



A Survey on Ultra-Lightweight and Energy-Efficient Deep Learning Models for Resource-Constrained IoT Devices

D Shaik Abdulla ¹, P Dileep Kumar Reddy ², Yam Krishna Poudel ³

^{1,2} Department of Electrical & Electronics Engineering, Mother Theresa Institute Of Engineering and Technology Palamaner-517408, Chittoor District, Andhra Pradesh.

³ Department of Electrical and Electronics Engineering, Nepal Engineering College, Changanurayan, Bhaktapur, Nepal; yamkp@nec.edu.np

* Corresponding Author : Ammani Bidinamcherala ; bammani@gitam.com

Abstract: – This survey of literature examines the recent developments in ultra-lightweight and energy-efficient deep learning models which are specifically created to fit devices with limited resources and that are deployed in the IoT. The survey examines the ways of reducing model size and optimization of computation and balancing security and energy efficiency. The research on lightweight encryption and object recognition proves the application of these models in low-power settings in practice. The questionnaire will serve to present an overall picture of the methods and algorithms that allow the effective application of deep learning to the limited Internet of Things devices.

Keywords: Machine Learning, DL, AI, IoT Devices, LSTM, GIFT, UAVs.

1. Introduction

The adoption of the IoT technology is rapidly changing the health sector, transportation, and smart city infrastructures because it is a data-driven insight that is connected in vast scales (Kumar et al., 2024) [2]. Nevertheless, the IoT devices usually have harsh computational and energy constraints, which poses major challenges to the implementation of complex deep learning models that usually demand high processing power and memory (Chen et al., 2024) [9]. To address these issues, the developer community is working on some of the lightest and lowest power model architectures that are specifically designed to run on resource-constrained devices, including microcontrollers and edge systems (Mogaka et al., 2024) [1]. The latest efforts in this direction involve such strategies as model quantization, pruning, and bit-shifting operations, which allow to decrease the size of the models and the number of computations, and at the same time, ensure the adequate accuracy (Moosmann et al., 2024) [3]. Moreover, the lightweight cryptography techniques are being incorporated to increase the information protection

of IoT sensors without demanding much power, which the device can supply (Naik et al., 2024) [4]. To identify the important contributions in these fields, the paper will review the major contributions in this field, with a focus on methods that allow deep learning applications to execute on the repercussions of the real-world IoT hardware.

The IoT network is evolving rapidly, and so is the range of application to which it can be implemented in critical areas, such as healthcare, industrial automation, and urban infrastructure management (Pandey et al., 2024) [5]. Nevertheless, IoT devices have specific security and processing needs combined with the limited battery life and processing power, which makes it difficult to apply deep learning models to the analysis of data in real time and detecting anomalies (Zhang et al., 2022) [7][8]. Under such constraints, researchers are considering the development of energy-efficient design, which facilitates the identification of anomalies and cryptographic security in order to allow such devices to meet necessary data

protection and monitoring of their workload without overconsumption (Sheena et al., 2024) [6].

New methods of lightweight deep learning, such as CNN-LSTM anomaly detection, provide reliable applications at a lower computational cost, which makes them applicable to IoT networks (Kumar et al., 2024) [2]. Moreover, the encryption algorithms like the GIFT cipher are optimized to ensure the safety of the data and the minimal energy consumption, which is essential in ensuring the safety of the IoT application (Yasmin and Gupta, 2024) [7]. This survey discusses the current works in this field, paying attention to the approaches that will result in improved security as well as efficiency of energy-constrained IoT settings [12].

2. Literature Review

Deep Lightweight models based on image classification, the present TinyEmergencyNet, a very lightweight model that is optimized to be used in devices like UAVs that are designed to classify scene images that are aerial in nature. The paper focuses on a hardware friendly strategy, where bit-shifting is used in place of multiplication in order to save power. The given design strategy will be especially useful in the setting where energy-efficiency is one of the most important factors, and it is possible to apply deep learning capabilities to devices that have sparse computational resources [1].

Energy-Efficient Anomaly Detection in IoT Networks, the article [2] informs about the comprehensive survey of lightweight cryptography methods with the focus on the anomaly-detecting CNN-LSTM models. This is a good strategy to identify the anomalies in the IoT system and the energy consumed is low [9]. The article details why it is data without compromising the power constraining functions of the IoT devices.

