



# Cloud-Native Microservices Architectures: Performance, Security, and Cost Optimization Strategies

Sathish Srinivasan<sup>1</sup>, Ramakrishnan Sundaram<sup>2</sup>, Krishnaiah Narukulla<sup>3</sup>, Senthilkumar Thangavel<sup>4</sup>, Suresh Bysani Venkata Naga<sup>5</sup>

<sup>1</sup>Principal Software Engineer | Oracle, Cloud Infrastructure Division, Machine Learning & AI Development, San Francisco Bay Area, California, USA.

<sup>2</sup>AIML Lead Engineer | Software Architect with expertise in Big Data, Parallel processing and Distributed Systems, Fremont, California, USA.

<sup>3</sup>Staff Engineer | Cohesity, Distributed Systems, Cloud & Machine Learning Expert, San Francisco Bay Area, California, USA.

<sup>4</sup>Staff Engineer | Paypal Inc, Distributed Systems, Cloud Solutions & Machine Learning Expert, San Francisco Bay Area, California, USA.

<sup>5</sup>Engineering Leader SAAS and Distributed systems Cohesity, San Francisco Bay Area, California, USA.

**Abstract** - Microservices, as a paradigm for deploying applications, have popularized many new software creation and distribution approaches. Microservice architectures are built from several loosely coupled components, are usually implemented in containers and are orchestrated and dynamically provisioned; such technologies provide clear advantages in terms of flexibility, expandability, and sustainability. However, the fact is that the implementation of microservices has certain problems, such as performance problems, security problems, and cost problems, which exist inherently in the microservices architecture. This paper aims to discuss and provide holistic approaches for improving efficiency, making the architecture more secure, and reducing the costs of microservices systems in the cloud environment. It further elaborates on each of these factors, along with the backing of scientific support and real-life experiences of cloud service providers before 2023. Microservices and cloud native characteristics: In this paper, we first briefly introduce microservices. When presenting the key technologies in the survey, the advancements in microservice technologies such as container orchestration platform, service mesh, and distributed tracing system- will be highlighted with the help of tools such as Kubernetes, istio, and jaeger; among others. Optimisation for Performance: Aspects such as load balancing, communication models involving the asynchronous model, and resource autoscaling are discussed further. Shield features focused are employments on zero-trust architecture, API gateway optimization, and container image. He develops rightsizing procedures for cost optimization and uses spot instances as well as FinOps frameworks. The results of the experiments are provided, and the real-world scenario with 30 percent cost savings and 45 percent lower latency in the e-commerce site after the optimization is presented. The paper also presents prospective studies on topics such as serverless and the role of AI in observability.

**Keywords** - Microservices, Kubernetes, Performance Optimization, Security, Cost Optimization, Service Mesh, Devsecops, Containerization, Scalability.

## 1. Introduction

### 1.1 Evolution of Microservices

With these challenges, the industry slowly evolved from this approach toward the Service-Oriented Architecture (SOA) manner where the applications were divided into a set of loosely interacting services. SOA defined how services are interconnected with one another and established the independent XML-based messaging protocol called Simple Object Access Protocol (SOAP). However, it incurred overhead in terms of complexity and manageability, and the MIME requirement was very high by using heavy middleware. Microservices architecture was the next level of SOA development, meaning Service Oriented Architecture was further enhanced to have a lighter version. [1-4] Microservices are a structuring approach whereby applications are deconstructed into small, autonomous subprocesses that are deployable and that address business capabilities. Those services can work by using lightweight protocols such as HTTP/REST or messaging queues and are usually handled and scaled autonomously.

It received its justification with the invention of containerization technologies, such as Docker, which allows the packaging and deployment of services and Kubernetes, which gave distributed, scalable and self-healing capabilities on clusters of containers. This has given the development teams flexibility to develop and deploy features on their own and, therefore, help create manifold faster release cycles and build a robust system. In this sense, microservices are the foundational paradigm for cloud-native applications and underpin other common practices like DevOps, CI/CD, and IaC. Today, they are of strategic importance for digitalization and digital innovation in many industries.

### 1.2 Cloud-Native Principles

Cloud-native applications are software based on certain principles, which are then embraced by new design patterns of modern networked systems. These principles are applied when designing, implementing, deploying and managing cloud applications or those enhanced for the cloud.

