# Delving into Crispness: Guided Label Refinement for Crisp Edge Detection

Yunfan Ye, Renjiao Yi, Zhirui Gao, Zhiping Cai*, Kai Xu*
National University of Defense Technology

arXiv:2306.15172v1 [cs.CV] 27 Jun 2023

*Abstract*—Learning-based edge detection usually suffers from predicting thick edges. Through extensive quantitative study with a new edge crispness measure, we find that noisy human-labeled edges are the main cause of thick predictions. Based on this observation, we advocate that more attention should be paid on label quality than on model design to achieve crisp edge detection. To this end, we propose an effective Canny-guided refinement of human-labeled edges whose result can be used to train crisp edge detectors. Essentially, it seeks for a subset of over-detected Canny edges that best align human labels. We show that several existing edge detectors can be turned into a crisp edge detector through training on our refined edge maps. Experiments demonstrate that deep models trained with refined edges achieve significant performance boost of crispness from $17.4\%$ to $30.6\%$. With the PiDiNet backbone, our method improves ODS and OIS by $12.2\%$ and $12.6\%$ on the Multicue dataset, respectively, without relying on non-maximal suppression. We further conduct experiments and show the superiority of our crisp edge detection for optical flow estimation and image segmentation.

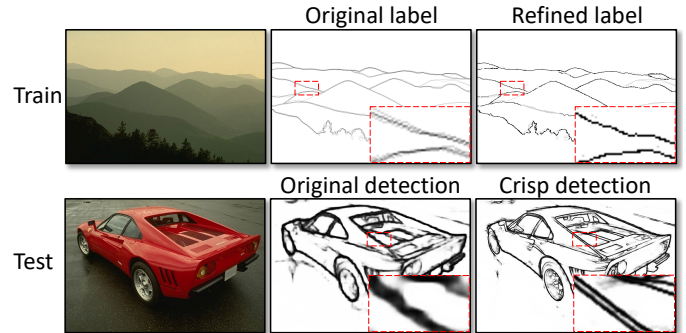*Index Terms*—Edge Detection, Crisp Edge, Label Quality, Canny-guided Refinement.



Fig. 1: We propose a guided refinement strategy for human-labeled edge maps. The refined edge maps can be used to train crisp edge detection (bottom-right) without post-processing.

## I. INTRODUCTION

Edge detection is a longstanding vision task for detecting object boundaries and visually salient edges from images. As a fundamental problem, it benefits various downstream tasks ranging from object proposal [1], [2], optical flow [3], semantic segmentation [4], [5] and image inpainting [6], [7] among others.

Traditional methods extract edges based on local features such as gradients [8], [9]. Recently, deep learning methods have achieved significant performance boost through learning more global and contextual features. Representatives include HED [10], RCF [11] and BDCN [12], which are based on the VGG [13] architecture with multi-layer deep supervision. Although these methods attain good accuracy, they all suffer from the issue of predicting thick edges and thus heavily rely on the post-process of non-maximal suppression (NMS). On the other hand, NMS finds difficulty in handling thick edges spatially close to each other due to ambiguity.

There have been some works studying how to detect *crisp* edges. Wang et al. [14] propose a refinement architecture to progressively increase the resolution of feature maps to generate crisp edges. Deng et al. [15] believe the widely used weighted cross-entropy loss [10] is the main cause of thick

edges and propose the Dice loss to tackle the problem. Compared to cross-entropy loss, however, Dice loss lacks precise localization [5] and tends to yield weaker edge response [15]. SEAL [16] and STEAL [17] refine human labels based on active contours during training, leading to crisper edges in the task of semantic edge detection.

In this paper, we observe through extensive empirical study that the noise of human labels is the main cause of thick edges by deep learning methods. It is very often that human labels deviate from the true edges due to imperfect human line drawing. We find that the more deviation the training labels have the thicker the predictions are made. Based on this observation, we advocate that more attention, if not all, should be paid on label quality rather than on model or loss design. Therefore, we propose an iterative Canny-guided refinement of human-labeled edges which can be used to learn crisp edge detection. We show that several existing edge detectors can be turned into a crisp edge detector through training on the refined edge maps, even for those using cross-entropy loss. Examples can be seen in Figure 1.

To study edge crispness quantitatively, we propose the first crispness measure of edge detection based on the ratio of edge pixel values after and before NMS. Based on the measure, we conducted experiments on the BSDS500 dataset [18] and observe the performance change of a detector trained on Canny labels perturbed with synthetic offsets of varying levels. The results verified that thick edges are mainly caused by noisy edge labels, and can be further aggravated when trained with cross-entropy loss.

To realize effective edge refinement, for each training image, we first overlap the human-labeled edge maps with the

corresponding over-detected Canny edge map and take the intersection of the two. This results in an edge map with edge segments coinciding with the true edges. We then correct the edge map by completing the segments with a pre-trained edge inpainting network. The inpainting network takes as input the gray-scale image of the input, an edge map to be completed, and a completion mask, and outputs a completed edge map. The mask is the difference between the current input edge map and the human labels, indicating which parts of the edge map need to be completed. Such completion essentially corrects the deviated portion of an edge through "guessing". The completed edges, after intersecting with Canny, are then fed into the inpainting network for the next-round refinement. The above process repeats until the change of the completion mask is negligible, implying that all human-labeled edges are collected and corrected. In essence, the interactive refinement seeks for a subset of the over-detected Canny edges aligning with human labels.

In experiments, we apply our edge refinement on several public datasets (including BSDS500 [18], Multicue [19] and BIPED [20]) and train several edge detectors (e.g., PiDiNet [21]) on the refined edge maps. The results demonstrate that the models trained with refined edges achieve significant performance boost of crispness from 17.4% to 30.6% on the public datasets. Furthermore, it brings 12.2% and 12.6% improvement for ODS and OIS, respectively on the Multicue dataset without any post-processing, compared with the original model. Besides the above experiments, we further applied our crisp edge detection to other vision tasks such as optical flow estimation and semantic segmentation. By integrating with edge-based algorithms Epicflow [3] and BNF [4], we show that crisp edge detection is essential to those vision tasks and can boost their original performance.

Our main contributions are as follows:

- Through quantitatively studying the impact of label noise on the crispness of edge detection, we observe that less label noise leads to crisper edges. In doing so, we propose the first crispness measure.
- We propose a novel Canny-guided iterative edge refinement method to correct human-labeled edges.
- We conduct extensive experiments of the proposed method and verify that corrected labels lead to significant improvement of edge crispness for existing detectors, and benefit other vision tasks.

## II. RELATED WORK

*1) Edge detecion:* The task of edge detection aims to extract object boundaries and visually salient edges from natural images. Traditional edge detectors focus on computing the image gradients to generate edges such as Sobel [8] and Canny [9]. They often suffer from noisy pixels without semantic understanding. While learning-based methods start to train detectors to create semantic contours like [18], [22]. Recently, edge detection is booming with the success of Convolutional Neural Networks (CNNs). Based on VGG [13] architecture, networks like [10]–[12] can achieve better precision performance with more semantic information. Efforts have also been

made to design lightweight architectures for efficient edge detection including [20], [21], [23].

While both precision performance and inference efficiency are improved, the issue of thick edges for result edge maps remains a problem. To evaluate the crispness of edges, [14] observe the degree of performance change by decreasing the maximal tolerant distance during the benchmark, and [15], [24] evaluate the edge maps for each model twice, one of the standard process and one of removing the post-processing. Such strategies to evaluate the edge crispness are all indirect and can not be measured by a certain value, which is inconvenient for comparison and study. To obtain crisper edges, previous works have made efforts on loss functions [15], [17], [24], refining labels during the training process [16], [17] and introducing extra modules [14], [24]. Though they can effectively obtain crisper edges, there are still some drawbacks including the balance of extra loss functions, timing decisions during the training process, or the introduction of extra modules and costs. In this paper, we explore a simple yet effective method for pre-processing the noisy labels to generate crisp edges concisely and propose the first measurable crispness evaluation metric.

*2) Image inpainting:* The task of image inpainting aims to fill missing regions in images with existing content. Driven by the progress of CNNs, the first deep learning method [25] adopts an encoder-decoder architecture and is improved in [26] by introducing a series of dilated convolution layers [27]. Some two-stage works attempt to improve the final inpainting results by introducing extra information like edges in the missing regions [6], [7], [28]. Downstream works including 3D Photography [29] can also benefit from inpainted edges. Inspired by those works, we also adopt the edge inpainting network in our proposed edge refinement method. Different from previous works, we apply the inpainting locally and adaptively based on the specific human labels guided by Canny.

*3) Edges for vision tasks:* As a fundamental vision task, edge detection can benefit kinds of mid-level and higher-level vision tasks. There are commonly two schemes to leverage edge detection, directly adopting edge maps as one of the inputs, or making edge detection play as an auxiliary task during the learning process to impose edge-preserving consistency.

For the first scheme, edge maps and the specific vision task are decoupled. For example, Edge-Boxes [1] and MCG [2] generate object bounding box proposals using edges, which can provide a sparse yet informative representation to indicate objects' positions. EpicFlow [3] targets at large displacements and computes optical flow by leveraging edges to interpolate matches of adjacent video frames, based on the observation that motion discontinuities appear most of the time at image edges. BNF [4] introduces boundary information to enhance semantic segment coherence and improve object localization based on the initial semantic segmentation.

For the second scheme, edge detection serves as an auxiliary task and is integrated with the specific vision task. Image inpainting [6], [7] leverage hallucinated edges of the missing region to avoid over-smoothed blurry and can reproduce filled regions of edge-preserving details. Moreover, by integrating

specifically designed edge detection branches, vision tasks including depth completion [30], depth estimation [31], stereo matching [32], semantic segmentation [5], instance segmentation [33] and glass-like object segmentation [34] can all benefit from extra edge-preserving consistency. In this paper, we adopt EpicFlow [3] and BNF [4] to show the superiority of crisp edge detection.

