



# Troubleshooting image segmentation models with human-in-the-loop

Haotao Wang<sup>1</sup> · Tianlong Chen<sup>1</sup> · Zhangyang Wang<sup>1</sup> · Kede Ma<sup>2</sup>

Received: 17 July 2021 / Revised: 8 October 2021 / Accepted: 21 October 2021 /  
Published online: 27 January 2022

© The Author(s), under exclusive licence to Springer Science+Business Media LLC, part of Springer Nature 2021

## Abstract

Image segmentation lays the foundation for many high-stakes vision applications such as autonomous driving and medical image analysis. It is, therefore, of great importance to not only improve the accuracy of segmentation models on well-established benchmarks, but also enhance their robustness in the real world so as to avoid sparse but fatal failures. In this paper, instead of chasing state-of-the-art performance on existing benchmarks, we turn our attention to a new challenging problem: how to efficiently expose failures of “top-performing” segmentation models in the real world and how to leverage such counterexamples to rectify the models. To achieve this with minimal human labelling effort, we first automatically sample a small set of images that are likely to falsify the target model from a large corpus of web images via the maximum discrepancy competition principle. We then propose a weakly labelling strategy to further reduce the number of false positives, before time-consuming pixel-level labelling by humans. Finally, we fine-tune the model to harness the identified failures, and repeat the whole process, resulting in an efficient and progressive framework for troubleshooting segmentation models. We demonstrate the feasibility of our framework using the semantic segmentation task in PASCAL VOC, and find that the fine-tuned model exhibits significantly improved generalization when applied to real-world images with greater content diversity. The code is available at [https://github.com/VITA-Group/Troubleshooting\\_Image\\_Segmentation](https://github.com/VITA-Group/Troubleshooting_Image_Segmentation).

---

Editors: Bo Han, Tongliang Liu, Quanming Yao, Mingming Gong, Gang Niu, Ivor W. Tsang, Masashi Sugiyama.

---

✉ Haotao Wang  
[htwang@utexas.edu](mailto:htwang@utexas.edu)

Tianlong Chen  
[tianlong.chen@utexas.edu](mailto:tianlong.chen@utexas.edu)

Zhangyang Wang  
[atlaswang@utexas.edu](mailto:atlaswang@utexas.edu)

Kede Ma  
[kede.ma@cityu.edu.hk](mailto:kede.ma@cityu.edu.hk)

<sup>1</sup> University of Texas at Austin, Austin, United States

<sup>2</sup> City University of Hong Kong, Hong Kong, China

**Keywords** Model troubleshooting · Semantic segmentation

## 1 Introduction

Image segmentation (i.e., pixel-level image labelling) has recently risen to explosive popularity, due in part to its profound impact on many high-stakes vision applications, such as autonomous driving and medical image analysis. While the performance of segmentation models, as measured by excessively reused test sets (Everingham et al., 2010; Lin et al., 2014), keeps improving (Badrinarayanan et al., 2017; Chen et al., 2018; Yu et al., 2018), two scientific questions have arisen to capture the community’s curiosity, and motivate the current work:

- Q1 Do “top-performing” segmentation models on existing benchmarks generalize to the real world with much richer variations?
- Q2 Can we identify and rectify the trained models’ sparse but fatal mistakes, without incurring significant workload of human labelling?

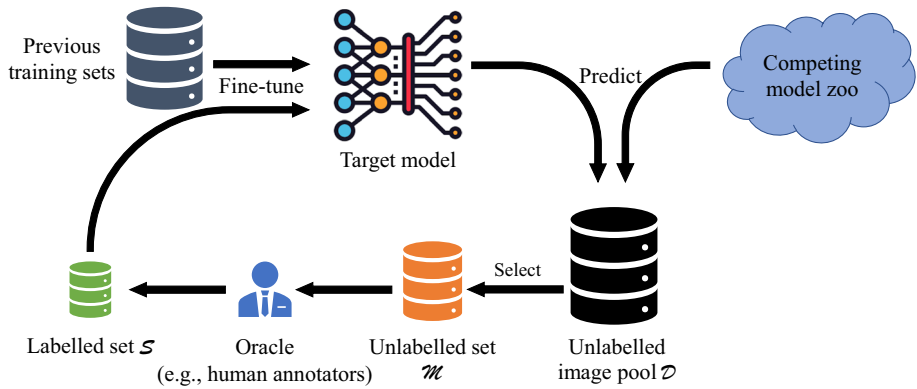
*The answer to the first question* is conceptually clearer, by taking reference to a series of recent work on image classification (Hendrycks et al., 2019; Recht et al., 2019). A typical test set for image classification can only include a maximum of ten thousands of images because human labelling (or verification of predicted labels) is expensive and time-consuming. Considering the high dimensionality of image space and the “human-level” performance of existing methods, such test sets may only spot an extremely small subset of possible mistakes that the model will make, suggesting their insufficiency to cover hard examples that may be encountered in the real world (Wang et al., 2020). The existence of natural adversarial examples (Hendrycks et al., 2019) also echos such hidden fragility of the classifiers to unseen examples, despite the impressive accuracy on existing benchmarks.

While the above problem has not been studied in the context of image segmentation, we argue that it would only be much *amplified* for two main reasons. First, segmentation benchmarks require pixel-level dense annotation. Compared to classification databases, they are much more expensive, laborious, and error-prone to label,<sup>1</sup> making existing segmentation datasets even more restricted in scale. Second, it is much harder for segmentation data to be class-balanced in the pixel level, making highly skewed class distributions notoriously common for this particular task (Bischke et al., 2018; Kervadec et al., 2019). Besides, the “universal” background class (often set to cover the distracting or uninteresting classes Everingham et al., 2010) adds additional complicity to image segmentation (Mostajabi et al., 2015). Thus, it remains questionable to what extent the impressive performance on existing benchmarks can be interpreted as (or translated into) real-world robustness. If “top-performing” segmentation models make sparse yet catastrophic mistakes that have not been spotted beforehand, they will fall short of the need by high-stakes applications.

*The answer to the second question* constitutes the main body of our technical work. In order to identify sparse failures of existing segmentation models, it is necessary to expose

---

<sup>1</sup> According to Everingham et al. (2010) and our practice, it can easily take ten times as long to segment an object than to draw a bounding box around it.



**Fig. 1** Proposed framework for troubleshooting segmentation models (Color figure online)

them to a much larger corpus of real-world labelled images (on the order millions or even billions). This is, however, implausible due to the expensiveness of dense labelling in image segmentation. The core question boils down to: how to efficiently decide what to label from the massive unlabelled images, such that a small number of annotated images maximally expose corner-case defects, and can be leveraged to improve the models.

In this paper, we introduce a two-stage framework with human-in-the-loop for efficiently troubleshooting image segmentation models (see Fig. 1). The first stage automatically mines, from a large pool  $\mathcal{D}$  of unlabelled images, a small image set  $\mathcal{M}$ , which are the most informative in exposing weaknesses of the target model. Specifically, inspired by previous studies on model falsification as model comparison (Ma et al., 2018; Wang et al., 2020; Wang & Simoncelli, 2008), we let the target model compete with a set of state-of-the-art methods with different design methodologies, and sample images by MAXimizing the Discrepancy (MAD) between the methods. To reduce the number of false positives, we propose a weakly labelling method of filtering  $\mathcal{M}$  to obtain a smaller refined set  $\mathcal{S}$ , subject to segmentation by human subjects. In the second stage, we fine-tune the target model to learn from the counterexamples in  $\mathcal{S}$  without forgetting previously seen data. The two stages may be iterated, enabling progressive troubleshooting of image segmentation models. Experiments on PASCAL VOC (Everingham et al., 2010) demonstrate the feasibility of the proposed method to address this new challenging problem, where we successfully discover corner-case errors of a “top-performing” segmentation model (Chen et al., 2017), and fix it for improved generalization in the wild.

## 2 Related work

### 2.1 MAD competition

The proposed method takes inspiration from the MAD competition (Wang & Simoncelli, 2008; Wang et al., 2020) to efficiently spot model failures. Previous works focused on performance evaluation. We take one step further to also fix the model errors detected in the MAD competition. To the best of our knowledge, our work is the first to extend the MAD idea to image segmentation, where labeling efficiency is more desired since pixel-wise

human annotation for image segmentation is much more time-consuming than image quality assessment (Wang & Simoncelli, 2008; Wang et al., 2021), image classification (Wang et al., 2020) and image enhancement (Cao et al., 2021) tasks previously explored.

