

Article

Not peer-reviewed version

Time Series Prediction of Backend Server Load via Deep Learning and Attention Mechanisms

[Le Liu](#)*

Posted Date: 12 November 2025

doi: 10.20944/preprints202511.0872.v1

Keywords: load forecasting; deep learning; time series modeling; sensitivity analysis



Preprints.org is a free multidisciplinary platform providing preprint service that is dedicated to making early versions of research outputs permanently available and citable. Preprints posted at Preprints.org appear in Web of Science, Crossref, Google Scholar, Scilit, Europe PMC.

Copyright: This open access article is published under a Creative Commons CC BY 4.0 license, which permit the free download, distribution, and reuse, provided that the author and preprint are cited in any reuse.

Disclaimer/Publisher's Note: The statements, opinions, and data contained in all publications are solely those of the individual author(s) and contributor(s) and not of MDPI and/or the editor(s). MDPI and/or the editor(s) disclaim responsibility for any injury to people or property resulting from any ideas, methods, instructions, or products referred to in the content.

Article

Time Series Prediction of Backend Server Load via Deep Learning and Attention Mechanisms

Le Liu

University of California, San Diego, San Diego, USA; liuleny@outlook.com

Abstract

This study proposes a deep learning-based time series modeling method for backend server load prediction, aiming to address the limitations of traditional approaches in handling complex nonlinear features and multidimensional dependencies. The method begins with normalization and outlier processing of server operation data to ensure stability and consistency, followed by feature transformation through linear mapping and embedding layers, and constructs a dynamic modeling framework with attention mechanisms to capture both local short-term fluctuations and global long-term trends. A decoding structure generates the prediction output, trained with mean squared error optimization. Experiments are conducted on large-scale public cluster logs, including CPU, memory, disk I/O, and network traffic, with comparative analysis across multiple models using metrics such as MSE, RMSE, MAE, and R^2 . Results show that the proposed method achieves superior accuracy and robustness compared with baseline methods, effectively reflecting real server behavior. Further sensitivity experiments on hyperparameters, environment conditions, and data perturbations validate the adaptability and robustness of the model under varying constraints. Overall, the approach enriches theoretical perspectives on load prediction and demonstrates strong practical potential for enhancing system performance and resource management efficiency.

Keywords: load forecasting; deep learning; time series modeling; sensitivity analysis

1. Introduction

In the context of rapid digitalization and extensive networking, backend servers act as the central hub of information systems[1]. They undertake critical functions such as computation, storage, and service response. The load status of these servers directly determines the efficiency and stability of the entire system. With the continuous expansion of business scale and the growth of access requests, server loads show highly dynamic and complex patterns. Peaks and sudden bursts occur frequently[2]. If load changes cannot be perceived and predicted in time, serious problems such as system delay, response failure, or even downtime may arise. Therefore, developing intelligent and efficient methods for load prediction has become a key research direction for ensuring stable backend system operation[3].

Time series forecasting provides an essential tool for load modeling and scheduling optimization. Backend server loads are usually represented as multidimensional and multi-granularity dynamic sequences. They include CPU utilization, memory usage, disk I/O, and network traffic, while also being affected by user behavior, changes in business types, and external conditions. These sequences often combine nonlinear, nonstationary, and periodic features. Traditional statistical approaches struggle to capture such hidden patterns, leading to poor prediction accuracy and limited generalization. By contrast, deep learning can extract and model latent structures from load sequences, offering stronger capability to represent complex dynamics and thus improving forecasting results[4].

Deep learning shows unique advantages in handling large-scale and complex data. In time series modeling, architectures based on recurrent neural networks, convolutional networks, and attention mechanisms can effectively capture both short- and long-term dependencies [5-7], as well as global

dynamic features. For backend server load data, which are high-dimensional, noisy, and influenced by multiple factors, deep learning can automatically extract hidden features and represent them through multilayer nonlinear mappings[8]. This capability overcomes the limitations of traditional models related to feature dependence and dimensionality. As a result, deep learning provides a new perspective for load prediction, making forecasts more consistent with the complexity of real operating environments.

The importance of backend load prediction goes beyond improving the performance of a single system. It also contributes to optimizing resource utilization and reducing energy consumption[9]. Accurate forecasts support elastic resource scheduling, preventing both idle resources and overload, which improves throughput and service quality. At the same time, effective prediction mechanisms reduce unnecessary energy use, lower the operating costs of data centers, and promote green computing and sustainable development. In cloud and edge computing environments, such prediction also supports resource isolation and fair scheduling under multi-tenant architectures, thereby safeguarding user experience and business continuity.

