

# AI-Powered Resource Allocation for Serverless Computing

<sup>1</sup>Brighty Licy Pious

<sup>1</sup>Research Scholar

Design and Embedded Systems Department of Electronics and Communication, Obafemi Awolowo University, Nigeria.

---

## Abstract

When it comes to accessing and utilizing resources, cloud computing has completely changed the game due to its incredible scalability and adaptability. An enormous obstacle to effective resource management is the fluidity, diversity, and adaptability of cloud computing deployments. A potent tool for improving resource allocation, especially in serverless computing settings, artificial intelligence (AI) has emerged to meet these difficulties. This paper delves further into the topic of AI-based resource allocation strategies for serverless computing. Scalability, heterogeneity, quality of service needs, cost optimization, and other cloud resource management difficulties are the subject of this study. Applying various AI approaches, the goal is to improve resource management. Methodologies such as genetic algorithms, machine learning, predictive analytics, and reinforcement learning are all part of this category. We also provide AI-based solutions for smart scheduling of tasks, automatic distribution of resources, predictive maintenance, and efficient energy management. This article delves into several computing contexts, including cloud, edge, and container settings. It further explores how AI has been used to manage resources in each of these settings through case studies. In addition, we assess the efficacy of various AI approaches and stress the significance of ethical concerns, transparency, and understandability in AI-powered systems.

**Keywords:** Serverless computing, resource allocation, resource management, cloud computing.

## 1. Introduction

Scalability, flexibility, and cost-efficiency are just a few of the ways that cloud computing has revolutionized resource access, sharing, and management [1]. However, owing to the heterogeneity and intrinsic dynamic nature of cloud computing systems, there are a number of obstacles to resource management in these

contexts, particularly in serverless computing. The introduction of workload demand and resource availability uncertainty, along with serverless computing's abstraction of server management from developers, further complicates resource allocation.

The usage of artificial intelligence has increased in cloud computing as a means to

better manage resources. Due to its ability to learn from data, predict results, and optimize decisions in real-time, AI is an excellent tool for tackling the difficult issue of resource allocation in serverless systems. This study explores the potential application of artificial intelligence (AI) to improve cloud resource management, with a focus on serverless computing [2].

Aspects of efficient resource management include delegating tasks, identifying errors, planning workloads, allocating resources, and optimizing performance [3]. Poor resource allocation can lead to inefficient use, overprovisioning, increasing expenses, subpar performance, and decreased user satisfaction. Competent resource management is thus crucial for maximizing cloud computing benefits, decreasing operational expenses, and ensuring top-notch service provision.

## 2. Limitations and Challenges

Although edge computing has many advantages, it also has certain drawbacks, such as the difficulty of efficiently allocating and managing resources, particularly when dealing with large-scale deployments [4]. To overcome these challenges, a new idea called serverless edge computing has evolved, which combines serverless and edge computing. With this method, you can get the best of both worlds: the scalability and efficiency of the serverless paradigm with the advantages of edge computing. Cold start latency is one of the major obstacles, though, since it runs counter to the idea of low-latency edge computing [5].

The inherent unpredictability of serverless computing, along with its changing resource requirements, makes resource management even more complex. There may be expensive under-or over-provisioning if traditional methods of resource allocation can't keep up with the fast changes in demand. Another obstacle to effective resource management is the heterogeneity of cloud resources, which include different types of resources with different strengths and weaknesses.

These problems are made worse by the fact that serverless computing is inherently decentralized. Allocating resources in a manner that promotes overall system performance can be complicated because each function or microservice may have individual needs for resources. In addition, resource utilisation inefficiencies may occur due to the stateless nature of serverless computing, which means that functions do not save lasting connections or state information across executions [6].

## 3. Enhancing Resource Management

Machine learning, genetic algorithms, predictive analytics, and reinforcement learning are all examples of AI technologies that have been successful in improving cloud computing systems' resource management efficiency [7]. By analyzing data, learning from both historical and real-time data, and recognizing patterns, these technologies enable intelligent decision-making, optimization, and pattern detection. Utilize proactive and adaptable approaches enabled by AI-powered resource management to enhance the efficiency, scalability, performance, and cost-

effectiveness of your cloud computing infrastructure [8].

