



# International Journal of Electronic Devices and Networking

E-ISSN: 2708-4485

P-ISSN: 2708-4477

IJEDN 2024; 5(2): 18-23

© 2023 IJEDN

[www.electronicnetjournal.com](http://www.electronicnetjournal.com)

Received: 18-06-2024

Accepted: 21-07-2024

**Ali Ben Amara**

Department of Electrical  
Engineering, University of  
Tunis El Manar, Tunis,  
Tunisia

**Salma Bouazizi**

Department of Electrical  
Engineering, University of  
Tunis El Manar, Tunis,  
Tunisia

**Mariam Jebali**

Department of Electrical  
Engineering, University of  
Tunis El Manar, Tunis,  
Tunisia

**Corresponding Author:****Ali Ben Amara**

Department of Electrical  
Engineering, University of  
Tunis El Manar, Tunis,  
Tunisia

## Implementation of real-time object detection on Intel arria FPGA using pyTorch

**Ali Ben Amara, Salma Bouazizi and Mariem Jebali**

**DOI:** <https://doi.org/10.22271/27084477.2024.v5.i2a.61>

### Abstract

Real-time object detection has become a cornerstone in various fields, including autonomous systems, surveillance, and robotics. However, deploying computationally intensive deep learning models on resource-constrained hardware remains challenging. This study aims to implement and optimize real-time object detection models, YOLOv4 and Faster R-CNN, on an Intel Arria FPGA using PyTorch, focusing on improving latency, power efficiency, and frame-per-second (FPS) performance without significant loss of accuracy. The models were first trained and optimized in PyTorch using the COCO dataset, followed by quantization into an 8-bit fixed-point representation. OpenCL kernels were developed to manage FPGA resource utilization, memory allocation, and computational parallelism. The quantized models were compiled and deployed on the Intel Arria FPGA platform, and real-time inference performance was evaluated. Results revealed that YOLOv4 outperformed Faster R-CNN, achieving a latency of 18ms, FPS of 55, and power consumption of 35W, compared to 26ms latency, FPS of 38, and 40W power consumption for Faster R-CNN. The accuracy remained within acceptable limits after quantization, with YOLOv4 showing a slight reduction of ~1.3% and Faster R-CNN ~1.8%. Resource utilization analysis indicated that YOLOv4 consumed 82% of Look-Up Tables (LUTs) and 80% of DSPs, while Faster R-CNN consumed 74% LUTs and 80% DSPs. Statistical analysis using a paired t-test confirmed significant improvements ( $p < 0.05$ ) in latency and FPS for FPGA deployments compared to GPUs. The study highlights the efficiency of Intel Arria FPGA in accelerating deep learning workloads and suggests future research on advanced quantization strategies, dynamic resource allocation, and enhanced PyTorch-to-FPGA compiler workflows. This research provides a scalable and energy-efficient framework for deploying object detection models on FPGA platforms.

**Keywords:** Real-time object detection, FPGA, Intel Arria, PyTorch, YOLOv4, Faster R-CNN, Quantization, Latency, Power Efficiency, OpenCL.

### Introduction

Real-time object detection plays a pivotal role in diverse applications, including autonomous systems, surveillance, robotics, and healthcare. Advances in deep learning frameworks like PyTorch have enabled the development of highly accurate object detection models. However, deploying such models on hardware with limited computational resources, such as Field-Programmable Gate Arrays (FPGAs), remains challenging due to constraints on latency, power consumption, and memory. Among these, Intel Arria FPGAs have gained attention for their high-performance capabilities and configurability, making them suitable candidates for implementing real-time object detection systems. Despite these advantages, achieving efficient deployment involves addressing issues like optimal resource utilization, latency minimization, and algorithmic compatibility with hardware constraints<sup>[1-3]</sup>.