More importantly, load prediction serves not only technical but also strategic purposes. Forecast results provide evidence for enterprise decision-making [10]. They guide capacity planning and disaster recovery strategies and enhance the ability of systems to handle high concurrency and sudden events. This research also drives the development of intelligent operations. By combining prediction with automated scheduling[11], it is possible to build a closed-loop optimization system that minimizes manual intervention and moves toward autonomous and intelligent operation [12]. In summary, research on deep learning-based time series prediction of backend server loads holds significant theoretical value and broad application prospects. It is of great importance for building efficient, stable, and sustainable digital infrastructure[13].

2. Proposed Approach

In this study, the time series prediction task of backend server load is modeled as a mapping problem from a multidimensional sequence to a real-valued sequence. The proposed approach employs a deep learning framework for backend server load prediction, as illustrated in Figure 1. At the data preprocessing stage, we apply normalization and outlier processing to the raw server operation data, ensuring stable and consistent inputs for the downstream modeling pipeline. For feature transformation, the model leverages linear mapping and embedding layers to encode multi-dimensional operational metrics such as CPU, memory, disk I/O, and network traffic. This enables effective representation of both local and global system states and supports robust learning across heterogeneous data sources. In the dynamic modeling module, we integrate spatiotemporal graph neural network principles, drawing on the multi-level performance forecasting methods advanced by Xue et al. [14]. By modeling servers and their interactions as nodes and edges in a spatiotemporal graph, the framework captures both the complex dependencies between system components and their temporal evolution, significantly enhancing predictive capability. The core prediction network further employs collaborative evolution strategies inspired by Li et al.[15], allowing intelligent agents within the system to adaptively evolve and coordinate forecasting decisions as operational patterns shift. This agent-based mechanism improves both adaptability and scalability in large-scale microservice environments. Additionally, a reinforcement learning-driven scheduling algorithm, as described by Zhang et al.[16], is applied to dynamically adjust the model's attention and computational resource allocation in response to workload fluctuations and multi-tenant demands. This optimization strategy ensures efficient resource utilization and maintains prediction accuracy under high concurrency. Let the input time series be $X = \{x_1, x_2, \dots, x_T\}$, where each $x_t \in R^d$ represents a d-dimensional system metric collected at time t. The goal is to learn a nonlinear function $f(\cdot)$ that, given the historical sequence, can predict the load value at future moments, namely:

$$\hat{y}_{t+h} = f(x_1, x_2, \dots, x_t)$$

where \hat{y}_{t+h} represents the predicted server load at the h th step in the future. The core of this problem lies in how to effectively model long-term dependencies and local dynamics to ensure that the prediction results can reflect the complex changes in the server operating environment. The overall model architecture is shown in Figure 1.

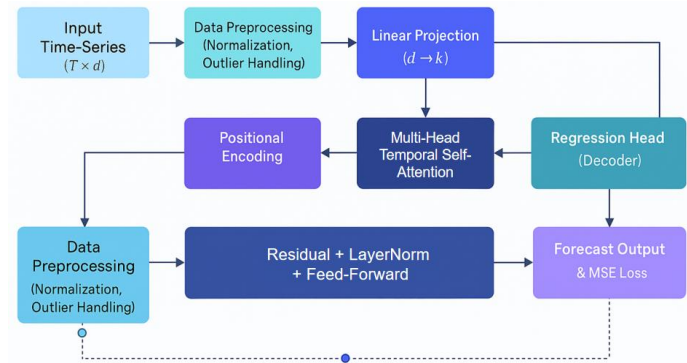


Figure 1. Overall model architecture.

To capture the correlation between input features, the original sequence is first linearly transformed and normalized. For each input vector x_t , it is mapped to a low-dimensional embedding space through the weight matrix $W \in R^{d \times k}$:

$$z_t = Wx_t + b$$

where $z_t \in R^k$ is the representation in the embedding space, and b is the bias term. This step unifies the input data into a representation space suitable for deep network modeling. The embedded sequence $\{z_1, z_2, \dots, z_T\}$ is then fed into the temporal modeling module to extract global and local dependencies.

