



Greylag Goose Optimization-Driven EALSTM for Accurate HVAC Chiller Energy Prediction

Doaa Sami Khafaga^{1,*}

¹Department of Computer Sciences, College of Computer and Information Sciences, Princess Nourah bint Abdulrahman University, P.O. Box 84428, Riyadh 11671, Saudi Arabia

Email: dskhafga@pnu.edu.sa

Abstract

Forecasting the energy consumption of heating, ventilation, and air conditioning (HVAC) chillers is vital for enhancing building efficiency, reducing operating costs, and supporting sustainability goals. However, the task remains challenging due to nonlinear system dynamics, strong dependence on weather conditions, and the scarcity of high-quality real-world datasets. In this work, we employ the *Chiller Energy Data* from Kaggle, which contains 13,561 cleaned records collected between August 2019 and June 2020, incorporating ten operational and meteorological features. Six baseline models, namely the Evolutionary Attention-based Long Short-Term Memory (EALSTM), Bidirectional LSTM (BiLSTM), standard LSTM, Gated Recurrent Unit (GRU), Temporal Convolutional Network (TCN), and Artificial Neural Network (ANN), are first benchmarked to assess their forecasting capability. To further improve predictive accuracy, we integrate EALSTM with ten meta-heuristic optimization algorithms, focusing on the Greylag Goose Optimization Algorithm (GGO) and comparing it with alternatives such as Harris Hawks Optimization (HHO), Artificial Physics Optimization (APO), Simulated Annealing Optimization (SAO), Grey Wolf Optimizer (GWO), and others. The optimized GGO+EALSTM framework achieves state-of-the-art performance with a mean squared error of 6.83×10^{-6} and an R^2 value of 0.98, reflecting a 96% reduction in error relative to simple feedforward models and significant improvements over other recurrent networks and optimizer-enhanced variants. The main contributions of this study include a structured benchmarking of neural architectures for chiller forecasting, the first systematic comparison of ten meta-heuristic optimizers applied to deep learning in this domain, and a visualization-based error analysis that strengthens interpretability and supports practical deployment. These results establish optimization-enhanced EALSTM as a robust and generalizable framework for HVAC energy forecasting, paving the way toward more efficient, reliable, and sustainable building energy management.

Keywords: HVAC energy forecasting; Chiller energy consumption; Evolutionary attention-based LSTM; Meta-heuristic optimization; Greylag Goose Optimization

1 Introduction

Heating, ventilation, and air conditioning (HVAC) systems—particularly chillers—represent a substantial portion of building energy consumption. In commercial and industrial buildings, chillers alone account for more than 40% of total energy use [1], [2]. Consequently, improving the accuracy of energy demand

predictions is essential to enhance energy efficiency, reduce operational costs, and mitigate environmental impact by lowering carbon emissions.

Forecasting chiller energy consumption is inherently challenging due to the nonlinear dynamics of HVAC systems, the strong influence of weather variability, and the scarcity of comprehensive, high-quality real-world datasets. These constraints hinder model generalization and highlight the necessity of robust modeling strategies capable of capturing complex system behavior even under limited data availability.

Recent advances in machine learning (ML) and deep learning (DL) have shown considerable promise in the field of energy forecasting. Recurrent neural networks (RNNs), convolutional neural networks (CNNs), and Transformer-based architectures are increasingly used to capture temporal and spatial dependencies within energy data [3], [4], [5]. In particular, long short-term memory (LSTM) networks and their variants have demonstrated strong performance for sequence modeling tasks in energy prediction [6], [7]. Attention mechanisms have further enhanced forecasting performance by improving feature weighting in dynamic conditions [8]. Nevertheless, the majority of existing works either focus on general energy demand prediction or specific HVAC subsystems, and only a limited number of studies benchmark multiple recurrent and neural architectures specifically for chiller energy forecasting.

In addition to the modeling approaches, optimization algorithms have emerged as a critical tool for enhancing the predictive power of ML models. Meta-heuristic optimization methods such as Particle Swarm Optimization (PSO), Genetic Algorithm (GA), Grey Wolf Optimizer (GWO), and Whale Optimization Algorithm (WOA) have been successfully employed for hyperparameter tuning and weight adjustment in various energy forecasting applications [9], [10], [11], [12]. More recently, novel bio-inspired optimizers like the Greylag Goose Optimization Algorithm (GGO) have been introduced, offering improved exploration–exploitation balance and robustness in avoiding local minima [13]. However, a systematic comparison of these optimization approaches when integrated with recurrent forecasting models for HVAC chiller systems remains absent in the literature.

This study aims to address the above gaps by conducting a comprehensive evaluation of six baseline ML and DL models on a real-world chiller dataset sourced from Kaggle, comprising 13,561 cleaned data samples with operational and meteorological features. The models benchmarked include Evolutionary Attention-based LSTM (EALSTM), Bidirectional LSTM (BiLSTM), standard LSTM, Gated Recurrent Unit (GRU), Temporal Convolutional Network (TCN), and Artificial Neural Network (ANN). We further propose an optimization-enhanced forecasting framework that integrates EALSTM with the Greylag Goose Optimization Algorithm (GGO), and systematically compare its performance against nine other meta-heuristic optimizers under identical conditions.

The contributions of this paper are fourfold. First, we establish a rigorous benchmarking of recurrent and neural forecasting models on a chiller energy dataset, evaluated using a comprehensive suite of metrics including mean squared error (MSE), root mean squared error (RMSE), mean absolute error (MAE), mean bias error (MBE), correlation coefficient (r), coefficient of determination (R^2), relative RMSE (RRMSE), Nash–Sutcliffe Efficiency (NSE), and Willmott's Index (WI). Second, we demonstrate that the GGO-enhanced EALSTM achieves superior forecasting performance compared to both baselines and other optimization-augmented models. Third, we provide detailed error visualizations, including residual analysis and temporal prediction plots, to enhance interpretability and practical applicability. Finally, our work establishes a methodological framework for future research on optimized deep learning models for HVAC energy forecasting, particularly under data-constrained conditions.

By combining structured benchmarking, optimization-enhanced modeling, and in-depth error analysis, this study contributes to advancing the state of the art in data-driven building energy management and sustainability.

2 Related Work

The accurate prediction and optimization of energy consumption in heating, ventilation, and air conditioning (HVAC) systems has attracted extensive research attention in recent years, with diverse approaches ranging from conventional machine learning models to advanced deep learning and meta-heuristic optimization frameworks. This literature can be thematically grouped into three broad categories: machine learning for load forecasting, deep learning models for HVAC prediction, and optimization-driven or reinforcement learning strategies for energy management.

Machine Learning for Load Forecasting

Early research has highlighted the value of machine learning models in predicting building cooling loads by leveraging external weather conditions, building characteristics, and operational parameters. Almeida et al. [14] conducted a comparative analysis of several machine learning and ensemble learning algorithms for cooling energy prediction at the European Central Bank building in Frankfurt. Their findings demonstrated the superiority of the Random Forest model across different feature sets, with consistently lower error metrics compared to alternative methods such as support vector regression, gradient boosting, and recurrent networks. These works underscore the importance of both model selection and feature engineering in achieving reliable forecasts.

Deep Learning Models for HVAC Prediction

With the growing availability of operational data, deep learning approaches have become increasingly popular for HVAC load forecasting. Heidarykiany and Ababei [15] investigated the sensitivity of long short-term memory (LSTM) networks to architectural and training parameters in predicting daily residential HVAC energy usage. Their analysis demonstrated that even minimalistic LSTM models, when carefully optimized, can yield strong predictive performance, highlighting the importance of hyperparameter tuning. Building on this direction, Zhang [16] combined Gaussian Process Regression with meta-heuristic algorithms such as the Weevil Damage Optimization Algorithm (WDOA) and Improved Manta-Ray Foraging Optimizer (IMRFO) to achieve highly accurate cooling load prediction, reporting an R^2 value of 0.99.

Optimization and Reinforcement Learning Strategies

Another line of research has focused on optimization-driven control strategies for HVAC energy management. Peng et al. [17] proposed a deep reinforcement learning (DRL) framework that integrates a CNN-LSTM prediction model with an enhanced deep deterministic policy gradient algorithm, achieving significant improvements in both prediction accuracy and real-time control efficiency. In parallel, Wang et al. [18] introduced a DQN-based coordination method for jointly managing HVAC systems and energy storage, reducing building operating costs while maintaining occupant comfort. Beyond HVAC, related studies have also explored sustainable energy management in broader contexts. For example, Hsu et al. [19] demonstrated that workload allocation optimization in data centers can reduce cooling-related energy use by more than 50%, emphasizing the relevance of optimization strategies for HVAC in computational environments. Complementing these works, Mengru et al. [20] investigated transfer learning approaches to overcome data scarcity in chiller systems, showing that multi-source transfer models reduced prediction error by up to 26.6%, thereby providing technical support for applications with limited operational data.

Synthesis

Taken together, the literature reveals a clear trajectory in the evolution of HVAC energy prediction research. Initial works demonstrated the benefits of machine learning models enhanced by feature engineering, followed by deep learning approaches that leveraged temporal and nonlinear dependencies more effectively. Recent studies have increasingly integrated meta-heuristic optimization and reinforcement learning to address the challenges of parameter sensitivity, nonlinearity, and real-time control. Despite this progress, systematic benchmarking of advanced recurrent models alongside a structured evaluation of multiple optimizers, particularly in the context of real-world chiller datasets, remains underexplored. The present study seeks to address this gap by evaluating six baseline neural models and comparing the performance of ten state-of-the-art optimization algorithms, with particular emphasis on the combination of evolutionary attention-based LSTM and Greylag Goose Optimization.

3 Dataset and Preprocessing

The dataset employed in this study is the *Chiller Energy Data* available on the Kaggle repository, which contains operational and weather-related variables for a commercial building located in Singapore. The temporal span of the dataset extends from August 18, 2019, to June 1, 2020, thus covering nearly ten months of continuous operation. After data cleaning and refinement, the final dataset consists of 13,561 usable records. The removal process included eliminating outliers, filtering corrupted entries, and addressing missing values to ensure reliability and consistency of the input data.