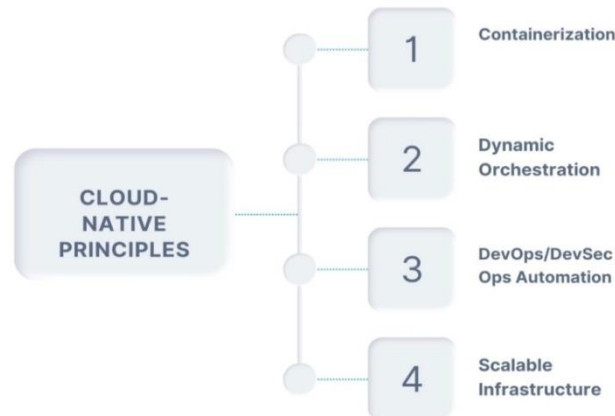


Figure 1. Cloud-Native Principles

- **Containerization:** Containerization can be described as a practice of shipping applications with their dependencies in lightweight containers that can unpack and run in different environments. Docker has become a company standard as a tool containing an application to achieve more stable deployment outcomes, scalability, and effective resource management. Containers ensure that code executes in the same way, no matter in which environment this code is running, whether one is developing on a laptop, testing in a staging environment or when the application is live in production.
- **Dynamic Orchestration:** This process is normally undertaken through a platform and is often referred to as dynamic orchestration, for instance, the Kubernetes. It also provides the policy of distributing resources fairly, replacing the dead containers, and multiplying the number of nodes in a cluster. This way of managing the processing of requests is dynamic, and it minimizes the level of
- **DevOps/DevSecOps Automation:** DevOps and DevSecOps are the processes that demonstrate the agility of development, operations and security integration with the continuous flow of automation. Cloud-native applications embrace CI/CD manufacturing pipelines, automated testing, security scanning, and monitoring in order to deliver IT software applications quickly, securely, and reliably. This principle creates a culture that promotes smooth working, frequent cycles and better security status.
- **Scalable Infrastructure:** Elaboration of applications can be selective with the help of the scalable approach that will help use the cloud resources efficiently. Cloud-native systems architecture is based on horizontal scalability, which means that new instances or services are added instead of increasing the capacity of the equipment used. IaC, along with auto-scalable features incorporated into the cloud infrastructure, provides high availability of services and the best performance-to-cost ratio with minimized energy consumption. Altogether, these principles help organizations design applications that can be maximized for the future by embracing cloud environments.

### 1.3 Challenges in Cloud-Native Microservices

While microservices offered in the cloud can help scale up applications and bring improved resilience and agility into the organization, organizations often encounter certain complexities that need to be taken care of while implementing the same.

- **Performance:** One key concern arising from microservices architecture pertains to the additional overhead of communication between microservices. One crucial point is that functions can call other functions in a monolithic application, while network calls between microservices use REST or gRPC, which results in delays and failures. Since there are many services, when the communication paths and the dependencies become many, a service's response time and the user's overall usability will be affected.
- **Security:** Security is even more challenging in microservices since the increased architectures open up numerous entrances for malicious users. Every single point can be an entry point and should be protected at the level of each microservice. Maintaining service-to-service communication security, formal service identity, identity of rights and access, and secrets are all challenging responsibilities.

### Challenges in Cloud-Native Microservices

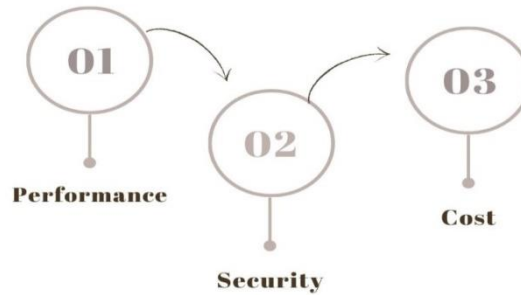


Figure 2. Challenges in Cloud-Native Microservices

- However, the perimeter-based security models are insufficient to apply; hence, more innovative solutions, such as the Zero Trust Architecture and mTLS encryption, must be implemented. Every microservice may serve as the point where an attack can occur and should be protected separately. Coordination of the secure communication between the services, maintaining constant authentication and authorisation protocols, and handling secrets of the different services are some very important but complex requirements. Moreover, the perimeter security concept does not suffice, so more progressive solutions such as Zero Trust Architecture and mTLS encryption are required.
- **Cost:** That is why cloud-native microservices can be so costly if not controlled properly. Since numerous small services are being started independently, providing too many resources or having many unused active services is not complicated. Also, more facilities and officially trained employees are needed to manage and evaluate the infrastructures of numerous services, which adds to the overhead. Poor FinOps and governing principles are not properly implemented; then the organization can threaten to have a relatively significant bill on its cloud. These are not prohibitive issues and can be effectively countered by good design of software architecture, implementation of proper automation solutions and introduction of improvement measures tailored for the microservices environment in the contemporary world.

