

Source-free Depth for Object Pop-out

Zongwei Wu^{1,2,3} Danda Pani Paudel^{1,4} Deng-Ping Fan^{1*} Jingjing Wang⁵ Shuo Wang¹
Cédric Demonceaux^{2,6} Radu Timofte³ Luc Van Gool^{1,4}

¹ CVL, ETH Zurich ² University of Burgundy, CNRS, ICB ³ Computer Vision Lab, CAIDAS & IFI, University of Würzburg
⁴ INSAIT, Sofia University ⁵ AUST ⁶ University of Lorraine, CNRS, Inria, Loria

Abstract

Depth cues are known to be useful for visual perception. However, direct measurement of depth is often impracticable. Fortunately, though, modern learning-based methods offer promising depth maps by inference in the wild. In this work, we adapt such depth inference models for object segmentation using the objects’ “pop-out” prior in 3D. The “pop-out” is a simple composition prior that assumes objects reside on the background surface. Such compositional prior allows us to reason about objects in the 3D space. More specifically, we adapt the inferred depth maps such that objects can be localized using only 3D information. Such separation, however, requires knowledge about contact surface which we learn using the weak supervision of the segmentation mask. Our intermediate representation of contact surface, and thereby reasoning about objects purely in 3D, allows us to better transfer the depth knowledge into semantics. The proposed adaptation method uses only the depth model without needing the source data used for training, making the learning process efficient and practical. Our experiments on eight datasets of two challenging tasks, namely salient object detection and camouflaged object detection, consistently demonstrate the benefit of our method in terms of both performance and generalizability. The source code is publicly available at <https://github.com/Zongwei97/PopNet>.

1. Introduction

The 3D knowledge of the scene is long known to be complementary to the task of visual perception [12, 16, 56, 67, 91]. Often in practice, though, visual perception needs to be carried out using only 2D images. Given multiple images, 3D geometry may be recovered using the structure-from-motion techniques [26, 41, 55, 86]. Such inversion, however, is not compatible when only a single image is available. Under such circumstances, image inversion

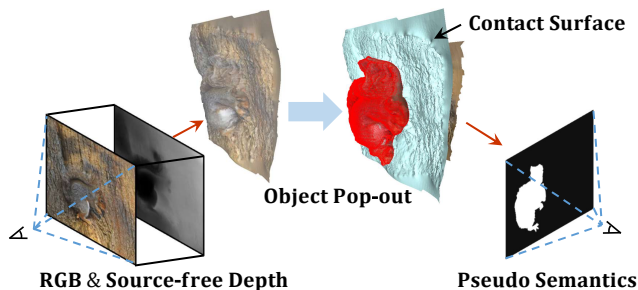


Figure 1. Depth to semantics conversion using the **object pop-out** prior. For input RGB and source-free depth pair, we learn contact surface. The obtained contact surface is then used to separate objects and backgrounds to derive pseudo semantics for supervision.

to depth map is usually done using learning-based methods [13, 45, 50, 51, 65], which have shown unparalleled success in the recent years. Unfortunately, the learning-based methods may not offer high-quality depth maps due to the generalization deficiency across domains [68, 73].

Despite poor generalization, the knowledge gained in one domain is shown to be useful in other close-by domains. This utility is harnessed by performing the so-called domain adaptation (DA) [2, 3, 6, 53, 61, 89]. In fact, it has been shown recently that DA methods can efficiently transfer knowledge using only the prediction models, *i.e.*, without requiring access to the data where the model is trained – also known as the source-free domain adaptation (SDA) [25, 28, 39, 71, 77]. The SDA methods are of gripping interest due to their efficiency and privacy promises.

The most existing SDA methods make one or both of these implicit assumptions: (a) a similar (as that of the source) supervising task at the target [28, 75, 77]; (b) task of discrete (and known) label space [27, 29, 36, 71]. The former assumption not only makes the source and target domains easier to compare but also potentially keeps the two domains closer. The latter assumption allows performing SDA by self-training where the discrete labels facilitate reasoning about the model’s confidence. The self-training is then performed by boosting the confidence at the target

*Corresponding Author: Deng-Ping Fan (dengpfan@gmail.com)

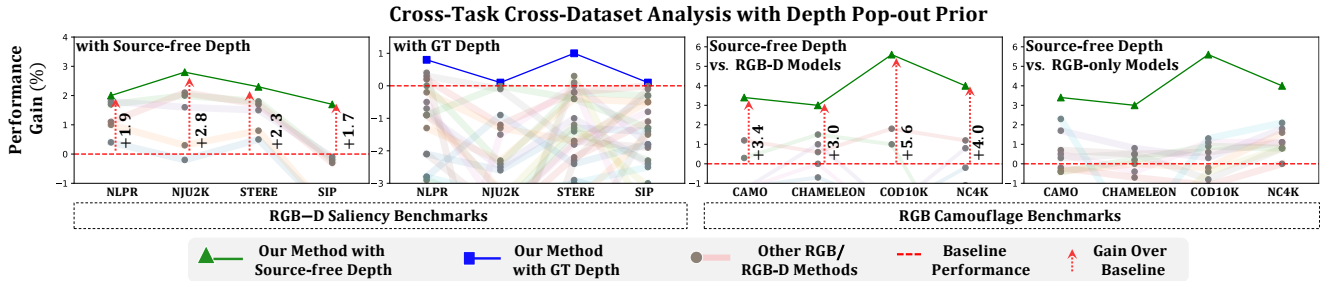


Figure 2. **The performance gained** in F-measure using our method over the established baselines. Our method ($\blacktriangle, \blacksquare$) significantly improves baselines on **8 datasets** and **2 tasks** (each task on 4 datasets– SOD: left two; COD: right two). We also compare **24 methods** (\bullet), where our method offers state-of-the-art results despite their task specialization. Note that all methods are connected using lines to illustrate their performance fluctuations across datasets. Please, refer Tables 1 & 2 and Section 4.3 for more details and discussions.

for some picked reliable examples.

In this work, we aim for source-free transfer of depth knowledge for object detection. Such transfer is desired to assist in locating the object by depth cues and to exploit the depth knowledge despite the domain gap. The addressed problem setting differs from the standard SDA in terms of (a) the source and target tasks’ difference; (b) and the continuous label space of the depth. These differences (with standard SDA) make our task at hand very challenging, which is addressed for the first time in this paper, up to our knowledge. To address such a challenging problem, we rely on the “pop-out” prior, which allows us to reason about the object’s location directly in the 3D space. The “pop-out” is a simple composition prior that assumes objects reside *on* the background surface. A graphical illustration of the used “pop-out” prior, using the results obtained by our method, is given in Figure 1.

The pop-out prior for image composition was successfully used by *Kang et al.* in [24]. An early work of *Treisman* has provided an in-depth study of such prior in [62]. In this work, we rely on the same compositional foundation of these works and exploit the pop-out before transferring the depth of knowledge across domains. Although our motivation comes from these early works, our experimental setup largely differs from theirs. We differ in terms of not only depth knowledge transfer across domains (without source data) but also in target supervision using only semantics.

The proposed method exploits the source-free depth to map it into a space where objects in depth stand out better against the background. This mapping is used for object and background separation using a learned contact surface between them. Such separation allows us to derive the semantic masks which can be directly compared against the ground truth for supervision. Using this supervision at the target, we can minimize the domain and task gap between the source and target. The overall framework of our method that performs cross-task cross-domain knowledge transfer by using the intermediate representation of the pop-

out space is shown in Figure 3. As can be seen, we first leverage an object popping network to encourage the object to jump out from the source-free depth. Then, we introduce another network *i.e.*, the segmentation with contact surface, to localize the object and predict the contact surface. These learning modules are jointly trained in an end-to-end manner, transferring the source-free depths into intermediate representations which are adapted to the target task, *i.e.*, object detection. To evaluate the proposed method, we conducted exhaustive experiments on eight datasets of two challenging tasks, namely salient object detection and camouflaged object detection. In both tasks, our method significantly improves the established baselines and offers state-of-the-art results at the same time, whose overview can be seen in Figure 2. The major contributions are as follows:

- Our problem of transferring source-free depth knowledge across domains and tasks is practical and novel.
- Our method relies on our object pop-out prior for visual understanding, which is simple and effective.
- Results of our method in two different tasks are significantly better than the baselines and existing models.

