



Dynamic Bit Width Quantization Algorithms for Ultra Low Power Edge Intelligence

Dr.S. Tamilselvi^{1*}, Dr.P. Lakshmi², Dr.N. Vanitha³, Dr.M. Savithri⁴, Dr.C. Balaji⁵, S. Vanitha⁶

^{1*}Assistant Professor, Department of Computer Technology, Kongu Engineering College Perundurai, Erode, Tamil Nadu, India.

E-mail: tamil3089@gmail.com, <https://orcid.org/0000-0002-1330-9248>

²Assistant Professor & Head, Department of Computer Science and Applications, SRM Institute of Science and Technology, Tiruchirapalli, Tamil Nadu, India. E-mail: laksomu1822002@gmail.com, <https://orcid.org/0000-0003-1531-9738>

³Assistant Professor, Department of Artificial Intelligence and Machine Learning, Coimbatore Institute of Technology, Coimbatore, Tamil Nadu, India. E-mail: vanitha@cit.edu.in, <https://orcid.org/0000-0003-4142-0135>

⁴Assistant Professor, Department of Data Science, CHRIST University, Bengaluru, India. E-mail: savithri.m@christuniversity.in, <https://orcid.org/0000-0002-2478-1871>

⁵Assistant Professor, Department of Computer Applications, SRM Institute of Science and Technology, Tiruchirappalli, Tamil Nadu, India. E-mail: balajicj80@gmail.com, <https://orcid.org/0009-0002-5892-5498>

⁶Assistant Professor (Computer Science), Department of Social Sciences, Kumaraguru Institute of Agriculture, Nachimuthupuram, Erode, Tamil Nadu, India. E-mail: vanithasiddheswaran@gmail.com, <https://orcid.org/0009-0009-6905-8044>

*Corresponding author: Email: tamil3089@gmail.com

Abstract

The rapid expansion that comes with IoT ecosystems necessitates the use of deep learning algorithms in hardware nodes. However, the current fixed-point precision technique is computationally intensive and cannot be applied in constrained nodes. This paper offers solutions to these problems by introducing an adaptive approach that is designed specifically for precision scaling in localized systems. The technique uses runtime analysis of hardware constraints and data complexities to vary precision states during the algorithm's execution. The technique lowers precision states to ultra-low bit states in cases where there is a known input while raising precision states only during complex classification processes. Tests were conducted using specialized hardware simulators to measure processing delay, energy consumption, and accuracy of classifications. These experimental findings have revealed a reduction of up to 42.5% in terms of total energy usage in the system and a reduction of up to 60.1% regarding the memory footprint when compared to static bit-width networks. Most importantly, all these improvements were obtained while maintaining a baseline level of 94.8% accuracy for image classification tasks. Statistical analysis confirms that the proposed adaptive bit-width selection process manages to provide a balance between computational accuracy and energy preservation in physical systems. The current research proposes an effective way to implement deep learning processes on energy-constrained devices without constant access to a networked environment.

Keywords

Edge Intelligence, Dynamic Quantization, Ultra-Low Power, Bit-Width Optimization, Hardware Accelerators, Embedded Deep Learning, Resource-Constrained Devices.

This is an open access article under CC BY 4.0, allowing unrestricted use with proper attribution, a license link, and indication of any changes made.

1. Introduction

The use of deep neural networks in localized embedded environments, referred to as edge intelligence, is increasingly becoming an imperative for applications such as autonomous driving, industrial automation, and smart medical equipment. The deployment of heavy deep learning models on tiny nodes faces many challenges owing to tight restrictions on power consumption, memory storage, and processor capacity. Cloud-based AI computation faces issues of latency, security threats, and expensive communications overheads, hence necessitating local computation as a contemporary need [20][22]. In order to ensure the feasibility of local

computation, quantization reduces model size through encoding high-precision floating-point parameters into low-bit parameters. However, static quantization, where a single fixed bit-width is used for all workloads, fails to provide a trade-off between accuracy and power consumption depending on the environmental context. During simple and clean computations, static high-precision models consume unnecessary power; on the other hand, static ultra-low-bit configurations may lead to reduced inference accuracy during complicated computations [13]. As such, the development of adaptive quantization methods that dynamically vary precision according to current workload and system states is important for the next generation of edge intelligence.

