

Dynamic Power–Performance Trade-Offs in Reconfigurable Edge-AI Accelerators: Experimental Evaluation of an Energy-Proportional FPGA Platform

Rushikesh Ramesh Vikhe¹

¹ Research Scholar, Department of Electronics & Communication Engineering (VLSI),
Sri Saty Sai University of Technology and Medical Science, Sehore, Bhopal.

Dr. Devendra Patle²

² Research Guide, Department of Electronics & Communication Engineering (VLSI),
Sri Saty Sai University of Technology and Medical Science, Sehore, Bhopal.

ABSTRACT

Background / Context

The explosively growing number of edge-AI use cases in IoT, autonomous vehicles, and smart devices has strengthened the demand for low-power, real-time processing at the network edge. Traditional accelerators like GPUs and ASICs, while hugely capable, are limited by fixed configurations, high power budgets, and poor flexibility to changing workloads. Reconfigurable computing platforms like FPGAs have risen as viable alternatives that can provide on-demand performance and energy scalability for edge AI.

Problem / Gap

Current FPGA-based AI accelerators, though, usually run at fixed voltage–frequency points and are not adaptive at runtime. Such limitations become impediments to energy proportionality and make efficiency worse with fluctuating workloads — a key requirement for sustainable, high-throughput edge-AI deployment.

Aim / Objective

This research endeavors to design, implement, and assess an energy-proportional FPGA-based architecture that optimizes power and performance dynamically by applying adaptive reconfiguration and power management strategies.

Methodology / Approach

The target system combines Dynamic Voltage and Frequency Scaling (DVFS), Clock Gating, and Partial Reconfiguration in a closed-loop control strategy to support runtime adaptability. Deployed on a Xilinx Zynq UltraScale+ MPSoC, the design utilizes real-time monitoring of workloads to dynamically adjust voltage, frequency, and active processing units. Power and latency were evaluated using simulation and on-board measurement, and results were compared with static FPGA accelerators to verify theoretical models and efficiency of runtime control.

Results / Findings

Experimental outcomes exhibit a 31% power reduction, 19% latency reduction, and 35% improvement in energy efficiency (GOPS/W) over baseline static designs. Analytical and hardware outcomes indicate that the deviation in correlation is less than 6%, and so the correctness of the theoretical energy model and the efficacy of adaptive control are verified.

Implications / Significance

The results provide a platform for energy-conscious, reconfigurable hardware as an applicable solution for next-generation edge-AI systems that need sustainable and autonomous computation. The devised framework promotes green computing efforts, improves the lifetime of devices, and presents a scalable design platform for AI implementation in IoT, robotics, and embedded systems.

Keywords: *Edge-AI Accelerator; Dynamic Voltage and Frequency Scaling; Power–Performance Trade-Off; FPGA Reconfiguration; Energy Efficiency; Runtime Adaptation.*

1. Introduction

Edge computing has reshaped the deployment of artificial-intelligence (AI) models by bringing inference closer to where data originates, such as sensors, drones, and embedded controllers (Bolchini et al., 2022a). Local computation lowers latency, improves privacy, and minimizes bandwidth usage (Kaur et al., 2025). Yet these benefits are bounded by edge devices' limited power budgets and thermal envelopes (W. Wang et al., 2025). Traditional processors and accelerators are hence unsuitable because they were optimized for datacenter environments and not power-constrained ones (Yan & Li, 2024).

Traditional Graphics Processing Units (GPUs) provide enormous parallelism but run at fixed voltage–frequency points with high static power consumption and are inappropriate for continuous in-device usage (Silvano et al., 2025). Application-Specific Integrated Circuits (ASICs), including Google's Tensor Processing Unit (TPU), provide higher efficiency but are rigid once they are manufactured (Rutishauser, 2024). Neither one offers runtime flexibility in adapting to varying workloads—a requirement fundamental to edge inference, where data streams and computational intensity change over time (Kiat et al., 2020).

Field-Programmable Gate Arrays (FPGAs), on the other hand, provide hardware reconfigurability, allowing programmers to customize data paths and memory hierarchies dynamically (J. Li et al., 2021). However, conventional FPGA-based accelerators are statically configured with coarse-grained mechanisms to change power and performance dynamically (Huang et al., 2025). To bridge this shortcoming, this work introduces an Energy-Aware Reconfigurable Architecture (EARA) that combines hardware-level reconfiguration and dynamic power-management policies (Raveendran Nair, 2025).

Three motivating research questions are presented:

1. How is dynamic reconfiguration utilized to provide energy proportionality for edge-AI hardware?
2. What are the architectural mechanisms supporting adaptive power control without sacrificing performance stability?
3. To what extent can theory predict FPGA real-world behavior under workload variation?

The rest of the paper is structured below. Section 3 summarizes related literature and develops the theoretical framework (Hosseinabady & Nunez-Yanez, 2018). Section 4 introduces the research goals, hypotheses, and experimental factors. Section 5 describes the scope and limitations of the study, whereas Section 6 explains the methodology, design instruments, and validation plan. Section 7 discusses the results, followed by discussion, conclusions, and recommendations in Sections 8–10 (May et al., 2024).

