

# Hierarchical Consistency Learning for Test-time Adaptation in Camouflage Perception

Mingfeng Zha, Tianyu Li, Guoqing Wang, *Member, IEEE*, Yunqiang Pei, Chaofan Qiao, Jiening Zhang  
Yang Yang, *Senior Member, IEEE*, and Heng Tao Shen, *Fellow, IEEE*

**Abstract**—Camouflaged object detection (COD) aims to localize targets that exhibit minimal perceptual differences from backgrounds through physical attributes. Existing methods, constrained by the static train-then-freeze paradigm, suffer from domain rigidity and annotation dependency, limiting their adaptability to scene variations and unseen camouflage patterns. To overcome these, we propose the hierarchical consistency learning (HCL) framework, which integrates test-time adaptation for dynamic representation recalibration. Specifically, we design the hierarchical representation reconstruction (HRR) to alleviate feature entanglement by synergizing spatial reconstruction with dual-stream frequency-domain decomposition, enhancing robustness against appearance homogenization. The pixel and spectrum inference provide structural and contextual priors. We further introduce task affinity guidance (TAG) to propagate knowledge across branches via channel-wise affinity, aligning local discriminative cues and mitigating semantic drift. To ensure semantic invariance, we formulate the prototype consistency calibration (PCC), which aggregates region features into compact prototypes and establishes prototype-feature similarity. This imposes implicit and hierarchical constraints that bridge task and representation gaps. Extensive experiments across four camouflaged and four underwater object benchmarks, under three degradation settings, demonstrate that our method consistently outperforms state-of-the-art approaches, highlighting its robustness and generalization under distribution shifts.

**Index Terms**—Camouflage perception, Test-time adaptation, Consistency representation.

## I. INTRODUCTION

Existing segmentation and detection frameworks have made notable progress in perceiving salient entities, yet they struggle with concealed targets. When objects with highly similar appearances blend into surrounding contexts without clear boundaries, they may critically mislead downstream analysis, *e.g.*, medical image diagnosis [1]. Camouflaged object detection (COD) focuses on identifying and localizing such

This work was supported in part by the National Natural Science Foundation of China under grant U23B2011, 62102069, U20B2063 and 62220106008, the Key R&D Program of Zhejiang under grant 2024SSYS0091, the Sichuan Science and Technology Program under Grant 2024NSFTD0034.

Mingfeng Zha, Tianyu Li, Guoqing Wang, Yunqiang Pei, Chaofan Qiao, Jiening Zhang and Yang Yang are with the Center for Future Media and School of Computer Science and Engineering, University of Electronic Science and Technology of China, Chengdu 611731, China. (Email: zhamf1116@gmail.com)

Heng Tao Shen is with the School of Computer Science and Technology, Tongji University, Shanghai 201804, China, with the Center for Future Media and School of Computer Science and Engineering, University of Electronic Science and Technology of China, Chengdu 611731, China, and with the Peng Cheng Laboratory, Shenzhen 518066, China.

Corresponding author: Tianyu Li. Email: cosmos.yu@hotmail.com

Project page: <https://winter-flow.github.io/project/HCL>

Manuscript received June 24, 2025; revised June 24, 2025.

hidden targets, providing critical visual cues for robust scene understanding and reliable decision-making.

COD presents fundamental challenges due to the strong visual entanglement between targets and backgrounds, *e.g.*, texture mimicry and chromatic assimilation. In Figure 1 (a), traditional approaches mainly adhere to a *train-then-freeze* paradigm, where models are trained on offline datasets and deployed with fixed parameters. However, this paradigm exhibits significant limitations in open-world or degenerate scenarios (Figure 1 (b)): 1) **Domain Rigidity**: Pre-trained models lack the flexibility to adapt to diverse test-time scenes (*e.g.*, variations in illumination), resulting in feature misalignment and poor pixel-level discrimination. 2) **Annotation Dependency**: Heavy reliance on limited labeled training data restricts generalization, particularly in long-tail or evolving camouflage scenarios, where exhaustive annotation is infeasible. Test-time adaptation (TTA) [2], [3] offers a *train-then-adapt* paradigm that actively leverages intrinsic signals during inference. Rather than treating testing as a passive process, TTA enables models to dynamically recalibrate feature sensitivity, amplifying subtle yet crucial cues for camouflaged object perception. This shift allows models to self-optimize in the presence of distribution shifts. Motivated by this, we propose the hierarchical consistency learning (HCL) framework to formulate a task-specific TTA strategy. This raises three crucial questions: 1) *Why introduce representation reconstruction and hierarchical modeling?* 2) *Why establish knowledge sharing?* 3) *Why enforce prototype consistency?*

We answer the first question. In COD, visual ambiguity arises from intrinsic multi-scale feature entanglement, where local texture similarities and global attribute alignment blur the distinction between targets and backgrounds. We aim to provide self-supervised signals and facilitate regional generalization and association through reconstruction. Existing approaches often rely on single-modality reconstruction in the spatial domain (*i.e.*, RGB pixel space), attempting to differentiate targets through pixel-level autoencoding. However, this strategy suffers from three critical limitations: 1) Highly sensitive to local perturbations (*e.g.*, lighting variations), making it prone to overfitting on noise; 2) Decoupling low-level details and high-level semantics; 3) Lacking robustness against occlusions and deformations that break spatial continuity and hinder global semantic reasoning. Therefore, we propose a hierarchical feature disentanglement strategy through multi-spectral joint reconstruction. Low-frequency reconstruction operates on global image attributes by decomposing signals in the spectrum space, thereby suppressing high-frequency

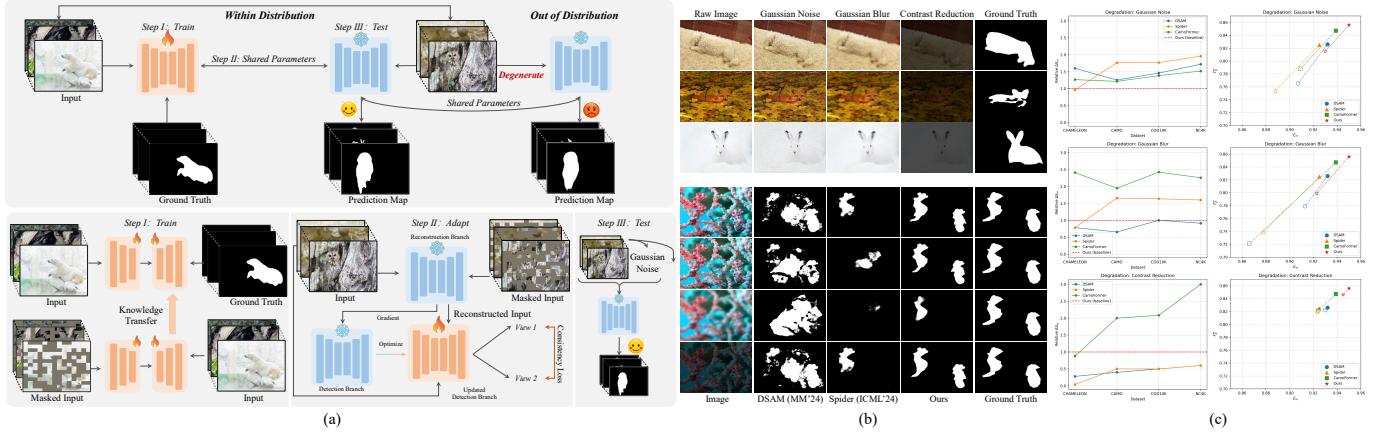


Fig. 1. (a) Traditional methods follow the *train-fix-test* paradigm (top), while our method adopts the *train-adapt-test* paradigm to dynamically perceive the scene (bottom). (b) Three testing sample degradation strategies (top) and qualitative comparisons (bottom). (c) Left: Performance drop  $\Delta S_m$  of comparison methods, relative to our method. Right: Performance variation ( $F_\beta^w$  and  $E_m$  on the NC4K) before (solid markers) and after degradation (hollow markers). Best view by zooming in.

noise and modeling coarse target-background discrepancies. In contrast, high-frequency reconstruction sharpens local details by enhancing edge contours and subtle texture residuals, compensating for the spatial ambiguity. These two branches are jointly optimized alongside spatial reconstruction under complementary constraints, enabling the model to disentangle complex visual patterns across scales.

We answer the second question. While MAE [4] implicitly models global structural priors through masked region reconstruction, which can reveal latent inconsistency, its pixel-level optimization objective significantly diverges from the high-level semantic discriminability required for detection tasks. When relying solely on reconstruction loss to adjust model parameters during test time, the reconstruction branch may overfit to local restoration, while the main detection branch, lacking direct supervision, gradually drifts away from the discriminative representations. This leads to a decoupling between self-supervised signals and detection objectives. To bridge this gap, we introduce affinity learning, which enforces consistency across the feature spaces of the reconstruction and detection branches. Rather than treating the two as independent, we transform the structural-semantic correlations learned during joint training into dynamic constraints for TTA. Concretely, the detection-relevant cues implicitly encoded in the reconstruction features serve as semantic anchors that guide parameter updates. When encountering novel degradation, structural anomalies captured by the reconstruction branch, e.g., edge discontinuity or illumination inconsistency, are leveraged through an affinity matrix to correct the local feature responses of the detection branch.

We answer the third question. The discriminative core of camouflaged objects often lies in global semantic relationships rather than isolated local feature differences. Guiding the model to ignore superficial perturbations and focus on high-level semantic invariance is crucial. Although reconstruction enforces the model to learn structured representations of targets, the process may introduce local noise or semantic deviations. To address this, we introduce prototype consistency,

which aggregates features from regions into compact prototype vectors. These prototypes, extracted from both reconstructed and original images, are fused and aligned with features in the embedding space. This alignment encourages the model to suppress local disturbances and focus on consistent semantic cues shared across samples. By shifting the optimization objective from pixel-level detail recovery to robust semantic space representation, the model becomes more robust to low-level variations and better at identifying true semantic shifts. Furthermore, prototype-feature similarity dynamically guides decoding, balancing local reconstruction with global semantic coherence, thereby enhancing the model's generalization to unseen camouflage patterns.

Technically, we introduce the HRR to enforce pixel-level and multi-frequency reconstruction while imposing domain consistency constraints. To further enhance knowledge transfer, we propose the TAG, which constructs channel-wise affinity maps to propagate local knowledge from the reconstruction branch to the main detection branch. To prevent the model from collapsing into trivial details, we formulate the PCC, which feeds the reconstructed image into the detection model to generate predictions. We employ entropy estimation and edge-guided mechanisms to produce confidence maps, which emphasize regions of high uncertainty and object boundaries. We then perform variational fusion of prototypes from the original and reconstructed images, followed by metric consistency computation to ensure robust and reliable predictions. In Figure 1 (c), our method yields promising results on the original benchmark and under diverse conditions.

In summary, our main contributions are as follows:

- We revisit existing frameworks for COD task and propose a sample-specific TTA strategy to handle data distribution shifts. This strategy requires no additional data and can be seamlessly integrated into other methods.
- We introduce three customized components: the HRR, which enables automatic adaptation to diverse scenes and imaging conditions; the TAG, which serves as a constraint to guide attention toward structurally inconsistent regions

- and refine features; and the PCC, which facilitates selective learning of decoupled and compact representations.
- Extensive experiments on eight benchmarks and three distribution shift settings validate the superiority of the proposed method and the effect of components.

## II. RELATED WORK

### A. Salience Perception

In contrast to camouflage perception, salience perception aims to identify visually distinctive content. Based on feature selection strategies, existing methods can be broadly categorized into handcrafted and learning-based approaches. The former leverages expert prior knowledge such as gradient or geometric cues, while the latter adopts data-driven paradigms to generate high-dimensional latent representations with greater robustness. Learning-based methods evolved beyond conventional 2D natural scenes to diverse downstream scenarios, including underwater, mirror, and panoramic environments [5]–[7], and have progressed from single-modality frameworks to multi-modal collaboration (*e.g.*, depth and thermal data). Furthermore, some works focus on modeling relative saliency differences among targets, *i.e.*, saliency ranking [8]. Since fully supervised learning generally requires large-scale, multi-scene datasets to obtain generalizable representations, recent efforts explored data-efficient alternatives such as self-supervised learning [9]. In this work, we propose to discover discriminative regions during inference by leveraging the image itself as a supervisory signal.

### B. Camouflage Perception

The development of COD has progressed rapidly. Based on supervised learning signals, we categorize it into four types: 1) Fully supervised, which utilizes edge [10], texture [11], frequency [12], depth [13], uncertainty [14], or text [15] guidance, or employs a coarse-to-fine progressive approach to mine potential clues [16], [17]; 2) Weakly supervised [18], which leverages manually annotated scribbles or points as ground truth and iteratively optimizes the labels; 3) Semi-supervised [19], which explores the relationships and consistencies between labeled and unlabeled data; 4) Zero/Few-shot learning [20], [21], which aims to transfer knowledge to unseen scenarios using only a limited number of samples. Additionally, some works [22]–[25] leverage features and semantic priors from pre-trained or large language/vision/multi-modal foundation models, introducing efficient fine-tuning mechanisms that achieve competitive performance with minimal learnable parameters. Furthermore, some studies investigate at the instance level [26], video level [27], and across different scenarios [28], [29]. Our work is most closely related to [30], but it overlooks scene adaptation during the testing phase by constructing a multi-task framework for pixel reconstruction and detection that is susceptible to noise interference. Pang *et al.* [31] proposed an open-world setting based on CLIP [32], utilizing text prompts. In contrast, our HCL does not rely on additional data or foundation models to facilitate the remapping of the model to address distribution shifts.

