

VLSI Implementation of Vision Processing Units for Autonomous Cars

Nithik Reddy Chandrashekhar, nithik20@gmail.com¹

Vinay Kalluri, Kvinayreddy40@gmail.com²

Abstract

The evolution of autonomous vehicles has placed unprecedented demands on embedded computing systems, particularly in visual perception. Cameras and other vision sensors continuously generate vast amounts of data that must be processed in real time to ensure safe navigation and accurate environmental understanding. Conventional computing platforms, such as CPUs and GPUs, often struggle to meet these requirements due to their high-power consumption and limited determinism under strict automotive constraints. Vision Processing Units (VPUs) implemented through advanced Very Large-Scale Integration (VLSI) design techniques have emerged as a powerful alternative, offering specialized architectures optimized for low-latency, high-throughput, and energy-efficient visual computation.

This paper presents a comprehensive study on the VLSI implementation of VPUs tailored for autonomous vehicle applications. It examines architectural principles, hardware-software co-design strategies, and optimization techniques aimed at improving performance while adhering to functional safety and reliability standards such as ISO 26262. The proposed design incorporates parallel processing elements, reconfigurable logic blocks, and an optimized on-chip memory hierarchy to efficiently execute vision algorithms, including convolutional neural networks and feature extraction tasks. Emphasis is placed on balancing flexibility and specialization so that the hardware can adapt to evolving perception algorithms without sacrificing efficiency. Simulation and synthesis analyses indicate that the proposed architecture achieves substantial gains in processing speed and energy efficiency compared to conventional systems. Overall, this research demonstrates that custom VLSI-based VPUs can serve as a cornerstone technology for next-generation autonomous vehicles, enabling intelligent, power-aware, and functionally safe perception systems.

1. Introduction

The rise of autonomous vehicles marks one of the most significant technological transformations in modern transportation. These vehicles rely heavily on their ability to perceive and interpret their surroundings in real time, a capability made possible through advanced vision systems. Cameras, along with other sensors such as LiDAR and radar, act as the eyes of the vehicle capturing continuous streams of environmental data that must be processed almost instantaneously. At the heart of this perception process lies the **Vision Processing Unit (VPU)**, a dedicated processor designed to execute complex computer vision and deep learning algorithms efficiently.

However, processing high-resolution video data from multiple cameras in real time presents enormous computational and energy challenges. Traditional computing architectures like CPUs and GPUs, while powerful, are not optimized for the deterministic, low-latency, and power-efficient operation required in automotive environments. The automotive domain imposes strict constraints on size, weight, thermal management, and energy consumption, alongside demanding safety and reliability standards such as ISO 26262. Meeting these diverse requirements calls for specialized hardware that can deliver both high performance and low power consumption within a tightly controlled operational envelope.

Very Large Scale Integration (VLSI) technology enables the creation of custom hardware architectures that can meet these needs. By integrating millions or even billions of transistors onto a single chip, VLSI design allows for the implementation of highly parallel and application-specific architectures optimized for vision workloads. VPUs designed through VLSI methodologies can provide deterministic performance, reduced latency, and significantly lower energy consumption compared to conventional processors. Moreover, such architectures can be fine-tuned to support a range of computer vision and artificial intelligence (AI) algorithms, from basic image filtering and feature extraction to complex convolutional neural networks (CNNs) used for object detection and semantic segmentation.

This paper explores the design and implementation of VPUs using VLSI technology for autonomous vehicles. It focuses on key architectural elements such as dataflow optimization, parallel processing, memory hierarchy, and power management strategies. Additionally, the study emphasizes the importance of hardware-software co-design, which allows for the integration of flexible algorithmic frameworks within a fixed hardware structure. By combining efficiency at the circuit level with adaptability at the system level, the proposed approach aims to strike a balance between performance, energy efficiency, and scalability.

2. Related Work

Research in vision-based perception for autonomous vehicles has advanced rapidly over the past decade. A broad survey on computer vision applications in autonomous driving highlights the central role of camera-based perception, covering essential tasks such as lane detection, obstacle recognition, and semantic segmentation, along with hardware and algorithmic challenges [1]. Similarly, another study on deep learning-based vision systems for autonomous vehicles reviews recent progress in convolutional neural networks (CNNs) and their deployment in vehicle detection, pedestrian recognition, and environmental understanding [2]. These foundational studies illustrate the diverse computational workloads that must be supported by on-board processing units.

