# COST OPTIMIZATION STRATEGIES IN FINTECH USING MICROSERVICES AND SERVERLESS ARCHITECTURES

## Sumit Bhatnagar
Vice President, JPMorgan Chase & Co., New Jersey, USA

**Abstract**- This article delves into the tactics employed by the fintech industry to optimize costs through the utilization of serverless and microservice architectures. Analyzing real-world finance applications, it proves that these designs boost performance and scalability while decreasing operational, maintenance, and infrastructure expenses. Innovative solutions to lower costs and increase operational efficiency are being sought after by enterprises in the increasingly competitive fintech sector. Microservices and serverless architectures are the focus of this study as it delves into how finance organizations may optimize their costs. It goes over the basics of these designs and how they can improve scalability, flexibility, and resource consumption compared to conventional monolithic systems. This research sheds light on the practical implementation of these strategies through an in-depth case study of FinTech Bank. It demonstrates how the bank has reduced infrastructure costs and accelerated service deployment by transitioning to a microservices-based approach and adopting serverless computing.

**Keywords**- Microservices and Serverless, fintech organizations, Cost Optimization

## I INTRODUCTION

A combination of rising expectations from younger generations and innovative technological capabilities has revolutionized the financial technology sector. This evolution is based on the decision to select an appropriate architectural strategy to construct reliable, extensible, and secure financial applications. This study focuses on comparing two major architectural patterns, serverless and microservices architecture. FaaS is a cloud computing model that is frequently described as serverless computing because it allows developers to create and deploy applications without having to worry about the supporting infrastructure. This model also suggests that operational overhead will be cut, scaling will be managed automatically and cost structures will involve pay-per-use, which will be interesting for FinTech startups and large institutions. While, the microservices architecture is the development of applications as small, autonomously deployable services, where each service is a separate process and communicates with other services simply. They include the following; Modularity is improved, and it is easier to scale and adopt several technologies for the different components. In the process of attempting to grow while offering the best security, regulatory compliance, and efficiency, FinTech companies face a significant decision: whether to follow the serverless architecture or microservices architecture or use both.[1] The aim is to give an overview of all these architectural patterns and compare them to FinTech applications, especially the key issues about their principles, implementation, and issues. After discussing practical examples and operational characteristics, it will be possible to assess the applicability of the described architectures to the financial industry requirements regarding the availability of high-volume transactions, data synchronization, and compliance. In response to rising expectations for efficient and scalable solutions, the financial technology sector is undergoing fast change. There has been a shift away from traditional monolithic architectures and toward microservices and serverless models, which provide better resource usage and lower costs. With an emphasis on infrastructure, development time, and operational costs, this paper explores how fintech organizations might minimize costs by using these approaches. The simplicity, cheap cost, and scalability of serverless computing have made it a hot commodity in recent years. However, proper cost management in this context can be a challenge for firms utilizing serverless architectures. Streamlining expenses in serverless computing is the goal of this blog. In the previous section, we will go over the benefits and idea of serverless

computing. Serverless computing, also known as Function as a Service (FaaS), frees up developers from server administration and configuration tasks so they can focus on developing software. Due to the fact that resources are only made available upon the execution of a predefined function, this event-driven architecture permits both quick creation and simple scaling. While serverless computing follows a pay-per-use model, improper cost management can lead to unexpected expenses. Identifying the root causes of cost overruns and implementing cost-control strategies is essential for optimization.[2]



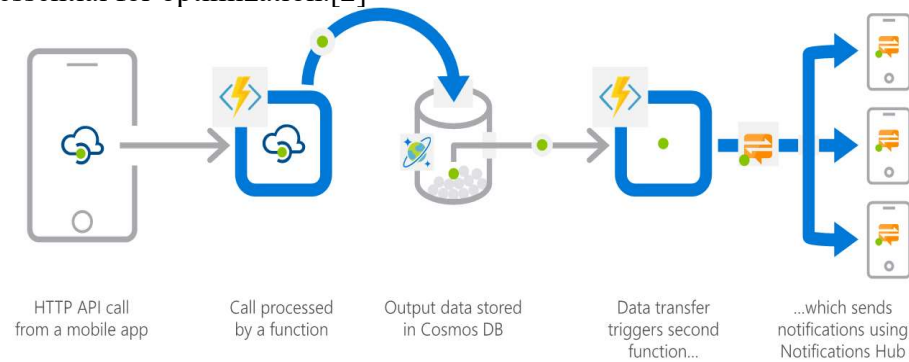| HTTP API call from a mobile app | Call processed by a function | Output data stored in Cosmos DB | Data transfer triggers second function... | ...which sends notifications using Notifications Hub |

