

---

# [Re] Interactive Two-Stream Decoder for Accurate and Fast Saliency Detection

---

Anonymous Author(s)

Affiliation

Address

email

## Reproducibility Summary

1

2 *This article describes our attempts to reproduce the work presented in the paper - Interactive Two-Stream Decoder for*  
3 *Accurate and Fast Saliency Detection which published in CVPR2020, as part of the ML Reproducibility Challenge*  
4 *2020. Our work consists of three parts: (1) evaluation using the provided pre-trained models from the original paper;*  
5 *(2) evaluation using our retrained models; (3) parameters analysis.*

### 6 **Scope of Reproducibility**

7 There are two main objectives of our report: (1) evaluation using the provided *pre-trained* models from the original  
8 paper; (2) evaluation using our *retrained* models

### 9 **Methodology**

10 We used publicly available source code provided by the authors. Minor changes were made to the source code in order  
11 to load the model weight properly. The reproducibility experiments followed the protocol as described in the original  
12 paper. We performed the training on a machine with two GTX2080TI GPUs. The training and validation took about 5  
13 and 7.5 hours for ResNet-based and VGG-based models, respectively.

### 14 **Results**

15 The reproducibility using authors' pretrained model was of no success and we recorded the discrepancy of as high as  
16 11.4%. While for the retrained models, the difference was observed to be as high as 8.6%. Therefore, we carried out  
17 statistical test to further confirm the reproducibility of the investigated work. We failed to reject the null hypothesis,  
18 which implies that there is no significant difference between the reported results and our results.

### 19 **What was easy**

20 The pretrained models have been shared by authors. One can try out the models with little effort as minor changes to  
21 the name of the parameter keys of the provided PyTorch models are expected.

### 22 **What was difficult**

23 We observed some differences between the reported results. We spent longer time to look into the issue and also had  
24 discussion with the authors. We repeated the experiments several times for each model to ensure the correctness of  
25 our results. Despite all the challenges, we finally come to the conclusion that the investigated work is reproducible, as  
26 evidenced by the statistical test.

### 27 **Communication with original authors**

28 We raised and discussed some issues via GitHub. We also discussed some technical details and also about the result  
29 discrepancy with the authors via emails.

## 30 **1 Introduction**

31 Saliency detection is a mechanism to detecting visually attentive objects or regions in an image. To better discover  
32 the visually distinctive objects or regions, the saliency detection methods are expected to understand both the global  
33 and local features Qin et al. [2019]. The global features are essential to provide high-level semantic information for  
34 identifying salient objects while the local features are exploited for features refinement. The investigated paper –  
35 Interactive Two-Stream Decoder for Accurate and Fast Saliency Detection (ITSD) Zhou et al. [2020b] explores the  
36 correlation between saliency (global) and contour (local) information by introducing interactive connections between  
37 the saliency and contour streams in the decoder. An adaptive contour loss is also introduced in the saliency stream to  
38 weight more on near boundary pixels of the salient objects, thus improves the network capability to handling the salient  
39 object boundary. The method was evaluated on six publicly available saliency datasets and encouraging results were  
40 reported.

41

## 42 **2 Scope of reproducibility**

43 The investigated work is considered lightweight and able to achieve real-time inference speed at about 88 FPS based on  
44 our evaluation. Furthermore, the reported results on six benchmark datasets are also quite encouraging, which achieve  
45 comparable performance as compared to several SOTA methods. Our reproduction of work consists of the following:

- 46 • Reproduce the reported results using the pretrained models as made available by the authors
- 47 • Reproduce the reported results by training the models following the exact protocol described in the original  
48 paper
- 49 • Evaluation of various parameters as reported in the original paper

## 50 **3 Methodology**

51 The authors have made the source code associated with the paper publicly available on GitHub as well as links to  
52 download their pretrained VGG-based and ResNet-based models Zhou et al. [2020a]. In order to properly load the  
53 model weights, a minor modification is expected to make to the name of the parameter keys in the given PyTorch  
54 models as the parameter names are inconsistent with the model definition. To retrain the ITSD models, we followed the  
55 training protocol as described in the original paper. The VGG and ResNet, which pretrained on ImageNet dataset, were  
56 employed as the backbone of the model. Data augmentation such as random flipping, random cropping and scaling was  
57 applied on the DUTS-TR dataset with the aim to increase variability in the dataset. The models were trained on two  
58 GTX2080TI GPUs using SGD at the learning rate of 0.01 for the first 20k iterations and decayed by a factor of 0.1 for  
59 the remaining 5k iterations.