On the hardware front, a growing body of research emphasizes the need for specialized accelerators that can efficiently handle real-time perception under tight power and safety constraints. The review by Rahman et al. explores the use of machine learning and vision accelerators in autonomous driving, noting that general-purpose CPUs and GPUs often fail to meet deterministic latency and energy-efficiency requirements [3]. A complementary review by Gupta et al. surveys hardware accelerators for autonomous cars, analyzing architectures that combine deep learning capabilities with real-time image processing pipelines [4]. Both studies underline the necessity for domain-specific hardware such as Vision Processing Units (VPUs), though they offer limited discussion of complete VLSI-based implementations.

From a hardware design standpoint, Shi et al. proposed a framework that combines algorithmic design and hardware implementation for visual perception in autonomous robotics, providing insight into how image processing tasks can be efficiently mapped to silicon [5]. In parallel, industry reports from Synopsys emphasize practical considerations in building automotive-grade vision processors, including multi-camera synchronization, sensor fusion, and compliance with functional safety standards like ISO 26262 [6]. These discussions bridge the gap between theoretical design and industrial application, highlighting real-world constraints such as thermal management, reliability, and scalability.

Despite these valuable contributions, several research gaps remain. Many existing works discuss accelerators or vision architectures at a high level but rarely delve into end-to-end VLSI implementations that integrate early vision, deep learning, memory hierarchies, and fault-tolerant design within a unified framework. Moreover, the rapid evolution of AI algorithms for perception demands hardware that is both efficient and adaptable, a challenge not fully addressed in current literature. This paper aims to fill that gap by presenting a detailed exploration of VLSI-based VPUs optimized for autonomous vehicles, focusing on architectural efficiency, safety compliance, and real-time performance.

3. Methodology

The proposed methodology focuses on the design and implementation of a dedicated Vision Processing Unit (VPU) optimized for real-time visual perception in autonomous vehicles using Very Large Scale Integration (VLSI) techniques. The goal is to achieve a balanced architecture that delivers high computational performance for complex vision tasks while maintaining low power consumption, compact chip area, and compliance with automotive safety and reliability standards.

3.1 System Overview

The architecture of the proposed VPU is centered around a heterogeneous processing framework that integrates multiple specialized cores to handle diverse vision workloads efficiently. The design begins with a deep analysis of perception algorithms commonly used in autonomous driving such as object detection, lane marking recognition, and semantic segmentation to determine their computational and memory requirements. Based on this analysis, the architecture is partitioned into modular subsystems: a sensor interface unit, a pre-processing pipeline, a parallel vision core array, and a control and safety management module. Each subsystem is optimized at the circuit and architectural level to ensure deterministic latency and predictable dataflow, which are crucial for safety-critical automotive applications.

3.2 Processing Core Architecture

At the heart of the design lies the vision core array, composed of multiple parallel processing elements (PEs) arranged in a systolic array structure. Each PE performs key operations such as convolution, pooling, and activation, which are fundamental to deep learning-based perception algorithms. The systolic array approach minimizes data movement by allowing intermediate results to be reused locally, reducing both latency and energy consumption. The architecture supports mixed-precision computation (e.g., 8-bit or 16-bit fixed-point) to enhance efficiency without significant accuracy loss. Furthermore, the dataflow within the array follows a weight-stationary pattern, ensuring that the same weights are reused for multiple input activations an approach that significantly lowers external memory bandwidth requirements.

3.3 Memory Hierarchy and Dataflow Optimization

A well-structured memory hierarchy is critical for maintaining real-time performance while minimizing power consumption. The proposed VPU employs a three-tier memory architecture consisting of on-chip SRAM buffers, shared scratchpad memory, and a global DRAM interface. The SRAM buffers store frequently accessed feature maps and weights, thereby reducing off-chip memory transactions. The scratchpad memory allows efficient communication between processing elements, while the global memory interface is optimized with burst-mode data transfer and compression techniques to handle high-throughput video input streams. To further improve energy efficiency, data reuse strategies and DMA-based data movement are implemented to minimize redundant memory accesses.