The recorded variables comprise one temporal feature (timestamp) and nine operational and meteorological inputs. Specifically, the operational parameters include chilled water flow rate (L/s), cooling water temperature (°C), building load (RT), and chiller energy consumption (kWh). In addition, weather-related attributes are represented by outside air temperature (°F), dew point (°F), humidity (%), wind speed (mph), and atmospheric pressure (inHg). These features together capture both the physical state of the HVAC system and the external climatic conditions, providing a comprehensive foundation for predictive modeling.

The processed dataset exhibits heterogeneous distributions across different variables. For instance, the chiller load and energy consumption display a right-skewed distribution, with relatively few high-demand peaks compared to the majority of normal operating conditions. Similarly, meteorological variables such as humidity and temperature demonstrate seasonal and diurnal variations, reflecting the tropical climate of Singapore. Such variability introduces additional complexity into the modeling task, as learning algorithms must effectively generalize across diverse operating conditions.

To prepare the dataset for model training, several preprocessing steps were applied. First, all continuous variables were normalized to a [0,1] range using min-max scaling, which ensured numerical stability during model optimization and prevented features with larger scales from dominating the learning process. Second, the dataset was divided into training, validation, and test subsets in a time-aware manner, thereby avoiding data leakage across temporal splits and ensuring that future data was never used in training phases. The allocation of samples followed a chronological split to maintain realistic forecasting conditions. Third, feature engineering was performed by generating lagged versions of key variables, enabling models to incorporate temporal dependencies explicitly. Lag features were particularly critical for capturing the delayed impact of weather fluctuations on building cooling loads and energy consumption.

The cross-correlation analysis (Fig. 1) investigates the dynamic relationship between building load (RT) and chilled water flow rate. Strong positive correlation at lower lags confirms the physical dependency of water flow on building cooling demand. Negative correlation peaks at intermediate lags suggest delayed oscillatory

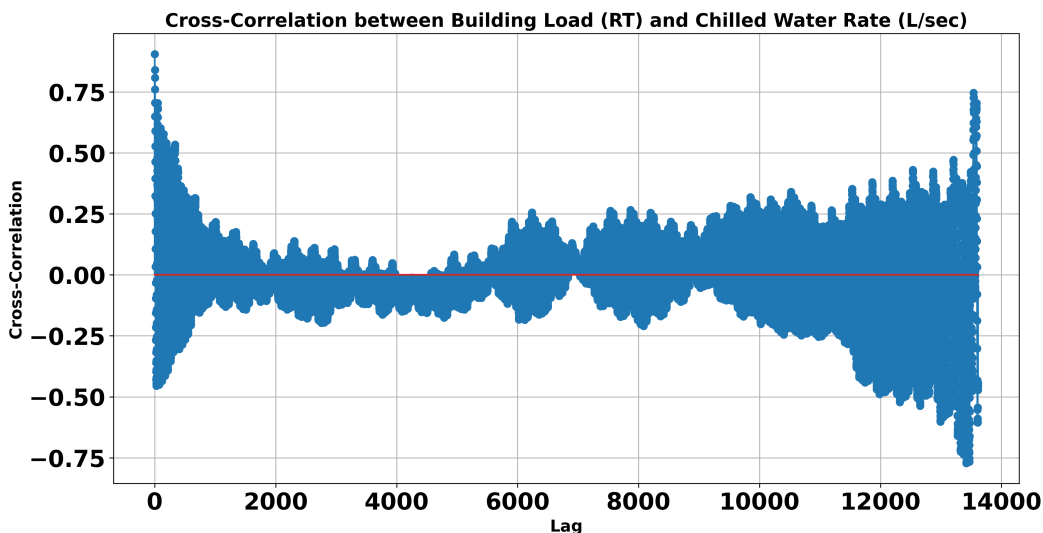


Figure 1: Cross-correlation between building load (RT) and chilled water rate. High positive correlation at small lags highlights direct dependency, while alternating negative correlations reflect delayed system responses.

patterns, likely reflecting operational control adjustments. This reinforces the inclusion of lagged features in predictive modeling.

The correlation heatmap of features (Fig. 2) reveals interdependencies among building, chiller, and environmental variables. Building load, chilled water rate, and chiller energy consumption exhibit the strongest correlations ($r > 0.85$), confirming their joint role in cooling demand. External weather factors, such as outside temperature and humidity, show moderate to strong correlations, indicating their indirect impact on system load. Negative correlation between humidity and energy consumption highlights its inverse influence on thermal loads.

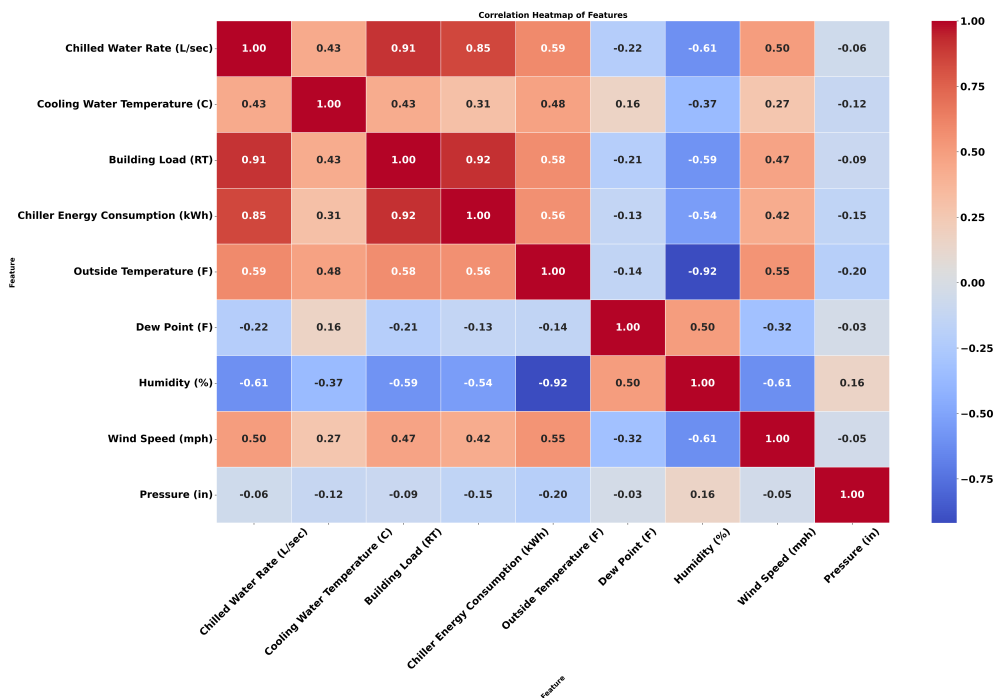


Figure 2: Correlation heatmap of dataset features. Strong associations are observed between building load, chilled water rate, and chiller energy use, while weather variables provide secondary but significant influences.

Figure 3 illustrates the normalized building load (RT) over the observation period. The series displays daily cyclicity and seasonal fluctuations, with higher demand during warmer months and sporadic spikes indicating peak cooling requirements. Outlier dips correspond to possible sensor errors or temporary equipment shutdowns. These time-series dynamics highlight the necessity of temporal models that capture both short-term fluctuations and seasonal variations.

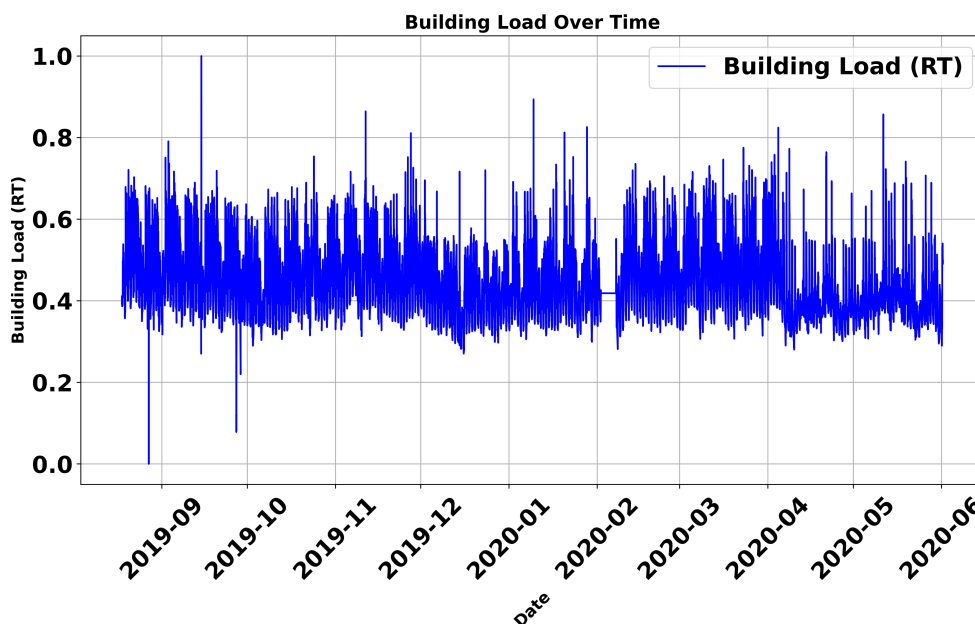


Figure 3: Time-series evolution of normalized building load (RT) from Aug 2019 to May 2020. Seasonal and daily cycles dominate, with occasional anomalies reflecting operational irregularities.

Residual analysis (Fig. 4) evaluates the deviation of predicted versus fitted building load values. The residuals are centered near zero, indicating unbiased predictions overall. However, minor heteroscedasticity is visible, with wider spread at higher fitted loads, suggesting increasing uncertainty under peak demand. This diagnostic emphasizes the importance of model calibration to reduce variance at load extremes.

