

UNIVERSITÀ DEGLI STUDI DI SIENA

DIPARTIMENTO DI INGEGNERIA DELL'INFORMAZIONE E SCIENZE MATEMATICHE



**UNIVERSITÀ
DI SIENA**
1240

**Design and Evaluation of FPGA Implementations of
Deep Learning Models: From Real-Time Object
Detection to Graph Convolutional Networks**

Seyed Hani Hozhabr

PhD in Information Engineering and Science

Supervisor

Prof. Roberto Giorgi

Examination Committee

Prof.

Prof.

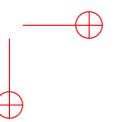
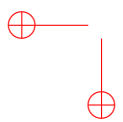
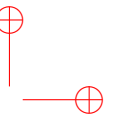
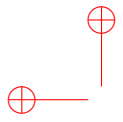
Prof.

Thesis reviewers

Prof.

Prof.

SIENA, 15/11/2025



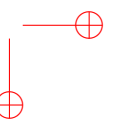
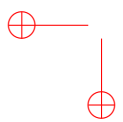
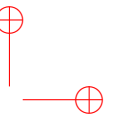
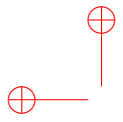
Abstract

This thesis investigates the design and evaluation of FPGA-based implementations for accelerating deep learning workloads, with a particular focus on real-time object detection and Graph Convolutional Networks (GCNs). The first part of the work presents the first dedicated and comprehensive survey of real-time object detection on FPGAs, analyzing state-of-the-art architectures, acceleration techniques, and optimization strategies across more than a decade of research. To enable fair comparison among heterogeneous implementations, the study introduces pixel throughput as a resolution-independent performance metric and provides an in-depth classification of algorithmic and hardware design approaches, together with a discussion of the challenges and emerging trends shaping next-generation FPGA-based detection systems.

Building on the insights gained from this survey, the second contribution evaluates a practical deployment of Tiny YOLOv2 on an AMD Xilinx Zynq UltraScale+ MPSoC. The implemented heterogeneous system achieves 30 FPS real-time performance while consuming only 1.1 W, outperforming GPU-based implementations (15 FPS at 13.5 W) and CPU baselines (2.5–8 FPS at 2.1–5.3 W) by large margins in both energy efficiency and performance–power ratio. A detailed profiling study using the Vitis AI Profiler reveals key execution bottlenecks—including CPU–DPU synchronization stalls and underutilized DDR bandwidth—offering quantitative insights into system-level inefficiencies and guiding future optimization efforts for FPGA-based MPSoC deployments.

The final part of the thesis introduces a lightweight, fully streaming FPGA accelerator for Graph Convolutional Networks, designed for mid-range and power-constrained devices. The architecture integrates scalable sparse–dense (SpMM) units, an efficient systolic array for dense transformations, and a dataset-aware 16-bit fixed-point quantization strategy that preserves full classification accuracy. Implemented on a Kintex-7 FPGA, the accelerator achieves 1.77–3.47 ms inference latency on the Cora and Citeseer datasets while consuming only 0.305–0.553 W, delivering up to 1,082 graphs/J, over 50× greater energy efficiency than contemporary GPU implementations, while requiring just 54–55 DSPs. These results highlight its suitability for real-time, embedded, and edge-AI applications where energy and resources are severely limited.

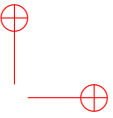
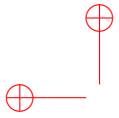
Overall, this dissertation provides a unified and experimentally validated body of methodologies for the FPGA acceleration of deep learning workloads. Through systematic analysis, empirical bottleneck characterization, and novel architectural design, it advances the state of the art in real-time and energy-aware deep learning on reconfigurable platforms and establishes a foundation for future research on efficient heterogeneous AI computing.



Contents

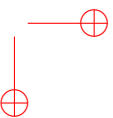
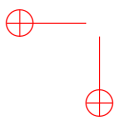
Abstract	i
1 Introduction	1
1.1 Motivation	2
1.2 Research Objectives	3
1.3 Thesis Contributions	3
1.4 Thesis Organization	4
2 Real-Time Object Detection on FPGAs	5
2.1 Introduction	5
2.2 Fundamentals and Background	10
2.2.1 OBJECT DETECTION OVERVIEW	10
2.2.2 Soft and Hard Real-time Constraints	20
2.2.3 Evaluation Metrics and Datasets	21
2.3 FPGA Basics and Applications in Object Detection	25
2.3.1 FPGA Overview	25
2.3.2 Applications of FPGAs in Computer Vision-Object Detection	25
2.3.3 FPGA Platforms for Accelerating Object Detection	29
2.4 FPGA-based Designs	29
2.4.1 FPGA architecture design approaches	31
2.4.2 Hardware Acceleration Techniques for FPGA-based Object Detection	34
2.4.3 Impact and Trade-offs of Acceleration Techniques in FPGA-based Object Detection	46
2.5 Results Analysis and Optimization	47
2.5.1 Case Studies: Adoption of Acceleration Techniques in FPGA-Based Real-Time Object Detection	48
2.5.2 COMPARATIVE ANALYSIS OF FPGA IMPLEMENTATIONS	55
2.5.3 OPTIMIZATION STRATEGIES FOR REAL-TIME PERFORMANCE	55
2.6 Challenges and Future Directions	60
2.6.1 Existing Challenges	60
2.6.2 Possible Future Research Trends	60

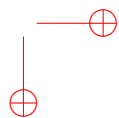
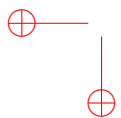
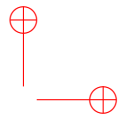
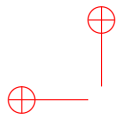
2.7	Final Remarks	62
3	Real-time Object Detection on FPGA-based Heterogeneous MPSoCs: A Preliminary Analysis of the Execution Bottlenecks	63
3.1	Introduction	63
3.2	Case Study: Accelerated YOLOv2 on Zynq Ultrascale+	64
3.2.1	Algorithm Pipeline Overview	64
3.2.2	Pipeline Design Overview	64
3.3	Primary Implementation Results	65
3.4	Execution Profiling and Bandwidth Analysis	67
3.5	Final Remarks	68
4	Design and Validation of a Lightweight, Energy-Efficient FPGA Streaming Architecture for Graph Convolutional Networks (GCNs)	69
4.1	Introduction	69
4.2	Literature Review	70
4.2.1	Symmetry-Aware and Sparse Data Optimizations	71
4.2.2	Quantization and Resource Efficiency	71
4.2.3	Dataflow and Architectural Flexibility	71
4.2.4	Software–Hardware Co-Design Frameworks	71
4.2.5	Comparative Insights and Research Gaps	72
4.3	Motivation	72
4.4	Architecture	72
4.4.1	Overview	72
4.4.2	Sparse-Dense Matrix Multiplication (SpMM) Module	75
4.4.3	Dense-Dense Matrix Multiplication (DDMM) Module	76
4.4.4	Parallel-to-Serial Converter (P2S_Conv) Module	76
4.5	Data-Aware Quantization and Scaling Methodology	78
4.5.1	Profiling-Guided Fixed-Point Configuration	78
4.6	Evaluation Results	79
4.6.1	Experimental Setup	79
4.6.2	Primary Experimental Results	79
4.6.3	Design Space Exploration for SpMM Parallelism	81
4.6.4	Cross-Platform Performance and Energy Efficiency on Benchmark Datasets	81
4.6.5	Comparative Performance Analysis	83
4.7	Final Remarks	85
5	Conclusion and Future Work	87
5.1	Conclusion	87
5.2	Future Work	88
A	Some General Definitions in The Context of Object Detection Metrics	91



Contents

B Overview of FPGA Architecture and Design Approaches	95
B.0.1 FPGA Architecture	95
B.0.2 FPGA Design Approaches	95
C Roofline Model	99
List of Acronyms	101
Bibliography	103





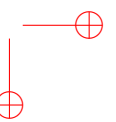
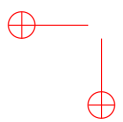
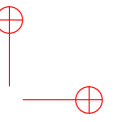
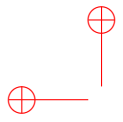
List of Figures

2.1	Increasing number of publications related to real-time object detection and FPGA over the years, especially after the introduction of one-stage object detection models like YOLO [1] and SSD [2], based on data extracted from Google Scholar. . .	6
2.2	Organization of the paper.	9
2.3	Classification of object detection techniques: Non-DNN-based and DNN-based detectors, along with instances of models for each category. It also shows which models may show real-time performance and which are more suitable for FPGA implementation. The focus of this study is the overlapping area (green zone.) . .	11
2.4	Typical structure of a computer vision system.	27
2.5	An example of FPGA co-processing architecture: images are captured by the CPU, and FPGA is used to accelerate the task. In comparison with the typical structure of a CV system shown in Figure 2.4, here, two tasks, i.e., the tasks performed in blocks 3 and 4, are done on the FPGA side. The white blocks represent the tasks run in the software. The green blocks indicate hardware units implemented on FPGA, and the light yellow blocks denote hardware units within the system.	28
2.6	Example of an FPGA accelerator structure with a single computation engine architecture. In this approach, common computation resources are temporarily shared across different layers. A software-based controller is also deployed to control and schedule the operations. (CONV=CONVOLUTION, NONLIN=NON LINEar layer, POOL=POOLing layer) (Figure adapted from [3]).	32
2.7	Example of an FPGA accelerator structure with a streaming architecture. Each block is individually designed and optimized to perform the required calculations of each layer. (Figure adapted from [3])	33
2.8	Example of applying structured and unstructured pruning on a fully connected layer. Removed (pruned) neurons and weights are shown by white circles and dotted lines, respectively.	37

2.9	Illustration of Quantization-Aware Training (QAT) (Left) and Post-Training Quantization (PTQ) (Right) procedures. In QAT, a pre-trained model is quantized and then fine-tuned. In PTQ, a pre-trained model is calibrated using calibration data to do quantization based on the calibration result. (Figure adapted from [4].)	38
2.10	An illustration of performing knowledge distillation with a teacher-student framework. (Figure adapted from [5].)	40
2.11	An example of adopting a systolic array to perform matrix multiplications in CNNs. (Figure adapted from [6])	43
2.12	Possible impact of model-related (top) and implementation-related (bottom) acceleration techniques on key metrics in FPGA-based object detection systems, focusing on throughput and accuracy. The corresponding section where each technique is discussed in this work is also indicated.	46
2.13	Data reuse scheme proposed in [7], minimizing memory accesses by reusing both input data and weights, with intermediate results stored in line buffers. $K \times K \times T_i$ defines the size of the input sliding cube in pixels, and T_0 denotes the number of weight blocks concurrently convolved with the input sliding cubes in each iteration. (Figure adapted from [7]).	49
2.14	An architecture for accelerating YOLOv3 on FPGA. It is based on a pipelined multicore processing architecture consisting of a Control Unit to manage all operations, a Compute Engine with multiple layers of 2D arrays of PEs, and an on-chip memory unit optimized for parallel data access (Figure adapted from [8].)	50
2.15	The effect of applying pruning (p), quantization (Q), and knowledge distillation (KD) techniques one after the other on the size and accuracy of the YOLOv4 model (the numbers extracted from [9]).	52
2.16	Overall accelerator architecture, consisting of a SIMD array, activation, pooling, and control units along with memories to store input, weight, and output values (Figure adapted from [10]).	53
2.17	Overall hardware block diagram of a proposed FPGA accelerator, with two separate DMA modules, each equipped with a dedicated DMA descriptor buffer. The data router will rearrange the pixels and weights to maximize data reuse capability (Figure adapted from in [11]).	54
2.18	Comparison of the reviewed works based on the achieved pixel throughput.	58
2.19	Comparison of the reviewed works based on their reported power efficiency (note: for some works power efficiency data is not available).	58
3.1	Pipeline design of the real-time object detection system, centered around the DPU inference engine, provided by Vitis AI, as the core of the accelerator design on the Kria KV260 board, featuring a Zynq UltraScale+ MPSoC.	64
3.2	Timeline analysis of the application using Vitis AI Profiler, illustrating the interaction between the application running on the PS and the DPU, data transfer between the PS and PL through DDR memory, and the model execution on the DPU. The long "Execution wait" periods can be inferred as a performance bottleneck, caused by inefficient task scheduling.	67

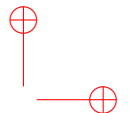
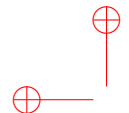
List of Figures

3.3	DDR bandwidth (read/write) utilization in Vitis Analyzer, where each color represents data transfer between a PL module (various DPU ports and the post-processing unit) and the DDR controller via high-performance memory ports. Considering that each port theoretically supports up to 2400 MB/s of total bandwidth (read + write), the observed memory bandwidth utilization suggests potential for optimization.	68
4.1	Overall streaming architecture of the proposed GCN accelerator.	73
4.2	Sparsity analysis of input matrices for Cora, Citeseer, and Pubmed datasets. The blue bars represent the sparsity of the Adjacency Matrix (\mathbf{A}), which drives the Aggregation phase design. The orange bars represent the sparsity of the Feature Matrix (\mathbf{X}), influencing the efficiency of the Combination phase.	74
4.3	Overall block diagram of the Sparse Matrix Multiplication (SpMM) Module. The module consists of an FSM-based control unit coordinating data flow between the memory blocks and the core <i>MAC_Sparse</i> units.	75
4.4	Block diagram of the Dense-Dense Matrix Multiplication (DDMM) Module. The module features a 1D Systolic Array for processing data received from the P2S Conversion stage, using an internal Memory Bank for weight storage.	77
4.5	Energy Efficiency (<i>EE</i>) comparison across CPU, GPU, and the proposed FPGA accelerator, demonstrating the superior <i>EE</i> achieved by the dedicated hardware solution.	82
4.6	Comparison of normalized throughput with SOTA on Cora dataset.	85
A.1	Definition of Intersection over Union (IoU), calculated by dividing the intersection of the predicted bounding box (Bbox) and ground truth (GT) Bbox by the union of them.	91
B.1	An overview of basic FPGA architecture, including programmable logic blocks, Block RAMs (<i>BRAMs</i>), DSP blocks, input/output blocks (<i>I/O</i>), and programmable interconnects (shown by black lines). (Figure adapted from [12]).	96
C.1	The Roofline model.	99



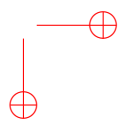
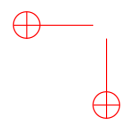
List of Tables

2.1	Comparison of recent related review studies with the present work, focusing on the main topics discussed.	7
2.2	A summary of the discussed models from all categories, including Non-DNN-based Detectors (NDD), DNN-based Two-stage Detectors (DTD), DNN-based One-stage Detectors (DOD), and Transformer-based Detectors (TBD) (Part I). . .	12
2.3	A summary of the discussed models from all categories, including Non-DNN-based Detectors (NDD), DNN-based Two-stage Detectors (DTD), DNN-based One-stage Detectors (DOD), and Transformer-based Detectors (TBD) (Part II). . .	13
2.4	A summary of the commonly-used datasets in the context of object detection . .	22
2.5	A summary of commonly used metrics for evaluating object detection models. . .	23
2.6	Comparing FPGAs with some possible alternatives in different criteria.	26
2.7	A list of some FPGA boards used in recent works for implementing Real-time object detection models, sorted by number of LUTs/ALMs.	30
2.8	A summary of two FPGA architecture design approaches for implementing object detection models.	31
2.9	Commonly used Hardware Acceleration Techniques for FPGA-based Object Detection.	35
2.10	The recent advanced works of FPGA-based real-time object detection implementations (part 1/2).	56
2.11	The recent advanced works of FPGA-based real-time object detection implementations (part 2/2).	57
3.1	A Summary of Resource Utilization of Implementation Results.	65
3.2	Comparative performance and efficiency of the developed OD with Tiny YOLOv2 receiving a live video stream from an attached camera @ 30 FPS as input on CPU, GPU, and KV260 Platforms, highlighting FPS, power consumption (W), and Power Delay Product (PDP) (W.Sec).	66
4.1	Computational complexity under different execution orders [13]	74
4.2	Characteristics of Citation Network Datasets for GCN Evaluation.	79
4.3	FPGA Resource Utilization and Core Performance Metrics.	80



List of Tables

4.4	Numerical Precision and Accuracy Preservation with 16-bit Fixed-Point Quantization.	80
4.5	Design Space Exploration for SpMM Parallelism Scaling with Cora.	81
4.6	Comparison of latency and Energy Efficiency (EE) across platforms on Cora and Citeseer datasets.	82
4.7	Performance and Hardware Resource Utilization of State-of-the-Art FPGA-based Graph Convolutional Network (GCN) Accelerators. The comparison highlights latency (μs), power consumption (W), accuracy (%), frequency (MHz), and the number of consumed resources, using the Cora and Citeseer graph datasets. . . .	84
A.1	Prediction result category in object detection models based on the label provided by the dataset (real label).	92



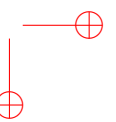
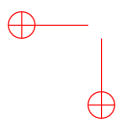
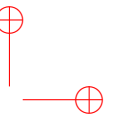
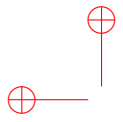


Acknowledgements

I would like to express my gratitude to Prof. Roberto Giorgi, my supervisor, for his guidance, continuous support, and mentorship throughout my doctoral studies. His scientific insight, encouragement, and commitment to excellence have shaped this work and influenced my growth as a researcher.

I gratefully acknowledge Xilinx/AMD for their technical support and for the donation of hardware platforms that were essential to the experimental components of this thesis. This research has been partially supported by the European Commission through the EDGE-ME Project, the AXIOM H2020 Project (Grant 645496), the TERAFLUX Project (Grant 249013), and the HiPEAC Project (Grant 101069836). Additional support was provided by the European Union—NextGenerationEU under the PNRR M4C2-Inv1.4 Italian Research Center on High-Performance Computing, Big-Data and Quantum Computing (cascade funding project EDGE-ME, MUR Grant CN0000013), as well as by the University of Siena, Campora-ES, and the Tuscany Region under the HIPERAIHL Project. Their contributions enabled the successful completion of this work.

Most importantly, I extend my heartfelt appreciation to my parents for their unconditional love, sacrifices, and steadfast belief in my path, and to my wife, whose patience, constant encouragement, and unwavering support sustained me throughout the most challenging stages of this journey. This thesis would not have been possible without them.



Chapter 1

Introduction

The contemporary landscape of computer science is defined by the pervasive integration of deep learning (DL) models across complex systems. From advanced driver-assistance systems (ADAS) and sophisticated industrial automation to large-scale network modeling and bio-informatics, these models, particularly Convolutional Neural Networks (CNNs) and Graph Neural Networks (GNNs), have demonstrated state-of-the-art accuracy [14–19]. However, translating this algorithmic prowess into reliable, real-world deployment remains constrained by architectural limitations inherent in the hardware that must support them. Modern DL models, characterized by millions of parameters and billions of operations per inference, demand computational resources far exceeding the capabilities of conventional, low-power embedded platforms [20–22].

This deployment challenge is multi-faceted, involving severe restrictions on power consumption, the absolute necessity for low and deterministic latency (particularly for hard real-time systems), and the requirement to efficiently process data of varying structural complexity. General-purpose architectures, such as Central Processing Units (CPUs), are typically bottlenecked by sequential instruction processing and insufficient on-chip memory for comprehensive feature reuse. Simultaneously, power-intensive Graphics Processing Units (GPUs) frequently fail to meet the stringent power budgets of battery-operated or thermally limited edge devices, and their non-deterministic execution often precludes their use in safety-critical hard real-time applications.

Addressing this critical architectural disparity requires a fundamental shift toward hardware specialization. This thesis, "Design and Evaluation of FPGA Implementations of Deep Learning Models: From Real-Time Object Detection to Graph Convolutional Networks," champions the Field-Programmable Gate Array (FPGA) as the optimal foundational technology for high-performance, energy-efficient DL inference at the edge. The FPGA's inherent reconfigurability allows for the creation of customized dataflow architectures and on-chip memory hierarchies that perfectly match the computational and memory access patterns of a target DL model. Furthermore, modern Heterogeneous Multi-Processor Systems-on-Chips (MPSoCs), which tightly couple programmable logic (PL) with embedded processor systems (PS), offer a fertile ground for sophisticated hardware/software co-design, enabling the fine-grained scheduling and control necessary for complex real-time applications.

The scope of this research is deliberately dual-pronged and expansive, encompassing two antithetical yet representative deep learning workloads, thereby ensuring the generalizability and robustness of the proposed architectural solutions:

- **Dense, Time-Critical Workloads (Real-Time Object Detection):** The acceleration of CNN-based Object Detection (OD) models (e.g., YOLO and SSD) is governed by dense, compute-intensive matrix operations (convolutions). The core challenge here is achieving maximum

throughput and deterministic latency within heterogeneous MPSoC environments, requiring rigorous analysis of execution bottlenecks such as DDR memory bandwidth and inter-system synchronization overhead.

- **Sparse, Irregular Workloads (Graph Convolutional Networks):** The acceleration of GCNs, designed to process complex, non-Euclidean data, is fundamentally defined by highly irregular memory access patterns stemming from Sparse-Dense Matrix Multiplication (SpMM). This necessitates radically different architectural strategies—specifically, specialized streaming engines designed for maximal energy efficiency on resource-constrained FPGAs, where the traditional metric of peak throughput is subordinated to minimizing the power-delay product.

By systematically investigating these contrasting computational regimes, this thesis develops and validates a body of advanced hardware/software co-design methodologies that provide a comprehensive framework for realizing adaptable, energy-optimal, and real-time deep learning inference across the entire spectrum of AI applications at the edge.

1.1 Motivation

The rapid convergence of advanced Artificial Intelligence (AI) models, particularly deep neural networks, and the proliferation of edge computing platforms has created a critical demand for real-time, energy-efficient processing across diverse application domains. This demand is particularly pronounced in safety-critical and high-throughput scenarios, such as autonomous vehicles, robotics, industrial automation, and large-scale graph analysis.

The architectural foundation of these cutting-edge AI workloads presents significant challenges to conventional computing platforms. Specifically:

- **Computer Vision (CV) Workloads:** Real-time Object Detection (OD) models, such as the YOLO (You Only Look Once) and SSD (Single Shot MultiBox Detector) families, rely heavily on Convolutional Neural Networks (CNNs). While highly accurate, these models are computationally intensive and memory-bound. Achieving the low and deterministic latency required by hard real-time systems—where missing a deadline is unacceptable—is infeasible on power-hungry General-Purpose Processors (CPUs) or Graphics Processing Units (GPUs).
- **Graph Workloads:** Graph Convolutional Networks (GCNs) have become essential for processing complex, non-Euclidean data structures like social networks and knowledge graphs. The core challenge here is the sparse matrix-matrix multiplication (SpMM), which results in irregular memory access patterns and poor data locality. This irregularity leads to severe performance and energy bottlenecks on conventional, dense-optimized hardware.

This research is fundamentally motivated by the need to bridge the gap between the increasing complexity of state-of-the-art AI algorithms and the stringent power, latency, and throughput constraints of embedded and edge-AI systems. The Field-Programmable Gate Array (FPGA) emerges as the ideal platform to resolve this conflict. FPGAs offer high parallelism, the flexibility to design custom data paths, and the ability to implement fine-grained precision control, making

them superior to both CPUs and GPUs in terms of power efficiency and deterministic latency, especially for demanding real-time inference tasks.

1.2 Research Objectives

The overarching goal of this thesis is to systematically investigate, design, and validate novel hardware and software co-design methodologies that enable energy-efficient and real-time execution of complex deep learning models on resource-constrained FPGA platforms.

To achieve this goal, the research is structured around the following specific objectives:

- **Systematic Analysis of Real-Time Challenges:** To conduct a comprehensive survey and comparative analysis of existing FPGA-based solutions for real-time Object Detection, identifying key acceleration techniques (e.g., pruning, quantization, pipelining) and determining the architectural trade-offs essential for meeting hard and soft real-time constraints.
- **Bottleneck Identification and Optimization in Heterogeneous Systems:** To evaluate the practical deployment of energy-efficient Computer Vision models (specifically Tiny YOLOv2) on FPGA-based Heterogeneous Multi-Processor Systems-on-Chips (MPSoCs). This includes profiling the execution pipeline to precisely locate and analyze system bottlenecks, such as DDR memory bandwidth underutilization and inefficient CPU-DPU (Deep Learning Processor Unit) task scheduling, to guide subsequent optimization efforts.
- **Lightweight Architecture Development for Sparse Workloads:** To propose and validate a novel, lightweight, streaming FPGA architecture explicitly tailored for accelerating GCNs. This architecture must effectively mitigate the challenges of sparse computations by integrating scalable SpMM units and a compact dense transformation stage while ensuring high energy efficiency and preserving full classification accuracy via dataset-aware fixed-point quantization.

1.3 Thesis Contributions

This thesis advances the state-of-the-art in energy-efficient and real-time AI acceleration by delivering the following original contributions:

- **The First Dedicated Survey on Real-Time OD on FPGAs:** This work provides the first comprehensive, dedicated review focusing specifically on the implementation and optimization of "real-time Object Detection on FPGAs". It introduces pixel throughput as a fair metric for comparative analysis and synthesizes knowledge on appropriate one-stage CNN architectures, hardware acceleration techniques, and optimization strategies necessary for industrial-grade, low-latency deployments.
- **Quantifiable Performance and Efficiency Benchmarks for MPSoCs:** Through the implementation and evaluation of a Tiny YOLOv2-based real-time OD system on an AMD Xilinx Zynq UltraScale+ MPSoC, this research provides quantitative, cross-platform performance comparisons. It demonstrates that the FPGA-based system achieves significantly superior

energy efficiency (up to 22.7x better Power Delay Product) compared to high-end GPU and CPU baselines, establishing a robust benchmark for edge-AI systems.

- **Preliminary Analysis of Heterogeneous MPSoC Bottlenecks:** A crucial contribution is the preliminary, empirical analysis of the execution bottlenecks within the heterogeneous MPSoC environment using the Vitis AI Profiler. This analysis pinpoints specific architectural and software integration challenges (e.g., inefficient parallelism and memory underutilization) that are critical for achieving peak performance in FPGA-based co-designed systems.
- **A Lightweight, Energy-Efficient GCN Streaming Accelerator:** The thesis proposes and validates a novel lightweight and energy-efficient FPGA accelerator specifically for GCN workloads. Key innovations include:
 - A unified, fully streaming dataflow architecture that removes instruction-driven overhead, directly mapping GCN operations into hardware.
 - Demonstrable high energy efficiency, achieving up to 1,082 graphs/J on the Cora dataset, exceeding modern GPU baselines by more than 50x, using a remarkably low DSP count, suitable for mid-range, embedded FPGA devices.
 - The validation of a dataset-aware 16-bit fixed-point quantization scheme that successfully preserves full classification accuracy while significantly minimizing resource utilization.

1.4 Thesis Organization

This thesis is organized into five chapters that collectively investigate the design and evaluation of FPGA-based acceleration techniques for deep learning workloads. Chapter 1 introduces the motivation, research objectives, and contributions, emphasizing the challenges of deploying CNN- and GCN-based models under stringent real-time and energy constraints. Chapter 2 provides a comprehensive and structured survey of real-time object detection on FPGAs, covering background fundamentals, FPGA architectures, acceleration techniques, comparative design analyses, and optimization strategies, culminating in an overview of existing challenges and future research trends. Chapter 3 presents a practical case study on heterogeneous MPSoCs, detailing the implementation of Tiny YOLOv2 on a Xilinx Zynq UltraScale+ device and delivering an empirical profiling-based analysis of execution bottlenecks, including CPU–DPU synchronization inefficiencies and DDR bandwidth underutilization. Chapter 4 introduces and evaluates a lightweight, energy-efficient FPGA streaming architecture for Graph Convolutional Networks, discussing its design rationale, sparse–dense compute modules, quantization strategy, and experimental performance across benchmark datasets. Finally, Chapter 5 offers concluding remarks and outlines future research directions emerging from the survey, MPSoC analysis, and GCN accelerator development.



Chapter 2

Real-Time Object Detection on FPGAs

2.1 Introduction

Real-time object detection is a computer vision task that aims at identifying and localizing objects in images or video sequences with acceptable inference speed, frame rate, and accuracy, satisfying the requirement of the application of interest. Autonomous driving [14, 15], robotics [16, 17], and healthcare monitoring [18, 19] are some examples in which real-time object detection can play a vital role.

Similar to classical real-time theory, we can define hard and soft constraints for real-time object detection. In hard real-time scenarios, the system fails if a specific deadline is not met, whereas in soft real-time situations, missing a deadline is undesirable but can be tolerated [20].

Achieving high power efficiency and, at the same time, performance satisfying hard constraints typically leads to the use of FPGA platforms [21, 22].

The role of Machine Learning Advancements in machine learning, particularly in deep learning [23], have led to the development of highly accurate and optimized object detection algorithms. Contemporary real-time object detectors commonly leverage Convolutional Neural Network (CNN) architectures to achieve an acceptable balance between accuracy and speed.

Recently, Transformer-based object detectors [24, 25] have gained considerable attention owing to their simplified, end-to-end architectures and impressive performance. However, their high computational demands limit their practical application [26], especially on resource-constrained platforms.

The role of FPGAs FPGAs are widely regarded as suitable platforms for implementing object detection systems as stand-alone target devices or, in heterogeneous systems, as accelerators alongside a processor [27]. They have gained remarkable attention for implementing real-time object detection systems due to their ability to offer low and deterministic latency, along with high throughput—critical factors for real-time systems [7].

Moreover, FPGAs can directly receive images from external imaging sources and process them without the need for processor intervention. This results in improved system performance and ensures lower and more predictable end-to-end latency. This is particularly significant as real-world object detection systems commonly rely on input image data from one or multiple imaging sources [28, 29].

Furthermore, the introduction of one-stage CNN-based object detection models [1, 2] in 2016, characterized by simpler architectures and higher speed compared to their two-stage counter-



parts [30–32], has facilitated the development of real-time systems, particularly on resource-constrained devices like FPGAs. Over time, with the introduction of more advanced one-stage object detection models [33,34], these architectures have managed to strike a balance between accuracy and speed. Figure 2.1 illustrates the increasing contribution of “real-time” systems in the context of implementing object detection on FPGAs over the last decades. This figure highlights that the introduction of these detection models has marked a significant turning point in attention to this topic.

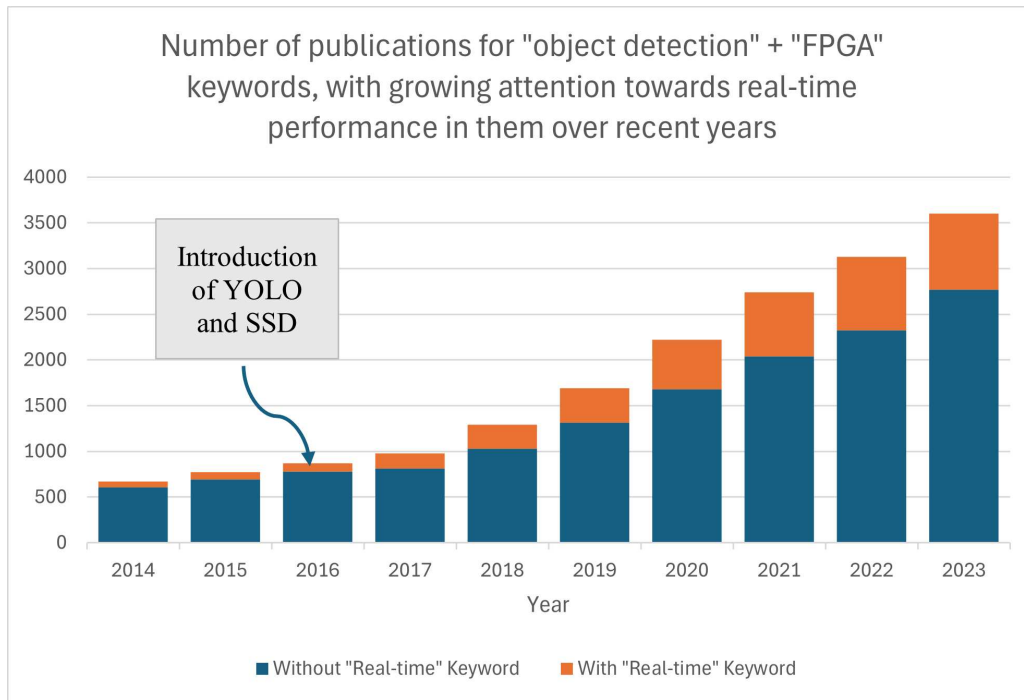


Figure 2.1: Increasing number of publications related to real-time object detection and FPGA over the years, especially after the introduction of one-stage object detection models like YOLO [1] and SSD [2], based on data extracted from Google Scholar.

Other related surveys This review examines implementations of real-time object detection on FPGAs, while other recent surveys partially address the topic as a whole as illustrated in Table 2.1. The survey literature is quite large, so we divided the contributions in three main categories discussed below: i) surveys on general object detection; ii) surveys on object detection on resource-constrained platforms; iii) surveys on FPGA-based object detection. In those categories, we discuss below and highlight in Table 2.1 the surveys of the last three years (2022, 2023, 2024), assuming that they already include comprehensively previous works and we explain how we extended such reviews.

Surveys on General Object Detection Numerous research works provide reviews of object detection algorithms, covering traditional and deep learning-based methods, commonly used data-

Table 2.1: Comparison of recent related review studies with the present work, focusing on the main topics discussed.

Year	Ref.	Main Surveyed Topics (Corresponding to ours)				Highlights
		Object Detection (OD) Methods	Real-time Performance	FPGA Implementation	HW Acceleration/Optimization Techniques	
2024	Ours	YES	YES	YES	YES	Focuses mainly on suitable OD models, HW acceleration methods, and optimization strategies to achieve real-time performance on FPGAs
2024	Mittal [35]	YES	YES	NO	NO	Focuses on the design and evaluation of lightweight object detection models based on deep learning, specifically tailored for edge devices
2023	Amjoud et al. [36]	YES	NO	NO	NO	Provides an overview of the current state of deep learning-based OD models
2023	Kamath et al. [37]	YES	NO	NO	NO	Provides a systematic review on deep learning-based OD on resource-constrained devices, emphasizing lightweight models
2023	R. Kaur et al. [38]	YES	NO	NO	NO	Reviews the evolution of OD, emphasizing DNN-based models while discussing key frameworks, architectures, datasets, and evaluation metrics
2023	Setyanto et al. [39]	YES	YES	NO	NO	Reviews OD methods and simplification strategies suitable for edge computing, focusing on compression techniques
2023	Zou et al. [40]	YES	YES	NO	NO	Reviews the development of OD with a discussion on detection speedup helpful for real-time performance
2022	Huang et al. [41]	YES	YES	NO	NO	Discusses CNN-based OD models suitable for implementation on Raspberry Pi while achieving real-time performance
2022	J. Kaur et al. [42]	YES	YES	NO	NO	Discusses various OD techniques, applications, and related tools and datasets, with a focus on existing challenges
2022	Zaidi et al. [43]	YES	YES	NO	NO	Reviews OD methods, focusing on lightweight models suitable for real-time performance on edge devices
2022	Zeng et al. [27]	YES	NO	YES	YES	Reviews some software and hardware optimization methods for the implementation of OD models on FPGAs

sets, typical performance metrics, and more [36, 38, 40, 42, 44–49].

The work by Amjoud et al. [36] provides an overview of the current state of object detection based on deep learning, covering two-stage anchor-based detectors, one-stage anchor-based detectors, anchor-free detectors, and transformer-based detectors. It also evaluates existing models across major object detection datasets such as Pascal VOC [50] and MS-COCO [51].

J. Kaur et al. [42] review key techniques, datasets, and tools for object detection, with a

focus on the advancements and challenges within the field. It covers various detection methods, noting how deep neural networks have improved performance, while challenges like small object detection, video analysis, and dataset limitations persist. The paper also explores future directions, such as combining detection models. They highlighted key research opportunities, such as designing efficient networks and leveraging multi-source information to enhance robustness in detection tasks.

The paper by R. Kaur et al. [38] reviews the evolution of object detection, focusing on deep convolutional neural networks and their applications across various fields. It compares deep learning methods with other object detection techniques and discusses key frameworks, architectures, datasets, and evaluation metrics.

Zou et al. [40] provide a comprehensive review of popular detectors, key technologies, speedup methods, datasets, and metrics spanning 20 years of object detection history, along with discussing promising future directions.

Surveys on Object Detection on Resource-Constrained Platforms Implementing these algorithms on resource-constrained devices is a practical and interesting topic. Extensive research has been conducted in this domain, covering various aspects of this issue [37, 39, 41, 43, 52–54]. From these works, several important design aspects emerge that need to be taken care of. They include the design of appropriate algorithms, techniques to simplify architectures, power efficiency (e.g., optimizing memory access, a primary source of power consumption), edge-specific metrics, such as cost, model size versus resource availability, and common metrics like accuracy, latency, and throughput.

Huang et al. [41], examine the challenges and advancements in enabling real-time object detection on resource-constrained edge devices. It explores methods of data processing on such devices, the development of efficient deep learning models, and the importance of input size reduction in neural networks for edge scenarios. The authors analyze existing approaches, comparing traditional machine learning methods and deep learning techniques for object positioning and classification, and evaluate their suitability for edge devices like Raspberry Pi. The paper also highlights trade-offs between computational efficiency, accuracy, and power consumption in these contexts.

Zaidi et al. [43] review recent developments in deep learning for object detection, covering benchmark datasets, evaluation metrics, and lightweight models for edge devices. It emphasizes the progress in creating faster, more accurate detectors to achieve real-time performance suitable for embedded applications.

The work by Kamath et al. [37] explores recent trends in deep learning-based object detection for resource-constrained devices. The study identifies key research areas, techniques, devices, and applications in this domain, based on a systematic literature review of 167 studies. It highlights the need for lightweight models that perform well on constrained devices, with a focus on the transportation industry.

Setyanto et al. [39] discuss various approaches for deploying object detection on near-edge devices, such as autonomous vehicles. They emphasize the challenges posed by limited computational resources and the need for efficient model compression techniques, such as network pruning and quantization, to maintain accuracy while reducing computational demand.

Mittal's review [35] provides an exploration of deep learning-based lightweight object de-

tection models designed for edge devices, addressing the increasing need for accurate, efficient, and low-latency detection systems. The study presents a taxonomy of lightweight detection algorithms, delves into key backbone architectures, and evaluates their performance on widely used datasets such as MS COCO [51] and Pascal VOC [55]. Additionally, it highlights the challenges and opportunities in the field, outlines real-world applications, compares state-of-the-art models, and offers insights into optimization strategies to enhance the performance of object detection models on edge platforms.

I. Introduction	B. Hardware Acceleration Techniques for FPGA-based Object Detection
II. Fundamentals and Background	C. Impact and Trade-offs of Acceleration Techniques in FPGA-based Object Detection
A. Object Detection Overview	V. Results Analysis and Optimization
B. Soft and Hard Real-time Constraints	A. Case Studies: Adoption of Acceleration Techniques in FPGA-based Real-time Object Detection
C. Evaluation Metrics and Datasets	B. Comparative Analysis of FPGA Implementations
III. FPGA Basics and Applications in Object Detection	C. Optimization Strategies for Real-time Performance
A. FPGA Overview	VI. Challenges and Future Direction
B. Applications of FPGAs in Computer Vision-Object Detection	A. Existing Challenges
C. FPGA Platforms for Accelerating Object Detection	B. Possible Future Research Trends
IV. FPGA-based Designs	VII. Conclusion
A. FPGA Architecture Design Approaches	

Figure 2.2: Organization of the paper.

Surveys on FPGA-based Object Detection Despite the proven benefits of FPGAs in developing object detection systems [7], especially when real-time performance is required [56–58], relatively little attention has been devoted to reviewing the related accomplished works.

Zeng et al. [27] partially address this gap by examining and comparing existing works in this domain. They also provide a summary of some software and hardware optimization techniques for the implementation of FPGA-based accelerators for object detection. However, in this survey, there is no emphasis on the real-time related aspects, which are connected more closely and examined in our work, besides analyzing also more recent works.

Our contribution However, to the best of our knowledge, this work marks the first dedicated review toward the implementation and optimization of “real-time object detection on FPGAs”. As illustrated in Figure 2.1, the contribution of real-time systems in this field is on the rise, highlighting the growing importance of focusing on this topic. Our discussion covers appropriate object detection algorithms, acceleration techniques, and optimization strategies for soft and hard real-time systems. Furthermore, we conduct a comprehensive examination and comparison of state-of-the-art works in this domain from various perspectives. Additionally, we delve into the existing challenges and propose potential solutions for achieving real-time object detection on FPGAs.

Organization of this survey As illustrated in Figure 2.2, the rest of the paper is organized as follows. Section 2.2 provides background information about object detection, the hard as well as soft real-time concepts. It also discusses evaluation metrics and datasets commonly used in

assessing real-time object detection systems. Section 2.3 provides an overview of FPGAs, exploring fundamental concepts, advantages, and common design approaches. It further delves into how FPGAs are applied in computer vision, particularly in object detection, and highlights specific platforms that enable accelerated processing for these tasks. Then, Section 2.4 delves into the implementation of real-time object detection systems using FPGAs, covering common architectures, acceleration techniques, and optimization strategies to achieve real-time performance. A review of state-of-the-art related works and their comparative analysis are also provided in that section. Section 2.6 addresses existing challenges in achieving real-time object detection systems on FPGAs, alongside potential directions for future research. The paper concludes in section 2.7.

2.2 Fundamentals and Background

2.2.1 OBJECT DETECTION OVERVIEW

Object detection is a computer vision task in which one or several objects in an image or a video frame can be identified and localized [40]. In other words, object detection algorithms perform two primary tasks:

- **Classification:** by assigning a label to every detected object based on the obtained probability that the object belongs to a certain category.
- **Localization:** by providing bounding boxes (bounding box detection) or highlighting significant landmarks (landmark detection) indicating the locations of detected objects within the input image.

Nowadays, object detection is widely applied in various real-world applications, such as video surveillance [59, 60], autonomous driving [14, 61], and healthcare monitoring [18, 62] to name a few. Also, object detection forms the foundation for numerous other computer vision tasks, such as object tracking [63], instance segmentation [64–66], and image captioning [67, 68].

The progress in object detection can be divided into two major periods: before and after the development of Deep Neural Networks (*DNNs*) [40]. Figure 2.3 depicts a classification of object detection techniques. This figure also illustrates which categories of models are more popular for real-time systems and which are better suited for FPGA implementation based on algorithm complexity and model size. The study specifically focuses on models that exhibit both characteristics, represented by the overlapping green area, which are related to CNN-based one-stage detectors.

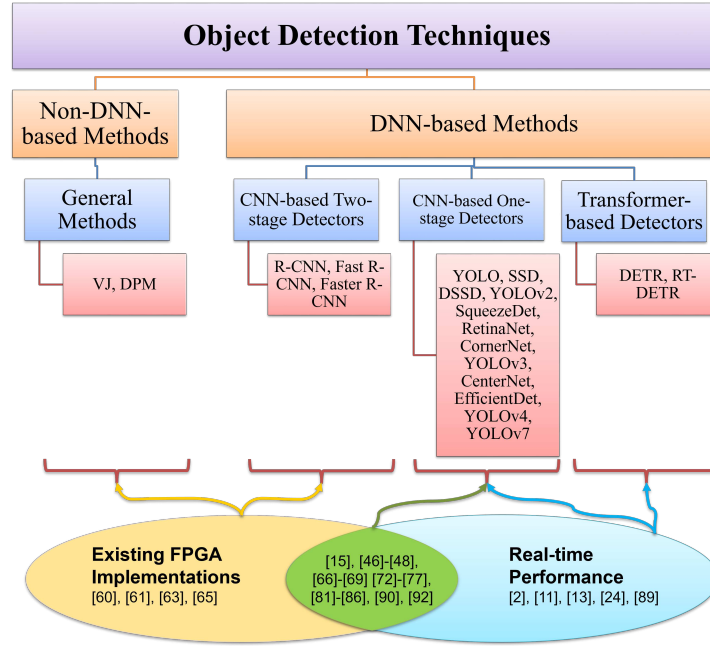


Figure 2.3: Classification of object detection techniques: Non-DNN-based and DNN-based detectors, along with instances of models for each category. It also shows which models may show real-time performance and which are more suitable for FPGA implementation. The focus of this study is the overlapping area (green zone.)

