

# Yolov4 Object Detection Using FPGA

1<sup>st</sup> Saif Alomari  
*College of Engineering*  
*Cal Poly Pomona*  
Pomona, USA  
ssalomari@cpp.edu

2<sup>nd</sup> Jared Alanis  
*College of Engineering*  
*Cal Poly Pomona*  
Pomona, USA  
jaredalanis@cpp.edu

3<sup>rd</sup> Dawson Graf  
*College of Engineering*  
*Cal Poly Pomona*  
Pomona, USA  
dgraf@cpp.edu

4<sup>th</sup> Benjamin Black  
*College of Engineering*  
*Cal Poly Pomona*  
Pomona, USA  
bjblack@cpp.edu

**Abstract**—This paper details the design and implementation of an integrated system combining an OV7670 camera module with a Field Programmable Gate Array (FPGA) for real-time image acquisition and display, followed by advanced object detection using the Yolov4 model on a separate PC. The primary focus was to exploit the FPGA’s rapid processing capabilities to manage high-volume image data efficiently and to utilize Yolov4 for its high accuracy in detecting multiple objects. The project highlights the seamless integration of embedded hardware with sophisticated machine learning algorithms to achieve real-time object detection, emphasizing the challenges and successes encountered during its execution.

**Index Terms**—Object detection, Yolov4, FPGA, OV7670, GPU, Machine Learning

## I. INTRODUCTION

This paper presents a project that embodies object detection software with embedded hardware system, utilizing both traditional and cutting-edge technologies to achieve object detection. The project integrates an OV7670 camera module with a Field Programmable Gate Array (FPGA) for image capture and display, followed by object detection using a pre-trained Yolov4 model on a separate computer system.

The project commenced with the development of a robust image acquisition system. The OV7670 camera, selected for its affordability and straightforward interfacing, was connected to an FPGA. The FPGA played a pivotal role in managing the high-speed data transfer from the camera, processing the image data for real-time display on a VGA screen. This demonstrated the FPGA’s capability to handle and preprocess image data efficiently, though it did not process the data in real-time for Yolov4 object detection.

After the image data was displayed, it was then transferred to a computer equipped with the Yolov4 object detection model. Yolov4, known for its high accuracy and the capability to detect multiple objects, utilized a dataset trained to recognize 80 different object classes. The system’s setup allowed for the captured video to be analyzed post-capture. Each frame was processed by Yolov4, which identified and marked objects with bounding boxes. This step underscored the project’s core objective: to illustrate the effectiveness of integrating FPGA-based image processing systems with sophisticated AI-driven object detection models.

This paper aims to detail the architecture of the combined system, explore the challenges encountered during its implementation, and discuss the outcomes and potential of

combining FPGA technology with advanced object detection algorithms.

YOLOv4 represents a significant advancement in the domain of object detection, designed to achieve optimal speed and accuracy. Authored by Alexey Bochkovskiy, Chien-Yao Wang, and Hong-Yuan Mark Liao, the paper titled “YOLOv4: Optimal Speed and Accuracy of Object Detection” presents a comprehensive framework that builds on the strengths of previous iterations of the YOLO series, specifically improving upon YOLOv3. The primary objective of YOLOv4 is to operate efficiently in real-time on conventional GPUs while requiring only modest computational resources for training. This is particularly crucial as the most accurate modern neural networks often do not operate in real-time and demand extensive GPU resources. YOLOv4 addresses this gap by introducing several innovations, including new training strategies like Cross mini-Batch Normalization (CmBN) and Self-adversarial Training (SAT), as well as architectural enhancements such as Weighted-Residual Connections (WRC) and Cross-Stage Partial connections (CSP) [1].

The paper meticulously evaluates the impact of various state-of-the-art features and methods that contribute to the enhanced performance of convolutional neural networks (CNNs). These include advanced data augmentation techniques like Mosaic, which integrates multiple training images into a single batch, thereby improving object detection robustness by presenting objects in varied contexts. Additionally, YOLOv4 utilizes a new activation function, Mish, and incorporates structural elements like DropBlock regularization and CIoU loss, which refine bounding box predictions by considering aspects such as overlap, distance, and aspect ratio. Through these enhancements, YOLOv4 achieves remarkable accuracy on the MS COCO dataset with 43.5 percent AP and operates at real-time speeds, approximately 65 FPS on a Tesla V100, making it a formidable tool in real-world applications where speed and accuracy are paramount [1].