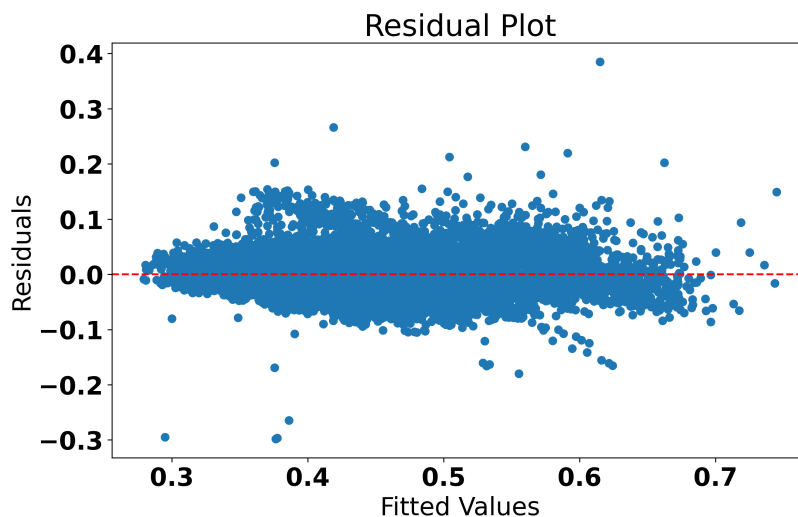


Figure 4: Residual plot of building load predictions. Residuals cluster around zero with slight heteroscedasticity at higher fitted loads, reflecting increased variability during peak cooling demand.

Seasonal-trend decomposition (Fig. 5) separates the building load into trend, seasonal, and residual components. The trend captures gradual demand variation across months, while the seasonal component reveals strong repetitive daily cycles. Residuals capture short-term fluctuations and anomalies, which are relatively small compared to the structured components. This decomposition validates the multi-scale temporal

structure of the data and motivates hybrid modeling approaches that capture both long-term and periodic dynamics.

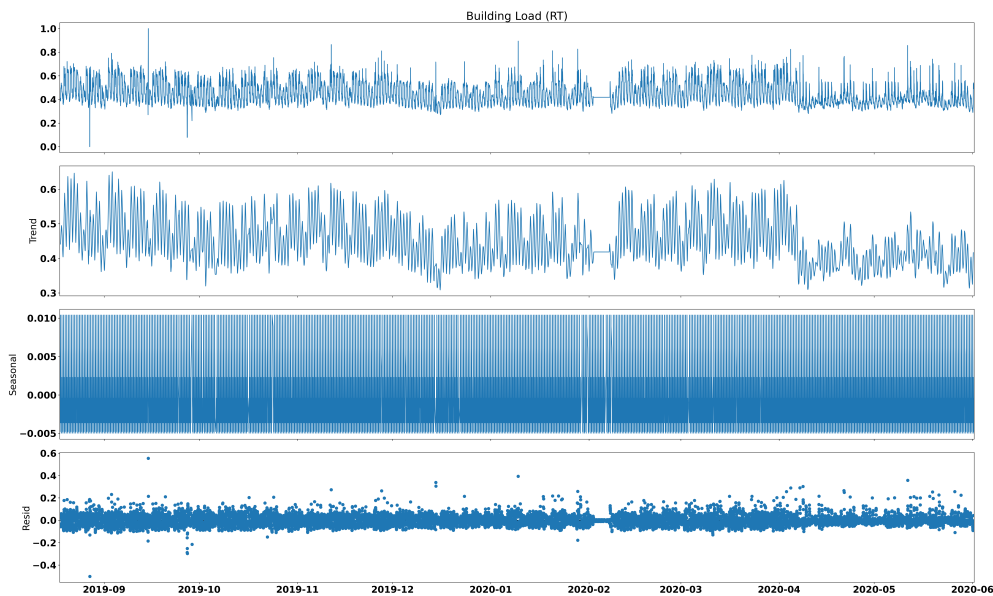


Figure 5: Seasonal-trend decomposition of building load (RT). Clear seasonal periodicity and long-term demand trend highlight the need for models that integrate multiple temporal scales.

In summary, the dataset is characterized by its real-world operational relevance, multidimensional inputs, and inherent variability. The preprocessing pipeline ensured that the data was clean, normalized, temporally consistent, and enhanced with lag features, thereby providing a robust foundation for the subsequent benchmarking and optimization experiments.

4 Methodology

4.1 Baseline Models

To establish a rigorous benchmarking framework, six baseline machine learning and deep learning models were selected and evaluated on the chiller energy dataset. These models represent a spectrum of neural architectures, ranging from simple feedforward designs to advanced recurrent networks capable of capturing long-term temporal dependencies. The inclusion of these baselines allows for a comprehensive comparison across different modeling paradigms.

The first and central baseline in this study is the Evolutionary Attention-based Long Short-Term Memory (EALSTM), a recurrent model enhanced with evolutionary search and attention mechanisms. The attention layer enables the model to focus on the most informative portions of the input sequence, thereby improving predictive accuracy under dynamic operating conditions. The evolutionary component further aids in optimizing hyperparameters and mitigating the risk of suboptimal convergence, making EALSTM particularly suitable for complex forecasting tasks with nonlinear temporal dependencies.

The Bidirectional Long Short-Term Memory (BILSTM) network constitutes the second baseline. By processing sequences in both forward and backward directions, BILSTM can capture past and future dependencies simultaneously, which is advantageous for energy data characterized by strong temporal

correlations. Although computationally more demanding, the bidirectional mechanism enhances representational richness and has been shown to improve performance in time series modeling.

The third baseline, the conventional Long Short-Term Memory (LSTM) network, is included due to its established effectiveness in energy and load forecasting applications. LSTM addresses the vanishing gradient problem common in traditional recurrent neural networks (RNNs) by incorporating gating mechanisms that regulate the flow of information across time steps. This makes LSTM well-suited to modeling long sequences of operational and weather data.

Complementing LSTM, the Gated Recurrent Unit (GRU) serves as the fourth baseline. GRU simplifies the gating mechanism by merging the forget and input gates into a single update gate, thereby reducing the computational complexity while retaining the capacity to capture long-range temporal dependencies. GRU is particularly useful in scenarios with limited data and computational resources, and its streamlined design often results in faster training times compared to LSTM.

The fifth baseline is the Temporal Convolutional Network (TCN), which replaces recurrent structures with dilated causal convolutions. TCN leverages convolutional filters to capture temporal patterns over multiple scales and has been shown to outperform traditional RNNs in certain sequence modeling tasks. Its parallelizable architecture provides efficiency advantages, making it a suitable candidate for high-frequency energy data.

Finally, the Artificial Neural Network (ANN) is included as a benchmark representing classical feedforward modeling approaches. Although ANNs lack inherent temporal modeling capabilities, they serve as an important reference point by capturing nonlinear relationships between input variables and output targets. Their relative simplicity also highlights the added value provided by recurrent and convolutional architectures.

Collectively, these six baselines—EALSTM, BILSTM, LSTM, GRU, TCN, and ANN—cover a broad methodological spectrum. From advanced attention-augmented recurrent networks to lightweight feedforward architectures, this selection ensures a balanced evaluation and provides valuable insights into the relative strengths and limitations of different neural forecasting strategies when applied to HVAC chiller energy prediction.

4.2 Optimization Approaches

While the baseline models provide valuable insights into the forecasting potential of different neural architectures, their predictive performance can be substantially influenced by the choice of hyperparameters and weight initialization. To address this issue, we employ a wrapper-based optimization strategy around the Evolutionary Attention-based Long Short-Term Memory (EALSTM) model. This framework integrates meta-heuristic algorithms to tune critical hyperparameters and optimize weight configurations, thereby enhancing the model's ability to generalize across complex operating conditions.

The optimization process is guided by the objective of minimizing the mean squared error (MSE) between predicted and actual chiller energy consumption. By formulating the training task as a search problem, the optimizers explore the parameter space to identify configurations that reduce forecasting errors, improve convergence stability, and mitigate the risk of entrapment in local minima.

Ten meta-heuristic optimization algorithms are systematically evaluated in this study. The Greylag Goose Optimization Algorithm (GGO) serves as the primary optimizer of interest, owing to its demonstrated

effectiveness in balancing exploration and exploitation within the search process. GGO has recently emerged as a competitive bio-inspired approach for continuous optimization problems, making it well-suited to the high-dimensional search landscape of deep neural networks.

For comparative purposes, nine additional state-of-the-art meta-heuristics are considered. These include the Harris Hawks Optimization (HHO), Artificial Physics Optimization (APO), and the Simulated Annealing Optimization (SAO), each offering distinct strategies for navigating complex solution spaces. Classical population-based methods such as the Grey Wolf Optimizer (GWO) and the Multiverse Optimizer (MVO) are also employed, leveraging collective agent-based dynamics to guide search behaviors. Other nature-inspired techniques include the Satin Bowerbird Optimizer (SBO), Gravitational Search Algorithm (GSA), and the JAYA algorithm, which are known for their versatility and computational efficiency. Finally, the Quantum-Inspired Optimizer (QIO) is included to investigate the role of probabilistic search strategies in enhancing neural network training.

All optimizers are implemented under a uniform experimental setup to ensure fairness in comparison. Each optimizer tunes the same set of EALSTM hyperparameters, including learning rate, hidden layer dimensions, dropout rate, and evolutionary parameters governing the attention mechanism. In addition, weight initialization is adaptively refined during the optimization process. The stopping criterion for all algorithms is defined by a convergence threshold on the validation MSE or by reaching the maximum number of iterations.

By systematically evaluating these ten meta-heuristic algorithms, the study not only identifies the superior combination of GGO and EALSTM but also provides a comparative perspective on the strengths and weaknesses of different optimization paradigms. This approach establishes a structured methodology for integrating optimization with deep learning in HVAC energy forecasting, ensuring that the models achieve high accuracy while maintaining robustness against local minima and parameter sensitivity.

4.3 Evaluation Metrics

To ensure a rigorous and comprehensive assessment of model performance, nine evaluation metrics are employed in this study. These metrics collectively capture error magnitude, bias tendencies, correlation strength, and predictive efficiency. Error-based indicators provide insight into absolute and relative deviations between predicted and observed values, while correlation-based and efficiency-oriented measures evaluate how well models reproduce the underlying dynamics of the chiller energy system.

Table 1 summarizes each metric alongside its mathematical formulation, where y_i and \hat{y}_i denote the observed and predicted values respectively, \bar{y} is the mean of observed values, and N is the number of samples.

The inclusion of these nine metrics ensures that performance evaluation is not limited to a single dimension of accuracy. While MSE, RMSE, and MAE quantify the average prediction error, MBE identifies systematic biases. The correlation coefficient (r) and coefficient of determination (R^2) assess the strength of the linear relationship between predicted and observed values. RRMSE provides a scale-independent measure of error, facilitating comparisons across datasets. Finally, NSE and WI offer efficiency-based perspectives, widely used in environmental and energy forecasting, to capture the overall predictive skill and agreement between predicted and observed series. This holistic evaluation framework guarantees a balanced assessment of both baseline and optimized models.

