

Comprehensive Performance Evaluation of YOLOv12, YOLO11, YOLOv10, YOLOv9 and YOLOv8 on Detecting and Counting Fruitlet in Complex Orchard Environments

Ranjan Sapkota^{a,*}, Zhichao Meng^c, Martin Churuvija^b, Xiaoqiang Du^c, Zenghong Mab^c, Manoj Karkee^b

^aCornell University, Ithaca, Ithaca, 14850, NY, USA

^bCenter for Precision & Automated Agricultural Systems, 24106 N Bunn Rd, Prosser, 99350, WA, USA

^cSchool of Mechanical Engineering, Zhejiang Sci-Tech University, Hangzhou, 310018, China

Abstract

This study systematically performed an extensive real-world evaluation of the performances of all configurations of YOLOv8, YOLOv9, YOLOv10, YOLO11(or YOLOv11), and YOLOv12 object detection algorithms in terms of precision, recall, mean Average Precision at 50% Intersection over Union (mAP@50), and computational speeds including pre-processing, inference, and post-processing times immature green apple (or fruitlet) detection in commercial orchards. Additionally, this research performed and validated in-field counting of the fruitlets using an iPhone and machine vision sensors. Among the configurations, YOLOv12l recorded the highest recall rate at 0.90, compared to all other configurations of YOLO models. Likewise, YOLOv10x achieved the highest precision score of 0.908, while YOLOv9 Gelan-c attained a precision of 0.903. Analysis of mAP@0.50 revealed that YOLOv9 Gelan-base and YOLOv9 Gelan-e reached peak scores of 0.935, with YOLO11s and YOLOv12l following closely at 0.933 and 0.931, respectively. For counting validation using images captured with an iPhone 14 Pro, the YOLO11n configuration demonstrated outstanding accuracy, recording RMSE values of 4.51 for Honeycrisp, 4.59 for Cosmic Crisp, 4.83 for Scilate, and 4.96 for Scifresh; corresponding MAE values were 4.07, 3.98, 7.73, and 3.85. Similar performance trends were observed with RGB-D sensor data. Moreover, sensor-specific training on Intel Realsense data significantly enhanced model performance. YOLOv11n achieved highest inference speed of 2.4 ms, outperforming YOLOv8n (4.1 ms), YOLOv9 Gelan-s (11.5 ms), YOLOv10n (5.5 ms), and YOLOv12n (4.6 ms), underscoring its suitability for real-time object detection applications.

1. Introduction

Object detection is a fundamental task in computer vision that enables automated systems to identify, classify, and locate objects within images or video streams Zou et al. (2023); Zhao et al. (2019). This capability is critical for a wide range of applications, from autonomous driving and surveillance to robotics, healthcare, and agricultural automation Sapkota et al. (2024d). By recognizing objects in real time, detection systems can drive rapid decision-making and improve operational safety and efficiency Wu et al. (2017). In complex, real-world environments, systems must contend with challenges such as variations in lighting, scale, occlusions, and cluttered backgrounds. Overcoming these hurdles is essential for developing robust and reliable detection algorithms that work effectively under diverse conditions. In addition, automated object detection reduces human intervention, cuts operational costs, and enhances precision, which in turn accelerates technological innovation Haleem et al. (2021); Manakitsa et al. (2024).

1.1. Evolution of Object Detection and the Emergence of YOLO

The evolution of object detection (as depicted in 1) has been driven by the quest for greater accuracy, speed, and robustness

in real-world scenarios. Early methods relied on handcrafted features and classical machine learning techniques. Techniques such as sliding window methods Lampert et al. (2008) and the Viola–Jones detector Castrillón et al. (2011) used features like Haar-like descriptors, histograms of oriented gradients, and local binary patterns to capture essential visual information. Although these methods laid the foundation for automated detection, they were limited by the reliance on manual feature engineering and exhaustive search processes, which made them computationally intensive and less adaptable to the intricacies of natural scenes. The advent of deep learning, particularly the introduction of convolutional neural networks (CNNs), revolutionized the field by enabling automated, hierarchical feature extraction Khan et al. (2020). Models such as R-CNN, Fast R-CNN, and Faster R-CNN significantly improved detection accuracy by learning complex representations directly from data. Simultaneously, the Single Shot MultiBox Detector (SSD) streamlined the process by integrating region proposal and detection in a single network pass, setting the stage for real-time performance Liu et al. (2016). Among these transformative developments, the YOLO (You Only Look Once) series emerged as a breakthrough approach by unifying detection and classification into one efficient, end-to-end framework, dramatically reducing computation time while maintaining competitive accuracy Redmon et al. (2016).

The initial iterations of the YOLO series marked a departure

*Ranjan Sapkota

Email address: ranjan.sapkota@example.com (Manoj Karkee)

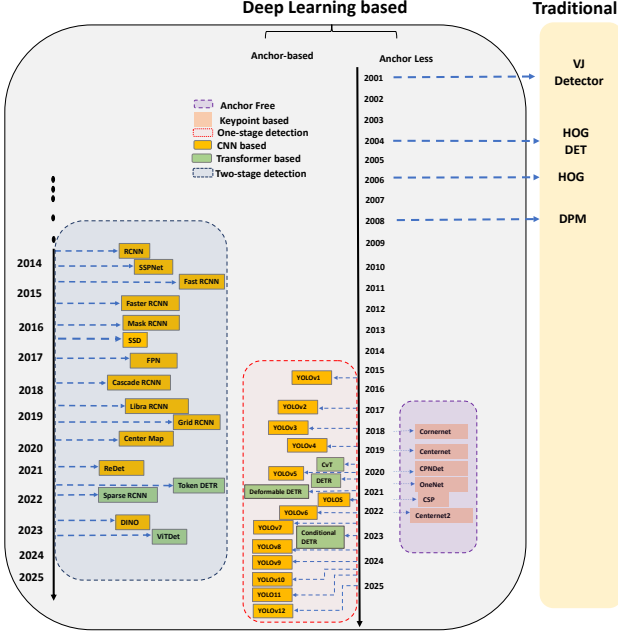


Figure 1: Timeline diagram depicting the evolution of YOLO algorithms from YOLOv1’s grid-based detection to YOLOv12’s attention-centric architecture Sapkota et al. (2024d)

from traditional multi-stage detectors. YOLOv1 introduced a grid-based approach that enabled the model to simultaneously predict bounding boxes and class probabilities, thereby establishing a new paradigm in real-time detection Redmon et al. (2016). Building on this success, YOLOv2 improved image resolution and expanded the range of detectable classes, enhancing both accuracy and versatility. YOLOv3 further refined the approach by incorporating multi-scale predictions and a deeper network architecture, which enabled more effective detection of small objects and complex features Redmon and Farhadi (2018). As the series progressed, YOLOv4 and YOLOv5 introduced architectural innovations such as Cross-Stage Partial (CSP) networks and advanced data augmentation techniques that enhanced feature representation and inference speed Bochkovskiy et al. (2020); Jocher et al. (2020). Subsequent versions like YOLOv6 and YOLOv7 continued to optimize the balance between accuracy and computational efficiency through innovations like dynamic label assignment and refined architectural blocks Li et al. (2022); Wang et al. (2023). Together, these early YOLO models laid a robust foundation for real-time object detection by addressing many limitations of earlier methods and setting a high standard for speed and precision.

1.2. Recent YOLO Iterations (YOLOv8 to YOLOv12)

- **YOLOv8 Architecture:** YOLOv8 builds on the advances of its predecessors by adopting a more efficient and anchor-free design. Its backbone leverages a refined version of CSP-based modules to extract hierarchical features effectively, while its neck integrates a streamlined Feature Pyramid Network that enhances multi-scale fea-

ture representation. The decoupled head architecture in YOLOv8 separates the objectness, classification, and regression tasks, thereby improving precision without incurring significant computational overhead. This design not only simplifies the detection pipeline but also ensures that the network is more adaptable to varying object sizes and aspect ratios. The overall architecture demonstrates a balance between speed and accuracy, making YOLOv8 suitable for real-time applications, as illustrated in Figure 2a.

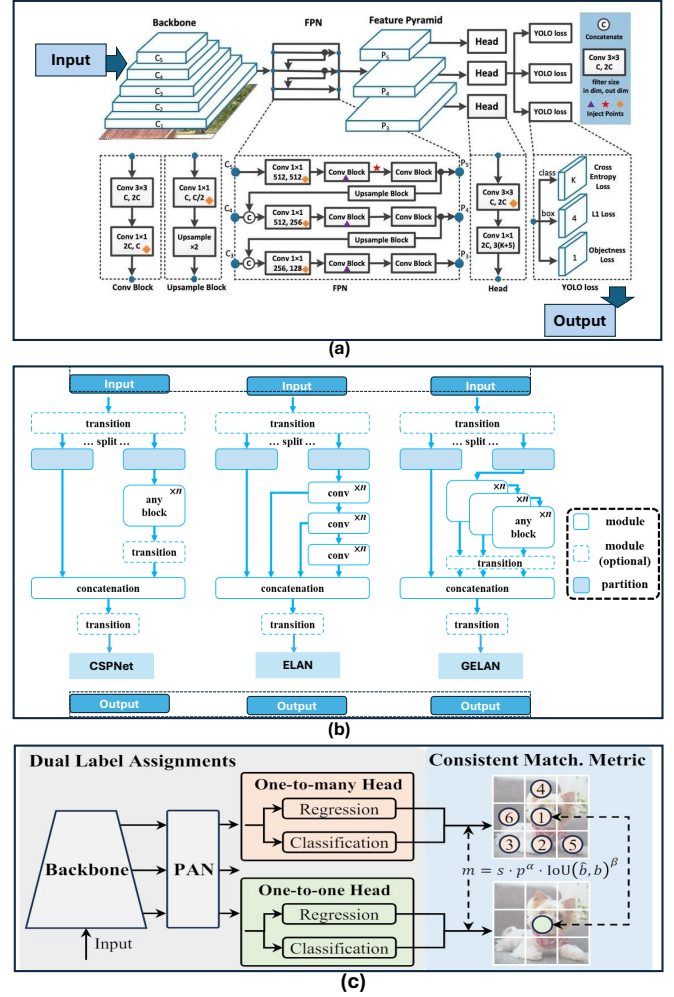


Figure 2: YOLOv8, YOLOv9 and YOLOv10 Architecture Diagram
(a) YOLOv8 employs a CSP-based backbone, anchor-free decoupled head, and streamlined FPN for efficient multi-scale feature extraction. (b) YOLOv9 integrates programmable gradient information with GELAN for robust feature aggregation. (c) YOLOv10 adopts a dual assignment strategy, lightweight heads, and spatial-channel decoupled downsampling to further boost overall speed and accuracy significantly.