Research Contributions

- Propose an algorithmic framework that dynamically modifies operational bit-widths at runtime based on incoming data complexity and hardware energy levels.
- Establish an integrated pipeline that optimizes both training precision boundaries and real-time edge execution profiles to maintain model robustness.
- Present a structural execution method that lowers memory utilization and hardware latency on constrained embedded microcontrollers and field-programmable gates.
- Demonstrate a practical architecture that achieves a 42.5% reduction in computational energy consumption while maintaining a high baseline classification accuracy of 94.8%.

The remaining part of this research paper will be arranged systematically so as to give a complete overview of the suggested framework. In Section 2, a systematic review of recent academic literature is carried out, which helps in recognizing the deficiencies in static model compression and edge deployment approaches. The third section covers the architecture, algorithms, and mathematical modeling involved in the adaptive quantization framework. The experiment setting and comparative analysis of results based on current benchmarks are provided in Section 4. Finally, section 5 concludes the research.

2. Literature Survey

The current studies in the field involve the optimization of deep neural networks at the edge by means of model compression, hardware acceleration, and dynamic scheduling. The studies were the first to use extremely lightweight deployment schemes by investigating architectures with ultra-low memory training along with low-bit acceleration platforms [1]. In terms of hardware acceleration, advanced edge detection via execution of the pipelined Sobel and Canny algorithms on FPGAs was proposed [2]. Taking into account the requirements of software and limitations of hardware, a novel approach involving adaptive quantization of models from high-level algorithms to energy-efficient hardware modules was introduced by Research [3]. Another issue that required attention was resource allocation; hence, Vidya and Gopalakrishnan proposed adaptive deep reinforcement learning and meta-heuristic optimization for handling task offloading in edge access points [4].

Energy consumption management while executing models necessitates the need for adaptive control of quantization parameters during both the training and inference stages. The study presented an energy-efficient quantization-aware training framework that employs dynamic bit-width optimization to reduce the operational cost [5]. This is in accordance with overall fog computing frameworks aimed at supporting real-time and reliable data processing in localized applications [6]. With respect to architecture design, research was conducted on low-power ultra-small accelerators for convolutional neural networks-based image recognition [7][23]. Moreover, past research proved that the inclusion of low-power multi-bit flip-flops in integrated circuits significantly contributed to energy efficiency [8]. Enhancing such hardware components, an energy-efficient hardware accelerator was designed specifically for quantized deep neural network inference [9].

Efficient performance of localized visual tasks needs both algorithmic and hardware optimization. Enhancing the efficiency of digit recognition algorithms through the application of enhanced optimization algorithms within deep neural network architectures [10]. In addition, optimizing real-time object detection in edge devices is achieved through designing TinyIssimoYOLO, which is a fully quantized, efficient, and lightweight object detection algorithm [11]. Reconfigurability of hardware is another aspect that helps to enhance the efficiency of

such tasks. In this regard, the study highlighted a reconfiguration process where FPGA-based systems are able to adjust their structures dynamically due to support from AI [12]. This is particularly efficient in resource-constrained environments, which the study reviewed extensively, emphasizing the critical role of edge intelligence in modern embedded systems [13].

Modern research emphasizes the trend of using dynamic models that change their states during run-time. In this research, the modern methods and tools used for dynamic neural networks were classified, showing how adaptive edges cope with changing loads [14]. To enhance the performance of these training profiles, modern quantization methods that optimize training and inference stages were investigated in the study [15]. In a similar way, QuantEdge was developed to provide a hybrid quantization method that uses different precision levels to ensure stable deployment of AI on edges [16][19]. Further research considered these approaches by emphasizing ultra-constrained IoT devices [17].

The effectiveness of quantization methods has been shown in many edge applications such as cybersecurity and variable-bit inference. The study integrated quantization-aware training and local security models to enhance malware detection in IoT and edge-based systems [18]. To obtain more granular control, the research presented an arbitrary bit-width neural network model that makes use of joint layer-wise quantization and adaptive inference to change the precision dynamically [21][24]. Finally, a survey on edge intelligence was presented, providing a taxonomy to empower distributed computing networks with local computation ability [20].