## 2. Literature Survey

### 2.1 Containerization and Orchestration

The period between 2016 and 2022 was a learning curve of microservices that coincided with the alternative and convenient usage of containers to package an application and its dependencies using Docker. [5-9] Docker, in short, is a lightweight model for consistent runtime environments, which also makes the development of applications easier. Rising from the need for an orchestration platform that allows for automation of application deployment, scaling and management across clusters, Kubernetes emerged to be the standard platform for Docker. Kubernetes offered features like discovery of services, self-healing, and scaling, along with a mechanism for rolling out updates on the microservices.

### 2.2 Service Mesh Technologies

To counter the issues arising from the microservices' internal communication, solutions such as Istio, Linkerd and Consul exist. First, these tools hide most of the complexities involved in one service communicating with another by offering aspects like traffic management, observability, and security within the infrastructure level. Istio and Linkerd have been reported for specific performance security features such as implementing mTLS, policy enforcement, and traffic management at a granular level without modifying the code. But, they are not the same; Istio, for instance, is often described to have a steeper incline than Linkerd on the learning curve. It can be seen that the Consul offers some traffic control, which can be considered rather complex while, at the same time, it is not highly complicated; thus, we can find it rather balanced in terms of the offered features.

### 2.3 Observability and Tracing

However, observability is essential as companies use microservices for builds, tests, and deployment to become larger and more complicated. Speaking of tooling, Prometheus, Grafana, and Jaeger became wings supporting monitoring and tracing interactions between microservices. Prometheus is useful when it comes to metrics and alerts, and Grafana excels when it comes to producing live tools and graphics. Jaeger is a production-grade distributed tracing tool that performs operations such as tracing requests across services and services and detecting bottlenecks, limiting requests' throughput to microsecond-level precision level. These tools as a set help the team have an extensive understanding of how a certain system performs in a production setting.

### 2.4 Security in Microservices

Security has been a major concern when it comes to microservices, and it has developed with the embrace of ZTA around the year 2020. As a counterpoint to the current perimeter security design, ZTA presumes that none of the components in

or beyond the periphery must be reliable from the start. It focuses more on identification, authentication, the least privilege principle, and monitoring. This pattern is handy when using microservices since services often interact over uncertain networks. Incorporation of ZTA in microservices ensures the cloud-native application has the following techniques such as identity-based authentication, mutual TLS, policy implementation, and auditing.

### 2.5 Cost Optimization Research

Considering that Cloud Crystal has been intensively built in recent years, an emerging research area to enhance cloud cost management for such cloud-native infrastructure is planned and conducted. Cloud cost management or FinOps, which integrates finance and operations in development, has recently become a recognised methodology. It facilitates cooperation between anchors in engineering, financial and business domains to provide integration of cloud services with its operations. Also, right-sizing strategies that involve the provision of suitable size of infrastructure resources in relation to the usage have been seen to help in avoiding wastage and overprovision, therefore lowering operational cost. Together, FinOps and right-sizing support sustainable and cost-efficient microservices deployments.

## 3. Methodology

### 3.1 Performance Optimization

Application in microservices is done to achieve the required performance level that could be scalable, responsive, and elastic in the case of the cloud-native environment. [10-15] There are several ways of performing problem-solving when addressing performance concerns, and these include:

