

Zoom In and Out: A Mixed-scale Triplet Network for Camouflaged Object Detection

Youwei Pang^{1†}, Xiaoqi Zhao^{1†}, Tian-Zhu Xiang³, Lihe Zhang^{1*} and Huchuan Lu^{1,2}

¹Dalian University of Technology, China ²Peng Cheng Laboratory, China

³Inception Institute of Artificial Intelligence, UAE

{lartpang, zxq}@mail.dlut.edu.cn, tianzhu.xiang19@gmail.com, {zhanglihe, lhchuan}@dlut.edu.cn

Abstract

The recently proposed camouflaged object detection (COD) attempts to segment objects that are visually blended into their surroundings, which is extremely complex and difficult in real-world scenarios. Apart from high intrinsic similarity between the camouflaged objects and their background, the objects are usually diverse in scale, fuzzy in appearance, and even severely occluded. To deal with these problems, we propose a mixed-scale triplet network, **ZoomNet**, which mimics the behavior of humans when observing vague images, i.e., zooming in and out. Specifically, our ZoomNet employs the zoom strategy to learn the discriminative mixed-scale semantics by the designed scale integration unit and hierarchical mixed-scale unit, which fully explores imperceptible clues between the candidate objects and background surroundings. Moreover, considering the uncertainty and ambiguity derived from indistinguishable textures, we construct a simple yet effective regularization constraint, uncertainty-aware loss, to promote the model to accurately produce predictions with higher confidence in candidate regions. Without bells and whistles, our proposed highly task-friendly model consistently surpasses the existing 23 state-of-the-art methods on four public datasets. Besides, the superior performance over the recent cutting-edge models on the SOD task also verifies the effectiveness and generality of our model. The code will be available at <https://github.com/lartpang/ZoomNet>.

1. Introduction

Camouflaged objects are often “seamlessly” integrated into the environment by changing their appearance, coloration or pattern to avoid detection, such as chameleons, cuttlefishes and flatfishes. This is mainly due to their self-protection mechanism in the harsh living environment.

[†]These authors contributed equally to this work.

^{*}Corresponding author.

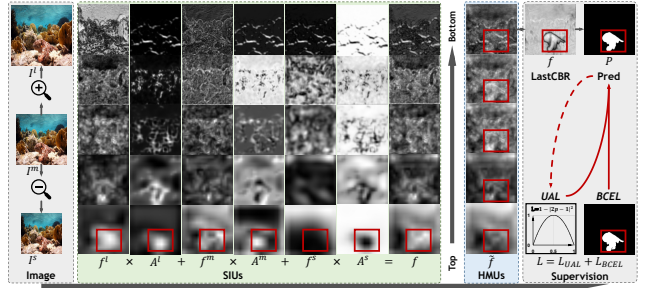


Figure 1. Illustration of ZoomNet. Based on *zoom* strategy, our model distills the *differentiated* features at different “zoom” scales. Then we design SIUs to screen and aggregate scale-specific features, and HMUs to reorganize and enhance mixed-scale features. Under the supervision of BCEL and the proposed UAL, the model produces the accurate and reliable camouflaged object prediction. Note that BCEL is computed based on ground truth while UAL is not. f : feature map; A : attention map. LastCBR: the last “Conv-BN-ReLU” layer before the prediction. $l/m/s$: Different input scales. The whiter region denotes the larger activation response.

Broadly speaking, camouflaged objects also refer to the objects that are extremely small in size, highly similar to the background, or heavily obscured. They subtly hide themselves in the surroundings, making them difficult to be found, e.g., soldiers wearing camouflaged uniforms and lions hiding in the grass. Camouflaged object detection (COD) is far more complex and challenging than traditional salient object detection or other object segmentation. Recently, it has attracted ever-growing research interest from the computer vision community and facilitates many valuable real-life applications, such as search and rescue [8], species discovery [36], and medical image analysis [10, 11, 68].

Recently, numerous deep learning-based methods have been proposed and achieved significant progress. Nevertheless, they are still struggled to accurately and reliably detect camouflaged objects, due to visual insignificance of camouflaged objects, and high diversity in scale, appearance

and occlusion. By observing our experiments, it is found that the current COD detectors are susceptible to distractors from background surroundings. Thus it is difficult to excavate discriminative and subtle semantic cues for camouflaged objects, resulting in the inability to clearly segment the camouflaged objects from the chaotic background and the predictions of some uncertain (low-confidence) regions. Taking these into mind, in this paper, we summarize the COD issue into two aspects: 1) *How to accurately locate camouflaged objects under conditions of inconspicuous appearance and various scales?* 2) *How to suppress the obvious interference from the background and infer camouflaged objects more reliably?* Intuitively, to accurately find the vague or camouflaged objects in the scene, humans may try to refer to and compare the changes in the shape or appearance at different scales by zooming in and out (re-scaling) the image. This specific behavior pattern of human beings motivates us to identify camouflaged objects by mimicking the zooming in and out strategy.

With this inspiration, in this paper, we propose a mixed-scale triplet network, *ZoomNet*, which significantly improves the existing camouflaged object detection performance. **Firstly**, for accurate object location, we employ scale space theory [20,21,48] to imitate zooming in and out strategy. Specifically, we design two key modules, *i.e.*, the scale integration unit (SIU) and the hierarchical mixed-scale unit (HMU). As shown in Fig. 1, our model extracts differentiated camouflaged object features at different “zoom” scales using the triplet architecture, then adopts SIUs to screen and aggregate scale-specific features, and utilizes HMUs to further reorganize and enhance mixed-scale features. Thus, our model is able to mine the accurate and subtle semantic clues between objects and background under the mixed scales, and produce accurate predictions. Besides, we use the shared weight strategy, which achieves a good balance of efficiency and effectiveness. **Secondly**, it is related to reliable prediction in complex scenarios. Although the object is accurately located, the indistinguishable texture and background will easily bring negative effects to the model learning, *e.g.* predicting uncertain/ambiguity regions, which greatly reduces the detection performance and cannot be ignored. This can be seen in Fig. 6 (Row 3 and 4) and Fig. 1 in the *supp.* To this end, we design an uncertainty-aware loss (UAL) to guide the model training, which is only based on the prior knowledge that a good COD prediction should have a clear polarization trend. Its GT-independent characteristic makes it suitable for enhancing the GT-based BCE loss. This targeted enhancement strategy can force the network to optimize the prediction of the uncertain regions during the training process, enabling our *ZoomNet* to distinguish the uncertain regions and segment the camouflaged objects reliably.

Our contributions can be summarized as follows: 1)

For the COD task, we propose a mixed-scale triplet network, *ZoomNet*, which can credibly capture the objects in complex scenes by characterizing and unifying the scale-specific appearance features at different “zoom” scales and the purposeful optimization strategy. 2) To obtain the discriminative feature representation of camouflaged objects, we design SIUs and HMUs to distill, aggregate and strengthen the scale-specific and subtle semantic representation for accurate COD. 3) We propose a simple yet effective optimization enhancement strategy, UAL, which can significantly suppress the uncertainty and interference from the background without increasing extra parameters. 4) Our model greatly surpasses recent 23 state-of-the-art methods under seven metrics on four COD datasets. Furthermore, it shows good generalization in the SOD task and the superior performance compared with the existing SOD methods.

2. Related Work

Camouflaged Object. The study of camouflage has a long history in biology. This behavior of creatures in nature can be regarded as the result of natural selection and adaptation. In fact, in human life and other parts of society, it also has a profound impact, *e.g.*, arts, popular culture, and design. More details can be found in [42]. In the field of computer vision, research on camouflaged objects is often associated with salient object detection (SOD), which mainly deals with those salient and easily observed objects in the scene. In general, saliency models are designed for the general observation paradigm (*i.e.*, finding visually prominent objects). They are not suitable for the specific observation (*i.e.*, finding concealed objects). Therefore, it is necessary to establish models based on the essential requirements and specific data of the task to learn the special knowledge.

Camouflaged Object Detection (COD). Different from the traditional SOD task, the COD pays more attention to the undetectable objects (mainly because of too small size, occlusion, concealment or self-disguise). Due to the differences in the attributes of the objects of interest, the goals of the two tasks are different. The difficulty and complexity of the COD far exceed the SOD due to the high similarity between the object and the environment. Some valuable attempts have been made in recent years. Recent works [16,29,59] construct the multi-task learning framework in the prediction process of camouflaged objects and introduce some auxiliary tasks like classification and edge detection. Some uncertainty-aware methods [17,56] are proposed to model and cope with the uncertainty in data annotation or COD data itself. In the other two methods [31,43], contextual feature learning also plays an important role. There are also a number of bio-inspired methods, such as [9,53]. They capture camouflaged objects by imitating the behavior process of hunters or changing the viewpoint of the scene. Although our method can also be at-

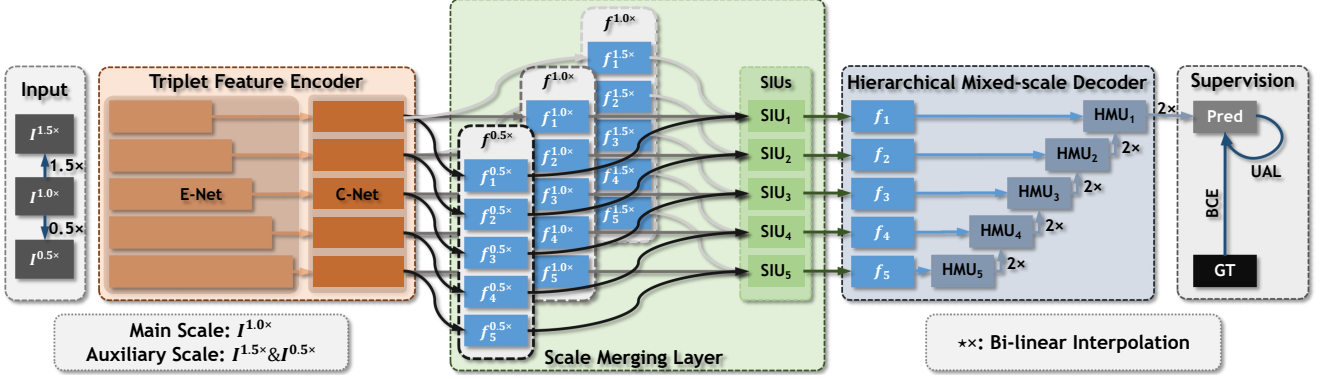


Figure 2. Overall framework. The shared triplet feature encoder is used to extract multi-level features corresponding to different input “zoom” scales, which is composed of E-Net and C-Net for extracting and compressing features, respectively. At different levels of the scale merging layer, SIUs are adopted to screen and aggregate the critical cues from different scales. Then the fused features are gradually integrated through the top-down up-sampling path in the hierarchical mixed-scale decoder. HMUs further enhance the feature discrimination by constructing a multi-path structure inside the features. Finally, a probability map of the camouflaged object corresponding to the input image can be obtained. In the training stage, the binary cross entropy and the proposed UAL are used as the loss function.

tributed to the last category, ours is different from the above methods. Our method simulates the behavior of humans to understand complex images by zooming in and out strategy. The proposed method explores the scale-specific and imperceptible semantic features under the mixed scales for accurate predictions, with the supervision of BCE and our proposed uncertainty-aware loss. Accordingly, our method achieves a more comprehensive understanding of the scene, and accurately and robustly segments the camouflaged objects from the complex background, which even can be transferred to the SOD task effectively and smoothly.

Scale Space Integration. The scale-space theory aims to promote an optimal understanding of image structure, which is an extremely effective and theoretically sound framework for addressing naturally occurring scale variations. Its ideas have been widely used in computer vision, including the image pyramid [2] and the feature pyramid [19]. Due to the structural and semantic differences at different scales, the corresponding features play different roles. However, the commonly-used inverted pyramid-like feature extraction structures [13, 38, 65] often cause the feature representation to lose too much texture and appearance details, which are unfavorable for dense prediction tasks [27, 39] that emphasize the integrity of regions and edges. Thus, some recent CNN-based COD methods [9, 31, 43, 59] and SOD methods [15, 23, 34, 35, 66, 67, 69] explore the combination strategy of inter-layer features to enhance the feature representation. These bring some positive gains for accurate localization and segmentation of objects. However, for the COD task, the existing approaches overlook the performance bottleneck caused by the ambiguity of the structural information of the data itself that makes it difficult to be fully perceived at a single scale.