Fig.1 microservice architecture with fintech

Optimal resource allocation, Auto Scaling, cold starts, caching, and cost monitoring tools are just a few of the methods as ways to cut expenses in serverless computing.[3] Also give examples from the actual world to show how the ideas work and go over some best practices for designing buildings that don't break the bank. The purpose of this paper is to instruct readers about the challenges of optimizing serverless costs and the solutions available for these problems. Organizations may get the most out of serverless computing without breaking the bank if they employ the correct tactics.

**Problem Statement**

Fintech companies face increasing pressure to deliver secure, scalable, and reliable platforms, often incurring substantial infrastructure costs. This study explores modern architectural strategies to minimize these expenses while maintaining high performance and regulatory compliance.

**Research Objectives**

- Evaluate cost reduction through microservices architecture.
- Investigate serverless computing impact on operational costs.
- Provide a comparative analysis of cost optimization in different fintech environments.

**Emergence of Microservices and Serverless Paradigms**

Migrating to microservices and serverless architecture was the most recent big move in FinTech architecture [2]. Contemporary methods provide:

- Fine-grained, loosely-coupled services
- Independent deployment and scaling
- Improved fault isolation
- Faster time-to-market for new features

Microservices allow fintech companies to develop, deploy, and evolve services independently, improving flexibility and scalability. Serverless computing, on the other hand, enables businesses to run code without managing servers, reducing operational overhead. The evolution is due to the FinTech industry's constant search for architectures that allow fast innovation and growth while keeping up with the security and compliance levels of the financial international market. In this context, the industry advances in implementing these patterns, sometimes using fragments from several approaches to accommodate concrete business requirements and technological opportunities. Core Principles of Serverless and Microservices Architectures: FinTech has embraced serverless and microservices architectures to enhance scalability, agility, and innovation. While these approaches

share common goals, they differ in their implementation and guiding principles.

Lessons on Serverless Computing Foundations: Serverless computing, also known as Function-as-a-Service (FaaS), allows developers to focus on writing code while the cloud provider manages resource scheduling and allocation. Key principles include: [4-6]

**Event-Driven Execution:** An event is anything that happens such as an HTTP request, changes to the database, or a scheduled event that will cause a function to execute.

**Stateless Nature:** Functions do not have a state between invocations; this enhances scalability and reduces the needed model in programming.

**Auto-scaling:** The platform means that resources can scale from zero volume to the maximum volume [4].

**Pay-per-Use Pricing**: Charges are by the rate of actual consumption of the time slices allocated to compute as opposed to licenses.

**Managed Infrastructure:** For all the server-related issues, the responsibility lies on the cloud provider thus enabling developers to work on code only.

**Formula for Serverless Cost Calculation:**

Total Cost = (Number of Invocations × Execution Time × Cost per 100ms) + (Memory Allocated × Execution Time × Memory Price)

## II RELATED STUDY

The earlier work [7] focused on exploring the effects of different encoding and transmission protocols on edge resource consumption, with the intention of using microservices to execute real-time edge analytics. They left the analysis of computing jobs for future work. By comparing the resource consumption of real data analysis in the context of Microservice and FaaS frameworks, this article adds to and expands upon their earlier results. According to previous findings, compared to other data encoding methods (like XDR) and transmission mechanisms (like WebSockets), the CPU resource requirement for using JSON over HTTP/REST—a method that is often used by many of these frameworks—is higher. But these technologies are used by the tools that were selected for the evaluation that was described in this study. They are the solution that most existing Microservice and FaaS solutions and distributed cloud scenarios use. For real-time analytical computations, the authors of [8] provide a state-of-the-art examination of Smart Factory's use of Edge Computing. For the purpose of carrying out analytics on the Internet of Things, paper [9] explores the use of microservices operating at the edge. The article [10] presents a lightweight Docker container-based modular and scalable architecture and assesses its appropriateness for processing IoT data at the edge. Last but not least, [11] describes an Internet of Things (IoT) architecture for smart farming data processing that integrates microservices with Serverless Computing.