For temporal modeling, this approach applies an attention-based mechanism to effectively capture complex feature interactions across different time steps. To ensure adaptive and context-sensitive modeling, the attention distribution is guided by multi-agent reinforcement learning strategies, which allow the model to dynamically adjust focus based on changing resource states and workload patterns [17]. This enables the system to maintain robust temporal feature extraction even under fluctuating operational conditions. Agent-based scheduling and coordination techniques are also employed to optimize the aggregation and utilization of temporal information in prediction, improving both the responsiveness and efficiency of the overall modeling process [18]. This coordinated approach helps the model prioritize critical periods and better capture the temporal dynamics relevant to system performance. Additionally, by incorporating graph neural network and transformer integration, the model is able to capture subtle dependencies and identify anomalous patterns within the time series data [19]. This integration enhances both the expressiveness and reliability of the attention mechanism, ensuring more accurate and interpretable temporal representations.

The attention weight between time positions is calculated as:

$$\alpha_{ij} = \frac{\exp\left(\frac{(z_i W_Q)(z_j W_K)^T}{\sqrt{k}}\right)}{\sum_{j=1}^T \exp\left(\frac{(z_i W_Q)(z_j W_K)^T}{\sqrt{k}}\right)}$$

$$\alpha_{ij} = \frac{\exp\left(\frac{(z_i W_Q)(z_j W_K)^T}{\sqrt{k}}\right)}{\sum_{j=1}^T \exp\left(\frac{(z_i W_Q)(z_j W_K)^T}{\sqrt{k}}\right)}$$

where $W_Q, W_K \in R^{k \times k}$ is the projection matrix of the query and key, respectively. Through weighted summation, the new representation is obtained:

$$h_i = \sum_{j=1}^T \alpha_{ij} (z_j W_V)$$

where $W_V \in R^{k \times k}$ is the projection matrix of the value. This mechanism can dynamically assign the importance of different time steps, thereby effectively capturing long-range dependencies.

In the output layer, to achieve prediction, a linear decoding structure is used to map the time series representation to the final target value [20]. The specific form is:

$$\hat{y}_{t+h} = \sigma(h_t W_o + b_o)$$

where $W_o \in R^{k \times 1}$, b_o is the decoding parameter, and $\sigma(\cdot)$ is the nonlinear activation function. The training goal of the entire model is to minimize the error between the predicted value and the true value, which is usually constrained by the mean square error function:

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

where N is the total number of samples, y_i is the actual value, and \hat{y}_i is the predicted value. This optimization objective guides the model to continuously adjust parameters, thereby learning a more accurate load forecast mapping function in complex dynamic environments.

3. Performance Evaluation

3.1. Dataset

The dataset used in this study comes from the publicly available Google Cluster Trace Dataset. It is derived from real logs of large-scale distributed computing clusters. The dataset records detailed load conditions of thousands of servers over continuous days. It contains key indicators such as CPU utilization, memory usage, disk I/O, and network traffic. These data comprehensively reflect the dynamic operation of backend servers in multi-task and multi-tenant environments. Due to its large scale and high sampling frequency, the dataset provides abundant support for load modeling and time series prediction.

The dataset has notable characteristics in both dimensionality and time span. It includes sudden peaks within short time windows and also preserves stable patterns across long cycles. This multi-granularity nature offers multiple perspectives for studying load prediction. It enables models to capture both short-term fluctuations and long-term trends. In addition, the dataset is anonymized. It contains no user or task identifiers, ensuring privacy and security during the research process.

The dataset was chosen because it realistically reflects server load patterns in production environments. Its scale is sufficient to meet the training requirements of deep learning models. Using this dataset makes it possible to evaluate prediction methods in complex and dynamic settings. This provides solid data support for intelligent scheduling and resource management in backend systems.

3.2. Experimental Results

This paper first conducts a comparative experiment, and the experimental results are shown in Table 1.

Table 1. Comparative experimental results.

Method	MSE	RMSE	MAE	R ²
Transformer[21]	0.0521	0.2283	0.1635	0.918
LSTM[22]	0.0734	0.2709	0.1846	0.892
BiLSTM[23]	0.0657	0.2563	0.1762	0.901
iTransformer[24]	0.0456	0.2137	0.1524	0.927
Ours	0.0389	0.1972	0.1411	0.936

From the results in Table 1, it can be observed that traditional recurrent neural network methods show limited performance in backend server load prediction. LSTM presents relatively high error values, indicating that it still struggles with long-term dependency modeling. Although BiLSTM improves bidirectional information capture to some extent, leading to lower MSE, RMSE, and MAE compared with LSTM, it still cannot fully adapt to the complex dynamics of server load sequences. The R² value also fails to reach an ideal level.

