

Photo-Sketching: Inferring Contour Drawings from Images

Mengtian Li¹ Zhe Lin² Radomír Měch² Ersin Yumer³ Deva Ramanan^{1,4}
¹Carnegie Mellon University ²Adobe Research ³Uber ATG ⁴Argo AI



Figure 1: Automatic contour drawing generation for images in the wild. Different from traditional edge or boundary detectors [48], our method predicts the most salient contours in images and reflects imperfections in ground truth human drawing, e.g. the ceiling is not perfectly straight in the right example as a novice drawer would draw it. Right photo by [ostap25 – stock.adobe.com](https://www.istock.com/ostap25).

Abstract

Edges, boundaries and contours are important subjects of study in both computer graphics and computer vision. On one hand, they are the 2D elements that convey 3D shapes, on the other hand, they are indicative of occlusion events and thus separation of objects or semantic concepts. In this paper, we aim to generate contour drawings, boundary-like drawings that capture the outline of the visual scene. Prior art often cast this problem as boundary detection. However, the set of visual cues presented in the boundary detection output are different from the ones in contour drawings, and also the artistic style is ignored. We address these issues by collecting a new dataset of contour drawings and proposing a learning-based method that resolves diversity in the annotation and, unlike boundary detectors, can work with imperfect alignment of the annotation and the actual ground truth. Our method surpasses previous methods quantitatively and qualitatively. Surprisingly, when our model fine-tunes on BSDS500, we achieve the state-of-the-art performance in salient boundary detection, suggesting contour drawing might be a scalable alternative to boundary annotation, which at the same time is easier and more interesting for annotators to draw.

1. Introduction

Edge-like visual representation, appearing in form of image edges, object boundaries, line drawings and pictorial scripts, is of great research interest in both computer vision and computer graphics. Automatic generation of such representation enables us to understand the geometry of the scene [40], and perform image manipulation in this sparse space [11]. This paper studies such representation in the form of *contour drawing*, which contains object boundaries, salient inner edges such as occluding contours, and salient background edges. These sets of visual cues convey 3D perspective, length and width as well as thickness and depth [45]. Contour drawings are usually based on real-world objects (immediately observed or from memory), and therefore, can be considered as an expression of human vision. Its counterpart in machine vision is edge and boundary detection. Interesting, the set of visual cues is different in contour drawings and in image boundaries (Fig 2). Comparing to image boundaries, contour drawings tend to have more details inside each object (including occluding contours and semantically-salient features such as eyes, mouths, etc.) and are made of strokes that are loosely aligned to pixels on the image edges. We propose a contour generation algorithm to output contour drawings given input images. This generation process involves identifying salient boundaries and is connected with the salient boundary detection in computer

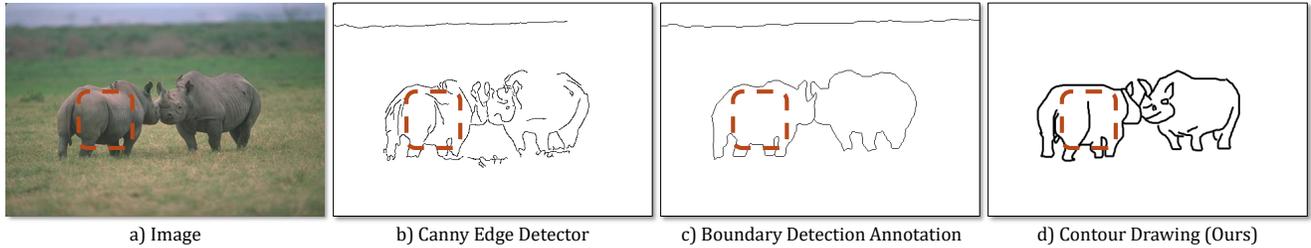


Figure 2: a) Which visual cues would you draw when sketching out an image? b) Traditional edge detectors [6] only capture high frequency signals in the image without image understanding. c) Boundary detectors are usually trained on edges that derived from closed segment annotations and therefore, they do not include salient inner boundaries by definition [37, 2]. d) In contrast, our contour drawing (ground truth is shown here) contains both the occluding contours and salient inner edges. For example, the dashed box in the top row contains a open contour ending in a *cusp* [28, 18].

vision. In fact, we will show that our contour generation algorithm can be re-purposed to perform salient boundary detection and achieve the best performance on standard benchmark. Another element involved contour drawing generation is to adopt proper artistic style. Fig 1 shows our method successfully captures the style and itself is a style transfer application. Moreover, contour drawing is an intermediate representation between image boundary and abstract line drawing. Our study of contour drawing paves the road for machine’s understanding and generation of abstract line drawings [15, 9].

What types of edge-like visual representation are studied in existing work? In non-photorealistic rendering, 2D lines that convey 3D shapes are widely studied. The most important of such might be the *occluding contours*, regions where the local surface normal is perpendicular to the viewing direction, and *creases*, edges along which the dihedral angle is small [28]. It is noted by DeCarlo *et al.* [10] that important details are missing if only those edges are rendered, and their solution is to add *suggestive contours*, regions where occluding contour would appear with minimal change in viewpoint. As a result of having clear mathematical definition, these edge-like representation can be directly computed using methods in differential geometry given the 3D model. In computer vision, a different set of visual cues are defined and are inferred from the image alone without knowledge of the 3D world, namely the *image edges* and *boundaries*. Image edges correspond to sharp changes in image intensity due to changes in albedo, surface orientation, or illumination [18]. Boundaries, as formally defined by Martin *et al.* [36], are contours in the image plane that represents a change in pixel ownership from one object or surface to another. Nonetheless, this definition ignores the fact that the contour can also appear on a smooth surface of the same object, for example the *cusp* in Fig 2. Since much progress has been driven by datasets, in practise, the boundaries are “defined” by the seminal benchmark of BSDS300 [36] and BSDS500 [2]. Interestingly, despite their popular-

ity, these datasets were originally designed and annotated to be a segmentation dataset. This means that boundaries are derived from closed segments annotated by humans [37], and yet not all boundaries form closed shapes.

