

Additional Material

1. Loss Formulation

Score Distillation Sampling (SDS). During distillation, the frozen Diffusion Transformer (DiT) guides the Retinex-UNet to align with its denoising prior. The SDS gradient is computed as:

$$\nabla_{\theta} \mathcal{L}_{\text{SDS}} = \mathbb{E}_{t, \epsilon} \left[w(t) (\epsilon_{\theta}(z_t; t, y) - \epsilon) \frac{\partial \hat{f}}{\partial \theta} \right], \quad (1)$$

where z_t is the noisy latent of \hat{f} at timestep t . This enables direct color-prior transfer without iterative diffusion.

Inter-Component Residual (ICR) Regularization. To reduce coupling between reflectance (R) and illumination (L), we apply two lightweight penalties that require no architectural change. The first penalizes cross-component similarity:

$$\mathcal{L}_{\text{icr-cos}} = \mathbb{E} \left[\frac{\langle \text{vec}(R), \text{vec}(L) \rangle}{\|\text{vec}(R)\|_2 \|\text{vec}(L)\|_2 + \epsilon} \right], \quad (2)$$

where $\epsilon > 0$ ensures numerical stability. The second discourages co-occurring strong gradients:

$$\mathcal{L}_{\text{icr-grad}} = \|\nabla R \odot \nabla L\|_1. \quad (3)$$

Together, these losses suppress halos, ghosting, and color drift while preserving structural fidelity.

Overall Objective. The Retinex-UNet is optimized with the joint loss:

$$\mathcal{L}_{\text{UNet}} = \mathcal{L}_{\text{SDS}} + \lambda_{\text{icr}} (\mathcal{L}_{\text{icr-cos}} + \mathcal{L}_{\text{icr-grad}}), \quad (4)$$

where λ_{icr} controls the regularization strength. Inference requires only a single UNet forward pass to produce (R, L) and \hat{f} , achieving stable, artifact-free color reconstruction.

2. Training Data Construction via Physical Prior

To provide the Diffusion Transformer (DiT) with a robust, physics-guided prior, we synthesize underwater images from the NYU-Depth V2 dataset [SHKF12] using a physical image-formation model, as illustrated in Figure 1. Specifically, for each RGB channel $c \in \{R, G, B\}$, we apply:

$$I_c^{\text{uw}} = I_c e^{-\eta_c d} + B_c (1 - e^{-\eta_c d}), \quad (5)$$

where $\eta_c d \in [0, 5]$ and $B_c \in [0, 1]$ denote the wavelength-dependent attenuation and background light, respectively. As shown in Figure 1, this process generates synthetic underwater training images from terrestrial data (top), while the model is evaluated on real

underwater scenes (bottom). Several visual examples of the synthesized data are further illustrated in Figure 2, demonstrating the diverse color casts and visibility degradations captured by the physical prior.

To verify the fidelity of the synthetic dataset, we analyze its color statistics in the CIELAB space. Following [PZB23, QT09], we randomly sample 2000 images each from LSUI (real underwater), MIT67 (terrestrial), and our synthesized dataset, computing the mean and standard deviation of the a and b channels. As summarized in Table 2, our synthesized data exhibit close chromatic alignment with real underwater images while remaining distinct from the terrestrial domain, confirming that the physical-prior formulation effectively models underwater color characteristics.

