


## Article

## YOLOv8-WBF: Ensemble Learning for Reliable Detection of Endangered Medaka (Oryzias)

Rahmatullah R. <sup>1, </sup>, Armin Lawi <sup>1,2,3</sup>, Muhammad Haerul <sup>1</sup>, Iman Mustika Ismail <sup>1</sup>, Irma Andriani <sup>4</sup>, Andi Iqbal Burhanuddin <sup>5</sup>, and Mario Köppen <sup>6</sup>

<sup>1</sup> Information Systems Study Program, Faculty of Mathematics and Natural Sciences, Hasanuddin University, Indonesia

<sup>2</sup> Data Science and Artificial Intelligence Research Group, Hasanuddin University, Indonesia

<sup>3</sup> B.J. Habibie Institute of Technology, Parepare, Indonesia

<sup>4</sup> Department of Biology, Faculty of Mathematics and Natural Sciences, Hasanuddin University, Indonesia

<sup>5</sup> Department of Fishery, Faculty of Fishery and Marine Sciences, Hasanuddin University, Indonesia

<sup>6</sup> Department of Creative Informatics, Faculty of Computer Science and Systems Engineering, Kyushu Institute of Technology, Japan

\* Correspondence : armin@unhas.ac.id)

**Abstract:** Reliable detection of Medaka (*Oryzias*) fish is essential for ecological monitoring and conservation, particularly for tracking population trends of endangered species. This study evaluates the performance of a state-of-the-art deep learning model (YOLOv8) and an ensemble approach using Weighted Box Fusion (WBF) on a manually annotated dataset of Medaka images collected from online sources. Models were trained and validated using 5-fold cross-validation, and performance was assessed using COCO metrics, including mean Average Precision (mAP), precision, recall, and bounding box regression error. The YOLOv8-WBF ensemble achieved a mAP@0.5:0.95 of 0.578, representing an 8% improvement over the best single model. It also enhanced bounding box localization and classification reliability, particularly for small and visually challenging fish instances. These accuracy gains came at the expense of computational efficiency, with inference requiring approximately five times more operations than a single YOLOv8 model. While less suited for real-time deployment, the ensemble approach offers more reliable detection for offline ecological workflows, where accuracy is prioritized over speed. By reducing missed detections of rare or occluded fish, this work contributes to more robust biodiversity monitoring and provides a baseline for developing optimized ensemble and lightweight detection models in aquatic conservation.

**Keywords:** Medaka (*Oryzias*); Deep Learning; Object Detection; YOLOv8; Weighted Box Fusion; Ensemble Learning; Ecological Monitoring; Biodiversity Conservation

Received:

Revised:

Accepted:

Published:

**Citation:** . Title. *Journal Not Specified* 2025, 1, 0. <https://doi.org/>

**Copyright:** © 2025 by the authors. Submitted to *Journal Not Specified* for possible open access publication under the terms and conditions of the Creative Commons Attribution (CC BY) license (<https://creativecommons.org/licenses/by/4.0/>).

## 1. Introduction

Medaka fish (*Oryzias*) are small freshwater fish valued as ornamental species and are significant for biodiversity studies. They are species that are under danger of extinction as declared by the International Union for Conservation of Nature (IUCN), which makes protecting them an ecological priority. Identifying Medaka apart is challenging since they have subtle morphological differences and small sizes. Not only that, but the variability of aquatic environments also makes it even more challenging. Traditionally, medaka fish have been sacrificed in taxonomy studies due to their genetic value for lineage studies [1]. Their body colour serves as a social signal and reflects environmental conditions [2], making them potential indicators of ecosystem health. Research conducted using traditional

methods, such as direct capture and then being put in an experimental tank to understand its anatomy and internal structure, can harm the organism and limit its application in long-term ecological and genetic studies. Reliable detection technologies can facilitate non-invasive conservation without the need for traditional conservation methods. This technique protects species and ecosystems while enabling effective ecological monitoring in conservation efforts.

Deep learning has evolved beyond digit identification to sophisticated object detection, facilitating applications in autonomous vehicles, medical diagnostics, agricultural automation, and environmental monitoring [3,4]. In ecological research, object detection has been utilized for monitoring insects in agriculture [5,6] and detecting wildlife in natural habitats [7–9], where precise identification is frequently challenging due to visual similarities among species and the complexity of water environments. These issues are also present in the medaka fish species, such as *Oryzias javanicus* and *Oryzias celebensis*, which have subtle morphological distinctions that complicate the reliable identification of endemic fish.