Different from them, we mimic the zoom strategy to synchronously consider differentiated relationships between object and background at multiple scales, thereby fully perceiving the camouflaged objects and confusing scenes. Besides, we also further explore the fine-grained feature scale space between channels.

3. Proposed Method

In this section, we first elaborate on the overall architecture of the proposed ZoomNet, and then present the details of each module and the uncertainty-aware loss.

3.1. Overall Architecture

The overall architecture of the proposed ZoomNet is illustrated in Fig. 2. Inspired by the zoom strategy from human beings when observing confusing scenes, we argue that different zoom scales often contain their specific information. Aggregating the differentiated information on different scales will benefit exploring the inconspicuous yet valuable clues from confusing scenarios, thus facilitating COD. To implement it, intuitively, we resort to the image pyramid. Specifically, we customize an image pyramid based on the single scale input to identify the camouflaged objects. The scales are divided into a main scale (*i.e.* the input scale) and two auxiliary scales. The latter is obtained by re-scaling the former to imitate the operation of zooming in and out. We utilize the shared triplet feature encoder to extract features on different scales and feed them to the scale merging layer. To integrate these features that contain rich scale-specific information, we design a series of scale integration units (SIUs) based on the attention-aware filtering mechanism. Thus, these auxiliary scales are integrated into the main scale, *i.e.*, information aggregation of “zoom in and out” op-

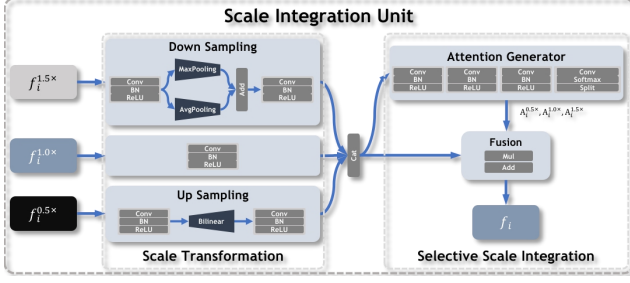


Figure 3. Illustration of the scale integration unit (SIU).

eration. This will largely enhance the model to distill critical and informative semantic cues for capturing difficult-to-detect camouflaged objects. After that, we construct hierarchical mixed-scale units (HMUs) to gradually integrate multi-level features in a top-down manner to enhance the mixed-scale feature representation. It further increases the receptive field range and diversifies feature representation within the module. The captured fine-grained and mixed-scale clues promote the model to accurately segment the camouflaged objects in the chaotic scenes. Besides, to overcome the uncertainty in the prediction caused by the inherent complexity of the data, we design an uncertainty-aware loss (UAL) to assist the BCE loss, enabling the model to distinguish these uncertain regions and produce an accurate and reliable prediction.

3.2. Triplet Feature Encoder

We start by extracting deep features through a shared triplet feature encoder for the group-wise inputs, which consists of the feature extraction and the channel compression networks, *i.e.* E-Net and C-Net. For the trade-off between efficiency and effectiveness, the main scale and the two auxiliary scales are empirically set to 1.0 \times , 1.5 \times and 0.5 \times . E-Net is constituted by the commonly-used ResNet-50 [14] that is removed the structure after “layer4”. C-Net is cascaded to further optimize computation and obtain a more compact feature. For more details about it, please see the *supp*. Thus, three sets of 64-channel feature maps corresponding to three input scales are produced, *i.e.*, $\{f_i^k\}_{k=1}^5$, $k \in \{0.5, 1.0, 1.5\}$. Next, these features are fed successively to the scale merging layer and the hierarchical mixed-scale decoder for subsequent processing.

3.3. Scale Merging Layer

We design an attention-based SIU to screen (weight) and combine scale-specific information, as shown in Fig. 3. Several such units make up the scale merging layer. Through filtering and aggregation, the expression of different scales is self-adaptively highlighted. Before scale integration, the features $f_i^{1.5}$ and $f_i^{0.5}$ are first resized to be consistent resolution with the main scale feature $f_i^{1.0}$. Specif-

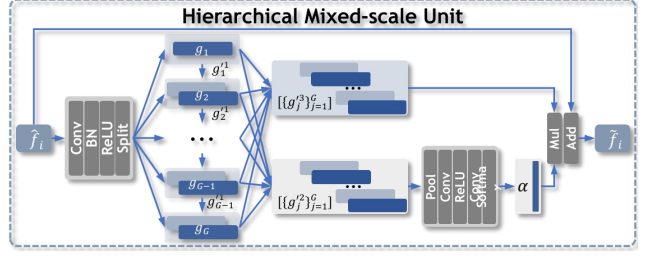


Figure 4. Hierarchical mixed-scale unit (HMu). We adopt group-wise interaction and channel-wise modulation to explore the discriminative and valuable semantics from different channels. Note that each group of features is executed sequentially from top to bottom. The latter one integrates part of the features of the previous one before the feature transformation.

ically, for $f_i^{1.5}$, we use a hybrid structure of “max-pooling + average-pooling” to down-sample it, which helps to preserve the effective and diverse responses for camouflaged objects in high-resolution features. For $f_i^{0.5}$, we directly up-sample it by the bi-linear interpolation. Then, these features are fed into the “attention generator”, and a three-channel feature map is calculated through a series of convolutional layers. After a softmax activation layer, the attention map A^k ($k \in \{0.5, 1.0, 1.5\}$) corresponding to each scale can be obtained and used as respective weights for the final integration. The process is formulated as:

$$\begin{aligned} A_i &= \text{softmax}(\Psi([\mathcal{U}(f_i^{0.5}), f_i^{1.0}, \mathcal{D}(f_i^{1.5})], \phi)), \\ f_i &= A_i^{0.5} \cdot \mathcal{U}(f_i^{0.5}) + A_i^{1.0} \cdot f_i^{1.0} + A_i^{1.5} \cdot \mathcal{D}(f_i^{1.5}), \end{aligned} \quad (1)$$

where $\Psi(\star, \phi)$ indicates the stacked “Conv-BN-ReLU” layers in the attention generator, and ϕ means the parameters of these layers. $[\star]$ represents the concatenation operation. \mathcal{D} and \mathcal{U} refer to the hybrid pooling and bi-linear interpolation operations mentioned above, respectively. Note that some operations before and after the sampling operation are not shown in Equ. 1 for simplicity but can be seen in Fig. 3. These designs aim to selectively aggregate the scale-specific information to explore subtle but critical semantic cues at different scales, boosting the feature representation.

3.4. Hierarchical Mixed-scale Decoder

After SIUs, the auxiliary-scale information is integrated into the main-scale branch. Similar to the multi-scale case, different channels also contain differentiated semantics. Thus, it is necessary to excavate valuable clues contained in different channels. To this end, we design HMUs to conduct information interaction and feature refinement between channels, which strengthen features from coarse-grained group-wise iteration to fine-grained channel-wise modulation in the decoder, as depicted in Fig. 4. The input \hat{f}_i of the HMu_i contains the multi-scale fused feature

Table 1. Comparisons of different methods on COD datasets. The best three results are highlighted in **red**, **green** and **blue**. “—”: Not available; *: Using more datasets.

Model	CAMO					CHAMELEON					COD10K					NC4K				
	$S_m \uparrow$	$F_{\beta}^w \uparrow$	MAE \downarrow	$F_{\beta} \uparrow$	$E_m \uparrow$	$S_m \uparrow$	$F_{\beta}^w \uparrow$	MAE \downarrow	$F_{\beta} \uparrow$	$E_m \uparrow$	$S_m \uparrow$	$F_{\beta}^w \uparrow$	MAE \downarrow	$F_{\beta} \uparrow$	$E_m \uparrow$	$S_m \uparrow$	$F_{\beta}^w \uparrow$	MAE \downarrow	$F_{\beta} \uparrow$	$E_m \uparrow$
Salient Object Detection / Medical Image Segmentation																				
NLDF [28]	0.665	0.495	0.123	0.564	0.790	0.798	0.652	0.063	0.714	0.893	0.701	0.473	0.059	0.539	0.819	0.738	0.586	0.083	0.656	0.846
PiCANet [24]	0.701	0.510	0.125	0.573	0.799	0.765	0.552	0.085	0.618	0.846	0.696	0.415	0.081	0.489	0.788	0.758	0.570	0.088	0.640	0.835
BASNet [37]	0.615	0.434	0.124	0.503	0.727	0.847	0.771	0.044	0.795	0.894	0.661	0.432	0.071	0.486	0.749	0.695	0.546	0.095	0.610	0.785
CPD [50]	0.716	0.556	0.113	0.618	0.796	0.857	0.731	0.048	0.771	0.923	0.750	0.531	0.053	0.595	0.853	0.787	0.645	0.072	0.705	0.866
PoolNet [23]	0.730	0.575	0.105	0.643	0.819	0.845	0.690	0.054	0.749	0.933	0.740	0.506	0.056	0.575	0.844	0.785	0.635	0.073	0.699	0.865
EGNet [63]	0.732	0.604	0.109	0.670	0.820	0.797	0.649	0.065	0.702	0.884	0.736	0.517	0.061	0.582	0.854	0.777	0.639	0.075	0.696	0.864
F3Net [46]	0.711	0.564	0.109	0.616	0.780	0.848	0.744	0.047	0.770	0.917	0.739	0.544	0.051	0.593	0.819	0.780	0.656	0.070	0.705	0.848
SCRN [51]	0.779	0.643	0.090	0.705	0.850	0.876	0.741	0.042	0.787	0.939	0.789	0.575	0.047	0.651	0.880	0.830	0.698	0.059	0.757	0.897
CSNet [12]	0.771	0.641	0.092	0.705	0.849	0.856	0.718	0.047	0.766	0.928	0.778	0.569	0.047	0.634	0.871	0.750	0.603	0.088	0.655	0.793
SSAL [61]	0.644	0.493	0.126	0.579	0.780	0.757	0.639	0.071	0.702	0.856	0.668	0.454	0.066	0.527	0.789	0.699	0.561	0.093	0.644	0.812
UCNet [60]	0.739	0.640	0.094	0.700	0.820	0.880	0.817	0.036	0.836	0.941	0.776	0.633	0.042	0.681	0.867	0.811	0.729	0.055	0.775	0.886
MINet [35]	0.748	0.637	0.090	0.691	0.838	0.855	0.771	0.036	0.802	0.937	0.770	0.608	0.042	0.657	0.859	0.812	0.720	0.056	0.764	0.887
ITSD [72]	0.750	0.610	0.102	0.663	0.830	0.814	0.662	0.057	0.705	0.901	0.767	0.557	0.051	0.615	0.861	0.811	0.679	0.064	0.729	0.883
PraNet [10]	0.769	0.663	0.094	0.710	0.837	0.860	0.763	0.044	0.789	0.935	0.789	0.629	0.045	0.671	0.879	0.822	0.724	0.059	0.763	0.888
Camouflaged Object Detection																				
ANet_SRM [16]	0.682	0.484	0.126	0.541	0.722	—	—	—	—	—	—	—	—	—	—	—	—	—	—	—
SINet [9]	0.745	0.644	0.092	0.702	0.829	0.872	0.806	0.034	0.827	0.946	0.776	0.631	0.043	0.679	0.874	0.808	0.723	0.058	0.769	0.883
SLSR [29]	0.787	0.696	0.080	0.744	0.854	0.890	0.822	0.030	0.841	0.948	0.804	0.673	0.037	0.715	0.892	0.840	0.766	0.048	0.804	0.907
MGL-R [59]	0.775	0.673	0.088	0.726	0.842	0.893	0.812	0.031	0.833	0.941	0.814	0.666	0.035	0.710	0.890	0.833	0.739	0.053	0.782	0.893
PFNet [31]	0.782	0.695	0.085	0.746	0.855	0.882	0.810	0.033	0.828	0.945	0.800	0.660	0.040	0.701	0.890	0.829	0.745	0.053	0.784	0.898
UJSC* [17]	0.800	0.728	0.073	0.772	0.873	0.891	0.833	0.030	0.847	0.955	0.809	0.684	0.035	0.721	0.891	0.842	0.771	0.047	0.806	0.907
MirrorNet [53]	0.785	0.719	0.077	0.754	0.850	—	—	—	—	—	—	—	—	—	—	—	—	—	—	—
C ² FNet [43]	0.796	0.719	0.080	0.762	0.864	0.888	0.828	0.032	0.844	0.946	0.813	0.686	0.036	0.723	0.900	0.838	0.762	0.049	0.795	0.904
UGTR [56]	0.784	0.684	0.086	0.735	0.851	0.888	0.794	0.031	0.819	0.940	0.817	0.666	0.036	0.711	0.890	0.839	0.746	0.052	0.787	0.899
Ours	0.820	0.752	0.066	0.794	0.892	0.902	0.845	0.023	0.864	0.958	0.838	0.729	0.029	0.766	0.911	0.853	0.784	0.043	0.818	0.912