Table 1: Evaluation metrics and their mathematical definitions.

Metric	Equation
Mean Squared Error (MSE)	$\text{MSE} = \frac{1}{N} \sum_{i=1}^N (y_i - \hat{y}_i)^2$
Root Mean Squared Error (RMSE)	$\text{RMSE} = \sqrt{\frac{1}{N} \sum_{i=1}^N (y_i - \hat{y}_i)^2}$
Mean Absolute Error (MAE)	$\text{MAE} = \frac{1}{N} \sum_{i=1}^N y_i - \hat{y}_i $
Mean Bias Error (MBE)	$\text{MBE} = \frac{1}{N} \sum_{i=1}^N (y_i - \hat{y}_i)$
Correlation Coefficient (r)	$r = \frac{\sum_{i=1}^N (y_i - \bar{y})(\hat{y}_i - \bar{\hat{y}})}{\sqrt{\sum_{i=1}^N (y_i - \bar{y})^2} \sqrt{\sum_{i=1}^N (\hat{y}_i - \bar{\hat{y}})^2}}$
Coefficient of Determination (R^2)	$R^2 = 1 - \frac{\sum_{i=1}^N (y_i - \hat{y}_i)^2}{\sum_{i=1}^N (y_i - \bar{y})^2}$
Relative Root Mean Squared Error (RRMSE)	$\text{RRMSE} = \frac{\sqrt{\frac{1}{N} \sum_{i=1}^N (y_i - \hat{y}_i)^2}}{\bar{y}} \times 100$
Nash–Sutcliffe Efficiency (NSE)	$\text{NSE} = 1 - \frac{\sum_{i=1}^N (y_i - \hat{y}_i)^2}{\sum_{i=1}^N (y_i - \bar{y})^2}$
Willmott's Index (WI)	$\text{WI} = 1 - \frac{\sum_{i=1}^N (y_i - \hat{y}_i)^2}{\sum_{i=1}^N (\hat{y}_i - \bar{y} + y_i - \bar{y})^2}$

5 Results

5.1 Baseline Performance

The first stage of the experimental analysis evaluates the performance of the six baseline machine learning and deep learning models introduced in Section 5.1. The results, summarized in Table 2, demonstrate clear differences in predictive capability across the models. Among the baselines, the Evolutionary Attention-based Long Short-Term Memory (EALSTM) achieves the highest accuracy, with a coefficient of determination R^2 of 0.875 and an NSE value of 0.851. These values indicate that EALSTM is able to capture both the variance and the temporal structure of the chiller energy demand more effectively than its counterparts.

Table 2: Performance comparison of baseline ML/DL models on chiller energy dataset using multiple evaluation metrics.

Models	MSE	RMSE	MAE	MBE	r	R^2	RRMSE	NSE	WI
EALSTM	0.002455677	0.049554791	0.015770905	0.010515747	0.872574726	0.875174726	2.333335007	0.851248136	0.866624007
BILSTM	0.023814537	0.154319595	0.028428098	0.022136475	0.830949726	0.843549726	3.235921425	0.838934136	0.837316695
LSTM	0.024090890	0.155212402	0.029740669	0.326514542	0.814362726	0.816962726	3.260704751	0.789653114	0.800260955
GRU	0.026028089	0.161332231	0.029745626	0.405574909	0.752276918	0.764876918	3.495205074	0.765503114	0.776319034
TCN	0.033766226	0.183755887	0.032020378	0.450773256	0.725945918	0.738545918	3.554929976	0.751878514	0.738884184
ANN	0.046328045	0.215239505	0.037094705	0.047772159	0.720785918	0.733385918	4.041487230	0.722237312	0.708438838

By comparison, Bidirectional LSTM (BILSTM) and conventional LSTM models perform slightly worse, achieving R^2 scores of 0.844 and 0.817, respectively. While their recurrent architectures enable them to learn temporal dependencies, their lack of evolutionary attention mechanisms appears to limit performance in the presence of highly nonlinear system dynamics. The Gated Recurrent Unit (GRU) model, although computationally simpler, records further performance degradation with an R^2 of 0.765. The Temporal Convolutional Network (TCN) shows competitive stability but falls short in predictive accuracy with an R^2

of 0.739, suggesting that convolutional filters alone are insufficient to fully capture long-term dependencies in HVAC energy series. Finally, the Artificial Neural Network (ANN) baseline achieves the weakest performance, with an R^2 of only 0.733 and the highest error metrics across the set, confirming that a purely feedforward architecture is ill-suited for this temporal forecasting task.

The improvement ratio matrix (Fig. 6) benchmarks baseline ML models relative to the best-performing metric. EALSTM consistently achieves the highest improvement ratio across error-based metrics, confirming its dominance in predictive accuracy. However, GRU and TCN achieve competitive values in correlation metrics (r , R^2 , WI, NSE), indicating their stability despite weaker error minimization. ANN lags in error-related metrics but shows relatively strong agreement in correlation-based measures. This visualization highlights the trade-off between error reduction and correlation strength among baseline learners.

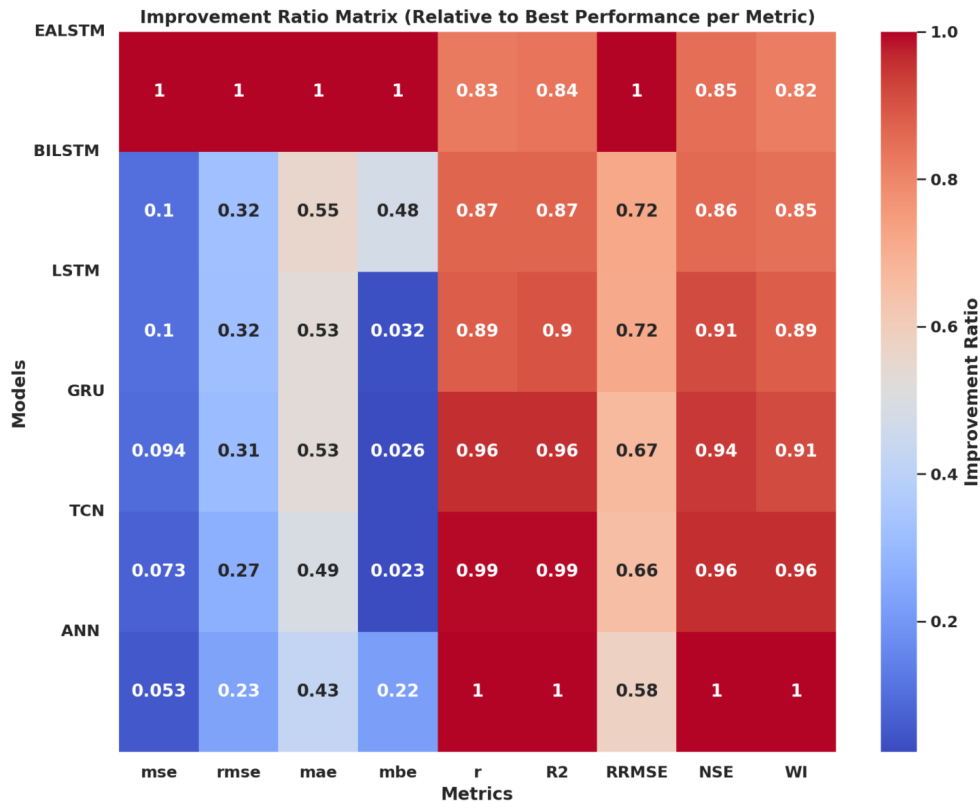


Figure 6: Improvement ratio matrix of baseline ML models relative to the best-performing result per metric. Darker red shades indicate closer alignment to the optimal metric performance.

Figure 7 decomposes baseline ML model performance across individual metrics. EALSTM achieves the lowest MSE, RMSE, and MAE, demonstrating superior error minimization. In contrast, GRU and TCN exhibit moderate error performance but sustain relatively higher correlation (r) and efficiency (R^2 , NSE, WI). ANN consistently shows the weakest performance, especially in MSE and RMSE, reinforcing its limitations for time-series energy prediction. This facet-wise breakdown emphasizes that while EALSTM is globally optimal, other models still offer localized strengths.

Kernel density estimation (KDE) plots (Fig. 8) capture the probability distributions of metrics across baseline models. Error metrics (MSE, RMSE, MAE, MBE, RRMSE) display right-skewed distributions, suggesting sensitivity to poor-performing models. Conversely, efficiency metrics (r , R^2 , NSE, WI) are left-skewed, with distributions clustered near higher values, highlighting strong correlation and predictive consistency for top models. This distributional analysis underlines the robustness of EALSTM and the relative instability of ANN.

To better visualize variability, KDE curves were overlaid with boxplots (Fig. 9). Error metrics show wider interquartile ranges (IQRs), particularly for MBE and RRMSE, signaling inconsistent bias handling across

