

RHCNet: Residual-Guided Hierarchical Calibration Network for Robust Underwater Object Detection

Supplementary Material

6. Supplement to The Text

In this supplementary material, we present additional visual results and extended discussions that were not included in the main text. We provide a detailed overview of our visual analysis of characteristic distributions for different types of underwater degradation, along with a comprehensive qualitative comparison, which further substantiates the effectiveness of the proposed RHCNet.

7. Feature Distribution Visualization Analysis

As illustrated in Fig. 6, compared with the baseline RetinaNet, RHCNet exhibits stronger feature consistency and higher object sensitivity within object regions under complex underwater scenarios, demonstrating its superior robustness in handling background interference and semantic ambiguity.

To further evaluate the discriminative feature representation capability of the model, Fig. 7 presents the feature space distributions of the baseline and our proposed method over one thousand test images from the DUO [24] dataset, where the baseline features appear widely scattered. In contrast, our method demonstrates markedly higher semantic compactness, reflecting its superior ability to capture and represent discriminative features.

8. Comprehensive Qualitative Comparison

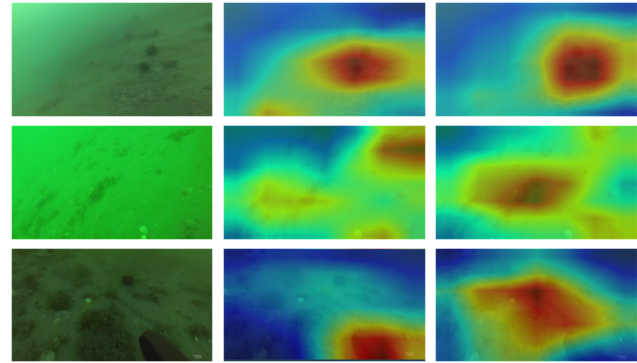
Fig. 8, Fig. 9, and Fig. 10 present qualitative comparison results for three typical complex underwater scenarios on the DUO [24] dataset.

1) Similar Foreground and Background Scenes. RHCNet effectively suppresses false activations in background regions during qualitative comparisons, resulting in more accurate object localization.

2) Object boundary degradation Scenes. Detection boxes generated by RHCNet more accurately cover object areas, markedly reducing instances of missed detections outside the box or redundant coverage within it caused by boundary blurring.

3) Low-Visibility Underwater Scenes. RHCNet maintains stable object coverage even under low-light and highly turbid conditions, demonstrating superior detection confidence and localization consistency compared with other state-of-the-art methods.

Fig. 11, Fig. 12, and Fig. 13 illustrate the corresponding results on the UTDAC dataset, encompassing degradation types consistent with those in the DUO dataset, where



(a) Raw image (b) Baseline (c) Ours

Figure 6. Visual comparison of feature heatmaps for the raw image, Baseline and Ours in underwater images.

RHCNet consistently surpasses six state-of-the-art counterparts in terms of object localization accuracy and boundary compactness, thereby further validating its strong generalization capability under diverse and challenging underwater conditions.

9. More Implementation Details

9.1. Data Augmentation

During training, a concise yet effective data augmentation strategy was adopted, where images were randomly rescaled across multiple resolutions with the shorter side sampled between 800 and 1920 pixels while maintaining the original aspect ratio, thereby substantially improving the model’s capacity to capture multi-scale object variations particularly beneficial in underwater scenarios with significant object size diversity; moreover, no additional preprocessing techniques such as histogram equalization, CLAHE contrast enhancement, or style transfer were utilized, ensuring that all observed performance gains arise primarily from the proposed architecture and its enhanced feature representation capability.

9.2. Evaluation Metrics

To ensure fairness and rigor in the evaluation process, this paper uniformly employs the official COCO object detection metrics for performance assessment across all experiments. Details are as follows.

Average Precision (AP): This metric represents the average accuracy achieved across a range of IoU thresholds

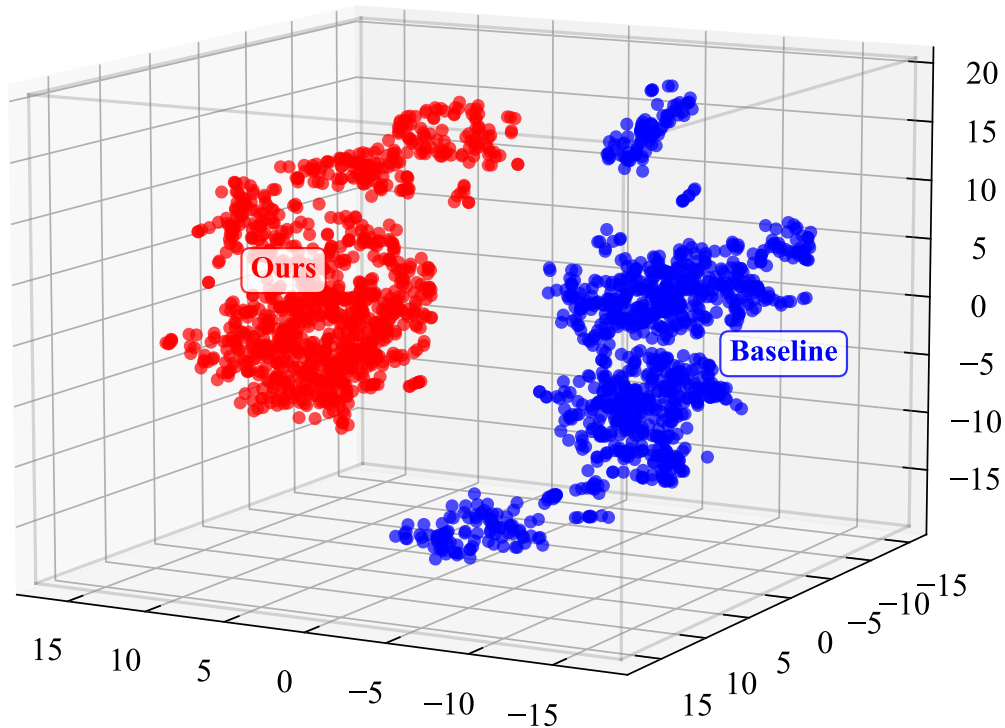


Figure 7. Scatter plot comparing feature distributions of our RHCNet and baseline methods in three-dimensional feature space.