We discuss the works that do not rely on Deep Neural Networks (here referred to as *Non-DNN-based Methods*) and those that rely on DNNs (referred to here as *DNN-based Methods*).

In the following, the object detection models listed in Tables 2.2 and 2.3 are discussed in detail. We provide references to the original work that introduced the model, along with its main features, best-achieved performance, and, if available, references to subsequent FPGA-based implementations. For detailed descriptions of the performance metrics used in the table and the following parts, please refer to Section 2.2.3.

The table shows that CNN-based one-stage object detection models achieve a good balance between accuracy and speed, making them excellent choices for developing real-time object detection systems.

On the other hand, the introduction of DETR [24] marks the beginning of a new era in object detection models in which Transformers are deployed. These models can capture global dependencies and contextual understanding in input images, improving object detection accuracy [100]. However, their large size and high memory requirements often render them unsuitable for scenarios such as implementing real-time object detection systems on resource-constrained devices like FPGAs [101].

Table 2.2: A summary of the discussed models from all categories, including Non-DNN-based Detectors (NDD), DNN-based Two-stage Detectors (DTD), DNN-based One-stage Detectors (DOD), and Transformer-based Detectors (TBD) (Part I).

Detection Model	Year	Detector Category	Backbone Network	#Param (M)	Required GFLOPs	Best Reported Performance				Dataset	Hardware Platform (software implementation)	FPGA Platform (hardware implementation)	Highlights
						Accuracy	Latency (ms)	FPS	Size				
VJ [69]	2001	NDD	N/A	N/A	N/A	90% (detection rate)	N/A	15	384*288	MIT	CPU (Intel Pentium III)	Yes [70,71]	Fast, simple, good for face detection
DPM [72]	2008	NDD	N/A	N/A	N/A	34% AP	N/A	N/A	N/A	VOC	CPU	Yes [73]	Objects are represented as a collection of parts, Effective for Complex Objects
R-CNN [74]	2014	DTD	AlexNet	N/A	N/A	58.5% mAP	47000	N/A	227*227	VOC	GPU (NVIDIA Titan Black)	N/A	Uses CNN to extract features from regions proposed by region proposal part, improved accuracy compared to traditional methods, multi-stage pipeline, slow
Fast R-CNN [31]	2015	DTD	AlexNet	N/A	N/A	66.9% mAP	2000	N/A	1000*600	VOC	GPU (NVIDIA K40)	N/A	Faster than R-CNN by using a shared convolutional layer to process the entire image., multi-stage pipeline
Faster R-CNN [32]	2015	DTD	VGG-16	N/A	N/A	69.9% mAP	200	5	1000*600	VOC	GPU (NVIDIA K40)	Yes [75]	Faster than Fast R-CNN by using a better selective search algorithm, end-to-end network
YOLO [1]	2016	DOD	Modified GoogleNet	N/A	40.19	63.4% mAP	25	45	448*448	VOC	GPU (Geforce GTX Titan X)	Yes [76]	Introduced the concept of dividing the image into a grid and predicting bounding boxes, real-time performance with moderate accuracy
SSD [2]	2016	DOD	ResNet-101	N/A	N/A	74.3% mAP	N/A	46	300*300	VOC	GPU (NVIDIA Titan X)	Yes [11,56, 57,77,78]	Real-time object detector, Utilized a set of default bounding boxes at different aspect ratios and scales
DSSD [79]	2017	DOD	ResNet-101	N/A	N/A	81.5% mAP	N/A	6.4	513*513	VOC	GPU (NVIDIA Titan X)	N/A	An extension of SSD, Improved accuracy by adopting a deconvolution layer to extract more semantic information, Slower than SSD
YOLOv2 [80]	2017	DOD	DarkNet-19	193	34.90	76.8% mAP	N/A	67	416*416	VOC	GPU (Geforce GTX Titan X)	Yes [81-83]	Detection of a large number of object categories, Addressed some limitations of YOLO, such as localization accuracy and small object detection
YOLOv2-tiny [80]	2017	DOD	DarkNet-19	60.5	6.97	57.1% mAP	N/A	207	416*416	VOC	GPU (Geforce GTX Titan X)	Yes [7,10, 58,84,85]	Lightweight and faster alternative to YOLOv2, optimized for real-time performance on resource-constrained devices
SqueezeDet [15]	2017	DOD	SqueezeNet	26.8	77.2	80.4% mAP	N/A	32.1	1242*375	KITTI	GPU (Geforce GTX Titan X)	N/A	Designed for efficient object detection on embedded systems by reducing model size and complexity
RetinaNet [86]	2017	DOD	ResNet-101-FPN	N/A	N/A	39.1% mAP	90	5	600*600	COCO	GPU(NVIDIA M40)	N/A	Addressed the imbalance between foreground and background classes with a focal loss function
CornerNet [87]	2018	DOD	Hourglass	N/A	N/A	42.2% mAP	224	N/A	511*511	COCO	GPU (Titan X)	N/A	Detected objects as paired keypoint locations (top-left and the bottom-right corners)
YOLOv3 [88]	2018	DOD	Darknet-53	237	65.86	55.3% mAP	29	35	416*416	COCO	GPU (Geforce GTX Titan X)	Yes [8,89]	Introduced multi-scale detection and feature pyramid network to improve accuracy
YOLOv3-tiny [88]	2018	DOD	Darknet-53	33.8	5.56	33.1% mAP	4.5	220	416*416	COCO	GPU (Geforce GTX Titan X)	Yes [90-93]	Efficient and faster variant of YOLOv3, offering a good balance between speed and accuracy on embedded devices

Table 2.3: A summary of the discussed models from all categories, including Non-DNN-based Detectors (NDD), DNN-based Two-stage Detectors (DTD), DNN-based One-stage Detectors (DOD), and Transformer-based Detectors (TBD) (Part II).

Detection Model	Year	Detector Category	Backbone Network	#Param (M)	Required GFLOPs	Best Reported Performance				Dataset	Hardware Platform (software implementation)	FPGA Platform (hardware implementation)	Highlights
						Accuracy	Latency (ms)	FPS	Size				
CenterNet [94]	2019	DOD	Hourglass	N/A	N/A	47% mAP	340	N/A	511*511	COCO	GPU (NVIDIA Tesla P100)	Yes [95]	Predicting the center point and regressing to the object size and orientation
EfficientDet [96]	2020	DOD	EfficientNet	6.6	6.1	39.6% AP	20	50	640*640	COCO	GPU (Titan V)	N/A	Utilized a compound scaling method to efficiently scale the network, balanced model efficiency and accuracy.
YOLOv4 [33]	2020	DOD	CSPDarknet-53	64.4	43.5	41.2% mAP	N/A	96	416*416	COCO	GPU (Tesla V100)	Yes [97]	Introduced significant architectural improvements over YOLOv3
YOLOv4-tiny [98]	2020	DOD	CSPDarknet-53	6.1	6.9	21.7% mAP	N/A	371	416*416	COCO	GPU (Tesla V100)	Yes [92]	Optimized for real-time detection with high speed, providing improved accuracy over previous tiny models
DETR [24]	2020	TBD	ResNet-50	41	86	42% mAP	N/A	28	800*1200	COCO	GPU (V100)	N/A	Transformer encoder-decoder, streamlining the detection pipeline, issues with processing small objects, too large
YOLOv7 [34]	2022	DOD	E-ELAN	36.9	51.2	51.2% mAP	N/A	161	640*640	COCO	GPU (Tesla V100)	N/A	Outperforms all YOLO versions in speed and detection accuracy
YOLOv7-tiny [34]	2022	DOD	E-ELAN	6.2	13.8	35.2% mAP	N/A	161	416*416	COCO	GPU (Tesla V100)	Yes [99]	Enhanced speed and accuracy for real-time detection, featuring optimized architecture for edge devices
RT-DETR [26]	2023	TBD	ResNet-50	42	136	53.1% mAP	8.8	108	1024*1024	COCO	GPU (T4)	N/A	Good speed-accuracy trade-off, eliminates NMS, Supports flexible speed tuning without retraining, still too large

Non-DNN-based Detectors

Early works on object detection tasks were performed using handcrafted features. Viola and Jones introduced the first real-time detector, called *VJ detector*, for human faces in the early 2000s [69]. Using sliding windows over the image to find out which windows contain a face and deploying some techniques to speed up the algorithm, this detector was about 15 times faster than other detectors at that time with comparable accuracy. However, in addition to the long training time, the VJ detector was limited to performing binary classification.

In 2005, Dalal and Triggs proposed a more accurate detector, which was capable of detecting more object classes, based on the Histogram of Oriented Gradients (*HOG*) [102]. HOG serves as a feature descriptor, widely used for extracting features from input images in computer vision, especially in object detection applications [102, 103]. Similar to methods such as Scale-Invariant Feature Transform (*SIFT*) [104], which focuses on the structure of objects, HOG counts intensity gradient orientation occurrences in localized regions of an image. In contrast to its equivalent descriptors [104, 105], HOG deals with a uniformly spaced array of cells forming a dense grid and takes advantage of overlapping local contrast normalization to enhance detection performance.

An extension of the HOG detector [102] is the Deformable Part-Based Model (*DPM*), pro-

posed by Felzenszwalb et al. [72] in 2008. DPM is one of the well-known non-DNN object detectors, as it won several Pascal VOC detection challenges [50]. Following the “divide and conquer” approach, DPM tries to decompose objects into parts during training and performs ensemble-based inference on these components. It means that the problem of detecting an object can be decomposed into detecting constituted parts of that object. Later, several improvements were made by Girshick on the original DPM detector to enhance the accuracy as well as the speed of detection [74, 106–108].

Efforts have been made to develop object detection systems on FPGAs using these types of models [70, 71, 73]. However, in addition to offering low speed (up to 11.75 FPS, to the best of our knowledge), since these detectors rely on fixed, handcrafted features, they often fail to deliver high accuracy and robustness, particularly when detecting diverse objects in complex backgrounds. As a result, they are not well-suited to achieve real-time performance for FPGA-based object detection systems.

DNN-based Detectors

Object detection has witnessed a remarkable breakthrough after introducing deep neural networks [109]. The ability to learn robust and high-level representations of images without any need for handcrafted features (such as SIFT [104] and HOG [102]), which had limited accuracy of non-DNN-based object detectors, motivated many researchers to leverage DNNs in vision tasks [110].

DNN-based object detectors can be divided into two categories: CNN-based and Transformer-based detectors [111]. The architectures of CNN-based detectors have progressed from “two-stage detectors” to “one-stage detectors”.

Object detectors incorporate a *backbone* network as a feature extractor to derive features from input images. VGG [112], GoogleNet [113], EfficientNet [114], and DenseNet [115] are some CNN architectures that can be adopted as backbone networks of object detectors. Recently, it has been proven that in addition to CNNs, Transformers can also be deployed as backbone networks in object detectors [116, 117].

2.a: Two-stage Detectors The Region-based Convolutional Neural Network (*RCNN*), introduced by Girshick et al. [30, 74], marked the initial breakthrough in utilizing DNNs for object detection. Based on this promising model, a new branch of object detectors emerged called “Two-Stage Detectors.” As the name implies, these detectors are made up of two major stages:

1. Region proposal: to propose a set of candidate regions or bounding boxes likely to contain objects of interest [118].
2. Classification and refinement: to perform classification for determining the presence of objects within each proposed region. If any object is detected in this stage, the corresponding bounding box coordinates are refined for more accurate localization.

Several two-stage detectors have been proposed [31, 32, 66, 119, 120], but reviewing them is beyond the scope of this survey. The existing implementations typically require more computational resources than one-stage detectors, particularly for the region proposal task, making them less suitable for real-time object detection. Furthermore, considering their complex architectures with a vast array of parameters, deploying these detectors on resource-constrained platforms such

as FPGAs is a challenging task [39]. In an attempt to implement these types of detectors on FPGAs, An et al. [75] implemented the Faster R-CNN algorithm [32] on an Arria-10 GX FPGA board. However, the results indicated a high latency of 153.6 ms, which is typically considered unsuitable for real-time performance.

2.b: One-stage Detectors Poor speed and complex architectures of two-stage detectors motivated researchers to develop detectors capable of performing object classification and localization in one single-stage, called one-stage detectors. Although they usually suffer from poor accuracy, especially in detecting small objects, compared to two-stage detectors [40], one-stage detectors feature simpler architectures, higher speed, and adaptability to various scales [39]. These attributes make them suitable for many real-world vision tasks, especially when it comes to leveraging resource-constrained devices.

YOLO family: In contrast to two-stage detectors, which treat object detection as a classification task, Redmon et al. [1] approached it as a regression problem by proposing the first version of You Only Look Once (*YOLO*) architecture in 2016.

In YOLO, the entire image is passed through a CNN, where objects are classified and localized directly. Initially, input images are divided into several cells, each responsible for predicting multiple bounding boxes along with their corresponding confidence scores. Every bounding box is defined by its coordinates, (x, y) representing the center of the box and (w, h) denoting its width and height, respectively. For each bounding box, YOLO predicts class probabilities for every category of objects. In addition, the confidence score indicates how confident the model is that a bounding box contains an object of interest. Then, during the post-processing, redundant and low-confidence bounding boxes are removed, keeping only the most confident and non-overlapping predictions. YOLO provides a set of bounding boxes, each containing a class label and a confidence score, as output.

Li et al. [76] showed the feasibility of deploying YOLO on FPGAs for developing real-time object detection systems. They achieved about 44x inference speedup on a ZYNQ7035 compared to an Intel Core i5-6200U CPU, demonstrating promising results for real-time performance.

While it can detect multiple objects at a high speed, making it suitable for real-time vision applications, YOLO suffers from poor localization accuracy, especially for small objects [80]. So far, many different versions and variants of YOLO have been developed, each improving performance and efficiency. Among all, YOLOv2 [80], YOLOv3 [88], YOLOv4 [33], and YOLOv7 [34] are briefly discussed in the following. Although several diverse and newer versions of YOLO have been introduced at the time of writing this paper [121, 122], it was not feasible to review all of them due to the extensive volume of content. Therefore, we have restricted our discussion to the versions introduced by the main YOLO research group that have been considered for implementation on FPGAs thus far.

Unlike its predecessor, which utilized GoogleNet [113] as its backbone network, YOLOv2 adopts a less complex and lighter architecture called DarkNet-19, with 19 convolutional and 5 max-pooling layers [123]. Leveraging the Batch Normalization technique [124], which accelerates the optimization process, and anchor idea¹ [32], the generalization capability was improved in this version. These enhancements made the model more powerful in predicting objects of

¹Anchor boxes are predefined boxes with various aspect ratios and scales. They are utilized to align the predicted bounding boxes with the actual objects present in the image.

varying scales and shapes. Overall, YOLOv2 significantly enhanced the performance of YOLO, achieving advancements in speed (by at least 30%), accuracy (approximately 22% on the Pascal VOC dataset [50]), flexibility, and its capability to handle a wider range of object categories [80].

Numerous successful efforts have deployed YOLOv2 on FPGAs to develop real-time object detection systems [7, 10, 58, 81–85], demonstrating its significant capability in this area. For instance, Nakahara et al. [81] deployed a simplified YOLOv2 model on a ZCU102 FPGA board, achieving a throughput of 35.71 FPS, a promising processing speed in various fields such as video surveillance. Other recent works are reviewed in this study.

To further improve detection accuracy and speed, YOLOv3 [88] adopts Darknet-53 [125], a deeper (53 convolutional layers) and more complex architecture than Darknet-19 [80, 123]. Also, unlike YOLOv2 [80], in which all anchor boxes have the same size, YOLOv3 uses anchor boxes in different scales and aspect ratios. Overall, YOLOv3 shows better detection accuracy, especially for small objects, compared to its ancestors.

Like YOLOv2, YOLOv3 is also popular in developing real-time object detection systems using FPGAs [8, 89–92]. Wang et al. [89], demonstrated that deploying this model on FPGAs could lead to higher throughput by 6.5x and lower power consumption by 5x compared to running the same model on a GeForce GTX1080 GPU. In another effort, Yu et al. [90] showed that a tiny YOLOv3 model could be deployed to develop object detection systems with low latency, a crucial parameter for achieving real-time performance, even on low-end FPGAs.

In YOLOv4 [33], many techniques are integrated to improve both accuracy and speed. By combining some methods, such as Weighted-Residual-Connections (*WRC*), Cross-Stage-Partial-connections (*CSP*), and Cross mini-Batch Normalization (*CmBN*), YOLOv4 achieve improved results in detection accuracy (up to 43.5% mAP on MS COCO [51]) at a speed of 30 frames per second (*FPS*) and higher on a GPU-based platform (see Tables 2.2 and 2.3). Using *WRC*, it becomes feasible to amplify the influence of crucial features acquired from preceding layers onto the current layer. This enhancement results in improved performance compared to using simple residual connections. *CSP* aims to lower computational complexity by splitting the feature map of the base layer into two sections and then merging them after further computations are executed on only one branch. *CmBN*, a variation of standard Batch Normalization [124], normalizes the model's activations based on the assumption that each batch consists of four mini-batches. This method gathers statistics only between mini-batches within a single batch, making YOLOv4 more efficient in both training and testing.

YOLOv4 is a practical choice for real-time FPGA-based object detection systems [92, 97], thanks to its efficient design and optimal balance between accuracy and speed. This study offers an in-depth analysis of the results obtained from recent advancements in this field.

YOLOv7 [34] outperforms all previous YOLO versions in speed and detection accuracy. The remarkable results are attained through optimizing model architecture, employing advanced training strategies, and utilizing efficient feature extraction techniques. These innovations enhance accuracy without significantly increasing inference costs. Additionally, this model can use parameters and computation more effectively by leveraging some proposed techniques such as compound scaling. The introduced compound scaling aims to allow developers to customize key model attributes, such as width, depth, and resolution, based on the desired application requirements and the characteristics of the target computing device while maintaining the initial model properties.

YOLOv7 is another viable option for developing real-time object detection systems on FP-

GAs [99], offering state-of-the-art accuracy, an optimized design for faster inference, and advanced techniques like compound model scaling that ensure compatibility with FPGA resource constraints.

SSD: Early YOLO models suffered from poor accuracy compared to two-stage detectors. The Single Shot Multibox Detector (*SSD*), proposed by Liu et al. [2], was the first one-stage detector that achieved accuracy comparable to two-stage detectors while operating at real-time speeds on GPUs.

The backbone network in SSD is based on VGG-16 [112] with a few customizations. In particular, some fully connected layers are replaced with convolutional ones, and more multi-scale convolutional layers are added at the end of the network to improve detection performance.

In SSD, the predictions are made based on features extracted at different convolutional layers. This strategy helps the model detect objects of different scales and sizes more efficiently since every layer can provide different semantic information. In addition, a fast Non-Maximum Suppression (*NMS*) technique [126] is deployed at every stage to remove redundant bounding boxes. This leads to reduced computation compared to using NMS only at the final stage. (NMS works by selecting the bounding box with the highest confidence score and suppressing all other boxes with significant overlap, ensuring only the most accurate detections are retained.)

SSD's lightweight architecture, single-stage detection pipeline enabling fast inference, scalability across various input resolutions, and compatibility with FPGA-friendly optimizations such as quantization and pruning make it a widely adopted model for FPGA-based real-time object detection. As discussed later in this study, the favorable results achieved in recent works [11, 56, 57, 77, 78], strongly support its popularity in this domain.

Later, Deconvolutional SSD (*DSSD*) [79] was introduced to improve the accuracy of SSD in detecting small objects by adopting a deconvolution layer² to extract more semantic information. However, this improved accuracy comes at the cost of slower inference speed than SSD.

CenterNet: One problem with keypoint-based detectors is that they generate many incorrect object bounding boxes, mainly because they do not examine the cropped regions [94]. Duan et al. [94] proposed a new keypoint-based object detection named *CenterNet* to address this issue.

In this detector, a third keypoint, i.e., the center of each object, is also predicted, leading to detection accuracy improvement up to 4.9% on the MS COCO dataset [51]. In addition, two customized pooling layers are introduced to provide more precise and recognizable information about the top-left and bottom-right corners, as well as the center of each object.

solovyev et al. [95] achieved promising results (about 19 FPS) by deploying CenterNet on a Cyclone V FPGA, demonstrating the potential of this model towards developing real-time object detection on FPGAs.

In the following, other one-stage detectors listed in Table 2.2 and 2.3 are briefly discussed. While these models demonstrate potential for deployment in real-time object detection systems based on their characteristics and performance, as summarized in the table, to the best of our knowledge, no FPGA implementations of them have been reported in recent publications.

SqueezeDet: SqueezeDet [15] is a lightweight, fully CNN-based object detector proposed by the developer of SqueezeNet DNN architecture [127] in 2017. Initially designed for autonom-

²Unlike convolution layers, where the input image is downsized through convolution with a kernel, deconvolution layers perform the reverse process.

ous vehicles, SqueezeDet aimed to achieve reasonable accuracy at real-time speed on resource-constrained devices.

In addition, this model addresses some other concerns related to such devices, such as model size and power efficiency. It also adopted SqueezeNet as its backbone network, with the architecture consisting of a single forward-pass neural network, making it an effectively lightweight and fast single-stage detector

RetinaNet: Recognizing the significant class imbalance as a primary factor leading to lower accuracy in one-stage detectors compared to two-stage ones, Lin et al. [86] introduced a new loss function, named “Focal Loss”, to tackle this issue. Focal Loss can address the class imbalance bottleneck by reducing the loss contribution from well-classified examples, enabling the model to prioritize and focus on difficult-to-classify instances.

Based on the proposed loss function, they introduced Retina-Net which concentrates more on misclassified examples during the training phase, resulting in a remarkable accuracy improvement. In addition, this model takes advantage of ResNet [128] (ResNet-50 and ResNet-101) followed by a Feature Pyramid Network (FPN) [129] as its backbone network. Deploying FPN helps Retina-Net detect objects of different scales and dimensions more accurately (40.8% AP on COCO [51]).

CornerNet: The main idea of CornerNet [87] is finding two corners of each object, as *keypoints*, instead of dealing with anchor boxes to perform object detection. A single CNN is utilized to predict two separate heatmaps and one embedding vector for each predicted corner, covering both the top-left and bottom-right corners. Heatmaps are binary masks that indicate the locations of a class’s corners, and embeddings are numerical vectors that group the corners associated with each object.

Law and Deng [87] also introduced a novel pooling layer, called *Corner Pooling*, to enable the proposed model to localize corners more accurately. Compared to its contemporary one-stage detectors, CornerNet shows a competitive accuracy of 42.2% AP on the MS COCO dataset [51].

EfficientDet: In 2020, Tan et al. [96] introduced several optimization methods to improve the efficiency of existing object detectors. These methods include a novel FPN, called Weighted Bi-directional Feature Pyramid Network (BiFPN), and a compound scaling method. While BiFPN effectively enables the model to perform multi-scale feature fusion, the compound scaling method can simultaneously and consistently scale the depth, width, and resolution of different parts of the model. By deploying EfficientNet [114] as the backbone network and integrating these optimization techniques, EfficientDet [96] was introduced.

This model can achieve remarkable accuracy (55.1% AP on MS COCO [51]) with significantly fewer parameters (up to 9x fewer) and reduced computation (up to 42x fewer FLOPs) compared to previous state-of-the-art object detectors. These advancements make EfficientDet an appealing choice for resource-constrained platforms.

Overall, CNN-based one-stage object detectors, particularly the YOLO family and SSD models, play a crucial role in real-time object detection on FPGAs due to their efficient architectures and high-speed processing. Unlike two-stage detectors that separate object proposal and classification stages, one-stage detectors integrate these tasks, enabling faster inference times—a key requirement for real-time applications. The combination of their streamlined architectures and FPGA capabilities results in high throughput, low latency, and power-efficient designs, making them highly attractive for real-time object detection systems. Deploying these models on FPGAs

supports applications like autonomous vehicles, robotics, and surveillance systems, where speed and accuracy are paramount.

2.c: Transformer-based detectors In addition to the above-mentioned deep learning-based objection models, i.e., CNN-based one-stage and two-stage detectors, Transformer-based object detectors have recently attracted remarkable attention in the computer vision field, especially for object detection tasks [130]. However, despite demonstrating excellent performance results, as shown in Tables 2.2 and 2.3, transformer-based object detectors are not well-suited for FPGA implementations due to their high computational complexity and memory demands. To the best of our knowledge, no FPGA-based real-time object detection implementations of these models have been reported.

From an architectural perspective, these detectors can be considered a subset of one-stage detectors. However, due to the distinct model structure, we review this category of detectors separately.

Transformers are a type of deep-learning model originally proposed for performing sequence-to-sequence tasks in Natural Language Processing (*NLP*) [131]. They are taking advantage of the self-attention mechanism, enabling them to identify and weigh the significance of different parts of the input data. As a result, Transformers can discover long-term dependencies and complex patterns within a sequence.

The recent achievements of transformers in NLP have prompted researchers to investigate their capabilities in the field of computer vision as well [132]. In this context, some studies have explored the integration of Transformers with CNN-based architectures [24, 133], whereas some endeavors have focused on constructing the entire model solely using Transformers and without incorporating any CNN network for feature extracting [117, 134, 135]. In the following, some of these efforts are briefly reviewed.

DETR: In 2020, Carion et al. [24] proposed the first end-to-end object detection model that leverages Transformers, named DEtection TRansformer (*DETR*). They introduced a novel approach for object detection, treating it as a direct set prediction task. Also, the proposed method simplifies detection by eliminating the need for hand-designed parts such as non-maximum suppression or anchor generation.

DETR utilizes a transformer encoder-decoder architecture. By reasoning about object relations and global image context, DETR outputs the final set of predictions in parallel using a fixed small set of learned object queries.

It achieves accuracy and runtime performance comparable to Faster RCNN [136] on the MS COCO dataset [51]. However, the practical application of DETRs is constrained by their high computational cost [26]. DETRs also show poor performance in detecting small objects [137]. Numerous modifications have been introduced thus far to tackle the challenges associated with DETR. [25, 138, 139].

RT-DETR: In 2023, Zhao et al. tried to reduce the computation complexity of DETR [24] to make it suitable for real-time object detection [26]. Based on those efforts, they proposed the first real-time end-to-end object detector, called Real-Time DEtection TRansformer (*RT-DETR*).

They improved the encoder part of DETR to be able to process multi-scale features. In addition, an Intersection over Union (*IoU*)-aware query selection component was proposed to enhance the initialization of object query.

RT-DETR outperforms other real-time detectors and end-to-end detectors of comparable size, excelling in both speed and accuracy (up to 54.3% AP on COCO [51]) and 108 FPS) to achieve state-of-the-art performance. Furthermore, the proposed detector allows for adaptable adjustment of inference speed by utilizing different decoder layers without any need for retraining. This feature is extremely useful for the practical implementation of real-time object detection systems.

2.d: Object Detectors based on Vision Transformers ViT: Prior to the introduction of Vision Transformer (*ViT*) by Dosovitskiy et al. [117] in 2020, attention mechanisms were primarily utilized for vision tasks in conjunction with CNNs. In computer vision, “attention” refers to examining the relationships between pairs of input image tokens or patches. This mechanism allows the model to prioritize relevant input features, aiding in capturing more informative representations of the input image.

Although some researchers tried to replace convolutional architectures with Transformers [134, 135], those models were not efficient enough for deploying on existing hardware accelerators. That was mainly because of utilizing customized attention patterns. However, ViT demonstrated that a pure Transformer could effectively perform image processing by treating it as sequences of patches without relying on CNNs [117]. The straightforward and scalable approach in ViT proves remarkably effective, particularly when combined with pre-training on extensive datasets.

As a result, Vision Transformer not only meets but often surpasses the performance of state-of-the-art CNNs [117]. This breakthrough opened the door to more widespread exploration of Transformer-based models for various computer vision tasks.

In an effort, for instance, developing a Transformer-based object detector by combining ViT and DETR [24] showed promising results [140].

Swin Transformer: Swin Transformer [116], introduced by Liu et al. in 2021, is a vision Transformer, basically designed as a backbone network for computer vision tasks. To address the challenges of applying Transformers to computer vision tasks—such as the varied scale of visual entities and higher image resolutions compared to texts—the authors introduced a hierarchical Transformer using **Shifted windows**.

Swin Transformer enhances efficiency by confining self-attention to non-overlapping local windows while enabling cross-window connections. Its hierarchical design approach accommodates various scales and maintains linear computational complexity relative to image size. Swin Transformer demonstrates its versatility across image classification, semantic segmentation, and object detection tasks, achieving state-of-the-art performance on the MS COCO dataset. [51].

The results show the potential of Transformer-based models as effective vision backbones, especially with the hierarchical design and the *shifted window* approach, which also benefits all-MLP (multi-layer perceptron) architectures.

2.2.2 Soft and Hard Real-time Constraints

Real-time systems are those in which the accuracy of performance is determined by both the logical results of computation and the timeliness of producing those results [141, 142]. Today, they play a key role in many applications, such as robotics, manufacturing, healthcare, and autonomous driving systems.

Real-time systems ensure that tasks are executed within a precise and predictable time frame, making these systems safe, predictable, and reliable [20]. The predictability of real-time systems can be measured by evaluating the existing latency and its variation between iterations, also known as jitter [20].

Another important feature of real-time systems is their ability to manage real-time and non-real-time tasks to prevent system failure [143]. For example, in a self-driving system equipped with an object detection system, the system must be able to prioritize sending a warning signal to the central control system when detecting an obstacle rather than displaying it on the screen.

Real-time systems can be classified into two main categories as follows:

- **Soft real-time systems:** They are characterized by possible performance degradation rather than complete failure when response time constraints are not met [144].
- **Hard real-time systems:** These are systems in which failure is inevitable if response-time constraints are not met.

More precisely, in hard real-time systems, missing a deadline may have dangerous consequences, such as human injury, equipment damage, or even death [143]. This makes it crucial to ensure deterministic and predictable behavior. In autonomous vehicles, for instance, timely processing of input data is crucial for collision avoidance and safe navigation. These systems must process data from various sensors and analyze received images in real time to detect pedestrians, other vehicles, and obstacles in the path. This example highlights the significance of a hard real-time system for executing object detection, which is crucial for ensuring the safety of both passengers and pedestrians.

Designing real-time systems across different domains involves addressing various considerations [145–148]. As highlighted earlier, contemporary approaches heavily rely on CNN models in the context of real-time object detection systems, which is our focus in this study. A crucial challenge in implementing these systems, particularly on embedded platforms, is achieving a delicate trade-off between accuracy and speed to satisfy real-time requirements [26, 149, 150].

To that end, a range of techniques can be employed, including algorithmic and hardware implementation strategies, such as quantization and model pruning [37, 151, 152]. These approaches aim to optimize computational efficiency while preserving adequate accuracy for timely and reliable object detection. These methods are investigated in more detail in this study.

2.2.3 Evaluation Metrics and Datasets

Common Datasets

Datasets are crucial in developing object detection models, including training and evaluation phases. During the training process, models learn to recognize patterns and features associated with different objects based on many example data provided by a dataset. Utilizing a diverse, extensive, and well-structured dataset for training can significantly improve the model's capability to detect unseen objects during inference [48].

Datasets are also used to evaluate the performance of an object detection model and compare it with others [110].

Below, we briefly discuss the commonly used datasets in the field of object detection. Table 2.4 summarizes the key characteristics of these datasets.

Table 2.4: A summary of the commonly-used datasets in the context of object detection

Dataset	# of classes	# of images	# of annotated objects	Examples of object types	Tasks
ImageNet [109]	20K	14M	1M (with Bbox)	animals, everyday objects, natural scenes, plants, people, food, abstract concepts	image classification, object detection
MS COCO [51]	91	328K	2.5M	animals, person, everyday objects, vehicles, food, sport equipment	object detection, image segmentation
Pascal VOC 07 [50]	20	9963	26640	person, animals, vehicles, and indoor objects	object detection, semantic segmentation, image classification
Pascal VOC 12 [55]	20	11K	27450	person, animals, vehicles, and indoor objects	object detection, semantic segmentation, image classification

ImageNet ImageNet [109] is a large-scale dataset commonly used in object classification and detection models. It has also been considered the main benchmark in the ImageNet Large Scale Visual Recognition Challenge (ILSVRC) [153] for object detection and image classification. It contains over 14 million images annotated in one of two ways: image-level or object-level. The former indicates whether or not an object class exists in the image, while the latter provides a bounding box and class label for every object instance in the image.

MS COCO Microsoft Common Objects in Context (*MS COCO*) dataset, introduced in 2014 by Lin et al. [51], is a large visual dataset widely used in computer vision tasks, such as object detection and image segmentation. It contains 2.5 million labeled instances in 328K images of 91 different object categories, including everyday objects and humans. It is primarily prepared for the detection and segmentation of objects that appear in their natural surroundings, containing more classes and instances compared to PASCAL VOC [55]. After years, this dataset remains one of the most popular choices in the field of computer vision, with its usage continuing to grow. In 2023 alone, it was referenced in over 2200 articles [154].

Pascal VOC The Pascal Visual Object Classes (*VOC*) dataset [50, 55] is widely used in object detection, semantic segmentation, and image classification tasks. Two more popular versions of this dataset are VOC 2007 [50] and VOC 2012 [55]. The former contains 20 object classes, including persons, animals, vehicles, and some indoor stuff, in 9,963 images with 24,640 annotated objects. The latter's number of categories has remained unchanged, while it contains more data: over 11K images with 27,450 Region of Interest (*ROI*) annotated objects.

Performance Metrics for Real-time Object Detection

Various parameters, including accuracy, processing throughput, model complexity, and inference time, can be employed to measure the efficiency and performance of object detection models and systems. Table 2.5 shows some commonly used metrics for evaluating object detection models. Considering this study’s primary focus, this section only discusses metrics that are more relevant to assessing real-time object detection systems, particularly those that are more important when FPGA implementation is concerned. To better understand other useful metrics in this context, such as Intersection over Union (*IoU*) and Average Precision (*AP*), it may be helpful to review Appendix A.

Table 2.5: A summary of commonly used metrics for evaluating object detection models.

Metric	Measure of
Intersection over Union (IoU)	Localization accuracy
mean Average Precision (mAP)	Detection (localization and classification) accuracy
Frames Per Second (FPS)	Processing throughput
Number of required Floating-point Operations (FLOPs)	Computational complexity
Number of trainable parameters	Computational complexity and memory footprint
Latency	Inference time

One popular evaluation metric in the context of object detection is mean Average Precision (*mAP*). Several object detection algorithms, including Faster R-CNN [32] and YOLO [1], employ mAP for assessing their models. It serves as a standard metric in various benchmark challenges like Pascal VOC [50, 55]. mAP evaluates detection accuracy across multiple classes and indicates the average of AP (see Appendix A) across all classes. Equation 2.1 shows how mAP is obtained for N classes of objects, where the AP_k is the AP of class k .

$$mAP := \frac{1}{N} \sum_{k=1}^N (AP_k) \quad (2.1)$$

Another key metric for evaluating object detection systems is the number of frames the model can process, which is, in fact, a measure of the speed and efficiency [155]. This factor, called frame rate, is measured as Frames Per Second (*FPS*) and is directly related to the term “real-time” in real-time object detection systems. When discussing “real-time” object detection models, it is crucial to first provide a clear and precise definition of what they are.

According to the discussion provided in section 2.2.2, the word “real-time” has a specific definition and does not mean simply being “fast”. However, an inaccurate definition of real-time object detection systems and models is sometimes used. This definition may focus only on speed or processing throughput and consider meeting a certain frame rate (e.g., 30 FPS) sufficient to classify a system as real-time. However, it should be evaluated according to specific criteria depending on the requirements of the application. For example, a system that detects objects at 30 FPS for a self-driving car traveling at speeds of 110 km/h or higher cannot meet the criteria for a hard real-time system. This is because, at such a speed, the distance to objects changes by approximately 1 meter (for stationary objects) in each frame interval. Therefore, although speed, or processing throughput, is an important factor in assessing object detection systems, it should not be considered the only criterion for categorizing a system as real-time or non-real-time.

Although the frame rate metric can provide valuable insight into real-time object detection system performance, it alone may not offer a complete picture, particularly when comparing systems that process images of varying resolutions. To address this limitation, we introduce pixel throughput, as defined in Equation 2.2, to evaluate the performance while processing video streams or images. Pixel throughput is calculated as the product of the achieved frame rate (FPS) and the number of pixels per frame, capturing both the temporal and spatial aspects of system performance.

$$\langle Pixel Throughput \rangle = \langle Frame Rate \rangle \times \langle Pixel Per Frame \rangle \quad (2.2)$$

This metric, measured in Mega Pixels Per Second (*MPPS*), allows for a more fair comparison between works that evaluate implementations that use different image sizes.

There are some other key parameters that can be used to evaluate object detection systems, which are particularly important when it comes to achieving real-time performance on resource-constrained devices. They include computational complexity, number of parameters, and latency [155–157].

The computational complexity of an object detection model can be represented by the number of Floating-point Operations (*FLOPs*) required for inference. Although higher FLOPs can result in higher accuracy, it directly affects the number of required computational resources and the model’s speed [155].

The number of trainable parameters of a model can also show the model’s complexity as well as memory footprint [155]. Memory footprint means how much memory we need to store all the parameters of a model. This metric can directly affect the power consumption and processing performance of the system. In fact, when the model parameters exceed the capacity of available on-chip memory of the target device, resorting to off-chip memory becomes unavoidable, a common scenario in real-world object detection systems. As the rate of memory access increases, system speed decreases while power consumption rises. [156].

Generally speaking, latency refers to the duration it takes for a system to generate output after input is received [156]. It quantifies the delay in the response of a digital system. Typically measured in milliseconds (*ms*) [26,58], A greater latency value indicates a slower system performance. Likewise, in an object detection system, latency indicates the time required for the model to analyze an image and provide information about the identified objects. Especially in real-time object detection systems, it is a crucial metric in system evaluation and a critical parameter for improvement if needed [58, 156].

2.3 FPGA Basics and Applications in Object Detection

2.3.1 FPGA Overview

Basics

Field-Programmable Gate Arrays (*FPGAs*) are popular devices for implementing digital hardware circuits [158]. Advances in process technology have significantly boosted the logic capacity of FPGAs, making them more attractive for larger and more computationally intensive designs [159].

For a brief discussion on FPGA architecture and design approaches, refer to Appendix B.

Why FPGAs?

The hardware reconfigurability of FPGAs, combined with their capacity for optimization based on specific needs, makes them suitable for a wide range of applications, including digital signal processing [160], medical imaging [161], and communication encoding [162].

Thanks to their reprogrammable and versatile characteristics, FPGAs enable developers to customize and adapt systems to meet specific application requirements [163]. Also, unlike processors, which typically support one or multiple predefined data types, FPGAs offer the flexibility to accommodate various custom data types with different sizes within a single design. This capability allows FPGA designers to adopt mixed-precision designs tailored to the specific requirements of the target application.

Moreover, FPGAs provide true parallelism, making them well-suited for implementing sophisticated and compute-intensive algorithms while achieving satisfactory latency and throughput [164]. In addition, regarding power efficiency, FPGAs are considered appropriate platforms for achieving high performance per watt [165, 166].

Furthermore, FPGAs can simultaneously receive input data from diverse and multiple sources, each potentially working with different data types and communication standards [167]. It is an important feature, particularly in numerous real-world applications like real-time object detection, where connectivity with other sensors and cameras is often necessary [168].

An essential aspect to consider is the comparison of FPGAs with other hardware platforms. Table 2.6 provides a comprehensive comparison of FPGAs with Application-Specific Integrated Circuits (ASICs), Graphics Processing Units (GPUs), and Central Processing Units (CPUs) across various criteria. In summary, FPGAs can offer strong parallel processing capabilities, flexibility, and low latency, making them suitable for, e.g., inference phases of CNN-based object detection but with higher initial costs. ASICs, on the other hand, excel in performance, power efficiency, and latency but lack flexibility and have high development costs, making them ideal only for mass production and specific tasks. GPUs are effective for parallel processing and training phases of AI but consume more power and exhibit longer latency. CPUs are versatile, cost-effective for general-purpose tasks, and have a mature development ecosystem but are generally less efficient for parallel processing and AI training compared to GPUs and ASICs.

2.3.2 Applications of FPGAs in Computer Vision-Object Detection

Computer Vision (*CV*) tasks typically involve techniques for capturing, processing, interpreting, and comprehending digital images or videos [169]. One of the most demanding CV tasks in this

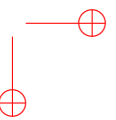
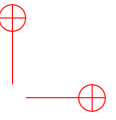


Table 2.6: Comparing FPGAs with some possible alternatives in different criteria.

Criteria	FPGAs	ASICs	GPUs	CPUs
General Performance	Good for parallel processing (more powerful than GPUs in some cases, like irregular data accesses and multiple-stream management), Limited computational and memory resources, totally customizable architecture	Superior performance in memory usage, power efficiency, latency, and throughput, Highly optimized pipelines for efficient data processing, fixed architecture	Good for parallel processing, high power consumption, fixed architecture	Good for sequential tasks, fixed architecture
Cost of the hardware	Higher initial cost but Longer life cycle compared to GPUs, cost savings due to additional capabilities integration	Cost-effective only in mass production, high initial development cost	Wide-spectrum of costs	Typically cost-effective for general-purpose computing tasks
Flexibility/Programmability	Reprogrammable for different functionalities, Adaptable to different tasks, More flexible compared to ASICs	Not reconfigurable after manufactured, suitable only for specific tasks	Not reconfigurable at the hardware level, offer software-level flexibility, less adaptable compared to FPGAs	Not reconfigurable at the hardware level, offer software-level flexibility
Power	Can achieve reasonable power efficiency	Superior power efficiency	Less power efficient	Varied power efficiency depending on workload and design
I/O Interface	Can receive input data from multiple and diverse sources	Fixed I/O types after production	Typically communicates with the host through PCIe interface, Less versatile in input sources compared to FPGAs	Typically versatile in terms of input and output sources
Development Tools	Less complex and faster development compared to ASICs, Lacks strong tools for development of AI-based systems, Developing high-performance architectures is not straightforward	Requires more complex development workflows, Extended development time	User-friendly environments with extensive libraries, Faster and more straightforward code development compared to FPGAs and ASICs	Extensive and mature development ecosystem
Latency	Can deliver deterministic and low latency, Can directly connect with peripheral hardware components, which improve latency	Can offer the lowest latency than other platforms	Exhibits longer latency due to communication with the host through PCIe interface	Typically can offer lower latency than GPUs due to a single memory system
Suitability for developing phases of CNN-based object detection	Primarily for the inference phase due to their limited computational and memory resources	Deployed for both training and inference phases	Primarily leveraged for the training phase	Used for both training and inference phases, Generally slower compared to GPUs



field, which is also our focus in this work, is Object Detection (*OD*).

