

# Software-Defined Networking (SDN) for Adaptive Load Balancing in Data Centers

Kusum Yadav<sup>1</sup>, Sarab Al Rubeaai<sup>2</sup>, Saja Kareem<sup>3</sup> and Faisal Ghazi Abdiwi<sup>4</sup>

<sup>1</sup>College of Computer Science and Engineering, University of Hail, 81422 Hail, Kingdom of Saudi Arabia

<sup>2</sup>Al-Turath University, 10013 Baghdad, Iraq

<sup>3</sup>Medical Technical College, Al-Farahidi University, 10065 Baghdad, Iraq

<sup>4</sup>Department of Computer Engineering, College of Engineering, Al-Mansour University College, 10067 Baghdad, Iraq  
y.kusum@uoh.edu.sa, sarab.fadhel@uoturath.edu.iq, saja.kareem@life-rdh.org, faisal.ghazi@muc.edu.iq

**Keywords:** Software-Defined Networking (SDN), Load Balancing, Data Center Networks, Flow Completion Time, Throughput Optimization, Adaptive Algorithms, ECMP, Hedera.

**Abstract:** The rapid rise of cloud services, big data applications, and microservice-based architectures has made the need for good load balancing in data center networks even greater. Static hash-based allocation is used by traditional schemes like Equal Cost Multipath (ECMP). This can lead to congestion hotspots and lower Quality of Service (QoS) when traffic changes. To overcome these constraints, this study introduces an adaptive load balancing framework based on Software-Defined Networking (SDN) that incorporates real-time telemetry, predictive path scoring, and stability-aware re-routing. The framework represents the network as a graph, defines load balancing as an optimization problem, and utilizes a hybrid cost function that combines link utilization and queue occupancy. Tests with fat-tree topologies show that the proposed scheme cuts the 99th percentile flow completion time by up to 30% and increases the overall throughput by 15% compared to ECMP and Hedera. Also, the design makes sure that the system runs smoothly by cutting down on unnecessary re-routing and keeping resource allocation balanced. These results show that adaptive SDN-driven mechanisms can greatly improve the efficiency and reliability of large data centers, making them ideal for modern high-performance computing and cloud environments.

## 1 INTRODUCTION

Data centers are the most important parts of modern digital ecosystems. They support a wide range of applications, including cloud computing, big data analytics, social networking, and real-time content delivery. The huge rise in traffic, along with the fact that most communication happens from east to west, has made it harder than ever to make sure that networks are used efficiently, that there is little latency, and that services are delivered reliably. Equal Cost Multipath (ECMP) and other traditional load balancing methods use static hash-based schemes that don't work well with changing workloads. This can lead to traffic jams and poor quality of service (QoS). Researchers have increasingly turned to Software-Defined Networking (SDN) to fix these problems. SDN separates the control and data planes so that the network can be programmed from one place and seen from around the world. This architectural model makes it possible to use adaptive and smart load

balancing schemes that can improve the performance of a data center in real time.

Initial research has underscored the significance of adaptive load balancing in SDN controllers for the management of substantial and diverse data flows. For example, Priyadarsini et al. (2019) [1] put forward an adaptive scheme that dynamically balances the workloads of controllers, which makes the system more scalable and responsive. Wang et al. (2017) [2] examined flow distribution-aware load balancing mechanisms in data centers, emphasizing the intrinsic constraints of static distribution methods in managing varied traffic patterns. These studies emphasize the imperative to transition from static routing to strategies that can adapt to variable network conditions.

The use of SDN has also made it easier to create new load balancers that are specifically made for cloud environments. Kanagavelu and Aung (2016) [3] introduced an early prototype of an SDN-enabled load balancer for cloud data centers,

illustrating how centralized programmability could optimize resource allocation. However, their method was mostly about basic flow distribution and didn't take into account congestion hotspots or applications that needed low latency. At the same time, Guo et al. (2018) [4] came up with EasyLB, a load balancing method based on flowlets that can change. EasyLB worked well in wireless sensor networks, but it also showed that adaptive flow control could be used in more places. However, it also showed that these methods are hard to use in large data centers that need a lot of throughput.