### 60 **3.1 Model descriptions**

61 The ITSD employs an encoder-decoder network architecture. The backbone of the encoder can be either of a pretrained  
62 VGG or ResNet. In order to keep the model lightweight, two measures are adopted: (1) the fully-connected (FC) layers  
63 in the standard VGG network is truncated with the aim to reduce the model size while there is no such need for the  
64 ResNet-based encoder since it has no FC layer; (2) channel pooling is applied to each of the encoded feature maps. To  
65 better recover the saliency map, the decoder of the network employs a two-stream structure, which comprises of five  
66 cascaded Feature Correlation Module (FCF) modules. The FCF provides interactive connection (fusion stream) which  
67 allows the respective saliency and contour information to be transmitted across the stream, and therefore encourages the  
68 network to learn the correlation. To generate the final saliency map, all of the intermediate features from the saliency  
69 stream are first up-sampled to the input image size followed by concatenation operation. An adaptive contour loss  
70 (ACT) which leverages on the contour information is proposed to provide adaptive weight to the saliency stream, which  
71 allows the model to pay more attention to those pixels near to the salient object boundary.

72

### 73 3.2 Datasets

74 Following the original paper, six commonly used publicly available benchmark datasets, namely, DUTS, SOD, PASCAL-  
75 S, ECSSD, HKU-IS and DUTS-OMRON were employed in our reproduction work. The number of images for each  
76 dataset is summarized in Table 1.

Datasets	Number of Images	Remarks
DUTS-TR	10553	Training
DUTS-TE	5019	Testing
SOD	300	
PASCAL-S	850	
ECSSD	1000	
HKU-IS	4447	
DUTS-OMRON	5168	

Table 1: Training and testing datasets

### 77 3.3 Hyperparameters

78 As our primary object was to evaluate the reproducibility of the investigated work, we employed the default parameters  
79 reported in the original paper, *i.e.*, input image size =  $288 \times 288$ ,  $\lambda = 1$ ,  $m = 4$ , and batch size = 8, to train the ITSD  
80 models. In additional, we also carried out further analysis on various parameters as reported in the paper.

### 81 3.4 Experimental setup

82 All the experiments conducted in this work followed the exact protocol as described in the original paper, which we  
83 hoped to mitigate the possible differences in some software libraries/packages, in order to reproduce the results as  
84 similarly to the reported values. We used the recommended settings from the original paper, as also described in Secs. 3  
85 & 3.3 for the training. Furthermore, we used the source code as made publicly available by the authors to perform both  
86 the training and the evaluation, which ensure there exists no bias in the implementation. All the training and evaluation  
87 were performed on the same machine with two Nvidia GTX2080TI GPUs, except for the inference speed evaluation  
88 which was done using single GPU. Our source code is also made publicly available at GitHub.

### 89 3.5 Computational requirements

90 The ResNet-based and VGG-based ITSD models had different requirements on the computational expectations.  
91 According to our evaluations, the times required to perform the training and validation were around 5 and 7.5 hours,  
92 respectively for ResNet-based and VGG-based models. The overall required time to perform training and validation  
93 was observed to be shorter for the ResNet-based model, which may be attributed to its efficient bottleneck architecture.

## 94 4 Results

95 In this section, we provide our evaluation results obtained from: (1) the pretrained models as shared by the authors; (2)  
96 the retrained ITSD models according to the described protocol and (3) evaluation of various parameters. We presented  
97 the evaluation results using two commonly used performance metrics, *i.e.*,  $F_\beta$  and  $MAE$ .

### 98 4.1 Performance evaluation using pretrained models

99 As can be observed from Table 2 and 3, the evaluation results in terms of  $F_\beta$  were close to the reported values, with the  
100 neglectable differences of up to 1%. However, the degree of differences in terms of  $MAE$  was quite noticeable, with  
101 the difference as high as 11.4% was observed.

### 102 4.2 Performance evaluation using retrained models

103 Table 4 and 5 summarize our evaluation results using the retrained ITSD models. It can be observed that the degree of  
104 differences can be as high as 2.4% in  $F_\beta$  and 8.6% in MAE. In order to quantitatively justifying the reproducibility of

Datasets	$F_\beta$			MAE		
	Our results	Reported	Difference (%)	Our results	Reported	Difference (%)
ECSSD	0.946	0.947	0.11	0.034	0.035	2.86
DUTS-TE	0.885	0.883	0.23	0.039	0.041	4.88
DUT-OMRON	0.819	0.824	0.61	0.059	0.061	3.28
PASCAL-S	0.876	0.871	0.57	0.064	0.071	9.86
SOD	0.872	0.880	0.91	0.092	0.095	3.16
HKU-IS	0.936	0.934	0.21	0.029	0.031	6.45