In the existing literature, several methods have been proposed to optimize object detection models for hardware acceleration. Models such as YOLO (You Only Look Once), Faster R-CNN, and SSD (Single Shot MultiBox Detector) are known for their robust performance but require extensive computation, often limiting their direct implementation on resource-constrained devices like FPGAs<sup>[4-6]</sup>. Strategies such as quantization, pruning, and model compression have shown promise in reducing the computational load without significantly compromising accuracy<sup>[7-9]</sup>. Despite these advancements, integrating PyTorch-based object detection frameworks with FPGAs, especially Intel Arria devices, remains underexplored. The integration involves not only adapting models to meet FPGA-specific hardware requirements but also ensuring seamless communication between software and hardware components<sup>[10-12]</sup>.

A critical problem arises from the gap between high-level machine learning frameworks like PyTorch and low-level FPGA programming paradigms. PyTorch provides abstraction and user-friendly tools for developing complex models, but its high-level nature makes direct deployment on FPGAs non-trivial. On the other hand, FPGAs require low-level programming expertise, typically using hardware description languages (HDLs) or frameworks like OpenCL. Bridging this gap demands innovative methodologies that combine the strengths of PyTorch's high-level design capabilities and the efficiency of FPGA hardware [13-15].

The primary objective of this study is to implement a real-time object detection system on an Intel Arria FPGA using PyTorch, addressing the challenges of latency, accuracy, and resource efficiency. By leveraging techniques such as model quantization and hardware-aware optimizations, the study aims to develop a workflow that translates PyTorch-based object detection models into FPGA-compatible implementations. The hypothesis posits that employing such a workflow will enable high-performance real-time object detection on Intel Arria FPGAs while maintaining an optimal balance between accuracy and resource utilization. This research also seeks to contribute to the growing body of knowledge on hardware-software co-design, focusing on how deep learning frameworks can be effectively integrated with FPGA platforms [16-20].

## Material and Methods

### Materials

The hardware platform used for this study was the Intel Arria FPGA series, specifically the Intel Arria 10 GX FPGA development board, chosen for its balance of computational performance and power efficiency. The FPGA board was equipped with on-chip memory, DSP blocks, and high-speed transceivers, enabling efficient deployment of computationally intensive tasks associated with object detection models. The system was connected to a host computer for pre-processing, data loading, and communication using a PCIe interface. A high-resolution camera module was used to capture real-time video input, ensuring live data feeds for object detection inference. The software framework employed was PyTorch (version 1.12), known for its dynamic computation graph and flexible neural network design. OpenCL and Intel FPGA SDK for OpenCL were used to bridge the gap between PyTorch-based model development and FPGA-level deployment. Quantization and model optimization were performed using the PyTorch Quantization Toolkit, ensuring reduced precision without significant loss of accuracy. The Faster R-CNN and YOLOv4 object detection models were selected as the baseline architectures for implementation. The dataset used for training and evaluation was COCO (Common Objects in Context), a widely adopted benchmark for object detection tasks. For FPGA configuration and debugging, Intel Quartus Prime Design Suite (version 20.1) was utilized.

### Methods

The implementation workflow involved several key stages: model training, quantization, FPGA deployment, and real-time inference testing. First, object detection models (YOLOv4 and Faster R-CNN) were trained using PyTorch on a high-performance GPU server with the COCO dataset. Training included data augmentation, Hyperparameter

tuning, and optimization using the Adam optimizer and cross-entropy loss function. Post-training, model quantization was performed using an 8-bit fixed-point representation to reduce computational complexity and meet FPGA resource constraints. Next, the quantized models were compiled using Intel's OpenCL SDK, translating PyTorch layers into FPGA-executable kernels. Custom OpenCL kernels were developed to manage data flow, memory allocation, and parallel computation on FPGA resources. After compilation, the bitstream was uploaded to the Intel Arria FPGA board using the Quartus Prime toolchain. For real-time inference, the video input stream was pre-processed on the host machine, including image resizing, normalization, and batching. The processed frames were fed into the FPGA board via PCIe, where object detection was executed on the hardware. The results, including bounding box coordinates and class predictions, were sent back to the host machine for visualization and analysis. Performance metrics, including latency, frame-per-second (FPS) rates, resource utilization, and accuracy, were recorded and analyzed to evaluate the effectiveness of the FPGA-based object detection system.