781 from 0.50 to 0.95, incremented by 0.05. It serves to evalu-
782 ate a model’s overall detection capability.

783 AP_{50}/AP_{75} : These denote the accuracy under condi-
784 tions of $IoU=0.50$ and $IoU=0.75$, respectively. The former
785 measures coarse-grained localization capability, while the
786 latter emphasizes precise detection.

787 $AP_S/AP_M/AP_L$: These denote detection accuracy for
788 small object areas (less than 32^2 pixels²), medium object
789 areas (between 32^2 pixels² and 96^2 pixels²), and large ob-
790 ject areas (greater than 96^2 pixels²), respectively. They are
791 employed to analyse a model’s robustness across different
792 object scales.

793 All the aforementioned metrics were uniformly tested on
794 datasets such as DUO [24], UTDAC [28], and COCO [21]
795 to validate the model’s robustness and stability in complex
796 underwater environments and general detection scenarios.

797 10. Additional Ablation Study

798 To validate the effectiveness of each component within
799 RHCNet, we provide extended ablation results and visual-
800 ized performance trend analyses based on Tab. 2 in the main
801 text. As shown in Fig. 14, the metrics AP, AP_{50} , AP_{75} ,

AP_S , AP_M , and AP_L on the DUO [24] dataset exhibit a
802 stable upward trend with increasing module stacking. 803

804 To further validate the effectiveness of each module, we
805 conducted additional ablation experiments on the UTDAC
806 [28] dataset, with results presented in Tab. 4. The over-
807 all findings demonstrate that the complete model (Case 7)
808 achieves optimal performance across all evaluation met-
809 rics, thereby corroborating the positive contributions of the
810 LAM, RGFE, PAM, and CGCA modules in enhancing the
811 model’s perceptual capabilities and detection robustness.

812 As shown in Fig. 15, our method exhibits a stable up-
813 ward trend across the AP, AP_{50} , AP_{75} , AP_S , AP_M , and
814 AP_L metrics on the UTDAC dataset as modules are stacked. 815

816 11. Future Work and Limitation

817 Although the proposed RHCNet demonstrates outstanding
818 performance across multiple evaluation metrics and ex-
819 hibits strong robustness in complex underwater scenarios,
820 several areas warrant further investigation alongside ex-
821 isting limitations. Firstly, RHCNet may still exhibit in-
822 sufficient feature representation when processing severely

Table 4. In the UTDAC [28] dataset, the ablation experiments of the key modules in RHCNet (\checkmark = MODULE RETAINED, BLANK = MODULE ABLATED).

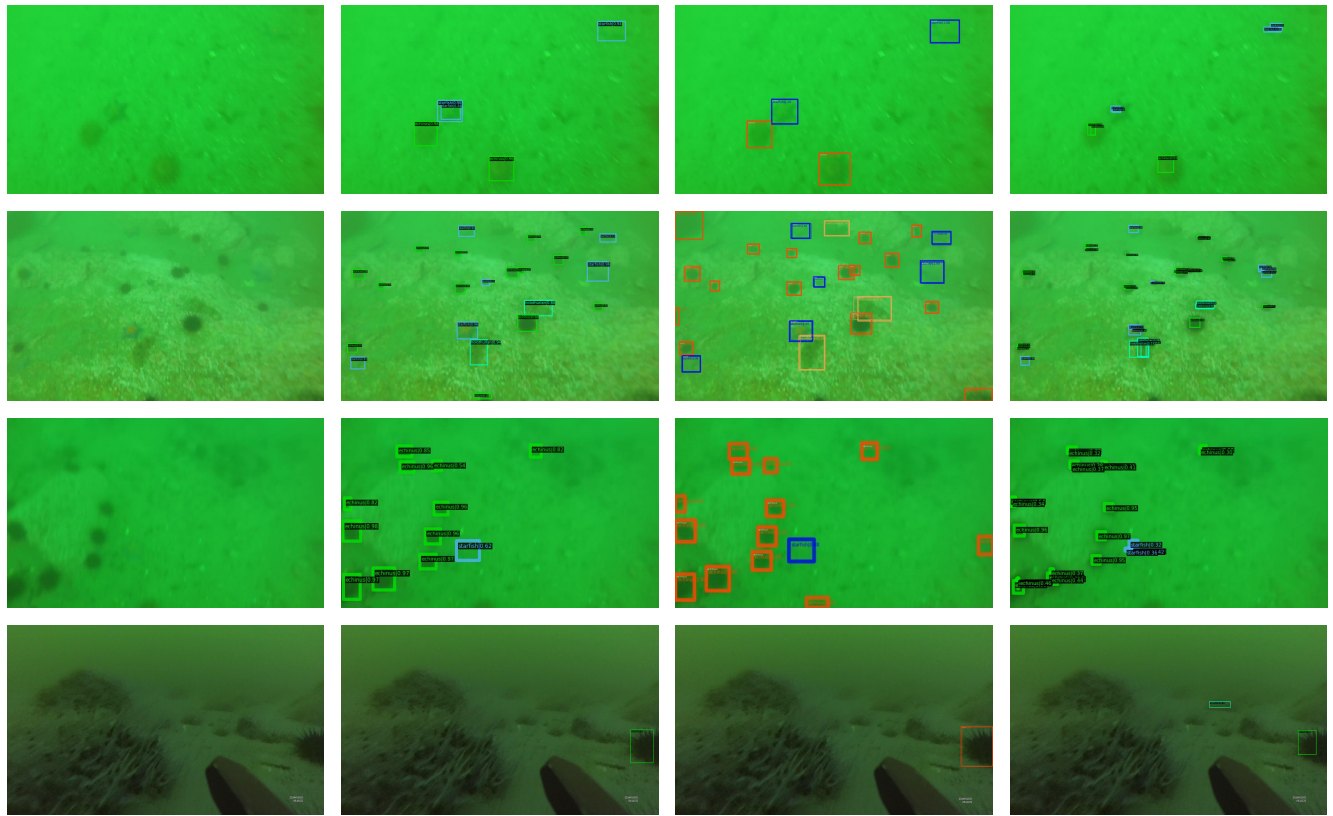
Cases	LAM	RGFE	PAM	CGCA	AP	AP ₅₀	AP ₇₅	AP _S	AP _M	AP _L
RetinaNet					46.02	70.53	48.36	19.57	40.34	51.23
1		\checkmark	\checkmark	\checkmark	48.96	83.42	53.76	23.20	45.07	54.71
2	\checkmark		\checkmark	\checkmark	48.45	83.11	53.64	22.68	44.82	54.19
3			\checkmark	\checkmark	47.13	82.62	52.49	22.07	43.90	53.37
4	\checkmark	\checkmark		\checkmark	48.84	83.26	53.31	22.75	44.83	54.08
5	\checkmark	\checkmark	\checkmark		47.67	82.81	52.64	22.18	44.01	53.69
6	\checkmark	\checkmark			47.54	81.75	52.13	21.08	43.87	53.34
7	\checkmark	\checkmark	\checkmark	\checkmark	50.87	85.86	55.37	24.82	46.23	56.42