## III. STUDY OF CRISP EDGE MECHANISMS

In order to find the cause of thick edges, we do a comprehensive study of crisp edge mechanisms. We first review the loss functions of current methods, as well as some crisp label refinement works. Our assumption is that misalignment in human annotations is the main reason for thick edges in deep-learning methods, instead of the weighted binary cross entropy (W-BCE) loss as some prior works [5], [15], [24] claim. To verify our assumption. We conduct experiments on different levels of misalignment simulating human-labeled edge maps, with qualitative and quantitative comparisons. Furthermore, due to the lack of an evaluation metric for edge crispness, we also propose the first measurable crispness evaluation metric, and prove its effectiveness with several experiments.

### A. Brief Reviews of Crisp Edges

*1) A Brief Review of Loss Functions:* The task of edge detection is a pixel-wise binary classification problem. Since the numbers of edge and non-edge pixels are highly imbalanced, with the majority of pixels being non-edges. The conventional cross-entropy loss would degenerate in this case by predicting all pixels as non-edges. HED [10] applied balancing weights to normal binary cross entropy loss named W-BCE loss to solve this problem. Based on it, RCF [11] proposed an annotator-robust loss function. For the $i$th pixel in the $j$th edge map with value $p_i^j$, let $y_i$ be the ground truth edge probability of $i$th pixel, the annotator-robust loss is calculated as:

$$l_i^j = \begin{cases} \alpha \cdot \log\left(1 - p_i^j\right), & if \ y_i = 0, \\ 0, & if \ 0 < y_i < \eta, \\ \beta \cdot \log p_i^j, & otherwise, \end{cases} \quad (1)$$

in which

$$\begin{aligned} \alpha &= \lambda \cdot \frac{|Y^+|}{|Y^+| + |Y^-|}, \\ \beta &= \frac{|Y^-|}{|Y^+| + |Y^-|}, \end{aligned} \quad (2)$$

where $Y^+$ and $Y^-$ denote the number of edge and non-edge pixels in the ground truth edge maps. $\lambda$ is a parameter to balance positive and negative samples. $\eta$ is a threshold to filter out the less confident edge pixels in ground truths, to avoid confusing the network. For an edge map, the final loss is $\mathcal{L} = \sum_i^j l_i^j$.

Previous works often hold a common conjecture that, although the W-BCE loss makes the network trainable, it also causes the issue of thick edges by predicting too many edge pixels [5], [14], [15]. Dice loss is proposed in [15], which measures the overlap between predictions and ground truths. It is insensitive to the number of edge/non-edge pixels, alleviating the class-imbalance problem and helping the neural

network predict crisper edges. However, Dice loss lacks precise localization [5] and tends to yield weaker edge response compared with W-BCE loss. Currently, there is no perfect solution in terms of loss functions for crisp edge detection.

*2) Brief Reviews of Label Alignment:* Another assumption of predicting thick edges is the misalignment in edge maps annotated by different people. With inconsistent ground truth edge maps for the same data, the networks are confused and predict thick edges covering most annotated edge pixels. Therefore another direction to produce crisp edges is to refine the edge labels by removing misalignment. [35] proposes to use dense conditional random fields (dense-CRF) to refine object contour in pre-processing. Its motivation is to improve the segmentation labels, and works not well for edge detection. SEAL [16] optimizes labels by formulating a computationally expensive bipartite graph min-cost assignment problem, while STEAL [17] infers true object boundaries via a level set formulation which preserves connectivity and proximity. Both of them refine misaligned semantic contours during training and can make the contours of segmentation predictions crisper. However, such active alignment strategies heavily rely on the quality of network predictions, and have to decide the timing of label refinement during the training process.

Training with better-annotated labels can indeed make predicted edges crisper. The phenomenon indicates that the cross-entropy loss may not be the main reason leading to thick edges. Noisy labels, which are unavoidable in human annotations, are actually the main factor preventing crisp edges, the W-BCE loss just aggravates the situation. Several experiments are conducted to verify the assumption in the following sections.

### B. Crispness Evaluation

Although learning-based edge detection has made significant progress in precision, the quality of predicted edge maps is not satisfying without post-processing. Previous edge detection methods often focus on the "correctness" (the precision of edge pixel localization) instead of the "crispness" (the width of result edges). In the post-processing, a common Non-maximum Suppression (NMS) step is applied for all methods to ensure thin edges before evaluating on ground truth. However, NMS may introduce new problems such as connecting nearby edges. Predicting crisp edges directly by networks would be the best solution.

Before experiments, we define a crispness metric to measure how crisp the edge map is. There are several works trying to evaluate the crispness of predicted edge maps before, where [14] captures the crispness by decreasing the maximal permissible distance when matching ground-truth edges during the benchmark. [15], [24] examine the crispness by evaluating the edge maps before and after NMS, since predictions of high-crispness tend to perform better without the aid of NMS. However, both of them evaluate the crispness by observing the performance changes after some operations, which is indirect and can not be measured by a crispness value. In this section, to better explore crisp edge mechanisms, we propose a conceptually simple yet effective evaluation metric, termed as $Crispness$. The crispness for each edge map is calculated
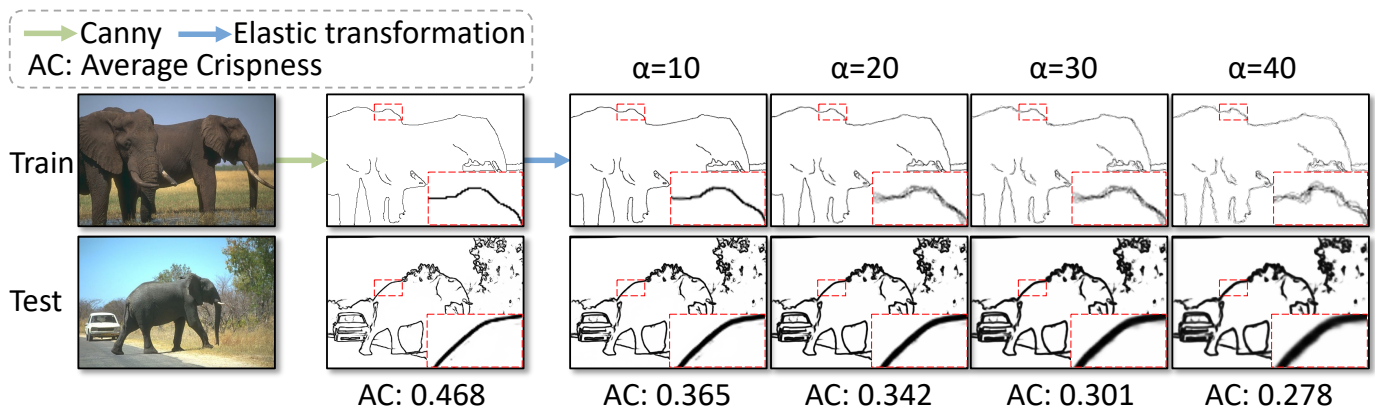
Fig. 2: Impact of noisy labels. Zoom-in is recommended for better comparisons.

as the ratio of the sum of pixel values after NMS, to the sum of pixel values before NMS:

$$Crispness = \frac{\sum_{i=1}^{w \times h} NMS(E)}{\sum_{i=1}^{w \times h} E},$$ (3)

where $E$ means the predicted edge map of size $w \times h$ and NMS represents the standard non-maximum suppression for edge detection. We take the average crispness (AC) when evaluating on a benchmark dataset.

The proposed crispness metric has several characteristics. First, it is lightweight and time-saving with about 100 FPS when evaluating on the BSDS500 dataset [18]. Second, a crispness value is calculated and ranges from 0 to 1, for more straightforward comparisons. For an edge map where the edges are at a single-pixel width, the crispness would be close or equal to 1. The crisper the edge is, the larger the crispness would be.

### C. Impact of Noisy Labels

To verify the assumption of noisy labels as the main reason for thick edges, we conduct experiments to measure the impact of noisy labels (label misalignment) on final crispness.

We apply the Canny [9] detector to images to simulate real ground truth with no offsets, since the edge pixels produced by Canny are all single-pixel width and precisely lie in the positions where gradients change most sharply. Then, we transform the Canny edge maps by moving pixels locally around using random displacement fields, called elastic transformation [36], which is widely used for data augmentation in computer vision tasks [37], [38]. The parameter $\alpha$ controls the strength of the displacement, higher values mean that pixels are moved further. We perform elastic transformation for each $\alpha$ several times, and take their average to simulate human annotations from different people.