### C. Test-time Adaptation

TTA adapts model parameters during the testing phase to dynamically align with new data distributions, facilitating the transition from *i.i.d.* modeling to out-of-distribution generalization. Sun *et al.* [2] proposed adapting models via online self-supervised learning, enabling real-time updates. Wang *et al.* [33] further improved efficiency by dynamically adjusting parameters via entropy minimization. To tackle error accumulation, Niu *et al.* [3] introduced anti-forgetting mechanisms with sample selection and regularization. Jang *et al.* [34] enhanced pseudo-label reliability by leveraging prototype matching. For lifelong adaptation, Brahma *et al.* [35] designed a probabilistic framework to balance new knowledge integration with prior preservation. Zhao *et al.* [36] tackled complex distribution shifts by combining statistical correction and sample reweighting. Ma *et al.* [37] boosted robustness through graph-structured label refinement and model averaging. Recently, some works introduced TTA into specific tasks, *e.g.*, video segmentation [38]. We introduce TTA into the COD and propose the HCL framework, which incorporates tailored components to mitigate the train-test gap, particularly under significant distribution shifts. Unlike [39], which relies on multiple foundation models and adapts only at the test stage through prompting, the HCL framework integrates seamlessly into the entire pipeline without requiring manual design (*e.g.*, prompts) or prior external knowledge.

## III. PROPOSED METHOD

### A. Overall Architecture

In Figure 2 and Algorithm 1, given an input  $\mathbf{I} \in \mathbb{R}^{3 \times H \times W}$ , we generate a random mask  $\mathbf{M}_{\text{spa}} \in \{0, 1\}^{H \times W}$  and extract multi-scale features  $\mathbf{X}^i \in \mathbb{R}^{C^i \times H^i \times W^i}$  using the encoder ( $C$ : Channel;  $H$ : Height;  $W$ : Width). We apply Fourier transform  $\mathbb{F}$  to project the input into the frequency domain, where mask  $\mathbf{M}_l$  (and  $\mathbf{M}_h$ ) and inverse transformation  $\mathbb{F}^{-1}$  are used to obtain the masked low-frequency component  $\hat{\mathbf{I}}_l$  and the masked high-frequency component  $\hat{\mathbf{I}}_h$ . The subsequent operations mirror those in the spatial domain. We alternately use  $\mathcal{L}_{\text{pix}}$  and  $\mathcal{L}_{\text{freq}}$  as reconstruction and constraint losses. The masked feature  $\mathbf{X}_{\text{rec}}^i$  transfers heterogeneous information to the unmasked feature  $\mathbf{X}^i$  via the TAG as a bridge. The reconstructed image  $\mathbf{I}_{\text{rec}}$  is then fed as an additional input to generate the prediction map  $\mathbf{O}_{\text{rec}}$  and confidence map  $\Phi$ , which are further used to derive the prototype  $\mathbf{P}_{\text{rec}}$ . We integrate  $\mathbf{P}_{\text{rec}}$  with the prototype  $\mathbf{P}$  from the original input, sequentially using fused prototype  $\mathbf{P}_{\text{fusion}}$  as an anchor to establish consistency across predictions and features.

### B. Preliminary

According to Parseval's Theorem, a signal  $x$  retains the same total energy in both the spatial and frequency domains, though with different emphasis. Based on the 2D Discrete Fourier Transform (DFT)  $\mathbb{F}$ , we transform  $x \in \mathbb{R}^{H \times W}$  into the Fourier spectrum,

$$\mathcal{F}(x)(u, v) = \sum_{h=0}^{H-1} \sum_{w=0}^{W-1} x(h, w) e^{-j2\pi(\frac{h}{H}u + \frac{w}{W}v)} \quad (1)$$

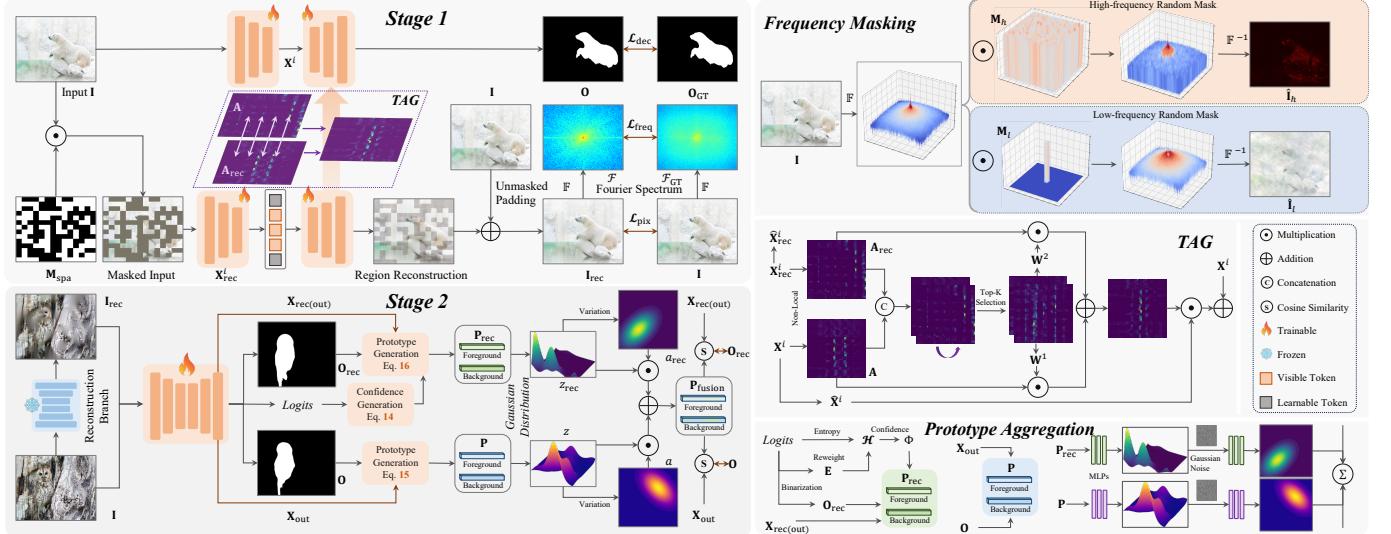


Fig. 2. Stage 1: The input image undergoes masked patch partitioning in the spatial domain, while threshold-based masking is applied to obtain high- and low-frequency components. Unmasked patches are encoded to extract visible tokens, and along with learnable tokens, are decoded in spatial and frequency branches for pixel- and spectrum-level reconstruction. The TAG enables cross-branch knowledge transfer, serving as an explicit constraint. Stage 2: The spatially reconstructed image is passed through the detection network to generate predictions. Uncertainty estimation and edge information are leveraged to construct the confidence map that refines pixel-wise weighting. Finally, prototype integration enforces cross-view consistency, forming a closed-loop from reconstruction to prediction.

where  $u$  and  $v$  represent the horizontal and vertical indices in the spectrum, respectively. To transform back to the spatial domain, we apply the inverse Fourier transform  $\mathbb{F}^{-1}$ ,

$$x(h, w) = \frac{1}{HW} \sum_{u=0}^{H-1} \sum_{v=0}^{W-1} \mathcal{F}(u, v) e^{j2\pi(\frac{uh}{H} + \frac{vw}{W})} \quad (2)$$

In practice, we use the Fast Fourier Transform (FFT) version to improve computational efficiency.

### C. Hierarchical Representation Reconstruction

We decompose the auxiliary reconstruction task into: spatial, high-frequency, and low-frequency, each accompanied by corresponding transformation constraints.

**Spatial Reconstruction.** To supervise the reconstruction process, we adopt the Mean Square Error (MSE) loss  $\mathcal{L}_{pix}$  on the output  $\mathbf{I}_{rec}$  from the pixel decoder, encouraging accurate recovery of the spatial regions masked by  $\mathbf{M}_{spa}$ ,

$$\mathcal{L}_{pix} = \frac{1}{HWC} \sum_{i=0}^{H-1} \sum_{j=0}^{W-1} \sum_{c=0}^{C-1} (\mathbf{I}_{rec}^{i,j,c} - \mathbf{I}^{i,j,c})^2 \quad (3)$$

**Frequency Reconstruction.** To avoid inherent reconstruction defects in the pixel domain, we employ Fourier transform to convert the raw image into the frequency domain, partitioning the spectrum into high- and low- frequency part based on the predefined threshold (the distance of spectral points from the center). For low-frequency reconstruction, we use a central mask  $\mathbf{M}_l \in \{0, 1\}^{H_l \times H_w}$  to randomly suppress low-frequency energy while retaining high-frequency,

$$\mathbf{M}_l(u, v) = \begin{cases} 1, & \text{unmasked} \\ 0, & \text{masked} \end{cases} \quad (4)$$

Therefore, we can obtain the degraded spectrum  $\hat{\mathcal{F}}$ ,

$$\hat{\mathcal{F}} = \sum_{c=0}^{C-1} \mathcal{F}_{low}^c \odot \mathbf{M}_l^c + \mathcal{F}_{high}^c \quad (5)$$

where  $\odot$  denotes the Hadamard product. We further utilize the inverse Fourier transform to generate the degraded image  $\hat{\mathbf{I}}_l$ . For high-frequency reconstruction, we perform similar operations and define the random mask  $\mathbf{M}_h$ , obtaining  $\hat{\mathbf{I}}_h$ . After generating features using the encoder, we reconstruct the high-frequency and low-frequency images from tokens, i.e.,  $\mathbf{I}_{rec(h)}$  and  $\mathbf{I}_{rec(l)}$ .

In spatial visual recognition, different regions contribute unequally to decision-making. Likewise, in frequency decoding, hard samples demand greater learning attention. To address this, we introduce focal frequency loss [40] to capture high-value clues and emphasize challenging frequency components,

$$\mathcal{L}_{freq} = \frac{1}{HW} \sum_{u=0}^{H-1} \sum_{v=0}^{W-1} \omega(u, v) \odot \gamma(\mathcal{F}(u, v), \mathcal{F}_{GT}(u, v))^2 \quad (6)$$

where  $\mathcal{F}(u, v)$  and  $\mathcal{F}_{GT}(u, v)$  represent the components of the reconstructed spectrum  $\mathbf{L}_{rec}$  (or  $\mathbf{H}_{rec}$ ) and the raw spectrum  $\mathbf{L}$  (or  $\mathbf{H}$ ) at position  $(u, v)$ , respectively.  $\omega$  and  $\gamma$  denote the weight parameter and distance metric, which can be formulated as,

$$\begin{aligned} \omega(u, v) &= \gamma(\mathcal{F}(u, v), \mathcal{F}_{GT}(u, v))^\beta \\ \gamma(\mathcal{F}, \mathcal{F}_{GT}) &= \sqrt{(\mathcal{R} - \mathcal{R}_{GT})^2 + (\mathcal{I} - \mathcal{I}_{GT})^2} \end{aligned} \quad (7)$$

where  $\mathcal{R}$  and  $\mathcal{I}$  denote the real and imaginary parts of the spectrum, respectively.  $\beta$  is a scaling factor, which is set to 1 by default.

**Transformation Consistency.** In theory, when  $\mathcal{L}_{pix} \rightarrow 0$  or  $\mathcal{L}_{freq} \rightarrow 0$ , the model can accurately predict the masked

regions, resulting in the reconstructed image that matches the raw image. However,  $\mathcal{L}_{\text{pix}}$  directly enforces pixel-level matching, ensuring the recovery of local details but is insensitive to global structures, *e.g.*, edge continuity. Conversely,  $\mathcal{L}_{\text{freq}}$  enhances the perception of global structures but is less sensitive to local pixel shifts. In other words, a sole focus on pixel matching may lead to blurriness, while an emphasis on spectral matching may introduce artifacts. Therefore, for each reconstruction branch, we apply the corresponding domain loss as a regularization term to achieve transformation consistency. We can formulate the total loss  $\mathcal{L}_{\text{HRR}}$  for the HRR,

$$\begin{aligned} \mathcal{L}_{\text{HRR}} = & \mathcal{L}_{\text{pix}(\text{rec})}(\mathbf{I}_{\text{rec}}, \mathbf{I}) + \mathcal{L}_{\text{freq}(\text{con})}(\mathcal{F}, \mathcal{F}_{\text{GT}}) + \\ & \sum_{k \in \{\mathbf{L}, \mathbf{H}\}} \lambda_k (\mathcal{L}_{\text{freq}(\text{rec})}(k_{\text{rec}}, k) + \mathcal{L}_{\text{pix}(\text{con})}(\mathbb{F}^{-1}(k), \mathbf{I})) \end{aligned} \quad (8)$$

where  $\mathcal{L}_{\sim(\text{rec})}$  and  $\mathcal{L}_{\sim(\text{con})}$  represent the reconstruction and consistency (constraint) losses, respectively.  $\lambda$  is the balancing hyperparameter. We empirically set  $\lambda_{\mathbf{L}}$  and  $\lambda_{\mathbf{H}}$  to 0.4 and 0.6.