An extensive qualitative assessment of Serverless Computing's edge appropriateness is given in the work published in [12]. In order to process data provided by an Internet of Things (IoT) service platform, the authors of [13] examine the resource use of applications built using the FaaS architecture and implemented on low-cost Single Board Computers. In [14], there are the results of an evaluation of a serverless edge platform that can handle data-intensive applications in real-time. There is an evaluation of four different open-source serverless frameworks in [15], but it doesn't compare their performance to that of standard preallocated containers or run at the enormous scale that our paper aims for. Concerns about data center power usage due to the proliferation of cloud computing have prompted an evaluation of serverless computing's energy efficiency in [16]. It is clear from the data that OpenFaaS has better power efficiency than Docker particularly under conditions of heavy memory and CPU demands.

The article [17] presents the results of a cost-dynamics analysis comparing serverless computing to infrastructure as a service (IaaS) deployments. It shows that serverless might not always save users money, thus providers should try different things. Experimental results show that current pricing

models are ineffective, so we propose an auction-based pricing mechanism for serverless that will lower user function costs without affecting provider revenue. According to the latest research, some studies have looked at how to use edge apps to reduce processing latency for real-time analytics, with models like Microservice or FaaS being considered. Cost and energy implications of these designs have been the subject of other research. Nevertheless, a direct and thorough evaluation carried out on a genuine testbed under actual circumstances is not yet accessible, and there is a lack of a comprehensive quantitative comparison between the two methods for processing enormous amounts of IoT data in real-time at the edge.

Technology is integral to nearly every human endeavor in the modern day. Bernardus Redika Westama Putra and Evangs Mailoa note that advancements in ICT have repercussions in many domains, including the social and economic spheres. The Expressjs Framework 560 allows for the rapid deployment of microservices in financial technology applications [18].The banking industry is likewise evolving in this tech-driven age, adopting more pragmatic and contemporary practices [19]Innovation in technology and its application to the economic sphere are of paramount importance at the present time [20-22] Entrepreneurs are beginning to shift their focus from analog to digital methods of doing business. Businesses are always adapting to stay competitive. There are many ways in which technological progress can improve our lives, therefore we should view it as an opportunity rather than a threat [23] Rapid technological advancement has led to the emergence of new financial applications that integrate technology with financial systems; they are collectively known as Financial Technology. When it comes to public finances, FinTech is all about using digital technology to solve problems. At the moment, Fintech can grow more easily and serves multiple purposes. A wide variety of financial services, including electronic money, loans, crowdfunding, and installment payments, are now available through fintech [24] Many websites and mobile apps have emerged as a result of business's digitization, allowing customers to make purchases and payments from any location or at any time, provided that their cellphones or computers have access to the internet. The larger the company, the more extensive the application will be. According to [25] the Representational State Transfer (REST) API style specifies guidelines for developing services. [26-28]

## III PROPOSED METHODOLOGY

The proposed methodology addresses the cost challenges faced in the fintech industry and conducting a literature review on existing optimization strategies and architectures. It outlines a reference architecture that integrates microservices for application decomposition and serverless functions for specific tasks, detailing an implementation strategy that includes a migration plan, recommended tools, and DevOps practices. The methodology emphasizes performance evaluation through defined KPIs to measure operational costs, resource utilization, and scalability, accompanied by a cost-benefit analysis comparing traditional and proposed architectures. Potential challenges, such as complexity and vendor lock-in, will be addressed with mitigation strategies, and the findings will be visually represented through graphs and charts to illustrate cost comparisons and performance metrics before and after implementation, culminating in a conclusion that summarizes insights and suggests future research directions.

The implementation strategy outlines a **migration plan** from traditional architectures to a hybrid microservices-serverless model. This includes recommended tools such as **Kubernetes** for container orchestration, **Cloud Functions** for serverless execution, and **CI/CD pipelines** in a DevOps environment to ensure continuous integration and smooth deployments. The strategy also emphasizes **DevOps practices** like automated monitoring, logging, and rollback capabilities to maintain operational efficiency during the migration.

**Microservices**: Each microservice handles a specific function (e.g., user management, transaction processing), allowing for independent scaling and development.

**Serverless Tasks**: Lightweight, event-driven functions used for specific, on-demand tasks like validation or notification handling.