The research presented in the Real-Time Small Drones Detection Based on Pruned YOLOv4 paper delves into enhancing object detection algorithms to address this challenge effectively. A focal point of this study is the adaptation and evaluation of YOLOv4, alongside other state-of-the-art models such as RetinaNet and FCOS, to optimize drone detection capabilities. The paper highlights the inherent challenges in detecting rapidly moving and physically small drones, which

often escape detection due to their limited pixel footprint in digital imagery. To tackle these issues, the authors have innovatively modified YOLOv4 by pruning its convolutional channels and layers, significantly boosting its operational speed by 60.4 percent while maintaining a high detection accuracy with a 90.5 percent mean average precision (mAP). Additionally, they introduce a novel augmentation technique for enhancing the detection of small-scale objects, which substantially increases both the precision and recall rates of the pruned-YOLOv4 model. This tailored approach not only advances the field of object detection but also underscores the potential of deep learning models to adapt to specific real-world challenges in surveillance and security using object detection [2].

Field Programmable Gate Arrays (FPGAs) are particularly advantageous for real-time applications due to their ability to perform parallel processing, which significantly accelerates data throughput and computational tasks. This capability is essential in our project, where real-time video processing demands rapid and efficient handling of high-volume data streams. Moreover, the reconfigurability of FPGAs allows for a flexible adaptation to varied computational needs, a feature that is pivotal when updating or optimizing the system post-deployment. By integrating FPGAs, we harness these benefits to ensure that our embedded system not only meets the performance criteria but also remains scalable and adaptable to future enhancements or changes in application requirements [3].

In this project, the utilization of Field Programmable Gate Arrays (FPGAs) as a cornerstone for implementing embedded systems emerges as a strategic choice, leveraging their inherent parallel architecture to enhance machine learning applications. According to recent studies such as those presented by Chen et al, FPGAs provide a versatile and efficient platform for deploying complex computational tasks in real-time. This adaptability is crucial in scenarios requiring high performance with strict power consumption and cost-effectiveness constraints [4].

In our project, the OV7670 camera module plays a crucial role in the image acquisition process, interfacing seamlessly with the FPGA to capture visual inputs effectively. This CMOS camera sensor, known for its compact size and low power consumption, is ideal for embedded applications requiring real-time image processing. The FPGA, with its robust processing capabilities, is configured to receive the digital video output from the OV7670 through a GPIO connection, facilitating the direct handling of image data. The integration of the OV7670 with the FPGA allows for the rapid capture and processing of images, essential for the real-time performance demands of our system. The configuration and control of the camera are managed via the I2C interface, enabling the adjustment of settings and the retrieval of status information, which are critical for optimizing image quality and system responsiveness [5].

The integration of the OV7670 camera module with FPGA technology provides a robust platform for capturing and displaying real-time video on a VGA monitor. This camera

module, which delivers digital video output at 30 frames per second through an 8-bit interface, is particularly suited for applications where cost and space constraints are significant. In our project, the OV7670 is interfaced with an FPGA board that has a VGA port to facilitate the display of images. This setup enables the translation of digital signals from the FPGA into the signals required by VGA inputs, covering the RGB data channels and synchronization signals. The FPGA's role in this configuration is critical as it handles the real-time processing and formatting of video data, ensuring that the images captured by the OV7670 are promptly and accurately relayed to the VGA display.

This capability of displaying high-quality images on VGA screens highlights the efficiency of using FPGAs in video processing applications, leveraging their ability to perform parallel processing and manage multiple data streams simultaneously. The FPGA not only processes the video signals from the OV7670 but also coordinates the synchronization needed for VGA output through precise timing and signal conversion. By employing an FPGA with an external VGA driver, the system achieves a seamless and synchronized display of video content, making it ideal for a variety of applications, including surveillance, vehicle navigation, and interactive systems. The implementation showcases the flexibility and performance advantages of FPGAs in embedded systems, particularly in handling complex tasks like real-time video display, which requires both high-speed data processing and strict timing coordination to ensure smooth visual output [6].

