

# Hybrid Model Deployment: Balancing Edge and Cloud Computation

<sup>1</sup>Dheeraj Vaddepally

<sup>1</sup>Independent Researcher, USA

Received: 13<sup>th</sup> Nov 2025 | Received Revised Version: 18<sup>th</sup> Dec 2025 | Accepted: 15<sup>th</sup> Jan 2026 | Published: 29<sup>th</sup> Jan 2026

Volume 08 Issue 01 2026 | Crossref DOI: [10.37547/tajet/v8i1-317](https://doi.org/10.37547/tajet/v8i1-317)

## Abstract

*During the past couple of years, rapid development of edge computing and cloud technologies has allowed for the implementation of hybrid models that offload computation to edge devices as well as cloud platforms. This article explores the design decisions in the implementation of hybrid models, specifically focusing on offloading processing to the cloud, and addresses the necessary trade-offs between latency and privacy. We begin by contrasting edge and cloud computing, highlighting the advantages of hybrid systems in enhancing scalability, real-time processing, and flexibility. Significant architectural concerns, such as model partitioning and offloading rules, are addressed in the context of the dynamic nature of edge and cloud environments. Latency is a fundamental concern that influences the effectiveness of hybrid systems, especially in applications involving real-time processing. We explore how to minimize latency through edge caching, adaptive algorithms, and local computation for enhanced system performance. Privacy becomes an issue when handling sensitive data on the edge and the cloud. In this paper, we present privacy-preserving mechanisms, such as data anonymization, encryption, and federated learning, to secure user information while leveraging the computational power of the cloud. By performance metric evaluation, such as latency, precision, and scalability, we compare hybrid model deployment with cloud-only and edge-only deployment. We concluded the paper by outlining challenges experienced in hybrid deployment, including network limitations and model complexity, and introduce future work ideas on further enhancing edge to cloud computation balance. This paper offers a thorough examination of deploying hybrid models and offers real-world architectural advice on how to maximize system performance without exacerbating latency and privacy concerns.*

**Keywords:** hybrid model deployment, edge computing, cloud computing, offloading processing, latency optimization, privacy-preserving techniques, model partitioning, real-time processing, federated learning, scalability, architectural decisions.

© 2026 Dheeraj Vaddepally. This work is licensed under a Creative Commons Attribution 4.0 International License (CC BY 4.0). The authors retain copyright and allow others to share, adapt, or redistribute the work with proper attribution.

**Cite This Article:** Vaddepally, D. (2026). Hybrid model deployment: Balancing edge and cloud computation. The American Journal of Engineering and Technology, 8(1), 178–185. <https://doi.org/10.37547/tajet/v8i1-317>.

## 1. Introduction

The growing adoption of Internet of Things (IoT) devices, autonomous systems, and machine learning models transformed the landscape of contemporary computing. With more distributed data generation and more real-time processing requirements, there is a need to innovate the deployment of machine learning models in heterogeneous computing environments. [1] Deployment of hybrid models that leverage both edge and cloud computing has

been suggested as a solution to meet such requirements. In this regard, computational burdens are intentionally offloaded between edge devices (e.g., smartphones, sensors, or IoT concentrators) and cloud infrastructure in order to achieve greater flexibility, scalability, and efficiency. Additional need for balance between edge and cloud computation arises to achieve some of the requirements of applications, e.g., real-time decision-

making, data privacy, and network efficiency. Edge computing is most suitable for latency-sensitive workloads, whereas cloud computing has virtually unlimited computational capacity under load. The hybrid approach enables enterprises to optimize both and have the best of both with balanced performance, providing bandwidth support, power limits, and hardware. Architectural selection becomes an important determinant in determining the compute task balance between edge and cloud. [1]

Cloud offloading processing has the potential to support computationally intensive operations beyond the capacity of edge devices but is accompanied by latency, data privacy, and security concerns. For the majority of real-time systems—autonomous vehicles, smart cities, and healthcare—latency minimization is critical to enable timely response. Cloud processing, however, as a single solution has the potential to introduce latency by virtue of network transmission and bandwidth limitations. [2] Besides that, privacy concerns arise with sensitive data processed or stored in edge devices to be communicated to the cloud for computation. Compliance frameworks such as GDPR make securing people's data a top priority, and thus ensuring privacy is an essential consideration for offloading for cloud computation. All these performance, latency, and privacy concerns should be well traded off while creating the hybrid model structures. This research is focused on the design of offloading machine learning models to edge and cloud environments.