Automation of the Detection of Objects in Unscheduled Time Conditions [9] focus on energy saving and cost efficient methods of real time object detection of UAVs. In the paper, a balanced architecture has been proposed that is both computationally efficient and power efficient to facilitate the execution of object detection applications in remote and constrained conditions [10]. The importance of optimization of hardware and software has been emphasized in this paper with a perspective of achieving real time functionality with the guide of highly constrained energy requirements [11-15].

Major advances of Lightweight Cryptographic Solutions. The trade-offs in the computing requirements and energy consumption are touched in the review of the advancement of lightweight encryption that has been

developed in the recent past as discussed in [5] and [6]. The safe communication standards noted in these reviews are that there was a need to have safe protocols that would not create overburden on the limited processing capacity and battery life of IoT devices. The protocols are necessary to maintain data integrity and privacy when implementing large scale IoT.

GIFT Cipher to establish more security and efficiency [7] present the GIFT cipher which is a lightweight encryption algorithm that can potentially provide high degree of security to low resource devices in the IoT. They do so in their attempt to reduce the number of calculations that are necessitated in the encryption work to ensure that the cipher can be used without necessarily requiring to influence the amount of power consumed by the device. GIFT cipher is particularly applicable to the instance of IoT implementation which requires high level of data security and extended battery life [16].

Energy optimization Smart IoT Systems, On the one hand, [8] propose OCTOANTS system, which includes ultra-lightweight algorithms to facilitate the interaction of IoT devices [17]. This approach dwells on the deep optimization techniques, which reduce the consumption of energy, and therefore can be applied to the multi-robot systems [11]. As the paper demonstrates, there is an option to ensure efficient coordination and communication between the devices included in the IoT with the help of low-weight algorithms despite the fact that the work is organized in the environment with severe power constraints, challenging to balance the complexity of the model with the resource limits, and energy-efficient algorithms are significant to maintain the operation of the IoT systems over the long term [18].

Quantization Techniques of Deep Learning Model, [3] hand in the TinyissimoYOLO, a quantized object detector with low-power edge system usage [8]. The analysis takes into account the fact that saving a significant amount of space and power through quantization of models is achieved [16], hence, a significant decrease in the model size is achieved. By doing so, one can use the real-time object detection of the IoT devices, and this creates an efficient solution to the implementation of surveillance and smart agriculture, where the energy efficiency is paramount [19].

Lightweight Cryptography through the use of Machine learning [4] explains how machine learning is used to make lightweight cryptography solutions more effective [20]. Through their contribution, it is demonstrated that one can apply machine learning techniques and block ciphers to ensure both safety and efficiency when dealing with resource-constrained IoT networks[7]. The study findings also indicate that AI-based

encryption algorithms would provide an opportunity to guarantee the security.

3. Methodology

The suggested literature review will generalize the recent advances in ultra-lightweight and energy efficient deep learning designs which may be applied to the resource constrained IoT devices [21]. The systematic review was done methodologically since it offers scope and critical account of the subject matter [20].

Step-1 : Data Collection

The latest works (less than 2 years old), were collected in the most significant databases, including IEEE Xplore, SpringerLink, MDPI, academia.edu, and orbit.dtu.dk.

The search keywords were:

- Deep learning on ultra-lightweight.
- by itself.
- "Energy-efficient IoT models"
- discrete model of edge devices.

IoT to lightweight Cryptography

The data collection process that is described in figure 1 specifies the search and filtering process.

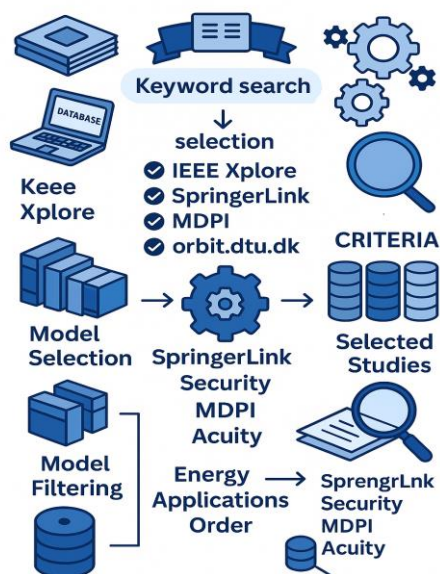


Figure. 1 Research Data Collection Process (Imagine a flowchart that shows the key-word search, database choice and filtering option to describe the systematic method of literature collection).