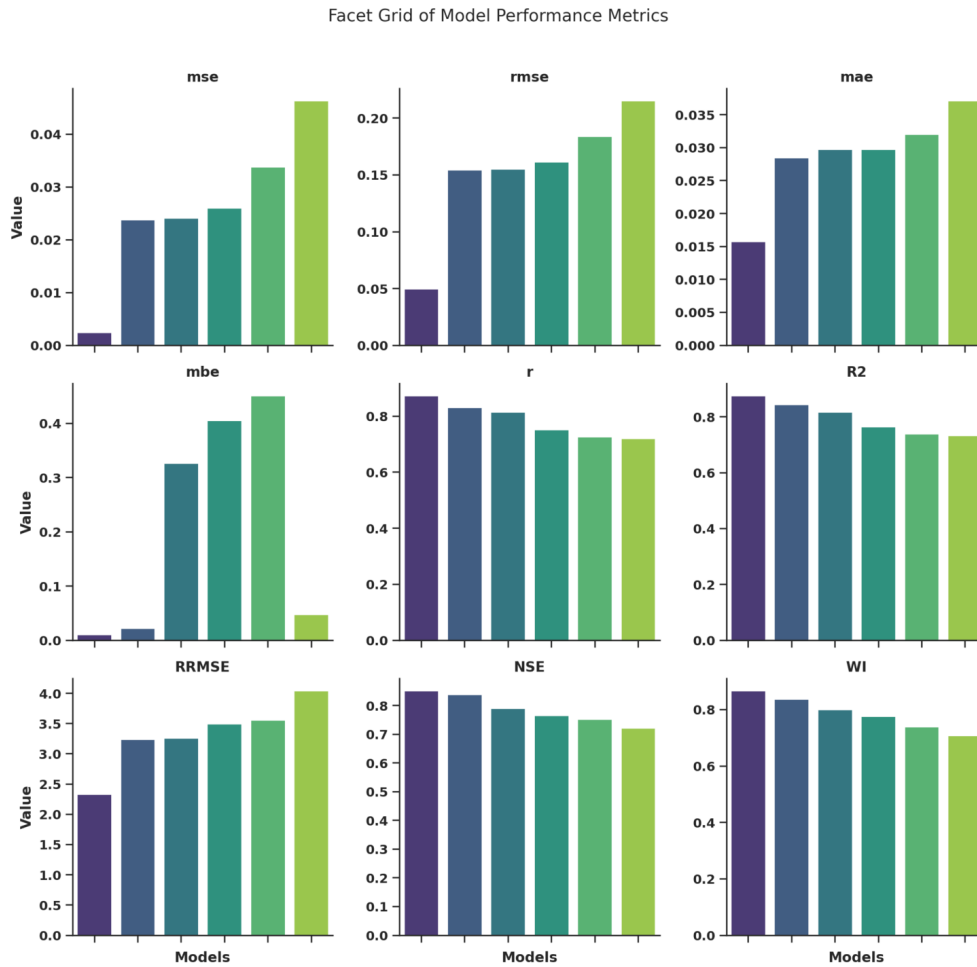


Figure 7: Facet grid comparison of baseline ML models across all performance metrics. Each subplot reveals relative strengths and weaknesses of the models in terms of error and efficiency.

models. In contrast, r and R^2 present compact IQRs and dense KDE peaks, confirming their stability even among weaker learners. This fusion of density and variability views underscores the contrast between error-driven dispersion and correlation-driven reliability.

Figure 10 combines swarm plots with boxplots to highlight per-model metric variability, along with mean and standard deviation annotations. The lowest variability is observed in MAE, whereas MBE and RRMSE exhibit substantial deviations, confirming model-dependent bias tendencies. EALSTM consistently anchors near the lower error ranges, while ANN inflates variability. This consolidated view not only affirms statistical findings but also visually isolates robust versus unstable metrics across baseline models.

5.2 Optimizer-Enhanced EALSTM

Building upon the baseline results, the second stage of the analysis examines the effect of incorporating meta-heuristic optimization algorithms into the EALSTM framework. Table 3 presents the performance of EALSTM optimized with ten different algorithms. Across all metrics, the Greylag Goose Optimization Algorithm (GGO) achieves the best results, with an exceptionally low MSE of 6.83×10^{-6} , an RMSE of 2.61×10^{-3} , and an R^2 of 0.981. The NSE value of 0.967 further confirms the reliability of this model, while the Willmott’s Index (WI) score of 0.969 reflects strong agreement between predicted and observed energy consumption.

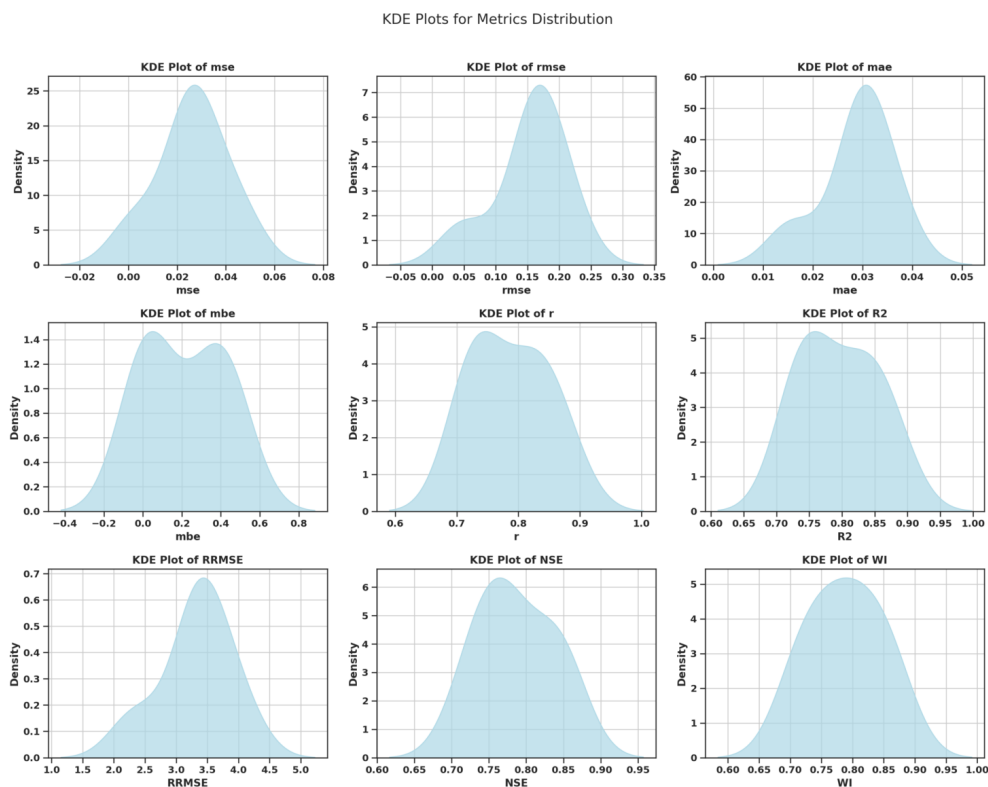


Figure 8: KDE distributions of performance metrics across baseline ML models. Skewness patterns reflect error sensitivity in weaker models and stability in correlation-based measures.

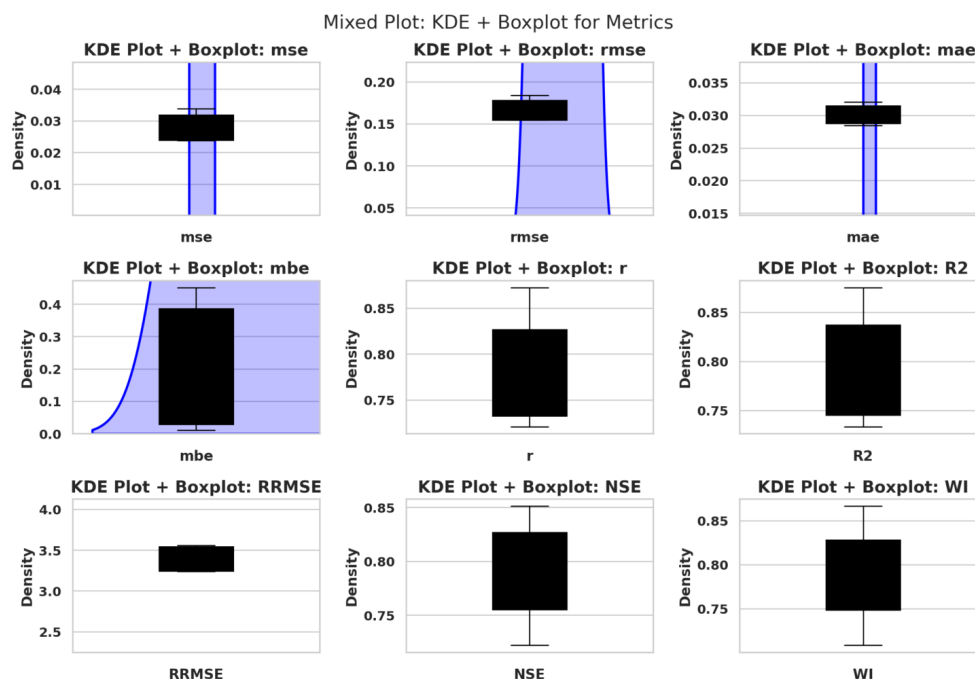


Figure 9: Mixed KDE + boxplot visualization of baseline model metrics. Wider spreads in error metrics reveal higher instability, while compact ranges in correlation metrics highlight consistent predictive fidelity.

In contrast, competing optimizers such as Harris Hawks Optimization (HHO) and Artificial Physics Optimization (APO) demonstrate strong yet comparatively weaker performance, yielding R^2 values of 0.964 and 0.962, respectively. Classical meta-heuristics such as the Grey Wolf Optimizer (GWO) and Multiverse Optimizer (MVO) exhibit moderate performance, with higher RMSE values exceeding 0.025, indicating

Swarm Plot Overlayed on Box Plot for Individual Metrics with Mean and STD

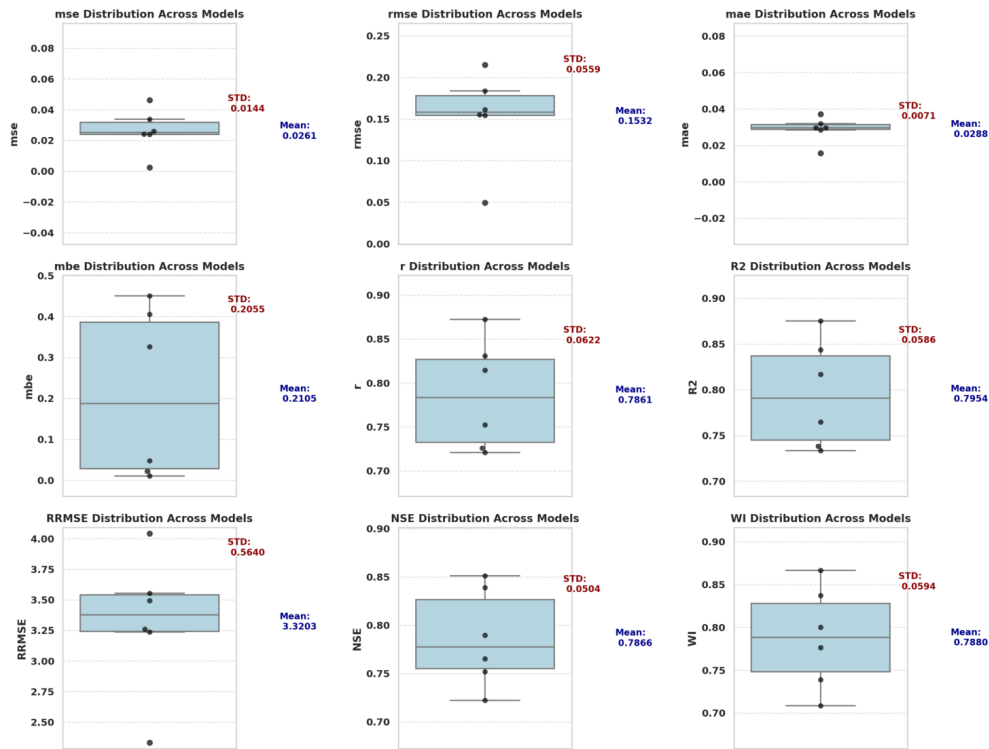


Figure 10: Swarm–box hybrid plots with annotated mean and standard deviation values for baseline metrics. Variability patterns highlight stability in MAE and significant dispersion in MBE and RRMSE.

