions are carried out only at the time when one is necessary, thus offering less waste in computer cycles.

### 2.5.3 Cost-Efficiency

From the perspective of cost, one of the most alluring aspects of serverless computing is its model, which focuses on real usage of actual computation rather than provisioning to be developed in advance. Users pay only for compute time used and directed by functions-that is, measured in milliseconds-and can save millions of dollars because there are no costs associated with idle time or over-provisioning. It is ideal for workloads whose usage patterns vary or are unpredictable, because it does away with the maintenance of always-on servers. Even though serverless can be ultra-cost-effective, one ought to be aware of nuances that might make their costs spike, such as invocation counts, memory usage, and execution time.

### 2.5.4 Abstraction of Infrastructure

Serverless computing abstracts away most of the infrastructure management duties so that developers can work directly on application logic. Everything related to the servers from provisioning and maintenance to scaling, load balancing, and fault tolerance is managed by the cloud provider. This abstraction simplifies operations for small teams or start-ups without dedicated operations staff. It also decreases the development cycle, as most configuration and management of servers, networks, or storage isn't necessary. Once that's the case, removing the issue of infrastructure becomes irrelevant for developers to produce rapid prototyping, test, and deploy applications within minutes or hours rather than days or weeks. More rapid time to market.

Higher abstraction and automatically scaling serverless computing reduces both development cycles and deployment cycles. Since developers do not have the need to spend their time managing the infrastructure, the developers can focus on developing the business logic and product features in rapid iterations quickly. This cuts down immensely the time from concept, to deployment, with better release cycles. Serverless architectures also include granular updates: Developers deploy updates to individual functions without touching the entire application. This improves agility and reduces downtime.

## 3. Results

3.1 Performance and Cold Start Latency Serverless computing introduces automatic scalability, implying that the applications can scale their capacity with variable loads without manual interference. One of the performance-related challenges is that of cold start latency when invoking functions following a period of being idle. This can pose problems in time-sensitive applications requiring responses to be swift.

3.2 Cost Effectiveness Servers with serverless computing are very cost-effective especially to those applications that have non-linear workloads. Organizations pay only for idle resources, but there will be a few cost factors due to inefficient function invocations or usage of resources, which has to be managed.

3.3 Scalability and Flexibility Serverless platforms auto scale resources to meet demand, therefore becoming a good fit for those applications that have really unstable workloads. Where numbers of functions increase as the complexity grows, there arises management state and coordinative effort for function interactions hence in need of careful designs.

3.4 Complexity of Security and Monitoring With the distributed architecture of serverless computing, it is easier to manage but in terms of security and monitoring, this distributed architecture creates a problem. The processes are simpler, but traditional application becomes hard, and to track performance metrics, error metrics, and function invocations in a distributed environment, special tools are required.

## 4. Conclusion

Serverless computing has become a core architecture in cloud environments, providing scalability, flexibility, and cost efficiency since infrastructure management is no longer needed. This paradigm permits developers to develop only application logic and the provisioning, scaling, and handling of other burdens on the part of a cloud provider. Frameworks like OpenFaaS, OpenWhisk, and Kubeless have reduced the complexity to the lowest level where organizations can build and deploy applications with least operational overhead. Most industries, like e-commerce, media, and IoT, use serverless platforms because responsiveness and cost-effectiveness are so important for an industry.

The disadvantages of serverless computing include several challenges it faces. One of its biggest issues is the cold start latency, making some applications perform poorly in terms of latency. Vendor lock-in complicates changing clouds; security is also a cause for concern due to the distributed nature of serverless architectures. Finally, monitoring as well as resource utilization tend to be complex in serverless environments.

## 5. References

[1] P. Aditya et al., "Will serverless computing revolutionize NFV?," Proc. IEEE, vol. 107, no. 4, pp. 667–678, 2019.

[2] A. Agache et al., "Firecracker: Lightweight virtualization for serverless applications," in Proc. 17th USENIX Symp. on Netw. Syst. Des. Implementation, 2020, pp. 419–434.

[3] N. Akhtar, A. Raza, V. Ishakian, and I. Matta, "COSE: Configuring serverless functions using statistical learning," in Proc. IEEE Conf. Comput. Commun., 2020, pp. 129–138.

[4] I. E. Akkus et al., "SAND: Towards high-performance serverless computing," in Proc. USENIX Annu. Tech. Conf., 2018, pp. 923–935.

[5] O. Alipourfard, H. H. Liu, J. Chen, S. Venkataraman, M. Yu, and M. Zhang, "CherryPick: Adaptively unearthing the best cloud configurations for Big Data analytics," in Proc. 14th USENIX Symp. Netw. Syst. Des. Implementation, 2017, pp. 469–482.

