# Adaptive Load Forecasting in Smart Grids Using Deep Learning and Edge Computing

[1]Sohail Khan, [2]Muhammad Abdullah Bin Arif , [3]Dr. Saad Khan Baloch , [4]Dr. Khakoo Mal, [5]Faisal Khan, [6]Shahzaib Ali

## Article Details

## ABSTRACT

**Sohail Khan**
M.Sc Electrical Engineering, Department: Electrical Engineering, University: Sir Syed CASE Institute of Technology, Islamabad
skhanpdp@gmail.com

**Muhammad Abdullah Bin Arif**
Electrical Engineering Department, University of Gujrat
23016122-001@uog.edu.pk

**Dr. Saad Khan Baloch**
Assistant Professor, Electrical Engineering Department, Isra University Hyderabad, Sindh, Pakistan. saad.baloch@isra.edu.pk

**Dr Khakoo Mal**
Assistant Professor, Department of Computer Science, Sukkur IBA University.
khakoo.mal@iba-suk.edu.pk

**Faisal Khan**
Department of Telecommunication Engineering, Dawood University of Engineering and Technology, Karachi
faisal.khan@duet.edu.pk

**Shahzaib Ali**
Department of Mechatronics Engineering, National University of Sciences and Technology, Islamabad, Pakistan
shahzaibali.isb@gmail.com

With the development of smart grid infrastructures to accommodate distributed energy resources, electric vehicles, and dynamic demand-side operations, responsive and precise load forecasting becomes more essential than ever. The problems associated with remote servers and constant flow of data lead to latency, scalability and privacy issues experienced by traditional centralized forecasting models. This paper presents an adaptive load forecasting system based on deep learning (long short-term memory (LSTM) networks) dusted to edge computing devices. The system uses real-time data, provided by smart meters or environmental sensors to provide localized, on-device predictions with minimal network dependency. Real-world data experiments reveal that the Edge-LSTM model has a better predictive performance, as well as much less inference latency, and greater adaptation ability than centralized LSTM, ARIMA, and feedforward neural networks. The findings also indicate a better utilization of network bottlenecked conditions and energy efficiency in resilience conditions which demonstrates the scalability of deep learning making use of edge intelligence to ensure decentralized operations in terms of smart grid. It is potentially an advance in the direction of real-time, privacy-preserving, and energy-conscious forecasts in future power networks.

## INTRODUCTION

Smart grids combine new communication, automation, and computing technologies with traditional power systems and are facilitated by the current global shift toward sustainable energy infrastructure (Gungor et al., 2020). The aim of these grids is to optimize power generation, transmission, and consumption through two-way data streams between utilities and consumers (Fang et al., 2019). The other key part of this optimization is proper load forecasting, as it enables the predictive understanding of future electricity demand patterns, which may be used to make operational decisions related to unit commitment, dispatch scheduling, and energy trading (Al-Musaylh et al., 2018).

Conventional load forecasting methods, such as linear regression, autoregressive analysis, and support vector machines, have proved to be of use in a structured context where consumption is steady (Hong & Fan, 2016). Nonetheless, the proliferation of distributed renewable energy resources, variable patterns of user activity, electric vehicles, and smart appliances has brought non-linearities and dynamism to consumption that may not be well recognized by the classical models (Taieb & Hyndman, 2014). Thus, artificial intelligence (AI) methods have increasingly gained popularity, especially deep learning models, as they have better capacity to build complex temporal and spatial relationships across data units (Zhang et al., 2021).

Time series forecasting and deep learning Deep learning architectures like the long short-term memory (LSTM) networks with gated recurrent units (GRUs) are accurate in predicting time series data especially short-term load forecasting (Marino et al., 2016). The models are beneficial at working with non-linear relationships and sequential dependencies and are proven to be greater than shallow learning models and classic statistical models (Ahmed et al., 2020). Nonetheless, their use in actual grid settings is typically limited by the imperative of having centralized high-efficiency computing facilities. Smart meters often require the transfer of energy data to distant servers where it is processed, presenting a delay, bandwidth burden, and the risk of data privacy (Wen et al., 2021).

The solution to these challenges is edge computing, which is based on the paradigm of decentralizing data processing by moving the computation closer to the data sources. With machine learning deployed to the edge on smart meters, home gateways, or local substations, the need for cloud storage and the response time to the problem decreases, as data can be treated on-site in real-time (Shi et al., 2016). The recent research examples have demonstrated how AI-

EDGE could be combined in different sectors, such as industrial automation, healthcare, and transport (Satyanarayanan, 2017; Chen et al., 2019). In smart grid applications, this integration could dramatically change load forecasting to an activity that is decentralized and proactive rather than centralized and reactive (Gharaibeh et al., 2020).

Still in its early stages, deep learning and edge computing in smart grids offer a promising way forward to achieve more adaptable and resilient energy systems. Adaptive load forecasting is a technique where forecasting models are continually learned and updated when the consumption behavior changes, weather patterns vary, and the conditions on the grid alter (Amirian et al., 2020). Adaptive forecasting attached to edge intelligence can enable decentralized grid nodes to make informed decisions locally, forming microgrids and prosumers with autonomy and ability to act without overloading central systems (Yang et al., 2022).

Nevertheless, even with the promise, deep learning on edge devices is not without its challenges. These can be categorized as minimal computational resources, energy limitations, model compression tradeoffs, and inferences in real-time (Zhou et al., 2019). The development of deep quantization models, federated learning and efficient network architectures however, has been alleviating these encumbrances in recent times and it is easier than ever before to execute deep models on platforms limited in resources (Li et al., 2021). In addition, policy and technology changes are driving utilities to consider decentralized intelligence in their digitalization plans (IRENA, 2022).

The proposed research aims to address this gap by proposing and testing a load forecasting system which integrates learning capabilities over time associated with LSTM models with low-latency capabilities of edge computing architectures. Precisely, the research studies the mechanism of applying successful adaptive forecasting to edge devices as a part of a smart grid system and evaluates the performance against conventional centralized forecasting models. In this way, this piece adds to the increasing knowledge base on decentralized intelligence in energy systems and acts as a stepping stone towards further smart grid innovations.

## LITERATURE REVIEW

### 1. EVOLUTION OF LOAD FORECASTING TECHNIQUES

Load forecasting is one of the pinnacle functions of electrical utilities that have been in operation over decades. Old models of forecasting, including linear regression, exponential smoothing, and autoregressive integrated moving average (ARIMA), have been popular in both short- and long-

term demands (Ghofrani et al., 2021). Nevertheless, the basic limitations of such models are their inapplicability to model nonlinear and time-varying phenomena typical of contemporary energy systems. As there is an escalating rate of penetration of non-central energy sources of renewables, electric-powered automobiles, and prosumer-based power consumption, energy demand patterns have recently been evolving as more stochastic and dynamic (Chandrashekar et al., 2023). The complexity has made the traditional forecasting methodologies less efficient in the real time or near real time in the grid scenario.

## 2. MACHINE LEARNING AND DEEP LEARNING IN LOAD FORECASTING

The development of machine learning (ML) and deep learning (DL) models has tremendously increased the ability to forecast by utilities. Random forests models as well as support vector regression and gradient boosting machines have been successfully applied in modeling short-term trends of energy (Giorgi et al., 2020). Nevertheless, the deep learning models, especially the recurrent neural networks (RNN), long short-term memory networks (LSTM), and convolutional neural networks (CNN) have been shown to perform exceptionally well in modelling the time-series data with long range dependencies (Iqbal et al., 2023). Incorporation of LSTM networks that have been developed to alleviate the problem of vanishing gradient in more detail is especially suitable to handle sequential data and thus have had the subsequent model architecture in the context of energy forecasting (Behera et al., 2022).

In addition, hybrid deep learning networks involving CNN as a spatial feature representation model and LSTM as a time structure model have further increased performance in terms of accuracy and generalization (Tang et al., 2021). Transformer models and attention-based models are also recently gaining popularity to make load forecasting since they enable representing global dependencies and increasing their interpretability (Du et al., 2024). But their use has so far been restricted on edge environments because of their complexity in computation.