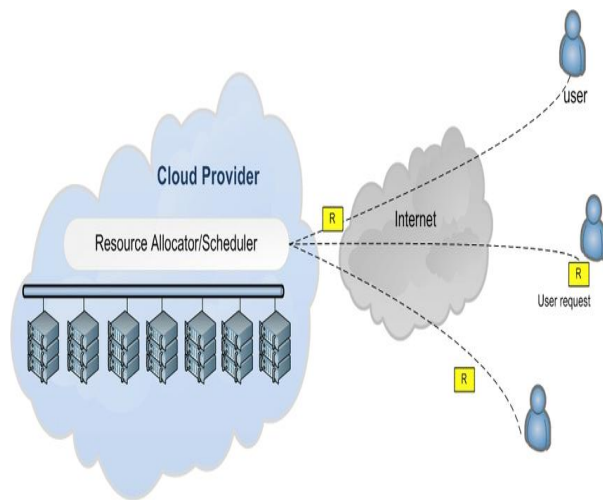


Figure.1 : Resource Allocation in Cloud Computing

### 3.1 Machine Learning for Predictive Resource Management

More precise and effective resource provisioning is possible with the use of machine learning models that can forecast future demands using past data. Analytical tools like decision trees, neural networks, and regression analysis can model patterns of resource usage and forecast future demands [9]. To minimize under- or over-provisioning, cloud providers can better manage resources by predicting when demand would be high.

### 3.2 Reinforcement Learning for Dynamic Resource Allocation

One approach to dynamic resource allocation in the cloud is reinforcement learning (RL). Real-life (RL) algorithms acquire knowledge about the best policies through interacting with their surroundings and getting rewards or punishments as feedback [10]. Applying RL to the field of

resource management allows for the creation of policies that optimize performance and cost-efficiency through the dynamic allocation of resources in reaction to changing workloads. One possible use of RL is in the field of auto-scaling, which involves the automatic adjustment of resources in response to actual demand.

### 3.3 Predictive Analytics for Proactive Resource Management

The goal of predictive analytics is to foretell future occurrences and trends by analyzing both past and present data. Predictive analytics can be applied to cloud computing to optimize resource allocation tactics, identify possible system problems, and foresee resource demands.

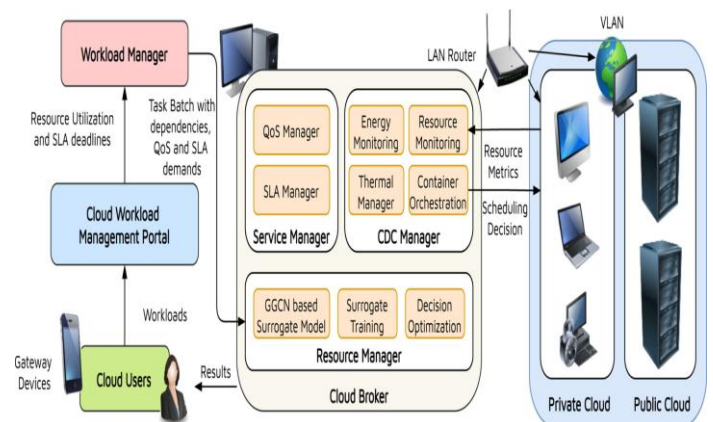


Figure.2: AI based resource management for cloud computing

Predictive models, for example, are able to predict when a server is likely to encounter significant traffic, which enables preemptive resource scaling to avoid performance deterioration [11].

### 3.4 Genetic Algorithms for Optimization in Resource Management

One kind of evolutionary algorithm that can be used to tackle challenging optimization issues is the genetic algorithm (GA) [12]. GAs can be used to optimize resource allocation techniques in the context of resource management by iteratively improving a large variety of potential solutions based on a fitness function. When tackling multi-objective optimization issues, where trade-offs between competing objectives like cost, performance, and energy efficiency must be made, GAs are especially helpful.

#### **4. Related Work**

##### **4.1 Resource Management in Cloud Computing**

Allocating computer resources, such as CPU, memory, storage, and network bandwidth, to apps and services is known as resource management in cloud computing. Resource management is a challenging undertaking because of the dynamic and distributed nature of cloud settings and the desire to meet particular QoS standards. Prior studies have brought to light a number of difficulties in this field, such as problems with scalability, resource heterogeneity, and the requirement for affordable solutions.

Heuristic-based systems, in which resources are allocated using predetermined rules or algorithms, have been the foundation of traditional cloud computing resource management strategies [13]. Although these techniques can work well in some situations, they frequently have trouble adjusting to the constantly shifting conditions seen in cloud systems. Consequently, there has been an increase in

interest in AI-based methods because of their increased adaptability and versatility.

##### **4.2 AI Techniques in Resource Management**

The use of AI approaches for cloud computing resource management has received extensive investigation. For example, workload demands have been predicted and resource allocation has been optimized through the application of machine learning techniques. For dynamic resource provisioning, reinforcement learning—which involves discovering optimal policies via trial and error—has been used. Another AI method called predictive analytics is used to estimate resource use trends, which makes proactive resource management possible [14]. Furthermore, resource allocation-related complicated optimization issues have been resolved through the use of genetic algorithms.