## 2.2 Differential testing

Our method is also loosely related to the cross-disciplinary field of software/system testing, especially the differential testing technique (McKeeman, 1998). By providing the same tests to multiple software implementations, differential testing aims to find bug-exposing test samples that lead to different results. Programmers can then dig into these test cases for potential bug fixing. While debugging software is very different from troubleshooting machine learning algorithms, a handful of recent work explored this idea to find pitfalls of deep learning-based autonomous driving systems (Pei et al., 2017; Tian et al., 2018; Zhang et al., 2018). Aiming to be fully automated without human intervention, these methods have to make strong assumptions such as the ground truth labels can be determined by majority vote<sup>2</sup> or are unchanged under some synthetic image transformations<sup>3</sup> (e.g., brightness and contrast change, or style transfer). Therefore, it is unclear how to generalize the results obtained in such noisy and often unrealistic settings, to the real world with both great content fidelity and diversity.

## 2.3 Adversarial examples

Introduced by Dalvi et al. (2004) and reignited by Szegedy et al. (2013), most adversarial attacks add small synthetic perturbations to inputs of computational models that cause them to make incorrect decisions. In image classification, Hendrycks et al. (2019) identified a set of natural images that behave like synthetic adversarial examples, which possess inherent transferability to falsify different image classifiers with the same type of errors. The selected counterexamples by the proposed framework might be treated as a new type of natural adversarial examples, that force the two models to make distinct predictions, therefore capable of fooling at least one model. Similar as natural counterexamples focused in this work, synthetic adversarial examples pose security risks of deploying machine learning algorithms in real-world applications. A large body of research (Chen et al., 2020; Madry et al., 2018; Zhang et al., 2019) delves deep into adversarial training, trying to defend against adversarial perturbations at the expensive cost of sacrificing the generalization on original test sets without perturbations (Schmidt et al., 2018; Tsipras et al., 2019; Zhang et al., 2019). This seems to suggest a trade-off between generalization to real-world benign examples and robustness to adversarial attacks. Wang et al. (2020). Recent research has shown that label noise (Luo et al., 2020; Wei et al., 2019; Wang et al., 2021; Xia et al., 2020) in training data may be the cause of adversarial vulnerability of deep neural networks (Sanyal et al., 2020). Readers of further interest are referred to Mohseni et al. (2021).

<sup>2</sup> As per (Hendrycks et al., 2019), machine learning algorithms with similar design philosophies tend to make common mistakes.

<sup>3</sup> In many areas of computer vision, methods trained on synthetic data cannot generalize to realistic data, and specialized techniques such as domain adaptation (Zhang et al., 2017; Zhao et al., 2019) have to be used to bridge the performance gap.

## 2.4 Semantic segmentation with deep learning

Fully convolutional network (Long et al., 2015) was among the first deep architectures adopted for high-quality segmentation. Skip connection (Ronneberger et al., 2015), recurrent module (Zheng et al., 2015), max index pooling (Badrinarayanan et al., 2017; Noh et al., 2015), dilated convolution (Chen et al., 2014; Chen et al., 2018; Yu and Koltun, 2016), and multi-scale training (Chen et al., 2016; Chen et al., 2018; Chen et al., 2020) are typical strategies to boost the segmentation performance. Conditional random fields (Ladický et al., 2010) used to dominate image segmentation before the advent of deep learning were also combined with convolutional networks to model spatial relationships (Zheng et al., 2015). We refer interested readers to Minaee et al. (2020) for a comprehensive survey of this field. Weakly supervised segmentation aims to learn pixel-wise dense predictions from training samples with partial annotations such as image-level tags, bounding boxes, or labeled pixels (Ahn et al., 2019; Gong et al., 2020; Ke et al., 2021; Lee et al., 2021; Zhou et al., 2020). Our method is orthogonal to those weakly-supervised segmentation methods. The core of our method is selecting the most informative errors of the current model and then fix them by finetuning. After selecting those corner-case samples, we can ask human annotators to provide any form of annotations. Such annotations are not limited to the pixel-wise dense labeling as done in our experiments, but can also be the image-level tags, object bounding boxes, or labeled points as done in weakly-supervised segmentation.

## 2.5 Active learning

Active learning aims to improve model performance by selecting data on which to learn. Our method generally falls into the category of disagreement-based active learning (Dasgupta et al., 2007; Hanneke, 2009; Hanneke, 2007). The different between our method and previous work is that we use a novel MAD competition to select the most controversial samples on which the models disagree on.

## 3 Proposed method

Suppose we have a target segmentation model  $f_t : \mathbb{R}^{h \times w} \mapsto \{1, \dots, c\}^{h \times w}$ , where  $h$  and  $w$  are the height and width of the input image, and  $c$  denotes the number of categories. Our goal is to efficiently identify and fix the failure cases of  $f_t$  encountered in the real world, while minimizing human labelling effort in this process. We start by constructing a large image database  $\mathcal{D}$ , whose collection may be guided by the keywords that represent the  $c$  categories. Rather than conducting large-scale subjective testing to obtain the ground truth segmentation map for each  $x \in \mathcal{D}$ , we choose to create a small subset of images  $\mathcal{M} \subset \mathcal{D}$ , which are strong candidates for revealing corner-case behaviors of  $f_t$ . To further reduce false positive examples in  $\mathcal{M}$ , we describe a method to gather a weak label for each  $x \in \mathcal{M}$  as an overall indicator of segmentation quality. Based on the labels, an even smaller set  $\mathcal{S} \subset \mathcal{M}$  can be obtained for dense annotation by humans. Last, we fine-tune  $f_t$  on the combination of  $\mathcal{S}$  and previously trained data, in an attempt to learn from the found failures without forgetting (Li & Hoiem, 2017).

### 3.1 Failure identification

#### 3.1.1 Constructing $\mathcal{M}$

Inspired by model falsification methodologies from computational vision (Wang & Simoncelli, 2008) and software engineering (McKeeman, 1998), we construct the set  $\mathcal{M} = \{x_i\}_{i=1}^{n_2}$  by sampling the most “controversial” images from the large-scale unlabelled database  $\mathcal{D} = \{x_i\}_{i=1}^{n_1}$ , where  $n_2 \ll n_1$ . Specifically, given the target model  $f_i$ , we let it compete with a group of state-of-the-art segmentation models  $\{g_j\}_{j=1}^m$  by maximizing the discrepancy (Wang et al., 2020) between  $f_i$  and  $g_j$  on  $\mathcal{D}$ :

$$\hat{x}^{(j)} = \arg \max_{x \in \mathcal{D}} d(f_i(x), g_j(x)), \quad j = 1, \dots, m, \quad (1)$$

where  $d(\cdot, \cdot)$  is a distance metric to gauge the dissimilarity between two segmentation maps. In practice, we use mean region intersection over union (mIoU) as  $d(\cdot, \cdot)$ .  $\hat{x}^{(j)}$  represents the most controversial image according to  $f_i$  and  $g_j$ , and therefore is the most informative in distinguishing between them. If the competing model  $g_j$  performs at a high level, and differs from the target model  $f_i$  in design,  $\hat{x}^{(j)}$  is likely to be a failure of  $f_i$ .