#### D. Task Affinity Guidance

Despite employing reconstruction as auxiliary supervision and implicitly updating the detection model parameters via gradient optimization, our approach may require more training iterations and could lead to suboptimal performance. To address this limitation, we explicitly establish feature associations. Unlike vanilla cross-attention mechanisms with high computational complexity ( $\mathcal{O}(H^2W^2)$ ), we construct a more efficient affinity matrix to reduce computational overhead and suppress noise. Specifically, we adopt channel-wise non-local associations with reduced complexity ( $\mathcal{O}(C^2)$ ) to capture long-range dependencies and enhance informative representations, ultimately generating the affinity map  $\mathbf{A}$  and the optimized feature  $\hat{\mathbf{X}}$ ,

$$[\hat{\mathbf{X}}, \mathbf{A}] = \text{NonLocal}(\mathbf{X}) \quad (9)$$

For  $\mathbf{X}_{\text{rec}}$ , the operation follows the same procedure, producing  $\hat{\mathbf{X}}_{\text{rec}}$  and  $\mathbf{A}_{\text{rec}}$ . Due to representation differences, we fuse  $\mathbf{A}$  and  $\mathbf{A}_{\text{rec}}$  to obtain the weight map  $\mathbf{W}$ ,

$$[\mathbf{W}^1, \mathbf{W}^2] = \text{Top-K}(\text{Self-Attention}(\text{Concat}(\mathbf{A}, \mathbf{A}_{\text{rec}}))) \quad (10)$$

where  $\text{Concat}$  denotes channel concatenation, while  $\text{Top-K}$  selects the top-k weights to filter out low-response elements. In Figure 6, we select the top 70% (or 80%) of elements to achieve a high signal-to-noise ratio. Thus, we can obtain the updated map that incorporates the reconstruction knowledge,

$$\mathbf{A} := \mathbf{W}^1 \odot \mathbf{A} + \mathbf{W}^2 \odot \mathbf{A}_{\text{rec}} \quad (11)$$

Correspondingly, the updated feature is,

$$\hat{\mathbf{X}} := \mathbf{X} + \mathbf{A} \odot \hat{\mathbf{X}} \quad (12)$$

#### E. Prototype Consistency Calibration

**Confidence Map Generation.** In practice, the representations of the reconstruction and the original perspective cannot be fully identical. Inevitably, noise disturbances and differentiated

modeling can lead to unreliability in certain regions of the predicted map  $\mathbf{O}_{\text{rec}}$ . To quantify the degree of reliability, we introduce uncertainty estimation. Unlike Monte Carlo strategy that requires multiple samples, we measure based on entropy  $\mathcal{H}$ , which demands only a single forward pass. For camouflaged targets, the boundary information is more critical than the main body. Therefore, we generate the edge map  $\mathbf{E}$  and utilize it as a spatial weighting factor. For any position  $(i, j)$  of output-layer feature  $\mathbf{X}_{\text{rec}(\text{out})} \in \mathbb{R}^{1 \times H \times W}$ , we apply sigmoid to obtain the probability  $\mathbf{p}$  for class  $m$ ,

$$\begin{aligned} \mathcal{H}(\mathbf{p}^{i,j}) = & - \sum_{m=0}^{M-1} w(i, j) \mathbf{p}^{i,j}(m) \log \mathbf{p}^{i,j}(m) \\ w(i, j) = & 1 + \alpha \cdot \mathbf{E}(i, j) \end{aligned} \quad (13)$$

where  $\alpha$  is a scaling factor to balance pixel-wise weights between edge and non-edge regions, and  $\mathbf{E}(i, j) \in \{0, 1\}$  denotes the edge indicator function:  $\mathbf{E}(i, j) = 1$  if pixel  $(i, j)$  lies on generated edge, and 0 otherwise (for COD task,  $M = 2$ ). In information theory, entropy serves as a quantitative measure of uncertainty, where higher entropy values correspond to increased information disorder and probabilistic ambiguity in the prediction confidence. Based on this, we further define the confidence map  $\Phi$ ,

$$\Phi^{i,j} = \frac{1}{H \times W} \left( 1 - \frac{\mathcal{H}(\mathbf{p}^{i,j})}{\sum_{i=0}^{H-1} \sum_{j=0}^{W-1} \mathcal{H}(\mathbf{p}^{i,j})} \right) \quad (14)$$

**Prototype Generation.** The prototype characterizes the overall representation of a region. We generate the raw prototype via masked average pooling (MAP),

$$\mathbf{P}^m = \frac{\sum_{i=0}^{H-1} \sum_{j=0}^{W-1} \mathbf{X}_{\text{out}}^{i,j} \mathbf{1}[\mathbf{O}^{i,j} = m]}{\sum_{i=0}^{H-1} \sum_{j=0}^{W-1} \mathbf{1}[\mathbf{O}^{i,j} = m]} \quad (15)$$

where  $\mathbf{1}$  is an indicator function. MAP treats features from all regions equally. However, the reconstructed features may contain errors (*e.g.*, blurriness and noise), particularly in incomplete or complex areas. Directly using MAP may lead to the smoothing of errors or even obscure the correct representations. We utilize  $\Phi$  to automatically filter reliable samples, achieving self-calibration,

$$\mathbf{P}_{\text{rec}}^m = \frac{\sum_{i=0}^{H-1} \sum_{j=0}^{W-1} \mathbf{X}_{\text{rec}(\text{out})}^{i,j} \Phi^{i,j} \mathbf{1}[\mathbf{O}_{\text{rec}}^{i,j} = m]}{\sum_{i=0}^{H-1} \sum_{j=0}^{W-1} \mathbf{1}[\mathbf{O}_{\text{rec}}^{i,j} = m]} \quad (16)$$

**Variational Prototype Fusion.** While naive fusion strategies like element-wise summation or channel-wise concatenation between prototypes  $\mathbf{P}^m$  and  $\mathbf{P}_{\text{rec}}^m$  appear straightforward, both suffer from critical limitations: 1) Direct summation induces linear amplification of noise artifacts: if  $\mathbf{P}^m$  contains misclassified background pixels due to blurred object boundaries and  $\mathbf{P}_{\text{rec}}^m$  exhibits restoration distortions (*e.g.*, occlusion), error patterns become compounded via accumulation; 2) Concatenation risks feature space misalignment, as  $\mathbf{P}^m$  primarily encodes fine-grained representations while  $\mathbf{P}_{\text{rec}}^m$  captures structural knowledge, creating distribution shifts that hinder semantic consistency when forcibly combined. We utilize variational fusion to alleviate through probabilistic modeling.

To further explore the intrinsic uncertainty, we project the samples into the latent Gaussian distribution  $\mathcal{N}$  to establish semantic associations. For  $\mathbf{P}^m$ , we have,

$$\begin{aligned} p(z|\mathbf{P}^m) &= \mathcal{N}(z; \mu, (\sigma)^2 \mathbf{I}) \\ \mu &= \mathcal{M}_\mu(\mathbf{P}^m), \quad \sigma = \mathcal{M}_\sigma(\mathbf{P}^m) \end{aligned} \quad (17)$$

where  $z$ ,  $\mu$ ,  $\sigma$ , and  $\mathbf{I}$  represent the reconstruction vector, mean, variance, and identity matrix, respectively. A larger  $\sigma$  indicates greater uncertainty.  $\mathcal{M}$  denotes the multilayer perceptron (MLPs). Similarly, for  $\mathbf{P}_{\text{rec}}^m$ , we can obtain  $z_{\text{rec}}$ .

Sampling latent variables from a distribution to generate attention weights results in non-differentiable gradients. To resolve this, we employ reparameterization trick, which decouples the stochasticity by introducing external noise, allowing gradient flow through a deterministic computational path. Specifically, we introduce standard Gaussian noise  $\eta$ ,

$$z = \mu + \eta \cdot \sigma, \quad \eta \in \mathcal{N}(0, 1) \quad (18)$$

Similar operation is performed for  $z_{\text{rec}}$ . Traditional attention mechanisms (like Softmax) generate fixed weights  $a^i$  through point estimates, which are essentially deterministic mappings and cannot reflect the model's confidence in the weight values,

$$\begin{aligned} \hat{a}^i &= \mathbf{W}^i z^i + \mathbf{b}^i, \quad i \in \{o, r\} \\ a^i &= \frac{\exp(\hat{a}^i)}{\sum_{j \in \{o, r\}} \exp(\hat{a}^j)} \end{aligned} \quad (19)$$

where  $o$  and  $r$  denote the original and reconstructed spatial images, respectively.  $\mathbf{W}$  and  $\mathbf{b}$  are learnable parameters used for projection. Instead of using deterministic weights, we leverage dual uncertainty modeling by treating the weights as probability distributions. 1) *Feature-level uncertainty*: The variances  $\sigma$  and  $\sigma_{\text{rec}}$  of the prototype features represent local confidence. 2) *Weight-level uncertainty*: The variance  $\sigma_a$  of the attention weights captures global confidence.

For each latent variable  $z^i$ , we employ shared MLPs to estimate the mean  $\mu_a$  and variance  $\sigma_a$ . Accordingly, we reformulate Eq. 19 as follows,

$$\begin{aligned} \hat{a}^i &\sim q(\hat{a}^i | z^i) = \mathcal{N}(\mu_a^i, (\sigma_a^i)^2 \mathbf{I}), \quad i \in \{o, r\} \\ a^i &= \frac{\exp(-\gamma(\sigma_a^i)^2)}{\sum_{j \in \{o, r\}} \exp(-\gamma(\sigma_a^j)^2)} \end{aligned} \quad (20)$$

where  $\gamma$  is a learnable temperature coefficient that controls the intensity of variance suppression on the weights. We can obtain the prototype  $\mathbf{P}_{\text{fusion}}^m$  after weighted fusion,

$$\mathbf{P}_{\text{fusion}}^m = \sum_{i \in \{o, r\}} a^i z^i \quad (21)$$

To prevent overfitting to noise, we enforce that the prototype distribution approaches the standard Gaussian prior through Kullback–Leibler (KL) divergence,

$$\mathcal{L}_{\text{KL}} = \sum_{i \in \{o, r\}} \frac{1}{2} ((\mu^i)^2 + (\sigma^i)^2 - \log((\sigma^i)^2) - 1) \quad (22)$$

Based on variational fusion, we reduce the weights in high-variance (low confidence) regions to suppress noise propagation. We dynamically adjust the fusion strategy based on input scenarios to achieve adaptive feature complementarity.

**Metric Consistency.** We use the similarity matrix  $\mathbf{S}$  between feature and prototype to estimate the class probability for each pixel. Essentially, this strategy is an implicit form of contrastive learning,

$$\mathbf{S}^{i,j} = \text{CosSim}(\mathbf{X}^{i,j}, \mathbf{P}^m) = \frac{\mathbf{X}^{i,j} \cdot \mathbf{P}^m}{\|\mathbf{X}^{i,j}\|_2 \cdot \|\mathbf{P}^m\|_2 + \epsilon} \quad (23)$$

where  $\text{CosSim}(\cdot)$  and  $\epsilon$  represent the cosine similarity and the minimum value, respectively. Similarly, we can obtain  $\mathbf{S}_{\text{rec}}$  for  $\mathbf{P}_{\text{rec}}$ . When  $\mathbf{P}_{\text{fusion}}$  is sufficiently accurate and robust, it follows that  $\mathbf{S} \rightarrow \mathbf{O}$  (or  $\mathbf{S}_{\text{rec}} \rightarrow \mathbf{O}_{\text{rec}}$ ). Base on cross entropy loss  $\mathcal{L}_{\text{CE}}$ , we have,

$$\mathcal{L}_{\text{pro}} = \mathcal{L}_{\text{CE}}(\mathbf{S}, \mathbf{O}), \quad \mathcal{L}_{\text{pro(rec)}} = \sum_{i=0}^{H-1} \sum_{j=0}^{W-1} \Phi^{i,j} \mathcal{L}_{\text{CE}}(\mathbf{S}_{\text{rec}}^{i,j}, \mathbf{O}_{\text{rec}}^{i,j}) \quad (24)$$

For the main detection branch, following [41], we employ structure loss  $\mathcal{L}_{\text{dec}}$ . Note that we apply supervision to the predictions at all stages of the decoder. We can derive the total loss  $\mathcal{L}_{\text{total}}$ ,

$$\mathcal{L}_{\text{total}} = \sum_{i \in \{\text{HRR, KL, pro, proe(rec), dec}\}} \lambda_i \mathcal{L}_i \quad (25)$$

In the TTA phase, we discard  $\mathcal{L}_{\text{dec}}$ .