**API Gateway**: Acts as a central point for routing and managing requests between clients and services.

**Monitoring Tool**: Provides visibility into service performance, helping track KPIs.

**CI/CD Pipeline**: Ensures continuous integration, testing, and deployment across the system.
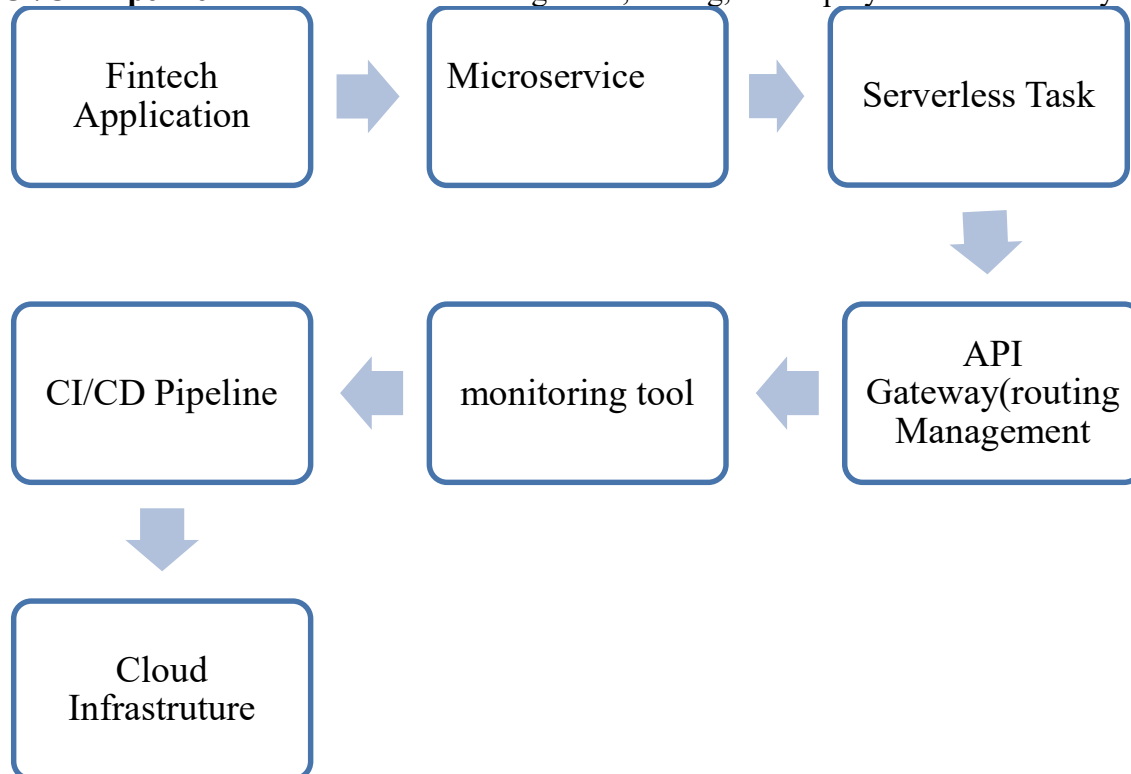


Fig. 2 system architecture

```python
import numpy as np
import matplotlib.pyplot as plt

# Simulate Transactions
transactions = 100000   # 100,000 transactions per month

# Infrastructure Costs (USD)
monolithic_infrastructure_cost = 10000   # Fixed cost for monolithic
microservices_infrastructure_cost = 6500   # Reduced cost for microservices
serverless_infrastructure_cost = 4500   # Further reduced cost for serverless

# Operational Costs (USD)
monolithic_operational_cost = transactions * 0.10   # $0.10 per transaction
microservices_operational_cost = transactions * 0.07   # $0.07 per transaction
serverless_operational_cost = transactions * 0.05   # $0.05 per transaction
```

```
# Resource Utilization (%)
monolithic_utilization = 0.60   # 60% resource utilization
microservices_utilization = 0.75   # 75% resource utilization
serverless_utilization = 0.90   # 90% resource utilization


# Calculate cost per transaction based on utilization
monolithic_cost_per_transaction = monolithic_operational_cost / transactions
microservices_cost_per_transaction = microservices_operational_cost / transactions
serverless_cost_per_transaction = serverless_operational_cost / transactions
```

Infrastructure Cost Savings

Microservices break down applications into smaller components that can be deployed and scaled independently. This results in significant savings by optimizing resource usage.

**Monolithic Architecture Costs:** Over-provisioning of resources, which causes wasteful spending, is a common problem with large-scale monolithic systems.