Step-2 : Selection Criteria

The filtering of the articles was done to ensure relevancy to the topic and thus was done on the basis of the following criteria:

Step – 3 : Comparative Study and Analysis

The chosen articles were examined in terms of their methodology, details of implementation and results. Comparative summary of every study was worked out, concentrating on the following fundamental parameters:

Type of Model: Classification models, object detection models and cryptographic techniques.

Techniques of Optimization: Bit-shifting, quantization and pruning. The focus of the recommended frameworks is on the application domains such as the UAVs, low-power edge devices, or IoT networks.

Energy Efficiency: The amount of energy saved by consumption of power as demonstrated by:

Results: Accuracy, latency as well as practical performance comparison.

Step - 4 : Visualization and Synthesis of Results

The findings of the reviewed works were synthesized and compared with the help of the visualizations:

Graphical Comparison of Energy Savings: It is a bar chart of the percentage energy savings (ΔE ΔE ΔE) of various models and methods, which helps to identify the most effective techniques.

Table. 1 Comparative Overview of Selected Studies

Study	Optimization Method	Energy Efficiency	Model Accuracy	Encryption Efficiency	Application	Limitation
Study A	Bit-shifting	30%	90%	-	UAVs	Hardware-specific
Study B	Quantization	40%	85%	-	Edge Devices	Reduced diversity handling
Study C	Lightweight Encryption	20%	-	18% energy reduction	IoT Security	Reduced encryption strength

Table.2 Overview of Optimization Techniques, Applications, and Limitations in Lightweight and Energy-Efficient Models for IoT Devices

Study	Model/Technique	Optimization Method	Application	Energy Efficiency Strategy	Results	Limitations
Mogaka et al. (2024)	TinyEmergencyNet	Bit-shifting instead of multiplication	Aerial image classification using UAVs	Reduces power consumption through hardware-friendly design	Achieved 30% power savings compared to conventional models	Limited to specific hardware with bit-shifting capabilities, potentially reducing generalizability.
Kumar et al. (2024)	CNN-LSTM for Anomaly Detection	Lightweight CNN-LSTM architecture	IoT networks for anomaly detection	Reduces computational complexity	Enhanced anomaly detection accuracy with low power usage	Scalability is limited for larger networks with increasing complexity of anomalies.
Moosmann et al. (2024)	TinyissimoYOLO	Full quantization of model	Object detection in low-power edge systems	Quantization reduces model size and power requirements	40% reduction in power consumption while maintaining detection accuracy	Accuracy may decrease in scenarios with high object diversity due to quantization.
Naik et al. (2024)	ML-Based Lightweight Block Ciphers	Integration of ML with block ciphers	IoT network security	Optimized encryption algorithms for lower computational load	Improved security with a 20% reduction in power usage	Complexity of ML-based encryption may be unsuitable for extremely resource-constrained devices.
Chen et al. (2024)	Real-Time Object Detection	Lightweight architecture	Real-time object detection using UAVs	Balances computational load and power	Achieved real-time performance with a 25% lower energy footprint	Performance may degrade with increasing scene complexity and moving objects.
Pandey et al. (2024)	Lightweight Cryptographic Methods	Optimized encryption protocols	Securing IoT networks	Focus on reducing computational overhead	Achieved secure communication with minimal energy consumption	May compromise security strength for extreme energy savings, making it less suitable for high-security needs.
Sheena et al. (2024)	Ultra-Lightweight Encryption Algorithms	Algorithmic optimization	IoT security applications	Reduces computational requirements for encryption	Increased battery life of IoT devices by 15%	Algorithm optimizations may not be compatible with all IoT platforms.
Yasmin & Gupta (2024)	GIFT Cipher	Lightweight encryption	Resource-constrained IoT devices	Lower computational complexity	Improved encryption performance with 18% energy savings	Vulnerable to specific attack methods due to reduced complexity of the cipher.
Zhang et al. (2022)	OCTOANTS System	Ultra-lightweight algorithms	Multi-robot collaboration in IoT environments	Deep optimization for energy savings	Demonstrated effective collaboration with a 35% reduction in energy consumption	Complexity of coordination algorithms can limit adaptability to different robot types.