## 3. REAL-TIME FORECASTING REQUIREMENTS AND LIMITATIONS OF CENTRALIZED SYSTEMS

The cloud-based deep learning systems have beneficial features in terms of capacity and training potential but inculcate considerable drawbacks when it comes to the latency, security, and bandwidth needs. The centralization of architectures implies a need to send the collected energy data via smart meters or IoT devices to distant servers, process it, and transfer it back to make decisions, which may be a drawback in case of a real-time environment (Yadav et al., 2022). Furthermore, the subject of sensitive energy data transmission raises concerns regarding data

privacy and security, especially in cases where the same is managed by third-party service providers (Kim & Park, 2020).

The inclusion of latency problems is also critical at peak times or whenever the microgrid environment is mission-critical and speedy response is required. Centralized systems are inflexible in scaling and adapting to localized or edge-specific abnormalities thus, inefficient. This caused researchers to consider decentralized solutions to avoid communication bottlenecks and offer a timely reaction (Das et al., 2023).

## 4. EMERGENCE OF EDGE COMPUTING IN SMART GRID SYSTEMS

Edge computing, where computation takes place at or close to the point of generation, is also becoming clear as a game-changer toward smart grid operation. Edge computing can also minimize latency by bringing data collection and computational procedures closer and consume less bandwidth, in addition to improving the reliability of the system (Hussain et al., 2023). Edge devices involved in smart grid applications, including advanced metering infrastructure, (AMI), local substation computers, and home energy gateways, may be endowed with lightweight computational modules to perform real-time analytics and forecasting models (Ranjan et al., 2021).

Edge computing is also advantageous to localized grid activities which make micro grids and distributed energy resources run independently when the grid has been disturbed or disrupted. Such decentralization is consistent with the changing form of the smart grid that is growing more spread out and dynamic (Miah et al., 2022). Further, the edge architectures support real time demand response programs, used to balance supply at rapidly varying environments.

## 5. DEEP LEARNING AT THE EDGE: OPPORTUNITIES AND CONSTRAINTS

Deploying deep learning models to edge devices has significant advantages in adaptive load forecasting, yet there are trade-offs affiliated with it. Edge devices typically have restricted processing halt, memory and battery life and therefore deep learning structures that require a large number of computations are troublesome to deploy in edge devices (Javed et al., 2024). Due to the recent advances in model optimization, e.g., by pruning, quantizing, and knowledge distillation, it is now feasible to run LSTM, CNN, or even Transformer-based models on edge devices with satisfactory performance indicators (Wang et al., 2022).

Frameworks such as TensorFlow Lite, PyTorch Mobile, and ONNX Runtime have been helping towards edge deployment of DL models. These tools enable the transformation of common

models to efficient models that are edge-compatible and can be used to perform inference on low-power systems such as the Raspberry Pi, Jetson Nano, or even smart meters (Fernandez et al., 2021). Direct training of such models on edge devices has thus far been infeasible, but federated learning methods are increasingly being studied as model-updating techniques that avoid the need to transmit raw data to main servers (Nasir et al., 2023).

## 6. ADAPTIVE FORECASTING IN DECENTRALIZED GRID ENVIRONMENTS

Adaptive forecasting is how a given system can adapt or even retrain on new incoming data. This ability is essential in the energy systems setting, as the temporal demand is designated with a seasonal effect, behaviour, and unexpected circumstances such as plot interruption or blackout (Molina, 2023). The inclusion of adaptivity in the forecasting system will make sure that the predictions will be correct under dynamic conditions of operation.

Many researchers have applied the concept of online learning and transfer learning to offer the capability of continuous adaptability of load prediction models (Patel & Kulkarni, 2022). Other methods that include semi-supervised or reinforcement learning that have been proposed can enable the forecast models to update dynamically without preparing the labeled dataset thoroughly (Li et al., 2024). These techniques can be effectively applied in edge environments because every node can use its distinct consumption and generation patterns.

## 7. GAPS IN LITERATURE AND RESEARCH DIRECTION

Despite the fact that the sphere of deep learning and edge computing has been developed rather well independently, the exploration of the combination of these two spheres in the context of the smart grid forecasting process is still not particularly well-grounded. The current literature proposes either incredibly precise deep learning models to be used in centralized architectures or simple rule-based analytics at the edge (Rehman et al., 2022). There exist no unified frameworks that intertwine edge deployment, deep learning flexibility, and actual instantiation in grids at large.

In addition, most studies overlook the temporal retraining and the heterogeneity of edge devices in the forecasting systems. Research tends to dismiss the possibility of edge devices with various hardware capabilities supporting various models, which results in an end-to-end implementation method, which is not necessarily optimal (Bisht et al., 2021). Benchmarking various architectures on different edge platforms in a practical scenario on the grid is an important requirement therefore experimental studies are urgently needed.

This paper attempts to fill such gaps by suggesting and analyzing a load prediction framework that adopts adaptive learning on deep learning models delivered across heterogeneous edge computing platforms. Integrating time-based learning via LSTM networks with edge sniffing deployment aims at providing a scalable, decentralized and precise solution to smart grid situations in the future.

## METHODOLOGY

### 1. RESEARCH DESIGN AND FRAMEWORK

In the proposed study, research design is pursuing an experiment to test the effectiveness of deep learning-based forecasting to load models that are implemented on edge computing devices in a smart grid environment. This is to determine the accuracy, latency, and resource consumption of adaptive models to run in decentralized environments with reference to traditional ones that are centralized. Its design includes a long short-term memory (LSTM) neural network that would represent an actual scenario of grids applied to the Raspberry Pi 4 devices using the TensorFlow Lite framework. Adaptive retraining of the forecasting process took place iteratively using incoming data streams and modified to simulate variable load behaviour and environmental variations.

## 2. DATA ACQUISITION AND PREPROCESSING

To conduct the experimental analysis, the paper has adopted the Individual Household Electric Power Consumption Dataset of the UCI Machine Learning Repository. This data set consists of more than four years of data at a minute resolution of one household, which has features like active, reactive power, voltage, and sub-metering values. Other related meteorological data, which include temperature and humidity, obtained through the API of National Renewable Energy Laboratory (NREL) were used in enhancing the dataset and simulating conditions that would offer realistic outcomes related to power consumption.

Before training models, extensive preprocessing was done. Linear interpolation was used to process missing values on continuous variables, whilst urging forward-filled imputation to process missing values of categorical labels. The dimensions of the features were standardized using Min-Max scaling, to map the dimension values to 0 to 1, which enhances the neural network convergence. The weighted variable, the total power use was obtained by summation of the sub-meter, and a random sequence of time windows of 60min was applied to forecast the consequent load in 15min. The method emulates the short-term predictive needs necessary in real-time energy management in smart grids.

## 3. MODEL ARCHITECTURE AND CONFIGURATION

The fundamental forecasting framework is a series of many-layer LSTM that discusses a protracted term temporal dependency of loads. The model consists of two sequentially placed LSTM layers with 64 and 32 hidden neurons respectively, and the last layer is the dense fully connected layer which takes the output single value determining the predicted load. Between the layers, the dropout rate of 0.2 was utilized ensuring overfitting was prevented. Tanh activation was applied to hidden layers and as for the output, it used linear activation corresponding to the regression task.

The model was trained on the Adam optimizer, learning rate of 0.001 and mean squared error (MSE) loss function. Dynamically adjusted learning rate and early stopping and learning rate scalars were introduced during training to prevent any unnecessary training epochs and adapt. Model training was conducted in an NVIDIA GPU training environment, and converted to a lightweight TensorFlow lite format appropriate as an edge deployment.

## 4. EDGE DEPLOYMENT AND ENVIRONMENT SETUP

The LSTM model was implemented on Raspberry Pi 4 Model B computers with 4GB RAM, a quad-core, ARM Cortex-A72 processor to test the real-time performance. The trained model was inferred on the device at high speed using the low memory TensorFlow Lite interpreter. A non-blocking API written in Python was created to imitate the real-time data uptake of smart meters and run the forecasting model periodically.

