

Exploring Depth Contribution for Camouflaged Object Detection

Mochu Xiang¹ Jing Zhang² Yunqiu Lv¹ Aixuan Li¹ Yiran Zhong³ Yuchao Dai^{1,*}
¹ Northwestern Polytechnical University ² Australian National University ³ SenseTime

Abstract

Camouflaged object detection (COD) aims to segment camouflaged objects hiding in the environment, which is challenging due to the similar appearance of camouflaged objects and their surroundings. Research in biology suggests depth can provide useful object localization cues for camouflaged object discovery. In this paper, we study the depth contribution for camouflaged object detection, where the depth maps are generated with existing monocular depth estimation (MDE) methods. Due to the domain gap between the MDE dataset and our COD dataset, the generated depth maps are not accurate enough to be directly used. We then introduce two solutions to avoid the noisy depth maps from dominating the training process. Firstly, we present an auxiliary depth estimation branch (“ADE”), aiming to regress the depth maps. We find that “ADE” is especially necessary for our “generated depth” scenario. Secondly, we introduce a multi-modal confidence-aware loss function via a generative adversarial network to weigh the contribution of depth for camouflaged object detection. Our extensive experiments on various camouflaged object detection datasets explain that the existing “sensor depth” based RGB-D segmentation techniques work poorly with “generated depth”, and our proposed two solutions work cooperatively, achieving effective depth contribution exploration for camouflaged object detection.

1. Introduction

As a key example of evolution by natural selection, camouflage is widely adopted by the preys in the wild to reduce their possibility of being detected by their predators [6]. Camouflaged object detection (COD) is the technique that segments the whole scope of the camouflaged object. It has wide applications in a variety of fields, such as military (e.g. military camouflaged pattern design [28]), agriculture (e.g. pest identification [41]), medicine (e.g. polyp segmentation [13]) and ecological protection (e.g. wildlife protection [49, 72]). Due to both scientific value and application value, COD deserves well exploration.

Compared with the generic object detection [29, 60] or

segmentation techniques [4, 63], COD is more challenging as the foreground objects usually share a very similar appearance to their surroundings. The visual cues for object identification, e.g. texture, contrast, edge, color, and object size, are vulnerable to attack from the basic camouflage strategies, e.g. background matching and disruptive coloration [57, 68]. Although some recent deep learning based studies [8, 12, 40, 74] have shown favorable performance, the misleading information in camouflage hinders the network from learning the most discriminative features of camouflage. We argue that more visual perceptual knowledge about camouflaged objects can be beneficial.

Research in biology suggests that depth provides 3D geometric information that makes the observer more sensitive to the true boundary and thus enables camouflage less effective [1, 35]. Accordingly, combining depth with RGB images can serve as a new way to solve the challenges in camouflaged object detection. Further, the visual system of many species operates in the real 3D scenes, and a lot of recent works for object segmentation [15, 32, 58, 77] have integrated depth map as another modal on top of RGB images to extract features about the 3D layout of the scene and object shape information. As far as we know, there exist no deep camouflaged object detection models exploring the depth contribution. In this paper, we present the first depth-guided camouflaged object detection network to study the contribution of depth for the COD task.

As there exists no RGB-D camouflaged object detection dataset, we generate depth map of the COD training dataset [12, 40] with existing monocular depth estimation method [59]. With the generated depth, a straightforward solution to achieve RGB-D COD is through multi-modal fusion strategies [15, 54]. However, the conventional monocular depth estimation models are trained on natural images, where there may not exist any camouflaged objects. The domain gap between the monocular depth estimation training dataset and COD training dataset leads to less accurate (or noisy) depth maps as shown in Fig. 1 “Depth”. Directly training with the noisy depth map may lead to an over-fitting model that generalizes poorly on the testing dataset [26].

To effectively use the depth information, we propose two main solutions, namely auxiliary depth estimation branch

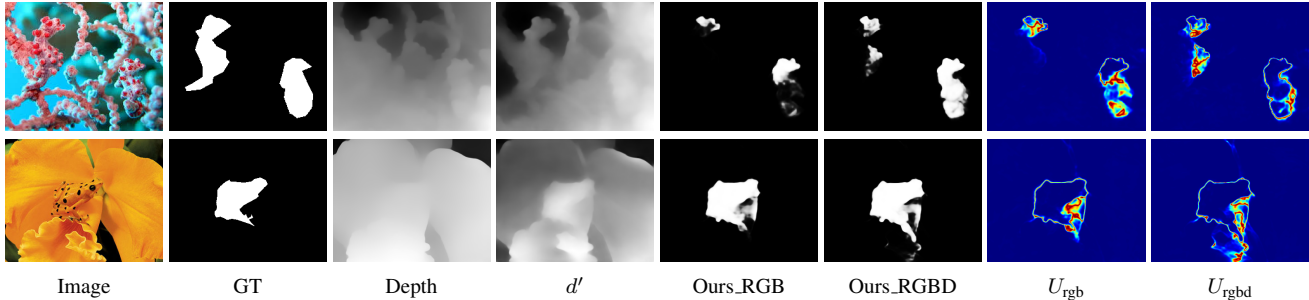


Figure 1. The auxiliary depth estimation branch (“ADE”) (for depth regression (d')) and the generative adversarial network [25] (for uncertainty estimation (U_{rgb} and U_{rgbD})) work cooperatively to produce effective camouflage maps (Ours_RGB and Ours_RGBD).

(“ADE”) and a multi-modal confidence-aware loss function via a generative adversarial network (GAN) [25]. The former aims to regress the depth maps supervised by the generated depth maps from existing MDE method, which is proven especially necessary for our “generated depth” scenario. With the latter, the stochastic attribute of the network makes it possible to estimate the confidence of model prediction, which will serve as the weight for our multi-modal loss function. The basic assumption is that compared with a less confident prediction, a highly confident prediction should contribute more to model updating. We show the estimated uncertainty maps U_{rgb} and U_{rgbD} (the inverse confidence map) of each modal (RGB and RGB-D in this paper) in Fig. 1, which validate the multi-modal confidence-aware learning, leading to better camouflage maps from the RGB-D branch (“Ours_RGBD”).

Our main contributions can be summarized as: 1) We advocate the contribution of depth for camouflaged object detection, and propose the first depth-guided camouflaged object detection network; 2) We introduce an auxiliary depth estimation branch, where the RGB-D feature fusion is achieved through multi-task learning instead of multi-modal fusion within a single task, which is proven especially effective in our “generated depth” scenario; 3) We present multi-modal confidence-aware loss function via generative adversarial network [25] to effectively weigh the contribution of depth for model updating.

2. Related Work

Camouflaged Object Detection Camouflage is a defense mechanism for animals to change their salient signatures and become invisible in the environment [6]. In early time, researchers have developed extensive methods using hand-crafted features, *e.g.* edge, brightness, color, gradient, texture, to detect the camouflaged object [2, 19, 42, 56, 67, 73]. However, these algorithms are far from practical applications because the well-performed camouflage is skilled in breaking the low-level features. Recent research in COD [8, 12, 16, 40, 74, 75, 79] mainly focus on using deep neural

networks to extract high-level semantic features to discriminate the concealed object from the complex scenarios. In addition, [35] claims that cuttlefish makes use of low-level cues to disrupt the perception of visual depth and therefore disguise their 3D shape. As described in [52], seeing the environment in 3D makes the presence of the object easier to be discerned. [1] verifies that the real 3D structure information provided by the depth is advantageous to overcome the disruption caused by edge enhancement. Based on these discussions, we argue that a better understanding of the 3D layout of the scene can be beneficial for more effective camouflaged object detection. In this paper, we aim to explore depth contribution for camouflaged object detection.