823 degraded images (e.g., intense light refraction, significant
824 suspended particle occlusion), indicating that the model’s
825 generalization capability under extreme conditions requires
826 enhancement. Secondly, the K-means cluster-guided fea-
827 ture calibration module depends on unsupervised cluster-
828 ing priors. It remains sensitive to the predefined number
829 of clusters, indicating a potential benefit from incorporat-
830 ing adaptive or learnable clustering strategies. Thirdly, the
831 model’s performance relies heavily on large-scale, well-
832 annotated training datasets, yet obtaining high-quality un-
833 derwater labels is both costly and labor-intensive. There-
834 fore, future work should investigate strategies such as self-
835 supervised representation learning, semi-supervised fine-
836 tuning, or synthetic-to-real domain adaptation to enhance
837 performance in data-scarce or weakly supervised settings.

838 12. Broader Impacts

839 The proposed RHCNet framework holds substantial signif-
840 icance beyond underwater target detection, demonstrating
841 strong potential for broader visual perception and cross-
842 domain detection tasks. In marine resource exploration,
843 RHCNet enhances automated recognition of marine or-
844 ganisms and submerged infrastructure, thereby improving
845 both the efficiency and precision of oceanographic surveys.
846 Within ecological conservation and environmental monitor-
847 ing, it enables accurate tracking of coral reef dynamics,
848 fish populations, and endangered species, offering essen-
849 tial technological support for the sustainable development
850 of marine ecosystems. In military and engineering contexts,
851 RHCNet enables autonomous perception and positioning
852 for underwater unmanned systems, thereby enhancing mis-
853 sion reliability and operational safety. More broadly, its
854 residual attention mechanism and hierarchical calibration
855 strategy provide a generalized solution to visual recognition
856 challenges under ambiguous conditions, extending its appli-
857 cability to low-light surveillance, medical imaging analysis,
858 and remote sensing interpretation. Nevertheless, the large-

scale deployment of intelligent underwater perception sys-
859 tems may raise ethical and privacy concerns, particularly in
860 defense and security domains; hence, future research should
861 not only emphasize technical optimization but also incor-
862 porate considerations of social responsibility and sustain-
863 ability, ensuring that technological advancement promotes
864 scientific innovation while upholding environmental ethics
865 and human values.
866

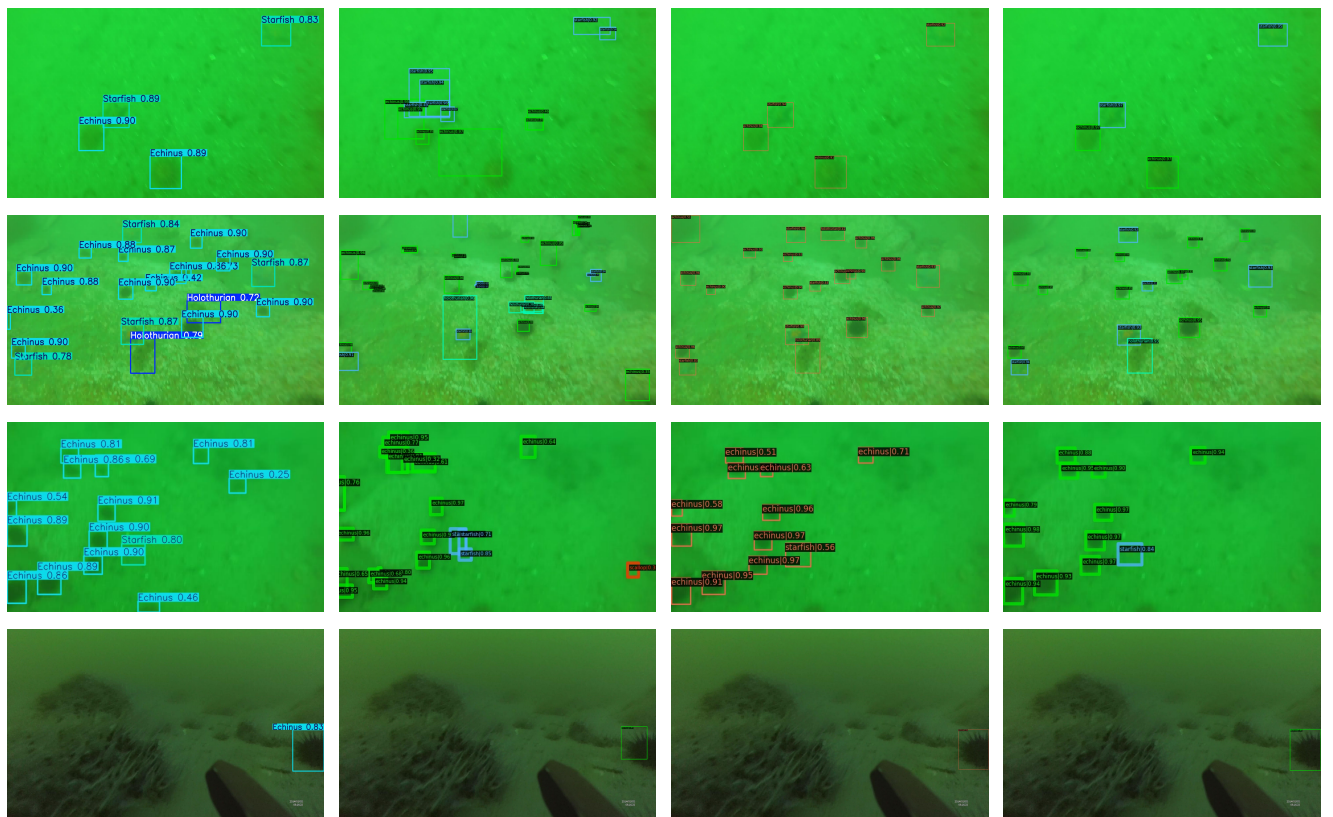


(a) Input

(b) ERLNet [7]

