



ISSN: 2617-6548

URL: www.ijirss.com



## Novel load prediction in microservice architecture using attention mechanism-based deep LSTM networks

 Snehal Chaflekar<sup>1\*</sup>,  Rajendra Rewatkar<sup>2</sup>

<sup>1,2</sup>Computer Science & Engineering (AIDS), Datta Meghe Institute of Higher Education and Research, Sawangi (Meghe), Wardha, Maharashtra, India.

Corresponding author: Snehal Chaflekar (Email: [snehalchaflekar@gmail.com](mailto:snehalchaflekar@gmail.com))

### Abstract

Load balancing in microservice architecture is essential for optimizing resource utilization and maintaining high availability. Traditional load balancing algorithms like First-Come-First-Serve (FCFS) and Round Robin often lead to inefficiencies due to their inability to account for server capabilities and varying request sizes. Machine Learning (ML) offers a promising solution by predicting future load patterns and distributing requests more effectively. In this study, we propose an innovative, novel attention mechanism-based Long Short-Term Memory (LSTM) network for web server load prediction. Our methodology involves detailed data preprocessing, sequence creation, and scaling to prepare the NASA HTTP dataset for model training. The attention mechanism enhances the LSTM network's ability to focus on relevant input sequences, significantly improving predictive accuracy. Compared to traditional algorithms such as linear regression, polynomial regression, L2 regularization, decision tree regression, XGBoost, and ARIMA, our model achieves the lowest Mean Squared Error (MSE) of 187,293.59 and Root Mean Squared Error (RMSE) of 432.77, with a strong R-squared value of 0.8532. This superior performance highlights the model's effectiveness in capturing both short-term and long-term dependencies in the data. This novel predictive model can be used to integrate into dynamic and efficient load balancing frameworks. Accurate future load predictions from AMDLN in the microservices environment optimize resource utilization and save infrastructure costs by providing accurate future load predictions for scaling up and scaling down of microservices.

**Keywords:** Attention mechanism, Cloud computing, Load prediction, Microservices.

**DOI:** 10.53894/ijirss.v8i3.6751

**Funding:** This study received no specific financial support.

**History: Received:** 10 March 2025 / **Revised:** 14 April 2025 / **Accepted:** 16 April 2025 / **Published:** 06 May 2025

**Copyright:** © 2025 by the authors. This article is an open access article distributed under the terms and conditions of the Creative Commons Attribution (CC BY) license (<https://creativecommons.org/licenses/by/4.0/>).

**Competing Interests:** The authors declare that they have no competing interests.

**Authors' Contributions:** Both authors contributed equally to the conception and design of the study. Both authors have read and agreed to the published version of the manuscript.

**Transparency:** The authors confirm that the manuscript is an honest, accurate, and transparent account of the study; that no vital features of the study have been omitted; and that any discrepancies from the study as planned have been explained. This study followed all ethical practices during writing.

**Acknowledgment:** Special thanks to Priyadarshini Bhagwati College of Engineering, Nagpur, especially the management and faculty of IT Dept. for the continuous support and encouragement, and allowing me to balance my role as a professor while pursuing my doctoral studies.

**Publisher:** Innovative Research Publishing

## 1. Introduction

In today's rapidly evolving technological landscape, where agility, scalability, and resilience are paramount, microservice architectures have emerged as the dominant paradigm for building modern applications. This architectural style deconstructs applications into a collection of small, independent services that communicate over a network, offering significant advantages over traditional monolithic architectures. These advantages include improved scalability, where individual services can be scaled independently to handle varying workloads, enhanced flexibility, allowing for easier adaptation to changing requirements, and greater resilience, as failures in one service are less likely to impact the entire application.

### 1.1. The Crucial Role of Load Balancing in Microservices

The distributed nature of microservices introduces complexities in managing and balancing the load across these independent services. Load balancing, the process of efficiently distributing workloads across multiple servers or instances, is essential to ensure optimal performance, high availability, and efficient resource utilization in microservice environments. Effective load balancing prevents overload on individual servers, improves response times, and enables seamless horizontal scaling by adding or removing instances as needed.

To illustrate the importance of load balancing in microservices, consider a hypothetical scenario where a sudden surge in user traffic overwhelms a specific service. Without proper load balancing, this service could become a bottleneck, leading to increased latency, degraded performance, and potentially even service outages. However, with effective load balancing, the incoming traffic would be distributed across multiple instances of the service, ensuring that no single instance becomes overloaded and maintaining the overall performance and availability of the application.

### 1.2. Challenges of Traditional Load Balancing Techniques and the Rise of Cloud Platforms

Traditional load balancing techniques, such as First-Come-First-Serve (FCFS), Round Robin, and Least Connections, often struggle to handle the dynamic nature and complexities of microservice workloads. These algorithms typically rely on simplistic metrics, such as the number of active connections or the order of request arrival, without considering crucial factors like server capabilities, request sizes, or the varying processing times of different microservices. This can lead to uneven load distribution, performance bottlenecks, and suboptimal resource utilization.

For instance, in a microservice environment with heterogeneous server capabilities, where some servers are more powerful than others, a simple Round Robin approach could assign computationally intensive tasks to less capable servers, leading to performance degradation. Similarly, if different microservices have varying processing requirements, a Least Connections algorithm might overload servers handling computationally intensive tasks while leaving servers handling simpler tasks underutilized.

The rise of cloud platforms, such as Amazon Web Services (AWS), Microsoft Azure, and Google Cloud Platform (GCP), has further emphasized the need for advanced load balancing solutions. Cloud environments offer on-demand access to vast computing resources, enabling dynamic scaling and cost optimization. However, the dynamic nature of cloud resources and the fluctuating workloads of microservices demand load balancing strategies that can adapt to changing conditions and optimize resource allocation in real-time. Efficient load balancing plays a crucial role in cost optimization within cloud environments. By evenly distributing workloads and preventing over-utilization of resources, load balancing minimizes wasted capacity and reduces overall infrastructure costs.

### 1.3. Need for Advanced Load Balancing Solutions

Leading companies like Netflix, Amazon, and Uber have successfully adopted microservice architectures to achieve exceptional scalability, agility, and resilience in their applications. Netflix, for instance, has transitioned its entire streaming platform to microservices, enabling it to handle millions of concurrent users and deliver a seamless streaming experience. This transition has allowed Netflix to scale its services independently, deploy new features rapidly, and maintain high availability even during peak usage periods.

The limitations of traditional load balancing techniques and the increasing complexity of microservice deployments have highlighted the need for more advanced load balancing solutions. These solutions should be able to:

- Dynamically adapt to changing workloads: Microservice workloads can fluctuate significantly, requiring load balancing solutions that can adjust to these changes in real time. This can be achieved through techniques such as predictive load balancing, where machine learning models are used to forecast future workloads and proactively allocate resources.
- Account for server capabilities: Different servers may have varying processing power, memory, and other resources. Load balancing solutions should consider these differences when distributing workloads. This can be achieved through techniques such as weighted load balancing, where servers are assigned weights based on their capabilities, and requests are distributed proportionally to these weights.

