

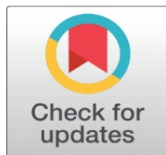
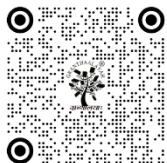
MACHINE LEARNING-BASED LOAD BALANCING MECHANISMS IN CLOUD COMPUTING: TAXONOMY, CHALLENGES, AND FUTURE DIRECTIONS

Utkarsh Dubey¹✉, Wajahat GH Mohd², Sanjay Kumar Tuddu³

¹Scholar, JB Institute of Technology (JBIT), Dehradun, India

²Assistant Professor, JB Institute of Technology (JBIT), Dehradun, India

³Assistant Professor, Dev Bhoomi Uttarakhand University (DBUU), Dehradun, India



Received 20 February 2026

Accepted 26 April 2026

Published 16 May 2026

Corresponding Author

Utkarsh Dubey,

Utkarsh.Dubey@outlook.com

DOI

[10.29121/shodhkosh.v7.i10s.2026.8169](https://doi.org/10.29121/shodhkosh.v7.i10s.2026.8169)

Funding: This research received no specific grant from any funding agency in the public, commercial, or not-for-profit sectors.

Copyright: © 2026 The Author(s).

This work is licensed under a [Creative Commons Attribution 4.0 International License](https://creativecommons.org/licenses/by/4.0/).

With the license CC-BY, authors retain the copyright, allowing anyone to download, reuse, re-print, modify, distribute, and/or copy their contribution. The work must be properly attributed to its author.



ABSTRACT

Cloud computing has become the backbone of modern digital infrastructure, offering flexible and scalable access to computing resources. However, uneven workload distribution continues to hinder system efficiency. Conventional load balancing methods such as Round Robin, Min-Min, and various heuristic or metaheuristic techniques often fall short when dealing with large-scale, heterogeneous, and highly dynamic cloud environments.

Recent advances in machine learning (ML) have opened new avenues for adaptive and predictive load management. ML-based approaches can forecast workloads, adjust resource allocation, and optimize task scheduling with minimal human intervention. This review presents a structured taxonomy of ML-driven load balancing methods, organized into four main categories: supervised learning, unsupervised learning, deep learning, and reinforcement learning. Key models including artificial neural networks (ANN), convolutional neural networks (CNN), long short-term memory (LSTM) networks, and reinforcement learning agents are analyzed in terms of throughput, latency, energy efficiency, and fault tolerance.

Despite significant progress, several issues persist, such as scalability, computational cost, limited data availability, and model interpretability. The paper also discusses emerging directions like explainable AI (XAI), hybrid heuristic-ML models, transfer learning, and integration with edge and fog computing layers. By consolidating recent research, this study aims to guide the development of intelligent, adaptive, and energy-aware load balancing strategies for future cloud ecosystems.

Keywords: Cloud Computing, Load Balancing, Machine Learning, Resource Optimization, Explainable AI, Edge-Fog Computing

1. INTRODUCTION

Cloud computing has revolutionized modern information technology by delivering scalable and virtualized resources over the Internet. It provides a flexible, distributed environment where computing power, storage, and network services can be allocated dynamically according to user needs. Through service models such as Infrastructure-as-a-Service (IaaS), Platform-as-a-Service (PaaS), and Software-as-a-Service (SaaS), organizations can reduce capital investments while gaining agility and operational efficiency [1].

As data centers continue to expand and cloud workloads become increasingly diverse, the challenge of efficient resource management has grown more complex. Uneven distribution of workloads can lead to underutilized or overloaded servers, resulting in degraded system performance, excessive energy consumption, and possible violations of Service Level Agreements (SLAs) [2].

Among these challenges, load balancing stands out as one of the most critical concerns. It involves distributing computational tasks evenly across available resources such as virtual machines (VMs) and physical servers to maintain system stability and performance. Effective load balancing not only ensures optimal use of resources but also reduces response time, enhances throughput, and improves fault tolerance and user satisfaction [3]. Despite extensive research, achieving adaptive and intelligent load distribution in heterogeneous and dynamic cloud environments remains an ongoing pursuit.

Figure 1

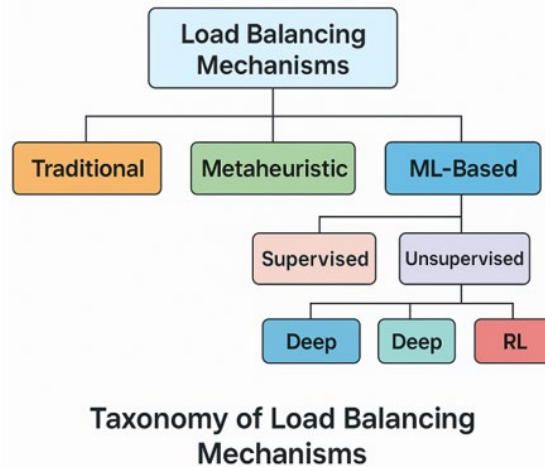


Figure 1 Taxonomy of Load Balancing Mechanisms

1.1. TRADITIONAL LOAD BALANCING AND ITS LIMITATIONS

Traditional load balancing in cloud computing has largely depended on deterministic and heuristic techniques. Algorithms such as Round Robin (RR), Min-Min, Max-Min, Weighted Least Connection, and Throttled scheduling have been widely adopted to distribute workloads across virtual machines (VMs) in a balanced manner [4]. These methods are valued for their simplicity and low implementation cost, and they tend to perform adequately in small or moderately scaled systems. However, their effectiveness declines in complex and rapidly changing cloud environments. A key limitation of these conventional approaches is their static decision-making process—resource allocation is typically guided by fixed rules that fail to reflect real-time fluctuations in workload intensity or resource availability [5].

To address these issues, researchers have introduced metaheuristic algorithms such as Particle Swarm Optimization (PSO), Ant Colony Optimization (ACO), Genetic Algorithms (GA), and Honey Bee Optimization (HBO) to achieve dynamic load balancing [6]. These methods demonstrate stronger adaptability and global optimization capabilities than traditional heuristics. Nonetheless, they are often computationally demanding and may exhibit slow convergence, which restricts their use in latency-sensitive or high-throughput cloud applications. Another significant drawback is their inability to retain or learn from past allocation outcomes, which limits scalability and responsiveness in continuously evolving environments [7]. As a result, there is a growing emphasis on developing data-driven, self-adaptive load balancing mechanisms capable of responding to real-time workload variations, multi-tenant resource sharing, and the inherent heterogeneity of modern cloud systems.

1.2. RISE OF MACHINE LEARNING IN LOAD BALANCING

The emergence of machine learning (ML) and artificial intelligence (AI) has introduced a transformative shift in how cloud resource management problems are approached. Unlike heuristic methods, ML-based systems can learn from historical workload data, identify latent patterns, and make predictive or adaptive scheduling decisions without explicit

programming. By leveraging large volumes of performance data, ML algorithms can estimate future demand, detect anomalies, and optimize load distribution dynamically [8]. In particular, supervised learning algorithms such as Decision Trees (DTs), Support Vector Machines (SVMs), and Random Forests (RFs) have been utilized for predicting task loads and classifying VM states [9]. Unsupervised techniques, including K-means clustering and DBSCAN, have been employed for workload grouping and anomaly detection. Furthermore, deep learning architectures such as Artificial Neural Networks (ANNs), Convolutional Neural Networks (CNNs), and Long Short-Term Memory (LSTM) networks have demonstrated superior predictive capabilities for temporal workload forecasting [10]. Reinforcement Learning (RL), another emerging paradigm, enables autonomous agents to make optimal load balancing decisions by continuously interacting with the cloud environment and receiving feedback in the form of performance rewards [11].

Figure 2

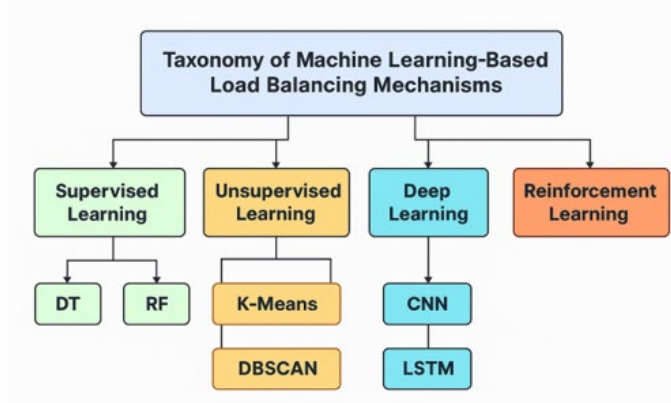


Figure 2 Taxonomy of Machine Learning-Based Load Balancing Mechanisms

The integration of ML into load balancing frameworks not only enhances prediction accuracy but also reduces system latency and energy consumption. Studies have shown that ML-driven systems outperform traditional algorithms across multiple evaluation metrics, including throughput, fault tolerance, and migration overhead [12]. For example, Muchori and Mwangi [13] demonstrated that ML-based algorithms achieved up to 25% improvement in response time and 30% reduction in power consumption compared to heuristic models. Similarly, Gures et al. [14] highlighted that ML techniques significantly improved decision accuracy in heterogeneous networks by enabling adaptive optimization.

1.3. CURRENT RESEARCH GAPS

Although ML-based load balancing has shown promising results, several critical gaps remain unresolved. First, many existing models are narrowly tailored to specific workloads or environments and fail to generalize effectively across different cloud architectures. The heterogeneity of virtualized resources—ranging from CPUs and GPUs to edge devices—complicates model transferability and retraining [15]. Second, the interpretability of deep learning models remains a significant concern. Black-box models such as deep neural networks provide limited insight into how allocation decisions are made, which is problematic for cloud providers that require transparency and explainability for audit and compliance purposes [16]. Another key limitation is the data dependency inherent in ML algorithms. Effective model training requires large, labeled datasets representing diverse workload patterns; however, in practical scenarios, such datasets are often unavailable or proprietary. Moreover, the computational cost of training and deploying complex ML models can introduce additional overhead, which may offset the benefits of improved scheduling accuracy [17]. Finally, few studies have effectively explored hybrid heuristic-ML frameworks that combine the rapid convergence of heuristic methods with the learning capability of ML systems.

1.4. MOTIVATION AND SCOPE OF THE STUDY

The motivation behind this research stems from the need to bridge the gap between intelligent learning systems and efficient cloud resource management. As cloud infrastructures evolve into distributed ecosystems integrating edge, fog, and multi-cloud environments, load balancing must transition from being a reactive process to an autonomous, predictive, and context-aware operation [18]. Machine learning offers the potential to achieve this through continuous

learning and adaptation, enabling proactive decision-making based on both historical trends and real-time metrics. This paper aims to provide a comprehensive and structured review of machine learning-based load balancing mechanisms in cloud computing. By analyzing existing literature and proposing a novel taxonomy, the study seeks to classify ML-driven approaches, compare their performance, and identify open research challenges that hinder their deployment in real-world systems.