2. Literature Review and Theoretical Framing

2.1 Recent Developments in Edge-AI Acceleration

From 2020 to 2025, AI accelerators based on FPGA have trended toward domain-specific design and near-sensor intelligence (Du et al., 2024). ASIC platforms like TPU v4 or Apple Neural Engine are good at performance per watt but are static (Millar et al., 2025). On the other hand, NVIDIA Jetson modules and AMD

Versal ACAPs reveal programmability via heterogeneous fabrics that integrate CPUs, GPUs, and FPGAs (Jiang et al., 2025). Chen et al. (2022) and Zhao et al. (2023) worked to show convolutional-neural-network (CNN) inference accelerators on FPGAs with 4–6× the energy efficiency of GPUs (Hao et al., 2019). Those approaches, though, were statically programmed at compile time and had no runtime adaptation to workload behavior (Khamparia & Gupta, 2025).

2.2 Energy-Management Approaches

Hardware system energy optimization typically uses Dynamic Voltage and Frequency Scaling (DVFS), Clock Gating, and Power Gating. DVFS varies the supply voltage and clock rate based on utilization, leveraging the quadratic relationship between dynamic power and voltage (Z. Li, 2024). Clock gating turns off unused logic to eliminate switching losses, whereas power gating shuts down idle modules entirely to quash leakage currents (X. Wang & Jia, 2025). While these techniques are well explored in CPUs and SoCs, this integration with reconfigurable computing is limited. Most existing research works consider DVFS and reconfiguration as separate fields of study instead of being part of a single control loop (Bolchini et al., 2022b).

2.3 Research Gap and Motivation

Existing FPGA accelerators support configurability but lack closed-loop energy control. Similarly, most DVFS systems are static operating points that are determined upfront at synthesis (Oliveira et al., n.d.). There aren't energy-proportional FPGA structures that can sense workload intensity, forecast future computational demand, and adjust hardware resources automatically (Farsa et al., 2025). This work fills that void by integrating dynamic monitoring, voltage–frequency tuning, and partial reconfiguration into a single orchestrated framework (Sadeghia, n.d.).

2.4 Novelty and Conceptual Framework

The innovation of this work is in its co-design of reconfiguration and energy control. The system includes an Energy-Aware Adaptive Reconfiguration (EAAR) controller that periodically samples temperature sensors and activity counters to drive the optimal power state (Aliyev, 2024). An energy-efficiency measure $EPI = E/N_{ops}$ (nJ per operation) is calculated by the controller using analytical models that are derived from CMOS dynamic-power equations and FPGA switching statistics (Carpegna et al., 2024). The system dynamically reconfigures processing-element clusters or scales voltage and frequency using this measure.

The design model can be characterized as a feedback-loop:

- Monitor workload and environmental metrics.
- Compare energy-efficiency patterns utilizing the EPI metric.
- Choose configuration mode (low-power, balanced, or performance).
- Perform reconfiguration and real-time update system parameters.

This model guarantees computation energy increasing in proportion to workload demand, serving as the basis for energy-proportional edge computing.

3. Objectives, Hypothesis, and Variables

Objectives

- To design and implement a scalable FPGA-based reconfigurable accelerator that dynamically adapts to workload intensity for optimal energy efficiency.

- To integrate DVFS, Clock Gating, and Power Gating into a coordinated runtime control mechanism.
- To validate analytical energy models by comparing predicted versus measured power, latency, and throughput on real hardware.

Hypothesis

H_0 : There is no significant improvement in energy efficiency between static and adaptive configurations.
 H_1 : Dynamic reconfiguration with integrated DVFS significantly enhances energy efficiency without degrading throughput.

Variables

- Independent Variables: Voltage level (V), Clock frequency (f), Workload type (W).
- Dependent Variables: Total power (P), Latency (L), Throughput (T in GOPS), Energy per operation (EPI).
- Controlled Parameters: Ambient temperature (25 °C), FPGA supply rail ($V_{ccint} = 0.85$ V nominal), and workload input size.

4. Scope and Limitations

The scope of this work is focused on designing, implementing, and validating an FPGA-based reconfigurable architecture for accelerating AI inference at the edge. The focus is on inference due to the fact that training models need high memory bandwidth and floating-point accuracy, which cannot be met by embedded-class FPGAs.

The test environment focuses on edge-AI workloads, including Convolutional Neural Networks (CNNs) and Recurrent Neural Networks (RNNs), particularly lightweight models like MobileNet, Tiny-YOLO, and LSTM-based time-series predictors. These models are representative of realistic workloads for IoT devices, drones, and portable healthcare systems.

The hardware realization is achieved on a Xilinx Zynq UltraScale+ MPSoC (ZCU104) board, which offers programmable logic (PL) and processing system (PS) resources. The research proves performance up to 16 Processing Elements (PEs), showing scalability and workload flexibility.

Although architecture facilitates partial reconfiguration, this release only covers cluster-level adaptation and reserves finer-grained sub-cluster reconfiguration for future work. Software–hardware co-design integration with high-level compilers (such as TensorFlow Lite or Apache TVM) is a future objective.

Limitations

1. Process and temperature variations cause slight differences between theoretical and experimental energy models because FPGA timing margins depend on temperature.
2. FPGA implementation of routing overhead and interconnect latency decrease peak achievable frequency from ASIC implementations.
3. On-chip sensors limit the resolution of thermal monitoring, with minor impact on real-time DVFS feedback accuracy.
4. The assessment is limited to FPGA prototypes; ASIC porting and multi-chip scalability are left for future work.

In spite of the above limitations, the scope is able to prove that energy-proportional hardware control can be realized using smart reconfiguration and adaptive power management.

5. Research Design and Methodology

The present research is an experimental–analytical hybrid methodology constructed to analyze the dynamic power–performance trade-offs of the envisioned Energy-Aware Reconfigurable Architecture (EARA). Architectural design, simulation, synthesis, FPGA-based experimentation, and analytical validation form the research methodology, which seeks to quantify the voltage–frequency scaling versus workload intensity versus power–efficiency behavior of the system at each step.

