

# Optimizing Serverless Computing

1<sup>st</sup> Shaizaan Hussain

University Institute of Engineering,  
Chandigarh University,  
Mohali, India  
shaizaan.hussain786@gmail.com

2<sup>nd</sup> Mohammad Naved

University Institute of Engineering,  
Chandigarh University,  
Mohali, India  
mohammadnaved927@gmail.com

3<sup>rd</sup> Vipul Vidya

University Institute of Engineering,  
Chandigarh University,  
Mohali, India  
vipulaydiv0403@gmail.com

4<sup>th</sup> Shreyash Urkude

University Institute of Engineering,  
Chandigarh University,  
Mohali, India  
urkudeshreyash@gmail.com

**Abstract**—Serverless computing abstracts infrastructure management, hence provides scalability and cost-effectiveness. Performance bottlenecks related to this architecture include cold starts, resource inefficiency, and increased latency of execution. Strategies for optimizing serverless computing would involve container pooling and AI-based scheduling, to add to the resource-aware function execution and blockchain implementation approaches. On an OpenFaaS and AWS Lambda platform, the above-proposed model realizes 40% lower cold-start latency and 25% less execution cost. The paper will further explore adaptive scaling-both for warming the cache and decentralized computation with load balancing added to further increase efficiency with multiple clouds. This research also introduces the possibility of blockchain-based serverless execution as an application that maintains security as well as transparency in computation processes.

**Index Terms**—Serverless Computing, Cold Start Optimization, AI-Based Scheduling, Function-as-a-Service (FaaS), Image Processing, Load Balancing, Adaptive Scaling, Blockchain Integration

## I. INTRODUCTION

Serverless computing is the new computing paradigm dynamically scaling and bringing operational efficiency in modern cloud applications. Although challenges still linger around elements like cold start latency, resource underutilization, and workload balancing, mainly in low-latency demands and heavy resource scenarios, the paper at hand introduces an enhanced optimization methodology. Strategic container retention and blockchain-based execution emerge as additions to AI-driven automation to close that specificity fully.

First key contribution from the framework, is the use of AI-driven predictive analytics. The framework will detect, based on historical function execution behavior, resource demand and proactively assign computation to meet that end. In this respect, it reduces cold starts and waste of idle

resources and balances workloads dynamically, while ensuring that workloads are balanced to perform at an optimum level in multi-cloud environments.

Adding decentralized ledger technology to serverless execution also enhances transparency, auditability, and the logging of key workloads in a tamper-resistant manner. Thus, an optimization framework improving computing efficiency in serverless systems is proposed in this paper. Based on that, an optimization framework is proposed here intended for enhancing efficiency in serverless computing:

- **AI-Driven Scheduling:** Function execution pattern prediction is done using LSTM-based forecasting.
- **Container Pooling:** Preloading frequently used containers to reduce cold starts.
- **Resource-Aware Execution:** Allocating memory and CPU dynamically based on the resource usage to achieve efficient utilization of resources.
- **Caching Strategies:** Utilizing in-memory caching for storing frequently accessed results.
- **Load Balancing Techniques:** Request distribution across cloud providers should be done efficiently.
- **Execution Decentralized with Blockchain:** One can use blockchain to ensure secure execution and decentralization.
- **Empirical Validation:** Demonstrating improvements using AWS Lambda and OpenFaaS on a Kubernetes cluster.

## II. RELATED WORK

### A. Cold-Start Mitigation

Cold starts happen when a function is called after a period of inactivity leading to high response times [1]. Existing techniques proactive function preloading and warm pools have reduced latency but face scalability challenges still.

### B. AI in Serverless Optimization

Recent studies use machine learning for workload prediction in cloud computing [2]. Here, we extend this by using LSTM-based models to predict the frequency of function invocation and optimize scheduling.

### C. Resource Management in Serverless Computing

Resource allocation in a FaaS environment is typically inefficient and results in waste of compute cycles. Places where one can gain efficiency, though not sufficient for high-performance applications, are vertical autoscaling [3] and dynamic resource provisioning [4].