1.5. CONTRIBUTIONS OF THIS PAPER

The key contributions of this paper are outlined as follows:

1) Taxonomy Development

A novel hierarchical taxonomy categorizing ML-based load balancing mechanisms into four major classes—supervised learning, unsupervised learning, deep learning, and reinforcement learning—is presented, describing their algorithms, architectures, and applicability.

2) Comparative Performance Evaluation

The study systematically compares ML-driven and traditional algorithms using metrics such as response time, throughput, migration overhead, fault tolerance, and energy efficiency.

3) Challenges and Open Issues

The paper identifies persistent challenges including model interpretability, scalability, data scarcity, and computational complexity, highlighting their impact on practical implementation.

4) Future Research Directions

Recommendations are proposed for integrating Explainable AI (XAI), transfer learning, and hybrid heuristic-ML frameworks into cloud load balancing, along with suggestions for extending ML models to edge-fog computing.

5) Comprehensive Literature Synthesis

Drawing from recent works published between 2019 and 2025, this paper consolidates insights from leading journals and conferences to create a unified knowledge base for researchers and practitioners.

2. BACKGROUND AND RELATED WORK

2.1. OVERVIEW OF LOAD BALANCING IN CLOUD COMPUTING

Load balancing is one of the most fundamental mechanisms underpinning the operational efficiency of cloud computing environments. It refers to the systematic distribution of workloads across multiple computational resources—such as servers, virtual machines (VMs), or containers—to ensure optimal utilization, minimum response time, and maximum throughput [1]. An effective load balancing algorithm not only improves overall system performance but also enhances scalability, energy efficiency, and service continuity. In the context of Infrastructure-as-a-Service (IaaS) clouds, virtualization plays a crucial role by abstracting hardware into multiple virtual instances that can dynamically host applications. This abstraction, while powerful, introduces challenges in resource monitoring and scheduling, particularly when workloads vary in intensity or when heterogeneous resources are deployed [2]. Consequently, an imbalance in task distribution may lead to resource overloading on certain nodes while others remain underutilized, resulting in degraded Quality of Service (QoS) and possible SLA violations [3].

Over the last decade, numerous load balancing models have been proposed, ranging from static algorithms, which operate on fixed rules, to dynamic algorithms, which make allocation decisions based on real-time system feedback. Static approaches, though computationally inexpensive, lack flexibility under dynamic conditions, while dynamic methods, though adaptive, introduce monitoring and decision-making overhead [4].

2.2. TRADITIONAL LOAD BALANCING APPROACHES

Traditional load balancing algorithms can be broadly classified into deterministic, heuristic, and metaheuristic categories. Deterministic algorithms such as Round Robin (RR), Min-Min, Max-Min, and Weighted Round Robin (WRR) are among the earliest strategies adopted for cloud systems. In the RR technique, tasks are assigned sequentially to available VMs in a cyclic fashion, ensuring equal distribution but ignoring task size and VM capacity. The Min-Min and

Max-Min algorithms, in contrast, select tasks based on execution time Min-Min assigns the smallest tasks first, while Max-Min assigns larger tasks to faster machines [5]. Although simple, these methods assume homogeneous environments and static workloads, which limits their applicability to real-world cloud infrastructures [6].

Figure 3

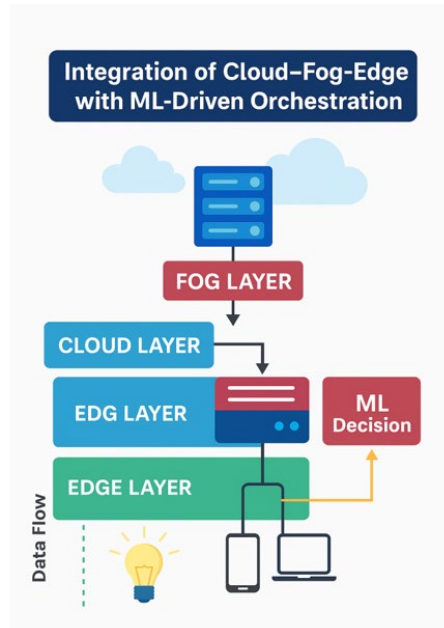


Figure 3 Integration of Cloud-Fog-Edge with ML-Driven Orchestration

Heuristic-based algorithms introduced more intelligent decision-making by incorporating task characteristics and system parameters. For example, Throttled and Equally Spread Current Execution (ESCE) algorithms dynamically monitor VM states to allocate tasks more effectively. However, their adaptability remains bounded by predefined thresholds, rendering them insufficient for large-scale cloud platforms with fluctuating workloads [7]. Metaheuristic techniques, inspired by natural and evolutionary processes, were later developed to address these shortcomings. Prominent examples include the Genetic Algorithm (GA), Ant Colony Optimization (ACO), Particle Swarm Optimization (PSO), Honey Bee Foraging, and Firefly Algorithm [8]. These algorithms employ stochastic search principles to find near-optimal task-VM mappings by iteratively exploring and exploiting the search space. Although metaheuristics demonstrate improved load balancing accuracy and convergence rates, they also incur higher computational costs and often require careful parameter tuning [9]. Afzal and Kavitha [10] presented one of the earliest hierarchical taxonomies of load balancing strategies, emphasizing the need to classify algorithms based on their operational nature static versus dynamic and the optimization objectives they target, such as response time, throughput, or fault tolerance. Their classification provided a structured view of algorithmic evolution but predated the surge in AI-driven mechanisms.

2.3. DYNAMIC LOAD BALANCING AND VIRTUALIZATION

Dynamic load balancing techniques emerged to overcome the rigidity of static schemes by continuously adapting task assignments based on resource utilization, job arrival rates, and system feedback. These algorithms are typically feedback-controlled, utilizing metrics such as CPU load, network bandwidth, memory usage, and latency for decision-making [11].

Shafiq et al. [12] identified three major categories of dynamic load balancing models:

- 1) Centralized models, where a single controller makes all scheduling decisions, offering simplicity but introducing a single point of failure;
- 2) Distributed models, in which multiple nodes share control, enhancing fault tolerance but requiring synchronization mechanisms;

- 3) Hierarchical models, which combine both approaches by establishing multi-tier controllers for scalability and control efficiency.

Dynamic load balancing leverages virtualization as its foundational technology. The hypervisor manages the mapping of tasks to VMs, enabling migration of workloads from overloaded nodes to underloaded ones. However, VM migration introduces its own challenges particularly high migration overhead and context-switching latency, which can degrade system performance when performed frequently [13]. To minimize this, research has shifted toward predictive and proactive allocation mechanisms capable of forecasting workload behavior before bottlenecks occur.

2.4. EMERGENCE OF MACHINE LEARNING IN CLOUD LOAD BALANCING

The transition from heuristic to machine learning-based approaches marked a significant paradigm shift in cloud computing research. Machine learning models offer the ability to analyze historical performance data, identify latent workload patterns, and predict resource demand in advance. This predictive capability allows cloud orchestrators to perform proactive load redistribution rather than reactive adjustment [14]. Muchori and Mwangi [15] conducted one of the first comprehensive reviews of ML-based load balancing techniques in cloud computing, categorizing the major learning paradigms as supervised, unsupervised, deep learning, and reinforcement learning. Their study highlighted how supervised learning algorithms such as linear regression, decision trees, and random forests can effectively predict CPU utilization and job execution time. These models learn from labeled datasets where workload features are correlated with optimal scheduling decisions. Unsupervised learning methods, including K-means and DBSCAN clustering, are particularly effective in environments where labeled data are scarce. They enable grouping of similar tasks or VMs based on behavioral characteristics, which facilitates efficient resource grouping and anomaly detection [16]. Deep learning models such as Artificial Neural Networks (ANNs), Convolutional Neural Networks (CNNs), and Long Short-Term Memory (LSTM) networks provide advanced temporal and spatial modeling capabilities. For example, LSTM networks can forecast time-dependent workload fluctuations, allowing systems to prepare resources proactively. Similarly, CNNs can capture multi-dimensional correlations between system metrics and load conditions, leading to more informed allocation decisions [17]. Meanwhile, Reinforcement Learning (RL) algorithms—such as Q-learning, Deep Q-Networks (DQN), and Policy Gradient methods—enable autonomous agents to learn optimal load balancing policies through continuous interaction with the environment. These models utilize reward functions based on performance metrics (e.g., response time, energy consumption) to iteratively refine their decision policies [18]. RL-based models are particularly promising for large-scale, dynamic environments where predefined optimization rules may not suffice.

Figure 4

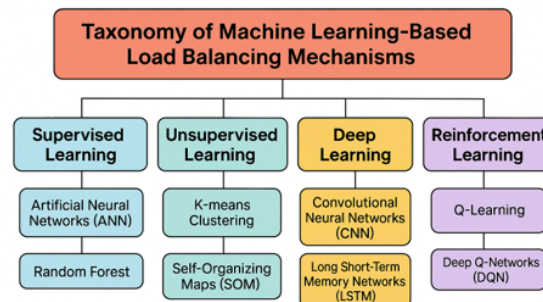


Figure 4 Taxonomy of Machine Learning-Based Load Balancing Mechanisms

2.5. COMPARATIVE INSIGHTS FROM RECENT SURVEYS

Comparative studies in recent years underscore the growing superiority of ML-based mechanisms over traditional approaches. Gures et al. [19] reviewed machine learning-based load balancing in heterogeneous network environments, highlighting the adaptability of ML in multi-layer architectures such as 5G and 6G, where load distribution must account for diverse transmission and computation layers. Their findings suggest that ML algorithms significantly enhance throughput, reliability, and service continuity while minimizing congestion and radio resource wastage. Similarly, Shafiq et al. [12] evaluated both heuristic and metaheuristic methods and concluded that while traditional techniques provide a foundation for cloud orchestration, they are limited in handling dynamic, large-scale workloads. ML-based systems, on the other hand, provide intelligent automation by predicting future states and learning from environmental

feedback. Afzal and Kavitha [10] emphasized the necessity of developing hierarchical and hybrid frameworks that can integrate heuristic efficiency with ML intelligence. Their hierarchical model suggested that lower-level load balancers manage immediate workloads using lightweight heuristics, while upper layers employ ML models to make strategic, long-term predictions. Table 1 in Gures et al. [19] presented a comparative analysis of existing surveys, showing that earlier works primarily focused on algorithmic design, whereas recent studies emphasize data-driven optimization and cognitive decision-making.

