# To be Critical: Self-Calibrated Weakly Supervised Learning for Salient Object Detection

Yongri Piao ⓘ, Jian Wang, Miao Zhang ⓘ, Zhengxuan Ma, and Huchuan Lu ⓘ, *Senior Member, IEEE*

*Abstract*—Weakly-supervised salient object detection (WSOD[1]) aims to develop saliency models using image-level annotations. Despite of the success of previous works, explorations on an effective training strategy for the saliency network and accurate matches between image-level annotations and salient objects are still inadequate. In this work, 1) we propose a self-calibrated training strategy by explicitly establishing a mutual calibration loop between pseudo labels and network predictions, liberating the saliency network from error-prone propagation caused by pseudo labels. 2) we prove that even a much smaller dataset (merely $1.8\%$ of ImageNet) with well-matched annotations can facilitate models to achieve better performance as well as generalizability. This sheds new light on the development of WSOD and encourages more contributions to the community. Comprehensive experiments demonstrate that our method outperforms all the existing WSOD methods by adopting the self-calibrated strategy only. Steady improvements are further achieved by training on the proposed dataset. Additionally, our method achieves $94.7\%$ of the performance of fully-supervised methods on average. And what is more, the fully supervised models adopting our predicted results as "ground truths" achieve successful results ($95.6\%$ for BASNet and $97.3\%$ for ITSD on F-measure), while costing only $0.32\%$ of labeling time for pixel-level annotation.

*Index Terms*—Salient object detection, Weakly supervised learning, Deep learning.

## I. INTRODUCTION

SALIENT object detection (SOD) aims to segment objects in an image that visually attract human attention most. It plays an important role in many computer vision and robotic vision tasks [1], such as image segmentation [2] and visual tracking [3]. Recently, deep learning based methods [4], [5], [6], [7], [8], [9], [10], [11] have proved its superiority and achieved remarkable progress. Success of those methods, however, heavily relies a large number of highly accurate pixel-level annotations, which are time-consuming and labor-intensive to collect. A trade-off between testing accuracy and training annotation cost has long existed in the SOD task.

To alleviate this predicament, several attempts have been made to explore different weakly supervised formats, such as noisy label [12], [13], scribble [14], [15] and image-level annotation (*i.e.*, classification label). Image-level annotation

[1]In this paper, we denote weakly supervised salient object detection methods using image-level labels as WSOD for convenience.



Previous work MSW [14]

Our model without the proposed SC

Our model with the proposed SC

data

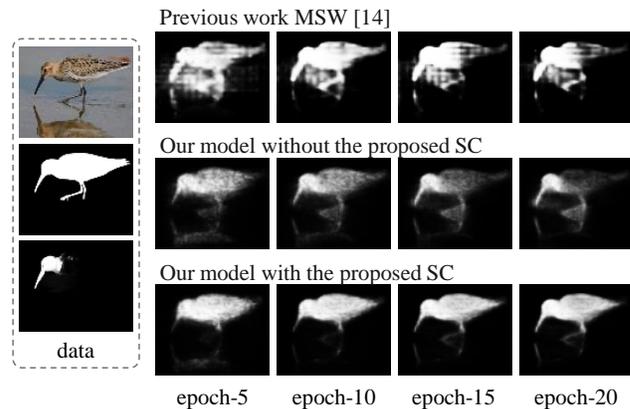epoch-5    epoch-10    epoch-15    epoch-20

Fig. 1. The visual saliency predictions during the training process of different models, in which **SC** represents our proposed self-calibrated training strategy. Column **data** represents image, ground truth and pseudo label, noting that the ground truth is just for exhibition and not used in our framework.

based WSOD methods usually adopt a two-stage scheme, which leverages a classification network to generate pseudo labels and then trains a saliency network on these labels. In this paper, we focus on this most challenging problem of developing WSOD by only using image-level annotation.

Some pioneering works [16], [17], [18] pursue accurate pseudo labels to train a saliency network and achieve good performance. However, given the fact that pseudo labels are still a far cry from the ground truths, the error remaining unaddressed in the pseudo labels can propagate to the generated predictions. This is consistent with the fact that as the number of epochs increases, the parameters of the model are updated and the prediction curve goes from underfitting to optimal to overfitting. Interestingly, we observe that the relatively good results containing global representations of saliency can be predicted at the early training process (*e.g.*, epoch-5), while the predictions are more prone to error at the latter training process (*e.g.*, epoch-20), as shown in the first two rows in Figure 1. This inspires us to go one step further exploring how this global representation can be evolved as the model is properly trained.

Moreover, previous works adopt existing large-scale datasets, *e.g.*, ImageNet [19] and COCO [20], to perform WSOD. However, an observable fact should not be ignored that there is an inherent inconsistency between image classification and SOD task. For example, many classification labels do not match the salient objects in both single-object and multi-object cases in ImageNet, as illustrated in Figure 2. Such cross-domain inconsistency caused by those mismatched samples impairs the generalizability of models and prevents WSOD methods from achieving optimal results.
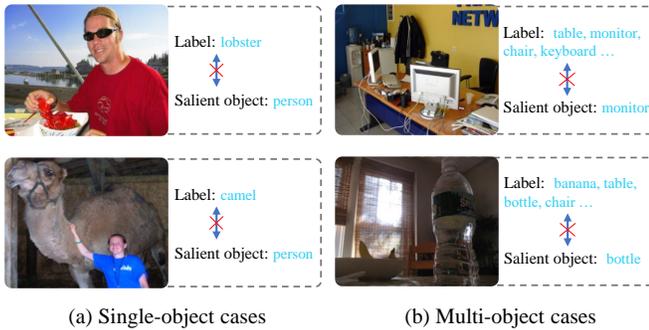
| (a) Single-object cases | (b) Multi-object cases |

Fig. 2.   Cross-domain inconsistency between ImageNet dataset and salient object detecion. (a) and (b) represent single-object and multi-object cases, respectively.

In this work, our core insight is that we can design a self-calibrated training strategy and exploit saliency-based image-level annotations to address the aforementioned challenges. To be specific, we **1)** aim to calibrate our network with progressively updated labels to curb the spread of errors in low-quality pseudo labels during the training process. **2)** develop reliable matches for which image-level annotations are correctly corresponding to salient objects. The source code will be released upon publication. Concretely, our contributions are as follows:

- We propose a self-calibrated training strategy to prevent the network from propagating the negative influence of error-prone pseudo labels. A mutual calibration loop is established between pseudo labels and network predictions to promote each other.
- We open up a fresh perspective on that even a much smaller dataset (merely 1.8% of ImageNet) with well-matched image-level annotations allows WSOD to achieve better performance. This encourages more existing data to be correctly annotated and further paves the way for the booming future of WSOD.
- Our method outperforms existing WSOD methods on all metrics over five benchmark datasets, and meanwhile achieves averagely 94.7% performance of state-of-the-art fully supervised methods. We also demonstrate that our method retains its competitive edge on most metrics even without our proposed dataset.
- We extend the proposed method to other fully supervised SOD methods. Our offered pseudo labels enable these methods to achieve comparatively high accuracy (95.6% for BASNet [21] and 97.3% for ITSD [22] on F-measure) while being free of pixel-level annotations, costing only 0.32% of labeling time for pixel-level annotation.