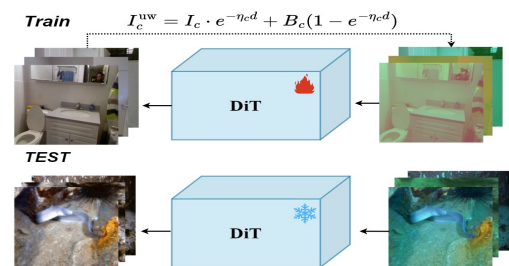


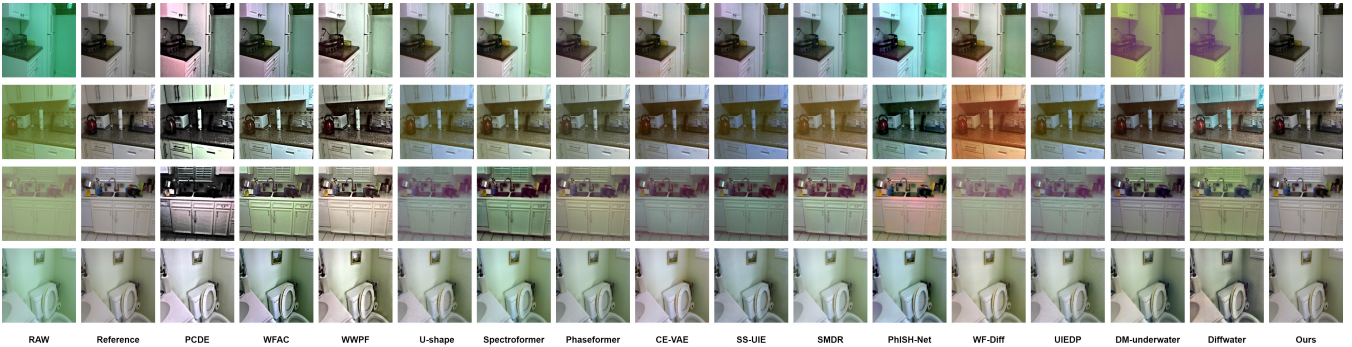
Figure 1: Training and testing pipeline. The Train set is synthesized from terrestrial images using a physical prior, while testing is performed on real underwater images.



Figure 2: Examples of synthetic underwater images generated from terrestrial images using the physical prior. The results exhibit diverse color casts and illumination degradations, enriching the training domain.

Table 1: PSNR and SSIM results on the synthetic underwater dataset. Higher is better.

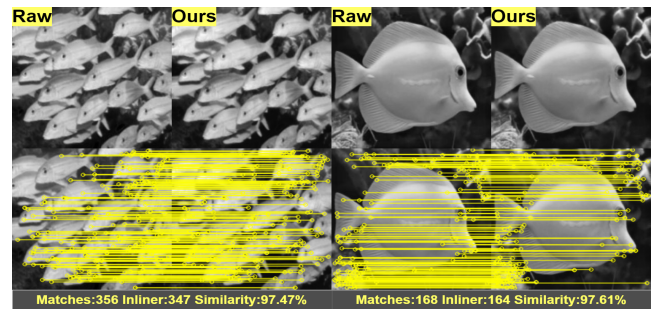
Metric	PCDE [ZZ ²³]	WFAC [ZLL ²⁵]	WWPF [ZZZ ²⁴]	U-shape [PZB ²³]	Spectroformer [KMM ²⁴]	Phaseformer [KNK ²⁵]	CE-VAE [PM ²⁵]	SS-UIE [PB ²⁵]	SMDR [ZG ²⁴]	PhISH-Net [CSB ²⁴]	WF-Diff [ZCDH ²⁴]	UIEDP [DOS ²¹]	DM-underwater [TK ²³]	Diffwater [GX ²³]	RetiDiff (Ours)
PSNR ↑	21.02	23.70	23.29	26.16	25.71	30.92	27.82	28.16	27.25	26.48	27.02	26.08	23.69	23.12	39.63
SSIM ↑	0.71	0.73	0.81	0.79	0.81	0.85	0.86	0.86	0.83	0.81	0.78	0.78	0.72	0.69	0.97

**Figure 3:** Qualitative comparison on the synthetic underwater dataset. **RetiDiff** effectively restores color distributions to match terrestrial references, while other methods retain domain-specific color bias or structural artifacts.**Table 2:** Statistical validation of synthesized images in CIELAB space (a, b channels). The close alignment between our synthesized dataset and real underwater images validates the chromatic fidelity achieved through the physics-guided data generation process.