- Handle diverse request types: Microservices may handle a variety of request types, each with different processing requirements. Load balancing solutions should be able to route requests to the appropriate servers based on their type and complexity. This can be achieved through techniques such as content-based routing, where requests are routed based on their content, such as the URL or HTTP headers.
- Integrate with cloud platforms: Cloud platforms offer a variety of services and tools that can be used to enhance load balancing. Load balancing solutions should be able to integrate with these platforms to provide a seamless experience. This can include integration with cloud monitoring services, auto-scaling features, and serverless computing platforms.

#### *1.4. Introducing AMDLN: An AI-Powered Solution for Load Balancing in Microservices*

In this study, we introduce AMDLN, an Attention Mechanism-based Deep LSTM Network designed for accurate and efficient load prediction in microservice architectures. This model leverages the strengths of Long Short-Term Memory (LSTM) networks, renowned for their ability to handle time-series data, and enhances their capabilities with an attention mechanism to focus on the most relevant parts of the input sequence. This approach aims to achieve superior load prediction accuracy, enabling proactive and adaptive load balancing to optimize resource utilization and enhance the performance of microservice applications.

AMD LN addresses the limitations of traditional load balancing techniques by leveraging the power of artificial intelligence (AI) to predict future workloads and dynamically adjust resource allocation. This AI-powered approach enables more efficient and effective load balancing, leading to improved performance, reduced costs, and enhanced resilience in microservice environments.

## **2. Literature Review**

In the domain of load prediction for cloud environments, traditional regression models and neural networks have been explored for their effectiveness and simplicity. Jaradat [1] provides an overview of various load prediction techniques in cloud environments, highlighting the application of simple regression models and basic neural networks. The study discusses linear and polynomial regression models, which are foundational in predicting load based on historical data. These models are appreciated for their simplicity and ease of implementation. However, their limitations are evident in handling non-linear patterns and interactions within the data, which is a common scenario in cloud computing environments.

The use of simple neural networks, particularly feed-forward neural networks, is also discussed in Jaradat [1]. These networks are used to predict load by learning from the historical load data. The key advantage of neural networks over simple regression models is their ability to capture non-linear relationships. However, they require a more extensive computational effort and a larger dataset for training to avoid overfitting and to generalize well on unseen data. The study demonstrates that while simple neural networks perform better than basic regression models, their performance is still limited compared to more sophisticated deep learning techniques.

ARIMA, SARIMA, LSTM, CNN. For more advanced time series forecasting in cloud load prediction, models like ARIMA, SARIMA, LSTM, and CNN have shown significant promise. Calheiros et al. [2] utilized the ARIMA model to predict workload in cloud environments. ARIMA (Auto Regressive Integrated Moving Average) is known for its strength in capturing temporal dependencies in the data. The study found that ARIMA could effectively model and predict short-term load, providing insights into future demands. However, the model's performance diminishes when dealing with non-linearities and sudden spikes in load.

Expanding on ARIMA, Goswami and Kandali [3] introduced SARIMA (Seasonal ARIMA) to handle seasonal patterns in load data. This model incorporates seasonal differencing to improve prediction accuracy for data with periodic fluctuations. The results indicated that SARIMA outperforms ARIMA in scenarios with evident seasonal trends, making it suitable for environments where load patterns repeat over regular intervals.

LSTM (Long Short-Term Memory) networks have gained attention for their ability to handle long-term dependencies and non-linear relationships in time series data. Liu et al. [4] proposed a hybrid model combining LSTM with ARIMA to leverage the strengths of both models. The ARIMA component captures linear temporal dependencies, while the LSTM network addresses non-linear patterns. This hybrid approach demonstrated superior performance in terms of prediction accuracy and robustness compared to using ARIMA or LSTM alone.

Similarly, Chen et al. [5] explored a combination of ARIMA and LSTM for short-term load forecasting. Their results echoed the findings of Liu et al. [4], showing that the hybrid model could effectively predict load with higher accuracy and lower error rates.

Sun and Zhuang [6] introduced a load forecasting algorithm based on the Kalman Filter and Adaptive Neuro-Fuzzy Inference System (ANFIS). This approach combines the strengths of statistical and fuzzy logic methods to enhance prediction accuracy. The Kalman Filter helps in estimating the state of the system dynamically, while ANFIS captures the fuzzy and uncertain nature of the data. The study demonstrated that this combination could provide more reliable predictions compared to traditional methods.

For handling spatial and temporal dependencies, Xu et al. [7] proposed a CNN-LSTM hybrid model. The Convolutional Neural Network (CNN) component extracts spatial features from the data, which are then processed by the LSTM to capture temporal dependencies. This hybrid model showed excellent performance in predicting load in cloud environments, particularly in capturing complex patterns and interactions in the data.

Zhang et al. [8] employed an improved deep learning approach combining CNN and LSTM for load forecasting in cloud computing environments. Their model demonstrated significant improvements in prediction accuracy and computational

efficiency. The study highlighted the importance of combining spatial and temporal modeling techniques to enhance the performance of load prediction systems.

Lastly, Xie et al. [9] developed a hybrid model integrating ARIMA with Triple Exponential Smoothing for real-time prediction of Docker container resource load. This model aimed to balance the strengths of ARIMA in capturing linear dependencies and exponential smoothing in handling level, trend, and seasonality in the data. The results showed that this hybrid model could effectively predict load with higher accuracy and adaptability to real-time changes in the cloud environment.

In the quest for more accurate and efficient load prediction models, researchers have explored various hybrid and other advanced techniques that leverage the strengths of multiple algorithms. This section reviews such methodologies and their applications in cloud environments.

Xu et al. [10] proposed a fusion model for CPU load prediction in cloud computing, combining multiple predictive techniques to enhance accuracy. This model integrates statistical methods with machine learning algorithms, specifically blending ARIMA with Support Vector Regression (SVR) and Bayesian networks. The results indicated a significant improvement in prediction accuracy compared to using individual models, demonstrating the effectiveness of hybrid approaches in handling diverse patterns in load data.

Feng et al. [11] introduced a load prediction optimization framework based on machine learning techniques in cloud computing environments. Their approach integrates Decision Trees and Gradient Boosting Machines to predict future load. This method benefits from the robustness and interpretability of Decision Trees while leveraging the predictive power of ensemble methods like Gradient Boosting. The study showed that this hybrid model outperforms traditional single-model approaches in terms of both accuracy and computational efficiency.

Rasheduzzaman et al. [12] tackled workload prediction using Google Cluster Trace data. They employed a mixture of Hidden Markov Models (HMM), Genetic Algorithms (GA), and Elman networks to predict workload. The HMM captures the temporal dependencies, GA optimizes the model parameters, and Elman networks handle non-linear relationships. This comprehensive approach demonstrated high predictive accuracy, highlighting the importance of combining different algorithmic strengths for load prediction.