f_i from the SIU_i and the feature \tilde{f}_{i+1} from the HMU_{i+1} :

$$\hat{f}_i = f_i + \mathcal{U}(\tilde{f}_{i+1}). \quad (2)$$

Group-wise Iteration. We adopt 1×1 convolution to extend the channel number of feature map \hat{f}_i . The features are then divided into G groups $\{g_j\}_{j=1}^G$ along the channel dimension. Feature interaction between groups is carried out in an iterative manner. Specifically, the first group $\{g_1\}$ is split into three feature sets $\{g_1^k\}_{k=1}^3$ after a convolution block. Among them, the g_1^1 is adopted for information exchange with the next group, and the other two are used for channel-wise modulation. In the j^{th} ($1 < j < G$) group, the feature g_j is concatenated with the feature g_{j-1}^1 from the previous group along the channel, followed by a convolution block and a split operation, which similarly divides this feature group into three feature sets. It is noted that the output of the group G with the similar input form to the previous groups only contains g_G^2 and g_G^3 . Such an iterative mixing strategy strives to learn the critical clues from different channels and obtain a powerful feature representation. From another perspective, the iterative structure in HMU can be equivalent to a kernel pyramid structure.

Channel-wise Modulation. The features $\{g_j^2\}_{j=1}^G$ are concatenated and converted into the feature modulation vector α by a small convolutional network, which is employed to weight another concatenated feature $\{g_j^3\}_{j=1}^G$. The weighted feature is then processed by a convolutional layer, which is defined as:

$$\tilde{f}_i = \mathcal{A}(\hat{f}_i + \mathcal{N}(\mathcal{T}(\alpha \cdot \{g_j^3\}_{j=1}^G))), \quad (3)$$

where \mathcal{A} , \mathcal{N} and \mathcal{T} represent the activation layer, the normalization layer and the convolutional layer, respectively.

Based on five cascaded HMUs and several stacked convolutional layers, a single-channel logits map is obtained. The final confidence map \mathbf{P} that highlights the camouflaged objects is then generated by a sigmoid function.

3.5. Loss Functions

The binary cross entropy loss (BCEL) is widely used in various image segmentation tasks and its mathematical form is $l_{BCEL}^{i,j} = -\mathbf{g}_{i,j} \log \mathbf{p}_{i,j} - (1 - \mathbf{g}_{i,j}) \log(1 - \mathbf{p}_{i,j})$, where $\mathbf{g}_{i,j} \in \{0, 1\}$ and $\mathbf{p}_{i,j} \in [0, 1]$ denote the ground truth and the predicted value at position (i, j) , respectively. As shown in Fig. 6, due to the complexity of the COD data, if trained only under the BCEL, the model produces serious ambiguity and uncertainty in the prediction and fails to accurately capture objects, of which both will reduce the reliability of COD. To force the model to enhance “confidence” in decision-making and increase the penalty for fuzzy prediction, we design a strong constraint as the auxiliary of the BCEL, *i.e.*, the uncertainty-aware loss (UAL).

In the final probability map of the camouflaged object, the pixel value range is $[0, 1]$, where 0 means the pixel belongs to the background, and 1 means it belongs to the camouflaged object. Therefore, the closer the predicted value is to 0.5, the more uncertain the determination about the property of the pixel is. To optimize it, a direct way is to use the ambiguity as the supplementary loss for these difficult samples. To this end, we first need to define the ambiguity measure of the pixel x , which maximizes at $x = 0.5$ and

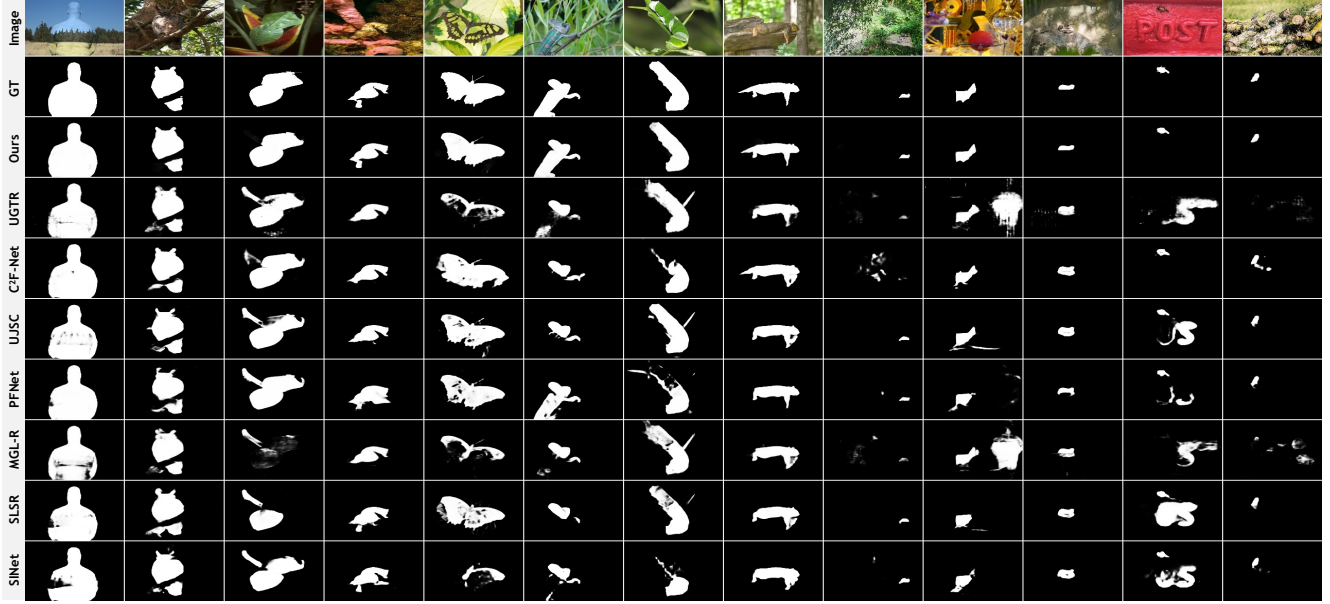


Figure 5. Visual comparisons of some recent COD methods and ours on different types of samples. Please zoom in for more details.

minimizes at $x = 0$ or $x = 1$. And as a loss, the function should be smooth and continuous with only a finite number of non-differentiable points. For brevity, we empirically consider two forms, $\Phi_{pow}^\alpha(x) = 1 - |2x - 1|^\alpha$ based on the power function and $\Phi_{exp}^\alpha(x) = e^{-(\alpha(x-0.5))^2}$ based on the exponential function. Besides, inspired by the form of the weighted BCE loss, we also try to use $\omega = 1 + \Phi_{pow}^2(x)$ as the weight of BCE loss to increase the loss of hard pixels. After massive experiments (Sec. 4.3), the proposed UAL is formulated as $l_{UAL}^{i,j} = 1 - |2p_{i,j} - 1|^2$. Finally, the total loss function can be written as:

$$L = L_{BCEL} + \lambda \times L_{UAL}, \quad (4)$$

where λ is the balance coefficient and we design three adjustment strategies of λ , i.e., a fixed constant value, an increasing linear strategy, and an increasing cosine strategy in Sec. 4.3. The different forms and corresponding results are listed in the supp. From the results, we find that the increasing strategies, especially “cosine”, do achieve better performance. So, the cosine strategy is used by default.

4. Experiments

4.1. Experiment Setup

Datasets. We use four COD datasets, CAMO [16], CHAMELEON [41], COD10K [9] and NC4K [29]. CAMO consists of 1,250 camouflaged and 1,250 non-camouflaged images. CHAMELEON contains 76 hand-annotated images. COD10K includes 5,066 camouflaged, 3,000 background and 1,934 non-camouflaged images. NC4K is

another large-scale COD testing dataset including 4,121 images from the Internet. Following the data partition of [9, 29, 31, 59], we use all images with camouflaged objects in the experiments, in which 3,040 images from COD10K and 1,000 images from CAMO are used for training, and the rest ones for testing. Besides, we also show the performance on five SOD datasets in the supp.

Evaluation Criteria. For COD and SOD, we use seven common metrics for evaluation based on [32, 33], including S-measure [6] (S_m), weighted F-measure [30] (F_β^ω), mean absolute error (MAE), F-measure [1] (F_β), E-measure [7] (E_m), precision-recall (PR) curve and F_β -threshold curve (F_β curve). The curves can be found in the supp.

Implementation Details. The proposed ZoomNet is implemented with PyTorch. As the settings in recent methods [9, 29, 31, 59], the encoder is initialized with the parameters of ResNet-50 pretrained on ImageNet, and the remaining parts are randomly initialized. SGD with momentum 0.9 and weight decay 0.0005 is chosen as the optimizer. The learning rate is initialized to 0.05 and follows a linear warm-up and linear decay strategy. The entire model is trained for 40 epochs with a batch size of 8 in an end-to-end manner on an NVIDIA 2080Ti GPU. During training and inference, the main scale is 384×384 . Random flipping and rotating are employed to augment the training data.

4.2. Comparisons with State-of-the-arts

COD is an emerging field, so we introduce some methods for salient object detection and medical image segmentation for comparison. The results of all these methods come from existing public data or are generated by models

Table 2. Ablation study on the COD10K-Test. SIU: Scale integration unit; HMU: Hierarchical mixed-scale unit with g groups; UAL: Uncertainty-aware loss; ⑤: The simple extension of baseline ① with the similar number of parameters and FLOPs to ④.

Model	GFLOPs	Params. (M)	HMU	SIU	UAL	$S_m \uparrow$	$F_\beta^\omega \uparrow$	MAE \downarrow	$F_\beta \uparrow$	$E_m \uparrow$
①	41.885	30.453				0.812	0.637	0.039	0.692	0.898
②	78.560	30.696	$g = 6$			0.820	0.654	0.037	0.706	0.897
③	166.821	30.453		✓		0.837	0.682	0.034	0.731	0.912
④	203.496	28.794	$g = 6$	✓		0.843	0.694	0.032	0.743	0.909
⑤	207.647	41.975				0.815	0.632	0.040	0.689	0.896
⑥	173.616	30.819	$g = 2$	✓	✓	0.835	0.721	0.030	0.758	0.902
⑦	188.556	31.600	$g = 4$	✓	✓	0.836	0.723	0.029	0.761	0.905
⑧	203.496	32.382	$g = 6$	✓	✓	0.838	0.729	0.029	0.766	0.911
⑨	218.436	33.163	$g = 8$	✓	✓	0.836	0.726	0.029	0.763	0.907

that are retrained based on the code released by the authors. **Quantitative Evaluation.** Tab. 1 shows the detailed comparison results. It can be seen that the proposed model consistently and significantly surpasses recent methods on all datasets without relying on any post-processing tricks. Compared with the recent best COD method UJSC, although it introduces extra SOD data for training and has suppressed other existing methods, our method still shows the obvious performance improvement on these datasets. Especially, our approach has more advantages on the metrics F_β^ω , MAE, and F_β . On four datasets, the proposed method averagely outperforms the second-best method C²F-Net by 19.3% in terms of MAE and the average gains in terms of F_β^ω and F_β are 4%. Besides, PR and F_β curves shown in the *supp.* also demonstrate the effectiveness of the proposed method. The flatness of the F_β curve reflects the consistency and uniformity of the prediction. Our curves are almost horizontal, which can be attributed to the effect of the proposed UAL. It drives the predictions to be more polarized and reduces the ambiguity.

