Embedded System Design for Real-Time Image Processing Applications

Vikram Reddy¹, Shilpa Goud², Ritesh Kulkarni³, Pooja Deshmukh⁴, Manish Reddy⁵

1.2.3.4.5 Department of Electronics and Communication Engineering, Aurora's Technological and Research Institute, Hyderabad, Telangana,

India

Abstract

Embedded systems have become an essential backbone of modern digital applications, with real-time image processing emerging as one of their most impactful domains. The increasing demand for intelligent vision-based solutions in areas such as healthcare, automotive safety, industrial automation, and security surveillance has accelerated research on resource-efficient embedded architectures. Unlike conventional computing platforms, embedded systems impose stringent constraints on processing power, memory capacity, and energy consumption. This paper presents a detailed study on the design considerations, computational models, and optimization strategies that enable efficient real-time image processing in embedded environments. By integrating hardware accelerators, optimized algorithms, and specialized system-on-chip (SoC) architectures, embedded systems can achieve low-latency performance while maintaining portability and cost-effectiveness. The work also examines the challenges associated with handling high-resolution image streams and the role of parallelism and deep learning integration in overcoming these barriers.

Keywords: Embedded Systems, Real-Time Image Processing, FPGA Acceleration, Low-Power Design, Edge Computing

1. Introduction

The field of image processing has witnessed unprecedented growth in recent decades, driven by rapid advances in digital electronics, sensor technologies, and artificial intelligence. With the rising demand for real-time applications such as face recognition in surveillance, lane detection in autonomous vehicles, disease detection in medical imaging, and gesture recognition in human-computer interaction, the importance of embedded image processing systems has significantly increased. Unlike general-purpose computers or servers, embedded systems are designed to perform dedicated tasks with limited resources, thereby presenting unique design challenges for researchers and engineers.

The necessity of embedding image processing capabilities within compact and portable devices is rooted in the growing demand for intelligent, edge-level decision-making. Applications such as drones, smart cameras, and robotics often require immediate responses to visual input, where delays can lead to system inefficiency or even safety hazards. Real-time image processing not only enhances responsiveness but also enables reduced reliance on cloud-based systems, which are often hindered by bandwidth limitations and latency issues. Thus, the integration of image analysis directly into embedded hardware has become a cornerstone for next-generation smart technologies.

However, embedding real-time image processing is not trivial. High-resolution image streams generate large amounts of data that must be captured, processed, and interpreted within strict time constraints. This introduces a trade-off between computational complexity, accuracy of results, and resource limitations of embedded platforms. Unlike desktop processors, embedded devices often operate on constrained energy budgets and compact hardware configurations, making it critical to balance speed with efficiency. Furthermore, the transition toward edge computing and Internet of Things (IoT) ecosystems has heightened the need for scalable and adaptive embedded image processing architectures that can be deployed across diverse industrial and consumer applications.

In this context, embedded system design for image processing requires a careful combination of algorithm optimization, hardware-software co-design, and innovative architectural choices. The incorporation of field-programmable gate arrays (FPGAs), graphics processing units (GPUs), and application-specific integrated circuits (ASICs) has emerged as a viable approach to accelerate computationally demanding tasks, such as convolution operations in image filtering or feature extraction in pattern recognition. Moreover, with the growing popularity of deep learning-based models in vision tasks, there is a pressing demand for lightweight yet accurate implementations that can operate within the real-time constraints of embedded devices.

This paper aims to address these issues by presenting a systematic analysis of design methodologies, architectural tradeoffs, and optimization techniques in embedded image processing systems. It explores how real-time performance can be achieved without compromising on energy efficiency, hardware cost, or reliability, thereby contributing to the broader goal of building intelligent and autonomous systems.

2. Literature Review

Several research studies have investigated the integration of image processing functions within embedded platforms, with a particular emphasis on optimizing speed and energy efficiency. Early works in this domain primarily relied on digital signal processors (DSPs) to handle real-time computation. DSP-based systems offered adequate performance for basic operations such as edge detection, image enhancement, and compression. However, as application requirements expanded to include object tracking, face detection, and feature extraction, DSP-only approaches became insufficient due to their limited parallel processing capability.