PiDiNet [21] is adopted as the edge detection backbone in the experiment. We choose BSDS500 dataset [18] to explore the impact of noisy labels. To be specific, we train PiDiNet on the training set of BSDS500 with labels all generated by Canny, and its different elastic transformation versions as comparing groups, then evaluate the average crispness of the different network versions on the testing set. We use the

commonly-used W-BCE loss in Equation 1 to train networks. The results are demonstrated in Figure 2. We can see that the smaller the $\alpha$ is, the crisper and better the result edges are, both qualitatively and quantitatively. It means larger annotation offsets (misalignment) cause thicker edges, and the crispness is close to the one training from human annotations when $\alpha=40$. The edge map predictions by the network trained from Canny labels are much crisper than other groups, which means networks trained with W-BCE loss can still get crisp edges. The label misalignment, instead of the loss function, is the real reason for thick edges.
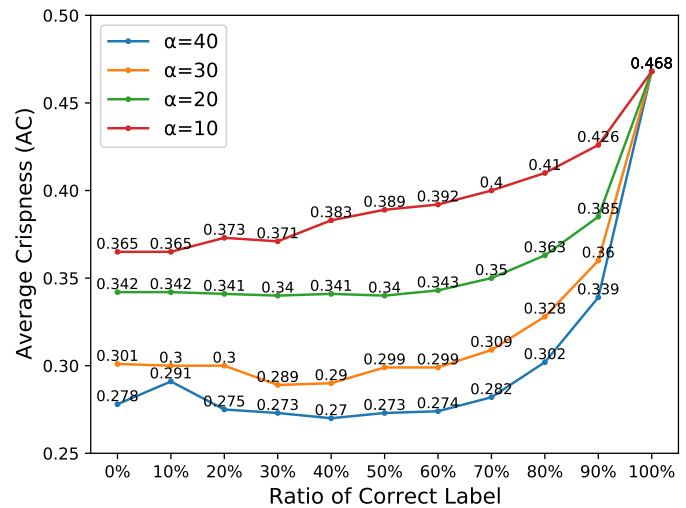


Fig. 3: Average Crispness for different percentage of correct training labels.

In real scenarios, it is common that only part of human annotations are noisy while some others are carefully annotated. Thus, we further conduct experiments with different percentages of elastic-transformed Canny labels with a fixed $\alpha$. The original Canny results are treated as "correct" labels. As shown in Figure 3, except for the line $\alpha = 10$, when we gradually increase the percentage of correct labels, the AC changes slowly before reaching 80%, while the AC increases dramatically between 90% and 100%. For the line $\alpha = 10$, which means the label offsets are very small and most labels can be considered correct enough, the AC gradually increases
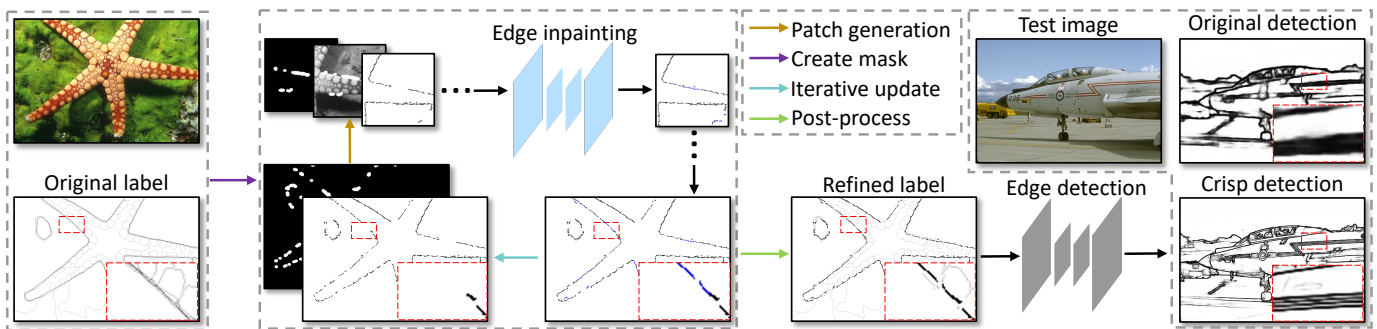
Fig. 4: The pipeline of our iterative edge label refinement and crisp edge detection.

in pace with the ratio of correct labels. Such interesting phenomenons indicate that the result edges can be crisp enough only when most of the training labels (>90%) have no offsets. It explains why most previous works consider W-BCE loss to be the main cause of thick edges, since it is hard or even impossible for human annotations to have "correct" labels of more than 90% in edge detection. We come to the conclusion that the issue of thick edges generated by neural networks is mainly caused by noisy human annotations, and aggravated by the W-BCE loss. With the premise of such a conclusion, we try to propose a general and simple edge label refinement method as a pre-processing before training edge detection networks. With less-noisy labels, the performance of the trained network would be much better in both precision and crispness. Moreover, since changes only happen in original labels, our method can be easily applied to any edge detection backbone, leading to a boost of crispness.

## IV. CRISP EDGE REFINEMENT PIPELINE

### A. Patch-based Edge Refinement

We aim to propose a technically simple but effective method that can directly pre-process and refine noisy human labels without introducing any extra modules and complicated loss combinations, enabling other backbone networks to generate crisper edges by training with refined labels. The label refinement pipeline is easy to use and can be integrated into any edge detection network, with minimal changes to the original labels. Based on the premise that Canny edges lie in the exact positions where gradients change most sharply, bringing no ambiguity when training with W-BCE loss, unlike human-annotated edges, we naturally develop the label refinement method facilitated by Canny edges. Specifically, we apply the Canny detector of a very low threshold to the multiple blurred images of the input image with different blur strategies, and fuse their results as the final Canny edges. Edges in the obtained Canny edge map are over-detected, which means the real crisp ground truth edges can be considered as a subset of these edges. Our task is to iteratively refine this subset of edges as completely as possible without any offsets, and get the crisp ground truth edge map at last.

To do so, we first calculate the overlap (Hadamard product) of human labels annotated by different people and the over-detected Canny edge map as the initial edge map. Apparently, most human-annotated edge maps have offsets to the

corresponding Canny edges, and cannot match Canny edge pixels precisely. As a result, their overlap would have many discontinuities. One straightforward approach is to inpaint the discontinuous edges through adversarial networks. Specifically, we build our edge inpainting model upon image inpainting methods in [6], [7], [39]. The key difference between our model and those recent two-stage approaches is that, their inpainting masks and the available contexts are static, while our method inpainting edges locally around each edge discontinuity adaptively. In other words, our inpaint masks are dynamic during the iterative label refinement pipeline.

We propose an adaptive patch-based edge refinement strategy to iteratively inpaint the missing edges from initial edge maps and compute the overlap with Canny edges, as illustrated in Figure 4. We first compute the Hadamard product of over-detected Canny edges from the blurred input image and human labels as initial edge, where there would be several discontinuous edges waiting for refinement (see Section V-C2 for more ablations of robustness). In order to decide which pixels to inpaint on the inital (refined) edge map, we check the neighborhood of each edge pixel in human labels, if there are no edges nearby on the inital (refined) edge map, we consider the current position as an inpainting pixel. At last, we dilate the binary mask of inpainting pixels as the inpainting mask, examples are provided in Figure 5.
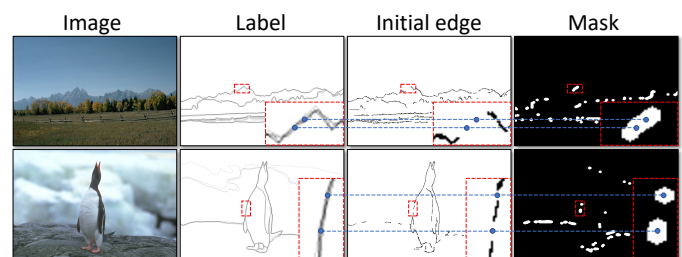


Fig. 5: Examples for generating the inpainting masks for initial edges. We check each position of edge pixels in the human label, and treat this position as an inpainting pixel (blue points are examples of some inpainting pixels) if there are no edges around the position on the initial (refined) edge map. The inpainting mask is further generated by dilating all inpainting pixels.

In the edge inpaint module, we adaptively inpaint the edge map by inputting the gray-scale image, the edge map of current iteration and the inpainting mask. The edge inpainting network is applied for each image patch and the inpainting process lasts
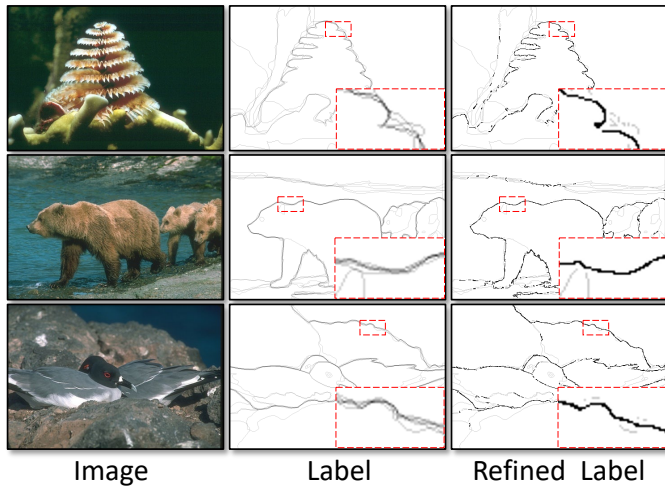
Fig. 6: Examples from BSDS500 dataset about the raw human annotations and our refined labels.

for several iterations until no new edge is created in the most recent iteration.

In this work, we employ the number of connected areas in the inpainting mask as the indicator for termination. Sometimes during the refinement process, one big connected area in the inpainting mask could be occasionally divided into two smaller ones if only the center part has been refined. However, such extreme cases are rare even in each image. Since the refinement method only inpaints edge pixels inside the inpainting mask, it can be guaranteed that the inpainting mask is always non-increasing by nature. Therefore, as the iterative refinement proceeds, the overall number of connected areas in the inpainting mask will gradually decrease until convergence. To avoid any potential infinite loop, there is also a maximum number of allowed iterations ($I_{max}$ in Algorithm 1).

At last, we set the values of those human annotated pixels that are not in the refined label to be smaller than $\eta$ (in Equation 1), so that the subsequent training ignores these unconfident pixels. See Algorithm 1 for detailed pipeline and formulations. Qualitative comparisons between original human annotations and our refined labels are illustrated in Figure 6.