Depth-Guided Segmentation For segmentation tasks, many early work has shown that depth information can reduce the ambiguity of RGB features in complex scenarios [7, 18, 27, 53, 61, 62, 65, 70]. In recent years, the deep multi-modal fusion models for segmentation have achieved significant performance improvement. [17] proposes a simple early fusion strategy by concatenating the RGB image and depth at the input layer, forming four-channel 3D-aware data. [46] on the contrary, perform multi-modal fusion at the output layer, leading to a late fusion model. To obtain more accurate and robust segmentation, a variety of models delicately design strategies for cross-level fusion, where they merge the complementary information in the middle level of the network [3, 18, 44, 45, 51, 55, 69].

Confidence-aware Learning Confidence-aware learning (or uncertainty-aware learning) aims to estimate the uncertainty representing the quality of the data (aleatoric uncertainty) or the awareness of true model (epistemic uncertainty) [36]. In this paper, as we mainly focus on modeling the quality of depth for camouflaged object detection, we only discuss the aleatoric uncertainty modeling strategies. [38] designs a network that yields a probabilistic distribution as output in order to capture such uncertainty. [64] employs a teacher-student paradigm to distill the aleatoric uncertainty, where the teacher network generates multiple predicative samples by incorporating aleatoric uncertainty

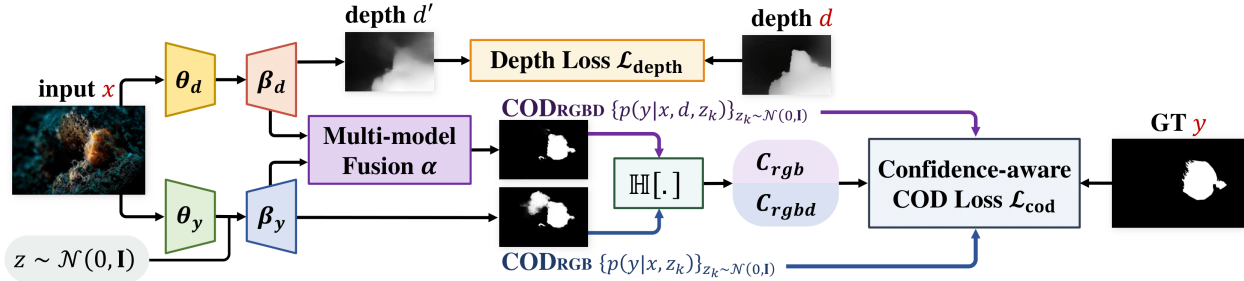


Figure 2. Our depth contribution exploration network takes RGB image x as input to achieve three main tasks: 1) RGB image based camouflaged object detection; 2) auxiliary depth estimation and 3) RGB-D camouflaged object detection. In addition, the proposed multi-modal confidence-aware loss function explicitly evaluates depth contribution with confidence of prediction (C_{rgb} and C_{rgbd}) as indicator, where the confidence is obtained with the proposed probabilistic model via generative adversarial network [25]. We only show the camouflage generator and the discriminator is introduced in Section 3.2.

for the student network to learn. [39] uses an adversarial perturbation technique to generate additional training data for the model to capture the aleatoric uncertainty.

Uniqueness of our solution Different from the existing camouflaged object detection methods [12, 40, 47, 74] that rely only on the RGB images, we introduce the first depth-guided camouflaged object detection framework, with depth contribution estimation solutions to adaptively fuse model predictions from both the RGB branch and the RGB-D branch. Although RGB-D data related segmentation models [15, 54] are widely studied, we claim that their multi-modal learning strategies based on “sensor depth” fails to explore the contribution of “generated depth”. On the contrary, they usually lead to worse performance compared with training only with RGB images (see Table 3). Our solution aims to produce a pixel-wise confidence map of each modal prediction, which is then treated as the weight to achieve multi-modal confidence-aware learning, leading to a better generated-depth contribution exploration.

3. Our Method

The original RGB image based COD training dataset is defined as $D = \{x_i, y_i\}_{i=1}^N$, where x_i and y_i are the input RGB image and the corresponding camouflage map, and i indexes the images, N is the size of D . In this paper, we study the contribution of depth for camouflaged object detection with a confidence-aware network shown in Fig. 2. To start our pipeline, we first generate the depth map d with existing monocular depth estimation method [59], leading to our RGB-D COD training dataset $D = \{x_i, d_i, y_i\}_{i=1}^N$. Then we introduce an auxiliary depth estimation branch and a multi-modal confidence-aware learning loss function to effectively explore depth contribution for COD.

3.1. Initial Depth Generation

As there exists no RGB-D based camouflaged object detection dataset, we generate pseudo depth maps from existing monocular depth estimation methods. We tried three state-of-the-art monocular depth estimation methods (MiDaS [59], Monodepth2 [24] and FrozenPeople [43]) to generate depth maps for our training dataset.

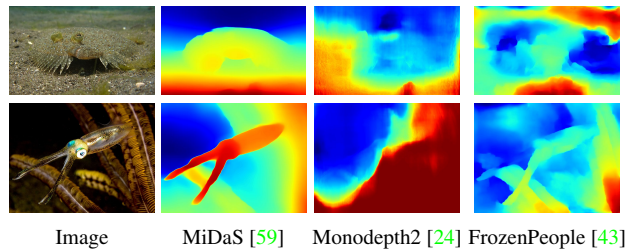


Figure 3. Camouflage images and the generated depth maps. Three methods are used to generate depth maps. Overall, MiDaS [59] has the best cross dataset performance and the strongest generalization ability. MonoDepth2 [24] is trained only on KITTI dataset, which performs the worst in our scenario.

Trained on 10 different RGB-D datasets, MiDaS [59] provides robust results across diverse scenes. Targeting autonomous driving, Monodepth2 [24] has a good performance on the KITTI benchmark under the self-supervision, but its ability of domain adaption is to be investigated. Since there are a large number of images with humans as camouflaged objects in our camouflage training dataset, we also tied FrozenPeople [43] with the single image configuration to generate depth for the characters in the scene.

Both MiDaS and Monodepth2 generate depth maps in the form of disparity, or inverse-depth, FrozenPeople generates depth values directly. We show the generated depth maps for our camouflaged object detection dataset in Fig. 3.

Due to the visually better performance of MiDaS depth maps¹, we adopt MiDaS depth maps in our experiments.

3.2. Depth Contribution Exploration

Due to the domain gap, the generated depth map from monocular depth estimation methods may not be very accurate as shown in Fig. 3. Directly training with the less-accurate depth map cannot improve the model performance as the network will over-fit on the less-accurate depth map, leading to poor generalization ability (see Table 3). Instead of directly using the depth as input, we design an auxiliary depth estimation branch, where the RGB-D feature fusion is achieved through multi-task learning (camouflaged object detection and monocular depth estimation) instead of multi-modal fusion within a single task. Four main modules are included in our framework: 1) a RGB COD module to generate camouflage map from the RGB image; 2) a auxiliary depth estimation module to regress the depth map of image x with d as supervision; 3) a multi-modal fusion module to aggregate the intermediate features from the above two modules; 4) a multi-modal confidence-aware learning module via generative adversarial network (GAN) [25] to adaptively weight the contribution of depth for model updating.

RGB camouflaged object detection: We build the RGB camouflaged object detection network upon ResNet50 [30] backbone, which maps the input image x to four stages features $f_{\theta_y}(x) = \{s_k\}_{k=1}^4$ of channel size 256, 512, 1024 and 2048 respectively, where θ_y is the parameter set of the backbone network. To obtain enlarged receptive field with less memory consumption, we feed each s_k to a multi-scale dilated convolution block [4] to obtain new backbone features $f_{\theta_y}(x) = \{s'_k\}_{k=1}^4$ of channel size $C = 32$. We further adopt decoder from [59] for high/low level feature aggregation. Specifically, we define prediction from the deterministic RGB branch as $f_{\beta_y}(f_{\theta_y}(x))$, where β_y is the parameter set of the decoder (excluding the multi-scale dilated convolutional blocks). As we intend to model confidence of each modal (RGB and RGB-D), we change the RGB camouflaged object detection network to a probabilistic network via generative adversarial network (GAN) [25]. In this way, we re-define model prediction of our stochastic prediction network as $f_{\beta_y}(f_{\theta_y}(x), z)$, where z is the latent variable. To fuse z with the backbone feature $f_{\theta_y}(x)$, we first tile z to the same spatial size of s'_4 , and then concatenate the tiled z with s'_4 , and feed it to another 3×3 convolutional layers of channel size $C = 32$. We use this new feature instead of s'_4 in the decoder parametered with β_y . Following the standard practice of GAN, z is assumed to follow the standard normal distribution: $z \sim \mathcal{N}(0, \mathbf{I})$.