Anupama et al. [13] proposed a hybrid model for resource utilization prediction in cloud computing environments. Their model integrates Linear Regression with Neural Networks, aiming to capture both linear trends and non-linear relationships in the data. The study found that this hybrid approach provides better accuracy and generalization capabilities compared to using either Linear Regression or Neural Networks alone.

Song, et al. [14] presented an optimized XGBoost-based sparrow search algorithm for short-term load forecasting. XGBoost, known for its gradient boosting framework, is combined with the Sparrow Search Algorithm to enhance its optimization capabilities. The study demonstrated that this hybrid approach significantly improves load prediction accuracy and reduces computational complexity, making it suitable for real-time applications.

Jindal et al. [15] focused on performance modeling for cloud microservice applications. They employed a hybrid approach combining queueing theory models with machine learning techniques to predict system performance under varying load conditions. This model helps in understanding the behavior of microservices and optimizing their deployment to ensure efficient load handling and resource utilization.

Load balancing is crucial in cloud computing environments to ensure efficient resource utilization and to avoid overloading any single resource. Various studies have explored different approaches to achieve effective load balancing.

Manjunath et al. [16] conducted a comprehensive survey on load prediction techniques in cloud environments. The paper categorizes various approaches into statistical methods, machine learning models, and hybrid techniques, providing a detailed comparison of their strengths and weaknesses. This survey helps in identifying the most suitable prediction models for different types of cloud workloads and paves the way for further research in improving these techniques.

Prevost et al. [17] investigated the prediction of cloud data center network loads using stochastic and neural models. Their approach combines stochastic modeling with neural networks to predict future loads and balance them accordingly. This hybrid method improves the accuracy of load predictions and helps in distributing the load more evenly across the network, thereby enhancing the overall performance of the cloud infrastructure.

Xiaolong et al. [18] proposed a server load prediction algorithm based on CM-MC (Cumulative Moving-Median Control) for cloud systems. This algorithm predicts server load and adjusts resource allocation dynamically to maintain load balance. The study found that the CM-MC algorithm effectively reduces server load variance and improves resource utilization efficiency.

Jena et al. [19] explored the hybridization of meta-heuristic algorithms for load balancing in cloud computing environments. Their approach integrates Particle Swarm Optimization (PSO) with Genetic Algorithms (GA) to optimize load distribution among virtual machines. The hybrid PSO-GA model showed significant improvements in balancing load and minimizing response time compared to traditional load balancing techniques.

Negi et al. [20] developed CMODLB (Clustering-based Multi-Objective Dynamic Load Balancing), an efficient load balancing approach that combines supervised (artificial neural networks), unsupervised (clustering), and soft computing (interval type-2 fuzzy logic system) techniques. This approach clusters virtual machines into underloaded and overloaded groups and balances the load based on these clusters. The study demonstrated that CMODLB enhances load balancing efficiency and reduces energy consumption in cloud environments.

Senthilkumar and Chitra [21] proposed a novel hybrid heuristic-metaheuristic load balancing algorithm for resource allocation in IaaS-cloud computing. This approach combines heuristic methods with meta-heuristic algorithms like Ant

Colony Optimization (ACO) and PSO to dynamically allocate resources and balance load. The results showed that this hybrid model achieves better load distribution and resource utilization compared to standalone heuristic or meta-heuristic methods.

Muchori and Mwangi [22] conducted a review of machine learning load balancing techniques in cloud computing. They examined various machine learning algorithms, such as reinforcement learning, deep learning, and hybrid models for their effectiveness in load balancing. The review highlighted that hybrid approaches combining multiple machine learning techniques often provide superior performance in terms of accuracy and efficiency.

Kumar and Prashar [23] introduced a bio-inspired hybrid algorithm for effective load balancing in cloud computing. Their approach integrates Artificial Bee Colony (ABC) optimization with PSO to achieve load balance. The study found that this bio-inspired hybrid model significantly improves load balancing efficiency and reduces the makespan of tasks in cloud environments.

Geeta and Kamakshi Prasad [24] presented a multi-objective cloud load-balancing approach using hybrid optimization techniques. This method combines Genetic Algorithms with Simulated Annealing to optimize multiple objectives such as load balance, energy consumption, and response time. The results indicated that this multi-objective hybrid approach provides a balanced trade-off between different performance metrics, ensuring efficient load balancing.

Mousavi, et al. [25] proposed a load balancing algorithm for resource allocation in cloud computing. Their approach uses a combination of heuristic and rule-based methods to dynamically allocate resources and balance the load. The study demonstrated that this algorithm effectively balances the load and improves resource utilization in cloud environments.

The attention mechanism has become a pivotal innovation in machine learning, especially in the field of natural language processing, and its applications are expanding to other domains such as load prediction in cloud computing. This section explores the seminal work and its subsequent applications in load prediction.

Vaswani et al. [26] introduced the transformer architecture with the paper "Attention Is All You Need," which revolutionized the use of attention mechanisms in deep learning. This architecture relies solely on self-attention mechanisms, dispensing with recurrent and convolutional layers entirely. The transformer model allows for greater parallelization and has significantly improved performance in various tasks. Although originally designed for natural language processing, the principles of the attention mechanism have been applied to load prediction and other domains, demonstrating its versatility.

Load Prediction Using Attention Following the breakthrough by Vaswani et al. [26] attention mechanisms have been adapted for load prediction in cloud computing environments. These adaptations leverage the ability of attention mechanisms to focus on relevant parts of the input data, thereby improving prediction accuracy and efficiency.

Zhang et al. [27] applied an improved Variational Mode Decomposition (VMD) combined with an attention mechanism for resource load prediction in cloud computing. This approach uses the attention mechanism to enhance the feature extraction process, focusing on the most relevant time series data for accurate load forecasting. The results demonstrated significant improvements in prediction accuracy compared to traditional models, highlighting the potential of attention-based methods in load prediction.

Xie et al. [28] proposed a container load prediction algorithm based on convolutional neural networks (CNN) and self-attention mechanisms. The model integrates CNNs to capture spatial features and the self-attention mechanism to focus on the most critical temporal dependencies. This hybrid approach enhances the model's ability to predict container loads accurately, making it suitable for dynamic cloud environments where resource demand fluctuates frequently.

Remaining Papers This section includes papers that have not been covered in the previous classifications but offer valuable insights into load prediction and balancing techniques in cloud computing.

Di et al. [29] developed a host load prediction model using a Bayesian approach for Google's compute cloud. This model uses Bayesian inference to predict future host loads, providing a probabilistic framework that accounts for uncertainty and variability in load patterns. The Bayesian model demonstrated robust performance in predicting load and assisting in effective resource allocation.

Huo et al. [30] investigated edge cloud load prediction using a dual-recurrent spatio-temporal graph neural network (DR-STGNN). This model combines recurrent neural networks (RNN) with graph neural networks (GNN) to capture both temporal and spatial dependencies in edge cloud environments. The DR-STGNN showed superior performance in load prediction, effectively managing the complexity of data generated in edge computing scenarios.