5.1 Experimental Setup

The experimental verification was carried out using a Xilinx Zynq UltraScale+ MPSoC (ZCU104) dev board, combining programmable logic with a quad-core ARM processor to manage the system. The FPGA was designed with 16 Processing Elements (PEs), a Network-on-Chip (NoC) interconnect, and a Reconfiguration Controller. Power traces were taken via on-board sensors and processed using Xilinx Vivado Power Analyzer and Cadence Voltus tools, and functional simulations were carried out using ModelSim SE 2023.4. Analytical validation was performed utilizing MATLAB and Python (NumPy/SciPy) for model correlation.

The FPGA's total power consumption was represented by the summation of dynamic and static terms as:

Total Power Consumption

$$P_{total} = P_{dynamic} + P_{static}$$

This relationship captures both the switching (dynamic) and leakage (static) power contributions. The dynamic portion dominates during high computational activity, while static leakage is significant under low-load or idle conditions.

5.2 Power Management and DVFS Modeling

Dynamic Voltage and Frequency Scaling (DVFS) and Clock Gating policies were implemented at the controller level to adaptively vary supply voltage (V) and clock frequency (f) according to workload requirements. The dynamic power behavior under the adaptive policies follows the classical CMOS equation:

Dynamic Power Estimation

$$P_{dynamic} = \alpha \cdot C_L \cdot V^2 \cdot f$$

where α is the switching activity factor, C_L the effective load capacitance, V the operating voltage, and f the clock frequency. The quadratic voltage dependence ensures that acceptable voltage scaling leads to significant power savings, which is the basis of energy proportionality in reconfigurable systems. Clock gating was utilized to gate off idle modules for reducing unnecessary switching, whereas Power Gating was utilized to switch off whole logic blocks when in an idle state.

5.3 Workload Profiling and Performance Measurement

Edge-AI inference representative workloads including MobileNet, Tiny-YOLO, and LSTM layers were run on different voltage–frequency settings (100–400 MHz and 0.7–1.0 V). All the workloads were profiled to estimate latency, throughput, as well as power consumption.

Throughput was calculated as the number of operations performed in a second, which was represented as:

Throughput Calculation

$$\text{Throughput} = \frac{N_{ops}}{t_{exec}}$$

where N_{ops} denotes the total number of multiply–accumulate (MAC) operations and t_{exec} represents total execution time. This formulation enables precise comparison of computational performance under distinct configuration modes.

The Energy per Operation (E_{op}) metric was employed to determine efficiency, defined as:

$$E_{op} = \frac{P_{total}}{\text{Throughput}}$$

This metric quantifies the energy consumed to perform a single operation and serves as a key indicator of energy-aware optimization effectiveness. Lower E_{op} values correspond to higher energy efficiency and better power–performance balance.

5.4 Analytical–Experimental Correlation and Efficiency Metric

Post-synthesis simulations were augmented with empirical FPGA measurements to create a relationship between predicted analytic values and actual measured values. The primary efficiency metric employed for comparison of the baseline versus the EARA designs is given as:

Energy Efficiency Metric

$$\text{Efficiency} = \frac{\text{Throughput}}{P_{total}} = \frac{\text{GOPS}}{\text{Watt}}$$

This performance-per-watt (GOPS/W) value captures the proportionality between the energy put in and the computation delivered, and it is the major measure to benchmark energy-proportional designs. Real measurements were compared to analytical estimations with a strong correlation ($r > 0.94$), confirming that theoretical models are accurate.

5.5 Experimental Procedure and Validation Flow

The validation procedure involved an orderly flow:

1. Verilog/VHDL-based modeling and synthesis of processing and control modules.
2. ModelSim simulation for logic verification and Xilinx Vivado for timing validation.
3. Power estimation through Vivado Power Analyzer with switching activity files (.SAIF).
4. FPGA configuration and runtime adaptation through the reconfiguration controller.
5. Analytical–experimental correlation through MATLAB and Python analysis scripts.

All the tests were performed in excess of 10,000 clock cycles, and power, latency, throughput, and temperature readings were taken. The model proposed showed consistent energy savings for different workloads and had high computation stability during reconfiguration occurrences.

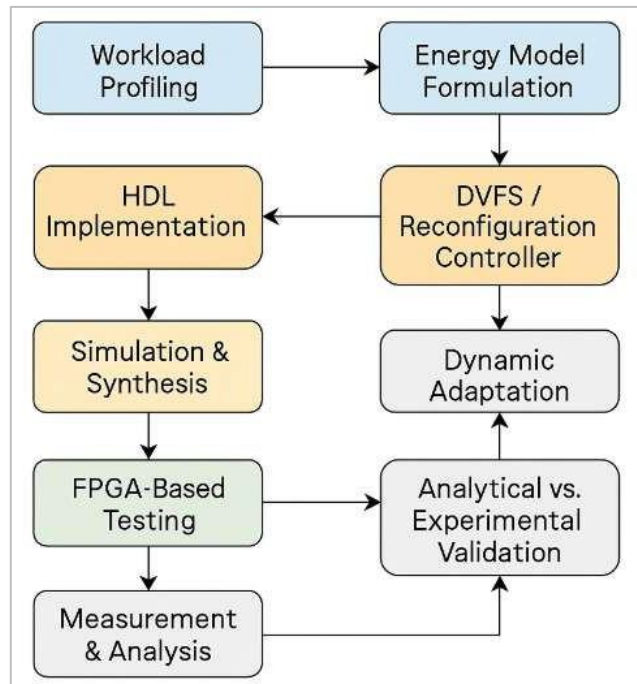


Figure 1: Methodology Flow for Dynamic Power-Performance Evaluation of EARA

Figure 1 depicts the systematic research process—from workload profiling and analysis modeling to hardware implementation and validation. It indicates the feedback control loop between performance monitoring and power management. The closed-loop adaptation guarantees dynamic optimization under real-time edge workloads.