Auxiliary depth estimation: Given the generated depth map d_i from [59], we design auxiliary depth estimation branch to extract depth related features from the input image

x . Specifically, we have a separate encoder with ResNet50 backbone for the auxiliary depth estimation branch with parameter set θ_d . Similarly, multi-scale dilated convolutional blocks are used for larger receptive field to generate the depth backbone features $f_{\theta_d}(x) = \{s'_k\}_{k=1}^4$ of the same channel size $C = 32$, and we adopt the same decoder structure in the RGB COD module to produce a one-channel depth map $d' = f_{\beta_d}(f_{\theta_d}(x))$.

Multi-modal Fusion: Given intermediate backbone features $f_{\theta_y}(x, z)$ from the RGB COD branch and the auxiliary depth estimation branch $f_{\theta_d}(x) = \{s'_k\}_{k=1}^4$, the goal of the ‘‘Multi-modal Fusion’’ branch is to effectively fuse features for COD and features for MDE to output a RGB-D camouflage map $f_\alpha(f_\theta(x), d, z^d)$, where $\theta = \{\theta_d, \theta_y\}$, α is the parameter set of the multi-modal fusion module, and z^d is the latent variable for the RGB-D camouflage branch, and we also have $z^d \sim \mathcal{N}(0, \mathbf{I})$. We concatenate each level of intermediate feature from RGB branch and depth branch, then we feed it to a 3×3 convolutions layer of channel size $C = 32$, and defined the fused feature as $f_\theta^{rgb-d} = f_\theta(x, d, z, z^d)$. We design a COD decoder within the multi-modal fusion module, which share the same structure as [59]. We define the RGB-D COD prediction as f_α for simplicity.

Multi-modal Confidence-aware Learning: Recall that the RGB COD model produces RGB image based camouflage map $p(y|x, z) = \mathcal{N}(f_{\beta_y}(f_{\theta_y}(x), z), \epsilon_{rgb})$, where $\epsilon_{rgb} \sim \mathcal{N}(0, \sigma_{rgb}^2)$ and σ_{rgb}^2 is the variance of the labeling noise, representing the inherent noise level from a generative model’s perspective. Similarly, with the ‘‘Multi-modal Fusion’’ module, we obtain the RGB-D camouflage map $p(y|x, d, z^d) = \mathcal{N}(f_\alpha, \epsilon_{rgb-d})$, where $\epsilon_{rgb-d} \sim \mathcal{N}(0, \sigma_{rgb-d}^2)$ also represents the inherent labeling noise. Directly training the network with above two camouflage predictions and depth regression leads to inferior performance (see ‘‘A_D’’ in Table 2), as there exists no explicit constraint for depth quality evaluation. Although the proposed auxiliary depth estimation module achieves better depth exploration than other multi-modal learning strategies widely used in the existing RGB-D segmentation models [15, 54] (compare ‘‘A_D’’ in Table 2 with multi-modal learning strategies in Table 3), there exists no explicit solution to weight the contribution of depth. Towards explicit depth contribution exploration, we first estimate confidence of each modal predictions, and use it as the weight within our multi-modal confidence-aware loss function.

Prediction confidence estimation: Given stochastic prediction $p(y|x, z)$ from the RGB COD branch, and $p(y|x, d, z^d)$ from the RGB-D branch, we first perform multiple iterations of sampling in the latent space $z \sim \mathcal{N}(0, \mathbf{I})$ and $z^d \sim \mathcal{N}(0, \mathbf{I})$, and obtain a sequence of predictions $\{p(y|x, z_k)\}_{z_k}$ and $\{p(y|x, d, z_k^d)\}_{z_k^d}$, where z_k and z_k^d represent the random sample of the latent variable for the RGB

¹The ‘‘depth map’’ represents ‘‘inverse-depth map’’ in the following.

COD branch and RGB-D COD branch respectively. For an ensemble based framework, the predictive uncertainty [36] (the inverse confidence) captures the overall uncertainty of model prediction, which is defined as entropy of the mean prediction ($\mathbb{H}[\cdot]$ in Fig. 2). In this way, we obtain the confidence map of each modal prediction as:

$$\begin{aligned} C_{\text{rgb}} &= 1 - \mathbb{H}[\mathbb{E}_{z \sim \mathcal{N}(0,1)} p(y|x, z)], \\ C_{\text{rgb-d}} &= 1 - \mathbb{H}[\mathbb{E}_{z^d \sim \mathcal{N}(0,1)} p(y|x, d, z^d)]. \end{aligned} \quad (1)$$

We show the inverse confidence maps (the uncertainty map) U_{rgb} and $U_{\text{rgb-d}}$ for RGB and RGB-D modal in Fig. 1.

Discriminator: As a generative adversarial network [25], we design an extra fully convolutional discriminator of the same structure as in [33], and define it as $g_\gamma(\cdot)$, which is used to distinguish model predictions and the ground truth annotations. Specifically, the discriminator takes the concatenation of model prediction (or ground truth) and image as input, and identify it as “fake” (or real) with an all zero feature map $\mathbf{0}$ (or all one feature map $\mathbf{1}$).

Train the model with multi-modal confidence-aware loss function: Given the COD outputs $p(y|x, z)$ and $p(y|x, d, z^d)$ from the RGB and RGB-D branch, and the corresponding confidence maps C_{rgb} and $C_{\text{rgb-d}}$, we introduce a multi-modal confidence-aware loss function to explicitly weight the depth contribution. The basic assumption is that if the model is confident about the prediction from one modal, then that modal should contribute more to the overall loss function, and vice versa. Inspired by [22], our multi-modal confidence-aware loss function is defined:

$$\mathcal{L}_{\text{cod}} = \omega_{\text{rgb}} \mathcal{L}_{\text{rgb}} + \omega_{\text{rgb-d}} \mathcal{L}_{\text{rgb-d}}, \quad (2)$$

where ω_{rgb} and $\omega_{\text{rgb-d}}$ are the pixel-wise confidence-aware weights, defined as: $\omega_{\text{rgb}} = (C_{\text{rgb}})/(C_{\text{rgb}} + C_{\text{rgb-d}})$ and $\omega_{\text{rgb-d}} = (C_{\text{rgb-d}})/(C_{\text{rgb}} + C_{\text{rgb-d}})$ respectively. $\mathcal{L}_{\text{rgb}} = \mathcal{L}_{\text{rgb}}(f_{\beta_y}(f_{\theta_y}(x), z), y)$ and $\mathcal{L}_{\text{rgb-d}} = \mathcal{L}_{\text{rgb-d}}(f_\alpha, y)$ are the task related loss functions, which are structure-aware loss functions [71] in this paper.