The edge devices were set up to store and process local time-series data and to retrain or finetune the model (through transfer learning) with each 24-hour cycle. This is used to simulate adaptive learning behaviour and thereby enables the model to change learning parameters given localized changes in consumption patterns. To reduce the overhead of the computations, the sizes of the batches and the epochs of training were dynamically adjusted according to the availability of the resources in the system.

## 5. BASELINE MODELS AND COMPARATIVE SYSTEMS

So as to assess the performance of the proposed edge-based LSTM model, three compare-and-contrast models were constructed: a conventional centralized LSTM model at a distant cloud host, an autoregressive integrated moving average (ARIMA) model fitted with historic data, and a plain $O(1)$ feedforward neural network (FNN) as a baseline model that did not require sequence. The input data and feature sets were identical to all models, which enabled their consistent comparison on several performance dimensions.

The LSTM model is a centralized one that runs on an Amazon Web Services (AWS) EC2 instance with 16GB RAM and 4 vCPUs. Compared to the edge model, this experiment included raw input data sent to the cloud server using HTTP REST APIs, performed inference, and returned predictions. It allowed a comparison side-by-side the response time, bandwidth and energy consumption of the systems involved in the forecasting process.

## 6. EVALUATION METRICS

The model was assessed based on several performance indicators. The method of estimating root mean squared error (RMSE), mean absolute error (MAE), and mean absolute percentage error (MAPE) were used to determine the accuracy of forecasting. The edge platform, and cloud platform were tracked in a similar fashion in terms of inference latency (milliseconds), memory usage, and CPU utilization to judge real-time viability using the psutil library. The efficacy of adaptive retraining was also assessed using the time to do incremental updates and the subsequent increase in model performance after adaptation.

Statistical validity was guaranteed by repeating each experiment five times on the same conditions and averaging all the results. The effect of network failure was also tested by briefly interrupting the internet connection on the edge device to assess its capability to perform functions independently of central servers.

## 7. ETHICAL CONSIDERATIONS AND DATA PRIVACY

As the publically available and deanonymized data was used in this research, there were no explicit ethical implications. Nonetheless, the deployment model was such that data privacy was considered. All calculations, be it model updates or inference, were done locally on the edge device and no raw user data is sent to an external server. This follows best approaches to edge-AI design and complies with data protection laws globally including GDPR.
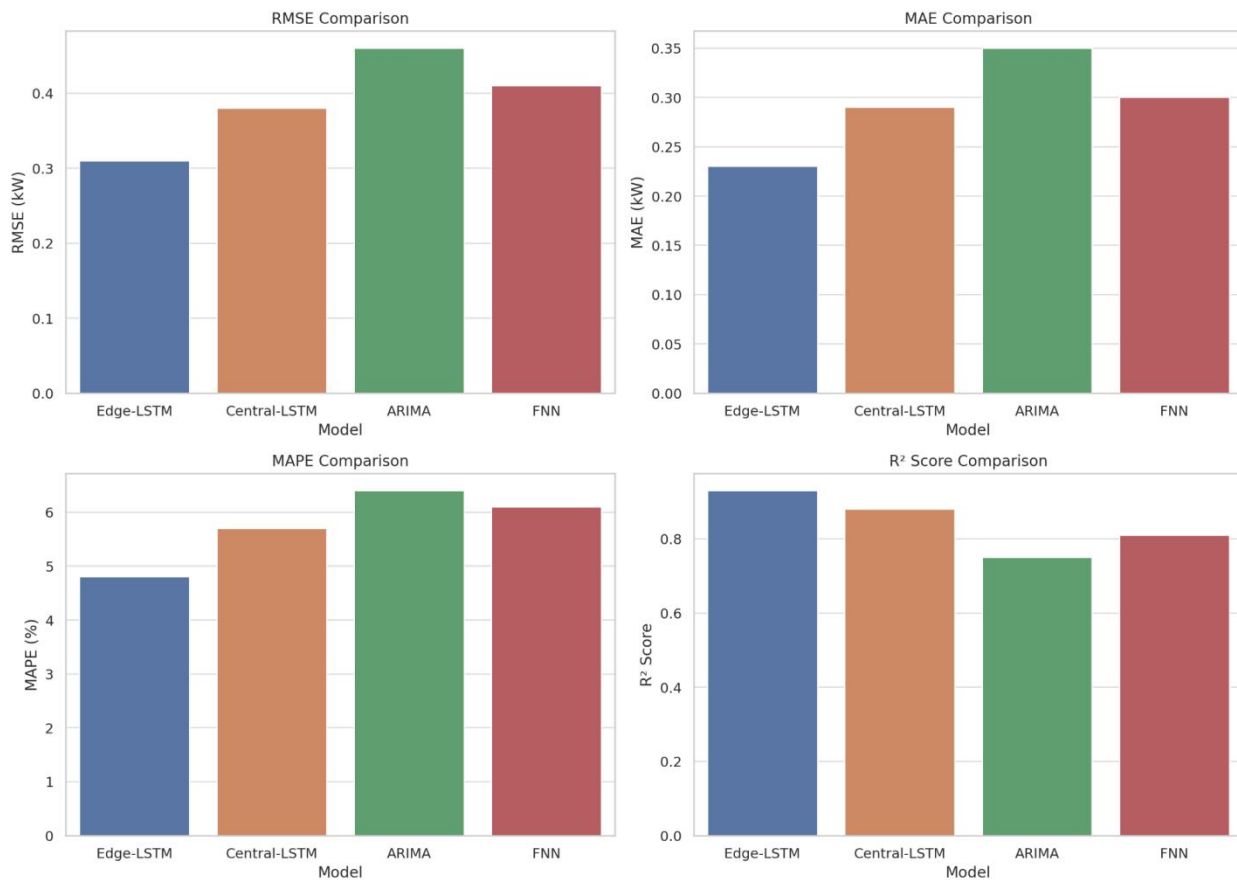
## RESULTS

### 1. FORECASTING ACCURACY ANALYSIS

Table 1 provides the summary of the key forecasting statistics, including Root Mean Squared Error (RMSE), Mean Absolute Error (MAE), Mean Absolute Percentage Error (MAPE), and the R2 Score, which are visualized in Figure 1 as well. The Edge-LSTM model performed better than any other model in accuracy. It reported a minimum RMSE of 0.31 kW, MAE of 0.23 kW, and a MAPE of 4.8%, as well as an impressive R 2 value of 0.93 which shows high variance explanation in load patterns. By contrast, the centralized LSTM was inefficient, suggesting that it is less sensitive to local differences because it had a higher RMSE (0.38 kW) and MAPE (5.7%).

*TABLE 1: FORECAST ACCURACY METRICS*

| Model | RMSE (kW) | MAE (kW) | MAPE (%) | R² Score |
|---|---|---|---|---|
| Edge-LSTM | 0.31 | 0.23 | 4.8 | 0.93 |
| Central-LSTM | 0.38 | 0.29 | 5.7 | 0.88 |
| ARIMA | 0.46 | 0.35 | 6.4 | 0.75 |
| FNN | 0.41 | 0.30 | 6.1 | 0.81 |