To avoid identifying different instantiations of the same underlying root cause (Pei et al., 2017) and to encourage content diversity of the candidate images, we describe a “content-aware” method for constructing  $\mathcal{M}$ . We first partition  $\mathcal{D}$  into  $c$  overlapped subgroups  $\{\mathcal{D}_{i,k}\}_{k=1}^c$  based on  $f_i$ 's predicted maps, where  $x \in \mathcal{D}_{i,k}$  if at least one pixel in  $f_i(x)$  belongs to the  $k$ -th category. In other words, images with the same predicted objects are put into the same subgroup. This is to guarantee the content diversity of selected corner-case samples, since we want to expose and fix as much weakness of the model as possible, instead of only focusing on a handful of most dominant weakness. After that, we add a content constraint by restricting the size of predicted pixels in the  $k$ -th category, i.e.,  $\sum \mathbb{1}[f_i(x) == k]/(h \times w)$ , within the range of  $[p_k, q_k]$ . This allows excluding images of exceedingly large (or small) object sizes, which may be of less practical relevance. Moreover, instead of focusing on the most controversial example defined in Eq. (1), we look at top- $n_3$  images in  $\mathcal{D}_{i,k}$  with  $n_3$  largest distances computed by

$$\left\{ \hat{x}_i^{(j,k)} \right\}_{i=1}^{n_3} = \arg \max_{\{x_i\}_{i=1}^{n_3} \in \mathcal{D}_{i,k}} \sum_{i=1}^{n_3} d(f_i(x_i), g_j(x_i)), \quad (2)$$

$$j = 1, \dots, m, \quad k = 1, \dots, c.$$

We then repeat this procedure, but with the roles of  $f_i$  and  $g_j$  reversed. That is, we partition  $\mathcal{D}$  into  $c$  subgroups  $\{\mathcal{D}_{j,k}\}_{k=1}^c$  according to  $g_j$ , and solve the maximization problem over  $\mathcal{D}_{j,k}$ . Finally, we gather all candidate images to arrive at the set  $\mathcal{M} = \{x_i\}_{i=1}^{n_2}$ , where  $n_2 \leq 2m c n_3 \ll n_1$ .<sup>4</sup>

<sup>4</sup> We have  $n_2 \leq 2m c n_3$  because the same images may be optimal in different problems specified in Eq. (2).

---

**Algorithm 1:** The proposed framework for efficiently troubleshooting segmentation models
 

---

**Input:** An unlabelled image set  $\mathcal{D}$ , a target model  $f_t^{(0)}$  and the dataset  $\mathcal{S}^{(0)}$  on which it is pre-trained, a group of competing models  $\{g_j\}_{j=1}^m$ , the maximum number  $r$  of fine-tuning rounds, hyper-parameters  $n_3, n_4$

**Output:** Improved  $f_t^{(r)}$

```

1  $\mathcal{S} \leftarrow \emptyset$ 
2 for  $j \leftarrow 1$  to  $m$  do
3   | Compute segmentation predictions  $\{g_j(x), x \in \mathcal{D}\}$ 
4 end
5 for  $i \leftarrow 0$  to  $r - 1$  do
6   |  $\mathcal{M}^{(i+1)} \leftarrow \emptyset$ 
7   | Compute segmentation predictions  $\{f_t^{(i)}(x), x \in \mathcal{D}\}$ 
8   | for  $j \leftarrow 1$  to  $m$  do
9     | Compute the distances  $\{d(f_t^{(i)}(x), g_j(x)), x \in \mathcal{D}\}$ 
10    | Divide  $\mathcal{D}$  into  $c$  subgroups according to  $f_t^{(i)}$ 
11    | Filter images by the content constraint
12    | Select top- $n_3$  images by solving Eq. equation 2 to form  $\mathcal{M}_j$ 
13    |  $\mathcal{M}^{(i+1)} \leftarrow \mathcal{M}^{(i+1)} \cup \mathcal{M}_j$ 
14    | Reverse the roles of  $f_t^{(i)}$  and  $g_j$ , and repeat Steps 10 to 13
15  | end
16  | Source weak human scores for  $\mathcal{M}^{(i+1)}$ 
17  | Select top- $n_4$  images with the lowest quality scores and collect pixel-level labels
    | from humans to form  $\mathcal{S}^{(i+1)}$ 
18  |  $\mathcal{S} \leftarrow \mathcal{S} \cup \mathcal{S}^{(i+1)}$ ;
19  | Fine-tune  $f_t^{(i)}$  on the combination of  $\mathcal{S}$  and  $\mathcal{S}^{(0)}$ 
20 end

```

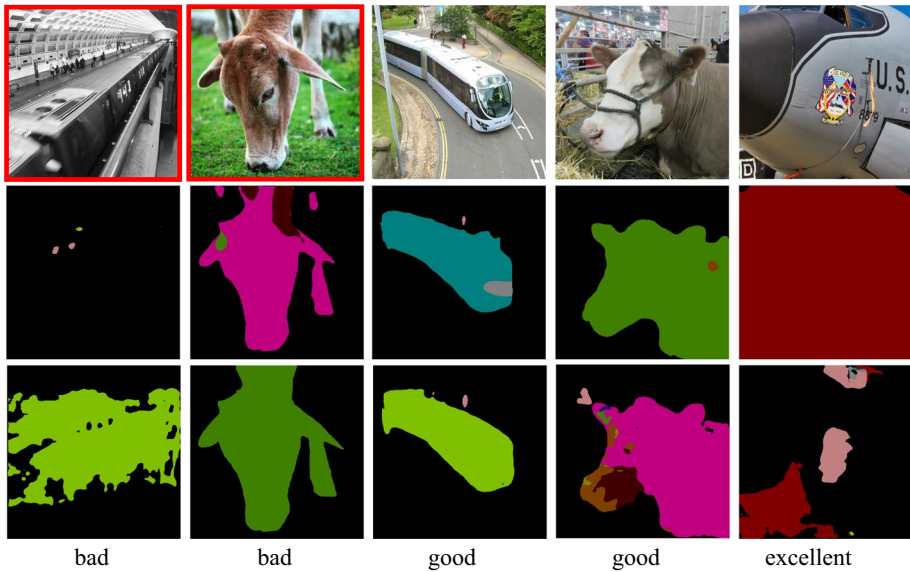
---

### 3.1.2 Constructing $\mathcal{S}$

Although images in the candidate set  $\mathcal{M}$  have great potentials of being counterexamples of  $f_t$ , some false positives may be included (see Fig. 2), especially when  $g_j$  is inferior to  $f_t$ . In addition, no data screening is involved in the construction of  $\mathcal{D}$ , increasing chances of including images that are out-of-domain (e.g., falling out of the  $c$  categories and/or containing inappropriate content). In view of these, we reduce false positives in  $\mathcal{M}$  via a weakly labelling strategy. For each  $x \in \mathcal{M}$ , we ask human subjects to give a discrete score on an absolute category rating (ACR) scale to indicate the segmentation quality of  $f_t(x)$ . The labels on the scale are “bad”, “poor”, “fair”, “good”, and “excellent” (see Fig. 8). We then rank all images in  $\mathcal{M}$  by the mean opinion scores, and choose top- $n_4$  images with the smallest scores to form the counterexample set  $\mathcal{S} = \{x_i\}_{i=1}^{n_4}$ . Finally, we seek pixel-level segmentation results from human annotators for each image in  $\mathcal{S}$  (see Fig. 9).

### 3.2 Model rectification

The labelled images in  $\mathcal{S}$  give us a great opportunity to improve  $f_t$  by learning from these failures. In order to be resistant to catastrophic forgetting (McCloskey & Cohen, 1989), the fine-tuning can also include labelled data previously used to train  $f_t$ . We may



**Fig. 2** Purpose of weakly labelling. First row: Candidate images from  $\mathcal{M}$ . Counterexamples selected into  $\mathcal{S}$  are highlighted with red rectangles, and the rest are false positives. Second and third rows: Predictions by the target and competing models, respectively. The last row: Discrete quality scores provided by human annotators on the segmentation results. See Fig. 5 for the color legend (Color figure online)

iterate through the whole procedure of failure identification and model rectification several rounds, leading to a progressive two-stage framework for efficiently troubleshooting segmentation models with human-in-the-loop.

When the iterative setting is enabled, the size of  $\mathcal{S}$  is growing:  $\mathcal{S} = \bigcup_{i=1}^r \mathcal{S}^{(i)}$ , where  $\mathcal{S}^{(i)}$  is the counterexample set created in the  $i$ -th fine-tuning round. Denoting the initial target model by  $f_t^{(0)}$ , we fine-tune  $f_t^{(i-1)}$  on the combination of  $\mathcal{S}$  accumulated in the previous  $i-1$  rounds and  $\mathcal{S}^{(0)}$  used for pre-training. We summarize the proposed framework in Algorithm 1.