A typical solution is to use efficient designs like FPGA-based systems to achieve real-time performance and meet latency constraints. In the following, we analyze the architecture of such systems. For a brief overview of FPGA architecture, refer to Appendix B.

In such systems, input data is received using one or more imaging sensors [170]. These sensors can be categorized into two types based on how their output is generated: frame-based and event-based [171]. The former can generate output, i.e., images, at specific intervals or frame rates, while in the latter, each pixel can independently respond to local alterations in light intensity that exceeds a given threshold [171]. Although the predominant focus in OD research relies on frame-based sensors, there is increasing attention to the development of algorithms for event-based computer vision [172] [171] [173]. Depending on the desired application, OD algorithms extract useful and meaningful information from the received image(s) to enable the system to make decisions or just to visualize the manipulated input data.

Figure 2.4 illustrates a typical and abstract block diagram of a computer vision system, where raw visual data is captured from cameras or sensors in the *Data Acquisition unit*. The data is then pre-processed and prepared (in the *Data Pre-processing unit*) for the next module, i.e., *Feature Extraction unit*. Data normalization, data enhancement, and color space conversion are some tasks that may be executed during data pre-processing [174, 175]. Depending on the application of interest, relevant patterns and features are then extracted to be analyzed and processed by the main OD algorithm in the next module, i.e., *Data Processing and Analysis unit*. The results can be interpreted and used in the subsequent unit to make a decision or trigger an appropriate action (*Decision Making and Control*). For instance, if an object is identified, it might trigger a robotic arm to pick it up. Finally, the results should be prepared to be transmitted to another system or utilized locally (*Output Preparation unit*).

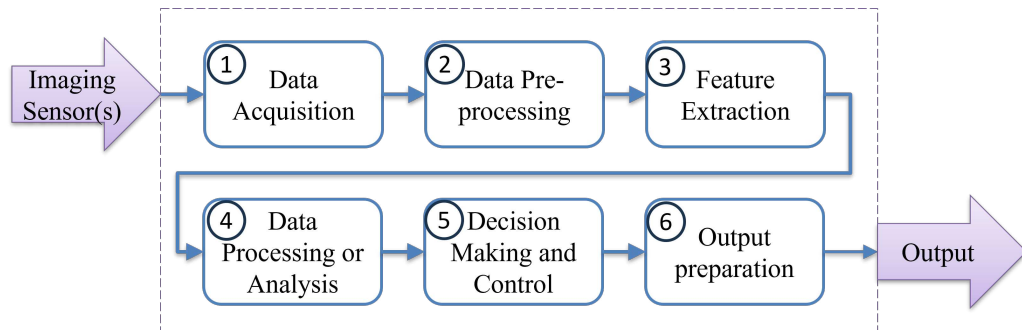


Figure 2.4: Typical structure of a computer vision system.

FPGAs or other co-processors can be integrated into the workflow shown in Figure 2.4 to accelerate the processing, resulting in an arrangement similar to Figure 2.5.

Object detection includes a broad range of applications and tasks, such as image classification, medical image analysis, and facial recognition, to name a few [176–178]. As discussed in section 2.2.1, deep-learning approaches, specifically using deep Convolutional Neural Networks, are increasingly used for accomplishing these tasks. However, they require high computational complexity and power consumption [41]. With all these in mind, as FPGAs can offer high per-

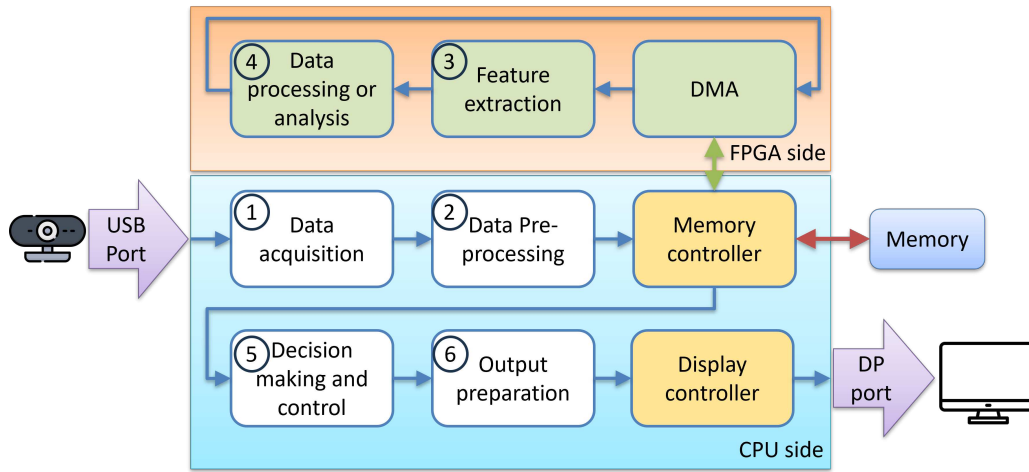


Figure 2.5: An example of FPGA co-processing architecture: images are captured by the CPU, and FPGA is used to accelerate the task. In comparison with the typical structure of a CV system shown in Figure 2.4, here, two tasks, i.e., the tasks performed in blocks 3 and 4, are done on the FPGA side. The white blocks represent the tasks run in the software. The green blocks indicate hardware units implemented on FPGA, and the light yellow blocks denote hardware units within the system.

formance and determinism as well as low latency systems, they are one of the popular platforms frequently adopted in OD systems [179, 180]. In addition, since FPGAs offer true parallelism, they are suitable for the implementation of inherently parallel, sophisticated, and compute-intensive algorithms in a way that the required latency and throughput needed in many OD tasks like real-time object detection are satisfied [181–183].

Depending on the type of deployed OD algorithm, the system specifications, and the intended application requirements, three main use cases for FPGAs in OD systems may be considered.

1. In the first use case, FPGAs can work as accelerators or co-processors alongside the main processor [179, 184], as shown in Figure 2.5. In this case, some parts of an OD system capable of being accelerated or optimized on hardware, specifically the feature extraction and the processing modules, can be offloaded to an FPGA to make the system more efficient. The input data/image is considered to be acquired through the CPU in this use case. Therefore, this architecture is commonly used when devices like “GigE Vision” feed the system [185] and USB3 cameras as their acquisition logic can optimally be executed in CPUs [184, 186]. The data can be exchanged between the FPGA and the host memory via Direct Memory Access (DMA).
2. In the second scenario, similar to the first one, the FPGA and CPU collaborate as a co-designed, heterogeneous OD system. However, in this case, the acquisition logic is implemented on the FPGA side. This architecture is particularly adopted in OD systems equipped with cameras, such as “MIPI cameras” [187], for which acquisition logic is easily implemented on FPGAs [188]. In this scenario, the data pre-processing unit is typically

implemented in the FPGA, resulting in reduced data transfer and improved overall performance [188].

3. Another use case is an OD system based only on an FPGA [93]. In a more complex scenario, a cluster of FPGAs can be considered. This architecture is more suitable for applications in which communication with connected devices through FPGAs is easier, and hardware implementation of the entire system units can result in better or at least similar to their software counterparts, providing a system with low latency and significantly high performance, in term of both power and throughput [93].

2.3.3 FPGA Platforms for Accelerating Object Detection

Although there are several FPGA suppliers worldwide, the FPGA market is dominated by AMD (formerly Xilinx) and Intel (formerly Altera) [183]. These two companies provide a broad range of FPGA devices for applications ranging from aerospace to data centers [193, 194].

Over the past decade, following significant breakthroughs in AI and its applications across various fields, such as object detection, there has been a growing trend to enhance the performance of AI-based systems using hardware accelerators.

In response, FPGA manufacturing companies have introduced a wide range of FPGA devices to facilitate the design and deployment of FPGA-based AI systems [193, 194]. Table 2.7 lists commonly used FPGA boards in recent works for the implementation of real-time object detection models, detailing their available resources. They continuously introduce high-performance FPGAs enhanced with AI capabilities and tools tailored for this new era. Key characteristics of these FPGA devices include better performance per watt, higher memory bandwidth to alleviate bottlenecks in memory-bound AI-based object detection models, and dedicated units for executing compute-bound AI models more efficiently [193–195].

In addition, these FPGA manufacturers offer various tools and frameworks to bridge the gap between AI model development and FPGA design flow. By utilizing these frameworks, designers can seamlessly develop object detection models using popular libraries like TensorFlow and PyTorch, evaluate the preliminary models, optimize them for the target FPGA device, and compile the final models for integration into the FPGA design [193, 194].

2.4 FPGA-based Designs

This section begins by examining two primary hardware architectures—single computation engine and streaming architectures—commonly considered for FPGA-based object detection systems. It then explores various acceleration techniques, categorized into model-related and implementation-related approaches, that can enhance performance across key metrics, including throughput, accuracy, and power consumption.

Following this, the section assesses the impact of each acceleration technique on system performance, particularly within real-time object detection contexts.

Table 2.7: A list of some FPGA boards used in recent works for implementing Real-time object detection models, sorted by number of LUTs/ALMs.

Board	FPGA Part	Number of Available Resources				Reference
		LUTs/ALMs ^a	FFs	DSPs	Memory Blocks (Size)	
ZC702	ZYNQ XC7Z020-1CLG484	53,200	106,400	220	140 (36Kb)	[97, 189]
Pynq-Z1	ZYNQ XC7Z020-1CLG400C	53,200	106,400	220	140 (36Kb)	[190]
Ultra96	Zyng UltraScale+ MPSoC ZU3EG A484	70,560	141,120	360	216 (36Kb)	[11, 191]
ZC706	ZYNQ XC7Z045 FFG900 - 2 SoC	218,600	437,200	900	545 (36Kb)	[57, 58, 78]
ZCU104	Zyng UltraScale+ MPSoC ZU7EV	230,400	460,800	1,728	312 (36Kb)	[192]
KU040	XCKU040-1FBVA676	242,400	484,800	1,920	600 (36Kb)	[92]
ZCU102	Zyng UltraScale+ XCZU9EG-2FFVB1156	274,080	548,160	2,520	912 (36Kb)	[8]
VCT07	Virtex 7 XC7VX485T-2FFG1761C	303,600	607,200	2,800	2,060 (36Kb)	[7]
Arria-10	Arria 10 GX1150	427,200	1,708,800	1,518	2,713 (20Kb)	[56, 82, 83, 85]
DE10-PRO	Stratix 10 GX2800	933,120	3,732,480	5,760	11,721 (20Kb)	[77]

^a For Arria 10 and Stratix 10 FPGA parts, the number of Adaptive Logic Modules (ALMs) is reported. Each ALM block has an eight-input adaptable LUT, two adders, and four FFs.

Table 2.8: A summary of two FPGA architecture design approaches for implementing object detection models.

Architecture Type	Highlights	Related Works
Single computation engine	Uses one single computation block to perform all parts of the model, prioritizes flexibility over customization, uses time-shared computation resources, adaptable to various object detection models, less efficient, bottleneck in memory bandwidth (Figure 2.6)	[57,78,82,196–198]
Streaming	Uses specialized blocks tailored for each part of the model, can take advantage of pipeline design, less adaptable to different models, faster, more resource consumer (Figure 2.7)	[7, 58, 97, 199]

2.4.1 FPGA architecture design approaches

The hardware architectures employed in FPGA-based object detection systems can be categorized into two main types, mirroring the classification outlined in [3]: *single computation engine* and *streaming architecture* (see Table 2.8).

Single Computation Engine Architecture

The single computation engine design approach, also known as the *one-size-fits-all*, prioritizes flexibility over customization [3]. It takes advantage of a single computation engine, often in the form of a systolic array of Processing Elements (*PEs*) [200], to perform all layers of an object detection system sequentially [3,57].

Figure 2.6 shows an example of adopting a single computation engine accelerator in which one operation unit is configured and deployed for performing all layers one after the other in time (not in space). The data is transferred between the FPGA and host memory via the DMA unit. There is also a “control unit” on the FPGA side, responsible for managing all executions. This unit receives instructions from the host side, where the central operational controller exists.

Single computation engine-based architectures temporally share common computation resources across different layers [58]. In such a scenario, the unit is time-shared, but we save resources compared to a pipelined design [196].

The computation engine is configurable to accommodate the specific characteristics of each layer. This adaptive approach reduces resource utilization and enables the implementation of any model, provided it utilizes layers supported by the engine [197, 198].

This structure offers the advantage of easy adaptation to various object detection models [151]. Nevertheless, this flexibility comes at the cost of reduced efficiency and varying and inconsistent performance when employing different models [201]. In a related experiment, Guo et al. [197] implemented YOLO [1] and Face Alignment models on Angel-Eye, a CNN accelerator featuring a single computation engine architecture, using Zynq XC7Z045. Compared

Streaming Architecture

On the other hand, streaming architectures are typically formed with a unique hardware block for each part of the targeted object detection model [7, 58], developing the architecture in space (rather than in time) (Figure 2.7).

Each computational block is optimized individually to exploit the inherent parallelism within its corresponding layer [3]. As depicted in Figure 2.7, in this design approach, once fetched from memory, image data passes sequentially through all dedicated hardware units for full processing until the final result is obtained and written back to memory. A pipelined architecture, where heterogeneous computational blocks are interconnected, can facilitate this data flow. This configuration enables the simultaneous execution of different layers, significantly enhancing overall system performance [7, 199].

However, it is crucial to design the PE for maintaining uniform processing times across individual layers to avoid idle states [151]. In this architecture, individual PEs are tailored and fine-tuned for each layer, creating an optimized structure specific to a particular object detection model. However, this specialization makes it less adaptable than a single computation engine architecture. Additionally, transitioning it to a new FPGA device can be challenging or even unfeasible due to the difficulty in customizing the implemented structure, especially when altering its size [3] [151]. Although pipelined accelerators can offer faster processing speeds due to reduced memory transfers, they incur higher hardware costs than one-size-fits-all configurable engines. This is because all individual layers must be simultaneously mapped in the hardware, making it less versatile and more resource-intensive [198].

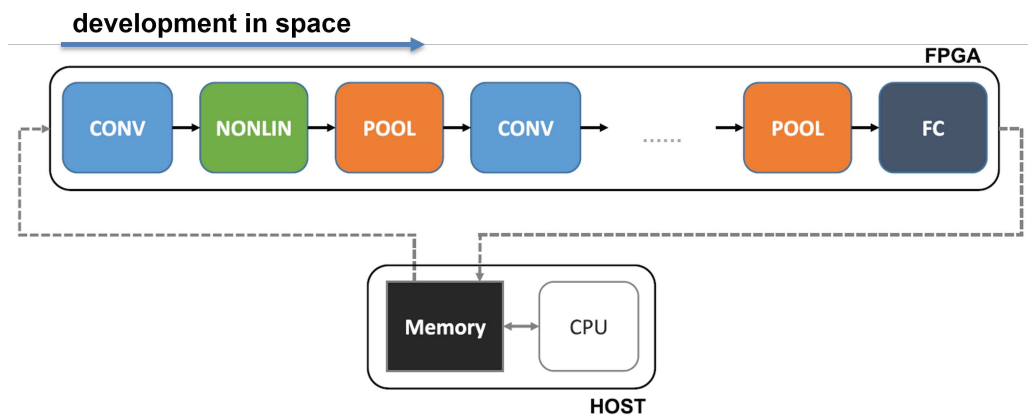


Figure 2.7: Example of an FPGA accelerator structure with a streaming architecture. Each block is individually designed and optimized to perform the required calculations of each layer. (Figure adapted from [3])

More Discussion

As mentioned, both architectures rely on PEs as their fundamental building blocks to execute computational tasks, primarily consisting of Multiply-ACcumulator (MAC) units [203]. In scenarios where the results of various input channels of a CNN need to be accumulated without mul-

tiplication, costly multipliers remain unused. To improve resource utilization efficiency, particularly in streaming architectures, a specialized Convolutional PE (*Conv-PE*) could be employed [151]. The Conv-PE comprises a few multipliers based on the kernel size, succeeded by an adder tree to sum the intermediate results obtained from each input channel, typically followed by a non-linearity unit [204, 205]. Moreover, using Conv-PEs reduces latency by requiring fewer pipelined stages, helping object detection systems satisfy real-time constraints.

On the other hand, in a one-size-fits-all architecture, frequently reading or writing parameters and intermediate calculation results to external memory creates a bottleneck in memory access and bandwidth. This results in reduced system performance and increased power consumption. In real-time systems, in particular, to compensate for it, additional buffers can be utilized to store weights and input data on available on-chip memory [206].

2.4.2 Hardware Acceleration Techniques for FPGA-based Object Detection

In the context of real-time object detection, given their heavy reliance on CNN-based models, as discussed in section 2.2.1, there is significant interest in designing and deploying suitable hardware accelerators to enhance system performance.

Generally speaking, an accelerator refers to a specialized hardware component designed to execute a specific set of tasks more efficiently in terms of performance and power consumption compared to a general-purpose processor or CPU [151, 207]. The development of floating point co-processors was one of the first attempts to adopt accelerators [208]. Since then, designers are increasingly taking advantage of hardware accelerators to enhance digital systems' power efficiency and performance. However, they have gained more attention in recent years, mostly due to AI breakthroughs, especially in DNNs [156].

Demands for powerful hardware accelerators have increased as AI-based applications have become more prevalent in various fields. Enhancing power efficiency, improving performance, and addressing model complexity of modern AI models are the most important reasons why developing accelerators is being paid remarkable attention in this era [49, 151]. Recent trends to solve more challenging tasks more accurately in different fields, such as computer vision and natural languages, through AI-based applications, have made researchers introduce more complex and larger models year by year. Consequently, such models cannot be exclusively run on general-purpose processors, as they lack the computational power necessary for both training and inference within reasonable time frames. It becomes a more important issue when developing real-time applications.

Furthermore, as AI models continue to grow in size, the need for increased memory access and data movement also rises. Memory access is notably more energy-intensive compared to arithmetic computations [156, 230]. Considering the limited capacity of on-chip storage in general-purpose processors and the need for external memory access, power efficiency is the most important reason for developing hardware accelerators. In addition to incorporating specialized hardware features to accelerate intensive computations, AI accelerators can be designed to reduce memory access and provide larger on-chip caches for improved performance.

FPGAs are highly regarded as suitable platforms for accelerating object detection because they can perform parallel and pipelined computations efficiently while maintaining high power

Table 2.9: Commonly used Hardware Acceleration Techniques for FPGA-based Object Detection.

1. Model-related Techniques for Hardware Acceleration	FPGA-related references	2. Implementation-related Techniques for Hardware Acceleration	FPGA-related references
a. Pruning	[57, 82, 209–212]	a. Data Reuse	[8, 11, 78, 191, 213, 214]
b. Quantization	[7, 8, 10, 11, 57, 58, 82, 191]	b. Mathematical-based Optimization	[78, 215–218]
c. Knowledge Distillation	[9, 219–221]	c. Systolic Arrays	[8, 222, 223]
d. Hardware-aware Neural Architecture Search	[10, 190, 224–226]	d. Code Modification	[7, 11, 57, 82, 97, 99]
		e. Roofline-based Optimization	[57, 82, 209, 217, 227, 228]
		f. Pipelining	[7, 8, 57, 58, 83, 90, 229]

efficiency. Moreover, FPGAs’ architectural flexibility and reconfigurability enable the implementation of custom logic units tailored to specific tasks in different object detection models. This flexibility can also make it possible to apply various optimization techniques, both at the hardware and software levels, aimed at improving system performance through architectural enhancements and object detection model simplifications. Therefore, FPGAs have emerged as the optimal hardware platform, offering both speed and power efficiency for implementing complex object detection models and accelerating them at the edge. Readers interested in a brief overview of FPGA architecture and common design approaches can refer to Appendix B.

This section delves into some commonly adopted techniques for designing FPGA-based hardware accelerators. Table 2.9 shows the methods discussed in the following two subsections, divided into two main categories: Model-related and implementation-related techniques. The former includes the methods adopted to prepare and optimize object detection models for FPGA implementation, while the latter explains techniques applicable during hardware architecture design and FPGA implementation. It is crucial to note that accuracy and execution time are two pivotal parameters in real-time systems, particularly in hard real-time scenarios [61]. Therefore, achieving a balance between these factors is of utmost importance.

Model-related Techniques for Hardware Acceleration

Some typical techniques that can be used to make the implementation feasible in hardware or allow for more efficient execution of the target application are discussed in detail below, namely:

a) Pruning; b) Quantization; c) Distillation; d) Hardware-aware Neural Architecture Search.

Pruning Pruning is typically performed in software before deploying the model on hardware. However, performance can be further enhanced if the pruning process accounts for the capabilities and custom hardware features of the target FPGA, as highlighted in some of the works discussed below.

Pruning is defined as identifying and eliminating neurons, kernels, weights, and channels that have minimal or negligible impact on the final accuracy of an AI model to reduce network complexity [231]. It offers several advantages, including reducing computational load and required memory and improving accuracy per parameter and per operation [157, 209]. To compensate for the possible effect of this on the accuracy, the pruned model, also known as a sparse network [231], needs to be retrained [232]. To avoid slow convergence during the retraining of a sparse network, pruning should be conducted incrementally, with each group of layers pruned in a separate stage [231].

Pruning methods can be broadly classified into two categories [4]: *unstructured pruning* and *structured pruning*. Figure 2.8 illustrates the distinction between these approaches in a fully connected layer.

In unstructured pruning [233–235], weights with low sensitivity are selectively removed throughout the network. This method allows for aggressive pruning, removing a significant portion of neural network parameters with insignificant accuracy loss. Wang et al. (2020) [82] apply unstructured pruning to YOLOv2 [80], demonstrating its effectiveness in enabling real-time object detection on FPGAs. Detailed results are provided in Sections 2.5.1 and 2.5.2. However, unstructured pruning leads to sparse matrix operations, which are challenging to accelerate and are often memory-bound [236].

Structured pruning [196, 237, 238], on the other hand, involves removing a group of parameters, such as the entire kernel. This approach alters the input and output shapes of layers and weight matrices, allowing for dense matrix operations to continue. However, aggressive structured pruning often results in considerable accuracy degradation. Achieving state-of-the-art performance during training and inference with high levels of pruning remains an open challenge [239].

Pruning may result in data sparsity and inefficient load balancing, mainly because the target device characteristics are not considered. To address these issues, Ramhorst et al. [210] introduce an FPGA resource-aware structured pruning algorithm capable of capturing, during training, the underlying mapping of network weights to computational and memory resources.

To enhance inference speed, Plochaet et al. [211] propose a pruning method suitable for FPGA-based AI accelerators in which the hardware constraints are considered.

Sui et al. [212] concentrate on row pruning to introduce a hardware-friendly pruning technique. They eliminate all rows except one for each convolution kernel and skip all zero calculations during FPGA deployment.

When developing real-time object detection systems on FPGAs, pruning can be adopted to enhance throughput, reduce latency, optimize resource utilization, and improve power efficiency, albeit with a potential trade-off in accuracy [57, 82] (cf. Section 2.5.1).

Quantization Data quantization is another commonly employed method to decrease the size of CNN-based object detection models. This approach involves replacing conventional 32-bit

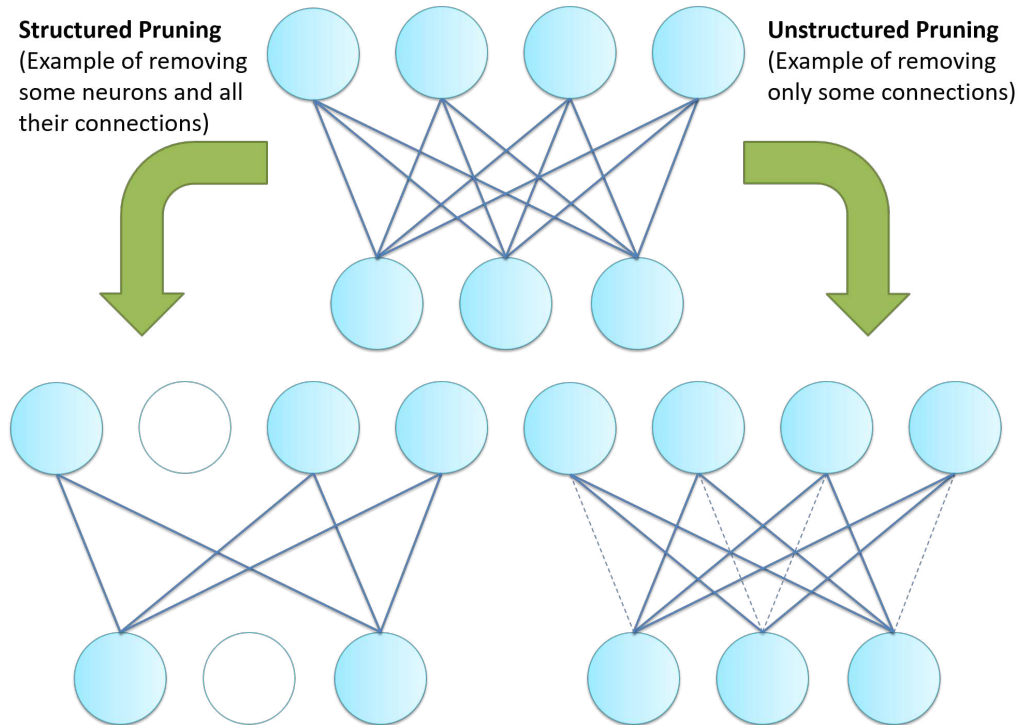


Figure 2.8: Example of applying structured and unstructured pruning on a fully connected layer. Removed (pruned) neurons and weights are shown by white circles and dotted lines, respectively.

floating-point weights and activation data with simpler representations, such as lower-bit floating-point or fixed-point numbers.

Therefore, to achieve this objective, FPGAs give the designer the maximum flexibility for choosing an arbitrary number for bits in the arithmetic operations and data representation.

Quantization can be deployed in both training and inference phases. However, despite its great achievements in training [240, 241], the majority of recent research on quantization has primarily concentrated on inference [4]. Often combined with pruning, quantization plays a crucial role in achieving an efficient hardware implementation. It reduces the required memory capacity and bandwidth while simplifying arithmetic operations. Quantization can be applied to both weights and activations. However, adjusting the bit width of weights generally has a smaller impact on accuracy compared to modifying the bit width of activations [203].

Data quantization can be applied *uniformly* [242, 243] or *non-uniformly* [244–247]. In the former, quantization levels are uniformly spaced, whereas in the latter, they do not necessarily need to be uniformly spaced. Uniform quantization is more popular due to its simplicity and efficient mapping to hardware [4]. For instance, in a successful effort to develop FPGA-based real-time object detection, Nguyen et al. [7] shows the effectiveness of adopting uniform quantization for activations. This approach enhances the design's speed and power efficiency by eliminating the need for external memory access. Detailed results are provided in Sections 2.5.1 and 2.5.2.

However, the dynamic range of data across various layers in a CNN tends to be large. Con-

sequently, using uniform quantization with a fixed-point data format for all layers may result in significant performance degradation [197].

On the other hand, non-uniform quantization offers the potential for higher accuracy within a fixed bit-width framework. This is because it allows for more effective capture of distributions, emphasizing crucial value regions and identifying optimal dynamic ranges [4].

One possible and hardware-efficient approach is using mixed-precision quantization, in which different bit precision is considered for each layer of CNN based on the sensitivity of that layer to quantization. This can address the accuracy loss problem in object detection systems, which commonly happens in low-precision quantization, particularly below 8-bit [248, 249]. There exist successful examples of adopting this approach in developing efficient real-time object detection systems on FPGAs [8, 10] (cf. Section 2.5.1). However, finding an optimum solution to decide on the best precision of each layer is an open challenge [4].

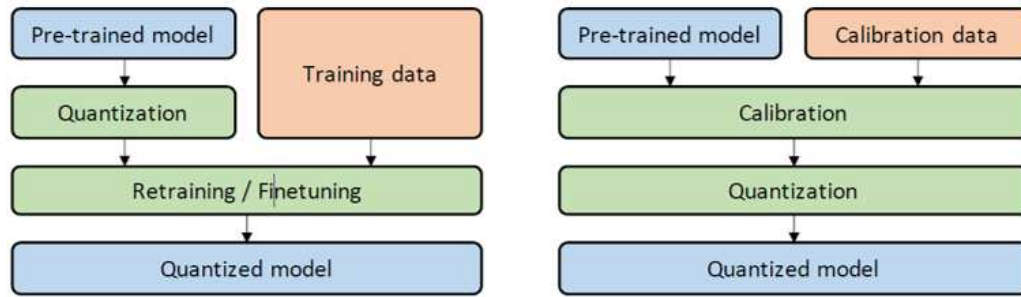


Figure 2.9: Illustration of Quantization-Aware Training (QAT) (Left) and Post-Training Quantization (PTQ) (Right) procedures. In QAT, a pre-trained model is quantized and then fine-tuned. In PTQ, a pre-trained model is calibrated using calibration data to do quantization based on the calibration result. (Figure adapted from [4].)

Data quantization can lead to a drop in model accuracy, which is not acceptable in some applications, such as hard real-time object detection systems. To mitigate this issue, the model parameters should be tuned or adjusted when applying quantization. In this regard, Quantization Aware Training (QAT) and Post-Training Quantization (PTQ) are two main approaches commonly used to achieve a more accurate and efficient quantization [4]. A pre-trained model undergoes quantization in the former, followed by fine-tuning using training data. This process adjusts parameters and aims to recover any accuracy degradation caused by quantization [4]. However, the retraining procedure can be time-consuming, especially when it comes to low-precision quantization. In PTQ, on the other hand, there is no need to model retraining. A pre-trained model is calibrated using calibration data, which is typically a small subset of training data. This calibration process computes the clipping ranges and scaling factors. Clipping ranges define the upper and lower bounds within which input data are constrained while scaling factors adjust the dynamic range of real-valued input data to adapt them to the desired output range. Subsequently, the model is quantized based on the calibration results [250]. Even though the process of obtaining the desired model is faster in the PTQ approach, this often results in lower accuracy compared to QAT [4], which may not be desirable in hard real-time applications, like object detection, where accuracy is critical. Figure 2.9 shows the overall procedure taken in QAT and PTQ.

In implementing object detection models on FPGAs, quantization is employed in various forms with different levels of precision, especially when real-time performance is required [7, 10, 57, 58, 191] (cf. Sections 2.5.1 and 2.5.2). Like Pruning, the primary goals of deploying this technique in such systems are improving throughput, latency, resource utilization, and power efficiency, with possible consequences in accuracy drop.

Knowledge Distillation Knowledge distillation is the process of transferring knowledge from a large and complex model or ensemble of models to a single smaller model, which can be feasibly deployed under real-world constraints.

The deployment of large object detection models poses a significant challenge, particularly for edge devices like FPGAs, which have restricted memory and computational resources. To address this challenge, a model compression method was initially proposed [219] to transfer knowledge from a large model into a smaller model without sacrificing performance. This process of training a small model from a larger one was formalized as a “Knowledge distillation” framework by Hinton et al. [220].

In knowledge distillation, as depicted in Figure 2.10, a small “student” model is trained to emulate a large “teacher” model. By leveraging the knowledge from the teacher, the student model aims to improve its accuracy.

Based on the discussions thus far, simplified models obtained through pruning and quantization techniques can be considered student models. By deploying knowledge distillation, we can bridge the accuracy gap between the simplified and original models, bringing them closer in performance.

In [221], for instance, a distillation technique tailored for quantized models was proposed, based on which quantized student networks attain accuracy levels comparable to their full-precision teacher model counterparts, with lower inference time.

In the context of real-time FPGA-based object detection, knowledge distillation is deployed to enhance the system’s accuracy, especially after applying model compression techniques such as pruning and quantization [9] (cf. Sections 2.5.1).

Hardware-aware Neural Architecture Search When considering resource-constrained platforms like FPGAs for developing DNN-based systems, optimizing the DNN model to be streamlined and compact is crucial. This optimization should encompass reducing the number of parameters and computations while maintaining acceptable accuracy. However, achieving an efficient network tailored to the characteristics of the target hardware is a challenging and time-consuming endeavor [251].

Hardware-aware Network Architecture Search *Hardware-aware NAS* or (*HW-NAS*) seeks to automate the process of discovering the most optimal architectures and configurations of a DNN, tailored specifically for a given hardware platform [224]. This process aims to strike a balance between accuracy and performance, ensuring acceptable tradeoffs [224]. Furthermore, through the consideration of the distinctive attributes of the target device and the implementation of multi-objective optimization algorithms, HW-NAS can generate hardware architectures that are inherently more compatible with the target system and more efficient in terms of performance and resource utilization [225].

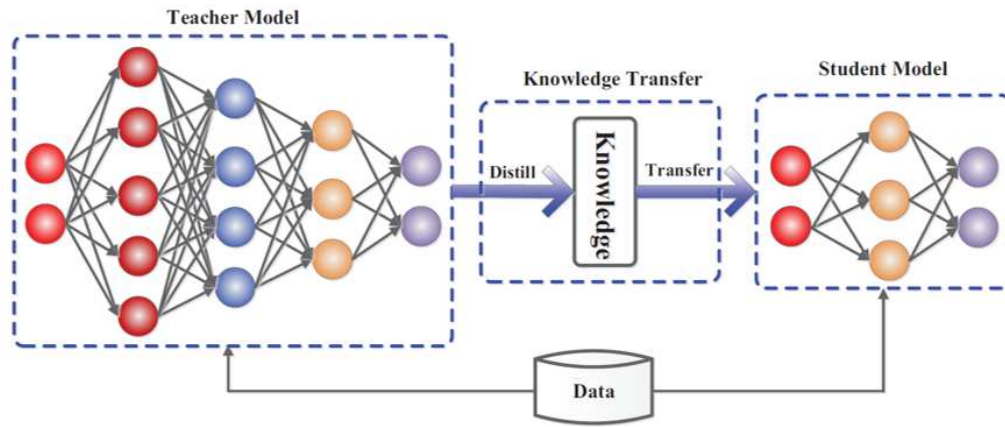


Figure 2.10: An illustration of performing knowledge distillation with a teacher-student framework. (Figure adapted from [5].)

HW-NAS has notably influenced image classification and object detection tasks by consistently achieving state-of-the-art results [225]. Using reinforcement learning [252], Jiang et al. [226] introduced an FPGA-aware NAS technique, called FNAS, tailored to identify architectures that can meet specified inference latency requirements. By employing a performance abstraction model, it can estimate neural network latency without the need for extensive training. FNAS effectively eliminates networks that do not align with the constraints in the search space, boosting search efficiency by a remarkable factor of 11.13 \times .

HW-NAS can effectively optimize all critical metrics of FPGA-based object detection systems, including throughput, accuracy, latency, resource utilization, and power efficiency. By leveraging the flexibility of FPGA designs alongside the efficiency of HW-NAS, the optimization process for object detection models becomes significantly more streamlined and impactful. This synergy enables the development of low-latency models, which are essential for achieving real-time performance [190] (cf. Section 2.5.1).

Implementation-related Techniques for Hardware Acceleration

Data Reuse An object detection model often requires access to a large amount of data, usually stored in external memory when the target platform is a resource-constrained device. Particularly with complex and large state-of-the-art object detection models, this strong dependency on access to external memory and data movement can lead to a bottleneck for performance, power, and computation efficiency [156, 253].

It is possible to reduce memory access by reusing pre-fetched or intermediate data, such as feature maps, weights, and convolution internal results, multiple times [253]. This can be achieved by leveraging the available on-chip memories of the target device or optimizing the algorithms. Also, convolution operations can be accelerated by maximizing data reuse [254]. This method involves utilizing pixels computed simultaneously at the same spatial position across various Output Feature Maps (OFMs). As a result, all pixels from Input Feature Maps (IFMs) and kernel data are accessed only once and stored in an on-chip BRAM of the FPGA (cf. Sec-

tion 2.3.1) until they are reused.

By efficiently employing loop transformation and BRAM-based on-chip buffers to leverage data locality, Beric et al. [255] demonstrated that it is possible to achieve an 11× speedup compared to the standard implementation on similar FPGA resources while also being more power-efficient due to reduced memory access.

Data reuse methods can be categorized into temporal and spatial reuse [156], both of which are critical for optimizing FPGA-based real-time object detection systems. In temporal reuse, one group of data is utilized multiple times by a single computational unit, allowing on-chip buffers to store a small set of required data, which minimizes memory access latency and power consumption. In spatial reuse, one set of data is simultaneously employed by different processing units, enabling parallelism and improving throughput without requiring additional buffering. These techniques are essential for maximizing the efficiency of FPGA implementations in real-time object detection applications. There are plenty of works in the context of FPGA-based real-time object detection adopting data reuse to achieve higher performance in terms of throughput and latency along with gaining more power efficiency [7, 8, 191] (cf. Section 2.5.1).

Mathematical-based Optimization To enhance performance by reducing computational complexity or memory access, various computation or transformation techniques can be employed, such as the Winograd transform [256], Fast Fourier Transform (FFT) [257].

Since its initial application in accelerating convolution execution in 2016 [258], the Winograd technique has garnered considerable attention in CNN implementations. Regarding FPGA-based implementation, there have been significant research achievements in neural network optimization based on the Winograd algorithm [215–218].

Winograd aims to reduce the number of multiplications using some pre-calculated values, considering the fact that the trainable parameters of a CNN model will remain fixed after training. In fact, it transforms overlapping kernels into non-overlapping ones to reduce the computation complexity of CNNs [256, 258]. As an example, Equation 2.3 shows that when convolving an input feature map I_n and a kernel W_n , both of size 1×3 , the traditional approach requires 6 multiplications and 4 summations.

$$\begin{bmatrix} I_0 & I_1 & I_2 \\ I_1 & I_2 & I_3 \end{bmatrix} \begin{bmatrix} W_0 \\ W_1 \\ W_2 \end{bmatrix} = \begin{bmatrix} O_0 \\ O_1 \end{bmatrix} \quad (2.3)$$

However, adopting the Winograd approach, the number of multiplication and summation will be 4 and 12, respectively, as shown in equations (2.4) and (2.5).

$$\begin{aligned} M_1 &= (I_0 - I_2)W_0 \\ M_2 &= (I_1 + I_2) \frac{W_0 + W_1 + W_2}{2} \\ M_3 &= (I_2 - I_1) \frac{W_0 - W_1 + W_2}{2} \\ M_4 &= (I_1 - I_3)W_2 \end{aligned} \quad (2.4)$$

The key insight in these calculations is that since the weights are fixed, the summations of those weights in Equation 2.4 can be pre-computed, so we can avoid 4 additions. Additionally, division by 2 only requires a shifter rather than any arithmetic logic.

$$\begin{bmatrix} O_0 \\ O_1 \end{bmatrix} = \begin{bmatrix} M_1 + M_2 + M_3 \\ M_2 - M_3 - M_4 \end{bmatrix} \quad (2.5)$$

Therefore, the actual number of multiplications and adders required to implement the convolution in equation (2.4) are 4 and 8, respectively. This shows that employing this method reduces the number of multipliers from 6 to 4 for the same calculation compared to the traditional method. It is important to note that the hardware implementation of multiplication is typically more resource-intensive than that of summation, which justifies the increase in the number of required adders from 4 in Equation 2.3 to 8 for the implementation of Equation 2.4. Furthermore, in FPGAs, addition can be performed faster than multiplication [253].

However, the Winograd method limits the reuse of the weights because the pre-computed weight sets in Equation 2.4 are just used once while moving the convolution window. This issue can be partially addressed by adopting appropriate considerations during the design, such as integration of the pooling with convolution, to leverage the benefits of the Winograd method in reducing power consumption and increasing throughput [151].

Cai et al. [78] demonstrate the effectiveness of employing this technique to develop an efficient FPGA-based object detection system, optimizing both speed and power consumption. Further details about this effort are provided in Section 2.5.1.

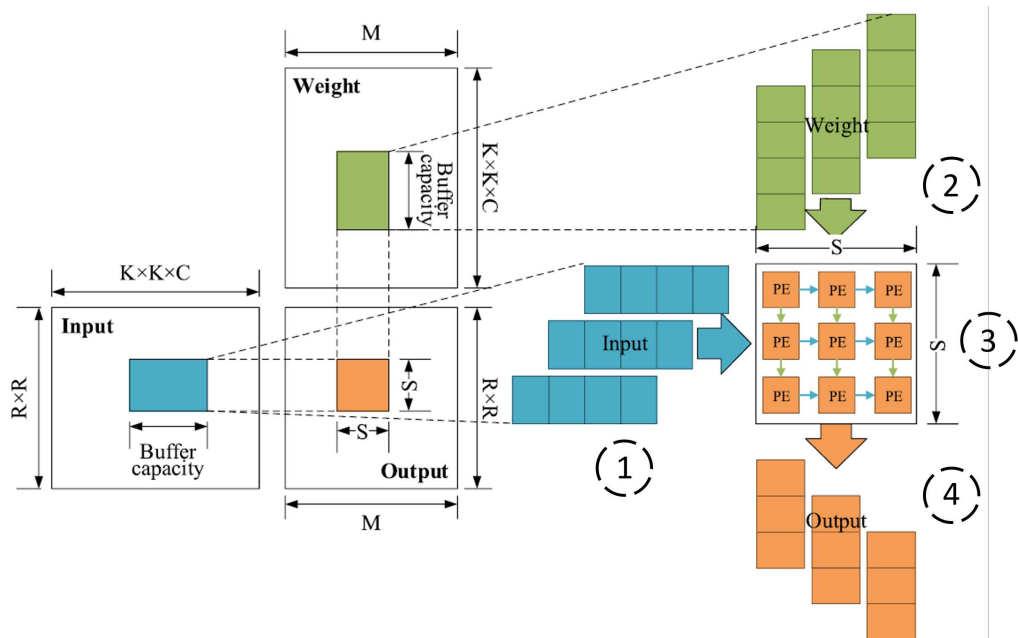
Another possible acceleration technique to perform convolution operations is using FFT [156, 257]. This method is similar to the Winograd transform and can reduce the number of required multipliers. The principle behind this approach is that convolution in the time domain can be transformed into multiplication in the frequency domain. Initially, we compute the FFT of both the weight and input. Subsequently, we obtain the output by taking the inverse FFT of their element-wise multiplication in the frequency space. As a result, there is a possible reduction in the number of multiplications for each input channel, strengthened from the order of $RSPQ$ to $PQ \log_2(PQ)$, where $P \times Q$ represents the output size and $R \times S$ signifies the filter size [156].

However, the advantages of this approach diminish with larger kernel sizes [203]. Additionally, combining FFT with sparsity, which often yields higher benefits, can be challenging [156]. Given that modern CNNs predominantly utilize small kernel sizes, FFT has declined in significance in modern hardware accelerators [151].

Systolic Arrays First introduced by Kung and Leiserson [200], systolic arrays are parallel computing architectures comprising a network of Processing Elements (*PEs*) organized in a regular grid, typically with a two-dimensional structure. In these arrays, data flows rhythmically between neighboring PEs, enabling efficient and synchronized computations. Some of the key features of the systolic array design include regularity, reconfigurability, and scalability [259].

When data is fetched from memory, it is sequentially passed from one PE to another, enhancing data reuse while minimizing memory access—a key advantage for power-efficient, real-time systems. Furthermore, systolic arrays demonstrate exceptional adaptability to a variety of DNN-based models and are capable of achieving high levels of parallelism and clock frequencies [206]. These attributes make them an ideal choice for developing real-time object detection systems.