---

### Algorithm 1 Our HCL Framework

---

- 1: **Input:** Image  $\mathbf{I}$
  - 2: **Stage 1: Hierarchical Reconstruction**
  - 3:  $\mathbf{M}_{\text{spa}} \leftarrow \text{SpaMask}(\mathbf{I}), [\mathbf{M}_l, \mathbf{M}_h] \leftarrow \text{FreqMask}(\mathbb{F}(\mathbf{I}))$
  - 4:  $\mathbf{X}_{\text{rec}}, \mathbf{X} \leftarrow \text{Enc}(\mathbf{I} \odot \mathbf{M}_{\text{spa}}, \mathbf{I})$   $\triangleright$  Encoding Process
  - 5:  $[\mathbf{X}_{\text{rec(l)}}, \mathbf{X}_{\text{rec(h)}}], \mathbf{X} \leftarrow \text{Enc}(\mathbb{F}^{-1}(\mathbb{F}(\mathbf{I} \odot [\mathbf{M}_l, \mathbf{M}_h])), \mathbf{I})$
  - 6:  $\mathbf{X} \leftarrow \text{TAG}(\mathbf{X}_{\text{rec}}, \mathbf{X}), \mathbf{X} \leftarrow \text{TAG}([\mathbf{X}_{\text{rec(l)}}, \mathbf{X}_{\text{rec(h)}}], \mathbf{X})$
  - 7:  $\mathbf{O}, \mathbf{I}_{\text{rec}} \leftarrow \text{Dec}(\mathbf{X}_{\text{rec}}, \mathbf{X})$   $\triangleright$  Decoding Process
  - 8:  $\mathbf{O}, [\mathbf{I}_{\text{rec(l)}}, \mathbf{I}_{\text{rec(h)}}] \leftarrow \text{Dec}(\mathbf{X}_{\text{rec}}, \mathbf{X})$
  - 9: **Stage 2: Consistency Calibration**
  - 10:  $\mathbf{I}_{\text{rec}} \leftarrow \text{FrozenNet}(\mathbf{I})$   $\triangleright$  Reconstruction
  - 11:  $[\mathbf{X}_{\text{rec(out)}}, \mathbf{O}_{\text{rec}}, \Phi] \leftarrow \text{TrainableNet}(\mathbf{I}_{\text{rec}})$   $\triangleright$  Detection
  - 12:  $[\mathbf{X}, \mathbf{O}] \leftarrow \text{TrainableNet}(\mathbf{I})$
  - 13:  $\mathbf{P}_{\text{rec}}, \mathbf{P} \leftarrow \text{Prototype}([\mathbf{X}_{\text{rec(out)}}, \mathbf{O}_{\text{rec}}, \Phi], [\mathbf{X}, \mathbf{O}])$
  - 14:  $\mathbf{P}_{\text{fusion}} \leftarrow \text{VariationalFusion}(\mathbf{P}_{\text{rec}}, \mathbf{P})$
  - 15:  $\mathbf{S}_{\text{rec}}, \mathbf{S} \leftarrow \text{CosSim}(\mathbf{P}_{\text{fusion}}, [\mathbf{X}_{\text{rec(out)}}, \mathbf{X}_{\text{out}}])$
  - 16: **if NOT INFERENCE then**
  - 17:  $\mathcal{L}_{\text{total}} \leftarrow \mathcal{L}_{\text{HRR}} + \mathcal{L}_{\text{KL}} + \mathcal{L}_{\text{pro}} + \mathcal{L}_{\text{pro(rec)}} + \mathcal{L}_{\text{dec}}$
  - 18: **else**
  - 19:  $\mathcal{L}_{\text{total}} \leftarrow \mathcal{L}_{\text{HRR}} + \mathcal{L}_{\text{KL}} + \mathcal{L}_{\text{pro}} + \mathcal{L}_{\text{pro(rec)}}$
  - 20: **end if**
  - 21: **Output:** Prediction Map  $\mathbf{O}$
- 

## IV. EXPERIMENTS

### A. Datasets

We conduct experiments on four standard COD benchmark datasets. CAMO [59] and COD10K [42] contain 1,000 and 3,040 training pairs, with 250 and 2,026 testing pairs, respectively. CHAMELEON [60] and NC4K [46] are used exclusively for testing, consisting of 76 and 4,121 images. Following [30], [52], we aggregate the training pairs from CAMO

TABLE I

QUANTITATIVE COMPARISON ON NORMAL BENCHMARK DATASETS. FOUNDATION MODELS INCLUDE, *e.g.*, SAM. EXTRA DATA INCLUDE, *e.g.*, DEPTH MAP AND MORE PIXELS. BEST PERFORMANCE IN **BOLD**, SECOND IN UNDERLINE.  $\ddagger$  REPRESENTS DATA IS UNAVAILABLE.  $\uparrow$  INDICATES HIGHER VALUES ARE BETTER, WHILE  $\downarrow$  INDICATES THE OPPOSITE. \*: HIGHER/MULTI-SCALE INPUT RESOLUTION. ‘OURS-R’ AND ‘OURS-P’ DENOTE RESNET50 AND PVT-v2 AS BACKBONE, RESPECTIVELY.

Method	Foundation Model	Extra Data	TTA	CHAMELEON (76)				CAMO (250)				COD10K (2026)				NC4K (4121)			
				$S_m \uparrow$	$E_m \uparrow$	$F_\beta^w \uparrow$	MAE $\downarrow$	$S_m \uparrow$	$E_m \uparrow$	$F_\beta^w \uparrow$	MAE $\downarrow$	$S_m \uparrow$	$E_m \uparrow$	$F_\beta^w \uparrow$	MAE $\downarrow$	$S_m \uparrow$	$E_m \uparrow$	$F_\beta^w \uparrow$	MAE $\downarrow$
PraNet [1] MICCAI’20	$\times$	$\times$	$\times$	0.860	0.898	0.763	0.044	0.769	0.833	0.663	0.094	0.789	0.839	0.629	0.045	0.822	0.876	0.724	0.059
SINet [42] CVPR’20	$\times$	$\times$	$\times$	0.872	0.946	0.806	0.034	0.751	0.771	0.606	0.100	0.771	0.806	0.551	0.051	0.810	0.873	0.772	0.057
MGL [43] CVPR’21	$\times$	$\times$	$\times$	0.893	0.923	0.813	0.030	0.775	0.847	0.673	0.088	0.814	0.865	0.666	0.035	$\ddagger$	$\ddagger$	$\ddagger$	$\ddagger$
PFNet [44] CVPR’21	$\times$	$\times$	$\times$	0.882	0.931	0.810	0.033	0.782	0.852	0.695	0.085	0.800	0.868	0.660	0.040	0.829	0.887	0.745	0.053
UGTR [45] ICCV’21	$\times$	$\times$	$\times$	0.888	0.940	0.794	0.031	0.785	0.859	0.686	0.086	0.818	0.850	0.667	0.035	$\ddagger$	$\ddagger$	$\ddagger$	$\ddagger$
LSR [46] CVPR’21	$\times$	$\times$	$\times$	0.893	0.938	0.839	0.033	0.793	0.826	0.725	0.085	0.793	0.868	0.685	0.041	0.839	0.883	0.779	0.053
SINet-v2 [47] TPAMI’21	$\times$	$\times$	$\times$	0.888	0.942	0.816	0.030	0.820	0.882	0.743	0.070	0.815	0.887	0.680	0.037	0.847	0.903	0.769	0.048
PreyNet * [48] MM’22	$\times$	$\checkmark$	$\times$	0.895	0.951	0.844	0.028	0.790	0.842	0.708	0.077	0.813	0.891	0.697	0.034	$\ddagger$	$\ddagger$	$\ddagger$	$\ddagger$
ZoomNet * [41] CVPR’22	$\times$	$\checkmark$	$\times$	0.902	0.958	<u>0.845</u>	<u>0.023</u>	<u>0.820</u>	<u>0.892</u>	<u>0.752</u>	<b>0.066</b>	0.838	0.888	0.729	<u>0.029</u>	0.853	0.896	0.784	0.043
PopNet [13] ICCV’23	$\times$	$\checkmark$	$\times$	<b>0.917</b>	<b>0.965</b>	<b>0.875</b>	<b>0.020</b>	0.808	0.859	0.744	0.077	<b>0.851</b>	<b>0.910</b>	<b>0.757</b>	<b>0.028</b>	0.861	0.909	<b>0.802</b>	0.042
FEDER [49] CVPR’23	$\times$	$\times$	$\times$	0.907	0.964	$\ddagger$	0.025	0.807	0.873	$\ddagger$	0.069	0.823	0.900	$\ddagger$	0.032	0.846	0.905	$\ddagger$	0.045
Ours-R	$\times$	$\times$	$\checkmark$	0.880	0.949	0.840	0.025	<b>0.827</b>	<b>0.898</b>	<b>0.765</b>	<u>0.067</u>	0.842	0.898	0.744	0.030	<b>0.865</b>	<b>0.915</b>	0.795	<b>0.042</b>
SAM [50] ICCV’23	$\checkmark$	$\times$	$\times$	$\ddagger$	$\ddagger$	$\ddagger$	$\ddagger$	$\ddagger$	$\ddagger$	$\ddagger$	$\ddagger$	0.778	0.800	0.701	0.050	0.765	0.778	0.696	0.078
HitNet * [51] AAAI’23	$\times$	$\checkmark$	$\times$	<b>0.922</b>	<b>0.970</b>	<b>0.903</b>	<b>0.018</b>	0.844	0.902	0.801	0.057	<b>0.868</b>	0.932	<b>0.798</b>	0.024	0.870	0.921	0.825	0.039
FPNet [52] MM’23	$\times$	$\times$	$\times$	0.914	0.960	<u>0.868</u>	<u>0.022</u>	0.851	0.912	0.802	0.056	0.851	0.909	0.755	0.028	$\ddagger$	$\ddagger$	$\ddagger$	$\ddagger$
FSPNet [53] CVPR’23	$\times$	$\times$	$\times$	0.908	0.943	0.851	0.023	0.856	0.899	0.799	0.050	0.851	0.895	0.735	0.026	0.879	0.915	0.816	0.035
EVP [54] CVPR’23	$\times$	$\times$	$\times$	0.871	0.917	0.795	0.036	0.846	0.895	0.777	0.067	0.843	0.907	0.742	0.032	0.874	$\ddagger$	$\ddagger$	$\ddagger$
CamoFormer [55] TPAMI’24	$\times$	$\times$	$\times$	0.910	0.957	0.866	0.022	<u>0.872</u>	<u>0.929</u>	<u>0.831</u>	<u>0.046</u>	<b>0.869</b>	0.932	<u>0.786</u>	<u>0.023</u>	<b>0.892</b>	<b>0.939</b>	<b>0.847</b>	0.030
VSCode [24] CVPR’24	$\times$	$\checkmark$	$\times$	$\ddagger$	$\ddagger$	$\ddagger$	$\ddagger$	$\ddagger$	$\ddagger$	$\ddagger$	$\ddagger$	0.836	0.892	0.768	0.060	0.847	0.913	0.744	0.028
Spider [56] ICML’24	$\times$	$\times$	$\times$	0.906	0.951	0.848	0.025	0.855	0.908	0.799	0.053	0.856	0.917	0.756	0.028	0.880	0.925	0.825	0.036
DSAM [57] MM’24	$\checkmark$	$\checkmark$	$\times$	0.853	0.924	0.785	0.045	0.832	0.913	0.794	0.061	0.846	0.921	0.760	0.033	0.871	0.932	0.826	0.040
Ours-P	$\times$	$\times$	$\checkmark$	0.893	<u>0.962</u>	0.853	0.023	<b>0.873</b>	<b>0.940</b>	<b>0.840</b>	<b>0.040</b>	0.860	<b>0.933</b>	0.781	<b>0.022</b>	<b>0.891</b>	<b>0.950</b>	<b>0.856</b>	<b>0.026</b>

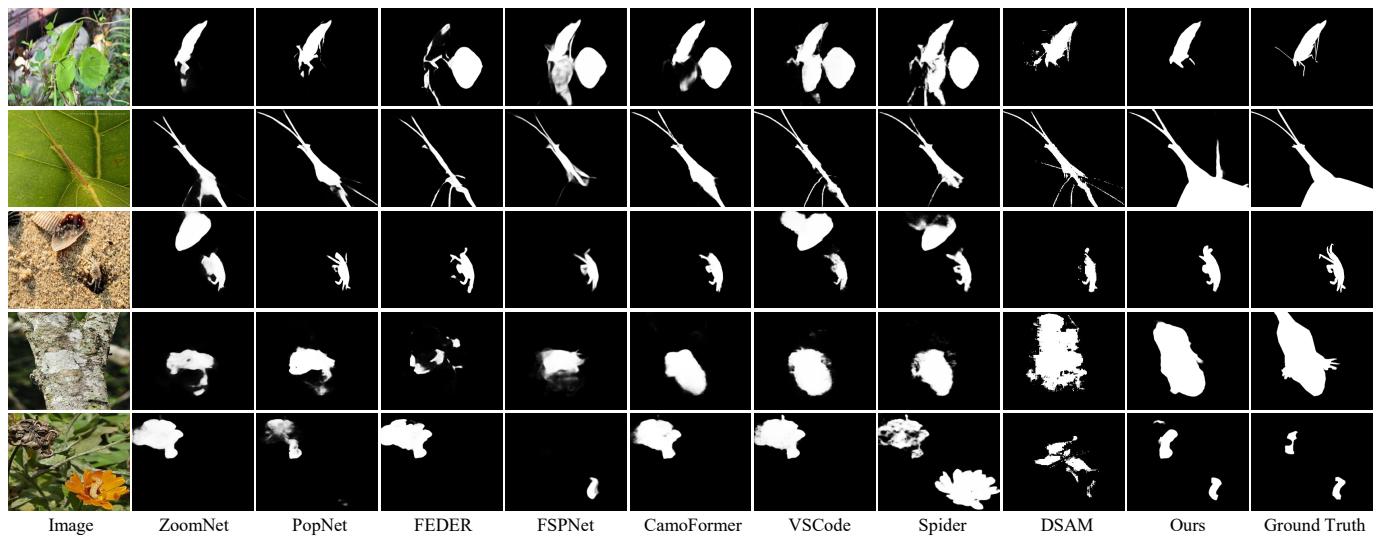


Fig. 3. Qualitative comparison on normal scenarios. Best view by zooming in.

and COD10K as the full training set, while the remaining data is assigned as the testing set. For the underwater benchmarks, we select MAS3K [61], RMAS [62], UFO120 [63], and RUWI [64]. MAS3K consists of 1,769 training pairs and 1,141 testing pairs with foreground objects. RMAS and UFO120 contain 3,014/500 and 1,500/120 training/testing pairs, respectively. Following [65], we set the training and testing sets of RUWI to 525 and 175 pairs, respectively. To further evaluate the model’s robustness under distribution shifts, inspired by [66], [67], we adopt three common degradation strategies: Gaussian blur (GB), Gaussian noise (GN), and contrast reduction (CR). These perturbations simulate real-world variations and assess the model’s adaptability in challenging conditions.