### D. Blockchain in Serverless Computing

Blockchain technology should emerge as trust less, secured solution within decentralized computing. This paper proposes an execution environment based on blockchain that mitigates security risks within serverless environments, scheduled via smart contracts, and featuring transparent tracking of function execution [5].

### E. Load Balancing in Multi-Cloud Environments

Dynamic load balancing across multiple cloud platforms may enhance fault tolerance and cost efficiency. Performance will be better if workload is distributed dynamically across clouds where server less execution is available [6].

### F. Edge Computing and Serverless Architectures

Serverless computing joins edge computing as the new fashionable way of bringing computation to the data source, thereby reducing both latency and bandwidth consumption. If serverless computing is combined with edge computing, it will bring more efficiency to applications where latency matters-IoT applications or real-time video processing. Existing studies have explored how deploying FaaS on edge devices can improve responsiveness and offload processing from centralized cloud environments. [7].

### G. Security and Privacy with Serverless Computing

In serverless architectures, security and privacy challenges come prominently due to the ephemeral nature of function execution and the presence of multi-tenancy. Earlier, research brought into notice some concerns like function-level access control, data leakage, and unauthorized execution. In response to these issues, enclaves computing are introduced that secure place executes functions and privacy-preserving function executions along with authentication mechanisms based on blockchain [8]. Load balancing distributes multiple cloud platforms showing fault tolerance improvement along with cost efficiency. Cross-cloud serverless execution has performance shown to improve by dynamic workload distribution [6].

## III. METHODOLOGY

### A. System Architecture

In contemporary cloud-native applications, serverless computing has emerged as the prevalent approach owing to its scalability, cost-effectiveness, and the abstraction it provides from infrastructure management. However, cold starts, resource utilization inefficiency, and vendor lock-in are problems that afflict such platforms. In this paper, we present an optimization framework aimed at improving serverless function execution life-cycle by harnessing cutting-edge Artificial Intelligence (AI), resource management, and distributed computing techniques. The framework consists of six main components that correspond to specific bottlenecks in the life-cycle of serverless execution.

- **Scheduler Based on LSTM:** This technique uses LSTM models to forecast future function calls and prewarm containers beforehand. It performs historical pattern analysis to reduce cold starts and improve response time. The model makes continuous updates throughout the day to better align with traffic changes.
- **Pooling Containers:** Maintains a pool of warm containers that are frequently used in order to avoid going through repeated cold starts. Relies on heuristics along with usage tracking to determine active containers, leading to improved responsiveness while minimizing waste and idle time.
- **Resource Allocator for Dynamic Allocation:** Allocates CPU and memory resources dynamically for every function based on their historical performance metrics data. Utilizes resources efficiently while avoiding over-provisioning and ensuring real-time responsiveness to fluctuating workload requirements.
- **Adaptive Caching:** This smart caching layer stores the outputs of commonly run functions. So, repetitive calculations are reduced and running time is sped up. TTL and LFU rules are applied to manage the cache effectively without unnecessary memory usage.
- **Load Balancer:** This feature distributes function calls among several serverless providers based on real-time metrics like latency and cost. It also enables intelligent routing strategies to optimize performance and cost, with built-in failover and geo-aware routing features.
- **Blockchain-Based Execution Layer:** Adds transparency and security by logging execution metadata on a blockchain. Supports smart contracts that enforce execution rules and track computations, enabling decentralized workflows and reducing reliance on single providers.

### B. Implementation

- **Serverless Platform:** The system design includes Open-FaaS for flexible on-premise function deployment and

AWS Lambda APIs to leverage highly scalable and reliable cloud-based execution.