- **YOLOv9 Architecture:** YOLOv9 introduces significant improvements through the incorporation of programmable gradient information (PGI) and the Generalized Effi-

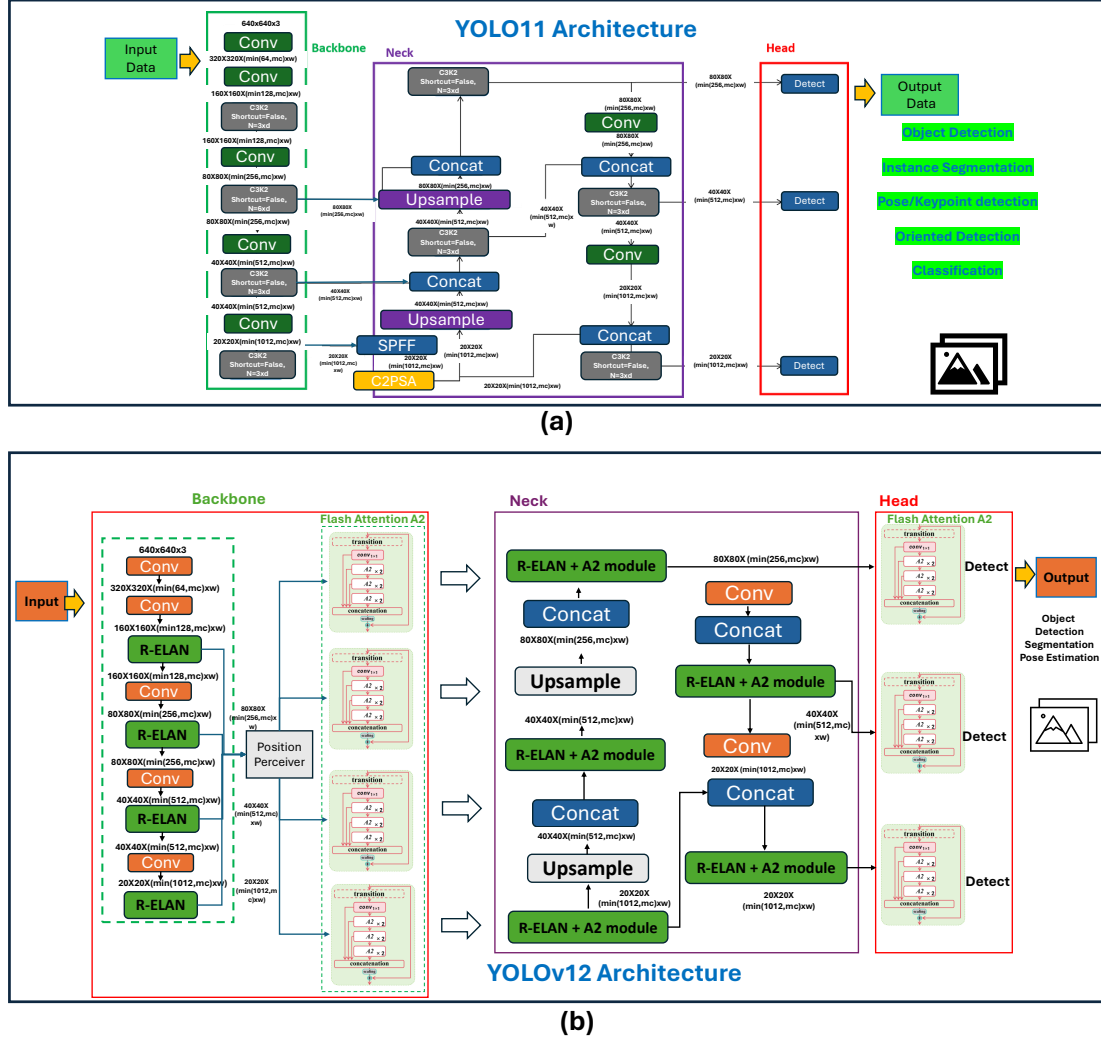


Figure 3: YOLO11 and YOLOv12 Architecture Diagram

(a) YOLOv11 features a refined architecture with advanced C3k2 blocks, SPPF, and C2PSA modules that significantly enhance multi-scale feature extraction and spatial attention for improved detection accuracy. (b) In contrast, YOLOv12 advances this design further with an attention-centric framework that integrates Area Attention modules and R-ELAN blocks to optimize feature fusion, dramatically boosting inference speed for state-of-the-art real-time object detection.

cient Layer Aggregation Network (GELAN). The PGI mechanism ensures that robust and reliable gradients are maintained across the network, which aids in efficient weight updates and better convergence Wang et al. (2024). GELAN further enhances the model by providing flexible and efficient multi-scale feature aggregation, which is essential for accurately detecting objects in complex scenes. This architectural refinement allows YOLOv9 to address the challenges of information loss in deep networks and to maintain high detection accuracy, particularly for smaller objects. Detailed architectural components of YOLOv9 are depicted in Figure 2b.

- **YOLOv10 Architecture:** Building upon the innovations of YOLOv9, YOLOv10 further streamlines the object detection pipeline by eliminating the need for non-maximum suppression through a novel dual assignment strategy

Wang et al. (2025). This approach reduces redundancy by directly assigning labels during training via one-to-many and one-to-one strategies, significantly cutting down on inference time. YOLOv10 also introduces lightweight classification heads and incorporates spatial-channel decoupled downsampling, which minimizes information loss during feature extraction. Additionally, the rank-guided block design optimizes parameter utilization, ensuring that the network remains efficient even under the demands of real-time detection tasks. The comprehensive architectural layout of YOLOv10 is shown in Figure 2c.

- **YOLO11 Architecture:** YOLOv11 further refines the object detection process by focusing on enhanced feature extraction and robust detection capabilities. It introduces the C3k2 block, a replacement for the previously used C2f block, that significantly improves gradient flow and com-

putational efficiency Sapkota et al. (2024e). YOLOv11 also integrates the Spatial Pyramid Pooling-Fast (SPPF) module to better capture multi-scale contextual information and employs the Convolutional block with Parallel Spatial Attention (C2PSA) to refine spatial feature representation. These enhancements collectively allow YOLOv11 to perform more accurate detection, particularly in scenarios involving occlusion and complex backgrounds Sapkota and Karkee (2024c). The architecture, highlighting these critical components, is detailed in Figure 31.

- **YOLOv12 Architecture:** Representing the state-of-the-art in the YOLO series, YOLOv12 adopts an attention-centric design that pushes the boundaries of real-time object detection Tian et al. (2025). Central to its innovation is the introduction of the Area Attention (A^2) module, which dynamically adjusts the receptive field to capture both global and local contextual cues with minimal computational expense. In parallel, the Residual Efficient Layer Aggregation Network (R-ELAN) refines feature aggregation by incorporating residual connections that ensure stable gradient flow and efficient information fusion. YOLOv12 also leverages additional optimizations, including FlashAttention and adaptive MLP ratios, to further enhance inference speed while maintaining high detection accuracy. These cumulative architectural improvements enable YOLOv12 to outperform previous iterations across multiple metrics, as comprehensively illustrated in Figure 3b.

1.3. Objectives

This study provides a comprehensive evaluation of the latest YOLO versions: YOLOv12, YOLO11, YOLOv10, YOLOv9, and YOLOv8, targeting fruitlet detection in commercial apple orchards by examining 26 configurations across these models and utilizing a commercial orchard's dataset of RGB images from an iPhone 14 across four apple varieties: Scifresh, Scilate, Honeycrisp, and Cosmic Crisp. Each image was annotated with a complete count of visible and occluded apples directly observed in the field, providing a comprehensive dataset for validation. The study also validates the top-performing models using machine vision sensor. The specific contributions of this study are:

- **Model Training and Configuration:** Comprehensive evaluation of the latest YOLO object detection models implemented across 26 configurations: YOLOv8 (5 configurations), YOLOv9 (6 configurations), YOLOv10 (6 configurations), specifically optimized for detecting green fruitlets in commercial apple orchards, YOLO11 (5 configurations), and YOLOv12(5 Configurations).
- **Comprehensive Metrics** Detailed examination of detection precision metrics, computational efficiency, and processing speeds at 3 steps (preprocess, inference and post-process) of the deep learning workflow.

- **Validation Across Varieties:** In-field counting accuracy validation using four apple varieties not included in the training set, to test generalizability and robustness of the models under varied agricultural conditions.
- **Integration of Smartphone Technology:** Utilization of high-resolution RGB images from an iPhone 14 Pro Max for adaptability of advanced smartphone imaging in field setting.

2. Methods

This study was conducted across commercial apple orchards in Prossers and Naches, Washington State, USA, focusing on evaluating datasets of apple fruitlets before thinning from four varieties: Scifresh, Scilate, Cosmic Crisp, and Honeycrisp as described by Figure 4. RGB images were acquired using a machine vision sensor IntelRealsense 435i (Intel Corporation, California, USA). These images were then manually labeled to prepare consistent datasets for training all configurations of YOLOv8, YOLOv9, and YOLOv10 models as illustrated by Figure 4. A total of 22 model configurations were examined: 5 for YOLOv8, 6 for YOLOv9, 6 for YOLOv10, 5 for YOLO11, and 5 for YOLOv12 (Figure 4c) each trained under standardized hyperparameter settings on the same computational system to ensure consistency and comparability throughout the process. The trained models' performance was then evaluated using the prepared datasets, followed by in-field counting validation using additional RGB images captured with an Apple iPhone 14 Pro Max smartphone (Apple Inc., California, USA). This validation process aimed to assess each model's fruit counting accuracy against manually counted ground truths, which included scenarios with occluded apples. Furthermore, this validation step was designed to rigorously test the highest-performing models in terms of speed and accuracy using images from different sensors. The details of this comprehensive methodology are illustrated in Figure 4, providing a visual reference for the experimental setup and data collection approach employed in this study.

2.1. Study Site and Data Acquisition

The image data collection for this study was conducted at Allan Brothers Orchard in Prosser, Washington State, USA, focusing on datasets of immature green apples for training models. RGB images were acquired using an Intel RealSense 435i camera during the month of June 2024. For validation, additional datasets were collected from different locations and times: Cosmic Crisp apple images from the ROZA experiment station at WSU IAREC in Prosser, Honeycrisp apple images from an orchard in Naches owned by Allan Brothers Fruit Company. Furthermore, All images were captured prior to the fruitlet thinning process conducted by orchard workers. The details of each machine vision sensor used in data collection of this study are:

- **Intel RealSense D435i:** This camera features a 2-megapixel RGB sensor capable of capturing high-quality images. Operating at a resolution of 1280×720 pixels,