In contrast, the Transformer structure demonstrates stronger modeling ability in time series forecasting. Its multi-head self-attention mechanism captures global dependencies more effectively. As a result, it achieves lower MSE, RMSE, and MAE than LSTM and BiLSTM. This shows its higher sensitivity to long-term trends and global patterns in load sequences. However, Transformer still has certain limitations in modeling local patterns and short-term fluctuations, leaving room for improvement in prediction accuracy.

iTransformer further enhances predictive performance in this experiment. It outperforms the standard Transformer on all error metrics, and its R² reaches 0.927. This indicates that with structural improvements, the model can balance both global dependencies and local variations more effectively. Such an advantage is especially significant under highly volatile and bursty load conditions, highlighting the modeling potential of iTransformer.

Finally, the proposed method in this study achieves the best performance across all metrics. Compared with other methods, Ours reports the lowest MSE, RMSE, and MAE, while R² reaches 0.936. This demonstrates that the model can capture complex nonlinear load patterns with greater precision. These results confirm the effectiveness and superiority of the proposed method for backend server load prediction. They also emphasize its potential in intelligent scheduling and resource optimization, providing strong support for building efficient and stable computing environments.

This paper also presents an experiment on the sensitivity of the upper limit of available memory to MSE, and the experimental results are shown in Figure 2.

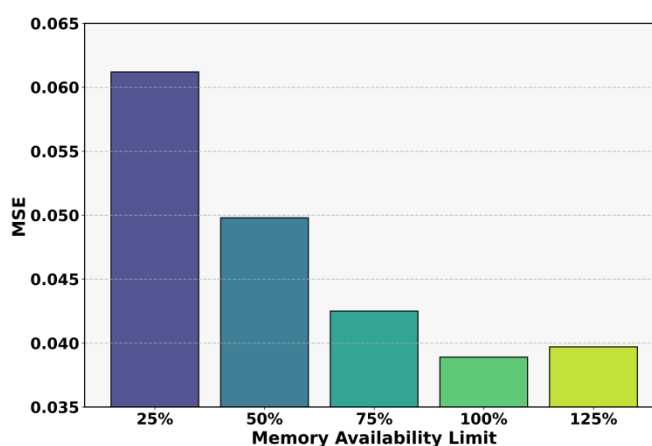


Figure 2. Sensitivity experiment of the memory available upper limit to MSE.

From Figure 2, it can be observed that when the upper limit of available memory is low, the prediction model shows significantly higher MSE. Under the condition of a 25% memory limitation, the model error exceeds 0.06. This indicates that insufficient memory constrains feature storage and intermediate representation computation, which reduces the ability to capture complex load patterns. It shows that hardware resource constraints have a direct impact on prediction performance, especially in high-intensity computation and large-scale sequence modeling scenarios.

As the memory limit increases, model performance improves, and MSE shows a decreasing trend. Under 50% and 75% memory conditions, prediction accuracy improves significantly. This indicates that additional computational resources support more stable parameter updates and more effective feature utilization. This phenomenon highlights the advantage of deep learning models under sufficient resources, where their ability to model nonlinear patterns and long-term dependencies in load sequences is better realized.

When the memory limit reaches 100%, the model achieves the lowest MSE of about 0.0389, showing the best performance. This suggests that without resource constraints, the model can achieve optimal time series modeling and effectively capture the complex dynamics of server loads. The result also confirms a clear positive correlation between resource allocation and prediction performance, providing practical evidence for backend server load management.

It is worth noting that when the memory limit further increases to 125%, the MSE does not continue to decrease but instead slightly rises. This implies that excessive resource allocation does not continuously improve prediction performance. The extra memory may not be fully utilized and may even introduce noise and additional computational overhead. This trend suggests that in practical deployment, computational resources should be allocated reasonably according to model characteristics and task requirements to achieve an optimal balance between performance and cost.

This paper further presents an experiment on the sensitivity of network jitter and packet loss rate to R^2 . The experimental results are shown in Figure 3.

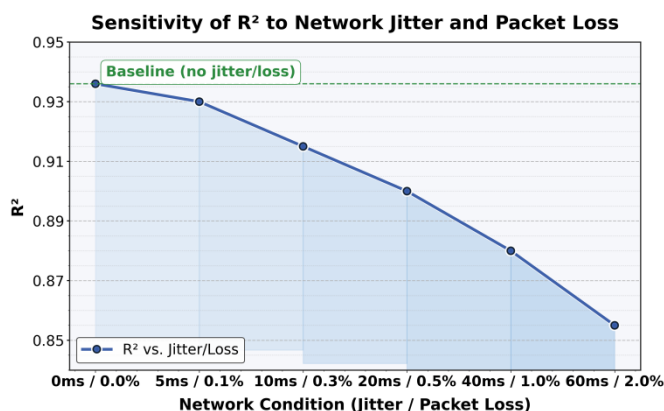


Figure 3. Experiment on the Sensitivity of Network Jitter and Packet Loss Rate to R^2 .