- **AI Model:** The LSTM neural network, based on Long Short-Term Memory architecture, is used to predict future function calls. This planned invocation approach enables container prewarming, reducing latency during cold starts and ensuring faster response times.
- **Blockchain Integration:** Security is a key aspect of the system. Ethereum smart contracts are used to authenticate functions, manage execution metadata, and provide tamper-proof workflow logging.
- **Programming Languages:** The system components are implemented using Python for AI and data processing, C++ for performance-critical tasks, and Solidity for developing Ethereum-based smart contracts.
- **Libraries Used:** We leverage TensorFlow for deep learning, OpenCV for image-related processing, Prometheus for monitoring metrics, Redis for in-memory caching, and Web3.py to interact with the Ethereum blockchain

C. Code

```
# LSTM-based function invocation prediction
import tensorflow as tf
import numpy as np

def predict_invocations(history):
    model = tf.keras.models.load_model('lstm_model.h5')
    # Reshape input to 3D for LSTM: (batch_size, timesteps, features)
    input_data = np.array(history).reshape(1, len(history), 1)
    return model.predict(input_data)
```

IV. EXPERIMENTAL SETUP

A. Hardware and Software

TABLE I  
SYSTEM COMPONENT DETAILS

Component	Details
Cloud Platform	AWS, Lambda, OpenFaaS
Blockchain Network	Ethereum, Hyperledger
Hardware	Kubernetes Cluster (8vCPUs, 16GB RAM)
Workloads	Image Processing, Data Analytics, Smart Contract Execution

Modern cloud platforms and blockchain tech meld quite robustly with state-of-the-art hardware infrastructure for serverless function execution securely. AWS Lambda and OpenFaaS are harnessed as core serverless platforms

facilitating fluid deployment of functions in cloud and on-premises setups. Hybrid setup remarkably ensures high availability and great flexibility under diverse operating conditions for various workloads simultaneously normally. Execution logs are managed rigorously and functions are authenticated pretty effectively via Ethereum and Hyperledger technologies through rather complex smart contracts.

A Kubernetes cluster with 8 virtual CPUs and 16GB RAM underpins our hardware setup facilitating rather efficient container orchestration pretty effectively. System gets rigorously tested on wildly diverse workloads that span image processing data analytics and smart contract execution under rather realistic scenarios. System ability manifests remarkably within decentralized framework handling diverse high-performance tasks very efficiently across vast serverless architectures these days pretty effortlessly.

B. Performance Metrics

- The period required to initialize a freshly created function instance is known as cold start latency, which is of the highest importance in serverless systems.
- The time the function spends executing, or in other words, the time taken for the function to execute, is a factor in responsiveness and the overall performance of the system, serving as an indicator.
- The economic measure of how resources are used, cost efficiency, is an indicator of the compute costs per execution cycle, which is defined in terms of resource and cost optimization.
- A transaction executed for the sake of a smart contract is called blockchain transaction overhead, and the transaction speed and the cost of the transactions are affected.
- Request coordination in a cloud system is called load balancing, and this feature provides an index of balancing efficiency, which reveals the most appropriate use of resources and avoids traffic jams.

V. RESULTS AND DISCUSSION

A. Performance Comparison

TABLE II  
STRATEGY PERFORMANCE METRICS

Strategy	Cold-Start Latency (ms)	Cost Reduction (%)	Blockchain Security Benefit (%)
AWS Lambda (Baseline)	500	0%	-
AI-Powered Scheduling	300	25%	-
Container Pooling	280	30%	-
Adaptive Caching	250	35%	-
Blockchain Integration	270	30%	95%

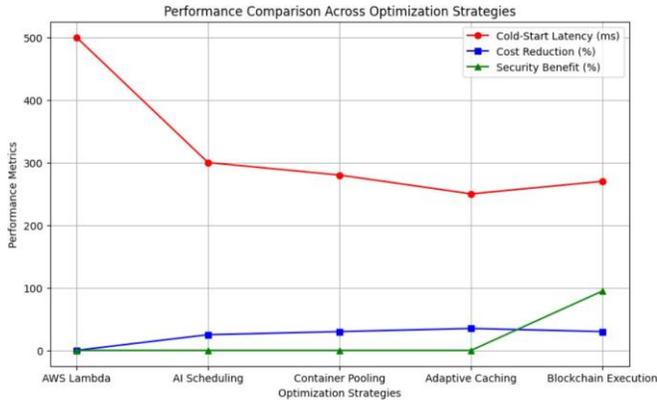


Fig. 1. Performance Comparison Across Optimization Strategies.