Meng et al. [31] explored load balancing techniques in cloud computing environments, focusing on a comparative analysis of different load balancing algorithms. The study provided a comprehensive overview of static and dynamic load balancing methods, highlighting their advantages and limitations. The findings emphasized the importance of adaptive load balancing strategies to handle varying workloads and optimize resource utilization.

As reviewed by Chaflekar et al. [32] and Chaflekar et al. [33], there exist job scheduling and load balancing challenges in cloud computing. There is a need to propose a priority-based approach for efficient and equitable resource utilization.

### **3. Methodology**

This section provides a comprehensive overview of the methodology employed in this study, encompassing the dataset used, data preprocessing techniques, and the proposed AMDLN model.

#### **3.1. Dataset**

The study utilizes the NASA-HTTP dataset, which comprises two months of HTTP requests to the NASA Kennedy Space Center WWW server, covering July and August 1995. The dataset includes detailed information such as host, timestamp, request, HTTP reply code, and bytes, providing valuable insights into traffic patterns and server load. This dataset is chosen for its comprehensive nature, capturing a wide range of web traffic characteristics, and its availability for public use, enabling reproducibility and further research.

### 3.2. Data Preprocessing

The initial phase of the methodology involves filtering and analyzing the raw data to extract relevant information and prepare it for model training. This process includes:

- **Data Cleaning:** Raw log files are parsed to extract key fields such as timestamp, request type, and response size. Regular expressions are used for accurate extraction, and timestamps are converted from string format to Python datetime objects to facilitate time-based operations.
- **Data Aggregation:** The extracted data is aggregated to calculate the total number of requests and bytes transferred per hour, transforming the raw log data into a time series format suitable for analysis and model training. Let:

$$R_h = \sum_{i=1}^N \delta(r_i, h)$$

where  $R_h$  represents the total requests in hour  $h$ ,  $r_i$  is an individual request, and  $\delta(r_i, h)$  is an indicator function:

$$\delta(r_i, h) = \{1, \text{if request } r_i \text{ belongs to hour } h \ 0, \text{otherwise}$$

Similarly, total bytes transferred per hour is computed as:

$$B_h = \sum_{i=1}^N \delta(r_i, h) \cdot s_i$$

where  $s_i$  is the response size (in bytes) of request  $r_i$ .

- **Statistical Analysis:** Key statistics such as minimum, maximum, mean, and total bytes are calculated to gain insights into the data distribution:

$$\mu_B = \frac{1}{H} \sum_{h=1}^H B_h$$

$$\sigma_B^2 = \frac{1}{H} \sum_{h=1}^H (B_h - \mu_B)^2$$

where  $\mu_B$  is the mean bytes transferred per hour,  $\sigma_B^2$  is the variance, and  $H$  is the total number of hours analyzed.

- **Visualization:** Temporal analysis is conducted using time series plots, where  $B_h$  and  $R_h$  are plotted over time to identify traffic patterns. The smoothed trend of requests over time is represented as:

$$R_h^{(smooth)} = \frac{1}{k} \sum_{j=-k}^k w_j R_{h+j}$$

where  $w_j$  are weights of a moving average filter and  $k$  defines the smoothing window size.

This comprehensive data filtering and analysis process ensures that the data is clean, relevant, and suitable for training the proposed AMDLN model.

### 3.3. Data Preparation for Model Training

The filtered and analyzed data is further processed to prepare it for training the AMDLN model. This process ensures that the model receives well-structured inputs, enabling it to effectively capture temporal dependencies and predict future load patterns. The primary steps in data preparation include feature scaling and sequence creation, both of which are crucial for improving model stability and performance.

#### 3.3.1. Feature Scaling

To ensure that all input features contribute equally to the training process, we normalize the data using Min-Max Scaling, which rescales each feature to a fixed range between 0 and 1. This prevents features with larger magnitudes from dominating the learning process and improves model convergence.

Let  $x_i$  be the original value of a feature at time step  $i$ , with  $x_{min}$  and  $x_{max}$  representing the minimum and maximum observed values of that feature across the dataset. The scaled value  $x_{scaled,i}$  is computed as:

$$x_{scaled,i} = \frac{x_i - x_{min}}{x_{max} - x_{min}}$$

This transformation ensures that all feature values are mapped to the interval  $[0,1]$ , preserving relative relationships between values.

#### 3.3.2. Justification for Feature Scaling

- It ensures numerical stability in optimization by preventing large gradient updates.
- It accelerates convergence in gradient-based learning algorithms.
- It helps in handling different feature scales, particularly when combining heterogeneous data sources (e.g., request counts vs. byte sizes).

If the data follows a non-uniform distribution, an alternative normalization method, such as Z-score normalization (standardization), may be applied:

$$x_{standardized,i} = \frac{x_i - \mu}{\sigma}$$

where  $\mu$  is the mean and  $\sigma$  is the standard deviation of the feature. However, Min-Max Scaling is preferred in this study as it retains the original data distribution within a bounded range.

### 3.3.3. Time Series Sequence Creation

To capture temporal dependencies in web request patterns, the time series data is transformed into sequences of fixed length  $T$ . Each input sequence consists of  $T$  consecutive hourly observations, and the corresponding target is the web server load at time step  $T + 1$ . The input sequences are constructed as follows:

$$X_t = \{x_{t-T}, x_{t-T+1}, \dots, x_t\}, \quad y_t = x_{t+1}$$

where:

- $X_t$  represents the input sequence consisting of the last  $T$  observations.
- $y_t$  is the target value, i.e., the predicted load at the next time step  $t + 1$ .

For a dataset containing  $N$  observations, the total number of sequences that can be extracted is:

$$N_{seq} = N - T$$

This ensures that the model receives a continuous stream of past observations, allowing it to learn underlying trends and seasonality in web traffic.

#### Why Sequence Creation?

- Encapsulates short-term dependencies (recent trends) and long-term dependencies (seasonal patterns).
- Converts independent time points into structured data suitable for LSTM-based models.
- Allows the model to learn dynamic shifts in request volume over time.

### 3.3.4. Sliding Window Approach

To enhance learning, a sliding window technique is applied to generate overlapping training sequences. Given a stride of  $S$ , the sequence extraction follows:

$$X_{t+k} = \{x_{t+k-T}, \dots, x_{t+k}\}, \quad y_{t+k} = x_{t+k+1}, \quad k = 0, S, 2S, \dots$$

where  $S$  is the step size that determines how frequently the window is moved forward.

- If  $S = 1$ , we extract overlapping sequences (maximizing data usage).
- If  $S > 1$ , we generate fewer sequences but introduce more variability.

### 3.3.5. Feature Representation

To enhance the model's ability to differentiate request types, each data sample includes:

- Raw traffic metrics: Request count, bytes transferred.
- Temporal features: Hour of day, day of week, weekend/weekday binary encoding.
- Lagged values: Past values at different time intervals.

The final input vector is defined as:

$$X_t = [x_t, x_{t-1}, \dots, x_{t-T}, \text{Hour}_t, \text{DayOfWeek}_t, \text{Weekend}_t]$$

where  $X_t$  is a multi-dimensional vector capturing both request trends and contextual time features?