### B. Upscaling Strategy for Crisp Edge Detection

Inspired by the multi-scale strategy in [11] which averages the resulting edge maps of input image pyramids to improve the performance, we propose an upscaling strategy to further make edges crisper. The original image and its up-scaled version are fed into the backbone edge detection network to get corresponding edge maps. Then, we downscale the up-scaled resulting edge map and restore it to its original size. A Hadamard multiplication is applied to the two edge maps to get a crisper result. This process is illustrated in Figure 7. With a fixed receptive field size, there is relatively less available context on the up-scaled input image compared to the image at its original size. Therefore the network is prone to generate more detailed edges. For images at the original or smaller sizes, there are more contextual information and the networks tend to create more meaningful and contextual but

---

**Algorithm 1** Iterative Edge Label Refinement.

**Input:**
   The original image $X$; The original edge label $Y$; The patch size $S$; The threshold $\eta$ as in Equation 1; The maximum iteration number $I_{max}$.

**Output:**
   The well-refined new edge label $Y'$.

1: Generate the initial edge $E$ by multiplying the over-detected Canny edge map $C$ and the original label filtered by $\eta$:
   $C = Canny(Blur(X), thres\_low)$,
   $Y[Y < \eta]] = 0$,
   $E = C \odot Y$;

2: Calculate the mask area $M$ where edge discontinuities need to be inpainted in $E_i$:
   $M = Create\_mask(Y, E)$;

3: **for** round $i$ in $\{1, ..., I_{max}\}$ **do**

4:     Generate patches $P = [p_1, p_2, ..., p_n]$ of patch size $S$ with the same center of each connected area in the mask $M$: $P = Create\_patches(M, S)$ ;

5:     **for** each patch $p_n$ in $P$ **do**

6:         Inpainting the edge discontinuities by edge patch $E(p_n)$ and the grayscale of original image patch $X_{gray}(p_n)$ guided by mask patch $M(p_n)$. Then update the patch guided by $C$:
   $E'(p_n) = Edge\_inpainting(E(p_n), X_{gray}(p_n), M(p_n))$,
   $E(p_n) = C \odot E'(p_n)$;

7:     **end for**

8:     Terminate the iterative refinement process when the current number of connected areas $N_{connect}^n$ of mask $M$ is no less than it of last iteration:
   **if** $N_{connect}^n >= N_{connect}^{n-1}$ **then**
       **break**;
   **end if**

9: **end for**

10: Post-processing to set a value smaller than $\eta$ in positions where any positive pixels in $Y$ that is negative in $E$:
   $Y' = Post\_process(Y, E, \eta)$.

---

thicker boundaries. Through the upscaling strategy, we take the advantage of different input sizes by fusing the edge maps.

This strategy has many merits. First, compared to the standard NMS, the upscaling strategy is end-to-end. While NMS has trouble separating adjacent edges that are glued together in the raw predictions, our method can clearly generate two separated edges. The upscaling strategy is a quick and effective way to refine thick edge predictions of edge detection networks to be crisper. This strategy introduces more computational and time costs to the whole pipeline, thus considering both the accuracy and speed, we set the upscaling factor as 1.5 for all cases. See Section V for more details.

## V. EXPERIMENTS

### A. Datasets and Implementation Details

*1) Benchmark Datasets:* We evaluate the proposed crisp edge method on three widely used datasets, BSDS500 [18],
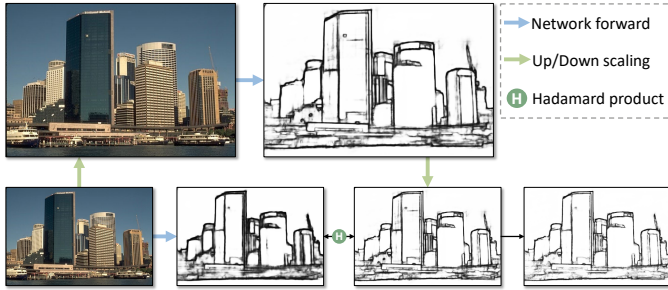
Fig. 7: The pipeline of crisp edge algorithm with the upscaling strategy.

Multicue [19], and BIPED [20]. BSDS500 consists of three splits of 200, 100 and 200 images for training, validation and testing, respectively. Each image is labeled by 4 to 9 annotators and the final edge ground truth is computed by taking the average of them. We adopt the training and validation sets of 300 images for training with data augmentation of flipping (2×), and rotation (4×). Multicue consists of images from 100 challenging natural scenes. Each scene contains a sequence of images from the left and right views, captured by a stereo camera. Only the last frame of every left-view sequence is labeled by annotators. Each of these 100 images is annotated by several people as well. We randomly split the 100 images into training and evaluation sets, consisting of 80 and 20 images respectively. The same as BSDS500 dataset, the average edge map of several annotations is computed as the ground truth. Images in the Multicue dataset are at higher resolutions, which are $1280 \times 720$. Thus we augment each image by flipping (2×), cropping (3×), and rotation (8×), leading to a training set that is 48 times larger than the original dataset. BIPED contains 250 annotated images of outdoor scenes, splitting into a training set of 200 images and a testing set of 50 images. All images are carefully annotated by experts in the computer vision field to get crisp edge labels at single-pixel width. The resolution of BIPED is also $1280 \times 720$, and images are augmented with flipping (2×), cropping (3×), and rotation (8×).

*2) Backbone networks:* Our crisp edge detection pipeline can be integrated with any edge detection backbones, here we adopt two most recent edge detectors, PiDiNet [21] and DexiNed [20], to validate the effectiveness of the proposed method. For the edge inpainting network, we adopt an existing architecture [7] to fill the edges in missing regions. The network takes the grayscale image, the incomplete edge map and the inpainting mask as inputs, and outputs an edge map filled with predicted edges in missing areas.

*3) Evaluation Metrics:* For general edge detection, the predicted edge maps are in gray-scale with pixel values in [0, 1], while the ground truth are binary in nature with each pixel labeled as either 0 (non-edge) or 1 (edge). To evaluate the precision, recall, and F-score, an optimal threshold that binarizes the predicted edge map is needed.

Following prior works, we compute the F-scores of Optimal Dataset Scale (ODS) and Optimal Image Scale (OIS), which are two strategies of determining the optimal thresholds to binarize predicted edge maps. ODS employs a fixed threshold for all images in the dataset while OIS chooses an optimal threshold for each image. F-scores are computed by $F = \frac{2 \cdot P \cdot R}{P + R}$, where $P$ denotes precision and $R$ denotes recall. Considering both accuracy and time cost, $P$ and $R$ are calculated for each threshold from 0.01 to 0.99 for every image. If multiple ground truth labels are given, such as the BSDS500 [18] and Multicue [19], $P$ is the number of edges that can match all the labels, divided by the number of edges in the binarized prediction result; $R$ is the sum of the number of edges that can match each label, divided by the sum of the number of edges from each label. For ODS and OIS, the maximum allowed distances between corresponding pixels from predicted edges and ground truths are set to 0.0075 for all experiments.

ODS and OIS emphasize more on the precision of edge pixel localization, while the main focus of this paper is improving the crispness with comparable performance. It is the reason for introducing a new metric Crispness in this paper. Noises on human labels are unavoidable, and there is no way to get true edge labels without any noise. Therefore, we apply ODS and OIS as a reference for accuracy performance and also report the proposed crispness measure as the main metric. Crispness is averaged among all images in the test dataset ("AC" in tables). The proposed method leads to significantly higher Crispness where ODS and OIS are comparable.

*4) Implementation Details:* Our method is implemented in the Pytorch [40] environment. While testing on PiDiNet and DexiNed backbones, we keep their original settings for a fair comparison. W-BCE loss in Equation 1 is adopted for training with Adam optimizer [41], at an initial learning rate of 0.005, decaying in a multi-step way. $\lambda$ is set to 1.1 and the threshold $\eta$ is set to 0.3 for all datasets while computing W-BCE losses. The architecture of the edge inpainting model is the same as [7] and trained on the Place2 [42] dataset. All images are resized to $256 \times 256$ for training with a batch size of 8, optimized by Adam optimizer. Therefore, we set the patch size as $256 \times 256$ as well while applying the locally patch-based edge refinement.

### B. Evaluations

In this section, we report the statistical evaluations of our label refinement method compared to two backbone methods PiDiNet and DexiNed. Experiments are conducted on BSDS500, Multicue and BIPED datasests. The quantitative evaluations on three datasets are in Table I. The analysis is organized in two aspects, after and before NMS, with and without the edge refinement method. Qualitative comparisons are presented in Figure 8.

As expected, we ("-R") achieve significant performance boosting on average crispness in all cases when training with our refined labels, both qualitatively and quantitatively. For the backbone of PiDiNet, our method improves the average crispness significantly by 17.4%, 22% and 28% on BSDS, Multicue and BIPED datasets. For the backbone of DexiNed, our method also improves the average crispness by 8%, 11.8% and 13.8% on BSDS, Multicue and BIPED datasets, respectively.