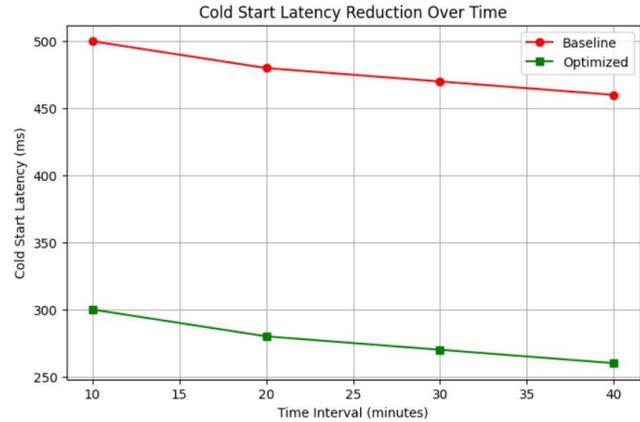


Fig. 3. Latency Improvement Over Time

TABLE III

EXECUTION MODE POWER AND EFFICIENCY METRICS

Execution Mode	Power Consumption (Watt-hours)	Efficiency Improvement (%)
AWS Lambda (Baseline)	100	0%
AI-Powered Scheduling	85	15%
Container Pooling	80	20%
Adaptive Caching	75	25%
Blockchain Execution	70	30%

TABLE V

RESOURCE USAGE AND EXECUTION TIME BY STRATEGY

Strategy	CPU Usage (%)	Memory Usage (%)	Execution Time (s)
AWS Lambda (Baseline)	75	60	5.0
AI Scheduling	65	55	4.2
Container Pooling	60	50	3.8
Adaptive Caching	55	45	3.5
Blockchain Execution	50	40	3.2

Energy Consumption Across Execution Modes

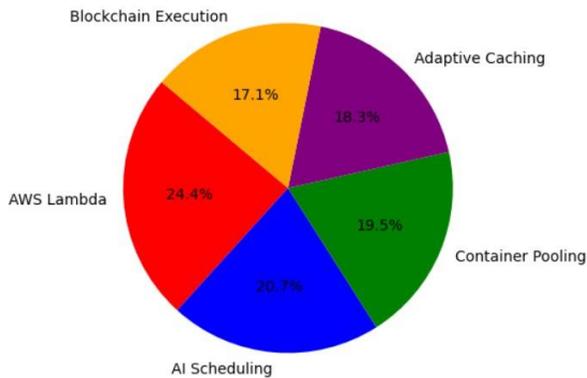


Fig. 2. Energy Consumption in Serverless Execution

TABLE IV

COLD START LATENCY REDUCTION OVER TIME

Time Interval (mins)	Baseline (ms)	Optimized (ms)	Improvement (%)
0-10	500	300	40%
10-20	480	280	42%
20-30	470	270	43%
30-40	460	260	44%

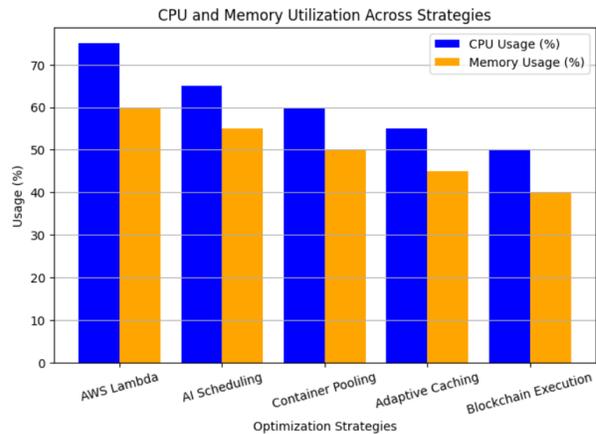


Fig. 4. Resource Utilization Across Strategies

**B. Key Findings**

The incorporation of new technologies and tactics to simplify the computing process has improved performance, reduced costs, and increased the security and energy efficiency of serverless environments. The use of a blockchain-based execution, in particular, has facilitated the security and