(c) DJENet [32]

(d) UDMD [38]



(e) YOLOv11 [15]

(f) SqNet [3]

(g) CIDNet [42]

(h) Ours

Figure 8. Qualitative evaluation of underwater scenarios with similar foreground and background in the DUO [24] dataset.

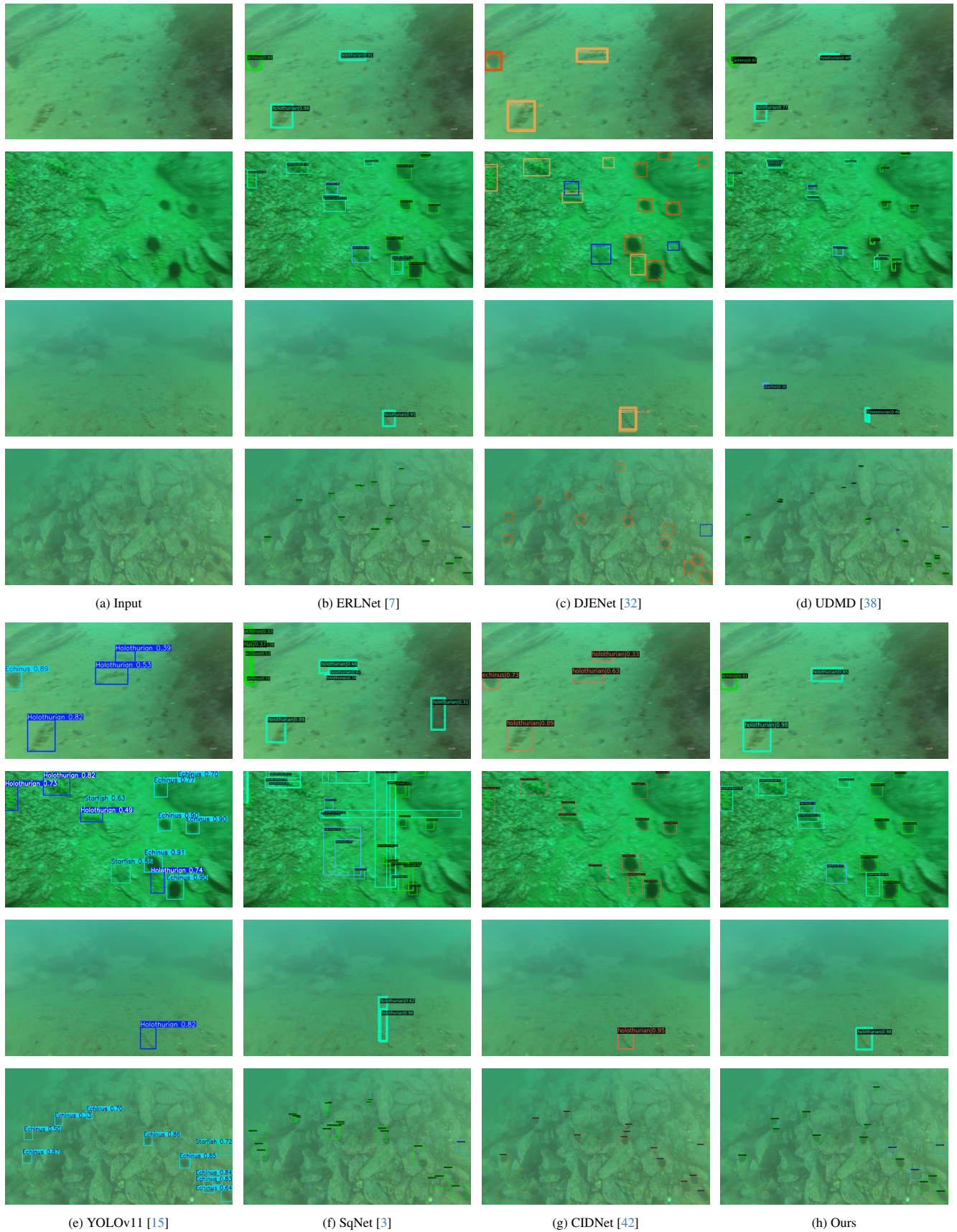


Figure 9. Qualitative evaluation of underwater scenarios with object edge degradation in the DUO [24] dataset.