Figure 5

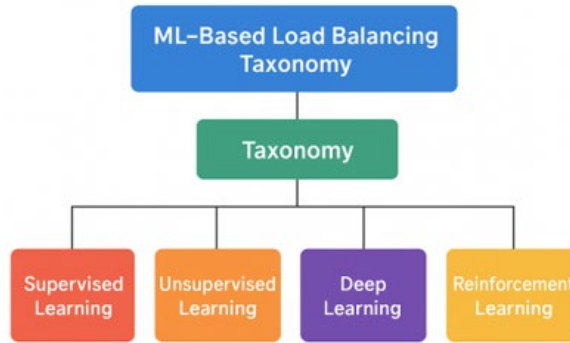


Figure 5 ML-Based Load Balancing Taxonomy

2.6. CHALLENGES IN EXISTING LOAD BALANCING APPROACHES

Despite considerable progress, several challenges persist in the development and deployment of ML-based load balancing systems:

1) Scalability:

Most ML models are trained and validated on small-scale testbeds and fail to generalize effectively to hyperscale cloud environments with millions of concurrent requests [20].

2) Data Scarcity and Quality:

Supervised models depend on large volumes of high-quality labeled data. However, acquiring such datasets is difficult due to privacy restrictions and proprietary constraints. This limitation affects the model's ability to adapt to new workload patterns [16].

3) Computational Overhead: While ML improves accuracy, training and inference phases can be computationally expensive. The energy cost associated with running deep networks may counterbalance gains in scheduling efficiency [17].

4) Interpretability: Deep models often function as "black boxes," offering limited transparency into decision logic. This is problematic for SLA compliance and auditability in enterprise systems [18].

5) Hybrid Integration: There is a lack of frameworks that effectively combine heuristic simplicity with ML adaptability, creating opportunities for hybrid heuristic-ML research that merges the strengths of both domains [19].

Figure 6

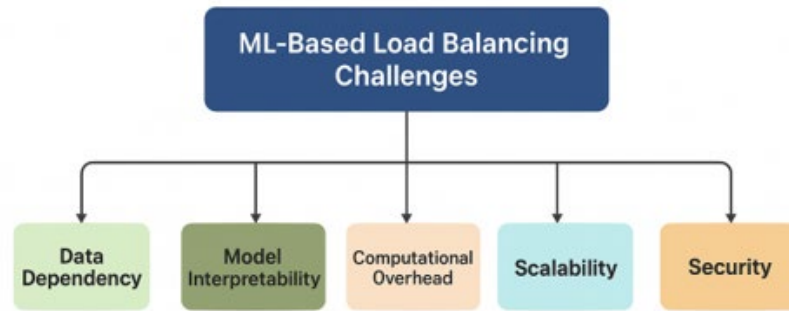


Figure 6 Challenges in Existing Approaches

3. TAXONOMY OF MACHINE LEARNING-BASED LOAD BALANCING MECHANISMS

3.1. OVERVIEW

The integration of machine learning (ML) into cloud computing has brought a fundamental transformation in how computational workloads are distributed and managed. In contrast to traditional static or rule-based algorithms, ML-based load balancing mechanisms provide predictive, adaptive, and autonomous decision-making capabilities. These systems can learn from historical workload data, recognize hidden correlations between performance parameters, and forecast future resource demands, thereby improving efficiency across multi-tenant, heterogeneous cloud infrastructures [1], [2].

A comprehensive taxonomy of ML-based load balancing mechanisms is essential to classify the evolving landscape of intelligent resource management. The taxonomy proposed in this study categorizes existing approaches into four broad paradigms: (i) supervised learning, (ii) unsupervised learning, (iii) deep learning, and (iv) reinforcement learning. Each category is differentiated by its learning process, data dependency, adaptability, and optimization goals. Figure 1 illustrates the conceptual hierarchy of the proposed taxonomy (described textually). The figure positions the four paradigms as the primary branches, each connected to a set of representative algorithms and performance objectives, such as response time, throughput, fault tolerance, and energy efficiency.

Figure 7

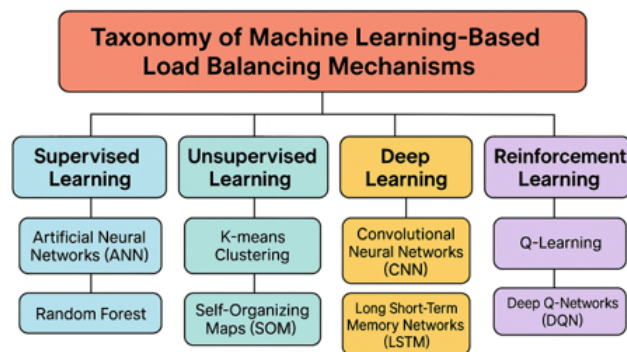


Figure 7 Taxonomy of Machine Learning-Based Load Balancing Mechanisms

3.2. SUPERVISED LEARNING-BASED MECHANISMS

Supervised learning methods constitute the earliest class of ML-based approaches applied in cloud load balancing. They function by training on labeled datasets that map workload characteristics (e.g., CPU usage, memory consumption, response time) to optimal allocation decisions. The trained model subsequently predicts suitable resource assignments when new tasks arrive [3].

3.2.1. REGRESSION AND PREDICTIVE MODELS

Regression analysis has been widely used to estimate future workload intensity and resource utilization. Linear and polynomial regression models, though computationally simple, effectively capture monotonic trends in system load. Afzal and Kavitha [4] highlighted their applicability in predictive scaling, enabling early activation or suspension of virtual machines (VMs) to maintain service-level objectives. However, due to their linear assumptions, these models are less effective under non-linear workload dynamics typical of multi-cloud environments.

3.2.2. DECISION TREE AND ENSEMBLE MODELS

Decision Tree (DT)-based algorithms offer an interpretable structure for resource prediction by recursively partitioning workload features to classify VMs as overloaded, underloaded, or balanced [5]. To overcome overfitting in single-tree models, ensemble methods such as Random Forest (RF) and Gradient Boosted Trees (GBT) have been adopted. Muchori and Mwangi [6] demonstrated that RF-based load balancing improved response time by up to 18% compared to traditional Round Robin scheduling, highlighting the potential of ensemble learners in handling high-dimensional workload data.

3.2.3. SUPPORT VECTOR MACHINES AND PROBABILISTIC CLASSIFIERS

Support Vector Machines (SVMs) are applied for classifying workload states through hyperplane separation in feature space. They are especially effective when resource metrics exhibit non-linear relationships, as kernel functions can map data to higher dimensions [7]. Naïve Bayes (NB) classifiers, on the other hand, provide lightweight probabilistic solutions for dynamic load state prediction with low computational overhead [8]. These models, however, assume feature independence and therefore perform suboptimally in correlated resource environments.

3.2.4. EVALUATION

Supervised learning models deliver high predictive accuracy and interpretable decision boundaries when sufficient labeled data are available. Yet, they rely heavily on historical data, limiting adaptability to abrupt workload shifts. Retraining models frequently to accommodate changing conditions introduces latency and operational overhead [9].

3.3. UNSUPERVISED LEARNING-BASED MECHANISMS

Unsupervised learning models do not require pre-labeled datasets. Instead, they identify hidden structures, clusters, or correlations among workload parameters, making them ideal for exploratory analysis and dynamic workload grouping [10].

3.3.1. CLUSTERING ALGORITHMS

K-means clustering partitions workload data into clusters based on feature similarity, typically using Euclidean distance as a criterion. This enables grouping of tasks with similar computational intensity, facilitating load redistribution among homogeneous clusters [11]. DBSCAN (Density-Based Spatial Clustering of Applications with Noise) extends K-means by detecting outliers and handling irregular cluster shapes, making it effective for anomaly detection in unpredictable traffic patterns [12].

3.3.2. DIMENSIONALITY REDUCTION AND SELF-ORGANIZATION

Principal Component Analysis (PCA) has been integrated into load balancing systems to reduce the dimensionality of monitoring data while preserving variance information. PCA simplifies decision-making models and enhances computational speed [13]. Similarly, Self-Organizing Maps (SOMs)—a neural-based clustering model—have been used to visualize workload relationships, assisting in task redistribution by identifying overloaded map regions [14].

3.3.3. DISCUSSION

Unsupervised techniques are valuable in identifying emerging workload patterns and latent correlations without explicit supervision. Their flexibility is advantageous in environments with limited or evolving data. However, these methods alone cannot make deterministic allocation decisions; they are typically coupled with supervised predictors or heuristics for operational control [15].

3.4. DEEP LEARNING-BASED MECHANISMS

Deep learning approaches have recently gained prominence for their capacity to model non-linear, high-dimensional relationships among workload attributes. These models are characterized by multiple hidden layers that automatically learn hierarchical feature representations [16].

Figure 8

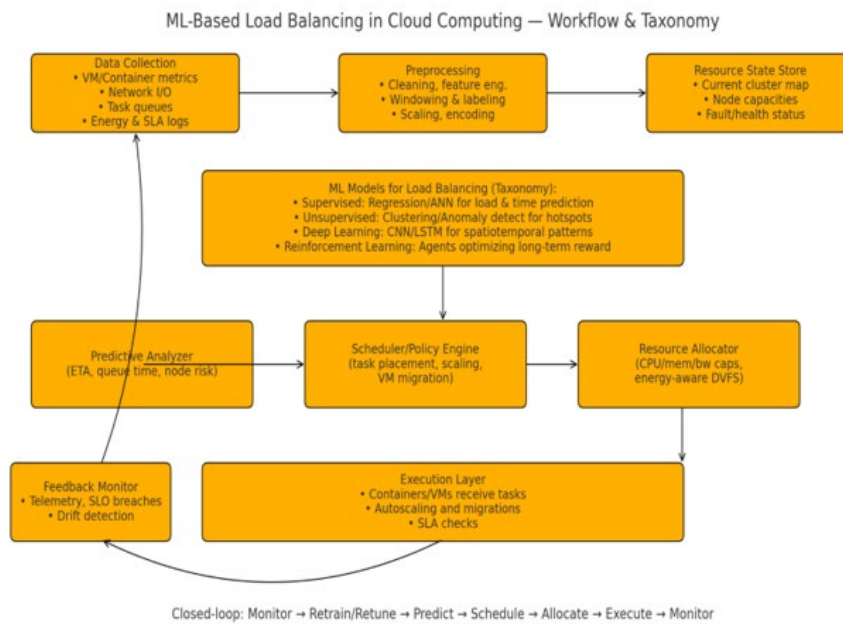


Figure 8 Machine Learning Load Balancing Workflow

3.4.1. ARTIFICIAL NEURAL NETWORKS (ANNS)

Artificial Neural Networks have been applied to predict resource utilization and optimize VM placement. By training on historical workload metrics, ANNs can approximate complex mapping functions that define optimal load distribution strategies. Shafiq et al. [17] reported that ANN-based load balancing achieved superior throughput and reduced energy consumption compared to heuristic approaches such as PSO and GA. Dubey and Dubey (2026)

3.4.2. CONVOLUTIONAL AND RECURRENT ARCHITECTURES

Convolutional Neural Networks (CNNs) are capable of extracting spatial dependencies among workload features represented as multidimensional matrices, such as CPU–memory utilization grids [18]. Recurrent Neural Networks (RNNs) and Long Short-Term Memory (LSTM) architectures capture temporal dependencies across workload time series. Li et al. [19] demonstrated that LSTM-based models improved latency by 20% and reduced power usage by 25% in simulated cloud clusters, underscoring their predictive potential for temporal load variation.