Categories	a-mean	a-std	b-mean	b-std
Underwater Images	-19.54	12.58	-1.64	14.54
Land Images	2.61	3.05	4.06	4.95
Synthesized Images	-17.15	10.12	0.79	11.95

4. SIFT-Based Structural Consistency

To further validate the structural fidelity of **RetiDiff**, we employ SIFT-based keypoint consistency analysis [Low04, AAHB12]. The SIFT algorithm extracts and matches local features between the input underwater image I^{uw} and the reconstructed output \hat{I} . As illustrated in Figure 4, **RetiDiff** achieves dense and accurate feature correspondences, where high inlier ratios and strong match similarity indicate that spatial structures are well preserved during enhancement rather than distorted by color correction. Quantitatively, Table 3 reports average matches, inliers, and inlier ratios across UIEB, LSUI, and TEST-U45, showing $\sim 97\%$ keypoint similarity, which corroborates the model’s strong structural consistency across diverse scenes.

**Figure 4:** SIFT-based keypoint matching between grayscale maps of the input (left) and reconstructed image (right). Yellow lines denote correspondences; overlays report match, inlier, and similarity statistics. **RetiDiff** achieves $\sim 97\%$ keypoint similarity, evidencing structural preservation.

3. Results on the Synthetic Dataset

To further verify the color reconstruction ability of **RetiDiff**, we evaluate all methods on our physics-based synthetic underwater dataset, where the ground truth terrestrial counterparts are available. This setting enables direct observation of how well each model restores the color distribution of underwater images toward that of land scenes.

As shown in Figure 3, traditional CNN-, GAN-, and transformer-based approaches exhibit residual color casts or contrast imbalance, failing to recover the natural color statistics of terrestrial domains. In contrast, **RetiDiff** produces balanced illumination and chromatic consistency closely aligned with the reference land images, demonstrating its ability to reconstruct the true color distribution rather than merely enhancing visual contrast.

We also provide quantitative PSNR and SSIM results in Table 1 to complement perceptual analysis. **RetiDiff** achieves the highest fidelity scores, confirming its capacity to recover both accurate color distributions and structural consistency under synthetic-to-real generalization.