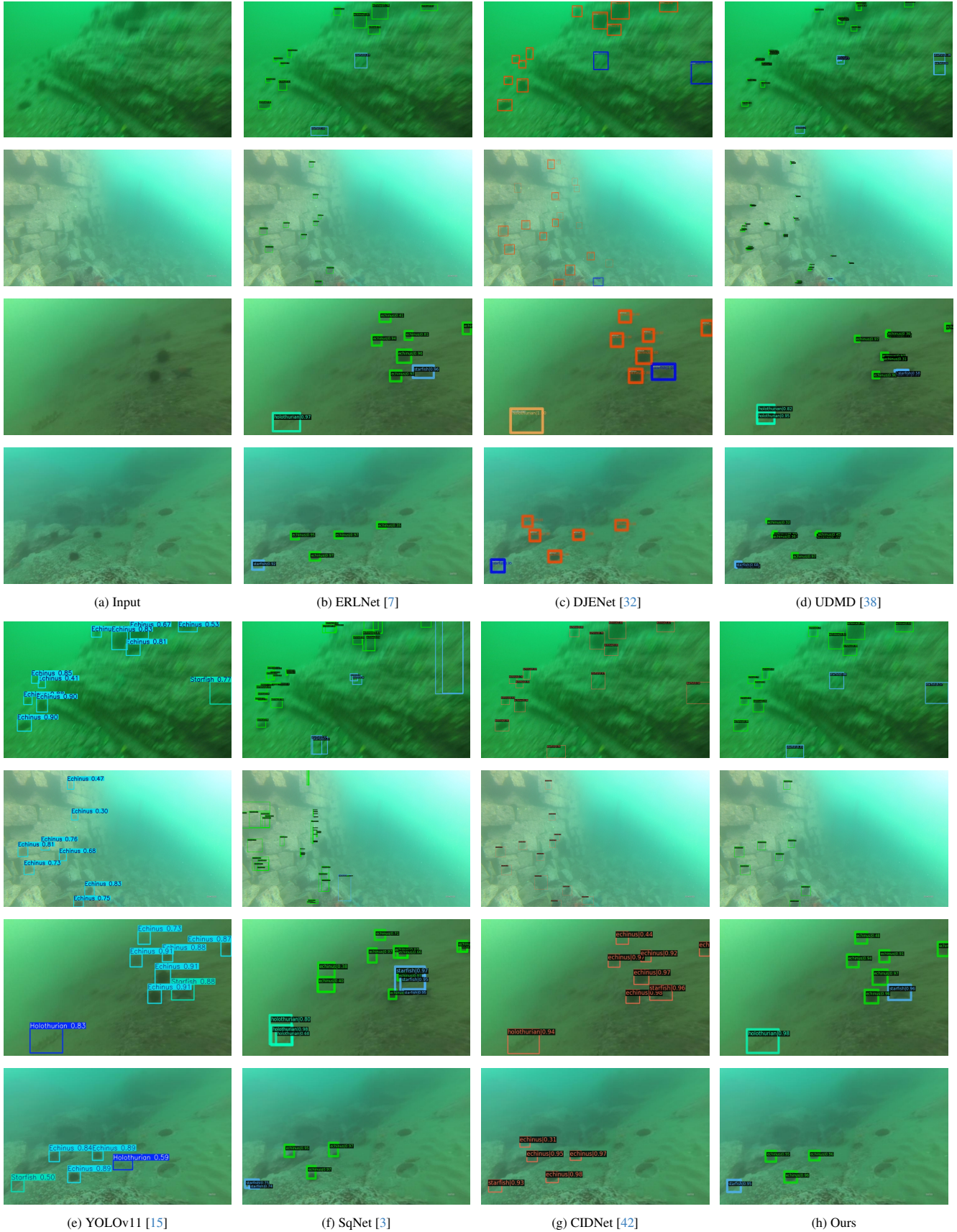


Figure 10. Qualitative evaluation of low-visibility underwater scenes in the DUO [24] dataset.

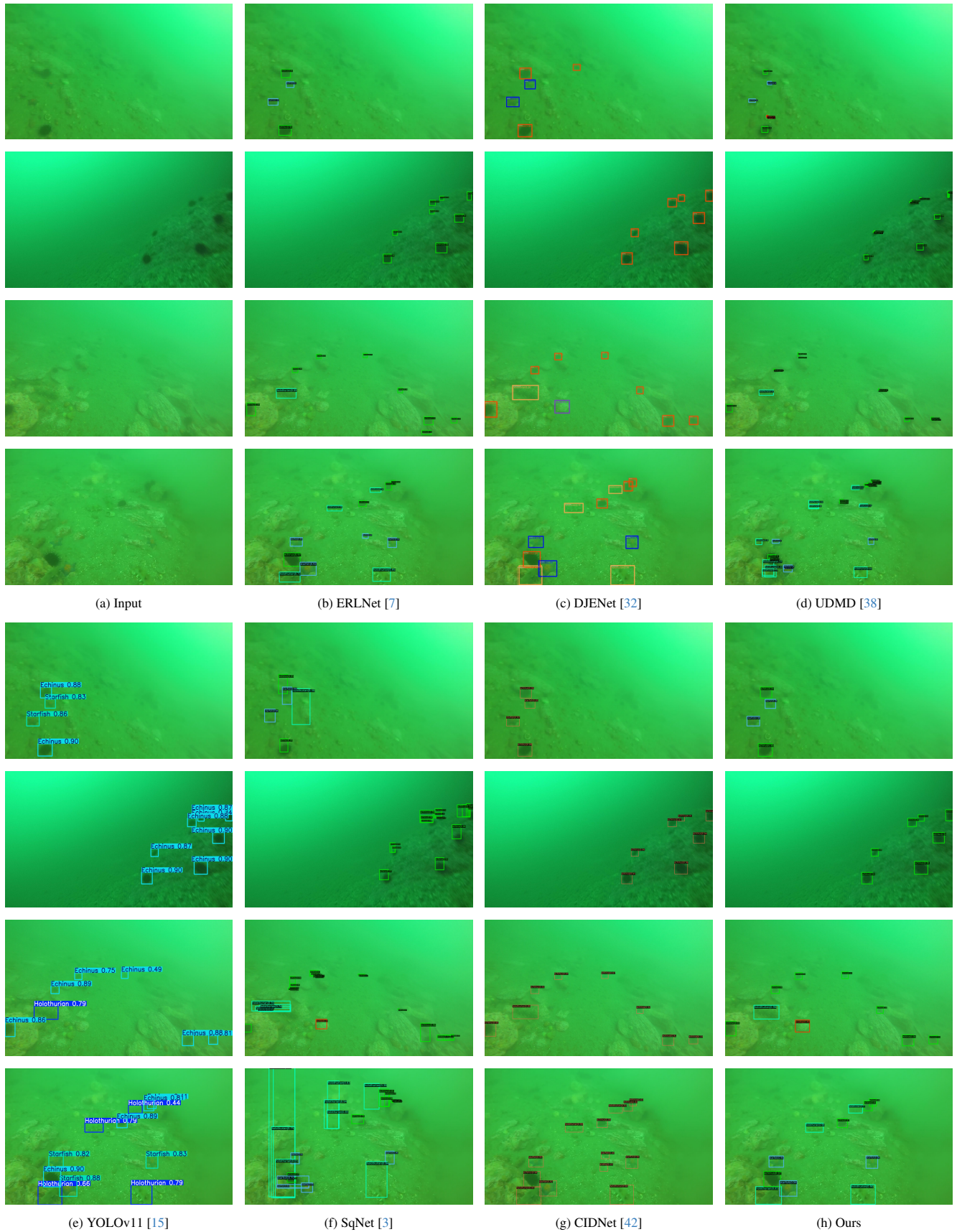


Figure 11. Qualitative evaluation on underwater scenes with similar foreground and background in the UTDAC [28] dataset.