The other related line of research revolves around the representation of sketch. Most work study the relationship between the strokes themselves without a reference object or image [20, 15, 42]. While some work on sketch-based image retrieval [41, 17] and sketch generation [20, 44] indeed models the correspondence between sketch and image, the type of drawings used are far too simple or abstract, and does not contain edge-level correspondence, making it unsuitable for training a generic scene sketch generator. A comparison is summarized in Fig 3 and Tab 1.

To accommodate our research on contour drawings, we collect a dataset containing 5000 drawings (Sec 2). The challenge for training a contour generator is to resolve the diversity among the contours for the same image obtained from multiple annotators. We address it by proposing a novel loss that allows the network to converge to an implicit consensus, while retaining details (Sec 3). Our contour generator can be applied to salient boundary detection. By simply fine-tuning on BSDS500, we achieve the state-of-the-art performance (Sec 4). Finally, we show our dataset can be expanded in a cost free way with a sketch game (Sec 5). Our code and dataset are available online ¹.

2. Collecting Contour Sketches

We create our novel task with the the popular crowdsourcing platform Amazon Mechanical Turk [5]. To collect drawings that are roughly boundary aligned, we allow the Turkers to trace over a faded background image. In order to obtain high-quality drawings, we design a labeling interface with a detailed instruction page including many positive and negative examples. The quality control is realized through manual inspection by treating drawings of

¹<http://www.cs.cmu.edu/~mengtial/proj/sketch>

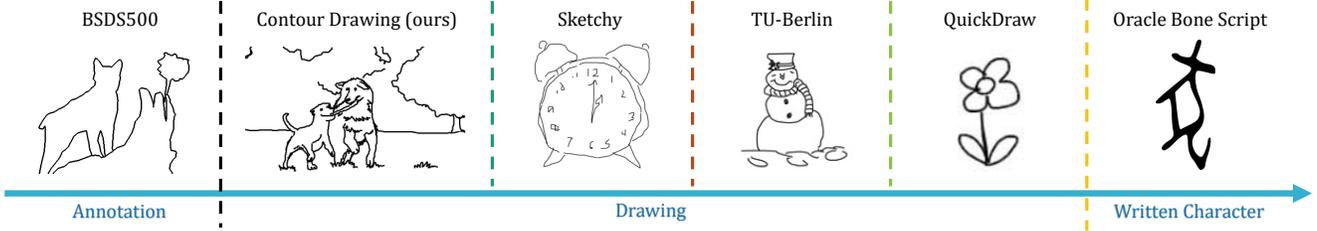


Figure 3: Comparison with other edge-like representations. Here we order examples [2, 41, 16, 20] by their level of abstraction. Our contour drawing cover more detailed internal boundaries than a boundary annotation, while having much better alignment to actual image contours, and much more complexity compared to other drawing-based representations.

Dataset	Edge Aligned	Multiple Obj	With Image	Vec Graphics	Stroke Order
BSDS500 [2]	✓	✓	✓	✗	✗
Contour Drawing (ours)	Roughly	✓	✓	✓	✓
Sketchy [41]	✗	✗	✓	✓	✓
TU-Berlin [16]	✗	✗	✗	✓	✓
QuickDraw [20]	✗	✗	✗	✓	✓

Table 1: Dataset comparison. Our proposed contour drawing dataset is different from prior work in terms of boundary alignment, multiple objects, corresponding image-sketch pairs, vector graphics encoding, and stroke ordering annotation.

the following types as rejection candidates: (1) missing inner boundary, (2) missing important objects, (3) with large misalignment with original edges, (4) the content not recognizable, (5) drawing humans with stick figures, (6) shaded on empty areas.

Finally, we collect 5000 high-quality drawings on a dataset of 1000 outdoor images crawled from Adobe Stock [1] and each image is paired with exactly 5 drawings. In addition, we have 1947 rejected submissions, which will be used in setting up an automatic quality guard as discussed in Sec 5.

3. Sketch Generation

In this section, we propose a new deep learning-based model to generate contour sketches from a given image and evaluate it against competing methods in both objective and subjective manner. A unique aspect of our problem here is that each training image is associated with multiple ground truth sketches drawn by different annotators.

3.1. Previous Methods

Early methods of line-drawing generation focus on human faces, where they build explicit models to represent facial features [4, 7]. Other work focuses on generating the style but leaving the task of deciding which edge to draw to the user [26, 47]. More recently, Song *et al.* [20] used LSTM to sequentially generate the stroke for simple doodles of several strokes. However, our contour drawings on average contain 44 strokes and around 5,000 control points,

way beyond the capacity of existing sequential models.