The results show that across UIEB, LSUI, and TEST-U45

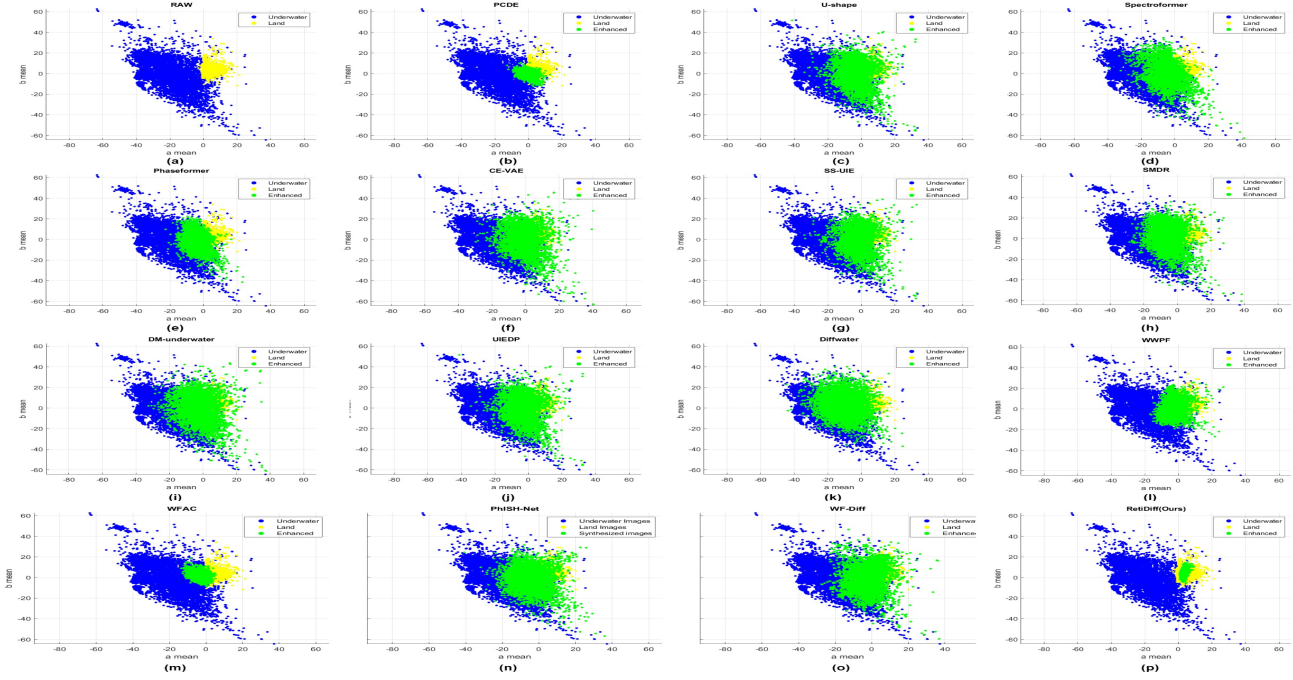


Figure 5: Distributions of underwater (blue), enhanced (green), and terrestrial (yellow) images in the CIELAB color space. Each point denotes one image based on the mean a - b values. RetiDiff shifts the chromatic distribution of enhanced images toward the terrestrial cluster, reducing color bias caused by underwater absorption.

Table 3: SIFT-based keypoint consistency between input I^{uv} and output \hat{I} . We report the average number of inliers, total matches, and inlier ratio (Average similarity).

Metrics	UIEB	LSUI	TEST-U45
Average inliers	278.86	327.45	292.42
Average matches	287.10	336.27	301.22
Average similarity (%)	97.13	97.38	97.08

76 datasets, **RetiDiff** maintains $\sim 97\%$ keypoint similarity, confirming
 77 that its color reconstruction process preserves geometric structures
 78 and fine spatial details.

79 **5. CIELAB Color Distribution Analysis**

80 To visually evaluate the effectiveness of **RetiDiff** in restoring nat-
 81 ural color tones, we analyze the distribution of enhanced images
 82 in the CIELAB color space. Each image is represented by the
 83 mean values of its a and b channels, which describe chromatic-
 84 ity along the red-green and yellow-blue axes. As shown in Fig-
 85 ure 5, underwater images (blue) cluster toward negative a and b
 86 values due to wavelength-dependent attenuation, while the enhanced
 87 results (green) shift markedly toward the terrestrial domain (yel-
 88 low). This movement indicates that **RetiDiff** effectively corrects
 89 underwater color casts and reconstructs color distributions consis-

90 tent with those of land images, demonstrating its capacity for phys-
 91 ically grounded and visually natural color restoration.

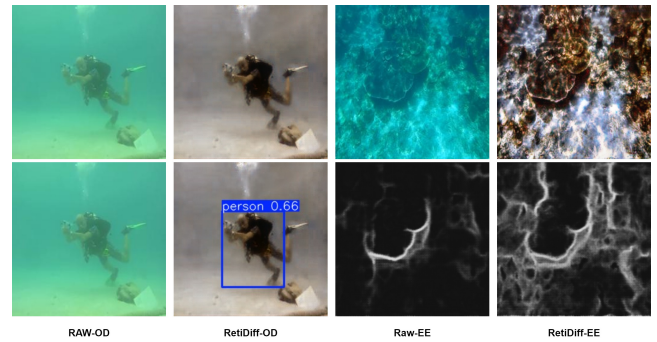


Figure 6: Applicability of **RetiDiff** in downstream tasks. Left: Ob-
 ject detection (OD) under YOLOv7 [WBL23], where the human tar-
 get undetected in the raw image becomes detectable after RetiDiff
 enhancement. Right: Sobel edge extraction (EE), where RetiDiff re-
 stores fine details and continuous contours in reef textures.

92 **6. Applicability Analysis**

93 To further demonstrate the practical value of **RetiDiff**, we evalu-
 94 ate its impact on downstream visual tasks, including object detec-