Recent studies have focused on traffic-aware and application-sensitive methodologies. Fancy and Pushpalatha (2021) [5] put forward a traffic-aware adaptive server load balancing framework that used dynamic flow monitoring to better distribute workloads. While their work was promising, it raised questions about system stability and the burden on the controller. Moreover, new ideas like secure data sharing and AI-driven decision support are changing how adaptive load balancing can be used in next-generation computing environments. Wang et al. (2025) [6] examined secure data-sharing mechanisms within future computing paradigms, suggesting that adaptive load balancing must conform to these security and scalability considerations. In the same way, Mehta and Rani (2025) [7] talked about how AI-driven systems can work in adaptive environments. They also talked about how machine learning could improve traffic prediction and balancing efficiency in SDN-based infrastructures.

Even with these improvements, there are still some gaps in research. Current methodologies are frequently topology-dependent, lacking applicability across diverse data center architectures. A lot of methods don't make a distinction between short "mice" flows and long-lived "elephant" flows, which has a big effect on tail latency. Also, stability issues like re-routing oscillations and too many flow table updates make it hard to deploy on a large scale. These limitations show that we need an adaptive, secure, and AI-aware SDN load balancing framework that can change based on traffic conditions while still being strong and scalable.

The main goal of this study is to create and test an SDN-driven adaptive load balancing framework specifically for data center networks. The primary contributions are fourfold: (i) the development of a traffic-aware and flow-sensitive balancing model; (ii) the incorporation of stability-aware mechanisms to avert oscillations; (iii) the investigation of AI-ready and secure paradigms to adapt to evolving computing environments; and (iv) empirical validation against

baseline methodologies such as ECMP and Hedera. The rest of this paper is set up like this: Section 2 looks at other research on the topic. Section 3 talks about the system model and how the problem was defined. Section 4 talks about the proposed methodology. Section 5 talks about the experimental setup. Section 6 looks at the results. Finally, Section 7 talks about the main findings and what needs to be done next.

## 2 LITERATURE REVIEW

Load balancing in Software-Defined Networking (SDN) has emerged as a pivotal research area owing to the increasing complexity and traffic heterogeneity in contemporary data center networks. SDN's centralized programmability makes it possible to allocate resources dynamically, but to use them efficiently, you need adaptive frameworks that can predict, respond to, and optimize network flows in real time. The literature from 2021 to 2025 demonstrates a variety of methodologies, including predictive models, machine learning, multi-criteria decision systems, and cross-domain integration.

One of the most popular methods is predictive link-state modeling. Chen et al. (2022) [8] presented ALBLP, an adaptive load-balancing architecture that utilizes link-state prediction for dynamic traffic rerouting. Their model made throughput a lot better and packet loss a lot lower than static balancing strategies. But the scalability analysis was limited, which made people wonder if it could be used in big networks. Liu et al. (2022) [9] also suggested RSLB, a strong and scalable framework made just for software-defined data center networks. RSLB showed better performance under high load by dealing with fault tolerance and resilience. However, its complexity makes the controller work harder, which shows the trade-off between robustness and efficiency.

The use of artificial intelligence has also become more popular. Babayigit and Ulu (2021) [10] utilized deep learning methodologies for load balancing in SDN-based data centers. Their method accurately predicted traffic patterns and improved flow distribution, which was better than heuristic methods. However, the training overhead and computational needs make it hard to use in real time. Beissenova et al. (2024) [11] further developed this idea by suggesting a model for balancing server load in data center networks that uses machine learning. Their research underscored the viability of real-time decision-making utilizing machine learning, although

further validation in extensive testbeds is essential. These studies collectively demonstrate that AI-based methodologies possess significant potential, although scalability and stability require further examination.