### B. Implementation Details.

We implement our model and conduct experiments using the PyTorch framework on NVIDIA A100 GPUs. We utilize PVT-v2 [68] and ResNet50 [69], both pre-trained on ImageNet, as backbone networks. For fair comparison, following [49], [53], we set the input resolution to 384×384. During training phrase, we leverage the AdamW optimizer with a batch size of 16, an initial learning rate of 0.001, and train for 200 epochs. Data augmentation strategies, including horizontal and vertical flipping, are applied to enhance generalization. During testing phrase, we apply the same settings for TTA stage as in the training phase. Moreover, we do not incorporate any post-processing techniques such as CRF to refine predictions. For

TABLE II

QUANTITATIVE COMPARISON ON DEGRADED BENCHMARK DATASETS. ATTR. INDICATES THE APPLIED DEGRADATION TYPE. \*: FOLLOWING THE SETTINGS IN [13], [57], WE USE [58] TO GENERATE DEPTH MAPS FOR DEGRADED TESTING IMAGES.

Method	Attr.	CHAMELEON (76)				CAMO (250)				COD10K (2026)				NC4K (4121)			
		$S_m \uparrow$	$E_m \uparrow$	$F_\beta^w \uparrow$	MAE $\downarrow$	$S_m \uparrow$	$E_m \uparrow$	$F_\beta^w \uparrow$	MAE $\downarrow$	$S_m \uparrow$	$E_m \uparrow$	$F_\beta^w \uparrow$	MAE $\downarrow$	$S_m \uparrow$	$E_m \uparrow$	$F_\beta^w \uparrow$	MAE $\downarrow$
DSAM * [57] MM'24	GN	0.805	0.894	0.714	0.059	0.768	0.863	0.704	0.087	0.808	0.897	0.700	0.043	0.828	0.907	0.765	0.054
	w/ HCL	0.837	0.914	0.758	0.051	0.800	0.892	0.756	0.075	0.833	0.907	0.741	0.038	0.855	0.921	0.797	0.048
Spider [56] ICML'24	GN	0.877	0.927	0.800	0.034	0.765	0.828	0.668	0.089	0.810	0.881	0.686	0.038	0.831	0.888	0.754	0.052
	w/ HCL	0.892	0.940	0.831	0.030	0.812	0.873	0.726	0.072	0.838	0.901	0.724	0.033	0.864	0.897	0.788	0.045
CamoFormer [55] TPAMI'24	GN	0.872	0.922	0.805	0.036	0.810	0.871	0.731	0.070	0.833	0.902	0.724	0.032	0.854	0.909	0.788	0.043
	w/ HCL	0.896	0.941	0.844	0.030	0.853	0.896	0.788	0.059	0.850	0.918	0.761	0.027	0.878	0.921	0.823	0.038
Ours (Baseline)	GN	0.802	0.887	0.740	0.043	0.753	0.840	0.709	0.075	0.785	0.868	0.686	0.038	0.807	0.875	0.760	0.042
Ours (w/ HCL)	GN	0.863	0.941	0.796	0.031	0.822	0.893	0.762	0.064	0.834	0.911	0.737	0.029	0.866	0.930	0.815	0.035
DSAM * [57] MM'24	GB	0.813	0.896	0.731	0.059	0.794	0.883	0.738	0.082	0.813	0.896	0.707	0.042	0.839	0.913	0.779	0.051
	w/ HCL	0.835	0.911	0.765	0.052	0.812	0.895	0.772	0.072	0.835	0.905	0.738	0.038	0.852	0.922	0.800	0.047
Spider [56] ICML'24	GB	0.866	0.922	0.785	0.039	0.759	0.808	0.647	0.091	0.802	0.869	0.668	0.044	0.824	0.878	0.739	0.057
	w/ HCL	0.888	0.936	0.826	0.033	0.803	0.849	0.708	0.077	0.828	0.887	0.710	0.036	0.851	0.899	0.790	0.047
CamoFormer [55] TPAMI'24	GB	0.787	0.837	0.659	0.058	0.759	0.811	0.651	0.088	0.789	0.859	0.653	0.044	0.813	0.866	0.721	0.058
	w/ HCL	0.832	0.872	0.700	0.046	0.797	0.852	0.705	0.072	0.815	0.891	0.696	0.035	0.861	0.903	0.778	0.044
Ours (Baseline)	GB	0.786	0.874	0.725	0.047	0.752	0.835	0.699	0.078	0.771	0.843	0.645	0.048	0.795	0.862	0.744	0.046
Ours (w/ HCL)	GB	0.842	0.943	0.773	0.035	0.815	0.896	0.758	0.064	0.827	0.907	0.724	0.029	0.856	0.923	0.799	0.038
DSAM * [57] MM'24	CR	0.846	0.923	0.773	0.046	0.828	0.908	0.787	0.064	0.840	0.917	0.750	0.034	0.868	0.930	0.822	0.041
	w/ HCL	0.856	0.929	0.783	0.044	0.839	0.917	0.801	0.062	0.851	0.929	0.763	0.032	0.879	0.937	0.830	0.040
Spider [56] ICML'24	CR	0.905	0.951	0.845	0.025	0.850	0.904	0.790	0.055	0.850	0.914	0.746	0.029	0.877	0.923	0.820	0.037
	w/ HCL	0.913	0.955	0.853	0.024	0.861	0.910	0.799	0.054	0.858	0.920	0.758	0.028	0.883	0.928	0.827	0.036
CamoFormer [55] TPAMI'24	CR	0.888	0.936	0.830	0.028	0.852	0.908	0.799	0.054	0.844	0.911	0.745	0.028	0.877	0.924	0.822	0.035
	w/ HCL	0.895	0.945	0.847	0.025	0.861	0.917	0.818	0.050	0.855	0.920	0.763	0.025	0.884	0.932	0.835	0.032
Ours (Baseline)	CR	0.822	0.901	0.754	0.041	0.825	0.888	0.775	0.056	0.802	0.874	0.715	0.033	0.830	0.889	0.783	0.037
Ours (w/ HCL)	CR	0.868	0.945	0.804	0.028	0.863	0.935	0.825	0.044	0.848	0.924	0.761	0.024	0.886	0.945	0.846	0.028

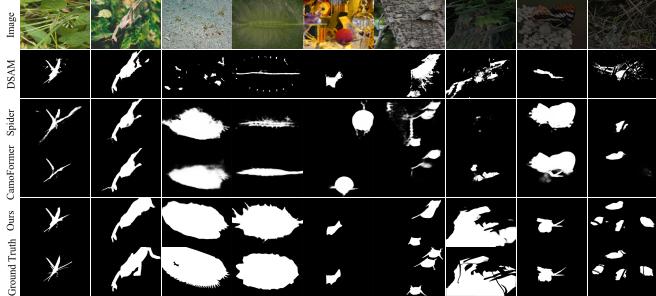


Fig. 4. Qualitative comparison on degraded camouflaged scenarios. (a)-(c): GN setting, (d)-(f): GB setting, (g)-(i): CR setting.

the reconstruction branch, we set the patch size to 16, and the depth of the encoder-decoder and the number of attention heads to 4, with the embedding dimension of 512 to reduce computational cost. To ensure a fair comparison, we compute the evaluation metrics using either the prediction results, pre-trained weights, or source codes provided by the projects.

### C. Evaluation Metrics.

We adopt four evaluation metrics: S-measure ( $S_m$ ), mean E-measure ( $E_m$ ), weighted F-measure, ( $F_\beta^w$ ), and Mean Absolute Error (MAE). Note that for the first three metrics, higher values indicate better performance, whereas for MAE, lower values are preferable.

### D. Comparison with State-of-the-Art Methods

1) *Quantitative Comparison.* In Table I, without relying on foundation models and incorporating depth maps, our method

outperforms DSAM by 2.0%, 1.8%, 3.0%, and 1.4% across four metrics on the most challenging NC4K dataset. We analyze that SAM, though pre-trained on large-scale natural images, tends to exhibit biases in camouflaged scenarios and struggles to adapt effectively with limited parameter updates and prior prompts. Moreover, our method also surpasses HitNet by 2.1%, 2.9%, 3.1%, and 1.3%. We argue that although HitNet improves representation capacity by increasing input resolution to preserve discriminative features, it lacks sensitivity to the diversity of camouflage patterns, thereby limiting the adaptability in complex settings. In Table II, under three interference conditions, CR shows the least performance degradation, followed by GN, with GB exhibiting the greatest decline. From the frequency perspective, CR preserves low-frequency structures essential for shape perception, GN adds high-frequency noise that can be partially suppressed, whereas GB removes critical high-frequency details, directly impairing boundary localization. Our method demonstrates consistent improvement of approximately 2%-8% on the NC4K. The higher gains in degraded scenes highlight the effectiveness of representation consistency and self-calibration. We integrate HCL into other methods and observe substantial performance improvements, indicating that train-test distribution shifts inherently exist and require additional customized components.

2) *Qualitative Comparison.* In Figure 3 and Figure 4, we present visualization results across various scenarios. Existing methods are prone to false negatives and false positives due to intrinsic physical characteristics and external disturbances, particularly in cases involving significant scale variations and multiple targets. In contrast, our approach effectively mitigates background clutter and imaging artifacts, accurately and comprehensively locating camouflaged regions.

TABLE III

QUANTITATIVE ABLATION OF PROPOSED COMPONENTS UNDER DIFFERENT CONDITIONS ON THE COD10K (TOP) AND NC4K (BOTTOM) DATASETS. C1, C2, C3, C4, AND C5 REPRESENT SPATIAL RECONSTRUCTION, LOW-FREQUENCY RECONSTRUCTION, HIGH-FREQUENCY RECONSTRUCTION, CONFIDENCE MAP-GUIDED AND VARIATIONAL FUSION, RESPECTIVELY.