2. Related Works

Source-free Adaptation: Knowledge transfer by domain adaptation without access to the source data has recently gathered vast interest [1, 25, 28, 39, 71, 77], due to the privacy, practicality, and efficiency reasons. We observe accessing the source data for transferring the depth knowledge learned by the off-the-shelf knowledge models [13, 51] is particularly impractical due to their multi-stage training on various datasets. The existing source-free adaptation methods either use generative [28, 29, 34], pseudo-label [25, 37, 64], or other customized [54, 76] approaches. In this work, we use the pseudo-label-based approach. However, using existing methods is not straightforward in our setting due to the tasks’ difference between source and target. With our pop-out technique that provides the pseudo semantics, the source-free depth is better transferred across tasks.

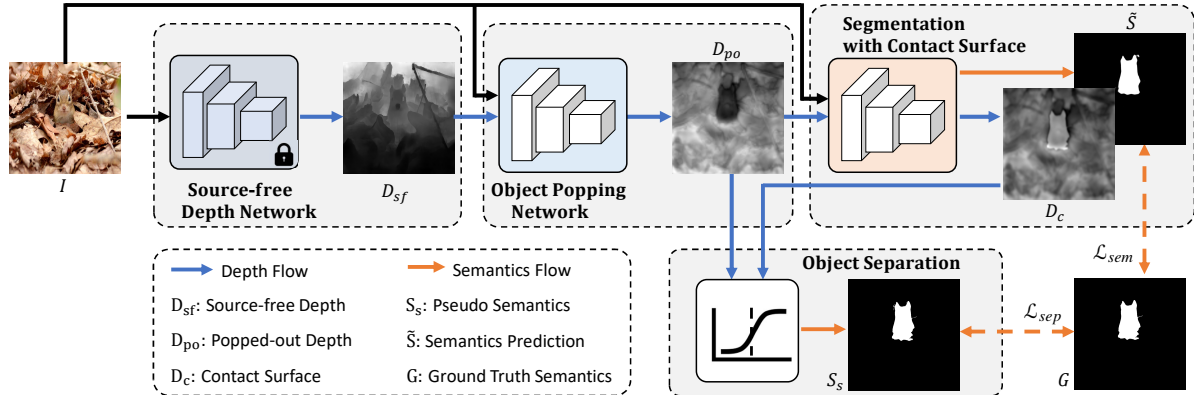


Figure 3. Our **proposed framework**, termed PopNet, is composed of a source-free network, an object popping network, a segmentation with contact surface, and an object separation module. The source-free depth network generates pseudo-depth in an off-the-shelf manner (Section 3.1). The object popping network converts the source-free depth into the popped-out depth of objects, bridging the cross-domain and cross-task gaps (Section 3.2). The segmentation network uses this depth to estimate the object’s mask and contact surface (Section 3.3). The object separation module then converts popped-out depth to the second mask of objects using the contact surface (Section 3.4). We compare both semantic masks against ground truth to supervise the whole pipeline in an end-to-end manner.

Salient Object Detection (SOD): Saliency detection aims to detect and segment the most prominent region within an image that visually attracts human attention [11]. A number of works have shown that saliency can be an auxiliary step for different vision tasks such as object tracking [93], object detection [58], *etc.* Conventional saliency works are unimodal, *i.e.*, they only require RGB images as input. In generic and common settings, RGB-based models [38, 70, 84] have already achieved very promising results. More recently, several works [10, 20, 49, 66, 68, 88, 92] exploit depth maps as additional clues to the 3D geometry since the depth can provide more truthful information on the object boundary as well as the scale awareness. These 3D features further improve the detection accuracy and performance in challenging scenarios [22, 33, 69, 80].

Camouflaged Object Detection (COD): Camouflage detection aims to find the preying object within an image. For computer vision society, primary works [9, 32, 47] often compare COD with SOD. A number of works [8, 9, 74, 79] have shown that simply extending saliency models [70] on COD will lead to undesired results, which is mainly caused by the nature of target object, *i.e.*, concealed or prominent. Hence, to constrain the attention on the concealed objects, several works come up with different perceptual systems that mimic human behavior vis-a-vis camouflaged objects, such as three-stage localize-segment-rank strategy [42]; iterative refinement [21], which is similar to repeatedly looking on the images; zooming into possible regions [47, 59]. Others [18, 40, 60, 82, 83, 90, 94, 95] deeply explore the texture difference with the help of the gradient [18], frequency [90], edges [60, 94], and probability [32, 74].

The latest psychological studies [5, 7] have shown that human perception can naturally benefit from the depth cues

to understand the scene: (A) the smooth variation within the object can contribute to alleviating the fake edges and preserving the object structure; (B) the depth discontinuity on the object boundary can make the segmentation easier. Inspired by these observations, we aim to explore the source-free depth for both SOD and COD tasks. To tackle the domain gap for source-free depth maps, we propose to jointly finetune the source-free depth together with the semantic network in an end-to-end manner, with both self-supervised loss and weak semantic supervision.

3. Proposed PopNet

Given an input RGB image I with size $I \in \mathbb{R}^{3 \times H \times W}$, where H and W are the spatial resolutions, *i.e.*, height and width, of the image, our objective is to predict the semantic mask $\hat{S} \in \mathbb{R}^{H \times W}$ for object detection. As shown in Figure 3, the input image I is firstly fed into a frozen-weight depth network to generate the source-free depth $D_{sf} \in \mathbb{R}^{H \times W}$ (Section 3.1). Then the mimicked multi-modal images are fed together into the depth popping network to compute the intermediate popped-out depth $D_{po} \in \mathbb{R}^{H \times W}$ (Section 3.2). This intermediate representation, as well as the RGB image I , is later processed by the segmentation network and transformed into a contact surface $D_c \in \mathbb{R}^{H \times W}$ and a semantics prediction $\hat{S} \in \mathbb{R}^{H \times W}$ (Section 3.3). On the one hand, the semantics prediction is directly supervised by the ground truth mask GT , denoted as G , which is similar to conventional segmentation supervision. On the other hand, we further explore the contact surface to pop the object out of the background by means of our object separation module (Section 3.4). This transfers the geometric cues into pseudo semantics and leads to another level of supervision.

3.1. Source-free Depth Network

In a practical setting, the GT depth is not always available. Therefore, we generate the source-free depth D_{sf} in an off-the-shelf manner to mimic the multi-modal input. We choose the state-of-the-art DPT model [50] with frozen weights as our source-free depth network, which offers us promising depth at the target. This choice is made upon its generalization capability as suggested in [51]. To obtain the depth map with the highest quality possible, by enhancing local details, we apply the boosting method [45] together with DPT. Despite the plausible results achieved by learning-based methods, the obtained source-free depth does not always offer high-quality geometric cues due to the domain gap. Therefore, we leverage the geometric and semantic priors to jointly finetune the source-free depth.

3.2. Object Popping Network

Network Architecture: We build a depth popping network to refine/smooth the source-free depth. The popping network follows the encoder-decoder design with skip connection as shown in Figure 4. In our case, we simply concatenate RGB and source-free depth at the input side to form a 4-channel input and feed them into the popping network. The encoder extracts semantic cues and generates five-scale outputs. Our decoder is composed of Conv2D, BN, ReLU, and upsampling layers. Following U-Net [52], we build a skip connection through simple addition.