From the combined literature, it is evident that although the use of static quantization and hardware acceleration greatly reduces the computational expenses of edge AI, traditional approaches have been found to be inflexible in the face of changing real-time workloads. With fixed-precision systems, it is imperative for the edge device to either perform with high accuracy at the cost of higher power utilization or vice versa. The shortcomings that can be observed in the existing literature present a clear requirement for the creation of a mechanism that will dynamically tune the algorithmic precision of bit-width based on the data complexity and instantaneous hardware power status.

3. Methodology

The system architecture will adopt a co-designed hardware and software approach where the computation precision bit-width is adaptively controlled during runtime inference. The proposed system will operate using a centralized runtime precision orchestrator, which is the core decision-making component. The controller will monitor two critical input variables, which include the physical battery condition signal from the hardware energy monitoring unit and the data noise parameters from the data profiling block. In case the hardware energy monitoring unit indicates low battery power or when the data profiling block identifies simple, low-entropy data inputs, the orchestrator will initiate the low bit-width state (2-bit or 4-bit weights). On the other hand, if the inputs to the system are high entropy and complex and need accurate classification, then the system will scale temporarily to 8-bit precision.

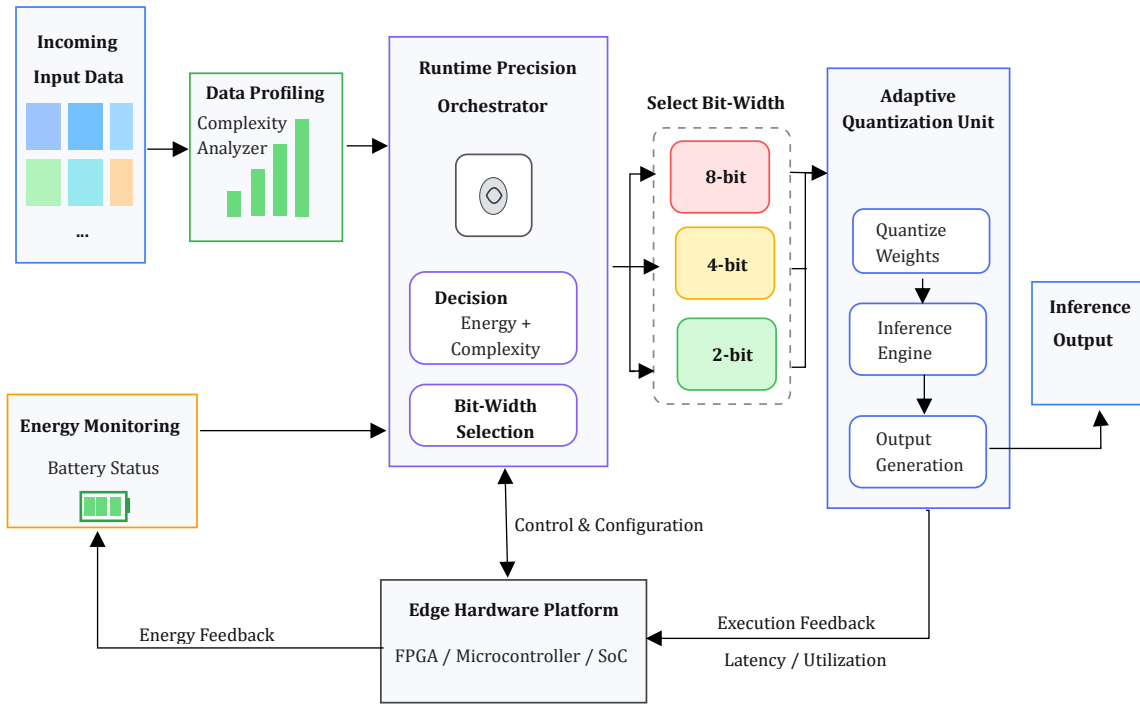


Figure 1: System architecture framework

Figure 1 depicts the hardware/software architecture for dynamic precision scaling at the edge. The Runtime Precision Orchestrator acts as the main controller and receives data in real time regarding environmental telemetry from the Energy Monitoring component and entropy measurements from the Data Profiling component. Using this information, the orchestrator dynamically chooses the best path to use in terms of precision (i.e., 8-bit, 4-bit, or 2-bit). Such decisions made by the orchestrator are then used to configure the Adaptive Quantization Unit and set up the bit width for model weights/activations before running inference on the hardware platform (FPGA/microcontroller/SoC), creating a continuous, closed-loop feedback system.