The YOLO architecture is one of the most widely used deep learning-based one-stage object detectors. Among its versions, YOLOv8 has shown a particularly good balance between speed and accuracy, as demonstrated in comparative benchmarks that evaluated YOLOv8 through YOLOv11 under real-world conditions [10]. Researchers have also proposed many variants to push YOLO's performance further, especially in challenging settings. For example, CEH-YOLO adds a high-order deformable attention (HDA) module to better highlight important spatial features, an Enhanced Spatial Pyramid Pooling-Fast (ESPPF) module for improved texture and color feature extraction, and a Composite Detection module to boost detection of small or overlapping underwater objects, along with using WIoU-v3 loss to improve bounding box regression under hard conditions [11]. [12] introduces the Softplus activation function to improve training stability, an AIFI module to strengthen intra-scale feature interactions (reducing false positives and missed detections), and lightweight neck convolution modules (GSConv, VoV-GSCSP) to reduce computational overhead while maintaining accuracy. SCoralDet focuses on underwater soft coral detection, using a Multi-Path Fusion Block (MPFB) to handle varied scales and lighting colors, lightweight modules for efficiency, and an Adaptive Power Transformation label assignment strategy to better align anchors with ground truth when coral structures are complex or blurred [13]. While these studies mainly alter the internal YOLO architecture to address trade-offs between accuracy and speed, in our work, we propose the use of an ensemble method that builds on YOLOv8 without modifying its core architecture by combining the strengths of multiple detection heads or models to improve reliability under environmental variability.

Despite its strong performance, YOLO models could still make misclassifications by detecting background regions as objects or producing duplicate overlapping predictions of the same object. By default, these redundant outputs are reduced through Non-Maximum Suppression (NMS), yet this technique has well-known limitations, particularly when objects overlap or when multiple plausible predictions exist. Recent studies have come up with a new ensemble method, Weighted Boxes Fusion (WBF), that provides a more accurate and robust alternative, since it merges bounding boxes based on confidence scores and spatial alignment rather than discarding valuable predictions outright [14]. The limitations of standard non-maximum suppression (NMS) underscore the need to investigate advanced fusion strategies in object detection. Ensemble methods are particularly effective because they combine multiple models or detection heads, thereby improving uncertainty management, reducing false detections, and increasing consistency in challenging scenarios. The primary contribution of this work is the development of an ensemble approach for YOLOv8 that incorporates Weighted Box Fusion (WBF) to improve detection reliability

in ecological monitoring applications, where accuracy is essential and misclassification significantly affect conservation outcomes.

In this work, we propose a novel Medaka fish dataset with images captured under diverse lighting colors to provide a realistic and challenging benchmark for ecological monitoring. Building on this resource MEDAKA- $\epsilon$ L, an ensemble-based detection framework for accurate and reliable identification of Medaka fish. By training 5 different YOLOv8n models trained across a 5-fold cross-validation of the dataset, which causes the model to learn its own unique feature of the Medaka fish. At the ensemble stage, predictions from multiple YOLOv8 models trained through cross-validation are combined using Weighted Boxes Fusion (WBF). It addresses challenges such as background misclassifications and redundant overlapping bounding boxes. Unlike traditional Non-Maximum Suppression (NMS), WBF merges bounding boxes based on confidence scores and spatial alignment, preserving valuable detections and reducing false positives that achieve higher accuracy. The key contributions of this work are as follows:

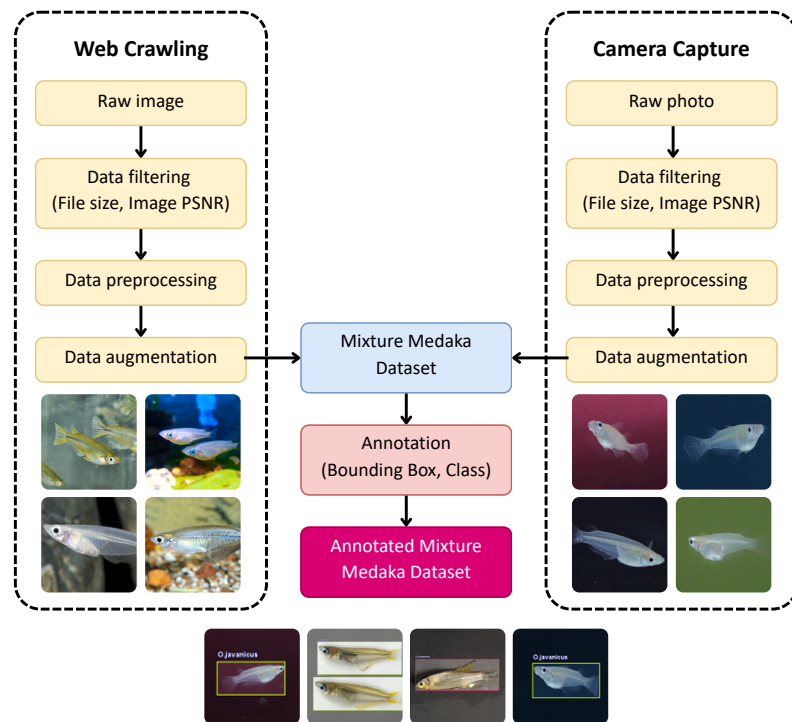
- Introduced the MEDAKA- $\epsilon$ L framework that uses 5 models trained across 5-fold cross-validation, and combines their outputs with Weighted Boxes Fusion (WBF) to improve detection accuracy.
- Proposed the Medaka dataset, a new collection of Medaka fish images that includes manually captured photos under different lighting colors, as well as annotated samples gathered from the internet.
- Developed an ensemble detection approach based on YOLOv8n models, where predictions are merged with WBF to keep valuable detections and reduce false positives.

## 2. Related Works

## 3. Materials and Methods

### 3.1. Data Gathering

Given the scarcity of publicly available Medaka fish imagery and the critical need for precise taxonomic identification in ecological monitoring, we constructed a comprehensive hybrid dataset called the *Medaka Fish Dataset*. This dataset, as shown in Figure 1, combines primary data that we collected through controlled laboratory imaging and *in situ* field observation, with curated internet-sourced images as our secondary data. This approach directly addresses the fundamental challenge specified in our introduction regarding the absence of dedicated datasets for endangered Medaka species detection while ensuring sufficient morphological diversity to support the robustness of our proposed model.



**Figure 1.** Representative samples from the Medaka dataset showing diversity in species, environmental conditions, and imaging scenarios. (a) *O. celebensis* specimens in various naturalistic settings. (b) *O. javanicus* specimens demonstrating morphological variation and environmental diversity.

Primary data acquisition combined (i) controlled laboratory imaging at the Genetics Laboratory, Faculty of Mathematics and Natural Sciences, Hasanuddin University, Indonesia and (ii) *in situ* observation in the Tanjung coastal freshwater-estuarine transition zone in Makassar, Indonesia. This dual setting was intentionally designed to balance morphological clarity with real-world environmental variability. In the laboratory setting, *O. celebensis* and *O. javanicus* were photographed using a Canon EOS M50 camera (24.1 MP, 6000×4000 px CMOS sensor) in modular glass aquaria with four different background colors. The camera operates at an aperture of  $f/2.5$  under diffuse LED lighting. This is done to help minimize light reflection and motion blur while maintaining natural colors. We also used four background colors (red, black, blue, and green) to (1) enhance contrast across pigmentation conditions, (2) avoid overfitting on single color-light pairs, and (3) approximate the natural substrate diversity.