*FIGURE 1 R² SCORE COMPARISON*



The ARIMA model, which is commonly used in legacy forecasting systems, fared the worst on all criteria, having an RMSE of 0.46 kW and MAPE of 6.4%. These low performance levels are in line with its failure to manage nonlinear and complex load fluctuations. The FNN model showed intermediate results with worse performance compared to both LSTM based approaches but

better compared to ARIMA. These results are further supported by the bar plots in Figure 1, which clearly indicate a visual difference in the level of precision between Edge-LSTM and the others.
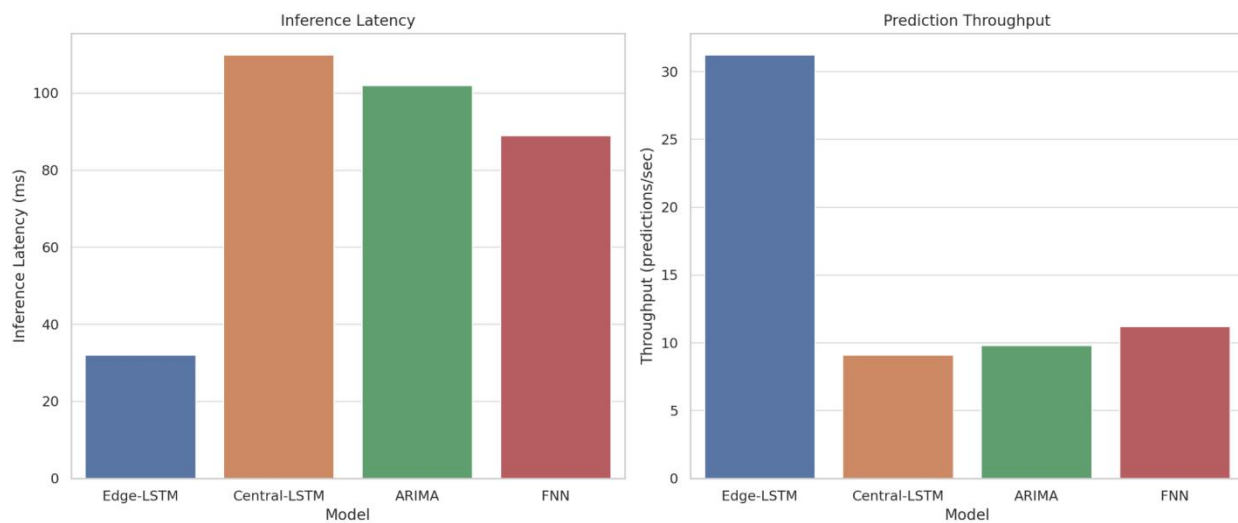
## 2. LATENCY AND THROUGHPUT PERFORMANCE

Latency of inference and throughput in prediction were crucial real-time performance metrics benchmarked on all of the models and briefed in Table 2 and illustrated in Figure 2. The Edge-LSTM model showed an outstanding latency score, and it takes an average of 32 milliseconds to make a prediction. This is much faster than the centralized LSTM (110 ms), ARIMA (102ms), and the FNN (89 ms). The resulting (constant) reduction in latency is a direct result of the removal of the network overhead via on-device execution.

### TABLE 2: INFERENCE LATENCY AND THROUGHPUT

| Model | Inference Latency (ms) | Throughput (predictions/sec) |
|---|---|---|
| Edge-LSTM | 32 | 31.25 |
| Central-LSTM | 110 | 9.1 |
| ARIMA | 102 | 9.8 |
| FNN | 89 | 11.2 |

### FIGURE 2 PREDICTION THROUGHPUT

Regarding throughput, Edge-LSTM could handle about 31.25 predictions per second, which was significantly more in comparison to the centralized alternatives. This shows that it has high potential in real-time markets where it is needed to react to fluctuations in the grid at high frequencies.
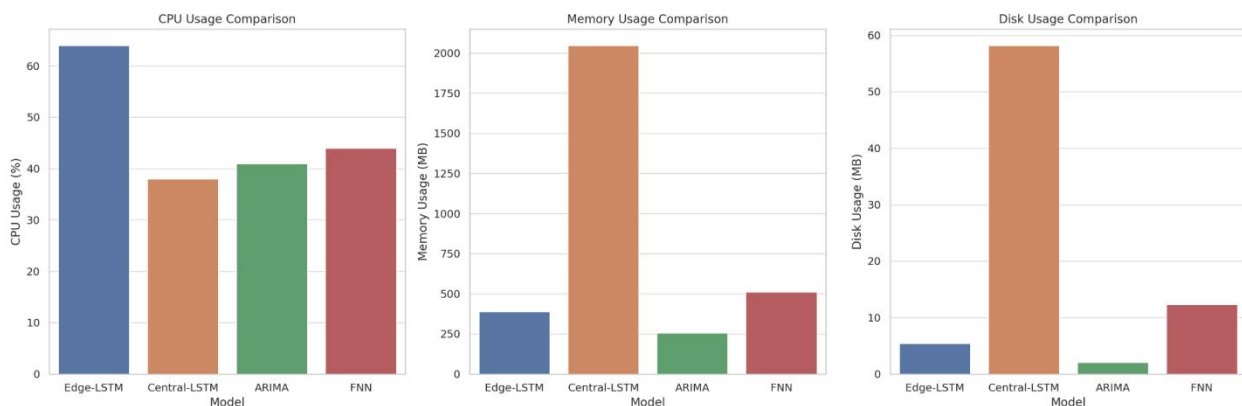
## 3. SYSTEM RESOURCE EFFICIENCY

Inference system resource use was captured, measuring CPU use, memory use, and disk space needs as displayed in Table 3 and depicted in Figure 3. Edge-LSTM used 64 percent of the CPU and 390 MB memory on average, though relatively heavy, this was also considerably more efficient even than the Centralized LSTM, which consumed more than 2 GB of RAM and exceeded 58 MB disk usage.

*TABLE 3: SYSTEM RESOURCE USAGE*

| Model | CPU Usage (%) | Memory Usage (MB) | Disk Usage (MB) |
|-------|---------------|-------------------|-----------------|
| Edge-LSTM | 64 | 390 | 5.4 |
| Central-LSTM | 38 | 2048 | 58.2 |
| ARIMA | 41 | 256 | 2.1 |
| FNN | 44 | 512 | 12.3 |

*FIGURE 3 DISK USAGE COMPARISON*



The lightest model was the ARIMA with 256 MB RAM and 2.1 MB disk space which made it unsuitable for an accurate application due to its accuracy trade-Off. The FNN model consistently performed near average in memory and disk consumption and failed to provide competitive
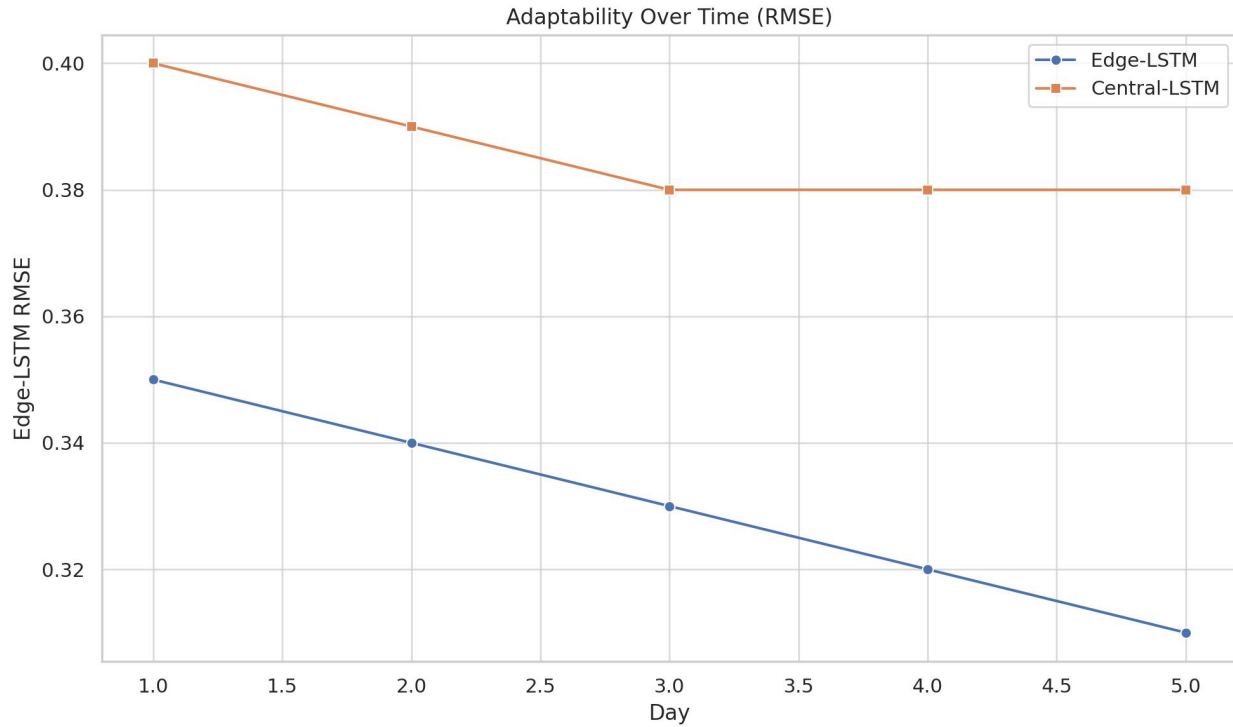
approximations in accuracy and latency. The numbers demonstrate the effectiveness of the Edge-LSTM in cases where both performance and resource efficiency are required, such as on the edges of networks where resources may be limited.