## II. RELATED WORK

### A. Salient Object Detection

Early SOD methods mainly focus on detecting salient objects by utilizing handcraft features and setting various priors, such as center prior [23], boundary prior [24] and so on [25], [26]. Recently, deep learning based methods demonstrate its advantages and achieve remarkable improvements. Plenty of promising works [27], [28], [29], [30], [31] are proposed and present various effective architectures. Among them, Hou *et al.*[27] present short connections to integrate the low-level and high-level features, and predict more detailed saliency maps. Wu *et al.*[28] propose a novel cascaded partial decoder framework and utilize generated relatively precise attention map to refine high-level features. in [29], [30], researchers propose to explore boundary of the salient objects to predict a more detailed prediction. Although appealing performance these methods have achieved, vast high-quality pixel-level annotations are needed for training their models, which are time-consuming and laborious.

### B. Weakly Supervised Salient Object Detection

For achieving a trade-off between labeling efficiency and model performance, researchers aim to perform salient object detect with low-cost annotations. To this end, WSOD is presented and achieves an appealing performance with image-level annotations only.

Wang *et al.*[16] design a foreground inference network (FIN) to predict saliency maps from image-level annotations, and introduce a global smooth pooling (GSP) to combine the advantages of global average pooling (GAP) and global max pooling (GMP), which explicitly computes the activation of salient objects. In [17], Li *et al.*also perform WSOD based on image-level annotations, they adopt a recurrent self-training strategy and propose a conditional random field based graphical model to cleanse the noisy pixel-wise annotations by enhancing the spatial coherence as well as salient object localization. Based on a traditional method MB+ [32], more accurate saliency maps are generated in less than one second per image. Zeng *et al.*[18] intelligently utilize multiple annotations (*i.e.*,, classification and caption annotations) and design a multi-source weak supervision framework to integrate information from various annotations. Benefited from multiple annotations and an interactive training strategy, a really sample saliency network can also achieve appealing performance. All the above methods target to train a classification network (on existing large-scale multiple objects dataset, *i.e.*,, ImageNet [19] or Microsoft COCO [20]) to generate class activation maps (CAMs) [33], then perform different refinement methods to generate pseudo labels. Supervised by these pseudo labels directly, a saliency network is trained and predicts the final saliency maps.

Different from the aforementioned works, we argue that: **1)** Developing an effective training strategy encourages more accurate predictions even under the supervision of inaccurate pseudo labels which would mislead the networks. **2)** Establishing accurate matches between classification labels and salient objects could facilitate the further development of WSOD.

## III. THE PROPOSED METHOD

In this section, we describe the details of our two-stage framework. As illustrated in the Figure 3, in the first training stage, we train a normal classification network based on the proposed saliency-based dataset, to generate more accurate pseudo labels. We then develop a saliency network using the pseudo labels in the second stage. A self-calibrated training
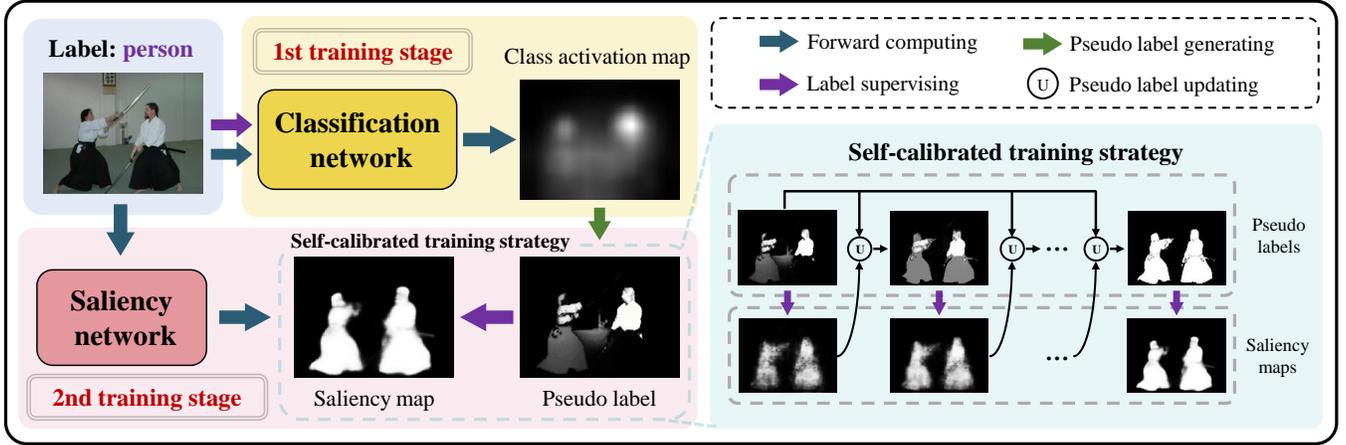
Fig. 3. Overall framework of our proposed method. In the first stage, classification labels are used to supervise classification network to generate CAMs and further produce pseudo labels. In the second stage, we train a saliency network with the above pseudo labels and propose a self-calibrated strategy to correct labels and predictions progressively.

strategy is proposed in this stage to immune network from inaccurate pseudo labels and encourage more accurate predictions.

### A. From Image-level to Pixel-level

Class activation maps (CAMs) [33] localize the most discriminative regions in an image using only a normal classification network and build a preliminary bridge from image-level annotations to pixel-level segmentation tasks. In this paper, we adopt CAMs following the same setting of [34], to generate pixel-level pseudo labels in the first training stage. To better understand our proposed approach, we will describe the generation of CAMs in a brief way.

For a classification network, we discard all the fully connected layers and apply an extra global average pooling (GAP) layer as well as a convolution layer as previous works do. In the training phase, we take images in classification dataset as input, and compute its classification scores $Cls$ as follows:

$$Cls = w_s{}^T * GAP(F^5) + b_s, \tag{1}$$

where $F^5$ represents the features from the last convolution block, $GAP(\cdot)$ denotes the global average pooling operation and $w_s^T$ as well as $b_s$ are the learnable parameters of the convolution layer. In the inference phase, we compute the CAMs of images in DUTS-Train dataset as follows:

$$C_{AM} = \sum_{k=1}^{N} Cls_k * Norm(Relu(w_s{}^T * F^5 + b_s)_k), \tag{2}$$

where $Relu(\cdot)$ and $Norm(\cdot)$ denote the relu activation function and normalization function, respectively. $w_s^T$ and $b_s$ are the shared parameters learnd in the training phase, $Cls_k$ represents the classification scores for category $k$ and $N$ represents the total number of categories. In this phase, multi-scale inference strategy is adopted, which rescales the original

image into four sizes and computes the average CAMs as the final output.

As Ahn *et al.*[34] have pointed out, CAMs mainly concentrate on the most discriminative regions and are too coarse to serve as pseudo labels. Various refinements have been conducted to generate pseudo labels. Different from [16], [18] using an clustering algorithm SLIC [35], a plug-and-play module PAMR [36] is adopt in our method. It performs refinement using the low-level color information of RGB images, which can be inserted into our framework flexibly and efficiently. Following the settings of [16], [18], we also adopt CRF [37] for a further refinement. Note that it is only used to generate pseudo labels in our method.