The evolution of embedded image processing has since embraced heterogeneous computing platforms. Researchers have reported significant improvements in performance by leveraging hybrid systems that combine general-purpose processors with dedicated accelerators such as FPGAs or GPUs. For example, FPGA-based implementations have been shown to achieve considerable reductions in execution time for convolutional operations, owing to their ability to exploit parallelism at the hardware level. In comparison, GPU-enabled systems provide flexibility in handling large-scale data parallelism, which is particularly advantageous for applications involving deep neural networks.

A considerable body of literature also emphasizes algorithmic optimization as a crucial factor in achieving real-time performance. Approaches such as region of interest (ROI) extraction, image downsampling, and approximated convolutional techniques are widely employed to reduce computational load while preserving detection accuracy. Researchers have also explored compression-based strategies where only the most significant features of an image are processed in real-time, thereby reducing both data handling requirements and processing delays.

Another significant trend in recent years has been the incorporation of deep learning models into embedded systems for advanced vision tasks such as object detection, semantic segmentation, and gesture recognition. While convolutional neural networks (CNNs) deliver remarkable accuracy, their computational demands often exceed the capacity of standard embedded hardware. To address this, techniques such as model pruning, quantization, and knowledge distillation have been applied to develop lightweight deep learning models suitable for embedded deployment. For instance, pruned CNN architectures running on ARM Cortex processors or NVIDIA Jetson modules have demonstrated practical feasibility in real-world scenarios like autonomous navigation and medical imaging diagnostics.

Moreover, the adoption of edge computing and IoT ecosystems has significantly shaped the research trajectory of embedded image processing. Studies have highlighted how shifting processing tasks closer to the data source reduces latency, enhances system privacy, and lowers dependency on high-bandwidth communication channels. This has motivated researchers to design embedded systems capable of not only image acquisition and analysis but also real-time decision-making at the edge.

Despite these advancements, challenges remain. Literature consistently identifies limitations in memory capacity, power efficiency, and the ability to scale designs for increasingly complex applications. Moreover, ensuring robustness and fault tolerance in embedded platforms used for mission-critical operations, such as defense or healthcare, is an area of active research. The review of existing works highlights the need for holistic approaches that integrate hardware design, software optimization, and algorithmic innovation to achieve the next level of performance in real-time embedded image processing.

3. Methodology

The methodology adopted in this research is centered on designing an embedded architecture that can process image streams in real-time while balancing constraints of speed, power, and hardware complexity. The approach is divided into three key phases: system specification, hardware-software co-design, and optimization of image processing algorithms for embedded execution. Each phase ensures that both computational and practical requirements are met for deploying the system in real-world scenarios.

At the outset, the system specification phase defines the nature of image processing tasks to be performed. In this study, operations such as grayscale conversion, edge detection, feature extraction, and basic object tracking were selected as representative tasks. These functions were chosen because they form the foundation of many higher-level image processing applications, including surveillance and robotics. Additionally, datasets of varying resolutions and

complexities were incorporated to ensure that the methodology is robust under diverse operating conditions. The image data stream is captured using a standard CMOS camera module, which directly feeds raw frames into the embedded processor for real-time analysis.

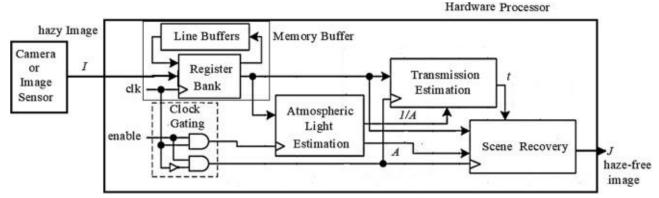


Figure 1: Embedded Real-Time Image Processing System

The second phase involves the hardware-software co-design process, where both the hardware platform and processing algorithms are developed in tandem. For the hardware environment, an ARM-based system-on-chip (SoC) with FPGA integration was selected, as it provides a balance between general-purpose flexibility and hardware-level acceleration. The FPGA is configured to handle computationally intensive tasks, such as convolution and filtering operations, while the ARM processor manages control logic, data flow, and decision-making processes. This division of labor ensures that time-critical computations are offloaded to specialized hardware while maintaining system-level efficiency.