## 4 Experiments

In this section, we use the semantic segmentation task defined in PASCAL VOC (Everingham et al., 2010) as a specific application to demonstrate the feasibility of our method. It is worth noting that the proposed framework can be applied to other segmentation tasks, such as those required in self driving (Cordts et al., 2016; Ess et al., 2009) and medical image analysis (Ronneberger et al., 2015).

### 4.1 Experimental setups

#### 4.1.1 Segmentation models

In our experiments, we choose the target model  $f_t$  to be the state-of-the-art DeepLabV3Plus (Chen et al., 2018) with DRN (Yu et al., 2017) as the backbone (termed



DeepLabV3P-DRN). We include five competing models: DeepLabV3Plus with ResNet101 (He et al., 2016) (termed DeepLabV3P-RN101), DeepLabV3 (Chen et al., 2017) with ResNet101 (termed DeepLabV3-RN101), DFN (Yu et al., 2018), Light-Weight RefineNet (Nekrasov et al., 2018) with ResNet50 and MobileNetV2 (Sandler et al., 2018) (termed LRefineNet-RN50 and LRefineNet-MNV2, respectively). Publicly available pre-trained weights on PASCAL VOC 2012 (Everingham et al., 2010) are used for all models. Following (Chen et al., 2017; Chen et al., 2018), images are cropped and resized to  $513 \times 513$  before inference.

#### 4.1.2 Constructing $\mathcal{D}$

We first briefly introduce the semantic segmentation database in PASCAL VOC (Everingham et al., 2010). It contains 1,464 images for training (denoted by  $\mathcal{S}^{(0)}$ ) and 1,149 for validation (denoted by  $\mathcal{T}^{(0)}$ ) with 20 scene categories (e.g., aeroplane, bicycle and person). We use the 20 class labels and combinations of them as keywords to crawl images from the Internet. No other constraints are imposed during data collection. As a result, the database  $\mathcal{D}$  includes a total of more than 40,000 images, which is much larger than PASCAL VOC training and validation sets.

#### 4.1.3 Constructing $\mathcal{M}$

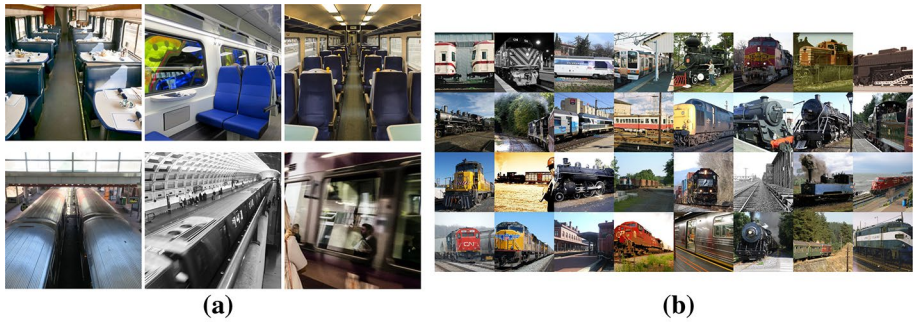
We empirically set the two parameters  $p_k$  and  $q_k$  in the content constraint to the first and the third quartiles of the  $k$ -th category sizes of the training images in PASCAL VOC, respectively. To strike a good balance between the human effort of weakly labelling and the content diversity of  $\mathcal{M}$ , we choose  $n_3 = 25$  in Eq. (2). Due to the existence of duplicated images, the actual sizes of  $\{\mathcal{M}^{(i)}\}$  range from 1,500 to 2,000 for different rounds.

#### 4.1.4 Constructing $\mathcal{S}$

We collect weak labels for images in  $\mathcal{M}$  from three volunteer graduate students, who have adequate knowledge of computer vision, and are told the detailed purpose of the study. Each subject is asked to give an integer score between one and five for each image to represent one of the five categories, with a higher value indicating better segmentation quality. All out-of-domain images are given a score of positive infinity, meaning that any subject is granted to eliminate an image without agreement from the other two. The mean opinion scores averaged across all subjects are used to rank images in  $\mathcal{M}$ . The hyper-parameter  $n_4$  in Algorithm 1 is set to 100. Representative examples in  $\mathcal{S}$  are shown in Figs. 3 and 4 along with images in  $\mathcal{T}^{(0)}$ . It is clear that images in  $\mathcal{S}$  are visually much harder.

We invite the same three students to provide ground truth segmentation results for images in  $\mathcal{S}$ , following the annotation guidance in PASCAL VOC, with the help of the online annotation tool Labelbox.<sup>5</sup> Our annotation process includes two stages with the goal of obtaining *consistent* segmentation maps. In the first stage, each subject is assigned one-third of the images to provide pixel-level labels. In the second stage, we carry out cross validation to improve the annotation consistency. Each student takes turn to check the segmentation maps completed by others, marking the positions and types of possible

<sup>5</sup> <https://labelbox.com/>



**Fig. 3** Visual comparison of train images in **a**  $S$  and **b** PASCAL VOC validation set (denoted by  $\mathcal{T}^{(0)}$ ). The images in  $S$  are visually much harder (Color figure online)



**Fig. 4** Visual comparison of horse images in **a**  $S$  and **b** PASCAL VOC validation set (denoted by  $\mathcal{T}^{(0)}$ ). The images in  $S$  are visually much harder (Color figure online)

annotation errors. During cross checking, the subjects can discuss with each other to reach an agreement on a small part of ambiguous annotation.

#### 4.1.5 Iterative model correction

We perform a total of  $r = 2$  fine-tuning rounds. As suggested in Chen et al. (2017), fine-tuning for each round is carried out by stochastic gradient descent with momentum 0.9 and weight decay  $5 \times 10^{-4}$ . We grid search the initial learning rate from  $\{1, 5, 10\} \times 10^{-5}$ , and choose the one with the highest mIoU on  $\mathcal{T}^{(0)}$ . The learning rate decay is guided by the polynomial policy. We set the mini-batch size and the maximum epoch number to 2 and 80, respectively.

#### 4.1.6 Model evaluation

How to reliably probe the generalization of computer vision models when deployed in the real world is by itself a challenging problem being actively investigated (Arnab et al., 2018; Hendrycks et al., 2019; Recht et al., 2019; Xia et al., 2020). Observing progress on  $\mathcal{T}^{(0)}$  is not a wise choice because this set may only contain few catastrophic failures of the target

**Table 1** Segmentation results in terms of mIoU on both  $\mathcal{T}^{(0)}$  and the unbiased test sets  $\{\mathcal{T}^{(i)}\}_{i=1}^{r+1}$ 

Fine-tuning round		0		1		2	
Test set		$\mathcal{T}^{(0)}$	$\mathcal{T}^{(1)}$	$\mathcal{T}^{(0)}$	$\mathcal{T}^{(2)}$	$\mathcal{T}^{(0)}$	$\mathcal{T}^{(3)}$
Competing models	DFN	0.8037	0.1349	0.8037	0.1054	0.8037	0.1365
	DeepLabV3-RN101	0.7795	0.1589	0.7795	0.1255	0.7795	0.1791
	DeepLabV3P-RN101	0.7843	0.2555	0.7843	0.1978	0.7843	0.1483
	LRefineNet-RN50	0.7710	0.1740	0.7710	0.1431	0.7710	0.2014
	LRefineNet-MNV2	0.7125	0.1325	0.7125	0.1194	0.7125	0.1678
Target model	DeepLabV3P-DRN	0.7887	0.1759	0.7827	0.2436	0.7828	0.4233