### B. Self-calibrated Training Strategy

In the second training stage, a saliency network is trained with the pseudo labels generated in the first training stage. As is mentioned above, the relatively good results containing global representations of saliency is gradually degraded as the training process continues. A straightforward method to tackle this dilemma is setting a validation set to pick the best result during the training process. However, we argue that it may lead to sub-optimal results because **1)** despite good saliency representations are learned at the early training stage, the predictions are coarse and lack detail as the loss function is still converging (as shown in the $2^{nd}$ row in Figure 1). **2)** the capability of networks to learn saliency representation is not fully excavated. **3)** we believe that WSOD should not use any pixel-level ground truth in the training process, even as a validation set. Following this main idea, we propose to establish a mutual calibration loop during the training process in which error-prone pseudo labels are recursively updated and calibrate network for better predictions in turn.

**Insight:** As is discussed in the Section I, under the supervision of noisy pseudo labels, the saliency network goes from optimal to overfitting. On the one hand, in our weakly supervised settings, this "overfitting" manifestes itself as the

**Algorithm 1** Self-calibrated training strategy

---

**Require:**   The images from DUTS-Train dataset, $I_n$;   The predictions of saliency network, $P_n$;   The original pseudo labels generated in the $1_{st}$ stage, $Y_1$.
**Ensure:**   the updated pseudo labels, $Y_{n+1}$.
 1: Performing $2_{nd}$ training stage, maximum epoch is $N$.
 2: **for** $n = 1$ to $N$ **do**
 3:    Refined predictions: $P_n^{'} = \text{PAMR}(P_n, I_n)$;
 4:    **if** $P_n^{'}(x, y) > 0.4$ **then**
 5:       $P_n^{'}(x, y) = 1$
 6:    **else**
 7:       $P_n^{'}(x, y) = 0$
 8:    **end if**
 9:    weighting factor $\lambda = (n/N)^{0.5}$;
10:    Updating pseudo labels: $Y_{n+1} = Y_1 * (1 - \lambda) + P_n^{'} * \lambda$;
11: **end for**

---

network being affected by the noisy pseudo labels and learning the inaccurate noise information in them, which heavily restricts the performance of WSOD. It is also worth to mention that this is fundamentally different from the "overfitting" in supervised learning, the latter means that the network learns the biased information in a less comprehensive training set. On the other hand, we conclude reasons of the optimal point before overfitting as: 1) Although many pseudo labels are noisy and inaccurate, the whole pseudo labels still describe general saliency cues. It can provide a roughly correct guidance for the saliency network. 2) Before the loss converges, the saliency network is prone to learn the regular and generalized saliency cues rather than the irregular and noisy information in pseudo labels. Such kind of robustness is also discussed in [38]. Motivated by the above analyses, we propose a self-calibrated training strategy to effectively utilize the robustness and tackle the negative overfitting.

To be specific, supervised by inaccurate pseudo labels $Y$, we take the predictions $P$ of the saliency network as saliency seeds. As is illustrated in Figure 3, coarse but more accurate seeds are predicted during the first few epochs regardless of the inaccurate supervision of error-prone pseudo labels. We take these seeds as correction terms to calibrate and update the original pseudo labels $Y$, while performing refinement again with PAMR. Detailed procedure is presentd in Algorithm 1, here we set a threshold to $0.4$ for the binarization operation on refined predictions $P^{'}$. We conduct self-calibrated strategy throughout the training process, that is, it is performed on each training batch. The loss function for this training stage can be described as:

$$
\begin{aligned}
L(P, Y) = &-\sum_{i=1} ((1 - \lambda)y_i + \lambda p_i^{'}) * \log p_i \\
&-((1 - \lambda)(1 - y_i) + \lambda(1 - p_i^{'})) * \log(1 - p_i),
\end{aligned} \tag{3}
$$

where $\lambda$ is the weighting factor that is illustrated in Algorithm 1. The intuition is that as the training process goes on, the saliency prediction is more accurate and larger weight should be given. $y_i$, $p_i$ and $p_i^{'}$ represent the elements of $Y$, $P$ and refined predictions $P^{'}$, respectively.

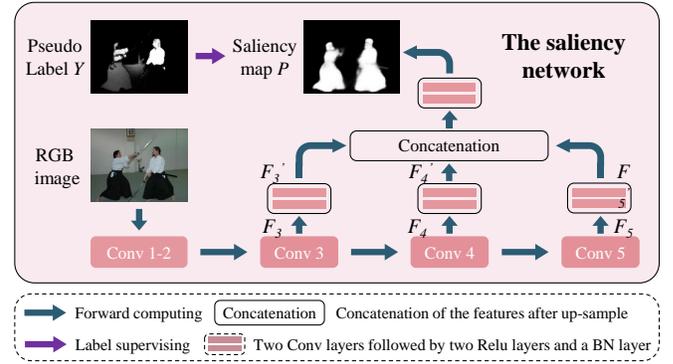As is illustrated in the Figure 3, equipped with our proposed



Fig. 4.   Detailed structure of our saliency network. We adopt a simple encoder-decoder architecture and take prediction $P$ as our final result.

self-calibrated training strategy, inaccurate pseudo labels are progressively updated, and in turn supervise the network. This mutual calibration loop finally encourages both accurate pseudo labels and predictions.

*C. Saliency Network*

As for the saliency network, we adopt a simple encoder-decoder architecture without any auxiliary modules, which is usually served as baseline for fully-supervised SOD methods [27], [28]. As illustrated in Figure 4, for an image from DUTS-Train dataset, we take features $F_3$, $F_4$ and $F_5$ from the encoder, to generate $F_3^{'}$, $F_4^{'}$ and $F_5^{'}$ through two convolution layers, and then adopt a bottom-up strategy to perform feature fusion, which can be denoted as:

$$
P = \sigma(Conv(Cat(Up(F_5^{'}), Up(F_4^{'}), F_3^{'}))), \tag{4}
$$

where $\sigma(\cdot)$ represents the sigmoid function, $Conv(\cdot)$ and $Cat(\cdot)$ denote the convolution and concatenation operation, respectively. $Up(\cdot)$ represents upsampling feature maps to the same size.

In the decoder, the number of output channels of all the middle convolution layers are set to 64 for acceleration. Note that our final prediction $P$ is predicted in an end-to-end manner in the test phase without any post-processing.

## IV. DATASET CONSTRUCTION

To explore the advantages of accurate matches between image-level annotations and corresponding salient objects, we establish a saliency-based classification dataset, which ensures all the classification labels correspond to the salient objects. Following this main idea, we relabel an existing widely-adopted saliency training set DUTS-Train [16] with well-matched image-level annotations, namely DUTS-Cls dataset. It fits with WSOD better than existing large-scale classification datasets due to the accurate matches, and facilitates the further improvements for WSOD.