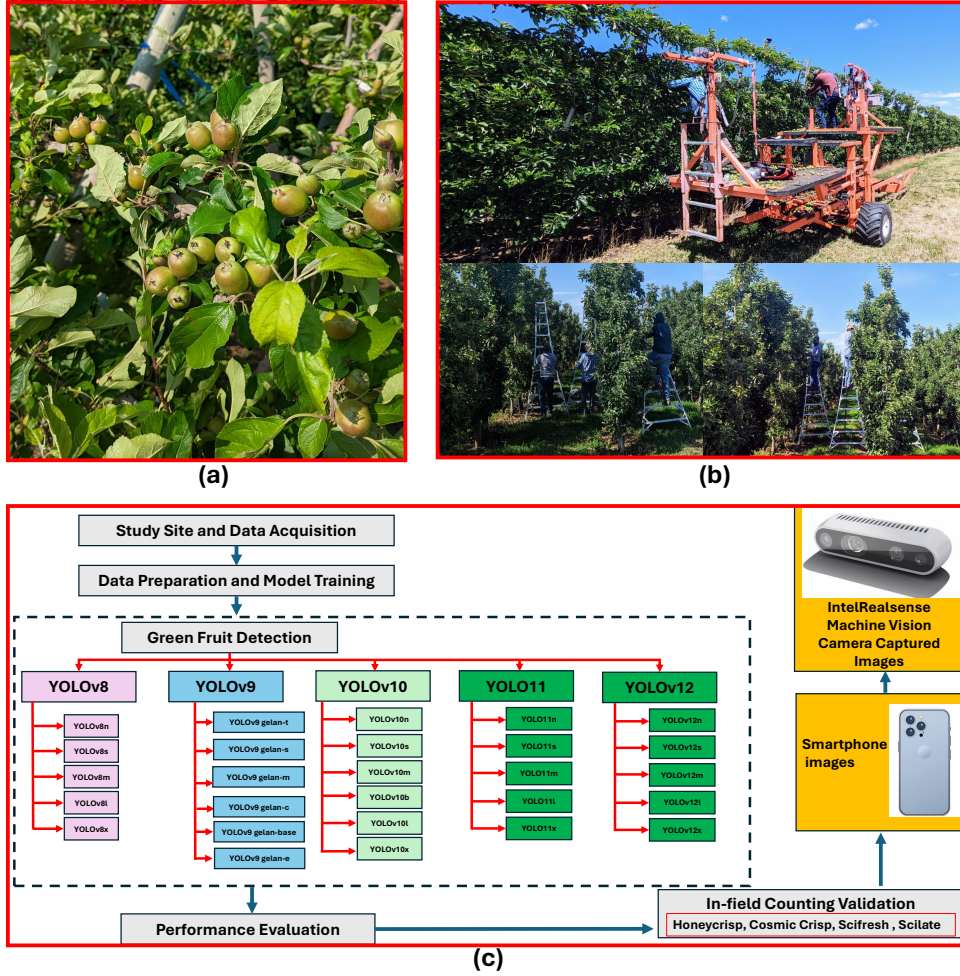


Figure 4: YOLO11 and YOLOv12 Architecture Diagram

Fruitlet thinning in commercial orchards: (a) High-density clusters of apple fruitlets on a Scilate apple tree during the peak thinning period in June 2022 (A commercial orchard in Prosser, WA), illustrating the typical overcropping seen in commercial orchards; (b) Top: Laborers utilizing an height-adjustable platform for efficiently thinning fruitlet in various parts of tree canopies; b) Bottom: A worker in a Scifresh apple orchard manually thinning excess fruitlets using an aluminum ladder, a common practice that highlights the labor-intensive nature of this crucial agricultural task. (c) Flow diagram of this study of comparison of YOLOv8, YOLOv9, YOLOv10, YOLO11, and YOLOv12 illustrating the study's methodology, including data collection, model training, and validation across multiple sensors and apple varieties in commercial orchards

it provides a comprehensive view with a 69.4° horizontal and 42.5° vertical field-of-view, ensuring broad coverage in various environments. The Intel RealSense D435i is designed to be compact and lightweight, making it highly effective for capturing RGB data in diverse settings.

- **Apple Iphone14 Pro Max:** The Apple iPhone 14 Pro Max features an advanced 48-megapixel RGB camera featuring a 24 mm lens with an f/1.78 aperture and second-generation sensor-shift optical image stabilization, ensuring high-resolution image capture with enhanced clarity. The camera's field of view spans 120° horizontally and 90° vertically, supported by a seven-element lens for minimized optical aberrations. It offers up to 15x digital zoom and 4K video recording capabilities. Additionally, the integration of Photonic Engine and Smart HDR 4 technology

optimizes performance in low-light conditions and dynamic range, making it ideal for detailed visual data collection in varying lighting environments.

2.2. Data Preparation and Model Training

A total of 1,147 images were manually annotated with bounding boxes using the online labeling platform provided by Roboflow. These images were allocated into training and validation at the ratio of 8:2 respectively, facilitated by Roboflow's distribution tools. No image preprocessing steps were undertaken in this study, as the objective was not to enhance any specific model but rather to compare and evaluate the performance of these models using raw data collected in natural orchard settings.

The computational analyses for this study were conducted on a high-performance workstation equipped with an Intel

Xeon(R) W-2155 CPU, featuring a base clock speed of 3.30 GHz across 20 cores. This setup provided substantial processing power necessary for handling intensive data processing tasks. The workstation was also outfitted with NVIDIA Corporation GP102 [TITAN Xp] graphics cards, enhancing its ability to perform complex image processing and machine learning tasks efficiently. The system included a substantial storage capacity of 7.0 TB, facilitating extensive data management and analysis. It operated under Ubuntu 20.04.6 LTS, a robust and stable 64-bit operating system, using GNOME version 3.36.8 for its graphical interface and X11 as its windowing system.

The analysis of YOLOv8, YOLOv9, YOLOv10, YOLO11, and YOLOv12 encompassed a total of 27 configurations (five for YOLOv8, six for YOLOv9, six for YOLOv10, five for YOLOv11, and five for YOLOv12). Each model configuration was rigorously tested for precision, recall, mAP@0.5, and image processing speed. For the training of all models, a consistent configuration was rigorously maintained to ensure uniformity and comparability across experiments. Each model was trained for 700 epochs, reflecting a substantial duration to adequately learn and adapt to the dataset’s complexities. The batch size was set at 8, optimizing the balance between memory usage and processing speed. Images were resized to a uniform resolution of 640x640 pixels to standardize the input data size across all models. The Stochastic Gradient Descent (SGD) optimizer was employed to update model weights effectively, favored for its efficiency in handling large-scale and complex data. Additionally, the configuration threshold for confidence was set at 0.25, and the Intersection over Union (IoU) threshold was maintained at 0.7, criteria chosen to optimize the balance between precision and recall during object detection tasks.

The hyperparameter settings outlined in table 1 provide a detailed framework for optimizing the training of YOLO models in the study. Key parameters such as the initial and final learning rates were set at 0.01, facilitating a controlled adjustment of learning throughout the training process. Momentum was maintained at 0.937 to ensure consistent updates across epochs, while a minimal weight decay of 0.0005 helped prevent overfitting. The training initiated with a warmup phase spanning 3 epochs to stabilize the learning parameters early in the training. Loss adjustments were specifically tuned, with box loss, class loss, and definition loss set at 7.5, 0.5, and 1.5 respectively, to balance the contributions of different components of the loss function. Image augmentation techniques such as hue, saturation, and value adjustments were precisely defined to enhance model robustness under varied lighting conditions typically found in orchard environments. Specific settings for geometric transformations like rotation, translation, and scaling were employed to simulate different orientations and sizes of objects, crucial for improving the model’s ability to generalize across diverse scenarios. The table 1 also highlights the use of flipping and mosaic data augmentation to further enrich the training dataset, ensuring comprehensive exposure to potential real-world variations.

Table 1: Hyperparameter Settings for YOLOv8, YOLOv9, YOLOv10, YOLOv11 and YOLOv12 used in training YOLO models for fruitlet detection in commercial apple orchards

Hyperparameter	Value	Description
Initial Learning Rate (lr0)	0.01	Sets the starting learning rate.
Final Learning Rate (lrf)	0.01	Determines the learning rate at the end of training.
Momentum	0.937	Controls the momentum for the SGD optimizer.
Weight Decay	0.0005	Helps in regularizing and preventing overfitting.
Warmup Epochs	3.0	Number of initial epochs for learning rate stabilization.
Box Loss Gain (box)	7.5	Weight of the bounding box loss component.
Class Loss Gain (cls)	0.5	Weight of the class prediction loss component.
Definition Loss Gain (dfl)	1.5	Weight of the definition loss component.

2.3. Performance Evaluation

To evaluate the detection capabilities of each configuration of YOLOv8, YOLOv9, YOLOv10, YOLO11, and YOLOv12 the metrics of precision, recall, and mean Average Precision at Intersection over Union (IoU) threshold of 0.50 (mAP@50) were employed. Precision was calculated as the ratio of true positive detections to the total predicted positives, given by equation 1, where TP denotes true positives and FP denotes false positives. Recall measured the ratio of true positive detections to the actual positives, formulated as equation 2 where FN represents false negatives. The mAP@50 was determined by averaging the precision across all recall levels for an IoU \geq 0.50. Additionally, the image processing speed for each model was analyzed and compared across three categories: preprocessing, inference, and postprocessing. These evaluations were systematically conducted across 22 configurations of the YOLO models: YOLOv8m, YOLOv8s, YOLOv8l, YOLOv8x, YOLOv8c for YOLOv8; YOLOv9 Gelan-e, YOLOv9 Gelan-c, YOLOv9 Gelan-s, YOLOv9 Gelan-t, YOLOv9 Gelan-m, and YOLOv9 Gelan for YOLOv9; YOLOv10m, YOLOv10s, YOLOv10l, YOLOv10x, YOLOv10c, YOLOv10d for YOLOv10, and YOLO11n, YOLO11s, YOLO11m, YOLO11l, and YOLO11x for YOLO11, and YOLOv12n, YOLOv12s, YOLOv12m, YOLOv12l, and YOLOv12x for YOLOv12 to evaluate their efficiency in detecting fruitlets in commercial orchards.

$$Precision = \frac{TP}{TP + FP} \quad (1)$$

$$Recall = \frac{TP}{TP + FN} \quad (2)$$

In addition to accuracy metrics, the model’s complexity and computational demand were evaluated by examining the number of convolutional layers, total parameters, and GFLOPs:

$$Parameters_{Model} = \text{Total trainable weights and biases} \quad (3)$$

$$GFLOPs = \frac{\text{Total floating-point operations}}{10^9} \text{ per image} \quad (4)$$

$$Layers_{Convolutional} = \text{Total number of convolutional layers in the model} \quad (5)$$

Additionally, the assessment involved analyzing preprocessing, inference, and postprocessing speeds for each model configuration, as these metrics are critical for real-time object detection systems. Preprocessing speed determines how quickly a model can prepare images for detection, inference speed measures the time taken to identify objects within images, and postprocessing speed reflects how swiftly the model finalizes the outputs after detection. Each of these stages is essential for efficient operation in agricultural applications, where timely and accurate detection can significantly impact decision-making and resource management. The computational efficiency of these models directly influences their practical utility in automated fruit detection systems, making this evaluation crucial for advancing agricultural technology solutions.

2.4. In-Field Counting Validation

Upon evaluating the performance metrics of each YOLOv8, YOLOv9, YOLOv10, YOLOv11, and YOLOv12 configuration for detecting fruitlets, the most effective model from each version was selected based on the highest accuracy achieved at mAP@0.5. These top-performing models from YOLOv8, YOLOv9, YOLOv10, YOLOv11, YOLOv12 were further validated to assess their detection capabilities in a real commercial orchard setting across four distinct apple varieties in Washington State. This validation process utilized images collected both by smartphone and machine vision sensor. Initially, images of Scifresh, Scilate, Cosmic Crisp, and Honeycrisp apples were captured using the smartphone. A total of 128 images, 32 per variety, were analyzed to evaluate the models' counting accuracy before the thinning process.