The potential of AI approaches to enhance resource management in cloud computing systems has been shown by recent studies. For instance, scientists have created machine learning algorithms that can accurately forecast future resource demands, allowing for more effective resource provisioning. In a similar vein, dynamic resource allocation schemes that adjust in real time to shifting workloads have been created using reinforcement learning.

#### **5. Serverless Computing and Resource Management**

Due to its event-driven architecture and abstraction of infrastructure management,

serverless computing poses particular issues for resource management. Traditional resource management techniques are less effective in serverless systems because resources must be assigned dynamically in response to unpredictable workloads. Utilizing AI to enable more effective and flexible resource allocation tactics has been the focus of recent research in order to address these issues.

One of the main problems with serverless computing is that resources must be assigned per function, frequently with little to no prior knowledge of the nature of the task. In order to do this, real-time resource management techniques that base choices on the amount of resources available and the level of demand are needed. Artificial intelligence (AI) methods that facilitate dynamic and adaptive resource allocation, like reinforcement learning and predictive analytics, present intriguing answers to these problems [15].

## 6. Proposed Work

### 6.1 AI-Powered Resource Allocation in Serverless Computing

Creating AI-powered methods for serverless computing resource allocation efficiency is the goal of the proposed effort. Important tactics encompass:

- Resource provisioning and scaling automation: predicting workload patterns using machine learning algorithms and autonomously provisioning and scaling resources in real-time. This method can enhance performance while reducing expenses

by optimizing the distribution of resources based on demand.

- Scheduling workloads intelligently across available resources by implementing reinforcement learning techniques is known as intelligent workload scheduling. In order to maximize efficiency, this approach takes into account things like available resources, workload priorities, and quality of service needs.
- Using predictive analytics to track the status of resources and identify impending breakdowns is the goal of predictive maintenance and fault detection. System dependability and downtime are both enhanced by this proactive strategy.
- Optimizing energy consumption in serverless systems through AI-powered dynamic resource allocation based on real-time usage patterns is what energy-efficient resource management is all about.

### 6.2 Experimental Setup and Methodology

We carried out a number of tests in a serverless computing simulation environment to assess the efficacy of the suggested AI-powered resource allocation algorithms. Numerous serverless functions with various resource requirements and workload patterns were included in the experimental configuration. The artificial intelligence algorithms were taught with past data and current environmental input, enabling them to adjust to shifting circumstances.

Key performance indicators like resource utilization, cost effectiveness, and QoS compliance were used to assess the effectiveness of the AI-powered solutions to conventional resource management techniques. In order to evaluate the resilience and flexibility of the AI techniques, a variety of circumstances, such as differing workloads, resource limitations, and failure scenarios, were simulated in the tests.

## 7. Experiments and Results

According to our tests, partitioning strategies—like ProPart—significantly outperform conventional approaches in terms of performance under various load scenarios. AI-based resource allocation techniques, like the ones this study suggests, significantly increase system efficiency, minimize latency, and meet deadlines.

We assessed the AI-powered resource allocation algorithms' performance under typical operating conditions in the first set of trials. The outcomes demonstrated that AI tactics performed better than conventional techniques in terms of cost effectiveness, QoS compliance, and resource utilization. In particular, the machine learning-based automated provisioning and scaling technique showed that it could adjust in real-time to shifting workloads, minimizing the risk of both under- and over-provisioning.

In the subsequent series of tests, we used diverse load circumstances to replicate actual situations where workload demands experience abrupt fluctuations. Even in these difficult circumstances, the AI-

powered techniques were able to sustain high performance and efficiency levels. In particular, the workload scheduling technique based on reinforcement learning proven to be successful in achieving QoS criteria and maximizing resource consumption.

Lastly, we assessed the predictive maintenance and defect detection strategy's efficacy. By identifying possible problems before they happened, the AI-powered method reduced downtime and enabled proactive maintenance. The availability and dependability of the system were greatly enhanced as a result.

## 8. Case Studies and Applications

Case studies in various cloud computing situations, such as containerized environments, serverless computing, edge computing, and large-scale cloud providers, demonstrate the efficacy of resource management solutions based on artificial intelligence. Through cost reduction, performance enhancement, and QoS compliance assurance, these studies show how AI may improve resource management.

### 8.1 Large-Scale Cloud Providers

AI-driven resource management techniques have been applied in large-scale cloud provider setups to improve the distribution of computer resources among several data centers. For instance, a top cloud provider achieved notable cost savings and enhanced performance by implementing a machine learning-based approach for automatic resource provisioning and scaling. By using an AI-driven strategy, the provider was able

to reduce both over- and under-provisioning by more effectively allocating resources based on actual demand.