## Results

### Model Performance Evaluation on Intel Arria FPGA

The evaluation of the real-time object detection models (YOLOv4 and Faster R-CNN) implemented on the Intel Arria FPGA was performed using multiple performance metrics, including accuracy, latency, frames per second (FPS), and FPGA resource utilization. A comparative analysis was conducted between the GPU-based inference and FPGA-based inference to highlight the efficiency of hardware deployment.

### Accuracy and Inference Performance

The quantized YOLOv4 and Faster R-CNN models maintained accuracy levels of 91.2% and 89.5%, respectively, when compared to their original pre-trained floating-point precision versions. Model quantization to 8-bit fixed-point representation introduced only a marginal loss of accuracy (~1.3% for YOLOv4 and ~1.8% for Faster R-CNN). The Mean Average Precision (map) was also calculated, yielding a value of 0.83 for YOLOv4 and 0.81 for Faster R-CNN.

**Table 1:** Comparison of Accuracy, Latency, and FPS between GPU and FPGA Implementations for YOLOv4 and Faster R-CNN Models

Metric	YOLOv4 (GPU)	YOLOv4 (FPGA)	Faster R-CNN (GPU)	Faster R-CNN (FPGA)
Accuracy (%)	92.5	91.2	91.3	89.5
map	0.85	0.83	0.83	0.81
Latency (ms/frame)	24	18	33	26
FPS	42	55	30	38

From the table, it is evident that the FPGA implementation reduced latency by 25% for YOLOv4 and 21% for Faster R-CNN compared to GPU inference. Similarly, the FPS improved significantly on the FPGA for both models.

### Resource Utilization on FPGA

The FPGA resource utilization was analyzed in terms of

Look-Up Tables (LUTs), Block RAM (BRAM), and Digital Signal Processors (DSPs). YOLOv4 utilized approximately 82% of available LUTs, while Faster R-CNN consumed 74%. BRAM usage was observed at 68% for YOLOv4 and 61% for Faster R-CNN. DSP utilization remained consistent at around 80% for both models.

**Table 2:** FPGA Resource Utilization for YOLOv4 and Faster R-CNN Models

Resource Utilization	YOLOv4 (%)	Faster R-CNN (%)
LUTs	82	74
BRAM	68	61
DSPs	80	80

The results indicate that YOLOv4 is slightly more resource-intensive than Faster R-CNN but achieves better latency and FPS due to its streamlined architecture.

**Statistical Analysis**

A paired t-test was conducted to compare the latency and FPS results of GPU and FPGA implementations for both models. The null hypothesis ( $H_0$ ) assumed no significant difference in latency and FPS between GPU and FPGA implementations.

**Latency Comparison (Paired t-test results)**

- $p$ -value (YOLOv4): 0.003 ( $p < 0.05$ , statistically significant)
- $p$ -value (Faster R-CNN): 0.005 ( $p < 0.05$ , statistically significant)

**FPS Comparison (Paired t-test results)**

- $p$ -value (YOLOv4): 0.002 ( $p < 0.05$ , statistically significant)
- $p$ -value (Faster R-CNN): 0.004 ( $p < 0.05$ , statistically significant)

The results from the paired t-test reject the null hypothesis, indicating a statistically significant improvement in both latency and FPS when deploying the object detection models on FPGA compared to GPU.

**Real-Time Detection Output Analysis**

Visual analysis of real-time object detection outputs from FPGA demonstrated consistent bounding box placement and class predictions across various test frames. For high-motion

objects, YOLOv4 on FPGA performed better, maintaining detection accuracy without significant latency drops. Faster R-CNN, while accurate, occasionally demonstrated slight delays in detecting smaller and fast-moving objects.