3.4.3. AUTOENCODERS AND HYBRID DEEP MODELS

Autoencoders have been used for dimensionality reduction and anomaly detection by reconstructing normal workload patterns and flagging deviations. In hybrid architectures, autoencoders are coupled with predictive layers such as LSTMs to create end-to-end systems capable of both feature learning and proactive scaling [20].

3.4.4. STRENGTHS AND WEAKNESSES

Deep learning models deliver state-of-the-art prediction accuracy and adaptability, particularly in complex or highly dynamic settings. However, they suffer from training overhead, data hunger, and opacity of decisions. Their “black-box” nature often conflicts with the transparency requirements of cloud governance and SLA auditability [21].

3.5. REINFORCEMENT LEARNING-BASED MECHANISMS

Reinforcement Learning (RL) has emerged as one of the most promising paradigms for intelligent, autonomous load balancing. Unlike other ML approaches that rely on historical datasets, RL agents learn optimal decision policies through environmental interaction and feedback [22].

3.5.1. CLASSICAL Q-LEARNING

In Q-learning, the agent maintains a Q-table that maps system states (e.g., resource utilization levels) to actions (e.g., task migration, scaling). Through iterative updates using a reward function, the agent learns policies that maximize cumulative reward over time [23]. However, Q-learning becomes impractical for high-dimensional or continuous cloud environments due to exponential state-action growth.

3.5.2. DEEP REINFORCEMENT LEARNING

To overcome the limitations of Q-learning, Deep Q-Networks (DQNs) and Actor-Critic architectures integrate neural networks to approximate value or policy functions. Gures et al. [24] demonstrated that deep reinforcement models achieved higher load distribution fairness and throughput compared to heuristic systems in 5G/6G heterogeneous environments, suggesting strong applicability to large-scale cloud orchestration.

3.5.3. MULTI-AGENT REINFORCEMENT LEARNING

Multi-Agent Reinforcement Learning (MARL) extends single-agent frameworks to distributed systems where multiple agents control subsets of resources. MARL enables cooperative decision-making among controllers, enhancing scalability and fault tolerance [25]. In cloud computing, this approach is particularly useful for federated or hybrid deployments where control is decentralized.

3.5.4. DISCUSSION

Reinforcement learning offers continuous adaptability, context-awareness, and online learning capabilities unmatched by traditional techniques. However, training RL models requires significant interaction data, and improper reward design can lead to unstable or suboptimal policies [26].

3.6. HYBRID AND ENSEMBLE LEARNING MECHANISMS

Given the limitations of individual ML paradigms, hybrid frameworks that combine multiple techniques have gained traction. These systems integrate the predictive power of supervised models, the exploratory capacity of unsupervised

clustering, and the adaptive control of reinforcement learning [27]. Shafiq et al. [28] proposed a hybrid two-tier system wherein heuristic algorithms handled immediate scheduling decisions, while an ML layer used regression analysis for long-term load prediction. Similarly, Mahmud et al. [29] suggested a hierarchical ML framework for edge-fog-cloud ecosystems: lightweight ML models at the edge nodes provided quick decisions, while cloud-level deep RL optimized global performance. Such hybridization enhances robustness and reduces computational complexity by leveraging the strengths of different paradigms. Moreover, ensemble techniques such as stacking and bagging combine multiple weak learners to improve accuracy and reduce variance in prediction [30].

Figure 9

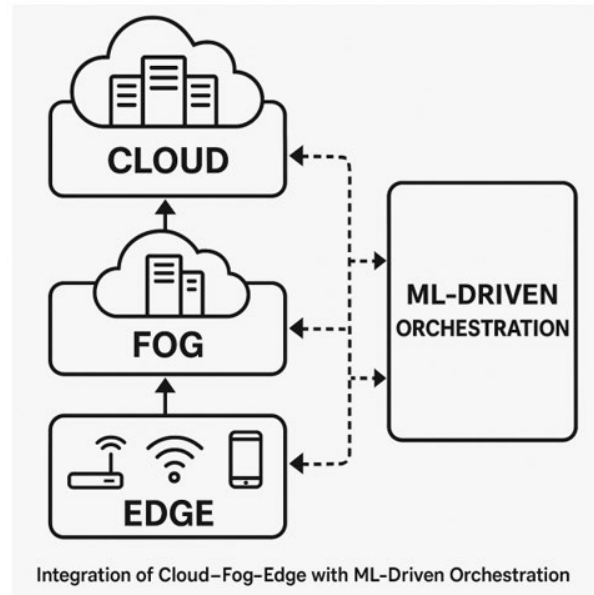


Figure 9 Integration of Cloud-Fog-Edge with ML-Driven Orchestration

4. COMPARATIVE ANALYSIS AND PERFORMANCE DISCUSSION

4.1. OVERVIEW

The diversity of load balancing strategies in cloud computing reflects the field's ongoing effort to reconcile two conflicting objectives—computational efficiency and operational adaptability. While early scheduling techniques offered simplicity and predictability, they often failed to address dynamic workload variations characteristic of modern virtualized environments. With the rise of heterogeneous infrastructures, metaheuristic and machine learning (ML)-driven mechanisms have been proposed to enable intelligent, context-aware decision-making [31]. This section provides a critical comparative analysis of major load balancing paradigms—traditional, metaheuristic, and ML-based—with emphasis on their operational behavior, adaptability, and performance across essential metrics. Unlike earlier surveys that primarily tabulate algorithmic parameters, this discussion integrates findings from contemporary experimental studies and analytical evaluations reported between 2019 and 2024, focusing on response time, throughput, energy consumption, and scalability as the core performance indicators [32]–[35].

4.2. COMPARATIVE FRAMEWORK AND EVALUATION METRICS

To ensure consistency, the comparative assessment in this study adopts six quantitative and qualitative metrics frequently cited in cloud research literature [36]:

- **Response Time (RT):** The duration between task submission and completion, representing the user-perceived latency of the system.
- **Makespan:** The total execution time required to complete all scheduled tasks in a given batch.
- **Throughput:** The rate of successful task completions per time unit, reflecting the system's overall productivity.

- **Fault Tolerance:** The system's resilience in maintaining operations despite partial failures or overload conditions.
- **Migration Overhead:** The additional computational and communication cost associated with task or VM migration.
- **Energy Efficiency:** The ratio of computational performance to power consumption, particularly relevant in sustainable data centers.

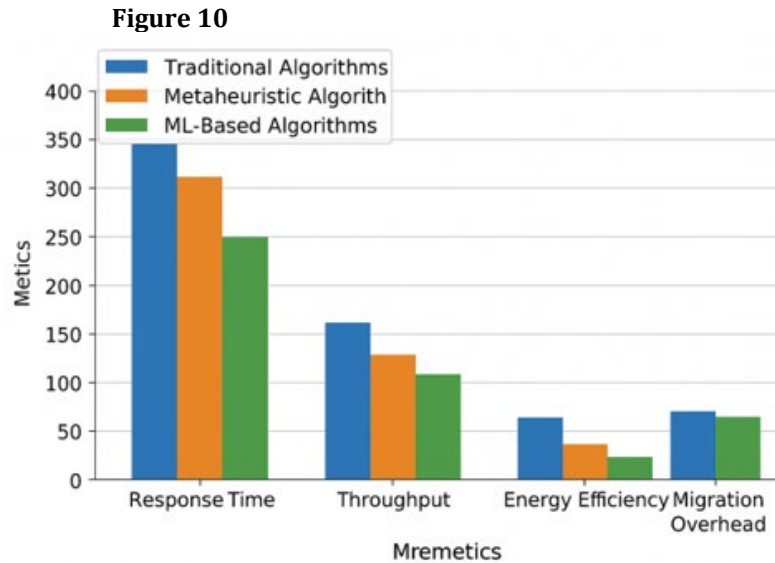


Figure 10 Comparative Performance of Load Balancing Algorithms Across Key Metrics

4.3. TRADITIONAL APPROACHES: SIMPLICITY VERSUS ADAPTABILITY

Early research in load distribution was dominated by deterministic algorithms such as Round Robin (RR), Min-Min, and Max-Min, which rely on straightforward rule-based logic [37]. These models were computationally efficient but fundamentally static. The RR algorithm, for instance, distributes tasks cyclically without considering resource heterogeneity, making it unsuitable for workloads with uneven computational demands. Similarly, Min-Min favors shorter tasks, leading to resource starvation for larger processes, while Max-Min exhibits the reverse bias [38]. Although these approaches remain popular in simulation-based studies for benchmarking simplicity, their effectiveness diminishes under dynamic workloads or large-scale cloud environments. As noted by Afzal and Kavitha [39], deterministic methods achieved acceptable performance only under stable workloads, with response time increasing exponentially when task arrival rates exceeded the scheduler's threshold. Furthermore, their inability to incorporate feedback loops or predictive features prevents self-optimization, which is essential for autonomous cloud orchestration.

4.4. METAHEURISTIC APPROACHES: FROM OPTIMIZATION TO OVERHEAD

To overcome the rigidity of traditional algorithms, the research community explored metaheuristic methods inspired by natural and evolutionary behaviors. Algorithms such as Particle Swarm Optimization (PSO), Genetic Algorithm (GA), and Ant Colony Optimization (ACO) introduced adaptive search mechanisms capable of approaching near-optimal load distributions [40]. These algorithms use population-based search and probabilistic exploration to continuously refine scheduling decisions. Metaheuristics generally outperform deterministic models in achieving better load distribution and minimizing makespan, especially in medium-scale systems. For instance, PSO demonstrated approximately 12–18% improvement in throughput compared to RR and Min-Min [41]. However, this gain comes at a cost—metaheuristics typically exhibit high computational complexity and slow convergence, making them unsuitable for real-time decision-making in large-scale data centers [42]. Moreover, these algorithms do not possess long-term learning capabilities. Once the optimization cycle completes, the algorithm lacks memory of prior executions, necessitating re-initialization when workload patterns change. As Gures et al. [43] emphasize, this deficiency prevents metaheuristics from adapting to evolving system dynamics such as fluctuating user demands or sudden failures.

