

International Journal of

INTELLIGENT SYSTEMS AND APPLICATIONS IN ENGINEERING

ISSN:2147-6799 www.ijisae.org Original Research Paper

Designing A Novel Cyber Resilient System for Smart Home Utilizing Recursive Feature Elimination

¹Deva Rugved, ²Dr. C. Madhusudhana Rao, ³Chinthala Bharath Raj, ⁴N. Tanushri,

Submitted: 14/03/2024 **Revised**: 29/04/2024 **Accepted**: 06/05/2024

Abstract: The proliferation of Internet of Things (IoT) devices in smart homes has introduced significant security challenges, making advanced intrusion detection systems (IDS) essential for itigating potential threats. In this study, proposed a comprehensive framework for designing and implementing an IDS specifically tailored for IoT-enabled smart homes. The approach integrates LightGBM (LGBM) for attack detection and employs feature reduction techniques to enhance model performance and efficiency. This paper evaluated IDS using both IoT datasets and tested three feature reduction methods: Recursive Feature Elimination (RFE), Principal Component Analysis (PCA), and Linear Discriminant Analysis (LDA). The results show that RFE provided the best balance of detection accuracy and computational efficiency. The model effectively identified Mirai and Gafgyt attacks, achieving high accuracy rates of 99% on the datasets. Compared to previous systems, model demonstrated superior performance with reduced computational overhead, enhancing both detection accuracy and operational efficiency. Privacy and security considerations were carefully addressed throughout the design process to ensure user data protection and compliance with relevant regulations. Future work will focus on expanding the model's capability to detect a broader range of attacks, improving scalability, and integrating advanced IoT technologies to further refine the IDS framework.

Keywords: IoT Intrusion Datasets, Feature Reduction, Light Gradient BoostingMachine (LGBM).

I. INTRODUCTION

The rapid proliferation of Internet of Things (IoT) devices in modern smart homes has revolutionized how this paper interact with our living environments. These interconnected devices offer unparalleled convenience, automation, and efficiency, transforming homes into intelligent ecosystems. However, this connectivity comes with significant security challenges. As IoT devices often lack robust security measures, they become vulnerable to various cyber threats, leading to potential breaches that could compromise the privacy and safety of users [1]. Consequently, ensuring the cyber resilience of smart homes has become a critical area of research [2].

Intrusion Detection Systems (IDS) have emerged as a primary defense mechanism against cyber threats in IoT environments. These systems monitor network traffic and device behavior to detect suspicious activities indicative of potential intrusions [3]. However, the effectiveness of IDS is

Student Department of Computer Science and Engineering INSTITUTE
OF AERONAUTICAL ENGINEERING

Hyderabad,India 21951A0534@iare.ac.in

²Department of Computer Science and Engineering INSTITUTE OF AERONAUTICAL ENGINEERING

Hyderahad India

³ Student Department of Computer Science and Engineering INSTITUTE OF AERONAUTICAL ENGINEERING

Hyderabad,India

ch.bharathraj@gmail.com

⁴ Student Department of Computer Science and Engineering INSTITUTE OF AERONAUTICAL ENGINEERING

Hyderabad,India

21951A05M5@iare.ac.in

often hindered by the high dimensionality of the data generated by IoT devices, which can lead to increased computational costs and reduced detection accuracy [4]. To address these challenges, feature reduction techniques are employed to streamline the dataset by selecting the most relevant features, thereby enhancing the performance and efficiency of IDS [5].

Current research in the field of IoT security has explored various methods for feature reduction, including techniques like Recursive Feature Elimination (RFE), Principal Component Analysis (PCA), and Linear Discriminant Analysis (LDA). For instance, Zhou et al.proposed a feature-selection approach utilizing CFS-BA heuristic algorithms for dimensionality reduction, which aimed to enhance detection accuracy while handling the complexity of high-dimensional data [1]. Similarly, Gupta et al implemented ensemble methods and deep learning techniques that indirectly benefit from feature reduction to improve detection rates [6]. While these methods have shown promise, there is still a gap in understanding which technique is most effective for optimizing intrusion detection in smart homes, particularly when applied to realworld IoT datasets.