From Figure 3, it can be observed that when the network environment remains stable, the model maintains high predictive performance, with R^2 close to 0.936. This indicates that the proposed method can effectively capture the temporal characteristics of server loads under ideal conditions. However, as network jitter and packet loss increase, R^2 gradually decreases. This shows that unstable network conditions weaken the model's ability to capture global dependencies. The phenomenon highlights the significant impact of data transmission quality on prediction performance.

When network jitter increases to 10 ms and packet loss reaches 0.3%, R^2 drops to about 0.915, showing a clear decline compared with the baseline. This indicates that when facing incomplete data or transmission delays, the model cannot fully maintain the integrity of temporal information, leading to interference in characterizing load dynamics. For backend servers, such performance

degradation means that prediction stability requires special attention in environments with frequent network fluctuations.

Under harsher network conditions, such as 40 ms jitter and 1% packet loss, R^2 falls to around 0.88. This further confirms the continuous negative impact of unstable network environments on prediction accuracy. Although deep learning models possess strong nonlinear modeling ability, they remain highly sensitive to input data quality. This result suggests that researchers should consider the robustness of input sequences when designing load prediction frameworks.

When network jitter reaches 60 ms and packet loss rises to 2%, R^2 drops to about 0.855, showing a clear loss compared with the baseline. This indicates that structural optimization of the model alone is not sufficient to resist the interference caused by network noise. It is necessary to combine enhancements at both the data level and system level, such as data reconstruction, fault-tolerant mechanisms, and multi-source information fusion, to maintain high predictive performance under complex network environments. This result also emphasizes the significance of this study, which explores robustness improvements of load prediction models under unstable conditions and provides stronger adaptability for intelligent backend system management.

4. CONCLUSIONS

This study focuses on the problem of backend server load time series prediction and builds a deep learning-based modeling framework. From input data preprocessing and embedding representation to attention-based temporal feature extraction and final prediction output, the framework forms a complete technical system. By characterizing load patterns in complex dynamic environments, the results show that the proposed method can more accurately capture nonlinear relationships and both short- and long-term dependencies, thereby improving prediction accuracy and stability. The research not only provides new ideas in methodological design but also lays the foundation for applying time series forecasting in high-load computing scenarios.

The significance of this work lies in its ability to address nonstationarity and multidimensional correlation issues that traditional statistical models struggle with. Through deep network structures, the predictions become more consistent with the complexity of real operating environments. The accuracy of load forecasting directly affects resource scheduling, fault tolerance, and performance optimization. Therefore, the proposed method plays an important role in ensuring the stability and service quality of large-scale computing platforms. Especially under sudden traffic and resource constraints, early awareness of load changes helps avoid bottlenecks and service interruptions.

Furthermore, this study has strong practical value for intelligent operation and maintenance in cloud computing, edge computing, and data centers. Accurate load prediction provides reliable support for elastic resource allocation, enabling dynamic adjustment of computing, storage, and network resources. This mechanism improves resource utilization and system throughput while reducing energy consumption, promoting the development of green computing. At the same time, the introduction of prediction models supports the formation of an intelligent decision loop, allowing systems to self-regulate based on prediction results. This reduces manual intervention and enhances automation.

In the future, related research can be extended to more complex and heterogeneous data environments. Examples include cross-task load prediction in multi-tenant architectures or joint modeling that incorporates more contextual information in large-scale distributed systems. In addition, integrating prediction models with anomaly detection, fault warning, and security protection functions will contribute to more robust and reliable intelligent operation systems. These explorations are not only significant for the sustainable development of backend systems but also drive the evolution and upgrading of intelligent computing infrastructure.