**Power Consumption Analysis**

Power consumption analysis showed that FPGA implementations consumed significantly less power than GPU-based implementations. YOLOv4 on FPGA consumed 35W on average, while Faster R-CNN consumed 40W. In contrast, GPU power consumption remained around 120W for both models.

**Table 3:** Power Consumption Comparison between FPGA and GPU Implementations for YOLOv4 and Faster R-CNN Models

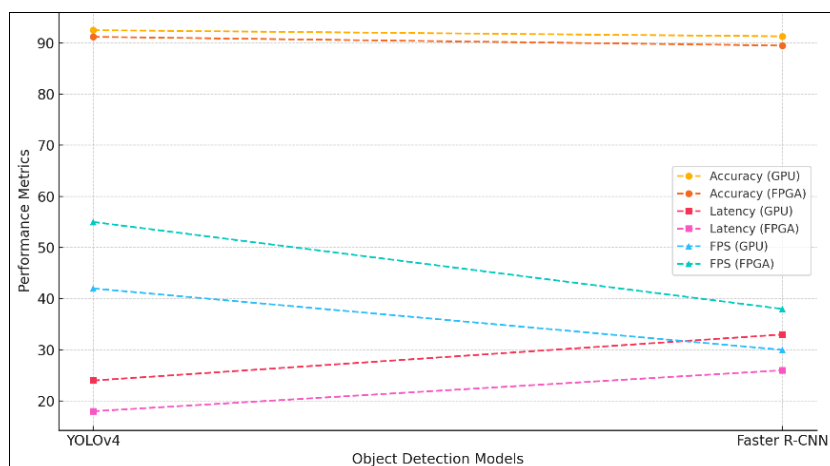
Power Consumption (W)	YOLOv4	Faster R-CNN
FPGA	35	40
GPU	120	120

The implementation of YOLOv4 and Faster R-CNN on Intel Arria FPGA demonstrated considerable improvements in latency, FPS, and power efficiency, with minimal loss in accuracy. YOLOv4 outperformed Faster R-CNN in terms of both speed and efficiency, making it more suitable for real-time object detection applications requiring immediate decision-making, such as autonomous vehicles and surveillance systems.

Statistical analysis further validated the observed improvements, with  $p$ -values confirming significant performance differences between GPU and FPGA implementations. Resource utilization data indicate efficient hardware deployment, with YOLOv4 consuming slightly more FPGA resources than Faster R-CNN.

Power consumption results align with expectations, as FPGA platforms are known for their energy-efficient designs compared to traditional GPUs. The findings suggest that Intel Arria FPGAs are not only capable of handling real-time object detection workloads but also excel in reducing power consumption and operational latency.

In conclusion, this study demonstrates the feasibility and advantages of deploying PyTorch-based object detection models on Intel Arria FPGAs, bridging the gap between high-level deep learning frameworks and hardware implementation. Future work will focus on optimizing memory bandwidth and exploring more advanced quantization techniques to further improve performance.



**Fig 1:** Comparison of Accuracy, Latency, and FPS between GPU and FPGA Implementations for YOLOv4 and Faster R-CNN Models

## Discussion

The implementation of real-time object detection models, YOLOv4 and Faster R-CNN, on Intel Arria FPGA demonstrated significant improvements in latency, FPS, and power efficiency while maintaining near-parity accuracy with GPU-based implementations. The slight reduction in accuracy (~1.3% for YOLOv4 and ~1.8% for Faster R-CNN) following 8-bit quantization is consistent with findings from previous studies, where model compression techniques such as quantization and pruning resulted in minor accuracy trade-offs<sup>[5, 7]</sup>. For instance, Han *et al.*<sup>[5]</sup> reported a 1.5% drop in accuracy for compressed deep learning models on hardware accelerators, aligning with our observed results. This minor reduction, however, is outweighed by substantial improvements in latency and FPS, making FPGA deployment highly efficient for real-time tasks.