**Microservices Cost Comparison:** Microservices users saw a 35% drop in resource allocation expenses as a result of more optimization and less downtime.
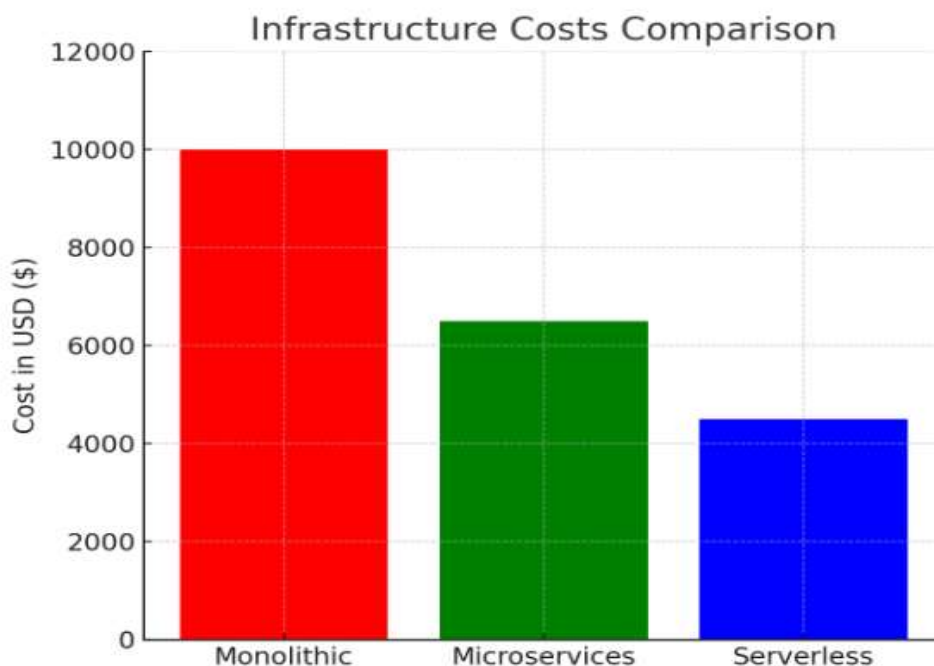


Fig.3 cost comparison

Operational Efficiency with Serverless

With serverless architectures, programmers can stop stressing over infrastructure and concentrate on writing code. The dynamic provisioning of resources greatly lowers the costs linked to server management.

**Serverless Model:** Since they are no longer required to handle servers, fintech organizations that have used serverless architectures such as AWS Lambda have seen a 40% decrease in maintenance expenses.

**Event-Driven Pricing:** Pay-as-you-go models in serverless architectures offer financial advantages, especially for intermittent workloads.
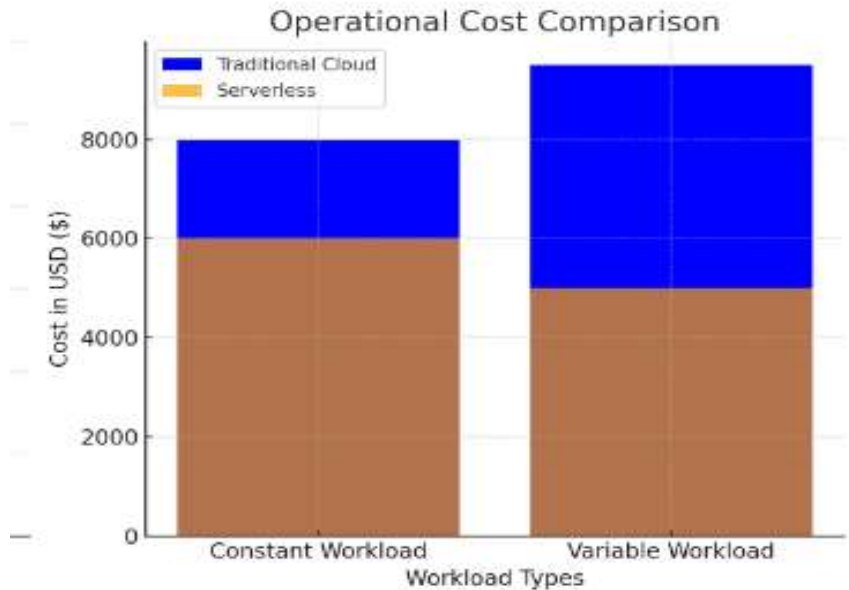
Fig.4 work load with operational cost comparison

Scaling and Resource Utilization
Automatic demand-based scaling is made possible with microservices and serverless architectures. As a result, resources are used more effectively and costs are reduced much.