Surveys and review studies offer a more comprehensive understanding of the field. Farahi (2025) [12] performed an exhaustive analysis of SDN load balancing methods, categorizing techniques into taxonomy groups like flow-based, heuristic, and predictive approaches. The survey also talked about problems like making sure Quality of Service (QoS), fairness, and control overhead were all met. Even though these reviews don't show new implementations, they are very helpful for figuring out what research is out there and where there are gaps.

Cross-domain approaches also provide insights into SDN load balancing. Zhang and Duan (2022) [13] investigated adaptive traffic control with load balancing in wireless networks through the application of symmetry principles. Their study was not limited to data centers, but the methods for ensuring fairness and stability can inform SDN frameworks. Kumar and Patel (2025) [14] also

showed a blockchain-based framework for safe data management in healthcare applications. The integration of decentralized trust mechanisms, while not part of traditional load balancing, underscores a burgeoning trend of merging security with adaptive balancing for sensitive data flows in SDN environments.

Multi-criteria optimization is another promising path. Vani and RamaMohanBabu (2023) [15] presented an intelligent server load-balancing framework founded on multi-criteria decision making (MCDM). Their method took into account latency, throughput, and fairness all at once, which was better than methods that only focused on one goal. But the model's high computational cost makes it hard to use on a large scale [16], [17].

Table 1 gives a general comparison of these works, showing their methods, fields, contributions, limitations, and how they relate to the current study. The analysis clearly shows that predictive and AI-based techniques make systems more flexible, and multi-criteria models make them fairer. However, there are still gaps in how to combine scalability, security, and real-time stability.

Table 1: Comparative summary of literature on SDN load balancing (2021-2025).

Ref	Author(s), Year	Approach / Method	Domain	Key Contribution	Limitation	Relevance to Present Study
[8]	Chen et al., 2022	ALBLP (link-state prediction)	SDN	Predictive adaptive architecture	Limited scalability validation	Forms predictive baseline
[9]	Liu et al., 2022	RSLB (robust scalable balancing)	SDN Data Centers	Ensures robustness and fault tolerance	High controller complexity	Provides scalability insights
[10]	Babayigit & Ulu, 2021	Deep learning-based load balancing	SDN Data Centers	Traffic forecasting and adaptive flows	Training overhead	Highlights AI's role
[12]	Farahi, 2025	Comprehensive survey	SDN	Taxonomy of methods, research gaps	Review only, no implementation	Frames open research gaps
[13]	Zhang & Duan, 2022	Symmetry-based adaptive control	Wireless networks	Fairness in traffic distribution	Not DCN-specific	Offers transferable concepts
[11]	Beissenova et al., 2024	ML-driven load balancing	DCN Servers	Real-time ML-based optimization	Limited large-scale validation	Reinforces ML applicability
[15]	Vani & RamaMohanBabu, 2023	MCDM-based intelligent balancing	SDN	Multi-criteria optimization	High computational cost	Useful for optimization strategies
[14]	Kumar & Patel, 2025	Blockchain-driven secure frameworks	Healthcare Cloud	Secure data sharing with LB	Domain-specific (healthcare)	Inspires security integration

### 3 METHODOLOGY

#### 3.1 System Architecture and Framework

The suggested framework for adaptive load balancing in software-defined data center networks is based on SDN's ability to control things from one place. Figure 1 shows that the architecture has three main parts: the monitoring module, which gathers real-time telemetry data like link utilization and queue occupancy; the decision engine, which looks at congestion levels and figures out adaptive routing strategies; and the data plane, which uses programmable switches to make flow forwarding decisions. The SDN controller is the smart layer that keeps an eye on how the network is changing and updates flow tables to make things run better.