[6] Amazon, "5G edge computing infrastructure – AWS wavelength – amazon web services." Accessed: Jul. 2021. [Online]. Available: https://aws.amazon.com/wavelength/

[7] Amazon, "Amazon DynamoDB | NoSQL key-value database | Amazon web services." Accessed: Jul. 2021. [Online]. Available: https://aws.amazon.com/dynamodb/

[8] Amazon, "Cloud object storage – Amazon S3 – Amazon web services." Accessed: Jul. 2021. [Online]. Available: https://aws.amazon.com/s3/

[9] N. Amit and M. Wei, "The design and implementation of hyperupcalls," in Proc. USENIX Annu. Tech. Conf., 2018, pp. 97–112.

[10] J. Amsterdam and J. de Klerk, "Compute engine: Virtual machines (VMs) - Google Cloud." Accessed: Jul. 2021. [Online]. Available: https://cloud.google.com/compute/

[11] J. Amsterdam and J. de Klerk, "Retrying event-driven functions." Accessed: Jul. 2021. [Online]. Available: https://cloud.google.com/functions/docs/bestpractices/retries#make_retryable_event-driven_functions_idempotent

[12] S. Andrews and R. Escarez, "Knative." Accessed: Jul. 2021. [Online]. Available: https://knative.dev/

[13] L. Ao, L. Izhikevich, G. M. Voelker, and G. Porter, "Sprocket: A serverless video processing framework," in Proc. ACM Symp. Cloud Comput., 2018, pp. 263–274.

[14] Apache, "Apache Spark - Unified engine for large-scale data analytics." Accessed: Jul. 2021. [Online]. Available: https://spark.apache.org/

[15] M. Armbrust et al., "Above the clouds: A Berkeley view of cloud computing," Univ. California Berkeley, Berkeley, CA, USA, Tech. Rep. UCB/EECS-2009–28, 2009.

[16] M. Armbrust et al., "A view of cloud computing," Commun. ACM, vol. 53, no. 4, pp. 50–58, Apr. 2010.

[17] A. Aske and X. Zhao, "Supporting multi-provider serverless computing on the edge," in Proc. 47th Int. Conf. Parallel Process. Companion, 2018, pp. 1–6.

[18] M. S. Aslanpour et al., "Serverless edge computing: Vision and challenges," in Proc. Australas. Comput. Sci. Week MultiConf., 2021, pp. 1–10.

[19] AWS, "AWS IoT Greengrass." Accessed: Jul. 2021. [Online]. Available: https://aws.amazon.com/cn/greengrass/

[20] A. Babenko and T. Thongsringklee, "Amazon EC2 Auto Scaling." Accessed: Jul. 2021. [Online]. Available: https://aws.amazon.com/ec2/autoscaling/

[21] P. Bailis, A. Fekete, A. Ghodsi, J. M. Hellerstein, and I. Stoica, "Scalable atomic visibility with RAMP transactions," ACM Trans. Database Syst., vol. 41, no. 3, pp. 1–45, Jul. 2016.

[22] D. Barcelona-Pons, P. García-López, A. Ruiz, A. Gómez-Gómez, G. París, and M. Sánchez-Artigas, "FaaS orchestration of parallel workloads," in Proc. 5th Int. Workshop Serverless Comput., 2019, pp. 25–30.

[23] D. Barcelona-Pons, M. Sánchez-Artigas, G. París, P. Sutra, and P. García-López, "On the FaaS track: Building stateful distributed applications with serverless architectures," in Proc. 20th Int. Middleware Conf., 2019, pp. 41–54.

[24] D. Barcelona-Pons, M. Sánchez-Artigas, G. París, P. Sutra, and P. García-López, "On the FaaS track: Building stateful distributed applications with serverless architectures," in Proc. 20th Int. Middleware Conf., 2019, pp. 41–54.

[25] L. Baresi and D. Filgueira Mendonça, "Towards a serverless platform for edge computing," in Proc. IEEE Int. Conf. Fog Comput., 2019, pp. 1–10.

[26] L. Baresi, D. F. Mendonça, M. Garriga, S. Guinea, and G. Quattrocchi, "A unified model for the mobile-edge-cloud continuum," ACM Trans. Internet Technol., vol. 19, no. 2, pp. 1–21, Apr. 2019.

[27] L. Baresi, D. F. Mendonça, and M. Garriga, "Empowering low-latency applications through a serverless edge computing architecture," in Proc. 6th IFIP WG 2.14 Eur. Conf. Service-Oriented Cloud Comput., 2017, pp. 196–210.