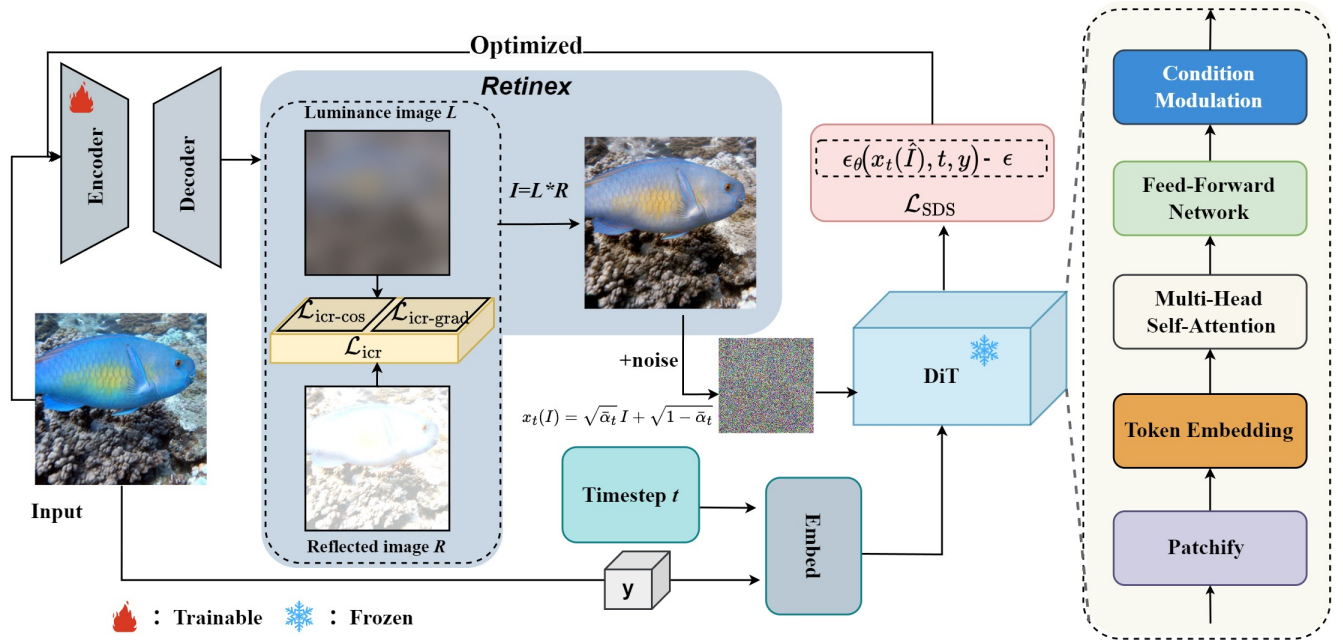


Figure 7: Architecture of RetiDiff.