Qualitative Evaluation. Visual comparisons of different methods on several typical samples are shown in Fig. 5. They present the complexity in different aspects, such as big objects (Col. 1), middle objects (Col. 2-8), small objects (Col. 9-13), occlusions (Col. 2 and 10), background interference (Col. 10-13), and indefinable boundaries (Col. 1, 2, 6-13). These results intuitively show the superior performance of the proposed method. In addition, it can be noticed that our predictions have clearer and more complete object regions and sharper contours.

4.3. Ablation Studies

In this section, we perform comprehensive ablation analyses on different components. Because COD10K is the most widely-used large-scale COD dataset, and contains various objects and scenes, all subsequent ablation experiments are carried out on it.

Effectiveness of SIUs and HMUs. In the proposed model, both the SIU and the HMU are very important structures. We install them one by one on the baseline model to eval-



Figure 6. Visual comparisons for showing the effects of the proposed components. B: Baseline; +S: +SIUs; +H: +HMUs; +S+H: +SIUs+HMUs; +S+H+L: +SIUs+HMUs+UAL.

Table 3. Comparisons of mixed and single scale input schemes on CAMO and COD10K. All models are based on ① in Tab. 2.

Input Scale	Combination Strategy	$S_m \uparrow$	$F_\beta^\omega \uparrow$	MAE \downarrow	$F_\beta \uparrow$	$E_m \uparrow$	AVG. Relative Improvement
1.0x	—	0.797	0.649	0.063	0.704	0.875	
0.5x	—	0.746	0.553	0.076	0.616	0.833	↓12.03%
0.5x, 1.0x	Addition	0.801	0.647	0.062	0.702	0.876	↑0.31%
0.5x, 1.0x	SIU	0.806	0.658	0.059	0.709	0.879	↑1.86%
1.5x	—	0.820	0.683	0.059	0.737	0.890	↑4.05%
0.5x, 1.5x	Addition	0.820	0.680	0.058	0.735	0.894	↑4.43%
0.5x, 1.5x	SIU	0.822	0.687	0.056	0.740	0.893	↑5.33%
1.0x, 1.5x	Addition	0.819	0.685	0.058	0.738	0.892	↑4.47%
1.0x, 1.5x	SIU	0.826	0.697	0.056	0.745	0.897	↑5.97%
0.5x, 1.0x, 1.5x	Addition	0.821	0.690	0.056	0.742	0.894	↑5.43%
0.5x, 1.0x, 1.5x	SIU	0.827	0.700	0.054	0.751	0.898	↑6.91%

uate their performance. The results are shown in Tab. 2. Our baseline ① and other models ② and ⑤ only use the inputs of the main scale. As can be seen, our baseline shows a good performance, probably due to the proper training setup and the more reasonable network architecture detailed in the *supp.* From ①-④, it can be seen that the two proposed modules make a significant contribution to the performance when compared to the baseline. Besides, the results in Fig. 6 show that the two modules can benefit each other and reduce their errors (e.g., Col. 1, 2, 5 and 7) to locate and distinguish objects more accurately. These components effectively help the model to excavate and distill the critical and valuable semantics and improve the capability of distinguishing hard objects. Under the cooperation between the proposed modules and loss functions, ZoomNet can completely capture the camouflaged objects of different scales and generate the predictions with higher contrast and consistency. In addition, in Tab. 2, the model ⑤ is a simple extension of the baseline model ① using some standard convolutional blocks, to make the similar number of parameters and FLOPs with ④. The model ④ still achieves better performance, which reflects the effectiveness of the proposed modules and the rationality of the design.

Number of Groups in HMUs. In Tab. 2, we also show the effects of different group numbers in the proposed HMU. It can be seen from the results that the best performance appears when the number of groups is equal to 6. Also, it achieves a good balance between performance and efficiency. So, in other experiments, we set the number of

Table 4. Different forms of the proposed UAL. Form 0 is our model without UAL. “—”: Unable to converge. Form 1.5 is used by default because of its balanced performance. Their curves are shown in the *supp.*

No.	Form	α	$S_m \uparrow$	$F_\beta \uparrow$	MAE \downarrow	$F_\beta \uparrow$	$E_m \uparrow$
0	χ	χ	0.843	0.694	0.032	0.743	0.909
1.1, 1.2, 1.3	$\Phi_{pow}^\alpha(x) = 1 - 2x - 1 ^\alpha$	1/8, 1/4, 1/2	—	—	—	—	—
1.4		1	0.834	0.716	0.029	0.757	0.903
1.5		2	0.838	0.729	0.029	0.766	0.911
1.6		4	0.834	0.727	0.029	0.763	0.903
1.7		8	0.833	0.725	0.029	0.760	0.900
2.1	$\Phi_{exp}^\alpha(x) = e^{-(\alpha(x-0.5))^2}$	1/8	0.844	0.698	0.032	0.744	0.907
2.2		1/4	0.845	0.700	0.032	0.746	0.911
2.3		1/2	0.843	0.701	0.031	0.746	0.908
2.4		1	0.842	0.713	0.030	0.754	0.908
2.5		2	0.839	0.720	0.030	0.761	0.908
2.6		4	0.839	0.706	0.032	0.752	0.909
2.7		8	0.841	0.698	0.032	0.745	0.910
3	BCE w/ $\omega = 1 + \Phi_{pow}^2(x)$	2	0.844	0.697	0.032	0.744	0.913

groups in each HMU to 6.

Mixed-scale Input Scheme. Our model is designed to mimic the behavior of “Zoom In&Out”. The feature expression is enriched by combining the scale-specific information from different scales. In Fig. 1, the intermediate features and attention maps show that our mixed-scale scheme plays a positive and important role in locating the camouflaged object. Considering that the objects in COD10K are mainly small objects [9], which may limit the role of $0.5\times$ input to some extent, we list average results on COD10K and CAMO in Tab. 3. The proposed scheme performs better than the single-scale one and simply mixed one. This verifies the rationality of such a design for the COD task.

Options of Setting λ . We compare three strategies and the results are listed in the *supp.*, in which the increasing cosine strategy achieves the best performance. This may be due to the advantage of its smooth change process. This smooth intensity warm-up strategy of UAL motivates the model to take advantage of UAL in improving the learning process and to mitigate the possible negative interference of UAL on BCEL due to the lower accuracy of the model during the early stage of training.

Forms of UAL. Different forms of UAL are listed in Tab. 4 and the corresponding curves are illustrated in the *supp.* As can be seen, Form 1.5 has a more balanced performance. Also, it is worth noting that, when approaching 0 or 1, the form which can maintain a larger gradient will obtain better performance in terms of F_β^ω , MAE and F_β . This may provide some reference for designing a better loss.

Effectiveness of UAL. The results of Fig. 6 intuitively show that the UAL greatly reduces the ambiguity caused by the interference from the background. Besides, we visualize the histogram maps of all results on CHAMELEON and the intermediate features from different stages in the decoder in the *supp.* In the stacked histogram map “w/o UAL”, a large number of pixels appear in the middle area, which corresponds to more visually blurred/uncertain pre-

dictions. Besides, in the corresponding feature visualization, especially in the region inside the red box, there is clear background interference due to the complex scenarios and blurred edges, which are extremely prone to yield false positive predictions. However, when UAL is introduced, it can be seen that the middle interval of “w UAL” is flatter than the one of “w/o UAL”, that is, most pixel values approach two extremes. And the feature maps become more discriminative and present a more compact and complete response in the regions of camouflaged objects.

4.4. Discussion on SOD and COD

It can be seen from the experiments in the *supp.* that our method not only performs well on COD, but also shows outstanding performance on SOD. Considering the difference between these two tasks, we may wonder *why our method consistently performs well on such two seemingly different tasks*. We attribute this to the generality and rationality of the designed structure. In fact, SOD and COD have a clear commonality, *i.e.*, the accurate segmentation has to depend on multi-scale and category-free discriminative features. By integrating rich scale-specific features, our model can extract critical and informative cues from scenes and objects, which helps precise localization and smooth segmentation of objects. In addition, the proposed UAL can mitigate the ambiguity of predictions caused by the inherent complexity of scenes. Although the objects in SOD are salient, it can also benefit from UAL due to the vagueness introduced by the CNN model itself in the detailed information recovery process. All of these components are built on the common demand of the two tasks, which provides a solid foundation for the performance.

5. Conclusion

In this paper, we propose the ZoomNet by imitating the behavior of human beings to zoom in and out on images. This process actually considers the differentiated expressions about the scene from different scales, which helps to improve the understanding and judgment of camouflaged objects. We first filter and aggregate scale-specific features through the scale merging layer to enhance feature representation. Next, in the hierarchical mixed-scale decoder, the strategies of grouping, mixing and fusion further mine the mixed-scale semantics. Lastly, we introduce the uncertainty-aware loss to penalize the ambiguity of the prediction. Extensive experiments verify the effectiveness of the proposed method in both the COD and SOD tasks with superior performance to existing state-of-the-art methods.

Acknowledgements This work was supported by the National Natural Science Foundation of China #61876202 and #61829102, the Liaoning Natural Science Foundation #2021-KF-12-10, and the Fundamental Research Funds for the Central Universities #DUT20ZD212.

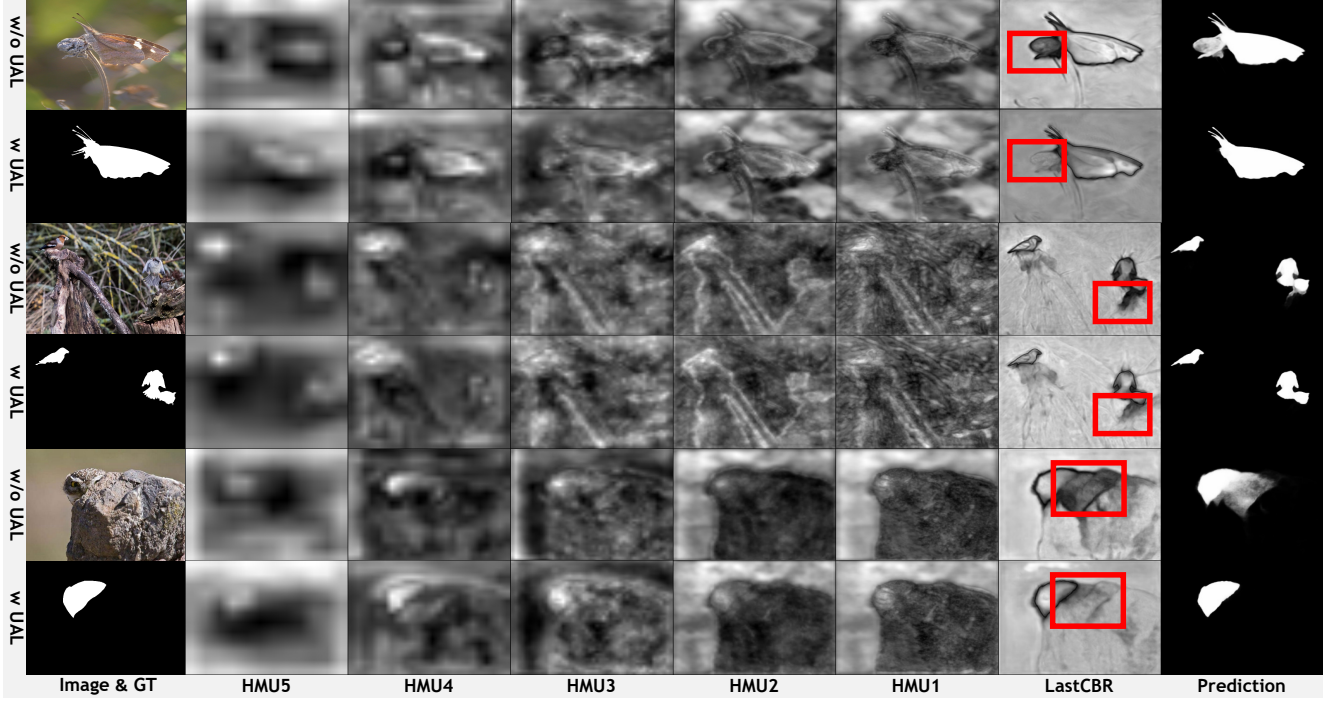


Figure 7. Visual comparison of intermediate feature maps from different stages of the decoder for showing the effects of the proposed UAL. Please zoom in for more details. HMU: Hierarchical mixed-scale unit; LastCBR: The last “Conv3 \times 3-BN-ReLU” structure before the layer generating the logits map.