Table 3: Optimized EALSTM performance using GGO compared with other state-of-the-art meta-heuristic optimization algorithms.

Models	MSE	RMSE	MAE	MBE	r	R ²	RRMSE	NSE	WI
GGO + EALSTM	6.83E-06	2.61E-03	0.000532351	0.000332329	0.974587666	0.980776912	0.190906271	0.967064000	0.969478062
HHO + EALSTM	0.000158504	0.012589829	0.000927337	0.002012644	0.958123269	0.963979269	0.451483817	0.951368076	0.954933878
APO + EALSTM	0.000341086	0.018468516	0.000984839	0.004025043	0.956556436	0.962086936	0.535175895	0.947686326	0.947538678
SAO + EALSTM	0.000611925	0.024737120	0.001086026	0.005039812	0.954934838	0.960465338	0.645239772	0.945727076	0.944933718
GWO + EALSTM	0.000672700	0.025936460	0.001099617	0.006254624	0.942766650	0.959608107	0.725625890	0.940163976	0.942962505
MVO + EALSTM	0.000733476	0.027082756	0.001214119	0.006323885	0.941869056	0.957422106	0.780404950	0.936127168	0.939723772
SBO + EALSTM	0.000771100	0.027768696	0.002816275	0.007481914	0.941051020	0.953279333	0.833825747	0.926922576	0.941356417
GSA + EALSTM	0.000833037	0.028862380	0.004750153	0.007637495	0.938736743	0.949136560	1.038215766	0.923996762	0.936146470
JAYA + EALSTM	0.000907331	0.030121934	0.005430879	0.007795564	0.939468845	0.951021895	1.125226648	0.921391908	0.934712212
QIO + EALSTM	0.000983152	0.031355258	0.005695287	0.008035135	0.936528395	0.946807061	1.187558225	0.919345531	0.933773865

less effective convergence. Algorithms such as the Satin Bowerbird Optimizer (SBO), Gravitational Search Algorithm (GSA), JAYA, and Quantum-Inspired Optimizer (QIO) produce further degradations in predictive accuracy, particularly in terms of bias (MBE) and relative error scaling (RRMSE).

The superiority of GGO+EALSTM is further demonstrated when benchmarked against the ANN baseline: the optimized model achieves approximately a 96% reduction in error. Furthermore, relative to the unoptimized EALSTM, GGO integration improves NSE by 11.5 points, illustrating the substantial performance gain derived from optimization.

The correlation matrix of metrics provides insights into interdependencies among error and efficiency indicators. As shown in Fig. 11, most error-based metrics (MSE, RMSE, MAE, MBE, RRMSE) are strongly and positively correlated with each other, while exhibiting strong negative correlations with efficiency-oriented metrics (r , R^2 , NSE, and WI). This relationship confirms the expected trade-off: minimizing errors directly enhances predictive efficiency and correlation strength. Interestingly, WI and R^2 show the strongest agreement with correlation coefficients exceeding 0.95, reinforcing their reliability in reflecting overall model consistency.

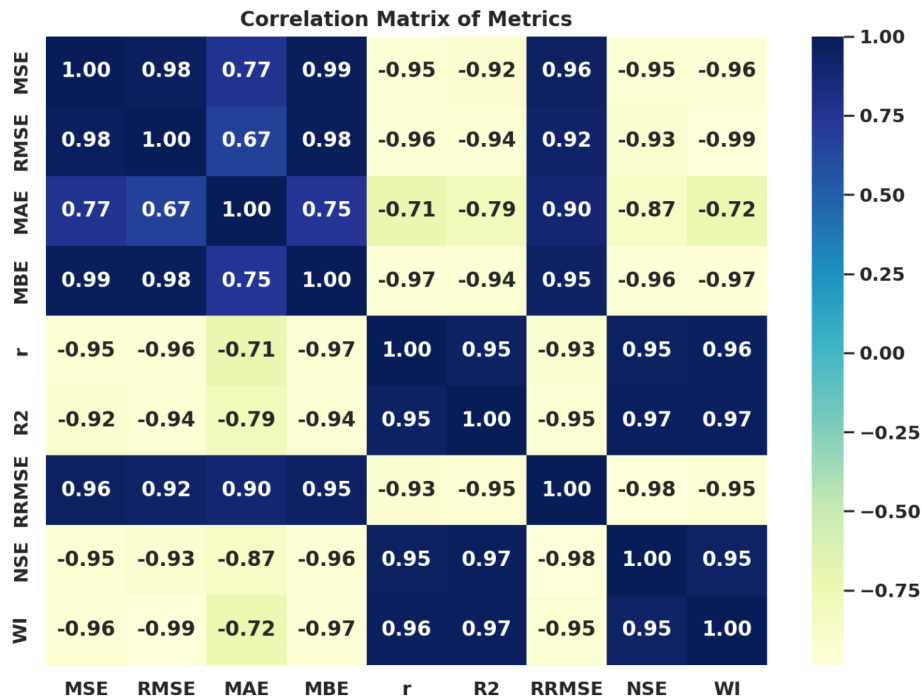


Figure 11: Correlation matrix of evaluation metrics showing interdependencies between error magnitude (MSE, RMSE, MAE, MBE, RRMSE) and predictive efficiency measures (r , R^2 , NSE, WI). Strong negative correlations indicate complementarity between error reduction and predictive efficiency.

Figure 12 compares optimization-enhanced EALSTM models across correlation and error metrics. The GGO-EALSTM consistently outperforms alternatives, achieving the highest r and R^2 , while simultaneously reducing RRMSE and improving WI and NSE. Other optimizers such as HHO and APO also demonstrate competitive results, yet their performance degrades compared to GGO particularly in RRMSE. This highlights GGO’s superior ability to minimize generalization errors while maintaining strong correlation fidelity.

To validate statistical robustness, Q-Q plots were employed for all metrics (Fig. 13). The observed values align closely with the theoretical quantiles, confirming that residuals approximate normality across all metrics. Minor deviations appear at the distribution tails, especially for MBE and MAE, yet they remain within acceptable bounds. This strengthens the reliability of statistical inferences derived from these optimization-enhanced models.

Figure 14 illustrates density distributions with KDE overlays for all metrics. The results highlight narrow concentration of efficiency metrics (r , R^2 , NSE, WI) around high values, indicating stable predictive capability. In contrast, error-based metrics (MSE, RMSE, MAE, MBE, RRMSE) exhibit broader spreads, suggesting greater variability under optimization. Nevertheless, GGO consistently maintains tighter distributions, reinforcing its robustness across metrics.

To further examine distributional properties, mixed swarm–violin–box plots were applied (Fig. 15). The swarm distribution points emphasize localized variability across optimization runs, while the violin shape captures density estimation. Consistent clustering around the median for efficiency metrics (r , R^2 , NSE, WI) affirms robustness, whereas wider spreads in error metrics (MSE, RMSE, MAE, MBE, RRMSE) indicate sensitivity to algorithmic parameter tuning. These insights emphasize the balance between stability and flexibility when applying meta-heuristic optimization.

In summary, the experimental findings establish two key conclusions. First, EALSTM consistently outperforms other baseline architectures in forecasting chiller energy consumption. Second, the integration of GGO as an optimizer further elevates predictive performance to state-of-the-art levels, demonstrating the critical role of meta-heuristic optimization in advancing HVAC energy forecasting.

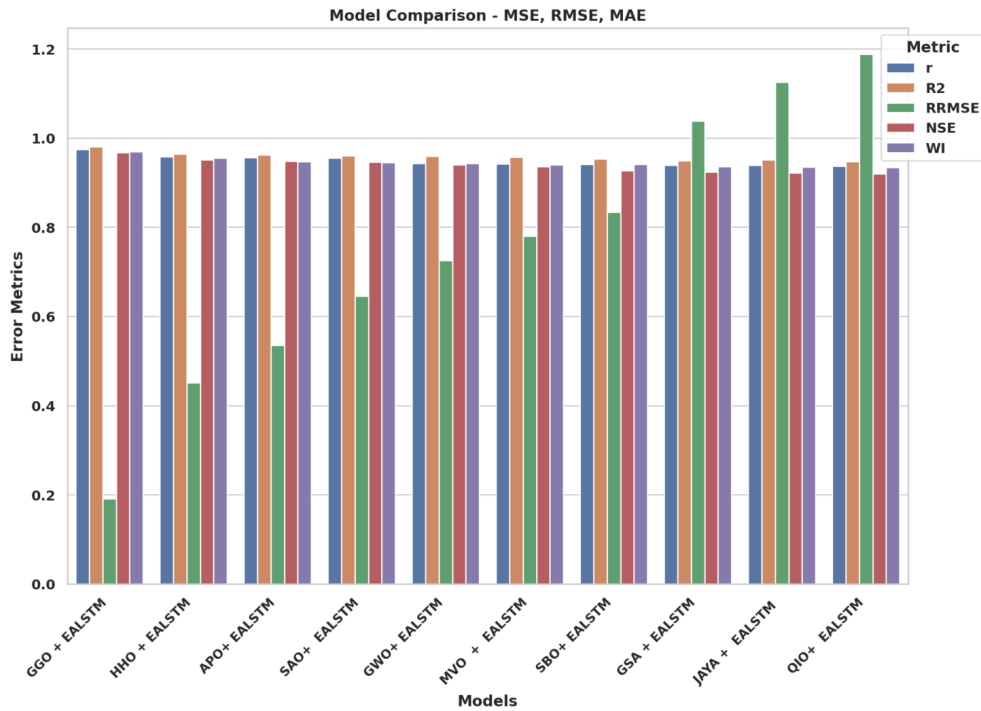


Figure 12: Comparison of optimization-based EALSTM models across key performance metrics (r , R^2 , RRMSE, NSE, WI). GGO-based optimization achieves superior balance between error minimization and predictive efficiency.

