

Received May 3rd, 2026; accepted May 29th, 2026. Date of publication June 30th, 2026
Digital Object Identifier: <https://doi.org/10.25047/jtit.v13i1.490>

Multivariate LSTM with SLO-Aware Loss for Virtual Machine Workload Prediction on Cloud Data Center

AGUS HARIYANTO¹, AHMAD FAHRIYANNUR ROSYADY¹, ADI SUCIPTO¹, BEKTI MARYUNI SUSANTO¹, SAPTA NUGRAHA², NICOLAS CHENU³

¹Jurusan Teknologi Informasi, Politeknik Negeri Jember, Jl. Mastrip Kotak Pos 164 Jember Jawa Timur, Indonesia

²Universitas Maritim Raja Ali Haji, Jl. Raya Dompok, Dompok, Kec. Bukit Bestari, Kota Tanjung Pinang, Kepulauan Riau, Indonesia

³Department of Industrial and Territorial Ecology, Polytech Ancecy-Chambery, Université Savoie Mont Blanc, France

CORRESPONDING AUTHOR: AGUS HARIYANTO (email: agus_hariyanto@polije.ac.id)

ABSTRACT One of the main requirements of cloud resource management is the ability to make an accurate prediction of virtual machine (VM) workload, which is also a prerequisite of auto-scaling and commitment to Service Level Objectives (SLOs). Conventional models for making predictions using symmetric loss functions, e.g., Mean Squared Error (MSE), have a major shortcoming that under-prediction errors are treated the same as over-prediction errors, while under-prediction is a lot more dangerous in terms of the operation of the system as it results in capacity shortage and SLO violation. In this research, a multivariate Long Short-Term Memory (LSTM) network is used to build the CPU workload prediction model and it is mixed with an SLO-aware loss, an asymmetric loss function that penalizes under-prediction ten times more severely than over-prediction. Four features input and a subset of 25, 000 rows of the Bitbrain GWA-T-12 fastStorage dataset were used as a platform for experiments and results are reproducible under the same conditions of using the fixed random seed. Two types of models were trained and compared: one with SLO-aware loss and one with standard MSE as baseline. Both models share the same architecture and hyperparameters. The main performance measure was the under-prediction rate, which is the most direct measure of SLO violation risk. The results reveal that the SLO-aware model has a rate of under-prediction of 0.04% as against 0.16% for MSE baseline, which is a four times reduction. These results provide a strong indication that SLO-aware loss effectively moves the model in the direction of conservative predictions that safeguarding SLO compliance. Thus, the design of loss function is a critical and impactful aspect of cloud VM workload prediction..

KEYWORDS: VM workload prediction, LSTM, SLO-aware loss, cloud computing

1. INTRODUCTION

The expansion of cloud computing services has necessitated the development of resource allocation algorithms that adjust to variations in workload [1], [2]. In virtualized systems, virtual machines (VMs) display dynamic and multidimensional load patterns encompassing CPU use, memory usage, and network traffic [3]. Accurate prediction of VM workload is essential for auto-scaling, load balancing, and VM scheduling techniques that seek cost effectiveness while adhering to Service Level Objectives (SLOs) [1], [4], [5].

In production cloud systems, SLO violations have serious operational and financial repercussions.

When demand is at its highest, under-provisioning causes latency spikes, request timeouts, and cascade failures that can spread throughout linked service chains [5]. These mistakes result in customer attrition, financial fines outlined in SLA contracts, and reputational harm that is hard to recover from in cutthroat cloud markets [1], [4]. On the other hand, the 30–40% resource utilization gap that is frequently seen in enterprise data centers is a result of systematic over-provisioning to prevent such failures, which represents wasted capital expenditure [2]. Rather than just minimizing symmetric error, accurate workload prediction that is directionally biased toward conservative over-estimation provides a principled middle ground: it

absorbs demand spikes without the blanket waste of static over-provisioning and without the SLO risk of reactive scaling that arrives too late [5], [6].

Deep learning techniques, especially Long Short-Term Memory (LSTM) networks [7], have been extensively used in workload forecasting because to their capacity to identify long-term temporal connections [8], [9], [6]. To increase prediction accuracy, a number of hybrid architectures have been put forth. Leka et al. [10] used a one-dimensional Convolutional Neural Network (CNN) with Long Short-Term Memory (LSTM) to model temporal dependencies and extract spatial characteristics from the Bitbrain dataset. Similarly, Shuvo et al. [11] proposed LSRU, a hybrid method that stacks a 1D convolutional layer on top of LSTM and GRU to capture both spatial and sequential features of VM workloads. Dang-Quang and Yoo [12] demonstrated that multivariate Bi-LSTM reduces prediction error by 46% compared to its univariate counterpart on the same Bitbrain fastStorage trace. Bai et al. [13] further advanced the field by integrating Temporal Convolutional Networks (TCN), Gated Recurrent Units (GRU), and a self-attention mechanism, showing superior generalizability across different cloud host clusters. Bhardwaj et al. [14] proposed a hybrid LSTM-RNN model that achieved a 58% RMSE reduction compared to conventional baselines on the GWA-T-12 Bitbrain dataset. Chen et al. [15] introduced EN-Beats, an ensemble learning approach that aggregates multiple N-Beats weak learners for multiple resource metric predictions, while Bansal and Kumar [16] explored an ensemble of CNNs, RNNs, and other deep learning models. Beyond accuracy, Mahbub et al. [17] raised a critical concern about the robustness of DL-based workload forecasting models, demonstrating their significant vulnerability to white-box adversarial attacks on the Bitbrain, Google, and Alibaba trace datasets.