### 3.3.6. Final Prepared Data

After preprocessing, the dataset consists of:

- $N_{seq}$  training sequences of shape  $(T, d)$ , where  $d$  is the number of features.
- Corresponding target values for each sequence.
- Normalized feature scales to improve training stability.

This data preparation process ensures that the AMDLN model is trained on structured, time-dependent input sequences, allowing it to learn complex traffic patterns and predict future server loads accurately.

## 3.4. Proposed Method: Attention Mechanism-Based Deep LSTM Networks (AMDLN)

The core of this study is the proposed Attention Mechanism-based Deep LSTM Networks (AMDLN), which integrates the capabilities of Bidirectional Long Short-Term Memory (BiLSTM) networks and an Attention Mechanism to dynamically predict web server load. This work draws inspiration from previous research on a Model-Agnostic ML-based Smart Load Balancer for Microservices (2024) [34], where ML models such as LSTM and GRU were employed for cost-efficient resource management. Building on earlier experience with secure cloud storage systems (2015) [35], it was hypothesized that newer ML models, when combined with systematic data preprocessing, could achieve superior cost-effectiveness and faster response times for load balancers in microservices environments. This hybrid approach enables the model to focus on the most relevant parts of the input sequence while capturing long-term dependencies from both past and future contexts.

### 3.4.1. Bidirectional LSTM Network

Unlike standard LSTMs, which process sequences only in a forward direction, Bidirectional LSTMs (BiLSTMs) process input sequences in both forward and backward directions, significantly enhancing the model's ability to understand both past and future temporal dependencies.

### 3.4.2. Mathematical Formulation of BiLSTM

Each BiLSTM cell consists of two LSTM networks: one processing the sequence from past to future and another from future to past. The forward and backward hidden states are concatenated to form the final representation:

$$\begin{aligned} \vec{h}_t &= LSTM_f(x_t, \vec{h}_{t-1}) \\ \overleftarrow{h}_t &= LSTM_b(x_t, \overleftarrow{h}_{t+1}) \\ h_t &= [\vec{h}_t; \overleftarrow{h}_t] \end{aligned}$$

where: -  $\vec{h}_t$  is the forward hidden state, -  $\overleftarrow{h}_t$  is the backward hidden state, -  $h_t$  is the concatenated final hidden state.

By stacking multiple BiLSTM layers, the model captures both short-term and long-term temporal dependencies in web server load variations.

### 3.4.3. Attention Mechanism

While BiLSTMs enhance feature extraction, they may still assign equal importance to all time steps. To address this, we integrate an Attention Mechanism, allowing the model to focus on the most relevant time steps.

### 3.4.4. Mathematical Formulation of Attention

The attention mechanism computes a weighted sum of hidden states, giving higher importance to more relevant time steps. The process is formulated as follows:

- Alignment Score: Computes a similarity measure between the hidden state  $h_t$  and the query state  $s_{t-1}$ :

$$score(h_t, s_{t-1}) = \tanh(W_a[h_t; s_{t-1}])$$

- Attention Weights: Softmax normalizes the scores to obtain weights  $\alpha_t$ :

$$\alpha_t = \frac{\exp(score(h_t, s_{t-1}))}{\sum_{k=1}^T \exp(score(h_k, s_{t-1}))}$$

- Context Vector: The final weighted sum of hidden states:

$$c_t = \sum_{i=1}^T \alpha_{t,i} h_i$$

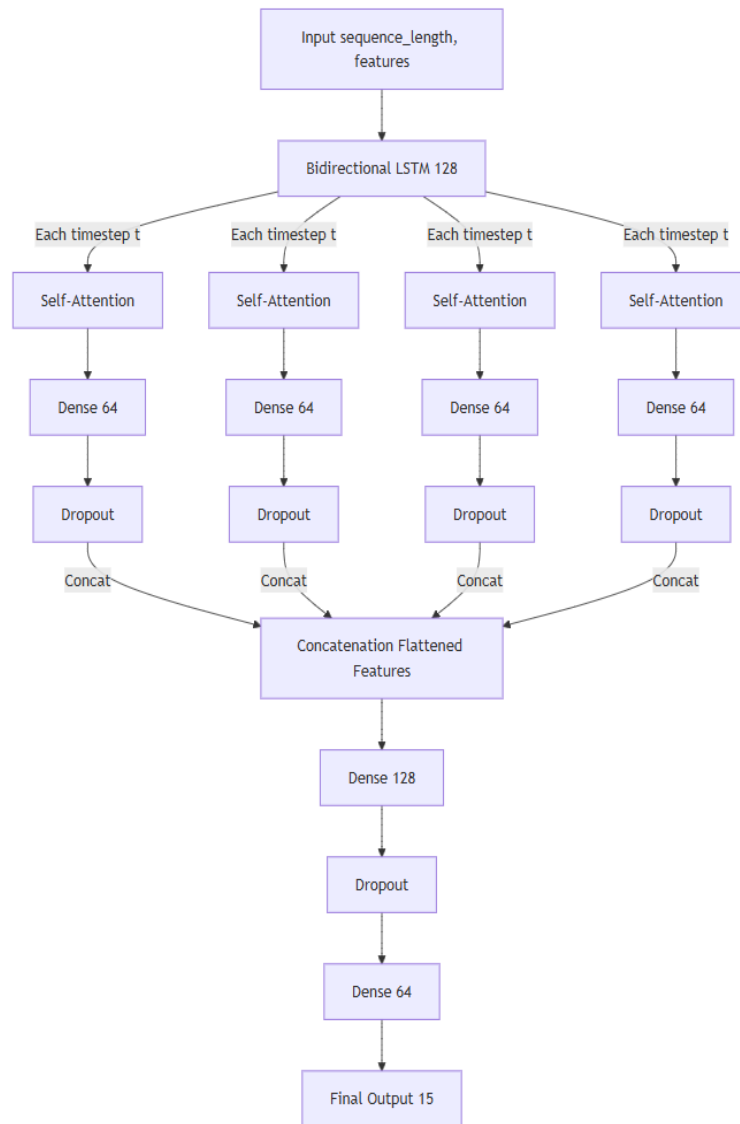
- Attention Output: The attention-enhanced hidden state is computed as:

$$\tilde{h}_t = \tanh(W_c[c_t; s_t])$$

### 3.4.5. AMDLN Model Architecture

The proposed AMDLN model integrates Bidirectional LSTMs with an Attention Mechanism to generate accurate demand predictions. Unlike conventional sequence models that aggregate all time steps into a single vector, AMDLN creates parallel branches, where each time step has its own attention-based dense layer, allowing the model to extract localized temporal dependencies before combining them into a final output.