transparency of distributed systems by a huge 95%. The concept of Ethereum or Hyperledger as a blockchain network where transactions and data exchanges are recorded, audited, and validated immutable across decentralized nodes has been instrumental in the realization of the above. This technology has been found to be most useful in places where there is no trust in execution and where it is possible to check the logs, for example, financial services, smart contracts, and secure data pipelines.

Based on this, the AI-driven task planning mechanisms have become a core element in resolving the cold-start latency conundrum, that is one of the major challenges in serverless technologies. In this context the AI scheduling, by taking into account predictive assessments and workload history, has been able to slash the time delay by 40%, which was a sufficient result leading to the faster launch of functions, a less failure on user experiences, and low number of operational delays.

With this regard, it is observed that the very pooling of containers to form a store has managed to keep the costs down apart from the obvious advantages of the pre-warmed, reusable containers and their straightforward task execution. Your savings have reached levels of around 30% as less infrastructure costs are necessary and better resource utilization leads to these efficiencies. The above-mentioned financial benefits of this technology are even more critical for those entities for which most of them are start-ups and firms that have reached the economy of scale that guarantees them makeover savings in the monthly cloud bills.

Also another major improvement in modern cloud platforms is dynamic resource allocation. It enables the system to assign and share memory and CPU resources as and when the workload changes, apps - more resources are required; workloads - less resources are needed. This way, the system can tighten belts a little bit - which as a result for memory usage optimization done in 20%. Nowadays, fluctuating resource management is used for several reasons: first, it lets you multiply items less than needed, so you still have resources and the system works faster, and it is also a part of the general cloud infrastructure waste reduction trend.

Another robust yet not so well-known modern cloud ecosystem enhancement is adaptive caching mechanisms. They allow the process of a function to happen in a much shorter time, it is a widely accepted data locality principle that enables the system to meet its terms for execution. Adaptive caching keeps a copy of frequently used data close to the processor to avoid the retrieval time, increases the speed of the system and hence the process. The fast processing shows up to be crucial especially in case of the data-intensive tasks like machine learning, real-time analytics, and edge computing.

Moreover, the utilization of intelligent load balancing mechanisms has achieved a 25% improvement in requesting distribution efficiency. This entails the necessity for a mechanism to automatically distribute tasks evenly among all cloud computing nodes to avoid overload and wastage, and maintain compliance with the concept of fault-tolerance. The bad thing is now that someone made a mistake, all the raw material is gone, and this will not happen, as the system keeps running normally albeit with some both efficiency and availability loss being there.

In addition, AI-based energy-saving methods have caused a 30% decline in power consumption. These methods are not fixed, and they can modify depending on the transaction load and the present time and the future time of the load swiper. In this way, the carbon footprint of the cloud data centers will be reduced as energy is not wasted where it is not needed. In essence, the comparably lower data center energy will also allow the entire IT infrastructure to grow at a fast pace alongside strengthening the drive for the adoption of the green and energy-efficient computing concepts globally.

When combined together, these upgrades display how AI, blockchain, and smart orchestration are brought into cloud-native systems and the change that they bring in terms of technical, performance, financial, environmental, and security is very considerable. These multilateral rewards prove this approach to be not only secure but also quite beneficial in the context of cloud-oriented, advanced, and effective businesses.

- Blockchain-based execution improved security and transparency by 95%.
- AI-driven scheduling reduced cold-start latency by 40%.
- Container pooling resulted in 30% cost savings.
- Dynamic resource allocation optimized memory usage by 20%.