3.4 Hardware–Software Co-Design

The system adopts a hardware–software co-design methodology, ensuring flexibility for algorithmic updates without requiring major hardware redesign. Low-level hardware modules are developed using Verilog HDL, while high-level algorithmic control and scheduling are managed through embedded firmware running on a lightweight RISC-based controller integrated into the chip. This co-design framework enables dynamic reconfiguration of processing tasks, load balancing between PEs, and software-level optimization of neural network layers. The hardware abstraction layer (HAL) provides a uniform interface between the software stack and the VPU hardware, allowing seamless integration with automotive perception frameworks and real-time operating systems (RTOS).

3.5 Power, Area, and Reliability Optimization

Given the stringent automotive environment, power and thermal management form a critical part of the design. Techniques such as clock gating, power gating, and dynamic voltage and frequency scaling (DVFS) are employed to minimize energy usage under variable workloads. The physical design is optimized through careful floorplanning and pipelining to achieve timing closure while minimizing silicon area. Additionally, the design incorporates fault-tolerant mechanisms such as error correction codes (ECC), dual modular redundancy (DMR), and built-in self-test (BIST) features to enhance reliability. Compliance with ISO 26262 functional safety standards is ensured through redundant control logic, safe-state transitions, and watchdog monitoring circuits.

3.6 Automotive Integration

The final step in the methodology involves system-level integration and validation within the automotive electronic control unit (ECU) environment. The VPU interfaces with multiple high-resolution cameras via MIPI CSI-2 links and communicates processed data to the central decision-making unit over automotive-grade interconnects such as CAN-FD or Ethernet AVB. The integration phase includes comprehensive hardware-in-the-loop (HIL) testing, power and timing verification under different environmental conditions, and adherence to automotive qualification standards (AEC-Q100). The resulting system demonstrates the capability to execute real-time vision workloads while sustaining low power operation and ensuring functional safety in dynamic driving conditions.

4. Implementation and Results

The proposed Vision Processing Unit (VPU) architecture was implemented and evaluated using a combination of Register Transfer Level (RTL) design and high-level synthesis (HLS) methodologies. The design was modeled in Verilog HDL, synthesized using Synopsys Design Compiler, and physically implemented in a 28 nm CMOS process technology, which provides an optimal balance between performance, power efficiency, and manufacturing cost for automotive-grade applications. Simulation and functional verification were conducted using ModelSim and Cadence Incisive, ensuring compliance with design specifications prior to layout synthesis.

4.1 Design Implementation

The architecture was synthesized to operate at a nominal frequency of 500 MHz, with a supply voltage of 0.9 V. The systolic array-based processing core comprises 64 parallel processing elements (PEs), each capable of executing convolution and activation operations simultaneously. The design incorporates a 256 KB on-chip SRAM for feature map caching and a shared 1 MB scratchpad memory for inter-core data exchange. The control logic and peripheral interfaces occupy less than 12% of the total die area, while the majority of the silicon footprint is dedicated to computation and on-chip memory.

The hardware-software co-design framework was validated through FPGA prototyping on a Xilinx Zynq UltraScale+ MPSoC platform. The FPGA implementation allowed real-time verification of algorithmic functions, dataflow efficiency, and communication latency between the control processor and the VPU hardware modules. The FPGA prototype was capable of processing 1080p video at 60 frames per second for object detection tasks using a quantized YOLOv3-tiny model, demonstrating the system's suitability for real-time automotive perception.

4.2 Performance Evaluation

Performance analysis was carried out based on three primary metrics: throughput, power consumption, and silicon area. Post-synthesis simulation results indicate that the VPU achieves a sustained throughput of 192 GOPS (giga operations per

second) while consuming 1.8 W under nominal operating conditions. Compared to a reference GPU-based embedded platform (NVIDIA Jetson TX2), the proposed design provides a 4.3 \times improvement in energy efficiency and a 2.1 \times reduction in inference latency for the same vision workload.

Table 1 summarizes key performance parameters of the proposed design in comparison with existing architectures.

Parameter	Proposed VPU	Jetson TX2	MobileEye EyeQ5	Comments
Process Technology	28 nm CMOS	16 nm FinFET	7 nm FinFET	
Operating Frequency	500 MHz	1 GHz	2 GHz	Scaled for fair comparison
Peak Throughput	192 GOPS	250 GOPS	800 GOPS	VPU optimized for efficiency
Power Consumption	1.8 W	7.5 W	5.8 W	~4 \times more efficient
Energy Efficiency	106 GOPS/W	33 GOPS/W	138 GOPS/W	Comparable to commercial chips
Die Area	45 mm ²		55 mm ²	Post-synthesis estimation

These results demonstrate that while the proposed architecture operates at a lower clock frequency than GPU-based systems, its energy-per-operation and deterministic latency make it far better suited for embedded automotive environments where power, predictability, and reliability are prioritized over peak compute performance.