Specifically, it considers offloading criteria to the cloud, including model partitioning methods and decision-making rules. The research also considers latency factors, analyzing mechanisms to reduce delay in hybrid systems, as well as privacy factors, analyzing mechanisms to protect sensitive data when offloading to the cloud. In this research, the goal of this paper is to present insights towards optimizing hybrid deployment for security and performance purposes.

## 2. Hybrid Model Deployment: An Overview

Hybrid model deployment includes the assignment of computational workload on both edge computing and

cloud computing systems utilizing the strengths of both models. Edge computing refers to the processing of data within devices that are geographically close to where the data are being generated, i.e., IoT sensors, smartphones, or other networked appliances. By processing data at the "edge" of the network, this approach reduces the quantity of data that must be transmitted over long distances to central servers, thereby reducing latency. Edge computing is ideally suited to latency-intensive applications like real-time monitoring, autonomous cars, and augmented reality, where instant processing and feedback are crucial. [3] However, edge devices are generally limited by processing capacity and storage space, which could restrict their ability to undertake more computationally intensive operations.

By contrast, cloud computing uses near-infinite processing resources through vast data centers. Data from edge devices can be forwarded to the cloud for advanced processing, analysis, or storage. Cloud provides elastic infrastructure to handle massive workloads, run complex machine learning models, and store massive datasets. The cloud can be endowed with immense benefits of processing, but the biggest drawback is the latency introduced due to the fact that data have to be conveyed over the internet, an issue that is a hindrance to applications involving real-time decision-making. In addition, cloud processing introduces latency problems when sensitive data are conveyed from edge devices to cloud servers. [4] Hybrid model combines the advantages of both edge and cloud computing by balancing computation between both environments. Deploying in a hybrid model divides computation such that some computations are computed locally within edge devices so as to maintain low latency while data is retained as private information, whereas computationally intensive computations are uploaded to the cloud. This can be achieved with variability in the balance between real-time processing needs and the capability of the cloud to perform high-level, high-complexity computation. For example, an autonomous vehicle can perform sensor data processing locally in real time to support rapid driving decisions and reserve intensive data analysis or machine learning model updates for the cloud. [5]