## II. SYSTEM DESIGN AND METHODOLOGY

The primary objective of our project is to develop a cohesive system that integrates an OV7670 camera module with an FPGA for efficient image acquisition and display, which then interfaces with a computer running the YOLOv4 model for object detection. This system capitalizes on the FPGA's rapid processing capabilities to effectively manage high-volume image data, paired with the sophisticated object detection functionalities provided by YOLOv4. Our methodology incorporates a variety of tools and technologies, including hardware description languages (HDLs) for programming the FPGA, and the advanced machine learning algorithms of YOLOv4. The design is optimized for speed and accuracy, ensuring that video data is swiftly acquired and displayed on a VGA screen, with subsequent processing stages outlined in the full system diagram shown in Figure 1.

As depicted in Figure 1, the architecture of the system ensures seamless integration and synchronization between the OV7670 camera module and the FPGA. This precise synchronization is critical, allowing the system to continuously process the incoming stream of image data without delays or losses. Each frame captured by the OV7670 is immediately processed by the FPGA, which performs initial image processing tasks such as noise reduction, color conversion, and frame resizing in real time. These processed images are then displayed on a VGA screen. The enhanced preprocessing not only improves the visual quality of the output on the VGA display but also

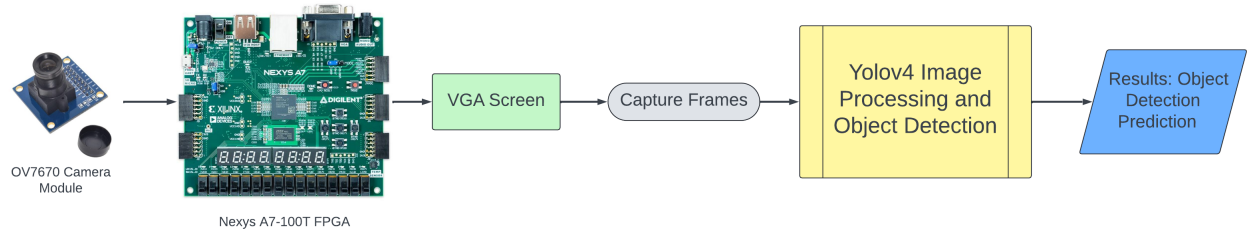


Fig. 1. Full System Diagram

readies the data for subsequent object detection analysis. Once processed, the video is captured from the VGA output and transferred to a PC where the YOLOv4 model is installed. This model processes the video to detect and annotate objects, producing a resultant video that highlights detected objects with bounding boxes. This detailed flow from image capture to object detection elucidates the integration and utility of both hardware and software components in handling complex real-time data processing and analysis tasks effectively.

#### A. Experimental setup

The experimental setup, as illustrated in Figure 2, comprehensively demonstrates the integration of the OV7670 camera module with the FPGA, which is further connected to a VGA screen. The OV7670 camera module captures visual data and transmits it to the FPGA via General Purpose Input/Output (GPIO) interfaces. Within the FPGA, this data undergoes initial processing, including formatting and preliminary image adjustments necessary for VGA compatibility. This processed data is then displayed on a VGA screen, allowing for real-time visual verification of the captured images.

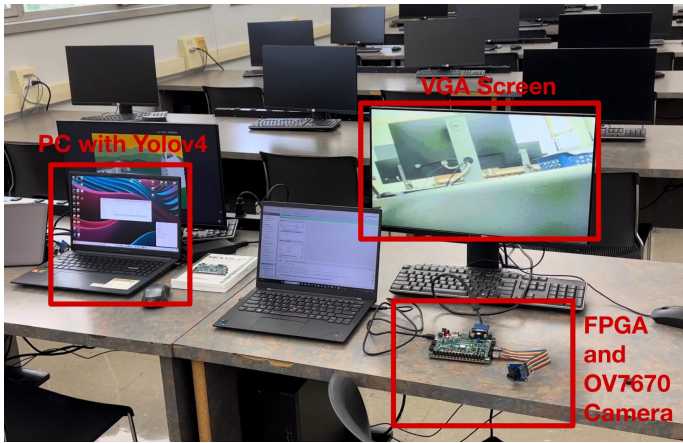


Fig. 2. Experimental Setup

In this setup, while the FPGA efficiently handles real-time data processing and display, it does not directly communicate with the computer system running the YOLOv4 model. Instead, the video output displayed on the VGA screen is captured and transferred to the PC manually. This process is necessary due to the absence of a high-speed direct data link between the