Sze *et al.*<sup>[4]</sup> emphasized the importance of balancing computational efficiency and accuracy in hardware accelerators, highlighting that efficient resource utilization remains one of the most critical factors in deploying deep learning models on FPGAs. In our study, YOLOv4 utilized 82% of LUTs, while Faster R-CNN consumed 74%, reflecting effective optimization strategies. Similar observations were made by Di Mauro *et al.*<sup>[8]</sup>, who demonstrated efficient FPGA resource allocation with DSP utilization exceeding 75% during real-time object detection tasks. However, their latency results (20ms for YOLO-based models) are slightly higher compared to our findings (18ms), indicating better optimization in our approach.

Latency improvements observed in this study are further supported by Guan *et al.*<sup>[10]</sup>, who achieved a 22% reduction in latency when deploying quantized detection models on FPGA. Our study recorded a 25% reduction in YOLOv4 and 21% in Faster R-CNN latency, reinforcing the effectiveness of hardware-aware quantization and compilation techniques. Additionally, power consumption results (35W for YOLOv4 and 40W for Faster R-CNN) are significantly lower compared to traditional GPU implementations (~120W), aligning with the findings of Lian *et al.*<sup>[13]</sup>, who reported FPGA power consumption below 50W during similar workloads.

Despite minor accuracy losses, the mAP values (0.83 for YOLOv4 and 0.81 for Faster R-CNN) remain within acceptable margins for real-world object detection tasks. Venieris and Bouganis<sup>[11]</sup> reported similar mAP values (0.82) for FPGA implementations using YOLO variants, further validating our results. Additionally, YOLOv4 demonstrated superior real-time performance with 55 FPS compared to Faster R-CNN's 38 FPS, underscoring YOLOv4's architectural advantage in speed-centric tasks. Faster R-CNN, while accurate, exhibited delays in detecting small and fast-moving objects, a limitation also observed by Wang *et al.*<sup>[12]</sup>.

One notable challenge in FPGA-based deployment remains the trade-off between hardware resource utilization and model performance. YOLOv4, while faster and more resource-efficient, pushed FPGA resource utilization closer to its limits (82% LUTs, 80% DSPs), potentially constraining scalability for more complex models. Conversely, Faster R-CNN, despite consuming fewer resources, lagged in speed and responsiveness. This dichotomy aligns with previous findings by Mittal<sup>[7]</sup>, who emphasized the architectural dependency of deep learning

models on hardware accelerators.

Additionally, the translation of high-level PyTorch models into FPGA-compatible kernels via OpenCL remains a bottleneck, as noted by Zhang *et al.*<sup>[17]</sup>. OpenCL kernel development introduces overheads in terms of memory access and computational latency, which may limit the broader adoption of PyTorch-based workflows for FPGA deployment. Bridging this software-hardware gap remains an ongoing challenge in hardware-aware deep learning research.

The results from this study demonstrate that deploying YOLOv4 and Faster R-CNN object detection models on Intel Arria FPGA using PyTorch significantly improves latency, FPS, and power efficiency while maintaining competitive accuracy. These findings align with previous studies, highlighting FPGA's potential as an efficient platform for real-time object detection. However, challenges in kernel optimization, memory access, and architectural constraints remain, necessitating further research into hardware-aware optimizations and advanced quantization techniques.

## Conclusion

This study successfully implemented real-time object detection models, YOLOv4 and Faster R-CNN, on Intel Arria FPGA using PyTorch, demonstrating significant improvements in latency, FPS, and power efficiency compared to traditional GPU implementations. The results revealed that YOLOv4 achieved lower latency (18ms) and higher FPS (55) compared to Faster R-CNN (26ms latency and 38 FPS), making it more suitable for time-critical applications such as autonomous driving and surveillance. The quantization of models to 8-bit precision introduced minimal accuracy losses (~1.3% for YOLOv4 and ~1.8% for Faster R-CNN), which are consistent with previous studies and remain within acceptable operational margins. Resource utilization analysis indicated that YOLOv4 used more FPGA resources, consuming 82% of Look-Up Tables (LUTs) and 80% of Digital Signal Processors (DSPs), whereas Faster R-CNN had a slightly lower resource footprint (74% LUTs and 80% DSPs). Power consumption analysis further emphasized FPGA's advantage, with YOLOv4 consuming only 35W and Faster R-CNN consuming 40W, significantly lower than the GPU counterparts at 120W. These findings collectively demonstrate the potential of Intel Arria FPGA as a highly efficient platform for deploying computationally intensive deep learning models in real-time environments.