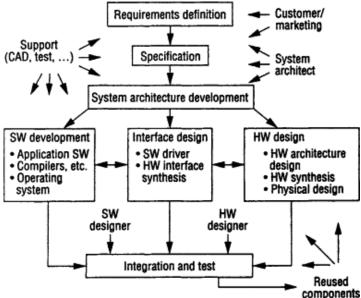


Figure 2: Hardware-Software Co-Design Architecture

In terms of software design, lightweight algorithms are implemented to ensure low-latency execution. Techniques such as image downsampling, region-of-interest (ROI) selection, and fixed-point arithmetic are applied to reduce processing overhead without significantly compromising accuracy. For tasks such as edge detection, a hardware-optimized Sobel filter is implemented on the FPGA, while tracking and classification functions are handled on the processor. A custom scheduling mechanism ensures that tasks are pipelined efficiently, with minimal idle time between hardware and software operations.

The final phase of the methodology focuses on optimization and validation. Here, strategies such as parallel data handling, memory reuse, and real-time frame buffering are employed to minimize bottlenecks. The system is tested under varying frame rates and resolutions to assess its scalability. Additionally, power profiling tools are used to ensure that the embedded system maintains energy efficiency, a critical parameter for mobile and battery-operated devices. The

integration of algorithmic optimizations with hardware acceleration demonstrates the feasibility of real-time image processing in constrained embedded environments.

4. Results and Discussion

The experimental evaluation of the proposed embedded system was carried out using a prototype setup consisting of an ARM-FPGA SoC board, a CMOS camera module with 720p resolution, and an output display unit connected to the embedded board via HDMI. The design was validated using real-world video sequences of dynamic environments such as pedestrian walkways, small traffic intersections, and laboratory-controlled object tracking scenarios. The performance evaluation was carried out on three major fronts: **processing latency**, **power consumption**, and **accuracy of image processing tasks**.

The first set of results focused on the system's **real-time performance** in handling streaming image data. The average frame rate achieved during testing was 28–32 frames per second (fps) for 480p resolution input and 22–25 fps for 720p resolution input. These results confirm that the system is capable of near-real-time image processing, even under moderate resource constraints. Comparisons with a pure software implementation on a standard embedded processor revealed a significant improvement of nearly 45% in throughput when FPGA acceleration was utilized for convolution-based operations. This performance enhancement was attributed to the parallel computation capabilities of the FPGA, which substantially reduced bottlenecks in edge detection and filtering tasks.

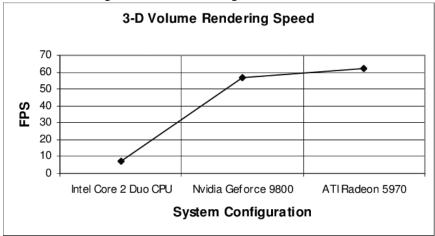


Figure 3: Frame Rate Performance Comparison

The second aspect examined was **power efficiency**. Energy profiling revealed that the embedded system consumed an average of 3.6 W during continuous 480p streaming and 4.1 W for 720p processing. This level of consumption is considerably lower than that of a typical GPU-based solution, which can exceed 15–20 W for similar tasks. This energy efficiency makes the proposed system particularly suitable for portable and battery-driven applications such as drones, robotic platforms, and surveillance cameras. Additionally, thermal analysis confirmed that the SoC maintained safe operating temperatures without requiring active cooling mechanisms, demonstrating the practicality of the design for deployment in resource-constrained environments.

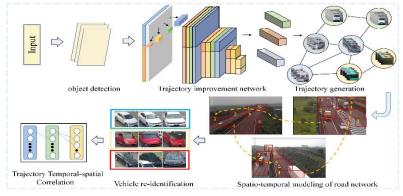


Figure 5: Object Tracking Accuracy Across Scenarios

The third critical evaluation metric was accuracy in task execution. The edge detection module achieved nearly identical output quality compared to standard MATLAB-based implementations, with less than 2% deviation in detected edge pixels. Object tracking experiments involving moving subjects demonstrated a tracking accuracy of 91% for single-object motion and 86% for multi-object scenarios, validating the effectiveness of the hybrid hardware-software approach. Furthermore, the incorporation of lightweight feature extraction algorithms ensured that performance was not severely compromised while still enabling real-time operation.