A critical but frequently overlooked limitation of these studies is the exclusive reliance on symmetric accuracy metrics — principally Root Mean Squared Error (RMSE) and Mean Absolute Error (MAE) — as the sole criteria for model evaluation. While RMSE and MAE quantify average prediction error magnitude, they are agnostic to the direction of error and therefore cannot distinguish between under-prediction and over-prediction of equal magnitude. In cloud resource management, however, this directional distinction is operationally decisive. Under-prediction rate — defined as the proportion of time steps in which the model's predicted workload falls below the actual workload ($\hat{y} < y$) — directly measures the frequency of conditions that may lead to resource capacity shortfalls and, consequently, SLO violations [1],[4][5]. A model with a low RMSE but a high under-prediction rate may still fail to protect SLOs in practice, because even a small but

systematic tendency to underestimate demand during peak periods is sufficient to trigger cascading service degradation [5]. Conversely, a model with a slightly higher RMSE but a very low under-prediction rate is operationally safer, as its predictions consistently provide at least the resources demanded by the actual workload. For this reason, this study adopts the under-prediction rate — rather than RMSE or MAE — as the primary evaluation metric, treating it as the direct operational proxy for SLO violation risk in the cloud VM workload prediction context.

However, the majority of prior studies use symmetric loss functions such as Mean Squared Error (MSE) or Mean Absolute Error (MAE) [18], [6]. Such functions treat under-prediction and over-prediction equivalently [19], whereas in cloud operations, the two error types have fundamentally different implications [4], [20]: over-prediction tends to cause excess resource allocation and cost inefficiency, but the service remains operational; under-prediction can cause capacity shortages during load spikes, leading to performance degradation and SLO violations [4], [5]. Research on asymmetric loss functions has been actively developed in the recent literature [19], with applications to time-series forecasting, financial forecasting, and network/cloud resource allocation [20], [21]. Nevertheless, the explicit integration of an SLO-aware loss into a multivariate LSTM architecture for real VM workload traces such as the Bitbrain dataset remains limited. Specifically, SLO-aware loss is an asymmetric variant of the squared error loss that assigns a heavier penalty factor $\alpha > 1$ to under-prediction errors ($\hat{y} < y$) while retaining the standard squared error for over-prediction cases ($\hat{y} \geq y$). This asymmetry directly encodes the operational cost difference in cloud environments: a model that over-predicts by a margin of ϵ incurs only a resource waste proportional to ϵ^2 , whereas a model that under-predicts by the same margin ϵ incurs a penalty of $\alpha \cdot \epsilon^2$, driving the optimizer to steer predictions toward a conservative regime. As a consequence, the trained model is structurally biased to over-estimate rather than under-estimate workload, reducing the risk of insufficient resource provisioning and thereby protecting SLO compliance, at the cost of a controlled degree of over-allocation. This paper provides empirical evidence of this behavior by directly comparing an SLO-aware-trained LSTM against an MSE-trained LSTM of identical architecture on the same dataset, enabling an isolated evaluation of the loss function's effect on under-prediction frequency.

Based on these considerations, the main research questions addressed in this study are: (i) how to design an LSTM-based VM workload prediction model that explicitly minimizes under-prediction so that it aligns with the goal of maintaining SLOs; (ii) how to leverage multivariate

information (CPU, memory, network) to reduce the under-prediction rate compared to univariate approaches by capturing inter-feature correlations that signal impending demand spikes [22], [23]; and (iii) how the proposed approach performs in terms of under-prediction rate when evaluated on a real workload trace dataset such as Bitbrain [3], [13], [10], [11]. This study contributes by applying an SLO-aware loss to a multivariate LSTM model and conducting an initial evaluation on a subset of the Bitbrain dataset, supported by an Optuna-based hyperparameter optimization framework [24].

II. METHOD

A. Multivariate LSTM Architecture

The proposed model takes as input a sliding window of length $T = 50$ time steps with four features per step: CPU usage (%), memory usage (KB), network received throughput (KB/s), and network transmitted throughput (KB/s). The output is the prediction of CPU usage at the next time step. The multivariate approach was chosen because VM resource features are mutually correlated, and their joint patterns provide richer signals for detecting conditions that precede demand spikes — the primary source of under-prediction risk and SLO violations [3], [22], [23]. For example, a model that tracks concurrent increases in network throughput and memory pressure is more likely to predict an impending CPU spike than a model that only tracks CPU utilization. The basic architecture consists of one or more LSTM layers [7] dropout for regularization, and a single dense layer that generates the regression output at the end. Hyperparameter optimization is used to automatically select the number of layers, units per layer, dropout rate, learning rate, and batch size (see subsection C below).

The information flow within each LSTM cell is governed by three gating mechanisms that selectively retain, update, and expose memory over time. Given input x_t at time step t and the previous hidden state h_{t-1} , the forget gate f_t , input gate i_t , candidate cell \tilde{c}_t , cell state c_t , output gate o_t , and hidden state h_t are computed as equations (1)–(6) [7], [10], [11]:

$$f_t = \sigma(W_f \cdot [h_{t-1}, x_t] + b_f) \quad (1)$$

$$i_t = \sigma(W_i \cdot [h_{t-1}, x_t] + b_i) \quad (2)$$

$$\tilde{c}_t = \tanh(W_c \cdot [h_{t-1}, x_t] + b_c) \quad (3)$$

$$c_t = f_t \odot c_{t-1} + i_t \odot \tilde{c}_t \quad (4)$$

$$o_t = \sigma(W_o \cdot [h_{t-1}, x_t] + b_o) \quad (5)$$

$$h_t = o_t \odot \tanh(c_t) \quad (6)$$

where σ denotes the sigmoid activation function, \tanh is the hyperbolic tangent, and \odot is the element-wise (Hadamard) product. W_f , W_i , W_c , W_o are learnable weight matrices and b_f , b_i , b_c , b_o are the corresponding bias vectors. The forget gate determines what proportion of the prior cell state to discard; the input gate controls the

magnitude of new information written to the cell; and the output gate modulates what portion of the cell state is exposed as the hidden state h_t , which serves as both the recurrent feedback and the temporal feature representation passed to the next layer.