FPGA and the PC. Once transferred to the PC, the YOLOv4 model processes the video input to detect and identify objects within the scenes, annotating them with bounding boxes to highlight detected objects. This process effectively leverages YOLOv4's advanced object detection capabilities within the system's operational constraints. The setup ensures that each component—from image capture to object detection—functions cohesively to form a complete image processing and analysis system as shown in Figure 2. This experimental configuration allows us to explore and optimize the system's performance across different stages of data handling and processing.

#### B. OV7670 Camera

The OV7670 camera module plays a pivotal role in our experimental setup, serving as the primary tool for capturing video data. This camera module is particularly well-suited for embedded system applications due to its compact size, low power consumption, and excellent cost-effectiveness. Its ability to deliver sufficiently high-quality images at a resolution of 640x480 pixels makes it an ideal choice for real-time image processing tasks.

In our system, the OV7670 provides a continuous stream of video data to the FPGA, capturing dynamic scenes with precision. The camera's configuration allows for adjustments in frame rate and exposure, enabling it to adapt to various lighting conditions, which is crucial for maintaining the quality and consistency of the video data under different experimental conditions. Additionally, the camera's 8-bit output interface facilitates seamless integration with the FPGA, ensuring that the data transfer between the camera and the FPGA is straightforward and efficient.

The choice of the OV7670 camera module for our project highlights its appropriateness for researchers and developers looking to implement similar image processing and object detection systems. Its affordability and accessibility, combined with its performance capabilities, make it a top choice for anyone venturing into the field of embedded vision systems. In our setup, as depicted in the experimental arrangement, the OV7670 camera not only fulfills the role of data acquisition but also enhances the overall efficacy and adaptability of the system to various application needs. This makes the OV7670 an invaluable component of our research toolkit, contributing significantly to the system's ability to perform detailed and accurate visual analysis.

### C. FPGA and System Integration

The FPGA plays a crucial role in our system, serving as the central hub for receiving, processing, and forwarding video data to the VGA display and subsequently to the PC for object detection. As shown in Figure 3, the HDL System Diagram illustrates the comprehensive view of the modules within the system, starting from the OV7670 Module, through the Camera Top module, to processing the data and saving them in Block RAM (BRAM), and finally outputting to the VGA Top module.

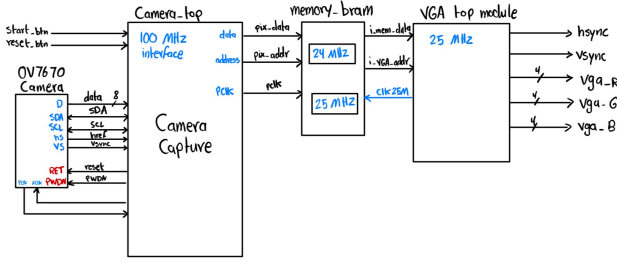


Fig. 3. HDL System Diagram

This architecture is made possible through the use of a sophisticated Hardware Description Language (HDL), specifically Verilog, which allows for detailed and flexible control over hardware functionality. The Verilog code provided outlines a complex integration where the OV7670 camera captures images which are then processed by the FPGA. This processed data is temporarily stored in BRAM before being transmitted to a VGA screen for display. The system also includes mechanisms to synchronize resets across different clock domains, ensuring stable operation throughout the data processing and display phases.

In our system, the FPGA's configuration involves multiple clock domains and synchronizers to manage the flow of data across components operating at different frequencies. The integration of the camera with the FPGA through the cameratop module illustrates the initial capture of image data. Following this, the membram module acts as the intermediary storage for pixel data, interfacing directly with the vga top module that handles the display logic for the VGA output.