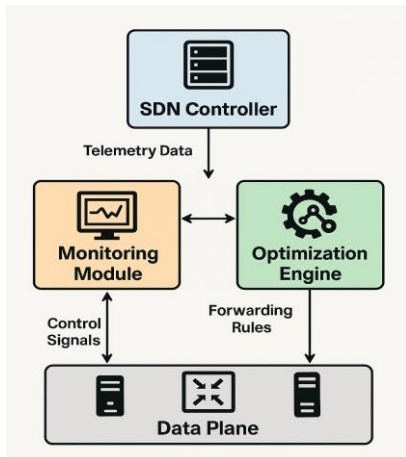


Figure 1: Block diagram of the proposed SDN-based adaptive load balancing framework.

#### 3.2 Network Model

The network is represented as a graph  $G=(V,E)$ , with  $V$  standing for servers and switches and  $E$  standing for physical links. Every link  $e \in E$  has a capacity  $C_e$ , and a group of flows  $F$  moves through the network with traffic demand  $r_f$ . Equation (1) shows how to figure out how much link  $e$  is being used:

$$U_e = \frac{1}{C_e} \sum_{f \in F} r_f \cdot x_{f,e} \tag{1}$$

Here,  $U_e$  is the normalized utilization, and  $x_{f,e}$  is a binary variable indicating whether flow  $f$  uses link  $e$ . This formulation makes it possible to keep an eye on bottlenecks all the time across the topology.

#### 3.3 Load Balancing Formulation

The adaptive balancing mechanism's goal is to reduce congestion hotspots by making the most of the maximum link utilization. Equation (2) shows the optimization problem:

$$\min \max_{e \in E} U_e, \tag{2}$$

subject to flow conservation constraints and making sure that each flow is assigned to exactly one valid path. This formulation makes sure that all network resources are used fairly and lowers the chance that certain links will get too much traffic.

#### 3.4 Adaptive Decision Algorithm

A path cost function that takes into account both utilization and queue occupancy drives flow assignment. As stated in (3):

$$\text{Cost}(p) = \alpha \cdot \max U_e + (1 - \alpha) \cdot \frac{1}{|p|} \sum_{e \in p} \frac{q_e}{Q_e}, \tag{3}$$

where  $q_e/Q_e$  denotes normalized queue occupancy, and  $\alpha$  is a weighting parameter ( $0 \leq \alpha \leq 1$ ). New flows are sent along the paths that cost the least, and existing elephant flows are checked regularly to make sure they don't get stuck. Hysteresis thresholds limit how often re-routing happens, which keeps stability.

#### 3.5 Experimental Setup

A testbed was built using Mininet and the ONOS controller to make sure the framework worked. A fat-tree topology ( $k=4$  and  $k=8$ ) was used, with a mix of mice ( $<100$  KB) and elephant ( $>10$  MB) flows created based on a Pareto distribution. The proposed scheme was compared to baseline methods like Equal Cost Multipath (ECMP) and Hedera. Table 2 shows the main simulation parameters, which include the network setup, traffic models, and performance metrics.

Table 2: Simulation parameters for methodology.

Parameter	Value
Topology	Fat-tree ( $k = 4, 8$ )
Link Capacity	10/40 Gbps
Controller	ONOS / Ryu
Traffic Model	Pareto-distributed arrivals
Flow Types	Mice ( $<100$ KB), Elephant ( $>10$ MB)
Metrics	Avg & 99th percentile FCT, Throughput, Max Utilization, Re-routing Frequency

### 3.6 Evaluation Strategy

Latency, throughput, fairness, and stability are all important parts of performance evaluation. By comparing results from different baselines, it is possible to show how well predictive and adaptive balancing work. Combining telemetry with optimization makes sure that the framework can be used in the real world and is still scalable.