When non-maximum suppression (NMS) is adopted on predicted results, for models using our label refinement pipeline
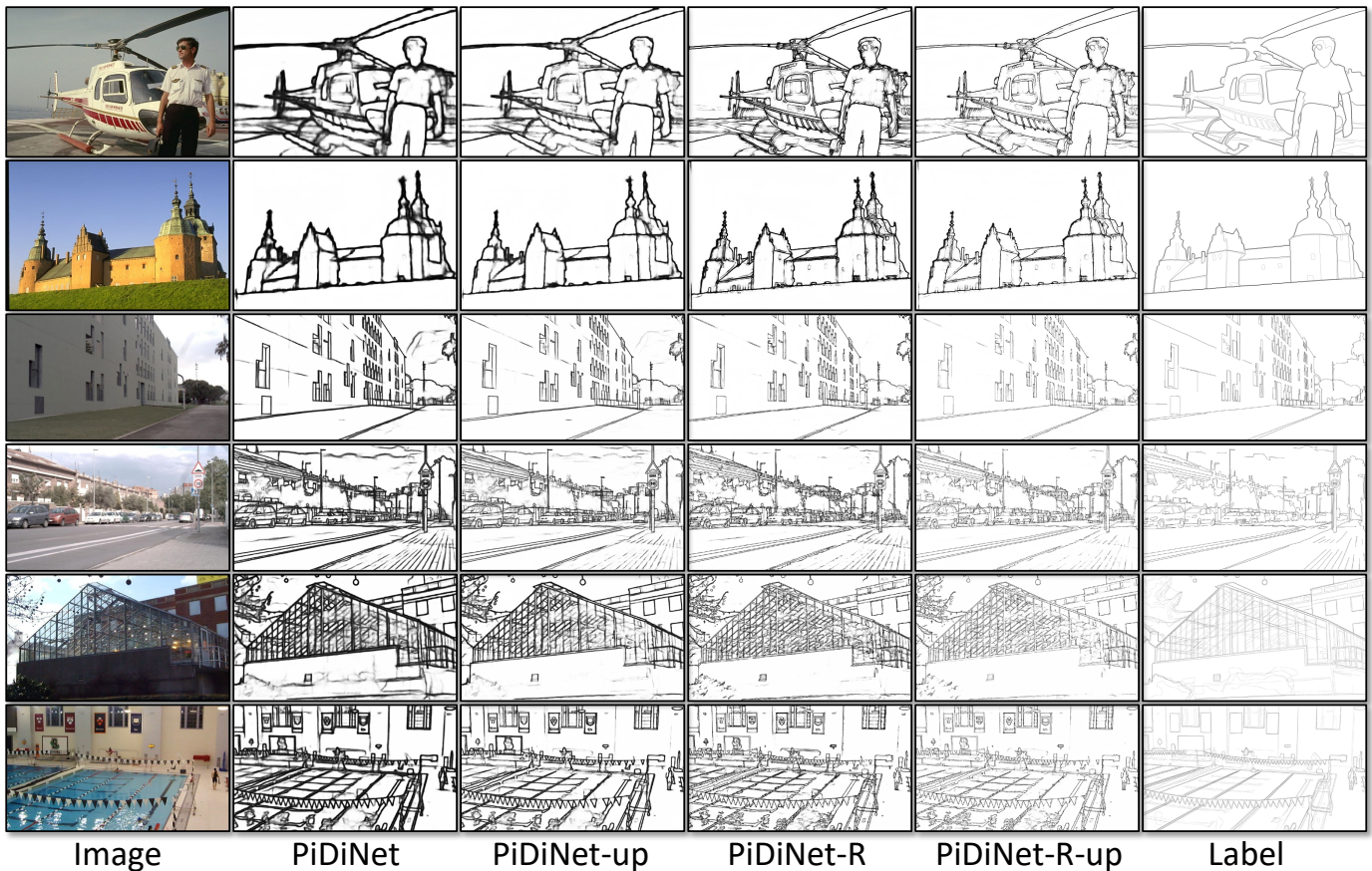
Fig. 8: Qualitative results on datasets of BSDS500 (the first two rows), BIPED (the second two rows) and Multicue (the last two rows) before NMS. Edge maps become much crisper with the aid of our edge refinement method ("-R") and upscaling strategy ("-up"). Zoom-in is recommended to see the crispness details.

| Dataset | Methods | after NMS | | before NMS | | AC |
|---|---|---|---|---|---|---|
| | | ODS | OIS | ODS | OIS | |
| BSDS 500 | PiDiNet | **0.780** | **0.794** | 0.763 | 0.776 | 0.250 |
| | PiDiNet-R | 0.771 | 0.782 | **0.771** | **0.780** | **0.424** |
| | DexiNed | 0.722 | 0.746 | 0.701 | 0.723 | 0.277 |
| | DexiNed-R | 0.719 | 0.744 | 0.708 | 0.730 | 0.357 |
| Multicue | PiDiNet | 0.874 | 0.878 | 0.764 | 0.778 | 0.204 |
| | PiDiNet-R | **0.892** | **0.907** | **0.886** | **0.904** | **0.424** |
| | DexiNed | 0.855 | 0.868 | 0.748 | 0.775 | 0.274 |
| | DexiNed-R | 0.863 | 0.877 | 0.849 | 0.863 | 0.392 |
| BIPED | PiDiNet | **0.868** | **0.876** | 0.842 | 0.852 | 0.232 |
| | PiDiNet-R | 0.863 | 0.873 | **0.863** | **0.873** | **0.512** |
| | DexiNed | 0.841 | 0.853 | 0.831 | 0.842 | 0.295 |
| | DexiNed-R | 0.858 | 0.868 | 0.855 | 0.868 | 0.433 |

TABLE I: Quantitative results for PiDiNet when training with original labels and refined labels. "-R" means training with refined labels and the same for other tables and figures.

("-R" in tables), F-scores of ODS and OIS are boosted on Multicue and BIPED dataset in terms of DexiNed backbone, and F-scores of ODS and OIS are boosted on Multicue in terms of PiDiNet backbone. We notice slight performance drops of around 1% for some conditions. The reason could be that not all human annotations are correct enough including the test datasets, and our refinement method actually changes some edge pixels with relatively big offsets, which may be unfavorable to the quantitative evaluation of our method.

Meanwhile, although NMS is a commonly used post-processing to get crisp edges, it is not differentiable, which

cannot be integrated into end-to-end training and testing. Furthermore, NMS is time-consuming. It runs at 100 fps on BSDS, whereas the inference of PiDiNet runs at 92 fps, which means adopting NMS doubles the running time, while our label refinement method does not lead to any extra time when testing. For high-resolution datasets such as BIPED, NMS takes even longer time and runs at only 25 fps. All time are reported on an RTX 2080 Ti GPU.

Without NMS ("before NMS" in tables), the performance of models training with original labels drops significantly, since thick edges tend to have high recall but very low precision, leading to low F-scores. With our refinement method, the performance boosts significantly because crisp edges predicted by networks already have a better balance between precision and recall.

For the upscaling strategy for crisp edge detection, as shown in Table II, this simple strategy ("-up" in the table) is effective in further improving the crispness. On BIPED dataset, applying the upscaling strategy to models trained with our refined labels achieve better performance on all metrics of ODS, OIS and AC, and combining with refined labels leads to a crispness improvement of 30.6%. Although on Multicue and BSDS datasets, combining label refinement and upscaling strategy leads to a slight performance drop (less than 1% when NMS is not adopted), it still provides a significant boost on

| Dataset | Methods | after NMS | | before NMS | | AC |
|---------|---------|-----|-----|-----|-----|-----|
| | | ODS | OIS | ODS | OIS | |
| BSDS 500 | PiDiNet | **0.780** | **0.794** | 0.763 | 0.776 | 0.250 |
| | PiDiNet-up | **0.780** | 0.793 | **0.773** | **0.787** | 0.323 |
| | PiDiNet-R-up | 0.774 | 0.785 | 0.771 | 0.780 | **0.522** |
| Multicue | PiDiNet | 0.874 | 0.878 | 0.764 | 0.778 | 0.204 |
| | PiDiNet-up | **0.890** | **0.900** | **0.880** | **0.890** | 0.284 |
| | PiDiNet-R-up | 0.874 | 0.893 | 0.863 | 0.887 | **0.555** |
| BIPED | PiDiNet | **0.868** | **0.876** | 0.842 | 0.852 | 0.232 |
| | PiDiNet-up | 0.860 | 0.869 | 0.850 | 0.864 | 0.318 |
| | PiDiNet-R-up | 0.862 | 0.872 | **0.860** | **0.870** | **0.624** |

TABLE II: Quantitative analysis for PiDiNet when applying the upscaling strategy ("-up") described in Section IV-B.

crispness. Comparing Table I and II, both label refinement and upscaling strategy work effectively for generating more accurate and crisper edge maps. As illustrated in Figure 8, combining both strategies also leads to significantly better (and crisper) visual results.

### C. Ablation Study

*1) Ablations for refinement strategies:* To demonstrate the effectiveness of our iterative Canny-guided edge refinement pipeline, we adopt the PiDiNet backbone and conduct several ablation studies on the BSDS500 dataset. The data splitting and augmentation are described in Section V-A.

The initial edge labels are computed by simply taking the overlap of the over-detected Canny edges and the original human labels. We find that by simply computing a Hadamard product with Canny edges for all training data, the average crispness is satisfactory enough, but the quantitative performance drops sharply without further optimization. It means, refining human labels by simply taking the overlap with Canny edges is not enough. After refining the initial edge once or iteratively by our pipeline, the predicted edges become both crisper and more accurate, which shows the effectiveness of each module in our designs.

Moreover, to verify the importance of Canny guidance, with all other settings unchanged, we conduct experiments without the Canny-guided strategy by simply dilating human labels to replace the over-detected Canny edge maps (Baseline-D), and directly generate refined training labels using NMS on the averaged labels from multiple annotators (Baseline-N). Qualitative and quantitative comparisons are shown in Figure 9 and Table III, respectively. We observe that the crispness drops significantly in all cases without Canny guidance, since Canny edges can be treated as, or close to true edges without any offsets and are essential for crisp edge detection, as studied in Section III-C. After Canny-guided refinement, the trained models can predict crisper and cleaner edges.

| Method | ODS | OIS | AC |
|--------|-----|-----|-----|
| without Canny guidance | | | |
| Baseline-D | 0.769 | 0.778 | **0.271** |
| Baseline-N | **0.775** | **0.788** | 0.252 |
| with Canny guidance | | | |
| Initial edge | 0.735 | 0.75 | 0.41 |
| Initial edge + Refine once | 0.753 | 0.766 | 0.414 |
| Patch-based refine iteratively | **0.771** | **0.782** | **0.424** |

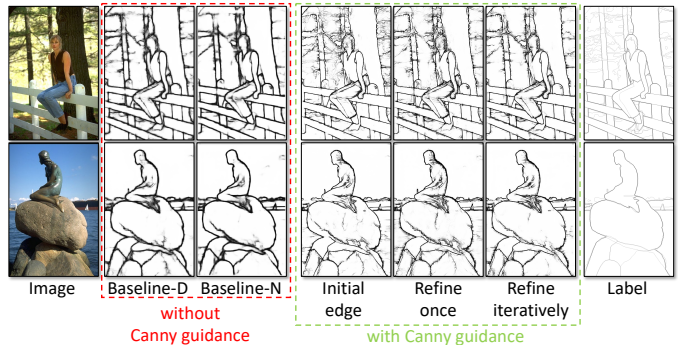TABLE III: Quantitative results for various refinement strategies.