In addition to Eq. 2, as a generative adversarial network with auxiliary depth estimation branch, we add the extra adversarial loss \mathcal{L}_{adv} and depth regression loss $\mathcal{L}_{\text{depth}}$ to \mathcal{L}_{cod} , and define it as our generator loss \mathcal{L}_{gen} :

$$\mathcal{L}_{\text{gen}} = \mathcal{L}_{\text{cod}} + \mathcal{L}_{\text{depth}} + \lambda \mathcal{L}_{\text{adv}}, \quad (3)$$

where λ is the hyper-parameter to weight the contribution of the adversarial loss, and empirically we set $\lambda = 0.1$. The adversarial loss \mathcal{L}_{adv} is defined as:

$$\begin{aligned} \mathcal{L}_{\text{adv}} &= \mathcal{L}_{\text{ce}}(g_\gamma(x, f_{\beta_y}(f_{\theta_y}(x), z)), \mathbf{1}) \\ &+ \mathcal{L}_{\text{ce}}(g_\gamma(x, f_\beta(f_\theta(x), d, z)), \mathbf{1}), \end{aligned} \quad (4)$$

with \mathcal{L}_{ce} as the binary cross-entropy loss. The extra adversarial loss aims to fool the discriminator to predict “real”

with generator camouflage prediction as input. The depth regression loss $\mathcal{L}_{\text{depth}}$ is defined as the weighted sum of the point-wise L_1 loss and the SSIM loss [23]:

$$\mathcal{L}_{\text{depth}} = (1 - \lambda) \frac{1}{n} \sum_{i=1}^n |d - d'| + \lambda \frac{1 - \text{SSIM}(d, d')}{2}, \quad (5)$$

and empirically we set $\lambda = 0.85$ [23].

Then the discriminator is updated via:

$$\begin{aligned} \mathcal{L}_{\text{dis}} &= \mathcal{L}_{\text{ce}}(\mathcal{L}_{\text{ce}}(g_\gamma(x, y), \mathbf{1}) \\ &+ \mathcal{L}_{\text{ce}}(\mathcal{L}_{\text{ce}}(g_\gamma(x, f_{\beta_y}(f_{\theta_y}(x), z)), \mathbf{0})) \\ &+ \mathcal{L}_{\text{ce}}(g_\gamma(x, f_\alpha), \mathbf{0})), \end{aligned} \quad (6)$$

where the discriminator aims to correctly distinguish model predictions from the ground truth, and this is what [25] defined as the “minimax game”. We show the complete learning pipeline of our method in Algorithm 1.

Algorithm 1 Depth Contribution Exploration for COD

Input:

(1) Training dataset $D = \{x_i, d_i, y_i\}_{i=1}^N$, where the depth map d_i is pre-computed from [59].

(2) The maximal number of learning epochs E .

Output: Model parameter $\theta = \{\theta_d, \theta_y\}$, $\beta = \{\beta_d, \beta_y\}$ and α for the camouflage generator, and γ for the discriminator;

- 1: **for** $t \leftarrow 1$ to E **do**
 - 2: Generate camouflage map $f_{\beta_y}(f_{\theta_y}(x), z)$ and f_α from the RGB and RGB-D COD branches.
 - 3: Obtain monocular depth $f_{\beta_d}(f_{\theta_d}(x_i))$ of the input x_i .
 - 4: Perform $T = 5$ iterations of sampling from the latent space: $z \sim \mathcal{N}(0, \mathbf{I})$, $z^d \sim \mathcal{N}(0, \mathbf{I})$, and obtain $\{p(y|x, z_k)\}_{z_k}$ and $\{p(y|x, d, z_k^d)\}_{z_k^d}$, where z_k and z_k^d represent the random samples of the latent variables.
 - 5: Compute confidence of each model C_{rgb} and $C_{\text{rgb-d}}$.
 - 6: Obtain the pixel-wise confidence weights $\omega_{\text{rgb}} = (C_{\text{rgb}})/(C_{\text{rgb}} + C_{\text{rgb-d}})$ and $\omega_{\text{rgb-d}} = (C_{\text{rgb-d}})/(C_{\text{rgb}} + C_{\text{rgb-d}})$;
 - 7: Define the multi-modal confidence-aware loss function $\mathcal{L}_{\text{cod}} = \omega_{\text{rgb}} \mathcal{L}_{\text{rgb}} + \omega_{\text{rgb-d}} \mathcal{L}_{\text{rgb-d}}$, and obtain generator loss \mathcal{L}_{gen} .
 - 8: Update $\theta = \{\theta_d, \theta_y\}$, $\beta = \{\beta_d, \beta_y\}$ and α via \mathcal{L}_{gen} .
 - 9: Compute discriminator loss \mathcal{L}_{dis} , and update γ with \mathcal{L}_{dis} .
 - 10: **end for**
-

4. Experimental Results

4.1. Setup

Datasets We train our method with benchmark COD training dataset, which is the combination of 3,040 images from COD10K training dataset [12] and 1,000 images from CAMO training dataset [40]. We then test model performance on four benchmark camouflage testing datasets, including CAMO [40] of size 250, CHAMELEON of size 76, COD10K testing set of size 2,026, and the newly released NC4K dataset [47] of size 4,121.

Table 1. Performance comparison with benchmark COD models. ‘‘BkB’’: the backbone model, and ‘‘R50’’ is the ResNet50 backbone [31], ‘‘R2_50’’ is the Res2Net50 backbone [21]. ‘‘Size’’: size of the training and testing images.

Method	BkB	Year	Size	CAMO [40]				CHAMELEON [66]				COD10K [12]				NC4K [47]			
				$S_\alpha \uparrow$	$F_\beta \uparrow$	$E_\xi \uparrow$	$\mathcal{M} \downarrow$	$S_\alpha \uparrow$	$F_\beta \uparrow$	$E_\xi \uparrow$	$\mathcal{M} \downarrow$	$S_\alpha \uparrow$	$F_\beta \uparrow$	$E_\xi \uparrow$	$\mathcal{M} \downarrow$	$S_\alpha \uparrow$	$F_\beta \uparrow$	$E_\xi \uparrow$	$\mathcal{M} \downarrow$
RGB COD Models																			
SINet [12]	R50	2020	352	.745	.702	.804	.092	.872	.827	.936	.034	.776	.679	.864	.043	.810	.772	.873	.057
SINet-V2 [11]	R2_50	2021	352	.820	.782	.882	.070	.888	.835	.942	.030	.815	.718	.887	.037	.847	.805	.903	.048
LSR [47]	R50	2021	352	.793	.725	.826	.085	.893	.839	.938	.033	.793	.685	.868	.041	.839	.779	.883	.053
MGL [76]	R50	2021	473	.775	.726	.812	.088	.893	.834	.918	.030	.814	.711	.852	.035	.833	.782	.867	.052
PFNet [48]	R50	2021	416	.782	.744	.840	.085	.882	.826	.922	.033	.800	.700	.875	.040	.829	.782	.886	.053
UJTR [75]	R50	2021	473	.785	.686	.859	.086	.888	.796	.918	.031	.818	.667	.850	.035	.839	.786	.873	.052
RGB-D COD Models																			
UCNet [77]	R50	2020	352	.729	.672	.785	.101	.869	.823	.924	.039	.738	.611	.825	.052	.784	.728	.849	.066
BBSNet [15]	R50	2020	352	.776	.689	.786	.093	.864	.768	.872	.049	.782	.633	.836	.050	.825	.745	.859	.062
JL-DCF [20]	R50	2020	352	.772	.670	.777	.102	.857	.734	.867	.052	.749	.581	.789	.053	.788	.713	.819	.072
SSF [78]	VGG16	2020	352	.748	.643	.782	.116	.866	.788	.904	.045	.729	.593	.778	.055	.770	.689	.827	.069
Ours	R50	2021	352	.794	.759	.853	.077	.885	.836	.941	.032	.801	.705	.882	.037	.837	.798	.899	.049
Ours	R50	2021	416	.798	.762	.860	.077	.889	.842	.947	.029	.807	.712	.881	.037	.842	.802	.901	.047
Ours	R50	2021	480	.803	.770	.858	.078	.891	.846	.939	.028	.820	.737	.894	.033	.845	.809	.903	.046
Ours	R2_50	2021	352	.819	.798	.881	.069	.895	.856	.951	.027	.829	.751	.903	.032	.855	.825	.910	.042

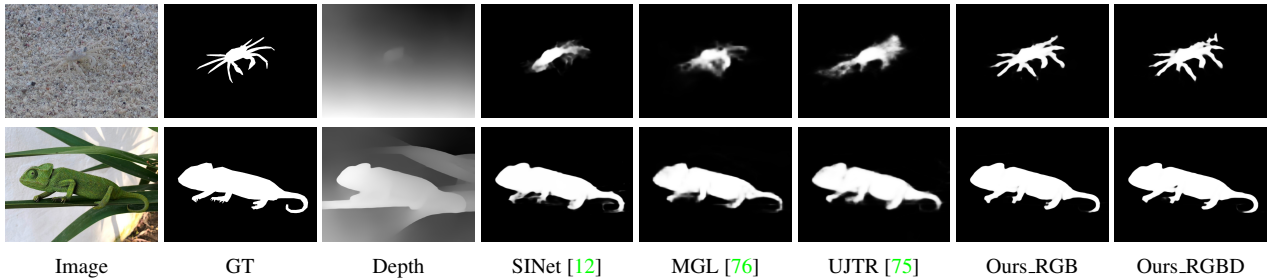


Figure 4. Visual comparison of our predictions with the benchmark techniques.