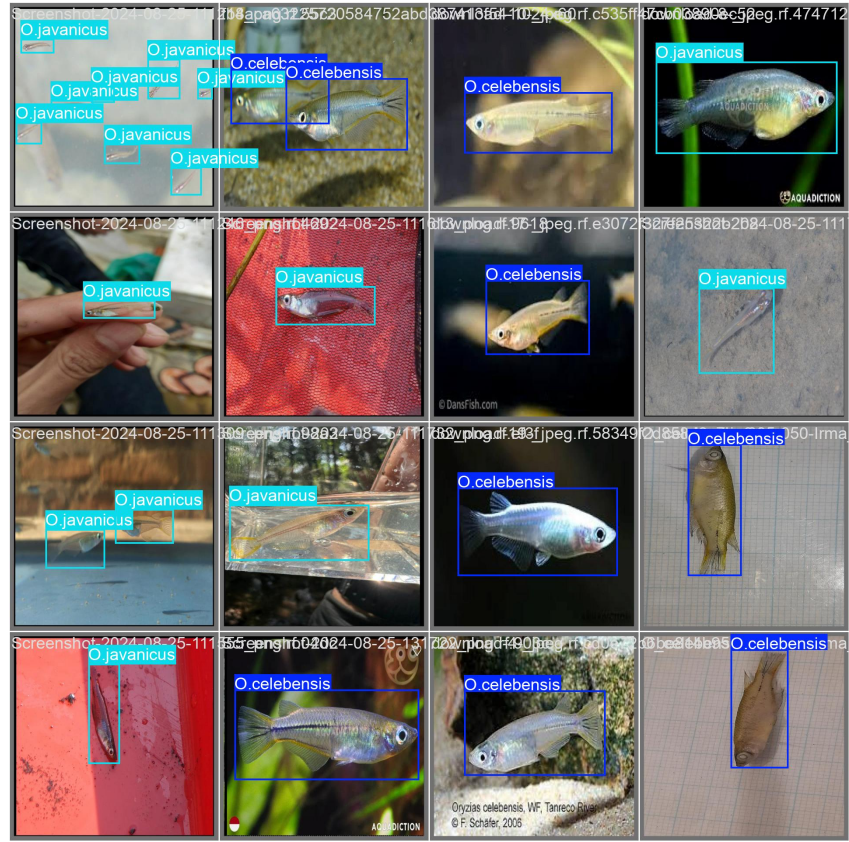
In addition to controlled laboratory imaging, we also added *in situ* observation images to increase the variety of our data. These field images contributed heterogeneous backgrounds, such as irregular substrates, floating particles, and fluctuating lighting. However, images affected by severe turbidity, glaring light reflections, occlusions, or unclear body contours were discarded. Nearly identical temporal sequences, excessive motion blur, and noticeable chromatic aberrations were also excluded. Only minor normalization was applied, consisting of orientation adjustments and removal of unusable frames. No denoising or color correction was performed in order to preserve the genuine visual variation.

In parallel, secondary images were curated from publicly available online sources. To ensure reliability, only specimens that were clear and taxonomically identifiable were retained. Also, images with artificial backgrounds, compression artifacts, or ambiguous species characteristics were excluded to maintain the integrity of our dataset.

Consequently, the combination of the two acquisition strategies yielded a comprehensive dataset that balances controlled imaging precision with ecological variability. In total,



1,511 high-resolution images ranging from 1920x1080 to 4096x3072 pixels of Medaka fish (*Oryzias* species) were compiled.



**Figure 2.** Examples of annotated Medaka images from the mixed dataset.

### 3.2. Data Filtering, Pre-processing, and Annotation

Subsequently, as illustrated in Figure 1, the collected raw images and raw photos were filtered based on their file size and PSNR values. The PSNR for each image was computed using Equation 2, where the Mean Squared Error (MSE) between an image and its reference was obtained using Equation 1.

$$MSE = \frac{1}{mn} \sum_{i=0}^{m-1} \sum_{j=0}^{n-1} [I(i, j) - K(i, j)]^2 \quad (1)$$

$$PSNR = 10 \log_{10} \frac{MAX_i^2}{MSE} \quad (2)$$

A higher PSNR indicates better image fidelity and lower noise. Hence, images falling below a predefined PSNR threshold were discarded to maintain dataset integrity. The remaining samples were then passed through a preprocessing pipeline.

There were two preprocessing steps applied in this research: (1) automatic orientation correction and (2) image resizing to a fixed resolution of 640x640 pixels. Both preprocessing steps follow the default configuration commonly adopted in YOLO-based object detection pipelines.