This appendix will introduce more details that cannot be expanded in the main text, while showing the performance on SOD.

A. Model Details

A.1. E-Net

E-Net is based on the feature extraction part of ResNet-50 [14] and the layers after the “layer4” are removed. We collect the feature maps before passing the first max-pooling layer and the output feature maps of “layer1”, “layer2”, “layer3” and “layer4” as the output feature maps of the E-Net. The numbers of channels corresponding to them are 64, 256, 512, 1024, and 2048, respectively.

A.2. C-Net

Following the setting of the method [67], in C-Net, we use an ASPP [3] simplified according to our needs as the feature compression layer corresponding to the “layer4” of E-Net and other layers are simply composed of an independent “Conv3 \times 3-BN-ReLU” (3×3 CBR) unit. The numbers of output channels of all levels are set to 64 in our models.

The ASPP layer is composed of five CBR branches. The kernel sizes and dilation rates of them are 1, 3, 3, 3, 1 and 1, 2, 5, 7, 1. All convolution operations use the padding to ensure that the input and output sizes are consistent. A

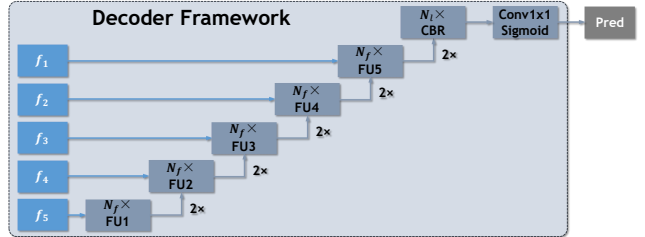


Figure 8. Illustration of the basic framework adopted by the decoder in our proposed method. FU: The fusion unit for fusing the up-sampled feature map from the previous FU and the shallower feature map f_i ($i = 1, 2, \dots, 5$). $2\times$: The bi-linear interpolation operation with a factor of 2. CBR: The “Conv3 \times 3-BN-ReLU” unit. Conv1 \times 1: The convolution operation with a kernel size of 1×1 . N_f and N_l : The numbers of FUs and the last CBR units.

global average pooling operation and an up-sampling operation are used before and after the second 1×1 CBR branch to capture the global context information and restore it to the original size. All results of the five branches are concatenated along the channel dimension and fused by a 3×3 CBR unit to obtain the output.

A.3. Decoder Framework

The decoder networks of our models in all experiments follow the same framework as shown in Fig. 8. Before being

```

1 class StackedCBRBlock(nn.Sequential):
2     def __init__(self, in_c, out_c, num_blocks=1, kernel_size=3):
3         super().__init__()
4         self.kernel_setting = dict(kernel_size=kernel_size, stride=1, padding=kernel_size // 2)
5         cs = [in_c] + [out_c] * num_blocks
6         self.channel_pairs = tuple(self.slide_win_select(cs, win_size=2, win_stride=1, drop_last=True))
7         for i, (i_c, o_c) in enumerate(self.channel_pairs):
8             self.add_module(name=f"cbr_{i}", module=CBR(i_c, o_c, **self.kernel_setting))
9
10    @staticmethod
11    def slide_win_select(items, win_size=1, win_stride=1, drop_last=False):
12        i = 0
13        while i + win_size <= len(items):
14            yield items[i: i + win_size]
15            i += win_stride
16        if not drop_last:
17            yield items[i: i + win_size]

```

Listing 1. Code of stacked CBR units.

fed into the fusion unit (FU), the up-sampled deeper feature map is directly added to the shallow feature map.

In our all experiments, N_f and N_l are set to 1. The numbers of input & output channels of the last 3×3 CBR unit are 64 and 32, respectively. The number of output channels of the “Conv1 \times 1” is 1 and a sigmoid layer is cascaded to convert the logits map to the prediction. In the decoder of the proposed ZoomNet, the FU is set to the HMU and the other layers remain the same.

A.4. Baseline Model

In the ablation study, we introduce a simple encoder-decoder network as our baseline model to evaluate the performance of different proposed components. It contains a feature extraction network “E-Net”, a simple multi-level feature compression convolutional network “C-Net”, and a basic convolutional decoder where the FU is set to the 3×3 CBR unit. In the following text, “CBR1-5” are used to refer to these five units.

A.5. Model ⑤

In Tab. 2 of the main text, based on the baseline model ①, we construct the model ⑤ with the similar amount of parameters and FLOPs to ④ to reflect the effectiveness of the method and the rationality of the design. For increasing the number of parameters and FLOPs, we made the following modifications to the baseline model ①:

- The number of output channels of all levels of C-Net: $64 \rightarrow 128$.
- The number of input/output channels of CBR1-5 units of the basic convolutional decoder: $64 \rightarrow 128$.
- The number of input channels of the last CBR unit: $64 \rightarrow 128$.
- The number of CBR units (N_f and N_l) of all levels of the basic convolutional decoder: $1 \rightarrow 3$.
- The kernel size of the convolution operation in all lev-

Table 5. Comparisons of the number of parameters and FLOPs based on <https://github.com/lartpang/MethodsCmp> corresponding to recent COD methods. All evaluations follow the inference settings in the corresponding papers.

Method	Ours	UGTR [56]	C ² F-Net [43]	UJSC [17]	PFNet [31]	MGL-R [59]	SLSR [29]	SINet [9]
Params.	32.382M	48.868M	28.411M	217.982M	46.498M	63.595M	50.935M	48.947M
FLOPs	203.496G	1.007T	26.167G	112.341G	53.222G	553.939G	66.625G	38.757
FPS	24.030	16.640	65.759	34.178	62.590	13.373	58.782	56.509

els of the basic convolutional decoder: $3 \rightarrow 5$.

To facilitate understanding, the corresponding code for the stacked CBR units used here is listed in List. 1.

Algorithm 1 The iteration struction in the HMU

Input: $\{g_j\}_{j=1}^G$: feature groups; $G \geq 2$: the number of groups; $C = 32$: the number of channels in a single feature group g_j ; \mathcal{S} : splitting operation; $\mathcal{T}_{C_o \times C_i}$: stacked CBR units with initial input and final output channel numbers of C_i and C_o as listed in List. 1; \mathcal{C} : concatenate operation;

Output: $\{g_j^2\}_{j=1}^G$: the feature set for generating the modulation vector α ; $\{g_j^3\}_{j=1}^G$: the feature set used to be modulated and generate the final output of the HMU;

```

1: for  $i \leftarrow 1, G$  do
2:     if  $i = 1$  then ▷ Group 1
3:          $g_i^1, g_i^2, g_i^3 \leftarrow \mathcal{S}(\mathcal{T}_{3C \times C}^i(g_i))$ ;
4:          $g_{prev}^1 \leftarrow g_i^1$ ;
5:     else if  $i = G$  then ▷ Group  $G$ 
6:          $g_i^2, g_i^3 \leftarrow \mathcal{S}(\mathcal{T}_{2C \times 2C}^i(\mathcal{C}(g_i, g_{prev}^1)))$ ;
7:     else ▷ Group  $i, 1 < i < G$ 
8:          $g_i^1, g_i^2, g_i^3 \leftarrow \mathcal{S}(\mathcal{T}_{3C \times 2C}^i(\mathcal{C}(g_i, g_{prev}^1)))$ ;
9:          $g_{prev}^1 \leftarrow g_i^1$ ;
10:    end if
11: end for

```

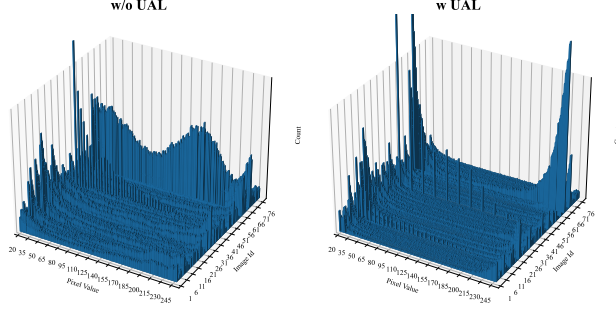


Figure 9. Visual comparison of histograms of all 76 prediction results on the CHAMELEON [41] dataset, which is a stack of the histogram of each prediction. A good result should embody a closely binarized histogram at both ends. For a more clear demonstration, only the interval with pixel values between 20 and 245 is counted here. It is best to zoom in for more details.

Table 6. Comparisons of different increasing strategies of λ . λ_{const} : A constant value and it is set to 1. t and T : The current and total number of iterations, respectively. λ_{min} and λ_{max} : The minimum and maximum values of λ , and they are set to 0 and 1 in our experiments. “Linear, $t_{min} \rightarrow t_{max}$ ”: The linearly increasing interval in the iterations is $[t_{min}, t_{max}]$. clip: Values outside the interval are clipped to the interval edges.

Strategy	λ	$S_m \uparrow$	$F_\beta \uparrow$	MAE \downarrow	$F_\beta \uparrow$	$E_m \uparrow$
Cosine	$\lambda_{min} + \frac{1}{2}(1 - \cos(\frac{t}{T}\pi))(\lambda_{max} - \lambda_{min})$	0.838	0.729	0.029	0.766	0.911
Linear, $t_{min} \rightarrow t_{max}$	clip($\lambda_{min} + \frac{t - t_{min}}{t_{max} - t_{min}}(\lambda_{max} - \lambda_{min}), \lambda_{min}, \lambda_{max}$)	0.834	0.723	0.029	0.760	0.908
Linear, $0.3T \rightarrow 0.7T$		0.832	0.719	0.030	0.758	0.904
Constant	λ_{const}	0.830	0.717	0.030	0.757	0.906

B. HMU: Perspective of Kernel Pyramid

The iteration structure of feature groups in HMU is actually equivalent to an integrated multi-path kernel pyramid structure with partial parameter sharing. In order to understand this intuitively, we highlight the feature information flow of different groups in the iterative structure in Fig. 11. Specifically, the 3×3 CBR unit corresponding to the feature group in the iteration structure can be split according to the output feature groups. As shown in the “Integrated Kernel Pyramid” on the left of Fig. 11, each original CBR unit with an output channel number of $3C$ is converted to three independent CBR units with a shared input. And the numbers of output channels of them are C . When we further decouple the integrated form on the left into the form on the right, we can clearly see that the information flow paths corresponding to different feature groups each form a multi-branch kernel pyramid structure and there are some shared parameters between these pyramids.

As mentioned in the main text of the paper, some of the channels in the output feature of each branch are used together to generate the modulation vector. It not only weights the channels inside each branch, but also weights different branches. If viewed from the aforementioned perspective of

Table 7. Comparison results of methods trained without CPD1K-TR on CPD1K-TE [71].

Model	$S_m \uparrow$	$F_\beta \uparrow$	MAE \downarrow	$F_\beta \uparrow$	$E_m \uparrow$
ZoomNet	0.759	0.537	0.011	0.578	0.843
C ² FNet	0.743	0.495	0.016	0.528	0.840
PFNet	0.722	0.460	0.017	0.494	0.819

the kernel pyramid, such an operation can be seen as a relative modulation of the different kernel pyramids contained in the iterative structure of the HMU.

Besides, in our HMU, C is set to 32. The number of channels of the final output feature of the HMU is the same as the input feature, both are 64. We also list the algorithm of the iteration structure in Alg. 1 to present the process more clearly and to complement the related statement in the main text.

C. More Comparisons

C.1. PR & F_β curves of COD Methods

In Fig. 12, we show the PR & F_β curves of different methods on four COD datasets. The red curve represents our method.