Evaluation Metrics We adopt four widely used metrics for performance evaluation, including 1) Mean Absolute Error (\mathcal{M}); 2) Mean F-measure (F_β); 3) Mean E-measure (E_ξ) [10]; 4) S-measure (S_α) [9]. Details of those evaluation metrics are introduced in the supplementary materials.

Implementation Details We train our model with Pytorch, where ResNet-50 is chosen as backbone, which is initialized with weights trained on ImageNet, and other newly added layers are initialized by default. We resize all the images and ground truth to 352×352 for both training and testing. The maximum epoch is 50. The initial learning rates is 2.5×10^{-5} . The whole training takes 17 hours with batch size 6 on one NVIDIA GTX 2080Ti GPUs.

4.2. Performance comparison

Quantitative comparison: We compare our results with the benchmark models in Table 1. As both training image size and backbone networks are important factors for model performance, for fair comparison, we train with different backbones and different training image sizes following the settings of existing solutions. Due to extra depth informa-

tion, we also generate benchmark RGB-D COD models by training RGB-D saliency detection models with our RGB-D COD dataset. Specifically, we re-train UCNet [77], BBSNet [15], JL-DCF [20] and SSF [78]) with our RGB-D camouflaged object detection dataset. The consistent better performance the RGB COD benchmark models compared with the re-trained RGB-D COD models explains that the existing sensor depth based RGB-D fusion strategies work poorly within our generated depth scenario. The main reason lies in the low quality depth map due to domain gap of monocular depth estimation dataset and our camouflaged object detection dataset. Table 1 shows that, with the proposed depth contribution exploration techniques, we obtain consistently improved performance compared with both the RGB COD models and the re-trained RGB-D COD models.

Qualitative comparison: In Fig. 4, we show predictions of benchmark techniques and our method, where we show prediction from both the RGB COD branch (‘‘Ours_RGB’’) and the RGB-D COD branch (‘‘Ours_RGBD’’). The better camouflage maps from both branches further explain effectiveness of our solutions.

Table 2. Performance of ablation study models.

Method	CAMO [40]		CHAMELEON [66]		COD10K [12]		NC4K [47]	
	$F_\beta \uparrow$	$\mathcal{M} \downarrow$	$F_\beta \uparrow$	$\mathcal{M} \downarrow$	$F_\beta \uparrow$	$\mathcal{M} \downarrow$	$F_\beta \uparrow$	$\mathcal{M} \downarrow$
Base	.751	.079	.815	.037	.696	.039	.789	.050
ADE	.735	.084	.826	.033	.699	.039	.794	.050
A_D	.746	.079	.835	.032	.702	.038	.793	.049
Ours	.759	.077	.836	.032	.705	.037	.798	.049

Running time comparison: The parameter number of our camouflage generator is 101M, which is larger than some RGB-D saliency detection models, *e.g.* UCNet [77] (62M), BBSNet [15] (49M), while smaller than JL-DCF [20] (144M). The extra parameters mainly come from our separate encoder with ResNet50 backbone for auxiliary depth estimation. At test time, it cost 0.2 second to process each image of size 352×352 , which is comparable with benchmark techniques.

4.3. Ablation Study

We thoroughly analyse the proposed solutions with extra experiments, and show results in Table 2. Note that, all the following experiments are trained with ResNet50 backbone with training/testing image size of 352.

Firstly, we train an RGB COD model with only the RGB image, and the performance is shown as “Base”. The network structure of “Base” is the same as our auxiliary depth estimation network, except that we use structure-aware loss function [71] for “Base”. Then, we introduce our auxiliary depth estimation branch to “Base”, and obtain “ADE”, where there exists no RGB and depth feature interaction. We further add the “Multi-modal Fusion” module to “ADE” to produce extra RGB-D camouflage map, and show performance as “A_D”. Note that, until now, all the discussed models are deterministic. We intend to explore the depth contribution with a generative adversarial network, and we then add latent variable z to “A_D” with an extra discriminator to obtain a generative adversarial network based RGB-D COD network, which is our final model, shown as “Ours”.

We observe competing performance of “Base” compared with the state-of-the-art camouflaged object detection networks, indicating effectiveness of the simple base model. As there exists no interaction of RGB feature and depth feature, the performance of “ADE” is similar to “Base”. At the same time, as the weight of each task is quite important for a multi-task learning framework (within “ADE”, we achieve simultaneous camouflaged object detection and monocular depth estimation), the performance of “ADE” can be further tuned with more effective task-relationship models [37]. Build upon “ADE” with extra “Multi-modal Fusion” block, we notice improved performance of “A_D” compared with “ADE”, which explains effectiveness of our “Multi-modal Fusion” block. Then, we add extra latent

variable z and a fully convolutional discriminator to “A_D”, and perform multi-modal confidence-aware learning. The improved performance of “Ours” explains effectiveness of our depth contribution exploration solutions.

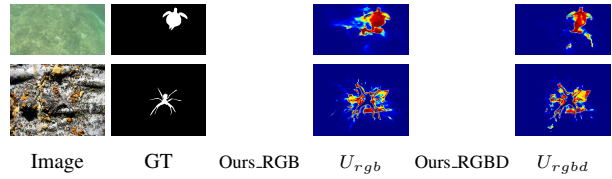


Figure 5. Predictive uncertainties deduced from predictions with different input sizes.

4.4. Training/Testing size and Backbone Analysis

Training/testing image size is usually an important factor for dense prediction tasks, and we observe it’s more important for camouflaged object detection, where higher resolution training images can always lead to better performance. In Table 1, we summarize the training/testing image sizes of existing methods, and train our model with different image sizes. We observe consistently improved performance with larger image sizes. The main reason is that lower resolution can be seen as adversarial attack for camouflaged object detection. Camouflaged objects usually share a similar appearance as the background, and it’s more difficult to discover camouflaged objects in a small size image than in a large size image. In this way, reducing the images sizes of a training dataset can be seen as improving the difficulty of the task, and training models with too many hard samples will limit its generalization ability.

We further notice that testing with different sizes also affects model predictions. Models perform best when we have the same training and testing data sizes. By associating predictions of a single input with a series of sizes, we can deduce the predictive uncertainty that lies intrinsically in the model and the task themselves (see Fig. 5).

Further, we notice backbone of SInet-V2 [11] is Res2Net50 [21]. For fair comparison, and also to investigate model performance w.r.t. backbone networks, we train extra model with Res2Net50 backbone. The consistent better performance of our model with Res2Net50 [21] backbone compared with SInet-V2 [11] validate our solution.

4.5. Multi-Modal Learning Discussion

We investigate the “sensor depth” based multi-modal learning strategies on our “generated depth” scenario.

Sensor depth based multi-modal learning framework for COD: With depth from monocular depth estimation method, following conventional multi-modal learning pipeline, we can train directly the RGB-D COD task.

Table 3. Performance of RGB-D COD with existing sensor depth based multi-modal learning framework.