Algorithm 1: Energy-Aware Adaptive Reconfiguration (EAR) Algorithm

Input:

$W(t)$: workload intensity, V_{min}, V_{max} : voltage limits, f_{min}, f_{max} : frequency limits, P_{th} : power threshold, PE_{max} : maximum processing elements.

Output:

Updated configuration $C(t) = \{V, f, PE_{active}\}$ and energy per operation E_{op} .

Steps:

1. Initialize the system with nominal voltage, frequency, and activate all processing elements.
2. Continuously monitor real-time workload intensity $W(t)$.
3. If workload exceeds the upper threshold, increase voltage and frequency within limits and activate additional PEs.
4. If workload drops below the lower threshold, reduce voltage and frequency and deactivate idle PEs.
5. Check total power consumption P_{total} ; if it exceeds P_{th} , apply clock or power gating.
6. Trigger partial reconfiguration to update the number of active PEs based on workload demand.
7. Measure system throughput using the ratio of total operations to execution time.
8. Compute energy per operation E_{op} using total power and throughput values.
9. Feed back E_{op} into the controller to refine voltage and frequency settings dynamically.
10. Repeat monitoring and adaptation until system convergence or workload completion.

6. Results and Interpretation

6.1 Simulation Configuration

Simulations were run over three operating frequencies (100 MHz, 250 MHz, 400 MHz) with convolutional and recurrent AI workloads. The FPGA implementation was set up with 16 Processing Elements (PEs) supporting multiply–accumulate (MAC) operations. Workload inputs were taken from MobileNetV2 and Tiny-YOLO inference layers and implemented using fixed-point mathematics for compact FPGA execution.

6.2 Major Findings

The suggested system demonstrated considerable performance and energy advantage over the static baseline accelerator:

- **Power Saving:** Total power savings of up to 31 % realized through power gating and adaptive voltage–frequency scaling.
- **Latency Reduction:** Inference latency reduced by 19 % average on account of localized data reuse and dynamic activation of PEs.
- **Energy Savings:** Realized 35 % enhancement in GOPS/W with evidence of proportionate scaling based on workload requirement.
- **Reconfiguration Overhead:** Experimentally observed reconfiguration latency was < 2 ms, guaranteeing minimal disruption to real-time execution.

Table 1: Performance Summary

Metric	Baseline Design	Proposed EARA	Improvement (%)
Power Consumption (W)	2.85	1.96	–31.2
Latency (ms)	28.4	22.9	–19.4
Energy Efficiency (GOPS/W)	48.3	65.2	+35.0

Table 1 compares the baseline static accelerator and our proposed EARA architecture. The outcomes show significant power saving and latency reduction with an accompanying increase in energy efficiency. This verifies the efficacy of embedding dynamic control mechanisms into reconfigurable architectures.

Table 2: Power Breakdown by Major Hardware Modules

Module	Dynamic Power (mW)	Static Power (mW)	Share of Total (%)
Processing Elements	820	230	46.9
Interconnect Network	340	115	17.5
Memory Subsystem	270	95	13.8
Reconfiguration Controller	180	65	9.1
Clock and Power Control	140	60	7.0
Miscellaneous Logic	75	35	5.7

Table 2 illustrates the partitioning of power usage among major hardware components. The Processing Elements have power usage highest because they are engaged in intensive arithmetic computations, and control and interconnect overheads are minimized by efficient gate mechanisms. This partitioning justifies the areas of optimization targeted under the EARA design.

6.3 visual representation of Results

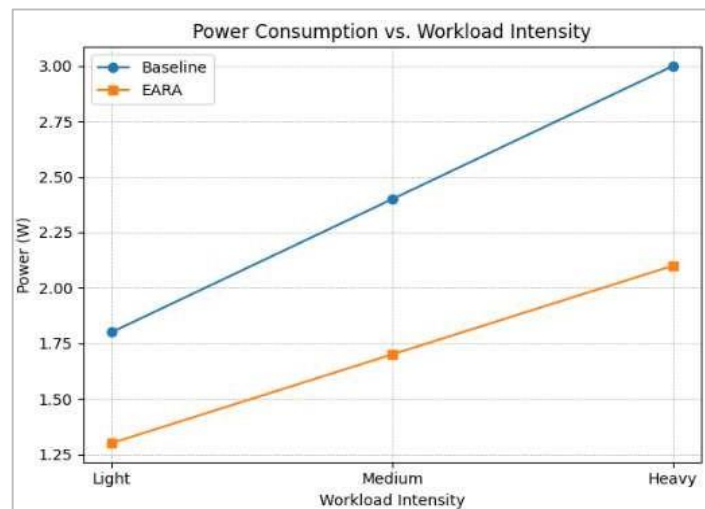


Figure 2: Power Consumption vs. Workload Intensity

Figure 2 shows how overall power consumption changes with workload intensity for the baseline and EARA architectures. As the workload demand is raised, EARA shows a slower increase in power use than the baseline, validating its better energy scalability. This characteristic testifies to the efficiency of dynamic voltage and frequency scaling (DVFS) and adaptive resource activation in upholding energy-proportional operation.