C.2. Comparisons of Param. & FLOPs

In Tab. 5, we list the number of parameters and FLOPs of existing COD methods and ours. Our method provides a performance-robust solution with the second-smallest amount of parameters for the COD task. But there may be still some redundancy in the design of the inference structure. The adopted explicit scale-independent design may bring additional inference cost. We will explore and improve this in future work.

C.3. Intermediate Feature Maps of the Decoder

We show the intermediate feature maps from different stages of the decoder in Fig. 7.

C.4. Effectiveness of UAL

In Fig. 9, we visualize the histogram maps of all results on CHAMELEON [41].

C.5. Different Forms of λ

The different adjustment functions of the coefficient λ and their results of UAL are list in Tab. 6.

C.6. Different Forms of UAL

The different forms of UAL are shown in Fig. 10.

C.7. Performance in More Complex Scenes

Actually, COD10K-TE is a very representative test dataset with rich and diverse scenarios and objects. Be-

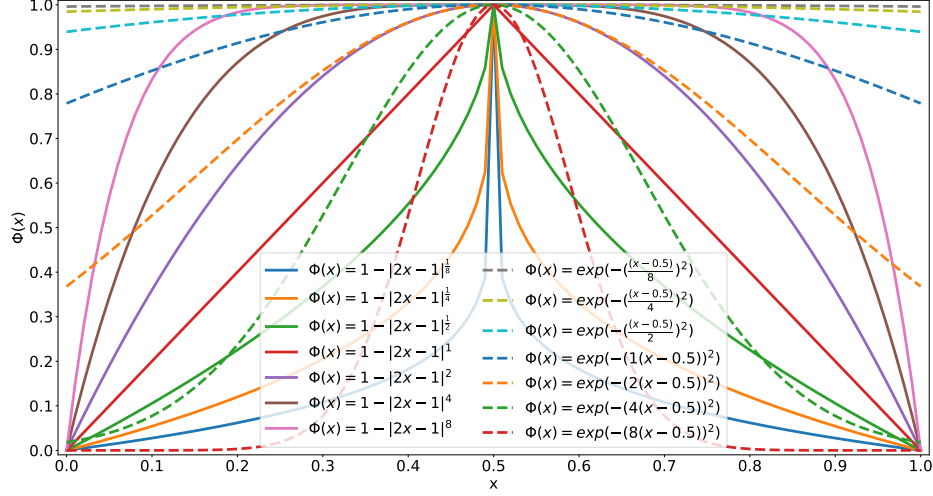


Figure 10. Curves of different forms of the proposed UAL.

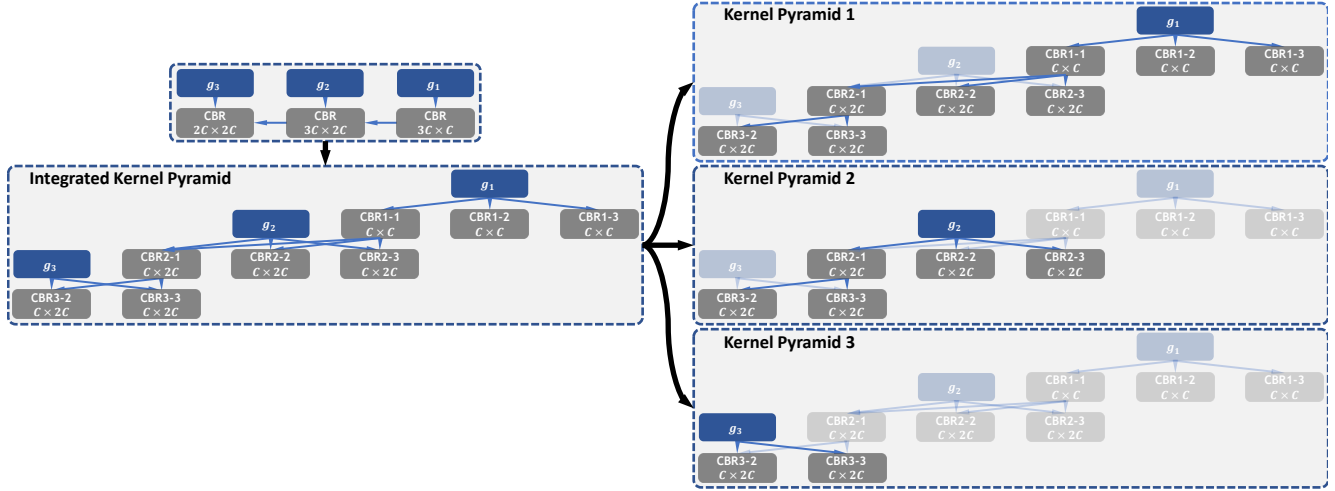


Figure 11. The iteration structure of feature groups in HMU can be regarded as an integrated kernel pyramid. Without loss of generality, we show the situation of the number of groups $G = 3$ in the figure. The actual final model is set to $G = 6$. The only difference lies in the number of repetitions of the kernel pyramid structure in the middle. “CBR1- j ”: The “Conv3 × 3-BN-ReLU” structure corresponding to the input feature group g_i and the j th output feature group. $C_o \times C_i$: The numbers of input and output channels of the CBR unit is C_i and C_o , respectively.

sides, there is also a very complex small-scale dataset CPD1K [71]. Tab. 7 shows the results of our method and some state-of-the-art competitors (all are trained without CPD1K-TR). The test results on CPD1K-TE can reflect the adaptability of the model to complex scenarios. The experiment shows the superior performance of our method in more complex scenarios.

D. Experiments on SOD

In order to show good generalization and further verify the rationality of the structural design, we evaluate the proposed model on the SOD task.

D.1. Datasets

Our experiment on SOD is based on the existing five SOD datasets, DUT-OMRON [55] (5168), DUTS [44] (10553 + 5017), ECSSD [54] (1000), HKU-IS [18] (4447) and Pascal-S [18] (850). We only use the training set of DUTS for training. During the test phase, we use the remaining data for inference.

D.2. Implementation Details

For a fair comparison on SOD, the proposed model is re-trained on DUTS [44] following the training strategies and techniques of [35, 46, 47, 51, 64]. The learning rate is initial-

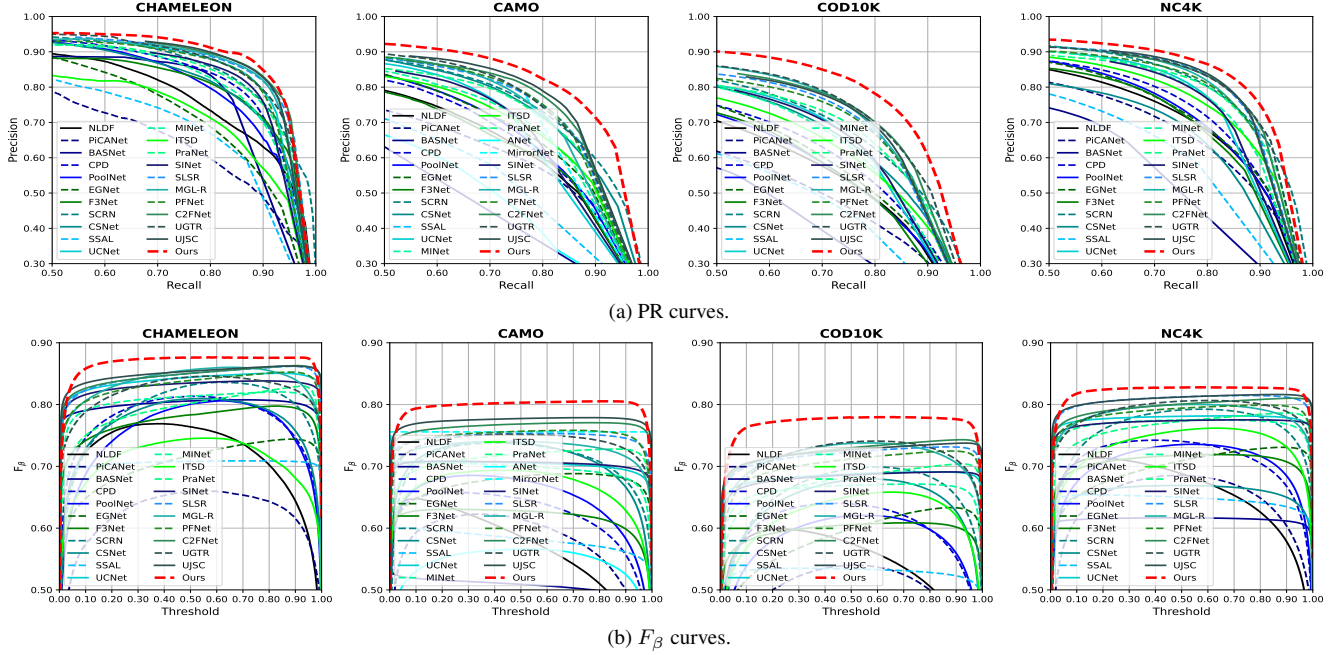


Figure 12. PR and F_β curves of the proposed model and recent SOTA algorithms over four COD datasets.

Table 8. More detailed comparison results on the SOD task. The best results are highlighted in **red**, **green** and **blue**. These results are based on the VGG [40], ResNet [14] and T2T-ViT [57] version of the corresponding method.