To be specific, we select and label images in DUTS-Train with image-level annotations, while discarding rare categories because only several images are contained. The proposed DUTS-Cls dataset contains 44 categories and 5959 images. As is illustrated in Figure 5, it reaches a relative equilibrium
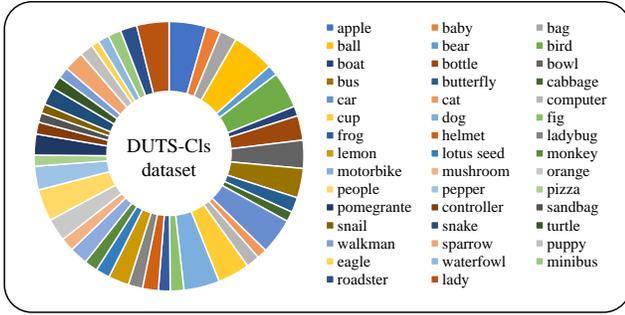
Fig. 5. Our introduced DUTS-Cls is a saliency-based dataset with image-level annotations, containing 44 categories and 5959 images, in which all the classification labels correspond to the most salient objects in images.

in terms of image numbers of each category and covers most common categories.

It is worth mentioning that labeling image-level annotations is quite fast, which only takes less than 1 seconds per image. Compared to about 3 minutes [39] for labeling a pixel-level ground truth, it takes less than $0.56\%$ of the time and labor cost for a sample. Annotating DUTS-Cls dataset (5959 samples) only costs $0.32\%$ of labeling time than annotating the whole DUTS-Train dataset (10553 samples) with pixel-level ground truth. This indicates that exploring WSOD with image-level annotation is quite efficient. Moreover, the DUTS-Cls dataset with well-matched image-level annotations offers a better choice for WSOD than ImageNet, and we genuinely hope it could contribute to the community and encourage more existing data to be correctly annotated at image level.

## V. EXPERIMENTS

### A. Implementation Details

We implement our method on the Pytorch toolbox with a single RTX 2080Ti GPU. The backbone adopted in our method is DenseNet-169 [40], which is same as the latest work [18]. During the first training stage, we train a classification network on our proposed DUTS-Cls dataset. In this stage, we adopt the Adam optimization algorithm [41], the learning rate is set to 1e-4 and maximum epoch is set to 20. In the second training stage, we only take the RGB images from DUTS-Train as our training set. In this stage, we use Adam optimization algorithm with the learning rate 3e-6 and maximum epoch 25. The batch size of both training stages is set to 20 and all the training and testing images are resized to $256 \times 256$.

**Hyperparameters setting.** For the weighting factor $\lambda$ of self-calibrated strategy, we conduct hyper-parameter experiments on ECSSD [42] dataset to pick the optimal value through F-measure [43]. According to the results (0.5 to 0.848, 0.6 to 0.853 and 0.7 to 0.849), we finally set the hyper-parameter $\lambda$ to 0.6.

### B. Datasets and Evaluation Metrics

For a fair comparison, we train our model on ImageNet and our proposed DUTS-Cls dataset respectively, the results are shown in Table I. We conduct comparisons on five following widely-adopted test datasets. ECSSD [42]: contains 1000 images which cover various scenes. DUT-OMRON [24]: includes 5168 challenging images consisting of single or multiple salient objects with complex contours and backgrounds. PASCAL-S [2]: is collected from the validation set of the PASCAL VOC semantic segmentation dataset [44], and contains 850 challenging images. HKU-IS [45]: includes 4447 images, many of which contain multiple disconnected salient objects. DUTS [16]: is the largest salient object detection benchmark, which contains 10553 training samples (DUTS-Train) and 5019 testing samples (DUTS-Test). Most images in DUTS-Test are challenging with various locations and scales.

To evaluate our method in a comprehensive and reliable way, we adopt four well-accepted metrics, including S-measure [46], E-measure [47], F-measure [43] as well as Mean Absolute Error (MAE).

### C. Comparison with State-of-the-arts

We compare our method with all the existing image-level annotation based WSOD methods: WSS [16], ASMO [17] and MSW [18]. To further demonstrate the effectiveness of our weakly supervised methods, we also compare the proposed method with nine state-of-the-art fully supervised methods including DSS [27], R$^3$Net [7], DGRL [48], BASNet [21], PFA [49], CPD [28], SCRN [50], ITSD [22] and MINet [51], all of which are trained on pixel-level ground truth and based on DNNs. For a fair comparison, we use the saliency maps provided by authors and perform the same evaluation code for all methods.

**Quantitative evaluation.** Table I shows the quantitative comparison on four evaluation metrics over five datasets. It can be seen that our method outperforms all the weakly supervised methods on all metrics. Especially, $31.0\%$ improvement on HKU-IS and $28.4\%$ on DUT-OMRON are achieved on MAE metric. Our method also improves the performance on two challenging datasets DUT-ORMON and PASCAL-S by a large margin, which indicates that our method can explore more accurate saliency cues even in complex scenes. **Additionally**, the proposed saliency-based dataset with well-matched image-level annotations enables our method to achieve better performance, while far less training samples (less than $1.45\%$ of the latest work MSW [18]) are required. To prove the effect of our method in a more objective manner, we also train our method on ImageNet dataset following the previous works. The results of "Ours-" shown in Table I demonstrate that our method can outperform existing methods on most metrics even without the proposed dataset thanks to the effective strategy. **Moreover**, we also compare our method with nine state-of-the-art fully supervised methods. It can be seen in Figure 7 that our method, even with the image-level annotations only and a simple baseline network without any auxiliary modules, can also achieve $94.7\%$ accuracy of fully supervised methods on average.

**Qualitative evaluation.** In Figure 6, we show the qualitative comparisons of our method with existing three WSOD methods as well as six state-of-the-art fully supervised methods. It can be seen that our method could discriminate salient objects from various challenging scenes (such as small objects case in

TABLE I

QUANTITATIVE COMPARISONS OF E-MEASURE ($E_s$), S-MEASURE ($S_\alpha$), F-MEASURE ($F_\beta$) AND MAE ($M$) METRICS OVER FIVE BENCHMARK DATASETS. THE SUPERVISION TYPE (**SUP.**) I INDICATES USING IMAGE-LEVEL ANNOTATIONS ONLY, AND I&C REPRESENTS DEVELOPING WSOD ON BOTH IMAGE-LEVEL ANNOTATIONS AND CAPTION ANNOTATIONS SIMULTANEOUSLY. **NUM.** REPRESENTS THE NUMBER OF TRAINING SAMPLES. - MEANS UNAVAILABLE RESULTS, OURS- AND OURS REPRESENT OUR METHOD TRAINED ON IMAGENET AND PROPOSED DUTS-CLS DATASET, RESPECTIVELY. THE BEST TWO RESULTS ARE MARKED IN **BOLDFACE** AND MAGENTA.