Table 2: Evaluation results using the pretrained ResNet-based model

Datasets	$F_\beta$			MAE		
	Our results	Reported	Difference (%)	Our results	Reported	Difference (%)
ECSSD	0.939	0.939	0.00	0.039	0.040	2.50
DUTS-TE	0.880	0.877	0.34	0.041	0.042	2.38
DUT-OMRON	0.805	0.813	0.98	0.059	0.063	6.35
PASCAL-S	0.871	0.871	0.00	0.067	0.074	9.46
SOD	0.860	0.869	1.04	0.106	0.100	6.00
HKU-IS	0.931	0.927	0.43	0.031	0.035	11.43

Table 3: Evaluation results using the pretrained VGG-based model

105 the results, we carried out statistical analysis using paired  $t$ -test at significance level of 0.05 to verify our null hypothesis  
106 if there was no significant difference between our evaluation results and the reported results in the article. The paired  
107  $t$ -test statistics value can be calculated as follows:

$$t = \frac{\bar{d}}{s/\sqrt{n_d}} \quad (1)$$

108 where  $\bar{d}$  is the mean differences,  $s$  is the standard deviation and  $n_d$  is the number of observations, *i.e.*, number of the  
109 employed benchmark datasets. The  $p$ -value obtained from the calculated  $t$ -score is given in the Table 6.

Datasets	$F_\beta$			MAE		
	Our results (best epoch result)	Reported	Difference (%)	Our results (best epoch result)	Reported	Difference (%)
ECSSD	0.945 (0.947)	0.947	0.21	0.035 (0.034)	0.035	0.00
DUTS-TE	0.883 (0.883)	0.883	0.00	0.040 (0.039)	0.041	2.44
DUT-OMRON	0.823 (0.824)	0.824	0.12	0.057 (0.055)	0.061	6.56
PASCAL-S	0.870 (0.872)	0.871	0.11	0.068 (0.066)	0.071	4.23
SOD	0.859 (0.875)	0.880	2.39	0.099 (0.091)	0.095	4.21
HKU-IS	0.935 (0.935)	0.934	0.11	0.029 (0.029)	0.031	6.45

Table 4: Evaluation results using the retrained ResNet-based model

Datasets	$F_\beta$			MAE		
	Our results (best epoch result)	Reported	Difference (%)	Our results (best epoch result)	Reported	Difference (%)
ECSSD	0.937 (0.941)	0.939	0.21	0.039 (0.037)	0.040	2.50
DUTS-TE	0.874 (0.879)	0.877	0.34	0.042 (0.041)	0.042	0.00
DUT-OMRON	0.801 (0.804)	0.813	1.48	0.062 (0.059)	0.063	1.59
PASCAL-S	0.869 (0.875)	0.871	0.23	0.068 (0.066)	0.074	8.11
SOD	0.863 (0.866)	0.869	0.69	0.103 (0.093)	0.100	3.00
HKU-IS	0.929 (0.932)	0.927	0.22	0.032 (0.030)	0.035	8.57