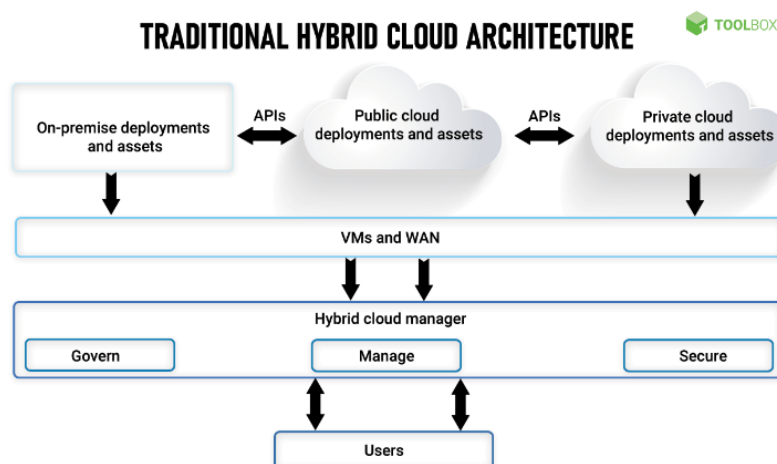


Fig. 1. Hybrid Cloud Architecture

Hybrid systems possess numerous advantages, particularly in scalability, efficiency, and real-time responsiveness. By the giving of workload selective offloading to the cloud, hybrid architectures can scale easily against changing workloads and use computing power efficiently. Hybrid systems also offer the type of adaptability real-time decision-making applications need to do local computation of sensitive data so that latency isn't experienced while still being able to leverage the scalable computer power of the cloud when needed. Other than that, hybrid deployment facilitates system efficiency through optimization of the usage of bandwidth so that transfer to the cloud would only occur for most critical information with less critical processes still left behind on the edge. This type of mechanism allows better privacy in that sensitive data will be processed in a local node and would not be transferred over to central platforms. [4] Generally, hybrid systems are offering a solid and scalable solution that is compromising on performance, latency, and privacy requirements across various applications.

### 3. Architectural Decisions for Offloading Processing

In hybrid model offloading in cloud, computation offloading to cloud must be accompanied by cautious design decisions. Those decisions are made based on varying parameters like the volume of data, model dimensions, network latency, and energy consumption. Proper balance between cloud and edge processing must be realized for effective system performance, confidentiality, and energy saving.

#### 3.1 Offloading Criteria

There are a number of factors that come into play when offloading computation work to the cloud compared to doing it locally on edge devices. One of the key ones is data large data can be expensive or difficult to transfer through the network, and so local processing to conserve bandwidth would be desirable. [6] Low latency tolerance small data is easily dealt with by the cloud where computation power exists to leverage for additional analysis.

Model complexity is also something to be worried about. Intricate models such as deep neural networks will require more computational power and memory than edge devices can support. Such computationally demanding operations on the nature of model training or inference, for example, will become more efficient by cloud outsourcing. Such less demanding computation such as pre-filtering raw data or merely simple feature extraction would be best left to edge devices.

Energy consumption is also a consideration, particularly in battery-powered devices such as IoT sensors or smartphones. Computation offloading reduces energy consumed in edge devices but at the expense of energy consumed to send data and wait for the reply. Effective offloading policies try to strike a balance between computation load and energy consumption and reduce processing time.

#### 3.2 Model Partitioning

Hybrid model deployment is one of the most significant architecture decisions and model partitioning—the

determination of how much of a model goes on the edge and how much goes to the cloud—is one of the most effective architecture decisions. It must be done in an intelligent manner such that operations necessary for real-time processing, such as feature extraction, are executed on the edge and computationally intensive operations such as inference or training of a model go to the cloud to process. In alternative hybrid architectures, computation is partitioned in the sense that the front-end part of the model executes on the edge and passes the mid-way result to the cloud for further processing. This can reduce data communication cost by transmitting relevant processed data instead of raw data. For instance, in computer vision tasks like image classification, feature extraction can be performed at the edge and leaving the computationally heavier task of final classification being run on the cloud. Models may be divided thoughtfully and systems can then experience a balance among energy consumption, latency, and computation overhead.

### **3.3 Decision Making Algorithms**

The degree and timing of the offloading to the cloud are managed by choice algorithms that differ in strategy as well as levels of sophistication. Heuristic algorithms use pre-specified pre-determined thresholds and rules, i.e., the availability of energy on devices or the status of network bandwidth, to determine offloading initiation. They are easy to use and perhaps do not learn in reactive manners to disparate conditions on-the-fly.

More complex ones are optimization-based algorithms, which try to reduce a cost function that has been developed, like latency or energy consumed, by making adaptive adjustments on offloading decisions. They are better suited for hybrid systems where network conditions or computation resources continuously change.[7]

Machine learning techniques have previously been used to develop more complex offloading decisions. Historical information can be utilized by machine learning algorithms to predict system behavior under different circumstances and make offloading decisions in advance that result in optimal efficiency in the long term. For example, a machine learning algorithm can predict when network congestion is most likely to occur and adjust the offloading strategy such that critical tasks are run locally to avoid delays.