3.2. Our Method

Naturally, the problem of generating contour drawing can be cast into an image translation problem or a classical boundary detection problem. Given the popularity of using conditional Generative Adversarial Networks (cGANs) to generate images from sketches or boundary maps, one might think the apparently easier inverse problem can be solved by reversing the image generation direction. However, none of the existing cGAN methods [23, 49, 33, 50] have shown results on such a task and our experiments show that they do not work on sketch generation out-of-the-box. We conjecture that the drawings are sparse and discrete representation compared to textured images. It might be easier to obtain gradients in the latter case. Also, our dataset has more than one target for each source image (1-to-many mapping). And modeling such diversity makes it difficult to optimize. On the other hand, classical boundary detection approaches linearly combines the different ground truths to form a single target per each input. This form of data augmentation bypasses the need to model the diverse outputs and results in a soft output as well, but it is not the case with multiple ground truth having edges not perfectly aligned. The soft representation no longer bears the meaning of boundary strength, but how well the edges are accidentally matched. Training on such data yields unreasonable output for both our method and existing boundary detection methods. Hence, our problem cannot be trivially

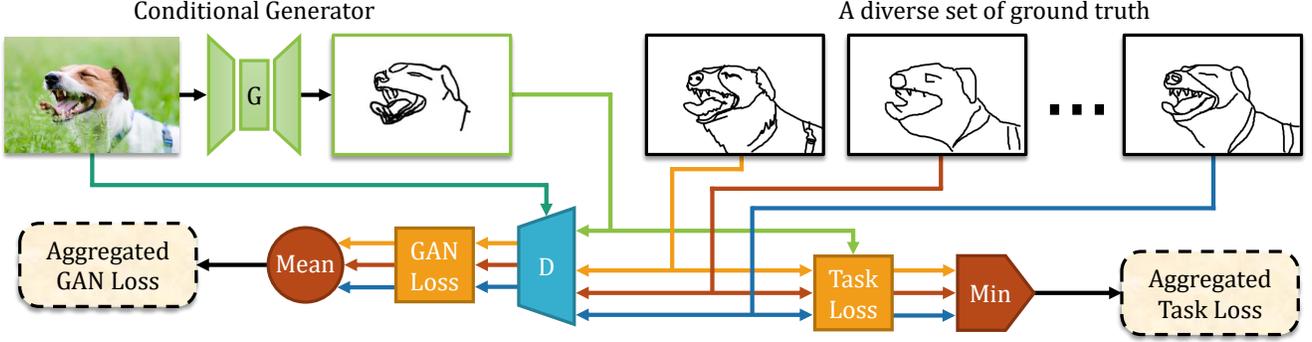


Figure 4: We train an image-conditioned contour generator with a novel MM-loss (Min-Mean-loss) that accounts for multiple diverse outputs encountered during training (top row). Training directly on the entire set of image-contour pairs generates conflicting gradients. To rectify this, we carefully aggregate the discriminator “GAN” loss and the regression “Task” loss. The discriminator averages the GAN loss across all image-contour pairs, while the regression-loss finds the minimum-cost contour to pair with this image (determined on-the-fly during learning). This ensures that the generator will not simply regress to the “mean” contour, which might be invalid. Photo by alexei_tm – [stock.adobe.com](https://www.adobe.com/stock).

solved by training or fine-tuning boundary detectors on the contour drawing dataset. Another issue with soft representation, as found in [22], is their poor correlation with the actual boundary strength. We share the same findings in our experiments, as it is difficult to find a single threshold for the final output that works well for all images. In this work, we use a different cGAN with a novel MM-loss (Fig 4).

3.2.1 Formulation

We leverage the power of the recently popular framework of adversarial training. In a Generative Adversarial Network (GAN), a random noise vector z is fed into the generation network G to generate an output image y . In the conditional setup (cGAN), the generator takes input an image x , and together with a z , it maps to a y . The generator G aims to generate “real” images conditioned on x , while there is another discriminator network D that is adversarially trained to tell the generated images from the actual ground truth target. Mathematically, the loss for such objective can be written as

$$\mathcal{L}_{cGAN}(x, y, z) = \min_G \max_D \mathbb{E}_{x, y} [\log D(x, y)] + \mathbb{E}_{x, z} [\log(1 - D(x, G(x, z)))] \quad (1)$$

As found by previous work [38, 23], the noise vector z is usually ignored in the optimization. Therefore, we do not include z in our experiments. We also followed the common approach in cGAN to include a task loss in addition to the GAN loss. This is reasonable since we have a target ground truth for us to compare with directly. For our contour generation task, we set the task loss to be L1 loss which

encourages sparsity required for contour outputs. The combined loss function now becomes

$$\mathcal{L}_c(x, y) = \lambda \mathcal{L}_{cGAN}(x, y) + \mathcal{L}_1(x, y), \quad (2)$$

where the non-negative constant λ adjusts the relative strength of the two objectives. Note that when $\lambda = 0$, the model reduces to a simple regression.

The above formulation assumes a 1-to-1 mapping between the two domains. However, we have multiple different targets $y_i^{(1)}, y_i^{(2)}, \dots, y_i^{(M_i)}$ for a same input x_i , making it a 1-to-many mapped problem. Note the number of targets M_i for each input may vary from examples to examples. If we ignore the fact of 1-to-many mapping, this is reduced to a regular 1-to-1 mapping problem: $(x_1, y_1^{(1)}), \dots, (x_1, y_1^{(M_1)}), \dots, (x_N, y_N^{(1)}), \dots, (x_N, y_N^{(M_N)})$, and those pairs are fetched in random order to train the network.