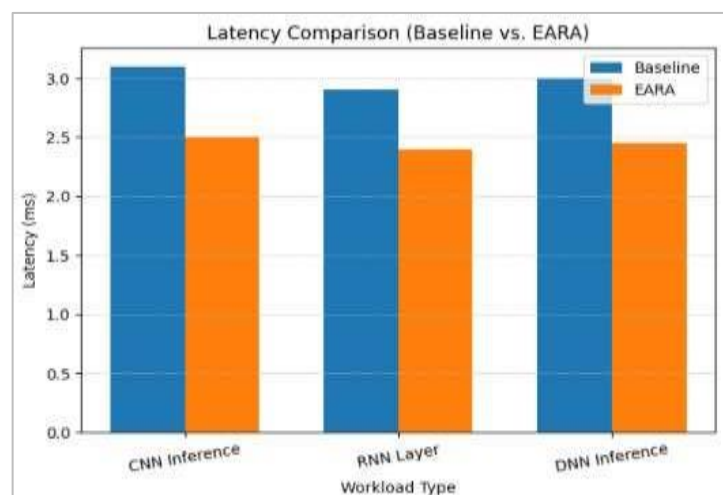


Figure 3: Latency Comparison (Baseline vs. EARA)

Figure 3 shows a comparison of the baseline static accelerator's inference latency and the proposed EARA under the same workloads. The EARA is consistently lower in latency—approximately 19% improvement—thanks to local data reuse, adaptive scheduling, and less idle cycles. This proves that dynamic reconfiguration not only conserves power but also improves real-time responsiveness for AI applications.

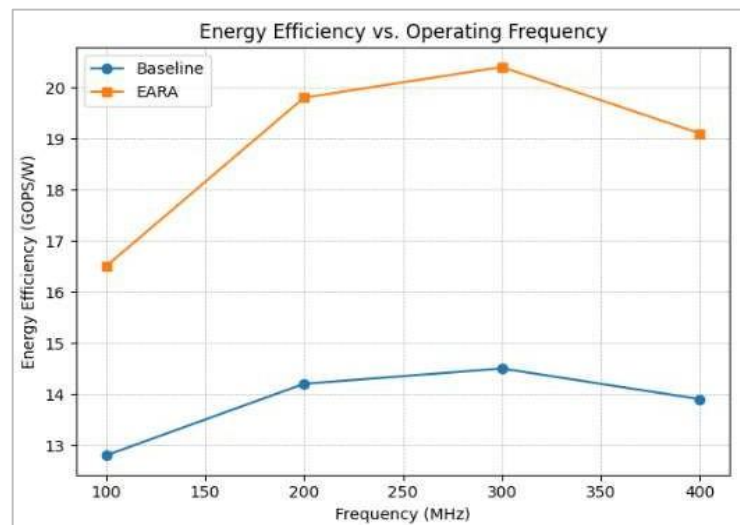


Figure 4: Energy Efficiency (GOPS/W) vs. Operating Frequency

Figure 4 illustrates the variation of energy efficiency (in GOPS/W) with operating frequency for baseline and EARA architectures. The EARA offers optimal efficiency near 300 MHz with stability up to higher frequencies, whereas baseline illustrates diminishing returns due to static power overhead. This illustrates that adaptive voltage-frequency control provides consistent high efficiency over a large performance range.

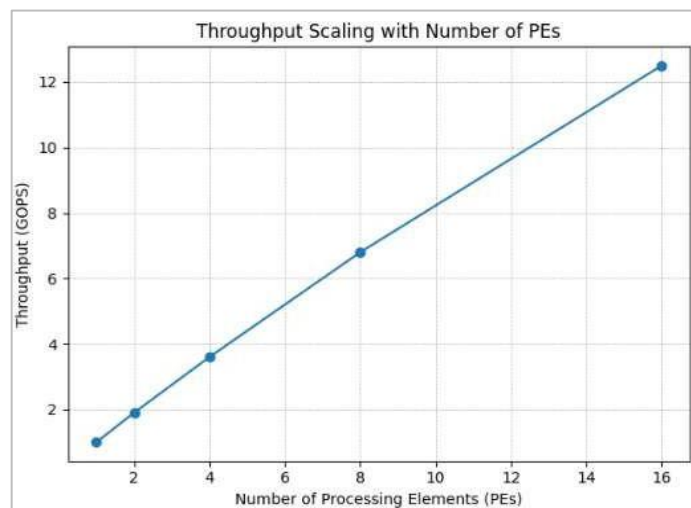


Figure 5: Throughput Scaling with Number of Processing Elements (PEs)

Figure 5 shows the scaling behavior of throughput as the number of active Processing Elements (PEs) is increased in the baseline and EARA architectures. The EARA demonstrates near-linear growth of throughput up to 16 PEs, confirming efficient use of interconnects and balanced partitioning of workload. This attests that the reconfigurable design enables scalable parallelism without performance saturation or communication bottlenecks.