### **3.4 Case Study**

A few examples of real-world application use cases reflect the power of hybrid model deployment patterns. Autonomous environmental monitoring drones are one such case. In such systems, edge-capable drones can perform video streams and detect important features (e.g., object detection or anomalies) in real-time. The more computationally intensive work of constructing predictive models from the data is delegated to cloud servers, where machine learning algorithms operate on the data and inform the drone of an update when it occurs. This hybrid design allows the drones to be useful without sacrificing real-time responsiveness. [8]

For example, in smart healthcare systems where patient vital signs such as blood pressure, oxygen saturation, and heart rate are monitored through wearable sensors. In these devices, edge processing is employed to process information locally and transmit only significant measurements or alerts to the cloud infrastructure for further processing. Offloading the complex health prognosis to the cloud and keeping life-critical real-time monitoring on-site, the medical personnel are able to respond in a timely fashion for emergency situations without taxing the system with duplicate data transfer.

## **4. Latency Considerations in Hybrid Model Deployment**

Latency, or the duration between data collection, processing, and action, is most likely to be the biggest determinant of hybrid model deployment success. Latency can be a determining factor for the system's performance as well as its ability to respond appropriately in applications such as real-time analytics and self-driving cars.

### **4.1 Impact of Latency on Edge Devices**

For real-time applications such as self-driving vehicles or augmented reality, even small boosts in latency will result in ineffectual or even dangerous operations. In self-driving cars, for example, slow sensor processing would cause long reaction times and put the car at risk to safety hazards. [9] For other types of real-time analytics applications such as factory automation or health management, latency will impact decision making and customer experience. Such systems usually need to process data in real time in order to monitor and react to shifting conditions, so low-latency processing becomes critical.

#### 4.2 Techniques to minimize Latency

There are several methods available to reduce latency in hybrid systems. One such method is edge caching, which consists of storing frequently accessed data or models locally on the edge device in order to reduce the amount of information that must be fetched from remote cloud servers. Edge devices can reduce wait time for cloud processing by locally storing important data.

Local processing is another approach, where the processing is done partly at the edge and only the most compute-intensive operations are offloaded to the cloud. In this

manner, the most time-critical processes are completed in real-time, and the cloud gets to do more detailed analysis whenever feasible.

Adaptive algorithms can also be used to minimize latency. Adaptive algorithms adjust adaptively in real-time the balance between edge and cloud computation based on factors such as network load, device battery life, or task urgency. Adaptive algorithms enable the ability to optimize performance and minimize latency by taking a smart decision on offloading versus local computation.

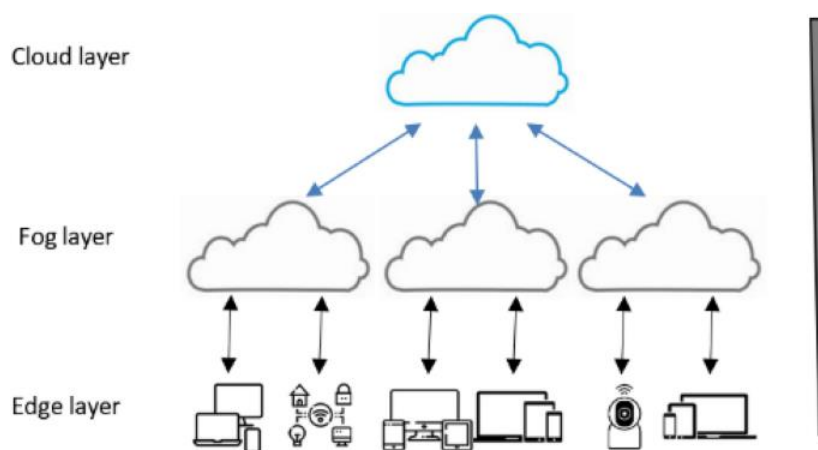


Fig. 2. Hybrid Model Offloading in Cloud

#### 4.3 Trade-offs Between Edge and Cloud Computation

While latency reduction is critical, it comes at the cost of edge computation vs. cloud computation trade-offs. Task execution at the edge can reduce latency, but edge devices lack the processing capacity of the cloud. Offloading tasks to the cloud offers more powerful computation but at the cost of increased latency due to network latency.