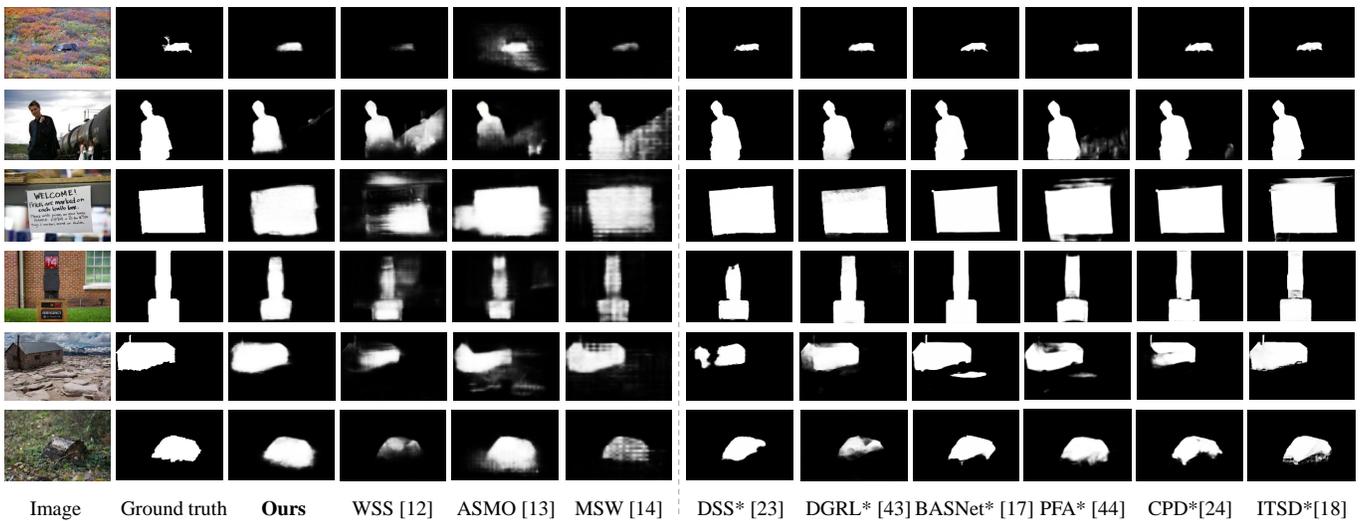| Methods | Sup. | ECSSD | | | | DUTS-Test | | | | HKU-IS | | | | DUT-OMRON | | | | PASCAL-S | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | $S_\alpha$ | $E_s$ | $F_\beta$ | $M$ | $S_\alpha$ | $E_s$ | $F_\beta$ | $M$ | $S_\alpha$ | $E_s$ | $F_\beta$ | $M$ | $S_\alpha$ | $E_s$ | $F_\beta$ | $M$ | $S_\alpha$ | $E_s$ | $F_\beta$ | $M$ |
| WSS [16] | I | .811 | .869 | .823 | .104 | .748 | .795 | .654 | .100 | .822 | .896 | .821 | .079 | .725 | .768 | .603 | .109 | .744 | .791 | .715 | .139 |
| ASMO [17] | I | .802 | .853 | .797 | .110 | .697 | .772 | .614 | .116 | - | - | - | - | .752 | .776 | .622 | .101 | .717 | .772 | .693 | .149 |
| MSW [18] | I&C | .827 | .884 | .840 | .096 | .759 | .814 | .684 | .091 | .818 | .895 | .814 | .084 | .756 | .763 | .609 | .109 | .768 | .790 | .713 | .133 |
| Ours- | I | .836 | .887 | .838 | .083 | .770 | **.830** | **.689** | .079 | .836 | .907 | .822 | .064 | .743 | .807 | .643 | .085 | .778 | .818 | .742 | .111 |
| Ours | I | **.858** | **.901** | **.853** | **.071** | **.776** | .829 | .688 | **.077** | **.850** | **.918** | **.835** | **.058** | **.766** | **.817** | **.667** | **.078** | **.781** | **.824** | **.749** | **.108** |



Fig. 6. Visual comparisons of our method with existing WSOD methods as well as six state-of-the-art fully supervised SOD methods (marked with *) in some challengling scenes.

Image   Ground truth   **Ours**   WSS [12]   ASMO [13]   MSW [14]   DSS* [23]   DGRL* [43]   BASNet* [17]   PFA* [44]   CPD*[24]   ITSD*[18]
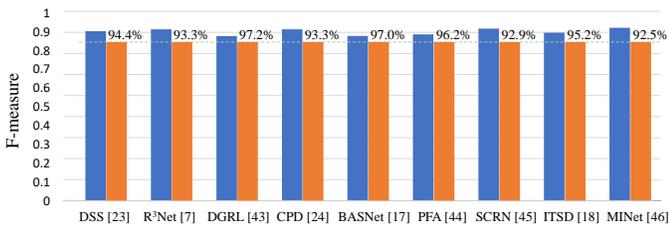


Fig. 7. Comparison of our method with 9 fully supervised methods on ECSSD dataset. The blue column represents the performance of each fully supervised methods and the orange one indicates ours. The corresponding data denote the percentages of performance of our method in different fully supervised methods.
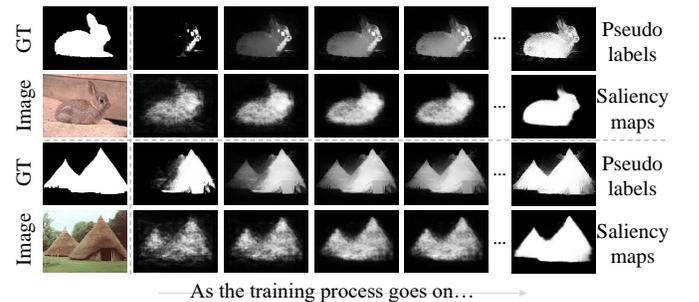


Fig. 8. Visual analysis of the effectiveness of our proposed self-calibrated strategy during the training process, noting that the ground truth is just for exhibition and not used in our framework.

the $1^{st}$ row and complex background cases in the $2^{nd}$ and $3^{rd}$ rows) and achieve more complete and accurate predictions. **Moreover**, compared with the fully supervised methods, our method also predicts comparable and even better results in some cases, such as the complete house and log in the $5^{th}$ and $6^{th}$ rows. But we would like to point out that our results also need to be improved in term of the boundary of the salient objects.

### D. Ablation Studies

**Effect of the self-calibrated strategy.** We conduct experiments on both ImageNet ($1^{st}$ and $3^{rd}$ rows) and DUTS-

TABLE II
QUANTITATIVE RESULTS OF ABLATION STUDIES. **DATASET** REPRESENTS DIFFERENT TRAINING SETS USED IN THE FIRST TRAINING STAGE.
**STRATEGY** DENOTES TRAINING STRATEGY USED IN THE SECOND STAGE, - INDICATES THE BASELINE MODEL WITHOUT ANY TRAINING STRATEGY
AND +SC REPRESENTS ADOPTING OUR PROPOSED SELF-CALIBRATED STRATEGY.

| Dataset | | Strategy | | ECSSD | | DUTS-Test | | HKU-IS | | DUT-OMRON | | PASCAL-S | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| ImageNet | DUTS-Cls | - | + SC | $F_\beta\uparrow$ | $M\downarrow$ | $F_\beta\uparrow$ | $M\downarrow$ | $F_\beta\uparrow$ | $M\downarrow$ | $F_\beta\uparrow$ | $M\downarrow$ | $F_\beta\uparrow$ | $M\downarrow$ |
| ✓ | | ✓ | | 0.776 | 0.121 | 0.642 | 0.094 | 0.773 | 0.090 | 0.568 | 0.111 | 0.694 | 0.140 |
| | ✓ | ✓ | | 0.836 | 0.096 | 0.675 | 0.085 | 0.822 | 0.075 | 0.648 | 0.083 | 0.735 | 0.126 |
| ✓ | | | ✓ | 0.838 | 0.083 | **0.689** | 0.079 | 0.822 | 0.064 | 0.643 | 0.085 | 0.742 | 0.111 |
| | ✓ | | ✓ | **0.853** | **0.071** | 0.688 | **0.077** | **0.835** | **0.058** | **0.667** | **0.078** | **0.749** | **0.108** |