Our method treats $(x_i, y_i^{(1)}, \dots, y_i^{(M_i)})$ as a single training example. To accommodate the extra targets in each training example, we propose a novel MM-loss (Min-Mean-loss) (Fig 4). Two different aggregate functions are used for the generator G and the discriminator D respectively. The final loss for each training example becomes

$$\mathcal{L}(x_i, y_i^{(1)}, \dots, y_i^{(M_i)}) = \frac{\lambda}{M_i} \sum_{j=1}^{M_i} \mathcal{L}_{cGAN}(x_i, y_i^{(j)}) + \min_{j \in \{1, \dots, M_i\}} \mathcal{L}_1(x_i, y_i^{(j)}), \quad (3)$$

The “mean” aggregate function asks the discriminator to learn from all modalities in the target domain and treat those

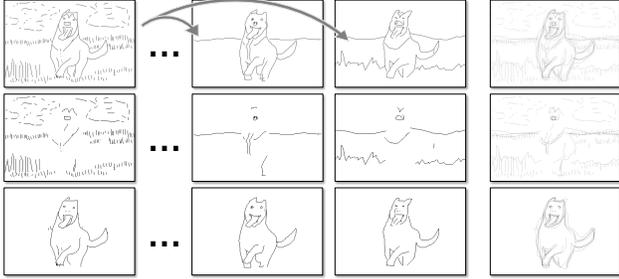


Figure 5: Finding consensus from diverse drawings. In row 1, we visualize 3 different ground truth drawings corresponding to the same image, followed by their overlay in the fourth column. We match strokes in one drawing to another, removing those strokes that could not be matched (row 2). The leftover matched strokes (row 3) are used for evaluation. Note our novel loss allows us to train on the original drawing (row 1) directly and it outperforms the results training on the consensus (row 3), as shown in Tab 2 second last row.

modalities with equal importance. The “min” aggregate function allows the generator to adaptively pick the most suitable modality to generate on-the-fly. Therefore, the problem of conflicting gradients caused by different modalities is greatly alleviated. In the diagnostic experiments (Tab 2), we find that training on the consensus drawing outperforms the baseline method, while training on the complete set of sketches with MM-loss outperforms training just on consensus. The “min” aggregation function might be reminiscent of the stochastic multiple-choice loss [30] which relies on a single target output but learns multiple network output branches to generate a diverse output. In our setting, we have a single stochastic output, but multiple ground-truth targets, and a part of the network (the discriminator) that still uses the set of all ground truths to back-propagate the gradients.

We use the standard encoder-decoder architecture (ResNet [21] based) that yields good performance on style translation tasks [25]. Unlike other pixel generation tasks, we find the skip connections between the encoder and the decoder make the performance drop. The reason might be that our targets contain mainly object boundaries instead of edges in textures, and removing the skip connections suppresses this low-level information. In many pixel-level prediction tasks, the skip connections are added to make pixel accurate predictions. However, we find that pixel accuracy is already encoded in the network itself since our output is sparse. This can be evidenced by the pixel accurate predictions of our same model applied to boundary detection (Sec 4). For the discriminator, we used a regular global GAN as opposed to PatchGAN [24] in related work. Although PatchGAN helps other networks to generate nice textures,

Method	F1-score	Precision	Recall
pix2pix [24] (baseline)	0.514	0.585	0.458
+ ResNet generator	0.561	0.620	0.512
+ our MM-loss	0.765	0.814	0.722
+ GlobalGAN	0.773	0.720	0.835
+ augmentation (final)	0.826	0.861	0.794
- train on consensus	0.802	0.915	0.714
- remove GAN loss	0.778	0.889	0.692

Table 2: Ablation study of our method on the validation set. The metrics are explained in Sec 3.3. We built up our model from a baseline method [24] and the final model uses the ResNet generator without skip connection, a global discriminator and our proposed MM-loss. Moreover, despite the inconsistency in the non-consensus strokes, training on the original drawings outperforms training on just consensus strokes (second last row). We conjecture that our MM-loss can resolve conflicting supervision on the fly. Also, the last row shows that by adopting adversarial training, we outperform pure regression.

it discourages the network to “think” globally, resulting in many broken edges for a single contour of the object. This problem is alleviated when using the global GAN. An ablation study is provided in Tab 2 with evaluation metric explained in the next subsection.

3.3. Evaluation

Quantitative Evaluation Boundary detection has a well-established evaluation protocol that matches predicted pixels to the ground truth under a given offset tolerance [36, 2]. Matching is done with min-cost bipartite assignment [19, 8]. To apply this approach to contour generation, we first need to reconcile the diverse drawing styles in the ground truth set. Hou *et al.* [22] propose a *consensus matching* evaluation of boundary detection that refines the ground-truth by matching pixels from one human annotation to another, removing those that are not unanimously matched across all annotators. We follow suit, but match at stroke level to ensure that strokes are not broken up in the final consensus drawing (Fig 5). In addition, since contour drawings are not exactly aligned with the image boundary, we double the standard offset tolerance used for boundary evaluation. The evaluation treats each ground truth pixel as an “object” in the precision-recall framework. We split the set of 1000 images with associated sketches into train-val-test sets of 800-100-100. The results are shown in Tab 3 and Fig 6. Pix2Pix and our method are trained on our dataset while HED is off-the-shelf. As explained earlier, boundary detection methods cannot work with diverse ground truths and imperfect alignment between the annotations and fine-