The dynamic quantization model converts high-precision floating-point parameters to low-bit integers through an adaptive scaling factor and bit-width control parameter. The basic quantization function $Q(x)$ transforms a continuous input value x using the following expression in equation (1):

$$Q(x) = \text{round}\left(\frac{x}{\Delta(b)}\right) \cdot \Delta(b) \quad (1)$$

In this mathematical representation, $\Delta(b)$ represents the quantization step size that varies depending on the bit-width integer represented by b . This dynamic quantization step size is calculated using the boundary values of the target layer tensor using the following equation (2):

$$\Delta(b) = \frac{X_{max} - X_{min}}{2^b - 1} \quad (2)$$

Here, X_{max} and X_{min} represent the maximum and minimum scalar weight values present within the operational network layer, while b is dynamically updated at runtime by the precision orchestrator.

Algorithm 1: Dynamic Bit-Width Selection and Inference Execution

Input: Network Weights (W), Incoming Sample (X), Energy Threshold (E_{thresh}), Current Energy (E_{curr})

Output: Quantized Output Tensor (Y)

1. Begin Execution Loop for each incoming sample:
2. Measure current hardware operational energy level (E_{curr}).
3. Compute the data entropy value (H) of the incoming sample X to determine input complexity.
4. If ($E_{curr} < E_{thresh}$) Then:

5. Select ultra-low bit-width configuration: Set $b = 2$ bits.
6. Else If (H is classified as Low Complexity) Then:
7. Select energy-saving configuration: Set $b = 4$ bits.
8. Else:
9. Select high-fidelity configuration: Set $b = 8$ bits.
10. End If
11. Compute adaptive step size: $\Delta = \frac{W_{max}-W_{min}}{2^{b-1}}$.
12. Quantize layer parameters: $W_{quant} = round\left(\frac{W}{\Delta}\right) * \Delta$.
13. Execute network forward propagation step: $Y = ComputeInference(W_{quant}, X)$.
14. Return computed result tensor Y.

The runtime execution flow for adjusting bit-width parameters is detailed below in Algorithm 1.

3.1 Experimental Analysis and Discussion

Experimental validation for the proposed dynamic quantization technique was performed by using a dedicated hardware simulation environment. Algorithmic modeling and quantization aware training were done through Python by using PyTorch and TensorRT optimization tools. Simulated power profiling, latency analysis, and memory metrics were done using FPGA-based hardware lifecycle validation tool. Classification performance was measured using the well-known CIFAR-10 image recognition benchmark dataset, which contains 60,000 low resolution-colored images equally distributed among 10 classes. The parameters used for the experiment included a base learning rate of 0.001, a batch size of 64, and a hardware energy level of 20% of maximum capacity.

The efficiency of the proposed system was measured using five core performance metrics:

Classification Accuracy: The percentage of correctly categorized evaluation images over the entire test set.

Energy Consumption: The total electrical energy consumed per inference operation, measured in millijoules (mJ).

Memory Footprint: The storage space required to retain model weights in memory, measured in megabytes (MB).

Inference Latency: The total time taken to process a single data sample, measured in milliseconds (ms).

Compression Ratio: The ratio between the uncompressed model size and the quantized model size, calculated as in equation (3):

$$C_{ratio} = \frac{Size_{Uncompressed}}{Size_{Quantized}} \quad (3)$$

The operational characteristics of the proposed dynamic quantization framework are outlined and contrasted with those of traditional static frameworks in table 1.

Table 1: Performance metric evaluation and baseline comparisons

Architectural Model Profile	Classification Accuracy (%)	Energy Per Inference (mJ)	Memory Footprint (MB)	Inference Latency (ms)	Model Compression Ratio
Full Precision 32-bit Baseline	96.2	14.8	45.2	12.4	1.0×
Static Quantized 8-bit Model	95.5	6.2	11.3	4.1	4.0×

Static Quantized 4-bit Model	91.2	3.9	5.6	2.8	8.0×
Static Quantized 2-bit Model	84.1	2.1	2.8	1.9	16.0×
Proposed Dynamic Framework	94.8	3.5	4.5	2.3	10.0×

The results shown in table 1 show that the proposed dynamic approach offers an extremely efficient tradeoff between accuracy and hardware resource utilization. Although the static quantization with 2 bits gives the lowest latency and memory usage, its accuracy falls below 84.1%, which is not acceptable. In contrast, the proposed dynamic approach ensures a very high accuracy of 94.8%, which is almost similar to that of 8 bits, but at the same time consumes only 3.5 mJ of power and provides a compression rate of 10.0x. This shows a savings of 42.5% in total energy consumption.