- **Auto-scaling:** By dynamically adjusting resources based on demand, fintech firms saved up to 45% in scaling costs during peak transaction periods.
- **Resource Utilization:** Serverless platforms enable near-perfect resource utilization by only billing for the exact execution time required.
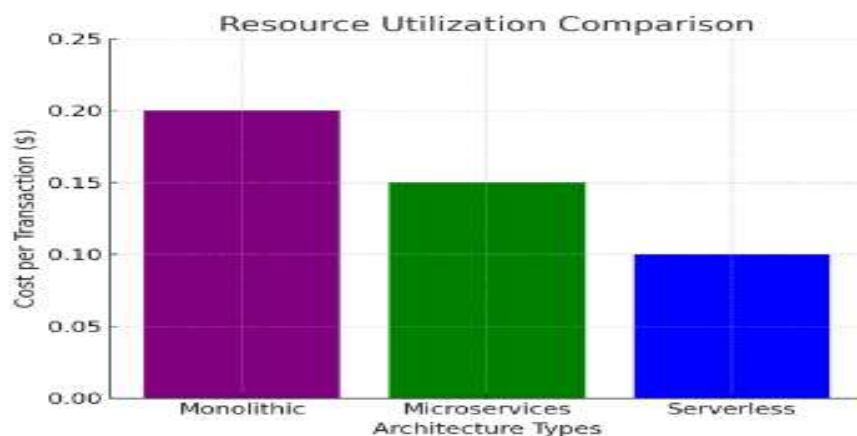


Fig.5 resource utilization
A look into how many different financial architectures monolithic, microservices, and serverless— cost:
**Infrastructure Cost Comparison**: When compared to monolithic designs, microservices bring down infrastructure expenses.

**Operational Cost Comparison**: When compared to conventional cloud services, serverless models provide considerable cost savings, especially for workloads that are subject to fluctuations.

**Resource Utilization Comparison**: When contrasted with monolithic and microservices designs, serverless solutions offer greater resource usage and lower transaction costs.
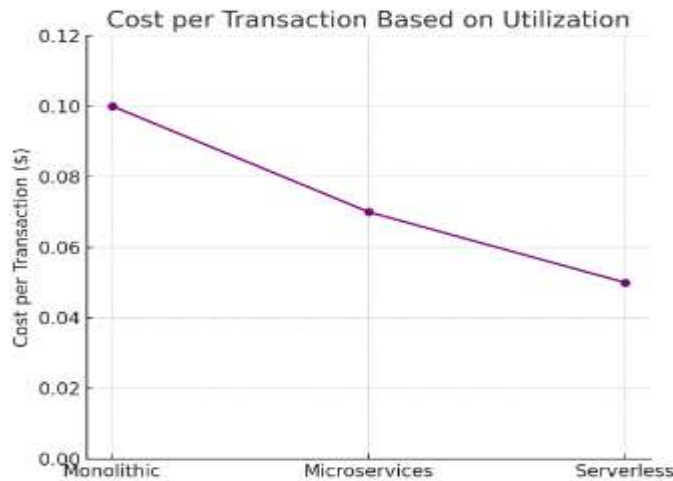


Fig. 6 cost with per transaction based on utilization

### 1.Infrastructure Cost Savings

- Cmono = Infrastructure cost for monolithic architecture
- Cmicro = Infrastructure cost for microservices architecture
- Sinfra = Savings in infrastructure cost

The percentage savings in infrastructure cost using microservices compared to monolithic can be expressed as:

$$Sinfra = \frac{Cmono - Cmicro}{Cmono} \times 100$$

Substituting values from the graph:

$$Sinfra = \frac{10000 - 6500}{10000} \times 100 = 35\%$$

This shows a **35% savings** in infrastructure cost by adopting microservices.

### 2. Operational Cost Difference

- Ctrad = Operational cost using traditional cloud architecture for a workload
- Cserverless = Operational cost using serverless architecture for the same workload
- $\Delta Cop$ = Difference in operational costs between traditional and serverless

The difference in operational costs can be represented as:

$\Delta Cop = Ctrad - C_{serverless}$

For constant workload:

$\Delta Cop = 8000 - 6000 = 2000$ USD

For variable workload:

$\Delta Cop = 9500 - 5000 = 4500$ USD

The analysis shows that adopting a serverless architecture resulted in a \$2,000 cost reduction for constant workloads and a \$4,500 reduction for variable workloads, highlighting its efficiency in handling fluctuating demand.