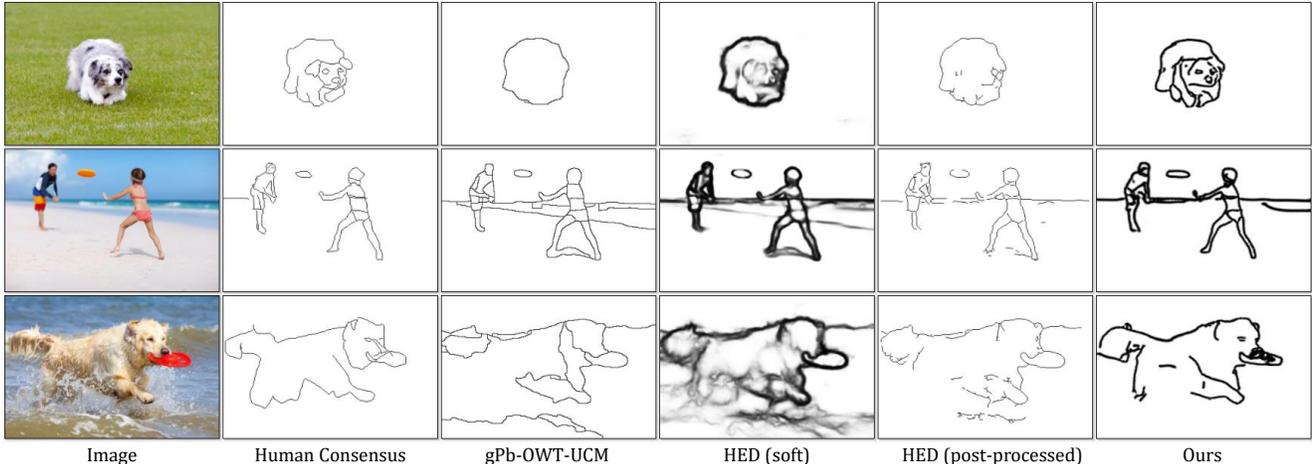


Figure 6: Qualitative results for contour drawing generation. Column 3 and 5 are results at their optimal thresholds for the entire test set (*i.e.* ODS, for readers familiar with BSDS evaluation [2]). Note that it is non-trivial to convert HED [48] soft output (column 4) to clean sketches (column 5) without artifacts (broken edges). The last row shows that our methods fails at places with partial occlusion (water splash). In contrast, all human annotators are not confused by such occlusion. Photos from top to bottom by martincp, BlueOrange Studio, and Phil Stev – stock.adobe.com.

Method	F1	Prec	Rec	Human
pix2pix [24]	0.536	0.625	0.469	4.0%
gPb-OWT-UCM [2]	0.697	0.634	0.774	4.5%
HED [48]	0.782	0.779	0.785	13.0%
Ours	0.822	0.879	0.773	19.5%

Table 3: We evaluate contour drawing generation using standard schemes for evaluating boundary detection, modified to allow for larger pixel offsets during matching. Our method outperforms strong baselines for image translation and boundary detection. In the last column, we measure accuracy with an A/B user study. Our generated drawings are able to fool significantly more human subjects. User studies are consistent with our quantitative evaluation, suggesting that boundary detection accuracy is a reasonably proxy for contour drawing generation.

tuning them on our dataset yields worse performance.

Perceptual Study Besides measuring the F1-score, we also evaluate results with A/B testing in a perceptual study (shown in the last column in Tab 3). For each algorithm, we present AMT Turkers [5] with a tuple consisting of an image, the generated drawing, and a random human drawing for that image for 5 seconds. Turkers are then asked to select the one drawn by a human. Past A/B tests for image generation tend to use a shorter presentation time of 1 second [24, 49], making it easier to fool a user.

In-the-Wild Testing We also test the generalizability of our method on arbitrary Internet images. Note that the re-

sults here are obtained by directly applying the top performing model on the sketch validation set, *without any tuning of the hyperparameters*. The qualitative results in Fig 1 show that our model has learned general representations for salient contours in the images without content bias and incorporates random perturbations present in human drawings. The generalization power to unseen contents suggests that our method can be applied to other tasks, for instance, salient boundary detection, which is discussed in the following section. Moreover, the generalization to arbitrary contents is also crucial for us to design a human-in-the-loop learning scheme for dataset expansion (Sec 5).

4. Application to Boundary Detection

For an algorithm to generate a contour drawing, it needs first to identify salient edges in an image and this implies that our contour sketch can be re-purposed for salient boundary detection.

Previous Methods Edge detectors are precursors to boundary detectors. Those methods [27, 6] are usually filter-based and closely related to intensity gradients in images. Since Martin *et al.* [36] first raised the problem of boundary detection, many efforts [2, 13, 39, 31, 14] have been devoted to building learning methods upon hand-crafted features. Benchmark performance was improved by a series of deep methods [43, 3, 48]. However, most of these methods merge the set of annotations to one before training ignoring annotation inconsistency issue. In addition, adversarial training has not yet been applied to boundary detection.