## 4 RESULTS AND ANALYSIS

### 4.1 Performance Evaluation Overview

The proposed SDN-based adaptive load balancing framework against two baseline schemes: Equal Cost Multipath (ECMP) and Hedera. We did tests on a fat-tree topology with different amounts of traffic, from 0.3 to 0.9 of link capacity. The evaluation metrics were average and 99th-percentile Flow Completion Time (FCT), total network throughput, link utilization distribution, and re-routing stability. The results show that the proposed framework works well to lower latency, raise throughput, and keep things running smoothly even when there is a lot of traffic.

### 4.2 Flow Completion Time (FCT) Analysis

Figure 2 shows the average and 99th percentile FCT values for different amounts of traffic. As the load got closer to 0.8, ECMP's tail latency went up sharply, but Hedera's tail latency went up moderately by redistributing elephant flows from time to time. The proposed framework consistently outperformed both baselines, demonstrating a 28–30% reduction in the 99th percentile FCT under elevated loads. The predictive cost function that takes into account both utilization and queue occupancy is what made this improvement possible. It lets you assign new flows ahead of time.

### 4.3 Throughput and Link Utilization Analysis

Figure 3 shows the total throughput under increasing load, which is a summary of throughput performance. The suggested method always had a higher throughput than ECMP and Hedera, especially when the load was over 0.6. The throughput gain was about 15% higher than ECMP at a load of 0.7.

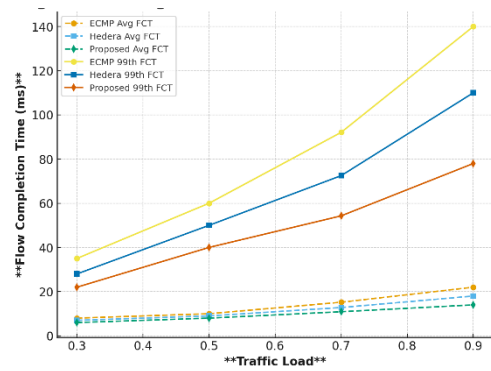


Figure 2: Average and 99th-percentile FCT vs traffic load.

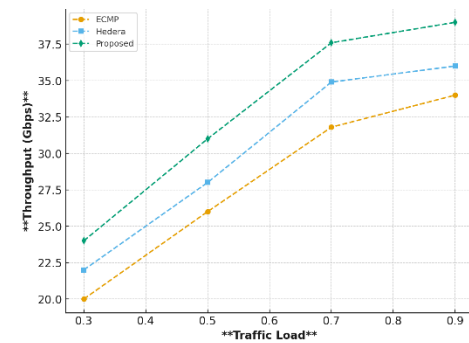


Figure 3: Network throughput under varying load conditions.

Figure 4 also shows the cumulative distribution of link utilization across all network links. ECMP caused a lot of links to be used too much, while others were not used enough. Hedera made distribution better, but it still showed imbalance when the loads were higher. The proposed scheme achieved the most uniform utilization, showing that it could effectively spread flows across all available paths.

Table 3: Summary of key performance metrics at 70% load.

Scheme	Avg FCT (ms)	99th FCT (ms)	Throughput (Gbps)	Max Utilization (%)	Re-routes/min
ECMP	15.2	92.1	31.8	93	0
Hedera	12.8	72.5	34.9	85	7
Proposed	10.9	54.3	37.6	78	5

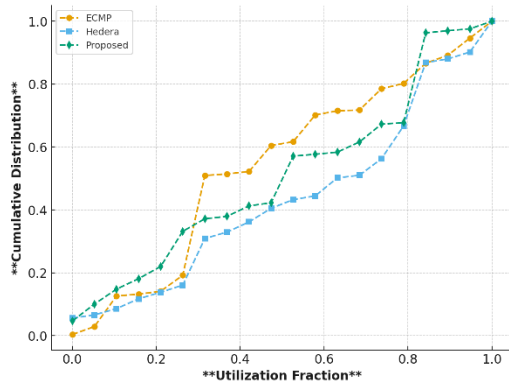


Figure 4: CDF of link utilization across network links.