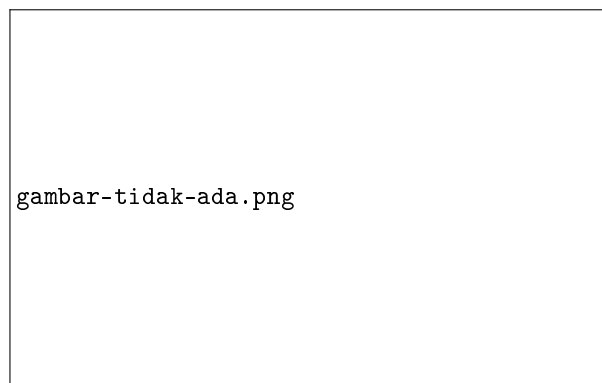
Finally, 1,139 filtered and preprocessed images were annotated using Roboflow, resulting in a total of 1,239 labeled instances: 641 of *Oryzias celebensis* and 598 of *Oryzias javanicus*, respectively. Figure 2 shows the examples of annotated Medaka images from our dataset. The annotation process was conducted by domain experts, who applied bounding boxes to each fish instance and assigned the corresponding species class. In most cases, each image

contained a single fish object (1,071 images), while 62 images included two to three objects, and only four images contained between four and five objects.

### 3.3. Yolo's Architecture V8

In this study, we used the YOLOv8 model as our foundation object detection model. Figure 3 shows the detailed architecture of YOLOv8 and its three main components, the backbone, neck, and head [15]. The first component, the backbone, extracts multi-scale features from input images using an advanced convolutional neural network (CNN) that improves upon the previous CSPLayer from YOLOv5, now called the C2f module [16]. This module applies depthwise separable convolution to optimize the balance between processing speed and feature extraction capabilities [15]. The second component, the neck, integrates and refines the extracted features through the combined adaptation of the Path Aggregation Network (PANet) and the Feature Pyramid Network (FPN). This structure performs multi-scale feature fusion by aggregating feature maps from three hierarchical levels of the backbone using the C2f module [17]. Finally, the head component is responsible for generating the final prediction, which includes bounding box coordinates, confidence scores, and class labels using an anchor-free approach. This approach replaces the previous version's anchor-based method to simplify the prediction process, reduce the number of hyperparameters, and improve generalization abilities to variations in object size and proportions [16]. Thus, by integrating these three components, YOLOv8 improves accuracy, inference speed, and flexibility for a range of object detection tasks.

The training process was carried out on the Medaka dataset with an input size of 640×640 pixels for 50 epochs using the AdamW optimizer with a initial learning rate of 0.01, a batch size of 16, and weight decay of 0.05. Model performance was evaluated using the mean Average Precision (mAP) and mean Average Recall (mAR) to measure detection capabilities for each object class.



**Figure 3.** Yolo V8 Architecture

### 3.4. *k*-fold validation

*k*-fold validation is a technique used to estimate the performance of a learning algorithm on a given dataset. It works by randomly dividing a dataset into *k* disjoint folds with approximately equal size. One fold is used as the validation set for the other *k* − 1 folds used to train the algorithm in each iteration. The algorithm's overall performance is then computed as the average of the evaluation metrics obtained across all *k* folds, and thus reflects performance estimation at the fold level [18].

Previous studies, such as [14,19], have shown that box ensembling indeed increases the performance of object detection models. However, these works mainly focused on ensembling different YOLO versions or model sizes, putting more attention on model diversity than data diversity. In contrast, we looked back at some traditional machine learning techniques, in which we identify opportunities to improve performance by focusing on

the dataset itself, especially in cases of data scarcity, through  $k$ -fold cross-validation-based ensembling learning models.

In this study, a YOLOv8 model was applied to the Medaka dataset using a 5-fold cross-validation scheme. The dataset was divided into five subsets, each comprising 16%, and one test subset comprising 20%. In each fold iteration, one 16% subset was used for validation, while the remaining four subsets were used for training, with the 20% subset remaining as the fixed test set. These folds were rotated alternately so that each 16% subset became the validation set once, resulting in five different fold configurations:

- Fold 1: 16% train, 16% train, 16% train, 16% train, 16% val, 20% test
- Fold 2: 16% train, 16% train, 16% train, 16% val, 16% train, 20% test
- Fold 3: 16% train, 16% train, 16% val, 16% train, 16% train, 20% test
- Fold 4: 16% train, 16% val, 16% train, 16% train, 16% train, 20% test
- Fold 5: 16% val, 16% train, 16% train, 16% train, 16% train, 20% test

### 3.5. Ensemble Methods

After applying 5-fold cross-validation, the predictions from the trained five YOLOv8 models were combined using ensemble techniques. We experimented with three ensemble methods commonly used in object detection: Weighted Box Fusion (WBF), Non-Maximum Suppression (NMS), and Non-Maximum Weighted (NMW). In our experiments, we evaluated these three methods with IoU threshold of  $\theta = 0.5, 0.75$ , and a range of  $0.50 : 0.95$  using MaxDets of 1, 10, and 100 across different object scales.