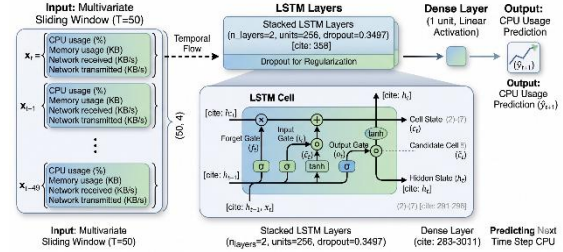


FIGURE 1. MULTIVARIATE LSTM ARCHITECTURE WITH SLO-AWARE LOSS FOR VM WORKLOAD PREDICTION

A Long Short-Term Memory (LSTM)-based architecture is shown in Fig. 1. for predicting CPU usage using multivariate time series data with a sliding window approach. The model receives an input sequence of length $T = 50$, where each time step includes four features: CPU usage (%), memory usage (KB), network received (KB/s), and network transmitted (KB/s), resulting in an input dimension of 50×4 . This sequential data is analyzed using stacked LSTM layers, which are made up of two layers with 256 hidden units each, to capture temporal patterns and system behavior across time. The model can learn both short-term and long-term dependencies thanks to the layered design, and a dropout mechanism with a rate of about 0.35 is used to lessen overfitting and enhance generalization. Three primary gating mechanisms are employed in each LSTM cell: the output gate, which governs the information transferred to the subsequent time step, the input gate, which decides how fresh information is assimilated, and the forget gate, which controls the erasure of irrelevant prior information. These gates enable the model to create a hidden state that reflects the acquired temporal properties and to maintain and update an internal memory (cell state). To generate the anticipated CPU consumption at the subsequent time step, the final hidden representation is subsequently transferred to a fully connected (dense) layer with a single neuron and a linear activation function. Overall, this architecture effectively models temporal relationships across multiple system metrics and generates continuous predictions of future CPU usage.

B. SLO-Aware Loss

The core methodological contribution of this study is the use of an asymmetric loss function designed to make the model more sensitive to under-prediction. The concept of asymmetric loss functions is rooted in the recent forecasting literature [19], and has been adapted to the context of network and cloud resource allocation [20]. Let \hat{y} be the predicted value and y the actual value; the SLO-aware loss is defined as in equation (7):

$$L(y, \hat{y}) = \alpha \cdot (y - \hat{y})^2 \text{ if } \hat{y} < y; (\hat{y} - y)^2 \text{ if } \hat{y} \geq y \quad (7)$$

where α is the under-prediction penalty factor. In this experiment, $\alpha = 10$, meaning that under-prediction errors are weighted ten times more heavily than over-prediction errors. Consequently, parameter optimization steers predictions toward a more conservative regime as a safeguard against potential SLO violations [25], [5].

C. Hyperparameter Optimization with Optuna

Hyperparameter search is performed using the Optuna framework [24] with the objective of minimizing `val_loss` (i.e., the SLO-aware loss on the validation set). Optuna was chosen because it supports a define-by-run API, sampling based on the Tree-structured Parzen Estimator (TPE), and automatic pruning to accelerate search convergence [24]. The search space comprises the number of LSTM layers (`n_layers`), number of units per layer (units), dropout rate, learning rate, and batch size. For each trial, the under-prediction rate on the validation set is also recorded as the operationally critical indicator of SLO risk; RMSE and MAE are additionally logged for contextual reference only.

D. Evaluation Metrics

The Under-prediction Rate, which is the percentage of test samples for which the projected value is less than the actual value ($\hat{y} < y$ is the main evaluation statistic used in this study. The operational condition that directly causes SLO breaches in cloud environments is this metric, which estimates the frequency of time steps during which the model's output would fall short of the actual workload need [1], [25], [5]. As stated in Section I, symmetric accuracy metrics like RMSE and MAE are insufficient as the only foundation for SLO-oriented model evaluation since they are unable to differentiate this directionally crucial failure pattern from over-prediction of comparable size. A model that more consistently makes conservative forecasts is shown by a reduced under-prediction rate, which lowers the danger of resource under-provisioning and safeguards SLO compliance.

RMSE (Root Mean Squared Error) and MAE (Mean Absolute Error) are recorded as secondary contextual metrics to explain each model's overall prediction performance and to make comparison with pertinent work that provides these standard benchmarks easier [13], [10], [11], [23]. They do not, however, function as the foundation for model selection or as the main criterion for assessing the effectiveness of the suggested strategy in this investigation. The study's main goal, which is to lower the chance of SLO violation through loss function design rather than to minimize average prediction error, is in line with the usage of under-prediction rate as the major indicator.