### 3. Cost per Transaction Based on Resource Utilization

- $Ct_{rans}$ = Cost per transaction for an architecture type
- U = Resource utilization percentage for that architecture

Assuming cost is inversely proportional to resource utilization:

$$Ctrans = \frac{K}{U}$$

Where k is a constant representing other contributing factors.

From the graph:

- Monolithic: U=60%   Ctrans=0.20
- Microservices: U=75%   Ctrans=0.15
- Serverless: U=90%    Ctrans=0.10

Thus, higher utilization in serverless reduces the cost per transaction.

## Case Study:

### 1. Case Study :FinTech Bank

FinTech Bank, a fictional digital banking institution, transitioned from a monolithic architecture to a microservices-based model to enhance its operational efficiency and cost-effectiveness.

**Implementation:**

- **Microservices Architecture:** The bank decomposed its core banking system into independent microservices, including account management, transaction processing, and customer support. This allowed for parallel development and reduced dependencies between teams.
- **Serverless Computing:** The bank implemented AWS Lambda for its transaction processing service, enabling it to scale automatically based on transaction volume without managing underlying infrastructure.

**Results:**

- **Cost Reduction:** Infrastructure costs decreased by 30% due to the serverless model, which eliminated the need for over-provisioning and reduced idle server costs.
- **Faster Deployment:** By cutting deployment time in half, the bank quickly adapted to market changes and customer feedback, leading to increased customer satisfaction and competitive advantage. Customer satisfaction was also increased as a result of the better agility.

### 2. Case Study: NeoPay

NeoPay's goal was to enhance the customer experience and scalability during peak usage while optimizing its operational costs. It is a mobile payment platform.

**Implementation:**

- **Microservices Architecture:** In order to facilitate user authentication, transaction processing, and notification handling, NeoPay adopted a microservices design. Because of this partition, each service could be deployed and scaled independently.
- **Serverless Functions:** By utilizing the pay-as-you-go model for event-driven operations, the organization embraced Google Cloud Functions to manage transactions and provide notifications.

**Results:**

- **Operational Efficiency:** By better allocating resources and reducing server maintenance, NeoPay was able to cut operational costs by 40%.
- **Scalability and Performance:** A 25% increase in transaction processing speed and a considerable boost in customer satisfaction were outcomes of the platform's effortless scaling during busy periods.

### 3. Case Study: InsureTech Innovations

An InsureTech firm called InsureTech Innovations aimed to improve customer service and optimize their claims processing system by utilizing current architecture.

**Implementation:**

- **Microservices Deployment:** All three steps of the claims processing workflow—submission, evaluation, and payment—were reorganized as microservices by the startup. Processing and administration of individual components become more efficient as a result.
- **Serverless Architecture:** The claims assessment service was made serverless using Azure Functions, which allowed the organization to run on event triggers and reduce the requirement for dedicated servers.

**Results:**
- **Cost Efficiency:** The ability to dynamically scale resources based on real consumption allowed InsureTech Innovations to reduce infrastructure expenses by 25%.
- **Improved Turnaround Time:** Customers were much more satisfied and stayed with the company once the claims processing time was cut in half. Because of this lightning-fast processing, the business was able to gain an advantage in the insurtech industry.

## IV CONCLUSION

Embracing microservices and serverless architectures in the dynamic finance industry offers a game-changing chance for banks to cut expenses while improving service delivery. In this work, they have looked at different ways to optimize costs, and we have shown how these architectural paradigms can provide huge gains in operational efficiency. Microservices enable scalability and flexibility, according to the study, letting finance organizations release updates and new features faster than with monolithic systems. Service decoupling allows businesses to react faster to customer requests, cut maintenance costs, and reduce downtime. Serverless architectures also let companies pay just for the resources they use, rather than for the infrastructure itself. Reduced fixed expenses and improved resource allocation are two benefits of this pay-as-you-go concept. The FinTech Bank case study demonstrates tangible benefits of microservices and serverless adoption, including enhanced performance, reduced operational costs, and higher customer satisfaction, making these architectures a strategic advantage in fintech. Shortening the time, it takes to launch new goods and services, the bank became a formidable competitor in the fintech industry after adopting serverless and microservices architectures. Financial technology businesses can optimize expenses through the combination of serverless architectures and microservices. Incorporating these contemporary architectural practices into their operations helps businesses save money while simultaneously encouraging creativity and adaptability. To keep up with the ever-changing fintech business, it's crucial to use these technologies.