Table 5: Evaluation results using the retrained VGG-based model

	$p$ -value	
	ResNet-based	VGG-based
$F_\beta$	0.2956	0.1050
MAE	0.4261	0.3276

Table 6: Statistical test at significance level of 0.05

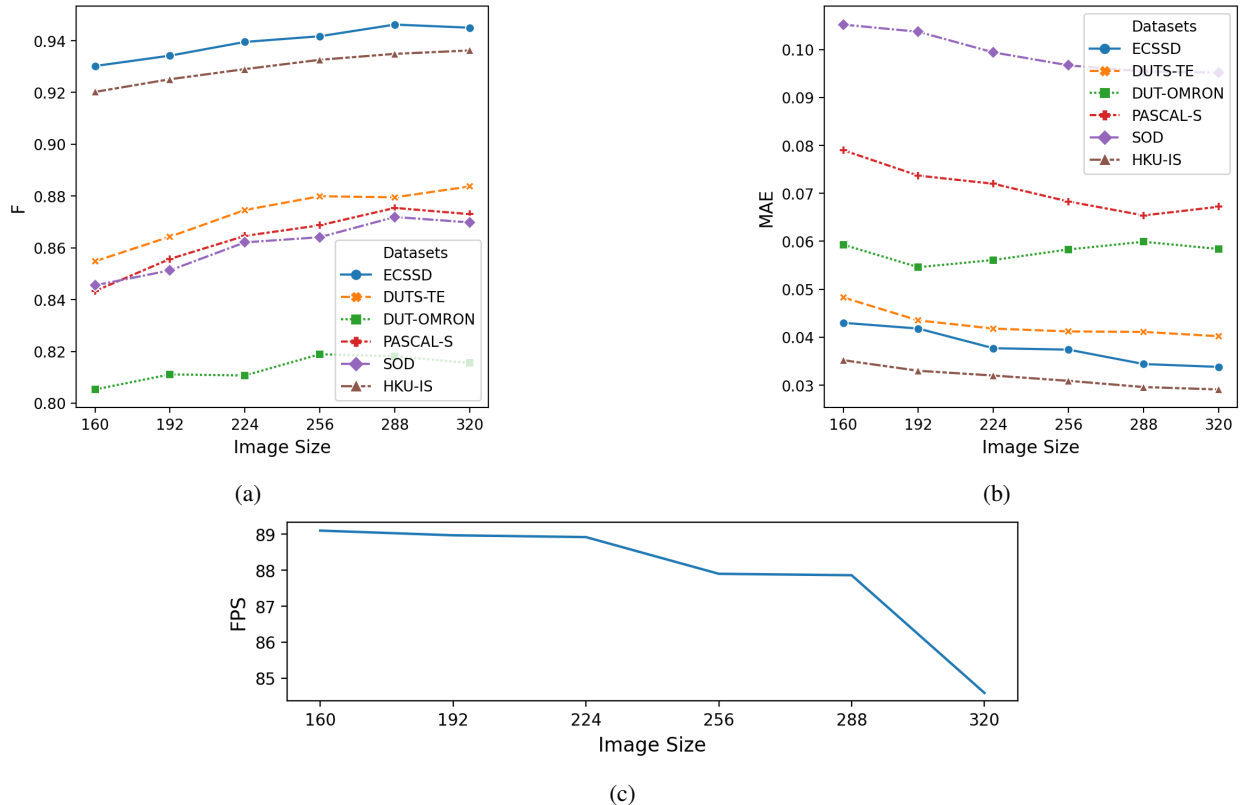


Figure 1: Performance of different input image sizes and inference speed. The recommended image size  $288 \times 288$  appear to be a better choice which balances the saliency and the inference performance. (a)  $F_\beta$ . (b) MAE (c) inference speed.

### 110 4.3 Additional results not present in the original paper

111 In this section, we provide further analysis on various parameters that had not been elaborated in the original paper to  
 112 provide more insights which may be deemed benefit for those working in the field. All the experiments conducted in  
 113 this section employed the ResNet-based ITSD model, with all the training parameters kept unchanged (see Section 3.3)  
 114 unless specified otherwise.

#### 115 4.3.1 Image sizes and inference speed

116 We have seen different input image sizes were employed for the saliency detection, such as Liu et al. [2019], Zhao et al.  
 117 [2019], Qin et al. [2019]. Ref. Zhou et al. [2020b] arbitrary chosen an input image size of  $288 \times 288$  which motivates  
 118 us to determine the impact of varying input image sizes on the performance of the ITSD model. We performed the  
 119 evaluation using the models which trained with input image sizes of  $160 \times 160$ ,  $192 \times 192$ ,  $224 \times 224$ ,  $256 \times 256$ ,  $288 \times$   
 120  $288$ ,  $320 \times 320$ . Generally, the performance was better by using larger input image size, as can be observed from the  
 121 evaluation results as shown in Fig. 1(a) and (b). In addition, we also provide the evaluation of inference speed with  
 122 different input image sizes on a single GPU, as can be observed from Fig. 1(c). The default input image size  $288 \times 288$   
 123 reported the inference speed ran at around 88 FPS while the inference speed for the input image size of  $160 \times 160$   
 124 achieved 89 FPS, at the expense of the saliency performance. The input image size of  $288 \times 288$  emerges to be a better  
 125 choice which balances the saliency performance and the inference speed.

#### 126 4.3.2 Hard examples factor $m$

127 The hard examples factor  $m$  in the saliency loss function serves to emphasize the boundary pixels so the model can  
 128 better handling the boundary of the salient regions. However, the effect of such parameter had not been discussed in the