Variant	HRR			TAG		PCC			Clean				Noise				Blur				Contrast				
	C1	C2	C3	C4	C5	$S_m \uparrow$	$E_m \uparrow$	$F_\beta^w \uparrow$	MAE $\downarrow$	$S_m \uparrow$	$E_m \uparrow$	$F_\beta^w \uparrow$	MAE $\downarrow$	$S_m \uparrow$	$E_m \uparrow$	$F_\beta^w \uparrow$	MAE $\downarrow$	$S_m \uparrow$	$E_m \uparrow$	$F_\beta^w \uparrow$	MAE $\downarrow$	$S_m \uparrow$	$E_m \uparrow$	$F_\beta^w \uparrow$	MAE $\downarrow$
I						0.818	0.889	0.747	0.029	0.785	0.868	0.686	0.038	0.771	0.843	0.645	0.048	0.802	0.874	0.715	0.033				
II	✓					0.827	0.898	0.755	0.027	0.794	0.874	0.695	0.036	0.780	0.855	0.657	0.044	0.812	0.884	0.720	0.031				
III		✓				0.822	0.894	0.753	0.028	0.793	0.870	0.688	0.037	0.781	0.850	0.651	0.045	0.808	0.885	0.718	0.032				
IV			✓			0.824	0.891	0.751	0.028	0.789	0.873	0.690	0.036	0.778	0.853	0.648	0.043	0.805	0.880	0.721	0.032				
V	✓	✓				0.830	0.903	0.759	0.027	0.801	0.878	0.693	0.035	0.783	0.862	0.670	0.042	0.822	0.889	0.723	0.030				
VI	✓		✓			0.831	0.902	0.754	0.026	0.798	0.880	0.698	0.034	0.785	0.864	0.665	0.041	0.817	0.890	0.725	0.030				
VII		✓	✓			0.826	0.899	0.755	0.025	0.796	0.875	0.697	0.034	0.784	0.860	0.669	0.040	0.815	0.888	0.724	0.031				
VIII	✓	✓	✓			0.839	0.907	0.761	0.024	0.807	0.887	0.705	0.033	0.793	0.875	0.684	0.037	0.834	0.901	0.732	0.029				
IX	✓	✓	✓	✓		0.844	0.914	0.765	0.024	0.814	0.895	0.718	0.032	0.803	0.888	0.700	0.035	0.840	0.900	0.741	0.027				
X	✓	✓	✓	✓	✓	0.848	0.920	0.771	0.023	0.820	0.902	0.725	0.031	0.810	0.896	0.711	0.032	0.839	0.911	0.748	0.026				
XI	✓	✓	✓	✓	✓	<b>0.860</b>	<b>0.933</b>	<b>0.781</b>	<b>0.022</b>	<b>0.834</b>	<b>0.911</b>	<b>0.737</b>	<b>0.029</b>	<b>0.827</b>	<b>0.907</b>	<b>0.724</b>	<b>0.029</b>	<b>0.848</b>	<b>0.924</b>	<b>0.761</b>	<b>0.024</b>				
XII						0.834	0.900	0.801	0.032	0.807	0.875	0.760	0.042	0.795	0.862	0.744	0.046	0.830	0.889	0.783	0.037				
XIII	✓					0.843	0.914	0.811	0.031	0.819	0.886	0.771	0.041	0.811	0.874	0.756	0.044	0.843	0.900	0.795	0.037				
XIV		✓				0.838	0.907	0.809	0.032	0.815	0.880	0.766	0.042	0.803	0.868	0.749	0.045	0.835	0.897	0.788	0.037				
XV			✓			0.840	0.909	0.805	0.031	0.810	0.883	0.764	0.042	0.804	0.869	0.751	0.045	0.837	0.893	0.790	0.038				
XVI	✓	✓				0.847	0.919	0.820	0.030	0.826	0.894	0.776	0.040	0.816	0.882	0.761	0.043	0.851	0.905	0.807	0.035				
XVII	✓		✓			0.845	0.921	0.815	0.030	0.824	0.897	0.774	0.041	0.820	0.877	0.763	0.044	0.853	0.907	0.803	0.034				
XVIII		✓	✓			0.842	0.912	0.817	0.030	0.820	0.885	0.772	0.040	0.809	0.876	0.759	0.042	0.841	0.902	0.792	0.036				
XIX	✓	✓	✓			0.859	0.930	0.835	0.028	0.845	0.909	0.784	0.039	0.829	0.894	0.771	0.042	0.870	0.920	0.828	0.031				
XX	✓	✓	✓	✓		0.870	0.935	0.838	0.027	0.842	0.918	0.796	0.038	0.837	0.905	0.783	0.041	0.876	0.925	0.837	0.029				
XXI	✓	✓	✓	✓	✓	0.878	0.941	0.844	0.026	0.852	0.923	0.802	0.036	0.843	0.913	0.787	0.040	0.880	0.931	0.833	0.028				
XXII	✓	✓	✓	✓	✓	<b>0.891</b>	<b>0.950</b>	<b>0.856</b>	<b>0.026</b>	<b>0.866</b>	<b>0.930</b>	<b>0.815</b>	<b>0.035</b>	<b>0.856</b>	<b>0.923</b>	<b>0.799</b>	<b>0.038</b>	<b>0.886</b>	<b>0.945</b>	<b>0.846</b>	<b>0.028</b>				

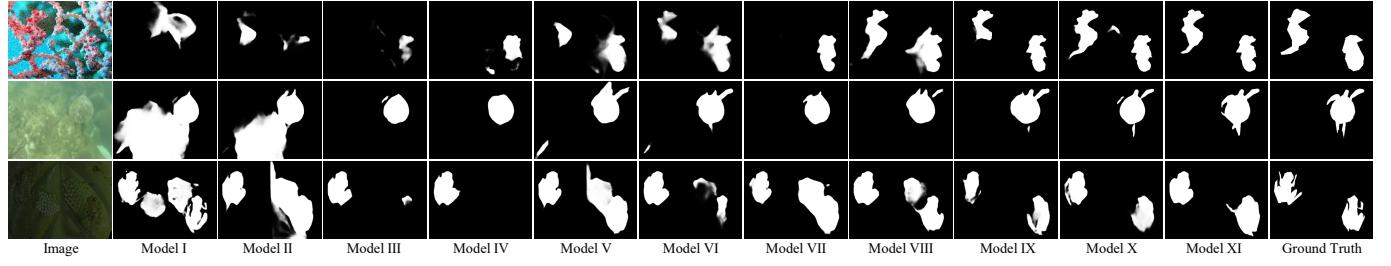


Fig. 5. Qualitative ablation of the proposed components. Variants corresponding to Table III.

### E. Ablation Study

We quantitatively and qualitatively validate the impact of the proposed components, method variants, and critical hyper-parameter settings.

1) *Effect of the crucial components.* In Table III, we integrate components through one-step process. We find that: 1) For the HRR, directly using spatial reconstruction (*i.e.*, vanilla MAE) does not yield significant results, while combining with frequency reconstruction achieves complementary effects, especially in degraded scenarios. 2) Integrating the TAG generally improves performance, although some cases show a decline. 3) The confidence map guidance and variational fusion components exhibit progressive coupling. We analyze: 1) Spatial and frequency domains emphasize different aspects of information recovery, and the variations in dynamic scenes and camouflage patterns necessitate more comprehensive clue extraction. 2) We utilize the TAG as a bridge to inject prior knowledge, enhancing detection representation when reconstruction results are ideal; conversely, it may weaken or even obscure when results are suboptimal. 3) The confidence map provides the robust prototype, while the variational component addresses internal instabilities and weight biases, transforming

TABLE IV  
QUANTITATIVE COMPARISON OF TTA STRATEGIES.

Method	COD10K-GB				NC4K-GB			
	$S_m \uparrow$	$E_m \uparrow$	$F_\beta^w \uparrow$	MAE $\downarrow$	$S_m \uparrow$	$E_m \uparrow$	$F_\beta^w \uparrow$	MAE $\downarrow$
TTT-Rot [2]	0.795	0.874	0.694	<b>0.033</b>	0.827	0.875	0.764	0.045
TENT [33]	0.782	0.845	0.700	0.035	0.816	0.867	0.752	0.041
UDA [70]	0.800	0.866	0.687	0.036	0.824	0.879	0.748	0.044
TTT-MIM [71]	0.792	0.879	0.695	0.034	0.821	0.890	0.762	0.043
Ours	<b>0.827</b>	<b>0.907</b>	<b>0.724</b>	<b>0.029</b>	<b>0.856</b>	<b>0.923</b>	<b>0.799</b>	<b>0.038</b>

point estimates into probability distributions to mitigate the turbulence. In Figure 5, as the proposed components are gradually integrated, the detection areas approach the ground truth. In Table II, embedding the HCL into other methods can also effectively boost performance by about 3%.

2) *Analysis of TTA strategies.* In Table IV, we consider: 1) Dynamically adjusting the statistics or parameters of the Batch Normalization layer to adapt to the distribution of the target domain [33], [70]; 2) Updating parameters through rotation predictions to construct consistency [2]. Our HCL shows the best performance. We argue that: 1) Both TENT and UDA struggle to handle confusable visual space and structural am-

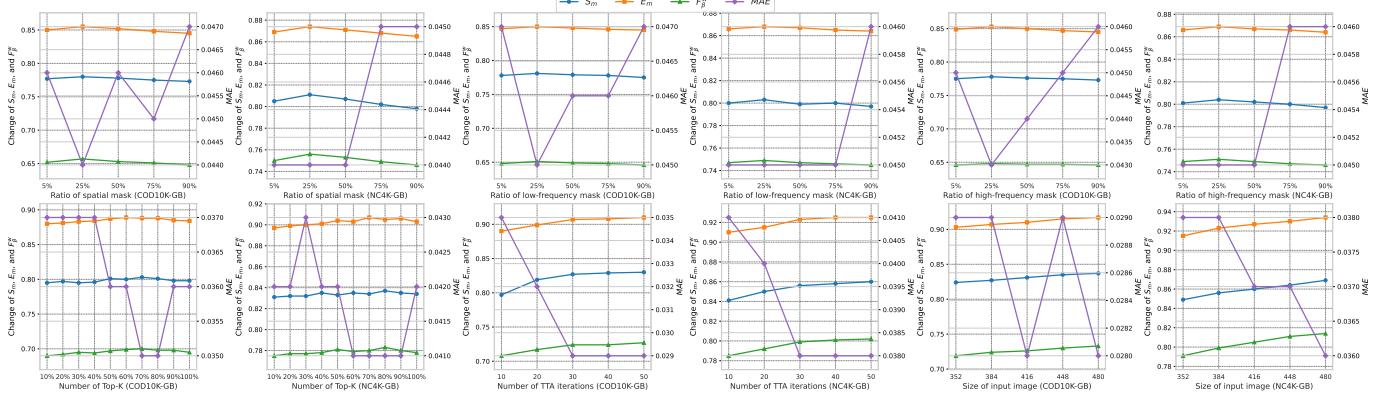


Fig. 6. Quantitative ablation on the ratios of spatial, low-frequency, and high-frequency masks (with optimal spatial ratio), as well as the effects of Top-K selection, TTA iterations, and input image resolution

TABLE V  
QUANTITATIVE COMPARISON OF MASK STRATEGIES.

Method	COD10K-GB				NC4K-GB			
	$S_m \uparrow$	$E_m \uparrow$	$F_\beta^w \uparrow$	MAE $\downarrow$	$S_m \uparrow$	$E_m \uparrow$	$F_\beta^w \uparrow$	MAE $\downarrow$
Round	0.820	<b>0.909</b>	0.713	0.030	0.849	<b>0.925</b>	0.789	0.039
Grid	<b>0.829</b>	0.900	0.715	<b>0.028</b>	0.853	0.920	0.791	<b>0.039</b>
Block	0.818	0.904	0.718	0.029	0.848	0.918	<b>0.792</b>	0.038
Random	0.827	0.907	<b>0.724</b>	0.029	<b>0.856</b>	0.923	<b>0.799</b>	<b>0.038</b>

TABLE VI  
QUANTITATIVE COMPARISON OF INTERACTIVE STRATEGIES.

Method	COD10K-GB				NC4K-GB			
	$S_m \uparrow$	$E_m \uparrow$	$F_\beta^w \uparrow$	MAE $\downarrow$	$S_m \uparrow$	$E_m \uparrow$	$F_\beta^w \uparrow$	MAE $\downarrow$
Spatial Cross	0.797	0.880	0.693	0.036	<b>0.840</b>	0.898	0.774	0.042
Channel Cross	0.795	<b>0.889</b>	0.690	0.037	0.833	0.900	<b>0.780</b>	0.042
TAG	<b>0.803</b>	0.888	<b>0.700</b>	<b>0.035</b>	0.837	<b>0.905</b>	<b>0.783</b>	<b>0.041</b>

biguity; 2) Due to the high coupling between foreground and background, camouflaged regions can be positioned anywhere while preserving characteristics. In other words, rotation is inadequate for capturing camouflage and degradation cues. In contrast, image reconstruction explores the intrinsic relationships between regions, and the representations learned through self-supervised learning align better with the detection branch.

3) *Analysis of mask strategies.* In Table V, random masking achieves superior performance compared to fixed geometric patterns (round, grid, block). We attribute to: 1) Fixed masking strategies consistently eliminate specific frequency components, creating persistent blind spots in feature space and learning. 2) Random masking dynamically varies the suppressed frequency bands across spatial locations, ensuring no single spectral component is entirely excluded during reconstruction. This encourages the model to holistically integrate multi-spectrum cues rather than over-relying on localized patterns. 3) The stochastic characteristic of random masks enhances robustness against diverse corruption types by preventing overfitting to artificial geometric artifacts.

4) *Analysis of interactive strategies.* In Table VI, our TAG outperforms other cross-attention-based strategies, which can

TABLE VII  
QUANTITATIVE COMPARISON OF PROTOTYPE INTEGRATION STRATEGIES.

Method	COD10K-GB				NC4K-GB			
	$S_m \uparrow$	$E_m \uparrow$	$F_\beta^w \uparrow$	MAE $\downarrow$	$S_m \uparrow$	$E_m \uparrow$	$F_\beta^w \uparrow$	MAE $\downarrow$
Sum	<b>0.823</b>	0.901	0.717	0.032	<b>0.848</b>	0.914	<b>0.795</b>	0.040
Concat	0.815	0.899	0.718	0.032	0.840	0.917	0.790	0.039
Cross Attention	0.820	0.902	0.715	0.031	0.847	0.915	0.792	0.039
Variational Fusion	<b>0.827</b>	<b>0.907</b>	<b>0.724</b>	<b>0.029</b>	<b>0.856</b>	<b>0.923</b>	<b>0.799</b>	<b>0.038</b>

be attributed to: 1) Instead of direct feature-level interaction, we refine internal features along the channel dimension and construct an affinity map, encouraging coherent and compact representations while mitigating noise entanglement; 2) A Top-K selection operator is introduced to filter out low-response values, enabling entropy compression and enhancing the quality of information propagation.

5) *Analysis of prototype integration strategies.* In Table VII, variational fusion demonstrates superior performance. We analyze that: 1) Unlike fixed fusion schemes, variational fusion learns distribution-based representations, allowing to dynamically assign confidence-aware weights to different features based on reliability; 2) By explicitly incorporating uncertainty, variational fusion suppresses noisy or redundant information and highlights informative cues, improving robustness in challenging conditions such as low saliency or distribution shifts; 3) Fusion in the probabilistic latent space facilitates smoother integration of complementary cues, enabling more consistent and discriminative representations than simple arithmetic operations or attention that may amplify noise.