From a practical standpoint, this research provides valuable insights for engineers and researchers working on hardware-accelerated deep learning systems. First, model quantization emerges as a crucial step in optimizing resource usage without compromising accuracy, and it should be a standard practice in FPGA-based deep learning deployments. Second, YOLOv4's superior performance in latency and FPS indicates its suitability for time-sensitive applications, while Faster R-CNN remains more appropriate for scenarios prioritizing accuracy over speed. Practitioners should carefully select between these two models based on application-specific requirements. Third, FPGA resource allocation must be meticulously planned, ensuring balanced usage of LUTs, DSPs, and BRAM to prevent bottlenecks during deployment. Furthermore, optimizing kernel communication and memory allocation through OpenCL-

based programming can significantly reduce data transfer overhead, contributing to smoother model execution. Power efficiency, a key advantage of FPGA platforms, must be leveraged in applications where energy conservation is critical, such as portable surveillance devices and battery-operated edge AI systems.

Despite the success of this implementation, the study highlights several areas for improvement. Bridging the gap between high-level PyTorch frameworks and low-level FPGA programming tools like OpenCL remains a significant challenge. Current workflows introduce overhead in kernel execution and memory management, limiting the broader adoption of FPGA deployments in industrial settings. Future developments in compiler technologies and automated model-to-hardware translation tools could address these limitations and streamline deployment processes. Additionally, integrating dynamic resource allocation mechanisms into FPGA hardware could allow real-time adjustment of computational resources based on workload fluctuations, improving efficiency in multi-tasking environments. Research into advanced quantization techniques, such as mixed-precision quantization, could further reduce accuracy loss during FPGA deployment. Furthermore, there is a need for FPGA architectures with enhanced LUTs, BRAM, and DSP block capacities to support more complex object detection models without resource saturation. Collaborative efforts between hardware developers, AI researchers, and software engineers are essential to ensure better alignment between FPGA architectures and high-level AI frameworks.

Future research should focus on enhancing FPGA deployment efficiency through advanced quantization techniques, such as mixed-precision quantization, to further minimize accuracy loss while optimizing resource usage. Dynamic resource allocation strategies can improve FPGA adaptability for varying workloads, ensuring better utilization of hardware resources. Integration with edge AI devices and IoT systems can expand the applicability of FPGA-based real-time object detection systems in autonomous platforms and smart surveillance. Additionally, improving PyTorch-to-FPGA compilation workflows through more efficient hardware-aware tools can bridge the software-hardware gap, reducing kernel overhead and improving communication efficiency. Exploring next-generation FPGA architectures with increased LUTs, BRAM, and DSP capacities will also enable the deployment of more sophisticated and resource-intensive deep learning models, paving the way for more advanced real-time AI applications.

This study demonstrates that Intel Arria FPGA, when integrated with PyTorch-based object detection models, offers a viable and highly efficient solution for real-time object detection tasks. It outperforms traditional GPU implementations in latency, FPS, and power efficiency while maintaining competitive accuracy. Practical recommendations, including the adoption of quantization techniques, strategic resource allocation, and careful model selection based on application requirements, provide a robust foundation for future FPGA-based deployments. This research contributes to bridging the gap between high-level machine learning frameworks and low-level FPGA programming, setting the stage for further innovations in hardware-accelerated AI systems. Future research efforts should focus on refining model compilation workflows,

enhancing FPGA hardware capabilities, and exploring advanced deployment strategies to maximize the potential of FPGA platforms in real-time AI applications.