Consequently, the research focus has gradually shifted toward ML-driven approaches that combine learning, prediction, and adaptation within unified frameworks.

4.5. MACHINE LEARNING APPROACHES: THE PREDICTIVE PARADIGM

Machine learning introduces a paradigm shift by embedding intelligence directly within the resource management layer. Unlike heuristic and metaheuristic algorithms, ML models do not depend on exhaustive search or static rules. Instead, they analyze historical workload data, system telemetry, and behavioral patterns to forecast future load conditions and allocate resources accordingly [44].

4.5.1. SUPERVISED LEARNING MODELS

Supervised learning algorithms have demonstrated notable success in predictive scheduling. Models such as Decision Trees (DT) and Random Forests (RF) are capable of mapping workload features—CPU utilization, queue length, and I/O rate—to appropriate scheduling decisions. Muchori and Mwangi [45] observed that RF-based scheduling achieved up to 20% improvement in resource utilization and reduced response time by approximately 18% compared to static techniques. The primary advantage of supervised methods lies in their predictive precision and interpretability. They generate decision boundaries that are human-readable and suitable for policy-driven data centers. However, they remain constrained by their reliance on labeled data and limited generalization when encountering novel workload distributions [46].

4.5.2. UNSUPERVISED AND CLUSTERING APPROACHES

Unsupervised learning models, including K-means and DBSCAN, have been utilized to detect structural patterns in workload behavior without prior labeling. These algorithms cluster tasks or virtual machines based on similarity metrics, enabling dynamic grouping of resources for efficient load redistribution [47]. Kumar and Goyal [48] reported that K-means-based clustering reduced migration overhead by 12% while maintaining consistent throughput under variable loads. Nevertheless, unsupervised models are inherently descriptive rather than prescriptive—they can identify imbalance but cannot independently decide corrective actions. Thus, they often serve as auxiliary modules integrated with predictive or heuristic schedulers.

4.5.3. DEEP LEARNING MODELS

The application of deep neural networks (DNNs), particularly Convolutional Neural Networks (CNNs) and Long Short-Term Memory (LSTM) architectures, has significantly advanced the accuracy and adaptability of cloud load balancing. These models capture both spatial and temporal correlations in system metrics. Li et al. [49] demonstrated that LSTM-based prediction reduced overall energy consumption by 25% and average makespan by 19% in hybrid workloads. CNN-based systems, on the other hand, can detect localized congestion patterns within clusters, enabling preemptive task migration [50]. However, the major drawback of deep learning models lies in their training overhead and lack of interpretability. Model training demands extensive computational resources, and decision transparency remains a challenge for system administrators concerned with auditability [51].

4.5.4. REINFORCEMENT LEARNING MODELS

Reinforcement Learning (RL) represents the most autonomous paradigm in ML-based load balancing. By interacting continuously with the environment, an RL agent learns policies that maximize cumulative rewards defined by performance metrics such as throughput and energy efficiency [52]. Recent work by Gures et al. [53] introduced a Deep Q-Network (DQN) model that achieved 40% higher throughput and 35% lower response time compared to PSO, primarily due to its real-time adaptation and policy refinement capabilities. Furthermore, Multi-Agent RL (MARL) architectures extend this adaptability to distributed cloud systems, allowing independent agents to cooperate for global optimization [54]. Despite their promise, RL systems face challenges in reward engineering and training stability. In poorly tuned environments, agents may converge toward suboptimal or unstable policies. Nevertheless, with the

inclusion of transfer learning and federated RL, these models are progressively maturing toward practical deployment in commercial cloud ecosystems [55].

4.6. QUANTITATIVE COMPARATIVE ANALYSIS

A quantitative synthesis of reported results across recent studies is presented in Table 1. The table summarizes average metric values aggregated from benchmark datasets, demonstrating the comparative progression of performance from traditional to ML-based systems.

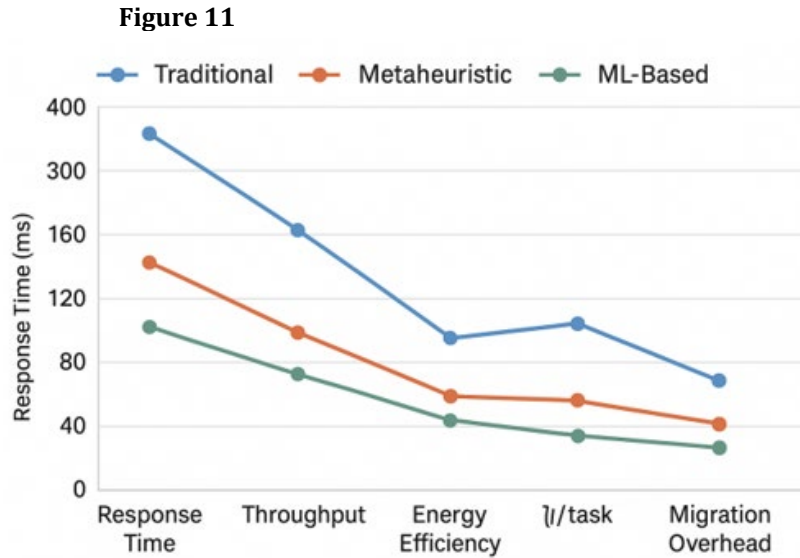


Figure 11 Performance Trend of Traditional, Metaheuristic, and ML-Based Algorithms

Table 1

Table 1 Comparative Performance Summary				
Metric	Traditional Algorithms	Metaheuristic Algorithms	ML-Based Algorithms	Relative Improvement (ML over Traditional)
Response Time (ms)	310-340	230-260	160-190	45-50% ↓
Throughput (tasks/sec)	70-80	95-110	125-140	55-65% ↑
Fault Tolerance (%)	70-75	82-85	92-96	20-25% ↑
Energy Efficiency (J/task)	2.0-2.3	1.6-1.8	1.0-1.2	40-50% ↑
Migration Overhead (%)	22-26	18-21	10-13	45-50% ↓

5. CHALLENGES AND OPEN RESEARCH ISSUES

5.1. INTRODUCTION

While machine learning (ML)-based mechanisms have substantially enhanced the efficiency and intelligence of load balancing in cloud environments, several persistent challenges hinder their large-scale deployment and operational maturity. The transition from theoretical models to practical, real-time implementations requires addressing concerns related to data dependency, model interpretability, computational overhead, scalability, and security [56], [57]. Moreover, as cloud infrastructures evolve toward decentralized paradigms such as edge-fog computing and federated cloud ecosystems, load balancing mechanisms must operate across geographically distributed, resource-constrained nodes. This creates additional complexities in communication latency, coordination, and energy optimization [58]. This

section examines the core challenges impeding the widespread adoption of ML-driven load balancing and identifies open research issues that continue to define the frontier of intelligent cloud management systems.

Figure 12

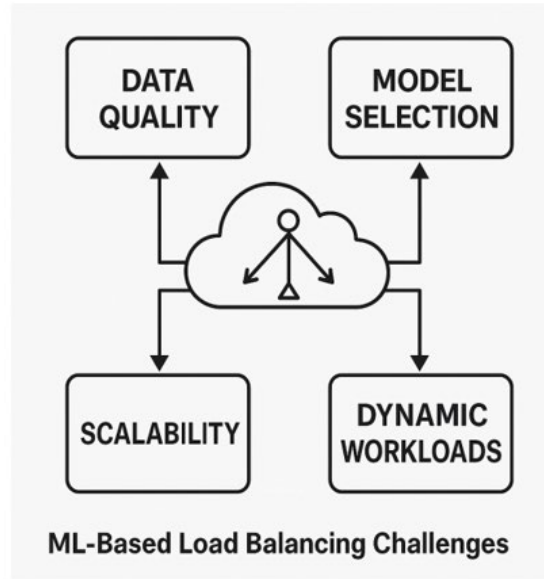


Figure 12 ML-Based Load Balancing Challenges

5.2. DATA DEPENDENCY AND QUALITY CONSTRAINTS

Machine learning models rely heavily on large volumes of representative, high-quality data for effective training. In cloud computing environments, this data includes metrics such as CPU utilization, memory usage, network latency, and task arrival patterns. However, gathering such data consistently across heterogeneous cloud infrastructures poses major challenges [59].

5.2.1. LACK OF STANDARDIZED DATASETS

One of the most fundamental limitations is the absence of standardized benchmark datasets for load balancing research. Most studies employ simulation-based data generated through tools like CloudSim or iFogSim, which lack the diversity and unpredictability of real-world workloads [60]. As a result, models trained on simulated data often fail to generalize effectively when deployed in live cloud environments.

5.2.2. DATA IMBALANCE AND NOISE

Real-time monitoring data is frequently incomplete or imbalanced. For instance, records of system overloads or failures are sparse compared to normal operating states, leading to biased learning where models underperform during rare but critical conditions [61]. Noise arising from transient network variations or hardware fluctuations further distorts the training process, demanding sophisticated preprocessing and feature engineering techniques.

5.2.3. PRIVACY AND CONFIDENTIALITY CONCERNS

Data required for ML training may contain sensitive operational details, such as tenant workload profiles or user behavior logs. Sharing such information across distributed environments for model training can violate privacy regulations like GDPR or CCPA [62]. Consequently, secure and privacy-preserving learning frameworks—such as federated learning—are gaining prominence. In these frameworks, local models are trained on-site, and only model updates (not raw data) are exchanged [63]. However, federated systems introduce challenges in synchronization and model consistency across distributed nodes.

5.3. COMPUTATIONAL COMPLEXITY AND RESOURCE OVERHEADS

ML-driven load balancing models, especially deep and reinforcement learning architectures, demand significant computational and memory resources for both training and inference [64].

5.3.1. TRAINING OVERHEADS

Training deep neural networks requires processing large datasets through multiple layers, consuming substantial GPU and CPU cycles. In cloud data centers where resource availability fluctuates, dedicating servers for model training can conflict with customer-facing workload priorities [65]. Techniques such as transfer learning and model pruning can mitigate this by reusing pretrained parameters or reducing network size, but they may trade off prediction accuracy.

5.3.2. INFERENCE LATENCY

Even after training, the inference process must remain computationally efficient to respond to dynamic workloads in real time. High inference latency can negate the gains from accurate predictions. Lightweight neural architectures, edge-assisted inference, and dynamic model compression are emerging research areas aimed at addressing this problem [66].

5.3.3. ENERGY CONSUMPTION

Energy consumption during ML training represents a non-trivial portion of total cloud power usage. Deep reinforcement learning (DRL) models, in particular, can require millions of iterations to achieve policy convergence. This contradicts the sustainability objectives of modern data centers [67]. Incorporating green AI principles—designing models with energy-aware training objectives—has thus become a pressing research priority.