TABLE III
THE EFFECTIVENESS OF OUR PROPOSED SELF-CALIBRATED STRATEGY
ON ECSSD DATASET. **+ SC** INDICATES SIMPLY APPLYING OUR
SELF-CALIBRATED STRATEGY DURING THE TRAINING PROCESS.

| Method | Strategy | $S_\alpha\uparrow$ | $E_s\uparrow$ | $F_\beta\uparrow$ | MAE $\downarrow$ |
|---|---|---|---|---|---|
| BSCA [52] | - | 0.846 | 0.884 | 0.814 | 0.084 |
| | + SC | **+0.007** | **+0.009** | **+0.018** | **-0.008** |
| MR [24] | - | 0.839 | 0.884 | 0.823 | 0.085 |
| | + SC | **+0.014** | **+0.010** | **+0.016** | **-0.009** |
| MSW [18] | - | 0.827 | 0.884 | 0.840 | 0.096 |
| | + SC | **+0.017** | **+0.012** | **+0.014** | **-0.014** |



Fig. 9. Visual analysis of effect of DUTS-Cls datset. $CAM_I$ and $CAM_D$ represent the CAMs generated by training on ImageNet and our DUTS-Cls dataset, respectively. Heatmap is adopted for better visualization.

Cls ($2^{nd}$ and $4^{th}$ rows) settings in Table II. It can be seen that the proposed self-calibrated strategy can not only enhance the performance of our method on ImageNet setting greatly, but also achieve great improvements even on the DUTS-Cls setting, especially on MAE metrics. **Besides**, the effectiveness of the proposed self-calibrated strategy can also be demonstrated by the visual results in Figure 8. It can be seen that the proposed strategy can keep and enhance the globally good representations during the training process, and predict accurate saliency maps even supervised by error-prone pseudo labels. **Moreover**, for a comprehensive evaluation, 1) We change the pseudo label by using two traditional SOD methods BSCA [52] and MR [24], and then train our model with and without the proposed strategy respectively, the results are shown in the first four rows in Table III. 2) We further apply our method on the lasted work MSW [18] by just adding our proposed strategy in the last two rows in Table III. These results strongly prove that the self-calibrated strategy can not only work well on our method, but also effective for other pseudo labels and other works.

**Effect of the DUTS-Cls dataset.** We introduce a saliency-based dataset with well-matched image-level annotations to offer a better choice for WSOD. The first two rows in Table II demonstrate that DUTS-Cls dataset encourages the baseline model to achieve remarkable improvements, compared to ImageNet dataset. And as is illustrated in the last two rows in Table II, it also proves its superiority by a steady improvement on most metrics even if good enough performance is already achieved by adopting the self-calibrated strategy. This is con-
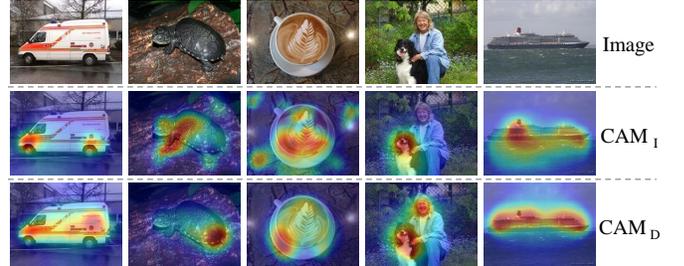
sistent with our argument that the cross-domain inconsistency does impede the performance of WSOD, and a saliency-based dataset can settle this matter better. **Additionally**, we visualize the CAMs trained on ImageNet (named $CAM_I$) and DUTS-Cls (named $CAM_D$) in Figure 9, it can be seen that $CAM_D$ have higher activation level within the salient objects trained on well-matched DUTS-Cls dataset. **Last but not least**, to further prove the effectiveness of the proposed DUTS-Cls dataset objectively, we also train the latest work MSW [18] on the DUTS-Cls dataset. As is shown in Figure 10, by simply replacing ImageNet with DUTS-Cls, considerable improvements are achieved in less training iterations. It is worth to mention that the DUTS-Cls dataset reaches less than 1.8% percent of ImageNet in terms of sample size. This strongly demonstrates the effectiveness and generalizability of the well-matched DUTS-Cls dataset for WSOD.

### E. Effectiveness on Unseen Category

The category number of classification dataset inevitably influences the performance of WSOD. Unlike ImageNet including 200 various categories, our proposed DUTS-Cls dataset only contains 44 categories. It is necessary to evaluate the effectiveness of our method as well as DUTS-Cls dataset on unseen categories.

To this end, we choose THUR [53] as the benchmark dataset for this experiment. THUR is a high-quality saliency dataset which consists of five categories including butterfly, coffee mug, dog, giraffe and airplane. The category airplane is unseen to our DUTS-Cls dataset but seen to ImageNet, while the category giraffe is unseen to both ImageNet and
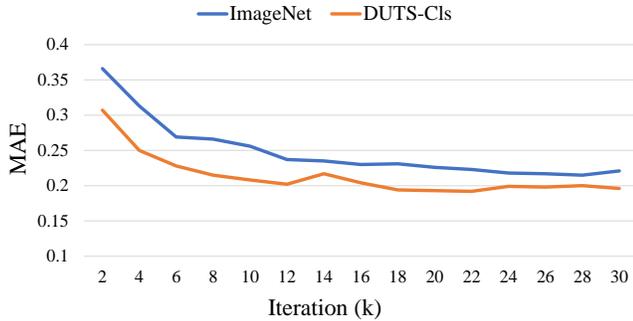
Fig. 10.   Experiments on the effect of our proposed DUTS-Cls dataset. We conduct experiments on the classification branch of the latest work MSW [18] for a fair comparison, the results are tested on the ECSSD dataset.
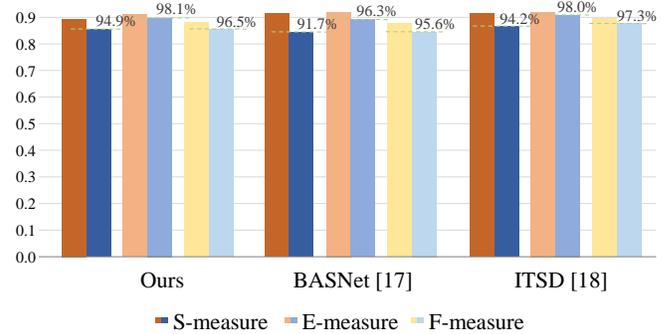


Fig. 11.   Comparisons of different methods trained on our offered labels (the right one) and ground truth (the left one) on ECSSD dataset. The number on each data pair denotes the corresponding percentage.