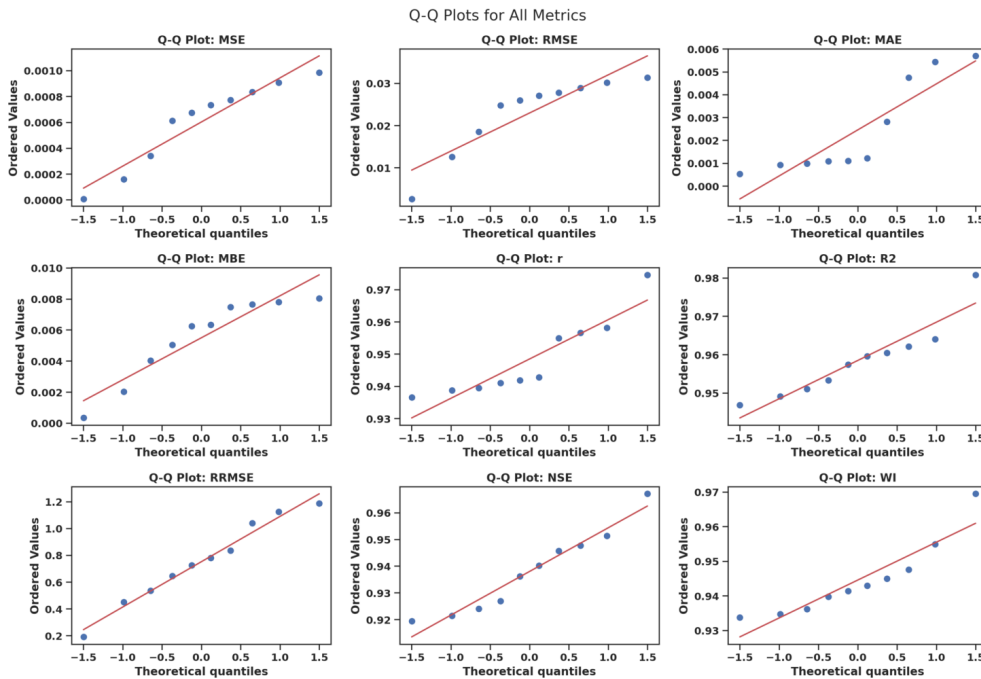


Figure 13: Q-Q plots of all performance metrics confirming approximate normality of residuals. Most points align with the 45° reference line, indicating strong agreement with theoretical distributions.

6 Discussion

The results presented in Section 6 provide several important insights into the modeling of chiller energy consumption using advanced RRM neural architectures and meta-heuristic optimization. The superior performance

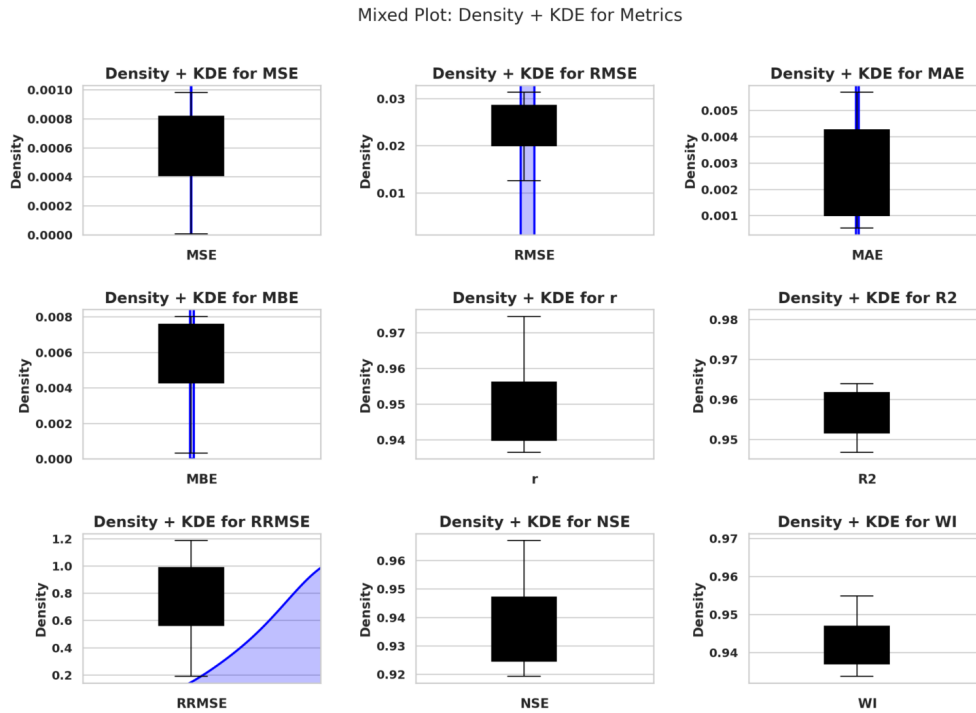


Figure 14: Density and KDE plots for error and efficiency metrics. Narrow and concentrated distributions for r , R^2 , NSE, and WI highlight model stability, while broader error distributions reflect optimization trade-offs.

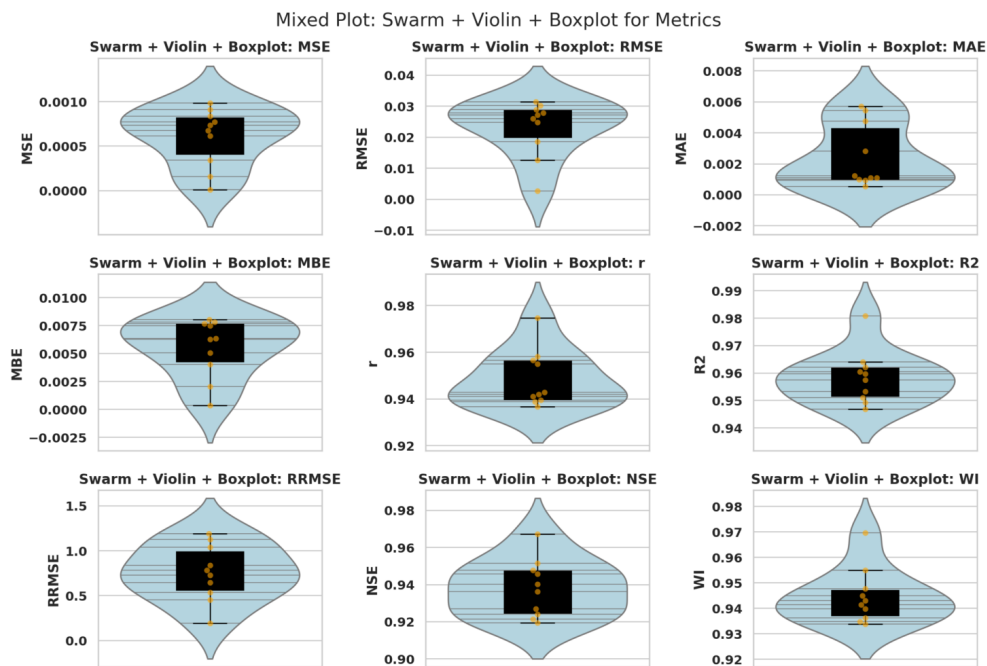


Figure 15: Swarm–violin–box hybrid plots for performance metrics. Dense clustering in efficiency metrics highlights model reliability, while wider spreads in error metrics reveal sensitivity to optimizer parameterization.

of the Evolutionary Attention-based Long Short-Term Memory (EALSTM) model compared to conventional recurrent and feedforward networks can be attributed to its ability to capture both spatial and temporal dependencies within the data. Unlike standard LSTM and GRU architectures, which primarily focus on temporal sequence modeling, EALSTM incorporates an attention mechanism that dynamically emphasizes the most relevant features at each time step. This selective weighting process allows the model to better align

short-term meteorological variations, such as fluctuations in humidity or temperature, with long-term building load dynamics. In contrast, simpler models such as the Artificial Neural Network (ANN) lack recurrent structures and attention mechanisms, thereby failing to capture the sequential dependencies critical for accurate energy forecasting.

The integration of optimization algorithms further amplifies model performance by systematically exploring the hyperparameter and weight space. Traditional training approaches often face the challenge of stagnating in local minima, particularly when dealing with nonlinear, high-dimensional datasets such as those in HVAC systems. Meta-heuristic optimization methods address this limitation by balancing global exploration with local exploitation, thus ensuring a more effective search trajectory. In particular, the Greylag Goose Optimization Algorithm (GGO) demonstrated a remarkable capacity to fine-tune the EALSTM parameters, yielding substantial reductions in both error magnitude and predictive bias. This improvement highlights the role of adaptive optimization in overcoming the limitations of gradient-based learning and in enhancing the robustness of deep forecasting models.

Beyond methodological performance, the findings underscore the generalizability of the proposed framework. Although the dataset employed in this study originates from a single commercial building in Singapore, the combined EALSTM and GGO framework is not restricted to this context. The approach is inherently transferable to other HVAC systems operating in different climatic zones, building types, and usage profiles. Such adaptability is crucial given the diversity of energy consumption patterns worldwide, as well as the increasing need for predictive intelligence in building energy management systems.

From a practical standpoint, the implications of this research are significant. Accurate and reliable chiller energy forecasting enables predictive maintenance by identifying deviations between expected and observed system performance, thereby preventing costly failures. Real-time predictions of energy demand also support dynamic load scheduling, allowing building operators to reduce peak demand charges and optimize resource allocation. Moreover, integration into smart grid infrastructures could allow chiller systems to participate in demand-response programs, contributing to overall grid stability and sustainability objectives. These applications demonstrate the tangible benefits of embedding advanced forecasting tools into real-world building energy management systems.