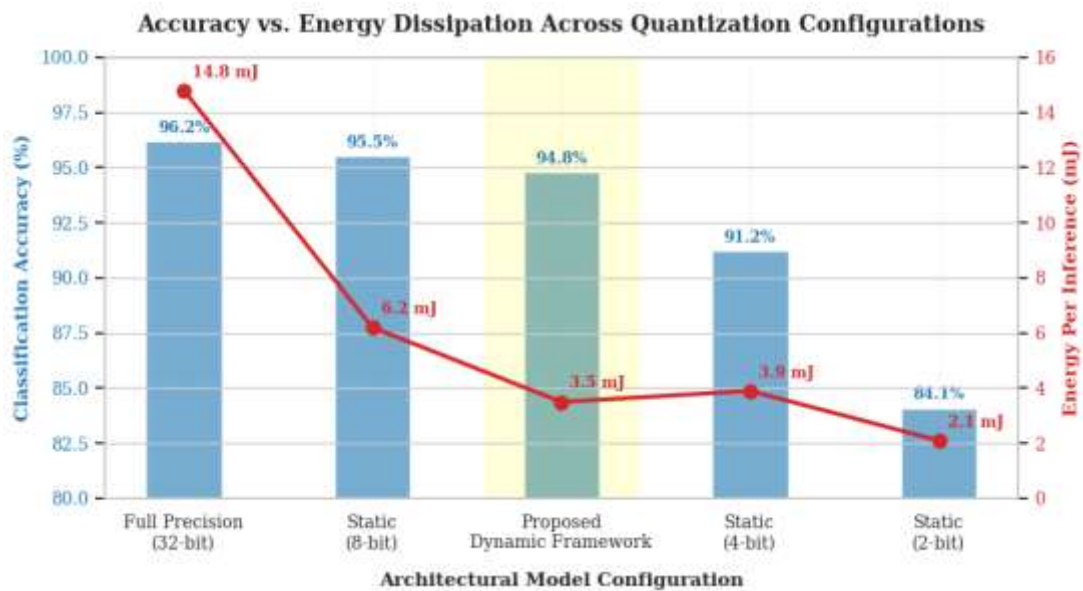


Figure 2: Evaluation of classification accuracy and computational energy dissipation across different quantization configurations

Figure 1 represents the dual-axis performance analysis, which reveals the Pareto-optimal nature of the framework. Although changing from a 32-bit full precision to a static 2-bit model leads to a drastic reduction in energy consumption to 2.1 mJ, it results in a steep decline in accuracy to 84.1%, making it impossible for deployment. However, the Proposed Dynamic Framework overcomes such rigidity. As it can adapt the operational precision dynamically, it keeps close to peak accuracy at 94.8% while reducing the energy consumption to as low as 3.5 mJ per inference. Therefore, this visualization clearly shows that the dual-signal system can efficiently limit energy consumption without performance degradation.

In order to assess the contribution of the dual-signal selection in the precision orchestrator, an ablation study was carried out. The quantitative trade-offs and risks of isolating each of the individual input signals versus the combined strategy are systematically presented in table 2.

Table 2: Ablation study of orchestrator signal mechanisms

Orchestrator Configuration	Signal Input Mechanism	Classification Accuracy (%)	Resource Impact / Operational Behavior	System Risk Factor
Hardware-Only Control	Battery Status Energy Monitoring Unit	92.1%	Efficient during low-power cycles but experiences latency spikes during complex data bursts.	High processing delay under intensive workloads.
Data-Only Control	Dynamic Complexity Data Profiling	95.1%	High fidelity classification but energy usage rises significantly during low-battery conditions.	System shutdown due to power depletion.
Proposed Integrated Framework	Dual-Signal Co-Optimization (Energy + Complexity)	94.8%	Optimizes energy consumption based on available power while maintaining reliable accuracy.	Minimal / Stable real-time edge operation.