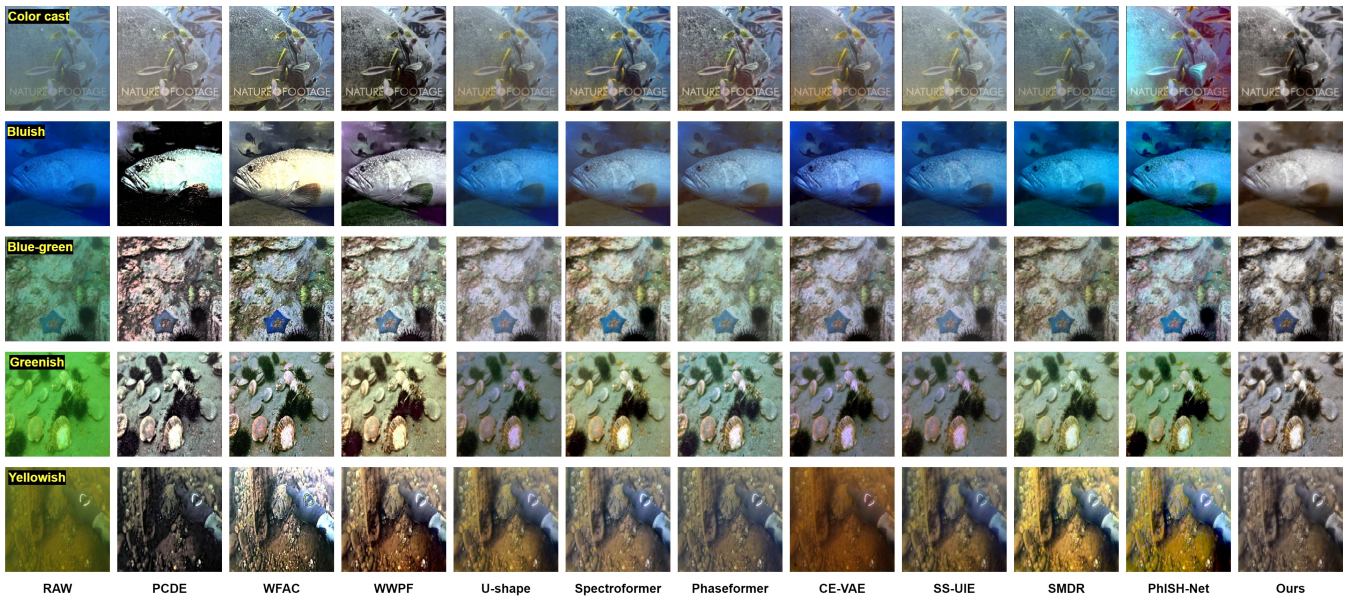


Figure 8: Qualitative comparison of the proposed method with state-of-the-art approaches on real underwater images.

95 tion and edge extraction. As illustrated in Figure 6, the left pair
 96 (RAW–OD vs. RetiDiff–OD) presents object detection results under
 97 YOLOv7 [WBL23]. In the raw underwater image, the human
 98 target is barely visible and fails to be detected due to low contrast
 99 and severe color degradation. After enhancement by **RetiDiff**, the
 100 restored image exhibits balanced chromaticity and higher local con-
 101 trast, enabling successful detection with tighter bounding boxes and
 102 higher confidence.

103 The right pair (RAW–EE vs. RetiDiff–EE) shows Sobel-based
 104 edge extraction on a reef scene. While the raw image produces
 105 weak and discontinuous edge responses, the RetiDiff-enhanced re-
 106 sult reveals finer structural details and more continuous contours,
 107 especially around coral textures and object boundaries. These re-
 108 sults confirm that **RetiDiff** not only restores perceptual realism but
 109 also provides feature-consistent enhancements beneficial to down-
 110 stream vision tasks in underwater environments.