TABLE IV
THE QUANTITATIVE RESULTS OF EFFECTIVENESS ON UNSEEN CATEGORY. **DATASET** REPRESENTS THE TRAINING SET USED IN THE FIRST TRAINING STAGE, THUR-PLANE AND THUR-GIRAFFE DENOTE THE SAMPLES OF PLANE AND GIRAFFE IN THUR DATASET, RESPECTIVELY.

| Method | Dataset | THUR | | THUR-plane | | THUR-giraffe | |
|--------|---------|------|------|------|------|------|------|
| | | $F_\beta\uparrow$ | $M\downarrow$ | $F_\beta\uparrow$ | $M\downarrow$ | $F_\beta\uparrow$ | $M\downarrow$ |
| MSW [18] | ImageNet | 0.624 | 0.104 | 0.716 | 0.079 | 0.547 | 0.088 |
| Ours | ImageNet | 0.676 | 0.089 | 0.788 | 0.055 | 0.550 | 0.088 |
| | DUTS-Cls | **0.689** | **0.082** | **0.809** | **0.050** | **0.588** | **0.073** |

DUTS-Cls dataset. As is illustrated in the $2^{nd}$ and $3^{rd}$ rows in Table IV, DUTS-Cls dataset encourages better predictions on the whole THUR dataset, and also outperforms ImageNet by a large margin on both airplane and giraffe categories. It proves the generalizability and effectiveness of the proposed DUTS-Cls dataset. Besides, the superiority of our method on unseen categories can be demonstrated in the $1^{st}$ and $2^{nd}$ rows of Table IV. Moreover, except the airplane and giraffe categories, our method also behaves well on other various unseen categories such as the cases shown in Figure 6. It further supports the effect of our method on unseen categories.

*F. Applications*

We extend our method to fully supervised methods by replacing manually labeled ground truth with our generated predictions on training set. To be specific, we infer predictions using our trained model on DUTS-Train dataset and adopt CRF for a further refinement.v It can be seen in Figure 11 that trained with our offered predictions as supervision, BASNet [21] and ITSD [22] achieve 95.6% and 97.3% of their fully supervised accuracy on F-measure without any pixel-level annotations. Additionally, our method also achieves 96.5% accuracy of its fully supervised accuracy on F-measure. These experiments indicate that our method can serve as an alternative to provide pixel-level supervisions for fully supervised SOD methods while maintaining comparatively high accuracy. This costs only 0.32% of pixel-level annotation time and labor.

## VI. CONCLUSION

In this paper, we propose a novel self-calibrated training strategy and introduce a saliency-based dataset with well-matched image-level annotations for WSOD. The proposed strategy establishes a mutual calibration loop between pseudo labels and network predictions, which effectively prevents the network from propagating the negative influence of error-prone pseudo labels. We also argue that cross-domain inconsistency exists between SOD and existing large-scale classification datasets, and impedes the development of WSOD. To offer a better choice for WSOD and encourage more contributions to the community, we introduce a saliency-based classification dataset DUTS-Cls to settle this matter well. Extensive experiments demonstrate the superiority of our method and effectiveness of our two ideas. In addition, our method can serve as an alternative to provide pixel-level labels for fully supervised SOD methods while maintaining comparatively high performance, costing only 0.32% of labeling time for pixel-level annotation.

## REFERENCES

[1] A. Borji, M.-M. Cheng, H. Jiang, and J. Li, "Salient object detection: A benchmark," *IEEE transactions on image processing*, vol. 24, no. 12, pp. 5706–5722, 2015. I

[2] Y. Li, X. Hou, C. Koch, J. M. Rehg, and A. L. Yuille, "The secrets of salient object segmentation," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2014, pp. 280–287. I, V-B

[3] S. Hong, T. You, S. Kwak, and B. Han, "Online tracking by learning discriminative saliency map with convolutional neural network," in *International conference on machine learning*, 2015, pp. 597–606. I

[4] N. Liu and J. Han, "Dhsnet: Deep hierarchical saliency network for salient object detection," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 678–686. I

[5] L. Wang, L. Wang, H. Lu, P. Zhang, and X. Ruan, "Saliency detection with recurrent fully convolutional networks," in *European conference on computer vision*. Springer, 2016, pp. 825–841. I

[6] M. Feng, H. Lu, and E. Ding, "Attentive feedback network for boundary-aware salient object detection," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 1623–1632. I

[7] Z. Deng, X. Hu, L. Zhu, X. Xu, J. Qin, G. Han, and P.-A. Heng, "R3net: Recurrent residual refinement network for saliency detection," in *Proceedings of the 27th International Joint Conference on Artificial Intelligence*, 2018, pp. 684–690. I, V-C

[8] Q. Zhang, R. Cong, C. Li, M.-M. Cheng, Y. Fang, X. Cao, Y. Zhao, and S. Kwong, "Dense attention fluid network for salient object detection in optical remote sensing images," *IEEE Transactions on Image Processing*, vol. 30, pp. 1305–1317, 2021. I

[9] Z. Chen, R. Cong, Q. Xu, and Q. Huang, "DPANet: Depth potentiality-aware gated attention network for RGB-D salient object detection," *IEEE Transactions on Image Processing*, 2020. I

[10] R. Cong, J. Lei, H. Fu, F. Porikli, Q. Huang, and C. Hou, "Video saliency detection via sparsity-based reconstruction and propagation," *IEEE Transactions on Image Processing*, vol. 28, no. 10, pp. 4819–4931, 2019. I

[11] R. Cong, J. Lei, H. Fu, Q. Huang, X. Cao, and C. Hou, "Co-saliency detection for RGBD images based on multi-constraint feature matching and cross label propagation," *IEEE Transactions on Image Processing*, vol. 27, no. 2, pp. 568–579, 2018. I

[12] T. Nguyen, M. Dax, C. K. Mummadi, N. Ngo, T. H. P. Nguyen, Z. Lou, and T. Brox, "Deepusps: Deep robust unsupervised saliency prediction via self-supervision," in *Advances in Neural Information Processing Systems 32*, 2019, pp. 204–214. I

[13] P. Siva, C. Russell, T. Xiang, and L. Agapito, "Looking beyond the image: Unsupervised learning for object saliency and detection," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2013, pp. 3238–3245. I

[14] J. Zhang, X. Yu, A. Li, P. Song, B. Liu, and Y. Dai, "Weakly-supervised salient object detection via scribble annotations," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 12 546–12 555. I

[15] S. yue Yu, B. Zhang, J. Xiao, and E. Lim, "Structure-consistent weakly supervised salient object detection with local saliency coherence," *ArXiv*, vol. abs/2012.04404, 2020. I

[16] L. Wang, H. Lu, Y. Wang, M. Feng, D. Wang, B. Yin, and X. Ruan, "Learning to detect salient objects with image-level supervision," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 136–145. I, II-B, III-A, III-A, IV, V-B, V-C, I

[17] G. Li, Y. Xie, and L. Lin, "Weakly supervised salient object detection using image labels," *arXiv preprint arXiv:1803.06503*, 2018. I, II-B, V-C, I

[18] Y. Zeng, Y. Zhuge, H. Lu, L. Zhang, M. Qian, and Y. Yu, "Multi-source weak supervision for saliency detection," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 6074–6083. I, II-B, III-A, III-A, V-A, V-C, V-C, I, III, V-D, V-D, 10, IV