Subsequently, an additional set of images was used for further validation. This included 32 images of Scifresh apples taken with an IntelRealsense machine vision camera. Notably, while the IntelRealsense camera images of Scifresh and Scilate apples served as the training dataset, validation was performed on varieties: Cosmic Crisp and Honeycrisp, which were not included in the training phase. For each image, a count of all visible and occluded apples fruitlets were manually performed directly in the field. This comprehensive manual counting provided a precise ground truth for each image sample across the different apple varieties, ensuring robust validation of the models' performance.

In our study assessing the performance of YOLO models for detecting fruitlets in commercial orchards, Root Mean Square Error (RMSE) and Mean Absolute Error (MAE) were employed as crucial metrics to quantify the accuracy of fruit counts predicted by the models compared to manually counted ground truths.

RMSE was calculated using Equation 3:

$$RMSE = \sqrt{\frac{1}{n} \sum_{i=1}^n (\text{predicted}_i - \text{actual}_i)^2} \quad (6)$$

Here, predicted_i denotes the number of fruits counted by the model, actual_i represents the manually counted fruits for each image, and n is the total number of images. This measure computes the square root of the average of the squared differences between the predicted and actual counts, thereby emphasizing larger errors and highlighting significant deviations in model performance.

Likewise, MAE was determined using Equation 4:

$$MAE = \frac{1}{n} \sum_{i=1}^n |\text{predicted}_i - \text{actual}_i| \quad (7)$$

In this equation, predicted_i and actual_i retain their previous definitions, with MAE calculating the average of the absolute differences between the predicted and actual counts. This metric, being less sensitive to large errors than RMSE, provides a straightforward indication of the average error magnitude per image, offering an intuitive measure of prediction accuracy across the dataset.

3. Results and Discussion

Figure 5 illustrated the detection outcomes of the fastest variants (nano) of YOLOv8, YOLOv9, YOLOv10, YOLOv11, and YOLOv12 on two representative images captured in commercial apple orchards. In the left image, the original photograph was displayed at the top, with the predictions from each model arranged sequentially below it. In this example, a green arrow on the upper portion of the image indicated an apple fruitlet partially visible in a complex scene, while another green arrow on the lower portion pointed to an apple fruitlet under severe occlusion. The green arrows denoted instances where the model correctly detected the fruitlets, whereas a red arrow at the same spatial location highlighted a missed detection. Specifically, YOLOv8 correctly identified the partially visible fruitlet at the top but failed to detect the heavily occluded fruitlet at the bottom. In contrast, YOLOv9 missed the top fruitlet yet successfully detected the bottom one. YOLOv10, however, failed to detect either fruitlet in this scenario. YOLOv11 detected the top fruitlet but again failed on the bottom one, and YOLOv12 also managed to identify only the top fruitlet, indicating a variable detection capability among the models.

In the right image, the original view was again presented at the top, followed by the predictions from the different YOLO models. A green dotted circle marked a region containing three apple fruitlets, with one fruitlet partially obscured by two others in the foreground. In this region, YOLOv8 detected only two fruitlets, as evidenced by the dotted circle, while YOLOv9 and YOLOv10 similarly failed to detect all three. Notably, YOLOv11 successfully identified all fruitlets within the occluded region, demonstrating improved robustness under challenging conditions. Conversely, YOLOv12, despite being the latest iteration, failed to detect the partially occluded fruitlet in that area. These results underscored the significant variability in detection performance across the YOLO variants, particularly in complex and occluded scenes. Overall, the findings revealed that while YOLOv11 exhibited superior detection performance

in certain challenging scenarios, models such as YOLOv10 and YOLOv12 still faced limitations. This variability emphasizes the need for further refinements in model architecture and training strategies to enhance robustness and accuracy in real-world agricultural applications.

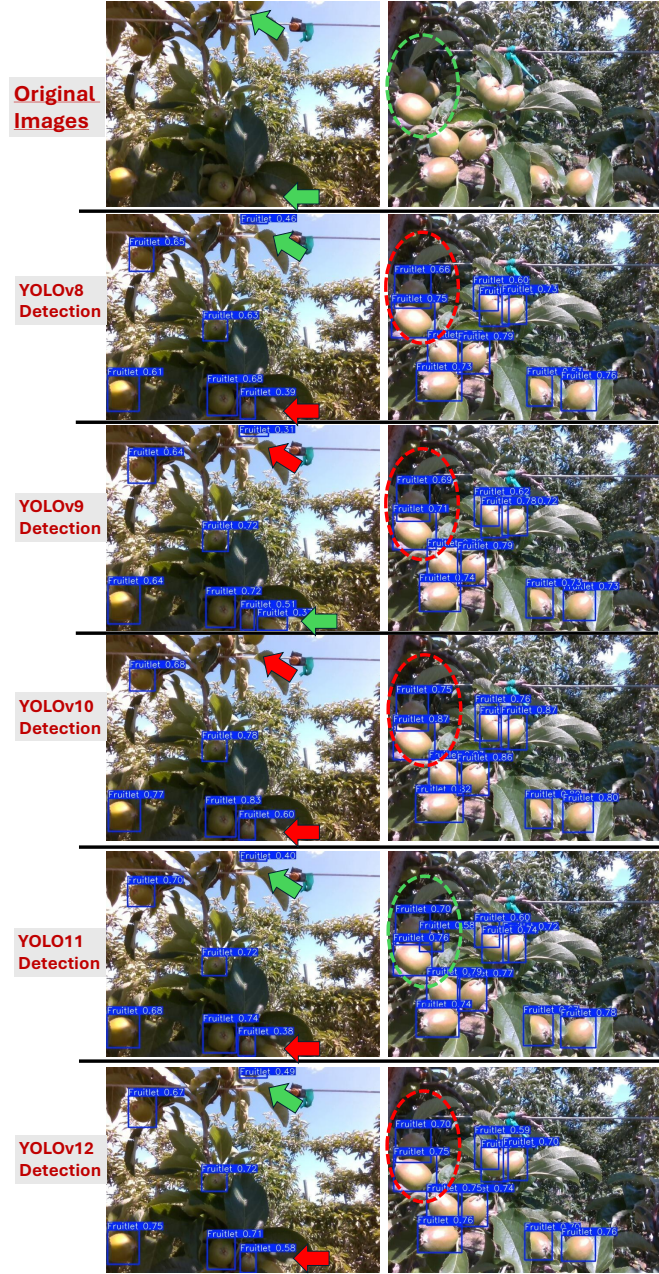


Figure 5: Illustration of Object Detection (Green Fruitlet) Results of YOLOv12, YOLOv11, YOLOv10, YOLOv9, YOLOv8

Detection results of YOLOv8, YOLOv9, YOLOv10, YOLOv11, and YOLOv12 are shown. The figure presents original images (top) and predictions (bottom) with green and red arrows indicating correct and missed detections.

3.1. Assessment of Detection Accuracy: Precision and Recall Metrics

Among all the models compared in this study across the configurations of YOLOv8 to YOLOv12, YOLOv10x achieved the highest precision score of 0.908, while YOLOv12l attained the highest recall score of 0.900. Within each model family, performance differences were evident. In the YOLOv12 family—which comprised four fully trained variants (YOLOv12n, YOLOv12s, YOLOv12m, and YOLOv12l, noting that YOLOv12x could not be completely trained due to memory constraints)—YOLOv12l demonstrated the highest precision at 0.87. In the YOLOv11 family, among the evaluated variants, YOLOv11l achieved the highest precision with a score of 0.899. YOLOv10x, as noted earlier, delivered the highest precision among all the model configurations compared in this study. Similarly, within the YOLOv9 series, the variant YOLOv9 gelan-c achieved the highest precision score of 0.903 among the six configurations tested, while in the YOLOv8 family—comprising five configurations—YOLOv8m obtained the highest precision, reaching a value of 0.897. These findings reflect heterogeneous performance across different YOLO architectures and underscore the critical importance of model selection and training optimization for specific detection tasks. Detailed precision and recall values for all the evaluated models are presented in Table 2, providing a comprehensive overview of detection performance across the various YOLO families. Overall, these results highlight the trade-offs between precision and recall inherent in different model architectures and emphasize the need for further optimization to achieve balanced performance in real-time object detection applications.

3.2. Evaluation of Detection Consistency: Mean Average Precision at IoU=0.50

The assessment of mean Average Precision at an Intersection over Union (IoU) threshold of 0.50 (mAP@0.50) across various YOLO model configurations in Figure 6 provides profound insights into their efficacy in detecting fruitlets within agricultural settings. Among all evaluated configurations, YOLOv9 models, particularly YOLOv9 Gelan-e and YOLOv9 Gelan-base, stand out with the highest mAP@0.50 scores of 0.935 both, respectively. However, all configurations of YOLOv9 achieved higher mAP@50 as depicted in Figure 10. These scores not only exceed those achieved by all configurations of YOLOv8, YOLOv10, YOLOv11 and YOLOv12 but also underscore YOLOv9’s superior precision in object detection tasks.

The newly included YOLOv12 and YOLOv11 series introduces impressive scores as well: YOLOv12l records an mAP@50 of 0.931, YOLOv12m at 0.928, YOLOv11n records an mAP@0.50 of 0.926, YOLOv11s at 0.933, YOLOv11m at 0.924, YOLOv11l at 0.932, and YOLOv11x at 0.922, indicating enhanced detection capabilities across these newer configurations. Within the YOLOv10 series, YOLOv10n leads with an mAP@0.50 of 0.921, closely followed by YOLOv10b and YOLOv10-M, which record scores of 0.919 and 0.917 respectively. Despite their strong performance, these figures remain slightly below the peak values presented by YOLOv9,

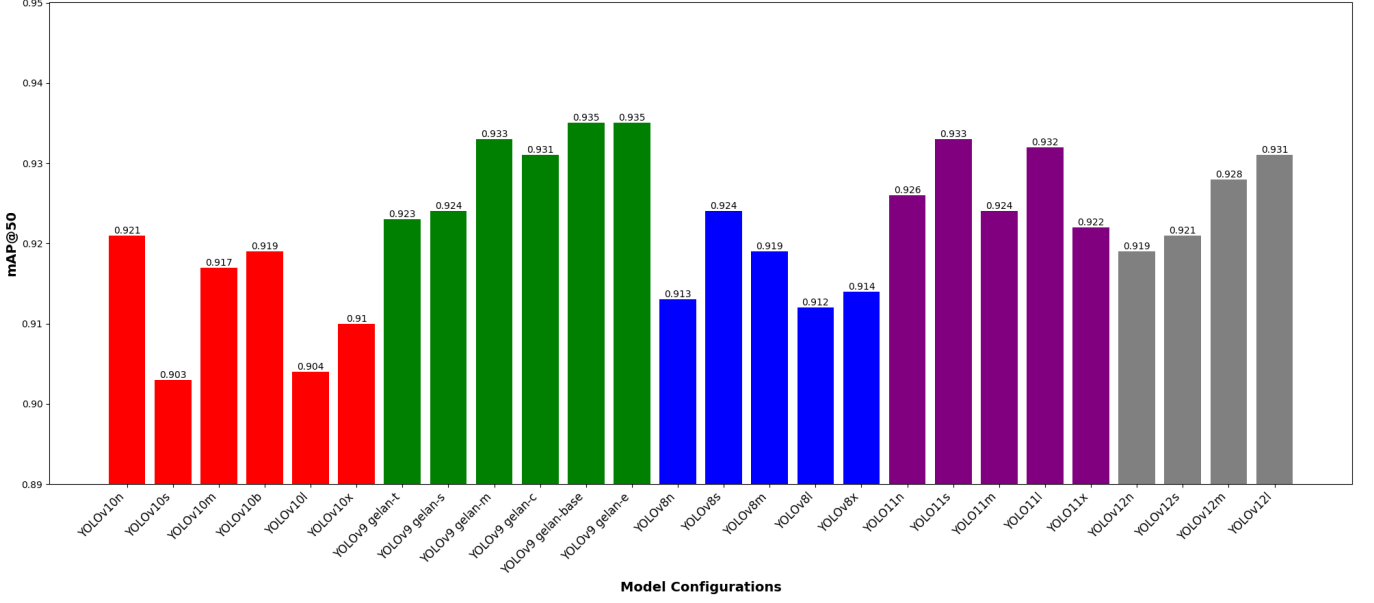


Figure 6: mAP@50 scores for all tested configurations of YOLOv8, YOLOv9, YOLOv10, YOLO11, and YOLOv12 models
Bar diagram showing mAP@50 scores for all Configurations of YOLOv8, YOLOv9, YOLOv10, YOLO11, and YOLOv12 for fruitlet detection in commercial orchards.