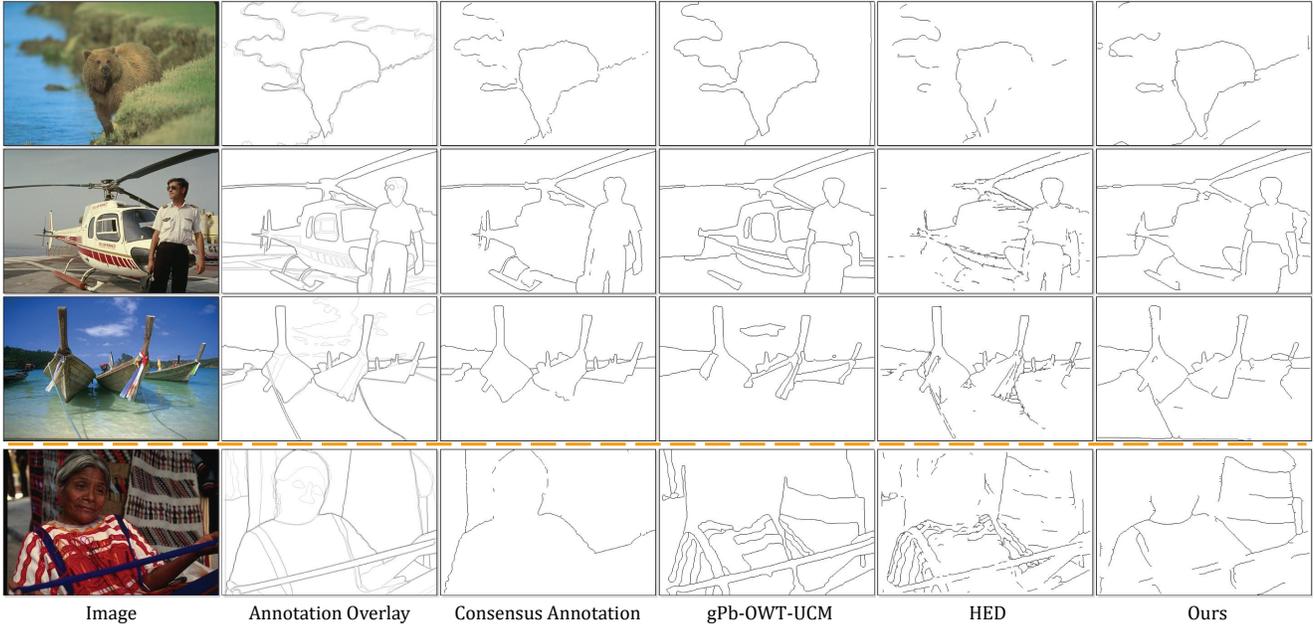


Figure 7: Qualitative results for salient boundary detection. Our method generates complete objects more frequently than competing methods and yet without over generating texture edges. The last row shows a challenging case where all methods fail.

Salient Boundaries Image boundaries may be somewhat ambiguous. This can be seen from the open-ended instructions used to annotate BSDS500 [35, 2] (divide the images into 2 to 20 regions), which often results in inconsistent boundaries labeled by annotators for the same image. Hou *et al.* [22] studied this inconsistency issue in BSDS500 with a series of experiments. They define *orphan* labels to be boundaries labeled by only a single annotator. They then ask human subjects if a particular algorithm’s false alarm (predicted boundary pixel that is not matched to any ground-truth) is stronger than a randomly-selected orphan. Subjects selected the false alarm 50.2% of the time. This seems worrisome as nearly half of the false alarms are stronger than a ground-truth boundary. When repeating this experiment with a consensus boundary pixel rather than an orphan boundary, subjects select the false alarm only 10.9% of the time. Motivated by this observation, Hou *et al.* suggest evaluating results using only consensus boundary pixels as the ground-truth.

Qualitatively, we find such boundary pixels correspond to salient contours on object boundaries. This criterion appears to be quite reliable in BSDS, but weak boundaries (of which not all annotators agree) may not be. Many of the weak boundaries may be artifacts caused by original segmentation-like annotation protocols in BSDS. As a result, the standard BSDS benchmark favors algorithms that tend to over-generate boundary predictions so as to have high recall of weak boundaries. Therefore, we focus on

Method	F1-score	Precision	Recall
gPb-OwT-UCM [2]	0.591	0.512	0.698
DeepContour [29]	0.615	0.540	0.714
DeepEdge [3]	0.651	0.608	0.700
HED [48]	0.665	0.580	0.778
RCF [34]	0.693	0.621	0.784
Ours (w/o pre-training)	0.637	0.627	0.646
Ours (final)	0.707	0.706	0.708

Table 4: Salient (consensus) boundary detection on the BSDS500 dataset, a standard benchmark set up by [22].

salient boundary detection and adopt the evaluation criteria proposed by [22].

Results The BSDS500 is also a dataset with 1-to-many mapping. We can directly apply our method to this task. When we train our method on BSDS500, we experiment with two settings: whether we pre-train on our contour drawing dataset. The results are summarized in Tab 4 and Fig 7. When our model is trained only on BSDS500, it performs worse than HED [48] & RCF [34], which are pre-trained on ImageNet. But after fine-tuning, our method notably outperforms HED & RCF by a sizeable margin. Interestingly, our fine-tuned model learns to generate contours with precise pixel alignment.

5. Cost-Free Data Expansion

The edge alignment in our contour drawing is not as perfect as in BSDS500, but it is this very imperfection that makes data collection much easier, and therefore much more scalable. In this “deep” era, both BSDS500 and our current dataset are considered “small”. Comparing to boundary annotations in BSDS500, our sketch annotations typically contain more details, but it is a much easier and more interesting task than annotating precise boundaries, since we only require loose alignment. During data collection, we frequently received comments like “this is really fun”, “I like this *game*” and “I enjoyed the task”. Motivated by such comments, we further extended the interface to a sketch drawing game with the goal of collecting large scale data for free.

Prior Work Von Ahn and Dabbish [46] first point out that games can be used to label images. Their game asks a pair of players to guess the label of the same image and extract the common input as the final label. Another game is built by Deng and his colleagues [12] in which they mask an image of an object and ask the player which fine-grained class (e.g. the species of a bird) is present. The player is allowed to erase blobs of the mask at some penalty to better inspect the image. Several other games, WhatsMySketch [15] and QuickDraw [20], are built with the idea of letting the player to draw the sketch in order for an artificial intelligence system to recognize it. DrawAFriend [32] lets players trace images of celebrities or their friends and send their finished drawings to their friend to guess the identity.