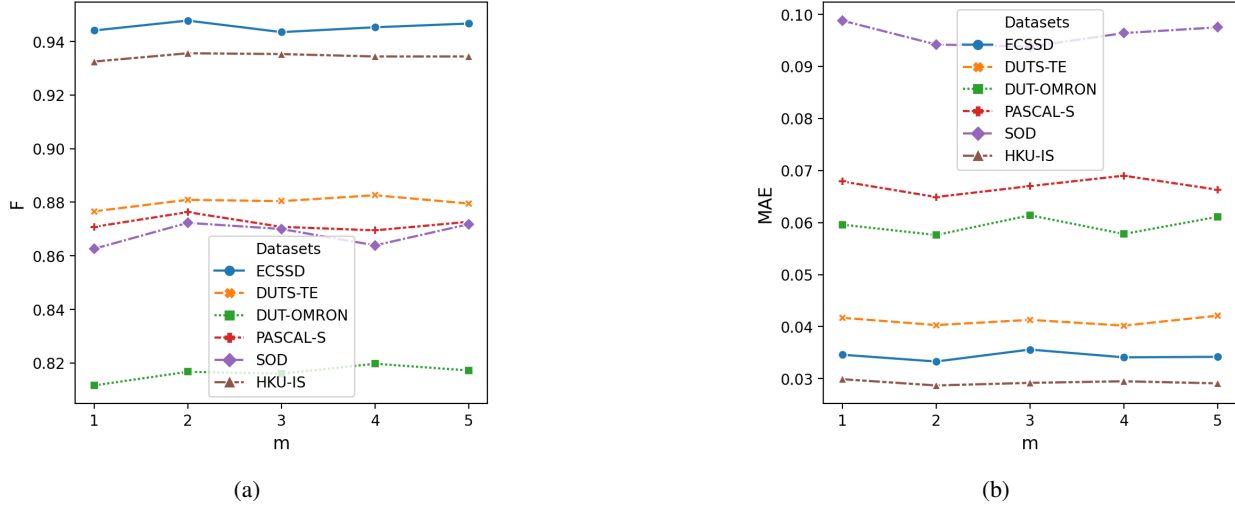


Figure 2: Effect of different hard example factor (a)  $F_\beta$ . (b) MAE

129 original paper. Fig. 2 shows our evaluation results of using various models respectively trained with  $m = \{1, 2, 3, 4, 5\}$ .  
 130 On average, some marginal gains were observed in both metrics with the hard examples factor  $m = 2$ .

### 131 4.3.3 Dilation and erosion

132 The contour map is obtained by computing the difference between the saliency maps produced by two morphological  
 133 operations, *i.e.*, dilation and erosion. The resultant from such operation is a contour map with the border width  
 134 determined by the morphological kernel size, which in turn affecting how the model learn about the hard example pixels  
 135 near the salient object boundary. In the original authors' implementation, morphological kernel size of  $5 \times 5$  was used.  
 136 In order to study the effect of different morphological kernel sizes, we considered  $1 \times 1, 3 \times 3, 5 \times 5, 7 \times 7, 9 \times 9$  in  
 137 our experiments. Fig 3 shows samples of contour maps generated using different morphological kernel sizes. Note that  
 138 no contour information is produced by using the morphological kernel size of  $1 \times 1$  and therefore the Eq. 12 (refer to  
 139 original paper) will be solely depend on the  $P^c$  to weight the saliency loss function. As can be observed from Fig. 4,  
 140 the morphological kernel size of  $3 \times 3$  can generally provide better performance than the default kernel size, but such  
 141 improvement are only marginal.

### 142 4.3.4 Lambda, $\lambda$

143 The  $\lambda$  in the total loss function serves to balance between the saliency and contour losses, which was set to  $\lambda = 1$  in the  
 144 original paper. In other words, one can determine how much emphasis to be placed on the contour loss by adjusting the  
 145  $\lambda$ . In our experiments, we performed evaluation for a range of values, *i.e.*,  $\lambda = [0.0, 1.6]$ , with a step size of 0.2. Based  
 146 on the average results from our evaluation, the  $\lambda = 0.6$  had shown to achieve better performance, which suggested  
 147 that the model is biased towards the saliency information as such information is more crucial, while complemented by  
 148 the contour information for further performance gain. The  $\lambda = 0$  totally suppresses the contour loss, which can be  
 149 considered as a special case of the total loss function, where only saliency information is considered.

### 150 4.3.5 Optimized vs. default parameters

	Optimized	Default
Input image size	$288 \times 288$	$288 \times 288$
Morphological kernel size	$3 \times 3$	$5 \times 5$
Hard examples factor, $m$	2	4
Lambda, $\lambda$	0.6	1