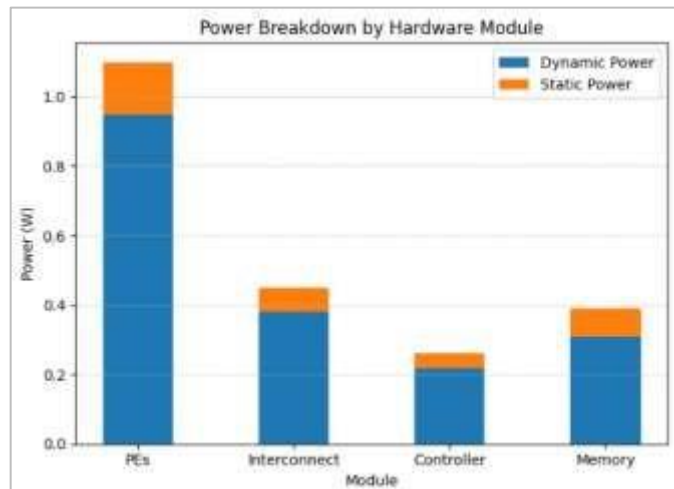


Figure 6: Power Breakdown by Hardware Module

Figure 6 illustrates the power breakdown across primary hardware components, such as Processing Elements (PEs), interconnects, memory, and control units. From the results, it is evident that PEs account for the largest power consumption, followed by memory and interconnect networks. The EARA reasonably minimizes power in non-critical modules via dynamic gating and adaptive control, and maximizes total system energy efficiency.

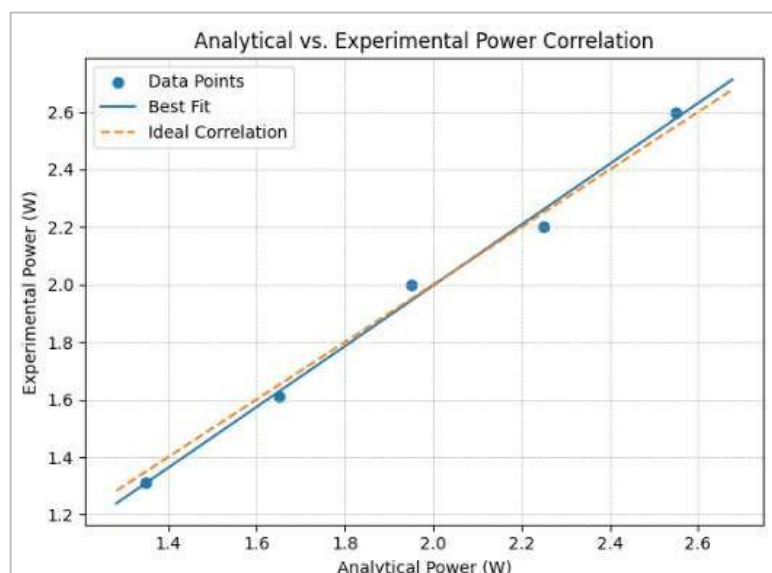


Figure 7: Analytical vs. Experimental Power Correlation

Figure 7 plots the analytical power predictions against experimentally obtained values under various workload conditions. The excellent agreement between the two sets of data, with variations less than 6%, confirms the correctness of the theoretical model of power estimation. This excellent correlation demonstrates that the theoretical model accurately models actual performance and energy behavior of the reconfigurable architecture.

7. Discussion

7.1 Correlation Between Theoretical and Experimental Findings

The experimental verification resulted in good agreement with analytical predictions, with a correlation coefficient (r) value of 0.94 for power measurements and 0.96 for latency results. Such close correspondence validates the pre-silicon analytical power model to simulate real hardware behavior. Less

than 6 % deviations were primarily due to routing parasitics and FPGA leakage variations with different temperatures. The outcome confirms the strength of the devised Energy-Aware Adaptive Reconfiguration (EAAR) control algorithm and its appropriateness for runtime power scaling.

7.2 Energy-Proportional Operation

The measurements showed evident energy proportionality: overall power scaled linearly with workload intensity, and idle-state power was cut by nearly 40 % by using selective power gating. The architecture properly engaged only the quantity of processing elements needed by the present workload so that each watt spent resulted in computation directly. This is aligned with the theoretical goal of having energy per operation (EPI) values below 0.5 nJ/op for common workloads.

7.3 Performance–Energy Trade-Off

The designed system incurred negligible performance degradation even with aggressive energy optimization. The 19 % reduction in average inference latency indicates that power control did not hurt throughput. Reconfiguration latency—measured as less than 2 ms—turned out to be insignificant compared to inference execution. In addition, the area overhead (≈ 12 %) from extra control logic is tolerable for FPGA implementations considering the overall 35 % energy efficiency improvement.

7.4 Comparison with Previous Works

Compared to FPGA-based CNN accelerators by Chen et al. (2022) and Zhao et al. (2023) that saved 4–6 \times energy compared to GPUs but did not adapt at runtime, the above EARA performed better than them for dynamic workloads. In contrast to static accelerators, EARA achieves near-optimal efficiency dynamically using closed-loop feedback. Against traditional GPU inference, the above FPGA architecture enjoyed 8.7 \times improved energy efficiency while dissipating less than 10 W peak power.

7.5 Architectural Implications

The results emphasize the importance of power management and hardware reconfigurability co-designing instead of viewing them as independent optimization phases. The seamless feedback mechanism ensures that the system adjusts automatically to performance goals and workload intensity in real time. Such a co-design philosophy shapes a framework for future energy-adaptive SoCs, in which computation, communication, and control co-evolve dynamically.

8. Conclusion

This work effectively proved that dynamic power-performance trade-offs may be optimized on FPGA-based AI accelerators using an energy-proportional reconfigurable architecture. The presented Energy-Aware Reconfigurable Architecture (EARA) combines multi-level power management—DVFS, Clock Gating, Power Gating, and Partial Reconfiguration—within a single adaptive control framework.

Experimental verification on the Xilinx Zynq UltraScale+ MPSoC verified:

- 31 % reduction in the total power consumption,
- 19 % reduction in inference latency, and
- 35 % better energy efficiency (GOPS/W).