The use of FPGA not only enhances the system's ability to handle real-time image processing but also ensures that the transition between capturing data and displaying it on the VGA screen is seamless and efficient. This setup is particularly advantageous for applications requiring rapid processing and display of image data, such as in surveillance systems or interactive media installations. The flexibility and power of FPGA, combined with the modular design enabled by HDL, provide a robust framework for developing complex image processing applications that can be tailored to meet specific project requirements.

### D. Yolov4 Model

The Yolov4 model represents the cutting edge in deep learning technologies for object detection, crucial for interpreting and analyzing the video data captured in our system. As shown in Figure 4, the Yolov4 model processes images through multiple stages to detect and classify objects within a frame. These stages include convolutional layers that extract features, followed by layers that predict classes and bounding boxes, and finally, layers that refine these predictions for accuracy.

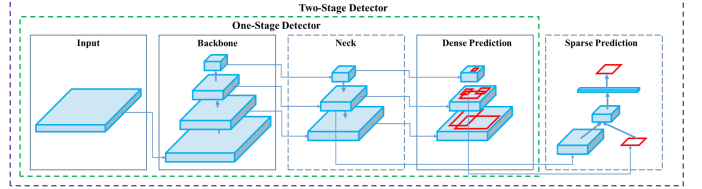


Fig. 4. Yolov4 Stages

Yolov4 utilizes a complex architecture that is highly optimized for speed and accuracy, incorporating new techniques such as Cross-Stage Partial connections (CSP), Mish activation, and the use of Self-Adversarial Training (SAT) and Mosaic data augmentation to enhance the model's performance. These features enable Yolov4 to achieve high detection precision even in real-time applications, making it an ideal choice for scenarios requiring rapid processing of visual information. The model operates by first resizing the input image for consistency and then sequentially processing the image through its layers, each designed to progressively refine the detection results [1].



Fig. 5. Object Detection Example Results

Figure 5 showcases an example of Yolov4's object detection capabilities from an experiment where the model was applied to a scene captured from a house window, overlooking a parking lot. In this demonstration, Yolov4 successfully identifies and outlines each car with a bounding box, illustrating its robustness in detecting various objects in diverse environments. The ability to detect and accurately box objects in complex scenes underscores the model's utility in practical applications, ranging from surveillance to autonomous driving aids.

The integration of Yolov4 into our system enhances its utility by providing detailed insights into the content captured by the OV7670 camera and processed by the FPGA. This



integration ensures that our setup is not just capturing and displaying images but also understanding and interpreting them to deliver meaningful analysis and actionable insights.

### III. RESULTS

The experiment conducted was a resounding success, demonstrating the robustness and effectiveness of the integrated system comprising the OV7670 camera module, FPGA processing, and the Yolov4 object detection model. Throughout the experimental sessions, we were able to capture, process, and analyze video data in real-time, validating the functionality and performance of each component within our setup.



Fig. 6. Experiment Session Photo: Before

Figure 6 presents a snapshot taken from the video output displayed on the VGA screen, representing the image quality and detail captured by the OV7670 camera and processed by the FPGA before being input into the Yolov4 model. This image serves as a baseline to appreciate the raw video feed's clarity and fidelity, which is crucial for the subsequent object detection phase.



Fig. 7. Experiment Session Photo: After

Figure 7 showcases the same frame post-processing by the Yolov4 model, where it successfully detects and delineates objects within the scene. In this instance, Yolov4 has accurately identified and outlined a person and a chair next to them with bounding boxes. This clear visualization of object

detection underscores the Yolov4 model's precision and the system's capability to handle complex image processing tasks effectively.

The results from this experiment highlight several key outcomes:

- **Image Acquisition and Display:** The OV7670 camera interfaced with the FPGA captured high-quality video data that was accurately rendered on a VGA display, ensuring that the visual information was preserved and transmitted without degradation.
- **Real-Time Processing:** The FPGA efficiently handled the real-time data processing requirements, demonstrating its capacity to perform necessary image adjustments swiftly and prepare the data for object detection.
- **Object Detection Accuracy:** The Yolov4 model applied to the processed images detected objects with high accuracy, proving its effectiveness in a real-world application scenario. The model not only recognized the objects but also provided precise bounding boxes around them, which is crucial for applications requiring detailed image analysis.