Considering its advantages, such as scalability and flexibility, systolic array architectures have been adopted to implement different computations on FPGAs [222, 223]. Figure 2.11 represents an example of using systolic array-based architectures for performing matrix multiplication. The regular and predictable architectures of systolic arrays make them highly suitable for hardware implementation. In addition, conventional approaches to data spatial reuse in FPGA computation units often encounter challenges such as significant fan-out and difficulty in meeting timing constraints. On the other hand, each PE in a systolic array architecture can seamlessly transmit result data to its neighboring PEs. These localized connections between PEs effectively diminish the necessity for transmitting extensive data over extended distances, thereby enhancing system performance and reducing latency.



1	Each element of the input matrix rows is fed and propagated horizontally across PEs on each clock cycle, with a one-clock-cycle delay for each new row.
2	Each element of the Weight matrix columns is fed and propagated vertically across PEs every clock cycle, with a one-clock-cycle delay for each new column.
3	Each PE performs partial multiplications and accumulates the results over successive clock cycles, processing data as it flows horizontally and vertically
4	Each row of the result exits vertically, with a one-clock-cycle delay from the previous row.

Figure 2.11: An example of adopting a systolic array to perform matrix multiplications in CNNs. (Figure adapted from [6])

Code Modification To maximize hardware usability by optimizing object detection models on the operator level, some hardware-specific code transformation techniques can be deployed. These techniques include loop optimization and operator fusion.

Loop optimization techniques, such as loop unrolling and loop tiling, can improve resource utilization efficiency and data reusability. This results in lower latency and increased power efficiency by reducing the frequency of off-chip memory accesses. [259]. Given that matrix multiplication serves as a main computational component within CNN-based object detection models, it can be conceptualized as nested loops. Therefore, loop modification techniques are frequently employed in hardware accelerators to optimize performance [209, 260–262].

when designing real-time object detection systems on FPGAs, techniques such as loop reordering, loop unrolling, and loop pipelining can be employed to accelerate execution. Additionally, loop tiling is often utilized to optimize memory allocation, improving overall system efficiency and performance [151, 263].

Loop reordering aims to minimize redundant memory accesses and maximize the utilization of on-chip memory. It involves rearranging the sequence of uncoupled multiplication and additional operations within the convolution process to optimize efficiency [263, 264]. Loop unrolling is a technique where multiple copies of the loop body are created to allow some or all iterations to execute in parallel, improving data access and throughput. This technique can be deployed to optimize hardware utilization, reducing computation time effectively and maximizing data reuse by fully leveraging internal memory [260, 261]. Loop pipelining is a technique that allows the overlap of iterations in a loop by starting a new iteration before the previous one completes, enhancing execution efficiency and throughput. However, whether pipelining is feasible depends on the hardware accelerator's architecture (cf. Section 2.4.1). The concept of loop tiling, or loop blocking, involves dividing the loop's iteration space into smaller blocks to enhance memory hierarchy efficiency [263]. By doing so, the loop's data fits within the cache until it is reused, thus reducing cache misses and enhancing performance significantly.

In the context of the FPGA-based real-time object detection, adopting these techniques is popular (cf. Section 2.5.1). Nguyen et al. [7] use loop reordering and tiling to optimize the data path of the deployed streaming architecture. To enhance data reuse, thereby optimizing memory access, Babu et al. [97] deploy loop tiling. Additionally, they leverage loop pipelining to enhance system throughput. For a different purpose, Fan et al. [57] employ loop unrolling and loop reordering to tackle the issue arising from the distinct kernel sizes of depth-wise and point-wise convolutions. This difference makes it difficult to directly utilize processing elements (*PEs*) designed for depth-wise convolution with point-wise convolution.

The operator fusion technique, also known as kernel or layer fusion technique, aims to reduce data exchanges between computation units and external memory by finding data dependencies between different parts of an object detection model [265, 266]. This process begins by partitioning the model into groups of operators to be fused. Within each group, intermediate feature maps are promptly accessible to subsequent units or layers that need the data. However, if none of the computation units in a group require the data at that time, it is transferred to external memory. It will then be reloaded into on-chip memory when another computation unit needs it. In designing FPGA-based object detection systems in particular, employing the operator fusion technique can offer significant advantages in improving the performance of the system and reducing memory access overhead [82, 85, 99]. For example, in their proposed FPGA design for real-time object

detection based on YOLOv2 [80], Xu et al. [85] merge different components of the model, including convolution, batch normalization, and activation function units, to minimize memory access. This approach enables them to achieve a high throughput of up to 71 FPS for the system (see Section 2.5.2 for further details).

Roofline-based Optimization When deploying an object detection model on a particular platform, such as an FPGA device, the achievable performance of the model is influenced by the compatibility between the model and the target platform [253]. In this regard, performance bottleneck analysis is a critical task that can be conducted using various methods [217,227]. Among these methods, the “Roofline Model” [267] stands out as particularly popular. The Roofline model is a visual performance model that enables the estimation of the gap between the actual performance (in terms of computational performance (FLOPS) and memory bandwidth (Bytes/s)) of a compute kernel compared to the peak theoretical performance of the underlying system.

A more detailed discussion on the Roofline Model is presented in AppendixC.

In recent years, the Roofline model has been deployed to optimize the CNN implementations on FPGAs in various works [209,228]. Some of these efforts have focused on using this technique combined with the High-Level Synthesis (HLS) design approach [268–270]. Calore et al. [271] optimized the neural network compiler through pipeline design and utilized the Roofline model to investigate the trade-off between memory and computational throughput, aiming to enhance FPGA performance. In particular, to implement real-time object detection models on FPGAs, many works have adopted this technique to find possible performance bottlenecks and optimize the design [57,272] (cf. Section 2.5.1).

Pipelining In designing FPGA accelerators for CNN-based object detection, pipelining is advantageous when two adjacent convolutional iterations can be performed without data dependency [7,57,58,229]. This allows the next convolutional operation to commence before the current one is fully completed, enhancing computational power.

Generally speaking, pipelining is a design technique aiming to enhance the maximum clock frequency and throughput of synchronous digital systems. This technique involves inserting registers or memory elements into the dataflow path to break down a large operation into several smaller ones [253]. In that way, each small operation requires less time, reducing the path length that a signal must traverse within a clock cycle and thereby enhancing the working frequency. Furthermore, smaller operations can be executed in parallel, leading to improved data throughput [253].

Pipelining proves particularly beneficial when processing a stream of data. In a pipelined circuit, various stages of the pipeline can handle different input stream data simultaneously within the same clock cycle. This concurrent processing enhances data processing throughput significantly [273], at the cost of more resource utilization and higher latency (cf. Section 2.4.3).

As a result, pipelining has become an indispensable operation in the design of most computational engines in the context of FPGA-based real-time object detection systems [7,27,83,90,253]. For example, Nguyen et al. [7] propose an FPGA-based streaming architecture tailored for real-time object detection showing high throughput up to 1877 GOPS, where all convolutional layers are fully pipelined (cf. Section 2.5.1).

2.4.3 Impact and Trade-offs of Acceleration Techniques in FPGA-based Object Detection

It is crucial to understand the potential effects of each technique, reviewed in Section 2.4.2, on system performance and efficiency when developing a real-time object detection system using FPGAs. This insight allows designers to select the most appropriate set of techniques, taking into account the requirements of the desired system and application, as well as the limitations and capabilities of the target FPGA device.

Metrics Methods (Section)	Throughput	Accuracy	Latency	Resource Utilization	Power Efficiency
Pruning (IV-B.1a)	↑	↓	↑	↑	↑
Quantization (IV-B.1b)	↑	↓	↑	↑	↑
Knowledge Distillation (IV-B.1c)	=	↑	=	=	=
HW-NAS (IV-B.1d)	↑	↑	↑	↑	↑
Metrics Methods (Section)	Throughput	Accuracy	Latency	Resource Utilization	Power Efficiency
Data Reuse (IV-B.2a)	↑	=	↑	↓	↑
Mathematical-based Optimization (IV-B.2b)	↑	↓	↑	↑	↑
Systolic Arrays (IV-B.2c)	↑	=	↑	↓	↑
Code Modification (IV-B.2d)	↑	=	↑	↓	↑
Roofline-based Optimization (IV-B.2e)	↑	=	=	↓	=
Pipelining (IV-B.2f)	↑	=	↓	↓	=

Figure 2.12: Possible impact of model-related (top) and implementation-related (bottom) acceleration techniques on key metrics in FPGA-based object detection systems, focusing on throughput and accuracy. The corresponding section where each technique is discussed in this work is also indicated.

Figure 2.12 highlights the possible impacts of different acceleration techniques on FPGA-

based object detection systems, focusing on key metrics including throughput, accuracy, latency, resource utilization, and power efficiency. Each technique can have a positive effect (indicated by green arrows), a negative effect (indicated by red arrows), or a neutral effect (indicated by an equal sign).

Regarding model-related techniques, pruning effectively increases throughput, reduces latency, improves resource utilization, and enhances power efficiency, although it may lead to a decrease in accuracy. Quantization shows similar benefits, improving throughput, latency, resource utilization, and power efficiency, with a potential drop in accuracy. Knowledge distillation, on the other hand, offers improvements in accuracy, making it a suitable compensator for previous techniques. Finally, Hardware-aware Neural Architecture Search (*HW-NAS*) can be employed to optimize all key metrics, leading to enhancements in throughput, accuracy, latency, resource utilization, and power efficiency. This positions HW-NAS as the most comprehensive acceleration technique for FPGA-based object detection systems, albeit with the caveat of its complex and time-consuming implementation process.

In the context of implementation-related techniques, data reuse enhances throughput and reduces latency by reusing intermediate results to avoid redundant computations. This technique tends to increase resource utilization due to higher memory requirements but improves power efficiency, with no impact on accuracy. Mathematical-based optimization methods can significantly improve throughput and power efficiency by simplifying calculations. However, these approaches might compromise accuracy because of approximations, although they can optimize resource utilization. Using systolic arrays can improve throughput and power efficiency by enabling parallel processing with a consistent data flow. While accuracy is maintained, this method may require more resources due to the dedicated hardware involved. Code modification, including techniques like loop unrolling, boosts throughput and power efficiency without affecting accuracy but may increase resource utilization and, in some cases, latency. Roofline-based optimization strikes a balance between computational load and memory access to enhance throughput, with no impact on accuracy but potentially higher resource demands and trade-offs in power efficiency. Pipelining improves throughput by allowing concurrent execution of multiple stages, thereby improving throughput. Although it does not affect accuracy, the additional hardware requirements for pipelining can lead to increased resource utilization and latency, with minimal effect on power efficiency.

It is important to note that these effects and results are not definitive and may vary with different implementations and deployments. However, the mentioned effects can be considered as the most anticipated outcomes. All in all, selecting the appropriate acceleration method involves weighing these trade-offs to achieve the desired performance while considering system constraints.

2.5 Results Analysis and Optimization

This section aims to review several case studies, illustrating how acceleration architectures and techniques discussed in Section 2.4 have been applied to develop real-time object detection systems, along with an analysis of their outcomes. Then, recent advancements and works in this field are compared across multiple metrics—such as throughput, accuracy, pixel throughput, and power efficiency—highlighting state-of-the-art results and revealing trends in FPGA-based real-

time object detection. Finally, optimization strategies specifically aimed at achieving real-time performance in FPGA-based object detection systems are briefly discussed.

2.5.1 Case Studies: Adoption of Acceleration Techniques in FPGA-Based Real-Time Object Detection

This section discusses examples of applying the aforementioned acceleration techniques in developing real-time object detection systems on FPGAs. The works are chronologically reviewed to illustrate the progress in this field. The results of adopting these techniques are presented and analyzed to demonstrate how they have contributed to enhancing the performance, accuracy, and efficiency of such systems. Additionally, an overview of the overall architecture proposed by some recent works is briefly examined to provide context on the approaches adopted in state-of-the-art object detection solutions.

Fan et al. [57] demonstrated the effectiveness of channel pruning in boosting the performance and hardware efficiency of FPGA-based real-time object detection systems. By pruning up to 10% of the channels of a customized SSD model [2], they achieved a model 3 times smaller with a minimal accuracy loss of just 1.8%. Also, by proposing a partial quantization scheme, they minimized the accuracy drop caused by quantization. This quantization scheme preserves 32-bit precision for certain components like Pre- and post-processing modules while quantizing most layers in the feature extractor to 8 bits. They also used 8-bit weights and activation in their proposed system. By employing the roofline model technique tailored for the target platform, i.e., a Xilinx Zynq ZC706 with a DRAM memory bandwidth of 1.2 GB/s, they improved overall performance by caching all intermediate results in the on-chip memory. By applying these techniques, they achieved 64.8 FPS, a latency of 15.43 ms, and a power consumption of 9.9 W using a single computation engine architecture. These results represent a more than sixfold improvement in throughput and latency, along with a 94% reduction in power consumption, compared to running the original model on a TITAN X Pascal GPU.

Nguyen et al. [7] implemented a real-time object detection system using YOLOv2 [80] on a VC707 FPGA, employing 1-bit weights and flexible 3-to-6-bit activations. This approach led to a 30-fold reduction in model size and a 5.4-fold decrease in activation size. They demonstrated that even with such low-bit quantization, an acceptable accuracy (51.68% mAP) with a high inference speed (66 FPS) is possible with a power consumption of 8.7 W. They developed an efficient streaming architecture tailored for real-time object detection, focusing on optimizing the data path, where all convolutional layers are fully pipelined. Additionally, they enhanced loop computations using loop reordering and tiling techniques, enabling the convolutional layers to run effectively on their resource-constrained FPGA.

They also introduced a data reuse scheme to minimize memory accesses and enhance computation performance by reusing both the input data and the weights (Figure 2.13). The input is processed in sliding cubes (with $K \times K \times T_i$ pixels) that move across the image in row passes. During each pass, the T_0 weight blocks are reused to perform convolutions, generating temporary output values that are stored in line buffers. These computations are performed in parallel, and the results are accumulated across multiple row passes. After each pass, the input sliding cube shifts, and new weight blocks are fetched. The convolutional outputs are accumulated in the line buffers, reducing the need for additional memory accesses. The weights are fetched once per row pass

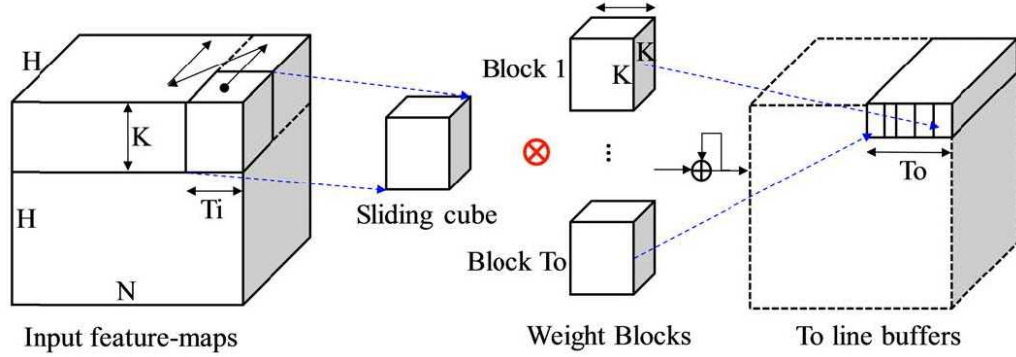


Figure 2.13: Data reuse scheme proposed in [7], minimizing memory accesses by reusing both input data and weights, with intermediate results stored in line buffers. $K \times K \times T_i$ defines the size of the input sliding cube in pixels, and T_0 denotes the number of weight blocks concurrently convolved with the input sliding cubes in each iteration. (Figure adapted from [7]).

(H times in total) and from on-chip SRAM, resulting in fewer memory accesses. Once all rows are processed, the final accumulated output is forwarded to the next layer, minimizing external memory access. The parameters T_i and T_0 are chosen to balance hardware cost and performance, ensuring efficient memory usage and parallel computation.

The results demonstrate approximately a 24% improvement in throughput and a 90% reduction in power consumption compared to running the same model with 32-bit floating-point precision on a GTX Titan X GPU.

Hao et al. [190] presented a co-search method that enhanced the performance of CNN-based object detection, leading to gains in accuracy, frame rate, and power efficiency. The proposed approach was tested on an object detection task from the DAC-SDC competition [274], where it outperformed the first-place winner using a PYNQ-Z1 board. This implementation could meet real-time requirements, delivering up to 29.7 FPS and achieving 68% accuracy. Their solution delivered a 6.2% improvement in IoU, reduced power usage by 40%, and was 2.5 times more power efficient compared to the state-of-the-art FPGA designs. Also, in comparison to GPU-based designs on Nvidia Jetson TX2, their method achieved nearly identical accuracy (approximately 69% IoU) while offering a power efficiency advantage of 3.1 to 3.8 times.

Wang et al. [82] also applied an unstructured pruning to YOLOv2 [80] and showed the effectiveness of this technique in implementing real-time object detection on FPGAs. Their approach leveraged hardware-aware optimizations, including model pruning and quantization, to reduce the computational workload and memory footprint of the YOLOv2 network. Using the roofline model as an optimization flow guidance to adopt a fine-grained unstructured pruning scheme, they achieved a 7x reduction in computational workload with only a 2.35% loss in accuracy. This technique, along with applying 8-bit quantization on the whole model and adopting a layer fusion approach, allowed for high-throughput object detection at 61.9 FPS using a single computation engine architecture. They also achieved a high power efficiency of 81.92 GOPS/W on an Arria-10 FPGA, demonstrating the viability of these techniques for real-time object detection.

Chang et al. [8] showed that adopting a mixed-precision quantization approach is a more ef-

fective way to minimize the accuracy loss than applying fixed bit-width quantization. By quantizing weights and activations to 5-8 bits and 7-8 bits, respectively, they developed a real-time object detection system based on YOLOv3 [88] using the ZCU102 FPGA platform. They achieved a throughput of 22 FPS with 49.7% mAP, reflecting only a 1.5% accuracy loss compared to the 32-bit floating-point model.

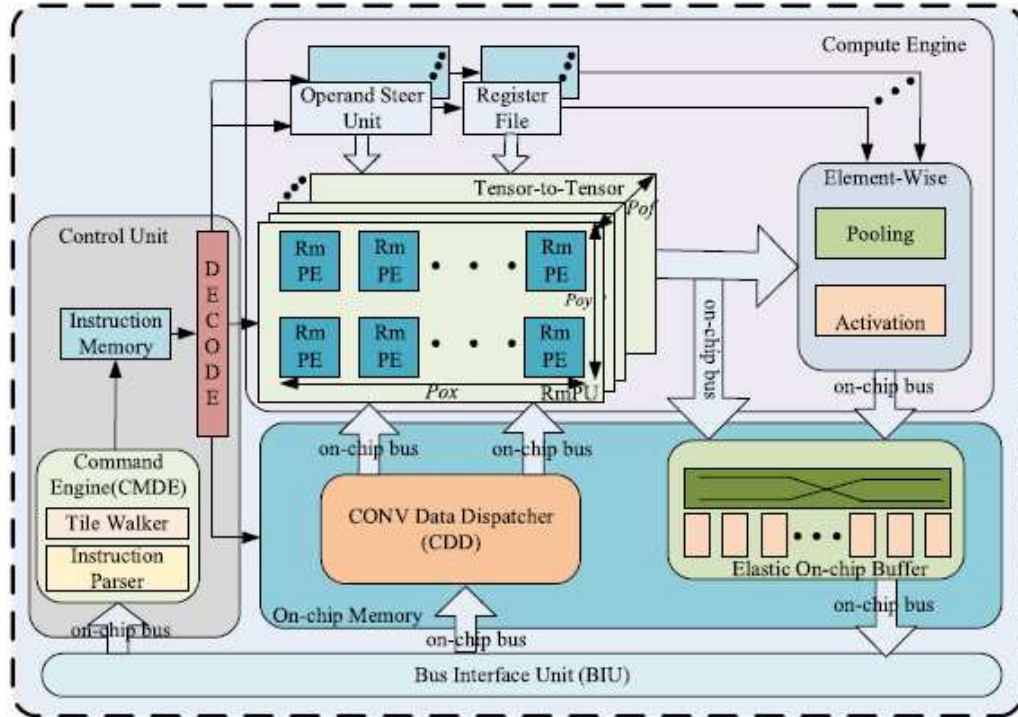


Figure 2.14: An architecture for accelerating YOLOv3 on FPGA. It is based on a pipelined multicore processing architecture consisting of a Control Unit to manage all operations, a Compute Engine with multiple layers of 2D arrays of PEs, and an on-chip memory unit optimized for parallel data access (Figure adapted from [8], .

As Figure 2.14 shows, the proposed architecture features a pipelined multicore unit composed of a computing engine, an elastic on-chip buffer, and a control unit.

The computing engine performs operations like convolution, activation, and pooling with mixed precision. This unit consists of Reconfigurable Microprocessing Elements (*RmPEs*) organized in a 3D array to adapt to various convolutional layer requirements. Each *RmPE* performs parallel MAC operations, using DSP and LUT-based adders. Multiple *RmPEs* form Reconfigurable Macro Processing Units (*RmPUs*), which combine into a 3D array for enhanced parallelism. The Control Unit manages configurations, including bit-widths, activation (e.g., ReLU), pooling types, and dataflow patterns, optimizing resource utilization. The unit supports flexible data processing by adapting to different CNN models through reconfigurable parameters for input/output widths and convolution array setup.

The on-chip buffer unit is a flexible memory system dynamically partitioned into physical banks, optimized for parallel data access, and configured by the Control Unit to match computational needs. It stores intermediate data for convolution, activation, and pooling tasks, reducing off-chip memory dependency and facilitating cross-layer data reuse.

The Control Unit orchestrates the pipelined multicore unit by breaking down convolution tasks into smaller, tile-based subtasks and coordinating the computing engine and on-chip buffer accordingly. It decodes instructions, queues them by type (control, memory access, and computation), and manages synchronized, parallel issuance, ensuring data dependencies are met for efficient execution.

Cai et al. [78] proposed an FPGA-based approach for real-time underwater object detection based on a single computation engine architecture, utilizing the Winograd algorithm to optimize various convolution operations within the network. They also improved data reuse through a ping-pong-based memory access approach. By adopting these techniques, they were able to achieve 33.14 FPS with MobileNetV3-SSDLite [2] when the accelerator was deployed on a Zynq XC7Z045 device running at 150 MHz. Compared to a CPU (Intel i7-8700), the proposed accelerator offers an 8.7× speedup, enhancing power efficiency by 60×.

Huang et al. [191] showed that employing an effective data reuse strategy in FPGA implementations can greatly enhance the performance of real-time object detection by optimizing memory access patterns and increasing data locality. Their approach employed buffering techniques to reduce the frequency of memory fetches, lowering latency and boosting throughput. This strategy proved particularly beneficial for deformable convolutions, where irregular memory access patterns typically limit data reuse. They modify the deformable convolution operation to enhance performance further by restricting adaptive offsets to fixed ranges, enabling more efficient data reuse. Additionally, they replaced the full 3×3 deformable convolutions with 3×3 depthwise deformable convolutions and 1×1 convolutions, akin to the depthwise separable convolution approach used in Xception [275], thereby streamlining computation and enhancing the efficiency of the FPGA implementation. Using their developed object detection model on an Ultra96 board, they achieved up to 32.2 FPS on the Pascal VOC dataset [50] while consuming only 5.6 W of power.

In their proposed workflow for implementing YOLOv4 on a Zynq UltraScale+ MPSoC FPGA, Liew et al. [9] demonstrated how leveraging knowledge distillation can help recover some of the accuracy lost due to pruning and quantization. As shown in Figure 2.15, the size of the baseline model (with 32-bit floating point precision) was reduced by 56% and 88% after applying channel-wise model pruning and then 8-bit quantization (INT8), respectively. This reduction in model size was crucial for achieving real-time performance at 33.34 FPS, though it came with an 8% decrease in accuracy. However, adopting knowledge distillation led to a recovery of some accuracy by about 2.5%.

Jain et al. [10] developed a real-time object detection system based on tiny-YOLOv2 [80], in which a mixed-precision quantization scheme (8-bit weight and 16-bit activation) was adopted. By taking advantage of QAT, they could not only achieve a satisfactory processing throughput (23 FPS) for real-time performance on the XC7Z035 FPGA but also recover the accuracy loss caused by applying quantization to reach a high accuracy of 57.1%. To find the best spatial and temporal unrolling factor for designing an power-efficient and high-speed architecture, they did a design space exploration using the ZigZag tool, an optimization framework [276].

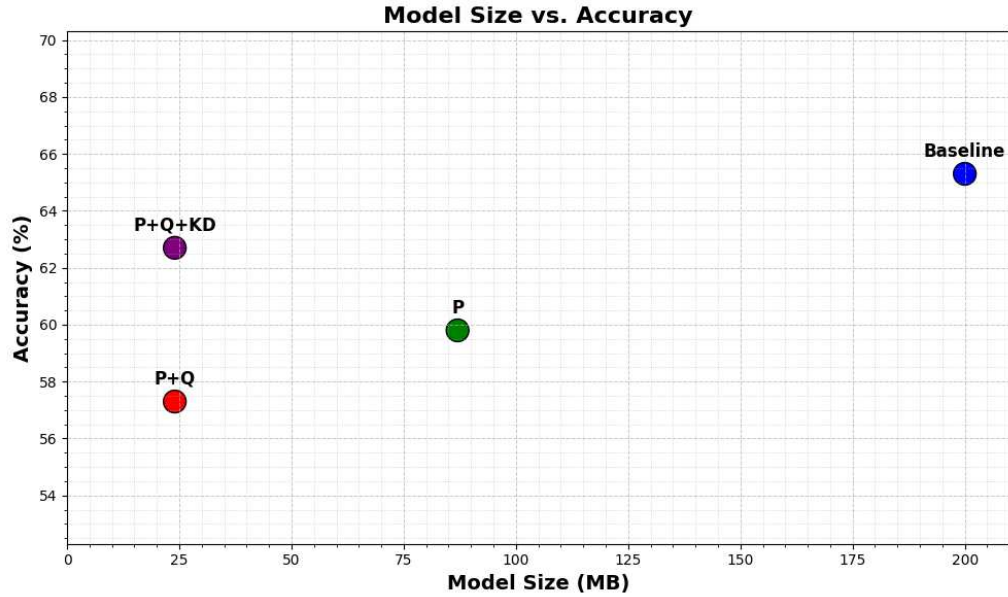


Figure 2.15: The effect of applying pruning (p), quantization (Q), and knowledge distillation (KD) techniques one after the other on the size and accuracy of the YOLOv4 model (the numbers extracted from [9]).

Subsequently, as Figure 2.16 shows, they developed an architecture using an array consisting of 32 Single Instruction Multiple Data (SIMD) units, where each SIMD unit consists of 14 parallel PEs. The architecture also includes a 2×2 max pooling unit with strides of 1 and 2, along with an activation unit that performs ReLU and Leaky ReLU. To minimize latency, dual-port memories are used to enable parallel read and write operations. A control unit is also implemented to manage data movement and oversee the operation of the SIMD arrays and other components.

Suh et al. [11] utilized a low-precision quantization approach to implement the SSD model, trained on a custom drone dataset on the Zynq ZU3EG FPGA. They introduced a uniform quantization scheme with power-of-two (POT) quantization boundaries, termed UniPOT, designed to streamline the model, reduce its size, and eliminate the need for multipliers by substituting them with shift registers. Their results showed that applying UniPOT quantization with 8-bit precision resulted in only a minimal accuracy drop (0.24% mAP) compared to the baseline 32-bit model. Their proposed design achieved 88.42% mAP, an power efficiency of 79 GOPS/W, and 158 GOPS throughput.

Figure 2.17 shows the overall proposed architecture. The design utilizes two separate DMA modules to handle read and write operations independently, each with its own descriptor buffer. The read DMA module loads image tiles and weights into designated input and weight buffers. A data router then rearranges the pixel and weight data for optimized reuse in MAC arrays, leveraging FIFOs to minimize buffer reads by shifting and reusing data in registers. This setup enables efficient pixel reuse within the register arrays, reducing memory access. Each Processing Element (PE) in the MAC array performs one multiply-accumulate (MAC) operation per cycle, benefiting

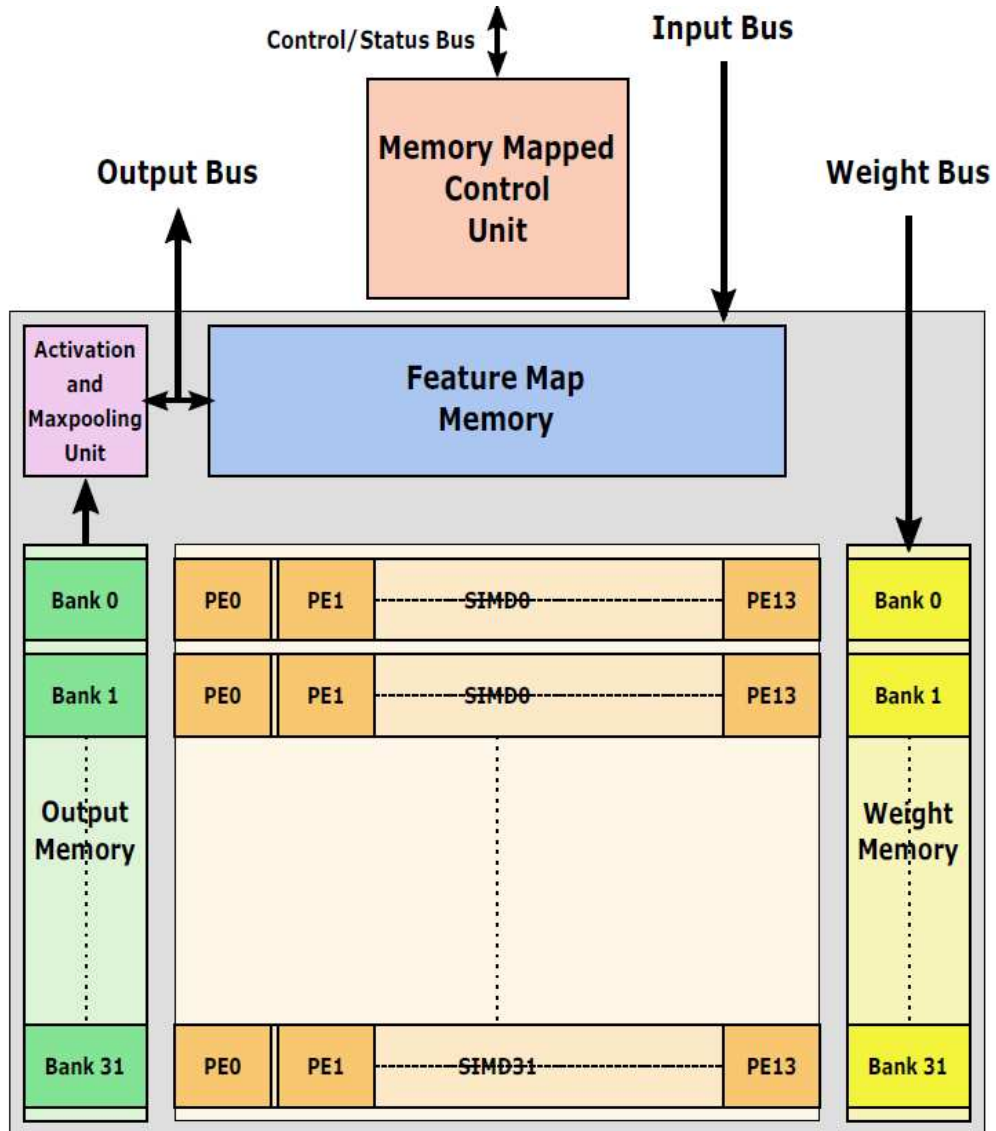


Figure 2.16: Overall accelerator architecture, consisting of a SIMD array, activation, pooling, and control units along with memories to store input, weight, and output values (Figure adapted from [10]).

from loop unrolling and tiling techniques to maximize efficiency.

Anupreetham et al. [77] developed a high-performance, real-time object detection system on a Stratix 10 FPGA, achieving low latency and high throughput by introducing a novel pipelined Non-Maximum Suppression (NMS) algorithm. This new NMS design reduced latency by removing the traditional dependency on sequential processing, which typically slows down the object detection pipeline. By allowing continuous dataflow between stages, the pipelined NMS enabled

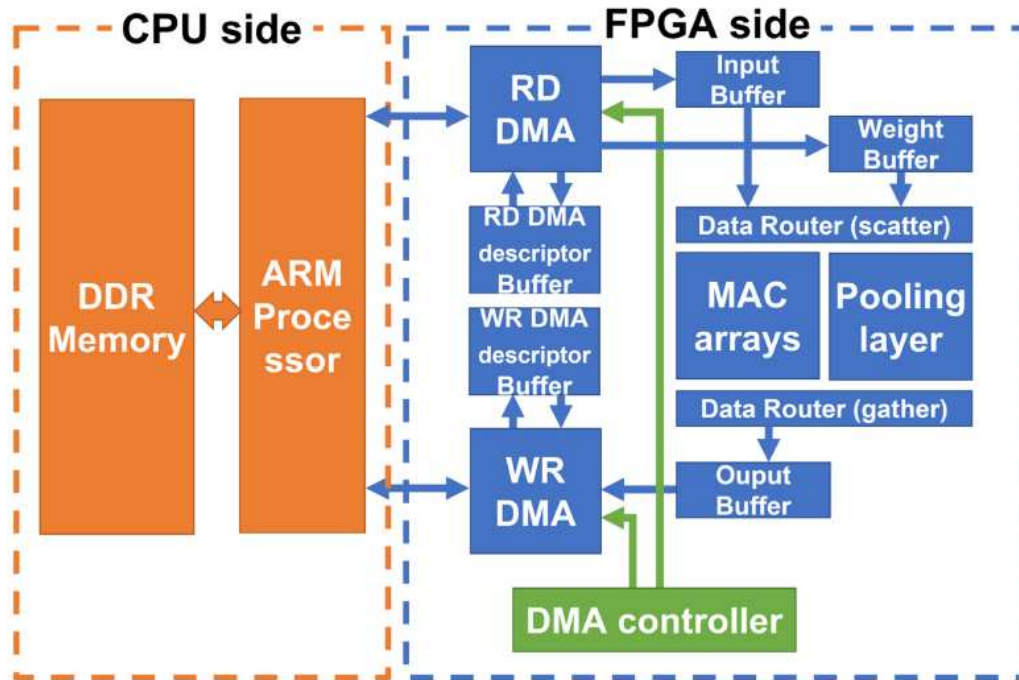


Figure 2.17: Overall hardware block diagram of a proposed FPGA accelerator, with two separate DMA modules, each equipped with a dedicated DMA descriptor buffer. The data router will rearrange the pixels and weights to maximize data reuse capability (Figure adapted from in [11]).

concurrent processing of feature extraction and SSD [2] detection stages, significantly decreasing end-to-end latency. Additionally, they enhanced the system's efficiency by incorporating a multi-threaded NMS module that processes multiple images in parallel, aligning with the throughput capacity of the CNN backbone accelerator and thus preventing stalls.

As a result, the system achieved an end-to-end latency of only 2.13 milliseconds, marking a substantial improvement over prior FPGA implementations. It delivered 5.3× higher throughput than previous solutions with comparable accuracy, i.e., 22.8% mAP on MS COCO [51], demonstrating a significant advancement in balancing speed, efficiency, and accuracy for FPGA-based real-time object detection.

D. Zhang et al. [93] introduced a fully hardware-accelerated end-to-end object detection system, designed entirely to run on FPGA devices. To meet real-time requirements and minimize system latency, all components, including post-processing algorithms, were implemented on hardware, eliminating the need for a CPU.

To enhance design flexibility and support various object detection models, they adopted a single computation engine architecture. However, to address the primary challenge of this approach—intensive memory access—they employed three dedicated data access units alongside on-chip buffers to improve data reuse. These data access units leveraged burst transfer mechanisms to optimize memory bandwidth utilization and minimize data transfer overhead. Furthermore, the proposed data buffering strategy allowed data transfers to occur in parallel with

computations, significantly boosting execution efficiency. A ping-pong buffer mechanism further enhanced data reuse within the design.

The system was evaluated using object detection models such as YOLOv2 and YOLOv3, demonstrating substantial performance improvements. Compared to state-of-the-art implementations, such as [7, 58, 77], it achieved up to 9x higher throughput and up to 5x lower latency.

When tested on a KC705 board, featuring a low-end XC7K325T FPGA, the system delivered an impressive throughput of 401 GOPS, consumed 12 W of power for the entire system, and achieved a frame rate of 65 FPS while running a quantized (8-bit) YOLOv3-tiny model.

2.5.2 COMPARATIVE ANALYSIS OF FPGA IMPLEMENTATIONS

Tables 2.10 and 2.11 present detailed information on recent studies, sorted based on the achieved pixel throughput, regarding the implementations of real-time object detection models on FPGAs. As shown in Table 2.10, the majority of these studies utilize two main detection model families: YOLO [1] and SSD [2] (cf. Section 2.2.1). These one-stage detectors are favored due to their balanced accuracy-speed performance, which is crucial for developing real-time object detection systems. Additionally, there are many lightweight variants of these models suitable for resource-constrained devices like FPGAs.

As a common technique employed in this context, different quantization regimes have been adopted based on the available resources on the target FPGA device and the accuracy required by the application of interest. The best-achieved accuracy, processing speed (in FPS), and latency, as the most relevant metrics for evaluating real-time object detection systems, are provided for each study. It is important to note that the image resolution can affect system performance—the larger the image size, the higher the achievable accuracy with lower speed. Therefore, looking at the achieved pixel throughput may help make a more meaningful and fair quantitative comparison between the existing works. For this reason, the reviewed works in Table 2.10 and 2.11 are sorted based on this metric. Additionally, Figure 2.18 shows a comparison of the reviewed works based on their pixel throughput performance. Another critical metric is the achieved computational power, or throughput, measured in Giga Operations Per Second (*GOPS*). While FPS reflects the processing throughput from a visual perspective, GOPS represents it from an architectural standpoint. Knowing the working Frequency (*Freq*) is also essential for a meaningful GOPS comparison. In addition, power consumption and power efficiency are critical factors, especially when developing systems for edge devices. Figure 2.19 shows the achieved power efficiency of the reviewed works (if reported in the corresponding manuscript) to provide a visual comparison and facilitate analysis.

Table 2.11 specifically highlights the deployed FPGA devices and the main development tools used in each work. All studies here have relied on either AMD (formerly Xilinx) or Intel (Altera) FPGA devices. This table also provides a report on resource utilization and mentions the primary acceleration techniques adopted in each study.

2.5.3 OPTIMIZATION STRATEGIES FOR REAL-TIME PERFORMANCE

To develop a real-time object detection system on resource-constrained devices such as FPGAs, it is essential to focus on key metrics like throughput, latency, and accuracy. A detailed review of

Table 2.10: The recent advanced works of FPGA-based real-time object detection implementations (part 1/2).

Reference	Year	Detection Model	Precision (W.A)	Freq. (MHz)	Accuracy (mAP)(%)	Image Resolution (Pixels)	Frame Rate (FPS)	Pixel Throughput (MPPS) ^{**}	Latency (ms)	Throughput (GOPs)	Power (W)	Power eff. ^{***} (GOPs/W)
Huang et al. [191]	2021	ShuffleNetV2 based	(4.8)	250	55.1	256×256	26.9	1.76	40	128	5.6	N/A
Cai et al. [78]	2021	SSDLite-Mobile-NetV3	(16,N/A)	150	56.8	300×300	33.14	2.98	N/A	N/A	9.34	N/A
Sub et al. [11]	2023	SSD-VGG16	(8.8)	200	76.2	300×300	34.18	3.07	29.25	138	2.4	57.5
Fan et al. [57]	2018	SSDLite-Mobile-NetV2	(8.8)	100	20.5	224×224	64.8	3.25	15.43	N/A	9.9	N/A
Jain et al. [10]	2022	Tiny YOLOv2	(8,16)	200	57.1	416×416	23	3.98	43.33	161.4	N/A	N/A
K. Kim et al. [189]	2023	SqueezeNet based	(8.8)	100	82.6	192×256	43.95	4.14	22.75	N/A	2.11	N/A
S. Kim et al. [56]	2021	SSDLite	(16,16)	200	78.6	224×224*	84.4	4.23	11.79	N/A	5.1	N/A
Li et al. [83]	2020	YOLOv2	(8,N/A)	200	73.6	416×416	25	4.32	N/A	740	27.2	27.2
M. Kim et al. [91]	2023	Tiny YOLOv3	(8.8)	100	81.19	320×320	76.75	4.5	13.03	95.08	2.2	43.16
Babu et al. [97]	2022	YOLOv4	N/A	100	N/A	416×416*	30	5.19	11.83	189.14	10.36	18.26
Chang et al. [8]	2021	YOLOv3	(2.8-2.8)	200	49.7	640×480	22	6.75	N/A	809	15.3	52.87
D. Zhang et al. [93]	2024	Tiny YOLOv3	(8.8)	200	58.1	416×416	65.7	11.36	N15.2	401	7.8	51.4
Nguyen et al. [7]	2019	Tiny YOLOv2	(1.6)	200	51.38	416×416	66.56	11.51	N/A	464.7	8.7	53.29
Hosseiny et al. [99]	2023	Tiny YOLOv7	(13,13)	100	34.1	640×640*	30	12.28	15	N/A	10.79	N/A
Xu et al. [85]	2021	Tiny YOLOv2	(8.8)	190	56	416×416	71	12.28	N/A	500	N/A	N/A
Wang et al. [82]	2020	YOLOv2	N/A	211	74.45	416×416	72.5	12.54	N/A	2,130	26	81.92
Pestana et al. [92]	2021	Tiny YOLOv3	(16,16)	143	31	768×576	32.4	14.33	30.9	180	3.87	46.51
J. Zhang et al. [58]	2021	Tiny YOLOv2	(8.8)	200	77.04	1280×384	43.7	21.47	27.78	464.5	10.25	45.3
Anupreetham et al. [77]	2024	SSDLite-Mobile-NetV1	N/A	400	22.8	640×640*	469	192	2.13	N/A	81	N/A

* It is not directly mentioned through the work; therefore, the image size is assumed based on the employed model or dataset.

** Pixel throughput, measured in Mega Pixels Per Second (MPPS), is calculated by multiplying the frame rate by the number of pixels per frame. (The table is sorted based on this parameter)

*** The term "energy efficiency" is also used in studies within this context to refer to GOPs/W.

Table 2.11: The recent advanced works of FPGA-based real-time object detection implementations (part 2/2).