The edge vs. cloud computation ratio depends on application requirements. For example, in use cases where latency is the top priority, such as emergency response systems or real-time gaming, edge processing is prioritized to the extent of minimizing delays. However, for computationally demanding operations such as training deep learning models, the use of the cloud's resources becomes inevitable at the cost of some extra latency. The goal is to achieve a balance that optimizes both latency minimization and computational efficiency.

#### 5. Privacy Considerations in Hybrid Model Deployment

In addition to latency, privacy is the foremost concern with hybrid model deployment, especially with sensitive information processed on edge devices and transmitted to the cloud. Computation split between the cloud and edge environments raises certain privacy issues that should be addressed carefully.

##### 5.1 Data Sensitivity on Edge Devices:

Edge devices usually handle extremely sensitive personal data, such as health readings from wearables, video feeds from home security cameras, or location data from mobile phones. Possession or processing of sensitive data on the devices poses privacy concerns, especially if the devices are compromised or data is accidentally leaked. For the

Internet of Things (IoT), billions of devices online process data at the edge, increasing the possibility of privacy violations. [10]

### **5.2 Edge Processing for Privacy:**

To minimize privacy risk, edge devices can perform up-front processing for anonymizing or encrypting sensitive data before it is transmitted to the cloud. Techniques like data anonymization aid in deleting or obscuring personally identifiable information to ensure the privacy of the user while not compromising the analysis of meaningful data. [4]

In some cases, edge data preprocessing can reduce the volume of raw data to be transmitted to the cloud, again enhancing privacy. For example, an edge device might transmit aggregated results or insights derived, but not raw sensor data, and hence restrict sensitive data exposure. Encryption is yet another prevalent technique employed to safeguard data in transit, such that even if intercepted, data cannot be read or altered.

### **5.3 Cloud Security & Encryption:**

As the cloud is outsourced for operations, privacy concerns convert to what is done to the cloud in maintaining sensitive data safety. For guarding privacy in systems based on cloud, providers implement different methods for security such as end-to-end encryption, storing data securely, and access control policies. It is the prime focus of cryptography to ensure data is transferred securely and stored safely in the cloud in a form that will prevent unauthorized individuals or entities from having access to that data.

Cloud platforms also utilize role-based access control and multi-factor authentication to further secure sensitive data and permit only authorized users or systems to access specific resources. However, even with these security features, the movement of sensitive data to the cloud necessarily introduces risks, particularly in multi-tenant or shared cloud environments where various users share the same infrastructure.

### **5.4 Privacy-Enhancing Technologies:**

There are several privacy-improving technologies that are on the rise to counter the data privacy challenges of hybrid systems. One of these is federated learning, a method

through which machine learning models can be trained on several edge devices without requiring raw data to be sent to the cloud. In federated learning, every edge device calculates model updates from local data, and the updates (rather than raw data) are transferred to the cloud for aggregation. This technique keeps sensitive data at the local end, thus maintaining privacy, but enables collaborative model enhancement.

Another method is differential privacy, a technique that renders individual data points un-reversible from aggregate results. Through the addition of noise to the data or results, differential privacy renders any one user's data unidentifiable uniquely, even when calculated by the cloud.

## **6. Performance Evaluation**

Hybrid model deployment system performance analysis is crucial to establish their scalability and efficiency in a range of applications. Various performance characteristics, benchmarking methods, and workload elasticity are employed for overall examination of hybrid systems.