III. EXPERIMENTAL SCENARIO

A. Dataset

The experiments use the Bitbrain Grid Workloads Archive dataset (GWA-T-12), specifically the `fastStorage` collection [3]. The original dataset comprises performance traces of 1,250 VMs connected to fast Storage Area Network (SAN) devices in the Bitbrains datacenter, with a total of 193,720 rows of multivariate resource utilization observations after aggregation. This dataset has been widely adopted as a benchmark in VM workload prediction studies [8], [22], [13], [10], [15], [14], [11], [23], [17]. For computational efficiency in this preliminary study, a capped subset of 25,000 rows was used, ordered by time (without shuffling), with a fixed random seed of 42 to ensure reproducibility of hyperparameter search and weight initialization. The final size of each subset is shown in Table 1.

Four resource metrics are selected as input features for the prediction model. CPU usage (%) directly represents the computational demand on the VM and is the primary prediction target. Memory usage (KB) is included because memory pressure frequently correlates with CPU-intensive workloads, particularly for in-memory applications and database servers common in the Bitbrains enterprise environment [3], [15]. Network received throughput (KB/s) reflects inbound data processing demand, which drives CPU load in web-serving and streaming applications. Network transmitted throughput (KB/s) captures outbound activity, which correlates with CPU usage in request-response workloads. The inclusion of all four features aligns with the multivariate prediction approach validated by Xu et al. [22] and Dang-Quang and Yoo [23], both of whom demonstrated that incorporating correlated resource metrics enables the model to better capture the inter-feature dynamics that precede CPU demand spikes — precisely the patterns whose mis-prediction produces under-prediction events on the Bitbrain `fastStorage` trace. All features are sampled at 5-minute intervals, consistent with the original dataset's collection frequency [3].

TABLE 1. SHAPE OF EACH DATA SUBSET

Subset	Shape (samples, time_steps, features)
Training	(19,960, 50, 4)
Validation	(2,495, 50, 4)
Test	(2,495, 50, 4)

B. Pre-processing

The pre-processing pipeline consists of the following steps. First, all Bitbrain CSV files are merged into a single multivariate parquet file. Second, four multivariate features are selected for CPU prediction. Third, all features are normalized using `MinMaxScaler` to the range [0, 1]. Fourth,

time-series sequences are constructed using a 50-step sliding window, following standard practice in deep-learning-based cloud workload forecasting [9], [22]. Fifth, the data is split into training, validation, and test sets with an 80%/10%/10% ratio in chronological order without shuffling, preserving the temporal structure consistent with standard time-series evaluation practice [26].

C. Experimental Configuration

The experiment configuration is summarized in Table 2. Note that the relatively small values of `n_trials`, `tune_epochs`, and `final_epochs` reflect the pilot-study nature of this experiment; follow-up experiments with a larger compute budget are outlined in the conclusion section.

TABLE 2. EXPERIMENTAL CONFIGURATION

Parameter	Value
<code>time_steps</code>	50
<code>Split (train/val/test)</code>	80% / 10% / 10%
<code>Loss function (Optuna)</code>	SLO-aware loss ($\alpha=10$)
<code>Loss function (baseline)</code>	MSE (standard)
<code>n_trials (Optuna)</code>	2
<code>tune_epochs</code>	2
<code>final_epochs</code>	2
<code>max_rows</code>	25,000
<code>random_seed</code>	42

D. Execution Steps

The execution pipeline proceeds as follows: (1) ingest all Bitbrain CSV files into a single multivariate parquet file; (2) select and normalize features using `MinMaxScaler`; (3) construct sliding-window time-series sequences; (4) split the data chronologically; (5) tune LSTM hyperparameters with Optuna using SLO-aware loss as the objective; (6) record the under-prediction rate as the primary SLO indicator for each trial, with RMSE, MAE, and MSE logged as secondary contextual metrics; (7) retrain two models using the winning hyperparameters from the Optuna objective — one with SLO-aware loss and one with standard MSE as a baseline comparison, keeping the architecture and number of epochs identical; and (8) evaluate both models on the test set with under-prediction rate as the primary criterion.

IV. RESULT AND DISCUSSION

A. Hyperparameter Optimization Results

Optuna executed two trials using SLO-aware loss as the optimization objective. Both trials produced a validation under-prediction rate of 0.00% — the primary evaluation metric — meaning neither trial produced any under-prediction on the validation set. The trial with the lowest `val_loss` (Trial 1, the Optuna objective winner) yielded the

configuration shown in Table 3, recording a validation RMSE of 0.018780 and MSE of 0.00035270 as secondary contextual metrics.

TABLE 3. BEST HYPERPARAMETERS FROM OPTUNA OBJECTIVE

Hyperparameter	Value
<code>n_layers</code>	2
<code>units</code>	256
<code>dropout</code>	0.3497
<code>lr</code>	0.000266
<code>batch_size</code>	128

Trial 0, which also achieved a validation under-prediction rate of 0.00%, used a shallower configuration (`n_layers=2`, `units=64`, `dropout=0.1468`, `lr=0.000205`, `batch_size=128`) and recorded a validation RMSE of 0.026515 and MSE of 0.00070304. Because both trials achieved an identical under-prediction rate of 0.00% on the validation set — meaning the primary evaluation criterion was tied — RMSE served as a tiebreaker criterion, and the Trial 1 configuration (lower RMSE: 0.018780 vs 0.026515) was selected as the winning hyperparameter set for both the SLO-aware and MSE baseline models.