Reference	FPGA/Board Model	Design Environment	Resource Utilization			Key Employed Techniques	
			BRAMs # (size)	LUTs/ALMs (K)	DSP FFs (K)		
Huang et al. [191]	Xilinx Ultra96	Vivado HLS ^a	216 (18K)	34	360	42	Data reusing, algorithm-hardware co-design
Cai et al. [78]	Xilinx XC7Z045	Vivado	446 (18K)	166	801	159	Wingrad-based, and memory access (using ping-pong operation) optimization
Suh et al. [11]	Xilinx ZU3EG	Vivado	102 (18K)	78	263	75	Low-precision quantization, loop tiling, model pruning
Fan et al. [57]	Xilinx ZC706	Vivado	311 (18K)	148	728	191	Partial quantization, resource sharing, channel pruning
Jain et al. [10]	Xilinx XC7Z035	N/A	462 (18K)	88	448	62	Design space exploration, quantization aware training, tiling strategy
K. Kim et al. [189]	Xilinx ZC702	Vivado	N/A	18	7	21	Model simplification and compression, parallel processing
S. Kim et al. [56]	Intel Arria 10	Quartus Prime	N/A	N/A	980	N/A	Arithmetic-based and memory access optimization, data reusing
Li et al. [83]	Intel Arria 10	OpenCL + RTL	N/A	N/A	410	N/A	Mixed-precision data, parallel design, pipelining
M. Kim et al. [91]	Xilinx Artix-7	Vivado	185 (18K)	50	240	58	Dynamic data reuse scheme, hardware-friendly quantization
Babu et al. [97]	Xilinx XC7Z020	Vivado HLS	115 (18K)	23	174	46	Code modification, memory access optimization
Chang et al. [8]	Xilinx ZCU102	Vivado	680 (18K)	185	2280	308	Mixed-precision quantization
D. Zhang et al. [93]	Xilinx XC7K325T	Vivado HLS	553 (18K)	136.5	687	N/A	Ping-pong buffer-based data reusing, memory access optimization
Nguyen et al. [7]	Xilinx VC707	Vivado HLS	1,026 (18K)	86	168	60	Binary weights, HW-aware quantization, data reusing
Hosseiny et al. [99]	Xilinx XC7Z100	Vivado HLS	1,092 (18K)	59	16	9	Code modification, pipelining, data reusing
Xu et al. [85]	Intel Arria-10 GX1150	OpenCL SDK	1849 (20K)	145	1092	N/A	Scalable pipeline design, layer fusion, SW-HW co-design
Wang et al. [82]	Intel Arria-10 GX1150	OpenCL	2,576 (20K)	188	76	N/A	Fine-grained pruning, Hardware/Software co-Design
Pestana et al. [92]	Xilinx XCKU040	N/A	384 (18K)	139	839	N/A	Post-training dynamic quantization
J. Zhang et al. [58]	Xilinx ZC706	N/A	256 (18K)	84	610	65	Resource sharing, pipelining, arithmetic-based optimization
Anupreetham et al. [77]	Intel Stratix 10 GX2800	Quartus Prime Pro	7,883 (20K)	N/A	5,009	N/A	Pipelining and parallel processing, code modification

^a Vivado offers a complete design environment for developing and implementing FPGA designs, while Vivado HLS facilitates high-level synthesis, converting C/C++ code into RTL code.

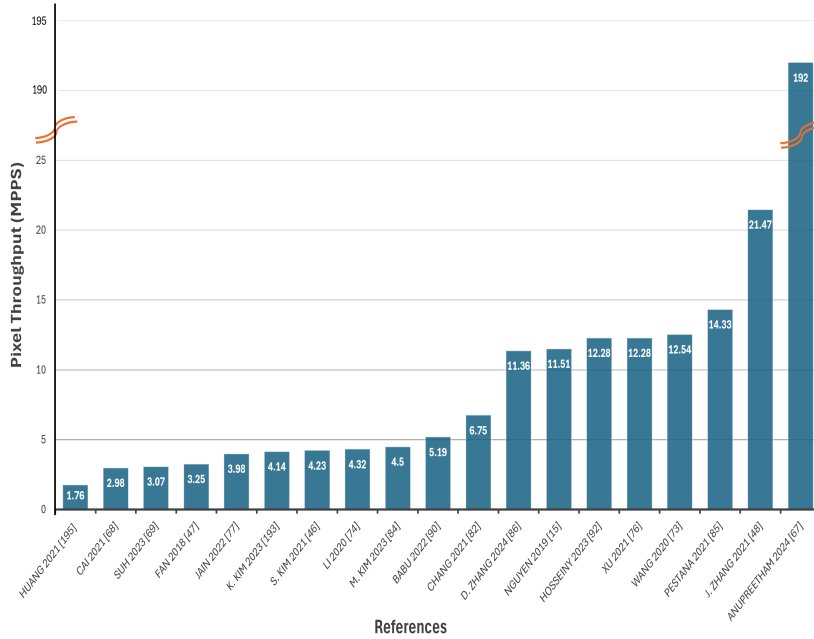


Figure 2.18: Comparison of the reviewed works based on the achieved pixel throughput.

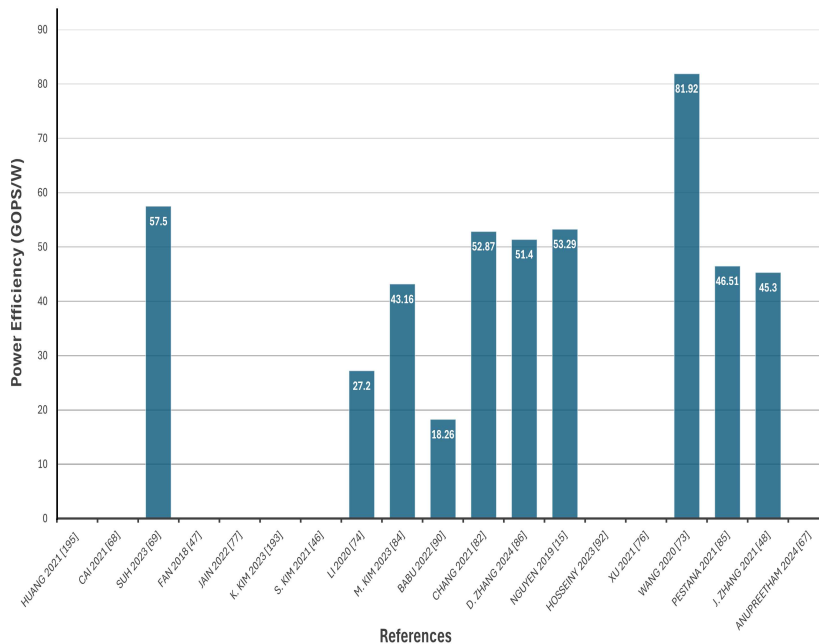


Figure 2.19: Comparison of the reviewed works based on their reported power efficiency (note: for some works power efficiency data is not available).

these metrics within the context of this paper is provided in Section 2.2.3.

As outlined in Section 2.2.1, deep neural network (DNN)-based models often provide the highest accuracy for object detection tasks. However, these models are computationally intensive, necessitating a careful balance between accuracy and execution time, as discussed by Castells et al. [61]. This balance becomes even more critical in hard real-time systems, where meeting stringent timing requirements is essential for safety and reliability.

Before adopting optimization techniques, discussed in Section 2.4.2, it is crucial to select the appropriate hardware architecture. Section 2.3.3 introduces two FPGA architectures popular for developing object detection systems: single computation engine and streaming architectures.

Among these, streaming architectures are more suitable for real-time object detection due to their reduced memory accesses, which enhance system speed and lower latency. These architectures are also optimized to perform specific tasks within each block, leveraging parallelism to boost performance in object detection models [3]. However, this level of optimization typically comes with higher hardware costs compared to one-size-fits-all configurable engines.

The rapid development of deeper and more powerful object detection models has shifted the focus from achieving high accuracy to optimizing the performance of these models on resource-constrained devices like FPGAs. Modern object detection models tend to be computationally demanding and require larger memory footprints, making performance optimization a primary challenge.

To address these challenges, the typical workflow for developing a real-time object detection system on FPGAs involves simplifying and optimizing the chosen model. In this regard, many works show the advantages of employing Binarized Neural Networks (*BNNs*) [277] and Ternary Neural Networks (*TNNs*) [278] to achieve real-time inference [7, 279–282]. BNNs [283] and TNNs [278] are two types of neural networks in which the weights and activations are quantized into two $(-1, 1)$ and three values $(-1, 0, 1)$, respectively.

However, compression and simplification techniques, such as pruning and quantization, can result in accuracy degradation. To mitigate this problem, retraining the compressed models and utilizing knowledge distillation [220] are effective strategies. Mishra et al. [284] propose a method, called *Apprentice*, which combines low-precision arithmetic with knowledge distillation to achieve nearly the same accuracy as the original model. Remarkably, this approach results in less than a 1% accuracy loss compared to a 32-bit floating-point model, even when employing ternary precision within the ResNet architecture [128] on the ImageNet dataset [153].

In addition, a real-world object detection system comprises diverse components, some of which are best suited for execution on a processor due to their sequential nature. Considering this, adopting a hardware-software co-design approach can facilitate the development of an optimal system [77, 82, 285, 286]. This approach may become more effective when combining hardware-aware NAS, discussed in section 2.4.2. Adopting a co-design approach can enable developers to design efficient object detection systems that satisfy real-time constraints, including latency, throughput, and accuracy, using FPGAs and FPGA-oriented system-on-chip devices [190, 226].

2.6 Challenges and Future Directions

2.6.1 Existing Challenges

Seeking to enhance their accuracy, object detection models are becoming increasingly larger and deeper with more complex architectures. Consequently, implementing such models on FPGAs is challenging due to their limited memory and computational resources. Furthermore, achieving real-time performance with these complex models is another significant issue. As discussed, model simplification and compression techniques can be employed to address this problem, albeit at the expense of a drop in accuracy. Although the knowledge distillation technique can mitigate the accuracy loss issue [220], efficiently applying this method is not an easy task [5]. Therefore, the first challenge is to make an object detection model hardware-friendly while maintaining acceptable model accuracy and speed to meet real-time constraints.

Developing FPGA-based object detection systems is a complex and multifaceted procedure. It begins with developing a suitable model or modifying an existing one based on the project requirements, leading up to the final hardware implementation. Consequently, this process requires knowledge from various areas (in both software and hardware) and the ability to work with multiple tools to perform each part of the development flow. For example, a designer may need to work with an AI framework such as PyTorch (working with Python), a high-level synthesis tool (working with C/C++), an FPGA design tool like Vivado (working with HDL), and some other software tools and programming languages depending on the project. The important point is that in each development phase, having comprehensive knowledge about the subsequent steps is needed to design an efficient system.

Many efforts have been made to provide appropriate tools to integrate this development process and reduce the need for in-depth knowledge in all parts, especially concerning hardware implementation [193, 194]. However, these tools are typically designed to work with some pre-defined object detection models for specific FPGA devices.

Consequently, the lengthy and challenging design process remains one of the main challenges in this field, directly impacting the time-to-market and, subsequently, the final price of the system.

2.6.2 Possible Future Research Trends

Naturally, addressing the existing challenges mentioned earlier remains a primary research focus. As these challenges evolve, the development of hardware-friendly models tailored to FPGA architectures, accounting for their computational and memory constraints, will likely remain a key direction. Such models should aim to maximize the inherent parallelism of FPGAs while minimizing resource utilization, leading to better trade-offs between speed, accuracy, and power efficiency. To this end, deeper exploration of co-design techniques, such as Hardware-aware Neural Architecture Search (HW-NAS) [224], can play a pivotal role. These approaches could pave the way for automatic model optimization, potentially enabling the deployment of more sophisticated detection models with less manual tuning.

A promising area is the creation of software tools that can seamlessly integrate the design, optimization, and deployment processes for FPGA-based object detection systems. Although tool development tends to be more commercially driven, incorporating open-source initiatives and academic collaborations could stimulate further innovation. Additionally, the emergence of

flexible hardware libraries and parameterizable soft cores capable of executing diverse object detection algorithms on a wide range of FPGA devices could standardize and accelerate the development cycle. Such efforts could emphasize modularity, enabling rapid adaptation to new FPGA platforms and model architectures without extensive re-engineering.

The trend toward using distributed FPGA clusters for real-time object detection is another area ripe for investigation. Distributed configurations could unlock higher performance and enable the design of fault-tolerant and scalable systems suitable for hard real-time requirements [253]. Future studies could explore novel frameworks for inter-node communication, computational task partitioning, and efficient scheduling strategies. As object detection models continue to grow in size and complexity, research on scalable data partitioning schemes, latency minimization, and load balancing across multiple FPGA nodes becomes increasingly vital. While some initial studies [287–289] have explored these areas, the field could benefit from more refined techniques to address the demands of emerging high-performance detection models.

The adoption of transformers in computer vision, particularly for object detection, has gained significant traction, with transformer-based models showing performance that rivals or surpasses that of conventional convolutional neural networks (CNNs) [100, 117]. However, implementing these models on FPGA platforms presents unique challenges due to their complex architecture, extensive parameter sets, and substantial computational and storage requirements [101]. Thus, more in-depth research is required to bridge the gap between transformer-based models' capabilities and the hardware constraints of FPGA devices. Future studies could explore the following directions to make FPGA-based transformer implementations feasible and efficient:

- **Hardware-Friendly Transformer Architectures:** Development of new transformer-based object detection models optimized specifically for FPGA deployment, focusing on reducing model complexity while maintaining accuracy.
- **Innovative Hardware Accelerator Designs:** Creation of specialized FPGA accelerators that match the unique structural characteristics of transformers, such as self-attention operations. These designs should emphasize parallelism and custom data paths to minimize latency and optimize resource usage. Furthermore, leveraging reconfigurable hardware properties to allocate resources for different transformer layers based on workload adaptively could improve overall system efficiency.
- **Enhanced Dataflow Techniques:** Efficient dataflow management can significantly reduce memory access and power consumption. Exploring new memory hierarchy designs, data compression techniques, and efficient caching strategies can enhance parameter sharing and data reuse across the transformer's layers. Furthermore, research could focus on dynamic memory allocation schemes tailored to transformer models to optimize on-chip memory utilization.

Ultimately, the future of FPGA-based real-time object detection will likely involve a combination of approaches, from the development of optimized hardware accelerators and software tools to advancements in model co-design. As object detection models become increasingly sophisticated, the field must evolve to address the growing demands for higher accuracy, lower latency, and better power efficiency. Researchers will need to push the boundaries of both algorithmic and

hardware architecture innovations to meet these challenges, setting the stage for breakthroughs in real-time object detection on FPGAs.

2.7 Final Remarks

Real-time object detection plays a crucial role in applications such as autonomous vehicles and robotics, demanding low latency, high throughput, and application-specific accuracy. Achieving these requirements necessitates efficient system architectures. FPGAs have emerged as a compelling choice for object detection due to their inherent parallelism, deterministic low latency, high throughput, and ability to process images directly from external sources.

This paper presents a comprehensive review of FPGA-based real-time object detection, covering key implementation and optimization strategies. We discuss the distinction between soft and hard real-time systems, advancements in object detection algorithms—particularly one-stage CNN architectures for their balance of speed and accuracy—and widely used evaluation metrics and datasets. Additionally, we analyze existing FPGA implementations, comparing their performance using pixel throughput as a fair metric.

Despite notable advancements, several challenges remain, including integrating increasingly complex models within FPGA constraints, optimizing resource utilization, and scaling designs for higher-resolution video streams. Overcoming these hurdles requires innovative hardware-software co-design approaches, enhanced toolchains, and hybrid architectures that combine FPGAs with other accelerators.

Moreover, FPGA-based real-time object detection faces challenges such as managing hardware resource limitations while maintaining real-time performance, navigating the intricate development process spanning hardware and software, and addressing the lack of flexible, integrated design tools. Techniques like model simplification and knowledge distillation help mitigate these issues but introduce accuracy trade-offs. Addressing these concerns necessitates advancements in hardware-aware model optimization, streamlined development frameworks, and modular FPGA architectures to enhance scalability and efficiency.

Future research directions may focus on optimizing hardware-friendly models, refining co-design methodologies like HW-NAS, and developing integrated software frameworks for efficient design. Exploring distributed FPGA clusters, modular hardware libraries, and specialized accelerators for transformer-based models could further enhance scalability, accuracy, and performance, enabling more efficient real-time systems across diverse applications.



Chapter 3

Real-time Object Detection on FPGA-based Heterogeneous MPSoCs: A Preliminary Analysis of the Execution Bottlenecks

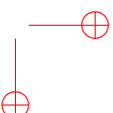
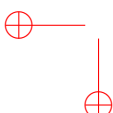
3.1 Introduction

Object Detection (*OD*) is a computer vision task that involves identifying and localizing objects within images or video sequences. In real-time OD systems, processing must be completed within a strict and predictable time frame to ensure reliability, safety, and consistency. Additionally, it is important to achieve an acceptable inference speed, frame rate, and accuracy that meets the requirements of the application of interest [290]. Nowadays, real-time OD is essential in various applications, including autonomous vehicles [14], and industrial automation [291].

To improve accuracy, object detection models are becoming progressively larger and more intricate, with deeper architectures. However, deploying such models on resource-constrained devices presents significant challenges, mainly due to limitations in memory and computational power [290]. The high computational demands and frequent memory access, required for the execution of such large, complex models, can lead to significantly increased power consumption and reduced processing speeds. Consequently, developing power-efficient, real-time OD systems remains a critical challenge.

To enhance system performance in terms of speed, throughput, and power efficiency, Hardware Accelerators (*HWAccs*) such as GPUs, FPGAs, ASICs, or other dataflow/multithreaded designs are used [292–296]. HWAccs are specialized components designed to execute specific tasks much more efficiently than general-purpose processors or CPUs. Heterogeneous Multiprocessor Systems-on-a-Chip (*MPSoCs*) have garnered significant attention in this context [297], as they integrate Programmable Logic (*PL*), or FPGA, with multiple hard processor cores on a single chip, enabling effective task offloading and optimization of system performance.

To leverage the advantages of heterogeneous MPSoCs, this work explores the development of real-time object detection using Tiny YOLOv2 [80] on an AMD Xilinx Zynq UltraScale+ platform. By utilizing both the Processing System (*PS*) and Programmable Logic (*PL*), the system achieves high performance while maintaining low power consumption. The evaluation includes a comparative analysis against CPU and GPU implementations. Additionally, profiling and memory analysis are conducted using the Vitis AI Profiler, provided by AMD Xilinx, to identify potential bottlenecks and optimization opportunities. The insights gained from this study contribute to improving real-time OD on other similar FPGA-based heterogeneous MPSoCs.



3.2 Case Study: Accelerated YOLOv2 on Zynq UltraScale+

3.2.1 Algorithm Pipeline Overview

Several effective approaches have been proposed for real-time object detection using deep learning techniques, including Single Shot MultiBox Detector (SSD) [2] and You Only Look Once (YOLO) [1]. Among them, YOLO offers one of the best trade-offs between accuracy and speed. It employs a single neural network to simultaneously predict object bounding boxes and class probabilities in a single pass.

In this work, we use Tiny YOLOv2 trained on the Pascal VOC dataset [50] as the main OD model. This model consists of nine convolutional layers, all of which are followed by Leaky ReLU activation except for the final layer, which uses a linear activation function. The model includes six max-pooling layers to progressively reduce spatial dimensions while preserving key features. The last two convolutional layers refine the extracted features, with the final 1×1 convolutional layer, mapping them to bounding box predictions. The region layer divides the image into a grid, assigns five anchor boxes per cell, and predicts object confidence scores, class probabilities, and bounding box coordinates. This streamlined architecture enables an efficient balance between speed and detection accuracy, making it well-suited for real-time applications on resource-constrained hardware [290].

3.2.2 Pipeline Design Overview

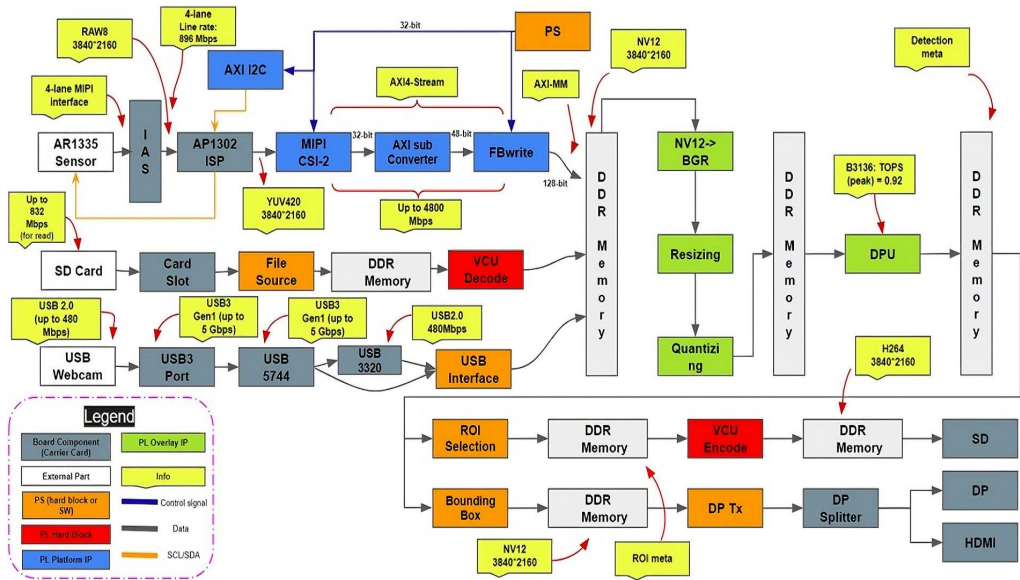


Figure 3.1: Pipeline design of the real-time object detection system, centered around the DPU inference engine, provided by Vitis AI, as the core of the accelerator design on the Kria KV260 board, featuring a Zynq UltraScale+ MPSoC.

To implement our real-time object detection system using Tiny YOLOv2, we leveraged AMD

Xilinx’s Smart Camera project [298] as the baseline design. The system operates on the cost-effective Kria KV260 board, equipped with a Zynq UltraScale+ MPSoC that integrates a quad-core ARM Cortex-A53 processor, dual ARM Cortex-R5F real-time processors, and FPGA fabric, making it well-suited for high-performance applications like real-time OD.

Figure 3.1 illustrates the pipeline design of the implemented system. The input video can be captured from different sources, including an SD card, a Mobile Industry Processor Interface (MIPI) single sensor device (MIPI camera), and a USB webcam. After applying some pre-processing functions such as image resizing and quantization, the processed video stream is fed into the Deep Learning Processor Unit (DPU), where the object detection model is executed. The inferred frames can then be stored on an SD card or displayed via a DisplayPort (DP) or High-Definition Multimedia Interface (HDMI).

3.3 Primary Implementation Results

Vivado and Vitis (both version 2021.1) are used to develop the hardware platform design on the PL side and the application software on the PS side, respectively. Vitis HLS 2021.1 is employed to implement PL accelerator kernels (excluding the DPU, shown as green units in Figure 3.1) using C/C++. The Tiny YOLOv2 model is deployed using Vitis AI 1.4.0, which facilitates model quantization from 32-bit floating-point to INT8 fixed-point format and compiles it for execution on the DPU.

Table 3.1 presents the resource utilization of the implemented hardware platform.

Table 3.1: A Summary of Resource Utilization of Implementation Results.

Resource	LUTs	BRAMs	DSPs	FFs
Available	117120	144	1248	234240
Utilization	75075	101.5	583	129141
Utilization (%)	64.10	70.49	46.71	55.13

The system delivers a throughput of approximately 30 Frames Per Second (*FPS*) when processing full-HD video streams captured from a MIPI camera or an SD card (both operating at 30 FPS).

To assess the performance achieved on the KV260 Zynq UltraScale+ MPSoC and conduct a comparative analysis, Table 3.2 presents a comparison of the obtained results with those from executing the same algorithm, Tiny YOLOv2, on an AMD Ryzen 5 5600H CPU and an NVIDIA GeForce RTX 3050 GPU.

The evaluation focuses on key metrics, including throughput (FPS), power consumption (W), and the power-delay product (PDP), with particular emphasis on relative PDP values as a measure of efficiency across different platforms.

As the table shows, the MPSoC-based system (KV260 board) stands out as the most efficient option for real-time object detection, delivering the highest throughput with the lowest power

Table 3.2: Comparative performance and efficiency of the developed OD with Tiny YOLOv2 receiving a live video stream from an attached camera @ 30 FPS as input on CPU, GPU, and KV260 Platforms, highlighting FPS, power consumption (W), and Power Delay Product (PDP) (W.Sec).

HW Platform	Processor	Quant. (#bits)	FPS	Power Consumption (W)	Power Delay Product (PDP) (W.Sec)	Efficiency (Relative PDP)	Comments
CPU	AMD Ryzen5 5600h	32	2.5	2.1	0.840	1	Single Core
CPU	AMD Ryzen5 5600h	32	8.0	5.3	0.663	1.26x	12 cores (with OpenMP)
GPU	NVIDIA GeForce RTX 3050	32	15	13.5	0.9	0.93x	12W by GPU and 1.5W by CPU are consumed.
KV260	Zynq Ultrascale+ MPSoC	8	30	1.1	0.037	22.7x	

consumption and the best power-delay efficiency. On the KV260 platform, the system is fed by 1920x1080 pixel video stream at 30 FPS using a MIPI camera, whereas other platforms receive 640x480 pixel video streams at 30 FPS via a USB camera. However, it should be mentioned that in all scenarios, the OD model receives resized video frames of 416×416 pixels. The KV260 platform achieves a throughput of 30 FPS, significantly outperforming the GPU's 15 FPS and the CPU's range of 2.5 to 8 FPS. It is important to note that when the GPU serves as the HWAcc in this system, the CPU-based video capture unit and data transfer to the GPU introduce latency, limiting overall FPS. However, on the KV260 board, the MIPI camera is directly connected to the PL side of the Zynq UltraScale+ MPSoC, leading to reduced data transfer latency and enhanced video stream acquisition performance. This capability ensures rapid video frame processing, essential for real-time object detection applications.

In terms of power consumption, the MPSoC-based system consumes just 1.1 watts, making it the most power-efficient platform among those tested. In contrast, the GPU uses 13.5 watts, while the CPU setups consume between 2.1 and 5.3 watts, depending on the number of cores utilized. Power consumption is measured using the "nvidia-smi" utility for the GPU and the "PowerTOP" tool for the CPU, both running on Ubuntu 20.04. Additionally, to further optimize CPU performance, OpenMP (Open Multi-Processing) is employed.

The MPSoC-based system also features the lowest Power Delay Product (PDP) of 0.036 W.sec, reflecting its superior energy efficiency. This metric, which combines performance and power efficiency and is calculated as $FPS^{-1} \times PowerConsumption$, underscores the KV260's advantage over the GPU (0.886 W.sec) and CPU (ranging from 0.562 to 0.852 W.sec, depending on the number of cores). Therefore, for the same power consumption, the KV260 platform delivers significantly better performance with reduced inference time compared to the other platforms.

3.4 Execution Profiling and Bandwidth Analysis

To identify performance bottlenecks and opportunities for optimization, we can profile our application using the Vitis AI Profiler, provided by AMD Xilinx. This tool collects runtime data from the DPU, ARM cores, and DDR memory controller during inference and presents the information in various formats, including text files and visual timelines, in Vitis Analyzer for further analysis. The profiler allows us to examine the DPU and PS components, including preprocessing, post-processing, and the DPU kernel during application execution. It displays the entire data pipeline on a unified timeline, highlighting areas for improvement, and captures DDR memory read/write accesses to offer a comprehensive system-wide view.

In our experiment, we use a MIPI camera as the input and a display through an HDMI port for the output. Figure 3.2 illustrates the profiling results (timeline trace), visualized in Vitis Analyzer 2021.1. It shows the application's profiling results, highlighting the interactions between the application running on the PS (host) and the DPU via Hardware Abstraction Layer (HAL) API calls¹, data transfer between the PS and PL through DDR memory, and layer-by-layer model execution on the DPU.

The profiling results indicate suboptimal parallelism, as evidenced by substantial CPU idle time while waiting for the DPU. A more efficient hardware-software co-design, leveraging improved task pipelining, could mitigate these stalls. Additionally, a hybrid control-flow/dataflow execution strategy with fine-grained task management may enhance CPU-DPU parallelism, leading to higher throughput.

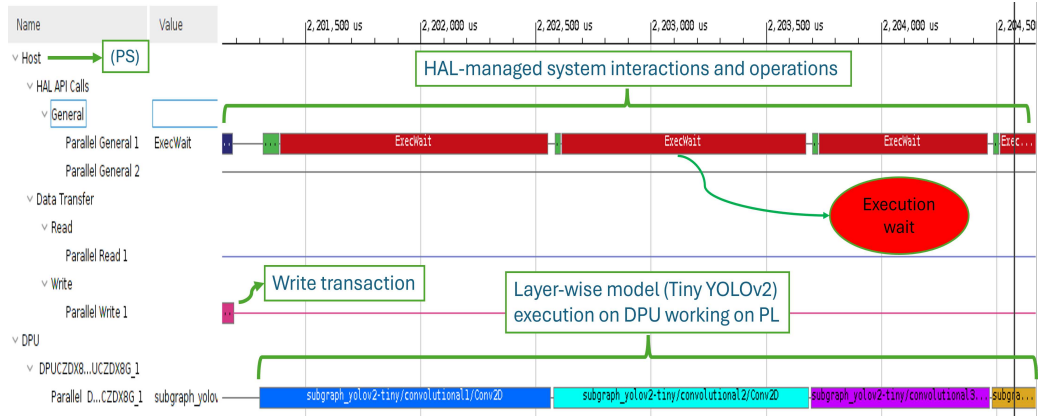


Figure 3.2: Timeline analysis of the application using Vitis AI Profiler, illustrating the interaction between the application running on the PS and the DPU, data transfer between the PS and PL through DDR memory, and the model execution on the DPU. The long "Execution wait" periods can be inferred as a performance bottleneck, caused by inefficient task scheduling.

Additionally, the memory bandwidth utilization analysis from Vitis Analyze, shown in Figure 3.3, reveals that the DDR bandwidth is not being fully utilized. Although each port connecting

¹The Hardware Abstraction Layer (HAL) API is a set of software functions that can provide a high-level interface for configuring and interacting with hardware blocks on an MPSoC, enabling applications to communicate with the hardware without dealing with low-level implementation details.

the DPU and post-processing module to DDR supports a data transfer capacity of approximately 2400 MB/s, the bandwidth evaluation in Vitis Analyzer indicates that, most of the time, only up to half of this capacity is used. This underutilization suggests a performance bottleneck in the data pipeline. To improve system performance and utilize DDR bandwidth more efficiently, strategies like optimizing data transfer, increasing parallelism between the DPU and post-processing stages, or enhancing memory access patterns with a hybrid execution approach could be considered

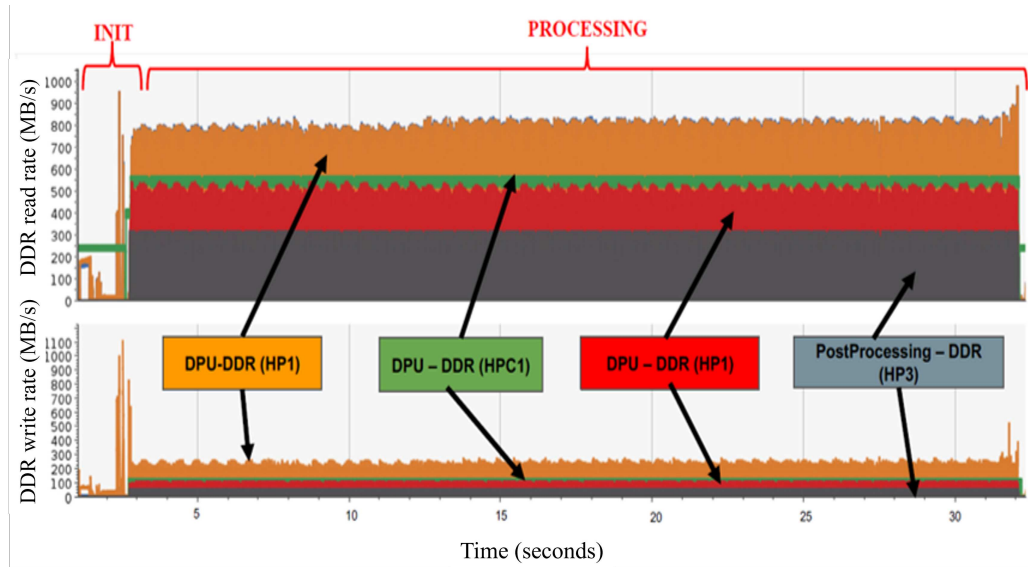


Figure 3.3: DDR bandwidth (read/write) utilization in Vitis Analyzer, where each color represents data transfer between a PL module (various DPU ports and the post-processing unit) and the DDR controller via high-performance memory ports. Considering that each port theoretically supports up to 2400 MB/s of total bandwidth (read + write), the observed memory bandwidth utilization suggests potential for optimization.

3.5 Final Remarks

This paper evaluates the efficiency of FPGA-based heterogeneous MPSoCs for real-time object detection by implementing Tiny YOLOv2 on an AMD Xilinx Zynq UltraScale+ MPSoC. The results show that the FPGA-based system achieves 30 FPS while consuming just 1.1 W of power, significantly outperforming CPU and GPU alternatives in both performance and power efficiency. A preliminary analysis using the Vitis Profiler identifies key bottlenecks, particularly in memory bandwidth utilization and task scheduling between the processing system and programmable logic. These findings suggest potential optimization through better data transfer, efficient hardware-software co-design, and hybrid control/dataflow strategies.



Chapter 4

Design and Validation of a Lightweight, Energy-Efficient FPGA Streaming Architecture for Graph Convolutional Networks (GCNs)

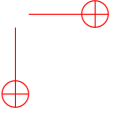
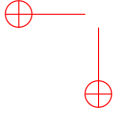
4.1 Introduction

Graph Convolutional Networks (GCNs) have emerged as a transformative paradigm for processing graph-structured data, establishing their dominance in domains ranging from social network analysis and recommendation systems to molecular property prediction and knowledge graph completion [299, 300]. The unique capability of GCNs to capture both node features and graph topology through learnable aggregation and transformation functions has positioned them as essential tools for machine learning tasks involving non-Euclidean data structures [13, 300, 301]. Recent advances in domains such as drug discovery, financial network analysis, and real-time decision-making systems have further amplified the demand for efficient GCN processing, particularly in resource-constrained environments where computational efficiency and low latency are paramount [300, 301].

Despite their mathematical elegance and widespread applicability, GCNs present significant computational challenges that distinguish them from traditional deep neural networks. The core GCN operation involves sparse matrix–matrix multiplication with irregular memory access patterns, creating bottlenecks in conventional computing architectures [302–304]. The adjacency matrix \tilde{A} exhibits extreme sparsity, often exceeding 99% in real-world graphs, leading to poor data locality, irregular workload distribution, and inefficient resource utilization on general-purpose processors [300, 302, 305]. Furthermore, the heterogeneous nature of graph data, where nodes exhibit vastly different degrees and connectivity patterns, exacerbates load imbalance issues and complicates parallelization strategies [300, 301, 306].

Field-Programmable Gate Arrays (FPGAs) have emerged as a compelling solution for GCN acceleration due to their inherent advantages in handling irregular computation patterns and their ability to implement custom memory hierarchies [13, 300, 303]. The reconfigurable nature of FPGAs enables domain-specific optimizations that directly address GCN computational bottlenecks, including specialized sparse matrix processing units, custom data formats for irregular access patterns, and fine-grained parallelism strategies [13, 301, 303]. Moreover, FPGA-based implementations offer superior energy efficiency compared to GPUs [307], making them particularly attractive for edge computing applications and large-scale deployment scenarios [300, 303, 306].

FPGA-based GCN accelerators have been explored through several complementary research directions, each addressing distinct computational challenges and introducing specific trade-offs. Some approaches focus on exploiting structural properties of graphs, such as adjacency matrix



symmetry and node degree distribution, to reduce computational redundancy [302,305]. Quantization-based techniques have also been explored to reduce arithmetic complexity and memory requirements. In particular, ternary quantization schemes have demonstrated strong effectiveness for resource-constrained implementations [303]. Another line of work emphasizes architectural flexibility and model generality, proposing overlay processors and dataflow architectures capable of supporting diverse GNN variants without requiring FPGA redeployment [13, 301].

However, existing FPGA-based GCN accelerators face several limitations that constrain their practical deployment. Specialized optimization techniques often sacrifice architectural generality, limiting accelerators to specific graph types or GCN variants [302,303,305]. Flexible overlay approaches, while supporting multiple models, introduce instruction interpretation overhead and complex scheduling mechanisms that reduce peak computational efficiency [13, 301]. Furthermore, most current implementations require sophisticated preprocessing algorithms or rely on graph-specific optimizations that compromise real-time processing capabilities and limit applicability to dynamic graph scenarios [306, 308].

This work addresses these limitations by introducing a lightweight and power-efficient streaming FPGA architecture that directly maps the core GCN operations into hardware while preserving broad model compatibility. The proposed design removes the overheads of instruction-driven overlays and avoids graph-specific preprocessing, enabling a purely dataflow-oriented execution that is well suited for resource-constrained devices. By centering the architecture around the fundamental sparse–dense and dense–dense computations of GCNs, and by employing scalable SpMM engines together with an efficient streaming datapath, the design achieves a favorable balance between computational throughput, resource utilization, and implementation simplicity.

The main contributions of this work are:

- a direct FPGA-based hardware realization of the GCN accelerator, eliminating instruction-interpretation and control overhead;
- a unified, fully streaming architecture that integrates sparse aggregation and dense transformation in a resource-efficient manner;
- a lightweight implementation tailored for mid-range FPGAs, demonstrating low power consumption and competitive performance on widely used citation datasets; and
- an analysis of architectural trade-offs that highlights the benefits of dataflow execution, resource scalability, and dataset-aware fixed-point quantization.

The remainder of this paper presents the related work (4.2), motivation (4.3), the proposed architecture (4.4), quantization methodology (4.5), evaluation results (4.6), and concluding remarks (4.7).

4.2 Literature Review

Graph Convolutional Networks (GCNs) have become essential for extracting features from non-Euclidean data structures, yet their unique sparsity and irregular computation patterns present substantial challenges to efficient hardware acceleration. In particular, the aggregation phase, dominated by sparse matrix-matrix multiplications, creates irregular memory access, low data reuse,

and workload imbalance. FPGA-based solutions have advanced through multiple perspectives, including symmetry exploitation, quantization, architectural flexibility, and software–hardware co-design.

4.2.1 Symmetry-Aware and Sparse Data Optimizations

Several architectures leverage graph symmetry to minimize redundant computation and memory traffic. Nair *et al.* process only the upper or lower triangular adjacency matrix, halving redundant edge accesses and improving memory reuse [302]. Their dynamically reconfigurable cores unify aggregation and transformation, reduce pipeline stalls, and achieve up to $110\times$ speedup over CPU and GPU baselines on Cora and Reddit. Jiang *et al.* present SSM-GCN, exploiting symmetric sparse matrix multiplication (SSpMM) through the SPCOO compression format, compatible with both symmetric and asymmetric matrices [305]. Their unified PE dynamically handles SpMM and SSpMM, reaching up to $318\times$ acceleration over CPU and $14\times$ over GPU, with DSP utilization up to 95% on the Alveo U50.

4.2.2 Quantization and Resource Efficiency

Quantization reduces arithmetic complexity and memory footprint. Chen *et al.* propose ATE-GCN, coupling asymmetrical ternary quantization with a unified PE array [303]. Specialized intervals fit the GCN weight distribution, keeping accuracy loss under 2% while reducing DSP usage up to $3\times$. Their system achieves $11\times$ latency reduction versus GPU on VCU118, combining ternary and SpMM modes at double frequency. Such quantization provides compute density but introduces portability trade-offs across different graph types.

4.2.3 Dataflow and Architectural Flexibility

Modern designs seek a balance between flexibility and performance. Tang *et al.* introduce Graph-OPU, an FPGA overlay processor with an integrated compiler and instruction set [13]. Supporting GCN, GAT, and GraphSAGE without redeployment, it unifies GEMM and SpMM via shared PEs and high-bandwidth memory, providing $1.45\times$ lower latency than prior overlays. Instruction abstraction, however, limits maximum efficiency versus static designs. Sarkar *et al.* present FlowGNN, a generic dataflow accelerator for real-time inference [301]. Abandoning preprocessing, FlowGNN pipelined message passing and node transformation using multi-queue multicasting, validated on Alveo U50 with up to $242\times$ speedup over CPU and $1.35\times$ over GPU. Generic architectures increase applicability but can penalize resource number efficiency.

4.2.4 Software–Hardware Co-Design Frameworks

Algorithmic and architectural co-optimization is another promising direction. Ran *et al.* combine attention-based graph sparsification with a two-stage accelerator for aggregation and combination phases [308]. Sparsified k -neighbor graphs improve locality, and their design speeds up inference by $7.3\times$ over CPU and $13.7\times$ over GPU, maintaining accuracy within 1%. Co-design demonstrates the benefits of integrated software preprocessing and hardware parallelism to address data irregularity.

4.2.5 Comparative Insights and Research Gaps

Recent works reveal a timeline from static symmetry-aware designs toward flexible overlays. Earlier systems optimized computational redundancy [302,308], while newer ones focus on compression, quantization, and reconfigurability [13,305]. Trade-offs remain between performance, flexibility, and regularity. No existing architecture fully unifies flexibility, low control overhead, and high compute density in resource-constrained FPGA fabrics. Overlays lose efficiency from instruction abstraction, while symmetry-based accelerators restrict model or graph applicability. Therefore, a gap persists for a streamlined, formula-centric FPGA design directly implementing

$$H^{(l+1)} = \sigma(\tilde{A}H^{(l)}W^{(l)})$$

that merges performance, efficiency, and broader GCN support.

4.3 Motivation

Existing works on FPGA-based acceleration of Graph Convolutional Networks (GCNs) primarily target high-end data center platforms characterized by advanced FPGA architectures such as UltraScale+ series, integration of High-Bandwidth Memory (HBM), and enhanced connectivity features. These designs typically leverage large FPGA fabrics and are tightly coupled with host systems to achieve superior performance within heterogeneous computing environments. However, such FPGA devices incur significantly higher costs due to their advanced hardware features, making them less suitable for cost-sensitive or resource-constrained applications.

In contrast, there is a noticeable gap in FPGA-based GCN accelerator designs optimized for embedded and edge computing contexts. Mid-range standalone FPGA devices, which offer a balanced trade-off between cost, power consumption, and computational capability, have not been adequately explored for efficient GCN acceleration. These devices often lack advanced memory subsystems like HBM and rely on more modest resources, demanding novel architectural approaches.

Therefore, there is a pressing need for scalable, lightweight, and energy-efficient GCN accelerator architectures tailored specifically for embedded edge platforms. Such architectures should maximize computational throughput and memory efficiency within the constraints of mid-range FPGA resources, enabling practical deployment of GCNs in power-sensitive and cost-restricted environments.

This work is motivated by bridging this gap: proposing an architecture that delivers power-efficient, streaming-based GCN acceleration on compact FPGAs suitable for edge computing, addressing key challenges of memory irregularity, workload imbalance, and computational efficiency without reliance on costly hardware features.

4.4 Architecture

4.4.1 Overview

As shown in Fig. 4.1, the accelerator employs a streaming architecture optimized for efficiency by adopting the following execution order:

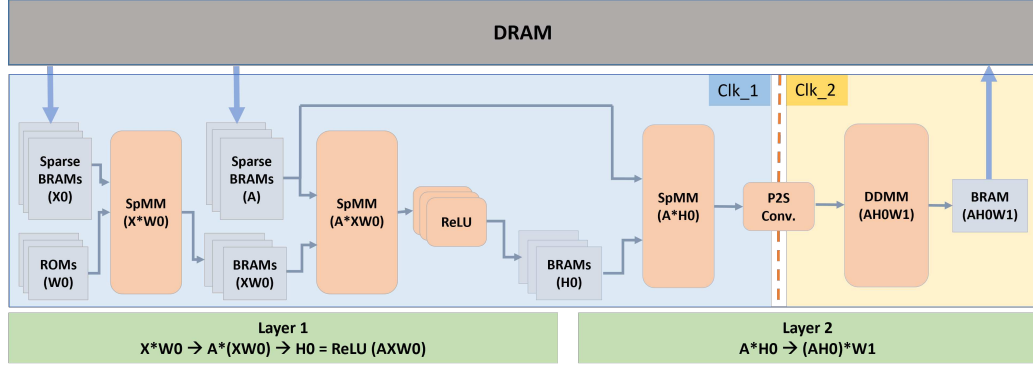


Figure 4.1: Overall streaming architecture of the proposed GCN accelerator.