CYBER RESILIENCE IN SMART HOME

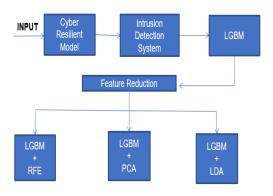


Fig 1 – Design and Architecture

The aim of the paper is to fill this gap by conducting a comprehensive evaluation of three feature reduction technique—RFE using LightGBM, PCA, and LDA—on an IoT intrusion detection dataset. By systematically applying these methods and assessing their impact on model performance as described in Fig 1, This Paper seek to identify the most effective technique for maintaining high detection accuracy while minimizing computational complexity. Our hypothesis is that RFE, due to its iterative nature and integration with LightGBM, will outperform PCA and LDA in terms of both feature selection and model accuracy [7][8].

| Dataset Name | BoTNeTIoT-L01 |
|----------------------------|--|
| Features | 23 |
| Total No.of Instances | 1048476 |
| Selected No.of Features | 23 (Features over a 10-second time window) |
| Total No.Of Classes | 2 (1 attack class + 1 benign class) |
| Predicted Attacks | Mirai, Gafgyt |

Table 1 – Dataset Description

The BoTNeTIoT-L01 dataset in Table 1 consists of 1,048,476 instances and comprises 23 features extracted over a 10-second time window. These features capture various network-related attributes, allowing for a comprehensive analysis of intrusion patterns. The dataset is labeled into two classes: a benign class representing normal network traffic and an attack class encompassing different types of malicious activity. The primary attacks detected in this dataset include Mirai and Gafgyt, which are among the most prevalent threats targeting IoT ecosystems.

II. RELATED WORK

Zhou et al. [1] proposed an intelligent intrusion detection system (IDS) framework that leverages feature selection and ensemble learning mechanisms. They heuristic CFS-BA algorithm introduced the dimensionality reduction to select the most relevant and distinct feature subsets while examining feature correlations. Their ensemble approach combined C4.5 and Random Forest (RF) with penalizing attribute algorithms, utilizing an average probability rule. They incorporated the probability distributions of base learners through voting techniques to enhance attack recognition performance. Their system was evaluated using the NSL-KDD, AWID, and CIC-IDS2017 datasets. However, their work did not address concerns related to time efficiency. Zhou et al.'s research does not address the time efficiency, scalability, or flexibility of their IDS system, and its complexity may make it difficult to develop and deploy in the real world.

Devprasad et al. [5] proposed a context-adaptive classification mechanism using hierarchy-based chi-square and bat algorithms. They tested their algorithm on the NSL-KDD and UNSW-NB15 datasets, employing Decision Tree (DT) and SVM algorithms. The ensemble classifier achieved a prediction accuracy of 89.43% and a false positive rate (FPR) of 3.215%. Despite these results, their model did not address execution time considerations.

Gupta et al. [6]. introduced a cost-sensitive network-based IDS designed to handle class imbalance through deep learning and ensemble algorithms. Their model, which was evaluated on CIDDS-001, NSL-KDD, and CICIDS2017, is divided into three stages: the first stage uses a deep neural network to distinguish normal from suspicious traffic, the second stage applies XGBoost for major attack classification, and the third stage employs Random Forest for minor attack classification. Their system demonstrated high detection rates, achieving 99% accuracy on NSL-KDD, 96% on CIDDS-001, and 92% on CICIDS2017. Although Gupta et al.'s study achieves great accuracy, its multi-stage strategy may restrict its scalability, flexibility to changing assault patterns, and practical implementation complexity.

Kumar et al. [4]. proposed a cyber-attack detection system utilizing Random Forest, K-Nearest Neighbors (KNN), and XGBoost algorithms. Evaluated on the BoT-IoT and DS2OS datasets, their system claimed detection rates between 90% and 100%. However, the authors acknowledged that not all threats could be detected, and potential for false alarms remained. They also noted that the XGBoost method, while accurate, does not address latency issues

Hadem et al. [3]. proposed an IDS using SVM within a Software-Defined Networking (SDN) framework, focusing on maximizing detection accuracy with minimal computational overhead and improved memory efficiency. Their proposed system recorded an accuracy of 95.98% on the full NSL-KDD dataset and 87.74% on a selected feature dataset

Htwe et al [7].. applied the Classification and Regression Tree (CART) method for their IDS architecture on the N-BaIoT dataset. They claimed that their classifier outperformed Naïve Bayes classifiers. However, their evaluation was limited to a single metric, accuracy, which may not provide a comprehensive assessment of model performance

Xu and Fan [8].proposed an IDS combining XGBoost with a logarithmic auto-encoder. Their model, evaluated on CICIDS2017 and UNSW-NB15 datasets, achieved accuracy scores of 99.92% and 95.11%, respectively. This study also evaluated the run-time performance of different classifiers, which adds a valuable dimension to their evaluation. While Xu and Fan's study demonstrates great accuracy, it may not be as generalizable to different datasets or attack situations, and it does not fully address the model's usefulness in dynamic, real-world contexts.