Table 7: Optimized vs. default parameters

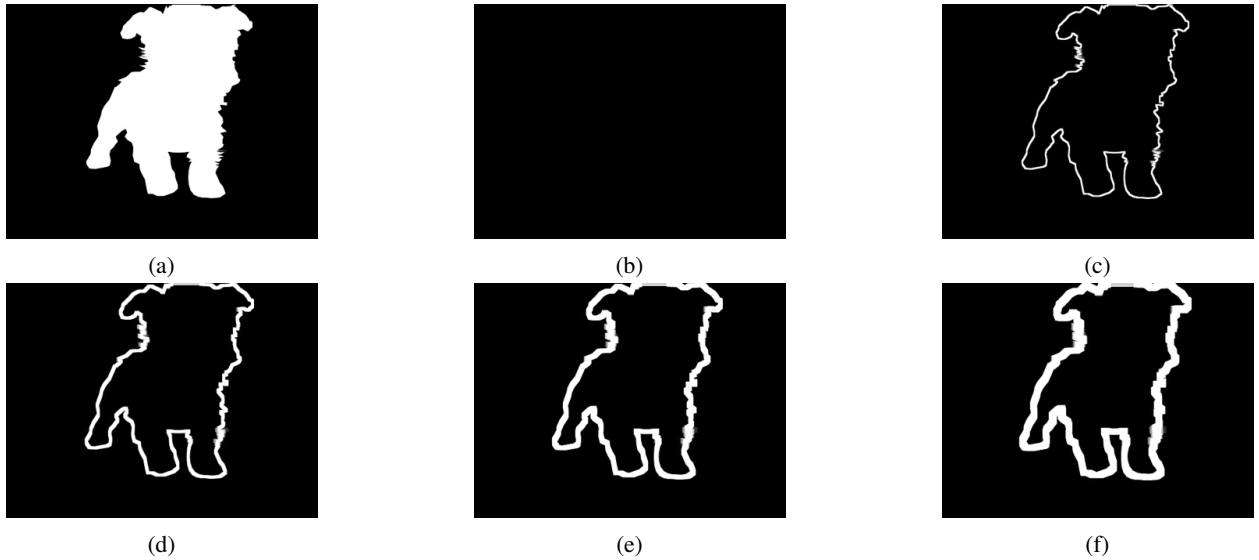
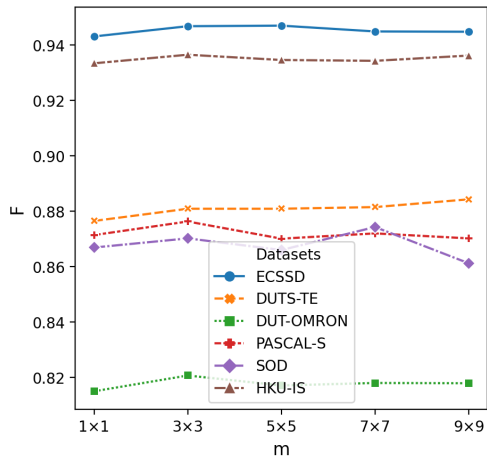
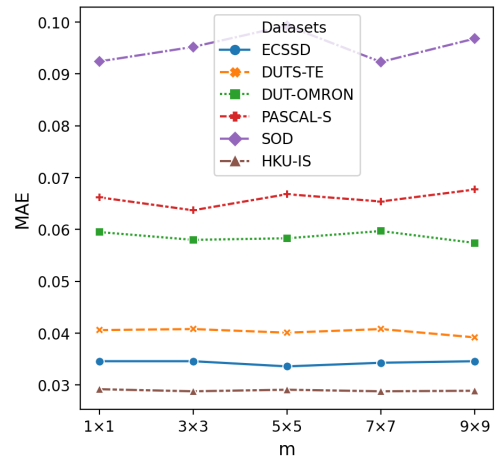


Figure 3: Contour maps generated using different morphological kernels (a) mask, (b)  $1 \times 1$ , (c)  $3 \times 3$ , (d)  $5 \times 5$ , (e)  $7 \times 7$ , (f)  $9 \times 9$

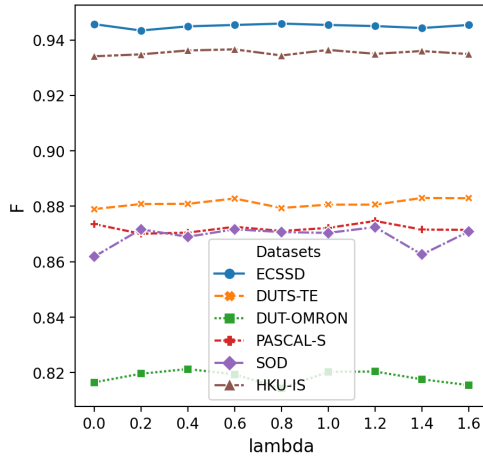


(a)

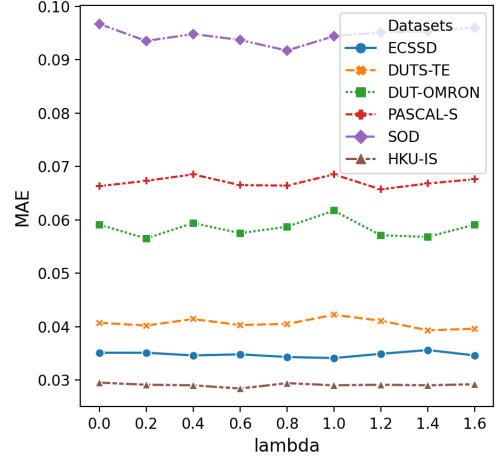


(b)