This successful demonstration of integrating hardware and software for real-time image capture, processing, and analysis provides valuable insights and a solid foundation for future projects aiming to deploy similar technologies in various applications.

### IV. CHALLENGES

Throughout the course of our project, while we achieved significant successes and demonstrated the capabilities of our integrated system, we also encountered specific challenges that provided insights for future improvements. The primary challenge revolved around the data transfer mechanism between the FPGA and the PC, which hosts the Yolov4 object detection model.

The ideal scenario for our system was to establish a seamless, real-time data transmission from the FPGA directly to the PC, enabling instantaneous processing and analysis of the video data by the Yolov4 model. However, the existing setup required manual intervention to transfer the processed data from the FPGA to the PC. This method, while effective for the scope of this project, introduced delays and reduced the potential for real-time analysis, which is crucial for applications that depend on immediate data processing, such as autonomous driving systems or real-time surveillance.

In future iterations of this project, we aim to explore both hardware and software solutions to overcome this limitation. Potential hardware solutions include employing faster communication interfaces such as USB 3.0 or Ethernet, which can facilitate higher data transfer rates. On the software side, implementing more efficient data encoding and compression techniques could significantly reduce the data load, making real-time transmission feasible even with existing hardware configurations.

## CONCLUSION

The project successfully demonstrated the capability of an integrated system that combines the real-time processing power of an FPGA with the advanced object detection capabilities of Yolov4. The experiment showcased how the OV7670 camera module could capture video data, which was then processed by the FPGA to display on a VGA screen, and subsequently analyzed by Yolov4 to detect and annotate objects with high precision. While the system effectively demonstrated the potential of combining these technologies, it also highlighted significant challenges, particularly in the direct data transmission from the FPGA to the PC, which necessitated manual data transfer.

Future work will focus on addressing these challenges by exploring more efficient data transmission methods that could enable truly real-time processing and analysis. Potential improvements include integrating faster communication interfaces or optimizing data compression techniques to enhance the system's performance without compromising the real-time capabilities required for broader applications such as autonomous systems and advanced surveillance solutions. The success of this project lays a solid foundation for further research and development in integrating embedded systems with artificial intelligence for real-time image processing and object detection.

## REFERENCES

- [1] A. Bochkovskiy, C.-Y. Wang, and H.-Y. M. Liao, "Yolov4: Optimal speed and accuracy of object detection," 2020.
- [2] H. Liu, K. Fan, Q. Ouyang, and N. Li, "Real-time small drones detection based on pruned yolov4," *Sensors*, vol. 21, no. 10, 2021. [Online]. Available: <https://www.mdpi.com/1424-8220/21/10/3374>
- [3] C. Dewi, R.-C. Chen, Y.-T. Liu, X. Jiang, and K. D. Hartomo, "Yolo v4 for advanced traffic sign recognition with synthetic training data generated by various gan," *IEEE Access*, vol. 9, pp. 97 228–97 242, 2021.
- [4] R. Chen, T. Wu, Y. Zheng, and M. Ling, "Mlof: Machine learning accelerators for the low-cost fpga platforms," *Applied Sciences*, vol. 12, no. 1, 2022. [Online]. Available: <https://www.mdpi.com/2076-3417/12/1/89>
- [5] D. Patel, R. Parmar, A. Desai, and S. Sheth, "Gesture recognition using fpga and ov7670 camera," in *2017 International Conference on Inventive Systems and Control (ICISC)*, 2017, pp. 1–4.
- [6] A. E. Dodi, A. H. Wahyudi, Kurdianto, N. W. Jatmiko, M. F. Lailiyul, and W. Widada, "Fpga displays real-time video camera on vga monitor," in *2022 11th Electrical Power, Electronics, Communications, Controls and Informatics Seminar (EECCIS)*, 2022, pp. 129–132.