#### 3.5.1. Non-Maximum Suppression (NMS)

In object detection, Non-Maximum Suppression (NMS) or to be exact, GreedyNMS, merges multiple detections into a single final detection. GreedyNMS selects the highest-scoring detection, thereafter discarding all detections whose overlap over a specified threshold  $\theta$ , and repeats this procedure on the remaining detections [20].

In addition to GreedyNMS, we apply SoftNMS with a sigma value of  $\sigma = 0.1$  in our experiments. SoftNMS differs from GreedyNMS in that it decreases the confidence scores of overlapping detections based on their IoU with the highest-scoring box, rather than discarding them outright. This approach allows detections that partially overlap with a high-confidence box to still contribute to the final results, improving recall without substantially affecting precision [21].

#### 3.5.2. Weighted Box Fusion (WBF)

Unlike Non-Maximum Suppression (NMS), which selects the single highest-confidence box and discards overlapping detections, WBF uses confidence scores of all proposed bounding boxes to construct the average boxes [22]. In practice, the final bounding box is computed as a weighted average of the coordinates from all boxes in each overlapping group, with weights given by the confidence scores assigned by each model.

#### 3.5.3. Non-Maximum Weighted (NMW)

Finally, we also employed Non-Maximum Weighted (NMW) to combine overlapping bounding boxes. In this method, each box in a group contributes to the final prediction proportionally to a weight calculated as the product of its confidence score and its IoU with the most confident box. The final bounding box is then obtained as the weighted average of all boxes in the group [23]. Unlike Weighted Boxes Fusion (WBF), which uses confidence scores from all models to compute the average, NMW incorporates both the confidence and spatial overlap (IoU) of boxes within a single prediction set, allowing it to exploit local agreement among overlapping detections.

### 3.6. Evaluation metrics

Finally, we evaluated the performance of each ensemble methods in the previous section using *Average Precision (AP)* and *Average Recall (AR)*.

**Average Precision (AP)** is calculated as the area under the Precision-Recall curve for each class, as defined in Equation (3):

$$AP = \int_0^1 \text{Precision}(r) dr \quad (3)$$

**Average Recall (AR)** is computed as the average recall over multiple IoU thresholds for each class, as defined in Equation (4):

$$AR = \frac{1}{T} \sum_{t=1}^T \text{Recall}_{\text{IoU}=t} \quad (4)$$

Higher AP and AR values indicate a model that is more accurate and better at detecting objects in the evaluation dataset.

## 4. Results

## 5. Discussion

**Author Contributions:** For research articles with several authors, a short paragraph specifying their individual contributions must be provided. The following statements should be used “Conceptualization, X.X. and Y.Y.; methodology, X.X.; software, X.X.; validation, X.X., Y.Y. and Z.Z.; formal analysis, X.X.; investigation, X.X.; resources, X.X.; data curation, X.X.; writing—original draft preparation, X.X.; writing—review and editing, X.X.; visualization, X.X.; supervision, X.X.; project administration, X.X.; funding acquisition, Y.Y. All authors have read and agreed to the published version of the manuscript.”, please turn to the [CRediT taxonomy](#) for the term explanation. Authorship must be limited to those who have contributed substantially to the work reported.

**Funding:** Please add: “This research received no external funding” or “This research was funded by NAME OF FUNDER grant number XXX.” and and “The APC was funded by XXX”. Check carefully that the details given are accurate and use the standard spelling of funding agency names at <https://search.crossref.org/funding>, any errors may affect your future funding.

**Institutional Review Board Statement:** In this section, you should add the Institutional Review Board Statement and approval number, if relevant to your study. You might choose to exclude this statement if the study did not require ethical approval. Please note that the Editorial Office might ask you for further information. Please add “The study was conducted in accordance with the Declaration of Helsinki, and approved by the Institutional Review Board (or Ethics Committee) of NAME OF INSTITUTE (protocol code XXX and date of approval).” for studies involving humans. OR “The animal study protocol was approved by the Institutional Review Board (or Ethics Committee) of NAME OF INSTITUTE (protocol code XXX and date of approval).” for studies involving animals. OR “Ethical review and approval were waived for this study due to REASON (please provide a detailed justification).” OR “Not applicable” for studies not involving humans or animals.