111 References

- 112 [AAHB12] ANCUTI C., ANCUTI C. O., HABER T., BEKAERT P.: En-
 113 hancing underwater images and videos by fusion. In *2012 IEEE Con-
 114 ference on Computer Vision and Pattern Recognition* (2012), pp. 81–88.
 115 doi:10.1109/CVPR.2012.6247661. 2
- 116 [CSB24] CHANDRASEKAR A., SREENIVAS M., BISWAS S.: Phish-
 117 net: Physics inspired system for high resolution underwater image en-
 118 hancement. In *2024 IEEE/CVF Winter Conference on Applications of
 119 Computer Vision (WACV)* (2024), pp. 1495–1505. doi:10.1109/
 120 WACV57701.2024.00153. 2
- 121 [DDSX*] DU D L. E., SI L., XU F., NIU J., SUN F.: Uiedp: under-
 122 water image enhancement with diffusion prior (2023). *arXiv preprint
 123 arXiv:2312.06240*. 2
- 124 [GXJ*23] GUAN M., XU H., JIANG G., YU M., CHEN Y., LUO T.,
 125 ZHANG X.: Diffwater: Underwater image enhancement based on condi-
 126 tional denoising diffusion probabilistic model. *IEEE Journal of Selected
 127 Topics in Applied Earth Observations and Remote Sensing* 17 (2023),
 128 2319–2335. 2
- 129 [KMM*24] KHAN R., MISHRA P., MEHTA N., PHUTKE S. S., VIP-
 130 PARTHI S. K., NANDI S., MURALA S.: Spectroformer: Multi-domain
 131 query cascaded transformer network for underwater image enhancement.
 132 In *Proceedings of the IEEE/CVF winter conference on applications of
 133 computer vision* (2024), pp. 1454–1463. 2
- 134 [KNK*25] KHAN M. R., NEGI A., KULKARNI A., PHUTKE S. S., VIP-
 135 PARTHI S. K., MURALA S.: Phaseformer: Phase-based attention mech-
 136 anism for underwater image restoration and beyond. In *2025 IEEE/CVF
 137 Winter Conference on Applications of Computer Vision (WACV)* (2025),
 138 IEEE, pp. 9618–9629. 2
- 139 [Low04] LOWE D. G.: Distinctive image features from scale-invariant
 140 keypoints. *International journal of computer vision* 60, 2 (2004), 91–
 141 110. 2
- 142 [PB25] PENG L., BIAN L.: Adaptive dual-domain learning for under-
 143 water image enhancement. In *Proceedings of the AAAI Conference on
 144 Artificial Intelligence* (2025), vol. 39, pp. 6461–6469. 2
- 145 [PM25] PUCCI R., MARTINEL N.: Ce-vae: Capsule enhanced variational
 146 autoencoder for underwater image enhancement. In *2025 IEEE/CVF
 147 Winter Conference on Applications of Computer Vision (WACV)* (2025),
 148 IEEE, pp. 2113–2123. 2
- 149 [PZB23] PENG L., ZHU C., BIAN L.: U-shape transformer for under-
 150 water image enhancement. *IEEE transactions on image processing* 32
 151 (2023), 3066–3079. 1, 2
- 152 [QT09] QUATTONI A., TORRALBA A.: Recognizing indoor scenes.
 153 In *2009 IEEE conference on computer vision and pattern recognition*
 154 (2009), IEEE, pp. 413–420. 1
- 155 [SHKF12] SILBERMAN N., HOIEM D., KOHLI P., FERGUS R.: Indoor
 156 segmentation and support inference from rgbd images. In *European con-
 157 ference on computer vision* (2012), Springer, pp. 746–760. 1
- 158 [TKI23] TANG Y., KAWASAKI H., IWAGUCHI T.: Underwater image
 159 enhancement by transformer-based diffusion model with non-uniform
 160 sampling for skip strategy. In *Proceedings of the 31st ACM International
 161 Conference on Multimedia* (2023), pp. 5419–5427. 2
- 162 [WBL23] WANG C.-Y., BOCHKOVSKIY A., LIAO H.-Y. M.: Yolov7:
 163 Trainable bag-of-freebies sets new state-of-the-art for real-time object
 164 detectors. In *Proceedings of the IEEE/CVF conference on computer vi-
 165 sion and pattern recognition* (2023), pp. 7464–7475. 3, 5
- 166 [ZCDH24] ZHAO C., CAI W., DONG C., HU C.: Wavelet-based fourier
 167 information interaction with frequency diffusion adjustment for under-
 168 water image restoration. In *Proceedings of the IEEE/CVF Conference
 169 on Computer Vision and Pattern Recognition* (2024), pp. 8281–8291. 2
- 170 [ZJZ*23] ZHANG W., JIN S., ZHUANG P., LIANG Z., LI C.: Under-
 171 water image enhancement via piecewise color correction and dual prior
 172 optimized contrast enhancement. *IEEE Signal Processing Letters* 30
 173 (2023), 229–233. doi:10.1109/LSP.2023.3255005. 2
- 174 [ZLL*25] ZHANG W., LIU Q., LU H., WANG J., LIANG J.: Underwater
 175 image enhancement via wavelet decomposition fusion of advantage con-
 176 trast. *IEEE Transactions on Circuits and Systems for Video Technology*
 177 35, 8 (2025), 7807–7820. doi:10.1109/TCSVT.2025.3545595.
 178 2
- 179 [ZZG*24] ZHANG D., ZHOU J., GUO C., ZHANG W., LI C.: Synergis-
 180 tic multiscale detail refinement via intrinsic supervision for underwater
 181 image enhancement. In *Proceedings of the AAAI conference on artificial
 182 intelligence* (2024), vol. 38, pp. 7033–7041. 2
- 183 [ZZZ*24] ZHANG W., ZHOU L., ZHUANG P., LI G., PAN X., ZHAO
 184 W., LI C.: Underwater image enhancement via weighted wavelet visual
 185 perception fusion. *IEEE Transactions on Circuits and Systems for Video
 186 Technology* 34, 4 (2024), 2469–2483. doi:10.1109/TCSVT.2023.
 187 3299314. 2