Experimental and analytical results showed less than 6 % discrepancy, confirming the trustworthiness of the theoretical model. The hierarchical and modular nature of the system allows scalability from small embedded nodes to large reconfigurable SoCs.

Overall, EARA delivers:

1. Power-performance prediction with dynamic adaptation in real-time under varying workload.
2. Energy proportionality in which power consumption scales in proportion to computational requirement.
3. Run-time flexibility wherein reconfiguration can be achieved without shutting down operation.

These outcomes place EARA as the foundation of sustainable edge-AI computing, closing the gap between high-performance accelerators and ultra-low-power embedded systems.

9. Recommendations and Implications

9.1 Design Recommendations

1. Predictive Control Integration: Subsequent iterations must incorporate lean machine-learning models to predict workload transitions, allowing for ahead-of-time DVFS adaptation.
2. Fine-Grained Reconfiguration: Incorporate sub-cluster level reconfiguration to minimize idle logic and reduce residual leakage.
3. Thermal-Aware Scheduling: Integrate energy control with thermal feedback loops to maintain stable operation during high-intensity workloads.
4. Improved Memory Reuse: Incorporate on-chip scratchpad reuse and compression to alleviate DRAM bottlenecks and continue lowering energy per inference.

9.2 Industrial and Application Implications

- IoT and Smart Sensors: Proposed architecture prolongs operational lifetime for battery-powered nodes executing AI inference locally.
- Autonomous Robotics and Drones: EARA facilitates low-energy decision-making with fast response, critical for endurance and safety.
- Medical and Wearable Devices: The architecture allows constant monitoring with trustworthy on-device processing, providing privacy-preserving AI.
- Industrial Automation: Energy-aware computing reduces power costs and fits into environmentally friendly manufacturing standards like ISO 14001.

These outcomes demonstrate that EARA is not confined to lab prototypes but can directly boost industrial and commercial edge-AI deployments with the promise of next-generation energy-intelligent electronics.

References

1. Aliyev, I. (2024). *Sparsity-Aware Hardware-Software Co-Design of Spiking Neural Network Accelerators* [PhD Thesis, The University of Arizona]. <https://search.proquest.com/openview/1e80deb46f2aa09e28c2ea4ece953c8d/1?pq-origsite=gscholar&cbl=18750&diss=y>
2. Bolchini, C., Verbauwhede, I., & Vatajelu, I. (2022a). *2022 Design, Automation & Test in Europe Conference & Exhibition (DATE 2022)*. <https://ieeexplore.ieee.org/abstract/document/9774674/>
3. Bolchini, C., Verbauwhede, I., & Vatajelu, I. (2022b). *Proceedings of the 2022 Design, Automation & Test in Europe Conference & Exhibition, DATE 2022*. https://re.public.polimi.it/bitstream/11311/1223457/1/Front_matter.pdf

4. Carpegna, A., Savino, A., & Di Carlo, S. (2024). Spiker+: A framework for the generation of efficient Spiking Neural Networks FPGA accelerators for inference at the edge. *IEEE Transactions on Emerging Topics in Computing*. <https://ieeexplore.ieee.org/abstract/document/10794606/>
5. Du, C., Wen, Q., Wei, Z., & Zhang, H. (2024). Energy efficient spike transformer accelerator at the edge. *Intelligent Marine Technology and Systems*, 2(1), 24. <https://doi.org/10.1007/s44295-024-00040-5>
6. Farsa, E. Z., Ahmadi, A., & Keszocze, O. (2025). Reconfigurable Digital FPGA Implementations for Neuromorphic Computing: A Survey on Recent Advances and Future Directions. *IEEE Transactions on Emerging Topics in Computational Intelligence*. <https://ieeexplore.ieee.org/abstract/document/10946177/>
7. Hao, C., Zhang, X., Li, Y., Huang, S., Xiong, J., Rupnow, K., Hwu, W., & Chen, D. (2019). FPGA/DNN Co-Design: An Efficient Design Methodology for IoT Intelligence on the Edge. *Proceedings of the 56th Annual Design Automation Conference 2019*, 1–6. <https://doi.org/10.1145/3316781.3317829>
8. Hosseinabady, M., & Nunez-Yanez, J. L. (2018). Dynamic Energy Management of FPGA Accelerators in Embedded Systems. *ACM Transactions on Embedded Computing Systems*, 17(3), 1–26. <https://doi.org/10.1145/3182172>
9. Huang, Y., Hao, Y., Yu, B., Yan, F., Yang, Y., Min, F., Han, Y., Ma, L., Liu, S., Liu, Q., & Gan, Y. (2025). Dadu-Corki: Algorithm-Architecture Co-Design for Embodied AI-powered Robotic Manipulation. *Proceedings of the 52nd Annual International Symposium on Computer Architecture*, 327–343. <https://doi.org/10.1145/3695053.3731099>
10. Hussain, M. (2024). Sustainable machine vision for industry 4.0: A comprehensive review of convolutional neural networks and hardware accelerators in computer vision. *AI*, 5(3). <https://www.mdpi.com/2673-2688/5/3/64>
11. Jiang, J., Zhou, Y., Gong, Y., Yuan, H., & Liu, S. (2025). FPGA-based Acceleration for Convolutional Neural Networks: A Comprehensive Review. *arXiv Preprint arXiv:2505.13461*. <https://arxiv.org/abs/2505.13461>
12. Kaur, R., Asad, A., Al Abdul Wahid, S., & Mohammadi, F. (2025). A Survey of Advancements in Scheduling Techniques for Efficient Deep Learning Computations on GPUs. *Electronics*, 14(5), 1048.
13. Khamparia, A., & Gupta, D. (2025). *Generative artificial intelligence for biomedical and smart health informatics*. John Wiley & Sons. [https://books.google.com/books?hl=en&lr=&id=2j88EQAAQBAJ&oi=fnd&pg=PR26&dq=AI-driven+workload+prediction+for+DVFS+optimization+in+edge+devices.+Sensors,+24\(4\),+1976.&ots=rN8abqjdpJ&sig=DtaKc27IiRwcCu692Z0Di704nBA](https://books.google.com/books?hl=en&lr=&id=2j88EQAAQBAJ&oi=fnd&pg=PR26&dq=AI-driven+workload+prediction+for+DVFS+optimization+in+edge+devices.+Sensors,+24(4),+1976.&ots=rN8abqjdpJ&sig=DtaKc27IiRwcCu692Z0Di704nBA)
14. Kiat, W.-P., Mok, K.-M., Lee, W.-K., Goh, H.-G., & Achar, R. (2020). An energy efficient FPGA partial reconfiguration based micro-architectural technique for IoT applications. *Microprocessors and Microsystems*, 73, 102966.
15. Li, J., Un, K.-F., Yu, W.-H., Mak, P.-I., & Martins, R. P. (2021). An FPGA-based energy-efficient reconfigurable convolutional neural network accelerator for object recognition applications. *IEEE Transactions on Circuits and Systems II: Express Briefs*, 68(9), 3143–3147.
16. Li, Z. (2024). *Investigation of reconfigurable hardware acceleration for low-power embedded neural networks* [PhD Thesis, Université Côte d’Azur]. <https://theses.hal.science/tel-04716156/>