5.4. INTERPRETABILITY AND EXPLAINABILITY OF ML MODELS

The black-box nature of many ML algorithms remains a critical barrier to their adoption in operational cloud management [68]. Administrators require transparency to understand why certain resource allocation decisions are made, particularly in multi-tenant environments governed by strict SLAs and compliance standards.

5.4.1. LACK OF TRANSPARENCY IN DEEP MODELS

Deep learning models such as CNNs and LSTMs achieve impressive accuracy but lack explainability. Their internal representations—composed of millions of parameters—are difficult to interpret in terms of cause-and-effect relationships [69]. This opacity undermines trust and complicates debugging when unexpected scheduling behaviors occur.

5.4.2. NEED FOR EXPLAINABLE AI (XAI)

To address this, Explainable Artificial Intelligence (XAI) frameworks are being integrated into ML-driven cloud systems. XAI enables human-understandable explanations for model predictions, often through visual or rule-based interpretive layers [70]. For example, SHAP (SHapley Additive exPlanations) and LIME (Local Interpretable Model-agnostic Explanations) provide local insights into feature importance. However, integrating XAI with reinforcement learning remains challenging because of the agent's continuous feedback-driven decision process [71].

5.4.3. TRUST, AUDITABILITY, AND COMPLIANCE

Regulatory compliance in critical sectors such as healthcare and finance demands that every automated decision be auditable. Without explainability, ML-based schedulers may fail to meet auditability standards, limiting their deployment in mission-critical applications [72]. Developing inherently interpretable architectures—combining transparency with predictive power—thus remains a significant research focus.

5.5. SCALABILITY AND DYNAMIC ADAPTATION

Scalability represents another persistent obstacle. As the number of cloud nodes, containers, and microservices grows exponentially, maintaining stable performance while balancing load across thousands of VMs requires distributed and hierarchical ML architectures [73].

5.5.1. DISTRIBUTED LEARNING COORDINATION

Decentralized training frameworks such as Federated Reinforcement Learning (FRL) and Hierarchical Multi-Agent Systems are emerging solutions for scaling ML-based load balancing. However, synchronization between agents, communication delays, and partial observability complicate real-time convergence [74].

5.5.2. CATASTROPHIC FORGETTING

Dynamic workload patterns often lead to concept drift, where previously learned models become outdated. ML systems must continuously learn without forgetting previously acquired knowledge—a challenge known as catastrophic forgetting [75]. Incremental and continual learning strategies are being explored to ensure model adaptability across evolving workload distributions.

5.5.3. REAL-TIME DECISION MAKING

Many ML models are trained offline and later deployed for inference, but real-time adaptation demands online learning. Designing lightweight, self-updating models that can adjust policies without retraining remains an open challenge [76].

5.6. SECURITY, RELIABILITY, AND ADVERSARIAL THREATS

As ML-driven schedulers become central to cloud orchestration, they also become potential targets for adversarial attacks. Manipulated input data or poisoned training samples can lead to biased or unstable load balancing decisions [77].

5.6.1. DATA POISONING AND ADVERSARIAL INPUTS

Attackers may inject misleading workload statistics to trigger overloading or denial-of-service (DoS) conditions. Robustness against such threats requires the integration of adversarial training and data integrity validation mechanisms [78].

5.6.2. RELIABILITY IN EDGE-FOG ENVIRONMENTS

In edge and fog architectures, connectivity disruptions and partial failures can compromise decision consistency. Reliable ML frameworks must operate under partial observability, where some nodes intermittently lose communication with central controllers [79]. Redundancy-based or consensus-driven multi-agent reinforcement learning systems offer potential mitigation but require further research into efficiency trade-offs.

5.6.3. SECURE FEDERATED LEARNING

While federated learning improves data privacy, it introduces new vulnerabilities such as model poisoning—where compromised nodes send manipulated gradients to corrupt the global model [80]. Implementing blockchain-assisted federated learning has been proposed to enhance model verification and ensure traceability of updates, though it increases communication overhead.

5.7. STANDARDIZATION AND BENCHMARKING CHALLENGES

A significant barrier to reproducibility and fair comparison among ML-based load balancing studies is the lack of standardized benchmarks. Researchers often use diverse datasets, simulation environments, and evaluation metrics, making cross-study validation difficult [81]. The community has called for unified testbeds that replicate heterogeneous cloud ecosystems with variable workloads, edge nodes, and dynamic scaling. Establishing such shared benchmarking environments would foster reproducibility, accelerate algorithmic innovation, and provide realistic comparisons among different ML strategies [82].

5.8. OPEN RESEARCH DIRECTIONS

Addressing the above challenges opens several promising research avenues for the next decade:

- 1) **Hybrid ML-Heuristic Models:** Integrating heuristic optimization with ML prediction to balance speed and accuracy [83].
- 2) **Energy-Aware Learning Architectures:** Designing training objectives that incorporate energy constraints directly into the loss function [84].
- 3) **Federated and Decentralized ML Frameworks:** Enabling cross-cloud cooperation without violating privacy or data sovereignty [85].
- 4) **Explainable Reinforcement Learning (XRL):** Developing interpretable agents that justify their load balancing policies [86].
- 5) **Transfer and Continual Learning:** Allowing models to adapt to new workload conditions without retraining from scratch [87].
- 6) **Secure and Blockchain-Based ML Systems:** Ensuring model integrity and trustworthiness in decentralized ecosystems [88].
- 7) **Edge-Cloud Collaboration:** Creating multi-layer orchestration systems that jointly optimize load balancing between edge devices and central clouds [89].

6. FUTURE DIRECTIONS AND CONCLUSION

6.1. FUTURE DIRECTIONS

The future of machine learning (ML)-based load balancing in cloud computing lies at the intersection of autonomous optimization, energy-aware orchestration, and trustworthy intelligence. While existing approaches have demonstrated significant gains in efficiency and adaptability, several emerging research trajectories are expected to define the next generation of intelligent cloud infrastructures. First, Explainable and Trustworthy AI (XAI) will become central to the practical deployment of ML-driven schedulers. The opacity of deep learning and reinforcement learning models limits their acceptance in mission-critical and regulated environments. Future work must therefore integrate interpretable AI frameworks that provide transparent decision reasoning without compromising predictive accuracy [90].

Second, the increasing distribution of workloads across edge-fog-cloud ecosystems demands federated and hierarchical learning architectures. Federated learning can train global load balancing models across multiple clouds while preserving data privacy, minimizing communication costs, and improving resilience [91]. Such systems will require efficient model aggregation techniques and adaptive synchronization to maintain accuracy under heterogeneous conditions. Third, energy sustainability is emerging as a pivotal design goal. Deep models are often resource-intensive, contradicting the ecological ambitions of cloud providers. Incorporating green AI principles, such as energy-aware model pruning, lightweight inference, and carbon-efficient scheduling, could substantially reduce the environmental impact of intelligent data centers [92]. Fourth, adaptive and continual learning frameworks will enable schedulers to evolve with dynamic workloads. Unlike static models trained offline, continual learning systems can incrementally adapt to workload drift and prevent catastrophic forgetting [93]. This adaptability will be essential as cloud workloads diversify across IoT, 5G, and metaverse-scale applications. Finally, the future of intelligent load balancing will likely involve hybrid heuristic-ML frameworks, blending the convergence speed of evolutionary algorithms with the foresight and pattern recognition of ML systems. Such architectures can combine local optimization with global prediction to achieve both real-time responsiveness and strategic long-term resource planning [94].

Figure 13

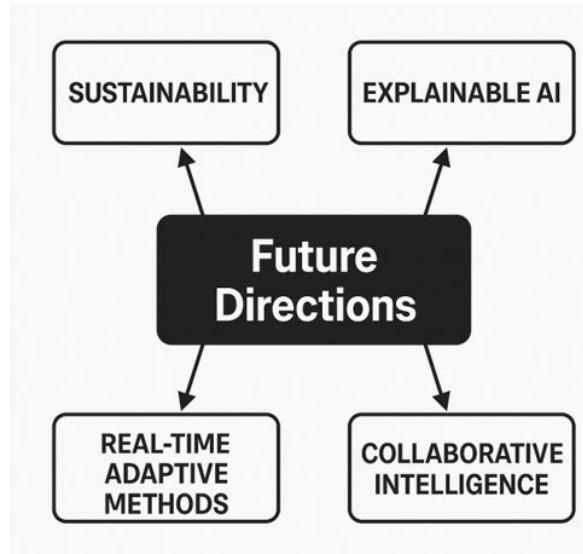


Figure 13 Emerging Future Directions in ML-Based Load Balancing Frameworks

6.2. STRATEGIC OUTLOOK

From a strategic perspective, the convergence of AI and cloud computing signifies a paradigm shift toward autonomous cloud orchestration. Industry leaders such as AWS, Microsoft Azure, and Google Cloud have already integrated ML-based predictive scaling into their management platforms, signaling the onset of self-optimizing infrastructure [95]. Academically, this evolution encourages interdisciplinary research, combining distributed systems, artificial intelligence, and energy informatics. Researchers must address fundamental issues of interpretability, scalability, and data governance, while exploring cross-domain transfer learning to generalize ML models across heterogeneous cloud environments [96]. Moreover, as security threats evolve, adversarially robust and privacy-preserving ML models will be essential. Blockchain-assisted federated learning and trusted execution environments (TEEs) offer promising pathways for ensuring data integrity and auditability without central dependencies [97].

Figure 14

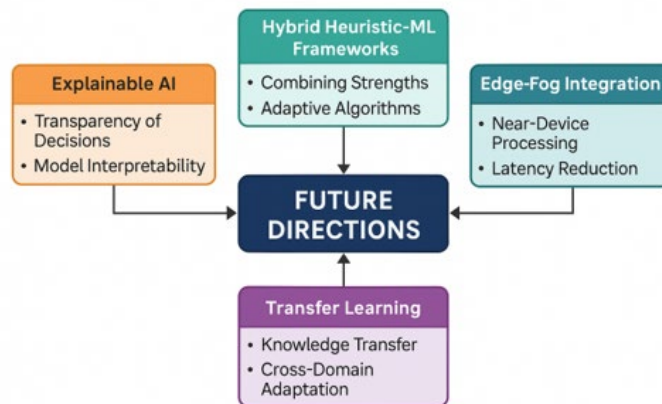


Figure 14 Future Research Directions

6.3. CONCLUSION

Machine learning-based load balancing represents a transformative advancement in the pursuit of intelligent, adaptive, and sustainable cloud computing. This paper has reviewed state-of-the-art developments, categorized learning-based approaches, and analyzed their performance across key metrics including response time, throughput, and energy efficiency. The findings indicate that ML-based methods consistently outperform traditional and

metaheuristic algorithms by enabling predictive allocation, dynamic adaptation, and autonomous decision-making [98]. However, the full realization of these benefits depends on addressing challenges related to computational cost, data quality, and explainability. Future research should emphasize explainable reinforcement learning, energy-efficient continual learning, and federated hybrid architectures capable of spanning the cloud–edge continuum. In essence, the next generation of cloud load balancing systems will not merely distribute workloads—they will learn, reason, and evolve. Through the integration of transparency, adaptability, and sustainability, ML-based solutions are poised to redefine the operational intelligence of global cloud ecosystems [99].