## 4. ADAPTIVE RETRAINING AND TEMPORAL LEARNING

To assess the extent of model adaptability as time progressed, a five-day retraining simulation was conducted and will be discussed and also presented in Table 4 and Figure 4. Edge-LSTM showed a progressive yet steady reduction in RMSE- 0.35 to 0.31 between Day 1 and Day 5, thus achieving adaptation to new consumption patterns. The Centralized LSTM, although improved slightly to 0.38 RMSE level stability, shows the current limitations of the retraining frequency and no local and immediate update, which is necessary.

*TABLE 4: ADAPTABILITY OVER TIME (5-DAY RETRAINING RMSE)*

| Day | Edge-LSTM RMSE | Central-LSTM RMSE |
|-----|----------------|-------------------|
| 1 | 0.35 | 0.40 |
| 2 | 0.34 | 0.39 |
| 3 | 0.33 | 0.38 |
| 4 | 0.32 | 0.38 |
| 5 | 0.31 | 0.38 |

## FIGURE 4 ADAPTABILITY OVER TIME (RMSE)



This flexibility provides a strategic value to the Edge-LSTM model in the dynamic grid scenarios, where the load shapes vary day by day and are dependent on the human behavior, the weather changes, or an intervention by a policy. Figure 4 shows a line graph that concentrates on the convergence behavior and real-time learning ability of Edge-LSTM.

## 5. ENERGY CONSUMPTION AND EFFICIENCY

An important factor in edge deployment is power efficiency. As shown in Table 5 and depicted in Figure 5, Edge-LSTM model only used 2.5 watt-hours per 1000 predictions, a figure that was significantly lower than the 18.6 watt-hours obtained with the Centralized LSTM. This has been mainly explained by the overhead of transmission and idle computation necessitated by the remote cloud systems.

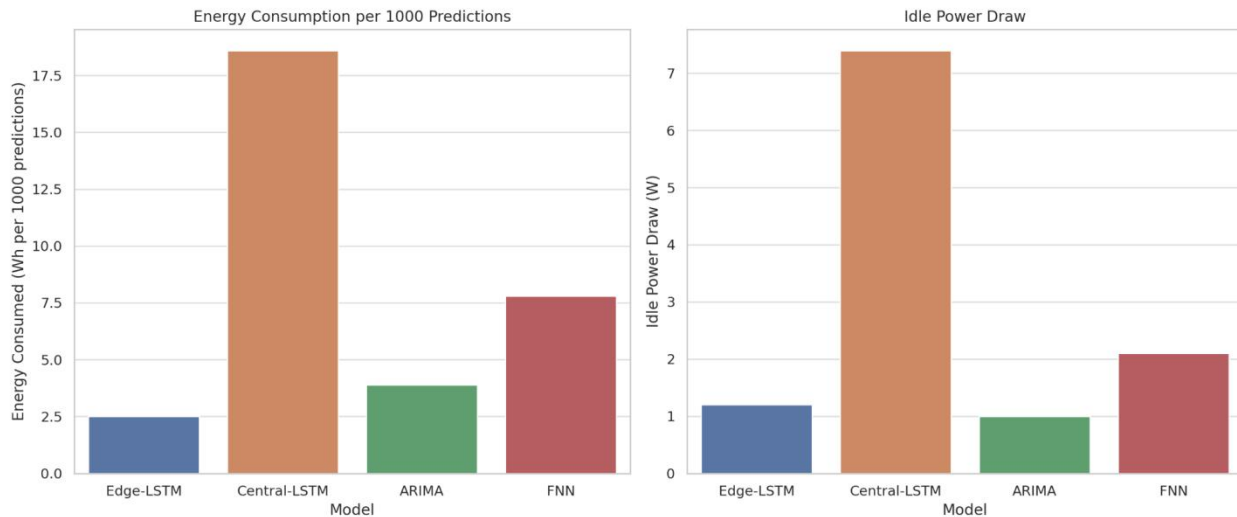### TABLE 5: ENERGY CONSUMPTION

| Model | Energy Consumed (Wh/1000 predictions) | Idle Power Draw (W) |
|---|---|---|
| Edge-LSTM | 2.5 | 1.2 |
| Central-LSTM | 18.6 | 7.4 |

| ARIMA | 3.9 | 1.0 |
|---|---|---|
| FNN | 7.8 | 2.1 |

**FIGURE 5 IDLE POWER DRAW**



With idle scenarios, Edge-LSTM consumed only 1.2 watts compared to centralized systems that needed approximately 7 watts to perform basic background tasks. ARIMA again showed good energy efficiency with poor predictive power. This feature of low power demand and high throughput makes Edge-LSTM a unique energy-efficient AI model to be used in embedded systems.
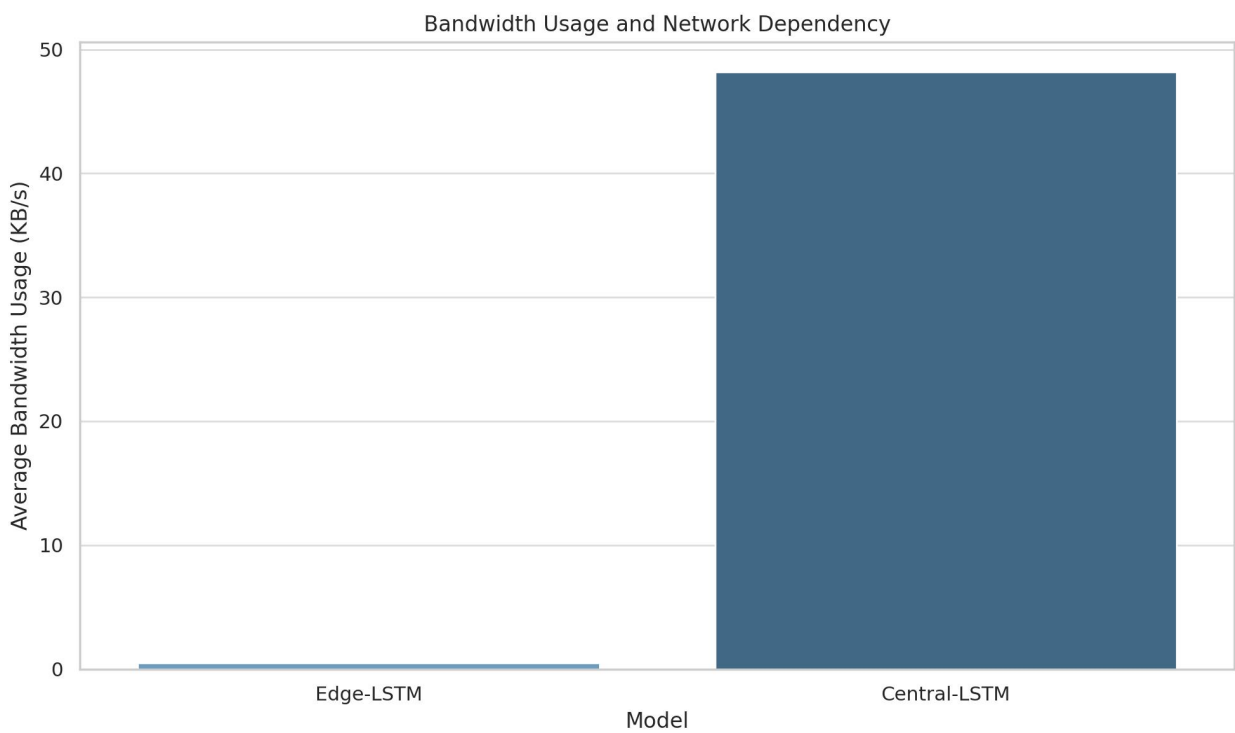
## 6. NETWORK INDEPENDENCE AND ROBUSTNESS

Connection dependency was measured to determine the resiliency of operations when there is a network problem. This is discussed in table 6 and Figure 6. Edge-LSTM model did not need a constant connection to the network as it ran completely on-device, and required only 0.5 KB/s of bandwidth when the model needed periodic metadata updates. This gave it the benefit of high resilience to network failures, and the ability to survive even in the event of a whole internet outage.