In the first GCN layer, the operations proceed as:

1. XW_0
2. $\hat{A}(XW_0)$
3. $H_0 = \text{ReLU}(\hat{A}XW_0)$

In the second layer, the sequence is:

1. $\hat{A}H_0$
2. $(\hat{A}H_0)W_1$
3. $H_1 = \text{Softmax}(\hat{A}H_0W_1)$

The normalized adjacency matrix \hat{A} and input feature matrix X , stored in Compressed Sparse Row (CSR) format, are fetched from off-chip DRAM into on-chip memory. Weight matrices are stored in on-chip ROM memory banks, with each column allocated to a separate memory block to enable parallel data reads.

To improve efficiency, the first layer begins with matrix multiplication XW_0 to reduce dimensionality, since the number of columns in weight matrix W_0 is typically much smaller than the number of features in X . This yields a smaller intermediate matrix, facilitating more efficient subsequent computations. Table 4.1 illustrates the significant reduction in computational complexity achieved by rearranging the matrix multiplication order from $((AX)W)$ to $A(XW)$, especially for highly sparse feature matrices.

In the second layer, however, because feature width is already reduced, the multiplication H_0W_1 does not significantly decrease the number of multiply-accumulate (MAC) operations. Therefore, performing the graph propagation step AH_0 before the weight transformation $(AH_0)W_1$ is more efficient. This ordering eliminates the need for additional memory blocks to store intermediate results, thereby reducing hardware footprint. The memory optimization strategy will be detailed later.

The softmax function is assumed to be executed on the host CPU, given its computational complexity involving exponentials and divisions, which are challenging and resource-intensive to

Table 4.1: Computational complexity under different execution orders [13]

Datasets	Operations (MAC Count)	
	Order 1 ((AX)W)	Order 2 (A(HX))
Cora	62.8M	1.33M
CiteSeer	198.0M	2.23M
PubMed	165.5M	18.6M

implement efficiently on an FPGA. This choice aligns with the accelerator design goals prioritizing system performance and energy efficiency.

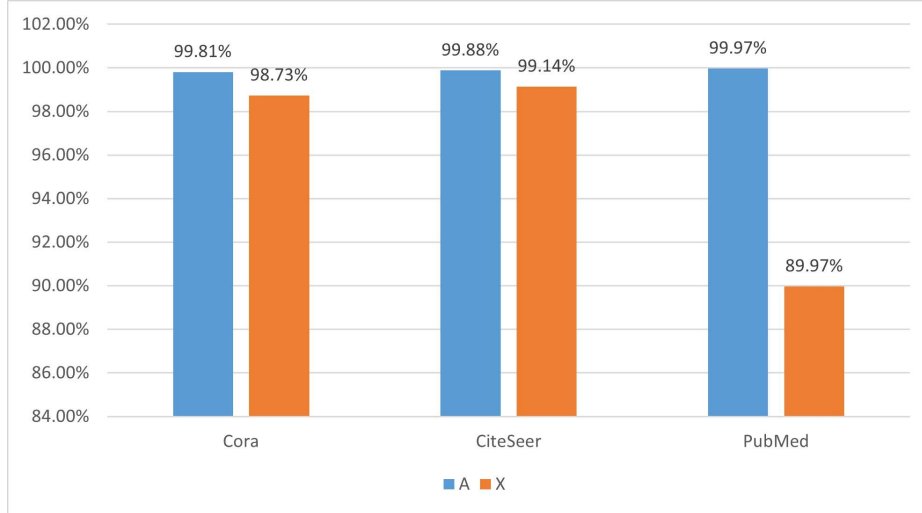


Figure 4.2: Sparsity analysis of input matrices for Cora, Citeseer, and Pubmed datasets. The blue bars represent the sparsity of the Adjacency Matrix (A), which drives the Aggregation phase design. The orange bars represent the sparsity of the Feature Matrix (X), influencing the efficiency of the Combination phase.

Considering the sparsity of matrices X and A in GCNs (see Fig. 4.2), the accelerator employs three Sparse-Dense Matrix Multiplication (SpMM) modules to efficiently compute XW_0 , $A(XW_0)$, and AH_0 , along with one Dense-Dense Matrix Multiplication (DDMM) module to compute $(AH_0)W_1$. To perform end-to-end system validation, the entire input matrices (i.e., \hat{A} and X) and output matrix are stored in the on-chip memory. Since the DDMM module processes inputs sequentially in row-major order, the PS2_Conv unit serializes the parallel output from the last SpMM module, which performs AH_0 .

In the following sections, each module is discussed in more detail.

4.4.2 Sparse-Dense Matrix Multiplication (SpMM) Module

The Sparse-Dense Matrix Multiplication (SpMM) module is a critical component of the accelerator, designed to efficiently perform sparse-dense matrix multiplications fundamental to Graph Convolutional Networks (GCNs). The input sparse matrices, i.e., \hat{A} and X , are encoded using the Compressed Sparse Row (CSR) format to optimize memory storage and access efficiency.

The CSR format represents the sparse matrix using three arrays: `row_ptr`, `col_idx`, and `values`. The `row_ptr` array stores the starting positions of each row's non-zero elements in the `col_idx` and `values` arrays. The `col_idx` array contains the column indices of these non-zero elements, and the `values` array holds the corresponding non-zero matrix entries. This format significantly reduces memory footprint by storing only non-zero entries and enables fast row-wise traversal, which is well-suited for streaming architectures.

Fig. 4.3 illustrates the block diagram of the SpMM module.

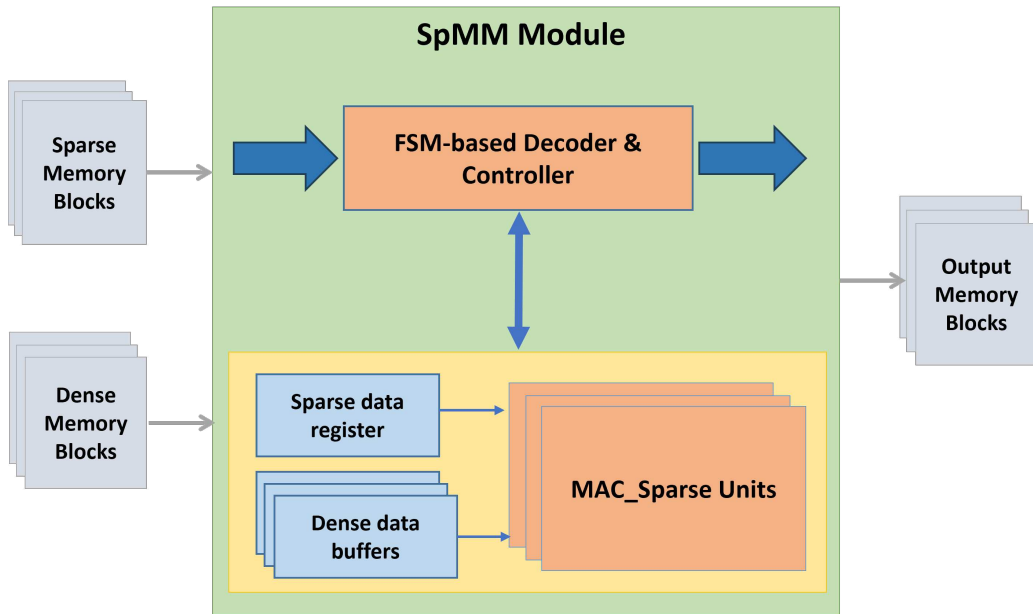


Figure 4.3: Overall block diagram of the Sparse Matrix Multiplication (SpMM) Module. The module consists of an FSM-based control unit coordinating data flow between the memory blocks and the core *MAC_Sparse* units.

During operation, the module iterates over each row of the input sparse matrix by reading `row_ptr` values to determine the number of non-zero elements per row. Each non-zero element is then fetched along with its column index and multiplied by the corresponding element from the dense feature matrix within the *MAC_Sparse* units. Partial products resulting from these multiplications are accumulated to form the dense output matrix rows. This computation follows an *inner product* dataflow, where output matrix rows are formed by accumulating products over the sparse matrix's non-zero elements multiplied by the dense matrix columns.

The architecture achieves parallelism by distributing the workload across a configurable number of *MAC_Sparse* units. Each *MAC_Sparse* unit contains a multiply-accumulate (MAC)

engine that processes assigned columns of the dense matrix, performing multiplications and accumulations in parallel for the non-zero elements of each sparse matrix row. By increasing the number of parallel `MAC_Sparse` units, the accelerator scales its throughput proportionally, effectively balancing hardware resource utilization and performance needs.

A finite state machine (FSM) orchestrates the entire computation procedure of the SpMM module. The FSM controls the sequencing of reading CSR arrays, fetching corresponding dense matrix data, enabling MAC operations, tracking progress via counters, and managing write-back of intermediate and final results to output memory. Beyond memory management, the FSM dynamically adapts to irregular sparsity patterns by processing varying numbers of non-zero elements per row, thus ensuring pipeline efficiency and preventing redundant computations.

In summary, by integrating CSR-based sparse data representation, parallel MAC units, inner product dataflow, and FSM-controlled computation and memory flow, the SpMM module delivers a scalable, resource-efficient, and high-throughput sparse-dense multiplication engine fundamental to the performance of the overall GCN accelerator.

4.4.3 Dense-Dense Matrix Multiplication (DDMM) Module

The Dense-Dense Matrix Multiplication (DDMM) module employs a one-dimensional systolic array architecture to efficiently perform matrix multiplication for dense matrices, which is essential for the later layers of Graph Convolutional Networks (GCNs).

As shown in Fig. 4.4, the weight parameters of the dense weight matrix are stored in dedicated memory banks, with each column assigned to a separate memory bank. Each memory bank is connected to a MAC unit within the systolic array, enabling parallel access to the weight values corresponding to that column.

The module accepts stream input vectors corresponding to activations or intermediate feature maps from the preceding pipeline stage. The input streams flow through the systolic array, where each MAC unit performs multiply-accumulate operations in a pipelined manner. The systolic structure facilitates sequential data propagation across MAC units, efficiently accumulating partial sums to compute each element of the output matrix.

Outputs are streamed out in a row-major order, supporting continuous data flow without intermediate buffering, thus minimizing latency and resource usage. This architecture is well-suited for FPGA implementations targeting embedded GCN acceleration.

4.4.4 Parallel-to-Serial Converter (P2S_Conv) Module

The *P2S_Conv* module bridges the data format gap between the parallel output of the Sparse-Dense Matrix Multiplication (SpMM) module and the serial input requirement of the Dense-Dense Matrix Multiplication (DDMM) module. It receives wide parallel data streams from the SpMM module and serializes these into a single-stream data flow suitable for subsequent processing. The architecture employs a *double buffering* (ping-pong) technique in conjunction with a dual-clock domain scheme, enabling concurrent data acquisition and serialization while maintaining data integrity and synchronization.

To prevent data loss and ensure continuous data flow without conflict, the module employs a double buffering technique, commonly known as ping-pong buffering. While one buffer is act-

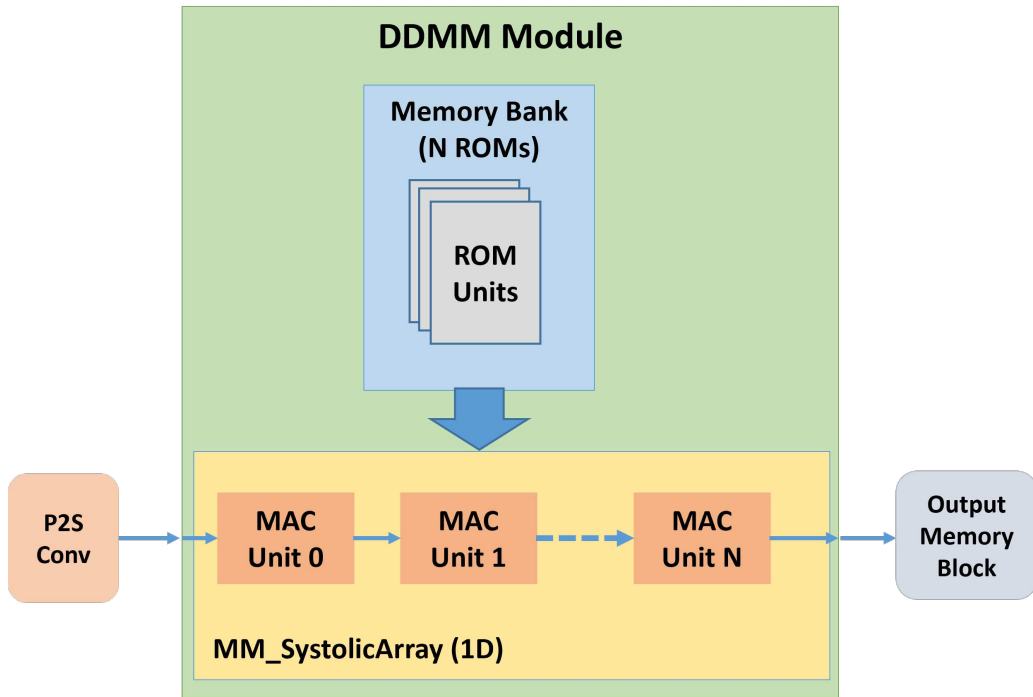


Figure 4.4: Block diagram of the Dense-Dense Matrix Multiplication (DDMM) Module. The module features a 1D Systolic Array for processing data received from the P2S Conversion stage, using an internal Memory Bank for weight storage.

ively being filled with new parallel data, the other buffer outputs serialized data to the DDMM module. This buffering strategy efficiently guarantees seamless handoff and synchronization between the domains or modules, thereby avoiding data conflicts that may arise if new data arrives before serialization of the previous batch is complete. Furthermore, to enhance performance and reduce serialization latency, the P2S_Conv module can operate at an increased clock frequency (twice the base system clock frequency in our case), after the initial data buffering. By operating the serialization logic at a higher clock frequency, the module achieves increased throughput and avoids data transfer conflicts.

Together, these design choices, i.e., double buffering combined with dynamic clock frequency scaling, not only ensure accurate and conflict-free data transfer but also maintain the streaming architecture's overall efficiency without introducing stalls or performance degradation.

Overall, the key fetures of *P2S_Conv* module includes:

- Converting wide parallel data (from SpMM) into a serialized data stream (for DDMM).
- Employing ping-pong double buffering to prevent data conflicts and losses.
- Incorporating a dual-clock domain design to enable concurrent parallel input loading and high-speed serial output generation.

- Utilizing an FSM-based control mechanism to ensure deterministic data flow and precise synchronization between buffers.
- Preserving data integrity and synchronization throughout the streaming pipeline.

4.5 Data-Aware Quantization and Scaling Methodology

The implementation of GCN accelerators on devices like FPGAs necessitates the use of reduced-precision numerical formats to minimize hardware resource utilization, latency, and power consumption. While software training environments typically rely on high-precision floating-point arithmetic, fixed-point representations are better suited for FPGA implementation due to their efficient mapping onto the device’s Digital Signal Processing (DSP) blocks and logic resources.

A critical challenge in adopting fixed-point formats is determining the optimal bit-width and scaling factor (the position of the binary point) for different intermediate tensors within the network. Suboptimal scaling leads to a trade-off between acceptable quantization error (precision loss due to limited range) and efficient resource utilization (wasted bit-width).

To address this challenge effectively, we employ a *data-aware, layer-wise static fixed-point quantization methodology* facilitated by comprehensive offline data distribution profiling. This approach is part of a software-hardware co-design flow that optimizes numerical precision while maintaining high inferencing accuracy for a given target dataset.

4.5.1 Profiling-Guided Fixed-Point Configuration

Our methodology utilizes a uniform 16-bit fixed-point representation throughout the accelerator for all major tensors, including the adjacency matrix (\mathbf{A}), feature matrix (\mathbf{X}), and weight matrices (\mathbf{W}). This fixed total width simplifies global data paths and control logic.

To maximize the effective utilization of this 16-bit width across various layers, which naturally exhibit different dynamic ranges, we perform a comprehensive analysis of value distributions during a software emulation phase. The methodology is structured as follows:

1. **Software Simulation and Range Analysis:** We execute the GCN model on target datasets (e.g., the CORA dataset) using a high-precision floating-point software environment. During this simulation, we monitor and record the maximum absolute value (V_{\max}) for all intermediate results generated after major computational steps, such as sparse matrix-vector multiplication ($\mathbf{A} \times \mathbf{X}$) and subsequent dense linear transformations ($\mathbf{A}\mathbf{X} \times \mathbf{W}_i$).
2. **Scaling Factor Determination:** Based on the observed V_{\max} for a given tensor, we determine the minimum required number of integer bits (I) to prevent overflow using the formula:

$$I = \lceil \log_2(V_{\max} + \epsilon) \rceil$$

where ϵ is a small margin. The remaining bits are allocated to the fractional part (F), such that the total width is fixed at 16 bits ($I + F = 16$). This approach ensures the dynamic range is captured without overflow while maximizing fractional precision.

3. **Hardware Parameterization:** The derived scaling factors and corresponding V_{\max} values are used as synthesis-time parameters for the hardware modules. This process configures the fixed-point arithmetic units statically before deployment, enabling rapid adaptation of the accelerator’s numerical behavior to specific dataset characteristics without requiring manual redesign of the HDL code.

The final implementation manages the fixed-point alignment through parameterized bit-shift operations determined entirely at synthesis time based on the input range parameters derived from profiling. This design choice ensures optimal static fixed-point precision for a given dataset, avoiding the area and latency overhead typically associated with run-time dynamic quantization logic.

4.6 Evaluation Results

4.6.1 Experimental Setup

We implement the proposed accelerator in SystemVerilog HDL on a Genesys 2 board equipped by a Kintex7 FPGA with AMD part number XC7K325T-2FFG900C. The complete design flow, encompassing synthesis, functional verification, physical implementation, and timing closure, is managed using the Xilinx Vivado 2021.1 Design Suite.

For validation of the end-to-end GCN inference pipeline, two canonical, small-to-medium scale graph datasets are employed: Cora and CiteSeer. Larger datasets, such as PubMed, could not be integrated for complete end-to-end validation under this memory model. Table 4.2 summarizes the key characteristics of the datasets used in this evaluation, including the number of nodes, edges, feature dimensions, classes, and the number of GCN layers.

Table 4.2: Characteristics of Citation Network Datasets for GCN Evaluation.

Dataset	Nodes	Edges	Features	Classes	Layers
Cora	2,708	10,556	1,433	7	2
Citeseer	3,327	9,104	3,703	6	2
Pubmed	19,717	88,648	500	3	2

4.6.2 Primary Experimental Results

This section details the key quantitative outcomes of the proposed GCN accelerator implementation, focusing on both resource utilization and core performance metrics derived from deployments on the Cora and CiteSeer benchmark datasets.

Hardware Utilization and Core Performance

The overall hardware utilization, as shown in Table 4.3, confirms a resource-optimized design point. The accelerator achieved efficient usage of key components, consuming only 55 and 54

Digital Signal Processing (DSP) components for the Cora and CiteSeer workloads, respectively. This low DSP footprint, alongside moderate Block RAM (BRAM) consumption (168 to 242 blocks), highlights the architecture’s focus on minimizing physical area overhead while prioritizing on-chip data locality to handle the sparse computations.

Table 4.3: FPGA Resource Utilization and Core Performance Metrics.

Dataset	Freq. (MHz)	LUT	FF	BRAM	DSP	Latency (ms)	Power (W)	EE (graphs/J)
Cora	150	1,446	2,823	168	55	1.77	0.305	~ 1,852
CiteSeer	145	1,501	3,168	242	54	3.47	0.518	~ 556

Throughput and Energy Efficiency

The implemented kernel delivered low execution latencies of 1.77 ms for Cora and 3.47 ms for CiteSeer, demonstrating suitability for most real-time applications [290]. A critical architectural objective was to achieve high Energy Efficiency (EE), defined as the ratio of throughput to average power consumption (graphs/Joule). The low measured power consumption (e.g., 0.305 W for Cora) directly translated into high energy efficiency, with the Cora implementation achieving approximately 1,852 graphs/J. This result unequivocally validates the architectural co-design strategy focused on high-throughput, low-power GCN acceleration.

Numerical Integrity and Final Accuracy

Despite the reduced precision required for hardware efficiency, the architecture successfully maintained numerical integrity. As illustrated in Table 4.4, the mean absolute error after quantization remained extremely low, measuring 1.66×10^{-4} for Cora and 1.25×10^{-4} for CiteSeer. This minimal error ensured that the final computed classification accuracy was preserved: 81.6% for Cora and 70.3% for CiteSeer. By matching the full-precision software baseline without any loss in accuracy, the 16-bit fixed-point scheme is definitively validated as a robust design choice for performance-critical GCN inference.

Table 4.4: Numerical Precision and Accuracy Preservation with 16-bit Fixed-Point Quantization.

Dataset	Max Absolute Error ($\times 10^{-4}$)	Mean Absolute Error ($\times 10^{-4}$)	Relative Quantization Error (MAE _{Normalized})	Classification Accuracy (Full-Precision / Fixed-Point)
Cora	7.12	1.66	0.0174%	81.6% / 81.6%
CiteSeer	10.20	1.25	0.0178%	70.3% / 70.3%

Note: Error metrics compare the 16-bit fixed-point hardware output (**H**) against the floating-point software reference (**S**). The consistently low relative error confirms the numerical integrity of the fixed-point arithmetic, ensuring no loss in classification accuracy.

4.6.3 Design Space Exploration for SpMM Parallelism

To evaluate the scalability of the proposed Sparse–Dense Matrix Multiplication (SpMM) module, a design space exploration is performed by varying the degree of parallelism, defined by the number of MAC_Sparse units, or Processing Elements (PEs), instantiated in the architecture. Each PE performs independent partial computations on non-zero elements of the sparse matrix, allowing the workload to be distributed across multiple computation units. This parameter directly affects the resource utilization, execution latency, and overall throughput of the accelerator.

Table 4.5 summarizes the results obtained from experiments conducted using the *Cora* dataset. The table reports the latency (T), the number of DSP blocks utilized (N_{DSP}), the achieved speedup (S), and the DSP efficiency (E_{DSP}) relative to the baseline configuration with $P = 2$. As shown, increasing the number of PEs leads to a significant reduction in latency—from 13.93 ms with 2 PEs to 1.77 ms with 16 PEs, corresponding to a speedup of $7.87\times$. Overall, these results demonstrate an effective trade-off between latency and resource utilization.

Table 4.5: Design Space Exploration for SpMM Parallelism Scaling with Cora.

P (#PEs/Layer)	T (Latency, ms)	N_{DSP} (#DSP Blocks)	S (Speedup)	E_{DSP} (DSP Efficiency)
2	13.93	13	1.00×	1.00×
4	6.98	19	1.99×	1.35×
8	3.50	31	3.98×	1.68×
16	1.77	55	7.87×	1.90×

Notes: **P** is the number of Processing Elements dedicated to the SpMM module per layer. **T** is the measured inference latency. **N_{DSP}** is the physical count of DSP blocks utilized. **S** is the speedup ratio relative to the baseline $P = 2$ configuration. **E_{DSP}** is the effective DSP utilization efficiency, calculated as the ratio of speedup to the normalized DSP count increase.

4.6.4 Cross-Platform Performance and Energy Efficiency on Benchmark Datasets

To rigorously evaluate the efficiency of our proposed FPGA-based GCN accelerator, we conducted a cross-platform comparison against state-of-the-art CPU and GPU implementations, specifically focusing on the widely adopted Cora and Citeseer benchmark datasets. The baseline performance data for the CPU (Intel i7-12700F) and GPU (Nvidia RTX 4090) implementations is sourced from a recent comparative study, the ATE-GCN paper [303].

Table 4.6 presents the end-to-end inference latency, corresponding power consumption, and resulting energy efficiency (EE) across the three platforms for both datasets. The results clearly demonstrate the performance and power advantages of the dedicated hardware solution:

- **Performance (Latency):** The GPU exhibits the lowest latency, achieving a $12.3\times$ speedup over the CPU on Cora (8.61 ms vs 0.70 ms). Our FPGA accelerator demonstrates

Table 4.6: Comparison of latency and Energy Efficiency (EE) across platforms on Cora and Citeseer datasets.

Platform	Dataset	Latency (ms)	Power (W)	EE (graphs/J)
CPU [303]	Cora	8.61	47.2	2.46
	Citeseer	9.98	N/A	N/A
GPU [303]	Cora	0.7	74.12	19.27
	Citeseer	1.13	N/A	N/A
FPGA (our)	Cora	1.67	0.553	1,082.83
	Citeseer	3.47	0.518	556.34

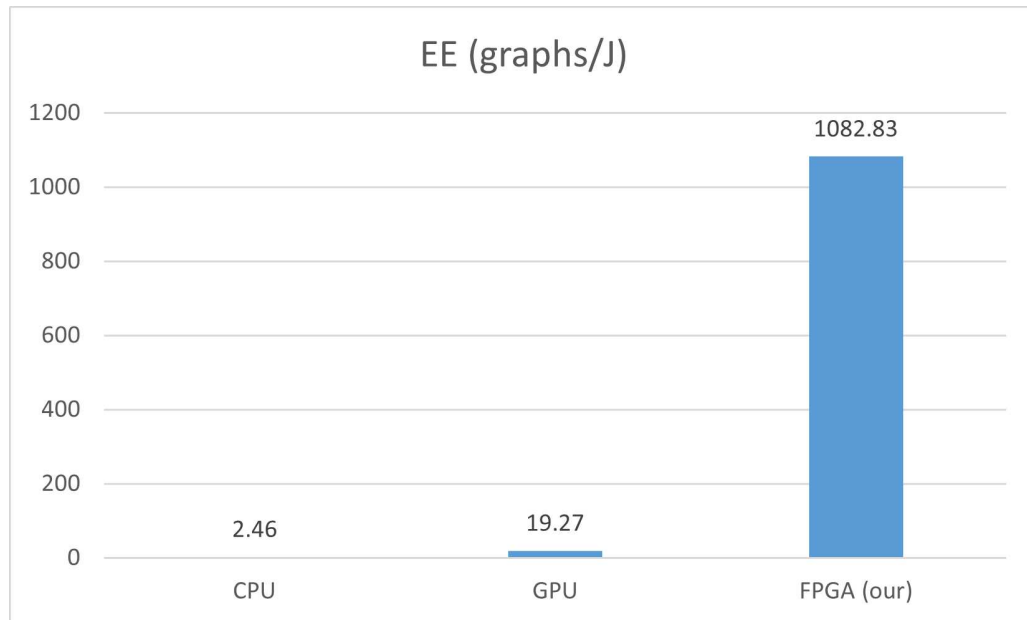


Figure 4.5: Energy Efficiency (*EE*) comparison across CPU, GPU, and the proposed FPGA accelerator, demonstrating the superior *EE* achieved by the dedicated hardware solution.

competitive latency, being approximately $2.38\times$ slower than the GPU on Cora (1.67 ms vs 0.70 ms)6.

- **Energy Efficiency (*EE*):** The *EE* metrics highlight the substantial advantage of the proposed FPGA-based design. On the Cora dataset, our FPGA design achieves an *EE* of 1,082.83 graphs/J, outperforming the GPU at 19.27 graphs/J by approximately $56.1\times$. This remarkable gain is primarily due to the drastically lower power consumption of our accelerator (0.553 W) compared to the CPU (47.2 W) and GPU (74.12 W).

This comparative evaluation confirms that while modern GPUs may excel in raw computa-

tional throughput, the proposed FPGA architecture offers a superior solution for energy-constrained environments. The ability of the dedicated FPGA hardware to maintain near-optimal energy efficiency, even when processing the irregular memory access patterns of the sparse aggregation phase, makes it a highly viable platform for real-time GCN inference, especially in edge computing applications where a high metric of graphs/Joule is paramount. The substantial advantage in Energy Efficiency (EE) for the proposed FPGA accelerator over commercial platforms is also visually summarized in Fig. 4.5.

4.6.5 Comparative Performance Analysis

We present a quantitative evaluation of our optimized GCN accelerator in comparison with state-of-the-art high-performance implementations, using the standard Cora and Citeseer datasets. Unlike existing architectures deployed on high-end FPGA platforms such as the AMD/Xilinx VCU118 and the Alveo U50, our design targets the embedded-class Kintex-7 platform, which imposes substantially stricter resource constraints.

As summarized in Table 4.7, our accelerator exhibits markedly lower resource utilization and power consumption—two key factors for edge-AI deployments on resource-constrained devices. The design requires at most 55 DSP slices (Cora), highlighting its suitability for platforms with tightly limited DSP budgets. This corresponds to an approximate $70\times$ reduction in DSP usage relative to the most resource-intensive architecture reported in the literature. The compact LUT and BRAM footprint further confirms the efficiency of the implementation on the constrained Kintex-7 fabric.

Measured core power consumption is exceptionally low, ranging from 0.553 W (Cora) to 0.518 W (Citeseer). Such low power demand is crucial for real-world edge-AI scenarios, where thermal limits and battery capacity preclude the deployment of power-hungry accelerators (e.g., compared to the 11.77 W reported for the VCU118 platform). In addition, despite adopting higher-precision quantization (16-bit fixed-point), our design achieves the highest reported accuracy among comparable works.

The comparatively higher absolute latency is a deliberate trade-off resulting from the architectural decision to minimize DSP usage. This design choice emphasizes feasibility on single-instance, resource-limited platforms rather than maximizing peak throughput.

To enable a fairer comparison of throughput across architectures operating at different clock frequencies and DSP budgets, we additionally report the normalized throughput metric defined as:

$$Norm.Throughput = \frac{1}{Latency(s) \times Freq.(s^{-1}) \times DSPs}$$

which incorporates both operating frequency and DSP utilization.

As shown in Fig. 4.6, our design achieves the second-highest performance after Graph-OPU [13], demonstrating that the tight integration of the dataflow architecture and processing elements enables highly effective computational extraction even from a constrained DSP pool.

Table 4.7: Performance and Hardware Resource Utilization of State-of-the-Art FPGA-based Graph Convolutional Network (GCN) Accelerators. The comparison highlights latency (μs), power consumption (W), accuracy (%), frequency (MHz), and the number of consumed resources, using the Cora and Citeseer graph datasets.

Work	Platform	Dataset	Quantization (#bits)	Freq. (MHz)	Latency (μs)	Resource Utilization (#)					Power (W)	Accuracy (%)
						LUT	FF	BRAM	URAM	DSP		
ATE-GCN [303]	VCU118	Cora	2-8	250	46.04	276,612	581,802	871	432	3,894	11.77	78.79
		Citeseer			75.48	N/A	N/A	N/A	N/A	N/A	N/A	62.42
SSM-GCN [305]	Alveo U50	Cora	8	225	64.29	184,000	156,000	386	0	1,152	N/A	N/A
		Citeseer			75.28	N/A	N/A	N/A	N/A	N/A	N/A	N/A
Graph-OPU [13]	Alveo U50	Cora	32	225	3.4	475,000	427,000	927	N/A	2,742	N/A	N/A
		Citeseer			5.09	N/A	N/A	N/A	N/A	N/A	N/A	N/A
Our work	Genesys 2	Cora	16	165	1670	1550	3041	168	N/A	55	0.553	81.6
		Citeseer		145	3380	1501	3168	242	N/A	54	0.518	70.3

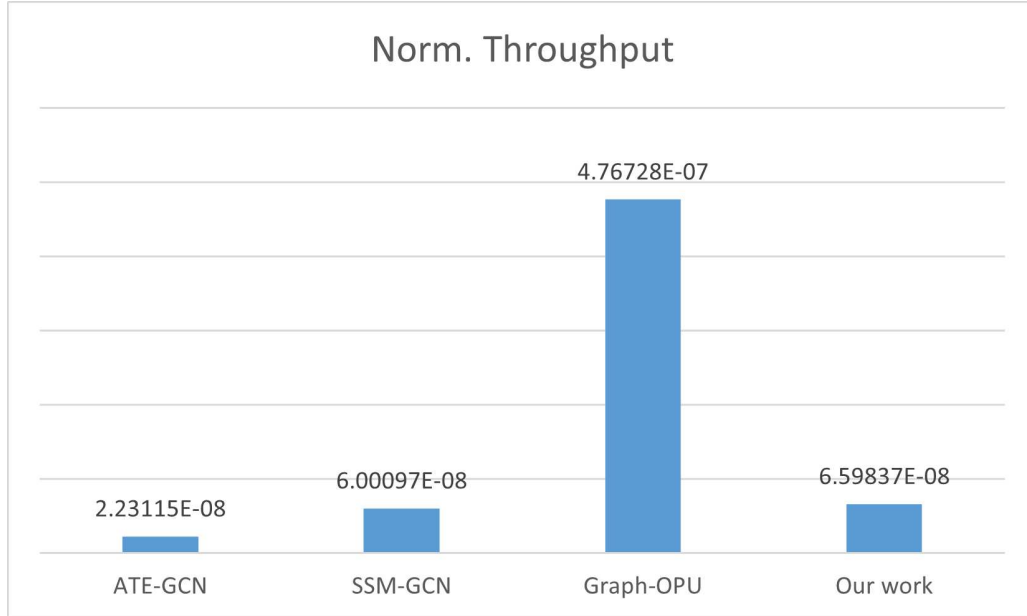


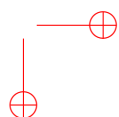
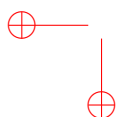
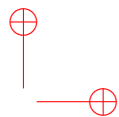
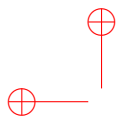
Figure 4.6: Comparison of normalized throughput with SOTA on Cora dataset.

4.7 Final Remarks

This paper presented a lightweight and energy-efficient FPGA accelerator for Graph Convolutional Networks (GCNs), targeting mid-range devices with tight resource and power constraints. The proposed design employs a unified streaming architecture, scalable SpMM engines, and a fixed-point quantization strategy to directly implement the GCN computation without instruction-driven overheads or preprocessing.

Experimental results on the Cora and Citeseer datasets demonstrate that the accelerator achieves 1.77–3.47 ms end-to-end latency while consuming only 0.305–0.553 W, resulting in energy efficiencies of up to 1,082 graphs/J, more than 50× higher than that of a modern GPU. Despite its compact footprint the architecture preserves full classification accuracy under 16-bit fixed-point arithmetic. A design-space exploration further confirms efficient scalability of the SpMM module under increasing parallelism.

Overall, this work shows that competitive GCN inference performance can be achieved on modest FPGA hardware through careful dataflow co-design and precision-aware implementation. Future work will extend the architecture to broader GNN variants, dynamic graph workloads, and multi-engine parallelism for higher throughput.





Chapter 5

Conclusion and Future Work

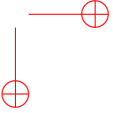
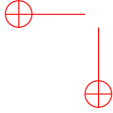
5.1 Conclusion

This dissertation investigated the design, optimization, and evaluation of FPGA-based hardware architectures for deep learning workloads, with a particular focus on real-time object detection (OD) and graph convolutional network (GCN) inference. Across heterogeneous platforms—from resource-constrained edge FPGA devices to modern MPSoCs—the thesis demonstrated how careful hardware–software co-design, dataflow-aware architecture development, and precision-aware techniques can substantially enhance performance and energy efficiency while maintaining model accuracy.

The first part of the thesis introduced a comprehensive and, to our knowledge, first dedicated survey on real-time object detection on FPGAs. By analyzing existing accelerator designs, architectural strategies, and optimization techniques, this survey established a unified categorization of FPGA-based OD accelerators and introduced pixel throughput as a fair and hardware-independent metric for cross-paper comparison. The survey highlighted clear trends toward dataflow architectures, hardware-friendly model adaptations, and the integration of quantization, pruning, and buffering optimizations to close the gap between algorithmic complexity and hardware constraints. This systematic review laid the scientific foundation for the subsequent contributions of the thesis.

The second part presented the first in-depth empirical study of real-time object detection on a heterogeneous FPGA-based MPSoC using Tiny YOLOv2 on an AMD Xilinx Zynq UltraScale+ platform. The implemented system achieved real-time performance with significantly higher energy efficiency compared to traditional CPU and GPU platforms. Through rigorous profiling using the Vitis AI Profiler, the work exposed execution bottlenecks rarely discussed in the literature, including insufficient CPU–DPU parallelism, DDR underutilization, and memory-bound pipeline stalls. These insights are critical for next-generation MPSoC-based accelerators and reinforce the need for hybrid execution pipelines that combine hardware parallelism with optimized software scheduling.

The final part of the thesis introduced a lightweight, streaming FPGA architecture for GCN inference. The proposed design combined sparse–dense matrix multiplication engines, a unified dataflow architecture, and a profiling-guided 16-bit fixed-point quantization scheme. Experiments on standard graph datasets demonstrated millisecond-level inference latency and energy efficiencies exceeding those of modern GPUs by more than an order of magnitude. This result confirms that well-designed FPGA accelerators can deliver competitive GNN performance while operating within the resource and power envelopes of embedded platforms. Furthermore, the design showed scalable parallelism, validating its applicability to future large-scale or multi-engine architectures.



Overall, the thesis shows that FPGA platforms—when paired with careful architecture co-design and precision-aware optimizations—can provide a highly efficient computing substrate for both CNN-based object detection and GNN workloads. The insights gained across the three complementary research directions provide a consolidated view of existing challenges, promising solutions, and open research opportunities for the broader community working at the intersection of machine learning and reconfigurable computing.

5.2 Future Work

Building on the outcomes of this research, several promising directions can further advance FPGA-based acceleration of deep learning models:

- **Hardware-Aware Neural Architecture Search (HW-NAS) for Real-Time OD**

Current OD models are not inherently optimized for the computational and memory characteristics of FPGAs. Future research should explore hardware-guided NAS pipelines that co-optimize model topology, precision, and hardware allocation. HW-NAS has the potential to generate FPGA-native OD architectures that maximize parallelism while minimizing on-chip memory pressure and external bandwidth demands.

- **Advanced Co-Design for MPSoCs**

The profiling results from the MPSoC study indicate substantial opportunities for improving task orchestration between PS and PL. Future work should investigate:

- fine-grained CPU–DPU pipelining,
- task-level parallelism and asynchronous execution models,
- optimized DDR access scheduling, and
- hybrid control/dataflow strategies.

Such optimizations could significantly improve throughput, especially for multi-stage OD pipelines and transformer-based models.

- **Streaming and Multi-Engine Architectures for GNNs**

The proposed GCN accelerator can be extended in several key directions:

- Support for broader GNN variants such as GraphSAGE, GAT, and message-passing transformers
- Dynamic graph processing, enabling inference on evolving or batched graph workloads
- Multi-engine parallelism to increase throughput for large-scale inference or multi-tenant applications

These extensions would transform the accelerator from a single-workload design into a flexible GNN hardware substrate.

- Efficient FPGA Implementations of Transformer-Based Models

Transformer architectures are rapidly becoming dominant in both vision and graph domains. However, their high parameter count and non-local attention mechanism pose significant challenges for FPGA deployment. Future work should explore:

- hardware-friendly transformer variants,
- low-cost attention approximations,
- hierarchical memory designs for attention maps,
- reconfigurable/folded datapaths for multi-head attention.

Achieving real-time transformer inference on mid-range FPGAs remains an open research problem with high practical relevance.

- Unified Toolflows and Parametric Hardware Libraries

Developing robust, open-source toolflows capable of bridging AI model development, quantization, and FPGA deployment will greatly reduce time-to-market and democratize FPGA-based ML acceleration. Parameterizable hardware libraries—especially for OD layers, attention blocks, and sparse operations—could standardize design reuse and simplify experimentation across hardware families.

- Distributed FPGA Clusters for Hard Real-Time Systems

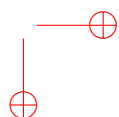
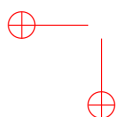
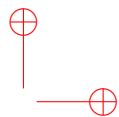
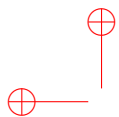
For demanding applications such as autonomous driving or industrial automation, future research may explore distributed FPGA clusters connected through high-bandwidth, low-latency links. Key topics include:

- load balancing and partitioning of CNN/GNN workloads,
- latency-aware scheduling,
- fault tolerance,
- multi-device synchronization.

This direction can unlock extreme scalability while preserving determinism.

Final Outlook

The convergence of machine learning and reconfigurable computing is accelerating rapidly. As deep learning models become more sophisticated and pervasive, the need for efficient, scalable, and application-specific hardware will only intensify. This dissertation contributes to this trajectory by providing both theoretical insights and practical implementations that demonstrate the capabilities of FPGAs in the era of modern AI. The future work directions outlined above serve as a roadmap toward even more powerful, flexible, and energy-efficient hardware architectures capable of enabling the next generation of intelligent systems.



Appendix A

Some General Definitions in The Context of Object Detection Metrics

Measuring the Intersection over Union (*IoU*) is a popular way to evaluate an object detector's localization accuracy. Given a ground truth Bounding box (*Bbox*), it calculates the ratio of the common, or overlapped, area between the ground truth Bbox and the predicted Bbox to the area of their union, as illustrated in Figure A.1. The greater the IoU ratio means the more detection accuracy. By setting a threshold, e.g., $IoU > 0.5$, as done in [50], it can be determined how precisely the object is localized.

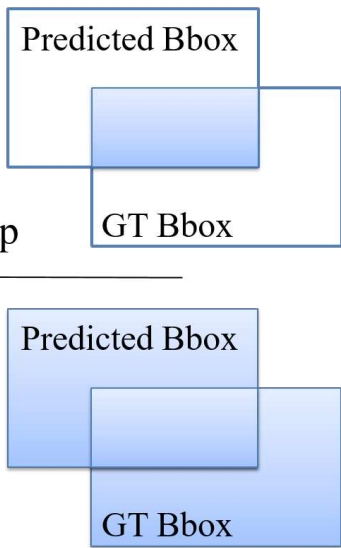
$$IoU = \frac{\text{Area of Overlap}}{\text{Area of Union}}$$


Figure A.1: Definition of Intersection over Union (IoU), calculated by dividing the intersection of the predicted bounding box (Bbox) and ground truth (GT) Bbox by the union of them.

Based on the calculated IoU and a pre-determined threshold, some other metrics, such as precision and recall, can be obtained. The IoU greater than or equal to the threshold value indicates that the prediction is correct. It should be noted that IoU is basically related to the localization aspect of object detection tasks and is not directly related to classification.

Given the real label and the predicted one for each object, the prediction result can be classified as shown in Table A.1. Accordingly, accuracy, precision, and recall values are defined as follows:

Table A.1: Prediction result category in object detection models based on the label provided by the dataset (real label).

Real Label	Predicted Label	Prediction Category
Positive	Positive	True Positive (TP)
Positive	Negative	False Negative (FN)
Negative	Positive	False Positive (FP)
Negative	Negative	True Negative (TN)

$$Accuracy := \frac{(TP + TN)}{(TP + FP + TN + FN)}$$

$$Precision := \frac{TP}{(TP + FP)}$$

$$Recall := \frac{TP}{(TP + FN)}$$

Accuracy determines how accurately the model makes predictions, showing the overall performance of the model. However, this metric has a significant drawback: it performs poorly with imbalanced data, where one class significantly outnumbers the others [309].

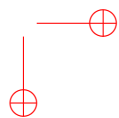
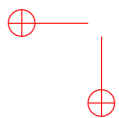
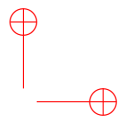
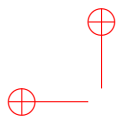
Two alternatives for accuracy are “precision” and “recall,” which focus on “true positive” predictions. The former represents how precisely the model can predict the positive class, measuring the quality of the detection task, while the latter, also known as “sensitivity”, refers to the proportion of “true positive” predictions to all positive instances in the dataset. From another point of view, if we need to minimize the false positive predictions, we should focus on improving precision, while recall is more important when the ability to predict all positives outweighs the detection accuracy.