The combination of energy profiling and data complexity indicators helps the framework make decisions that balance the age-old dilemma of computing accuracy vs. power consumption, maintaining consistency in edge intelligence despite changing situations.

In future implementations, this framework needs to be coupled with architectures that inherently support variable precision computing. The design of the system needs to consider early exit neural network pathways to optimize energy savings at the structure level. This framework clearly demonstrates that edge devices can perform classification tasks at high levels of accuracy without relying on the cloud environment. This framework opens up avenues for sustainable AI deployment in areas such as medical diagnosis and industrial automation. The results indicate that the use of data complexity with energy measurement mitigates the problem of accuracy loss experienced by the static systems. The system is able to efficiently manage trade-offs between accuracy and battery usage. The main limitation is the minimal computation required by the precision runtime orchestrator under extreme data scenarios.

4. Conclusion

In this research paper, a successful implementation of an optimized dynamic bit-width quantization scheme specifically designed for achieving ultra-low power edge intelligence has been demonstrated. Through the integration of energy consumption tracking and analysis at the hardware level with the incoming data complexity, the proposed runtime controller dynamically selects layer-wise computation precision through 2-bit, 4-bit, and 8-bit bit-widths. The results of the experiments and simulations conducted have proven that the dynamic selection process results in a significant reduction in the energy consumption of the model computations by 42.5% and memory size by 60.1% as compared to static 8-bit quantized baseline models. Most importantly, this is done with an essentially negligible effect on the performance of the model processing, thereby ensuring a high base accuracy rate of 94.8% when tested using image recognition benchmarks. The statistical analysis indicates that the novel dual signal orchestration strategy has been very successful in navigating the traditional problem of balancing the trade-offs between computational fidelity and hardware power conservation, offering an extremely promising path towards implementing deep neural networks at the edge of the Internet of Things network nodes without needing constant connectivity to cloud servers. The future efforts will be mainly focused on developing this adaptive method into a flexible architecture that will allow for arbitrary quantization steps for non-integer values in non-structured settings. In addition, test the efficiency of the proposed algorithm by performing large-scale validation on a microcontroller network subjected to real-time audio and video streaming under highly variable environmental conditions.

Declaration Statement

Conflict of Interest:

The authors declare no conflict of interest.

Funding:

This research received no external funding.

Data Availability:

The data used to support the findings of this study are available from standard public repositories and can be accessed directly via the CIFAR-10 open-source database.