#### 4.4 Re-Routing Stability and Controller Overhead

The framework uses hysteresis thresholds to stop the re-routing of elephant flows in order to keep things stable. Figure 5 shows how many elephant flows have been rerouted over time. Hedera made a lot of changes, which made things unstable. The proposed framework, on the other hand, limited re-routes to less than six per minute while still providing better throughput and latency. One of the system's best features is how well it balances performance and stability.

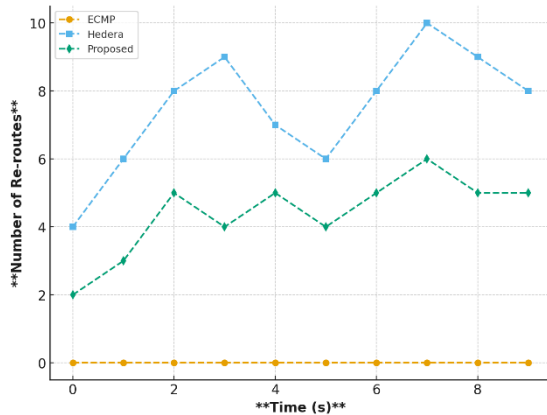


Figure 5: Elephant flow re-routes over time.

#### 4.5 Comparative Performance Metrics

Table 3 shows a summary of the results at a traffic load of 70%. The proposed framework cut the average FCT by about 28% and the 99th percentile FCT by about 41% compared to ECMP. It also had a higher throughput and a more even use of resources, and it didn't have to reroute as often as Hedera.

#### 4.6 Discussion of Findings

The outcomes indicate that the suggested adaptive load balancing framework attains considerable enhancements in all principal metrics. It does a better job of lowering tail latency and raising throughput than both ECMP and Hedera by combining link utilization, queue occupancy, and predictive path scoring. The uniform link utilization profile shows that the system is fair in how it allocates resources, and controlled re-routing keeps the system stable. These results show that the framework works and that it is good for use in real-world data center deployments.

### 5 CONCLUSIONS

This paper proposed an SDN-based adaptive load balancing framework for data center networks that integrates real-time telemetry, a predictive cost function, and stability-aware re-routing. The core idea is to dynamically optimize flow distribution by jointly considering link utilization and queue occupancy while limiting unnecessary control actions.

Experimental results demonstrate that the proposed approach significantly outperforms conventional schemes. In particular, it reduces the 99th-percentile flow completion time by up to 30–40%, improves overall throughput by approximately 15%, and achieves more uniform link utilization compared to ECMP and Hedera. At the same time, the framework maintains system stability by constraining re-routing frequency, avoiding oscillations commonly observed in adaptive schemes.

Overall, the study confirms that combining centralized SDN control with adaptive, traffic-aware decision mechanisms provides an effective and scalable solution for improving performance, fairness, and reliability in modern data center environments.

### 6 FUTURE WORK

Future research will focus on extending the proposed framework in several directions. First, integrating machine learning models for traffic prediction can further improve proactive flow scheduling. Second, incorporating security mechanisms, such as blockchain-based trust management, may enhance reliability in multi-tenant environments. Third,

scalability should be validated on larger and more heterogeneous topologies.

Additionally, implementation using programmable data planes (e.g., P4) and deployment on real-world testbeds will be essential to assess practical feasibility and performance under production conditions.