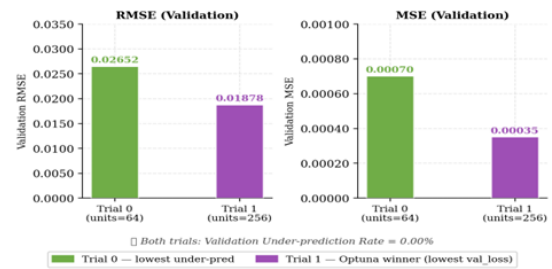


Figure 2. Optuna Trial Comparison on Validation Set

Figure 1 presents a side-by-side comparison of the two Optuna trials on the validation set in terms of RMSE (left) and MSE (right). This figure is included specifically to document the tiebreaker decision, not as evidence of the primary evaluation outcome. Both trials achieved an identical validation under-prediction rate of 0.00% — the primary evaluation metric of this study — meaning neither model under-predicted on any validation sample, and both are equally qualified from the SLO-protection standpoint. Since the primary criterion was tied, RMSE was used as a secondary tiebreaker: Trial 1 (256 units) achieved a lower validation RMSE of 0.01878 compared to 0.02652 for Trial 0 (64 units), and was therefore selected as the winning configuration. This tiebreaker role is the only context in which RMSE influences model selection in this study. It is also important to note that with only two trials, the Optuna search covers a very limited portion of the hyperparameter space; future work with 30–50 trials is expected to explore configurations that directly optimize under-prediction rate at larger scale.

A. Test-Set Performance

After retraining with the winning hyperparameters, both models — one trained with SLO-aware loss and one with standard MSE — were evaluated on the test set. The primary evaluation metric in this study is the under-prediction rate, which directly quantifies how frequently the model produces predictions that fall below the actual workload value — the condition that triggers potential SLO violations. Results are shown in Table 4.

TABLE 4. UNDER-PREDICTION RATE COMPARISON ON TEST SET

Model	Under-prediction Rate	Reduction
SLO-aware Loss (proposed)	0.04%	4× lower
MSE Baseline (comparison)	0.16%	—

Figure 2 presents the under-prediction rate — the central metric of this study — for both models on the test set. The SLO-aware model achieves an under-prediction rate of only 0.04%, compared to 0.16% for the MSE baseline, corresponding to a fourfold reduction in the frequency of predictions that fall below the actual workload value. In the context of cloud resource management, each under-prediction event represents a time step during which the provisioned capacity may be insufficient to serve the actual demand, potentially triggering an SLO violation [1], [9]. This fourfold difference is the direct consequence of the asymmetric loss formulation: by imposing a tenfold penalty on under-prediction errors during training, the SLO-aware model is structurally guided to favor over-estimation over under-estimation. The annotation '4× lower SLO risk' in the figure communicates this operational significance explicitly. Although both rates are low in absolute terms — a consequence of the right-skewed Bitbrain workload distribution and the limited training epochs as discussed in subsection C — the relative difference of 4× between the two models is consistent and meaningful, directly reflecting the design intent of the SLO-aware loss function.

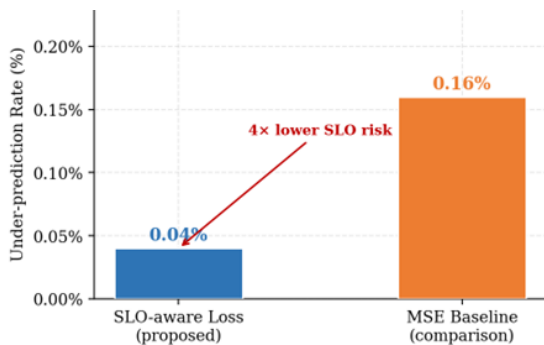


Figure 3. Under-prediction Rate on Test Set (Lower = Fewer Potential SLO Violations)

C. Discussion

Empirical validation of SLO-aware loss effectiveness. The central finding of this study is that the SLO-aware model achieves an under-prediction rate of 0.04% on the test set, compared to 0.16% for the MSE baseline — a fourfold reduction. This empirically confirms, rather than merely theoretically asserts, that asymmetric loss penalization directly reduces the frequency of under-prediction events. Importantly, both models were trained with identical architectures and hyperparameters, differing only in loss function. The observed difference in under-prediction behavior is therefore attributable solely to the loss formulation, providing a controlled and interpretable experimental result. The SLO-aware model is structurally biased by its training objective to avoid under-estimation, resulting in predictions that more consistently sit at or above the actual workload value — precisely the behavior required in SLO-critical cloud environments where capacity shortfalls carry heavier consequences than marginal over-allocation [25], [5].

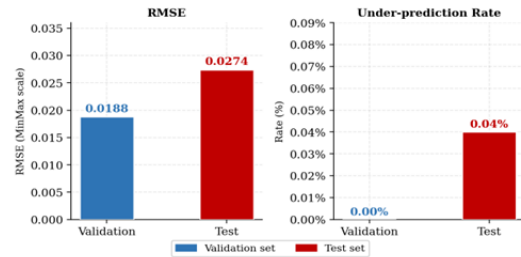


Figure 4. Validation vs. Test Performance — SLO-aware Model

Figure 3 examines the generalization behavior of the SLO-aware model by comparing the under-prediction rate between the validation set and the test set (right panel). On the validation set, the model records a perfect under-prediction rate of 0.00%, whereas on the test set this value rises slightly to 0.04%. Critically, even on unseen test data, the model maintains an extremely low under-prediction rate, demonstrating that the conservative prediction tendency instilled by the SLO-aware loss generalizes robustly beyond the training and validation phases. This is the key generalization finding of this paper: the asymmetric loss creates a durable directional bias — the model consistently prefers over-estimation over under-estimation — and this behavior transfers reliably to the test distribution. The RMSE panel (left) is shown for completeness to characterize overall model behavior, but it is not the primary evaluation criterion of this study.