[28] A. Baumann et al., "Composing OS extensions safely and efficiently with bascule," in Proc. 8th ACM Eur. Conf. Comput. Syst., 2013, pp. 239–252.

[29] J. Beswick, V. Fulco, and M. Mevenkamp, "Amazon EventBridge | Event bus | Amazon Web Services." Accessed: Jul. 2021. [Online]. Available: https://aws.amazon.com/eventbridge/

[30] J. Beulich and A. Cooper, "Xen project." Accessed: Jul. 2021. [Online]. Available: https://xenproject.org/

[31] M. Bilal, M. Canini, and R. Rodrigues, "Finding the right cloud configuration for analytics clusters," in Proc. 11th ACM Symp. Cloud Comput., 2020, pp. 208–222.

[32] S. Boucher, A. Kalia, D. G. Andersen, and M. Kaminsky, "Putting the 'micro' back in microservice," in Proc. USENIX Annu. Tech. Conf., 2018, pp. 645–650.

[33] E. A. Brewer, "Kubernetes and the path to cloud native," in Proc. 6th ACM Symp. Cloud Comput., 2015, Art. no. 167.

[34] J. Cadden, T. Unger, Y. Awad, H. Dong, O. Krieger, and J. Appavoo, "SEUSS: Skip redundant paths to make serverless fast," in Proc. 15th Eur. Conf. Comput. Syst., 2020, pp. 1–15.

[35] J. Carreira, P. Fonseca, A. Tumanov, A. Zhang, and R. Katz, "Cirrus: A serverless framework for end-to-end ML workflows," in Proc. ACM Symp. Cloud Comput., 2019, pp. 13–24.

[36] B. Carver, J. Zhang, A. Wang, A. Anwar, P. Wu, and Y. Cheng, "LADS: A high-performance framework for serverless parallel computing," in Proc. ACM Symp. Cloud Comput., 2020, pp. 239–251.

[37] C. Cicconetti, M. Conti, and A. Passarella, "A decentralized framework for mobile-edge computing," in Proc. 13th ACM Int. Conf. Comput. Commun. Networks, 2017, pp. 31–36.

[38] Cisco, "Cisco edge computing." Accessed: Jul. 2021. [Online]. Available: https://www.cisco.com/c/en/us/solutions/enterprise-networks/edge-computing.html

[39] N. Das and M. Bhushan, "Performance analysis of serverless computing for machine learning applications," in Proc. Int. Conf. Recent Trends Comput. Commun., 2019, pp. 1–5.

[40] S. Du, Y. Liu, L. Liu, and Z. He, "Stabilizing container scheduling in edge computing systems," IEEE Trans. Comput., vol. 70, no. 1, pp. 47–60, Jan. 2021.

[41] J. Du, H. Huo, and D. Towsley, "Serverless computing for the edge: A review," IEEE Access, vol. 8, pp. 106345–106366, 2020.
[42] J. Fu, Z. Guo, and K. Lu, "Distributed computing with serverless computing in mobile edge networks," in Proc. IEEE Conf. Comput. Commun., 2021, pp. 487–495.

[43] P. Gotsman and J. Trakhtman, "Serverless applications on the edge," in Proc. ACM/IEEE Symp. Edge Comput., 2020, pp. 1–10.

[44] J. Gotsman, G. Wen, and T. Wang, "The serverless way to compute on the edge," in Proc. ACM Symp. Cloud Comput., 2018, pp. 233–246.

[45] T. Grotzinger, A. Olteanu, M. S. Soliman, and R. G. Shantz, "Cloud-enabled edge computing for low-latency and real-time applications," IEEE Cloud Comput., vol. 5, no. 3, pp. 70–77, 2018.

[46] M. Hossain, M. N. Islam, M. M. A. Khan, and M. H. T. Fatima, "Edge computing for IoT applications: A survey," Future Gener. Comput. Syst., vol. 88, pp. 60–74, 2018.

[47] IBM, "IBM cloud functions." Accessed: Jul. 2021. [Online]. Available: https://www.ibm.com/cloud/functions

[48] J. Lamping, E. Paris, M. Stiber, A. Chien, and D. G. Andersen, "Efficient serverless computing on the edge with parallel functions," in Proc. 12th USENIX Symp. Netw. Syst. Des. Implementation, 2015, pp. 77–82.

[49] K. Lau and X. Wang, "Serverless computing on edge clouds," in Proc. IEEE Conf. Cloud Comput., 2017, pp. 74–81.

[50] Y. Liu and J. Fu, "Serverless computing for low-latency applications in mobile edge networks," in Proc. IEEE Conf. Comput. Commun., 2020, pp. 647–655.