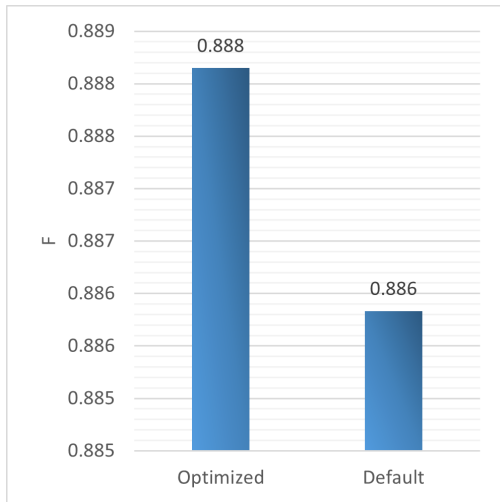
Figure 4: Effect of morphological kernel size (a)  $F_\beta$ . (b) MAE



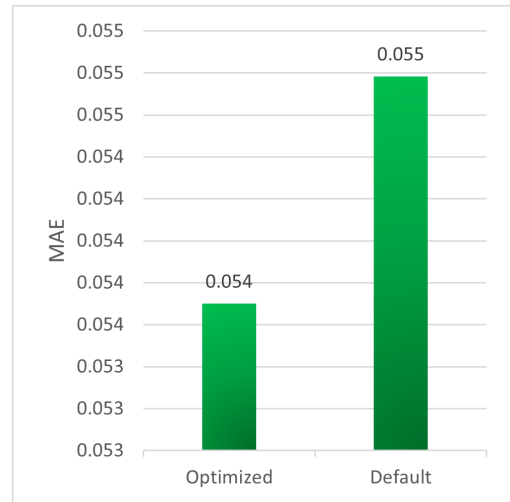
(a)



(b)

Figure 5: Effect of lambda (a)  $F_\beta$ . (b) MAE

(a)



(b)

Figure 6: Marginal performance improvement is observed with our optimized parameters as compared to the default parameters. (a)  $F_\beta$ . (b) MAE

151 Sections 4.3.1 - 4.3.4 show our evaluations on the various parameters of the ITSD model, and each individual parameter  
 152 has been shown to induce certain degree of improvement on the average performance metrics. As such, we carried  
 153 out experiment by using the optimized parameters, as summarized in Table 7. Fig. 6 presents the average results of  
 154 the comparison between the models trained with optimized and the default parameters, which only marginal gain was  
 155 observed. Our primary intention to carry out the experiments in Sections 4.3.1 - 4.3.4 was to study the implication of  
 156 each parameter towards the saliency detection performance. Therefore, each of the experiment conducted only evaluated  
 157 the parameter of concerned while other parameters were kept unchanged in order to better facilitate the analysis. One  
 158 should consider a proper parameter optimization approach such as grid search to exhaustively searches through the  
 159 parameter spaces to better estimate the best parameters' combination.