Comparison with related work. Prior studies on VM workload prediction — including CNN-LSTM [10], LSRU [11], TCN-GRU-attention [13], and multivariate Bi-LSTM [23] — focus primarily on minimizing symmetric accuracy metrics such as

RMSE and MAE using symmetric loss functions. These formulations implicitly treat under-prediction and over-prediction as equally undesirable, which does not reflect the asymmetric cost structure of SLO-governed cloud operations. The proposed SLO-aware loss addresses this gap directly by embedding the operational cost asymmetry into the training objective itself. Rather than relying on post-hoc thresholding or reactive scaling to compensate for under-prediction, the model learns proactively to avoid the under-estimation region. This study demonstrates that loss function design is a complementary and actionable dimension alongside architecture design in cloud workload prediction, and that even a simple multivariate LSTM can achieve strong SLO-protective behavior when equipped with an appropriate loss formulation.

Implications for cloud resource management. A fourfold reduction in under-prediction rate translates directly to a fourfold reduction in the expected frequency of SLO-violating time steps during real deployment. In production environments where SLO penalty clauses are financially significant — such as financial services, e-commerce, and enterprise cloud platforms [1], [5] — this improvement has concrete business value that extends well beyond a marginal improvement in RMSE. Cloud operators can deploy the SLO-aware model as part of a proactive auto-scaling pipeline, where the conservative prediction bias preemptively triggers resource scale-up before demand spikes materialize, rather than reacting after violations occur [4]. Future work can extend this by parameterizing α dynamically based on SLO tier or workload volatility, enabling fine-grained control over the aggressiveness of the SLO protection strategy.

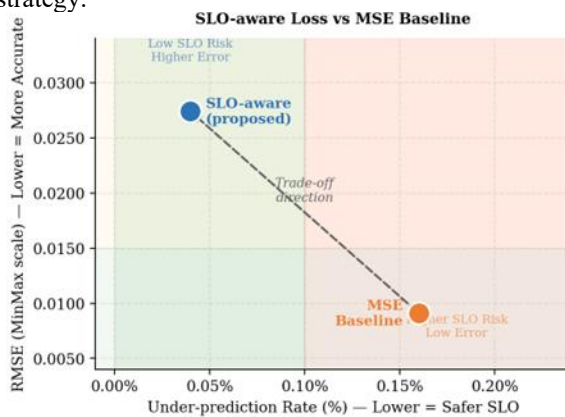


Figure 5. Accuracy-SLO Risk Trade-off: SLO-aware Loss vs MSE Baseline

Figure 4 provides a two-dimensional visualization of the relationship between the two models, with RMSE on the vertical axis and under-prediction rate on the horizontal axis. While this figure includes RMSE for spatial context, the horizontal axis — under-prediction rate — represents the primary dimension of interest in this

study. The green-shaded region on the left (under-prediction rate $< 0.10\%$) denotes the low SLO-risk zone, and the red-shaded region on the right represents higher SLO risk. The SLO-aware model (blue circle) sits firmly in the low SLO-risk zone, while the MSE baseline (orange circle) falls in the higher-risk zone. The dashed arrow illustrates that moving from the MSE baseline toward the SLO-aware model along the horizontal axis represents the primary objective of this work: reducing under-prediction risk. In SLO-critical deployments, the horizontal position of a model — not its RMSE — determines whether it is operationally safe, making the SLO-aware model the preferable choice for the target deployment context [1], [5].

Analysis of low absolute under-prediction rates in both models. A noteworthy observation is that both models achieved remarkably low absolute under-prediction rates (0.04% and 0.16%), which may appear surprising given the limited training budget of only 2 epochs. This behavior is attributable to two concurrent factors. First, the Bitbrain fastStorage workload trace exhibits a right-skewed distribution characteristic of enterprise computing environments: CPU utilization frequently remains at low-to-moderate levels, punctuated by occasional high-demand spikes [3], [13]. In such a distribution, a model that has not yet tightly fitted the training data tends to predict values clustered around the central tendency, which for right-skewed data sits below the mean but above many observations — naturally producing more over-predictions than under-predictions. Second, because both models were trained for only 2 epochs, neither has converged to a distribution-specific prediction surface; instead, both operate in a regime of partial underfitting that coincidentally aligns with the asymmetry of the workload distribution. As training epochs increase in future experiments, the absolute under-prediction rates of both models are expected to shift, and the relative advantage of SLO-aware loss over MSE is expected to become more pronounced as the models develop more refined prediction surfaces that can distinguish spike from baseline.

D. Limitations

Several limitations of this preliminary study should be acknowledged transparently. First, the dataset used comprises only 25,000 rows from the fastStorage trace, so generalization to the full dataset and other workload patterns still requires validation [13], [10]. Second, the compute budget (2 epochs and 2 Optuna trials) remains severely constrained; a larger budget would likely improve both models and may further reduce the under-prediction rate of the SLO-aware model. Third, the penalty factor $\alpha = 10$ was set in an ad-hoc manner; a systematic sensitivity analysis across different α values is needed to identify the optimal value that

minimizes under-prediction rate while containing over-allocation cost. Fourth, although the fixed random seed (42) improves reproducibility, cross-seed validation has not been performed to confirm the stability of the under-prediction rate findings across different random initializations.