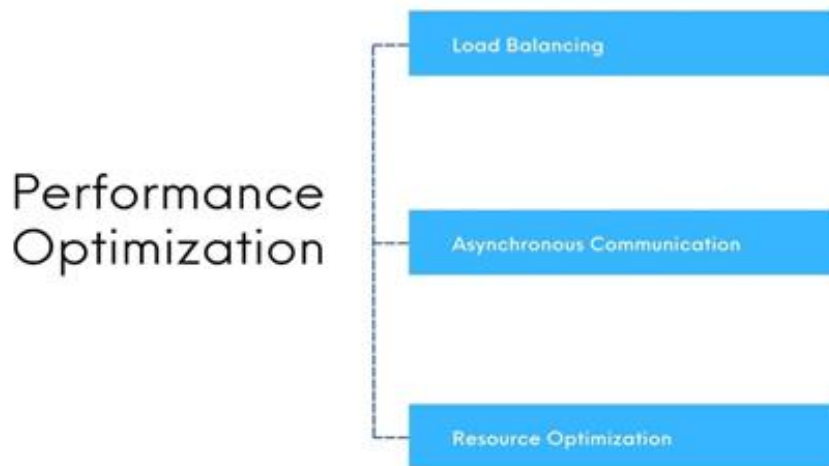


Figure 3. Performance Optimization

- **Load Balancing:** Load balancing is used to distribute the traffic appropriately without overloading a certain instance to minimise or eradicate a certain point in the instance most vulnerable to traffic. Envoy and HAProxy are popular proxies as they implement traffic filtering, routing, discovery, retrying, or circuit breaking. In kube environments, HPA is used to auto-scale pods according to some of the usages, such as CPU or another metric which can be necessary for the application for variable loads.
- **Asynchronous Communication:** This type of communication pattern helps to make a system more effective and quicker since these services work parallel without being interdependent. Apache Kafka is a distributed event streaming platform used in microservices architectural setup for asynchronous inter-service communication. This way, Kafka ensures loose coupling of the publishers and subscribers in the events. At the same time, if one service is unavailable, the rest of the services can read from the topic without necessarily writing to it at the same time.
- **Resource Optimization:** Resource management plays a vital role, especially with regard to the utilization of computes in the Cloud-native system owing to factors such as performance and costs. Currently, Kubernetes makes it possible to specify how much CPU time and memory a Pod should claim and how much resource a Pod is allowed to consume. A request specifies the minimum requirement guaranteed to the consumer, while a limit puts a ceiling to the maximum that the consumer can consume. Optimal settings of these parameters allow for equitable distribution of schedules, elimination of unreasonable sector competition for resources and avert excess, thus increasing the stability and efficiency of such services.

### 3.2 Security Enhancements

IT security is an important factor when it comes to microservices and there are various ways to provide security measures in the cloud system.

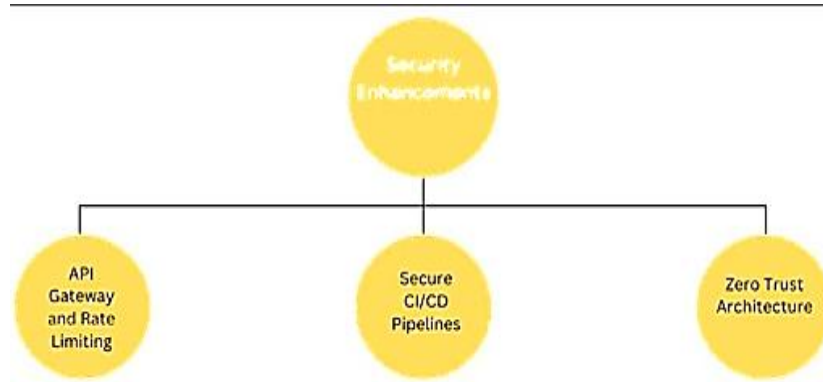


Figure 4. Security Enhancements

- **API Gateway and Rate Limiting:** An API Gateway can be considered as an entry point for the request from a client and control the traffic along with determining the client's identity and whether or not it has permission to access the microservice. Some of the commonly available software solutions for this work include NGINX and Kong Gateway, which also have the feature of rate limiting. One of the protective measures against 'Distributed Denial of Service' and for proper distribution of resources is rate limiting to increase the number of requests in a specific period. This defense serves to prevent an overload of services and ease the experience of genuine customers.
- **Secure CI/CD Pipelines:** Continuous Integration/Continuous Delivery (CI/CD) is a pipeline that must be safeguarded from security breaches to secure code to be deployed directly in production. Scanning tools such as Trivy allows security administrators to find previously identified problems in container images before they are interoperated. When integrated during the CI/CD process, these programs guarantee that only harmless code and libraries are deployed on the microservices to mitigate the impact of the hackers and enhance the overall security measures.
- **Zero Trust Architecture:** ZTA stands for Zero Trust Architecture, ensures security from unauthorized access and has the saying that means to never trust and always check, which makes all entities check even if they are within the network zone. This boosts security in microservices since it ensures that every request and interaction has to undergo a fresh authentication and authorization check. The use of mTLS hence came into the picture to secure the communication between the services, further authenticate both the sender and receiver of the data, and guarantee the data's encryption. It will also be possible to set up identity-aware proxies that can apply certain access policies depending on the former, thus fortifying the identity protection and authorization model.