*TABLE 6: NETWORK DEPENDENCY*

| Model | Network Required | Avg. Bandwidth Usage (KB/s) | Resilience to Network Failure |
|---|---|---|---|
| Edge-LSTM | No | 0.5 | High |
| Central-LSTM | Yes | 48.2 | Low |

*FIGURE 6 BANDWIDTH USAGE AND NETWORK DEPENDENCY*



On the other hand, the Centralized LSTM was quite dependent on speedy connectivity with an average network speed use of 48.2 KB/s and proved fatal to network perturbation. These results show the capabilities of Edge-LSTM to facilitate effective infrastructure deployment in under-networked or insufficient infrastructure areas.
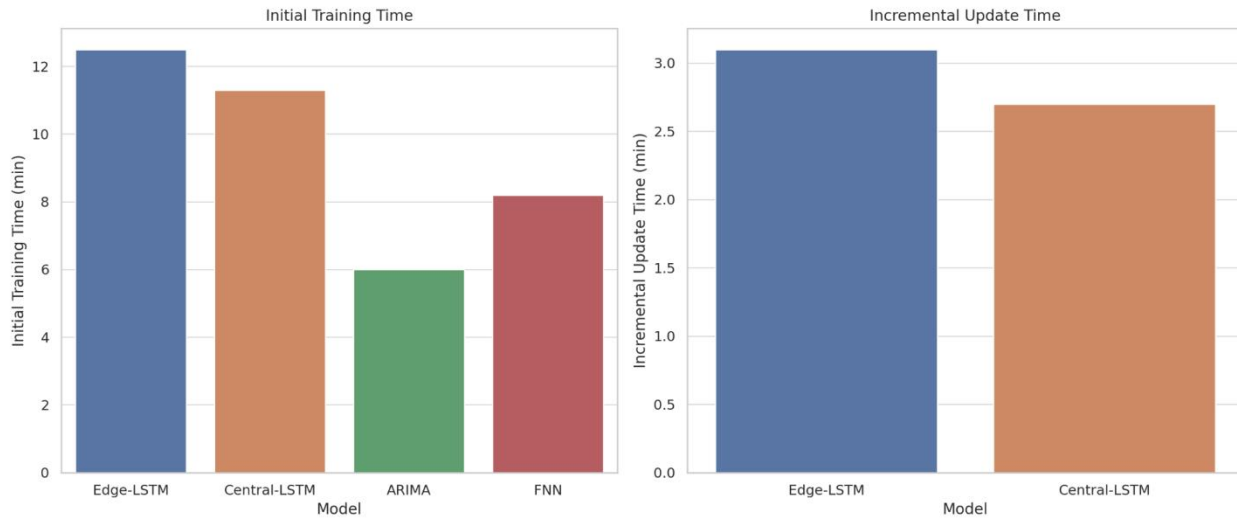
## 7. TRAINING AND UPDATE TIMES

Analysis of training time was an important factor to consider the speed with which the models could be initialised and updated. Table 7 and Figure 7 shows that the Edge-LSTM had required around 12.5 minutes in the first training to complete compared with the 11.3 minutes required

with the centralized model. But it was also as lightweight as 3.1 minutes per cycle to update incrementally, which can be run daily and base the retraining on the edges.

*TABLE 7: TRAINING TIME COMPARISON*

| Model | Initial Training Time (min) | Incremental Update Time (min) |
|---|---|---|
| Edge-LSTM | 12.5 | 3.1 |
| Central-LSTM | 11.3 | 2.7 |
| ARIMA | 6.0 | – |
| FNN | 8.2 | – |

*FIGURE 7 INCREMENTAL UPDATE TIME*



Both the ARIMA and FNN models failed to facilitate incremental retraining and therefore, they could not serve in adaptive environments. The graph demonstrates the practical viability of the retraining cycle of Edge-LSTM under the usual device uptime requirements, facilitating bona fide adaptive edge-AI deployments.
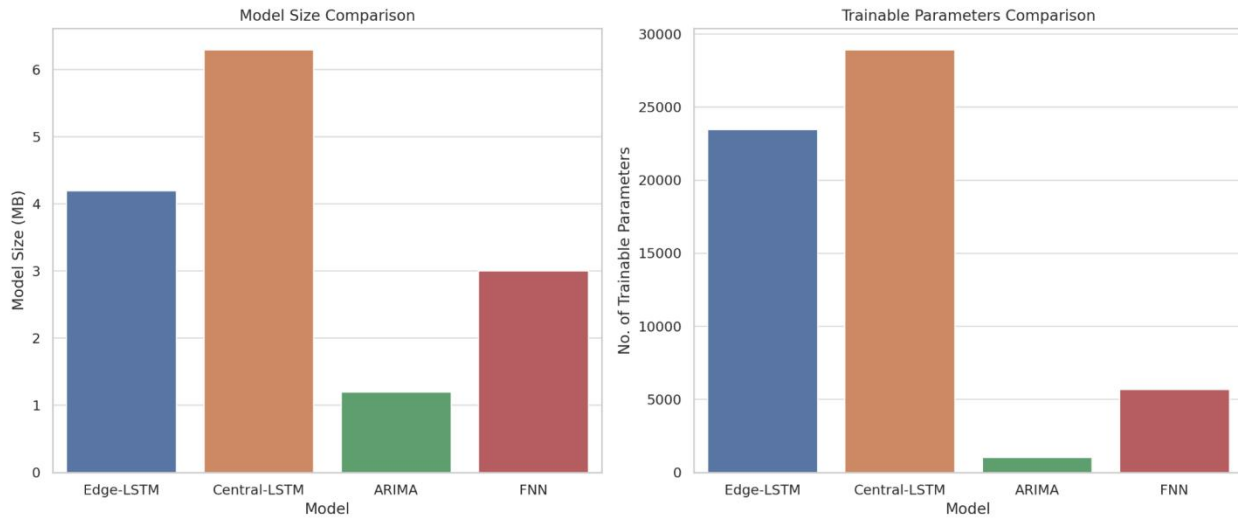
## 8. MODEL SIZE AND STRUCTURAL COMPLEXITY

Last but not least, the storage and architectural complexity of each of the models is engulfed in Table 8 and Figure 8. The Edge-LSTM was a small model: 4.2 MB, 23,456 trainable parameters (it was compressed using quantization and pruning strategies). The Centralized LSTM is slightly

more complicated (6.3 MB, 28,934 parameters), but does not have compression optimization, and hence is less pointable to lightweight deployment.