Le et al. [9]. developed a multiclass classification-based IDS for imbalanced datasets in Industrial IoT (IIoT). Their XGBoost classifier detected abnormal network traffic behavior, achieving attack detection rates of 99.9% and 99.87% on the TON_IoT and X-IIoTDS datasets, respectively. Their model demonstrated strong performance in detecting cyber-attacks. Le et al.'s work focuses on accurately identifying cyberattacks in IIoT, but it ignores possible problems with model scalability, generalizability to a variety of developing attack scenarios, and processing efficiency in real-time.

Awajan [10]. proposed a deep learning-based IDS for IoT devices that addresses five types of intrusions. The model, which relies on a deep neural network, showed an average accuracy of 93.74%. However, it requires retraining for each new IoT network, and the accuracy rate could be improved

III. CLASSIFICATION ALGORITHM

Light Gradient Boosting Machine (LGBM):

A histogram-based decision tree approach called Light Gradient Boosting Machine (LightGBM) significantly improves model execution times and reduces memory usage on computers, enhancing overall model efficiency. LightGBM is notably more optimized compared to other boosting ensemble decision tree algorithms [1]. It is recognized as a more advanced learning algorithm that is faster, more scalable, and more effective. The algorithm

excels in handling large data flows, such as those encountered in the design of intrusion detection models for IoT-enabled smart homes [2].

LightGBM employs leaf-wise tree growth, which differentiates it from many other boosting techniques. Unlike traditional methods that compute splits level-wise, LightGBM uses a pre-sorted approach and a histogram-based algorithm to determine the best split values [3]. The algorithm incorporates two novel methods: Exclusive Feature Bundling (EFB) and Gradient-based One-Side Sampling (GOSS). EFB addresses the limitations of conventional histogram-based methods by bundling features to reduce dimensionality without losing significant information. GOSS, on the other hand, enhances the model's accuracy by downsampling instances based on gradient sizes, focusing on large gradient samples while discarding smaller ones [4].

In comparison to uniform random sampling, GOSS improves accuracy by prioritizing large gradient samples that are undertrained, while small gradient samples are well-trained. Additionally, LightGBM's leaf-wise growth, as opposed to level-wise or depth-wise tree splitting used in other boosting algorithms, enables the decision tree to grow more efficiently [5]. Fig -2,This method of leaf-wise growth ensures that the model makes more precise splits, leading to better performance in various applications, including data science and intrusion detection systems [6].

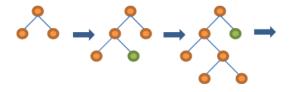


Fig -2 Light Gradient Boosting Machine

Recursive Feature Elimination:

Recursive Feature Elimination (RFE) is a potent feature selection method that enhances interpretability and model performance in machine learning. RFE selectively eliminates less relevant features repeatedly according to their impact on model correctness, in contrast to conventional approaches that use all features available [1]. The goals of this study are to improve model efficiency and forecast accuracy in a variety of areas by examining the ideas, methodology, and applications of RFE [2][3].

A significant development in feature selection methods is Recursive Feature Elimination (RFE), which provides a rationale for improving model performance and interpretability. It is a flexible tool for optimizing predictive models in practical applications due to its iterative character and interoperability with a wide

range of machine learning techniques [4][5]. The revolutionary power of RFE in promoting data-driven insights across domains and enabling well-informed decision-making is highlighted in this study [6][7].

Principal Computational Analysis:

This method lowers the dimensionality of high-dimensional data while maintaining its core structure [1][2]. PCA improves model performance, visualization, and understanding by converting data into a lower-dimensional space of principal components [3]. This study examines the benefits, applications, methods, and guiding principles of PCA in streamlining machine learning processes in a variety of fields [4][5].