Fig. 9: Ablations for various refinement strategies.

*2) Ablations for robustness:* Our iterative strategy starts the refinement with an initial edge generated by computing the Hadamard product of over-detected Canny edges and human labels. Therefore, it is natural to come up with a question that how robust such a seed could be. In most cases, human labels and the corresponding Canny edges can coincide within small regions, producing a relatively complete initial edge and benefiting subsequent refinement. However, there still could be extreme cases that the initial edge map is badly sparse caused by slight offsets across a large region on human labels.

We conduct experiments to test the robustness of the initial edge maps by randomly dropout different percentages of pixels (edges) on them. The original initial edge maps without any dropout are treated as the baseline version. As demonstrated in Figure 10, when the dropout ratio is no more than 50%, the recovered final edge map changes slightly compared with the original baseline, so as the quantitative performance change (∼1%) in Table IV; while if the initial edge maps are badly broken (dropout more than 70% edges), some discontinuous regions become hard to be inpainted. However, in such cases, the human labels themselves are too noisy to work as ground truths for supervised methods.
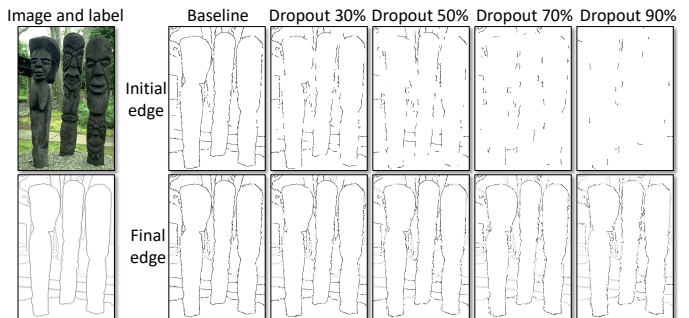


Fig. 10: Ablations for the robustness of initial edges by dropout various levels of pixels (edges).

| Method | ODS | OIS | AC |
|--------|-----|-----|-----|
| Baseline | 0.771 | 0.782 | 0.424 |
| Dropout 30% | 0.767 | 0.776 | 0.424 |
| Dropout 50% | 0.762 | 0.771 | 0.427 |
| Dropout 70% | 0.756 | 0.767 | 0.427 |
| Dropout 90% | 0.742 | 0.750 | 0.431 |

TABLE IV: Quantitative results of initial edges with different dropout levels.

*3) Ablations for other crisp edge methods:* We directly refine human labels to achieve crisp edge detection, which means our method can be easily integrated with other crisp edge detection works that focus on loss function designs [15], [24]. To thoroughly verify the compatibility of our method, we conduct experiments based on PiDiNet [21] with Dice loss [15] and tracing loss [24] following their settings, training with original labels and our refined labels, and evaluating after and before NMS.

By observing the qualitative and quantitative results in Figure 11 and Table V, several interesting conclusions can be drawn. (a) Compared with cross-entropy loss, Dice loss and tracing loss can generate crisper edge maps, but the improvement for AC is still less than only training with our refined labels. (b) Dice loss and tracing loss can be perfectly integrated with our method and further improve AC in all cases. (c) Dice loss and tracing loss can enjoy a free performance boost after integrating our method in most cases. Especially, Dice loss training with refined labels can achieve the best performance, since dice loss emphasizes the image-level similarity between two sets of edges, which can balance the discontinuities in refined labels. (d) After integrating our method with Dice loss and tracing loss, the predicted edge maps are almost crisp enough that the performance change after and before NMS is very slight ($<1\%$ in all cases). Such a phenomenon shows the potential to directly adopt the predicted results without NMS for downstream tasks, making the whole pipeline end-to-end.

| Dataset | Methods | after NMS | | before NMS | | AC |
|---|---|---|---|---|---|---|
| | | ODS | OIS | ODS | OIS | |
| BSDS 500 | PiDiNet-D | **0.782** | **0.798** | **0.774** | 0.788 | 0.306 |
| | PiDiNet-T | 0.78 | 0.796 | 0.77 | 0.787 | 0.333 |
| | PiDiNet-R | 0.771 | 0.782 | 0.771 | 0.78 | 0.424 |
| | PiDiNet-D-R | 0.779 | 0.794 | 0.773 | 0.787 | **0.546** |
| | PiDiNet-T-R | 0.773 | 0.79 | 0.772 | **0.789** | 0.545 |
| Multicue | PiDiNet-D | 0.871 | 0.876 | 0.77 | 0.779 | 0.208 |
| | PiDiNet-T | 0.866 | 0.871 | 0.742 | 0.754 | 0.217 |
| | PiDiNet-R | 0.892 | 0.907 | 0.886 | 0.904 | 0.424 |
| | PiDiNet-D-R | **0.899** | **0.912** | **0.894** | **0.909** | 0.463 |
| | PiDiNet-T-R | 0.884 | 0.903 | 0.876 | 0.894 | **0.559** |
| BIPED | PiDiNet-D | **0.873** | 0.877 | 0.852 | 0.858 | 0.34 |
| | PiDiNet-T | 0.872 | 0.878 | 0.845 | 0.852 | 0.296 |
| | PiDiNet-R | 0.863 | 0.873 | 0.863 | 0.873 | 0.512 |
| | PiDiNet-D-R | 0.869 | **0.879** | **0.867** | **0.877** | 0.696 |
| | PiDiNet-T-R | 0.863 | 0.874 | 0.862 | 0.873 | **0.704** |

TABLE V: Quantitative results for PiDiNet when training with original labels and refined labels based on dice loss and tracing loss. "-D" and "-T" means training with dice loss and tracing loss. "-D-R" and "-T-R" means training with our refined labels with dice loss and tracing loss, respectively.

## D. Crisp Edge Detection for Other Vision Tasks

*1) Optical Flow:* Optical flow estimation is the task of estimating per-pixel motion cues between consecutive video frames. The main difficulties include large displacements by fast-moving objects, occlusions, motion blur, and textureless surfaces. Based on the observation that motion boundaries often tend to appear at image edges, EpicFlow [3] computes sparse matches from DeepMatching [43] and leverages detected edges to perform sparse-to-dense interpolation relying on the edge-aware geodesic distance. The obtained dense

correspondences thus are robust to large displacements and motion discontinuities. Due to the nature of the EpicFlow algorithm, the quality of edge detection is important to the final results, where crisper edges are much cleaner and bring less ambiguity. To show the advantage of crisp edge detection, we replace the edge detection part in EpicFlow with PiDiNet trained on the BSDS500 dataset with original labels and refined ones, respectively. Detailed settings can be seen in Section V-A.

We adopt two widely-used optical flow datasets, Sintel [44] and FlyingChairs [45] in this paper. For the Sintel dataset, we conduct experiments on the training set of the "final" version that features realistic rendering effects such as motion, defocus blur and atmospheric effects. For the FlyingChairs dataset, since the whole dataset is relatively large, we only take the first 1k samples, which is enough for evaluations. As in previous works, we adopt Averaged Endpoint Error (AEE) as the evaluation metric (lower is better). Experimental results are summarized in Table VI, where PiDiNet trained with our refined labels achieves better AEE performance than the original thick one. We also present some qualitative examples of the original and crisp edges together with their optical flow estimations in Figure 12 and 13. Without fine-tuning the parameters under the same baseline of EpicFlow [3], we can observe that crisp edge maps reduce ambiguities and bring more accurate and clear results near object boundaries. Both qualitative and quantitative results show the advantage of crisp edge detection for optical flow estimation.

| | Sintel | | FlyingChairs | |
|---|---|---|---|---|
| Methods | AEE | AC | AEE | AC |
| Baseline EpicFlow [3] | 3.686 | \ | \ | \ |
| PiDiNet-EpicFlow | 3.627 | 0.254 | 2.734 | 0.264 |
| PiDiNet-R-EpicFlow | **3.602** | **0.445** | **2.718** | **0.448** |

TABLE VI: Optical flow estimation performance achieved with PiDiNet edge detection integrated by the EpicFlow [3] on Sintel and FlyingChairs dataset. "-R" means training with refined labels.

*2) Semantic Segmentation:* Semantic segmentation is an important high-level vision task that aims at densely classifying each pixel of an image with the semantic category. It achieves a significant performance boost with the help of fully convolutional networks [46]. However, the successive down-sampling layers generate low-resolution feature maps and fail to capture precise boundaries of objects. To address this issue, Bertasius et al. [4] proposed Boundary Neural Field (BNF) to integrate the initial predictions with edge cues by minimizing the global boundary-based energy. Therefore, an accurate and crisp edge map is important for BNF. To show the effectiveness of crisp edge detection, we adopt DeepLabV3+ [47] with various backbones (MobileNet [48] and ResNet50 [49]) to generate the initial segmentation results, then apply our edges maps and BNF to further improve the predictions. Edge maps are generated by PiDiNet trained on original and refined labels the same as V-D1.