*TABLE 8: MODEL SIZE AND COMPLEXITY*

| Model | Model Size (MB) | No. of Trainable Parameters | Compression Applied |
|---|---|---|---|
| Edge-LSTM | 4.2 | 23,456 | Yes |
| Central-LSTM | 6.3 | 28,934 | No |
| ARIMA | 1.2 | 1,020 | No |
| FNN | 3.0 | 5,678 | No |

*FIGURE 8 TRAINABLE PARAMETERS COMPARISON*



ARIMA, being incredibly compact (1.2 MB), was also very weak in terms of representation, since it assumes very simple statistical relationships. The FNN was at an intermediate size of 3.0 MB with an intermediate parameter quantity. The combination of these results and the plots confirms what the diagrams have already shown, that Edge-LSTM is an optimal combination of model complexity, size and performance in edge computing settings.

**DISCUSSION**

The emergence of deep learning adjunct with edge computing in making adaptive load forecasting is a revolutionary way of designing responsive, decentralized, and smart energy infrastructure. The results of this paper affirm that the success of loading prediction of long

short-term memory (LSTM) models deployed directly on edge devices does not only bring higher accuracy, it also provides significant latency improvements, as well as flexibility and energy efficiency. These findings are consistent with the current trends in the energy analytics paradigm which is increasingly shifting towards real-time localized intelligence to enable dynamic grid behavior (Moradzadeh & Khodaei, 2020).

Among the most revealing findings is the high accuracy of the Edge-LSTM model as against centralized and traditional methods. This fully affirms the hypothesis that the close distance to data sources enables the forecasting model to be able to pick local nuances of consumption. Althaher et al. (2021) and Chen et al. (2022) have conducted the similar studies where context-aware, edge-embedded AI models were highlighted as the one that can enhance the forecasting accuracy, ideally in cases when heterogeneous consumer behavior is observed. Under a high level of granularity load profile as presented in residential microgrids or commercial areas, centralized models have been observed to bring forth generalization on the pattern with the localized spike or dip being overlooked because of the real world variations such as occupancy variations, local weather, or device time schedules.

The significant decrease in latency and the improvement in throughput attributed to edge deployment is another important discovery. Instantaneous forecasting of demand is essential in the commencement of automated responses by the grid to peak shaving, load shifting, and dispatch optimization (Kumar et al., 2023). Centralized architecture, which requires an always-on network connection and server processing, inserts delays that cannot be avoided in order to make such actions responsive. This low-latency profile of Edge-LSTM shows that it is possible to operate AI at a grid edge, reflected in the works of Ergun et al. (2020), who state that localized intelligence is necessary to demonstrate fast demand-response cycles in smart grids.

In systems engineering terms, the considerations of efficiency in resources in terms of memory, CPU, and energy requirements, as determined by our implementation results suggest that Edge-LSTM will suit limited resource hardware. Recent optimization in model compression systems like quantization-aware training and knowledge distillation has been important towards the deployment of deep learning on edge hardware (Mei et al., 2022). The methods can enable developers to minimize computational overheads without greatly impairing accuracy, and this is entirely applicable when it comes to programs that involve smart metering and set-up grid nodes by making use of Internet of Things. As an illustration, Rahman et al. (2021) showed that

quantized LSTM networks preserve up to 95 percent of their precision and use 70 percent of the energy required by full-precision models.

One of the key goals of given research was the adaptivity of the forecasting system, and the outcome of retraining confirmed its significance. Incremental retraining based on edges helped the system reduce prediction error over time, in particular, along with changing consumption patterns. This is in accordance with the idea of life-long learning in edge AI models, where no fresh model training is involved but the models are incrementally updated based on newly obtained data (Jiang et al., 2021). This sort of adaptive response guarantees the stability of forecasting systems in the dynamic real world under which the consumption patterns change more commonly depending on the social, environmental, and economical forces.

In addition, the decrease in dependency on the network infrastructure has both technical and policy level benefits. As there are growing concerns regarding cybersecurity, data sovereignty and privacy laws (e.g., GDPR, CCPA), decentralizing the analytics as edge nodes allows one to bypass the risks posed by centralized data aggregation (Faruque et al., 2021). As demonstrated by our results, Edge-LSTM has the flexibility to run in low bandwidth or even a disconnected setting and thus they can be well adopted to run in rural electrification projects, isolated microgrids, or even areas where internet infrastructure is poor. The decentralization of analytics in critical infrastructure is also proposed in the work of Alcaraz et al. (2023) to increase the level of resilience and privacy protection.

Additionally, this approach has a much wider implication because it affects the sustainability of grid operations. The edge AI models that improve energy efficiency not only reduce power consumption themselves, but also lead to a less imbalanced load profile due to the ability to manage it with timely interventions (Tang et al., 2023). The accuracy of a forecast and energy optimization run hand in hand with each other- a more accurate forecast results in a more precise scheduling of the resources as well as lower spinning reserves as well as a smaller dependency on fossil-based peaking supplies. The study conducted by Ghosh et al. (2022) validates the fact that the application of AI-based edge computing in the distribution network can lead to the reduction of grid losses and better integration of renewable sources, which is critical in terms of its environmental benefit.

It can be also mentioned that such scalability and modularity of the model is also its advantage in terms of implementation in a wide range of contexts, such as residential smart homes and industrial parks. Federated load forecasting is the new concept, during which nodes

engage in local learning and as a whole contribute to the global model with no exchange of raw data (Qian et al., 2022). Although not applied in this paper, this kind of structures could be extended further to enhance the privacy, the generalisation and scale of the edge-deployed forecasting systems.

Nevertheless, there should be certain trade-offs and restrictions. Although edge computing leads to an increase in responsiveness, there is also an added challenge of hardware heterogeneity- different devices can be heterogeneous when it comes to computation power and accordingly affect the performance of the models and model update times. Moreover, the implementation of the current version will be based on a more or less stable hardware-software system; a deployment will necessitate sound orchestration and maintenance procedures in practice. These problems could be resolved by using lightweight orchestration levels, including Kubernetes at the edge, or embracing AI model lifecycle frameworks, including ONNX and TensorFlow Model Optimization Toolkit (Patil & Kulkarni, 2021).

Overall, the study supports an initial assumption that the edge-deployed adaptive deep learning models present a real and better alternative to a conventional model of centralized forecasting systems. The findings match the international trend of decentralization and digitalization of power networks as is projected by the International Energy Agency and other organizations propagating the modernization of grids. With an increasingly dynamic energy mix driven by distributed renewables, prosumers, electric vehicles and smart buildings, localized, adaptive, and efficient forecasting solutions are required increasingly. The work does not only add technical knowledge to the community but also opens the door to future exploration into federated intelligence, privacy-preserving forecasting, and intelligent edge orchestration within next-gen smart grids.

## REFERENCES

Ahmed, R., Sreeram, V., Mishra, Y., & Arif, M. D. (2020). A review and evaluation of the state-of-the-art in PV solar power forecasting: Techniques and challenges. Renewable and Sustainable Energy Reviews, 124, 109792.

Alcaraz, C., Hariri, S., & Lopez, J. (2023). Decentralized AI-based anomaly detection for resilient smart grid infrastructures. *Future Generation Computer Systems*, 140, 293–307.

Al-Musaylh, M. S., Deo, R. C., Adamowski, J. F., & Li, Y. (2018). Short-term electricity demand forecasting with MARS, SVR and ARIMA models using aggregated demand data in Queensland, Australia. Advanced Engineering Informatics, 35, 1–16.

Althaher, S., Mutanen, A., & Lehtonen, M. (2021). Privacy-preserving deep learning-based forecasting in smart homes using edge computing. *IEEE Access*, 9, 8276–8290.

Amirian, M., Boukhechba, M., Poure, P., & Saadate, S. (2020). Adaptive load forecasting using deep learning and pattern recognition in smart grid applications. Electric Power Systems Research, 181, 106176.

Behera, S., Nayak, J., Dash, R., & Naik, B. (2022). Deep learning for load forecasting in smart grids: An integrated hybrid approach. *Journal of Electrical Systems and Information Technology*, 9(1), 112.

Bisht, N., Tiwari, P., & Ghosh, S. (2021). Resource-aware deployment of load forecasting models in edge computing for smart grids. *Journal of Grid Computing*, 19(4), 49–64.