Model	Backbone	Year	DUT-OMRON					DUTS-TE					ECSSD					HKU-IS					PASCAL-S				
			$S_m \uparrow$	$F_\beta \uparrow$	MAE \downarrow	$F_\beta \uparrow$	$E_m \uparrow$	$S_m \uparrow$	$F_\beta \uparrow$	MAE \downarrow	$F_\beta \uparrow$	$E_m \uparrow$	$S_m \uparrow$	$F_\beta \uparrow$	MAE \downarrow	$F_\beta \uparrow$	$E_m \uparrow$	$S_m \uparrow$	$F_\beta \uparrow$	MAE \downarrow	$F_\beta \uparrow$	$E_m \uparrow$	$S_m \uparrow$	$F_\beta \uparrow$	MAE \downarrow	$F_\beta \uparrow$	$E_m \uparrow$
RAS [4]	VGG16	2018	0.814	0.695	0.062	0.731	0.860	0.839	0.740	0.059	0.779	0.889	0.893	0.857	0.056	0.887	0.931	0.887	0.843	0.045	0.875	0.940	0.793	0.735	0.106	0.790	0.846
MLMSNet [49]	VGG16	2019	0.809	0.681	0.064	0.710	0.848	0.862	0.761	0.049	0.792	0.907	0.911	0.871	0.045	0.890	0.944	0.907	0.859	0.039	0.878	0.950	0.845	0.785	0.075	0.814	0.893
PAGNet [45]	VGG16	2019	0.824	0.722	0.062	0.743	0.858	0.854	0.769	0.052	0.793	0.896	0.912	0.886	0.042	0.904	0.947	0.903	0.865	0.037	0.884	0.948	0.838	0.789	0.079	0.819	0.885
PICANet [24]	ResNet50	2018	0.832	0.695	0.065	0.729	0.876	0.869	0.755	0.051	0.791	0.920	0.917	0.867	0.046	0.890	0.952	0.904	0.840	0.043	0.866	0.950	0.852	0.779	0.078	0.812	0.899
BASNet [37]	ResNet34	2019	0.836	0.751	0.056	0.767	0.871	0.866	0.803	0.048	0.823	0.903	0.916	0.904	0.037	0.917	0.951	0.909	0.889	0.032	0.902	0.951	0.834	0.797	0.079	0.824	0.883
CPD [50]	ResNet50	2019	0.825	0.719	0.056	0.742	0.868	0.869	0.795	0.043	0.821	0.914	0.918	0.898	0.037	0.913	0.951	0.905	0.875	0.034	0.892	0.950	0.844	0.800	0.074	0.827	0.888
PoolNet [23]	ResNet50	2019	0.831	0.725	0.054	0.747	0.867	0.887	0.817	0.037	0.840	0.926	0.926	0.904	0.035	0.918	0.956	0.919	0.888	0.030	0.903	0.958	0.864	0.819	0.067	0.846	0.905
EGNet [63]	ResNet50	2019	0.841	0.738	0.053	0.760	0.878	0.887	0.816	0.039	0.839	0.927	0.925	0.903	0.037	0.918	0.955	0.918	0.887	0.031	0.902	0.958	0.850	0.804	0.076	0.833	0.892
HRS [58]	ResNet50	2019	0.772	0.645	0.066	0.690	0.841	0.829	0.746	0.051	0.791	0.899	0.883	0.859	0.054	0.894	0.934	0.882	0.851	0.042	0.883	0.941	0.799	0.744	0.091	0.792	0.866
SCRN [51]	ResNet50	2019	0.837	0.720	0.056	0.749	0.875	0.885	0.803	0.040	0.833	0.925	0.927	0.900	0.037	0.916	0.956	0.916	0.876	0.034	0.894	0.956	0.865	0.813	0.066	0.840	0.906
F3Net [46]	ResNet50	2020	0.838	0.747	0.053	0.766	0.872	0.888	0.835	0.035	0.852	0.927	0.924	0.912	0.033	0.925	0.955	0.917	0.900	0.028	0.910	0.958	0.857	0.823	0.064	0.843	0.901
GCPANet [5]	ResNet50	2020	0.839	0.734	0.056	0.756	0.869	0.891	0.821	0.038	0.841	0.929	0.927	0.903	0.035	0.916	0.955	0.920	0.889	0.031	0.901	0.958	0.864	0.819	0.063	0.840	0.906
LDf [47]	ResNet50	2020	0.839	0.752	0.052	0.770	0.869	0.892	0.845	0.034	0.861	0.930	0.924	0.915	0.034	0.927	0.954	0.919	0.904	0.028	0.913	0.958	0.859	0.829	0.062	0.851	0.905
DFI [22]	ResNet50	2020	0.840	0.738	0.055	0.762	0.877	0.887	0.817	0.039	0.840	0.928	0.927	0.906	0.035	0.920	0.957	0.919	0.890	0.031	0.903	0.961	0.864	0.824	0.066	0.849	0.907
GateNet [67]	ResNet50	2020	0.838	0.729	0.055	0.757	0.875	0.885	0.809	0.040	0.837	0.928	0.920	0.894	0.040	0.913	0.952	0.915	0.880	0.033	0.897	0.955	0.854	0.804	0.071	0.835	0.900
ITSd [72]	ResNet50	2020	0.840	0.750	0.061	0.768	0.880	0.885	0.824	0.041	0.840	0.929	0.925	0.910	0.034	0.921	0.959	0.917	0.894	0.031	0.904	0.960	0.859	0.823	0.066	0.843	0.910
MINet [35]	ResNet50	2020	0.833	0.738	0.056	0.757	0.869	0.884	0.825	0.037	0.844	0.927	0.925	0.911	0.033	0.923	0.957	0.919	0.897	0.029	0.909	0.960	0.854	0.818	0.066	0.841	0.901
VST [25]	T2T-ViT-14	2021	0.850	0.755	0.058	0.774	0.888	0.896	0.828	0.037	0.845	0.939	0.932	0.910	0.033	0.920	0.964	0.928	0.897	0.029	0.907	0.968	0.871	0.827	0.062	0.847	0.918
SAMNet [26]	Handcraft	2021	0.830	0.699	0.065	0.734	0.877	0.849	0.729	0.058	0.768	0.901	0.907	0.858	0.050	0.883	0.945	0.898	0.837	0.045	0.864	0.946	0.822	0.743	0.095	0.784	0.869
SGL-KRN [52]	ResNet50	2021	0.846	0.765	0.049	0.783	0.885	0.893	0.847	0.034	0.865	0.939	0.923	0.910	0.036	0.924	0.954	0.921	0.904	0.028	0.915	0.961	0.854	0.823	0.070	0.849	0.904
CTDNet [70]	ResNet50	2021	0.844	0.762	0.052	0.779	0.881	0.893	0.847	0.034	0.862	0.935	0.925	0.915	0.032	0.927	0.956	0.921	0.909	0.027	0.918	0.961	0.859	0.829	0.064	0.851	0.904
Auto-MSFNet [62]	ResNet50	2021	0.832	0.757	0.050	0.772	0.875	0.877	0.841	0.034	0.855	0.931	0.914	0.916	0.032	0.927	0.954	0.908	0.903	0.027	0.912	0.959	0.849	0.830	0.063	0.852	0.902
Ours	ResNet50	2021	0.841	0.755	0.053	0.771	0.872	0.900	0.854	0.033	0.866	0.936	0.935	0.926	0.027	0.933	0.963	0.931	0.918	0.023	0.923	0.967	0.869	0.844	0.057	0.860	0.917

ized to 0.05 and follows a linear warm-up and linear decay strategy. And the main scale is changed to 352×352 to achieve a trade-off between performance and speed. The model tends to converge after 50 epochs with a batch size of 22.

D.3. Comparisons with State-of-the-arts

We compare the proposed model with 22 existing methods. All the results are listed in Tab. 8 and shown in Fig. 13. Our model outperforms all these competitors, which shows that the proposed model can deal with the more general binary segmentation task.

E. Limitations and Future Work

Although our ZoomNet provides a powerful and effective solution for the COD task, some limitations still exist and are worth exploring further.

1. In the current work, the shared feature extraction structure explicitly collects complementary information from different scales on the image pyramid, which is designed to mimic the behavior of zooming in and out. But for human beings, the process of information extraction and integration should be implicit and internalized in the process of knowledge learning. Moreover, this explicit scale-independent design also brings the additional inference cost. Although our method has

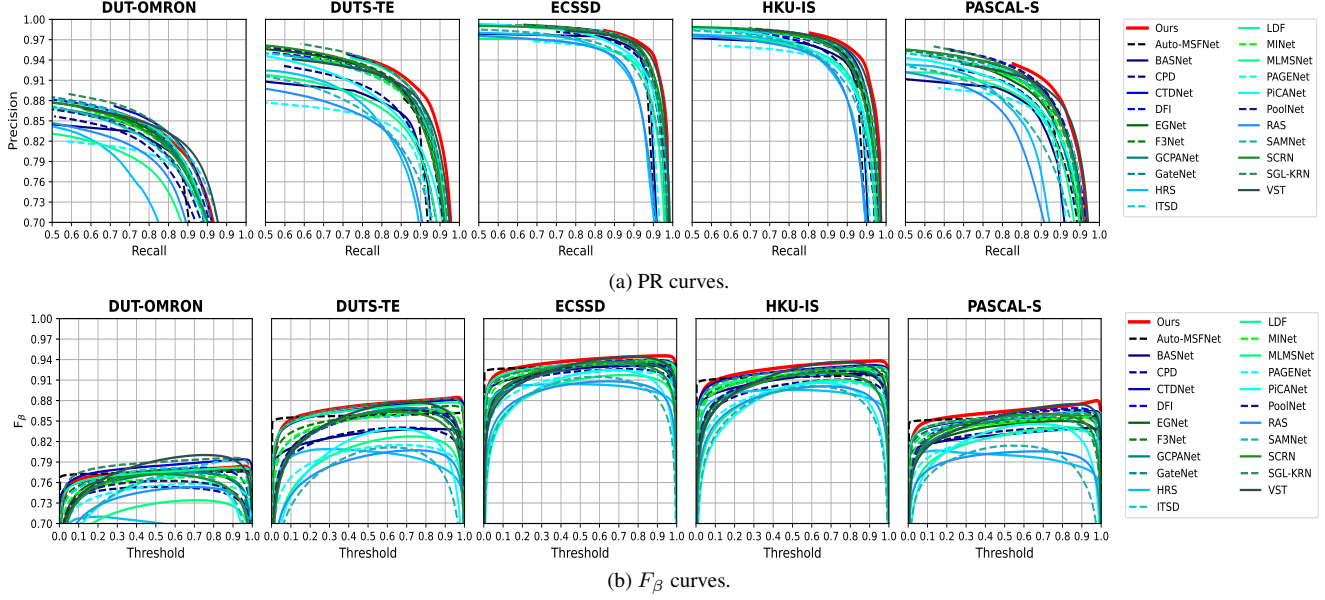


Figure 13. PR and F_β curves of the proposed model and recent SOTA algorithms over five SOD datasets.

achieved good performance on COD and SOD tasks, the inference speed is still slightly slower than the current fastest method, C2FNet [43].

2. Besides, there is still room for improvement in the way of mining effective clues from small-scale features in SIU.

In future work, we will try to further simplify the inference structure to make it more in line with the actual human decision-making process and optimize the ability of our method to extract contextual cues from small-scale features.

References

- [1] Radhakrishna Achanta, Sheila Hemami, Francisco Estrada, and Sabine Süsstrunk. Frequency-tuned salient region detection. In *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition*, number CONF, pages 1597–1604, 2009. 6
- [2] Edward Adelson, Charles Anderson, James Bergen, Peter Burt, and Joan Ogden. Pyramid methods in image processing. *RCA Eng.*, 29, 11 1983. 3
- [3] 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*, 40(4):834–848, 2017. 9
- [4] Shuhan Chen, Xiuli Tan, Ben Wang, and Xuelong Hu. Reverse attention for salient object detection. In *Proceedings of European Conference on Computer Vision*, pages 234–250, 2018. 13
- [5] Zuyao Chen, Qianqian Xu, Runmin Cong, and Qingming Huang. Global context-aware progressive aggregation network for salient object detection. In *AAAI Conference on Artificial Intelligence*, volume 34, pages 10599–10606, 2020. 13
- [6] Deng-Ping Fan, Ming-Ming Cheng, Yun Liu, Tao Li, and Ali Borji. Structure-measure: A new way to evaluate foreground maps. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 4548–4557, 2017. 6
- [7] 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*, pages 698–704, 2018. 6
- [8] Deng-Ping Fan, Ge-Peng Ji, Ming-Ming Cheng, and Ling Shao. Concealed object detection. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2021. 1
- [9] Deng-Ping Fan, Ge-Peng Ji, Guolei Sun, Ming-Ming Cheng, Jianbing Shen, and Ling Shao. Camouflaged object detection. In *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition*, pages 2777–2787, 2020. 2, 3, 5, 6, 8, 10
- [10] Deng-Ping Fan, Ge-Peng Ji, Tao Zhou, Geng Chen, Huazhu Fu, Jianbing Shen, and Ling Shao. Prananet: Parallel reverse attention network for polyp segmentation. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 263–273, 2020. 1, 5
- [11] Deng-Ping Fan, Tao Zhou, Ge-Peng Ji, Yi Zhou, Geng Chen, Huazhu Fu, Jianbing Shen, and Ling Shao. Inf-net: Automatic covid-19 lung infection segmentation from ct images. *IEEE Transactions on Medical Imaging*, 39(8):2626–2637, 2020. 1
- [12] Shang-Hua Gao, Yong-Qiang Tan, Ming-Ming Cheng, Chengze Lu, Yunpeng Chen, and Shuicheng Yan. Highly ef-