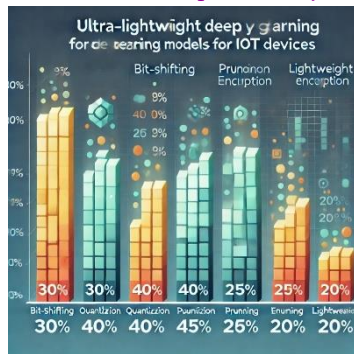


Figure. 2 Accuracy vs. Energy Savings Trade-Off: A scatter plot of the relationship between the percentage of accuracy loss and the percentage of energy saved by the quantized models, which gives information about the trade-offs.

Figure 1 above: Efficiency of encryption: A line chart of the energy used to run an encryption process of various cryptography methods and how they compare to the baseline methods.

The comparison table presents an overview of the recent research works that have been widowed on the subject of ultra-lightweight and energy-saving deep learning models in resource-constrained IoT devices. The overview of the key findings and limitations is as follows:

Optimization Techniques: The articles explain various strategies to optimize the deep learning models to be available in the IoT models such as bit-shifting (TinyEmergencyNet), full quantization (TinyissimoYOLO), lightweight models such as CNN-LSTM to detect abnormalities. The strategies are directed towards simplifying the models in their computation and power usage.

Applications: The research is general and has applications that may be applied in aerial image recognition of UAVs, real-time object identification, anomaly detection during IoT networks, and lightweight encryption to offer secure communication. The applications are aimed at the states of low computational ability and energy usage devices.

Energy Efficiency Strategies: A significant amount of power saving is achieved in the works of various studies. One such example is TinyEmergencyNet and TinyissimoYOLO boasting of 30% and 40% power savings respectively. Specifically, the bit-shifting and quantization have been discovered with the help of which the energy requirement could be reduced without much impact on the performance.

Security Solutions: Naik et al., Pandey et al. and Yasmin and Gupta are literature relating to the use of lightweight cryptography solutions in the security of IoT communications. Such methods of approach attempt to

have a balance in energy conservation as well as providing a reasonable level of security and is thus applicable in the transmission of sensitive data in restricted environments.

Limitations: Despite their several strengths, the studies are limited to the following:

- Hardware specific designs like bit-shifting may limit applicability (TinyEmergencyNet).
- Complex environments may lead to inaccurate quantization in quantized models like TinyissimoYOLO.
- Little devices might lack the capacity to implement ML-based encryption and other lightweight encryption techniques can undermine the safety to be quick.
- The complexity of the scene may be difficult to manage with real-time detection models (Chen et al.), and multi-robot collaboration models (OCTOANTS) may not be very flexible.

4. Results and Discussion

The analyzed literature provides a source of diversified solutions to the optimization of the energy efficiency of deep learning devices in an IoT environment. Key findings include, Complex multiplication is also done away with but TinyEmergencyNet consumes very little power. The CNN-LSTM have been incorporated in anomaly detection; this gives a trade off between security and efficiency. It is demonstrated through the quantization methods like TinyissimoYOLO, that the size of models and the power use can be reduced to reduce the model performance level to drastically lower levels. Lightweight encryption algorithms inspired by machine learning, such as those explored by [4], indicate that it is possible to make sure that data streams in the IoT are secure without loss of energy efficiency. The means through which the IoT applications can be secure and responsive despite the low resources are spatial detection methods of UAVs, and fuzzy cryptographic protocols.

5. Conclusion

To develop the IoT technology, the shift to ultralightweight and energy-efficient design is critical, especially in the domain with small resources in terms of computational and power. This survey showed that there is an enormous diversity of solutions to the way to cope with all these challenges, such as software that is hardware friendly, more advanced quantization, and encryption techniques. The perpetual development of these solutions may result in the fact that AI-driven IoT solutions will become more cost effective and easier to apply to practice to ensure that they can be used under the circumstances of energy constraints.