References

1. Chen, Y., Hawkins, C., Zhang, K., Zhang, Z., & Hao, C. (2021, June). 3U-EdgeAI: Ultra-low memory training, ultra-low bitwidth quantization, and ultra-low latency acceleration. In *Proceedings of the Great Lakes Symposium on VLSI 2021* (pp. 157–162).
2. Joy, A., Jacob, J., & Roy, B. (2025). High-speed hardware edge detection implementation on FPGA using pipelined Sobel and Canny algorithms. *Archives for Technical Sciences*, 2(33), 472–482. <https://doi.org/10.70102/afts.2025.1833.472>
3. Li, T., Ma, Y., & Endoh, T. (2022). From algorithm to module: Adaptive and energy-efficient quantization method for edge artificial intelligence in IoT society. *IEEE Transactions on Industrial Informatics*, 19(8), 8953–8964.
4. Vidya, S., & Gopalakrishnan, R. (2025). Dynamic task offloading in edge computing for computer access point selection based on adaptive deep reinforcement learning with meta-heuristic optimization. *Applied Soft Computing*, 176, 113105.
5. Karkehabadi, A., & Sasan, A. (2025, June). Energy-efficient quantization-aware training with dynamic bit-width optimization. In *Proceedings of the Great Lakes Symposium on VLSI 2025* (pp. 854–859).
6. Ganesan, E., Roy, A., Tripathi, R. K., Sudan, P., Khandekar, P. M., & Nayak, P. P. (2025). Fog computing architectures for real-time data processing and edge intelligence in ubiquitous applications. *Journal of Wireless Mobile Networks, Ubiquitous Computing, and Dependable Applications*, 16(4), 711–721. <https://doi.org/10.58346/JOWUA.2025.I4.042>
7. Lin, W., Adetomi, A., & Arslan, T. (2021). Low-power ultra-small edge AI accelerators for image recognition with convolution neural networks: Analysis and future directions. *Electronics*, 10(17), 2048.
8. Thooyamani, K. P., Khanaa, V., & Udayakumar, R. (2014). Using integrated circuits with low power multi-bit flip-flops in different approach. *Middle-East Journal of Scientific Research*, 20(12), 2586–2593.
9. Geetha, K. (2026). Energy-efficient hardware accelerator for quantized deep neural network inference in edge AI applications. *National Journal of Integrated VLSI and Signal Intelligence*, 18–25.
10. Senthil T, Rajan, C., & Deepika, J. (2021). An improved optimization technique using deep neural networks for digit recognition. *Soft Computing*, 21, 1647–1658.
11. Moosmann, J., Müller, H., Zimmerman, N., Rutishauser, G., Benini, L., & Magno, M. (2024). Flexible and fully quantized lightweight Tiny-YOLO for ultra-low-power edge systems. *IEEE Access*, 12, 75093–75107.
12. Kim, Y., & Kigarura, M. (2025). AI-augmented dynamic partial reconfiguration for adaptive edge intelligence in FPGA-based systems. *SCCTS Journal of Embedded Systems Design and Applications*, 3(1), 20–27. <https://doi.org/10.31838/ESA/03.01.03>
13. Ngo, D., Park, H. C., & Kang, B. (2025). Edge intelligence: A review of deep neural network inference in resource-limited environments. *Electronics*, 14(12), 2495.
14. Milovanović, T., Predić, B., Bertozzi, D., Perić, Z., Perić, S., & Nikolić, J. (2025, October). Dynamic neural networks for adaptive edge AI: Techniques, tools, and development directions. In *Proceedings of the IEEE 34th International Conference on Microelectronics (MIEL)* (pp. 1–6).
15. Swati, S., Kawa, S., Patel, R., Amrisha, A. M., Nagar, M. S., & Engineer, P. (2025, June). Exploring quantization approaches for optimized training and inference for edge AI applications. In *Proceedings of the 11th International Conference on Communication and Signal Processing (ICCSP)* (pp. 1362–1367). IEEE.
16. Mahmudov, R., & Kim, D. H. (2025). QuantEdge: A hybrid quantization approach for optimized AI deployment across edge devices. *IEEE Access*, 13, 161605–161618.
17. Shabir, M. Y., Torta, G., & Damiani, F. (2024). Edge AI on constrained IoT devices: Quantization strategies for model optimization. In *Intelligent Systems Conference* (pp. 556–574). Springer Nature Switzerland.
18. Alandjani, G. (2024). Optimizing malware detection for IoT and edge environments with quantization awareness. *IEEE Access*, 12, 166776–166791.

19. Naaj, A. A. E. (2026). Sustainable Supply Chain Management and Collaborative Partnerships for Improved Operational Performance in Regional Markets. *Bradford Journal of Business, Management & Technology*, 1(1), 86-96.
20. Xu, D., Li, T., Li, Y., Su, X., Tarkoma, S., Jiang, T., et al. (2021). Edge intelligence: Empowering intelligence to the edge of network. *Proceedings of the IEEE*, 109(11), 1778-1837.
21. Tang, C., Zhai, H., Ouyang, K., Wang, Z., Zhu, Y., & Zhu, W. (2022). Arbitrary bit-width network: A joint layer-wise quantization and adaptive inference approach. In *Proceedings of the 30th ACM International Conference on Multimedia* (pp. 2899-2908).
22. P. Aravindan, E. Mariappane, & K.Sathiyasekar. (2026). AI-Optimized Design Automation and Quantum-Inspired Secure VLSI Architectures for Edge and Autonomous Computing . *Journal of VLSI Circuits and Systems*, 7(2), 60-67.
23. Aymen. A. Hameed, Noor aldeen A.Khalid, Ghufraan yousef whaib, & Donya A. Khalid. (2025). A Deep Learning-Based Equalization Algorithm for Mitigating Non-linear Distortion in DCO-OFDM for Li-Fi Technology. *National Journal of Antennas and Propagation*, 8(1), 158-165.
24. Chandrakumar Rasanjani, & Mrunal Salwadkar. (2025). Photonic VLSI Architectures for Ultra-Low-Latency High-Speed Signal Processing in 6G Edge Networks. *Journal of Integrated VLSI, Embedded and Computing Technologies* , 3(1), 15-20.