**Informed Consent Statement:** Any research article describing a study involving humans should contain this statement. Please add “Informed consent was obtained from all subjects involved in the study.” OR “Patient consent was waived due to REASON (please provide a detailed justification).” OR “Not applicable” for studies not involving humans. You might also choose to exclude this statement if the study did not involve humans.

Written informed consent for publication must be obtained from participating patients who can be identified (including by the patients themselves). Please state “Written informed consent has been obtained from the patient(s) to publish this paper” if applicable.



**Data Availability Statement:** We encourage all authors of articles published in MDPI journals to share their research data. In this section, please provide details regarding where data supporting reported results can be found, including links to publicly archived datasets analyzed or generated during the study. Where no new data were created, or where data is unavailable due to privacy or ethical restrictions, a statement is still required. Suggested Data Availability Statements are available in section “MDPI Research Data Policies” at <https://www.mdpi.com/ethics>.

**Acknowledgments:** In this section you can acknowledge any support given which is not covered by the author contribution or funding sections. This may include administrative and technical support, or donations in kind (e.g., materials used for experiments). Where GenAI has been used for purposes such as generating text, data, or graphics, or for study design, data collection, analysis, or interpretation of data, please add “During the preparation of this manuscript/study, the author(s) used [tool name, version information] for the purposes of [description of use]. The authors have reviewed and edited the output and take full responsibility for the content of this publication.”

**Conflicts of Interest:** Declare conflicts of interest or state “The authors declare no conflicts of interest.” Authors must identify and declare any personal circumstances or interest that may be perceived as inappropriately influencing the representation or interpretation of reported research results. Any role of the funders in the design of the study; in the collection, analyses or interpretation of data; in the writing of the manuscript; or in the decision to publish the results must be declared in this section. If there is no role, please state “The funders had no role in the design of the study; in the collection, analyses, or interpretation of data; in the writing of the manuscript; or in the decision to publish the results”.

Abbreviations

The following abbreviations are used in this manuscript:

- MDPI    Multidisciplinary Digital Publishing Institute
- DOAJ    Directory of open access journals
- TLA    Three letter acronym
- LD    Linear dichroism

Appendix A

Appendix A.1

The appendix is an optional section that can contain details and data supplemental to the main text—for example, explanations of experimental details that would disrupt the flow of the main text but nonetheless remain crucial to understanding and reproducing the research shown; figures of replicates for experiments of which representative data are shown in the main text can be added here if brief, or as Supplementary Data. Mathematical proofs of results not central to the paper can be added as an appendix.

Table A1. This is a table caption.

Title 1	Title 2	Title 3
Entry 1	Data	Data
Entry 2	Data	Data

Appendix B

All appendix sections must be cited in the main text. In the appendices, Figures, Tables, etc. should be labeled, starting with “A”—e.g., Figure A1, Figure A2, etc.