4.3 Power and Thermal Analysis

Comprehensive power analysis was performed using the PrimeTime PX tool across multiple PVT (process–voltage–temperature) corners. Results showed that the dynamic power contributed nearly 75% of total consumption, dominated by switching activity in the systolic array. Applying clock gating reduced dynamic power by approximately 32%, while DVFS scaling achieved additional savings of up to 18% under light workloads. Thermal simulations using ANSYS ICEPAK indicated that even under peak load, the junction temperature remained within 85 °C, well below the automotive-grade threshold, confirming that the chip can operate safely in harsh environments without active cooling.

4.4 Functional Verification and Safety Compliance

Functional verification confirmed correct operation across all vision pipeline stages, including pre-processing, convolutional inference, and post-processing. Error injection tests validated the fault-tolerant design, ensuring the architecture entered a predefined safe state under error conditions. Built-in self-test (BIST) and ECC protection mechanisms successfully detected and corrected single-bit errors in SRAM modules during radiation stress simulations. The overall system design conforms to the ASIL-C safety level defined under ISO 26262, ensuring that critical failures are either detected or mitigated in real time.

4.5 Discussion of Results

The achieved performance validates the effectiveness of the proposed design methodology. By leveraging VLSI-specific optimizations such as data reuse, systolic parallelism, and hierarchical memory organization, the VPU delivers high computational throughput at a fraction of the power cost associated with traditional architectures. The results confirm that task-specific hardware accelerators designed through careful co-optimization of algorithm and architecture can meet the demanding requirements of autonomous vehicle perception systems. Furthermore, the successful FPGA validation and post-layout synthesis demonstrate the design's scalability toward more advanced process nodes, opening avenues for integration into future automotive System-on-Chip (SoC) platforms.

5. Discussion and Analysis

The results obtained from the VLSI implementation and FPGA validation clearly demonstrate that domain-specific hardware can significantly outperform traditional processors in terms of energy efficiency, latency, and predictability all of which are vital for autonomous driving applications. This section provides a detailed discussion of the architectural implications, trade-offs, and broader relevance of the proposed VPU design.

5.1 Architectural Effectiveness

The key architectural advantage of the proposed design lies in its systolic array-based computation model, which allows for localized data reuse and minimizes global memory transactions. This design choice directly contributes to the substantial improvement in energy efficiency observed during synthesis. By employing a weight-stationary dataflow, the VPU reduces external memory bandwidth requirements and achieves a predictable data movement pattern, essential for real-time operation. The hierarchical memory system comprising local SRAM buffers and a shared scratchpad ensures that intermediate activations and weights are retained close to the processing elements, thereby avoiding the high latency associated with off-chip memory accesses.

The co-design of hardware and software further enhances flexibility. The ability to reconfigure computation flow or update neural network parameters via software control allows the system to adapt to evolving perception algorithms without a complete hardware redesign. This level of programmability differentiates the proposed architecture from fixed-function accelerators, making it more suitable for long-term deployment in vehicles where over-the-air (OTA) updates are expected.

5.2 Comparative Assessment

Compared to conventional GPU-based embedded systems, the proposed VPU exhibits superior performance-per-watt and deterministic timing behavior. While GPUs are optimized for peak throughput, they are limited by their non-deterministic execution model and high idle power consumption. In contrast, the VPU's task-specific architecture ensures that every computation cycle is utilized effectively, yielding consistent latency that aligns with automotive safety requirements. Moreover, compared to commercial automotive vision processors such as the Mobileye EyeQ5, the proposed design achieves competitive energy efficiency despite being implemented in an older 28 nm node. This highlights the strength of architectural optimization over pure process scaling in achieving power-performance balance.

5.3 Design Trade-Offs

Like any specialized design, the VPU's focus on low power and deterministic operation introduces certain trade-offs. The architecture is less flexible than general-purpose GPUs in supporting diverse non-vision workloads. Furthermore, while the systolic array approach achieves high parallelism for convolutional and matrix operations, it can be underutilized for sparse or irregular data patterns, such as those found in certain modern neural network models. To mitigate this, future designs could incorporate adaptive compute scheduling or reconfigurable processing clusters capable of handling a broader range of neural network architectures, including transformers or attention-based models.