The “F1 score”, also known as the “balanced F-score” or “F-measure”, provides a means to evaluate the trade-off between recall and precision. Representing the harmonic mean of precision and recall, it ranges from 0 (indicating the worst value) to 1 (indicating the best value). The F1-score disregards variations in confidence values, limiting its utility to comparing object detectors solely at a predetermined confidence threshold level [310].

$$F1 := \frac{2}{\frac{1}{Precision} + \frac{1}{Recall}} = \frac{(2 * Precision * Recall)}{(Precision + Recall)} \quad (A.1)$$

Average precision (*AP*) has recently emerged as the most commonly utilized evaluation metric for detection tasks [40]. *AP* can be derived by calculating the area under the precision and recall curve (PR Curve). In other words, it indicates the average precision values across all recall values ranging from 0 to 1.

Widely utilized in computer vision, AP serves as a popular evaluation measure for assessing the prediction accuracy of object detection models [55]. AP can be assessed across various IoU threshold ranges. For instance, it can be computed for 10 IoU values ranging from 50% to 95%, with increments of 5%, typically denoted as “AP@50:5:95”. Additionally, it can be evaluated at specific IoU thresholds, commonly 50% and 75% denoted as “AP50” and “AP75” respectively [44].





Appendix B

Overview of FPGA Architecture and Design Approaches

B.0.1 FPGA Architecture

An FPGA consists of a set of programmable logic blocks, memory blocks, and specialized arithmetic units such as Digital Signal Processing (DSP) blocks. These components are interconnected by programmable interconnects, allowing for the development of highly flexible digital circuits. In addition, programmable input/output (I/O) blocks facilitate external connectivity [158]. Figure B.1 shows the basic architecture of an FPGA.

Some key factors in selecting an FPGA device as a target device for a specific application typically include the number of available I/O blocks and programmable logic blocks, the number and capabilities of fixed-function logic blocks (such as multipliers), and the total size of on-chip memory resources. Programmable logic blocks, which form the fundamental building blocks of FPGAs, directly impact the design's flexibility. Although there is no rigid standard for the architecture of these blocks, they typically include Flip-Flops (*FFs*), Look-Up Tables (*LUTs*), arithmetic carry logic, and multiplexers. LUTs are small, programmable memory blocks that can store truth tables of logic functions.

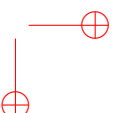
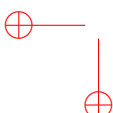
Another key factor in choosing FPGAs for a specific application is the available memory resources. There are two main types of on-chip memory within an FPGA device:

- Distributed RAM: In addition to forming various logic functions, LUTs can store data sets and act as “distributed” memory cells throughout the FPGA, as their name suggests. A k -input LUT can store 2^k bits.
- Block RAM (*BRAM*): BRAMs are built by dedicated SRAM memory blocks. They are typically used to store large amounts of data inside an FPGA. Additionally, supplementary peripheral circuitry enhances BRAM's reconfigurability for diverse applications and facilitates its connection to the programmable routing within the FPGA.

B.0.2 FPGA Design Approaches

The operation of all FPGA blocks and the configuration of the programmable interconnects are typically managed by a so-called “bitstream”, which is provided to internal dedicated (*SRAM*) cells [311]. This operation is done either once at boot time or even during the run time (partial reconfiguration).

The following outlines three primary FPGA design approaches, allowing developers to select one or a combination of them to realize their designs:



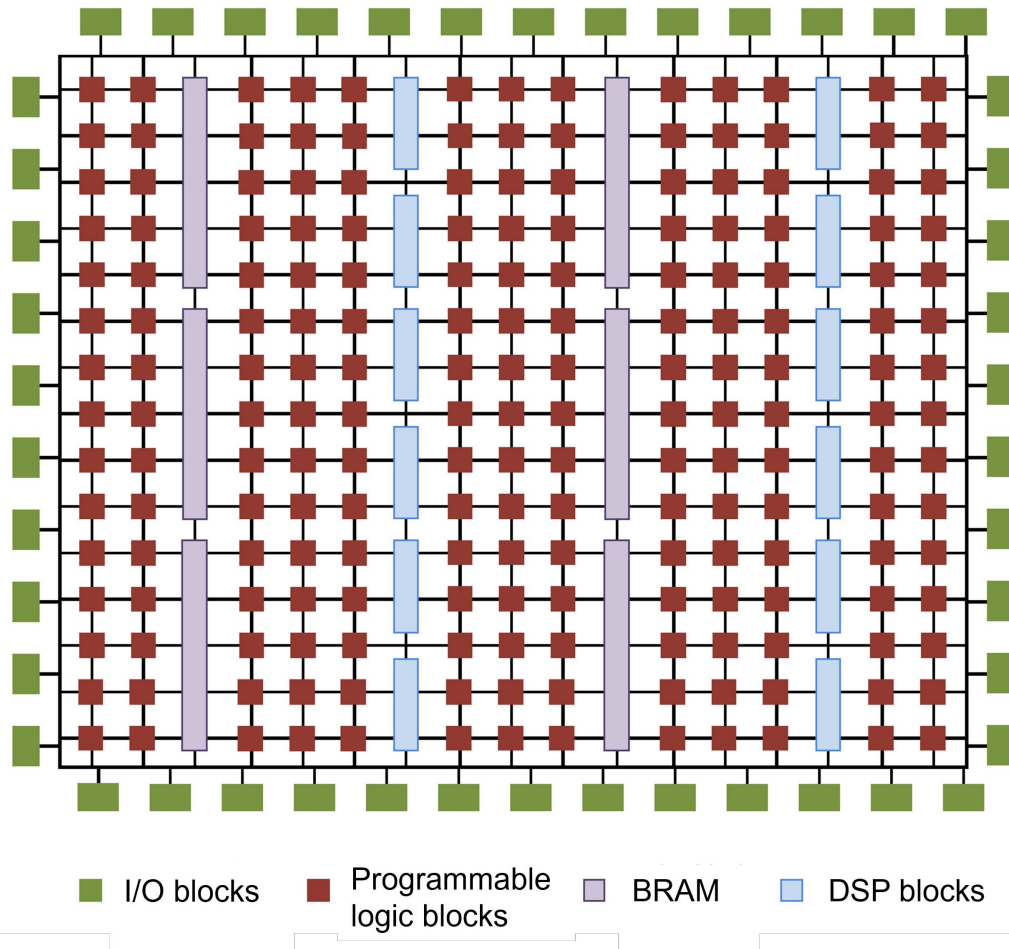


Figure B.1: An overview of basic FPGA architecture, including programmable logic blocks, Block RAMs (*BRAMs*), DSP blocks, input/output blocks (*I/O*), and programmable interconnects (shown by black lines). (Figure adapted from [12]).

1. Hardware Description Language (*HDL*) design: In this approach, designers describe the intended functionality using an HDL, such as Verilog or VHDL. Then, the HDL design is compiled through an intricate Computer-Aided Design (CAD) flow, generating a “bitstream” file [311]. This bitstream file, also referred to as the configuration file, is used to program the FPGA.
2. High-level design: High-level languages like C/C++, OpenCL, and SyCL can also be employed for FPGA design. In this approach, the high-level design can be translated into its corresponding HDL using available tools such as Vivado HLS [312], HLS compiler [313], and HDL Coder [314].
3. Model-based design (*MBD*): This approach can be considered a form of high-level design. Instead of using programming languages, high-level models or blocks are employed to

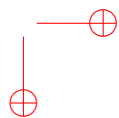
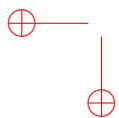
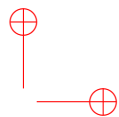
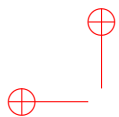
design and simulate the systems to be implemented. In this approach, once the system's behavior is verified, the design is translated into HDL code. Some of the most commonly used tools in MBD include Simulink [314], LabVIEW [315], and Model Composer [316].

Adopting each of these approaches can depend on various factors, including the designer's preference and the target application. For example, software engineers who are more accustomed to high-level languages might prefer the second approach. Additionally, for some applications, such as aerospace and automotive, due to the availability of ready-made models and the ease of behavioral simulation of developing systems, adopting MBD is very popular [317].

Each of these approaches has its own advantages and disadvantages, which are beyond the scope of this study. However, it is noteworthy that in all cases, having in-depth knowledge of hardware implementation is essential for designing an optimal system.

Additionally, helpful tools, frameworks, hardware libraries, and reusable reconfigurable components ("IP Cores") are continuously being introduced to facilitate the design of FPGA-based systems in various fields. One example of a very commonly used IP core for developing FPGA-based object detection systems is the Deep Learning Processor Unit (DPU) [195]. This programmable engine enables designers to implement many object detection models on FPGAs without requiring low-level design. In this regard, designers can use tools, such as the Vitis AI tool to select, prepare, optimize, and evaluate a detection model and compile the instructions used in the deployed DPU.

Further information about FPGA architecture and design could be found, e.g., in [318, 319].



Appendix C

Roofline Model

In this visual model, the peak computational performance provided by the hardware platform and the maximum off-chip memory bandwidth are the two critical factors for estimating the attainable performance [267].

The equation C.1, represents a “roofline”-type curve on the Cartesian plane where the X-axis is the Operational Intensity (*OI*) or Compute-To-Communication (*CTC*) measured in Floating Point Operations per Byte or just Operations per Byte [201] and the Y-axis represents the computational performance measured in Floating Point Operations per second or just Operations per second (Figure C.1).

$$\text{AchievablePerformance} = \min \begin{cases} \text{PeakComputationalPerformance} \\ \text{PeakMemoryBandwidth} \times \text{CTC} \end{cases} \quad (\text{C.1})$$

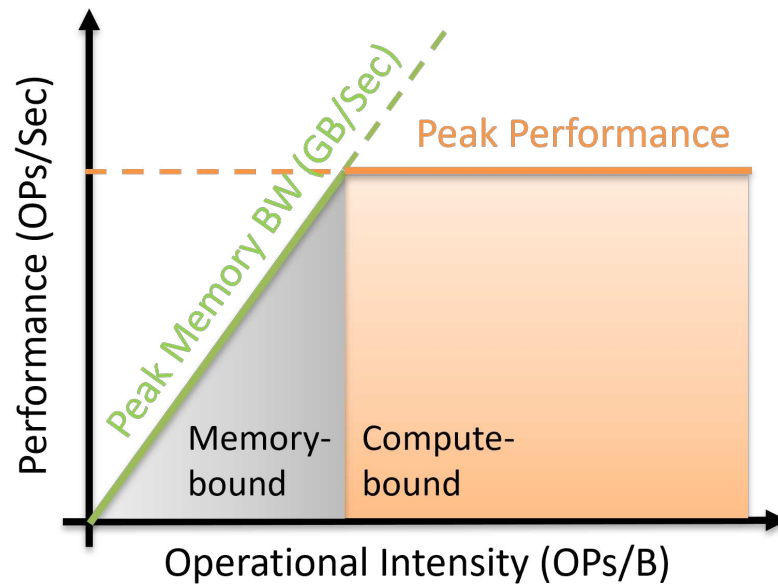


Figure C.1: The Roofline model.

In the Roofline model, an application’s performance can be represented as a point on a

Cartesian plane, which visualizes the distance between the actual performance and the obtainable one.

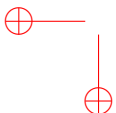
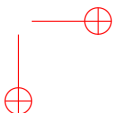
In an application, the CTC ratio, also known as operational intensity, represents the ratio of the total number of executed operations to the amount of data transferred to and from external memory and can be measured as operation per byte (*op/Byte*) or floating point operation per Byte (*Flops/Byte*). The performance of the application is calculated at run time while executing it.

The roofline constrains the achievable performance, delineating the figure into two distinct regions: *memory-bound* and *compute-bound* areas. The ridge point, also known as the machine balance point, highlighted in Figure C.1, is at the intersection between the diagonal and horizontal lines. Applications with CTC ratios on the left side of the ridge point are considered “memory-bound”, indicating that the system cannot efficiently utilize all computational resources due to limited off-chip communication. On the other hand, applications with a CTC ratio on the right side of the ridge point are referred to as “compute-bound”.

If the application is memory-bound, optimizing memory inefficiencies is often a fruitful strategy, focusing on factors such as memory access pattern, data locality, and cache reuse [267, 269, 320]. If the application is compute-bound, the bottleneck lies in the computational power or efficiency of the accelerator. In this case, applying some optimization techniques such as loop unrolling and loop reordering may help [267].



List of Acronyms

- AXI** Advanced eXtensible Interface
- BRAM** Block Random Access Memory
- CNN** Convolutional Neural Network
- COO** Coordinate List Format
- CPU** Central Processing Unit
- CSR** Compressed Sparse Row
- DDR** Double Data Rate (External Memory)
- DL** Deep Learning
- DPU** Deep Processing Unit
- DSP** Digital Signal Processing Block
- FPS** Frames Per Second
- FPGA** Field-Programmable Gate Array
- GCN** Graph Convolutional Network
- GNN** Graph Neural Network
- GPU** Graphics Processing Unit
- HLS** High-Level Synthesis
- IP** Intellectual Property (Hardware Core)
- ML** Machine Learning
- MPSoC** Multiprocessor System-on-Chip
- OD** Object Detection
- 
- 

OPS Operations Per Second

PL Programmable Logic

PS Processing System

PT Pixel Throughput

ReLU Rectified Linear Unit

SA Systolic Array

SDK Software Development Kit

SoC System-on-Chip

SpMM Sparse–Dense Matrix Multiplication

SRAM Static Random Access Memory

Bibliography

- [1] J. Redmon, S. Divvala, R. Girshick, and A. Farhadi, "You only look once: Unified, real-time object detection," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 779–788.
- [2] W. Liu, D. Anguelov, D. Erhan, C. Szegedy, S. Reed, C.-Y. Fu, and A. C. Berg, "Ssd: Single shot multibox detector," in *Computer Vision—ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11–14, 2016, Proceedings, Part I 14*. Springer, 2016, pp. 21–37.
- [3] S. I. Venieris, A. Kouris, and C.-S. Bouganis, "Toolflows for mapping convolutional neural networks on FPGAs: A survey and future directions," *ACM Computing Surveys (CSUR)*, vol. 51, no. 3, pp. 1–39, 2018.
- [4] A. Gholami, S. Kim, Z. Dong, Z. Yao, M. W. Mahoney, and K. Keutzer, "A survey of quantization methods for efficient neural network inference," in *Low-Power Computer Vision*. Chapman and Hall/CRC, 2022, pp. 291–326.
- [5] J. Gou, B. Yu, S. J. Maybank, and D. Tao, "Knowledge distillation: A survey," *International Journal of Computer Vision*, vol. 129, no. 6, pp. 1789–1819, 2021.
- [6] R. Xu, S. Ma, Y. Wang, Y. Guo, D. Li, and Y. Qiao, "Heterogeneous systolic array architecture for compact cnns hardware accelerators," *IEEE Transactions on Parallel and Distributed Systems*, vol. 33, no. 11, pp. 2860–2871, 2021.
- [7] D. T. Nguyen, T. N. Nguyen, H. Kim, and H.-J. Lee, "A high-throughput and power-efficient FPGA implementation of yolo cnn for object detection," *IEEE Transactions on Very Large Scale Integration (VLSI) Systems*, vol. 27, no. 8, pp. 1861–1873, 2019.
- [8] L. Chang, S. Zhang, H. Du, Y. Chen, and S. Wang, "A reconfigurable neural network processor with tile-grained multicore pipeline for object detection on fpga," *IEEE Transactions on Very Large Scale Integration (VLSI) Systems*, vol. 29, no. 11, pp. 1967–1980, 2021.
- [9] L.-Y. Liew and S.-D. Wang, "Object detection edge performance optimization on fpga-based heterogeneous multiprocessor systems," in *2022 IEEE International Conference on Consumer Electronics (ICCE)*. IEEE, 2022, pp. 1–6.
- [10] V. Jain, N. Jadhav, and M. Verhelst, "Enabling real-time object detection on low cost fpgas," *Journal of Real-Time Image Processing*, vol. 19, no. 1, pp. 217–229, 2022.
- [11] H.-S. Suh, J. Meng, T. Nguyen, V. Kumar, Y. Cao, and J.-S. Seo, "Algorithm-hardware co-optimization for energy-efficient drone detection on resource-constrained fpga," *ACM Transactions on Reconfigurable Technology and Systems*, vol. 16, no. 2, pp. 1–25, 2023.

- [12] B. Ronak and S. A. Fahmy, "Mapping for maximum performance on FPGA dsp blocks," *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems*, vol. 35, no. 4, pp. 573–585, 2015.
- [13] E. Tang, S. Li, R. Chen, H. Zhou, Y. Ma, H. Zhang, J. Yu, and K. Wang, "Graph-opu: A highly flexible fpga-based overlay processor for graph neural networks," *ACM Transactions on Reconfigurable Technology and Systems*, vol. 17, no. 4, pp. 1–33, 2024.
- [14] T. Turay and T. Vladimirova, "Toward performing image classification and object detection with convolutional neural networks in autonomous driving systems: A survey," *IEEE Access*, vol. 10, pp. 14 076–14 119, 2022.
- [15] B. Wu, F. Iandola, P. H. Jin, and K. Keutzer, "Squeezedet: Unified, small, low power fully convolutional neural networks for real-time object detection for autonomous driving," in *Proceedings of the IEEE conference on computer vision and pattern recognition workshops*, 2017, pp. 129–137.
- [16] H. Karaoguz and P. Jensfelt, "Object detection approach for robot grasp detection," in *2019 International Conference on Robotics and Automation (ICRA)*. IEEE, 2019, pp. 4953–4959.
- [17] S. Chatterjee, F. H. Zunjani, and G. C. Nandi, "Real-time object detection and recognition on low-compute humanoid robots using deep learning," in *2020 6th International Conference on Control, Automation and Robotics (ICCAR)*. IEEE, 2020, pp. 202–208.
- [18] A. Kaur, Y. Singh, N. Neeru, L. Kaur, and A. Singh, "A survey on deep learning approaches to medical images and a systematic look up into real-time object detection," *Archives of Computational Methods in Engineering*, vol. 29, pp. 1–41, 2021.
- [19] R. Yang and Y. Yu, "Artificial convolutional neural network in object detection and semantic segmentation for medical imaging analysis," *Frontiers in oncology*, vol. 11, p. 638182, 2021.
- [20] G. C. Buttazzo, *Hard real-time computing systems: predictable scheduling algorithms and applications*. Springer Science & Business Media, 2011, vol. 24.
- [21] L. Pezzarossa, M. Schoeberl, and J. Sparsø, "Reconfiguration in FPGA-based multi-core platforms for hard real-time applications," in *2016 11th International Symposium on Reconfigurable Communication-centric Systems-on-Chip (ReCoSoC)*. IEEE, 2016, pp. 1–8.
- [22] S. Saha, S. Ehsan, A. Stoica, R. Stolkin, and K. McDonald-Maier, "Real-time application processing for FPGA-based resilient embedded systems in harsh environments," in *2018 NASA/ESA Conference on Adaptive Hardware and Systems (AHS)*. IEEE, 2018, pp. 299–304.
- [23] Y. LeCun, Y. Bengio, and G. Hinton, "Deep learning," *nature*, vol. 521, no. 7553, pp. 436–444, 2015.
- [24] N. Carion, F. Massa, G. Synnaeve, N. Usunier, A. Kirillov, and S. Zagoruyko, "End-to-end object detection with transformers," in *European conference on computer vision*. Springer, 2020, pp. 213–229.
- [25] X. Zhu, W. Su, L. Lu, B. Li, X. Wang, and J. Dai, "Deformable detr: Deformable transformers for end-to-end object detection," *arXiv preprint arXiv:2010.04159*, 2020.
- [26] Y. Zhao, W. Lv, S. Xu, J. Wei, G. Wang, Q. Dang, Y. Liu, and J. Chen, "Detrs beat yolos on real-time object detection," *arXiv preprint arXiv:2304.08069*, 2023.
- [27] K. Zeng, Q. Ma, J. W. Wu, Z. Chen, T. Shen, and C. Yan, "FPGA-based accelerator for object detection: a comprehensive survey," *Journal of Supercomputing*, vol. 78, pp. 14 096–14 136, 8 2022.
- [28] T. I. Amosa, P. Sebastian, L. I. Izhar, O. Ibrahim, L. S. Ayinla, A. A. Bahashwan, A. Bala, and Y. A. Samaila, "Multi-camera multi-object tracking: a review of current trends and future advances," *Neurocomputing*, vol. 552, p. 126558, 2023.

- [29] L. Fei and B. Han, "Multi-object multi-camera tracking based on deep learning for intelligent transportation: A review," *Sensors*, vol. 23, no. 8, p. 3852, 2023.
- [30] R. Girshick, J. Donahue, T. Darrell, and J. Malik, "Region-based convolutional networks for accurate object detection and segmentation," *IEEE transactions on pattern analysis and machine intelligence*, vol. 38, no. 1, pp. 142–158, 2015.
- [31] R. Girshick, "Fast r-cnn," in *Proceedings of the IEEE international conference on computer vision*, 2015, pp. 1440–1448.
- [32] S. Ren, K. He, R. Girshick, and J. Sun, "Faster r-cnn: Towards real-time object detection with region proposal networks," *Advances in neural information processing systems*, vol. 28, 2015.
- [33] A. Bochkovskiy, C.-Y. Wang, and H.-Y. M. Liao, "Yolov4: Optimal speed and accuracy of object detection," *arXiv preprint arXiv:2004.10934*, 2020.
- [34] C.-Y. Wang, A. Bochkovskiy, and H.-Y. M. Liao, "Yolov7: Trainable bag-of-freebies sets new state-of-the-art for real-time object detectors," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2023, pp. 7464–7475.
- [35] P. Mittal, "A comprehensive survey of deep learning-based lightweight object detection models for edge devices," *Artificial Intelligence Review*, vol. 57, no. 9, p. 242, 2024.
- [36] A. B. Amjoud and M. Amrouch, "Object detection using deep learning, cnns and vision transformers: a review," *IEEE Access*, 2023.
- [37] V. Kamath and A. Renuka, "Deep learning based object detection for resource constrained devices: Systematic review, future trends and challenges ahead," *Neurocomputing*, vol. 531, pp. 34–60, 2023.
- [38] R. Kaur and S. Singh, "A comprehensive review of object detection with deep learning," *Digital Signal Processing*, vol. 132, p. 103812, 2023.
- [39] A. Setyanto, T. B. Sasongko, M. A. Fikri, and I. K. Kim, "Near-edge computing aware object detection: A review," *IEEE Access*, vol. 12, pp. 2989–3011, 2024.
- [40] Z. Zou, K. Chen, Z. Shi, Y. Guo, and J. Ye, "Object detection in 20 years: A survey," *Proceedings of the IEEE*, vol. 111, no. 3, pp. 257–276, 2023.
- [41] Z. Huang, S. Yang, M. Zhou, Z. Gong, A. Abusorrah, C. Lin, and Z. Huang, "Making accurate object detection at the edge: Review and new approach," *Artificial Intelligence Review*, vol. 55, no. 3, pp. 2245–2274, 2022.
- [42] J. Kaur and W. Singh, "Tools, techniques, datasets and application areas for object detection in an image: a review," *Multimedia Tools and Applications*, vol. 81, no. 27, pp. 38 297–38 351, 2022.
- [43] S. S. A. Zaidi, M. S. Ansari, A. Aslam, N. Kanwal, M. Asghar, and B. Lee, "A survey of modern deep learning-based object detection models," *Digital Signal Processing*, vol. 126, p. 103514, 2022.
- [44] R. Padilla, S. L. Netto, and E. A. Da Silva, "A survey on performance metrics for object-detection algorithms," in *2020 international conference on systems, signals and image processing (IWSSIP)*. IEEE, 2020, pp. 237–242.
- [45] M. Ahmed, K. A. Hashmi, A. Pagani, M. Liwicki, D. Stricker, and M. Z. Afzal, "Survey and performance analysis of deep learning based object detection in challenging environments," *Sensors*, vol. 21, no. 15, p. 5116, 2021.
- [46] Y. Xiao, Z. Tian, J. Yu, Y. Zhang, S. Liu, S. Du, and X. Lan, "A review of object detection based on deep learning," *Multimedia Tools and Applications*, vol. 79, pp. 23 729–23 791, 2020.
- [47] L. Du, R. Zhang, and X. Wang, "Overview of two-stage object detection algorithms," in *Journal of Physics: Conference Series*, vol. 1544. IOP Publishing, 2020, pp. 1–7.

- [48] X. Wu, D. Sahoo, and S. C. Hoi, "Recent advances in deep learning for object detection," *Neurocomputing*, vol. 396, pp. 39–64, 2020.
- [49] A. Dhillon and G. K. Verma, "Convolutional neural network: a review of models, methodologies and applications to object detection," *Progress in Artificial Intelligence*, vol. 9, no. 2, pp. 85–112, 2020.
- [50] M. Everingham, L. Van Gool, C. K. Williams, J. Winn, and A. Zisserman, "The pascal visual object classes (voc) challenge," *International journal of computer vision*, vol. 88, pp. 303–338, 2010.
- [51] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick, "Microsoft coco: Common objects in context," in *Computer Vision–ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6–12, 2014, Proceedings, Part V 13*. Springer, 2014, pp. 740–755.
- [52] Y. Chen, B. Zheng, Z. Zhang, Q. Wang, C. Shen, and Q. Zhang, "Deep learning on mobile and embedded devices: State-of-the-art, challenges, and future directions," *ACM Computing Surveys (CSUR)*, vol. 53, no. 4, pp. 1–37, 2020.
- [53] T. S. Ajani, A. L. Imoize, and A. A. Atayero, "An overview of machine learning within embedded and mobile devices—optimizations and applications," *Sensors*, vol. 21, no. 13, p. 4412, 2021.
- [54] R. Mishra, H. P. Gupta, and T. Dutta, "A survey on deep neural network compression: Challenges, overview, and solutions," *arXiv preprint arXiv:2010.03954*, 2020.
- [55] M. Everingham and J. Winn, "The pascal visual object classes challenge 2012 (voc2012) development kit," *Pattern Anal. Stat. Model. Comput. Learn., Tech. Rep.*, vol. 2007, no. 1-45, p. 5, 2012.
- [56] S. Kim, S. Na, B. Y. Kong, J. Choi, and I.-C. Park, "Real-time ssdlite object detection on FPGA," *IEEE Transactions on Very Large Scale Integration (VLSI) Systems*, vol. 29, no. 6, pp. 1192–1205, 2021.
- [57] H. Fan, S. Liu, M. Ferianc, H.-C. Ng, Z. Que, S. Liu, X. Niu, and W. Luk, "A real-time object detection accelerator with compressed ssdlite on FPGA," in *2018 International conference on field-programmable technology (FPT)*. IEEE, 2018, pp. 14–21.
- [58] J. Zhang, L. Cheng, C. Li, Y. Li, G. He, N. Xu, and Y. Lian, "A low-latency FPGA implementation for real-time object detection," in *2021 IEEE international symposium on circuits and systems (ISCAS)*. IEEE, 2021, pp. 1–5.
- [59] J. C. Nascimento and J. S. Marques, "Performance evaluation of object detection algorithms for video surveillance," *IEEE Transactions on Multimedia*, vol. 8, no. 4, pp. 761–774, 2006.
- [60] S. Jha, C. Seo, E. Yang, and G. P. Joshi, "Real-time object detection and tracking system for video surveillance system," *Multimedia Tools and Applications*, vol. 80, no. 3, pp. 3981–3996, 2021.
- [61] D. Castells-Rufas, V. Ngo, J. Borrego-Carazo, M. Codina, C. Sanchez, D. Gil, and J. Carrabina, "A survey of FPGA-based vision systems for autonomous cars," *IEEE Access*, vol. 10, pp. 132 525–132 563, 2022.
- [62] A. K. Tehrani and H. Rivaz, "Displacement estimation in ultrasound elastography using pyramidal convolutional neural network," *IEEE Transactions on ultrasonics, ferroelectrics, and frequency control*, vol. 67, no. 12, pp. 2629–2639, 2020.
- [63] K. Kang, H. Li, J. Yan, X. Zeng, B. Yang, T. Xiao, C. Zhang, Z. Wang, R. Wang, X. Wang *et al.*, "T-cnn: Tubelets with convolutional neural networks for object detection from videos," *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 28, no. 10, pp. 2896–2907, 2017.
- [64] B. Hariharan, P. Arbeláez, R. Girshick, and J. Malik, "Simultaneous detection and segmentation," in *Computer Vision–ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6–12, 2014, Proceedings, Part VII 13*. Springer, 2014, pp. 297–312.

- [65] J. Dai, K. He, and J. Sun, "Instance-aware semantic segmentation via multi-task network cascades," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 3150–3158.
- [66] K. He, G. Gkioxari, P. Dollár, and R. Girshick, "Mask r-cnn," in *Proceedings of the IEEE international conference on computer vision*, 2017, pp. 2961–2969.
- [67] A. Karpathy and L. Fei-Fei, "Deep visual-semantic alignments for generating image descriptions," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2015, pp. 3128–3137.
- [68] Q. Wu, C. Shen, P. Wang, A. Dick, and A. Van Den Hengel, "Image captioning and visual question answering based on attributes and external knowledge," *IEEE transactions on pattern analysis and machine intelligence*, vol. 40, no. 6, pp. 1367–1381, 2017.
- [69] P. Viola and M. Jones, "Rapid object detection using a boosted cascade of simple features," in *Proceedings of the 2001 IEEE Computer Society Conference on Computer Vision and Pattern Recognition. CVPR 2001*, vol. 1, 2001, pp. 1–9.
- [70] P. Irgens, C. Bader, T. Lé, D. Saxena, and C. Ababei, "An efficient and cost-effective FPGA based implementation of the viola-jones face detection algorithm," *HardwareX*, vol. 1, pp. 68–75, 2017.
- [71] F. N. Sitorus, Y. Manihuruk, and G. F. Panggabean, "Implementation of viola-jones algorithm for object detection using FPGA," in *2019 International Conference of Computer Science and Information Technology (ICoSNIKOM)*. IEEE, 2019, pp. 1–6.
- [72] P. Felzenszwalb, D. McAllester, and D. Ramanan, "A discriminatively trained, multiscale, deformable part model," in *2008 IEEE conference on computer vision and pattern recognition*. Ieee, 2008, pp. 1–8.
- [73] D. Tasson, A. Montagnini, R. Marzotto, M. Farenzena, and M. Cristani, "FPGA-based pedestrian detection under strong distortions," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, 2015, pp. 66–71.
- [74] R. Girshick, J. Donahue, T. Darrell, and J. Malik, "Rich feature hierarchies for accurate object detection and semantic segmentation," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2014, pp. 580–587.
- [75] J. An, D. Zhang, K. Xu, and D. Wang, "An opencl-based FPGA accelerator for faster r-cnn," *Entropy*, vol. 24, no. 10, p. 1346, 2022.
- [76] Z. Li and J. Wang, "An improved algorithm for deep learning yolo network based on xilinx zynq FPGA," in *2020 International Conference on Culture-oriented Science & Technology (ICCST)*. IEEE, 2020, pp. 447–451.
- [77] A. Anupreetham, M. Ibrahim, M. Hall, A. Boutros, A. Kuzhively, A. Mohanty, E. Nurvitadhi, V. Betz, Y. Cao, and J.-S. Seo, "High throughput FPGA-based object detection via algorithm-hardware co-design," *ACM Transactions on Reconfigurable Technology and Systems*, vol. 17, no. 1, pp. 1–20, 2024.
- [78] L. Cai, C. Wang, and Y. Xu, "A real-time fpga accelerator based on winograd algorithm for underwater object detection," *Electronics*, vol. 10, no. 23, p. 2889, 2021.
- [79] C.-Y. Fu, W. Liu, A. Ranga, A. Tyagi, and A. C. Berg, "Dssd: Deconvolutional single shot detector," *arXiv preprint arXiv:1701.06659*, 2017.
- [80] J. Redmon and A. Farhadi, "Yolo9000: better, faster, stronger," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 7263–7271.

- [81] H. Nakahara, M. Shimoda, and S. Sato, "A demonstration of fpga-based you only look once version2 (yolov2)," in *2018 28th International Conference on Field Programmable Logic and Applications (FPL)*. IEEE, 2018, pp. 457–4571.
- [82] Z. Wang, K. Xu, S. Wu, L. Liu, L. Liu, and D. Wang, "Sparse-yolo: Hardware/software co-design of an FPGA accelerator for yolov2," *IEEE Access*, vol. 8, pp. 116 569–116 585, 2020.
- [83] S. Li, Y. Luo, K. Sun, N. Yadav, and K. K. Choi, "A novel FPGA accelerator design for real-time and ultra-low power deep convolutional neural networks compared with titan x gpu," *IEEE Access*, vol. 8, pp. 105 455–105 471, 2020.
- [84] J. W. Yap, Z. bin Mohd Yussof, S. I. bin Salim, and K. C. Lim, "Fixed point implementation of tiny-yolo-v2 using OpenCL on FPGA," *International Journal of Advanced Computer Science and Applications*, vol. 9, no. 10, 2018.
- [85] K. Xu, X. Wang, X. Liu, C. Cao, H. Li, H. Peng, and D. Wang, "A dedicated hardware accelerator for real-time acceleration of yolov2," *Journal of Real-Time Image Processing*, vol. 18, pp. 481–492, 2021.
- [86] T.-Y. Lin, P. Goyal, R. Girshick, K. He, and P. Dollár, "Focal loss for dense object detection," in *Proceedings of the IEEE international conference on computer vision*, 2017, pp. 2980–2988.
- [87] H. Law and J. Deng, "Cornernet: Detecting objects as paired keypoints," in *Proceedings of the European conference on computer vision (ECCV)*, 2018, pp. 734–750.
- [88] J. Redmon and A. Farhadi, "Yolov3: An incremental improvement," *arXiv preprint arXiv:1804.02767*, 2018.
- [89] J. Wang and S. Gu, "FPGA implementation of object detection accelerator based on vitis-ai," in *2021 11th International Conference on Information Science and Technology (ICIST)*. IEEE, 2021, pp. 571–577.
- [90] Z. Yu and C.-S. Bouganis, "A parameterisable FPGA-tailored architecture for yolov3-tiny," in *Applied Reconfigurable Computing. Architectures, Tools, and Applications: 16th International Symposium, ARC 2020, Toledo, Spain, April 1–3, 2020, Proceedings 16*. Springer, 2020, pp. 330–344.
- [91] M. Kim, K. Oh, Y. Cho, H. Seo, X. T. Nguyen, and H.-J. Lee, "A low-latency fpga accelerator for yolov3-tiny with flexible layerwise mapping and dataflow," *IEEE Transactions on Circuits and Systems I: Regular Papers*, 2023.
- [92] D. Pestana, P. R. Miranda, J. D. Lopes, R. P. Duarte, M. P. Véstias, H. C. Neto, and J. T. De Sousa, "A full featured configurable accelerator for object detection with yolo," *IEEE Access*, vol. 9, pp. 75 864–75 877, 2021.
- [93] D. Zhang, A. Wang, R. Mo, and D. Wang, "End-to-end acceleration of the yolo object detection framework on FPGA-only devices," *Neural Computing and Applications*, vol. 36, no. 3, pp. 1067–1089, 2024.
- [94] K. Duan, S. Bai, L. Xie, H. Qi, Q. Huang, and Q. Tian, "Centernet: Keypoint triplets for object detection," in *Proceedings of the IEEE/CVF international conference on computer vision*, 2019, pp. 6569–6578.
- [95] R. A. Solovyev, D. V. Telpukhov, I. I. Romanova, A. G. Kustov, and I. A. Mkrtchan, "Real-time object detection with FPGA using centernet," in *2021 IEEE Conference of Russian Young Researchers in Electrical and Electronic Engineering (ElConRus)*. IEEE, 2021, pp. 2029–2034.
- [96] M. Tan, R. Pang, and Q. V. Le, "Efficientdet: Scalable and efficient object detection," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2020, pp. 10 781–10 790.

- [97] P. Babu and E. Parthasarathy, "Hardware acceleration for object detection using yolov4 algorithm on xilinx zynq platform," *Journal of Real-Time Image Processing*, vol. 19, no. 5, pp. 931–940, 2022.
- [98] C.-Y. Wang, A. Bochkovskiy, and H.-Y. M. Liao, "Scaled-yolov4: Scaling cross stage partial network," in *Proceedings of the IEEE/cvf conference on computer vision and pattern recognition*, 2021, pp. 13 029–13 038.
- [99] A. Hosseiny and H. Jahanirad, "Hardware acceleration of yolov7-tiny using high-level synthesis tools," *Journal of Real-Time Image Processing*, vol. 20, no. 4, p. 75, 2023.
- [100] Y. Li, N. Miao, L. Ma, F. Shuang, and X. Huang, "Transformer for object detection: Review and benchmark," *Eng. Appl. Artif. Intell.*, vol. 126, no. PC, feb 2024. [Online]. Available: <https://doi.org/10.1016/j.engappai.2023.107021>
- [101] T. Senoo, R. Kayanoma, A. Jinguji, and H. Nakahara, "A light-weight vision transformer toward near memory computation on an FPGA," in *International Symposium on Applied Reconfigurable Computing*. Springer, 2023, pp. 338–353.
- [102] N. Dalal and B. Triggs, "Histograms of oriented gradients for human detection," in *2005 IEEE computer society conference on computer vision and pattern recognition (CVPR'05)*, vol. 1. Ieee, 2005, pp. 886–893.
- [103] A. Suleiman and V. Sze, "An energy-efficient hardware implementation of hog-based object detection at 1080hd 60 fps with multi-scale support," *Journal of Signal Processing Systems*, vol. 84, pp. 325–337, 2016.
- [104] D. G. Lowe, "Distinctive image features from scale-invariant keypoints," *International journal of computer vision*, vol. 60, pp. 91–110, 2004.
- [105] S. Belongie, J. Malik, and J. Puzicha, "Matching shapes," in *Proceedings Eighth IEEE International Conference on Computer Vision. ICCV 2001*, vol. 1. IEEE, 2001, pp. 454–461.
- [106] P. F. Felzenszwalb, R. B. Girshick, D. McAllester, and D. Ramanan, "Object detection with discriminatively trained part-based models," *IEEE transactions on pattern analysis and machine intelligence*, vol. 32, no. 9, pp. 1627–1645, 2009.
- [107] R. Girshick, P. Felzenszwalb, and D. McAllester, "Object detection with grammar models," *Advances in neural information processing systems*, vol. 24, 2011.
- [108] R. B. Girshick, "From rigid templates to grammars: object detection with structured models," Ph.D. dissertation, University of Chicago, USA, 2012, aAI3513455.
- [109] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "Imagenet classification with deep convolutional neural networks," *Advances in neural information processing systems*, vol. 25, 2012.
- [110] Z.-Q. Zhao, P. Zheng, S.-t. Xu, and X. Wu, "Object detection with deep learning: A review," *IEEE transactions on neural networks and learning systems*, vol. 30, no. 11, pp. 3212–3232, 2019.
- [111] E. Arkin, N. Yadikar, X. Xu, A. Aysa, and K. Ubul, "A survey: object detection methods from cnn to transformer," *Multimedia Tools and Applications*, vol. 82, no. 14, pp. 21 353–21 383, 2023.
- [112] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," *arXiv preprint arXiv:1409.1556*, 2014.
- [113] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich, "Going deeper with convolutions," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2015, pp. 1–9.
- [114] M. Tan and Q. Le, "Efficientnet: Rethinking model scaling for convolutional neural networks," in *International conference on machine learning*. PMLR, 2019, pp. 6105–6114.

- [115] G. Huang, Z. Liu, L. Van Der Maaten, and K. Q. Weinberger, "Densely connected convolutional networks," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 4700–4708.
- [116] Z. Liu, Y. Lin, Y. Cao, H. Hu, Y. Wei, Z. Zhang, S. Lin, and B. Guo, "Swin transformer: Hierarchical vision transformer using shifted windows," in *Proceedings of the IEEE/CVF international conference on computer vision*, 2021, pp. 10012–10022.
- [117] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly *et al.*, "An image is worth 16x16 words: Transformers for image recognition at scale," *arXiv preprint arXiv:2010.11929*, 2020.
- [118] M. Carranza-García, J. Torres-Mateo, P. Lara-Benítez, and J. García-Gutiérrez, "On the performance of one-stage and two-stage object detectors in autonomous vehicles using camera data," *Remote Sensing*, vol. 13, no. 1, p. 89, 2020.
- [119] K. He, X. Zhang, S. Ren, and J. Sun, "Spatial pyramid pooling in deep convolutional networks for visual recognition," *IEEE transactions on pattern analysis and machine intelligence*, vol. 37, no. 9, pp. 1904–1916, 2015.
- [120] J. Dai, Y. Li, K. He, and J. Sun, "R-fcn: Object detection via region-based fully convolutional networks," *Advances in neural information processing systems*, vol. 29, 2016.
- [121] C.-Y. Wang, I.-H. Yeh, and H.-Y. M. Liao, "Yolov9: Learning what you want to learn using programmable gradient information," *arXiv preprint arXiv:2402.13616*, 2024.
- [122] A. Wang, H. Chen, L. Liu, K. Chen, Z. Lin, J. Han, and G. Ding, "Yolov10: Real-time end-to-end object detection," *arXiv preprint arXiv:2405.14458*, 2024.
- [123] "Darknet: Open source neural networks in c," <https://pjreddie.com/darknet/>, (Accessed on 04/05/2024).
- [124] S. Ioffe and C. Szegedy, "Batch normalization: Accelerating deep network training by reducing internal covariate shift," in *International conference on machine learning*. pmlr, 2015, pp. 448–456.
- [125] C.-Y. Wang, H.-Y. M. Liao, Y.-H. Wu, P.-Y. Chen, J.-W. Hsieh, and I.-H. Yeh, "Cspnet: A new backbone that can enhance learning capability of cnn," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition workshops*, 2020, pp. 390–391.
- [126] A. Neubeck and L. Van Gool, "Efficient non-maximum suppression," in *18th international conference on pattern recognition (ICPR'06)*, vol. 3. IEEE, 2006, pp. 850–855.
- [127] F. N. Iandola, S. Han, M. W. Moskewicz, K. Ashraf, W. J. Dally, and K. Keutzer, "Squeezenet: Alexnet-level accuracy with 50x fewer parameters and < 0.5 mb model size," *arXiv preprint arXiv:1602.07360*, 2016.
- [128] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 770–778.
- [129] T.-Y. Lin, P. Dollár, R. Girshick, K. He, B. Hariharan, and S. Belongie, "Feature pyramid networks for object detection," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 2117–2125.
- [130] S. Khan, M. Naseer, M. Hayat, S. W. Zamir, F. S. Khan, and M. Shah, "Transformers in vision: A survey," *ACM computing surveys (CSUR)*, vol. 54, no. 10s, pp. 1–41, 2022.
- [131] T. Wolf, L. Debut, V. Sanh, J. Chaumond, C. Delangue, A. Moi, P. Cistac, T. Rault, R. Louf, M. Funtowicz *et al.*, "Transformers: State-of-the-art natural language processing," in *Proceedings of the 2020 conference on empirical methods in natural language processing: system demonstrations*, 2020, pp. 38–45.