### **6.1 Performance Metrics**

The most important performance metrics to employ when testing hybrid models are power consumption, throughput, accuracy, and latency. Latency is the time taken between receiving data and displaying output, and is critical in applications that require timely outputs like autonomous vehicles and real-time computation. Throughput is the capacity of the system to handle large volumes of data within a specified time and is critical in applications like industrial control or video monitoring.

Accuracy is a numerical measure for the quality of model predictions or classifications, particularly for use in applications like image classification or predictive maintenance. Power consumption is another ultimate and significant measure, especially for edge devices that use limited battery power or power supplies. The greatest challenge to hybrid systems is achieving a balance between performance and energy efficiency. [7]

### **6.2 Benchmarking Hybrid Systems**

Benchmarking is necessary in order to compare hybrid models against fully edge and cloud-based models. Performance benchmarks tend to entail the simulation of workloads and comparing the performance of each system

when processing data. Pure edge computing restricts all data processing to devices, providing low latency but low computational power. Cloud-based models provide enormous computational power but increased latency due to network delay.

Hybrid systems get a balance in between wherein some of the processing is done locally and heavy computation is pushed to the cloud. Benchmarking hybrid systems includes the measurement of performance on different aspects, such as how well a hybrid system performs in terms of real-time processing and low latency as opposed to performing heavy computation through cloud computing.

### **6.3 Scalability and Adaptability**

Hybrid deployments are used to address scalability in heterogeneous applications through dynamic task distribution balancing between the cloud and the edge. Scalability is especially needed in those applications where workload varies from high to very low, for instance, e-commerce websites or IoT networks.

Hybrid system flexibility refers to its ability to manage its computing burden relative to actual situations like network overflow, device loading, or user loading. Adaptability in algorithms in software makes hybrid systems redirect additional work onto the cloud under edge device overflow or alter processing strategy if there is bandwidth constriction. Flexibility allows for hybrid deployment so that performance under different loads may be optimized.

## **7. Challenges And Future Direction**

Though deployment of the hybrid model has numerous advantages, issues are also present and need to be addressed in order to make the most out of its potential. They can be anything from technical constraints to areas of research in the future to take advantage of the hybrid computing impact.

### **7.1 Problem Areas of Hybrid Model Deployment**

Network bandwidth is one of the largest issues in hybrid. Hybrid systems are based on having reliable and fast network connections to perform computation offloading to the cloud. The hybrid system's performance would severely be affected in networks or locations with unstable network infrastructure or those lacking proper networks. Unreliable latency is a problem, especially where there needs to be real-time handling of applications. Changes in

network latency cause delays in processing data, which can be undesirable for the responsiveness and precision of the system.

A further primary challenge is model complexity. The more advanced the model with extra neural networks or intricate structures, the heavier computational power and memory it necessitates. Redistribution of such high-complexity models between cloud and edge poses further complexity in that how individual elements of the model are projected to be executed locally compared to the cloud should factor in both latency and local computation available.

### **7.2 Future Directions**

During the next few years, certain of these advances in hybrid computing would be able to solve the aforementioned problems. The way to proceed is the use of artificial intelligence to make intelligent offloading decisions and determine offloading. It is possible to configure machine learning architecture so it can automatically learn the optimum set of edge and cloud processing as a function of runtime variables such as model size, network condition, and battery level. These intelligent offloading methods have a large potential to contribute to efficiency and performance in hybrid deployments.

The second region with massive research opportunity lies in designing advanced privacy-resilient approaches. As an ever-increasing challenge, data privacy can be resolved through designing differential privacy and federated learning methods such that personal information will never be violated in confidence but even useful computations could be facilitated with hybrid platforms. This tech enables edge devices to carry out in-situ local computations while shielding sensitive raw information from getting processed in the cloud, preventing latency and breach worries.

5G/6G technology advances would also transform the application of hybrid models. Safer and faster network connections would remove most of the latency and bandwidth challenges, allowing hybrid systems to shift more workloads to the cloud without compromising performance. With the possibilities of ultra-low latency and vast bandwidth, 5G/6G networks will be able to host real-time, massive-scale hybrid applications, ranging from smart cities to autonomous vehicles.