Table 2: Precision and Recall for YOLO Configurations (YOLOv8, YOLOv9, YOLOv10, YOLO11, and YOLOv12) used in Fruitlet Detection. For F1-score metrics visualization, refer to Appendix Figure 9

Configuration	Precision	Recall
YOLOv8n	0.840	0.883
YOLOv8s	0.883	0.864
YOLOv8m	0.897	0.858
YOLOv8l	0.877	0.852
YOLOv8x	0.870	0.865
YOLOv9 gelan-e	0.875	0.891
YOLOv9 gelan-c	0.903	0.872
YOLOv9 gelan-s	0.877	0.881
YOLOv9 gelan-t	0.873	0.864
YOLOv9 gelan-m	0.866	0.899
YOLOv9 gelan-base	0.881	0.895
YOLOv10n	0.891	0.862
YOLOv10s	0.839	0.872
YOLOv10m	0.877	0.864
YOLOv10b	0.898	0.870
YOLOv10l	0.842	0.872
YOLOv10x	0.908	0.846
YOLO11n	0.897	0.868
YOLO11s	0.881	0.883
YOLO11m	0.858	0.897
YOLO11l	0.899	0.883
YOLO11x	0.891	0.866
YOLOv12n	0.864	0.881
YOLOv12s	0.880	0.840
YOLOv12m	0.860	0.883
YOLOv12l	0.871	0.900

indicating that specific enhancements in YOLOv9 have likely boosted its accuracy capabilities. In contrast, the YOLOv8 configurations show a wider spectrum of performance, with YOLOv8s achieving the highest mAP@0.50 within its group at 0.924, nearly matching the top-performing configurations of YOLOv10. Nonetheless, configurations such as YOLOv8l and YOLOv8x, with lower scores of 0.912 and 0.914 respectively. This variation may be attributed to an over-parameterization of the models for the given task. Considering the dataset’s simplicity and fewer categories relative to the more complex COCO dataset, the extensive parameters in these models might be excessive and potentially obstruct optimal performance.

3.3. Analysis of Computational Efficiency: Image Processing Speed

In assessing computational efficiency, particularly image processing speeds, YOLO11 emerged as the top performer, achieving remarkably low inference speeds across its variants, with an exceptional rate of just 2.4 ms. Moreover, the pre-processing speeds for YOLO11 were notably quick, registering at 0.1 ms for all variants, except for YOLO11l, which still performed impressively at 0.2 ms. This rapid preprocessing facilitates faster readiness of images for subsequent detection phases. Despite YOLO11’s superior performance, YOLOv12n, YOLOv12m, and YOLOv8x was also notable within its series for achieving the fastest preprocessing speed at merely 0.8, 0.3, and 0.9 ms respectively. While preprocessing speeds are expected to be consistent across models, the discrepancies observed, particularly in YOLOv12n, YOLOv12m, and YOLOv8x’s speed, are likely due to random variations rather than fundamental differences in model architecture.

Expanding on the theme of processing efficiency, YOLO11n configurations demonstrated excellence in inference speed. YOLO11n, in particular, recorded a speed of 2.4 ms, significantly outpacing the fastest models from the YOLOv12, YOLOv10, YOLOv9 and YOLOv8 series. In comparison, the quickest YOLOv9 model, YOLOv9 Gelan-s, logged an inference time of 11.5 ms, and the leading YOLOv10 model, YOLOv10-s, achieved 5.5 ms. These results underscore YOLOv8n’s superior capability in rapid image processing, highlighting its robustness and suitability for scenarios that demand high-speed, accurate object detection. This analysis reveals that while the YOLO11 series leads in low-latency performance, previous iterations like YOLOv8 still maintain competitive advantages, particularly in environments where quick decision-making is crucial. The detailed performance metrics for each model configuration across the YOLOv8, YOLOv9, YOLOv10, and YOLOv11 series are systematically presented in Table 2.

3.4. Field Validation of Counting Accuracy: RMSE and MAE Metrics

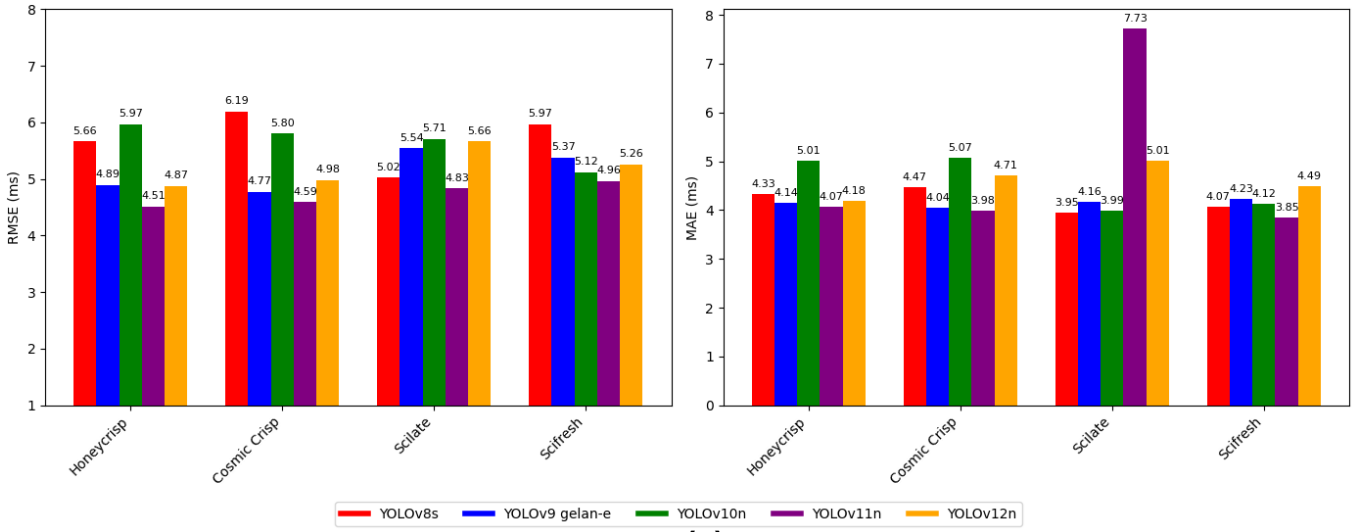
For the counting validation performed on four apple varieties using images collected with an Apple iPhone 14 smartphone, the top-performing configurations from each YOLO version (YOLOv8n, YOLOv9 Gelan-e, YOLOv10n, YOLOv11n, and YOLOv12n) were rigorously assessed for fruit counting accuracy. Among these, YOLOv11n demonstrated exceptional performance by consistently yielding lower error metrics compared to its predecessors. In particular, for Honeycrisp apples, YOLOv11n recorded an RMSE of 4.51 ms and an MAE of 4.07 ms, marking a significant improvement in detection precision and reliability. Similarly, for Cosmic Crisp apples, the configuration achieved an RMSE of 4.59 ms and an MAE of 3.98 ms. These results underscore the model’s robustness in handling subtle variations in fruit appearance and background complexity, indicating that YOLOv11n is highly effective in accurately detecting and counting fruitlets under challenging in-field conditions.

In contrast, for the Scilate variety, YOLOv11n produced an RMSE of 4.83 ms and a notably higher MAE of 7.73 ms, while for Scifresh apples it achieved an RMSE of 4.96 ms and an MAE of 3.85 ms. These differences highlight the impact of fruit variety and environmental factors on model performance, suggesting that detection accuracy can vary with the inherent complexity of the orchard scene. Figure 7a provides a comprehensive comparison of RMSE (left) and MAE (right) for green apple detection across the evaluated YOLO configurations, effectively illustrating these performance variations. Overall, the quantitative findings indicate that YOLOv11n offers superior accuracy and precision in in-field counting validation, thereby establishing its advanced capability to reliably detect and count green apple fruitlets. These insights pave the way for further enhancements in automated fruit detection and precision agriculture applications.

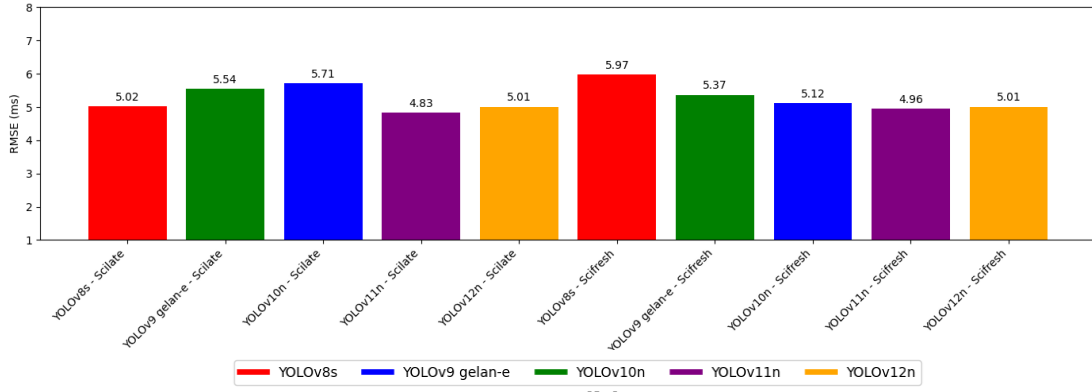
It is noteworthy that the counting validation was conducted on images collected using an Apple iPhone 14, which were not part of the training dataset; the models were exclusively

Table 3: Processing Speeds for YOLOv8, YOLOv9, YOLOv10, YOLO11 and YOLOv12 Configurations in Fruitlet Detection. In terms of speed, YOLO11n achieved the fastest inference at 2.4 ms, and four YOLO11 configurations recorded the fastest preprocessing speed at 0.1 ms. Meanwhile, YOLOv12m and YOLOv12l demonstrated the quickest postprocessing speeds. **Note***: Preprocessing refers to the initial stage where images are prepared for analysis, involving adjustments such as scaling, normalization, and augmentation to optimize them for detection. Inference refers to the core phase where the model analyzes the preprocessed images to detect and identify objects based on learned features and patterns. Postprocessing refers to the final stage that refines the outputs from inference, applying techniques like Non-Maximum Suppression (NMS) to eliminate redundant detections and finalize the list of detected objects.