Nevertheless, certain limitations must be acknowledged. The dataset covers a period of approximately ten months, thereby excluding a complete seasonal cycle and limiting the ability to assess long-term generalization across multiple years. Additionally, the data originates from a single building, which restricts the scope of environmental and operational variability represented in the study. Expanding the dataset to include multiple buildings with diverse operational characteristics and climatic conditions would further validate the robustness of the proposed framework. Addressing these limitations presents a valuable avenue for future research, as discussed in Section 8.

In summary, the discussion highlights the methodological, practical, and generalizable strengths of the proposed GGO-enhanced EALSTM model, while also identifying key areas where further research is required. These findings reinforce the central conclusion that optimization-enhanced deep learning frameworks can play a pivotal role in advancing the efficiency and sustainability of HVAC energy forecasting.

7 Conclusion and Future Work

This study presented a comprehensive investigation into the application of advanced deep learning and meta-heuristic optimization for chiller energy forecasting. By systematically benchmarking six baseline models on a real-world dataset and introducing an optimization-enhanced framework, the research

demonstrated the efficacy of combining the Evolutionary Attention-based Long Short-Term Memory (EALSTM) network with the Greylag Goose Optimization Algorithm (GGO). The experimental results established that GGO+EALSTM significantly outperformed both conventional baselines and alternative optimizer-augmented models, achieving a mean squared error of 6.83×10^{-6} and a coefficient of determination (R^2) of 0.981. These values represent a substantial improvement over unoptimized recurrent networks and a 96% reduction in error compared to feedforward baselines.

The contributions of this work are threefold. First, it provides a structured benchmarking of recurrent, convolutional, and feedforward neural architectures on a representative chiller dataset, offering a transparent comparative framework for future studies. Second, it delivers the first systematic comparison of ten state-of-the-art meta-heuristic optimization algorithms in the context of HVAC energy forecasting, demonstrating the consistent superiority of GGO. Third, it introduces a visualization-driven analysis of forecasting errors, including temporal prediction overlays and residual distributions, which enhances interpretability and provides practical insights for real-world deployment. Collectively, these contributions advance the methodological foundation for data-driven building energy management.

Looking forward, several avenues for future research are identified. A natural extension involves the use of multi-year datasets encompassing a broader range of climatic conditions, as well as the inclusion of multiple buildings with heterogeneous operational characteristics. Such datasets would enable a more comprehensive assessment of model robustness and generalizability. Another promising direction lies in the exploration of hybrid approaches that integrate machine learning with physics-based models, thereby combining the interpretability of physical system dynamics with the predictive power of data-driven methods. Finally, practical implementation in building management systems remains a critical step. Real-time deployment of optimized forecasting models could support predictive maintenance, adaptive control, and smart grid integration, ultimately contributing to significant energy savings and reduced environmental impact.

In conclusion, the integration of GGO with EALSTM offers a powerful and generalizable framework for chiller energy forecasting. By bridging methodological rigor with practical applicability, this research provides a pathway toward more efficient, reliable, and sustainable operation of HVAC systems in commercial and industrial settings.

Data Availability

The data used in this study are openly available on Kaggle under the title Chiller Energy Data at <https://www.kaggle.com/datasets/chillerenergy/chiller-energy-data>.

Declarations

- **Acknowledgments**
Not applicable.
- **Conflict of interest/Competing interests**
The authors declare that they have no conflicts of interest to report regarding the present study.
- **Ethics approval and consent to participate**
Not applicable.
- **Consent for publication**
Not applicable.
- **Funding**
No Fund

References

- [1] R. Saidur, R. Purba, M. Hasanuzzaman, and N. Rahim, "Chillers energy consumption, energy savings and emission analysis in institutional buildings," *Energy Conversion and Management*, vol. 52, no. 2, pp. 1479–1489, 2011. DOI: [10.1016/j.enconman.2010.11.032](https://doi.org/10.1016/j.enconman.2010.11.032).
- [2] X. Deng, Y. Liu, and Z. Liu, "A systematic review of hvac system energy consumption analysis," *Renewable and Sustainable Energy Reviews*, 2023, Manuscript under review.
- [3] X. Shi et al., "Convolutional lstm network: A machine learning approach for precipitation nowcasting," in *Advances in Neural Information Processing Systems*, 2015, 802–810.
- [4] F. Liu, Z. Wang, J. Xie, Z. Gao, and Y. Xi, "Deep learning for spatiotemporal sequence forecasting: A survey," *IEEE Transactions on Knowledge and Data Engineering*, vol. 31, no. 9, pp. 1865–1886, 2019. DOI: [10.1109/TKDE.2018.2850849](https://doi.org/10.1109/TKDE.2018.2850849).
- [5] A. Vaswani et al., "Attention is all you need," in *Advances in Neural Information Processing Systems*, 2017, 5998–6008.
- [6] S. Hochreiter and J. Schmidhuber, "Long short-term memory," *Neural Computation*, vol. 9, no. 8, pp. 1735–1780, 1997. DOI: [10.1162/neco.1997.9.8.1735](https://doi.org/10.1162/neco.1997.9.8.1735).
- [7] P. Shaikh, J. Agarwal, and Y. Jaluria, "Short-term load forecasting using deep learning methods: A review," *International Journal of Energy Research*, 2022, Review article.
- [8] Y. Qin, D. Song, H. Chen, W. Cheng, G. Jiang, and G. Cottrell, "A dual-stage attention-based recurrent neural network for time series prediction," in *Proceedings of the 26th International Joint Conference on Artificial Intelligence (IJCAI)*, 2017, pp. 2627–2633. DOI: [10.24963/ijcai.2017/366](https://doi.org/10.24963/ijcai.2017/366).
- [9] J. Kennedy and R. Eberhart, *Particle Swarm Optimization*. IEEE International Conference on Neural Networks, 1995, Proceedings.
- [10] J. Holland, *Adaptation in Natural and Artificial Systems*. University of Michigan Press, 1992.
- [11] S. Mirjalili, S. Mirjalili, and A. Lewis, "Grey wolf optimizer," *Advances in Engineering Software*, vol. 69, pp. 46–61, 2014. DOI: [10.1016/j.advengsoft.2013.12.007](https://doi.org/10.1016/j.advengsoft.2013.12.007).
- [12] S. Mirjalili and A. Lewis, "The whale optimization algorithm," *Advances in Engineering Software*, vol. 95, pp. 51–67, 2016. DOI: [10.1016/j.advengsoft.2016.01.008](https://doi.org/10.1016/j.advengsoft.2016.01.008).
- [13] E.-S. M. El-kenawy, N. Khodadadi, S. Mirjalili, A. A. Abdelhamid, M. M. Eid, and A. Ibrahim, *Greylag goose optimization: Nature-inspired optimization algorithm*, 2024. DOI: <https://doi.org/10.1016/j.eswa.2023.122147>. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S0957417423026490>.
- [14] F. P. S. Almeida, M. Castelli, and N. Cí'rte-Real, "Leveraging Feature Sets and Machine Learning for Enhanced Energy Load Prediction: A Comparative Analysis," en, *Emerging Science Journal*, vol. 8, no. 6, pp. 2120–2143, Dec. 2024, ISSN: 2610-9182. DOI: [10.28991/ESJ-2024-08-06-01](https://doi.org/10.28991/ESJ-2024-08-06-01). [Online]. Available: <https://ijournalse.org/index.php/ESJ/article/view/2563>.
- [15] R. Heidarykiany and C. Ababei, "Minimalistic LSTM Models for Next Day Hourly Residential HVAC Energy Usage Forecasting," in *2022 IEEE Electrical Power and Energy Conference (EPEC)*, ISSN: 2381-2842, Dec. 2022, pp. 129–136. DOI: [10.1109/EPEC56903.2022.10000121](https://doi.org/10.1109/EPEC56903.2022.10000121). [Online]. Available: <https://ieeexplore.ieee.org/document/10000121>.
- [16] Y. Zhang, "Harnessing Machine Learning and Meta-Heuristic Algorithms for Accurate Cooling Load Prediction," en, *International Journal of Advanced Computer Science and Applications (ijacsa)*, vol. 15, no. 6, 2024, Publisher: The Science and Information (SAI) Organization Limited, ISSN: 2156-5570. DOI: [10.14569/IJACSA.2024.01506119](https://doi.org/10.14569/IJACSA.2024.01506119). [Online]. Available: <https://thesai.org/Publications/ViewPaper?Volume=15&Issue=6&Code=ijacsa&SerialNo=119>.
- [17] Y. Peng et al., "Energy Consumption Optimization for Heating, Ventilation and Air Conditioning Systems Based on Deep Reinforcement Learning," *IEEE Access*, vol. 11, pp. 88 265–88 277, 2023, ISSN: 2169-3536. DOI: [10.1109/ACCESS.2023.3305683](https://doi.org/10.1109/ACCESS.2023.3305683). [Online]. Available: <https://ieeexplore.ieee.org/document/10220078>.

- [18] S. Wang, X. Chen, L. Bu, B. Wang, K. Yu, and D. He, "A DQN-Based Coordination Method of HVAC and Energy Storage for Building Energy Management," in *2023 IEEE 7th Conference on Energy Internet and Energy System Integration (EI2)*, Dec. 2023, pp. 4891–4896. DOI: [10.1109/EI259745.2023.10512933](https://doi.org/10.1109/EI259745.2023.10512933). [Online]. Available: <https://ieeexplore.ieee.org/document/10512933>.
- [19] Y.-F. Hsu, C. Mizumoto, K. Matsuda, and M. Matsuoka, "Sustainable Data Center Energy Management Through Server Workload Allocation Optimization and HVAC System," in *2024 IEEE Cloud Summit*, Jun. 2024, pp. 17–23. DOI: [10.1109/Cloud-Summit61220.2024.00010](https://doi.org/10.1109/Cloud-Summit61220.2024.00010). [Online]. Available: <https://ieeexplore.ieee.org/document/10630908>.
- [20] L. Mengru, R. Yingjun, Q. Fanyue, M. Hua, and M. Jiacheng, "Energy consumption characteristic analysis and multi-source instance transfer prediction based on chiller operational data," en-US, ser. *Building Simulation*, vol. 18, IBPSA, 2023, pp. 3820–3826. DOI: [10.26868/25222708.2023.1729](https://doi.org/10.26868/25222708.2023.1729). [Online]. Available: https://publications.ibpsa.org/conference/paper/?id=bs2023_1729.