Gaming Interface The fact that many data-collection games make use of sketches reinforces our observation that drawing itself is an interesting task. Therefore, we develop a game app for scalable data collection (Fig 8, demo video can be found on the project website). The challenge here is to set up a game reward system to provide user real-time feedback and also have automatic quality control mechanism for the collected data. In our game, we first process the image with a boundary detector or a contour generator (as described in Sec 3), and randomly sample points on the generated boundary map. We then define those points to be reward points, and when a player’s drawing matches a reward point, he or she receives the reward associated with the point. We also randomly sample points that are far enough from the generated boundary maps to be penalties points, *i.e.*, when a player stroke is too near a penalty point, he or she loses some of scores. Now the players get instant feedback on how well they draw and the total scores they obtained is a measure of the sketch quality. Note that all reward points are hidden to the players. Since both the reward points and penalty points are sparse, this will not limit the player to precisely follows the algorithm’s output. Then we set a cut-off score (as a percentage of the total available rewards) to perform final decision on whether to accept this



Figure 8: Our sketch game for automatic large scale data collection. The game implements real-time reward/penalty feedback and an automatic quality control mechanism.

sketch or not.

Evaluation To evaluate this reward system, or *AI score* for short, we conduct two experiments. For the first experiment, we set the sketches in our initial data collection phase (Sec 2) as ground truth, where 5000 of them are manually marked qualified and 1947 of them unqualified. When we use the AI score to rate those sketches, it can identify 90% of unqualified sketches, showing the capability of rejecting poor submissions. For the second experiment, we collected 100 additional sketches using the game interface, then we manually marked them as qualified or unqualified based on the standard in Sec 2. For this new testing, the AI score identifies 96% of the unqualified sketches.

In the future, we plan to release the game to the public to build a free sketch collection machine. As the collection process continues, we can keep update our sketch generation models, which allow us to generate more accurate reward points for the game, and make a never-ending sketch collection and generation system by closing the loop.

6. Conclusion

In this work, we examine the problem of generating contour drawings from images, introducing a dataset, benchmark criteria, and generation model. From a graphics perspective, this problem generates aesthetically pleasing line drawings. From a vision perspective, contour sketches allow for the exploration of open contours arising from geometric occlusion events. We show that such data can be used to learn low-level representations that can be fine-tuned to produce state-of-the-art results for salient boundary detection. Looking forward, contour sketches appear to be a scalable alternative for collecting geometric visual annotations (potentially through game interfaces).