## **8. Conclusion**

Hybrid model deployment provides a strong method of bringing together the best of edge and cloud computing with efficient, scalable, and flexible solutions for the majority of applications. Relieving edge devices of computation burdens between edge devices and cloud enables hybrid systems to take advantage of the best of two worlds: low-latency real-time processing at the edge and sheer computation power of the cloud for heavy processing. But this compromise must be specially architected for in the design, most notably when and how much offloading, depending on data size, model complexity, power demands, and privacy requirements. The two main drivers of hybrid system design and deployment are latency and privacy. Local processing, edge caching, and adaptive algorithms reduce latency, and privacy-preserving technologies like federated learning and differential privacy safeguard sensitive information. With the evolving technology, particularly with the advent of 5G/6G networks and advanced AI-based offloading techniques, deployment with hybrid models will be an even better solution for meeting the needs of contemporary computing.

While variable network bandwidth and model complexity remain issues, technological innovation and research can address these issues.

By enhancing the performance, security, and flexibility of hybrid systems, innovation in the future will open new doors for real-time applications across a wide array of industries, ranging from self-driving cars to IoT networks. Finally, the use of hybrid models is a crucial step towards developing smarter, more efficient, and privacy-respecting computing systems in a networked world.

## References

1. Alonso-Monsalve, S., García-Carballeira, F., & Calderón, A. (2018). A heterogeneous mobile cloud computing model for hybrid clouds. *Future Generation Computer Systems*, 87, 651-666.
2. Hosseinzadeh, M., Tho, Q. T., Ali, S., Rahmani, A. M., Souri, A., Norouzi, M., & Huynh, B. (2020). A hybrid service selection and composition model for cloud-edge computing in the internet of things. *IEEE Access*, 8, 85939-85949.
3. Nezami, Z., Zamanifar, K., Djemame, K., & Pournaras, E. (2021). Decentralized edge-to-cloud load balancing: Service placement for the Internet of Things. *Ieee Access*, 9, 64983-65000.
4. Dong, Y., Xu, G., Zhang, M., & Meng, X. (2021). A high-efficient joint 'cloud-edge' aware strategy for task deployment and load balancing. *IEEE Access*, 9, 12791-12802.
5. Zhang, W. Z., Elgendy, I. A., Hammad, M., Ilyasu, A. M., Du, X., Guizani, M., & Abd El-Latif, A. A. (2020). Secure and optimized load balancing for multitier IoT and edge-cloud computing systems. *IEEE Internet of Things Journal*, 8(10), 8119-8132.
6. Pal, S., Jhanjhi, N. Z., Abdulbaqi, A. S., Akila, D., Almazroi, A. A., & Alsubaei, F. S. (2023). A hybrid edge-cloud system for networking service components optimization using the internet of things. *Electronics*, 12(3), 649.
7. Simaiya, S., Lilhore, U. K., Sharma, Y. K., Rao, K. B., Maheswara Rao, V. V. R., Baliyan, A., ... & Alroobaea, R. (2024). A hybrid cloud load balancing and host utilization prediction method using deep learning and optimization techniques. *Scientific Reports*, 14(1), 1337.
8. Bulkan, U., Dagiuklas, T., Iqbal, M., Huq, K. M. S., Al-Dulaimi, A., & Rodriguez, J. (2018). On the load balancing of edge computing resources for on-line video delivery. *IEEE Access*, 6, 73916-73927.
9. Merseedi, K. J., & Zeebaree, S. R. (2024). The cloud architectures for distributed multi-cloud computing: a review of hybrid and federated cloud environment. *The Indonesian Journal of Computer Science*, 13(2).
10. Sharma, A. (2024). Optimizing Hybrid Cloud Architectures: A Comprehensive Study Of Performance Engineering Best Practices. *International Journal Of Engineering And Technology Research (Ijetr)*, 9(2), 288-299