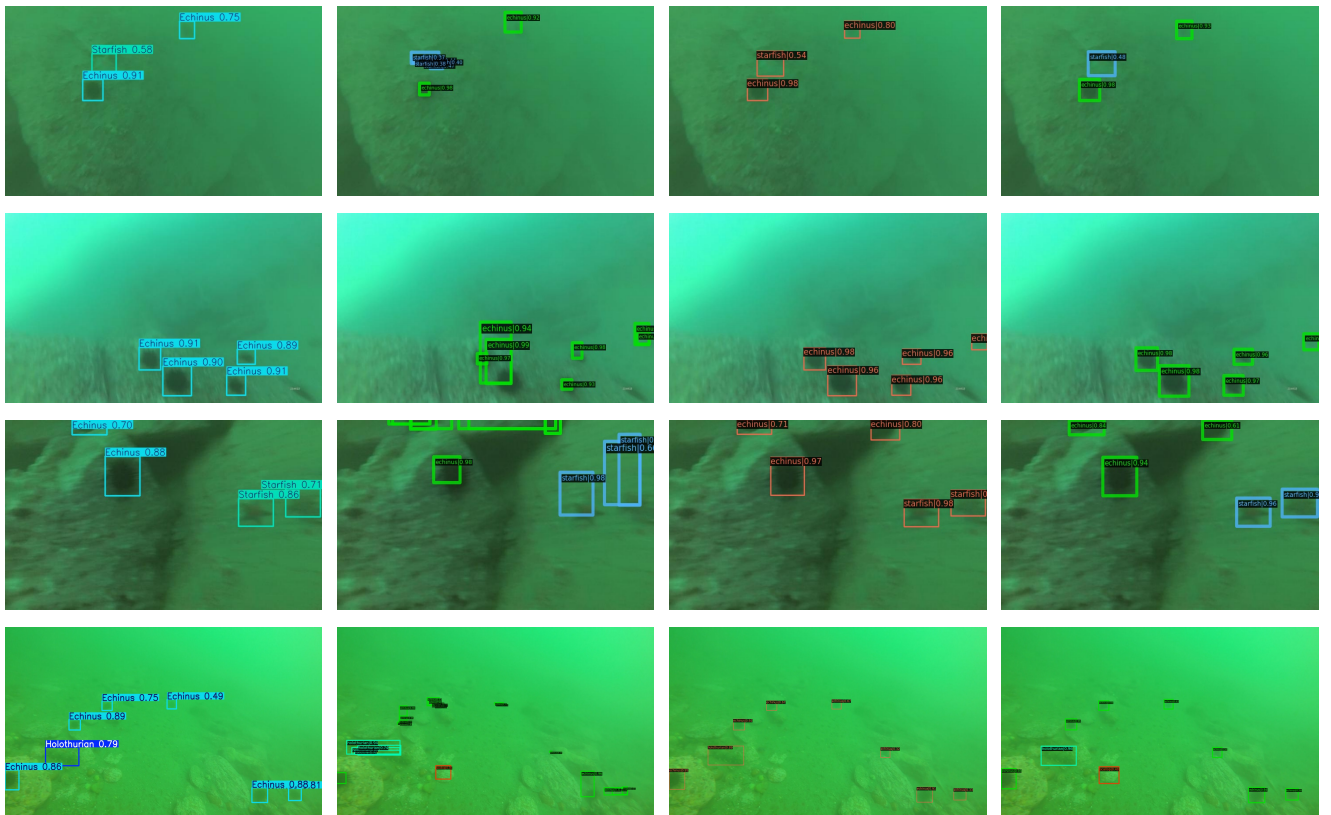
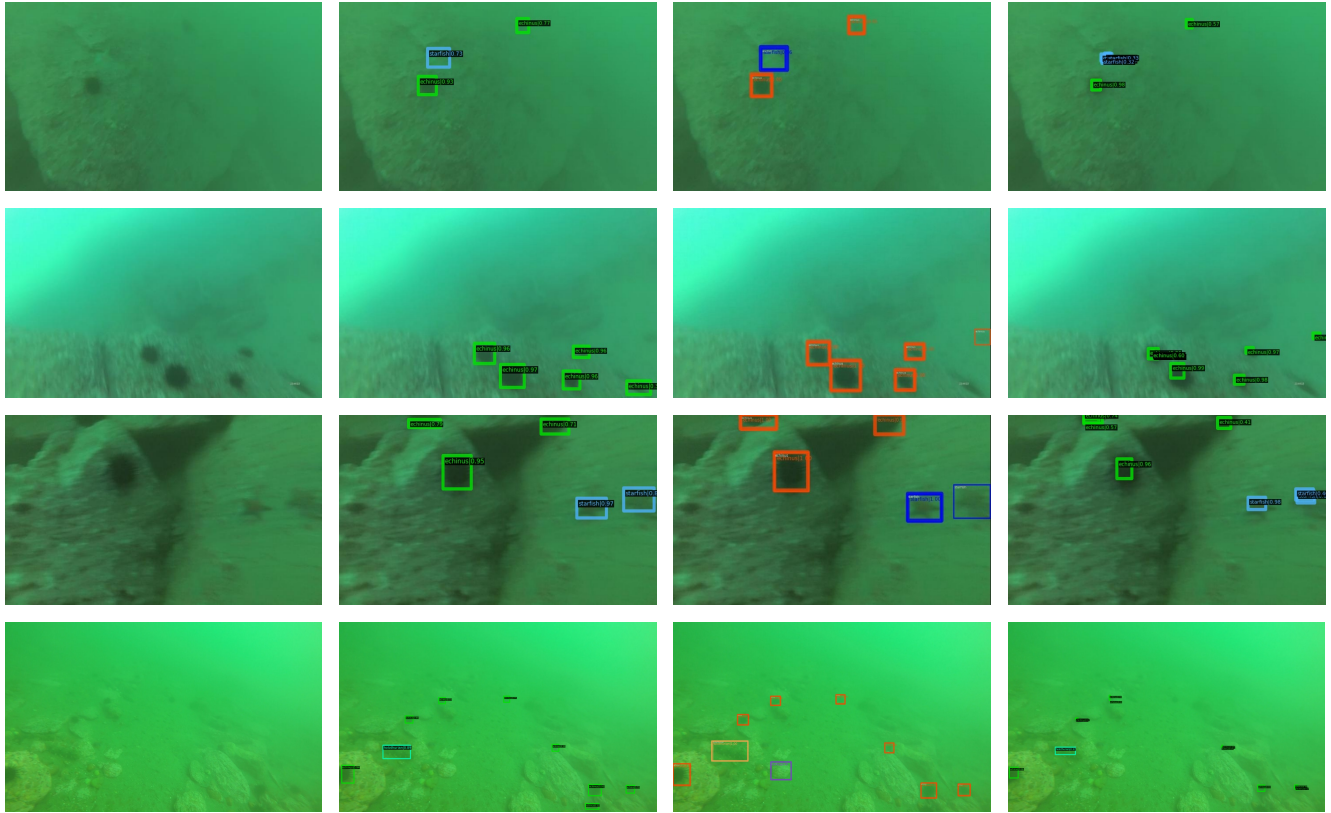