Principal Component Analysis (PCA) is a fundamental approach in machine learning for dimensionality reduction. It provides insights into data structure and improves the efficiency and interpretability of models [6][7]. Its use in a variety of fields, including financial analysis, bioinformatics, image and text processing, and more, highlights its adaptability and significance in turning high-dimensional data into useful insights [8][9].

Linear Discriminant Analysis:

A traditional machine learning method called Linear Discriminant Analysis (LDA) is employed for both supervised classification and dimensionality reduction [1][2]. LDA converts high-dimensional data into a lower-dimensional space by optimizing class separability and decreasing intra-class variation, improving model performance and interpretability [3][4]. This work investigates the benefits, uses, and methods of Linear Discriminant Analysis (LDA) for classification job optimization in many areas [5][6].

Linear Discriminant Analysis (LDA) is a fundamental methodology in machine learning that provides a systematic method for classifying and reducing dimensionality [7][8]. It is essential for activities ranging from pattern identification and medical diagnosis to natural language processing and financial forecasting because of its capacity to improve class separability and interpretability [9][10]. The revolutionary power of LDA in promoting data-driven decision-making and streamlining machine learning operations is highlighted in this research [1][10].

IV. METHODOLOGY

In this work, Fig-3, LightGBM and Recursive Feature Elimination (RFE) to accomplish feature reduction for an IoT intrusion detection dataset. Data preparation, model training, repeated feature removal, and model assessment

are the main processes in our methodology. Every stage aims to improve the intrusion detection system's effectiveness and performance.

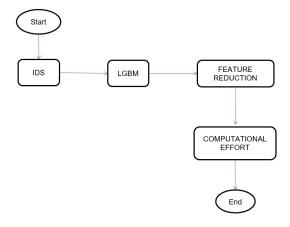


Fig-3 Workflow

1. Data Preprocessing and Exploration

- **1.1. Data Collection and Cleaning:** This Paper begins by collecting the IoT intrusion detection dataset and performing essential preprocessing steps. This includes addressing missing values, removing duplicates, and correcting data inconsistencies. Proper data preprocessing is crucial for ensuring the dataset's quality and suitability for model training [1].
- **1.2. Data Understanding:** Exploratory Data Analysis (EDA) is conducted to understand the dataset's structure and characteristics. This involves generating summary statistics, visualizing feature distributions, and analyzing feature correlations. EDA helps identify potential data issues and provides insights into the dataset's underlying patterns, setting the stage for effective feature selection [2].

2. Imparting Feature Importance to LightGBM

- **2.1. Training of the LightGBM Model:** This paper uses all of the characteristics in the dataset to build a LightGBM model. LightGBM is selected because to its efficaciousness and efficiency in managing extensive datasets and intricate feature interactions [3]. In order to determine which characteristics in the dataset have the most influence, this model computes feature significance scores.
- **2.2. Analysis of Feature Importance:** Based on the LightGBM model's feature importance scores, the features that most strongly influence the model's predictions are identified. The feature selection procedure has to be guided by the results of this study [4].

3. Implementing Recursive Feature Elimination (RFE)

3.1. Initial RFE Implementation: Recursive Feature Elimination (RFE) is applied by iteratively training the LightGBM model and removing the least important features

based on their importance scores. RFE helps in refining the feature set, focusing on the most significant features while removing those with lesser impact [5].

3.2. Iterative Feature Elimination: The RFE process is repeated iteratively. In each iteration, the LightGBM model is retrained, and the least important features are removed based on updated importance scores. This iterative approach ensures that only the most relevant features are retained, optimizing the model's performance and reducing dimensionality [6].

4.Implementing Principal Computational Analysis(PCA)

- **4.1. Initial PCA Implementation:** Principal Component Analysis (PCA) is used to reduce the dimensionality of the dataset by transforming features into principal components that capture the maximum variance:
 - Standardization: Standardize the data to have a mean of zero and variance of one [4].
 - Transformation: Apply PCA to project the data onto a lower-dimensional space [6].
 - Selection: Choose the number of principal components that retain the desired amount of variance [1].
- **4.2. Variance Preservation:** PCA's process involves preserving as much variance as possible:
 - Variance Calculation: Assess the explained variance ratio of each principal component [5].
 - Component Selection: Select principal components based on their ability to explain a significant portion of the total variance [2].
 - Dimensionality Reduction: Reduce the feature set to the selected principal components to simplify the

| Algorithms | Time | Accuracy |
|------------|------|----------|
| | | |
| LGBM | = | 99.99% |
| LGBM+ RFE | 68s | 99.99% |
| LGBM+PCA | 130s | 75% |
| LGBM+LDA | 111s | 99.99% |

model while maintaining data variability [3].