6) *Analysis of mask ratios.* In Figure 6, spatial and frequency masking both achieve optimal performance at 25% mask ratio, as moderate masking balances feature stability and contextual integrity. Excessive spatial masking ( $>25\%$ ) removes critical local details, while excessive frequency masking ( $>25\%$ ) disproportionately disrupts high-frequency textures or low-frequency structures, creating irrecoverable information gaps. Besides, frequency masking shows sharper performance decay due to the entangled spectral components increasing reconstruction ambiguity, whereas spatial masking preserves more structural coherence at higher ratios.

TABLE VIII  
QUANTITATIVE COMPARISON ON DEGRADED UNDERWATER BENCHMARK DATASETS. N/A INDICATES NO DEGRADATION APPLIED.

Method	Attr.	MAS3K (1141)				RMAS (500)				UFO120 (120)				RUWI (175)				
		$S_m \uparrow$	$E_m \uparrow$	$F_\beta^w \uparrow$	MAE $\downarrow$	$S_m \uparrow$	$E_m \uparrow$	$F_\beta^w \uparrow$	MAE $\downarrow$	$S_m \uparrow$	$E_m \uparrow$	$F_\beta^w \uparrow$	MAE $\downarrow$	$S_m \uparrow$	$E_m \uparrow$	$F_\beta^w \uparrow$	MAE $\downarrow$	
DualSAM [65]	CVPR'24	N/A	0.884	0.933	0.838	0.023	0.860	0.944	0.812	0.022	0.856	0.914	0.864	0.064	0.903	0.959	0.939	0.035
MASSAM [56]	IJCAI'24	N/A	0.887	0.938	0.840	0.025	0.865	0.948	0.819	0.021	0.861	0.914	0.864	0.063	0.894	0.961	0.941	0.035
DualSAM [65]	CVPR'24	GN	0.814	0.870	0.755	0.045	0.749	0.781	0.613	0.041	0.798	0.865	0.793	0.085	0.783	0.872	0.827	0.056
	w/ HCL	GN	0.853	0.901	0.799	0.032	0.785	0.833	0.651	0.033	0.827	0.892	0.831	0.075	0.832	0.903	0.869	0.045
MASSAM [56]	IJCAI'24	GN	0.833	0.866	0.748	0.042	0.733	0.766	0.587	0.047	0.784	0.850	0.779	0.083	0.773	0.853	0.812	0.053
	w/ HCL	GN	0.844	0.900	0.777	0.036	0.770	0.801	0.629	0.038	0.814	0.870	0.794	0.077	0.837	0.898	0.867	0.048
DualSAM [65]	CVPR'24	GB	0.795	0.857	0.746	0.050	0.740	0.769	0.615	0.044	0.778	0.848	0.758	0.082	0.795	0.863	0.834	0.051
	w/ HCL	GB	0.844	0.895	0.788	0.037	0.784	0.808	0.648	0.035	0.811	0.875	0.799	0.073	0.838	0.912	0.873	0.045
MASSAM [56]	IJCAI'24	GB	0.785	0.844	0.750	0.046	0.745	0.781	0.620	0.043	0.773	0.835	0.739	0.079	0.763	0.849	0.830	0.055
	w/ HCL	GB	0.839	0.875	0.770	0.037	0.769	0.800	0.649	0.037	0.810	0.874	0.786	0.071	0.803	0.891	0.878	0.047
DualSAM [65]	CVPR'24	CR	0.844	0.898	0.807	0.035	0.833	0.909	0.778	0.029	0.827	0.876	0.820	0.074	0.873	0.941	0.918	0.042
	w/ HCL	CR	0.867	0.921	0.819	0.031	0.840	0.930	0.790	0.026	0.844	0.899	0.840	0.071	0.888	0.945	0.938	0.038
MASSAM [56]	IJCAI'24	CR	0.849	0.908	0.803	0.036	0.828	0.911	0.772	0.031	0.820	0.869	0.823	0.073	0.865	0.935	0.911	0.044
	w/ HCL	CR	0.864	0.916	0.824	0.032	0.844	0.924	0.796	0.027	0.834	0.883	0.833	0.069	0.882	0.940	0.924	0.038

7) *Analysis of Top-K selection operator.* In Figure 6, selecting the top 70% (80%) values yields the best performance. When the threshold is lower, performance improves gradually; when it exceeds 70% (80%), a performance drop is observed. We argue that: 1) Retaining too few connections (low Top-K ratio) limits the model’s ability to capture sufficient contextual cues, especially long-range dependencies essential for complex scenes; 2) Retaining too many connections introduces noisy or less informative interactions, which may overwhelm the discriminative patterns and degrade the effectiveness of the feature refinement process. The optimal threshold achieves a balance between preserving informative associations and filtering out noise, optimizing both precision and generalization.

8) *Analysis of TTA iterations.* In Figure 6, we observe a turning point at 30 adaptation iterations: performance improves steadily when the iteration count is below 30, and then stabilizes beyond that. We analyze that: 1) During the early testing phase, a clear distribution mismatch exists between the test samples and the model’s learned feature space, as indicated by the significant discrepancy between prototypes  $P$  and  $P_{rec}$  distributions; 2) Through iterative refinement (e.g., progressive calibration guided by confidence map  $\Phi$ ), the model gradually aligns its internal parameters with the distribution of samples, reaching a relatively optimal state. Further updates yield diminishing returns, resulting in performance convergence.

9) *Analysis of input size.* In Figure 6, the performance generally exhibits a positive correlation with the input resolution. We attribute this to: 1) Higher-resolution inputs provide denser spectral sampling in the Fourier domain (expanded along the  $u$  and  $v$  dimensions), enabling the model to more precisely separate critical high-frequency details and structured low-frequency components; 2) The larger pixel space offers greater capacity to preserve discriminative features of the target, mitigating information loss during the encoding process.

#### F. Broader Impacts

We evaluate state-of-the-art methods on four underwater benchmarks under the same three degradation settings, and further integrate HCL to explore its performance gains.

1) *Quantitative Analysis.* In Table VIII, we present results under GN, GB, and CR conditions. Across all datasets, GB causes the most severe degradation, followed by GN, while CR shows the least impact. By integrating HCL, both DualSAM and MASSAM consistently achieve performance gains across metrics and datasets. Notably, for DualSAM on the MAS3K, HCL improves  $F_\beta^w$  by 4.2% and reduces MAE by 1.3% under GB, highlighting its ability to recover fine details. Under GN and CR, HCL still brings steady gains of 1.2%–4.4%

2) *Qualitative Analysis.* In Figure 7, we present qualitative comparisons under various challenging scenarios. Integrating HCL boosts comprehensive perception. In multi-object scenes (columns a, b, and e), it preserves complete instance integrity and avoids omissions by leveraging prototype consistency and fusion to suppress interference between adjacent objects. For objects with highly irregular shapes (columns c, d, and h), the variant captures long-range dependencies and preserves boundary details through non-local affinity modeling. In low-saliency settings (column f), where foreground-background contrast is weak, our hierarchical frequency reconstruction enhances subtle discriminative cues, effectively improving foreground-background separation. In camouflaged instances (columns f, g, and h), our method leverages task-guided reconstruction and entropy-based confidence weighting to highlight critical boundaries and suppress false detections.

3) *Why can HCL work under underwater conditions?* 1) HCL leverages spatial- and frequency-decoupled reconstruction for TTA, which is not limited to camouflage scenarios and can be extended to other binary segmentation tasks; 2) Existing underwater benchmarks and learning paradigms also suffer from train-test distribution shifts and fixed model parameters, making TTA essential, especially when degradation further amplifies distributional discrepancies.

#### G. Limitation and Future Work

In Figure 8, the HCL framework remains challenged in extreme scenarios: highly camouflaged targets (e.g., insects perfectly mimicking leaves) may evade detection as global feature homogenization overshadows subtle discriminative cues,

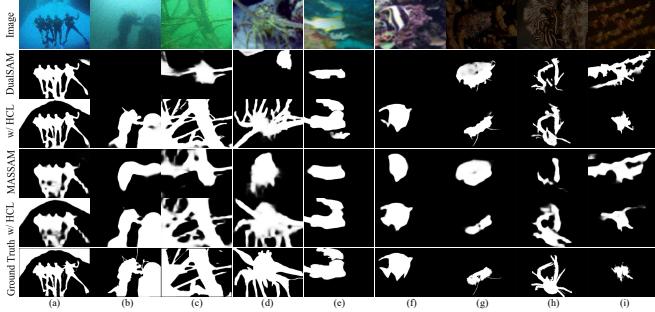


Fig. 7. Qualitative comparison on degraded underwater scenarios. (a)-(c): GN setting, (d)-(f): GB setting, (g)-(i): CR setting.

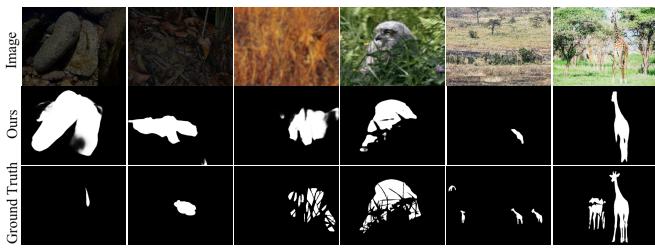


Fig. 8. Failure cases under degradation conditions.

while highly irregular backgrounds (*e.g.*, forest floors with fragmented shadows) disrupt frequency-based analysis, causing localization errors. These limitations reflect the inherent difficulty in balancing structural preservation and noise suppression under environmental uncertainty. We will focus on: 1) Integrating attention with dynamic sparsity to enhance feature disentanglement by focusing on sample-critical frequency bands; 2) Extending the HCL framework to instance-level and video-level scenarios.

## V. CONCLUSION

This paper tackles the challenge of COD under distribution shifts, where static models struggle to generalize. We propose HCL, a dynamic TTA framework that enables sample-specific self-calibration via spatial-frequency decoupled reconstruction. Key components such as TAG and entropy-guided prototype fusion enhance structural consistency and feature robustness. These findings highlight the importance of dynamic adaptation and suggest HCL's potential for broader binary segmentation tasks.