Structure Preserving: To supervise our popping network, we first guide the depth refinement with the help of generated pseudo-depth. We only leverage structural similarity since we aim to detect, preserve, and extract object structure from the intermediate representation. We use the following SSIM loss [13] for structural similarity.

$$\mathcal{L}_{dep} = SSIM(D_{po}, D_{sf}). \quad (1)$$

Local Depth Smoothing: In addition to the pseudo-depth supervision, we propose two losses to constrain the depth also by semantics. We assume that the objects' structure should be distinguishable from the background, *i.e.*, it should be smooth within the object region and sharp on the bounding pixels. Hence, we propose to leverage the weak semantic cues together with the geometric priors. Specifically, we first introduce a local smoothness loss. The locality is defined by the ground truth semantics G . Technically, we mask out background pixels with element-wise multiplication to suppress the inactive area through $D_{obj} = D_{po} \otimes G$. Let ∇_x and ∇_y be the Sobel operations. Then, our local loss \mathcal{L}_{loc} is expressed as:

$$\begin{aligned} \overrightarrow{n(p)} &= (-\nabla_x(D_{obj}(p)), -\nabla_y(D_{obj}(p)), 1); \\ \mathcal{L}_{loc} &= \sum_p \sum_{q \in \mathcal{N}(p)} 1 - \text{cosine}(\overrightarrow{n(p)}; \overrightarrow{n(q)}), \end{aligned} \quad (2)$$

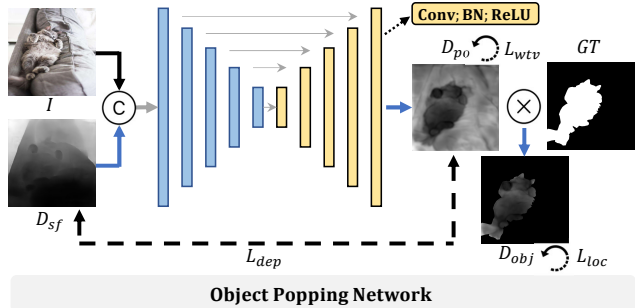


Figure 4. Our **object popping network** maps RGB-D inputs into popped-out depths. This network is supervised using a combination of structure preserving, local depth smoothing, and depth edge sharpening losses \mathcal{L}_{dep} , \mathcal{L}_{loc} , and \mathcal{L}_{wtv} , respectively.

where \vec{n} stands for the normal, p denotes the pixels within the object region, $\mathcal{N}(p)$ is the neighboring pixels, and cosine is the cosine similarity between two vectors. This way, our local loss only works on the object area, making the object structure consistent within the target region. Applying local smoothness loss reduces the depth noise at the object level.

Depth Edge Sharpening: In addition to the local smoothness, we also use edge sharpening. The edge sharpening loss is formulated as a weighted total variation. For this, we first compute the edge-aware weight $w(p)$ at any pixel p as,

$$w(p) = \begin{cases} w_0, & \text{if } \nabla_x(G(p))^2 + \nabla_y(G(p))^2 \neq 0, \\ w_0 + \gamma, & \text{otherwise,} \end{cases} \quad (3)$$

where w_0 is a pre-defined non-zero weight and γ is an additional weight for boundary pixels. In our setting, we choose w_0 as the normalized (by the image size) count of the boundary pixels, and we set $\gamma = 0.5$. We adopt the square form such that the large gradients play more important roles. Our weighted total-variation loss is given by,

$$\mathcal{L}_{wtv} = \sum_p \sum_{q \in \mathcal{N}(p)} w(p) \cdot \|D_{po}(p) - D_{po}(q)\|^2. \quad (4)$$

Our weighted total-variance loss differs from the conventional edge loss, due to our weighting function. More specifically, our weighting function relies on semantic boundaries, unlike the commonly used image gradients [14, 17]. Our motivation for using semantic boundaries instead of image gradients comes from our interest in performing object detection under challenging conditions, such as camouflaged objects. In such cases, image gradients may result in misleading weights. At first glance, our loss function seems to be similar to semantic-guided depth estimation methods [4, 35, 72]. However, [72] uses GT depths, and [4, 35] use multi-frames for supervision. Unfortunately, such supervision is not possible in our setup. Note that we exploit single-view source-free depth while only using semantic ground truth for supervision. Furthermore, we transfer the source-free depths' knowledge despite their domain

gap. Existing regularizations between depth and semantics are already exploited in our method, which is shown to be complimentary to our pop-out prior. We argue that our network architecture that exploits the pop-out prior is not trivial while simultaneously being generic and easy to use.

The total objective function \mathcal{L}_{pop} to supervise our object popping network is given by,

$$\mathcal{L}_{pop} = \mathcal{L}_{dep} + \lambda_1 \cdot \mathcal{L}_{loc} + \lambda_2 \cdot \mathcal{L}_{wtv}, \quad (5)$$

where λ_1 and λ_2 are the hyperparameters.

3.3. Segmentation with Contact Surface

The smoothness and edge losses encourage homogenizing the object structure, making it noticeable from the background. We now aim to further enlarge the object-background distance to make the object structure jump out. Specifically, we use an RGB-D segmentation network as shown in Figure 5. The main component of our segmentation network is a three-stream RGB-D network with some fusion design. In our setting, we choose [92] as our baseline since it is one of the SOTA RGB-D methods for saliency detection. We add a surface head to learn the depth of the contact surface D_c , which has the same resolution as the input depth D_{po} . Our surface head is composed of ConvLayer (Conv2D, BN, ReLU) and a Conv2D, which first decodes the feature maps and then transfers them into a 1-D map.

3.4. Object Separation

Using the previously discussed contact surface, in this section, we aim to separate the object from its background. At this point, we make an assumption that pixels in front of the contact surface belong to objects. The remaining pixels belong to the background. This assumption allows us to explicitly transfer the 3D knowledge into 2D semantics. Let the predicted depth of the contact surface be $D_c \in \mathbb{R}^{H \times W}$. We obtain the pseudo semantics S_s , using the popped-out depth D_{po} and surface’s depth D_c , as:

$$S_s = \text{sigmoid}(\sigma \cdot (D_{po} - D_c)), \quad (6)$$

where σ is a scalar value that controls the slope of the sigmoid function. In our experiments, we use $\sigma = 10$ to perform soft-thresholding, mimicking the desired hard one for binary outputs. Such soft thresholding facilitates the gradient back-propagation required for training. Finally, we minimize the gap between the pseudo semantics S_s and the GT semantics G with binary cross-entropy (BCE):

$$\mathcal{L}_{sep} = \text{BCE}(S_s, G). \quad (7)$$

3.5. Overall Loss Function

Both of our trainable modules, *i.e.*, *object popping* and *segmentation networks*, of our framework (Figure 3) are

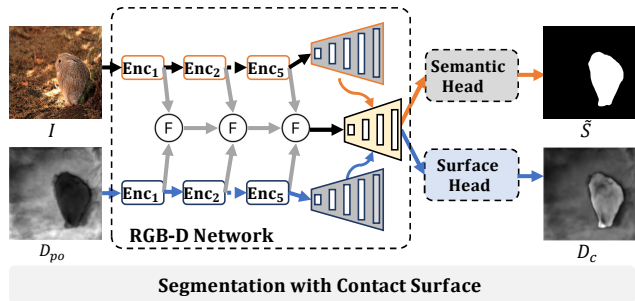


Figure 5. Our **segmentation network** uses a basic RGB-D three-stream network [92]. In addition to the conventional semantic head, we learn to predict pixel-wise **contact surface**. The contact surface is later used to transfer depth knowledge to semantics.

trained in an end-to-end manner. Therefore, the overall loss function consists of three parts: our depth popping loss \mathcal{L}_{pop} , our object separation loss \mathcal{L}_{sep} , and the conventional semantic loss \mathcal{L}_{sem} from the RGB-D baseline network. The total loss \mathcal{L}_{total} used for training is given by,

$$\mathcal{L}_{total} = \mathcal{L}_{pop} + \alpha_1 \cdot \mathcal{L}_{sep} + \alpha_2 \cdot \mathcal{L}_{sem}, \quad (8)$$

where α_1 and α_2 are hyperparameters.

Remarks: Our losses work in a complementary manner. The \mathcal{L}_{pop} plays the same role as a smooth filter as in image smoothing. It removes the noisy depth response due to the domain gap while preserving the object structure with the help of a weak semantic label. Hence, the smoothed background becomes less informative while the object region becomes uniform, making it easily detectable. These functionalities contribute to transferring the source-free depth into popped-out depth, as desired. Such a process brings objects above the background surface, despite their distance from the camera and other distracting surfaces. The \mathcal{L}_{sep} , on the other hand, fully benefits from the “pop-out” prior to segment the object from the background. With the help of the learned contact surface, this loss enlarges the foreground-background distance by pulling them in opposite directions. Such pulling results in a binary-like mask, which effectively bridges the gap between geometric knowledge and semantics. Finally, both depth-transferred and learned semantics are compared to the ground truths for supervised training.

4. Results

4.1. Experimental Setup

Dataset Preparation: To better illustrate the generalizability of our approach, we evaluate the effectiveness of our approach on both SOD and COD benchmarks. We choose four widely used RGB-D SOD datasets, *i.e.*, NLPR [48], NJUK [23], STERE [46], and SIP [10], as well as four COD datasets, *i.e.*, CAMO [30], CHAMELEON [57],

Table 1. Quantitative comparison on RGB-D SOD datasets. \uparrow (\downarrow) denotes that the higher (lower) is better. We use the Mean Absolute Error (M), max F-measure (F_m), S-measure (S_m), and max E-measure (E_m) as evaluation metrics. G.D. stands for GT Depth. **Bold** denotes the best performance.