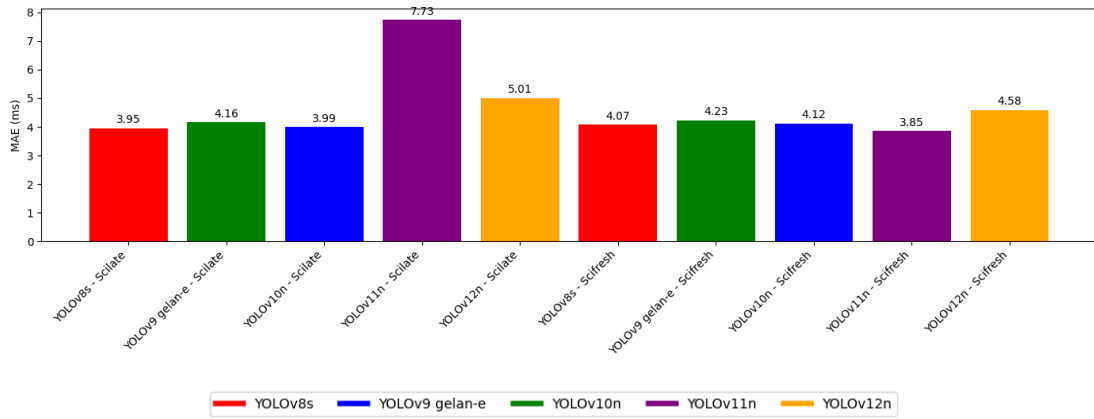
YOLO models	YOLO Configuration	Pre processing (ms)	Inference (ms)	Post processing (ms)
YOLOv8				
	YOLOv8n	1.3	4.1	2.3
	YOLOv8s	1.3	6.4	2.3
	YOLOv8m	1.3	11.2	2.1
	YOLOv8l	1.2	18.7	2.2
	YOLOv8x	0.9	24.8	2.3
YOLOv9				
	YOLOv9 Gelan-t	1.3	14.1	2.2
	YOLOv9 Gelan-s	1.3	11.5	2.2
	YOLOv9 Gelan-m	1.3	14.0	2.1
	YOLOv9 Gelan-c	1.3	17.0	2.0
	YOLOv9 Gelan-base	1.2	17.2	2.0
	YOLOv9 Gelan-e	1.1	33.5	1.9
YOLOv10				
	YOLOv10n	1.4	5.5	1.6
	YOLOv10s	1.3	7.7	1.6
	YOLOv10m	1.3	13.0	1.6
	YOLOv10b	1.3	16.7	1.5
	YOLOv10l	1.4	19.6	1.5
	YOLOv10x	1.2	26.5	1.5
YOLOv11				
	YOLOv11n	0.1	2.4	2.2
	YOLOv11s	0.1	5.0	2.5
	YOLOv11m	0.1	11.9	0.6
	YOLOv11l	0.2	14.6	1.2
	YOLOv11x	0.1	26.3	0.6
YOLOv12				
	YOLOv12n	0.8	4.6	0.7
	YOLOv12s	3.8	9.6	1.0
	YOLOv12m	0.3	15.9	0.5
	YOLOv12l	0.6	19.2	0.5



(a)



(b)



(c)

Figure 7: Illustration of In-Field Counting and Validation for Green Fruit Detection by YOLOv8, YOLOv9, YOLOv10, YOLOv11 and YOLOv12 configurations

(a) RMSE and MAE for in-field counting validation for Green Apple Detection Using iPhone 14 pro Images: Comparison of RMSE (left) and MAE (right) for green apple detection across YOLOv8s, YOLOv9 gelan-e, YOLOv10n, YOLOv11n, and YOLOv12n configurations, highlighting model accuracy and precision. (b) RMSE and MAE for Green Apple Detection Using a consumer grade Machine Vision Sensor Displays RMSE and (c) MAE for green apple detection across YOLOv8s, YOLOv9 gelan-e, YOLOv10n, YOLOv11n, and YOLOv12n assessed with Intel Realsense camera

trained on images captured by an Intel RealSense camera. Despite this domain gap, the top-performing configurations, specially YOLOv11n and YOLOv9 Gelan-e, demonstrated robust counting accuracy across four apple varieties. In particular, YOLOv9 Gelan-e excelled with Honeycrisp apples by achieving an RMSE of 4.89 and an MAE of 4.14, while for Cosmic Crisp, it registered the lowest RMSE of 4.77 and an MAE of 4.04. For the Scilate variety, the model produced an RMSE of 5.54, and for Scifresh, an RMSE of 5.37 was observed, with corresponding MAE values of 4.16 and 4.23, respectively. These performance metrics indicate that YOLOv9 Gelan-e consistently outperformed the YOLOv8, YOLOv10, YOLOv11, and YOLOv12 configurations in precision fruit counting under these novel imaging conditions.

The results clearly illustrate the models' ability to generalize to new data, even when the images are acquired with different sensors. Figure 7b provides a detailed visual distribution of RMSE and MAE values, capturing the in-field counting validation for green apple detection using iPhone 14 images. This evaluation not only confirmed the strong performance of YOLOv11n but also underscored the competitive robustness of YOLOv9 Gelan-e in challenging conditions. The observed discrepancies between model configurations emphasize the need for further exploration into domain adaptation and cross-sensor generalization techniques. Ultimately, these findings highlight the potential for training with high-quality machine vision sensor data to produce models that maintain high accuracy when deployed with more widely available smartphone cameras, thereby enhancing practical applications in precision agriculture.

In the evaluation of images captured by consumer-scale machine vision cameras for apple counting, the updated results highlight the exceptional accuracy of the YOLO11n configuration. Specifically, YOLO11n demonstrated remarkable performance in counting Scifresh apples, achieving an RMSE of 3.06 and an MAE of 2.33. This outperforms the previously noted accuracy of the YOLOv9 Gelan-e configuration, which led earlier assessments with an RMSE of 3.11 and an MAE of 2.46 for the same variety using Intel Realsense cameras. Figures 12 and 13 illustrate the distribution of RMSE and MAE across multiple configurations and apple varieties; Scilate, Scifresh, Honeycrisp, Cosmic Crisp captured with an Apple iPhone 14, with special attention to the performance on Scifresh apples captured using Intel Realsense cameras. These examples highlight the variable performance of each model and underscore the need for larger, more diverse training datasets.

3.5. Discussion

Figure 8 presents several examples of green apple detection on smartphone-captured images of two apple varieties, Honey Crisp and Cosmic Crisp, which were not part of the training set for any of the YOLO models evaluated. The figure is divided into five panels, each illustrating the detection performance of different YOLO configurations. For instance, Figure 8a depicts YOLOv8n's performance on Cosmic Crisp apples; despite the apples being clearly visible, albeit partially occluded as highlighted by red dotted regions, YOLOv8n failed to detect the ap-

ple in the lower region, where approximately 90% of the fruit was visible. In Figure 8b, the YOLOv9 Gelan-base model almost identified the majority of the green apples in the frame, yet it missed one apple, as indicated by red dotted circles. The apple in the upper region was severely occluded and confusing, which likely contributed to the detection failure. Similarly, Figure 8c illustrates YOLOv10n's performance on a densely clustered scene in Cosmic Crisp apples, where one red-dotted apple was not correctly detected.

Furthermore, Figure 8d shows the detection results of YOLOv11n on Cosmic Crisp apples, where the model missed two apples due to occlusion. In the final panel, Figure 8e, the YOLOv12 configuration is evaluated for green fruitlet detection; here, the red dotted regions indicate apple fruitlets that were not detected. Although some of the smaller, occluded fruitlets in complex scenes were partially recognized, a clearly visible fruitlet in the upper red dotted region was not identified by YOLOv12. These diverse outcomes underscore the variable detection capabilities of the different YOLO models when confronted with unseen data acquired from a smartphone. Overall, the results from this study, which involved 1,149 images, suggest that while current models show promise, there is a critical need to expand and diversify training datasets and to utilize more computationally intensive training environments. Such improvements would likely enhance model accuracy and generalizability in real-world fruit detection applications.

Recent studies in green fruit detection have proposed a variety of innovative approaches to overcome the challenges posed by complex orchard environments. Liu et al. (2022) introduced a model that replaces the standard Feature Pyramid Network (FPN) with a Residual FPN (RFPN) and integrates a double-layer Convolutional Block Attention (CBA) network, achieving an 85.3% segmentation accuracy for obscured green fruits. Similarly, Jia et al. (2022) developed Fast-FDM, a detection framework based on an optimized, anchor-free Foveabox model using EfficientNetV2-S as its backbone, combined with a weighted bi-directional feature pyramid network (BiFPN) and adaptive training sample selection (ATSS), resulting in a mean average precision (mAP) of 62.3% with reduced computational costs. Hussain et al. (2023) leveraged Mask R-CNN for instance segmentation and implemented orientation estimation algorithms to achieve green fruit segmentation APs of 83.4% overall and 91.3% for larger masks, thereby facilitating robust fruit thinning operations. In addition, Sun et al. (2022) proposed GHFormer-Net, which enhances small green apple and begonia fruit detection in low-light conditions by incorporating a gradient harmonizing mechanism into a RetinaNet-PVTv2 framework, achieving APs of up to 85.2% on benchmark datasets. These methods collectively underscore significant strides in balancing detection accuracy and computational efficiency, yet each presents inherent trade-offs in speed, robustness, or sensor adaptability. In contrast, the YOLO-based methods evaluated in our study offer a novel alternative with superior real-time performance. Notably, the YOLO11n configuration achieved an extraordinary inference speed of only 2.4 ms, a metric that has not yet been reported in the literature for green fruit detection. This excep-