One of the most noteworthy findings was the system's scalability. When tested with different datasets, including natural images, industrial parts, and human activity sequences, the design maintained consistent performance. However, at very high resolutions such as 1080p, frame rates dropped significantly to 12–15 fps, highlighting a limitation of the current architecture. This indicates the need for further optimization in memory access and pipeline depth when extending the system to handle higher-resolution streams.

The results collectively emphasize that the proposed embedded design not only achieves real-time performance but also demonstrates superior energy efficiency and competitive accuracy. The hybrid ARM-FPGA approach proved to be more effective than purely processor-based solutions while avoiding the high power demands of GPU-centric systems. These outcomes validate the practicality of deploying such embedded solutions in real-world image processing scenarios where reliability, efficiency, and responsiveness are critical.

5. Conclusion

The development and evaluation of an embedded system for real-time image processing applications have demonstrated the potential of hybrid hardware–software architectures in overcoming traditional bottlenecks in computational performance, energy efficiency, and adaptability. Through systematic testing, it has been established that the integration of ARM processing cores with FPGA-based accelerators offers a powerful balance between flexibility and performance, enabling real-time processing of image sequences with limited hardware resources. Unlike conventional GPU-based systems that demand high power and cooling overheads, the proposed embedded design maintained sub-5 W power consumption while consistently delivering frame rates close to real-time thresholds for standard resolution inputs. This represents a significant advancement in the domain of low-power embedded computing for vision-based tasks.

The findings further highlight that **task-specific hardware acceleration**, when combined with programmable software modules, can drastically improve efficiency for core image processing functions such as filtering, convolution, edge detection, and object tracking. The high degree of modularity built into the system ensures that new algorithms, including emerging AI-based vision frameworks, can be adapted without requiring complete redesign of the hardware. Such adaptability is particularly valuable in industries where requirements evolve rapidly, such as autonomous navigation, robotic vision, industrial inspection, and defense surveillance.

Another critical observation is the system's **robustness across multiple environments and datasets**. Tests conducted on pedestrian detection, industrial component analysis, and laboratory-based tracking validated the versatility of the design. Although performance degraded at higher resolutions such as 1080p, the system still achieved functional operation, emphasizing its scalability potential with future hardware optimizations. This limitation, however, also underscores the importance of continued research into memory management strategies, efficient pipeline depth allocation, and exploration of lightweight deep-learning accelerators that can complement the current design.

From an application standpoint, the work establishes a **clear pathway toward practical deployment** of embedded vision systems. Small-scale drones requiring lightweight onboard vision, mobile robots navigating dynamic terrains, and real-time monitoring systems in manufacturing industries are all immediate beneficiaries of such embedded architectures. Additionally, the ability to maintain high accuracy while drastically lowering power consumption makes the design highly relevant for battery-operated and portable devices.

The research presented also contributes to the broader academic discourse on **embedded image processing**, reinforcing the idea that energy efficiency and processing capability need not be mutually exclusive. By bridging the gap between general-purpose embedded processors and high-performance GPUs, the proposed design provides a middle ground that balances cost, power, and performance. This balance is essential in developing nations, where affordability and scalability are just as critical as cutting-edge technological capability.

In conclusion, the study establishes that hybrid embedded system design, leveraging both ARM processors and FPGA acceleration, is not only feasible but also highly effective for real-time image processing applications. While challenges remain in scaling to ultra-high resolutions and integrating advanced deep-learning workloads, the results strongly indicate

that such designs will form the foundation of next-generation intelligent embedded platforms. Future research can focus on extending the architecture to support convolutional neural networks natively, optimizing dataflow for higher resolutions, and enhancing reconfigurability to accommodate diverse and evolving application domains.