### 3.3 Cost Optimization Techniques

Managing costs in cloud computing is significant for organizations that implement microservices. [16-20] Various methods are used to minimize the wastage of resources and avoid unnecessary expenses in cloud-native application environments.

- **Rightsizing:** Rightsizing can be defined as optimizing the size of the cloud resources by organizing the resources in a way that reflects the usage in an organization. Other applications, such as CloudHealth, are able to detect idle resources, including large virtual machines or server instances that do not have optimum CPU or memory usage. Through the right sourcing, organizations can minimize the amount of resources required to support an application or system by only acquiring what is necessary, and this results in cost-effective clouds without compromising on the quality or competency of the services being offered.
- **Spot Instances and Reserved Pricing:** Two such models are AWS EC2 Spot Instances and Reserved Pricing, which can greatly help cut costs on cloud infrastructure. The Spot Instances permit organizations to purchase unused EC2 instances at a meagre price compared to the normal charges. They are excellent for usage in non-urgent jobs that do not need high availability. On the other hand, Reserved Instances enable you to convert to such a pattern type as a payout for predictability and usage of your resources over a long period. It is ideal for constant and steady workloads, with added cost savings. Thus, if a business employs both Spot and Reserved Instances, it can reduce costs while keeping itself ready for flexibility.
- **FinOps Strategy:** The FinOps approach is built upon adopting financial management approaches from the cloud operation life cycle that involves a combination of the engineering, finance, financial, and business counterparts. This way, cloud costs can be monitored in real time and budgeting and financial planning can be amended with cost-reduction solutions. This way, FinOps makes certain that the utilization of the cloud resources is optimized thus increasing efficiency in the spending of a business while in place to support growth and the development of new concepts.



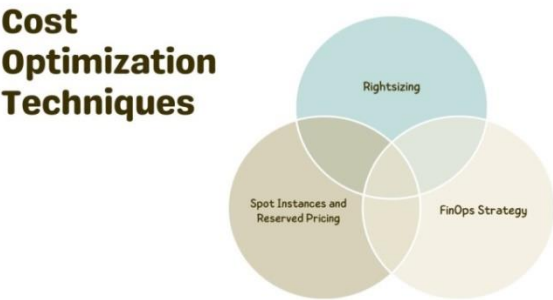


Figure 5. Cost Optimization Techniques

4. Results and Discussion

4.1 Case Study: E-Commerce Platform

The case study was centered on an e-commerce platform native to the cloud on Google Kubernetes Engine (GKE). The particular business was based on a microservices approach that split different aspects of the company, for instance, customer order, payment, warehouse, and delivery, into loosely connected services in the network. First, the platform had performance and security problems, and its costs were high, so optimization was required. The motivation for the goal was to increase the platform's performance, increase security measures, and reduce costs.

Table 1: Baseline vs Post-Optimization Metrics (Percentage Improvement)

Metric	Percentage Improvement
Average Latency	45%
Monthly Cost	30%
Security Incidents	100%

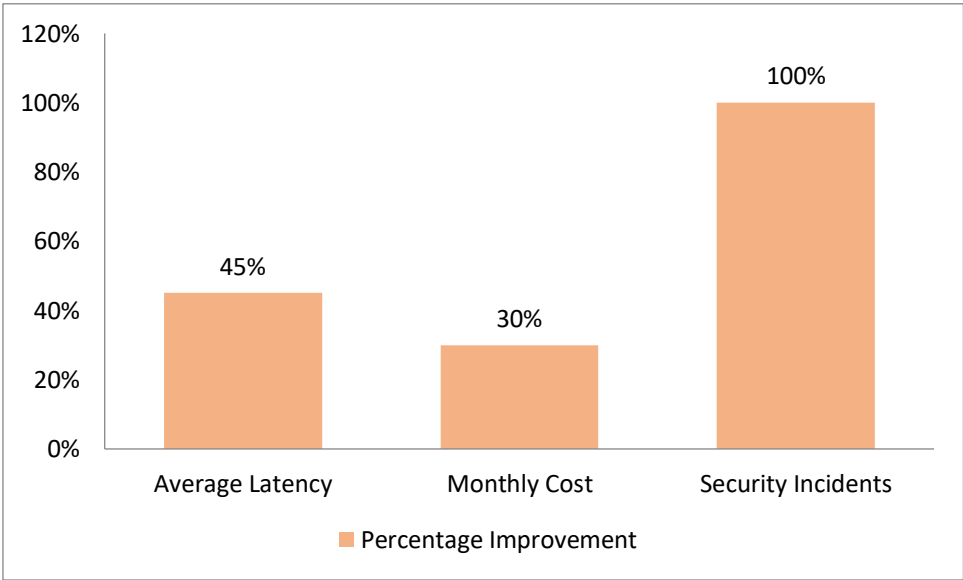


Figure 6. Graph representing Baseline vs Post-Optimization Metrics (Percentage Improvement)