**Figure 1.**  
AMDNLN Architecture Diagram.

- Input Layer: The model takes a sequence of  $T$  previous time steps as input:

$$X = \{x_{t-T}, x_{t-T+1}, \dots, x_t\} \text{ where } X \in R^{T \times d}$$

Here,  $d$  represents the number of input features at each time step.
- Bidirectional LSTM Encoder: The input sequence is processed through a Bidirectional LSTM network to extract temporal dependencies from both past and future contexts:

$$H = BiLSTM(X) \text{ where } H \in R^{T \times 2h}$$

The hidden state matrix  $H$  contains both forward ( $\vec{h}_t$ ) and backward ( $\overleftarrow{h}_t$ ) hidden representations concatenated at each time step:

$$h_t = [\vec{h}_t; \overleftarrow{h}_t] \text{ where } h_t \in R^{2h}$$

This enables the model to capture dependencies from both past and future data points.
- Branching per Sequence Step: After the BiLSTM, each time step  $t$  has its own separate processing branch, ensuring that each moment in the sequence is treated individually before being merged. The sequence branches are structured as follows:

  - Step Extraction: The  $t^{th}$  hidden state is extracted:
 
$$h_t' = H[t] \text{ where } h_t' \in R^{2h}$$
  - Self-Attention Mechanism: Self-attention is applied at each time step to prioritize relevant features:
 
$$A_t = Attention(h_t', h_t')$$
  - Dense Transformation: The attention-refined features are passed through a dense layer:
 
$$d_t = \sigma(W_d A_t + b_d) \text{ where } d_t \in R^k$$
  - Dropout Regularization: To prevent overfitting, dropout is applied:
 
$$\tilde{d}_t = Dropout(d_t)$$

Each of the  $T$  time steps follows this sequence of operations independently.

- Concatenation Layer: All time-step-specific representations are concatenated to form a final feature vector:

$$F = \text{Concatenate}([\tilde{d}_1, \tilde{d}_2, \dots, \tilde{d}_T]) \text{ where } F \in R^{T \times k}$$

- Final Dense Layers: The concatenated features are passed through fully connected layers to produce the final prediction:

$$y' = \sigma(W_f F + b_f) \text{ where } y' \in R^{128}$$

$$\hat{y} = W_o y' + b_o \text{ where } \hat{y} \in R^M$$

where  $M$  represents the number of output predictions (e.g., demand for different request types).

- Output Layer: The final output  $\hat{y}$  predicts the demand for each request type at time  $t + 1$ .

This architecture allows the model to capture detailed, time-step-specific interactions before aggregating them into a robust global prediction.

### 3.4.6. Loss Function and Optimization

The AMDLN model is trained using the Mean Squared Error (MSE) Loss:

$$L = \frac{1}{N} \sum_{i=1}^N (y_i - \hat{y}_i)^2$$

The Adam optimizer is used for adaptive learning rate adjustments.

### 3.4.7. Final Model Pipeline

The complete AMDLN model follows these steps:

- Preprocess web server logs and extract demand patterns.
- Normalize and segment data into fixed-length sequences.
- Train the Bidirectional LSTM with Attention Mechanism.
- Optimize using Adam with early stopping.
- Evaluate performance using MSE, RMSE, and  $R^2$ .

## 3.5. Analysis

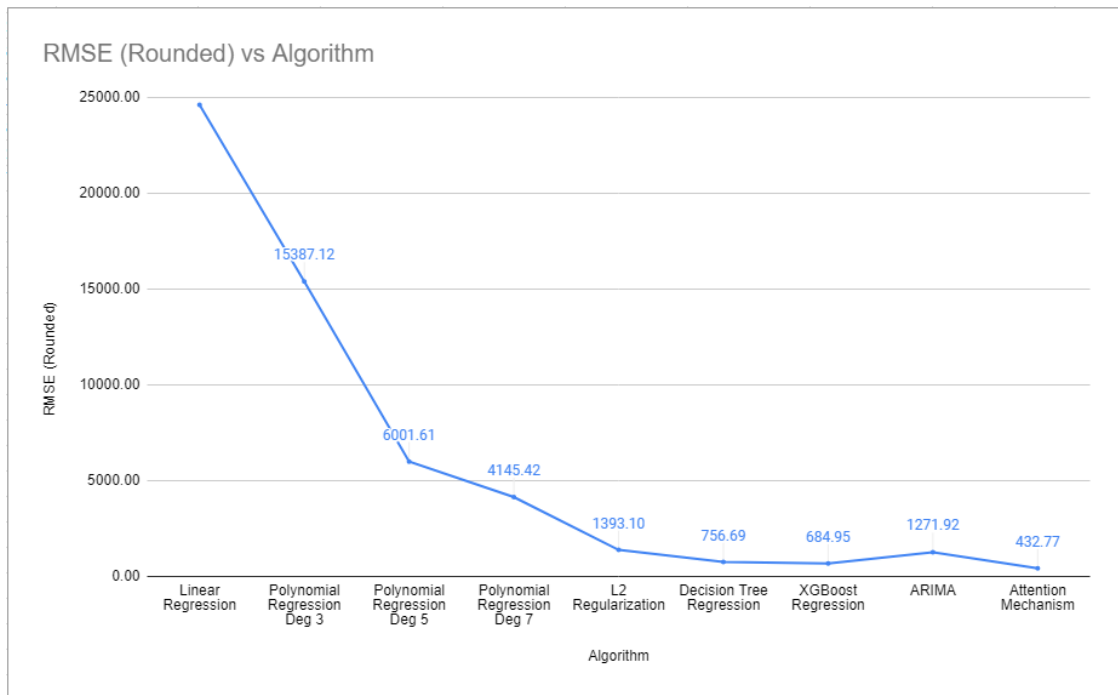
In this section, we analyze the performance of various algorithms used for predicting web server load. The results are evaluated based on Mean Squared Error (MSE), Root Mean Squared Error (RMSE), and R-squared (R2) metrics. The following table summarizes the performance of each algorithm:

### 3.5.1. Performance Comparison of Algorithms

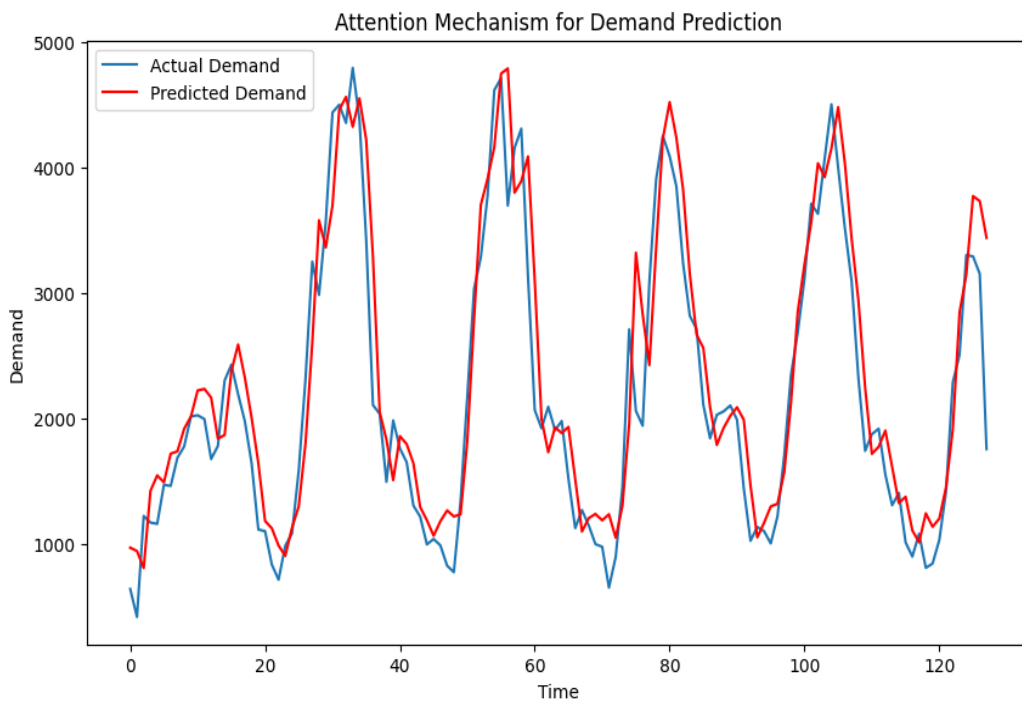
**Table 1.**

Performance Comparison of Algorithms in-terms of MSE, RMSE and R2 parameter.

Algorithm	MSE	RMSE	R2
Linear Regression	605,945,914.6	24,615.97	0.2026
Polynomial Regression Deg 3	236,763,394.9	15,387.12	0.6884
Polynomial Regression Deg 5	36,019,342	6,001.61	0.9526
Polynomial Regression Deg 7	17,184,536.62	4,145.42	0.9774
L2 Regularization	1,940,741.452	1,393.10	0.1100
Decision Tree Regression	572,574.9986	756.69	0.9992
XGBoost Regression	469,155.7802	684.95	0.9994
ARIMA	1,617,770.763	1,271.92	0.9799
Attention Mechanism	187,293.5947	432.77	0.8532



**Figure 2.** RMSE of all Algorithms (Lowest AMDLN).



**Figure 3.** AMDLN model Testing (Expected vs Predicted).

**Linear Regression:** Linear regression resulted in a high MSE of 605,945,914.6 and a low R2 of 0.2026, indicating that the model struggled to capture the complexity of the data.

**Polynomial Regression:** Increasing the polynomial degree significantly improved the model’s performance. Polynomial regression with a degree of 7 achieved an MSE of 17,184,536.62 and an R2 of 0.9774, demonstrating a strong fit to the data.

**L2 Regularization:** Ridge regression (L2 regularization) showed poor performance with an MSE of 1,940,741.452 and an R2 of 0.1100, indicating that regularization was not effective for this dataset.

**Decision Tree Regression:** Decision tree regression performed exceptionally well with an MSE of 572,574.9986 and an R2 of 0.9992, suggesting that it effectively captured the non-linear relationships in the data.

**XGBoost Regression:** XGBoost further improved the performance, achieving an MSE of 469,155.7802 and an R2 of 0.9994, highlighting its robustness and ability to handle complex patterns.

**ARIMA:** The ARIMA model, designed for time series data, achieved an MSE of 1,617,770.763 and an R2 of 0.9799, showing strong performance in capturing temporal dependencies.

Attention Mechanism: Our proposed attention mechanism-based LSTM network outperformed all other models with a significantly lower MSE of 187,293.5947 and an RMSE of 432.77. Although its R<sup>2</sup> value of 0.8532 is slightly lower compared to XGBoost and decision tree models, the attention mechanism provides a nuanced understanding of the data by focusing on relevant input sequences, making it highly effective for predictive accuracy.

#### 3.5.2. Comparative Analysis

While traditional models like linear regression and polynomial regression of lower degrees struggled to capture the complexity and non-linearity in the data, higher-degree polynomial regressions showed marked improvements. Decision tree regression and XGBoost demonstrated exceptional performance, with XGBoost slightly outperforming decision trees due to its advanced boosting techniques.

However, our attention mechanism-based LSTM network stands out due to its ability to leverage the sequential nature of time series data and focus on the most relevant parts of the input sequence. This capability not only improves predictive accuracy but also provides insights into which parts of the sequence are most critical for the prediction, a feature not present in other models.

In summary, while several models showed strong performance, our attention mechanism-based approach offers a comprehensive and highly effective solution for web server load prediction, combining predictive power with interpretability.

## 4. Conclusion and Future Scope

In this study, we explored various algorithms for predicting web server load, including linear regression, polynomial regression, L2 regularization, decision tree regression, XGBoost, and ARIMA. Our proposed method, an attention mechanism-based LSTM network, demonstrated superior performance with a significantly lower MSE and RMSE, highlighting its effectiveness in capturing both short-term and long-term dependencies in time series data. The ability of the attention mechanism to focus on relevant parts of the input sequence further enhanced the model's predictive accuracy.

### 4.1. Future Scope

Accurate load prediction is a critical step toward effective load balancing in cloud environments. Future work will focus on integrating these predictive models into load balancing frameworks, using machine learning to dynamically allocate resources based on predicted demand. This approach aims to optimize resource utilization, reduce latency, and improve overall system performance, ultimately leading to more efficient and responsive cloud services.

## References

- [1] E. M. A. Jaradat, "Load prediction techniques in cloud environment," *International Journal of Research In Science & Engineering*, pp. 1-10, 2022. <https://doi.org/10.55529/ijrise.21.1.10>
- [2] R. N. Calheiros, E. Masoumi, R. Ranjan, and R. Buyya, "Workload prediction using ARIMA model and its impact on cloud applications' QoS," *IEEE transactions on cloud computing*, vol. 3, no. 4, pp. 449-458, 2014. <https://doi.org/10.1109/TCC.2014.2350475>
- [3] K. Goswami and A. B. Kandali, "Electricity demand prediction using data driven forecasting scheme: ARIMA and SARIMA for real-time load data of assam," presented at the 2020 International Conference on Computational Performance Evaluation (ComPE), 2020.
- [4] X. Liu, X. Xie, and Q. Guo, "Research on cloud computing load forecasting based on LSTM-ARIMA combined model," presented at the 2022 Tenth International Conference on Advanced Cloud and Big Data, 2022.
- [5] S. Chen, R. Lin, and W. Zeng, "Short-term load forecasting method based on ARIMA and LSTM," presented at the 2022 IEEE 22nd International Conference on Communication Technology (ICCT), 2022.
- [6] J. Sun and Y. Zhuang, "The cloud computing load forecasting algorithm based on kalman filter and ANFIS," in *Proceedings of the 2016 4th International Conference on Machinery, Materials and Computing Technology*, Atlantis Press, 2016, pp. 565-569.
- [7] F. Xu, G. Weng, Q. Ye, and Q. Xia, "Research on load forecasting based on CNN-LSTM hybrid deep learning model," presented at the 2022 IEEE 5th International Conference on Electronics Technology, 2022.
- [8] K. Zhang, W. Guo, J. Feng, and M. Liu, "Load forecasting method based on improved deep learning in cloud computing environment," *Scientific Programming*, vol. 2021, no. 1, p. 3250732, 2021. <https://doi.org/10.1155/2021/3250732>
- [9] Y. Xie *et al.*, "Real-time prediction of docker container resource load based on a hybrid model of ARIMA and triple exponential smoothing," *IEEE Transactions on Cloud Computing*, vol. 10, no. 2, pp. 1386-1401, 2020. <https://doi.org/10.1109/tcc.2020.2989631>
- [10] D. Xu, S. Yang, and H. Luo, "A fusion model for CPU load prediction in cloud computing," *Journal of Networks*, vol. 8, no. 11, p. 2506, 2013. <https://doi.org/10.4304/jnw.8.11.2506-2511>
- [11] G. Feng, J. Xu, W. Jian, and Z. Liu, "Load prediction optimization based on machine learning in cloud computing environment," in *Proceedings of the 4th International Conference on Advanced Information Science and System*, 2022.
- [12] M. Rasheduzzaman, M. A. Islam, and R. M. Rahman, "Workload prediction on google cluster trace," *International Journal of Grid and High Performance Computing*, vol. 6, no. 3, pp. 34-52, 2014. <https://doi.org/10.4018/ijghpc.2014070103>
- [13] K. Anupama, B. Shivakumar, and R. Nagaraja, "Resource utilization prediction in cloud computing using hybrid model," *International Journal of Advanced Computer Science and Applications*, vol. 12, no. 4, pp. 1-10, 2021. <https://doi.org/10.14569/IJACSA.2021.0120447>
- [14] J. Song, L. Jin, Y. Xie, and C. Wei, "Optimized XGBoost based sparrow search algorithm for short-term load forecasting," presented at the 2021 IEEE International Conference on Computer Science, Artificial Intelligence and Electronic Engineering, 2021.
- [15] A. Jindal, V. Podolskiy, and M. Gerndt, "Performance modeling for cloud microservice applications," in *Proceedings of the 2019 ACM/SPEC International Conference on Performance Engineering*, 2019.