We conduct experiments on the Pascal Context validation set [50], and report three widely-adopted segmentation evaluation metrics, pixel accuracy (PA), mean pixel accuracy (MPA) and mean intersection over union (mIOU). Quantitative results
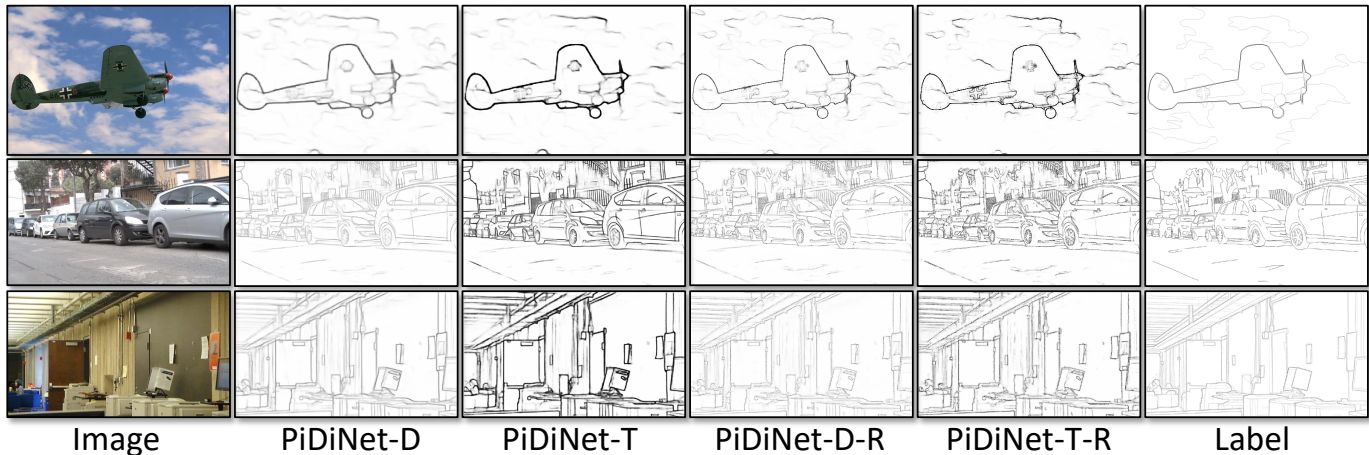
Fig. 11: Qualitative results on datasets of BSDS500 (the first row), BIPED (the second row) and Multicue (the last row) before NMS. Although dice loss ("-D") and tracing loss ("-T") can generate crisp edges, training with our refined labels ("-R") can bring a free performance boost and make edge maps even crisper.
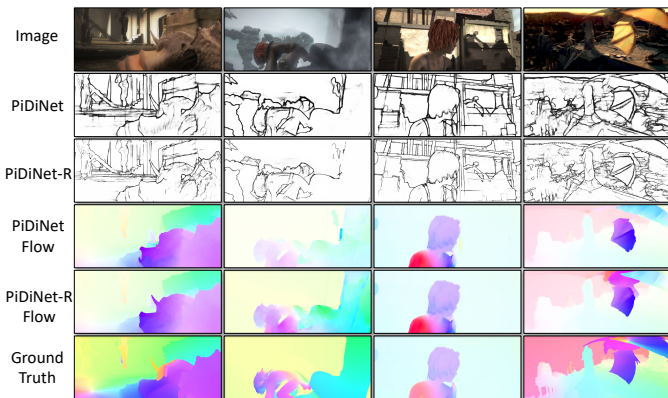


Fig. 12: Illustration of optical flow estimation examples on Sintel dataset. Crisp edge maps bring better localization near motion discontinuities.
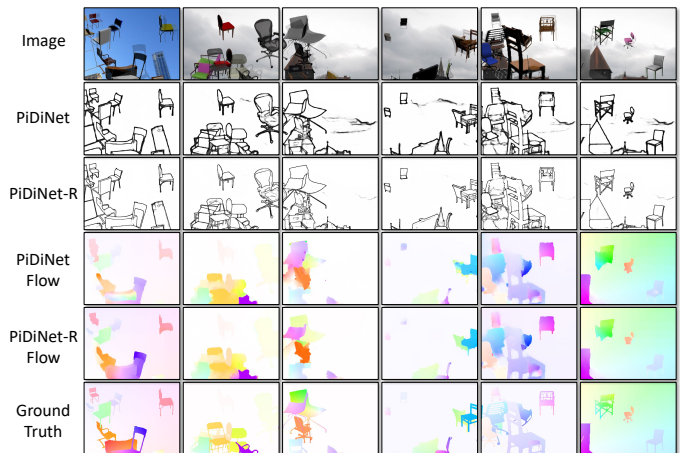


Fig. 13: Illustration of optical flow estimation examples on FlyingChairs dataset. Crisp edge maps yield more accurate results with the same baseline algorithm.

are summarized in Table VII, where PiDiNet trained with refined labels achieves the best performance on all metrics with all backbones, which proves the advantages of crisp edge detection. We illustrate some qualitative examples of the original and crisp edges and their refined segmentation results in Figure 14. These results illustrate that semantic segmentation can benefit from better-localized edge maps, and be further improved by crisper ones that are better at capturing the precise contour of semantic objects.

| Method | PA | MPA | mIOU | AC |
|---|---|---|---|---|
| DeepLabV3+ MobileNet [48] | 0.918 | 0.779 | 0.666 | \ |
| PiDiNet-BNF | 0.922 | 0.790 | 0.680 | 0.252 |
| PiDiNet-R-BNF | **0.923** | **0.803** | **0.684** | **0.444** |
| DeepLabV3+ ResNet50 [49] | 0.939 | 0.857 | 0.744 | \ |
| PiDiNet-BNF | 0.939 | 0.862 | 0.747 | 0.252 |
| PiDiNet-R-BNF | **0.940** | **0.871** | **0.749** | **0.444** |

TABLE VII: Semantic segmentation performance achieved with PiDiNet edge detection integrated by BNF [4] on Pascal Context validation set. "-R" means training with refined labels. We report the initial predictions and their refined results on two backbones of MobileNet and ResNet50.

## VI. CONCLUSIONS AND LIMITATIONS

In this paper, we explore the reason of predicting "thick" edges by deep learning edge detectors. We introduce the first crispness metric for edge detection and explore the impact of noisy human-labeled edges. We find that the issue of thick predictions is mainly caused by noisy human labels, and the problem is aggravated by the cross-entropy loss. Based on the observation, we propose an iterative Canny-guided refinement method to refine human labels. Our pipeline does not introduce any extra modules or loss functions in training and testing, and can be integrated with any edge detection backbones. Comprehensive experiments demonstrate the effectiveness of our method for improving the edge crispness compared to models trained from original labels. We also verify the superiority of crisp edge detection for other downstream vision tasks such as optical flow estimation and semantic segmentation.

*Limitations*: Although training with refined labels can generate crisp edges, the performance of our plug-and-play method still relies on the adopted edge detector and the quality of
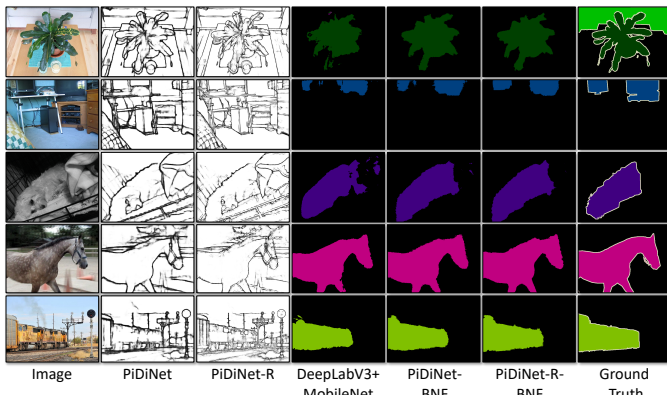
Fig. 14: Illustration of semantic segmentation examples on Pascal Context validation set. Crisp edge maps can bring more complete segmentation and better localization near object boundaries.
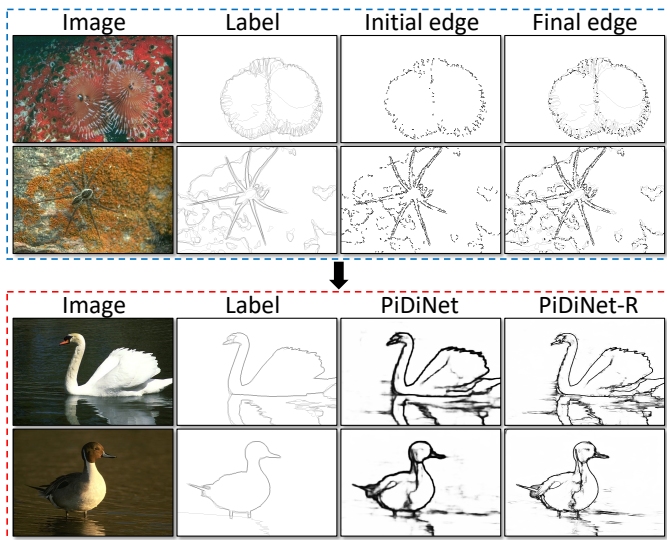


Fig. 15: Examples of some failure cases. The refined training labels and predicted results are inside the blue and red frames, respectively.

human labels. As shown in Figure 15, we present some failure cases of refined labels in the blue frame and prediction results in the red frame.

As demonstrated in the blue frame, failure cases often happen on images with rich texture and unclear boundaries, which is also challenging for human annotators. In such cases, the produced initial edges will be so sparse that the discontinuities are hard to be inpainted without enough semantic context, leading to unsatisfactory final edges.