G.D.	Public.	Dataset Metric	NLPR [48]				NJUK [23]				STERE [46]				SIP [10]			
			$M\downarrow$	$F_m\uparrow$	$S_m\uparrow$	$E_m\uparrow$	$M\downarrow$	$F_m\uparrow$	$S_m\uparrow$	$E_m\uparrow$	$M\downarrow$	$F_m\uparrow$	$S_m\uparrow$	$E_m\uparrow$	$M\downarrow$	$F_m\uparrow$	$S_m\uparrow$	$E_m\uparrow$
Performance of RGB-D Models Trained with Source-free Depth																		
\times	<i>MM</i> ₂₁ [85]	DFM-Net	.027	.909	.914	.944	.046	.903	.895	.927	.042	.906	.903	.934	.067	.873	.850	.891
\times	<i>TIP</i> ₂₂ [63]	DCMF	.027	.915	.921	.943	.044	.908	.903	.929	.041	.909	.907	.931	.067	.873	.853	.893
\times	<i>CVPR</i> ₂₂ [21]	SegMAR	.024	.923	.920	.952	.036	.921	.909	.941	.037	.916	.907	.936	.052	.893	.872	.914
\times	<i>CVPR</i> ₂₂ [47]	ZoomNet	.023	.916	.919	.944	.037	.926	.914	.940	.037	.918	.909	.938	.054	.891	.868	.909
\times	Ours	PopNet	.022	.925	.926	.956	.031	.931	.920	.949	.032	.922	.916	.947	.046	.911	.885	.926
Performance of RGB-D Models Trained with GT Depth																		
\checkmark	<i>TIP</i> ₂₁ [87]	BIANet	.032	.888	.900	.930	.056	.878	.867	.898	.048	.898	.895	.918	.091	.816	.802	.847
\checkmark	<i>TIP</i> ₂₁ [33]	HAINet	.024	.920	.924	.956	.037	.924	.911	.940	.040	.917	.907	.938	.052	.907	.879	.917
\checkmark	<i>TNNLS</i> ₂₁ [10]	D3Net	.029	.904	.911	.942	.046	.909	.899	.927	.044	.902	.906	.925	.063	.880	.860	.897
\checkmark	<i>ECCV</i> ₂₂ [31]	SPSN	.023	.917	.923	.956	.032	.927	.918	.949	.035	.909	.906	.941	.043	.910	.891	.932
\checkmark	Ours	PopNet	.019	.927	.932	.963	.030	.936	.924	.952	.033	.924	.917	.947	.040	.923	.897	.937

COD10K [9], and NC4K [42]. For SOD datasets, we conduct experiments with both GT depth and source-free depth. We follow the conventional learning protocol [19, 69, 92] and use 700 images from NLPR and 1,485 images from NJUK for training. The rest are used for testing. For the unimodal COD dataset, we compare with both RGB COD models and RGB-D SOD models retrained on the COD datasets with the same source-free depth D_{sf} . We follow the conventional training/testing protocol [8, 9, 21, 42, 47] and use 3,040 images from COD10K and 1,000 images from CAMO for training. The rest are used for testing.

Evaluation Metrics: We evaluate the performance with four generally-recognized metrics: Mean Absolute Error (M), max F-measure (F_m), S-measure (S_m), and max E-measure (E_m). All the object segmentation masks are trained or downloaded from the official resources. To make a fair comparison, we evaluate the prediction semantics with the standardized evaluation protocol as [92].

Implementation Details: Our model is implemented based on Pytorch with a V100 GPU. We use the Adam algorithm as an optimizer. The learning rate is initialized to $1e-4$ and is further divided by 10 every 60 epochs. We set the input resolution to 352×352 resolution for RGB and depth. A detailed comparison with higher resolution can be found in Table 3. During training, conventional data augmentation such as random flipping, rotating, and border clipping are adopted. The training takes around 6 hours for RGB-SOD tasks and 12 hours for COD tasks for 100 epochs.

4.2. Comparisons

Comparison with RGB-D SOD Models: We present in Table 1 the performance on SOD benchmarks with our source-free depth or with GT depth. It can be seen that our model with source-free depth achieves very competitive performance compared to many RGB-D models with GT depths. Our method with GT depth also outperforms the SOTA counterparts. The qualitative comparison can be found in Fig. 6. Note that we also retrain the SOTA uni-

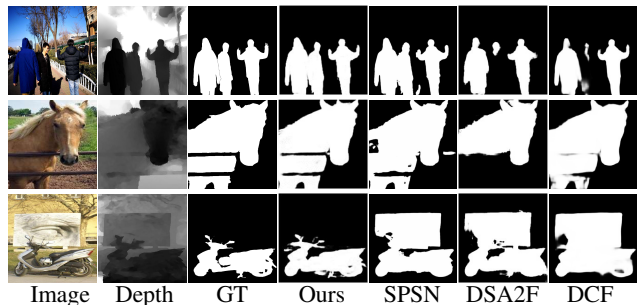


Figure 6. **Qualitative Comparison with GT Depth.** Our method outperforms all counterparts while dealing with multi-objects, large depth variation, and visually-mixed foreground-background.

modal COD models SegMAR [21] and ZoomNet [47] on the SOD dataset only with RGB images. We show that our model with source-free depth outperforms these counterparts, showing that our method can better generalize across different tasks with favorably better performance.

Comparison with COD Models: We present in Table 2 the performance of the most competitive SOTA methods, including task-specific COD models as well as retrained RGB-D SOD models with source-free depth. For a fair comparison, we retrain all RGB-D methods, SegMAR [21], and ZoomNet [47] in an end-to-end manner on the same resolution images as ours. Some RGB-only methods perform even better than many RGB-D methods, mainly due to their COD task-specific designs and the lack of ground-truth depths required for RGB-D methods. It is important to note that competing RGB-only COD methods do not perform favorably on SOD tasks. Please, refer to Table 1 & 3 to observe their poor cross-tasks generalization.

Towards Higher Resolution: Previous studies have shown that the image resolution may influence the model performance [43, 47, 73, 83]. For example, the current SOTA COD method ZoomNet [47] is with main scale 384^2 and implies the highest resolution of $(384 \times 1.5)^2 = 576^2$, as it operates on $0.5 \times, 1 \times,$ and $1.5 \times$ scales. To make a fair compari-

Table 2. Quantitative comparison on COD datasets. Pseudo stands for source-free depth used for RGB-D methods.

Pseudo Public.	Dataset Metric	CAMO [30]				CHAMELEON [57]				COD10K [9]				NC4K [42]				
		$M \downarrow$	$F_m \uparrow$	$S_m \uparrow$	$E_m \uparrow$	$M \downarrow$	$F_m \uparrow$	$S_m \uparrow$	$E_m \uparrow$	$M \downarrow$	$F_m \uparrow$	$S_m \uparrow$	$E_m \uparrow$	$M \downarrow$	$F_m \uparrow$	$S_m \uparrow$	$E_m \uparrow$	
Performance of RGB COD Models																		
✗	CVPR ₂₀ [9]	SINet	.099	.762	.751	.790	.044	.845	.868	.908	.051	.708	.771	.832	.058	.804	.808	.873
✗	CVPR ₂₁ [42]	SLSR	.080	.791	.787	.843	.030	.866	.889	.938	.037	.756	.804	.854	.048	.836	.839	.898
✗	CVPR ₂₁ [78]	MGL-R	.088	.791	.775	.820	.031	.868	.893	.932	.035	.767	.813	.874	.053	.828	.832	.876
✗	CVPR ₂₁ [44]	PFNet	.085	.793	.782	.845	.033	.859	.882	.927	.040	.747	.800	.880	.053	.820	.829	.891
✗	CVPR ₂₁ [32]	UJSC	.072	.812	.800	.861	.030	.874	.891	.948	.035	.761	.808	.886	.047	.838	.841	.900
✗	IJCAI ₂₁ [59]	C2FNet	.079	.802	.796	.856	.032	.871	.888	.936	.036	.764	.813	.894	.049	.831	.838	.898
✗	ICCV ₂₁ [74]	UGTR	.086	.800	.783	.829	.031	.862	.887	.926	.036	.769	.816	.873	.052	.831	.839	.884
✗	CVPR ₂₂ [21]	SegMAR	.080	.799	.794	.857	.032	.871	.887	.935	.039	.750	.799	.876	.050	.828	.836	.893
✗	CVPR ₂₂ [47]	ZoomNet	.074	.818	.801	.858	.033	.829	.859	.915	.034	.771	.808	.872	.045	.841	.843	.893
Performance of RGB-D Models Retained with Source-free Depth																		
✓	MM ₂₁ [80]	CDINet	.100	.638	.732	.766	.036	.787	.879	.903	.044	.610	.778	.821	.067	.697	.793	.830
✓	CVPR ₂₁ [19]	DCF	.089	.724	.749	.834	.037	.821	.850	.923	.040	.685	.766	.864	.061	.765	.791	.878
✓	ICCV ₂₁ [81]	CMINet	.087	.798	.782	.827	.032	.881	.891	.930	.039	.768	.811	.868	.053	.832	.839	.888
✓	ICCV ₂₁ [92]	SPNet	.083	.807	.783	.831	.033	.872	.888	.930	.037	.776	.808	.869	.054	.828	.825	.874
✓	TIP ₂₂ [63]	DCMF	.115	.737	.728	.757	.059	.807	.830	.853	.063	.679	.748	.776	.077	.782	.794	.820
✓	ECCV ₂₂ [31]	SPSN	.084	.782	.773	.829	.032	.866	.887	.932	.042	.727	.789	.854	.059	.803	.813	.867
✓	Ours	PopNet	.073	.821	.806	.869	.022	.893	.910	.962	.031	.789	.827	.897	.043	.852	.852	.908