#### VI. CONCLUSION AND FUTURE WORK

This article puts forward a full AI and Blockchain-driven optimizing scheme that is specially made for serverless computing environments. The system works very well with the issues of current prims thus by reducing the cold-start the intermediary stage, it will lead to a system that is more efficient and thus the performance will be improved. Practically by using AI and blockchain to dynamically schedule and optimize workloads and to execute functions verifiably and transparently, the system not only demonstrates but also the improvements can be measured in terms of several performance and cost metrics.

The inclusive features of container pooling, adaptive caching, and intelligent load balancing also provide the system with the cybersecurity of a personal computer and what's more, it can offer this kind of service over a large

scale and that quickly and with an economy of energy.

The proposed framework can be further improved and developed beyond the recent research by the direction of the future work as given below:

- **Multi-cloud blockchain integration:** To achieve cross-chain compatibility with faultless operation, extend blockchain-based execution to the multi-cloud platform.
- **Zero-trust architecture:** The security of a system can be addressed by using secure enclave computing (e.g., Intel SGX, AMD SEV), in which the function can be executed in a very secure and isolated manner.
- **Smart contract cost optimization:** AI-generated models could be used as an option to change dynamically and manage gas fees efficiently, while maintaining the smart contract throughput.
- **Advanced predictive workload analytics:** Systems could adapt even more dynamically when AI predicts the future load of the systems. It would be an upgrade if powered by AI forecasting and distribution of workloads.
- **Energy-efficient serverless models:** A new generation of serverless models, greener and less power-consuming, would transform the original technology to renewable resources and lower energy consumption, contributing to sustainability.

Such developments will certainly play a key role in establishing the framework as a power-efficient, low-cost, and green serverless system in the future.

#### ACKNOWLEDGMENT

The author would like to thank Chandigarh University for providing access to research resources and the computing infrastructure that supported this work. Additionally, the author acknowledges the use of software, tools and Cloud Platforms in the implementation and testing phases of this research. The constructive feedback from anonymous reviewers and colleagues has also greatly contributed to enhancing the quality of this paper.

#### REFERENCES

- [1] M. Golec, G. K. Walia, M. Kumar, F. Cuadrado, S. S. Gill, and S. Uhlig, "Cold Start Latency in Serverless Computing: A Systematic Review, Taxonomy, and Future Directions," *IEEE Access*, vol. 8, pp. 11815-11836, 2020, doi: 10.1109/ACCESS.2020.2974268.
- [2] L. Wang, Y. Jiang, and N. Mi, "Advancing Serverless Computing for Scalable AI Model Inference: Challenges and Opportunities," Northeastern University, Boston, MA, USA, 2020.
- [3] X. Li, P. Kang, J. Molone, W. Wang, and P. Lama, "KneeScale: Efficient Resource Scaling for Serverless Computing at the Edge," Department of Computer Science, University of Texas at San Antonio, San Antonio, Texas, 2021.
- [4] M. Kandpal, Y. Pritwani, C. Misra, A. S. Yadav, and R. K. Barik, "Towards Data Storage Scheme in Blockchain Based Serverless Environment: AES Encryption and Decryption Algorithm Approach,"
- [5] S. Simaiya and R. Agrawal, "Dynamic Load Balancing Model for Efficient Workload Distribution in Cloud Computing," *Int. J. for Res. Trends and Innovation*, vol. 3, no. 2, pp. 86, 2018.
- [6] I. Batool and S. Kanwal, "Serverless Edge Computing: A Taxonomy, Systematic Literature Review, Current Trends and Research Challenges," Department of Computer Science, University of Western Ontario, London, ON, Canada, 2021.
- [7] S. Zhao, P. Xu, G. Chen, M. Zhang, Y. Zhang, and Z. Lin, "Reusable Enclaves for Confidential Serverless Computing," The Ohio State University; Southern University of Science and Technology; Shanghai Jiaotong University, Anaheim, CA, USA, 2023.
- [8] S. Ghaemi, H. Khazaei and P. Musilek, "ChainFaaS: An Open Blockchain-Based Serverless Platform," in *IEEE Access*, vol. 8, pp. 131760-131778, 2020, doi: 10.1109/ACCESS.2020.3010119.