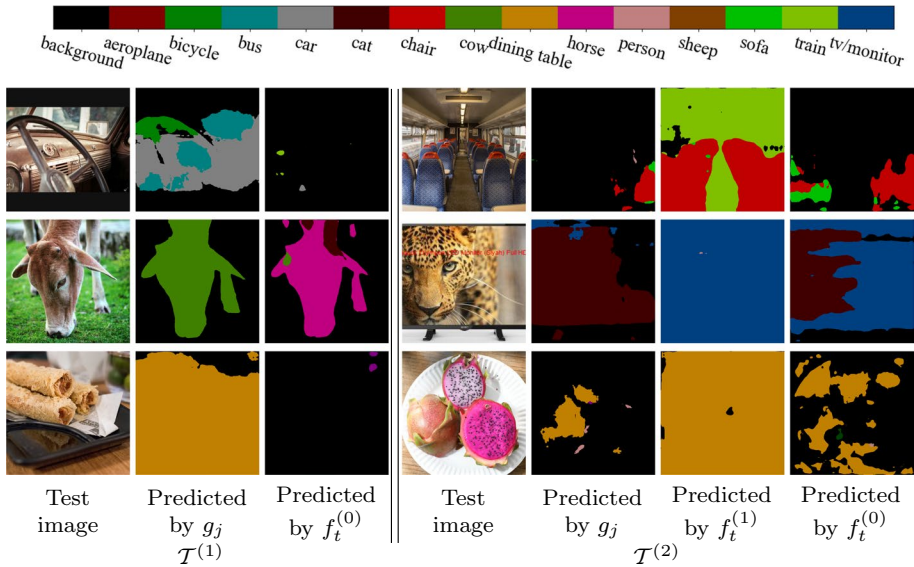
model  $f_t^{(i)}$ . It may also be unfair to employ the counterexample set  $\mathcal{S}^{(i+1)}$  to evaluate the relative progress of  $f_t^{(i)}$  against the set of competing models  $\{g_j\}_{j=1}^m$  due to the inclusion of the weakly labelling and filtering procedure. Inspired by the maximum discrepancy competition methodology for image classification (Wang et al., 2020), here we construct a new unbiased test set  $\mathcal{T}^{(i+1)} \subset \mathcal{M}^{(i+1)}$  to compare  $f_t^{(i)}$  with  $\{g_j\}_{j=1}^m$  by adding the following constraints. First, similar as the construction of  $\mathcal{S}$ , all out-of-domain images are filtered out. Second, to encourage class diversity, we retain a maximum of four images that contain the same main foreground object (i.e.,  $\mathcal{T}^{(i+1)}$  has at most four “car” images). Third, to keep evaluation fair for the competing models, images used to fine-tune the target model  $f_t^{(i)}$  are not included in  $\mathcal{T}^{(i+1)}$ . In our experiments, the size of  $\mathcal{T}^{(i+1)}$  is set to 30.

## 4.2 Main results

### 4.2.1 Quantitative results

We use a standard evaluation metric - mIoU to quantify the semantic segmentation performance. All results are listed in Table 1, where we find that, before the first round of fine-tuning, all models achieve competitive results on  $\mathcal{T}^{(0)}$ , implying close-to-saturation performance on PASCAL VOC. However, when tested on  $\mathcal{T}^{(1)}$ , the performance of all models drops significantly, indicating that many images in  $\mathcal{T}^{(1)}$  are able to falsify both the target model  $f_t^{(0)}$  and the associated competing model  $g_j$ . This also provides direct evidence that hard corner cases of existing segmentation models could be exposed. It is also proof-of-concept that the selection procedure is working as intended. Moreover, the top-1 model on  $\mathcal{T}^{(0)}$  does not necessarily perform the best on  $\mathcal{T}^{(1)}$ , conforming to the results in Wang et al. (2020).

After the first round of fine-tuning,  $f_t^{(1)}$  achieves noticeable improvements on  $\mathcal{T}^{(2)}$ , whereas all competing models experience different degrees of performance drops. This suggests that the target model begins to introspect and learn from its counterexamples in  $\mathcal{S}^{(1)}$ . After the second round of fine-tuning, the mIoU of  $f_t^{(2)}$  on  $\mathcal{T}^{(3)}$  is boosted by around 18%, surpassing all competing models by a larger margin than the previous round. This shows that our method successfully learns from and combines the best aspects of the competing models to fix its own defects, with approximately the same performance on  $\mathcal{T}^{(0)}$ . In our experiments, we only perform two rounds of fine-tuning due to limited computation and human resources, while we expect further performance gains under the proposed framework if tuned for more rounds.



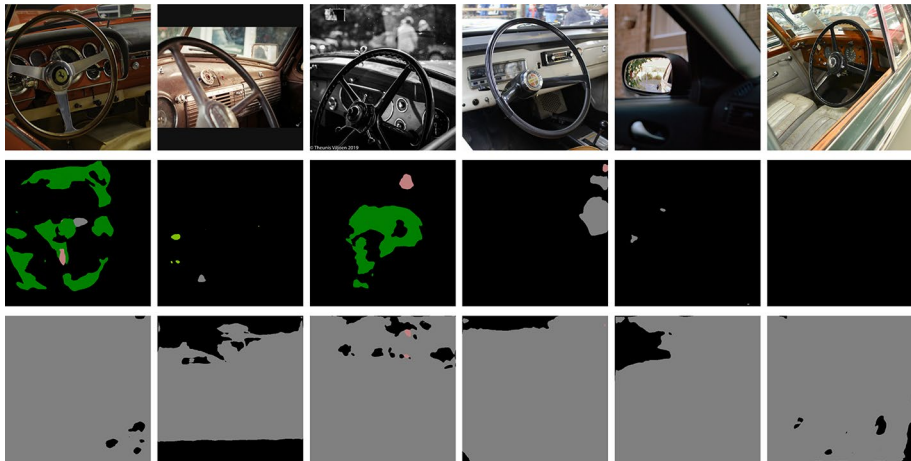
**Fig. 5** Left panel: Representative images in the test set  $\mathcal{T}^{(1)}$  before fine-tuning the target model  $f_t^{(0)}$ . From the corresponding predicted segmentation maps, we find that the competing models  $\{g_j\}$  successfully identify the failures of  $f_t^{(0)}$ . Right panel: Representative images in the test set  $\mathcal{T}^{(2)}$  after the first round of fine-tuning, where we see that  $f_t^{(1)}$  achieves noticeable performance improvements by learning from its failures in  $\mathcal{S}^{(1)}$  (Color figure online)

### 4.2.2 Qualitative results

We show representative test images before and after the first round of fine-tuning in Fig. 5. Before fine-tuning, the competing models can effectively find defects of the target model with incorrect semantics and/or boundaries. After fine-tuning, the target model does a clearly better job on the corner-case samples. We also visually compare  $f_t^{(0)}$  and  $f_t^{(1)}$  on  $\mathcal{T}^{(2)}$ , where we observe the improved robustness to unseen corner-case images. Another interesting finding is that our method can fix a general class of model mistakes by learning from a small set of valuable failures. Figure 6 shows such an example, where the target model only sees 10 “car” images in  $\mathcal{S}^{(1)}$ , and is able to generalize to images with similar unusual viewpoints.

### 4.2.3 Annotation cost

As described in Sect. 4.1, when setting  $n_3 = 25$ , the size of  $\mathcal{M}^{(i)}$  ranges from 1500 to 2000, and the size of  $\mathcal{S}^{(i)}$  is always 100 since we set  $n_4 = 100$ . In this section, we take the first round of finetuning for example to show the annotation cost. In the first round, we have  $|\mathcal{M}^{(1)}| = 1843$  and  $|\mathcal{S}^{(1)}| = 100$ . If we do not introduce the weakly labeling process, we will need to do pixel-wise annotation on all 1843 images. In contrast, with the weakly labeling process, we only need to give a global quality score for each of the 1843 images and then densely label the selected 100 worst-case images. According to our statistics, the average time cost of weakly labeling is around 7 seconds per image, and that for densely pixel-wise



**Fig. 6** Representative “car” images from  $\mathcal{M}^{(2)}$ . First row: Test images. Second row: Predictions by  $f_t^{(0)}$ . Third row: Predictions by  $f_t^{(1)}$  that is only exposed to a small set of “car” images in  $\mathcal{S}^{(1)}$ . The generalization of the target model on “car” images with unusual viewpoints is largely improved after the first round of fine-tuning. (All images shown in this figure are *not* in training set of  $f_t^{(1)}$ ) (Color figure online)

labeling is 135 seconds per image, using our own GUIs (see Appendix for the interface). To finish the first round of labeling,  $1843 \text{ s} \times 135 = 69.11 \text{ h}$  is needed without the weakly labeling phase. In contrast, only  $7 \text{ s} \times 1843 + 100 \text{ s} = 7.33 \text{ h}$  is needed if we add the weakly labeling phase. Adding the weakly labeling phase reduces the annotation cost by around ten times.