References

- [1] Adobe. Adobe stock (<https://stock.adobe.com/>), 2018. 3
- [2] P. Arbelaez, M. Maire, C. Fowlkes, and J. Malik. Contour detection and hierarchical image segmentation. *TPAMI*, 33(5):898–916, May 2011. 2, 3, 5, 6, 7
- [3] G. Bertasius, J. Shi, and L. Torresani. Deepedge: A multi-scale bifurcated deep network for top-down contour detection. *CVPR*, pages 4380–4389, 2015. 6, 7
- [4] S. E. Brennan. Caricature generator: The dynamic exaggeration of faces by computer. *Leonardo*. 3
- [5] M. D. Buhrmester, T. K. G. Kwang, and S. D. Gosling. Amazon’s mechanical turk: A new source of inexpensive, yet high-quality, data? *Perspectives on psychological science: a journal of the Association for Psychological Science*, 6 1:3–5, 2011. 2, 6
- [6] J. F. Canny. A computational approach to edge detection. *TPAMI*, PAMI-8:679–698, 1986. 2, 6
- [7] H. Chen, Y.-Q. Xu, H.-Y. Shum, S.-C. Zhu, and N.-N. Zheng. Example-based facial sketch generation with non-parametric sampling. In *ICCV*, volume 2, pages 433–438 vol.2, 2001. 3
- [8] B. V. Cherkassky and A. V. Goldberg. On implementing the push–relabel method for the maximum flow problem. *Algorithmica*, 19(4):390–410, Dec 1997. 5
- [9] F. Cole, A. Golovinskiy, A. Limpaecher, H. S. Barros, A. Finkelstein, T. A. Funkhouser, and S. Rusinkiewicz. Where do people draw lines? *TOG*, 27:88:1–88:11, 2008. 2
- [10] D. DeCarlo, A. Finkelstein, S. Rusinkiewicz, and A. Santella. Suggestive contours for conveying shape. *TOG*, 22(3):848–855, 2003. 2
- [11] T. Dekel, C. Gan, D. Krishnan, C. Liu, and W. T. Freeman. Sparse, smart contours to represent and edit images. *CVPR*, 2018. 1
- [12] J. Deng, J. Krause, L. Fei-Fei, C. Li, Y. W. Li, and F. fei Li. Fine-grained crowdsourcing for fine-grained recognition. *CVPR*, pages 580–587, 2013. 8
- [13] P. Dollár, Z. Tu, and S. J. Belongie. Supervised learning of edges and object boundaries. *CVPR*, 2:1964–1971, 2006. 6
- [14] P. Dollár and C. L. Zitnick. Fast edge detection using structured forests. *TPAMI*, 37:1558–1570, 2015. 6
- [15] M. Eitz, J. Hays, and M. Alexa. How do humans sketch objects? *TOG*, 31:44:1–44:10, 2012. 2, 8
- [16] M. Eitz, J. Hays, and M. Alexa. How do humans sketch objects? *SIGGRAPH*, 31(4):44:1–44:10, 2012. 3
- [17] M. Eitz, K. Hildebrand, T. Boubekeur, and M. Alexa. Sketch-based image retrieval: Benchmark and bag-of-features descriptors. *IEEE transactions on visualization and computer graphics*, 17(11):1624–1636, 2011. 2
- [18] D. A. Forsyth and J. Ponce. In *Computer Vision - A Modern Approach, Second Edition*, 2012. 2
- [19] A. V. Goldberg and R. Kennedy. An efficient cost scaling algorithm for the assignment problem. *Mathematical Programming*, 71(2):153–177, Dec 1995. 5
- [20] D. Ha and D. Eck. A neural representation of sketch drawings. *ICLR*, 2018. 2, 3, 8
- [21] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. *CVPR*, pages 770–778, 2016. 5
- [22] X. Hou, A. L. Yuille, and C. Koch. Boundary detection benchmarking: Beyond f-measures. *CVPR*, pages 2123–2130, 2013. 4, 5, 7
- [23] P. Isola, J.-Y. Zhu, T. Zhou, and A. A. Efros. Image-to-image translation with conditional adversarial networks. In *CVPR*, 2017. 3, 4
- [24] P. Isola, J.-Y. Zhu, T. Zhou, and A. A. Efros. Image-to-image translation with conditional adversarial networks. *CVPR*, 2017. 5, 6
- [25] J. Johnson, A. Alahi, L. Fei-Fei, C. Li, Y. W. Li, and F. fei Li. Perceptual losses for real-time style transfer and super-resolution. In *ECCV*, 2016. 5
- [26] H. W. Kang, W. He, C. K. Chui, and U. K. Chakraborty. Interactive sketch generation. *The Visual Computer*, 21(8-10):821–830, 2005. 3
- [27] J. Kittler. On the accuracy of the sobel edge detector. *Image and Vision Computing*, 1(1):37–42, 1983. 6
- [28] J. J. Koenderink. *Solid shape*. MIT press, 1990. 2
- [29] I. Kokkinos. Boundary detection using f-measure-, filter- and feature- (f3) boost. In *ECCV*, 2010. 7
- [30] S. Lee, S. P. S. Prakash, M. Cogswell, V. Ranjan, D. Crandall, and D. Batra. Stochastic multiple choice learning for training diverse deep ensembles. In *NIPS*, pages 2119–2127, 2016. 5
- [31] J. J. Lim, C. L. Zitnick, and P. Dollár. Sketch tokens: A learned mid-level representation for contour and object detection. *CVPR*, pages 3158–3165, 2013. 6
- [32] A. Limpaecher, N. Feltman, A. Treuille, and M. F. Cohen. Real-time drawing assistance through crowdsourcing. *TOG*, 32:54:1–54:8, 2013. 8
- [33] M.-Y. Liu, T. Breuel, and J. Kautz. Unsupervised image-to-image translation networks. In *NIPS*, pages 700–708, 2017. 3
- [34] Y. Liu, M.-M. Cheng, X. Hu, K. Wang, and X. Bai. Richer convolutional features for edge detection. *CVPR*, pages 5872–5881, 2017. 7
- [35] D. Martin, C. Fowlkes, D. Tal, and J. Malik. A database of human segmented natural images and its application to evaluating segmentation algorithms and measuring ecological statistics. In *ICCV*, volume 2, pages 416–423, July 2001. 7
- [36] D. R. Martin, C. C. Fowlkes, and J. Malik. Learning to detect natural image boundaries using local brightness, color, and texture cues. *TPAMI*, 26(5):530–549, May 2004. 2, 5, 6
- [37] D. R. Martin, C. C. Fowlkes, D. Tal, and J. Malik. A database of human segmented natural images and its application to evaluating segmentation algorithms and measuring ecological statistics. In *ICCV*, 2001. 2
- [38] M. Mathieu, C. Couprie, and Y. LeCun. Deep multi-scale video prediction beyond mean square error. *ICLR*, 2016. 4
- [39] X. Ren. Multi-scale improves boundary detection in natural images. In *ECCV*, 2008. 6
- [40] J. Revaud, P. Weinzaepfel, Z. Harchaoui, and C. Schmid. Epicflow: Edge-preserving interpolation of correspondences for optical flow. In *CVPR*, pages 1164–1172, 2015. 1

- [41] P. Sangkloy, N. Burnell, C. Ham, and J. Hays. The sketchy database: Learning to retrieve badly drawn bunnies. *SIG-GRAPH*, 2016. [2](#), [3](#)
- [42] R. G. Schneider and T. Tuytelaars. Sketch classification and classification-driven analysis using fisher vectors. *TOG*, 33(6):174, 2014. [2](#)
- [43] W. Shen, X. Wang, Y. Wang, X. Bai, and Z. Zhang. Deep-contour: A deep convolutional feature learned by positive-sharing loss for contour detection. *CVPR*, pages 3982–3991, 2015. [6](#)
- [44] J. Song, K. Pang, Y.-Z. Song, T. Xiang, and T. M. Hospedales. Learning to sketch with shortcut cycle consistency. *CVPR*, 2018. [2](#)
- [45] J. Sutherland. Gesture drawings. *American Artist*, 61:11, 1997. [1](#)
- [46] L. von Ahn and L. A. Dabbish. Labeling images with a computer game. In *CHI*, 2004. [8](#)
- [47] N. Xie, T. Zhao, F. Tian, X. Zhang, and M. Sugiyama. Stroke-based stylization learning and rendering with inverse reinforcement learning. In *IJCAI*, 2015. [3](#)
- [48] S. Xie and Z. Tu. Holistically-nested edge detection. In *ICCV*, 2015. [1](#), [6](#), [7](#)
- [49] J.-Y. Zhu, T. Park, P. Isola, and A. A. Efros. Unpaired image-to-image translation using cycle-consistent adversarial networks. In *ICCV*, 2017. [3](#), [6](#)
- [50] J.-Y. Zhu, R. Zhang, D. Pathak, T. Darrell, A. A. Efros, O. Wang, and E. Shechtman. Toward multimodal image-to-image translation. In *NIPS*. 2017. [3](#)