- [132] J. Beal, E. Kim, E. Tzeng, D. H. Park, A. Zhai, and D. Kislyuk, "Toward transformer-based object detection," *arXiv preprint arXiv:2012.09958*, 2020.
- [133] W. F. Hendria, Q. T. Phan, F. Adzaka, and C. Jeong, "Combining transformer and cnn for object detection in uav imagery," *ICT Express*, vol. 9, no. 2, pp. 258–263, 2023.
- [134] H. Wang, Y. Zhu, B. Green, H. Adam, A. Yuille, and L.-C. Chen, "Axial-deeplab: Stand-alone axial-attention for panoptic segmentation," in *European conference on computer vision*. Springer, 2020, pp. 108–126.
- [135] P. Ramachandran, N. Parmar, A. Vaswani, I. Bello, A. Levskaya, and J. Shlens, "Stand-alone self-attention in vision models," *Advances in neural information processing systems*, vol. 32, 2019.
- [136] X. Wang, R. Girshick, A. Gupta, and K. He, "Non-local neural networks," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 7794–7803.
- [137] T. Shehzadi, K. A. Hashmi, D. Stricker, and M. Z. Afzal, "2d object detection with transformers: a review," *arXiv preprint arXiv:2306.04670*, 2023.
- [138] Z. Dai, B. Cai, Y. Lin, and J. Chen, "Up-detr: Unsupervised pre-training for object detection with transformers," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2021, pp. 1601–1610.
- [139] Z. Yao, J. Ai, B. Li, and C. Zhang, "Efficient detr: improving end-to-end object detector with dense prior," *arXiv preprint arXiv:2104.01318*, 2021.
- [140] H. Song, D. Sun, S. Chun, V. Jampani, D. Han, B. Heo, W. Kim, and M.-H. Yang, "VidT: An efficient and effective fully transformer-based object detector," *arXiv preprint arXiv:2110.03921*, 2021.
- [141] J. A. Stankovic, "Misconceptions about real-time computing: A serious problem for next-generation systems," *Computer*, vol. 21, no. 10, pp. 10–19, 1988.
- [142] K. G. Shin and P. Ramanathan, "Real-time computing: A new discipline of computer science and engineering," *Proceedings of the IEEE*, vol. 82, no. 1, pp. 6–24, 1994.
- [143] J. A. Stankovic, "Real-time and embedded systems," *ACM Computing Surveys (CSUR)*, vol. 28, no. 1, pp. 205–208, 1996.
- [144] S. B. N. Meghanathan, "A survey of contemporary real-time operating systems," *Informatica*, vol. 29, no. 2, 2005.
- [145] B. P. Douglass, *Doing hard time: developing real-time systems with UML, objects, frameworks, and patterns*. Addison-Wesley Professional, 1999, vol. 1.
- [146] B. P. Douglass, *Real-time design patterns: robust scalable architecture for real-time systems*. Addison-Wesley Professional, 2003.
- [147] Q. Li and C. Yao, *Real-time concepts for embedded systems*. CRC press, 2003.
- [148] P. A. Laplante *et al.*, *Real-time systems design and analysis*. Wiley New York, 2004.
- [149] J. Huang, V. Rathod, C. Sun, M. Zhu, A. Korattikara, A. Fathi, I. Fischer, Z. Wojna, Y. Song, S. Guadarrama *et al.*, "Speed/accuracy trade-offs for modern convolutional object detectors," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 7310–7311.
- [150] H. Mao, S. Yao, T. Tang, B. Li, J. Yao, and Y. Wang, "Towards real-time object detection on embedded systems," *IEEE Transactions on Emerging Topics in Computing*, vol. 6, no. 3, pp. 417–431, 2016.
- [151] L. Baischer, M. Wess, and N. TaheriNejad, "Learning on hardware: A tutorial on neural network accelerators and co-processors," *arXiv preprint arXiv:2104.09252*, 2021.

- [152] A. Sateesan, S. Sinha, S. KG, and A. Vinod, "A survey of algorithmic and hardware optimization techniques for vision convolutional neural networks on FPGAs," *Neural Processing Letters*, vol. 53, no. 3, pp. 2331–2377, 2021.
- [153] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein *et al.*, "Imagenet large scale visual recognition challenge," *International journal of computer vision*, vol. 115, pp. 211–252, 2015.
- [154] "The latest in machine learning | papers with code," <https://paperswithcode.com/>, (Accessed on 04/07/2024).
- [155] I. Rodriguez-Conde, C. Campos, and F. Fdez-Riverola, "On-device object detection for more efficient and privacy-compliant visual perception in context-aware systems," *Applied Sciences*, vol. 11, no. 19, p. 9173, 2021.
- [156] V. Sze, Y.-H. Chen, T.-J. Yang, and J. S. Emer, *Efficient processing of deep neural networks*. Springer, 2020.
- [157] J. Lee, L. Mukhanov, A. S. Molahosseini, U. Minhas, Y. Hua, J. Martinez del Rincon, K. Dichev, C.-H. Hong, and H. Vandierendonck, "Resource-efficient convolutional networks: A survey on model-, arithmetic-, and implementation-level techniques," *ACM Computing Surveys*, vol. 55, no. 13s, pp. 1–36, 2023.
- [158] I. Kuon, R. Tessier, J. Rose *et al.*, "FPGA architecture: Survey and challenges," *Foundations and Trends® in Electronic Design Automation*, vol. 2, no. 2, pp. 135–253, 2008.
- [159] U. Farooq, Z. Marrakchi, H. Mehrez, U. Farooq, Z. Marrakchi, and H. Mehrez, "FPGA architectures: An overview," *Tree-Based Heterogeneous FPGA Architectures: Application Specific Exploration and Optimization*, pp. 7–48, 2012.
- [160] U. Meyer-Baese, *Digital Signal Processing with Field Programmable Gate Arrays*. Springer, 2014.
- [161] S. Coric, M. Leeser, E. Miller, and M. Trepanier, "Parallel-beam backprojection: an FPGA implementation optimized for medical imaging," in *Proceedings of the 2002 ACM/SIGDA tenth international symposium on Field-programmable gate arrays*, 2002, pp. 217–226.
- [162] R. Mehra, G. Saini, and S. Singh, "FPGA based high speed bch encoder for wireless communication applications," in *2011 International Conference on Communication Systems and Network Technologies*. IEEE, 2011, pp. 576–579.
- [163] Z. A. O. Nasri Sulaiman, M. Marhaban, and M. Hamidon, "Design and implementation of FPGA-based systems-a review," *Australian Journal of Basic and Applied Sciences*, vol. 3, no. 4, pp. 3575–3596, 2009.
- [164] T. El-Ghazawi, D. Bennett, D. Poznanovic, A. Cattle, K. Underwood, R. Pennington, D. Buell, A. George, and V. Kindratenko, "Is high-performance reconfigurable computing the next supercomputing paradigm?" in *Proceedings of the 2006 ACM/IEEE conference on Supercomputing*, 2006, pp. 71–es.
- [165] G. Lacey, G. W. Taylor, and S. Areibi, "Deep learning on FPGAs: Past, present, and future," *arXiv preprint arXiv:1602.04283*, 2016.
- [166] J. Cong, Z. Fang, M. Lo, H. Wang, J. Xu, and S. Zhang, "Understanding performance differences of FPGAs and gpus," in *2018 IEEE 26th Annual International Symposium on Field-Programmable Custom Computing Machines (FCCM)*. IEEE, 2018, pp. 93–96.
- [167] W. Wayne *et al.*, *FPGA-based system design*. Pearson Education India, 2004.

- [168] A. C. Sankaranarayanan, A. Veeraraghavan, and R. Chellappa, "Object detection, tracking and recognition for multiple smart cameras," *Proceedings of the IEEE*, vol. 96, no. 10, pp. 1606–1624, 2008.
- [169] R. Szeliski, *Computer vision: algorithms and applications*. Springer Nature, 2022.
- [170] V. Wiley and T. Lucas, "Computer vision and image processing: a paper review," *International Journal of Artificial Intelligence Research*, vol. 2, no. 1, pp. 29–36, 2018.
- [171] G. Gallego, T. Delbrück, G. Orchard, C. Bartolozzi, B. Taba, A. Censi, S. Leutenegger, A. J. Davison, J. Conrath, K. Daniilidis *et al.*, "Event-based vision: A survey," *IEEE transactions on pattern analysis and machine intelligence*, vol. 44, no. 1, pp. 154–180, 2020.
- [172] E. Yurtsever, J. Lambert, A. Carballo, and K. Takeda, "A survey of autonomous driving: Common practices and emerging technologies," *IEEE access*, vol. 8, pp. 58 443–58 469, 2020.
- [173] X. Zheng, Y. Liu, Y. Lu, T. Hua, T. Pan, W. Zhang, D. Tao, and L. Wang, "Deep learning for event-based vision: A comprehensive survey and benchmarks," *arXiv preprint arXiv:2302.08890*, 2023.
- [174] J. Beyerer, F. Puente León, C. Frese, J. Beyerer, F. Puente León, and C. Frese, "Preprocessing and image enhancement," *Machine Vision: Automated Visual Inspection: Theory, Practice and Applications*, pp. 465–519, 2016.
- [175] S. Krig and S. Krig, "Image pre-processing," *Computer Vision Metrics: Textbook Edition*, pp. 35–74, 2016.
- [176] M. Sen, I. Corretjer, F. Haim, S. Saha, S. S. Bhattacharyya, J. Schlessman, and W. Wolf, "Computer vision on FPGAs: Design methodology and its application to gesture recognition," in *2005 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR'05)-Workshops*. IEEE, 2005, pp. 133–133.
- [177] A. Gupta, "Current research opportunities for image processing and computer vision," *Computer Science*, vol. 20, pp. 387–410, 2019.
- [178] E. Elyan, P. Vuttipittayamongkol, P. Johnston, K. Martin, K. McPherson, C. Jayne, M. K. Sarker *et al.*, "Computer vision and machine learning for medical image analysis: recent advances, challenges, and way forward." *Artificial Intelligence Surgery*, vol. 2, 2022.
- [179] X. Feng, Y. Jiang, X. Yang, M. Du, and X. Li, "Computer vision algorithms and hardware implementations: A survey," pp. 309–320, 11 2019.
- [180] D. Bhowmik and K. Appiah, "Embedded vision systems: A review of the literature," in *Applied Reconfigurable Computing. Architectures, Tools, and Applications: 14th International Symposium, ARC 2018, Santorini, Greece, May 2-4, 2018, Proceedings 14*. Springer, 2018, pp. 204–216.
- [181] S. Jin, J. Cho, X. Dai Pham, K. M. Lee, S.-K. Park, M. Kim, and J. W. Jeon, "FPGA design and implementation of a real-time stereo vision system," *IEEE transactions on circuits and systems for video technology*, vol. 20, no. 1, pp. 15–26, 2009.
- [182] S. Mittal, S. Gupta, and S. Dasgupta, "FPGA: An efficient and promising platform for real-time image processing applications," in *National Conference On Research and Development In Hardware Systems (CSI-RDHS)*, 2008.
- [183] K. P. Seng, P. J. Lee, and L. M. Ang, "Embedded intelligence on FPGA: Survey, applications and challenges," *Electronics*, vol. 10, no. 8, p. 895, 2021.
- [184] A. Irwansyah, O. W. Ibraheem, J. Hagemeyer, M. Pormann, and U. Rueckert, "FPGA-based multi-robot tracking," *Journal of Parallel and Distributed Computing*, vol. 107, pp. 146–161, 2017.

- [185] “Gige vision,” <https://www.automate.org/vision/vision-standards/vision-standards-gige-vision>, (Accessed on 05/27/2024).
- [186] M. A. Altuncu, T. Guven, Y. Becerikli, and S. Sahin, “Real-time system implementation for image processing with hardware/software co-design on the xilinx zynq platform,” *International Journal of Information and Electronics Engineering*, vol. 5, no. 6, p. 473, 2015.
- [187] “Interface specifications for mobile products | mipi alliance,” <https://www.mipi.org/>, (Accessed on 06/25/2024).
- [188] J. Ahmad and A. Warren, “FPGA based deterministic latency image acquisition and processing system for automated driving systems,” in *2018 IEEE International Symposium on Circuits and Systems (ISCAS)*. IEEE, 2018, pp. 1–5.
- [189] K. Kim, S.-J. Jang, J. Park, E. Lee, and S.-S. Lee, “Lightweight and energy-efficient deep learning accelerator for real-time object detection on edge devices,” *Sensors*, vol. 23, no. 3, p. 1185, 2023.
- [190] C. Hao, X. Zhang, Y. Li, S. Huang, J. Xiong, K. Rupnow, W.-m. Hwu, and D. Chen, “FPGA/dnn co-design: An efficient design methodology for iot intelligence on the edge,” in *Proceedings of the 56th Annual Design Automation Conference 2019*, 2019, pp. 1–6.
- [191] Q. Huang, D. Wang, Z. Dong, Y. Gao, Y. Cai, T. Li, B. Wu, K. Keutzer, and J. Wawrzynek, “Codenet: Efficient deployment of input-adaptive object detection on embedded FPGAs,” in *The 2021 ACM/SIGDA International Symposium on Field-Programmable Gate Arrays*, 2021, pp. 206–216.
- [192] Z. Zhang, M. P. Mahmud, and A. Z. Kouzani, “Resource-constrained FPGA implementation of yolov2,” *Neural Computing and Applications*, vol. 34, no. 19, pp. 16 989–17 006, 2022.
- [193] “Amd together we advance_ai,” <https://www.amd.com/en.html>, (Accessed on 05/28/2024).
- [194] “Intel | data center solutions, iot, and pc innovation,” <https://www.intel.com/content/www/us/en/homepage.html>, (Accessed on 05/28/2024).
- [195] A. Xilinx, “Dpuczd8g for zynq ultrascale+ mpsocs. product guide,” 2022.
- [196] J. Yu, K. Guo, Y. Hu, X. Ning, J. Qiu, H. Mao, S. Yao, T. Tang, B. Li, Y. Wang *et al.*, “Real-time object detection towards high power efficiency,” in *2018 Design, Automation & Test in Europe Conference & Exhibition (DATE)*. IEEE, 2018, pp. 704–708.
- [197] K. Guo, L. Sui, J. Qiu, J. Yu, J. Wang, S. Yao, S. Han, Y. Wang, and H. Yang, “Angel-eye: A complete design flow for mapping cnn onto embedded FPGA,” *IEEE transactions on computer-aided design of integrated circuits and systems*, vol. 37, no. 1, pp. 35–47, 2017.
- [198] M. Reis, M. Véstias, and H. Neto, “Designing deep learning models on FPGA with multiple heterogeneous engines,” *ACM Transactions on Reconfigurable Technology and Systems*, vol. 17, no. 1, pp. 1–30, 2024.
- [199] A. Montgomerie-Corcoran, P. Toupas, Z. Yu, and C.-S. Bouganis, “Satay: a streaming architecture toolflow for accelerating yolo models on FPGA devices,” in *2023 International Conference on Field Programmable Technology (ICFPT)*. IEEE, 2023, pp. 179–187.
- [200] H. T. Kung and C. E. Leiserson, “Systolic arrays (for vlsi),” in *Sparse Matrix Proceedings 1978*, vol. 1. Society for industrial and applied mathematics Philadelphia, PA, USA, 1979, pp. 256–282.
- [201] E. Wang, J. J. Davis, R. Zhao, H.-C. Ng, X. Niu, W. Luk, P. Y. Cheung, and G. A. Constantinides, “Deep neural network approximation for custom hardware: Where we’ve been, where we’re going,” *ACM Computing Surveys (CSUR)*, vol. 52, no. 2, pp. 1–39, 2019.

- [202] R. Hameed, W. Qadeer, M. Wachs, O. Azizi, A. Solomatnikov, B. C. Lee, S. Richardson, C. Kozyrkis, and M. Horowitz, "Understanding sources of inefficiency in general-purpose chips," in *Proceedings of the 37th annual international symposium on Computer architecture*, 2010, pp. 37–47.
- [203] V. Sze, Y.-H. Chen, T.-J. Yang, and J. S. Emer, "Efficient processing of deep neural networks: A tutorial and survey," *Proceedings of the IEEE*, vol. 105, no. 12, pp. 2295–2329, 2017.
- [204] J. Qiu, J. Wang, S. Yao, K. Guo, B. Li, E. Zhou, J. Yu, T. Tang, N. Xu, S. Song *et al.*, "Going deeper with embedded FPGA platform for convolutional neural network," in *Proceedings of the 2016 ACM/SIGDA international symposium on field-programmable gate arrays*, 2016, pp. 26–35.
- [205] Y. Ma, Y. Cao, S. Vrudhula, and J.-s. Seo, "Optimizing the convolution operation to accelerate deep neural networks on FPGA," *IEEE Transactions on Very Large Scale Integration (VLSI) Systems*, vol. 26, no. 7, pp. 1354–1367, 2018.
- [206] X. Wei, C. H. Yu, P. Zhang, Y. Chen, Y. Wang, H. Hu, Y. Liang, and J. Cong, "Automated systolic array architecture synthesis for high throughput cnn inference on FPGAs," in *Proceedings of the 54th Annual Design Automation Conference 2017*, 2017, pp. 1–6.
- [207] C. Wang, W. Lou, L. Gong, L. Jin, L. Tan, Y. Hu, X. Li, and X. Zhou, "Reconfigurable hardware accelerators: Opportunities, trends, and challenges," *arXiv preprint arXiv:1712.04771*, 2017.
- [208] S. F. Anderson, J. G. Earle, R. E. Goldschmidt, and D. M. Powers, "The ibm system/360 model 91: Floating-point execution unit," *IBM Journal of research and development*, vol. 11, no. 1, pp. 34–53, 1967.
- [209] C. Zhang, P. Li, G. Sun, Y. Guan, B. Xiao, and J. Cong, "Optimizing FPGA-based accelerator design for deep convolutional neural networks," in *Proceedings of the 2015 ACM/SIGDA international symposium on field-programmable gate arrays*, 2015, pp. 161–170.
- [210] B. Ramhorst, V. Lončar, and G. A. Constantinides, "FPGA resource-aware structured pruning for real-time neural networks," in *2023 International Conference on Field Programmable Technology (ICFPT)*. IEEE, 2023, pp. 282–283.
- [211] J. Plochaet and T. Goedemé, "Hardware-aware pruning for FPGA deep learning accelerators," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 4481–4489.
- [212] X. Sui, Q. Lv, L. Zhi, B. Zhu, Y. Yang, Y. Zhang, and Z. Tan, "A hardware-friendly high-precision cnn pruning method and its FPGA implementation," *Sensors*, vol. 23, no. 2, p. 824, 2023.
- [213] G. Dinelli, G. Meoni, E. Rapuano, T. Pacini, and L. Fanucci, "Mem-opt: A scheduling and data re-use system to optimize on-chip memory usage for cnns on-board FPGAs," *IEEE Journal on Emerging and Selected Topics in Circuits and Systems*, vol. 10, no. 3, pp. 335–347, 2020.
- [214] X.-Q. Nguyen and C. Pham-Quoc, "An FPGA-based convolution ip core for deep neural networks acceleration," *REV Journal on Electronics and Communications*, vol. 12, no. 1-2, 2022.
- [215] S. Kala, J. Mathew, B. R. Jose, and S. Nalesh, "Uniwig: Unified winograd-gemm architecture for accelerating cnn on FPGAs," in *2019 32nd International Conference on VLSI Design and 2019 18th International Conference on Embedded Systems (VLSID)*. IEEE, 2019, pp. 209–214.
- [216] X. Wang, C. Wang, J. Cao, L. Gong, and X. Zhou, "Winonn: Optimizing FPGA-based convolutional neural network accelerators using sparse winograd algorithm," *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems*, vol. 39, no. 11, pp. 4290–4302, 2020.
- [217] C. Bao, T. Xie, W. Feng, L. Chang, and C. Yu, "A power-efficient optimizing framework FPGA accelerator based on winograd for yolo," *Ieee Access*, vol. 8, pp. 94 307–94 317, 2020.

- [218] Y. Wei, R. Wang, Y. Wang, F. Zhou, and N. Gou, "Optimizing FPGA-Based Target Detection: A Hybrid Winograd-GEMM Accelerator with Enhanced Resource Efficiency and Throughput," *Research Square pre-print*, 2024.
- [219] C. Buciluă, R. Caruana, and A. Niculescu-Mizil, "Model compression," in *Proceedings of the 12th ACM SIGKDD international conference on Knowledge discovery and data mining*, 2006, pp. 535–541.
- [220] G. Hinton, O. Vinyals, and J. Dean, "Distilling the knowledge in a neural network," *arXiv preprint arXiv:1503.02531*, 2015.
- [221] A. Polino, R. Pascanu, and D. Alistarh, "Model compression via distillation and quantization," *arXiv preprint arXiv:1802.05668*, 2018.
- [222] C. Zhang, G. Sun, Z. Fang, P. Zhou, and J. Cong, "Caffeine: Towards uniformed representation and acceleration for deep convolutional neural networks," in *Proceedings of the ACM Turing Award Celebration Conference-China 2023*, 2023, pp. 47–48.
- [223] U. Aydonat, S. O'Connell, D. Capalija, A. C. Ling, and G. R. Chiu, "An opencl™ deep learning accelerator on arria 10," in *Proceedings of the 2017 ACM/SIGDA International Symposium on Field-Programmable Gate Arrays*, 2017, pp. 55–64.
- [224] K. T. Chitty-Venkata and A. K. Somani, "Neural architecture search survey: A hardware perspective," *ACM Computing Surveys*, vol. 55, no. 4, pp. 1–36, 2022.
- [225] H. Benmeziane, K. E. Maghraoui, H. Ouarnoughi, S. Niar, M. Wistuba, and N. Wang, "A comprehensive survey on hardware-aware neural architecture search," *arXiv preprint arXiv:2101.09336*, 2021.
- [226] W. Jiang, X. Zhang, E. H.-M. Sha, L. Yang, Q. Zhuge, Y. Shi, and J. Hu, "Accuracy vs. efficiency: Achieving both through FPGA-implementation aware neural architecture search," in *Proceedings of the 56th Annual Design Automation Conference 2019*, 2019, pp. 1–6.
- [227] L. Lu, Y. Liang, Q. Xiao, and S. Yan, "Evaluating fast algorithms for convolutional neural networks on FPGAs," in *2017 IEEE 25th annual international symposium on field-programmable custom computing machines (FCCM)*. IEEE, 2017, pp. 101–108.
- [228] B. Liu, D. Zou, L. Feng, S. Feng, P. Fu, and J. Li, "An FPGA-based cnn accelerator integrating depthwise separable convolution," *Electronics*, vol. 8, no. 3, p. 281, 2019.
- [229] H. Zhang, J. Jiang, Y. Fu, and Y. Chang, "Yolov3-tiny object detection soc based on FPGA platform," in *2021 6th International Conference on Integrated Circuits and Microsystems (ICICM)*. IEEE, 2021, pp. 291–294.
- [230] S. Han, J. Pool, J. Tran, and W. Dally, "Learning both weights and connections for efficient neural network," *Advances in neural information processing systems*, vol. 28, 2015.
- [231] S. Han, H. Mao, and W. J. Dally, "Deep compression: Compressing deep neural networks with pruning, trained quantization and huffman coding," *arXiv preprint arXiv:1510.00149*, 2015.
- [232] Y. Guo, A. Yao, and Y. Chen, "Dynamic network surgery for efficient dnns," *Advances in neural information processing systems*, vol. 29, 2016.
- [233] Y. LeCun, J. Denker, and S. Solla, "Optimal brain damage," *Advances in neural information processing systems*, vol. 2, 1989.
- [234] S. Park, J. Lee, S. Mo, and J. Shin, "Lookahead: A far-sighted alternative of magnitude-based pruning," *arXiv preprint arXiv:2002.04809*, 2020.

- [235] X. Xiao, Z. Wang, and S. Rajasekaran, "Autoprune: Automatic network pruning by regularizing auxiliary parameters," *Advances in neural information processing systems*, vol. 32, 2019.
- [236] T. Gale, E. Elsen, and S. Hooker, "The state of sparsity in deep neural networks.(2019)," *arXiv preprint cs.LG/1902.09574*, 2019.
- [237] Z. Huang and N. Wang, "Data-driven sparse structure selection for deep neural networks," in *Proceedings of the European conference on computer vision (ECCV)*, 2018, pp. 304–320.
- [238] J.-H. Luo, J. Wu, and W. Lin, "Thinet: A filter level pruning method for deep neural network compression," in *Proceedings of the IEEE international conference on computer vision*, 2017, pp. 5058–5066.
- [239] D. Blalock, J. J. Gonzalez Ortiz, J. Frankle, and J. Gutttag, "What is the state of neural network pruning?" *Proceedings of machine learning and systems*, vol. 2, pp. 129–146, 2020.
- [240] F. Faghri, I. Tabrizian, I. Markov, D. Alistarh, D. M. Roy, and A. Ramezani-Kebrya, "Adaptive gradient quantization for data-parallel sgd," *Advances in neural information processing systems*, vol. 33, pp. 3174–3185, 2020.
- [241] B. Chmiel, L. Ben-Uri, M. Shkolnik, E. Hoffer, R. Banner, and D. Soudry, "Neural gradients are near-lognormal: improved quantized and sparse training," *arXiv preprint arXiv:2006.08173*, 2020.
- [242] Y.-H. Chen, T. Krishna, J. S. 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, 2016.
- [243] G. Zhou, J. Zhou, and H. Lin, "Research on nvidia deep learning accelerator," in *2018 12th IEEE International Conference on Anti-counterfeiting, Security, and Identification (ASID)*. IEEE, 2018, pp. 192–195.
- [244] S. Jung, C. Son, S. Lee, J. Son, J.-J. Han, Y. Kwak, S. J. Hwang, and C. Choi, "Learning to quantize deep networks by optimizing quantization intervals with task loss," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2019, pp. 4350–4359.
- [245] E. Park, J. Ahn, and S. Yoo, "Weighted-entropy-based quantization for deep neural networks," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 5456–5464.
- [246] E. Park, S. Yoo, and P. Vajda, "Value-aware quantization for training and inference of neural networks," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018, pp. 580–595.
- [247] F. Tung and G. Mori, "Clip-q: Deep network compression learning by in-parallel pruning-quantization," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 7873–7882.
- [248] P. Hu, X. Peng, H. Zhu, M. M. S. Aly, and J. Lin, "Opq: Compressing deep neural networks with one-shot pruning-quantization," in *Proceedings of the AAAI conference on artificial intelligence*, vol. 35, 2021, pp. 7780–7788.
- [249] H. Liu, S. Elkerdawy, N. Ray, and M. Elhoushi, "Layer importance estimation with imprinting for neural network quantization," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, pp. 2408–2417.
- [250] J. Fang, A. Shafiee, H. Abdel-Aziz, D. Thorsley, G. Georgiadis, and J. H. Hassoun, "Post-training piecewise linear quantization for deep neural networks," in *Computer Vision—ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part II 16*. Springer, 2020, pp. 69–86.

- [251] M. M. H. Shuvo, S. K. Islam, J. Cheng, and B. I. Morshed, "Efficient acceleration of deep learning inference on resource-constrained edge devices: A review," *Proceedings of the IEEE*, vol. 111, no. 1, pp. 42–91, 2022.
- [252] K. Arulkumaran, M. P. Deisenroth, M. Brundage, and A. A. Bharath, "A brief survey of deep reinforcement learning," *arXiv preprint arXiv:1708.05866*, 2017.
- [253] C. Wang and Z. Luo, "A review of the optimal design of neural networks based on FPGA," *Applied Sciences*, vol. 12, no. 21, p. 10771, 2022.
- [254] S. Moini, B. Alizadeh, M. Emad, and R. Ebrahimpour, "A resource-limited hardware accelerator for convolutional neural networks in embedded vision applications," *IEEE Transactions on Circuits and Systems II: Express Briefs*, vol. 64, no. 10, pp. 1217–1221, 2017.
- [255] A. Beric, J. van Meerbergen, G. de Haan, and R. Sethuraman, "Memory-centric video processing," *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 18, no. 4, pp. 439–452, 2008.
- [256] S. Winograd, *Arithmetic complexity of computations*. Siam, 1980, vol. 33.
- [257] H. J. Nussbaumer, *The Fast Fourier Transform*. Springer, 1982.
- [258] A. Lavin and S. Gray, "Fast algorithms for convolutional neural networks," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 4013–4021.
- [259] A. Nechi, L. Groth, S. Mulhem, F. Merchant, R. Buchty, and M. Berekovic, "FPGA-based deep learning inference accelerators: Where are we standing?" *ACM Transactions on Reconfigurable Technology and Systems*, vol. 16, no. 4, pp. 1–32, 2023.
- [260] Y. Ma, Y. Cao, S. Vrudhula, and J.-s. Seo, "Optimizing loop operation and dataflow in FPGA acceleration of deep convolutional neural networks," in *Proceedings of the 2017 ACM/SIGDA International Symposium on Field-Programmable Gate Arrays*, 2017, pp. 45–54.
- [261] Y. Ma, Y. Cao, S. Vrudhula, and J.-s. Seo, "An automatic rtl compiler for high-throughput FPGA implementation of diverse deep convolutional neural networks," in *2017 27th International Conference on Field Programmable Logic and Applications (FPL)*. IEEE, 2017, pp. 1–8.
- [262] C. Zhu, K. Huang, S. Yang, Z. Zhu, H. Zhang, and H. Shen, "An efficient hardware accelerator for structured sparse convolutional neural networks on FPGAs," *IEEE Transactions on Very Large Scale Integration (VLSI) Systems*, vol. 28, no. 9, pp. 1953–1965, 2020.
- [263] S. Mittal, "A survey of FPGA-based accelerators for convolutional neural networks," *Neural computing and applications*, vol. 32, no. 4, pp. 1109–1139, 2020.
- [264] J. Cheng, P.-s. Wang, G. Li, Q.-h. Hu, and H.-q. Lu, "Recent advances in efficient computation of deep convolutional neural networks," *Frontiers of Information Technology & Electronic Engineering*, vol. 19, pp. 64–77, 2018.
- [265] M. Alwani, H. Chen, M. Ferdman, and P. Milder, "Fused-layer cnn accelerators," in *2016 49th Annual IEEE/ACM International Symposium on Microarchitecture (MICRO)*. IEEE, 2016, pp. 1–12.
- [266] Y. Xing, S. Liang, L. Sui, Z. Zhang, J. Qiu, X. Jia, X. Liu, Y. Wang, Y. Shan, and Y. Wang, "Dnnvm: End-to-end compiler leveraging operation fusion on FPGA-based cnn accelerators," in *Proceedings of the 2019 ACM/SIGDA International Symposium on Field-Programmable Gate Arrays*, 2019, pp. 187–188.
- [267] S. Williams, A. Waterman, and D. Patterson, "Roofline: an insightful visual performance model for multicore architectures," *Communications of the ACM*, vol. 52, no. 4, pp. 65–76, 2009.
- [268] M. Siracusa, L. Di Tucci, M. Rabozzi, S. Williams, E. D. Sozzo, and M. D. Santambrogio, "A cad-based methodology to optimize hls code via the roofline model," in *Proceedings of the 39th International Conference on Computer-Aided Design*, 2020, pp. 1–9.

- [269] B. Da Silva, A. Braeken, E. H. D'Hollander, and A. Touhafi, "Performance modeling for FPGAs: extending the roofline model with high-level synthesis tools," *International Journal of Reconfigurable Computing*, vol. 2013, pp. 7–7, 2013.
- [270] P. Coussy, D. D. Gajski, M. Meredith, and A. Takach, "An introduction to high-level synthesis," *IEEE Design & Test of Computers*, vol. 26, no. 4, pp. 8–17, 2009.
- [271] E. Calore and S. F. Schifano, "Performance assessment of FPGAs as hpc accelerators using the FPGA empirical roofline," in *2021 31st International Conference on Field-Programmable Logic and Applications (FPL)*. IEEE, 2021, pp. 83–90.
- [272] G. Li, J. Zhang, M. Zhang, and H. Corporaal, "An efficient FPGA implementation for real-time and low-power uav object detection," in *2022 IEEE International Symposium on Circuits and Systems (ISCAS)*. IEEE, 2022, pp. 1387–1391.
- [273] "Pipelining," <https://www.intel.com/content/www/us/en/docs/oneapi-fpga-add-on/optimization-guide/2023-1/pipelining-001.html>, (Accessed on 03/29/2024).
- [274] X. Xu, X. Zhang, B. Yu, X. S. Hu, C. Rowen, J. Hu, and Y. Shi, "Dac-sdc low power object detection challenge for uav applications," *IEEE transactions on pattern analysis and machine intelligence*, vol. 43, no. 2, pp. 392–403, 2019.
- [275] F. Chollet, "Xception: Deep learning with depthwise separable convolutions," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 1251–1258.
- [276] L. Mei, P. Houshmand, V. Jain, S. Giraldo, and M. Verhelst, "Zigzag: Enlarging joint architecture-mapping design space exploration for dnn accelerators," *IEEE Transactions on Computers*, vol. 70, no. 8, pp. 1160–1174, 2021.
- [277] I. Hubara, M. Courbariaux, D. Soudry, R. El-Yaniv, and Y. Bengio, "Binarized neural networks," *Advances in neural information processing systems*, vol. 29, 2016.
- [278] H. Alemdar, V. Leroy, A. Prost-Boucle, and F. Pétrot, "Ternary neural networks for resource-efficient ai applications," in *2017 international joint conference on neural networks (IJCNN)*. IEEE, 2017, pp. 2547–2554.
- [279] L. Deng, P. Jiao, J. Pei, Z. Wu, and G. Li, "Gxnor-net: Training deep neural networks with ternary weights and activations without full-precision memory under a unified discretization framework," *Neural Networks*, vol. 100, pp. 49–58, 2018.
- [280] M. Ghasemzadeh, M. Samragh, and F. Koushanfar, "Rebnet: Residual binarized neural network," in *2018 IEEE 26th annual international symposium on field-programmable custom computing machines (FCCM)*. IEEE, 2018, pp. 57–64.
- [281] D. Przewlocka-Rus, S. S. Sarwar, H. E. Sumbul, Y. Li, and B. De Salvo, "Power-of-two quantization for low bitwidth and hardware compliant neural networks," *arXiv preprint arXiv:2203.05025*, 2022.
- [282] T. Xia, B. Zhao, J. Ma, G. Fu, W. Zhao, N. Zheng, and P. Ren, "An energy-and-area-efficient cnn accelerator for universal powers-of-two quantization," *IEEE Transactions on Circuits and Systems I: Regular Papers*, vol. 70, no. 3, pp. 1242–1255, 2022.
- [283] M. Courbariaux, I. Hubara, D. Soudry, R. El-Yaniv, and Y. Bengio, "Binarized neural networks: Training deep neural networks with weights and activations constrained to+ 1 or-1," *arXiv preprint arXiv:1602.02830*, 2016.
- [284] A. Mishra and D. Marr, "Apprentice: Using knowledge distillation techniques to improve low-precision network accuracy," *arXiv preprint arXiv:1711.05852*, 2017.

- [285] J. Haris, P. Gibson, J. Cano, N. B. Agostini, and D. Kaeli, "Secda: Efficient hardware/software co-design of FPGA-based dnn accelerators for edge inference," in *2021 IEEE 33rd International Symposium on Computer Architecture and High Performance Computing (SBAC-PAD)*. IEEE, 2021, pp. 33–43.
- [286] A. Ahmad, M. A. Pasha, and G. J. Raza, "Accelerating tiny yolov3 using FPGA-based hardware/software co-design," in *2020 IEEE international symposium on circuits and systems (ISCAS)*. IEEE, 2020, pp. 1–5.
- [287] C. Zhang, D. Wu, J. Sun, G. Sun, G. Luo, and J. Cong, "Energy-efficient cnn implementation on a deeply pipelined FPGA cluster," in *Proceedings of the 2016 International Symposium on Low Power Electronics and Design*, 2016, pp. 326–331.
- [288] Y. Fukushima, K. Iizuka, and H. Amano, "Parallel implementation of cnn on multi-FPGA cluster," *IEICE TRANSACTIONS on Information and Systems*, vol. 106, no. 7, pp. 1198–1208, 2023.
- [289] Y. Fukushima, K. Iizuka, and H. Amano, "Parallel implementation of vision transformer on a multi-FPGA cluster," in *2023 Eleventh International Symposium on Computing and Networking (CANDAR)*. IEEE, 2023, pp. 100–106.
- [290] S. H. Hozhabr and R. Giorgi, "A survey on real-time object detection on fpgas," *IEEE Access*, vol. 13, pp. 38 195–38 238, 2025.
- [291] H. M. Ahmad and A. Rahimi, "Deep learning methods for object detection in smart manufacturing: A survey," *Journal of Manufacturing Systems*, vol. 64, pp. 181–196, 2022.
- [292] Y. Hu, Y. Liu, and Z. Liu, "A survey on convolutional neural network accelerators: Gpu, FPGA and asic," in *2022 14th International Conference on Computer Research and Development (ICCRD)*. IEEE, 2022, pp. 100–107.
- [293] R. Giorgi, N. Bettin, S. Ermini, F. Montefoschi, and A. Rizzo, "An iris+voice recognition system for a smart doorbell," in *IEEE 8th MECCO*, Jun. 2019, pp. 419–422.
- [294] D. Theodoropoulos and et al., "The AXIOM project (agile, extensible, fast i/o module)," in *IEEE SAMOS*, Jul. 2015, pp. 262–269.
- [295] L. Verdoscia, R. Vaccaro, and R. Giorgi, "A clockless computing system based on the static data-flow paradigm," in *Proc. IEEE Int.l Workshop on Data-Flow Execution Models for Extreme Scale Computing (DFM-2014)*, Edmonton, Canada, Aug. 2014, pp. 30–37.
- [296] R. Giorgi, Z. Popovic, and N. Puzovic, "Implementing fine/medium grained tlp support in a many-core architecture," in *SAMOS 2009*. Springer, Jul. 2009, pp. 78–87.
- [297] G. Tatar and S. Bayar, "Real-time multi-task adas implementation on reconfigurable heterogeneous mp soc architecture," *IEEE Access*, vol. 11, pp. 80 741–80 760, 2023.
- [298] "Smartcamera—kria kv260 2021.1 documentation," [Online; accessed 2025-03-07]. [Online]. Available: <https://xilinx.github.io/kria-apps-docs>
- [299] A. Sahebi, M. Procaccini, and R. Giorgi, "Hashgrid: an optimized architecture for accelerating graph computing on fpgas," *Future Generation Computer Systems*, vol. 162, p. 107497, 2025.
- [300] M. Procaccini, A. Sahebi, and R. Giorgi, "A survey of graph convolutional networks (gcns) in fpga-based accelerators," *Journal of Big Data*, vol. 11, no. 1, p. 163, 2024.
- [301] R. Sarkar, S. Abi-Karam, Y. He, L. Sathidevi, and C. Hao, "Flowgnn: A dataflow architecture for real-time workload-agnostic graph neural network inference," in *2023 IEEE International Symposium on High-Performance Computer Architecture (HPCA)*. IEEE, 2023, pp. 1099–1112.

- [302] G. R. Nair, H.-S. Suh, M. Halappanavar, F. Liu, J.-s. Seo, and Y. Cao, "Fpga acceleration of gcn in light of the symmetry of graph adjacency matrix," in *2023 Design, Automation & Test in Europe Conference & Exhibition (DATE)*. IEEE, 2023, pp. 1–6.
- [303] R. Chen, J. Liu, S. Tang, Y. Liu, Y. Zhu, M. Ling, and B. Da Silva, "Ate-gcn: An fpga-based graph convolutional network accelerator with asymmetrical ternary quantization," in *2025 Design, Automation & Test in Europe Conference (DATE)*. IEEE, 2025, pp. 1–6.
- [304] A. Sahebi, M. Barbone, M. Procaccini, W. Luk, G. Gaydadjiev, and R. Giorgi, "Distributed large-scale graph processing on fpgas," *Journal of big Data*, vol. 10, no. 1, p. 95, 2023.
- [305] W. Jiang, Y. Liu, and H. Zhang, "Ssm-gcn: An fpga-based efficient gcn accelerator for symmetric sparse matrices," *Journal of Engineering Research*, 2025.
- [306] R. Hwang, M. Kang, J. Lee, D. Kam, Y. Lee, and M. Rhu, "Grow: A row-stationary sparse-dense gemm accelerator for memory-efficient graph convolutional neural networks," in *2023 IEEE International Symposium on High-Performance Computer Architecture (HPCA)*. IEEE, 2023, pp. 42–55.
- [307] S. H. Hozhabr and R. Giorgi, "Real-time object detection on fpga-based heterogeneous mpsoCs: A preliminary analysis of the execution bottlenecks," in *2025 14th Mediterranean Conference on Embedded Computing (MECO)*. IEEE, 2025, pp. 1–4.
- [308] S. Ran, B. Zhao, X. Dai, C. Cheng, and Y. Zhang, "Software-hardware co-design for accelerating large-scale graph convolutional network inference on fpga," *Neurocomputing*, vol. 532, pp. 129–140, 2023.
- [309] F. Thabtah, S. Hammoud, F. Kamalov, and A. Gonsalves, "Data imbalance in classification: Experimental evaluation," *Information Sciences*, vol. 513, pp. 429–441, 2020.
- [310] R. Padilla, W. L. Passos, T. L. Dias, S. L. Netto, and E. A. Da Silva, "A comparative analysis of object detection metrics with a companion open-source toolkit," *Electronics*, vol. 10, no. 3, p. 279, 2021.
- [311] A. Boutros and V. Betz, "FPGA architecture: Principles and progression," *IEEE Circuits and Systems Magazine*, vol. 21, no. 2, pp. 4–29, 2021.
- [312] D. O'Loughlin, A. Coffey, F. Callaly, D. Lyons, and F. Morgan, "Xilinx vivado high level synthesis: Case studies," 2014.
- [313] J. Cong, J. Lau, G. Liu, S. Neuendorffer, P. Pan, K. Vissers, and Z. Zhang, "FPGA hls today: successes, challenges, and opportunities," *ACM Transactions on Reconfigurable Technology and Systems (TRETS)*, vol. 15, no. 4, pp. 1–42, 2022.
- [314] J. C. T. Hai, O. C. Pun, and T. W. Haw, "Accelerating video and image processing design for FPGA using hdl coder and simulink," in *2015 IEEE Conference on Sustainable Utilization And Development In Engineering and Technology (CSUDET)*. IEEE, 2015, pp. 1–5.
- [315] N. Kehtarnavaz and S. Mahotra, *Digital Signal Processing Laboratory: LabVIEW-Based FPGA Implementation*. Universal-Publishers, 2010.
- [316] Xilinx, *UG1262*, 2nd ed., Xilinx, October 2019, available at https://www.xilinx.com/support/documentation/sw_manuals/xilinx2019_2/ug1262-model-composer-user-guide.pdf.
- [317] G. Nicolescu and P. J. Mosterman, *Model-based design for embedded systems*. Crc Press, 2018.
- [318] W. Wolf, *FPGA-Based System Design*, 1st ed. USA: Prentice Hall Press, 2009.
- [319] S. Hauck and A. DeHon, *Reconfigurable computing: the theory and practice of FPGA-based computation*. Elsevier, 2010.
- [320] M. Siracusa, E. Del Sozzo, M. Rabozzi, L. Di Tucci, S. Williams, D. Sciuto, and M. D. Santambrogio, "A comprehensive methodology to optimize FPGA designs via the roofline model," *IEEE Transactions on Computers*, vol. 71, no. 8, pp. 1903–1915, 2021.