5. Implementing Linear Discriminant Analysis (LDA)

- **5.1. Initial LDA Implementation:** Linear Discriminant Analysis (LDA) reduces dimensionality by maximizing class separability:
 - Standardization: If necessary, standardize the data before applying LDA [7].

- Transformation: Apply LDA to project the data into a lower-dimensional space that enhances class separability [10].
- Component Selection: Determine the number of components based on their ability to separate classes effectively [8].
- **5.2. Class Separability Enhancement:** LDA's process focuses on improving class distinction:
 - Separation Maximization: Calculate the discriminant components that best separate different classes [9].
 - Dimensionality Reduction: Reduce the feature set based on the selected discriminant components to improve classification performance [4].
 - Evaluation: Assess the effectiveness of the reduced dimensionality in enhancing class separability and overall model performance [6].

6. Interpretation of Selected Features

Feature Importance Interpretation: Finally, interpret the selected features to gain insights into their relevance for intrusion detection. Understanding the most important features provides valuable information about the intrusion detection mechanisms and enhances the interpretability of the model [10].

To implement this project, This paper followed these steps

- Split your dataset into training and testing sets. The training set is used to train your feature reduction techniques and models, while the testing set is used to evaluate their performance.
- Perform RFE, PCA, LDA on the features present in dataset for the feature reduction process.
- Train your machine learning models (e.g. LGBM) using the reduced feature sets from each technique.
 - Based on the evaluation results, select the feature reduction technique that yields the highest model performance or the most suitable outcome for your application.

V. RESULTS

Table- 2 Results Comparision

Table 2 presents the comparative performance of different feature reduction techniques applied to the Intrusion Detection System (IDS) using the LightGBM model. The results highlight the trade-offs between computational efficiency and classification performance.

The baseline LightGBM model achieved an accuracy of 99.99% without any feature reduction, serving as the

reference point for comparison. When Recursive Feature Elimination (RFE) was applied to LightGBM, the model maintained its accuracy at 99.99%, demonstrating that the removal of less significant features did not degrade performance. Additionally, RFE significantly improved computational efficiency, reducing execution time to 68 seconds, making it the most optimal method among the tested techniques.

Principal Component Analysis (PCA), on the other hand, led to a notable drop in classification accuracy, reducing it to 75%. Despite PCA's ability to reduce the dimensionality of the dataset, the transformation of features into new principal components appears to have removed critical information, leading to diminished performance. Furthermore, PCA exhibited the highest computational cost, with an execution time of 130 seconds, making it the least efficient approach in terms of both performance and time.

Linear Discriminant Analysis (LDA) maintained the model's high accuracy of 99.99%, similar to RFE. However, it required 111 seconds for execution, making it computationally more expensive than RFE while offering no additional improvement in classification performance.

From the comparative analysis, it is evident from Fig - 4 ,RFE is the most effective feature selection technique for this IDS application. It successfully retains model performance while substantially reducing computational effort. LDA, although preserving accuracy, incurs a higher computational cost, whereas PCA proves to be the least suitable method due to its significant accuracy drop and high processing time.

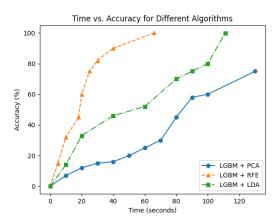


Fig - 4 Accuracy Comparison

VI. CONCLUSION AND FUTURE SCOPE

This paper sought to improve performance of a LightGBM classifier on an IoT intrusion detection dataset by employing feature reduction techniques such as Recursive Feature Elimination (RFE), Principal Component Analysis (PCA), and Linear Discriminant Analysis (LDA) . The objective was to determine the most efficient way to reduce the

number of characteristics while maintaining or enhancing classification accuracy. This paper tested various feature reduction methods to optimize feature selection, taking into account accuracy, interpretability, and computing efficiency. RFE was used to repeatedly remove less significant features [4][5], PCA was used to convert the data into a lower-dimensional space while conserving variance [1][6], and LDA was used to maximise class separability [2][7][8][9]. By selecting best out of these approaches with LightGBM, our work aims to balance performance and efficiency for improved IoT intrusion detection [10].