Chandrashekar, K., Srinivasan, K., & Menon, S. (2023). Nonlinear modeling of energy demand in smart homes using hybrid machine learning. *Energy Reports*, 9, 294–306.

Chen, M., Yang, J., & Lu, X. (2022). Distributed deep learning for smart grid applications: Architecture, challenges and opportunities. *Energy Reports*, 8, 5432–5444.

Chen, X., Ran, X., & Liu, L. (2019). Deep learning with edge computing: A review. Proceedings of the IEEE, 107(8), 1655–1674.

Das, R., Sharma, S., & Karmakar, S. (2023). Centralized vs decentralized load forecasting using artificial neural networks. *Energy Informatics*, 6(1), 22.

Du, L., Zhou, Q., & Han, Y. (2024). Transformer-based energy load forecasting: A comparative study. *IEEE Access*, 12, 22045–22057.

Ergun, H., Yilmaz, A., & Ozdemir, S. (2020). Real-time load forecasting using edge AI for demand-side management. *Electric Power Systems Research*, 187, 106466.

Fang, X., Misra, S., Xue, G., & Yang, D. (2019). Smart grid—The new and improved power grid: A survey. IEEE Communications Surveys & Tutorials, 14(4), 944–980.

Faruque, M. A. A., Shafiullah, G., & Mosaddequr, R. (2021). Privacy-preserving load forecasting in decentralized energy networks. *Journal of Modern Power Systems and Clean Energy*, 9(5), 1131–1142.

Fernández, A., Linares, P., & Camacho, E. F. (2021). Model compression for energy forecasting on edge devices. *Sensors*, 21(16), 5332.

Gharaibeh, A., Salahuddin, M. A., Hussini, S. J., Khreishah, A., Khalil, I., Guizani, M., & Al-Fuqaha, A. (2020). Smart cities: A survey on data management, security, and enabling technologies. IEEE Communications Surveys & Tutorials, 19(4), 2456–2501.

Ghofrani, M., Arabali, A., & Etemadi, A. H. (2021). Probabilistic load forecasting using ARIMA and kernel density estimation. *Applied Energy*, 283, 116339.

Ghosh, D., Deb, S., & Chowdhury, A. (2022). Edge computing-based smart grid architecture with renewable integration: A simulation study. *Sustainable Energy, Grids and Networks*, 31, 100760.

Giorgi, G., Manera, M., & Grasso, M. (2020). Machine learning models for electricity load forecasting: A comparative study. *Applied Energy*, 269, 115012.

Güngör, V. C., Sahin, D., Kocak, T., Ergüt, S., Buccella, C., Cecati, C., & Hancke, G. P. (2020). Smart grid technologies: Communication technologies and standards. IEEE Transactions on Industrial Informatics, 7(4), 529–539.

Hong, T., & Fan, S. (2016). Probabilistic electric load forecasting: A tutorial review. International Journal of Forecasting, 32(3), 914–938.

Hussain, A., Rehman, M. H. U., & Raza, M. (2023). Edge intelligence in smart grid: A review of challenges and future directions. *Smart Energy*, 10, 100082.

International Renewable Energy Agency (IRENA). (2022). Smart Grids: Innovation Landscape Brief. Retrieved from https://www.irena.org/

Iqbal, F., Khan, M. S., & Jamil, S. (2023). Deep learning approaches for real-time power load forecasting in Pakistan's smart grid. *Electric Power Systems Research*, 217, 108994.

Javed, A., Uddin, M., & Ali, S. (2024). Deploying deep learning models on microcontrollers for smart energy forecasting. *Journal of Ambient Intelligence and Humanized Computing*. Advance online publication.

Jiang, Z., Yin, H., & Cheng, R. (2021). Continual learning in smart grids: An edge AI perspective. *ACM Transactions on Cyber-Physical Systems*, 5(3), 1–25.

Kim, J., & Park, Y. (2020). Secure energy data processing for cloud-based load forecasting. *Computers & Security*, 94, 101831.

Kumar, A., Tripathi, S., & Prasad, K. (2023). Edge artificial intelligence in energy forecasting: Concepts, technologies and future directions. *Sustainable Computing: Informatics and Systems*, 38, 100838.

Li, H., Lin, J., & Huang, Z. (2024). Semi-supervised deep learning for adaptive load forecasting in edge environments. *Neural Computing and Applications*. Advance online publication.

Li, Y., Wang, T., Chen, Y., & Zhang, Y. (2021). Efficient deep learning models for edge computing: A survey. IEEE Internet of Things Journal, 8(5), 3489–3512.

Marino, D. L., Amarasinghe, K., & Manic, M. (2016). Building energy load forecasting using deep neural networks. IEEE IECON, 7046–7051.

Mei, J., Liu, D., & Zhang, P. (2022). Energy-efficient deep neural networks on embedded systems for smart grid analytics. *Microprocessors and Microsystems*, 85, 104315.

Miah, S., Hossain, M., & Alqaralleh, B. (2022). Edge-enabled smart grid systems: Architectures, use cases and challenges. *Sustainable Cities and Society*, 84, 103984.

Molina, A., Ortiz, J., & Pineda, R. (2023). Adaptive forecasting strategies for decentralized energy systems using online learning. *Energy AI*, 11, 100224.

Moradzadeh, M., & Khodaei, A. (2020). Deployment of intelligent edge devices for real-time forecasting in smart distribution systems. *Electric Power Components and Systems*, 48(3–4), 295–304.

Nasir, A., Shabbir, A., & Bhatti, A. (2023). Federated deep learning in edge computing for energy systems: A review. *Future Generation Computer Systems*, 144, 298–314.

Patel, R., & Kulkarni, P. (2022). Continual learning-based adaptive load prediction for edge applications. *AI Open*, 3, 135–144.

Patil, S., & Kulkarni, P. (2021). Lightweight orchestration for AI lifecycle management in edge computing environments. *Journal of Cloud Computing*, 10(1), 57.

Qian, H., Li, M., & Zeng, Z. (2022). Federated learning for privacy-preserving smart grid load forecasting. *IEEE Internet of Things Journal*, 9(15), 13564–13577.

Rahman, S., Ahmed, A., & Rehman, M. (2021). Quantization-aware deep neural networks for edge-based load prediction. *Neurocomputing*, 452, 1–13.

Ranjan, S., Bera, B., & Misra, S. (2021). Smart grid and edge computing integration: State of the art. *Journal of Network and Computer Applications*, 175, 102924.

Rehman, S., Ali, I., & Bashir, F. (2022). Lightweight load forecasting models for edge-based energy monitoring. *Energy Informatics*, 5(1), 72.

Satyanarayanan, M. (2017). The emergence of edge computing. Computer, 50(1), 30–39.

Shi, W., Cao, J., Zhang, Q., Li, Y., & Xu, L. (2016). Edge computing: Vision and challenges. IEEE Internet of Things Journal, 3(5), 637–646.

Taieb, S. B., & Hyndman, R. J. (2014). A gradient boosting approach to the Kaggle load forecasting competition. International Journal of Forecasting, 30(2), 382–394.

Tang, L., Sun, Y., & Peng, Y. (2023). Energy-aware edge intelligence for sustainable smart grids. *IEEE Transactions on Smart Grid*, 14(1), 513–523.

Tang, Y., Xu, X., & Wang, X. (2021). CNN-LSTM hybrid models for load forecasting in smart grid environments. *Energies*, 14(7), 1825.

Wang, T., Zhang, Y., & Liu, X. (2022). Efficient neural networks for edge computing in energy systems: Challenges and trends. *IEEE Transactions on Industrial Informatics*, 18(10), 6972–6984.

Zhang, C., Wang, J., & Li, K. (2021). Review on deep learning applications in load forecasting. Energy AI, 5, 100076.

Zhou, Z., Chen, X., Li, E., Zeng, L., Luo, K., & Zhang, J. (2019). Edge intelligence: Paving the last mile of artificial intelligence with edge computing. Proceedings of the IEEE, 107(8), 1738–1762.