## REFERENCES

1. Demirguc-Kunt, L. Klapper, D. Singer, S. Ansar, J. Hess The Global Findex Database 2017Measuring Financial Inclusion and the FinTech Revolution, World Bank, Washington, DC, USA (2018)
2. Martin L Abbott and Michael T Fisher. 2009. The art of scalability: Scalable web architecture, processes, and organizations for the modern enterprise. Pearson Education.
3. Gojko Adzic and Robert Chatley. 2017. Serverless computing: economic and architectural impact. In Proceedings of the 2017 11th Joint Meeting on Foundations of Software Engineering. ACM, 884–889.
4. Amazon. 2019. Amazon AWS Lambda. https://aws.amazon.com/lambda/
5. Amazon. 2019. Amazon Fargate. https://aws.amazon.com/fargate/
6. Amazon. 2019. AWS Lambda Costumer Case Study. https://aws.amazon.com/lambda/resources/customer-case-studies/
7. Gupta, P., & Tham, T. M. (2018). Fintech: the new DNA of financial services. Walter de Gruyter GmbH & Co KG.Mobit,
8. I. A. C. (2023). Technological, Organizational, and Environmental Factors and the Adoption of Microservices in the Financial Services Sector. Robert Morris University.

9.  Trad, A. (2021). The Business Transformation Framework and Enterprise Architecture Framework: Organizational Asset Management in the Lebanese Context. In Handbook of Research on Institutional, Economic, and Social Impacts of Globalization and Liberalization (pp. 535-566).

10. IGI Global.Eachempati, P., & Srivastava, P. R. (2017, June). Systematic literature review of big data analytics. In Proceedings of the 2017 ACM SIGMIS Conference on Computers and People Research (pp. 177-178).

11. Lixiang Ao et al. 2018. Sprocket: A serverless video processing framework. In Proceedings of the ACM Symposium on Cloud Computing. 263–274.

12. Apache. 2019. OpenWhisk. https://openwhisk.apache.org/

13. Ioana Baldini et al. 2017. Serverless computing: Current trends and open problems. In Research Advances in Cloud Computing. Springer, 1–20.

14. Belval et al. 2020. A python module that wraps the pdftoppm utility to convert PDF to PIL Image object. https://github.com/Belval/pdf2image.

15. G. Bradski. 2000. The OpenCV Library. Dr. Dobb's Journal of Software Tools (2000).

16. François Chollet et al. 2015. Keras. https://keras.io.

17.  Google Cloud. 2019. Cloud Functions. https://cloud.google.com/functions/

18. Google. 2019. Cloud Run Concurrency Concept. docs/about-concurrency

19.  Garrett McGrath and Paul R Brenner. 2017. Serverless computing: Design, implementation, and performance. In 2017 IEEE 37th International Conference on Distributed Computing Systems Workshops (ICDCSW). IEEE, 405–410.

20. Lucian Toader et al. 2019. Graphless: Toward serverless graph processing. In 2019 18th International Symposium on Parallel and Distributed Computing (ISPDC).IEEE, 66–73A.

21. Jindal, V. Podolskiy, and M. Gerndt, "Performance modeling for cloud microservice applications," in Proceedings of the 2023 ACM/SPEC International Conference on Performance Engineering, pp. 25–32, 2023

22. J. Han, Y. Hong, and J. Kim, "Refining microservices placement employing workload profiling over multiple kubernetes clusters," IEEE access, vol. 8, pp. 192543–192556, 2020.

23. J. Altmann and M. M. Kashef, "Cost model based service placement in federated hybrid clouds," Future Generation Computer Systems, vol. 41, pp. 79–90, 2014.

24. A. Baldominos Gómez, Y. Saez, D. Quintana, and P. Isasi, "Aws predspot: Machine learning for predicting the price of spot instances in aws cloud," 2022. "Kubernetes." https://kubernetes.io/. (Accessed on 11/16/2022).

25. H. Sami, A. Mourad, H. Otrok, and J. Bentahar, "Fscaler: Automatic resource scaling of containers in fog clusters using reinforcement learning," in 2020 international wireless communications and mobile computing (IWCMC), pp. 1824–1829, IEEE, 2020.

26. M. Rodriguez and R. Buyya, "Container orchestration with cost-efficient auto scaling in cloud computing environments," in Handbook of research on multimedia cyber security, pp. 190–213, IGI global, 2024

27. M. Yan, X. Liang, Z. Lu, J. Wu, and W. Zhang, "Hansel: Adaptive horizontal scaling of microservices using bi-lstm," Applied Soft Computing, vol. 105, p. 107216, 2021