Method	CAMO [40]		CHAMELEON [66]		COD10K [12]		NC4K [47]	
	$F_\beta \uparrow$	$\mathcal{M} \downarrow$	$F_\beta \uparrow$	$\mathcal{M} \downarrow$	$F_\beta \uparrow$	$\mathcal{M} \downarrow$	$F_\beta \uparrow$	$\mathcal{M} \downarrow$
Base	.751	.079	.815	.037	.696	.039	.789	.050
Early	.718	.092	.816	.037	.677	.042	.770	.056
Cross	.680	.102	.826	.035	.694	.040	.772	.060
Late	.689	.100	.808	.043	.681	.044	.752	.068

Table 4. Multi-modal learning performance for sensor depth.

Method	NJU2K [34]		SSB [50]		DES [5]		NLPR [53]		SIP [14]	
	$F_\beta \uparrow$	$\mathcal{M} \downarrow$	$F_\beta \uparrow$	$\mathcal{M} \downarrow$	$F_\beta \uparrow$	$\mathcal{M} \downarrow$	$F_\beta \uparrow$	$\mathcal{M} \downarrow$	$F_\beta \uparrow$	$\mathcal{M} \downarrow$
Base	.898	.036	.884	.037	.904	.022	.897	.024	.866	.049
Early	.904	.036	.887	.037	.924	.017	.904	.022	.882	.047
Cross	.905	.036	.890	.036	.928	.016	.904	.024	.886	.047
Late	.897	.040	.863	.053	.906	.020	.897	.024	.886	.048

Specifically, we introduce three different settings, namely early fusion model (“Early”), cross-level fusion model (“Cross”) and late fusion model (“Late”), and show model performance in Table 3. For “Early”, we concatenate x and d at the input layer, and feed it to a 3×3 convolutional layers to obtain a feature map of channel size 3, which is then fed to the RGB COD network to produce camouflage prediction (without the latent variable z as input). For “Cross”, we re-purpose the auxiliary depth estimation to depth based COD, and keep the “Multi-modal Fusion” module. In this way, we have camouflage prediction from the depth branch, the RGB branch and also the RGB-D branch, and prediction of the RGB-D branch is used for performance evaluation. For “Late”, we have two copies of RGB COD network that take RGB and depth as input to produce two copies of predictions. We concatenate the two predictions and feed it to a 3×3 convolutional layers to produce our final prediction. Table 3 shows that the existing “sensor depth” based multi-modal learning framework works poorly on our “generated depth” scenario, leading to inferior performance compared to RGB image based “Base” model.

Effectiveness of the multi-modal learning framework in Table 3 with “sensor depth”: Similar to camouflaged object detection, salient object detection [71] is also defined as class-agnostic binary segmentation task. We then analyse how the sensor depth based multi-modal learning framework performs for RGB-D salient object detection [15], and show performance in Table 4, where the depth are the raw “sensor depth” instead of the “generated depth”. Table 4 shows that the multi-modal learning framework in Table 4 performs well with sensor depth scenario, and the different conclusion with “sensor depth” and “generated depth” further explain the value of our work.

4.6. When is the Generated Depth More Helpful?

RGB image captures appearance information of a scenario, while depth captures geometric information, indicating the distance of objects to the camera. It’s usually claimed that extra data from different modal can lead to better understanding of the same scene. In this paper, we investigate how the generated depth can improve camouflaged object detection. Although the overall performance indicates that depth can be helpful if we use it properly, we still notice samples where the final camouflage map from the RGB-D branch is inferior to that from the RGB branch as shown in Fig. 6.

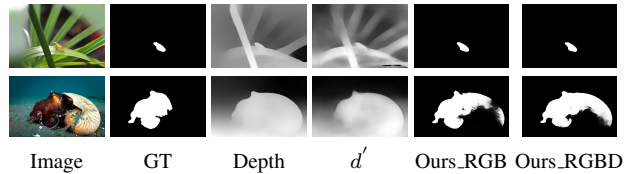


Figure 6. Scenarios when the RGB COD branch outperforms RGB-D COD branch.

Fig. 6 shows that depth can be a double-edged sword. In good cases, both RGB and RGB-D COD branches provide decent segmentation results, and we notice that depth predictions (d') in these scenarios tend to highlight the camouflaged objects better compared with the original generated depths (d). Further, when we can observe the camouflaged objects from the generated depth (d), the depth d can always benefit the RGB-D COD. In hard cases, especially when camouflaged objects dominate the view and have large depth variation, RGB-D COD prediction can be misled by the depth estimation task.

5. Conclusions

We tackle the problem of camouflaged object detection (COD) with the aid of depth, which is generated with the existing monocular depth estimation method [59]. To effectively explore depth contribution, we introduce an auxiliary depth estimation module, and a probabilistic learning network via a generative adversarial network, where the modal related confidence is estimated and serves as a weight for multi-modal confidence-aware learning. With extensive experimental results, we conclude that the existing “sensor depth” based multi-modal learning pipelines perform poorly with the “generated depth” setting, leading to inferior performance compared to the base model with only RGB image used. Instead, we suggest performing auxiliary depth estimation with an effective depth confidence estimation module to prevent the model from being dominated by the noisy depth, which has been proven effective in our “generated depth” based RGB-D COD task.