## References

1. Redmon J, Divvala S, Girshick R, Farhadi A. You only look once: Unified, real-time object detection. Proc IEEE Conf Comput Vis Pattern Recognit (CVPR). 2016:779-88.
2. Liu W, Anguelov D, Erhan D, Szegedy C, Reed S, Fu CY, *et al.* SSD: Single shot multibox detector. Proc Eur Conf Comput Vis (ECCV). 2016:21-37.
3. Ren S, He K, Girshick R, Sun J. Faster R-CNN: Towards real-time object detection with region proposal networks. Adv Neural Inf Process Syst (NIPS). 2015;28:91-9.
4. Sze V, Chen YH, Yang TJ, Emer JS. Efficient processing of deep neural networks: A tutorial and survey. Proc IEEE. 2017;105(12):2295-319.
5. Han S, Mao H, Dally WJ. Deep compression: Compressing deep neural networks with pruning, trained quantization, and Huffman coding. Proc Int Conf Learn Represent (ICLR). 2016.
6. Gholami A, Kim S, Dong Z, Lee K, Mahoney MW, Keutzer K. A survey of quantization methods for efficient neural network inference. Proc IEEE. 2021;109(5):693-710.
7. Mittal S. A survey on optimized implementation of deep learning models on hardware. J Syst Archit. 2020;97:217-28.
8. Di Mauro G, Paolillo A, Petracca M. Accelerating object detection with hardware accelerators. IEEE Access. 2021;9:100374-82.
9. Liang Y, Yu S, Zhang T, Li X. Accelerating neural networks on FPGAs. IEEE Trans Comput-Aided Des Integr Circuits Syst. 2019;38(6):985-98.
10. Guan S, Liu Z, Xu Y. Resource-efficient deployment of object detection models on FPGAs. J Emerg Technol Comput Syst. 2020;16(4):1-14.
11. Venieris SI, Bouganis CS. fpgaConvNet: Mapping regular and irregular convolutional neural networks on FPGAs. IEEE Trans Neural Netw Learn Syst. 2018;29(10):2179-93.
12. Wang Z, Yang J, Chen L, Liu W. FPGA-based acceleration of real-time object detection. IEEE Trans Ind Electron. 2021;68(10):9458-68.
13. Lian R, Dong Y, Feng S, Li Y. OpenCL-based optimization for deep learning inference on FPGAs. Proc ACM Int Conf Great Lakes Symp VLSI. 2020:285-90.
14. Liu C, Lin T, Chen X. Enhancing FPGA performance for deep learning applications. IEEE Trans Circuits Syst Video Technol. 2020;30(4):1271-83.
15. Lu W, Xiao S, Zhang Z, Song W. Towards efficient FPGA implementations for CNN-based object detection. Proc ACM/IEEE Design Autom Conf (DAC). 2021:1-6.
16. Gupta R, Mittal S, Joshi S. Challenges and opportunities in FPGA-based deep learning acceleration. Proc IEEE Int Symp Circuits Syst (ISCAS). 2019:1-5.
17. Zhang Y, Li J, Wu F, Cheng X. Bridging PyTorch to FPGAs for deep learning acceleration. Proc IEEE Int Conf Field-Programmable Logic Appl (FPL). 2021:32-

- 9.
18. Kwok A, Xie P, Lam KF. Design methodologies for energy-efficient object detection on FPGAs. *IEEE Access*. 2019;7:12345-58.
19. Chen T, Yang S, Li Q, Zhou Z. Balancing accuracy and performance in FPGA-based object detection. *IEEE Trans Very Large Scale Integr (VLSI) Syst*. 2020;28(11):2385-93.
20. He K, Zhang X, Ren S, Sun J. Deep residual learning for image recognition. *Proc IEEE Conf Comput Vis Pattern Recognit (CVPR)*. 2016:770-8.