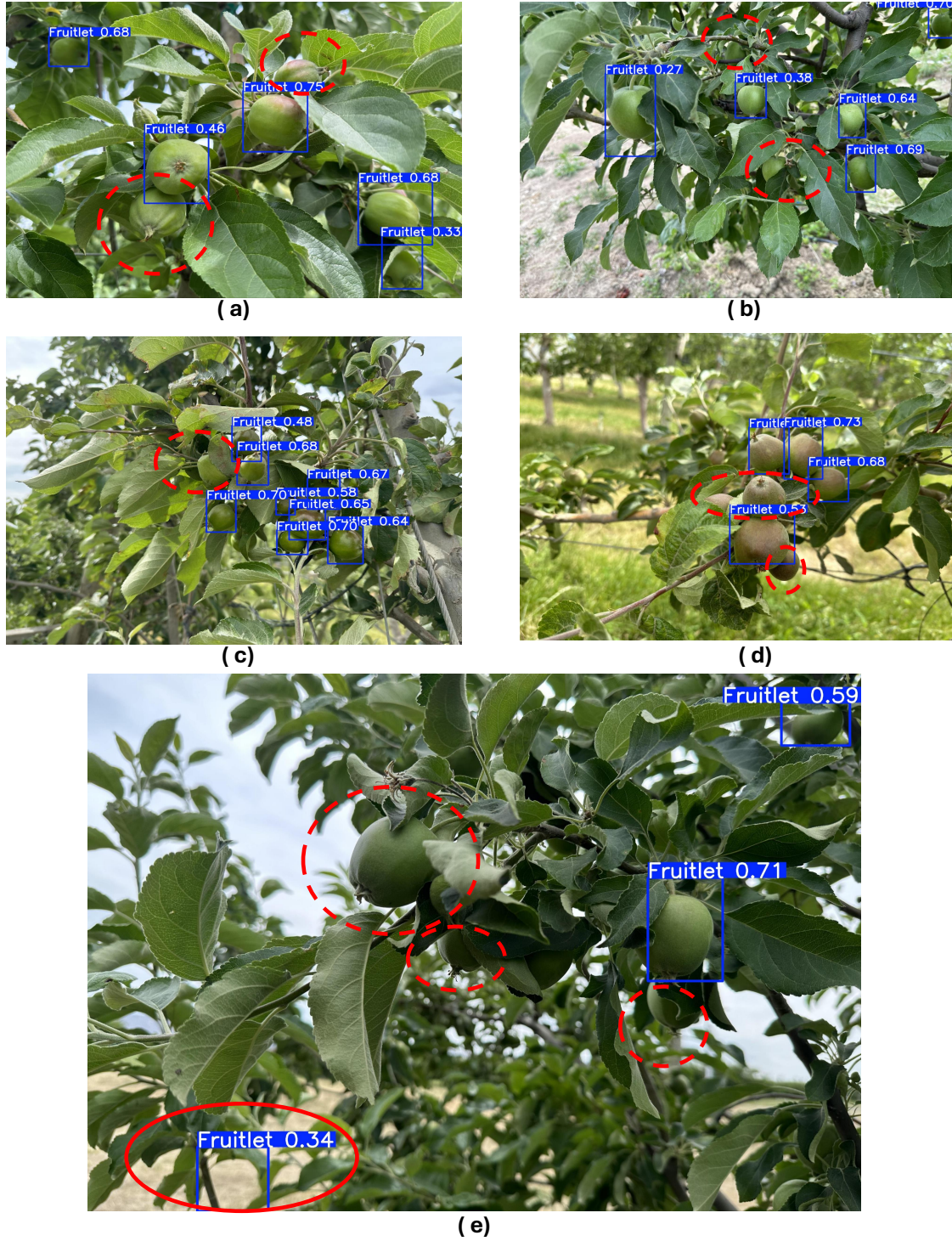


Figure 8: Detection results for green apple detection on iPhone 14 images from Honey Crisp and Cosmic Crisp varieties, which were not used during training: (a) Showing YOLOv8n missing a nearly 90% visible apple in a partially occluded region (red dotted area); (b) Illustrating YOLOv9 Gelan-base, which detected most apples but missed one in both upper and lower occluded areas; (c) Illustrating YOLOv10n failing to detect a clearly visible apple in a densely clustered scene; (d) Illustrating YOLOv11n missing two apples due to occlusion, while panel (e) Showing YOLOv12 failing to detect a clearly visible fruitlet despite partial detection of smaller, occluded fruitlets.

tional speed, coupled with competitive precision and recall metrics, highlights YOLO11n’s potential for rapid, high-accuracy detection in dynamic agricultural settings. By demonstrating robust performance even when deployed on images acquired

from smartphones, which were not part of the training set, our findings indicate that the YOLO series, particularly YOLO11n, can offer a more efficient and practical solution for automated fruit detection and counting. This performance advantage po-

sitions YOLO11n as a promising candidate for integration into precision agriculture systems, enabling faster decision-making and improved operational efficiency in real-world scenarios.

4. Conclusion and Future Suggestions

In this study, we conducted a comprehensive evaluation of various configurations of three state-of-the-art YOLO object detection models (YOLOv8, YOLOv9, YOLOv10, and YOLO11) to assess their effectiveness in detecting green apples before thinning in complex orchard environments using different sensors and conditions. The specific conclusions and future suggestions of this study are specifically summarized in following five points:

- **Model-Specific Performance:** YOLOv10x and YOLOv9 gelan-c demons demonstrated remarkable performance in precision and recall metrics, respectively. YOLOv10x achieved the highest precision score of 0.908 while YOLOv9 gelan-c achieved the second highest precision score of 0.903, showcasing superior accuracy in correctly identifying fruitlets without false positives. Meanwhile, YOLOv12l recorded the highest recall rate of 0.900, indicating its efficacy in capturing almost all fruitlets in the images.
- **Accuracy Across Configurations:** For mean Average Precision at a 50% Intersection over Union (mAP@0.50), analysis revealed that YOLOv9 configurations, particularly YOLOv9 Gelan-base and YOLOv9 Gelan-e, achieved the highest precision with scores of 0.935, demonstrating exceptional robustness in agricultural object detection. Furthermore, YOLO11s and YOLOv12l performed comparably well with mAP@0.50 values of 0.933 and 0.931, respectively, thereby underscoring the efficacy of these advanced YOLO configurations for precise detection in agricultural settings.
- **Counting Validation:** YOLO11n model demonstrated outstanding precision, achieving the lowest RMSE and MAE across all tested varieties, underscoring its robustness in fruit detection applications.

Specifically, in iphone 14 pro captured images, for YOLO11n recorded the best RMSE scores for each variety: 4.51 for Honeycrisp, 4.59 for Cosmic Crisp, 4.83 for Scilate, and 4.96 for Scifresh. These figures represent the model's consistent ability to accurately estimate fruit counts with minimal error, reflecting its sophisticated detection capabilities. Moreover, the MAE scores further illustrate YOLO11n's precision, with the lowest values recorded at 4.07 for Honeycrisp, 3.98 for Cosmic Crisp, and 3.85 for Scifresh. The exception was Scilate, where YOLO11n recorded a higher MAE of 7.73, suggesting an area for potential.

Likewise for a consumer grade RGB-D camera counting validation, YOLO11n consistently demonstrated the most accurate performance. For Scilate, YOLO11n achieved

the lowest RMSE of 4.83 and a significantly higher MAE of 7.73, indicating its precise but inconsistent counting ability under certain conditions. Conversely, for the Scifresh variety, YOLO11n again led with the lowest RMSE of 4.96 and an MAE of 3.85, showcasing its robustness in accurately detecting and counting fruitlets.

- **Sensor-Specific Training Impact:** The study demonstrated that models trained exclusively on data from a specific sensor, such as the Intel Realsense, exhibited significantly enhanced performance when evaluated on data from the same sensor. This finding underscores the critical importance of incorporating diverse sensor data during training to achieve robust model performance across varied deployment scenarios.
 - **Recommendations for Model Deployment:** When deploying YOLO models in automation and robotics, particularly for real-time agricultural tasks such as early-stage crop load management in apple orchards, the choice of model configuration becomes crucial. Each YOLO family has standout configurations optimized for speed and accuracy: YOLOv8n leads its family with an impressive inference speed of 4.1 ms, making it highly suitable for rapid processing needs. In the YOLOv9 series, YOLOv9 Gelan-s excels with a best inference speed of 11.5 ms. For the YOLOv10 configurations, YOLOv10n tops with a speed of 5.5 ms. For YOLOv12, YOLOv12n configuration achieved top inference speed at 4.6ms.
- However, surpassing these, YOLO11n from the YOLO11 series achieves an extraordinary low inference speed of only 2.4 ms, making it the optimal choice for environments where both high speed and precision are critical for efficient automation and real-time decision-making.

Author contributions statement

Ranjan Sapkota: Conceptualization, Data Curation, Methodology, Software, Literature Search, writing original draft, Vizualization. **Manoj Karkee:** Review, Editing and Overall Funding to Supervisory.

Acknowledgement

This work was supported by the National Science Foundation and the United States Department of Agriculture, National Institute of Food and Agriculture through the "Artificial Intelligence (AI) Institute for Agriculture" Program under Award AWD003473, and the AgAID Institute. , Astrid Wimmer, Randall Cason, Diego Lopez, Giulio Diracca, and Priyanka Upadhyaya for their invaluable efforts in data preparation and logistical support throughout this project. Special thanks to Dave Allan for granting orchard access. We also acknowledge the contribution of open-source platforms Roboflow (<https://roboflow.com/>), Ultralytics (<https://docs.ultralytics.com/models/yolo11/>), Hugging Face (<https://huggingface.co/>), and OpenAI (ChatGPT) for the

models and implementation assistance in our project through their open-source platform.

Declarations

The authors declare no conflicts of interest.

Our other agricultural object detection and image processing research can be found on Sapkota et al. (2023, 2024a,b); Sapkota and Karkee (2024b,c); Sapkota et al. (2024e,c); Meng et al. (2025); Sapkota and Karkee (2023); Sapkota et al. (2024f); Churuvija et al. (2025); Sapkota and Karkee (2024a); Khanal et al. (2023); Sapkota et al. (2025b,a)