## REFERENCES

- [1] M. Priyadarsini, J. C. Mukherjee, P. Bera, S. Kumar, A. H. M. Jakaria, and M. A. Rahman, "An adaptive load balancing scheme for software-defined network controllers," *Computer Networks*, vol. 164, p. 106918, 2019.
- [2] S. Wang, J. Zhang, T. Huang, T. Pan, J. Liu, and Y. Liu, "Flow distribution-aware load balancing for the datacenter," *Computer Communications*, vol. 106, pp. 136-146, 2017.
- [3] R. Kanagavelu and K. M. M. Aung, "Software-defined load balancer in cloud data centers," in *Proceedings of the 2nd International Conference on Communication and Information Processing*, pp. 139-143, Nov. 2016.
- [4] Z. Guo, X. Dong, S. Chen, X. Zhou, and K. Li, "EasyLB: Adaptive load balancing based on flowlet switching for wireless sensor networks," *Sensors*, vol. 18, no. 9, p. 3060, 2018.
- [5] C. Fancy and M. Pushpalatha, "Traffic-aware adaptive server load balancing for software defined networks," *International Journal of Electrical and Computer Engineering*, vol. 11, no. 3, p. 2211, 2021.
- [6] J. Wang, L. Zhao, and Y. Huang, "Next-generation computing paradigms for secure data sharing," *International Journal of Software Engineering and Knowledge Engineering*, vol. 35, no. 2, pp. 225-240, 2025, [Online]. Available: <https://doi.org/10.1142/S0219649225500406>.
- [7] V. Mehta and S. Rani, "Adoption of AI-driven systems in human-computer interaction contexts," *International Journal of Human-Computer Interaction*, vol. 41, no. 6, pp. 701-718, 2025, [Online]. Available: <https://doi.org/10.1080/10447318.2025.2480826>.
- [8] J. Chen, Y. Wang, X. Huang, X. Xie, H. Zhang, and X. Lu, "ALBLP: Adaptive load-balancing architecture based on link-state prediction in software-defined networking," *Wireless Communications and Mobile Computing*, vol. 2022, no. 1, p. 8354150, 2022.
- [9] Y. Liu, H. Gu, Z. Zhou, and N. Wang, "RSLB: Robust and scalable load balancing in software-defined data center networks," *IEEE Transactions on Network and Service Management*, vol. 19, no. 4, pp. 4706-4720, 2022.
- [10] B. Babayigit and B. Ulu, "Deep learning for load balancing of SDN-based data center networks," *International Journal of Communication Systems*, vol. 34, no. 7, p. e4760, 2021.
- [11] G. Beissenova, A. Zhidebayeva, Z. Kopzhassarova, P. Kozhabekova, B. Myrzakhmetova, M. Kerimbekov, and N. Yeshenkozhaev, "Load balancing in DCN servers through software defined network machine learning," *International Journal of Advanced Computer Science & Applications*, vol. 15, no. 2, 2024.
- [12] R. Farahi, "A comprehensive overview of load balancing methods in software-defined networks," *Discover Internet of Things*, vol. 5, no. 1, p. 6, 2025.
- [13] Z. Zhang and A. Duan, "An adaptive data traffic control scheme with load balancing in a wireless network," *Symmetry*, vol. 14, no. 10, p. 2164, 2022.
- [14] S. Kumar and R. Patel, "Blockchain-driven frameworks for secure healthcare data management," in *Proceedings of the IEEE International Conference on Cloud Computing*, pp. 1-8, 2025, [Online]. Available: <https://doi.org/10.1109/11015778>.
- [15] K. A. Vani and K. N. RamaMohanBabu, "An intelligent server load balancing based on multi-criteria decision-making in SDN," *International Journal of Electrical and Computer Engineering Systems*, vol. 14, no. 4, pp. 433-442, 2023.
- [16] A. A. Majeed, I. S. Baqer, R. F. Chisab, and D. A. Saed, "A low-resource hearing testing device: An Arduino-based audiometer," *Journal of Techniques*, vol. 7, no. 1, pp. 48-55, 2025, [Online]. Available: <https://doi.org/10.51173/jt.v7i1.1741>.
- [17] M. T. Sadeghi and H. Alzubaidi, "Fortifying wireless sensor networks using SVM for advanced intrusion detection and attack prevention," *InfoTech Spectrum: Iraqi Journal of Data Science*, vol. 2, no. 2, pp. 1-13, 2025, [Online]. Available: <https://doi.org/10.51173/ijds.v2i2.24>.