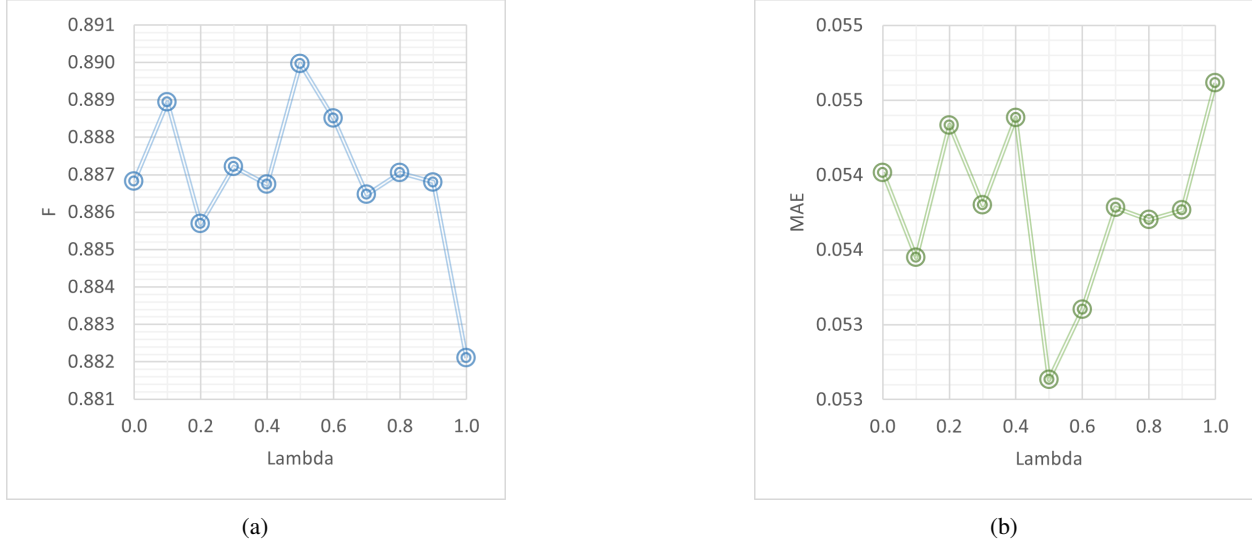


Figure 7: Constrained lambda. (a)  $F_\beta$ . (b) MAE

#### 160 4.4 Constrained lambda

161 Although the free parameter, *i.e.*,  $\lambda$  in the total loss function was analyzed in Section 4.3.4, it is still uncertain which of  
 162 the losses contribute more towards the optimization of the network. In order to further analyze the importance of each  
 163 loss term, we imposed constraint to the total loss function by considering a weighted sum rule, as given as follows:

$$L(P^s, P^c, G^s, G^c) = \sum_{i=0}^5 (1 - \lambda) L^s(P_i^s, P_i^c, G_i^s, G_i^c) + \lambda \sum_{j=1}^5 L^c(P_j^c, G_j^c) \quad (2)$$

s.t.  $0 \leq \lambda \leq 1$

164 By constraining the lambda would enable us to analyze the correlation between the saliency and contour losses, as can  
 165 be observed from average performance metrics in Fig. 7. The best saliency detection performance was achieved when  
 166 the  $\lambda = 0.5$ , which suggested that both the saliency and contour losses contributed equally and complemented each  
 167 other. While the worst performance was reported when the  $\lambda = 1.0$ , which implied that the contour loss alone appeared  
 168 to be suboptimal for optimizing the model.

#### 169 5 Discussion and conclusion

170 We had some issues to get the published code working initially but those issues were solved after some discussion with  
 171 the authors. One of the major issues was the inconsistent use of the name of the parameter keys in the given PyTorch  
 172 models, which resulted in failing to load the model weight. The issue can be solved by updating the name of the  
 173 parameter keys to be compatible with the defined models.

174  
 175 Despite of all the challenges we faced during the work reproduction, we eventually managed to reproduce the work of  
 176 the investigated paper, *i.e.*, ITSD. Our reproduction attempts were mainly comprised of two parts: (1) reproduction  
 177 using the pretrained models as shared by the authors; (2) reproduction by training the models. Our attempt using the  
 178 pretrained models was of no success to reproduce the exact results as reported in the original paper, as can be observed  
 179 in Table 2 and 3. After seeking clarification from the authors, one of the possible causes of such discrepancy may be  
 180 due to the environment settings. Another possible reason may be due to the best epoch results were reported in the  
 181 paper, while the shared models represented the final training artifacts. Our assumption was evidenced by carrying out  
 182 the experiments to capture the best epoch results, which we successfully obtained either similar or better performance  
 183 as compared to the reported results. Our second reproduction attempt was to perform the retraining of the ITSD models.  
 184 Despite differences in the experimental results, we carried out statistical test, *i.e.*, paired *t*-test, to further confirmed the

185 reproducibility of the presented work in the studied paper. Further to our reproduction attempts, we had also carried  
186 out experiments to perform evaluation on the various parameters which had not been elaborated in the original paper.  
187 We hope that our reproduction work can be beneficial to those working in the same direction, as well as serving as a  
188 recognition to the contributions of the original authors.

## 189 **References**

- 190 Jiang-Jiang Liu, Qibin Hou, Ming-Ming Cheng, Jiashi Feng, and Jianmin Jiang. A simple pooling-based design for  
191 real-time salient object detection. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2019.
- 192 X. Qin, Z. Zhang, C. Huang, C. Gao, M. Dehghan, and M. Jagersand. Basnet: Boundary-aware salient object detection.  
193 In *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 7471–7481, 2019. doi:  
194 10.1109/CVPR.2019.00766.
- 195 J. Zhao, J. Liu, D. Fan, Y. Cao, J. Yang, and M. Cheng. Egnets: Edge guidance network for salient object detection. In  
196 *2019 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 8778–8787, 2019. doi: 10.1109/ICCV.  
197 2019.00887.
- 198 H. Zhou, X. Xie, J. H. Lai, Z. Chen, and L. Yang. Itsd-pytorch, 2020a. URL [https://github.com/moothes/  
199 ITSD-pytorch](https://github.com/moothes/ITSD-pytorch).
- 200 H. Zhou, X. Xie, J. H. Lai, Z. Chen, and L. Yang. Interactive two-stream decoder for accurate and fast saliency detection.  
201 In *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 9138–9147, 2020b. doi:  
202 10.1109/CVPR42600.2020.00916.