V. CONCLUSION

This study proposed a multivariate LSTM-based VM workload prediction approach equipped with an SLO-aware loss that imposes asymmetric penalties on under-prediction. The primary research objective was to reduce the under-prediction rate — the direct measure of SLO violation risk — rather than to minimize general accuracy metrics. Experiments on a 25,000-row subset of the Bitbrain GWA-T-12 fastStorage dataset, conducted with a fixed random seed of 42, enabled a controlled comparison between the SLO-aware model and an MSE baseline trained with identical architecture and hyperparameters. The SLO-aware model achieved an under-prediction rate of 0.04% compared to 0.16% for the MSE baseline — a fourfold reduction — empirically confirming that embedding operational cost asymmetry into the loss function is an effective and targeted strategy for protecting SLO compliance in cloud VM workload prediction. Future research directions include scaling to the full Bitbrain dataset; allocating an adequate training budget with sufficient epochs and at least 30–50 Optuna trials; conducting a sensitivity study on the SLO-aware loss penalty parameter α to identify the Pareto-optimal point between under-prediction rate and over-allocation cost; comparing the under-prediction rate of the proposed approach against advanced hybrid architectures such as CNN-LSTM, TCN-GRU-attention, LSRU, and multivariate Bi-LSTM; implementing the model in an end-to-end auto-scaling simulation to measure the business impact of reduced SLO violations; and validating on other workload trace datasets such as Google Cluster Trace or Alibaba Cluster Trace to confirm generalizability.

REFERENCE

- [1] J. Dogani, R. Namvar, and F. Khunjush, "Auto-scaling techniques in container-based cloud and edge/fog computing: Taxonomy and survey," Sep. 01, 2023, *Elsevier B.V.* doi: 10.1016/j.comcom.2023.06.010.
- [2] S. Deng *et al.*, "Cloud-Native Computing: A Survey From the Perspective of Services," *Proceedings of the IEEE*, vol. 112, no. 1, pp. 12–46, Jan. 2024, doi: 10.1109/JPROC.2024.3353855.
- [3] S. Shen, V. Van Beek, and A. Iosup, "Statistical characterization of business-critical workloads hosted in cloud datacenters," in *Proceedings - 2015 IEEE/ACM 15th International Symposium on Cluster, Cloud, and Grid Computing, CCGrid 2015*, Institute of Electrical and Electronics Engineers Inc., Jul. 2015, pp. 465–474. doi: 10.1109/CCGrid.2015.60.
- [4] M. Masdari and A. Khoshnevis, "A survey and classification of the workload forecasting methods in cloud computing," *Cluster Comput.*, vol. 23, no. 4, pp. 2399–2424, Dec. 2020, doi: 10.1007/s10586-019-03010-3.
- [5] P. Souza *et al.*, "Predicting and Avoiding SLA Violations of Containerized Applications using Machine Learning and Elasticity," in *International Conference on Cloud Computing and Services Science, CLOSER - Proceedings*, Science and Technology Publications, Lda, 2022, pp. 74–85. doi: 10.5220/0011085100003200.
- [6] M. P. Yadav, N. Pal, and D. K. Yadav, "Workload prediction over cloud server using time series data," in *Proceedings of the Confluence 2021: 11th International Conference on Cloud Computing, Data Science and Engineering*, Institute of Electrical and Electronics Engineers Inc., Jan. 2021, pp. 267–272. doi: 10.1109/Confluence51648.2021.9377032.
- [7] K. Greff, R. K. Srivastava, J. Koutnik, B. R. Steunebrink, and J. Schmidhuber, "LSTM: A Search Space Odyssey," *IEEE Trans. Neural Netw. Learn. Syst.*, vol. 28, no. 10, pp. 2222–2232, Oct. 2017, doi: 10.1109/TNNLS.2016.2582924.
- [8] M. E. Karim, M. M. S. Maswood, S. Das, and A. G. Alharbi, "BHyPreC: A Novel Bi-LSTM Based Hybrid Recurrent Neural Network Model to Predict the CPU Workload of Cloud Virtual Machine," *IEEE Access*, vol. 9, pp. 131476–131495, 2021, doi: 10.1109/ACCESS.2021.3113714.
- [9] L. Nashold and R. Krishnan, "Using LSTM and SARIMA Models to Forecast Cluster CPU Usage," Jul. 2020, [Online]. Available: <http://arxiv.org/abs/2007.08092>
- [10] H. L. Leka, Z. Fengli, A. T. Kenea, A. T. Tegene, P. Atandoh, and N. W. Hundera, "A Hybrid CNN-LSTM Model for Virtual Machine Workload Forecasting in Cloud Data Center," in *2021 18th International Computer Conference on Wavelet Active Media Technology and Information Processing, ICCWAMTIP 2021*, Institute of Electrical and Electronics Engineers Inc., 2021, pp. 474–478. doi: 10.1109/ICCWAMTIP53232.2021.9674067.
- [11] Md. N. H. Shuvo, M. M. S. Maswood, and A. G. Alharbi, "LSRU: A Novel Deep Learning based Hybrid Method to Predict the Workload of Virtual Machines in Cloud Data Center," in *2020 IEEE Region 10 Symposium (TENSYP): 5-7 June 2020, Dhaka, Bangladesh*, IEEE, Jun. 2020.
- [12] N. M. Dang-Quang and M. Yoo, "Deep learning-based autoscaling using bidirectional long short-term memory for kubernetes," *Applied Sciences (Switzerland)*, vol. 11, no. 9, May 2021, doi: 10.3390/app11093835.
- [13] Y. Bai, L. Chen, Y. Lei, and H. Xie, "A Deep Learning Prediction Approach for Machine Workload in Cloud Computing," in *2023 5th International Conference on Data-Driven Optimization of Complex Systems, DOCS 2023*, Institute of Electrical and Electronics Engineers Inc., 2023. doi: 10.1109/DOCS60977.2023.10294685.
- [14] A. Bhardwaj, Y. R. Sai, and N. V. Patil, "Deep Learning Models for Workload Prediction in Virtual Machine Cloud Data Center," Institute of Electrical and Electronics Engineers (IEEE), Dec. 2025, pp. 1–6. doi: 10.1109/conecct65861.2025.11306757.
- [15] M. Chen, M. R. Read, P. Arroba, and R. Buyya, "EN-Beats: A Novel Ensemble Learning-based Method for Multiple Resource Predictions in Cloud." [Online]. Available: <http://gwa.ewi.tudelft.nl/datasets/gwa-t-12-bitbrains>
- [16] S. Bansal and M. Kumar, "Deep Learning-based Workload Prediction in Cloud Computing to Enhance the Performance," in *ICSCCC 2023 - 3rd International Conference on Secure Cyber Computing and Communications*, Institute of Electrical and Electronics Engineers Inc., 2023, pp. 635–640. doi: 10.1109/ICSCCC58608.2023.10176790.
- [17] N. I. Mahbub, M. D. Hossain, S. Akhter, M. I. Hossain, K. Jeong, and E. N. Huh, "Robustness of Workload