References

- [1] Wendy J Adams, Erich W Graf, and Matt Anderson. Disruptive coloration and binocular disparity: breaking camouflage. *Proceedings of the Royal Society B*, 286(1896):20182045, 2019. [1](#), [2](#)
- [2] Nagappa U Bhajantri and P Nagabhushan. Camouflage defect identification: a novel approach. In *International Conference on Information Technology*, pages 145–148, 2006. [2](#)
- [3] Hao Chen and Youfu Li. Progressively complementarity-aware fusion network for RGB-D salient object detection. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 3051–3060, 2018. [2](#)
- [4] Liang-Chieh Chen, George Papandreou, Iasonas Kokkinos, Kevin Murphy, and Alan L Yuille. Deeplab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected CRFs. *IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)*, 40(4):834–848, 2017. [1](#), [4](#)
- [5] Yupeng Cheng, Huazhu Fu, Xingxing Wei, Jiangjian Xiao, and Xiaochun Cao. Depth enhanced saliency detection method. In *Proceedings of international conference on internet multimedia computing and service*, pages 23–27, 2014. [8](#)
- [6] Anthony C Copeland and Mohan M Trivedi. Models and metrics for signature strength evaluation of camouflaged targets. In *Algorithms for Synthetic Aperture Radar Imagery IV*, volume 3070, pages 194–199, 1997. [1](#), [2](#)
- [7] Zhuo Deng, Sinisa Todorovic, and Longin Jan Latecki. Semantic segmentation of RGBD images with mutex constraints. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1733–1741, 2015. [2](#)
- [8] Bo Dong, Mingchen Zhuge, Yongxiong Wang, Hongbo Bi, and Geng Chen. Towards accurate camouflaged object detection with mixture convolution and interactive fusion. *arXiv preprint arXiv:2101.05687*, 2021. [1](#), [2](#)
- [9] Deng-Ping Fan, Ming-Ming Cheng, Yun Liu, Tao Li, and Ali Borji. Structure-measure: A new way to evaluate foreground maps. In *IEEE International Conference on Computer Vision (ICCV)*, pages 4548–4557, 2017. [6](#)
- [10] Deng-Ping Fan, Cheng Gong, Yang Cao, Bo Ren, Ming-Ming Cheng, and Ali Borji. Enhanced-alignment measure for binary foreground map evaluation. In *International Joint Conference on Artificial Intelligence (IJCAI)*, pages 698–704, 2018. [6](#)
- [11] Deng-Ping Fan, Ge-Peng Ji, Ming-Ming Cheng, and Ling Shao. Concealed object detection. *IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)*, 2021. [6](#), [7](#)
- [12] Deng-Ping Fan, Ge-Peng Ji, Guolei Sun, Ming-Ming Cheng, Jianbing Shen, and Ling Shao. Camouflaged object detection. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2777–2787, 2020. [1](#), [2](#), [3](#), [5](#), [6](#), [7](#), [8](#)
- [13] Deng-Ping Fan, Ge-Peng Ji, Tao Zhou, Geng Chen, Huazhu Fu, Jianbing Shen, and Ling Shao. PraNet: Parallel reverse attention network for polyp segmentation. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 263–273, 2020. [1](#)
- [14] Deng-Ping Fan, Zheng Lin, Zhao Zhang, Menglong Zhu, and Ming-Ming Cheng. Rethinking RGB-D salient object detection: models, datasets, and large-scale benchmarks. *IEEE Transactions on Neural Networks and Learning Systems (TNNLS)*, 2020. [8](#)
- [15] Deng-Ping Fan, Yingjie Zhai, Ali Borji, Jufeng Yang, and Ling Shao. BBS-Net: RGB-D salient object detection with a bifurcated backbone strategy network. In *European Conference on Computer Vision (ECCV)*, pages 275–292, 2020. [1](#), [3](#), [4](#), [6](#), [7](#), [8](#)
- [16] Zheng Fang, Xiongwei Zhang, Xiaotong Deng, Tiejiong Cao, and Changyan Zheng. Camouflage people detection via strong semantic dilation network. In *Proceedings of the ACM Turing Celebration Conference-China*, pages 1–7, 2019. [2](#)
- [17] Clément Farabet, Camille Couprie, Laurent Najman, and Yann LeCun. Learning hierarchical features for scene labeling. *IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)*, 35(8):1915–1929, 2012. [2](#)
- [18] David Feng, Nick Barnes, Shaodi You, and Chris McCarthy. Local background enclosure for RGB-D salient object detection. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2343–2350, 2016. [2](#)
- [19] Xue Feng, Cui Guoying, Hong Richang, and Gu Jing. Camouflage texture evaluation using a saliency map. *Multimedia Systems*, 21(2):169–175, 2015. [2](#)
- [20] Keren Fu, Deng-Ping Fan, Ge-Peng Ji, and Qijun Zhao. JL-DCF: Joint learning and densely-cooperative fusion framework for RGB-D salient object detection. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 3052–3062, 2020. [6](#), [7](#)
- [21] S. Gao, M. Cheng, K. Zhao, X. Zhang, M. Yang, and P. Torr. Res2net: A new multi-scale backbone architecture. *IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)*, 43(02):652–662, 2021. [6](#), [7](#)
- [22] Anjith George and Sebastien Marcel. Cross modal focal loss for RGBD face anti-spoofing. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 7882–7891, 2021. [5](#)
- [23] Clément Godard, Oisín Mac Aodha, and Gabriel J. Brostow. Unsupervised monocular depth estimation with left-right consistency. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017. [5](#)
- [24] Clément Godard, Oisín Mac Aodha, Michael Firman, and Gabriel J. Brostow. Digging into self-supervised monocular depth prediction. In *IEEE International Conference on Computer Vision (ICCV)*, pages 3828–3838, 2019. [3](#)
- [25] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In *Conference on Neural Information Processing Systems (NeurIPS)*, pages 2672–2680, 2014. [2](#), [3](#), [4](#), [5](#)
- [26] Chuan Guo, Geoff Pleiss, Yu Sun, and Kilian Q. Weinberger. On calibration of modern neural networks. In *International Conference on Machine Learning (ICML)*, pages 1321–1330, 2017. [1](#)
- [27] Saurabh Gupta, Pablo Arbelaez, and Jitendra Malik. Perceptual organization and recognition of indoor scenes from

- RGB-D images. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 564–571, 2013. 2
- [28] Joanna R Hall, Olivia Matthews, Timothy N Volonakis, Eric Liggins, Karl P Lymer, Roland Baddeley, Innes C Cuthill, and Nicholas E Scott-Samuel. A platform for initial testing of multiple camouflage patterns. *Defence Technology*, 2020. 1
- [29] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. Mask R-CNN. In *IEEE International Conference on Computer Vision (ICCV)*, pages 2961–2969, 2017. 1
- [30] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 770–778, 2016. 4
- [31] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 770–778, 2016. 6
- [32] Xinxin Hu, Kailun Yang, Lei Fei, and Kaiwei Wang. ACNet: Attention based network to exploit complementary features for RGBD semantic segmentation. In *IEEE International Conference on Image Processing (ICIP)*, pages 1440–1444, 2019. 1
- [33] Wei-Chih Hung, Yi-Hsuan Tsai, Yan-Ting Liou, Yen-Yu Lin, and Ming-Hsuan Yang. Adversarial learning for semi-supervised semantic segmentation. *arXiv preprint arXiv:1802.07934*, 2018. 5
- [34] Ran Ju, Yang Liu, Tongwei Ren, Ling Ge, and Gangshan Wu. Depth-aware salient object detection using anisotropic center-surround difference. *Signal Processing: Image Communication*, 38:115–126, 2015. 8
- [35] Emma J Kelman, Daniel Osorio, and Roland J Baddeley. A review of cuttlefish camouflage and object recognition and evidence for depth perception. *Journal of Experimental Biology*, 211(11):1757–1763, 2008. 1, 2
- [36] Alex Kendall and Yarin Gal. What uncertainties do we need in bayesian deep learning for computer vision? In *Conference on Neural Information Processing Systems (NeurIPS)*, pages 5580–5590, 2017. 2, 5
- [37] Alex Kendall, Yarin Gal, and Roberto Cipolla. Multi-task learning using uncertainty to weigh losses for scene geometry and semantics. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 7482–7491, 2018. 7
- [38] Lingkai Kong, Jimeng Sun, and Chao Zhang. SDE-Net: Equipping deep neural networks with uncertainty estimates. *International Conference on Machine Learning (ICML)*, pages 5405–5415, 2020. 2
- [39] Balaji Lakshminarayanan, Alexander Pritzel, and Charles Blundell. Simple and scalable predictive uncertainty estimation using deep ensembles. In *Conference on Neural Information Processing Systems (NeurIPS)*, pages 6402–6413, 2017. 3
- [40] Trung-Nghia Le, Tam V Nguyen, Zhongliang Nie, Minh-Triet Tran, and Akihiro Sugimoto. Anabranch network for camouflaged object segmentation. *Computer Vision and Image Understanding (CVIU)*, 184:45–56, 2019. 1, 2, 3, 5, 6, 7, 8
- [41] Simcha Lev-Yadun, Amots Dafni, Moshe A Flaishman, Moshe Inbar, Ido Izhaki, Gadi Katzir, and Gidi Ne’eman. Plant coloration undermines herbivorous insect camouflage. *BioEssays*, 26(10):1126–1130, 2004. 1
- [42] Shuai. Li, Dinei. Florencio, Wanqing Li, Yaqin Zhao, and Chris Cook. A fusion framework for camouflaged moving foreground detection in the wavelet domain. *IEEE Transactions on Image Processing (TIP)*, 27(8):3918–3930, 2018. 2
- [43] Zhengqi Li, Tali Dekel, Forrester Cole, Richard Tucker, Noah Snavely, Ce Liu, and William T Freeman. Learning the depths of moving people by watching frozen people. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 4521–4530, 2019. 3
- [44] Di Lin, Guangyong Chen, Daniel Cohen-Or, Pheng-Ann Heng, and Hui Huang. Cascaded feature network for semantic segmentation of RGB-D images. In *IEEE International Conference on Computer Vision (ICCV)*, pages 1311–1319, 2017. 2
- [45] Nian Liu, Ni Zhang, and Junwei Han. Learning selective self-mutual attention for RGB-D saliency detection. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 13756–13765, 2020. 2
- [46] Jonathan Long, Evan Shelhamer, and Trevor Darrell. Fully convolutional networks for semantic segmentation. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 3431–3440, 2015. 2
- [47] Yunqiu Lv, Jing Zhang, Yuchao Dai, Aixuan Li, Bowen Liu, Nick Barnes, and Deng-Ping Fan. Simultaneously localize, segment and rank the camouflaged objects. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021. 3, 5, 6, 7, 8
- [48] Haiyang Mei, Ge-Peng Ji, Ziqi Wei, Xin Yang, Xiaopeng Wei, and Deng-Ping Fan. Camouflaged object segmentation with distraction mining. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 8772–8781, 2021. 6
- [49] Melia G Nafus, Jennifer M Germano, Jeanette A Perry, Brian D Todd, Allyson Walsh, and Ronald R Swaisgood. Hiding in plain sight: a study on camouflage and habitat selection in a slow-moving desert herbivore. *Behavioral Ecology*, 26(5):1389–1394, 2015. 1
- [50] Yuzhen Niu, Yujie Geng, Xueqing Li, and Feng Liu. Leveraging stereopsis for saliency analysis. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 454–461, 2012. 8
- [51] Seong-Jin Park, Ki-Sang Hong, and Seungyong Lee. RDFNet: RGB-D multi-level residual feature fusion for indoor semantic segmentation. In *IEEE International Conference on Computer Vision (ICCV)*, pages 4980–4989, 2017. 2
- [52] Olivier Penacchio, P George Lovell, Innes C Cuthill, Graeme D Ruxton, and Julie M Harris. Three-dimensional camouflage: exploiting photons to conceal form. *The American Naturalist*, 186(4):553–563, 2015. 2
- [53] Houwen Peng, Bing Li, Weihua Xiong, Weiming Hu, and Rongrong Ji. RGBD salient object detection: A benchmark