Consequently, as illustrated in the red frame, due to the strong generalization ability of learning-based edge detection methods, the prediction results may still contain some salient noises (e.g. water reflection) that are not semantically meaningful. Our method may amplify such existing salient noises more depending on the performance of the edge detection backbone and its original predictions. To tackle this problem, training a better edge inpainting network can be an interesting direction for future work.

REFERENCES

[1] C. L. Zitnick and P. Dollár, "Edge boxes: Locating object proposals from edges," in *European conference on computer vision*. Springer, 2014, pp. 391–405.

[2] P. Arbeláez, J. Pont-Tuset, J. T. Barron, F. Marques, and J. Malik, "Multiscale combinatorial grouping," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2014, pp. 328–335.

[3] J. Revaud, P. Weinzaepfel, Z. Harchaoui, and C. Schmid, "Epicflow: Edge-preserving interpolation of correspondences for optical flow," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2015, pp. 1164–1172.

[4] G. Bertasius, J. Shi, and L. Torresani, "Semantic segmentation with boundary neural fields," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 3602–3610.

[5] T. Cheng, X. Wang, L. Huang, and W. Liu, "Boundary-preserving mask r-cnn," in *European conference on computer vision*. Springer, 2020, pp. 660–676.

[6] W. Xiong, J. Yu, Z. Lin, J. Yang, X. Lu, C. Barnes, and J. Luo, "Foreground-aware image inpainting," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 5840–5848.

[7] K. Nazeri, E. Ng, T. Joseph, F. Z. Qureshi, and M. Ebrahimi, "Edgeconnect: Generative image inpainting with adversarial edge learning," *arXiv preprint arXiv:1901.00212*, 2019.

[8] J. Kittler, "On the accuracy of the sobel edge detector," *Image and Vision Computing*, vol. 1, no. 1, pp. 37–42, 1983.

[9] J. Canny, "A computational approach to edge detection," *IEEE Transactions on pattern analysis and machine intelligence*, no. 6, pp. 679–698, 1986.

[10] S. Xie and Z. Tu, "Holistically-nested edge detection," in *Proceedings of the IEEE international conference on computer vision*, 2015, pp. 1395–1403.

[11] Y. Liu, M.-M. Cheng, X. Hu, K. Wang, and X. Bai, "Richer convolutional features for edge detection," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 3000–3009.

[12] J. He, S. Zhang, M. Yang, Y. Shan, and T. Huang, "Bi-directional cascade network for perceptual edge detection," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 3828–3837.

[13] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," *arXiv preprint arXiv:1409.1556*, 2014.

[14] Y. Wang, X. Zhao, and K. Huang, "Deep crisp boundaries," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 3892–3900.

[15] R. Deng, C. Shen, S. Liu, H. Wang, and X. Liu, "Learning to predict crisp boundaries," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018, pp. 562–578.

[16] Z. Yu, W. Liu, Y. Zou, C. Feng, S. Ramalingam, B. Kumar, and J. Kautz, "Simultaneous edge alignment and learning," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018, pp. 388–404.

[17] D. Acuna, A. Kar, and S. Fidler, "Devil is in the edges: Learning semantic boundaries from noisy annotations," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 11 075–11 083.

[18] P. Arbelaez, M. Maire, C. Fowlkes, and J. Malik, "Contour detection and hierarchical image segmentation," *IEEE transactions on pattern analysis and machine intelligence*, vol. 33, no. 5, pp. 898–916, 2010.

[19] D. A. Mély, J. Kim, M. McGill, Y. Guo, and T. Serre, "A systematic comparison between visual cues for boundary detection," *Vision research*, vol. 120, pp. 93–107, 2016.

[20] X. S. Poma, E. Riba, and A. Sappa, "Dense extreme inception network: Towards a robust cnn model for edge detection," in *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, 2020, pp. 1923–1932.

[21] Z. Su, W. Liu, Z. Yu, D. Hu, Q. Liao, Q. Tian, M. Pietikäinen, and L. Liu, "Pixel difference networks for efficient edge detection," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 5117–5127.

[22] P. Dollár and C. L. Zitnick, "Fast edge detection using structured forests," *IEEE transactions on pattern analysis and machine intelligence*, vol. 37, no. 8, pp. 1558–1570, 2014.

[23] J. K. Wibisono and H.-M. Hang, "Fined: Fast inference network for edge detection," *arXiv preprint arXiv:2012.08392*, 2020.

[24] L. Huan, N. Xue, X. Zheng, W. He, J. Gong, and G.-S. Xia, "Unmixing convolutional features for crisp edge detection," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2021.

[25] D. Pathak, P. Krahenbuhl, J. Donahue, T. Darrell, and A. A. Efros, "Context encoders: Feature learning by inpainting," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 2536–2544.

[26] S. Iizuka, E. Simo-Serra, and H. Ishikawa, "Globally and locally consistent image completion," *ACM Transactions on Graphics (ToG)*, vol. 36, no. 4, pp. 1–14, 2017.

[27] F. Yu and V. Koltun, "Multi-scale context aggregation by dilated convolutions," *arXiv preprint arXiv:1511.07122*, 2015.

[28] Y. Ren, X. Yu, R. Zhang, T. H. Li, S. Liu, and G. Li, "Structureflow: Image inpainting via structure-aware appearance flow," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2019, pp. 181–190.

[29] M.-L. Shih, S.-Y. Su, J. Kopf, and J.-B. Huang, "3d photography using context-aware layered depth inpainting," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 8028–8038.

[30] Y.-K. Huang, T.-H. Wu, Y.-C. Liu, and W. H. Hsu, "Indoor depth completion with boundary consistency and self-attention," in *Proceedings of the IEEE/CVF International Conference on Computer Vision Workshops*, 2019, pp. 0–0.

[31] S. Zhu, G. Brazil, and X. Liu, "The edge of depth: Explicit constraints between segmentation and depth," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 13 116–13 125.

[32] X. Song, X. Zhao, L. Fang, H. Hu, and Y. Yu, "Edgestereo: An effective multi-task learning network for stereo matching and edge detection," *International Journal of Computer Vision*, vol. 128, no. 4, pp. 910–930, 2020.

[33] A. Kirillov, E. Levinkov, B. Andres, B. Savchynskyy, and C. Rother, "Instancecut: from edges to instances with multicut," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 5008–5017.

[34] H. He, X. Li, G. Cheng, J. Shi, Y. Tong, G. Meng, V. Prinet, and L. Weng, "Enhanced boundary learning for glass-like object segmentation," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 15 859–15 868.

[35] J. Yang, B. Price, S. Cohen, H. Lee, and M.-H. Yang, "Object contour detection with a fully convolutional encoder-decoder network," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 193–202.

[36] P. Y. Simard, D. Steinkraus, J. C. Platt *et al.*, "Best practices for convolutional neural networks applied to visual document analysis." in *Icdar*, vol. 3, no. 2003. Edinburgh, 2003.

[37] P. T. Jackson, A. A. Abarghouei, S. Bonner, T. P. Breckon, and B. Obara, "Style augmentation: data augmentation via style randomization." in *CVPR workshops*, vol. 6, 2019, pp. 10–11.

[38] S. Zhao, Z. Liu, J. Lin, J.-Y. Zhu, and S. Han, "Differentiable augmentation for data-efficient gan training," *Advances in Neural Information Processing Systems*, vol. 33, pp. 7559–7570, 2020.

[39] G. Liu, F. A. Reda, K. J. Shih, T.-C. Wang, A. Tao, and B. Catanzaro, "Image inpainting for irregular holes using partial convolutions," in *Proceedings of the European conference on computer vision (ECCV)*, 2018, pp. 85–100.

[40] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimelshein, L. Antiga *et al.*, "Pytorch: An imperative style, high-performance deep learning library," *Advances in neural information processing systems*, vol. 32, 2019.

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

[42] B. Zhou, A. Lapedriza, A. Khosla, A. Oliva, and A. Torralba, "Places: A 10 million image database for scene recognition," *IEEE transactions on pattern analysis and machine intelligence*, vol. 40, no. 6, pp. 1452–1464, 2017.

[43] P. Weinzaepfel, J. Revaud, Z. Harchaoui, and C. Schmid, "Deepflow: Large displacement optical flow with deep matching," in *Proceedings of the IEEE international conference on computer vision*, 2013, pp. 1385–1392.

[44] D. J. Butler, J. Wulff, G. B. Stanley, and M. J. Black, "A naturalistic open source movie for optical flow evaluation," in *European conference on computer vision*. Springer, 2012, pp. 611–625.

[45] A. Dosovitskiy, P. Fischer, E. Ilg, P. Hausser, C. Hazirbas, V. Golkov, P. Van Der Smagt, D. Cremers, and T. Brox, "Flownet: Learning optical flow with convolutional networks," in *Proceedings of the IEEE international conference on computer vision*, 2015, pp. 2758–2766.

[46] J. Long, E. Shelhamer, and T. Darrell, "Fully convolutional networks for semantic segmentation," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2015, pp. 3431–3440.

[47] L.-C. Chen, Y. Zhu, G. Papandreou, F. Schroff, and H. Adam, "Encoder-decoder with atrous separable convolution for semantic image segmentation," in *Proceedings of the European conference on computer vision (ECCV)*, 2018, pp. 801–818.

[48] A. G. Howard, M. Zhu, B. Chen, D. Kalenichenko, W. Wang, T. Weyand, M. Andreetto, and H. Adam, "Mobilenets: Efficient convolutional neural networks for mobile vision applications," *arXiv preprint arXiv:1704.04861*, 2017.

[49] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 770–778.

[50] 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.