Further exploration could involve researching combined methods that combine RFE, PCA, and LDA to maximize their complimentary capabilities. Testing these strategies on bigger and different IoT datasets to ensure scalability and generalizability. Investigate other measures like recall, and F1-score to acquire a better understanding of model performance.

VII. REFERENCES

- [1] Y. Zhou, G. Cheng, S. Jiang, and M. Dai, "Building an efficient intrusion detection system based on feature selection and ensemble classier, Comput. Netw., vol. 174, Jun. 2020, Art. no. 107247.
- [2] Verma and V. Ranga, "Machine learning based intrusion detection systems for IoT applications," Wireless Pers. Commun., vol. 111, no. 4, pp. 2287– 2310, Apr. 2020.
- [3] P. Hadem, D. K. Saikia, and S. Moulik, "An SDN-based intrusion detection system using SVM with selective logging for IP traceback," Comput. Netw., vol. 191, May 2021, Art. no. 108015.
- [4] P. Kumar, G. P. Gupta, and R. Tripathi, "Toward design of an intelligent cyber attack detection system using hybrid feature reduced approach for IoT networks," Arabian J. Sci. Eng., vol. 46, no. 4, pp. 3749–3778, Apr. 2021.K.
- [5] D. Devprasad, S. Ramanujam, and S. B. Rajendran, "Context adaptive ensemble classi □cation mechanism with multi-criteria decision making for network intrusion detection," Concurrency Comput., Pract. Exper., vol. 34, no. 21, pp. 1–12, Jun. 2022.
- [6] N. Gupta, V. Jindal, and P. Bedi, "CSE-IDS: Using cost-sensitive deep learning and ensemble algorithms to handle class imbalance in networkbased intrusion detection systems," Comput. Secur., vol. 112, Jan. 2022, Art. no. 102499
- [7] S. H. Twe, Y. M. Thant, and M. M. S. Thwin, "Botnets attack detection using machine learning approach for IoT environment," J. Phys., Conf. Ser., vol. 1646, no. 1, pp. 1–8, 2020.

- [8] W. Xu and Y. Fan, "Intrusion detection systems based on logarithmic autoencoder and XGBoost," Secur. Commun. Netw., vol. 2022, pp. 1–8, Apr. 2022
- [9] T.-T.-H. Le, Y. E. Oktian, and H. Kim, 'XGBoost for imbalanced multiclass classification-based industrial Internet of Things intrusion detection systems,' Sustainability, vol. 14, no. 14, pp. 8707, Jul. 2022
- [10] Awajan, "A novel deep learning-based intrusion detection system for IoT networks," Computers, vol. 12, no. 2, pp. 34, Feb. 2023.
- [11] W. Xu and Y. Fan, "Intrusion detection systems based on logarithmic autoencoder and XGBoost," Secur. Commun. Netw., vol. 2022, pp. 1–8, Apr. 2022.
- [12] T.-T.-H. Le, Y. E. Oktian, and H. Kim, "XGBoost for imbalanced multiclass classification-based industrial Internet of Things intrusion detection systems," Sustainability, vol. 14, no. 14, pp. 8707, Jul. 2022.
- [13] M. Nakip and E. Gelenbe, "Online Self-Supervised Deep Learning for Intrusion Detection Systems," IEEE Transactions on Information Forensics and Security, vol. 19, pp. 5668–5683, 2024.

- [14] J. Wu, H. Dai, K. B. Kent, J. Yen, C. Xu, and Y. Wang, "Open Set Dandelion Network for IoT Intrusion Detection," arXiv preprint arXiv:2311.11249, 2023.
- [15] Yazdinejad, H. Haddadpajouh, A. Dehghantanha, R. M. Parizi, and G. Srivastava, "Cryptocurrency malware hunting: A deep Recurrent Neural Network approach," Applied Soft Computing, vol. 96, Art. no. 106658, Nov. 2020.
- [16] P. Fröhlich, E. Gelenbe, J. Fiołka, J. Chęciński, and M. Nowak, "Smart SDN Management of Fog Services to Optimize QoS and Energy," Sensors, vol. 21, no. 9, pp. 1–20, 2021.