- and algorithms. In *European Conference on Computer Vision (ECCV)*, pages 92–109, 2014. 2, 8
- [54] Yongri Piao, Wei Ji, Jingjing Li, Miao Zhang, and Huchuan Lu. Depth-induced multi-scale recurrent attention network for saliency detection. In *IEEE International Conference on Computer Vision (ICCV)*, pages 7254–7263, 2019. 1, 3, 4
- [55] Yongri Piao, Wei Ji, Jingjing Li, Miao Zhang, and Huchuan Lu. Depth-induced multi-scale recurrent attention network for saliency detection. In *IEEE International Conference on Computer Vision (ICCV)*, pages 7254–7263, 2019. 2
- [56] Thomas W Pike. Quantifying camouflage and conspicuousness using visual salience. *Methods in Ecology and Evolution*, 9(8):1883–1895, 2018. 2
- [57] Natasha Price, Samuel Green, Jolyon Troscianko, Tom Tregenza, and Martin Stevens. Background matching and disruptive coloration as habitat-specific strategies for camouflage. *Scientific reports*, 9(1):1–10, 2019. 1
- [58] Chen Qian, Hongsheng Li, and Gang Zeng. Bi-directional cross-modality feature propagation with separation-and-aggregation gate for RGB-D semantic segmentation. In *European Conference on Computer Vision (ECCV)*, pages 561–577, 2020. 1
- [59] René Ranftl, Katrin Lasinger, David Hafner, Konrad Schindler, and Vladlen Koltun. Towards robust monocular depth estimation: Mixing datasets for zero-shot cross-dataset transfer. *IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)*, 2020. 1, 3, 4, 5, 8
- [60] Joseph Redmon, Santosh Divvala, Ross Girshick, and Ali Farhadi. You only look once: Unified, real-time object detection. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 779–788, 2016. 1
- [61] Jianqiang Ren, Xiaojin Gong, Lu Yu, Wenhui Zhou, and Michael Ying Yang. Exploiting global priors for RGB-D saliency detection. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR) Workshop*, pages 25–32, 2015. 2
- [62] Xiaofeng Ren, Liefeng Bo, and Dieter Fox. RGB-(D) scene labeling: Features and algorithms. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2759–2766, 2012. 2
- [63] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-Net: Convolutional networks for biomedical image segmentation. In *International Conference on Medical image computing and computer-assisted intervention*, pages 234–241, 2015. 1
- [64] Yichen Shen, Zhilu Zhang, Mert R Sabuncu, and Lin Sun. Real-time uncertainty estimation in computer vision via uncertainty-aware distribution distillation. In *IEEE Winter Conference on Applications of Computer Vision (WACV)*, pages 707–716, 2021. 2
- [65] Nathan Silberman and Rob Fergus. Indoor scene segmentation using a structured light sensor. In *IEEE International Conference on Computer Vision (ICCV) Workshop*, pages 601–608, 2011. 2
- [66] Przemysław Skurowski, Hassan Abdulameer, Jakub Baszczyk, Tomasz Depta, Adam Kornacki, and Przemysław Kozie. Animal camouflage analysis: Chameleon database. In *Unpublished Manuscript*, 2018. 6, 7, 8
- [67] Ariel Tankus and Yehezkel Yeshurun. Convexity-based visual camouflage breaking. *Computer Vision and Image Understanding (CVIU)*, 82(3):208–237, 2001. 2
- [68] Gerald Handerson Thayer. *Concealing-coloration in the animal kingdom: an exposition of the laws of disguise through color and pattern*. Macmillan Company, 1918. 1
- [69] Abhinav Valada, Rohit Mohan, and Wolfram Burgard. Self-supervised model adaptation for multimodal semantic segmentation. *International Journal of Computer Vision (IJCV)*, pages 1–47, 2019. 2
- [70] Anzhi Wang and Minghui Wang. RGB-D salient object detection via minimum barrier distance transform and saliency fusion. *IEEE Signal Processing Letters (SPL)*, 24(5):663–667, 2017. 2
- [71] Jun Wei, Shuhui Wang, and Qingming Huang. F³Net: Fusion, feedback and focus for salient object detection. In *AAAI Conference on Artificial Intelligence (AAAI)*, pages 12321–12328, 2020. 5, 7, 8
- [72] Evan C Wilson, Amy A Shipley, Benjamin Zuckerberg, M Zachariah Peery, and Jonathan N Pauli. An experimental translocation identifies habitat features that buffer camouflage mismatch in snowshoe hares. *Conservation Letters*, 12(2):12614, 2019. 1
- [73] Feng Xue, Chengxi Yong, Shan Xu, Hao Dong, Yuetong Luo, and Wei Jia. Camouflage performance analysis and evaluation framework based on features fusion. *Multimedia Tools and Applications*, 75(7):4065–4082, 2016. 2
- [74] Jinnan Yan, Trung-Nghia Le, Khanh-Duy Nguyen, Minh-Triet Tran, Thanh-Toan Do, and Tam V Nguyen. Mirror-net: Bio-inspired adversarial attack for camouflaged object segmentation. *arXiv preprint arXiv:2007.12881*, 2020. 1, 2, 3
- [75] Fan Yang, Qiang Zhai, Xin Li, Rui Huang, Ao Luo, Hong Cheng, and Deng-Ping Fan. Uncertainty-guided transformer reasoning for camouflaged object detection. In *IEEE International Conference on Computer Vision (ICCV)*, pages 4146–4155, 2021. 2, 6
- [76] Qiang Zhai, Xin Li, Fan Yang, Chenglizhao Chen, Hong Cheng, and Deng-Ping Fan. Mutual graph learning for camouflaged object detection. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 12997–13007, 2021. 6
- [77] Jing Zhang, Deng-Ping Fan, Yuchao Dai, Saeed Anwar, Fatemeh Sadat Saleh, Tong Zhang, and Nick Barnes. UC-Net: Uncertainty inspired RGB-D saliency detection via conditional variational autoencoders. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 8582–8591, 2020. 1, 6, 7
- [78] Miao Zhang, Weisong Ren, Yongri Piao, Zhengkun Rong, and Huchuan Lu. Select, supplement and focus for RGB-D saliency detection. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 3472–3481, 2020. 6
- [79] Yunfei Zheng, Xiongwei Zhang, Feng Wang, Tiejong Cao, Meng Sun, and Xiaobing Wang. Detection of people with camouflage pattern via dense deconvolution network. *IEEE Signal Processing Letters (SPL)*, 26(1):29–33, 2018. 2