## REFERENCES

- [1] D.-P. Fan, G.-P. Ji, T. Zhou, G. Chen, H. Fu, J. Shen, and L. Shao, “Planet: Parallel reverse attention network for polyp segmentation,” in *International conference on medical image computing and computer-assisted intervention*. Springer, 2020, pp. 263–273.
- [2] Y. Sun, X. Wang, Z. Liu, J. Miller, A. Efros, and M. Hardt, “Test-time training with self-supervision for generalization under distribution shifts,” in *International conference on machine learning*. PMLR, 2020, pp. 9229–9248.
- [3] S. Niu, J. Wu, Y. Zhang, Y. Chen, S. Zheng, P. Zhao, and M. Tan, “Efficient test-time model adaptation without forgetting,” in *International conference on machine learning*. PMLR, 2022, pp. 16 888–16 905.
- [4] K. He, X. Chen, S. Xie, Y. Li, P. Dollár, and R. Girshick, “Masked autoencoders are scalable vision learners,” in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2022, pp. 16 000–16 009.
- [5] M. Zha, G. Wang, Y. Pei, T. Li, X. Tang, C. Li, Y. Yang, and H. T. Shen, “Heterogeneous experts and hierarchical perception for underwater salient object detection,” *IEEE Transactions on Image Processing*, 2025.
- [6] M. Zha, Y. Pei, G. Wang, T. Li, Y. Yang, W. Qian, and H. T. Shen, “Weakly-supervised mirror detection via scribble annotations,” in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 38, no. 7, 2024, pp. 6953–6961.
- [7] Y. Zhao, L. Zhao, Q. Yu, L. Sheng, J. Zhang, and D. Xu, “Distortion-aware transformer in 360 salient object detection,” in *Proceedings of the 31st ACM International Conference on Multimedia*, 2023, pp. 499–508.
- [8] F. Liu, Y. Liu, K. Xu, S. Ye, G. P. Hancke, and R. W. Lau, “Language-guided salient object ranking,” in *Proceedings of the Computer Vision and Pattern Recognition Conference*, 2025, pp. 29 803–29 813.
- [9] H. Guan, J. Lin, and R. W. Lau, “A contrastive-learning framework for unsupervised salient object detection,” *IEEE Transactions on Image Processing*, 2025.
- [10] Y. Sun, S. Wang, C. Chen, and T.-Z. Xiang, “Boundary-guided camouflaged object detection,” *arXiv preprint arXiv:2207.00794*, 2022.
- [11] J. Zhu, X. Zhang, S. Zhang, and J. Liu, “Inferring camouflaged objects by texture-aware interactive guidance network,” in *Proceedings of the AAAI conference on artificial intelligence*, vol. 35, no. 4, 2021, pp. 3599–3607.
- [12] Y. Sun, C. Xu, J. Yang, H. Xuan, and L. Luo, “Frequency-spatial entanglement learning for camouflaged object detection,” in *European Conference on Computer Vision*. Springer, 2024, pp. 343–360.
- [13] Z. Wu, D. P. Paudel, D.-P. Fan, J. Wang, S. Wang, C. Demonceaux, R. Timofte, and L. Van Gool, “Source-free depth for object pop-out,” in *Proceedings of the IEEE/CVF international conference on computer vision*, 2023, pp. 1032–1042.
- [14] Y. Zhang, J. Zhang, W. Hamidouche, and O. Deforges, “Predictive uncertainty estimation for camouflaged object detection,” *IEEE Transactions on Image Processing*, vol. 32, pp. 3580–3591, 2023.
- [15] H. Zhang, Y. Lyu, Q. Yu, H. Liu, H. Ma, D. Yuan, and Y. Yang, “Unlocking attributes’ contribution to successful camouflage: A combined textual and visual analysis strategy,” in *European Conference on Computer Vision*. Springer, 2024, pp. 315–331.
- [16] Y. Liu, D. Zhang, Q. Zhang, and J. Han, “Integrating part-object relationship and contrast for camouflaged object detection,” *IEEE Transactions on Information Forensics and Security*, vol. 16, pp. 5154–5166, 2021.
- [17] X. Zhang, J. Zhou, L. Yan, S. Zhong, and X. Zou, “Hunt camouflaged objects via revealing mutation regions,” *IEEE Transactions on Information Forensics and Security*, 2025.
- [18] H. Chen, D. Shao, G. Guo, and S. Gao, “Just a hint: Point-supervised camouflaged object detection,” in *European Conference on Computer Vision*. Springer, 2024, pp. 332–348.
- [19] X. Lai, Z. Yang, J. Hu, S. Zhang, L. Cao, G. Jiang, Z. Wang, S. Zhang, and R. Ji, “Camoteacher: Dual-rotation consistency learning for semi-supervised camouflaged object detection,” in *European Conference on Computer Vision*. Springer, 2024, pp. 438–455.
- [20] H. Li, C.-M. Feng, Y. Xu, T. Zhou, L. Yao, and X. Chang, “Zero-shot camouflaged object detection,” *IEEE Transactions on Image Processing*, vol. 32, pp. 5126–5137, 2023.
- [21] Z. Wang, Y. Li, Y. Yang, Y. Li, and G. Liu, “Few-shot camouflaged object segmentation,” in *2024 International Joint Conference on Neural Networks (IJCNN)*. IEEE, 2024, pp. 1–10.
- [22] J. Zhao, X. Li, F. Yang, Q. Zhai, A. Luo, Z. Jiao, and H. Cheng, “Focus-diffuser: Perceiving local disparities for camouflaged object detection,” in *European Conference on Computer Vision*. Springer, 2024, pp. 181–198.
- [23] H. Chen, P. Wei, G. Guo, and S. Gao, “Sam-cod+: Sam-guided unified framework for weakly-supervised camouflaged object detection,” *IEEE Transactions on Circuits and Systems for Video Technology*, 2024.
- [24] Z. Luo, N. Liu, W. Zhao, X. Yang, D. Zhang, D.-P. Fan, F. Khan, and J. Han, “Vscode: General visual salient and camouflaged object detection with 2d prompt learning,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024, pp. 17 169–17 180.
- [25] X. Zhang, Z. Yu, L. Zhao, D.-P. Fan, and G. Xiao, “Comprompter: reconceptualized segment anything model with multiprompt network for camouflaged object detection,” *Science China Information Sciences*, vol. 68, no. 1, p. 112104, 2025.
- [26] Z. He, C. Xia, S. Qiao, and J. Li, “Text-prompt camouflaged instance segmentation with graduated camouflage learning,” in *Proceedings of the 32nd ACM International Conference on Multimedia*, 2024, pp. 5584–5593.
- [27] W. Hui, Z. Zhu, S. Zheng, and Y. Zhao, “Endow sam with keen eyes: Temporal-spatial prompt learning for video camouflaged object

- detection,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024, pp. 19 058–19 067.
- [28] L. Wang, J. Yang, Y. Zhang, F. Wang, and F. Zheng, “Depth-aware concealed crop detection in dense agricultural scenes,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024, pp. 17 201–17 211.
- [29] J. Yang, Q. Wang, F. Zheng, P. Chen, A. Leonardis, and D.-P. Fan, “Plantcamo: Plant camouflage detection,” *CAAI Artificial Intelligence Research (AIR)*, 2025.
- [30] C. Hao, Z. Yu, X. Liu, J. Xu, H. Yue, and J. Yang, “A simple yet effective network based on vision transformer for camouflaged object and salient object detection,” *IEEE Transactions on Image Processing*, 2025.
- [31] Y. Pang, X. Zhao, J. Zuo, L. Zhang, and H. Lu, “Open-vocabulary camouflaged object segmentation,” in *European Conference on Computer Vision*. Springer, 2024, pp. 476–495.
- [32] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell, P. Mishkin, J. Clark *et al.*, “Learning transferable visual models from natural language supervision,” in *International conference on machine learning*. PMLR, 2021, pp. 8748–8763.
- [33] D. Wang, E. Shelhamer, S. Liu, B. Olshausen, and T. Darrell, “Tent: Fully test-time adaptation by entropy minimization,” *arXiv preprint arXiv:2006.10726*, 2020.
- [34] M. Jang, S.-Y. Chung, and H. W. Chung, “Test-time adaptation via self-training with nearest neighbor information,” *arXiv preprint arXiv:2207.10792*, 2022.
- [35] D. Brahma and P. Rai, “A probabilistic framework for lifelong test-time adaptation,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 3582–3591.
- [36] B. Zhao, C. Chen, and S.-T. Xia, “Delta: degradation-free fully test-time adaptation,” *arXiv preprint arXiv:2301.13018*, 2023.
- [37] J. Ma, “Improved self-training for test-time adaptation,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024, pp. 23 701–23 710.
- [38] W. Liu, X. Shen, H. Li, X. Bi, B. Liu, C.-M. Pun, and X. Cun, “Depth-aware test-time training for zero-shot video object segmentation,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024, pp. 19 218–19 227.
- [39] J. Hu, J. Lin, S. Gong, and W. Cai, “Relax image-specific prompt requirement in sam: A single generic prompt for segmenting camouflaged objects,” in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 38, no. 11, 2024, pp. 12 511–12 518.
- [40] L. Jiang, B. Dai, W. Wu, and C. C. Loy, “Focal frequency loss for image reconstruction and synthesis,” in *Proceedings of the IEEE/CVF international conference on computer vision*, 2021, pp. 13 919–13 929.
- [41] Y. Pang, X. Zhao, T.-Z. Xiang, L. Zhang, and H. Lu, “Zoom in and out: A mixed-scale triplet network for camouflaged object detection,” in *Proceedings of the IEEE/CVF Conference on computer vision and pattern recognition*, 2022, pp. 2160–2170.
- [42] D.-P. Fan, G.-P. Ji, G. Sun, M.-M. Cheng, J. Shen, and L. Shao, “Camouflaged object detection,” in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2020, pp. 2777–2787.
- [43] Q. Zhai, X. Li, F. Yang, C. Chen, H. Cheng, and D.-P. Fan, “Mutual graph learning for camouflaged object detection,” in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2021, pp. 12 997–13 007.
- [44] H. Mei, G.-P. Ji, Z. Wei, X. Yang, X. Wei, and D.-P. Fan, “Camouflaged object segmentation with distraction mining,” in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2021, pp. 8772–8781.
- [45] F. Yang, Q. Zhai, X. Li, R. Huang, A. Luo, H. Cheng, and D.-P. Fan, “Uncertainty-guided transformer reasoning for camouflaged object detection,” in *Proceedings of the IEEE/CVF international conference on computer vision*, 2021, pp. 4146–4155.
- [46] Y. Lv, J. Zhang, Y. Dai, A. Li, B. Liu, N. Barnes, and D.-P. Fan, “Simultaneously localize, segment and rank the camouflaged objects,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, pp. 11 591–11 601.
- [47] D.-P. Fan, G.-P. Ji, M.-M. Cheng, and L. Shao, “Concealed object detection,” *IEEE transactions on pattern analysis and machine intelligence*, vol. 44, no. 10, pp. 6024–6042, 2021.
- [48] M. Zhang, S. Xu, Y. Piao, D. Shi, S. Lin, and H. Lu, “Preynet: Preying on camouflaged objects,” in *Proceedings of the 30th ACM international conference on multimedia*, 2022, pp. 5323–5332.
- [49] C. He, K. Li, Y. Zhang, L. Tang, Y. Zhang, Z. Guo, and X. Li, “Camouflaged object detection with feature decomposition and edge reconstruction,” in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2023, pp. 22 046–22 055.
- [50] A. Kirillov, E. Mintun, N. Ravi, H. Mao, C. Rolland, L. Gustafson, T. Xiao, S. Whitehead, A. C. Berg, W.-Y. Lo *et al.*, “Segment anything,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2023, pp. 4015–4026.
- [51] X. Hu, S. Wang, X. Qin, H. Dai, W. Ren, D. Luo, Y. Tai, and L. Shao, “High-resolution iterative feedback network for camouflaged object detection,” in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 37, no. 1, 2023, pp. 881–889.
- [52] R. Cong, M. Sun, S. Zhang, X. Zhou, W. Zhang, and Y. Zhao, “Frequency perception network for camouflaged object detection,” in *Proceedings of the ACM International Conference on Multimedia*, 2023, pp. 1179–1189.
- [53] Z. Huang, H. Dai, T.-Z. Xiang, S. Wang, H.-X. Chen, J. Qin, and H. Xiong, “Feature shrinkage pyramid for camouflaged object detection with transformers,” in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2023, pp. 5557–5566.
- [54] W. Liu, X. Shen, C.-M. Pun, and X. Cun, “Explicit visual prompting for low-level structure segmentations,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 19 434–19 445.
- [55] B. Yin, X. Zhang, D.-P. Fan, S. Jiao, M.-M. Cheng, L. Van Gool, and Q. Hou, “Camoformer: Masked separable attention for camouflaged object detection,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2024.
- [56] X. Zhao, Y. Pang, W. Ji, B. Sheng, J. Zuo, L. Zhang, and H. Lu, “Spider: A unified framework for context-dependent concept understanding,” *arXiv preprint arXiv:2405.01002*, 2024.
- [57] Z. Yu, X. Zhang, L. Zhao, Y. Bin, and G. Xiao, “Exploring deeper! segment anything model with depth perception for camouflaged object detection,” in *Proceedings of the 32nd ACM International Conference on Multimedia*, 2024, pp. 4322–4330.
- [58] R. Ranftl, A. Bochkovskiy, and V. Koltun, “Vision transformers for dense prediction,” in *Proceedings of the IEEE/CVF international conference on computer vision*, 2021, pp. 12 179–12 188.
- [59] T.-N. Le, T. V. Nguyen, Z. Nie, M.-T. Tran, and A. Sugimoto, “Anabanch networks for camouflaged object segmentation,” *Computer vision and image understanding*, vol. 184, pp. 45–56, 2019.
- [60] P. Skurowski, H. Abdulameer, J. Błaszczyk, T. Depta, A. Kornacki, and P. Koziel, “Animal camouflage analysis: Chameleon database,” *Unpublished manuscript*, vol. 2, no. 6, p. 7, 2018.
- [61] L. Li, E. Rigall, J. Dong, and G. Chen, “Mas3k: An open dataset for marine animal segmentation,” in *International Symposium on Benchmarking, Measuring and Optimization*. Springer, 2020, pp. 194–212.
- [62] Z. Fu, R. Chen, Y. Huang, E. Cheng, X. Ding, and K.-K. Ma, “Masnet: A robust deep marine animal segmentation network,” *IEEE Journal of Oceanic Engineering*, 2023.
- [63] M. J. Islam, P. Luo, and J. Sattar, “Simultaneous enhancement and super-resolution of underwater imagery for improved visual perception,” *arXiv preprint arXiv:2002.01155*, 2020.
- [64] P. Drews-Jr, I. d. Souza, I. P. Maurell, E. V. Protas, and S. S. C. Botelho, “Underwater image segmentation in the wild using deep learning,” *Journal of the Brazilian Computer Society*, vol. 27, pp. 1–14, 2021.
- [65] P. Zhang, T. Yan, Y. Liu, and H. Lu, “Fantastic animals and where to find them: Segment any marine animal with dual sam,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024, pp. 2578–2587.
- [66] D. Hendrycks and T. Dietterich, “Benchmarking neural network robustness to common corruptions and perturbations,” *arXiv preprint arXiv:1903.12261*, 2019.
- [67] K. Huang, L. Fang, and C. Tian, “Learning to adapt using test-time images for salient object detection in optical remote sensing images,” *IEEE Transactions on Geoscience and Remote Sensing*, 2024.
- [68] W. Wang, E. Xie, X. Li, D.-P. Fan, K. Song, D. Liang, T. Lu, P. Luo, and L. Shao, “Pvt v2: Improved baselines with pyramid vision transformer,” *Computational Visual Media*, vol. 8, no. 3, pp. 415–424, 2022.
- [69] K. He, X. Zhang, S. Ren, and J. Sun, “Deep residual learning for image recognition,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 770–778.
- [70] M. J. Mirza, J. Micorek, H. Possegger, and H. Bischof, “The norm must go on: Dynamic unsupervised domain adaptation by normalization,” in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2022, pp. 14 765–14 775.
- [71] Y. Mansour, X. Zhong, S. Caglar, and R. Heckel, “Ttt-mim: Test-time training with masked image modeling for denoising distribution shifts,” in *European Conference on Computer Vision*. Springer, 2024, pp. 341–357.