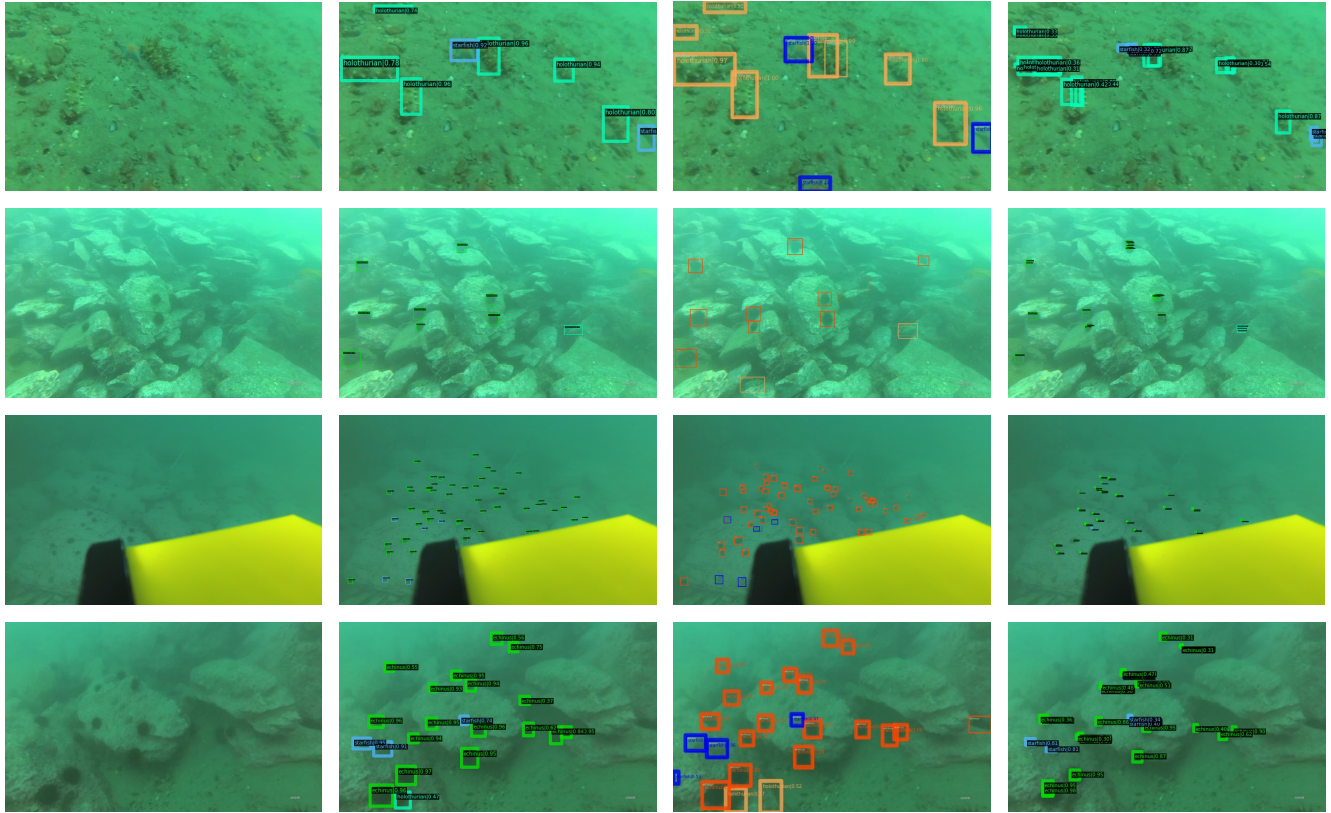