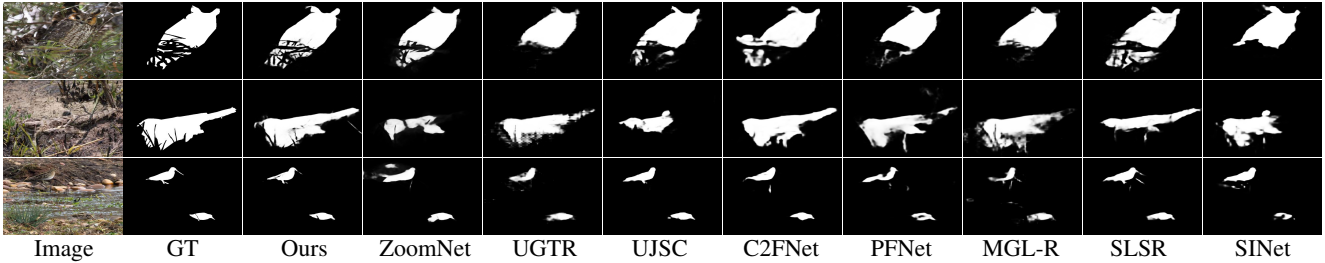


Figure 7. **Qualitative comparison.** Our method can better preserve the object structure compared to other counterparts, especially while dealing with objects with occlusion (two first rows). Our method performs favorably with multiple objects (last row). Better to zoom in.

son, we retrain the model with the same resolution (352^2 or 512^2) as ours. We show in Table3-ZoomNet* that the results deteriorate as expected. Compared to these counterparts, our method offers a good trade-off between accuracy and efficiency. More comparisons on SOD and COD benchmarks can be found in the [supplementary material](#).

Qualitative Comparisons: Figure 7 presents the output of our network on challenging cases. It can be seen while dealing with objects occluded by thin objects ($1^{st} - 4^{th}$ rows), our method can accurately reason about the segmentation masks closer to the GT. We also achieve better performance while dealing with multiple objects (last row). More discussions on multiple objects can be found in Table 6.

4.3. Ablation Study

Loss: In this section, we conduct experiments on analyzing the effectiveness of the proposed losses. The quantitative results of different loss combinations are provided in Table 4. It can be seen that each proposed loss behaves properly, *i.e.*, improving the performance compared to the baseline. More discussions and ablation studies on the hyperparameters can be found in the [supplementary material](#).

Table 3. End-to-end comparison with different resolutions on SOD benchmarks. Our method with source-free depth generalizes significantly better compared to SOTA COD models.

Model	Size	Flops (G)	NJUK [23]				SIP [10]			
			$M \downarrow$	$F_m \uparrow$	$S_m \uparrow$	$E_m \uparrow$	$M \downarrow$	$F_m \uparrow$	$S_m \uparrow$	$E_m \uparrow$
SegMAR [21]	352^2	67.3	.036	.921	.909	.941	.052	.893	.872	.914
ZoomNet [47]	352^2	167.8	.037	.926	.914	.940	.054	.891	.868	.909
Ours	352^2	228.8	.031	.931	.920	.949	.046	.911	.885	.926
SegMAR [21]	512^2	142.4	.035	.927	.914	.943	.050	.899	.878	.917
ZoomNet [47]	512^2	353.4	.036	.926	.915	.942	.052	.895	.873	.910
Ours	512^2	484.0	.031	.933	.922	.951	.044	.911	.890	.927

Table 4. Ablation study on the proposed losses.

\mathcal{L}_{dep}	\mathcal{L}_{loc}	\mathcal{L}_{wiv}	\mathcal{L}_{sep}	SIP [10]				NC4K [42]			
				$M \downarrow$	$F_m \uparrow$	$S_m \uparrow$	$E_m \uparrow$	$M \downarrow$	$F_m \uparrow$	$S_m \uparrow$	$E_m \uparrow$
-	-	-	-	.048	.903	.884	.922	.052	.832	.832	.893
✓	-	-	-	.046	.907	.889	.925	.051	.833	.839	.895
-	✓	-	-	.045	.908	.893	.929	.048	.837	.844	.898
-	-	✓	-	.046	.906	.891	.927	.050	.833	.841	.894
-	-	-	✓	.043	.914	.893	.933	.048	.840	.848	.900
✓	✓	-	-	.044	.911	.893	.928	.049	.837	.844	.897
✓	-	✓	-	.046	.909	.893	.927	.046	.840	.845	.898
✓	-	✓	✓	.040	.918	.897	.935	.045	.848	.849	.904
✓	✓	-	✓	.042	.916	.894	.931	.044	.850	.850	.906
✓	✓	✓	✓	.040	.923	.897	.937	.043	.852	.852	.908

Table 5. Generalization and cost over different RGB-D baselines.

Dataset Metric	Flops (G)	Param (M)	SIP [10]				NC4K [42]			
			$M \downarrow$	$F_m \uparrow$	$S_m \uparrow$	$E_m \uparrow$	$M \downarrow$	$F_m \uparrow$	$S_m \uparrow$	$E_m \uparrow$
HAINet [33]	363.2	59.8	.053	.899	.874	.919	.057	.809	.804	.872
+ Ours	373.7	72.5	.051	.910	.886	.923	.055	.814	.811	.878
SPNet [92]	149.0	150.4	.044	.911	.887	.914	.054	.828	.825	.874
+ Ours	159.5	163.1	.042	.917	.894	.932	.044	.851	.851	.905

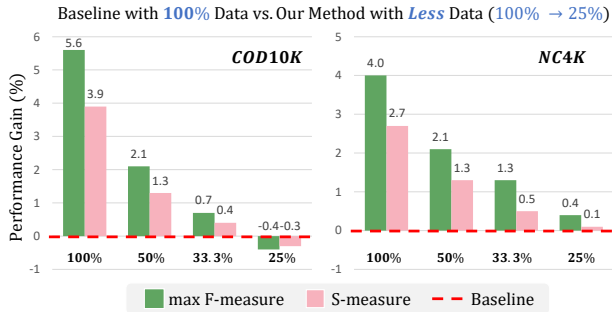


Figure 8. Benefit of our method on **reducing training data**. When our network is trained only with 25% data, the performance remains competitive compared to the baseline, *i.e.*, the absolute performance gains are +0.4% max F-measure and +0.1% S-measure on NC4K compared to the baseline trained with all data.

Object Popping Network as Plug-in: Our object popping network can be easily adapted with different encoders and with different existing RGB-D models. For example, with a ResNet-18 [15] encoder and convolution-based decoder, our popping only costs around 12.7M additional learning parameters or 48.7 MB model size. We show in Table 5 that our method can favorably improve performance over the baseline with less than an extra 10% GFlops.