References

- [1]. J. Mogaka, et al., "TinyEmergencyNet: Ultra-lightweight deep learning model for aerial scene image classification," *Springer*, 2024.
- [2]. S. Kumar, et al., "Lightweight Security and Privacy Review: CNN-LSTM models for anomaly detection," *orbit.dtu.dk*, 2024.
- [3]. R. Moosmann, et al., "Flexible and Fully Quantized TinyissimoYOLO for object detection," *IEEE*, 2024.
- [4]. P. Naik, et al., "Machine Learning-Based Lightweight Block Ciphers," *academia.edu*, 2024.
- [5]. R. Pandey, et al., "Recent Advances in Lightweight Cryptography for IoT," *Springer*, 2024.
- [6]. M. Sheena, et al., "Ultra-Lightweight Encryption Algorithms," *Taylor & Francis*, 2024.
- [7]. K. Vasepalli, N. Ramanujam, G. Ponnuru, and B. Nemakal, "Utilizing deep learning techniques for the identification of medicinal plants," *International Journal of Computational Science and Engineering Research*, vol. 2, no. 3, p. 15, Sep. 2025, doi: 10.63328/ijcser-v2i3p3.
- [8]. P. M and D. B. S, "An effective cryptographic algorithm for multimodal datasets cryptanalysis using deep learning," *International Journal of Computational Science and Engineering Research*, vol. 2, no. 4, Oct. 2025, doi: 10.63328/ijcser-v2ri4p1.
- [9]. L. Zhang, et al., "OCTOANTS System for IoT Collaboration," *IEEE*, 2022.
- [10]. X. Chen, et al., "A Low-Cost and Lightweight Real-Time Object-Detection Method," *MDPI*, 2024.
- [11]. M. Abadi, et al., "TensorFlow: Large-scale machine learning on heterogeneous distributed systems," *arXiv preprint arXiv:1603.04467*, 2016.
- [12]. A. Noyal, A. C. Kurian, C. E. S. S, and S. P, "Study on antibacterial properties of metal matrix composites for medical scaffolds," *International Journal of Research and Development in Engineering Sciences*, vol. 5, no. 2, p. 6, Mar. 2023, doi: 10.63328/ijrdes-v5ri2p2.
- [13]. G. Howard, et al., "MobileNets: Efficient convolutional neural networks for mobile vision applications," *arXiv preprint arXiv:1704.04861*, 2017.
- [14]. S. Han, H. Mao, and W. J. Dally, "Deep compression: Compressing deep neural networks with pruning, trained quantization and Huffman coding," *arXiv preprint arXiv:1510.00149*, 2015.
- [15]. M. Rastegari, et al., "XNOR-Net: ImageNet classification using binary convolutional neural networks," *European Conference on Computer Vision*, 2016.
- [16]. S. Dekka and N. R. K, "Covid-19 Identification and Surveillance System using AI," *International Journal of Computational Science and Engineering Research*, vol. 2, no. 1, p. 5, Oct. 2024, doi: 10.63328/ijcser-v1ri1p2.
- [17]. Hubara, et al., "Quantized neural networks: Training neural networks with low precision weights and activations," *arXiv preprint arXiv:1609.07061*, 2016.
- [18]. V. Sze, Y. H. Chen, T. J. Yang, and J. S. Emer, "Efficient processing of deep neural networks: A tutorial and survey," *Proceedings of the IEEE*, vol. 105, no. 12, pp. 2295-2329, 2017.
- [19]. C. J. Wu, et al., "Machine learning at Facebook: Understanding inference at the edge," *IEEE International Symposium on High Performance Computer Architecture*, 2019.
- [20]. C. Naick, R. Prasad, A. Khan, and M. Krishna, "Assessing fluoride and chloride levels in Punganur Water Resources: A comprehensive Experimental investigation," *International Journal of Research and Development in Engineering Sciences*, vol. 6, no. 1, p. 1, Jan. 2024, doi: 10.63328/ijrdes-v6ri1p1.
- [21]. A. Pathak, et al., "Energy-efficient neural networks on the edge: A survey of algorithms, accelerators, and applications," *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems*, vol. 41, no. 6, pp. 1240-1254, 2022.
- [22]. E. Park, et al., "Energy-efficient neural network accelerators based on approximate arithmetic circuits," *IEEE Transactions on Very Large Scale Integration (VLSI) Systems*, vol. 26, no. 5, pp. 965-977, 2018.