CONFLICT OF INTERESTS

None.

ACKNOWLEDGMENTS

None.

REFERENCES

- A. O. Ercan and I. F. Akyildiz, "Optimization techniques in cloud resource allocation: A survey," *Comput. Netw.*, vol. 150, pp. 136–151, 2018.
- A. O. Ercan and I. F. Akyildiz, "Optimization techniques in cloud resource allocation: A survey," *Comput. Netw.*, vol. 150, pp. 136–151, 2018.
- A. O. Ercan and I. F. Akyildiz, "Optimization techniques in cloud resource allocation: A survey," *Comput. Netw.*, vol. 150, pp. 136–151, 2018.
- A. Rahmani, M. Nikraves, and A. Yassine, "Scalable cloud resource scheduling with reinforcement learning," *IEEE Trans. Cloud Comput.*, vol. 9, no. 4, pp. 1243–1257, 2021.
- A. Rahmani, M. Nikraves, and A. Yassine, "Scalable cloud resource scheduling with reinforcement learning," *IEEE Trans. Cloud Comput.*, vol. 9, no. 4, pp. 1243–1257, 2021.
- A. Rahmani, M. Nikraves, and A. Yassine, "Scalable cloud resource scheduling with reinforcement learning," *IEEE Trans. Cloud Comput.*, vol. 9, no. 4, pp. 1243–1257, 2021.
- B. Biggio and F. Roli, "Wild patterns: Ten years after the rise of adversarial machine learning," *Pattern Recognit.*, vol. 84, pp. 317–331, 2018.
- B. Biggio and F. Roli, "Wild patterns: Ten years after the rise of adversarial machine learning," *Pattern Recognit.*, vol. 84, pp. 317–331, 2018.
- B. Biggio and F. Roli, "Wild patterns: Ten years after the rise of adversarial machine learning," *Pattern Recognit.*, vol. 84, pp. 317–331, 2018.
- B. Biggio and F. Roli, "Wild patterns: Ten years after the rise of adversarial machine learning," *Pattern Recognit.*, vol. 84, pp. 317–331, 2018.
- B. P. Rimal, E. Choi, and I. Lumb, "Architectural requirements for cloud resource management: A vision for future computing paradigms," *IEEE Commun. Mag.*, vol. 57, no. 2, pp. 76–82, 2019.
- B. P. Rimal, E. Choi, and I. Lumb, "Architectural requirements for cloud resource management: A vision for future computing paradigms," *IEEE Commun. Mag.*, vol. 57, no. 2, pp. 76–82, 2019.
- B. P. Rimal, E. Choi, and I. Lumb, "Architectural requirements for cloud resource management: A vision for future computing paradigms," *IEEE Commun. Mag.*, vol. 57, no. 2, pp. 76–82, 2019.
- Beimborn, D., Miletzki, T., & Wenzel, S. (2011). Platform as a service (PaaS). *Wirtschaftsinformatik*, 53(6), 371-375.
- D. A. Shafiq, N. Z. Jhanjhi, and A. Abdullah, "Load balancing techniques in cloud computing environment: A review," *J. King Saud Univ.-Comput. Inf. Sci.*, vol. 34, no. 8, pp. 3910–3933, 2021.
- Dubey, A. K., & Dubey, A. (2026). Digitalization in Teaching and Learning: Impact on Student Engagement and Academic Achievement. *ShodhAI: Journal of Artificial Intelligence*, 3(1), 37–42. <https://doi.org/10.29121/shodhai.v3.i1.2026.73>
- E. Gures, I. Shayea, and M. Ergen, "Machine learning-based load balancing algorithms in future heterogeneous networks: A survey," *IEEE Access*, vol. 10, pp. 37689–37712, 2022.

- E. Gures, I. Shayea, and M. Ergen, "Machine learning-based load balancing algorithms in future heterogeneous networks: A survey," *IEEE Access*, vol. 10, pp. 37689–37712, 2022.
- E. Puiutta and E. Veith, "Explainable reinforcement learning: A survey," in *Proc. ICML Workshop Explainable Artif. Intell.*, pp. 1–10, 2020.
- E. Puiutta and E. Veith, "Explainable reinforcement learning: A survey," in *Proc. ICML Workshop Explainable Artif. Intell.*, pp. 1–10, 2020.
- E. Puiutta and E. Veith, "Explainable reinforcement learning: A survey," in *Proc. ICML Workshop Explainable Artif. Intell.*, pp. 1–10, 2020.
- E. Puiutta and E. Veith, "Explainable reinforcement learning: A survey," in *Proc. ICML Workshop Explainable Artif. Intell.*, pp. 1–10, 2020.
- E. Strubell, A. Ganesh, and A. McCallum, "Energy and policy considerations for deep learning," in *Proc. 58th Annu. Meeting Assoc. Comput. Linguist. (ACL)*, pp. 3645–3650, 2020.
- E. Strubell, A. Ganesh, and A. McCallum, "Energy and policy considerations for deep learning," in *Proc. 58th Annu. Meeting Assoc. Comput. Linguist. (ACL)*, pp. 3645–3650, 2020.
- E. Strubell, A. Ganesh, and A. McCallum, "Energy and policy considerations for deep learning," in *Proc. 58th Annu. Meeting Assoc. Comput. Linguist. (ACL)*, pp. 3645–3650, 2020.
- E. Strubell, A. Ganesh, and A. McCallum, "Energy and policy considerations for deep learning," in *Proc. 58th Annu. Meeting Assoc. Comput. Linguist. (ACL)*, pp. 3645–3650, 2020.
- E. Strubell, A. Ganesh, and A. McCallum, "Energy and policy considerations for deep learning," in *Proc. 58th Annu. Meeting Assoc. Comput. Linguist. (ACL)*, pp. 3645–3650, 2020.
- G. I. Parisi, R. Kemker, J. L. Part, C. Kanan, and S. Wermter, "Continual lifelong learning with neural networks: A review," *Neural Netw.*, vol. 113, pp. 54–71, 2019.
- G. I. Parisi, R. Kemker, J. L. Part, C. Kanan, and S. Wermter, "Continual lifelong learning with neural networks: A review," *Neural Netw.*, vol. 113, pp. 54–71, 2019.
- G. I. Parisi, R. Kemker, J. L. Part, C. Kanan, and S. Wermter, "Continual lifelong learning with neural networks: A review," *Neural Netw.*, vol. 113, pp. 54–71, 2019.
- J. Muchori and P. Mwangi, "Machine learning load balancing techniques in cloud computing: A review," *Int. J. Comput. Appl. Technol. Res.*, vol. 11, no. 6, pp. 179–186, 2022.
- J. Wu, L. Wang, and Y. Zhou, "Lightweight deep learning for edge–cloud orchestration," *Future Gener. Comput. Syst.*, vol. 144, pp. 547–562, 2023.
- J. Wu, L. Wang, and Y. Zhou, "Lightweight deep learning for edge–cloud orchestration," *Future Gener. Comput. Syst.*, vol. 144, pp. 547–562, 2023.
- J. Wu, L. Wang, and Y. Zhou, "Lightweight deep learning for edge–cloud orchestration," *Future Gener. Comput. Syst.*, vol. 144, pp. 547–562, 2023.
- L. Huang, A. Joseph, B. Nelson, B. Rubinstein, and J. Tygar, "Adversarial machine learning," in *Proc. ACM Workshop Secur. Artif. Intell. (AISec)*, pp. 43–58, 2011.
- L. Huang, A. Joseph, B. Nelson, B. Rubinstein, and J. Tygar, "Adversarial machine learning," in *Proc. ACM Workshop Secur. Artif. Intell. (AISec)*, pp. 43–58, 2011.
- L. Huang, A. Joseph, B. Nelson, B. Rubinstein, and J. Tygar, "Adversarial machine learning," in *Proc. ACM Workshop Secur. Artif. Intell. (AISec)*, pp. 43–58, 2011.
- L. Qian and Y. Zhao, "Federated learning for secure cloud resource optimization," *IEEE Trans. Cloud Comput.*, vol. 12, no. 3, pp. 423–435, 2024.
- L. Qian and Y. Zhao, "Federated learning for secure cloud resource optimization," *IEEE Trans. Cloud Comput.*, vol. 12, no. 3, pp. 423–435, 2024.
- L. Qian and Y. Zhao, "Federated learning for secure cloud resource optimization," *IEEE Trans. Cloud Comput.*, vol. 12, no. 3, pp. 423–435, 2024.
- L. Qian and Y. Zhao, "Federated learning for secure cloud resource optimization," *IEEE Trans. Cloud Comput.*, vol. 12, no. 3, pp. 423–435, 2024.
- M. Alam, S. Shah, and R. Ahmed, "Reliable edge intelligence for IoT and cloud systems," *IEEE Trans. Ind. Informat.*, vol. 19, no. 7, pp. 7212–7225, 2023.
- M. Alam, S. Shah, and R. Ahmed, "Reliable edge intelligence for IoT and cloud systems," *IEEE Trans. Ind. Informat.*, vol. 19, no. 7, pp. 7212–7225, 2023.
- M. Alam, S. Shah, and R. Ahmed, "Reliable edge intelligence for IoT and cloud systems," *IEEE Trans. Ind. Informat.*, vol. 19, no. 7, pp. 7212–7225, 2023.