Pop-out Under Reduced Training Data: Here, we are interested in analyzing source-free depth’s benefits. Therefore, we conduct different experiments by reducing the training data. As shown in Figure 8(left), when both our PopNet and our baseline are trained with all data, our PopNet can lead to absolute improvements with 5.6% in max F-measure and with 3.9% in S-measure on COD10K dataset. While our PopNet is trained with only 25% data, it can still achieve competitive performance compared to our baseline trained with all data. Similar phenomena can be observed in the NC4K dataset as shown in Figure 8(right). To conclude, our method can efficiently explore the geometric prior and significantly reduce the required training data volume.

Gain over baseline: With depth cues, as shown in Figure 9, we boost the performance in 3125 over 4121 images (~75% cases). We also show in Table 6 that our network performs favorably over baselines with single or multiple objects. Our method may fail when the source-free depth is clueless. This mainly happens when the object is well concealed and fools the depth network. However, such cases are also challenging, even for humans.

Are RGB-D Methods Better Than RGB-only Methods?

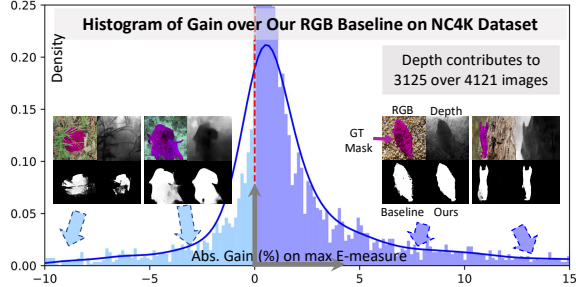


Figure 9. Histogram on the gain. Please zoom in for details.

Table 6. Multi-Object performance on NC4K [42] with size 512².

Obj. Nbr. (%) Metric	Single (92%)		Two (6%)		More (2%)		Overall	
	$M \downarrow$	$F_m \uparrow$	$M \downarrow$	$F_m \uparrow$	$M \downarrow$	$F_m \uparrow$	$M \downarrow$	$F_m \uparrow$
RGB Baseline	.054	.828	.067	.811	.091	.738	.056	.825
+ D_{sf}	.050	.842	.063	.821	.093	.746	.051	.839
Ours	.040	.864	.051	.847	.079	.767	.042	.861

Table 7. Performance with RGB- D_{sf} model vs. with RGB-only baseline. D_{sf} stands for source-free depth.

Dataset Metric	COD10K [9]				NC4K [42]			
	$M \downarrow$	$F_m \uparrow$	$S_m \uparrow$	$E_m \uparrow$	$M \downarrow$	$F_m \uparrow$	$S_m \uparrow$	$E_m \uparrow$
DASNet [88]	.041	.643	.793	.864	.055	.747	.830	.879
+ D_{sf}	.041	.642	.796	.858	.055	.743	.830	.874
SPNet [92]	.040	.743	.801	.867	.052	.846	.833	.883
+ D_{sf}	.037	.776	.808	.869	.054	.828	.825	.874

RGB-D methods are indeed better than RGB-only methods provided GT depth. However, only a few RGB-D methods can benefit from the source-free depth. For example, as shown in Table 7, DASNet [88] achieves poorer performance compared to RGB baseline when trained with source-free depth. Similarly, even for one of the SOTA RGB-D models, SPNet [92], the performance on NC4K dataset with RGB-only input is better than with additional source-free depth. Moreover, when provided with source-free depth, none of the existing RGB-D methods outperform the best performing RGB-only method (*e.g.*, ZoomNet [47]) on COD. The same observation was also made on SOD as well. This could be because of the domain gap coupled with the fusion design, among others. We also found it non-trivial to extend the most well-performing RGB-only methods to the RGB-D case. Note that our PopNet performs better than all existing RGB-only and RGB-D methods, with source-free or GT depth maps.

5. Conclusion

We demonstrate a successful case of cross-domain cross-task depth to semantics knowledge transfer using only the source model. In this paper, the source-free depth at the target offered by a given source model is used. The proposed method learns to transfer knowledge from depth to semantics using the objects’ pop-out prior. We facilitate our network to use such prior by designing a novel network archi-

ture. The designed network reasons about the objects by popping them out from the provided depth maps. This process is followed by separating objects from the background using the learned contact surface. We show that the joint learning of object pop-out and contact surface can be successfully supervised using the target semantics. Exhaustive experiments on SOD and COD benchmarks show the successful transfer of depth knowledge to the target, in terms of improved performance and generalization.

Acknowledgement The authors thank the anonymous reviewers and ACs for their tremendous efforts and helpful comments. This research is financed in part by the Conseil Régional de Bourgogne-Franche-Comté, Toyota Motor Europe (research project TRACE-Zurich), the Alexander von Humboldt Foundation, and the Ministry of Education and Science of Bulgaria (support for INSAIT, part of the Bulgarian National Roadmap for Research Infrastructure).

References

- [1] Peshal Agarwal, Danda Pani Paudel, Jan-Nico Zaech, and Luc Van Gool. Unsupervised robust domain adaptation without source data. In *WACV*, 2022. 2
- [2] Amir Atapour-Abarghouei and Toby P Breckon. Real-time monocular depth estimation using synthetic data with domain adaptation via image style transfer. In *IEEE CVPR*, 2018. 1
- [3] Adriano Cardace, Luca De Luigi, Pierluigi Zama Ramirez, Samuele Salti, and Luigi Di Stefano. Plugging self-supervised monocular depth into unsupervised domain adaptation for semantic segmentation. In *WACV*, 2022. 1
- [4] Po-Yi Chen, Alexander H Liu, Yen-Cheng Liu, and Yu-Chiang Frank Wang. Towards scene understanding: Unsupervised monocular depth estimation with semantic-aware representation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019. 4
- [5] Innes C Cuthill, Martin Stevens, Jenna Sheppard, Tracey Maddocks, C Alejandro Párraga, and Tom S Troscianko. Disruptive coloration and background pattern matching. *Nature*, 434(7029):72–74, 2005. 3
- [6] Jiahua Dong, Zhen Fang, Anjin Liu, Gan Sun, and Tongliang Liu. Confident anchor-induced multi-source free domain adaptation. *NeurIPS*, 2021. 1
- [7] Aliya El Nagar, Daniel Osorio, Sarah Zylinski, and Steven M Sait. Visual perception and camouflage response to 3d backgrounds and cast shadows in the european cuttlefish, *sepia officinalis*. *JEB*, 224(11), 2021. 3
- [8] Deng-Ping Fan, Ge-Peng Ji, Ming-Ming Cheng, and Ling Shao. Concealed object detection. *IEEE TPAMI*, 44(10):6024–6042, 2022. 3, 6
- [9] Deng-Ping Fan, Ge-Peng Ji, Guolei Sun, Ming-Ming Cheng, Jianbing Shen, and Ling Shao. Camouflaged object detection. In *IEEE CVPR*, 2020. 3, 6, 7, 8
- [10] 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 TNNLS*, 32(5):2075–2089, 2021. 3, 5, 6, 7, 8
- [11] Deng-Ping Fan, Jing Zhang, Gang Xu, Ming-Ming Cheng, and Ling Shao. Salient objects in clutter. *IEEE TPAMI*, 2023. 3
- [12] Jun Fu, Jing Liu, Haijie Tian, Yong Li, Yongjun Bao, Zhiwei Fang, and Hanqing Lu. Dual attention network for scene segmentation. In *IEEE CVPR*, 2019. 1
- [13] Clément Godard, Oisín Mac Aodha, Michael Firman, and Gabriel J Brostow. Digging into self-supervised monocular depth estimation. In *IEEE ICCV*, 2019. 1, 2, 4
- [14] Yuan-Chen Guo, Di Kang, Linchao Bao, Yu He, and Song-Hai Zhang. Nerfren: Neural radiance fields with reflections. In *IEEE CVPR*, 2022. 4
- [15] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *IEEE CVPR*, 2016. 8
- [16] Ian P Howard. *Perceiving in depth, volume 1: basic mechanisms*. Oxford University Press, 2012. 1
- [17] Junjie Hu, Mete Ozay, Yan Zhang, and Takayuki Okatani. Revisiting single image depth estimation: Toward higher resolution maps with accurate object boundaries. In *WACV*, 2019. 4
- [18] Ge-Peng Ji, Deng-Ping Fan, Yu-Cheng Chou, Dengxin Dai, Alexander Liniger, and Luc Van Gool. Deep gradient learning for efficient camouflaged object detection. *MIR*, 2023. 3
- [19] Wei Ji, Jingjing Li, Shuang Yu, Miao Zhang, Yongri Piao, Shunyu Yao, Qi Bi, Kai Ma, Yefeng Zheng, Huchuan Lu, et al. Calibrated RGB-D salient object detection. In *IEEE CVPR*, 2021. 6, 7
- [20] Wei Ji, Jingjing Li, Miao Zhang, Yongri Piao, and Huchuan Lu. Accurate RGB-D salient object detection via collaborative learning. In *ECCV*, 2020. 3
- [21] Qi Jia, Shuilian Yao, Yu Liu, Xin Fan, Risheng Liu, and Zhongxuan Luo. Segment, magnify and reiterate: Detecting camouflaged objects the hard way. In *IEEE CVPR*, 2022. 3, 6, 7
- [22] Wen-Da Jin, Jun Xu, Qi Han, Yi Zhang, and Ming-Ming Cheng. Cdnet: Complementary depth network for rgb-d salient object detection. *IEEE TIP*, 30:3376–3390, 2021. 3
- [23] Ran Ju, Ling Ge, Wenjing Geng, Tongwei Ren, and Gangshan Wu. Depth saliency based on anisotropic center-surround difference. In *IEEE ICIP*, 2014. 5, 6, 7
- [24] Hongwen Kang, Alexei A Efros, Martial Hebert, and Takeo Kanade. Image composition for object pop-out. In *IEEE ICCVW*, 2009. 2
- [25] Youngeun Kim, Donghyeon Cho, Kyeongtak Han, Priyadarshini Panda, and Sungeun Hong. Domain adaptation without source data. *IEEE TAI*, 2(6):508–518, 2021. 1, 2
- [26] Johannes Kopf, Xuejian Rong, and Jia-Bin Huang. Robust consistent video depth estimation. In *IEEE CVPR*, 2021. 1
- [27] Jogendra Nath Kundu, Akshay R Kulkarni, Suvaansh Bhambrri, Deepesh Mehta, Shreyas Anand Kulkarni, Varun Jampani, and Venkatesh Babu Radhakrishnan. Balancing discriminability and transferability for source-free domain adaptation. In *ICML*, 2022. 1