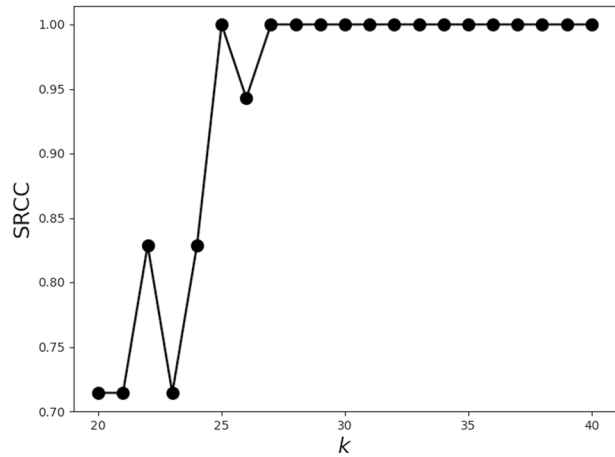
### 4.3 Ablation study

In order to show the testing results hold as we increase the size of  $\mathcal{T}^{(i)}$  (denoted as  $k$ ), we first evaluate the performance of the six models with different number of testing images, and conduct a global ranking among the six models. We use top- $k$  ranking to denote the global ranking evaluated with  $k$  testing images. Then we calculate the SRCC values between the top-40 ranking (as reference) and other top- $k$  rankings with  $k \in \{20, 21, \dots, 39\}$ . As shown in Fig. 7, the ranking results remain stable when  $k > 25$ . And our active finetuned target model always achieves the best performance on  $\mathcal{T}^{(1)}$  among the six models, for any  $k$  between 20 and 40.

### 4.4 Comparison with entropy-based method

Entropy has been used as an uncertainty metric to select hard samples (Joshi et al., 2009). The selected hard samples can be further used as training data to improve target model performance. In this section, we compare our method with the entropy-based method in Joshi et al. (2009) to show the advantage of our method than traditional entropy-based methods. Since our goal is to improve model performance on open-world images instead of fixed benchmark datasets, we need to tell which method can achieve models with better performance on open-world images, which can be efficiently done by the MAD competition in Wang et al. (2020).

**Fig. 7** The SRCC values between the top-40 ranking and other top- $k$  rankings



**Table 2** MAD competition results on two models achieved by our method and entropy-base method respectively

Method	Entropy-based	Ours
mIoU	19.65	26.53

For both methods, we conduct fine-tuning for one round and all experimental settings (including target model structure, human labelling budget, etc.) are kept identical. The only difference comes from the sampling strategy: our method samples according to Algorithm 1, while entropy-based method selects images with largest entropy. We then conduct MAD competition (Wang et al., 2020) on the two models obtained from each method to compare their performance on open-world images. Numerical results are shown in Table 2. As we can see, the model achieved by our method is winning the MAD competition with a noticeable mIoU advantage, showing that the model achieved by our method has better performance on open-world images than the model fine-tuned by entropy-based method.

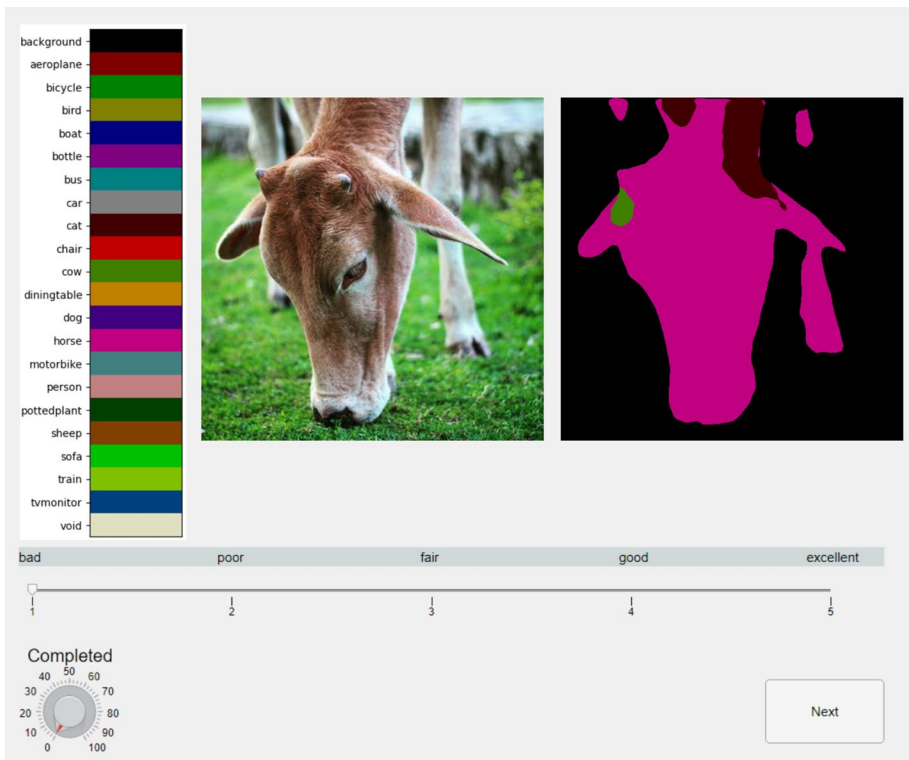
## 5 Conclusion

We have formulated a new computer vision problem that aims to efficiently expose the failures of top-performing segmentation models and learn from such failures for improved generalization in the real world. We proposed an efficient human-in-the-loop framework for troubleshooting segmentation models, and demonstrated its promise under realistic settings. We hope this work sheds light on a new line of research that requires both machine vision and human interaction. In the future, we plan to explore the idea of leveraging natural failures for improving model robustness in broader vision problems, such as video understanding (Tran et al., 2015), computer-aided diagnosis (Mansoor et al., 2015), and computational neuroscience (Golan et al., 2019). Moreover, while the current work remains to focus on improving the model's standard generalization (with an emphasis on

natural corner-case samples), our future study will investigate how this could be jointly optimized with adversarial robustness.

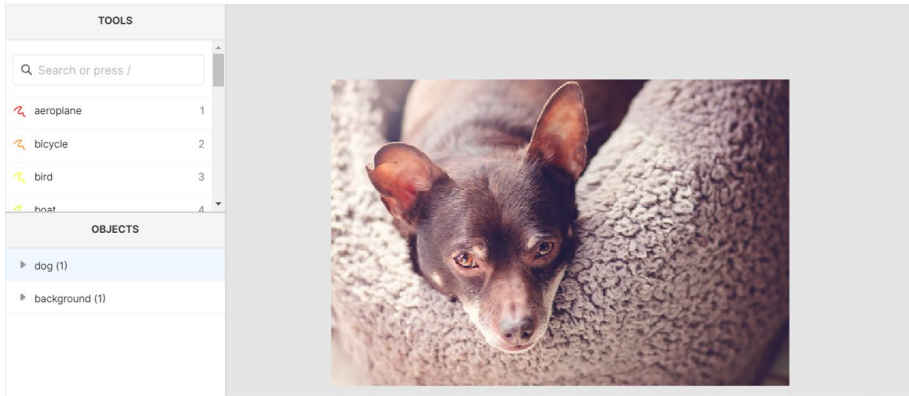
## Appendix: Subjective testing GUIs

Graphical user interfaces (GUIs) for our weakly and pixel-wise labelling experiments are shown in Figs. 8 and 9, respectively. For weakly labelling experiments, we build our own GUI: an image is rendered alongside the prediction made by the target model; a scale-and-slider applet is utilized to collect the absolute category rating score of that image as described in Sect. 3.1. For pixel-wise segmentation labelling experiments, we use Label-Box segmentation template<sup>6</sup> to build our GUI.



**Fig. 8** GUI for our weakly labelling experiments (Color figure online)

<sup>6</sup> <https://labelbox.com/product/image-segmentation>



**Fig. 9** LabelBox GUI for our pixel-level segmentation labelling experiments (Color figure online)

**Author Contributions** All authors whose names appear on the submission made substantial contributions to the conception or design of the work; or the acquisition, analysis, or interpretation of data; or the creation of new software used in the work; drafted the work or revised it critically for important intellectual content; approved the version to be published; and agree to be accountable for all aspects of the work in ensuring that questions related to the accuracy or integrity of any part of the work are appropriately investigated and resolved.