[19] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, "Imagenet: A large-scale hierarchical image database," in *2009 IEEE conference on computer vision and pattern recognition*. Ieee, 2009, pp. 248–255. I, II-B

[20] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick, "Microsoft coco: Common objects in context," in *European conference on computer vision*. Springer, 2014, pp. 740–755. I, II-B

[21] X. Qin, Z. Zhang, C. Huang, C. Gao, M. Dehghan, and M. Jagersand, "Basnet: Boundary-aware salient object detection," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 7479–7489. I, V-C, V-F

[22] H. Zhou, X. Xie, J.-H. Lai, Z. Chen, and L. Yang, "Interactive two-stream decoder for accurate and fast saliency detection," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 9141–9150. I, V-C, V-F

[23] Z. Jiang and L. S. Davis, "Submodular salient region detection," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2013, pp. 2043–2050. II-A

[24] C. Yang, L. Zhang, H. Lu, X. Ruan, and M.-H. Yang, "Saliency detection via graph-based manifold ranking," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2013, pp. 3166–3173. II-A, V-B, III, V-D

[25] W. Zhu, S. Liang, Y. Wei, and J. Sun, "Saliency optimization from robust background detection," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2014, pp. 2814–2821. II-A

[26] J. Kim, D. Han, Y.-W. Tai, and J. Kim, "Salient region detection via high-dimensional color transform," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2014, pp. 883–890. II-A

[27] Q. Hou, M.-M. Cheng, X. Hu, A. Borji, Z. Tu, and P. H. Torr, "Deeply supervised salient object detection with short connections," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 3203–3212. II-A, II-A, III-C, V-C

[28] Z. Wu, L. Su, and Q. Huang, "Cascaded partial decoder for fast and accurate salient object detection," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 3907–3916. II-A, II-A, III-C, V-C

[29] J. Su, J. Li, Y. Zhang, C. Xia, and Y. Tian, "Selectivity or invariance: Boundary-aware salient object detection," in *2019 IEEE/CVF Interna-tional Conference on Computer Vision, ICCV 2019*. IEEE, 2019, pp. 3798–3807. II-A, II-A

[30] J. Zhao, J. Liu, D. Fan, Y. Cao, J. Yang, and M. Cheng, "Egnet: Edge guidance network for salient object detection," in *2019 IEEE/CVF International Conference on Computer Vision, ICCV 2019*. IEEE, 2019, pp. 8778–8787. II-A, II-A

[31] J. Wei, S. Wang, and Q. Huang, "F$^3$net: Fusion, feedback and focus for salient object detection," in *The Thirty-Fourth AAAI Conference on Artificial Intelligence, AAAI 2020*. AAAI Press, 2020, pp. 12 321–12 328. II-A

[32] J. Zhang, S. Sclaroff, Z. Lin, X. Shen, B. Price, and R. Mech, "Minimum barrier salient object detection at 80 fps," in *Proceedings of the IEEE international conference on computer vision*, 2015, pp. 1404–1412. II-B

[33] B. Zhou, A. Khosla, A. Lapedriza, A. Oliva, and A. Torralba, "Learning deep features for discriminative localization," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 2921–2929. II-B, III-A

[34] J. Ahn, S. Cho, and S. Kwak, "Weakly supervised learning of instance segmentation with inter-pixel relations," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 2209–2218. III-A, III-A

[35] R. Achanta, A. Shaji, K. Smith, A. Lucchi, P. Fua, and S. Süsstrunk, "Slic superpixels compared to state-of-the-art superpixel methods," *IEEE transactions on pattern analysis and machine intelligence*, vol. 34, no. 11, pp. 2274–2282, 2012. III-A

[36] N. Araslanov and S. Roth, "Single-stage semantic segmentation from image labels," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 4253–4262. III-A

[37] P. Krähenbühl and V. Koltun, "Efficient inference in fully connected crfs with gaussian edge potentials," in *Advances in neural information processing systems*, 2011, pp. 109–117. III-A

[38] J. Fan, Z. Zhang, and T. Tan, "Employing multi-estimations for weakly-supervised semantic segmentation," in *2020 IEEE/CVF European Conference on Computer Vision, ECCV 2020*, vol. 12362, 2020, pp. 332–348. III-B

[39] Y. Niu, Y. Geng, X. Li, and F. Liu, "Leveraging stereopsis for saliency analysis," in *2012 IEEE Conference on Computer Vision and Pattern Recognition*, 2012, pp. 454–461. IV

[40] G. Huang, Z. Liu, and K. Q. Weinberger, "Densely connected convolutional networks," *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 2261–2269, 2017. V-A

[41] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," *arXiv preprint arXiv:1412.6980*, 2014. V-A

[42] Q. Yan, L. Xu, J. Shi, and J. Jia, "Hierarchical saliency detection," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2013, pp. 1155–1162. V-A, V-B

[43] R. Achanta, S. Hemami, F. Estrada, and S. Susstrunk, "Frequency-tuned salient region detection," in *2009 IEEE conference on computer vision and pattern recognition*, 2009, pp. 1597–1604. V-A, V-B

[44] M. Everingham, L. Van Gool, C. K. Williams, J. Winn, and A. Zisserman, "The pascal visual object classes (voc) challenge," *International journal of computer vision*, vol. 88, no. 2, pp. 303–338, 2010. V-B

[45] G. Li and Y. Yu, "Visual saliency based on multiscale deep features," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2015, pp. 5455–5463. V-B

[46] D.-P. Fan, M.-M. Cheng, Y. Liu, T. Li, and A. Borji, "Structure-measure: A new way to evaluate foreground maps," in *Proceedings of the IEEE international conference on computer vision*, 2017, pp. 4548–4557. V-B

[47] D.-P. Fan, C. Gong, Y. Cao, B. Ren, M.-M. Cheng, and A. Borji, "Enhanced-alignment measure for binary foreground map evaluation," *arXiv preprint arXiv:1805.10421*, 2018. V-B

[48] T. Wang, L. Zhang, S. Wang, H. Lu, G. Yang, X. Ruan, and A. Borji, "Detect globally, refine locally: A novel approach to saliency detection," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 3127–3135. V-C

[49] T. Zhao and X. Wu, "Pyramid feature attention network for saliency detection," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 3085–3094. V-C

[50] Z. Wu, L. Su, and Q. Huang, "Stacked cross refinement network for edge-aware salient object detection," in *Proceedings of the IEEE International Conference on Computer Vision*, 2019, pp. 7264–7273. V-C

[51] Y. Pang, X. Zhao, L. Zhang, and H. Lu, "Multi-scale interactive network for salient object detection," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 9413–9422. V-C

[52] Y. Qin, H. Lu, Y. Xu, and H. Wang, "Saliency detection via cellular automata," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2015, pp. 110–119. III, V-D

[53] M.-M. Cheng, N. J. Mitra, X. Huang, and S.-M. Hu, "Salientshape: group saliency in image collections," *The visual computer*, vol. 30, no. 4, pp. 443–453, 2014. V-E