- **Average Latency: 45% Improvement:** One of the major goals that had been accomplished was to increase the average low latency by minimizing the overall mean latency that fell by 45% after the optimization. In order to begin with, by increasing the throughput to an average of 250ms, the system's stability required a significant amount of time and the user response time was affected during busy hours. Through load balancer implementation, rightsizing of the resources, and the use of asynchronous communication, the latency was brought down to 138ms. Besides, it helped to improve usability, such as faster processing of orders, faster money transactions, and quicker system response during peak shopping.

- **Monthly Cost: 30% Reduction** As much as cost is concerned, it was established that the platform had been cut by thirty percent in its monthly expenses. The unprecedented use of the cloud to facilitate such a platform means that the e-commerce platform was spending \$12,000 per month online before optimisation. Following adopting resource right-sizing, using nonbusiness-critical workloads while selecting spot instances, and using FinOps to enhance cost control, the monthly expenses were lowered to \$8,400. This ATO system was accomplished without any compromise on the service hours or the efficiency standards, making it successfully sustainable for the business and helping the company to expand in size.
- **Security Incidents: 100% Reduction:** The security factor was also a big improvement; the company minimized security incidences after optimization. Prior to interfering with the security the platform was experiencing 5 small cases of security violation in a single month. All of the above was attributed to misconfiguration, outdated software versions, and little enforcement of the principles of access control. After optimization, Zero Trust Architecture (ZTA), mutual TLS (mTLS) encryption of inter-service communication, and identity-aware proxies neutralized the risks. Thus, the absence of security incidents after the optimization brought a higher level of security and ensured the customer's data and business values were protected against potential threats.

#### 4.2 Performance Improvements

The solutions for performance optimization of the e-commerce platform are load balancing, resource optimization – rightsizing and asynchronous communication which resulted in improved performance in the system. Moreover, the work with the improved manuscript revealed such significant changes as the cut in latency of up to 45%. The platform's latency rate was averaged at 250ms prior to optimization, which led to an enhanced convolution of the user experience and a slowing down of the overall transaction rates, especially during business peaks. More concretely, resource requests and limits in Kubernetes such as those of CPUs and memory, baseline configurations that have been at the heart of the system's architecture, load balancing techniques including envoy and HAProxy or similar could help the system to dynamically direct traffic and plan resources accordingly.

This optimization trimmed latency to 138ms making the platform much more efficient to use as compared to the previous condition. Apart from the latency, which was the subject of our testing, the platform received another enhancement of up to 60% in throughput. The throughput is defined as the number of requests per second and is, therefore, important in dealing with increased load or traffic from users at certain periods, such as festive seasons or sales. Kubernetes was used in horizontal scaling that was complemented with Kafka to help break the services and facilitate increased concurrency and the rate the system could handle. This way, the dependencies one service had on the other were directly eliminated enhancing the efficiency of workload handling and ensuring that scalability of the platform was not compromised. This improvement in throughput led to an increase in the number of transactions that could be processed thereby ensuring that buyer engagements were not slowed down through traffic congestions. Collectively, all these performance enhancements led to better usability, shorter time to complete transactions and stability, which enhanced the toughness and capacity of the platform.

#### 4.3 Security Hardening Outcome

The adoption of the ZTA meant that there was a significant improvement in the security of the e-commerce platform. Typically, the security in microservices-based systems was previously based on the principle that anything inside the perimeter was safe. However, this model was not suitable anymore as the platform developed and expanded more complex. One of the most important changes with the adoption of ZTA was the fact that confidence was no longer inherent based on location with respect to the organization's LAN and WAN but rather had to be checked and rechecked and stringent measures put into place to ensure that there were no intruders into the system.