**Funding** No funding was received for conducting this study.

**Data availability** The PASCAL VOC dataset is publicly available at <http://host.robots.ox.ac.uk/pascal/VOC/>.

**Code availability** The code is available at [https://github.com/VITA-Group/Troubleshooting\\_Image\\_Segmentation](https://github.com/VITA-Group/Troubleshooting_Image_Segmentation).

## Declarations

**Conflict of interests** Haotao Wang has been employed by Texas A&M University, University of Texas at Austin, Kwai Inc, NVIDIA Corporation, Amazon.com Services LLC. Tianlong Chen has been employed by Texas A&M University, University of Texas at Austin, Walmart Technology, Microsoft and Facebook. Zhangyang Wang have been employed by University of Illinois Urbana-Champaign, Texas A&M University, University of Texas at Austin, US Army Research Lab, Adobe, Microsoft, and Amazon. Kede Ma has been employed by University of Waterloo, New York University and City University of Hong Kong.

**Ethical approval** All procedures performed in studies involving human participants were in accordance with the ethical standards of the institutional and/or national research committee and with the 1964 Helsinki declaration and its later amendments or comparable ethical standards.

**Consent to participate** All human participants consented for participating in this study.

**Consent for publication** All contents in this paper are consented for publication.



## References

- Ahn, J., Cho, S., & Kwak, S. (2019). Weakly supervised learning of instance segmentation with inter-pixel relations. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition* (pp. 2209–2218).
- Arnab, A., Miksik, O., & Torr, P. H. (2018). On the robustness of semantic segmentation models to adversarial attacks. In *IEEE conference on computer vision and pattern recognition (CVPR)* (pp. 888–897).
- Badrinarayanan, V., Kendall, A., & Cipolla, R. (2017). SegNet: A deep convolutional encoder-decoder architecture for image segmentation. *IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)*, 39(12), 2481.
- Bischke, B., Helber, P., Borth, D., & Dengel, A. (2018). Segmentation of imbalanced classes in satellite imagery using adaptive uncertainty weighted class loss. In *IEEE international geoscience and remote sensing symposium (IGARSS)* (pp. 6191–6194).
- Cao, P., Wang, Z., & Ma, K. (2021). Debiased subjective assessment of real-world image enhancement. In: *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition* (pp. 711–721)
- Chen, W., Gong, X., Liu, X., Zhang, Q., Li, Y., & Wang, Z. (2020). FasterSeg: Searching for faster real-time semantic segmentation. In *International conference on learning representations*.
- Chen, T., Liu, S., Chang, S., Cheng, Y., Amini, L., & Wang, Z. (2020). Adversarial robustness: From self-supervised pre-training to fine-tuning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition* (pp. 699–708).
- Chen, L. C., Papandreou, G., Kokkinos, I., Murphy, K., & Yuille, A. L. (2014). Semantic image segmentation with deep convolutional nets and fully connected crfs. arXiv preprint [arXiv:1412.7062](https://arxiv.org/abs/1412.7062)
- Chen, L. C., Papandreou, G., Schroff, F., & Adam, H. (2017). *Rethinking atrous convolution for semantic image segmentation*. arXiv preprint [arXiv:1706.05587](https://arxiv.org/abs/1706.05587)
- Chen, L. C., Yang, Y., Wang, J., Xu, W., & Yuille, A. L. (2016). Attention to scale: Scale-aware semantic image segmentation. In *IEEE conference on computer vision and pattern recognition (CVPR)* (pp. 3640–3649).
- Chen, L. C., Zhu, Y., Papandreou, G., Schroff, F., & Adam, H. (2018). Encoder-decoder with atrous separable convolution for semantic image segmentation. In *European conference on computer vision (ECCV)* (pp. 801–818).
- Chen, L. C., Papandreou, G., Kokkinos, I., Murphy, K., & Yuille, A. L. (2018). DeepLab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected CRFs. *IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)*, 40(4), 834.
- Cordts, M., Omran, M., Ramos, S., Rehfeld, T., Enzweiler, M., Benenson, R., Franke, U., Roth, S., & Schiele, B. (2016). The cityscapes dataset for semantic urban scene understanding. In *IEEE conference on computer vision and pattern recognition (CVPR)* (pp. 3213–3223).
- Dalvi, N., Domingos, P., Sanghai, S., & Verma, D. (2004). Adversarial classification. In *ACM SIGKDD international conference on knowledge discovery and data mining (KDD)* (pp. 99–108).
- Dasgupta, S., Hsu, D. J., & Monteleoni, C. (2007). A general agnostic active learning algorithm. *Advances in Neural Information Processing Systems*, 20, 353.
- Ess, A., Müller, T., Grabner, H., & Van Gool, L. J. (2009). Segmentation-based urban traffic scene understanding. In *British machine vision conference (BMVC)*.
- Everingham, M., Van Gool, L., Williams, C. K., Winn, J., & Zisserman, A. (2010). The Pascal visual object classes (VOC) challenge. *International Journal of Computer Vision (IJCV)*, 88(2), 303.
- Golan, T., Raju, P. C., & Kriegeskorte, N. (2019). Controversial stimuli: Pitting neural networks against each other as models of human recognition. arXiv preprint [arXiv:1911.09288](https://arxiv.org/abs/1911.09288)
- Gong, C., Yang, J., You, J. J., & Sugiyama, M. (2020). Centroid estimation with guaranteed efficiency: A general framework for weakly supervised learning. *IEEE Transactions on Pattern Analysis and Machine Intelligence*.
- Hanneke, S. (2007). A bound on the label complexity of agnostic active learning. In *Proceedings of the 24th international conference on machine learning* (pp. 353–360).
- Hanneke, S. (2009). Adaptive rates of convergence in active learning. In *COLT. (Citeseer)*
- He, K., Zhang, X., Ren, S., & Sun, J. (2016). Deep residual learning for image recognition. In *IEEE conference on computer vision and pattern recognition (CVPR)* (pp. 770–778).
- Hendrycks, D., Zhao, K., Basart, S., Steinhardt, J., & Song, D. (2019). *Natural adversarial examples*. arXiv preprint [arXiv:1907.07174](https://arxiv.org/abs/1907.07174)
- Joshi, A. J., Porikli, F., & Papanikolopoulos, N. (2009). Multi-class active learning for image classification. In *IEEE conference on computer vision and pattern recognition (CVPR)* (pp. 2372–2379).