## References

1. Mahmudi, A.; et al. Title of the study (please update with actual title). *Journal Name* **2022**, *XX*, XXX–XXX. <https://doi.org/DOLifavailable>.
2. Ueda, R.; Ansai, S.; Takeuchi, H. Rapid body colouration changes in *Oryzias celebensis* as a social signal influenced by environmental background. *Scientific Reports* **2020**. Electronic supplementary material is available online, <https://doi.org/10.6084/m9.figshare.c.7358213>.
3. LeCun, Y.; Bengio, Y.; Hinton, G. Deep learning. *Nature* **2015**, *521*, 436–444. <https://doi.org/10.1038/nature14539>.
4. Zhao, Z.; Zheng, P.; Xu, S.; Wu, X. Object detection with deep learning: A review. *IEEE Transactions on Neural Networks and Learning Systems* **2019**, *30*, 3212–3232.
5. Tang, L.; et al. Deep learning-based insect detection for precision agriculture: A comprehensive review. *Computers and Electronics in Agriculture* **2023**.
6. Ciampi, S.; et al. Automatic insect detection in sticky trap images using deep learning models. *Ecological Informatics* **2023**, *73*, 101–123.
7. Roy, A.; et al. WilDect-YOLO: A one-stage detector for endangered wildlife monitoring. *Ecological Informatics* **2023**.
8. Wenhan, Z.; et al. Enhanced YOLOv5 for wildlife detection in complex environments. *Ecological Informatics* **2024**.
9. Sun, J.; Zhang, H. YOLOv7 improvements for robust wildlife monitoring. *Ecological Informatics* **2024**.
10. Sapkota, A.; Zhang, X.; Gong, Y.; Zhang, Q. Comprehensive Performance Evaluation of YOLO11, YOLOv10, YOLOv9 and YOLOv8 on Detecting and Counting Fruitlet in Complex Orchard Environments. *arXiv preprint arXiv:2407.12040* **2024**.
11. Chen, M.; Hou, Y.; Li, H.; Gao, J.; Liu, W.; Gong, W. CEH-YOLO: An efficient YOLO-based detection model for underwater ecological images. *Ecological Informatics* **2024**, *79*, 102513. <https://doi.org/10.1016/j.ecoinf.2024.102513>.
12. Yu, X.; Liu, Z.; Xu, H.; Chen, Z.; Wu, H. YOLO-SAG: An improved YOLOv8 model for underwater target detection. *Ecological Informatics* **2024**, *80*, 102570. <https://doi.org/10.1016/j.ecoinf.2024.102570>.
13. Wang, R.; Zhang, X.; Chen, Y.; Li, P.; Zhang, C. SCoralDet: Efficient real-time underwater soft coral detection with YOLO. *Ecological Informatics* **2024**, *80*, 102562. <https://doi.org/10.1016/j.ecoinf.2024.102562>.
14. Solovyev, R.; Wang, W.; Gabruseva, T. Weighted Boxes Fusion: Ensembling Boxes from Different Object Detection Models. *Image and Vision Computing* **2021**, *117*, 104–127. <https://doi.org/10.1016/j.imavis.2021.104127>.
15. Yaseen, M. What is YOLOv8: An In-Depth Exploration of the Internal Features of the Next-Generation Object Detector, [2408.15857 [cs]]. version: 1, <https://doi.org/10.48550/arXiv.2408.15857>.
16. Terven, J.; Cordova-Esparza, D. A comprehensive review of YOLO architectures in computer vision: From YOLOv1 to YOLOv8 and beyond. *Machine Learning and Knowledge Extraction* **2023**, *5*, 1680–1716. <https://doi.org/10.3390/make5040083>.
17. Wang, S.; Li, Y.; Qiao, S. ALF-YOLO: Enhanced YOLOv8 based on multiscale attention feature fusion for ship detection. *308*, 118233. <https://doi.org/10.1016/j.oceaneng.2024.118233>.
18. Wong, T.T. Performance evaluation of classification algorithms by *k*-fold and leave-one-out cross validation. *48*, 2839–2846. <https://doi.org/10.1016/j.patcog.2015.03.009>.
19. Khalili, S.; Shakiba, A. A face detection method via ensemble of four versions of YOLOs. In Proceedings of the 2022 International Conference on Machine Vision and Image Processing (MVIP). IEEE, pp. 1–4. <https://doi.org/10.1109/MVIP53647.2022.9738779>.
20. Hosang, J.; Benenson, R.; Schiele, B. Learning Non-maximum Suppression. In Proceedings of the 2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR). IEEE, pp. 6469–6477. <https://doi.org/10.1109/CVPR.2017.685>.
21. Bodla, N.; Singh, B.; Chellappa, R.; Davis, L.S. Soft-NMS – Improving Object Detection With One Line of Code, [1704.04503 [cs]]. <https://doi.org/10.48550/arXiv.1704.04503>.
22. Solovyev, R.; Wang, W.; Gabruseva, T. Weighted boxes fusion: Ensembling boxes from different object detection models. *107*, 104117, [1910.13302 [cs]]. <https://doi.org/10.1016/j.imavis.2021.104117>.
23. Zhou, H.; Li, Z.; Ning, C.; Tang, J. CAD: Scale Invariant Framework for Real-Time Object Detection. In Proceedings of the 2017 IEEE International Conference on Computer Vision Workshops (ICCVW), pp. 760–768. ISSN: 2473-9944, <https://doi.org/10.1109/ICCVW.2017.95>.

**Disclaimer/Publisher’s Note:** The statements, opinions and data contained in all publications are solely those of the individual author(s) and contributor(s) and not of MDPI and/or the editor(s). MDPI and/or the editor(s) disclaim responsibility for any injury to people or property resulting from any ideas, methods, instructions or products referred to in the content.