- ficient salient object detection with 100k parameters. In Andrea Vedaldi, Horst Bischof, Thomas Brox, and Jan-Michael Frahm, editors, *Proceedings of European Conference on Computer Vision*, pages 702–721, Cham, 2020. Springer International Publishing. 5
- [13] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Spatial pyramid pooling in deep convolutional networks for visual recognition. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 37(9):1904–1916, 2015. 3
- [14] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition*, pages 770–778, 2016. 4, 9, 13
- [15] Wei Ji, Ge Yan, Jingjing Li, Yongri Piao, Shunyu Yao, Miao Zhang, Li Cheng, and Huchuan Lu. Dmra: Depth-induced multi-scale recurrent attention network for rgb-d saliency detection. *IEEE Transactions on Image Processing*, pages 1–16, 2022. 3
- [16] Trung-Nghia Le, Tam V. Nguyen, Zhongliang Nie, Minh-Triet Tran, and Akihiro Sugimoto. Anabran network for camouflaged object segmentation. *Computer Vision and Image Understanding*, 184:45–56, 2019. 2, 5, 6
- [17] Aixuan Li, Jing Zhang, Yunqiu Lyu, Bowen Liu, Tong Zhang, and Yuchao Dai. Uncertainty-aware joint salient object and camouflaged object detection. In *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition*, 2021. 2, 5, 10
- [18] Guanbin Li and Yizhou Yu. Visual saliency based on multiscale deep features. In *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition*, pages 5455–5463, 2015. 12
- [19] Tsung-Yi Lin, Piotr Dollár, Ross Girshick, Kaiming He, Bharath Hariharan, and Serge Belongie. Feature pyramid networks for object detection. In *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition*, pages 2117–2125, 2017. 3
- [20] Tony Lindeberg. *Scale-Space Theory in Computer Vision*, pages 187–226. 01 1994. 2
- [21] Tony Lindeberg. Feature detection with automatic scale selection. *International Journal of Computer Vision*, 30(2):79–116, Nov 1998. 2
- [22] Jiang-Jiang Liu, Qibin Hou, and Ming-Ming Cheng. Dynamic feature integration for simultaneous detection of salient object, edge and skeleton. *IEEE Transactions on Image Processing*, pages 1–15, 2020. 13
- [23] 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 *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition*, pages 3912–3921, 2019. 3, 5, 13
- [24] Nian Liu, Junwei Han, and Ming-Hsuan Yang. Picanet: Learning pixel-wise contextual attention for saliency detection. In *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition*, pages 3089–3098, 2018. 5, 13
- [25] Nian Liu, Ni Zhang, Kaiyuan Wan, Ling Shao, and Junwei Han. Visual saliency transformer. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 4722–4732, October 2021. 13
- [26] Yun Liu, Xin-Yu Zhang, Jia-Wang Bian, Le Zhang, and Ming-Ming Cheng. SAMNet: Stereoscopically attentive multi-scale network for lightweight salient object detection. *IEEE Transactions on Image Processing*, 30:3804–3814, 2021. 13
- [27] Jonathan Long, Evan Shelhamer, and Trevor Darrell. Fully convolutional networks for semantic segmentation. In *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition*, pages 3431–3440, 2015. 3
- [28] Zhiming Luo, Akshaya Mishra, Andrew Achkar, Justin Eichel, Shaozi Li, and Pierre-Marc Jodoin. Non-local deep features for salient object detection. In *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition*, pages 6593–6601, Los Alamitos, CA, USA, jul 2017. IEEE Computer Society. 5
- [29] Yunqiu Lyu, Jing Zhang, Yuchao Dai, Aixuan Li, Bowen Liu, Nick Barnes, and Deng-Ping Fan. Simultaneously localize, segment and rank the camouflaged objects. In *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition*, 2021. 2, 5, 6, 10
- [30] Ran Margolin, Lihi Zelnik-Manor, and Ayellet Tal. How to evaluate foreground maps? In *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition*, pages 248–255, 2014. 6
- [31] Haiyang Mei, Ge-Peng Ji, Ziqi Wei, Xin Yang, Xiaopeng Wei, and Deng-Ping Fan. Camouflaged object segmentation with distraction mining. In *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition*, June 2021. 2, 3, 5, 6, 10
- [32] Youwei Pang. Pysodevaltoolkit: A python-based evaluation toolbox for salient object detection and camouflaged object detection. <https://github.com/lartpang/PySODEvalToolkit>, 2020. 6
- [33] Youwei Pang. Pysodmetrics: A simple and efficient implementation of sod metrics. <https://github.com/lartpang/PySODMetrics>, 2020. 6
- [34] Youwei Pang, Lihe Zhang, Xiaoqi Zhao, and Huchuan Lu. Hierarchical dynamic filtering network for rgb-d salient object detection. In *Proceedings of European Conference on Computer Vision*, pages 235–252, 2020. 3
- [35] Youwei Pang, Xiaoqi Zhao, Lihe Zhang, and Huchuan Lu. Multi-scale interactive network for salient object detection. In *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition*, pages 9410–9419, June 2020. 3, 5, 12, 13
- [36] Ricardo Pérez-de la Fuente, Xavier Delclòs, Enrique Peñalver, Mariela Speranza, Jacek Wierzbos, Carmen Ascaso, and Michael S Engel. Early evolution and ecology of camouflage in insects. *PNAS*, 109(52):21414–21419, 2012. 1
- [37] Xuebin Qin, Zichen Zhang, Chenyang Huang, Chao Gao, Masood Dehghan, and Martin Jagersand. Basnet: Boundary-aware salient object detection. In *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition*, pages 7479–7489, 2019. 5, 13

- [38] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster r-cnn: towards real-time object detection with region proposal networks. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 39(6):1137–1149, 2016. **3**
- [39] 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. **3**
- [40] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014. **13**
- [41] Przemysław Skurowski, Hassan Abdulameer, Jakub Błaszczak, Tomasz Depta, Adam Kornacki, and Przemysław Kozieł. Animal camouflage analysis: Chameleon database, 2017. <http://kgwisc.aei.polsl.pl/index.php/pl/dataset/63-animal-camouflage-analysis>. **6, 11**
- [42] Martin Stevens and Sami Merilaita. Animal camouflage: Current issues and new perspectives. *Philosophical transactions of the Royal Society of London. Series B, Biological sciences*, 364:423–7, 12 2008. **2**
- [43] Yujia Sun, Geng Chen, Tao Zhou, Yi Zhang, and Nian Liu. Context-aware cross-level fusion network for camouflaged object detection. In Zhi-Hua Zhou, editor, *International Joint Conference on Artificial Intelligence*, pages 1025–1031. ijcai.org, 2021. **2, 3, 5, 10, 14**
- [44] Lijun Wang, Huchuan Lu, Yifan Wang, Mengyang Feng, Dong Wang, Baocai Yin, and Xiang Ruan. Learning to detect salient objects with image-level supervision. In *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition*, pages 136–145, 2017. **12**
- [45] Wenguan Wang, Shuyang Zhao, Jianbing Shen, Steven CH Hoi, and Ali Borji. Salient object detection with pyramid attention and salient edges. In *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition*, pages 1448–1457, 2019. **13**
- [46] Jun Wei, Shuhui Wang, and Qingming Huang. F³net: Fusion, feedback and focus for salient object detection. In *AAAI Conference on Artificial Intelligence*, volume 34, pages 12321–12328, 2020. **5, 12, 13**
- [47] Jun Wei, Shuhui Wang, Zhe Wu, Chi Su, Qingming Huang, and Qi Tian. Label decoupling framework for salient object detection. In *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition*, pages 13022–13031, 2020. **12, 13**
- [48] A. Witkin. Scale-space filtering: A new approach to multi-scale description. In *International Conference on Acoustics, Speech and Signal Processing*, volume 9, pages 150–153, March 1984. **2**
- [49] Runmin Wu, Mengyang Feng, Wenlong Guan, Dong Wang, Huchuan Lu, and Errui Ding. A mutual learning method for salient object detection with intertwined multi-supervision. In *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition*, pages 8150–8159, 2019. **13**
- [50] Zhe Wu, Li Su, and Qingming Huang. Cascaded partial decoder for fast and accurate salient object detection. In *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition*, pages 3907–3916, 2019. **5, 13**
- [51] Zhe Wu, Li Su, and Qingming Huang. Stacked cross refinement network for edge-aware salient object detection. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 7264–7273, 2019. **5, 12, 13**
- [52] Binwei Xu, Haoran Liang, Ronghua Liang, and Peng Chen. Locate globally, segment locally: A progressive architecture with knowledge review network for salient object detection. *AAAI Conference on Artificial Intelligence*, 35(4):3004–3012, May 2021. **13**
- [53] Jinnan Yan, Trung-Nghia Le, Khanh-Duy Nguyen, Minh-Triet Tran, Thanh-Toan Do, and Tam V. Nguyen. Mirror-net: Bio-inspired camouflaged object segmentation. *IEEE Access*, 9:43290–43300, 2021. **2, 5**
- [54] Qiong Yan, Li Xu, Jianping Shi, and Jiaya Jia. Hierarchical saliency detection. In *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition*, pages 1155–1162, 2013. **12**
- [55] Chuan Yang, Lihe Zhang, Huchuan Lu, Xiang Ruan, and Ming-Hsuan Yang. Saliency detection via graph-based manifold ranking. In *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition*, pages 3166–3173, 2013. **12**
- [56] 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 *Proceedings of the IEEE International Conference on Computer Vision*, pages 4146–4155, October 2021. **2, 5, 10**
- [57] Li Yuan, Yunpeng Chen, Tao Wang, Weihao Yu, Yujun Shi, Zi-Hang Jiang, Francis E.H. Tay, Jiashi Feng, and Shuicheng Yan. Tokens-to-token vit: Training vision transformers from scratch on imagenet. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 558–567, October 2021. **13**
- [58] Yi Zeng, Pingping Zhang, Jianming Zhang, Zhe Lin, and Huchuan Lu. Towards high-resolution salient object detection. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 7233–7242, 2019. **13**
- [59] Qiang Zhai, Xin Li, Fan Yang, Chenglizhao Chen, Hong Cheng, and Deng-Ping Fan. Mutual graph learning for camouflaged object detection. In *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition*, 2021. **2, 3, 5, 6, 10**
- [60] 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 *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition*, pages 8579–8588, 2020. **5**
- [61] Jing Zhang, Xin Yu, Aixuan Li, Peipei Song, Bowen Liu, and Yuchao Dai. Weakly-supervised salient object detection via scribble annotations. In *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition*, 2020. **5**
- [62] Miao Zhang, Tingwei Liu, Yongri Piao, Shunyu Yao, and Huchuan Lu. *Auto-MSFNet: Search Multi-Scale Fusion Network for Salient Object Detection*, page 667–676. Association

tion for Computing Machinery, New York, NY, USA, 2021. 13

- [63] Jia-Xing Zhao, Jiang-Jiang Liu, Deng-Ping Fan, Yang Cao, Jufeng Yang, and Ming-Ming Cheng. Egnet: Edge guidance network for salient object detection. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 8779–8788, 2019. 5, 13
- [64] Ting Zhao and Xiangqian Wu. Pyramid feature attention network for saliency detection. In *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition*, pages 3080–3089, 2019. 12
- [65] Xiaoqi Zhao, Youwei Pang, Jiaxing Yang, Lihe Zhang, and Huchuan Lu. Multi-source fusion and automatic predictor selection for zero-shot video object segmentation. In *Proceedings of the ACM International Conference on Multimedia*, pages 2645–2653, 2021. 3
- [66] Xiaoqi Zhao, Youwei Pang, Lihe Zhang, Huchuan Lu, and Xiang Ruan. Self-supervised pretraining for rgb-d salient object detection. In *AAAI Conference on Artificial Intelligence*, 2022. 3
- [67] Xiaoqi Zhao, Youwei Pang, Lihe Zhang, Huchuan Lu, and Lei Zhang. Suppress and balance: A simple gated network for salient object detection. In *Proceedings of European Conference on Computer Vision*, pages 35–51, 2020. 3, 9, 13
- [68] Xiaoqi Zhao, Lihe Zhang, and Huchuan Lu. Automatic polyp segmentation via multi-scale subtraction network. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*. Springer, 2021. 1
- [69] Xiaoqi Zhao, Lihe Zhang, Youwei Pang, Huchuan Lu, and Lei Zhang. A single stream network for robust and real-time rgb-d salient object detection. In *Proceedings of European Conference on Computer Vision*, pages 646–662, 2020. 3
- [70] Zhirui Zhao, Changqun Xia, Chenxi Xie, and Jia Li. *Complementary Trilateral Decoder for Fast and Accurate Salient Object Detection*, page 4967–4975. Association for Computing Machinery, New York, NY, USA, 2021. 13
- [71] Yunfei Zheng, Xiongwei Zhang, Feng Wang, Tieyong Cao, Meng Sun, and Xiaobing Wang. Detection of people with camouflage pattern via dense deconvolution network. *IEEE Signal Processing Letters*, 26(1):29–33, Jan 2019. 11, 12
- [72] Huajun Zhou, Xiaohua Xie, Jian-Huang Lai, Zixuan Chen, and Lingxiao Yang. Interactive two-stream decoder for accurate and fast saliency detection. In *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition*, pages 9138–9147, 2020. 5, 13