Figure 12. Qualitative evaluation of underwater scenarios with object edge degradation in the UTDAC [28] dataset.

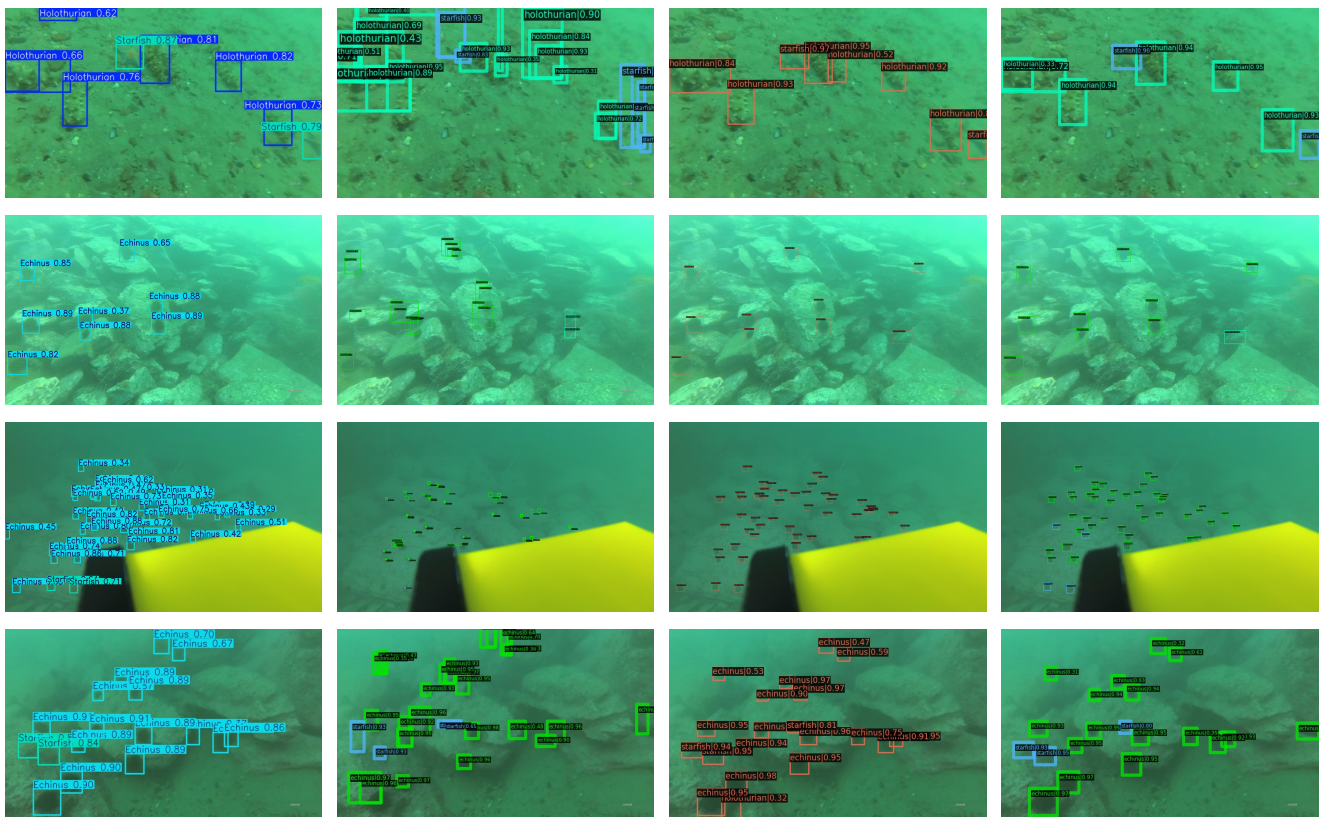


(a) Input

(b) ERLNet [7]

(c) DJENet [32]

(d) UDMD [38]



(e) YOLOv11 [15]

(f) SqNet [3]

(g) CIDNet [42]

(h) Ours

Figure 13. Qualitative evaluation of low-visibility underwater scenes in the UTDAC [28] dataset.

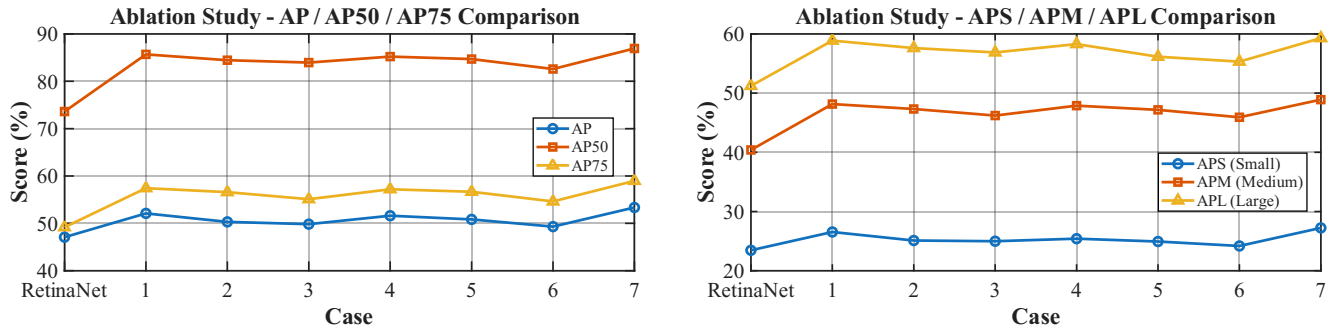


Figure 14. Line chart of ablation experiments on the DUO [24] dataset

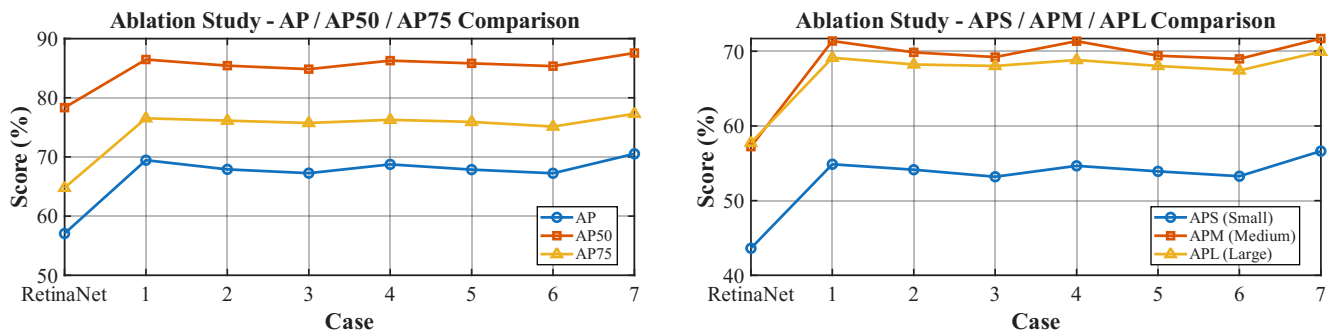


Figure 15. Line chart of ablation experiments on the UTDAC [28] dataset