References

- Bochkovskiy, A., Wang, C.Y., Liao, H.Y.M., 2020. Yolov4: Optimal speed and accuracy of object detection. *arXiv preprint arXiv:2004.10934*.
- Castrillón, M., Déniz, O., Hernández, D., Lorenzo, J., 2011. A comparison of face and facial feature detectors based on the viola-jones general object detection framework. *Machine Vision and Applications* 22, 481–494.
- Churuvija, M., Sapkota, R., Ahmed, D., Karkee, M., 2025. A pose-versatile imaging system for comprehensive 3d modeling of planar-canopy fruit trees for automated orchard operations. *Computers and Electronics in Agriculture* 230, 109899.
- Haleem, A., Javaid, M., Singh, R.P., Rab, S., Suman, R., 2021. Hyperautomation for the enhancement of automation in industries. *Sensors International* 2, 100124.
- Hussain, M., He, L., Schupp, J., Lyons, D., Heinemann, P., 2023. Green fruit segmentation and orientation estimation for robotic green fruit thinning of apples. *Computers and Electronics in Agriculture* 207, 107734.
- Jia, W., Wang, Z., Zhang, Z., Yang, X., Hou, S., Zheng, Y., 2022. A fast and efficient green apple object detection model based on foveabox. *Journal of King Saud University-Computer and Information Sciences* 34, 5156–5169.
- Jocher, G., Stoken, A., Borovec, J., Changyu, L., Hogan, A., Diaconu, L., Poznanski, J., Yu, L., Rai, P., Ferriday, R., et al., 2020. ultralytics/yolov5: v3. 0. Zenodo.
- Khan, A., Sohail, A., Zahoora, U., Qureshi, A.S., 2020. A survey of the recent architectures of deep convolutional neural networks. *Artificial intelligence review* 53, 5455–5516.
- Khanal, S.R., Sapkota, R., Ahmed, D., Bhattarai, U., Karkee, M., 2023. Machine vision system for early-stage apple flowers and flower clusters detection for precision thinning and pollination. *IFAC-PapersOnLine* 56, 8914–8919.
- Lampert, C.H., Blaschko, M.B., Hofmann, T., 2008. Beyond sliding windows: Object localization by efficient subwindow search, in: 2008 IEEE conference on computer vision and pattern recognition, IEEE. pp. 1–8.
- Li, C., Li, L., Jiang, H., Weng, K., Geng, Y., Li, L., Ke, Z., Li, Q., Cheng, M., Nie, W., et al., 2022. Yolov6: A single-stage object detection framework for industrial applications. *arXiv preprint arXiv:2209.02976*.
- Liu, M., Jia, W., Wang, Z., Niu, Y., Yang, X., Ruan, C., 2022. An accurate detection and segmentation model of obscured green fruits. *Computers and Electronics in Agriculture* 197, 106984.
- Liu, W., Anguelov, D., Erhan, D., Szegedy, C., Reed, S., Fu, C.Y., Berg, A.C., 2016. 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. pp. 21–37.
- Manakitsa, N., Maraslidis, G.S., Moysis, L., Fragulis, G.F., 2024. A review of machine learning and deep learning for object detection, semantic segmentation, and human action recognition in machine and robotic vision. *Technologies* 12, 15.
- Meng, Z., Du, X., Sapkota, R., Ma, Z., Cheng, H., 2025. Yolov10-pose and yolov9-pose: Real-time strawberry stalk pose detection models. *Computers in Industry* 165, 104231.
- Redmon, J., Divvala, S., Girshick, R., Farhadi, A., 2016. You only look once: Unified, real-time object detection, in: *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 779–788.
- Redmon, J., Farhadi, A., 2018. Yolov3: An incremental improvement. *arXiv preprint arXiv:1804.02767*.
- Sapkota, R., Ahmed, D., Churuvija, M., Karkee, M., 2024a. Immature green apple detection and sizing in commercial orchards using yolov8 and shape fitting techniques. *IEEE Access* 12, 43436–43452.
- Sapkota, R., Karkee, M., 2023. Creating image datasets in agricultural environments using dall. e: generative ai-powered large language model. *arXiv preprint arXiv:2307.08789*.
- Sapkota, R., Karkee, M., 2024a. Comparing yolov11 and yolov8 for instance segmentation of occluded and non-occluded immature green fruits in complex orchard environment. *arXiv preprint arXiv:2410.19869*.
- Sapkota, R., Karkee, M., 2024b. Integrating yolo11 and convolution block attention module for multi-season segmentation of tree trunks and branches in commercial apple orchards. *arXiv preprint arXiv:2412.05728*.
- Sapkota, R., Karkee, M., 2024c. Yolo11 and vision transformers based 3d pose estimation of immature green fruits in commercial apple orchards for robotic thinning. *arXiv preprint arXiv:2410.19846*.
- Sapkota, R., Meng, Z., Karkee, M., 2024b. Synthetic meets authentic: Leveraging llm generated datasets for yolo11 and yolov10-based apple detection through machine vision sensors. *Smart Agricultural Technology* 9, 100614.
- Sapkota, R., Paudel, A., Karkee, M., 2024c. Zero-shot automatic annotation and instance segmentation using llm-generated datasets: Eliminating field imaging and manual annotation for deep learning model development. *arXiv preprint arXiv:2411.11285*.
- Sapkota, R., Qureshi, R., Calero, M.F., Badgajar, C., Nepal, U., Poullose, A., Zeno, P., Vaddevolu, U.B.P., Khan, S., Shoman, M., et al., 2024d. Yolov12 to its genesis: a decadal and comprehensive review of the you only look once (yolo) series. *arXiv preprint arXiv:2406.19407*.
- Sapkota, R., Qureshi, R., Flores-Calero, M., Badgajar, C., Nepal, U., Poullose, A., Zeno, P., Bhanu Prakash Vaddevolu, U., Yan, P., Karkee, M., et al., 2024e. Yolov10 to its genesis: A decadal and comprehensive review of the you only look once series. Available at SSRN 4874098.
- Sapkota, R., Qureshi, R., Hassan, S.Z., Shutske, J., Shoman, M., Sajjad, M., Dharejo, F.A., Paudel, A., Li, J., Meng, Z., et al., 2024f. Multi-modal llms in agriculture: A comprehensive review. *Authorea Preprints*.
- Sapkota, R., Raza, S., Karkee, M., 2025a. Comprehensive analysis of transparency and accessibility of chatgpt, deepseek, and other sota large language models. *Preprints.org DOI 10*.
- Sapkota, R., Raza, S., Shoman, M., Paudel, A., Karkee, M., 2025b. Image, text, and speech data augmentation using multimodal llms for deep learning: A survey. *arXiv preprint arXiv:2501.18648*.
- Sapkota, R., Stenger, J., Ostlie, M., Flores, P., 2023. Towards reducing chemical usage for weed control in agriculture using uas imagery analysis and computer vision techniques. *Scientific reports* 13, 6548.
- Sun, M., Xu, L., Luo, R., Lu, Y., Jia, W., 2022. Ghformer-net: Towards more accurate small green apple/begonia fruit detection in the nighttime. *Journal of King Saud University-Computer and Information Sciences* 34, 4421–4432.
- Tian, Y., Ye, Q., Doermann, D., 2025. Yolov12: Attention-centric real-time object detectors. *arXiv preprint arXiv:2502.12524*.
- Wang, A., Chen, H., Liu, L., Chen, K., Lin, Z., Han, J., et al., 2025. Yolov10: Real-time end-to-end object detection. *Advances in Neural Information Processing Systems* 37, 107984–108011.
- Wang, C.Y., Bochkovskiy, A., Liao, H.Y.M., 2023. 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*, pp. 7464–7475.
- Wang, C.Y., Yeh, I.H., Mark Liao, H.Y., 2024. Yolov9: Learning what you want to learn using programmable gradient information, in: *European conference on computer vision*, Springer. pp. 1–21.
- Wu, B., Iandola, F., Jin, P.H., Keutzer, K., 2017. Squeezenet: 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*, pp. 129–137.
- Zhao, Z.Q., Zheng, P., Xu, S.t., Wu, X., 2019. Object detection with deep learning: A review. *IEEE transactions on neural networks and learning systems* 30, 3212–3232.
- Zou, Z., Chen, K., Shi, Z., Guo, Y., Ye, J., 2023. Object detection in 20 years: A survey. *Proceedings of the IEEE* 111, 257–276.

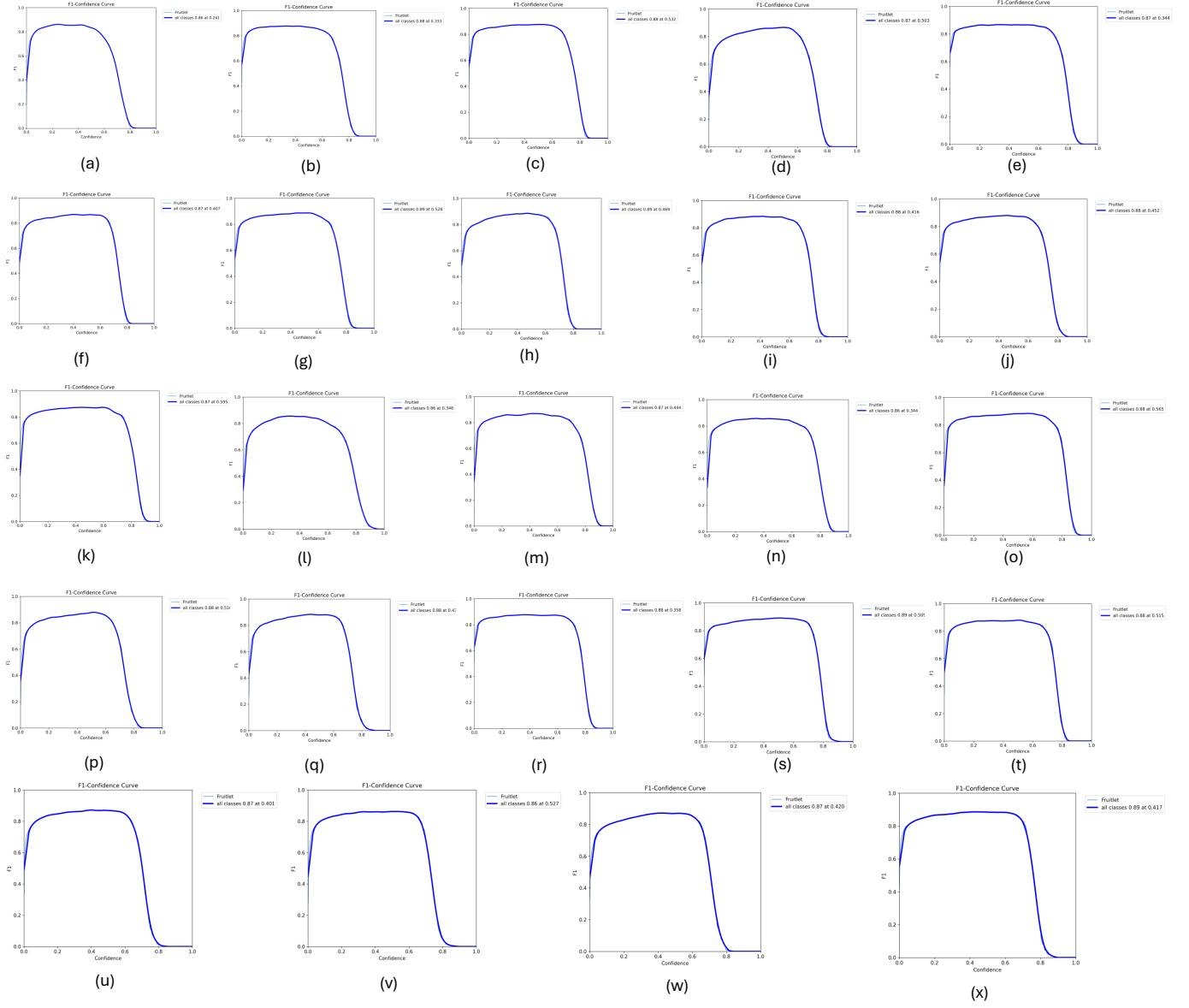


Figure 9: Illustrating the F1-scores of YOLO models (YOLOv8, YOLOv9, YOLOv10, YOLOv11, and YOLOv12 for fruitlet detection: Subfigure (a)–(e) illustrate the F1 scores for YOLOv8 configurations (YOLOv8n, YOLOv8s, YOLOv8m, YOLOv8l, YOLOv8x); (f)–(k) for YOLOv9 configurations (YOLOv9 Gelan-t, Gelan-s, Gelan-m, Gelan-c, Gelan-base, Gelan-e); (l)–(q) for YOLOv10 configurations; (r)–(v) for YOLOv11 configurations; and (w)–(z) for YOLOv12 configurations (YOLOv12n, YOLOv12s, YOLOv12m, YOLOv12l), highlighting comparative detection accuracy in complex orchard environments.