- [28] Jogendra Nath Kundu, Naveen Venkat, R Venkatesh Babu, et al. Universal source-free domain adaptation. In *IEEE CVPR*, 2020. 1, 2
- [29] Vinod K Kurmi, Venkatesh K Subramanian, and Vinay P Namboodiri. Domain impression: A source data free domain adaptation method. In *WACV*, 2021. 1, 2
- [30] Trung-Nghia Le, Tam V Nguyen, Zhongliang Nie, Minh-Triet Tran, and Akihiro Sugimoto. Anabranch network for camouflaged object segmentation. *CVIU*, 184:45–56, 2019. 5, 7
- [31] Minhyeok Lee, Chaewon Park, Suhwan Cho, and Sangyoun Lee. Spsn: Superpixel prototype sampling network for rgb-d salient object detection. In *ECCV*, 2022. 6, 7
- [32] Aixuan Li, Jing Zhang, Yunqiu Lv, Bowen Liu, Tong Zhang, and Yuchao Dai. Uncertainty-aware joint salient object and camouflaged object detection. In *IEEE CVPR*, 2021. 3, 7
- [33] Gongyang Li, Zhi Liu, Minyu Chen, Zhen Bai, Weisi Lin, and Haibin Ling. Hierarchical alternate interaction network for rgb-d salient object detection. *IEEE TIP*, 30:3528–3542, 2021. 3, 6, 8
- [34] Rui Li, Qianfen Jiao, Wenming Cao, Hau-San Wong, and Si Wu. Model adaptation: Unsupervised domain adaptation without source data. In *IEEE CVPR*, 2020. 2
- [35] Rui Li, Danna Xue, Shaolin Su, Xiantuo He, Qing Mao, Yu Zhu, Jinqiu Sun, and Yanning Zhang. Learning depth via leveraging semantics: Self-supervised monocular depth estimation with both implicit and explicit semantic guidance. *Pattern Recognition (PR)*, page 109297, 2023. 4
- [36] Xinhao Li, Jingjing Li, Lei Zhu, Guoqing Wang, and Zi Huang. Imbalanced source-free domain adaptation. In *ACM MM*, 2021. 1
- [37] Jian Liang, Dapeng Hu, and Jiashi Feng. Do we really need to access the source data? source hypothesis transfer for unsupervised domain adaptation. In *ICML*, 2020. 2
- [38] Jiang-Jiang Liu, Qibin Hou, Ming-Ming Cheng, Jiashi Feng, and Jianmin Jiang. A simple pooling-based design for real-time salient object detection. In *IEEE CVPR*, 2019. 3
- [39] Yuang Liu, Wei Zhang, and Jun Wang. Source-free domain adaptation for semantic segmentation. In *IEEE CVPR*, 2021. 1, 2
- [40] Yukang Lu, Dingyao Min, Keren Fu, and Qijun Zhao. Depth-cooperated trimodal network for video salient object detection. In *ICIP*. IEEE, 2022. 3
- [41] Xuan Luo, Jia-Bin Huang, Richard Szeliski, Kevin Matzen, and Johannes Kopf. Consistent video depth estimation. *ACM TOG*, 39(4):1–13, 2020. 1
- [42] 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 CVPR*, 2021. 3, 6, 7, 8
- [43] Daniel McKee, Zitong Zhan, Bing Shuai, Davide Modolo, Joseph Tighe, and Svetlana Lazebnik. Transfer of representations to video label propagation: Implementation factors matter. *arXiv preprint arXiv:2203.05553*, 2022. 6
- [44] Haiyang Mei, Ge-Peng Ji, Ziqi Wei, Xin Yang, Xiaopeng Wei, and Deng-Ping Fan. Camouflaged object segmentation with distraction mining. In *IEEE CVPR*, 2021. 7
- [45] S Mahdi H Miangoleh, Sebastian Dille, Long Mai, Sylvain Paris, and Yagiz Aksoy. Boosting monocular depth estimation models to high-resolution via content-adaptive multi-resolution merging. In *IEEE CVPR*, 2021. 1, 4
- [46] Yuzhen Niu, Yujie Geng, Xueqing Li, and Feng Liu. Leveraging stereopsis for saliency analysis. In *IEEE CVPR*, 2012. 5, 6
- [47] Youwei Pang, Xiaoqi Zhao, Tian-Zhu Xiang, Lihe Zhang, and Huchuan Lu. Zoom in and out: A mixed-scale triplet network for camouflaged object detection. In *IEEE CVPR*, 2022. 3, 6, 7, 8
- [48] Houwen Peng, Bing Li, Weihua Xiong, Weiming Hu, and Rongrong Ji. RGBD salient object detection: a benchmark and algorithms. In *ECCV*, 2014. 5, 6
- [49] Yongri Piao, Zhengkun Rong, Miao Zhang, Weisong Ren, and Huchuan Lu. A2dele: Adaptive and attentive depth distiller for efficient RGB-D salient object detection. In *IEEE CVPR*, 2020. 3
- [50] René Ranftl, Alexey Bochkovskiy, and Vladlen Koltun. Vision transformers for dense prediction. In *IEEE ICCV*, 2021. 1, 4
- [51] 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 TPAMI*, 44(3):1623–1637, 2022. 1, 2, 4
- [52] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *MICCAI*, 2015. 4
- [53] Suman Saha, Anton Obukhov, Danda Pani Paudel, Menelaos Kanakis, Yuhua Chen, Stamatios Georgoulis, and Luc Van Gool. Learning to relate depth and semantics for unsupervised domain adaptation. In *IEEE CVPR*, 2021. 1
- [54] Roshni Sahoo, Divya Shanmugam, and John Guttag. Unsupervised domain adaptation in the absence of source data. *arXiv preprint arXiv:2007.10233*, 2020. 2
- [55] Johannes L Schonberger and Jan-Michael Frahm. Structure-from-motion revisited. In *IEEE CVPR*, 2016. 1
- [56] Nathan Silberman, Derek Hoiem, Pushmeet Kohli, and Rob Fergus. Indoor segmentation and support inference from RGBD images. In *ECCV*, 2012. 1
- [57] Przemysław Skurowski, Hassan Abdulameer, J Błaszczczyk, Tomasz Depta, Adam Kornacki, and P Kozieł. Animal camouflage analysis: Chameleon database. *Unpublished manuscript*, 2(6):7, 2018. 5, 7
- [58] Hwanjun Song, Eunyoung Kim, Varun Jampan, Deqing Sun, Jae-Gil Lee, and Ming-Hsuan Yang. Exploiting scene depth for object detection with multimodal transformers. In *BMVC*, 2021. 3
- [59] Yujia Sun, Geng Chen, Tao Zhou, Yi Zhang, and Nian Liu. Context-aware cross-level fusion network for camouflaged object detection. In *IJCAI*, 2021. 3, 7
- [60] Yujia Sun, Shuo Wang, Chenglizhao Chen, and Tian-Zhu Xiang. Boundary-guided camouflaged object detection. In *IJCAI*, 2022. 3
- [61] Alessio Tonioni, Matteo Poggi, Stefano Mattoccia, and Luigi Di Stefano. Unsupervised domain adaptation for depth prediction from images. *IEEE TPAMI*, 42(10):2396–2409, 2019. 1