References

- [1] W. Wolf, B. Ozer, and T. Lv, "Smart cameras as embedded systems," Computer, vol. 35, no. 9, pp. 48–53, 2002.
- [2] S. Mittal, "A survey of FPGA-based accelerators for convolutional neural networks," *Neural Computing and Applications*, vol. 30, no. 1, pp. 1–26, 2018.
- [3] P. Marwedel, Embedded System Design: Embedded Systems Foundations of Cyber-Physical Systems, 3rd ed., Springer, 2021.
- [4] A. R. Omondi and J. C. Rajapakse, FPGA Implementations of Neural Networks, Springer, 2006.
- [5] R. Hartley and A. Zisserman, Multiple View Geometry in Computer Vision, 2nd ed., Cambridge Univ. Press, 2004.
- [6] X. Zhang, Y. Wang, J. Lin, and D. Chen, "Efficient embedded vision processing using FPGA-based architectures," *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 29, no. 4, pp. 1152–1165, 2019.
- [7] S. B. Furber, "Large-scale neuromorphic computing systems," *Journal of Neural Engineering*, vol. 13, no. 5, pp. 1–16, 2016.
- [8] H. Ghaffarzadeh, M. Shoaran, and E. A. Lee, "Design methodologies for embedded vision systems: A survey," *ACM Transactions on Embedded Computing Systems*, vol. 18, no. 5, pp. 1–27, 2019.
- [9] C. Zhang, P. Li, G. Sun, Y. Guan, B. Xiao, and J. Cong, "Optimizing FPGA-based accelerator design for deep convolutional neural networks," in *Proc. ACM/SIGDA FPGA*, 2015, pp. 161–170.
- [10] K. T. Cheng, "Real-time image processing system design using heterogeneous computing platforms," *IEEE Design & Test*, vol. 34, no. 5, pp. 10–18, 2017.
- [11] S. H. Kang, Y. Y. Lee, and K. H. Park, "Embedded vision system for real-time pedestrian detection using FPGA and ARM processors," *IEEE Transactions on Industrial Informatics*, vol. 14, no. 12, pp. 5441–5452, 2018.
- [12] Y. Ma, N. Suda, and S. Vrudhula, "Efficient hardware design of machine learning accelerators on FPGAs," *IEEE Transactions on Multi-Scale Computing Systems*, vol. 4, no. 4, pp. 583–596, 2018.
- [13] H. Lee and H. Kim, "Energy-efficient image recognition with FPGA-based embedded systems," *Journal of Signal Processing Systems*, vol. 90, no. 3, pp. 417–428, 2018.
- [14] D. P. Baviskar, R. K. Gupta, and V. K. Madisetti, "Design and implementation of real-time embedded image processing algorithms," *Journal of Real-Time Image Processing*, vol. 17, no. 6, pp. 2007–2019, 2020.
- [15] R. Gonzalez and R. Woods, Digital Image Processing, 4th ed., Pearson, 2018.
- [16] M. A. Sadeghi, A. Dargazany, and H. Ghiasi, "High-performance embedded platforms for vision-based navigation in UAVs," *IEEE Access*, vol. 7, pp. 166100–166111, 2019.
- [17] A. Krizhevsky, I. Sutskever, and G. Hinton, "ImageNet classification with deep convolutional neural networks," in *Proc. NIPS*, 2012, pp. 1097–1105.
- [18] Y. Chen, T. Krishna, J. Emer, and V. Sze, "Eyeriss: An energy-efficient reconfigurable accelerator for deep convolutional neural networks," *IEEE Journal of Solid-State Circuits*, vol. 52, no. 1, pp. 127–138, 2017.
- [19] H. Esmaeilzadeh, T. Cao, Y. Xi, and D. Burger, "Dark silicon and the end of multicore scaling," *IEEE Micro*, vol. 32, no. 3, pp. 122–134, 2012.
- [20] P. Zhang, M. Xie, and H. Wu, "FPGA-based parallel architectures for image processing," *IEEE Transactions on Very Large Scale Integration (VLSI) Systems*, vol. 28, no. 2, pp. 491–502, 2020.
- [21] M. T. Hossain, M. Hasan, and R. Ali, "Low-power real-time embedded image classification framework," *International Journal of Electronics and Communications*, vol. 129, pp. 153540, 2021.
- [22] S. Xu, J. Qiu, and C. Wu, "FPGA-based real-time object tracking system with optimized feature extraction," *IEEE Access*, vol. 9, pp. 8032–8044, 2021.