## 8.2 Edge Computing

AI-driven resource management techniques have been used in edge computing environments to solve the particular difficulties brought on by dispersed and decentralized infrastructures. For example, an edge computing network of a telecommunications business used a reinforcement learning-based approach for dynamic resource distribution. Even with workloads that were changing quickly, the organization was still able to achieve low-latency requirements and optimize resource utilization thanks to the AI-driven methodology.

## 8.3 Serverless Computing

AI-powered resource management techniques have been employed in serverless computing settings to address the issues posed by unpredictable and dynamic workloads. One of the top e-commerce platforms, for instance, used a proactive resource management approach based on predictive analytics in its serverless architecture. The platform was able to predict times of high demand and deploy resources appropriately thanks to the AI-driven strategy, which enhanced performance and reduced costs.

## 8.4 Containerized Environments

Artificial intelligence (AI)-driven resource management techniques have been applied to containerized settings to maximize the distribution of computing resources among several containers. A technology company

used a genetic algorithm approach to optimize resource allocation in its containerized infrastructure in a multi-objective manner. The company was able to manage conflicting objectives including cost, performance, and energy conservation thanks to the AI-driven approach, which significantly improved system performance as a whole.

## 9. Conclusion

To get the most out of cloud computing, serverless computing requires effective resource management. The intricacies of resource distribution in these settings can be effectively addressed by AI technologies. Cloud providers may greatly increase the effectiveness and dependability of their services by utilizing AI-based techniques like intelligent scheduling, automated provisioning, predictive maintenance, and energy-efficient administration.

Future studies should concentrate on incorporating AI technologies into the frameworks for resource management that are already in place, with a focus on coordination and optimization in real-time. Additionally, in order to guarantee that AI-powered systems are reliable and meet user expectations, ethical issues, openness, and explainability should be given top priority. AI will become more and more significant in determining how resource management is developed in the future as cloud computing develops.

## References:

[1] Obi, Ogugua Chimezie, et al. "Review of evolving cloud computing paradigms: security, efficiency, and innovations."

*Computer Science & IT Research Journal*  
5.2 (2024): 270-292.

[2] Mampage, Anupama, Shanika Karunasekera, and Rajkumar Buyya. "A holistic view on resource management in serverless computing environments: Taxonomy and future directions." *ACM Computing Surveys (CSUR)* 54.11s (2022): 1-36.

[3] Hameed, Abdul, et al. "A survey and taxonomy on energy efficient resource allocation techniques for cloud computing systems." *Computing* 98 (2016): 751-774.

[4] Haibeh, Lina A., Mustapha CE Yagoub, and Abdallah Jarray. "A survey on mobile edge computing infrastructure: Design, resource management, and optimization approaches." *IEEE Access* 10 (2022): 27591-27610.

[5] Zhao, Kongyange, et al. "Taming serverless cold start of cloud model inference with edge computing." *IEEE Transactions on Mobile Computing* (2023).

[6] Mampage, Anupama, Shanika Karunasekera, and Rajkumar Buyya. "A holistic view on resource management in serverless computing environments: Taxonomy and future directions." *ACM Computing Surveys (CSUR)* 54.11s (2022): 1-36.

[7] Khan, Tahseen, et al. "Machine learning (ML)-centric resource management in cloud computing: A review and future directions." *Journal of Network and Computer Applications* 204 (2022): 103405.

[8] Kanungo, Satyanarayan. "AI-driven resource management strategies for cloud

computing systems, services, and applications." *World Journal of Advanced Engineering Technology and Sciences* 11.2 (2024): 559-566.

[9] Matsunaga, Andréa, and José AB Fortes. "On the use of machine learning to predict the time and resources consumed by applications." *2010 10th IEEE/ACM International Conference on Cluster, Cloud and Grid Computing. IEEE*, 2010.

[10] Lin, Jinying, et al. "A review on interactive reinforcement learning from human social feedback." *IEEE Access* 8 (2020): 120757-120765.

[11] Wang, Thomas, Simone Ferlin, and Marco Chiesa. "Predicting CPU usage for proactive autoscaling." *Proceedings of the 1st Workshop on Machine Learning and Systems*. 2021.

[12] Goldberg, David E. "Genetic and evolutionary algorithms come of age." *Communications of the ACM* 37.3 (1994): 113-120.

[13] Houssein, Essam H., et al. "Task scheduling in cloud computing based on meta-heuristics: review, taxonomy, open challenges, and future trends." *Swarm and Evolutionary Computation* 62 (2021): 100841.

[14] Bravo, César, et al. "State of the art of artificial intelligence and predictive analytics in the E&P industry: a technology survey." *Spe Journal* 19.04 (2014): 547-563.

[15] Mahmood, M. Rezwani, et al. "A comprehensive review on artificial intelligence/machine learning algorithms

*for empowering the future IoT toward 6G era." IEEE Access 10 (2022): 87535-8756.*