- [62] Anne Treisman. Preattentive processing in vision. *ICVGIP*, 31(2):156–177, 1985. 2
- [63] Fengyun Wang, Jinshan Pan, Shoukun Xu, and Jinhui Tang. Learning discriminative cross-modality features for RGB-D saliency detection. *IEEE TIP*, 31:1285–1297, 2022. 6, 7
- [64] Qin Wang, Olga Fink, Luc Van Gool, and Dengxin Dai. Continual test-time domain adaptation. In *IEEE CVPR*, 2022. 2
- [65] Cho-Ying Wu, Jialiang Wang, Michael Hall, Ulrich Neumann, and Shuochen Su. Toward practical monocular indoor depth estimation. In *IEEE CVPR*, 2022. 1
- [66] Yu-Huan Wu, Yun Liu, Jun Xu, Jia-Wang Bian, Yu-Chao Gu, and Ming-Ming Cheng. Mobilesal: Extremely efficient rgb-d salient object detection. *IEEE TPAMI*, 44(12):10261–10269, 2022. 3
- [67] Zongwei Wu, Guillaume Allibert, Christophe Stolz, and Cédric Demonceaux. Depth-adapted CNN for RGB-D cameras. In *ACCV*, 2020. 1
- [68] Zongwei Wu, Guillaume Allibert, Christophe Stolz, Chao Ma, and Cédric Demonceaux. Modality-guided subnetwork for salient object detection. In *3DV*, 2021. 1, 3
- [69] Zongwei Wu, Shriarulmohzivarman Gobichettipalayam, Brahim Tamadazte, Guillaume Allibert, Danda Pani Paudel, and Cédric Demonceaux. Robust RGB-D fusion for saliency detection. *3DV*, 2022. 3, 6
- [70] Zhe Wu, Li Su, and Qingming Huang. Cascaded partial decoder for fast and accurate salient object detection. In *IEEE CVPR*, 2019. 3
- [71] Haifeng Xia, Handong Zhao, and Zhengming Ding. Adaptive adversarial network for source-free domain adaptation. In *IEEE ICCV*, 2021. 1, 2
- [72] Ke Xian, Jianming Zhang, Oliver Wang, Long Mai, Zhe Lin, and Zhiguo Cao. Structure-guided ranking loss for single image depth prediction. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020. 4
- [73] Mochu Xiang, Jing Zhang, Yunqiu Lv, Aixuan Li, Yiran Zhong, and Yuchao Dai. Exploring depth contribution for camouflaged object detection. *arXiv e-prints*, pages arXiv–2106, 2021. 1, 6
- [74] 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 ICCV*, 2021. 3, 7
- [75] Shiqi Yang, Joost van de Weijer, Luis Herranz, Shangling Jui, et al. Exploiting the intrinsic neighborhood structure for source-free domain adaptation. *NeurIPS*, 2021. 1
- [76] Shiqi Yang, Yaxing Wang, Joost van de Weijer, and Luis Herranz. Unsupervised domain adaptation without source data by casting a bait. *arXiv preprint arXiv:2010.12427*, 2020. 2
- [77] Shiqi Yang, Yaxing Wang, Joost van de Weijer, Luis Herranz, and Shangling Jui. Generalized source-free domain adaptation. In *IEEE ICCV*, 2021. 1, 2
- [78] Qiang Zhai, Xin Li, Fan Yang, Chenglizhao Chen, Hong Cheng, and Deng-Ping Fan. Mutual graph learning for camouflaged object detection. In *IEEE CVPR*, 2021. 7
- [79] Wei Zhai, Yang Cao, HaiYong Xie, and Zheng-Jun Zha. Deep texton-coherence network for camouflaged object detection. *IEEE TMM*, 2022. 3
- [80] Chen Zhang, Runmin Cong, Qinwei Lin, Lin Ma, Feng Li, Yao Zhao, and Sam Kwong. Cross-modality discrepant interaction network for RGB-D salient object detection. In *ACM MM*, 2021. 3, 7
- [81] Jing Zhang, Deng-Ping Fan, Yuchao Dai, Xin Yu, Yiran Zhong, Nick Barnes, and Ling Shao. RGB-D saliency detection via cascaded mutual information minimization. In *IEEE ICCV*, 2021. 7
- [82] Jing Zhang, Yunqiu Lv, Mochu Xiang, Aixuan Li, Yuchao Dai, and Yiran Zhong. Depth-guided camouflaged object detection. *CoRR*, abs/2106.13217, 2021. 3
- [83] Miao Zhang, Shuang Xu, Yongri Piao, Dongxiang Shi, Shusen Lin, and Huchuan Lu. Preynet: Preying on camouflaged objects. In *ACM MM*, 2022. 3, 6
- [84] Pingping Zhang, Dong Wang, Huchuan Lu, Hongyu Wang, and Baocai Yin. Learning uncertain convolutional features for accurate saliency detection. In *IEEE ICCV*, 2017. 3
- [85] Wenbo Zhang, Ge-Peng Ji, Zhuo Wang, Keren Fu, and Qijun Zhao. Depth quality-inspired feature manipulation for efficient RGB-D salient object detection. In *ACM MM*, 2021. 6
- [86] Zhoutong Zhang, Forrester Cole, Richard Tucker, William T Freeman, and Tali Dekel. Consistent depth of moving objects in video. *ACM TOG*, 40(4):1–12, 2021. 1
- [87] Zhao Zhang, Zheng Lin, Jun Xu, Wen-Da Jin, Shao-Ping Lu, and Deng-Ping Fan. Bilateral attention network for rgb-d salient object detection. *IEEE TIP*, 30:1949–1961, 2021. 6
- [88] Jiawei Zhao, Yifan Zhao, Jia Li, and Xiaowu Chen. Is depth really necessary for salient object detection? In *ACM MM*, 2020. 3, 8
- [89] Shanshan Zhao, Huan Fu, Mingming Gong, and Dacheng Tao. Geometry-aware symmetric domain adaptation for monocular depth estimation. In *IEEE CVPR*, 2019. 1
- [90] Yijie Zhong, Bo Li, Lv Tang, Senyun Kuang, Shuang Wu, and Shouhong Ding. Detecting camouflaged object in frequency domain. In *IEEE CVPR*, 2022. 3
- [91] Tao Zhou, Deng-Ping Fan, Ming-Ming Cheng, Jianbing Shen, and Ling Shao. Rgb-d salient object detection: A survey. *CVMJ*, pages 1–33, 2021. 1
- [92] Tao Zhou, Huazhu Fu, Geng Chen, Yi Zhou, Deng-Ping Fan, and Ling Shao. Specificity-preserving RGB-D saliency detection. In *IEEE ICCV*, 2021. 3, 5, 6, 7, 8
- [93] Zikun Zhou, Wenjie Pei, Xin Li, Hongpeng Wang, Feng Zheng, and Zhenyu He. Saliency-associated object tracking. In *IEEE ICCV*, 2021. 3
- [94] Hongwei Zhu, Peng Li, Haoran Xie, Xuefeng Yan, Dong Liang, Dapeng Chen, Mingqiang Wei, and Jing Qin. I can find you! boundary-guided separated attention network for camouflaged object detection. In *AAAI*, 2022. 3
- [95] Jinchao Zhu, Xiaoyu Zhang, Shuo Zhang, and Junnan Liu. Inferring camouflaged objects by texture-aware interactive guidance network. In *AAAI*, 2021. 3