- Ke, T. W., Hwang, J. J., & Yu, S. X. (2021). Universal weakly supervised segmentation by pixel-to-segment contrastive learning. arXiv preprint [arXiv:2105.00957](https://arxiv.org/abs/2105.00957)
- Kervadec, H., Bouchtiba, J., Desrosiers, C., Granger, E., Dolz, J., & Ayed, I.B. (2019). Boundary loss for highly unbalanced segmentation. In *International conference on medical imaging with deep learning (MIDL)* (pp. 285–296).
- Ladický, L., Sturgess, P., Alahari, K., Russell, C., & Torr, P. H. (2010). What, where and how many? Combining object detectors and CRFs. In *European conference on computer vision (ECCV)* (pp. 424–437).
- Lee, J., Yi, J., Shin, C., & Yoon, S. (2021). BBAM: Bounding box attribution map for weakly supervised semantic and instance segmentation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition* (pp. 2643–2652).
- Li, Z., & Hoiem, D. (2017). Learning without forgetting. *IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)*, 40(12), 2935.
- Lin, T. Y., Maire, M., Belongie, S., Hays, J., Perona, P., Ramanan, D., Dollár, P., & Zitnick, C. L. (2014). Microsoft COCO: Common objects in context. In *European conference on computer vision (ECCV)* (pp. 740–755).
- Long, J., Shelhamer, E., & Darrell, T. (2015). Fully convolutional networks for semantic segmentation. In *IEEE conference on computer vision and pattern recognition (CVPR)* (pp. 3431–3440).
- Luo, Y., Han, B., & Gong, C. (2020). A bi-level formulation for label noise learning with spectral cluster discovery. In *IJCAI* (pp. 2605–2611).
- Madry, A., Makelov, A., Schmidt, L., Tsipras, D., & Vladu, A. (2018). Towards deep learning models resistant to adversarial attacks. In *International conference on learning representations (ICLR)*.
- Ma, K., Duanmu, Z., Wang, Z., Wu, Q., Liu, W., Yong, H., et al. (2018). Group maximum differentiation competition: Model comparison with few samples. *IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)*, 42(4), 851.
- Mansoor, A., Bagci, U., Foster, B., Xu, Z., Papadakis, G. Z., Folio, L. R., et al. (2015). Segmentation and image analysis of abnormal lungs at CT: Current approaches, challenges, and future trends. *Radiographics*, 35(4), 1056.
- McCloskey, M., & Cohen, N. J. (1989). Catastrophic interference in connectionist networks: The sequential learning problem. In *Psychology of learning and motivation* (Vol. 24, pp. 109–165). Elsevier.
- McKeeman, W. M. (1998). Differential testing for software. *Digital Technical Journal*, 10(1), 100.
- Minae, S., Boykov, Y., Porikli, F., Plaza, A., Kehtarnavaz, N., & Terzopoulos, D. (2020). Image segmentation using deep learning: A survey. arXiv preprint [arXiv:2001.05566](https://arxiv.org/abs/2001.05566)
- Mohseni, S., Wang, H., Yu, Z., Xiao, C., Wang, Z., & Yadawa, J. (2021). Practical machine learning safety: A survey and primer. arXiv preprint [arXiv:2106.04823](https://arxiv.org/abs/2106.04823)
- Mostajabi, M., Yadollahpour, P., & Shakhnarovich, G. (2015). Feedforward semantic segmentation with zoom-out features. In *IEEE conference on computer vision and pattern recognition (CVPR)* (pp. 3376–3385).
- Nekrasov, V., Shen, C., & Reid, I. (2018). Light-weight RefineNet for real-time semantic segmentation. In *British machine vision conference (BMVC)*.
- Noh, H., Hong, S., & Han, B. (2015). Learning deconvolution network for semantic segmentation. In *IEEE international conference on computer vision (ICCV)* (pp. 1520–1528).
- Pei, K., Cao, Y., Yang, J., & Jana, S. (2017). DeepXplore: Automated whitebox testing of deep learning systems. In *Symposium on operating systems principles (SOSP)* (pp. 1–18).
- Recht, B., Roelofs, R., Schmidt, L., & Shankar, V. (2019). Do ImageNet classifiers generalize to ImageNet?. arXiv preprint [arXiv:1902.10811](https://arxiv.org/abs/1902.10811)
- Ronneberger, O., Fischer, P., & Brox, T. (2015). U-Net: Convolutional networks for biomedical image segmentation. In *International conference on medical image computing and computer-assisted intervention (MICCAI)* (pp. 234–241).
- Sandler, M., Howard, A., Zhu, M., Zhmoginov, A., & Chen, L.C. (2018). MobileNetV2: Inverted residuals and linear bottlenecks. In *IEEE conference on computer vision and pattern recognition (CVPR)* (pp. 4510–4520).
- Sanyal, A., Dokania, P. K., Kanade, V., & Torr, P. H. (2020). How benign is benign overfitting?. arXiv preprint [arXiv:2007.04028](https://arxiv.org/abs/2007.04028)
- Schmidt, L., Santurkar, S., Tsipras, D., Talwar, K., & Madry, A. (2018). Adversarially robust generalization requires more data. In *Advances in neural information processing systems (NeurIPS)* (pp. 5014–5026).
- Szegedy, C., Zaremba, W., Sutskever, I., Bruna, J., Erhan, D., Goodfellow, I., & Fergus, R. (2013). Intriguing properties of neural networks. In *International conference on learning representations (ICLR)*.

- Tian, Y., Pei, K., Jana, S., & Ray, B. (2018). DeepTest: Automated testing of deep-neural-network-driven autonomous cars. In *International conference on software engineering (ICSE)* (pp. 303–314).
- Tran, D., Bourdev, L., Fergus, R., Torresani, L., & Paluri, M. (2015). Learning spatiotemporal features with 3D convolutional networks. In *IEEE international conference on computer vision (ICCV)* (pp. 4489–4497).
- Tsipras, D., Santurkar, S., Engstrom, L., Turner, A., & Madry, A. (2019). Robustness may be at odds with accuracy. In *International conference on learning representations (ICLR)*.
- Wang, H., Chen, T., Wang, Z., & Ma, K. (2020). I am going MAD: Maximum discrepancy competition for comparing classifiers adaptively. In *International conference on learning representations (ICLR)*.
- Wang, Q., Han, B., Liu, T., Niu, G., Yang, J., & Gong, C. (2021). Tackling instance-dependent label noise via a universal probabilistic model. arXiv preprint [arXiv:2101.05467](https://arxiv.org/abs/2101.05467)
- Wang, Z., Wang, H., Chen, T., Wang, Z., & Ma, K. (2021). Troubleshooting Blind Image Quality Models in the Wild. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition* (pp. 16, 256–16, 265).
- Wang, H. N., Chen, T., Gui, S., Hu, T., Liu, J., & Wang, Z. (2020). Once-for-all adversarial training: In-Situ tradeoff between robustness and accuracy for free. *Advances in Neural Information Processing Systems*, 33, 7449.
- Wang, Z., & Simoncelli, E. P. (2008). Maximum differentiation (MAD) competition: A methodology for comparing computational models of perceptual quantities. *Journal of Vision*, 8(12), 8.1.
- Wei, Y., Gong, C., Chen, S., Liu, T., Yang, J., & Tao, D. (2019). Harnessing side information for classification under label noise. *IEEE Transactions on Neural Networks and Learning Systems*, 31(9), 3178.
- Xia, X., Liu, T., Han, B., Gong, C., Wang, N., Ge, Z., & Chang, Y. (2020). Robust early-learning: Hindering the memorization of noisy labels. In *International conference on learning representations*.
- Yu, F., & Koltun, V. (2016). Multi-scale context aggregation by dilated convolutions. In *International conference on learning representations (ICLR)*.
- Yu, F., Koltun, V., & Funkhouser, T. (2017). Dilated residual networks. In *IEEE conference on computer vision and pattern recognition (CVPR)* (pp. 472–480).
- Yu, C., Wang, J., Peng, C., Gao, C., Yu, G., & Sang, N. (2018). Learning a discriminative feature network for semantic segmentation. In *IEEE conference on computer vision and pattern recognition (CVPR)* (pp 1857–1866).
- Zhang, Y., David, P., & Gong, B. (2017). Curriculum domain adaptation for semantic segmentation of urban scenes. In *IEEE international conference on computer vision (ICCV)* (pp. 2020–2030).
- Zhang, H., Yu, Y., Jiao, J., Xing, E. P., Ghaoui, L. E., & Jordan, M. I. (2019). Theoretically principled trade-off between robustness and accuracy. In *International conference on machine learning (ICML)* (pp. 7472–7482).
- Zhang, M., Zhang, Y., Zhang, L., Liu, C., & Khurshid, S. (2018). DeepRoad: GAN-based metamorphic testing and input validation framework for autonomous driving systems. In *ACM/IEEE international conference on automated software engineering (ASE)* (pp 132–142).
- Zhao, S., Li, B., Yue, X., Gu, Y., Xu, P., Hu, R., Chai, H., & Keutzer, K. (2019). Multi-source domain adaptation for semantic segmentation. In *Advances in neural information processing systems (NeurIPS)* (pp. 7285–7298).
- Zheng, S., Jayasumana, S., Romera-Paredes, B., Vineet, V., Su, Z., Du, D., Huang, C., & Torr, P.H. (2015). Conditional random fields as recurrent neural networks. In *IEEE international conference on computer vision (ICCV)* (pp. 1529–1537).
- Zhou, L., Gong, C., Liu, Z., & Fu, K. (2020). SAL: Selection and attention losses for weakly supervised semantic segmentation. *IEEE Transactions on Multimedia*, 23, 1035.