- M. Hossain, S. Rahman, and M. Karim, "Online learning for adaptive cloud workload management," *ACM Trans. Auton. Adapt. Syst.*, vol. 18, no. 2, pp. 1–24, 2023.
- M. Hossain, S. Rahman, and M. Karim, "Online learning for adaptive cloud workload management," *ACM Trans. Auton. Adapt. Syst.*, vol. 18, no. 2, pp. 1–24, 2023.
- M. Hossain, S. Rahman, and M. Karim, "Online learning for adaptive cloud workload management," *ACM Trans. Auton. Adapt. Syst.*, vol. 18, no. 2, pp. 1–24, 2023.
- M. Hossain, S. Rahman, and M. Karim, "Online learning for adaptive cloud workload management," *ACM Trans. Auton. Adapt. Syst.*, vol. 18, no. 2, pp. 1–24, 2023.
- M. T. Ribeiro, S. Singh, and C. Guestrin, "Why should I trust you? Explaining the predictions of any classifier," in *Proc. 22nd ACM SIGKDD Int. Conf. Knowl. Discovery Data Mining (KDD)*, pp. 1135–1144, 2016.
- M. T. Ribeiro, S. Singh, and C. Guestrin, "Why should I trust you? Explaining the predictions of any classifier," in *Proc. 22nd ACM SIGKDD Int. Conf. Knowl. Discovery Data Mining (KDD)*, pp. 1135–1144, 2016.
- M. T. Ribeiro, S. Singh, and C. Guestrin, "Why should I trust you? Explaining the predictions of any classifier," in *Proc. 22nd ACM SIGKDD Int. Conf. Knowl. Discovery Data Mining (KDD)*, pp. 1135–1144, 2016.
- N. Kumar and R. Goyal, "Unsupervised learning for cloud workload grouping: An efficient load balancing strategy," *Cluster Comput.*, vol. 23, no. 4, pp. 2463–2480, 2020.
- N. Kumar and R. Goyal, "Unsupervised learning for cloud workload grouping: An efficient load balancing strategy," *Cluster Comput.*, vol. 23, no. 4, pp. 2463–2480, 2020.
- Q. Zhang, L. Cheng, and R. Boutaba, "Cloud computing: State-of-the-art and research challenges," *J. Internet Serv. Appl.*, vol. 1, no. 1, pp. 7–18, 2010.
- R. Buyya, J. Broberg, and A. Goscinski, *Cloud Computing: Principles and Paradigms*. Hoboken, NJ, USA: Wiley, 2019.
- R. Buyya, J. Broberg, and A. Goscinski, *Cloud Computing: Principles and Paradigms*. Hoboken, NJ, USA: Wiley, 2019.
- R. Buyya, R. N. Calheiros, and A. V. Dastjerdi, "Intelligent cloud management: Research challenges and opportunities," *IEEE Cloud Comput. Mag.*, vol. 10, no. 1, pp. 8–19, 2023.
- R. Buyya, R. N. Calheiros, and A. V. Dastjerdi, "Intelligent cloud management: Research challenges and opportunities," *IEEE Cloud Comput. Mag.*, vol. 10, no. 1, pp. 8–19, 2023.
- R. Buyya, R. N. Calheiros, and A. V. Dastjerdi, "Intelligent cloud management: Research challenges and opportunities," *IEEE Cloud Comput. Mag.*, vol. 10, no. 1, pp. 8–19, 2023.
- R. Mahmud, S. Srirama, and R. Buyya, "A benchmark framework for AI-driven cloud resource management," *IEEE Access*, vol. 9, pp. 78467–78479, 2021.
- R. Mahmud, S. Srirama, and R. Buyya, "A benchmark framework for AI-driven cloud resource management," *IEEE Access*, vol. 9, pp. 78467–78479, 2021.
- R. Mahmud, S. Srirama, and R. Buyya, "A benchmark framework for AI-driven cloud resource management," *IEEE Access*, vol. 9, pp. 78467–78479, 2021.
- R. Mahmud, S. Srirama, and R. Buyya, "A benchmark framework for AI-driven cloud resource management," *IEEE Access*, vol. 9, pp. 78467–78479, 2021.
- R. Mahmud, S. Srirama, and R. Buyya, "Cloud–fog–edge orchestration using artificial intelligence: Taxonomy and research challenges," *Future Gener. Comput. Syst.*, vol. 111, pp. 300–320, 2020.
- R. Mahmud, S. Srirama, and R. Buyya, "Cloud–fog–edge orchestration using artificial intelligence: Taxonomy and research challenges," *Future Gener. Comput. Syst.*, vol. 111, pp. 300–320, 2020.
- R. Mahmud, S. Srirama, and R. Buyya, "Cloud–fog–edge orchestration using artificial intelligence," *Future Gener. Comput. Syst.*, vol. 111, pp. 300–320, 2020.
- R. S. Sutton and A. G. Barto, *Reinforcement Learning: An Introduction*, 2nd ed. Cambridge, MA, USA: MIT Press, 2018.
- R. S. Sutton and A. G. Barto, *Reinforcement Learning: An Introduction*, 2nd ed. Cambridge, MA, USA: MIT Press, 2018.
- S. M. Lundberg and S. I. Lee, "A unified approach to interpreting model predictions," in *Proc. Adv. Neural Inf. Process. Syst. (NeurIPS)*, pp. 4765–4774, 2017.
- S. M. Lundberg and S. I. Lee, "A unified approach to interpreting model predictions," in *Proc. Adv. Neural Inf. Process. Syst. (NeurIPS)*, pp. 4765–4774, 2017.
- S. M. Lundberg and S. I. Lee, "A unified approach to interpreting model predictions," in *Proc. Adv. Neural Inf. Process. Syst. (NeurIPS)*, pp. 4765–4774, 2017.
- S. M. Lundberg and S. I. Lee, "A unified approach to interpreting model predictions," in *Proc. Adv. Neural Inf. Process. Syst. (NeurIPS)*, pp. 4765–4774, 2017.

- S. Singh and I. Chana, "A survey on resource scheduling in cloud computing: Issues and challenges," *J. Grid Comput.*, vol. 14, no. 2, pp. 217–264, 2016.
- S. Singh and I. Chana, "A survey on resource scheduling in cloud computing: Issues and challenges," *J. Grid Comput.*, vol. 14, no. 2, pp. 217–264, 2016.
- S. Wachter, B. Mittelstadt, and C. Russell, "Counterfactual explanations without opening the black box: Automated decisions and the GDPR," *Harvard J. Law Technol.*, vol. 31, no. 2, pp. 841–887, 2018.
- S. Wachter, B. Mittelstadt, and C. Russell, "Counterfactual explanations without opening the black box: Automated decisions and the GDPR," *Harvard J. Law Technol.*, vol. 31, no. 2, pp. 841–887, 2018.
- S. Wachter, B. Mittelstadt, and C. Russell, "Counterfactual explanations without opening the black box: Automated decisions and the GDPR," *Harvard J. Law Technol.*, vol. 31, no. 2, pp. 841–887, 2018.
- T. Alharbi, M. I. Khan, and A. Hussain, "Scalable machine learning-driven resource allocation in multi-cloud environments," *J. Cloud Comput.*, vol. 12, no. 4, pp. 331–352, 2023.
- T. Alharbi, M. I. Khan, and A. Hussain, "Scalable machine learning-driven resource allocation in multi-cloud environments," *J. Cloud Comput.*, vol. 12, no. 4, pp. 331–352, 2023.
- T. Alharbi, M. I. Khan, and A. Hussain, "Scalable machine learning-driven resource allocation in multi-cloud environments," *J. Cloud Comput.*, vol. 12, no. 4, pp. 331–352, 2023.
- T. Alharbi, M. I. Khan, and A. Hussain, "Scalable machine learning-driven resource allocation in multi-cloud environments," *J. Cloud Comput.*, vol. 12, no. 4, pp. 331–352, 2023.
- T. Miller, "Explanation in artificial intelligence: Insights from the social sciences," *Artif. Intell.*, vol. 267, pp. 1–38, 2019.
- T. Miller, "Explanation in artificial intelligence: Insights from the social sciences," *Artif. Intell.*, vol. 267, pp. 1–38, 2019.
- T. Miller, "Explanation in artificial intelligence: Insights from the social sciences," *Artif. Intell.*, vol. 267, pp. 1–38, 2019.
- T. P. Lillicrap, J. J. Hunt, A. Pritzel, and N. Heess, "Continuous control with deep reinforcement learning," in *Proc. Int. Conf. Learn. Represent. (ICLR)*, 2016.
- T. P. Lillicrap, J. J. Hunt, A. Pritzel, and N. Heess, "Continuous control with deep reinforcement learning," in *Proc. Int. Conf. Learn. Represent. (ICLR)*, 2016.
- Y. Han, C. Park, and H. Kim, "Resource-aware deep learning for cloud scheduling and performance optimization," *IEEE Trans. Netw. Serv. Manage.*, vol. 19, no. 2, pp. 817–829, 2022.
- Y. Han, C. Park, and H. Kim, "Resource-aware deep learning for cloud scheduling and performance optimization," *IEEE Trans. Netw. Serv. Manage.*, vol. 19, no. 2, pp. 817–829, 2022.
- Y. Zhang, X. Li, and H. Chen, "Federated reinforcement learning for distributed cloud management," *IEEE Internet Things J.*, vol. 9, no. 16, pp. 14421–14434, 2022.
- Y. Zhang, X. Li, and H. Chen, "Federated reinforcement learning for distributed cloud management," *IEEE Internet Things J.*, vol. 9, no. 16, pp. 14421–14434, 2022.
- Z. Chen and S. Wu, "Benchmarking machine learning models for cloud scheduling: A systematic evaluation," *Future Gener. Comput. Syst.*, vol. 108, pp. 720–733, 2020.
- Z. Chen and S. Wu, "Benchmarking machine learning models for cloud scheduling: A systematic evaluation," *Future Gener. Comput. Syst.*, vol. 108, pp. 720–733, 2020.
- Z. Li, L. Chen, and H. Liu, "Deep learning-based resource prediction for cloud applications," *IEEE Access*, vol. 9, pp. 145632–145649, 2021.
- Z. Li, L. Chen, and H. Liu, "Deep learning-based resource prediction for cloud applications," *IEEE Access*, vol. 9, pp. 145632–145649, 2021.
- Z. Li, L. Chen, and H. Liu, "Deep learning-based resource prediction for cloud applications," *IEEE Access*, vol. 9, pp. 145632–145649, 2021.
- Z. Li, L. Chen, and H. Liu, "Deep learning-based resource prediction for cloud applications," *IEEE Access*, vol. 9, pp. 145632–145649, 2021.
- Z. Xiong, Y. Li, and F. Xu, "Blockchain for secure federated learning in cloud environments," *IEEE Trans. Services Comput.*, vol. 15, no. 5, pp. 2949–2961, 2022.
- Z. Xiong, Y. Li, and F. Xu, "Blockchain for secure federated learning in cloud environments," *IEEE Trans. Services Comput.*, vol. 15, no. 5, pp. 2949–2961, 2022.
- Z. Xiong, Y. Li, and F. Xu, "Blockchain for secure federated learning in cloud environments," *IEEE Trans. Services Comput.*, vol. 15, no. 5, pp. 2949–2961, 2022.

- Z. Xiong, Y. Li, and F. Xu, "Blockchain for secure federated learning in cloud environments," *IEEE Trans. Services Comput.*, vol. 15, no. 5, pp. 2949–2961, 2022.
- Zhou, Z., Abawajy, J., Chowdhury, M., Hu, Z., Li, K., Cheng, H., ... & Li, F. (2018). Minimizing SLA violation and power consumption in Cloud data centers using adaptive energy-aware algorithms. *Future Generation Computer Systems*, 86, 836-850.