References

1. H. Chen and Z. Dang, "A hybrid deep learning-based load forecasting model for logical range," *Applied Sciences*, vol. 15, no. 10, 5628, 2025.
2. K. Aidi and D. Gao, "Temporal-spatial deep learning for memory usage forecasting in cloud servers," 2025.
3. Z. Lai, D. Zhang, H. Li, et al., "LightCTS: A lightweight framework for correlated time series forecasting," *Proceedings of the ACM on Management of Data*, vol. 1, no. 2, pp. 1-26, 2023.
4. J. Tang, S. Chen, C. Gong, et al., "LLM-PS: Empowering large language models for time series forecasting with temporal patterns and semantics," *arXiv preprint arXiv:2503.09656*, 2025.
5. H. Wang, "Temporal-semantic graph attention networks for cloud anomaly recognition," 2024.
6. W. Xu, M. Jiang, S. Long, Y. Lin, K. Ma and Z. Xu, "Graph neural network and temporal sequence integration for AI-powered financial compliance detection," 2025.
7. H. Yang, M. Wang, L. Dai, Y. Wu and J. Du, "Federated graph neural networks for heterogeneous graphs with data privacy and structural consistency," 2025.
8. S. Simaiya, U. K. Lilhore, Y. K. Sharma, et al., "A hybrid cloud load balancing and host utilization prediction method using deep learning and optimization techniques," *Scientific Reports*, vol. 14, no. 1, 1337, 2024.
9. R. Akter, M. G. Shirkoohi, J. Wang, et al., "An efficient hybrid deep neural network model for multi-horizon forecasting of power loads in academic buildings," *Energy and Buildings*, vol. 329, 115217, 2025.
10. K. Y. Chan, B. Abu-Salih, R. Qaddoura, et al., "Deep neural networks in the cloud: Review, applications, challenges and research directions," *Neurocomputing*, vol. 545, 126327, 2023.
11. J. Hu, B. Zhang, T. Xu, H. Yang and M. Gao, "Structure-aware temporal modeling for chronic disease progression prediction," *arXiv preprint arXiv:2508.14942*, 2025.
12. X. Wang, X. Zhang and X. Wang, "Deep skin lesion segmentation with Transformer-CNN fusion: Toward intelligent skin cancer analysis," *arXiv preprint arXiv:2508.14509*, 2025.
13. D. Gao, "High fidelity text to image generation with contrastive alignment and structural guidance," *arXiv preprint arXiv:2508.10280*, 2025.
14. Z. Xue, Y. Zi, N. Qi, M. Gong and Y. Zou, "Multi-level service performance forecasting via spatiotemporal graph neural networks," *arXiv preprint arXiv:2508.07122*, 2025.
15. Y. Li, S. Han, S. Wang, M. Wang and R. Meng, "Collaborative evolution of intelligent agents in large-scale microservice systems," *arXiv preprint arXiv:2508.20508*, 2025.
16. X. Zhang, X. Wang and X. Wang, "A reinforcement learning-driven task scheduling algorithm for multi-tenant distributed systems," *arXiv preprint arXiv:2508.08525*, 2025.
17. G. Yao, H. Liu and L. Dai, "Multi-agent reinforcement learning for adaptive resource orchestration in cloud-native clusters," *arXiv preprint arXiv:2508.10253*, 2025.
18. R. Zhang, "AI-driven multi-agent scheduling and service quality optimization in microservice systems," *Transactions on Computational and Scientific Methods*, vol. 5, no. 8, 2025.
19. Y. Zi, M. Gong, Z. Xue, Y. Zou, N. Qi and Y. Deng, "Graph neural network and transformer integration for unsupervised system anomaly discovery," *arXiv preprint arXiv:2508.09401*, 2025.
20. Ren, Y. (2024). Strategic Cache Allocation via Game-Aware Multi-Agent Reinforcement Learning. *Transactions on Computational and Scientific Methods*, 4(8).
21. Y. Nie, N. H. Nguyen, P. Sinthong, et al., "A time series is worth 64 words: Long-term forecasting with transformers," *arXiv preprint arXiv:2211.14730*, 2022.
22. K. Wang, J. Zhang, X. Li, et al., "Long-term power load forecasting using LSTM-informer with ensemble learning," *Electronics*, vol. 12, no. 10, 2175, 2023.
23. F. Liu and C. Liang, "Short-term power load forecasting based on AC-BiLSTM model," *Energy Reports*, vol. 11, pp. 1570-1579, 2024.
24. Y. Liu, T. Hu, H. Zhang, et al., "iTransformer: Inverted transformers are effective for time series forecasting," *arXiv preprint arXiv:2310.06625*, 2023.

Disclaimer/Publisher's Note: The statements, opinions and data contained in all publications are solely those of the individual author(s) and contributor(s) and not of MDPI and/or the editor(s). MDPI and/or the editor(s)

disclaim responsibility for any injury to people or property resulting from any ideas, methods, instructions or products referred to in the content.