- Forecasting Models in Cloud Data Centers: A White-Box Adversarial Attack Perspective,” *IEEE Access*, vol. 12, pp. 55248–55263, 2024, doi: 10.1109/ACCESS.2024.3385863.
- [18] S. Bhagavathiperumal and M. Goyal, “Workload Analysis of Cloud Resources using Time Series and Machine Learning Prediction,” in *2019 IEEE Asia-Pacific Conference on Computer Science and Data Engineering (CSDE)*, Aug. 2020. doi: 10.1109/CSDE48274.2019.9162385.
- [19] S. S. Rajpal, R. Mahadeva, A. K. Goyal, and V. Sarda, “Improving Forecasting Accuracy of Stock Market Indices Utilizing Attention-Based LSTM Networks with a Novel Asymmetric Loss Function,” *AI (Switzerland)*, vol. 6, no. 10, Oct. 2025, doi: 10.3390/ai6100268.
- [20] X. Zhou, X. Liu, D. Zhai, J. Jiang, and X. Ji, “Asymmetric Loss Functions for Noise-Tolerant Learning: Theory and Applications,” *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 45, no. 7, 2023, doi: 10.1109/TPAMI.
- [21] Y. S. Patel and J. Bedi, “MAG-D: A multivariate attention network based approach for cloud workload forecasting,” *Future Generation Computer Systems*, vol. 142, pp. 376–392, May 2023, doi: 10.1016/j.future.2023.01.002.
- [22] M. Xu, C. Song, H. Wu, S. S. Gill, K. Ye, and C. Xu, “esDNN: Deep Neural Network Based Multivariate Workload Prediction in Cloud Computing Environments,” *ACM Trans. Internet Technol.*, vol. 22, no. 3, Aug. 2022, doi: 10.1145/3524114.
- [23] N. M. Dang-Quang and M. Yoo, “Multivariate Deep Learning Model For Workload Prediction In Cloud Computing,” in *International Conference on ICT Convergence*, IEEE Computer Society, 2021, pp. 858–862. doi: 10.1109/ICTC52510.2021.9620931.
- [24] T. Akiba, S. Sano, T. Yanase, T. Ohta, and M. Koyama, “Optuna: A Next-generation Hyperparameter Optimization Framework,” in *Proceedings of the ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, Association for Computing Machinery, Jul. 2019, pp. 2623–2631. doi: 10.1145/3292500.3330701.
- [25] J. Ghanim, M. Issa, and M. Awad, “An Asymmetric Loss with Anomaly Detection LSTM Framework for Power Consumption Prediction,” Feb. 2023, doi: 10.1109/MELECON53508.2022.9842895.
- [26] R. J. Hyndman and G. Athanasopoulos, *Forecasting: Principles & Practice*. Melbourne, Australia: OTexts, 2021.

AGUS HARIYANTO is a researcher and lecturer at Information Technology Department, Politeknik Negeri Jember Indonesia . His research interests include computer network, security network, cloud computing, and internet of things. He has published papers on network management using cloud computing, network security implementation and internet of things implementation in applications.

AHMAD FAHRIYANNUR ROSYADY, was born in jember, East Java Indonesia, in 1992. He Received the Bachelor Degree from Bina Nusantara University Jakarta, in 2017 in Informatic Engineering and the Master Degree From Institut Teknologi Sepuluh Nopember, Surabaya Indonesia, in Technology Management of Information System. His Research Interest Include Information Technology, Internet Of things, Machine Learning, Artificial Intelegent..

ADI SUCIPTO.

BEKTI MARYUNI SUSANTO. was born in Yogyakarta Province, Indonesia, in 1984. He received the Bachelor degree from the Yogyakarta State University, Indonesia in 2010 in Electrical Engineering Education and the Master degree from the STMIK Nusa Mandiri Jakarta, Indonesia, in 2012, in Computer Science. His research interests include cloud computing, internet of things, and machine learning. He can be contacted at email: bekti@polije.ac.id..

SAPTA NUGRAHA.