- [16] R. Manjunath, C. Manaswini, and N. Bangar, "A survey on load prediction techniques in cloud environment," *International Journal of Engineering Research and Technology*, vol. 2, pp. 1-10, 2018.
- [17] J. J. Prevost, K. Nagothu, B. Kelley, and M. Jamshidi, "Prediction of cloud data center networks loads using stochastic and neural models," presented at the 2011 6th International Conference on System of Systems Engineering, 2011.
- [18] X. Xiaolong, Z. Qitong, M. Yiqi, and L. Xinyuan, "Server load prediction algorithm based on CM-MC for cloud systems," *Journal of Systems Engineering and Electronics*, vol. 29, no. 5, pp. 1069-1078, 2018. <https://doi.org/10.21629/JSEE.2018.05.17>
- [19] U. K. Jena, P. K. Das, and M. R. Kabat, "Hybridization of meta-heuristic algorithm for load balancing in cloud computing environment," *Journal of King Saud University-Computer and Information Sciences*, vol. 34, no. 6, pp. 2332-2342, 2022. <https://doi.org/10.1016/j.jksuci.2020.01.012>
- [20] S. Negi, M. M. S. Rauthan, K. S. Vaisla, and N. Panwar, "CMODLB: An efficient load balancing approach in cloud computing environment," *The Journal of Supercomputing*, vol. 77, no. 8, pp. 8787-8839, 2021. <https://doi.org/10.1007/s11227-020-03601-7>
- [21] G. Senthilkumar and M. Chitra, "A novel hybrid heuristic-metaheuristic load balancing algorithm for resource allocation in IaaS-cloud computing," presented at the 2020 Third International Conference on Smart Systems and Inventive Technology, 2020.
- [22] J. G. Muchori and P. M. Mwangi, "Machine learning load balancing techniques in cloud computing: a review," *International Journal of Computer Applications Technology and Research*, pp. 1-10, 2022. <https://doi.org/10.7753/ijcatr1106.1002>
- [23] R. Kumar and T. Prashar, "A bio-inspired hybrid algorithm for effective load balancing in cloud computing," *International Journal of Cloud Computing*, vol. 5, no. 3, pp. 218-246, 2016. <https://doi.org/10.1080/1206212X.2023.2260616>
- [24] K. Geeta and V. Kamakshi Prasad, "Multi-objective cloud load-balancing with hybrid optimization," *International Journal of Computers and Applications*, vol. 45, no. 10, pp. 611-625, 2023. <https://doi.org/10.1080/1206212X.2023.2260616>
- [25] S. Mousavi, A. Mosavi, and A. Várkonyi-Kóczy, "A load balancing algorithm for resource allocation in cloud computing," in *Proceedings of the 6th International Conference on Information Technology and Computer Science*, 2017, pp. 289-296.
- [26] A. Vaswani *et al.*, "Attention is all you need," *Advances in Neural Information Processing Systems*, vol. 30, pp. 1-10, 2017.
- [27] M. Zhang, Z. Fan, Y. Miao, and L. Yang, "Cloud computing resource load prediction based on improved VMD and attention mechanism," presented at the Journal of Physics: Conference Series, 2023.
- [28] X. Xie, C. Wang, D. Zhang, L. Zheng, and P. Gao, "Container load prediction algorithm based on convolution neural network and self-attention," *2023 International Conference on Pattern Recognition and Artificial Intelligence (PRAI)*, 2023.
- [29] S. Di, D. Kondo, and W. Cirne, "Host load prediction in a google compute cloud with a bayesian model," presented at the 2012 International Conference for High Performance Computing, Networking, Storage and Analysis, 2012.
- [30] X. Huo, X. Yu, and J. Wu, "Research and application of edge cloud load prediction based on DR-STGNN?," in *Proceedings of the 7th International Conference on Computer Science and Application Engineering*, 2023.
- [31] Y. Meng, J. Chen, and X. Liu, "A service load interval prediction method for cloud-edge collaborations based on type identification and SVMs," in *Proceedings of the 2023 IEEE 9th International Conference on Smart Cloud (SmartCloud)*, 2023, pp. 1258705-1258705-8.
- [32] S. Chaflekar, R. Rewatkar, K. T. V. Reddy, U. Pacharaney, M. Adak, and S. Khedkar, "Job scheduling approach in load balancing in cloud computing environment," in *AIP Conference Proceedings*, AIP Publishing, 2024, vol. 3188, no. 1.
- [33] S. Chaflekar, R. Rewatkar, and K. T. V. Reddy, "Load balancing in cloud environment": Issues and challenges," in *AIP Conference Proceedings*, AIP Publishing., 2024, vol. 3188, p. 1.
- [34] S. Chaflekar and R. Rewatkar, "Model agnostic predictive smart load balancer in microservices environment," presented at the 2024 2nd DMIHER International Conference on Artificial Intelligence in Healthcare, Education and Industry (IDICAIEI), Wardha, India, 2024.
- [35] S. Chaflekar and R. Raipure, "Application of cloud system on secure cloud storage," *International Research Journal of Engineering and Technology*, vol. 2, no. 6, pp. 105-109, 2015.