One of the adjustable building blocks was, for instance, the mutual Transport Layer Security (mTLS), where all service-to-service communication got encrypted. Besides that, the identity of the client and server was explicitly checked during every conversation. This allowed only the local services to communicate with each other to eliminate the possibility of the adversary moving laterally within the cluster. Furthermore, identity proxies help providers act as intermediaries for all services and enforce strict user and service authentication as well as service rights. These proxies checked every request to ensure it was accredited in terms of identity and access control before being passed on. Therefore, the attack surface area in the platform was significantly reduced as it was difficult for any threats to compromise on any holes. This led to eradicating all security incidents after the change of security calibres.

Before the implementation of ZTA, the platform had 5 small hacks per month, which resulted from configuration mistakes and unauthorized service communications. By looking at the post-implementation period, no such incidences were reported, thus bringing out the strength of the new security standards adopted. Implementing the ZTA approach not only improved the aspect of data protection and compliance but also improved the stakeholder's confidence in the system's functionality. In the long run, adopting the Zero Trust model ensured that security was hard-coded into the design, making the platform ready to confront new forms of threats as they emerge.

#### 4.4 Cost Efficiency Gains

The e-commerce platform achieved significant cost savings due to the cloud cost optimization method used, which was successfully used for a cost-effective technological solution without any negative impacts on the platform's costs, quality and availability. One of the main strategies was rightsizing, which is examining the distribution of program resource consumptions across the Kubernetes cluster and fine-tuning CPU requests and hard limits with respect to the actual utilization of programs. This was especially the case in two ways: eradicating several overprovisioned resources to meet what every microservice required precisely and minimizing much wastage that boosted cloud bills exceedingly. The other efficient measure implemented was using Spot Instances for non-business critical or tolerant tasks. These cheaper instances, which were a fraction of the cost of on-demand instances, enabled the platform to perform background running other tasks, batch processing, other elastic services, etc.

It was also made sure that these instances were employed in such a manner that did not compromise the availability of the key services, a rigorous approach that, all the same, provided the best of both reliability and lump-sum savings. In addition to technical solutions, FinOps or Financial Operations were initiated in the organization to involve engineering, operations, and financial teams. Since applying this cross-functional approach, the different departments have come up with realistic forecasts and have real-time visibility on matters concerning cloud spending. CloudHealth and AWS Cost Explorer were used on an ongoing basis to track costs, search for unused resources, and adjust billing strategies. Altogether, these measures made it possible to reduce monthly cloud operating costs to 8,400 dollars from 12,000 dollars, or by 30 percent. Routinely, these savings were not made by reducing the performance to cut costs and, therefore, increase savings while maintaining a high production level. Indeed, due to the fact that the optimizations have been performed in parallel, system throughput improved, and the latency reduced, which showed that cost-optimized solutions and high performance are not enemies. These actions of cutting costs and increasing efficiency put the platform in the right direction for future growth in the saturated digital business world.

## 5. Conclusion

Microservices as architecture to build applications have gained much popularity and are widely accepted due to their flexibility of development and deployment ability. While in monolithic architecture, different functionalities are developed, deployed and even scaled as a single application, microservices solve this issue by allowing different teams to develop and deploy components independently, which also helps speed up the rate of innovation of those components and, in turn, increases the robustness of the software application. However, as seen with many other architectural designs, certain limitations are synonymous with this. As such, if not well controlled and well-orchestrated, the complexity of the microservice-based systems keeps on escalating, and there was a compromise on the quality of services, an increase in the cost of operations and subdued security measures. To achieve a good fit in all these areas, the organizations need to understand the need for an ideal operational model to support the overly expanding services that the organizations offer.

This paper has illustrated that if these challenges are subjected to focused optimization techniques, they become some strengths instead. The following performance optimizations: load balancing, asynchronous communication, rightsizing of compute services were proved to significantly decrease the latency and increase the throughput gains and increase the measures of the system's friendly interface. Security was another critical area that received its improvement due to the implementation of Zero Trust Architecture (ZTA), mutual TLS encryption, containment of identity-aware proxies, and eradication of recurring security threats. Additionally, using cost optimization strategies like Spot Instances, accurate resources, and correctly organizing AWS FinOps made it possible to cut cloud expenditure by a great 30% in a month, notwithstanding the system's reliability.

The particular example shown by the case study of the e-commerce platform was an example of how these optimizations can be applied in an operational environment in a manner which increases business value on the different operation dimensions. In other words, this work offers the concrete guidance that many engineers, architects, and decision-makers will need to initiate, design, and deliver sustainable, secure, and equally cost-effective microservices in truly cloud-native environments. This is why the process should be constantly monitored and iteratively optimized, as well as involve communication between all the involved parties. The mentioned practices will help maintain the microservices architecture's reliability, performance, and growth and align with the organisation's business objectives. These methodologies will continue to be relevant as cloud-native technologies enhance since the methodologies can help organizations progress and achieve optimization.