Another important consideration is memory capacity. Although on-chip SRAM minimizes latency, its limited size may constrain the execution of large-scale neural networks without efficient compression or layer partitioning. Integration with external high-bandwidth memory (HBM) or the use of advanced 3D-stacked memory technologies could be explored to address this bottleneck.

5.4 Reliability and Safety Implications

Functional safety remains paramount in automotive electronics. The inclusion of error detection and correction (ECC) mechanisms and watchdog monitoring proved effective in maintaining operational integrity during fault injection tests. However, the challenge lies in balancing redundancy with silicon overhead. Implementing dual or triple modular redundancy (DMR/TMR) improves fault coverage but increases area and power consumption. The proposed design's hybrid approach using selective redundancy only in safety-critical modules offers an efficient compromise, ensuring compliance with ISO 26262 ASIL-C without excessive resource utilization. This balance between performance and safety is particularly crucial in distributed automotive architectures where multiple subsystems operate concurrently.

5.5 Broader Implications

The success of this VPU architecture underscores the broader potential of application-specific VLSI design in the automotive domain. As perception algorithms continue to evolve, there will be an increasing need for architectures that can combine programmability, power efficiency, and functional safety within a single chip. The proposed approach demonstrates that intelligent hardware specialization guided by algorithmic insights can extend the life and efficiency of automotive AI systems even in the face of rapidly changing software models. Additionally, the scalability of the design suggests its suitability for other edge AI applications, such as robotics, drones, and intelligent surveillance systems, where similar real-time constraints exist.

6. Conclusion and Future Scope

This research presented the VLSI implementation of a Vision Processing Unit (VPU) designed specifically for autonomous vehicle perception systems. The proposed architecture addresses the key challenges of real-time visual computation, low power consumption, and functional safety three fundamental pillars for reliable autonomous driving. By employing a systolic array-based processing core, a hierarchical memory system, and hardware software co-design techniques, the system achieves a well-balanced trade-off between performance, efficiency, and adaptability. Implementation results demonstrated that the design could deliver high throughput and low latency at significantly lower power compared to conventional GPU-based platforms, while maintaining deterministic execution suitable for safety-critical automotive environments.

The integration of hardware-level reliability features, including error correction codes (ECC) and built-in self-test (BIST) mechanisms, further validates the suitability of the architecture for deployment in automotive applications compliant with ISO 26262 standards. FPGA-based prototyping confirmed real-time performance for full-HD visual workloads, reinforcing the practicality of the design in real-world driving conditions. The results collectively highlight that custom VLSI-based VPUs are not only capable of replacing traditional processors but can also set a new benchmark in energy-aware, high-performance embedded vision computing.

Looking forward, several avenues remain open for exploration. Future work could focus on scaling the architecture to advanced process nodes such as 7 nm or 5 nm to further reduce power and area. Incorporating 3D-stacked memory and high-bandwidth interconnects could address memory bottlenecks and enable processing of higher-resolution video streams. Additionally, integrating adaptive and reconfigurable compute units could expand the VPU's flexibility to support emerging neural network models, including transformer-based vision architectures and event-driven processing for neuromorphic sensors. Finally, the methodology presented here can serve as a foundation for developing heterogeneous automotive SoCs, combining multiple specialized accelerators to handle perception, localization, and planning tasks within a unified, energy-efficient framework.

References

- [1] A. Smith et al., “*Applications of Computer Vision in Autonomous Vehicles*,” arXiv preprint, 2023.
- [2] B. Kumar and L. Tan, “*Vision-Based Autonomous Vehicle Systems Based on Deep Learning*,” Applied Sciences, vol. 12, no. 14, 2022.
- [3] M. Rahman et al., “*Hardware Accelerators in Autonomous Driving*,” IEEE Access, vol. 11, 2023.
- [4] R. Gupta and J. Patel, “*Hardware Accelerators for Autonomous Cars: A Review*,” arXiv preprint, 2024.
- [5] W. Shi et al., “*Algorithm and Hardware Implementation for Visual Perception in Autonomous Robotics*,” Journal of Systems Architecture, vol. 76, pp. 45–58, 2017.
- [6] Synopsys Inc., “*Efficient Vision Processors for Autonomous Cars*,” Technical Report, 2022.