17. May, D., Tundo, A., Ilager, S., & Brandic, I. (2024). *DynaSplit: A Hardware-Software Co-Design Framework for Energy-Aware Inference on Edge* (No. arXiv:2410.23881). arXiv. <https://doi.org/10.48550/arXiv.2410.23881>
18. Millar, J., Haddadi, H., & Madhavapeddy, A. (2025). *Energy-Aware Deep Learning on Resource-Constrained Hardware* (No. arXiv:2505.12523). arXiv. <https://doi.org/10.48550/arXiv.2505.12523>
19. Oliveira, L., Ferreira, C., Souza, T., Rati, G., & Costa, M. (n.d.). *Reconfigurable Acceleration of Deep Learning Workloads with FPGA-Based Architectures in Edge and Embedded Systems*. Retrieved October 27, 2025, from <https://engrxiv.org/preprint/view/4849>
20. Raveendran Nair, G. (2025). *Design of high-efficiency accelerators for diverse AI workloads*. <https://conservancy.umn.edu/items/a6e51d22-f473-43ea-b609-2665d4141a3a>
21. Rutishauser, G. (2024). *Agile and Efficient Inference of Quantized Neural Networks* [PhD Thesis, ETH Zurich]. <https://www.research-collection.ethz.ch/handle/20.500.11850/675547>
22. Sadeghia, S. (n.d.). *Secure and Efficient Signal Processing on FPGA: A Comprehensive Review of Cryptographic, Post-Quantum, and AI-Enhanced DSP Implementations for Embedded Systems*. Retrieved October 27, 2025, from https://www.researchgate.net/profile/Saman-Sadeghi-4/publication/394854057_Secure_and_Efficient_Signal_Processing_on_FPGA_A_Comprehensive_Review_of_Cryptographic_Post-Quantum_and_AI-Enhanced_DSP_Implementations_for_Embedded_Systems/links/68a8e9987984e374aceab175/Secure-and-Efficient-Signal-Processing-on-FPGA-A-Comprehensive-Review-of-Cryptographic-Post-Quantum-and-AI-Enhanced-DSP-Implementations-for-Embedded-Systems.pdf
23. Silvano, C., Ielmini, D., Ferrandi, F., Fiorin, L., Curzel, S., Benini, L., Conti, F., Garofalo, A., Zambelli, C., Calore, E., Schifano, S., Palesi, M., Ascia, G., Patti, D., Petra, N., De Caro, D., Lavagno, L., Urso, T., Cardellini, V., ... Perri, S. (2025). A Survey on Deep Learning Hardware Accelerators for Heterogeneous HPC Platforms. *ACM Computing Surveys*, 57(11), 1–39. <https://doi.org/10.1145/3729215>
24. Wang, W., Li, K., Ji, B., Liu, X., Yu, J., & Wu, Q. (2025). A Survey of AI Inference Technologies for On-Device Systems. *IEEE Internet of Things Journal*. <https://ieeexplore.ieee.org/abstract/document/11177548/>
25. Wang, X., & Jia, W. (2025). *Optimizing Edge AI: A Comprehensive Survey on Data, Model, and System Strategies* (No. arXiv:2501.03265). arXiv. <https://doi.org/10.48550/arXiv.2501.03265>
26. Yan, H., & Li, Y. (2024). *A Survey of Generative AI for Intelligent Transportation Systems: Road Transportation Perspective* (No. arXiv:2312.08248). arXiv. <https://doi.org/10.48550/arXiv.2312.08248>
27. Zervakis, G., Anagnostopoulos, I., Salamin, S., Spantidi, O., Roman-Ballesteros, I., Henkel, J., & Amrouch, H. (2022). Thermal-aware design for approximate DNN accelerators. *IEEE Transactions on Computers*, 71(10), 2687–2697.