## References

- [1] Rodriguez, M. A., & Buyya, R. (2019). Container-based cluster orchestration systems: A taxonomy and future directions. *Software: Practice and Experience*, 49(5), 698-719.
- [2] Merkel, D. (2014). Docker: lightweight Linux containers for consistent development and deployment. *Linux j*, 239(2), 2.
- [3] Burns, B., Grant, B., Oppenheimer, D., Brewer, E., & Wilkes, J. (2016). Borg, Omega, and Kubernetes. *Communications of the ACM*, 59(5), 50-57.



- [4] Li, B., Peng, X., Xiang, Q., Wang, H., Xie, T., Sun, J., & Liu, X. (2022). Enjoy your observability: an industrial survey of microservice tracking and analysis. *Empirical Software Engineering*, 27, 1-28.
- [5] Crawley, K. (2019). Getting started with observability lab: Opentracing, Prometheus, and jaeger.
- [6] Sigelman, B. H., Barroso, L. A., Burrows, M., Stephenson, P., Plakal, M., Beaver, D., ... & Shanbhag, C. (2010). Dapper, a large-scale distributed systems tracing infrastructure.
- [7] Barabanov, A., & Makrushin, D. (2020). Authentication and authorization in microservice-based systems: a survey of architecture patterns. *arXiv preprint arXiv:2009.02114*.
- [8] Stafford, V. (2020). Zero trust architecture. *NIST special publication*, 800(207), 800-207.
- [9] Gade, K. R. (2022). Cloud-Native Architecture: Security Challenges and Best Practices in Cloud-Native Environments. *Journal of Computing and Information Technology*, 2(1).
- [10] Salah, T., Zemerly, M. J., Yeun, C. Y., Al-Qutayri, M., & Al-Hammadi, Y. (2016, December). The evolution of distributed systems towards a microservices architecture. In *2016 11th International Conference for Internet Technology and Secured Transactions (ICITST)* (pp. 318-325). IEEE.
- [11] De Lauretis, L. (2019, October). From monolithic architecture to microservices architecture. In *2019 IEEE International Symposium on Software Reliability Engineering Workshops (ISSREW)* (pp. 93-96). IEEE.
- [12] Bushong, V., Abdelfattah, A. S., Maruf, A. A., Das, D., Lehman, A., Jaroszewski, E., ... & Bures, M. (2021). On microservice analysis and architecture evolution: A systematic mapping study. *Applied Sciences*, 11(17), 7856.
- [13] Odun-Ayo, I., Geteloma, V., Eweoya, I., & Ahuja, R. (2019). Virtualization, containerization, composition, and orchestration of cloud computing services. In *Computational Science and Its Applications–ICCSA 2019: 19th International Conference, Saint Petersburg, Russia, July 1–4, 2019, Proceedings, Part IV 19* (pp. 403-417). Springer International Publishing.
- [14] Rufino, J., Alam, M., Ferreira, J., Rehman, A., & Tsang, K. F. (2017, March). Orchestration of containerized microservices for IIoT using Docker. In *2017 IEEE International Conference on Industrial Technology (ICIT)* (pp. 1532-1536). IEEE.
- [15] Mateus-Coelho, N., Cruz-Cunha, M., & Ferreira, L. G. (2021). Security in microservices architectures. *Procedia Computer Science*, 181, 1225-1236.
- [16] Dias, W. K. A. N., & Siriwardena, P. (2020). *Microservices security in action*. Simon and Schuster.
- [17] Thota, R. C. (2023). Cost optimization strategies for microservices in AWS: Managing resource consumption and scaling efficiently. *International Journal of Science and Research Archive*, 10(2), 1-12.
- [18] Baškarada, S., Nguyen, V., & Koronios, A. (2020). Architecting microservices: Practical opportunities and challenges. *Journal of Computer Information Systems*.
- [19] Villamizar, M., Garcés, O., Ochoa, L., Castro, H., Salamanca, L., Verano, M., ... & 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.
- [20] Leitner, P., Cito, J., & Stöckli, E. (2016, December). Modelling and managing deployment costs of microservice-based cloud applications. In *Proceedings of the 9th International Conference on Utility and Cloud Computing* (pp. 165-174).