# FAIRS Soft Focus Generator and Attention for Robust Object Segmentation from Extreme Points

Ahmed H. Shahin[1][*], Prateek Munjal[2], Ling Shao[3], and Shadab Khan[2]

[1] Centre for Medical Image Computing, University College London, London, UK
a.shahin@cs.ucl.ac.uk
[2] Group42 Healthcare, Abu Dhabi, UAE
{prateekmunjal31,skhanshadab}@gmail.com
[3] Inception Institute of Artificial Intelligence, Abu Dhabi, UAE
ling.shao@inceptioniai.org

**Abstract.** Semantic segmentation from user inputs has been actively studied to facilitate interactive segmentation for data annotation and other applications. Recent studies have shown that extreme points can be effectively used to encode user inputs. A heat map generated from the extreme points can be appended to the RGB image and input to the model for training. In this study, we present FAIRS a new approach to generate object segmentation from user inputs in the form of extreme points and corrective clicks. We propose a novel approach for effectively encoding the user input from extreme points and corrective clicks, in a novel and scalable manner that allows the network to work with a variable number of clicks, including corrective clicks for output refinement. We also integrate a dual attention module with our approach to increase the efficacy of the model in preferentially attending to the objects. We demonstrate that these additions help achieve significant improvements over state-of-the-art in dense object segmentation from user inputs, on multiple large-scale datasets. Through experiments, we demonstrate our method's ability to generate high-quality training data as well as its scalability in incorporating extreme points, guiding clicks, and corrective clicks in a principled manner.

## 1 Introduction

Semantic segmentation has been one of the longstanding problems in computer vision. Segmentation algorithms produce masks to classify the pixels into foreground/background classes. These algorithms are used for a wide variety of tasks, ranging from typical applications in security [33], robotics [26], satellite imaging [4], medical imaging [30], to other interesting applications such as counting number of penguins in their arctic colonies [3]. Such algorithms require a large amount of ground truth labeled data for training, which is annotated with human

---

[*] Work done during an internship at Inception Institute of Artificial Intelligence.
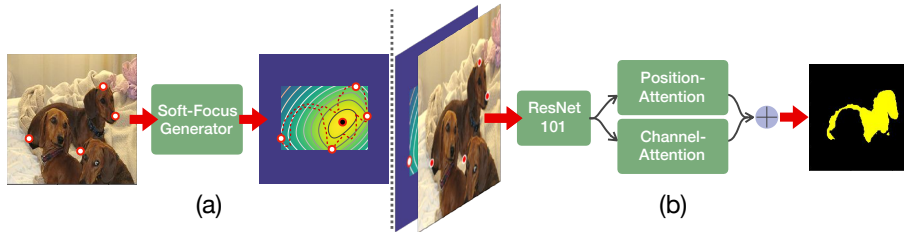
**Fig. 1.** Overview of our approach to generate object segmentation from extreme points using the proposed Soft-Focus Generator (SFG) module that results in a nearly-convex smoothly varying potential field using an $n$-ellipse formulation as shown in (a). (b) shows overview of the rest of the pipeline, we use ResNet-101 backbone with a dual attention module as proposed in [10] to produce a segmentation of the object of interest.

oversight and is therefore slow and expensive. To reduce the costs and accelerate the annotation process, methods to generate annotations from user inputs have been widely studied [27,22].

Several promising methods have been proposed that rely on user-provided cues such as bounding box [8], clicks [22,15], and scribbles [18]. These methods have worked to a varying degree of success on various datasets, and incorporating such cues from the users in a principled manner that works across datasets and conditions remains an open problem. One particular form of user clicks, called extreme points (EPs), have recently received significant attention owing to the study by [27], which showed that extreme points can be annotated more quickly than other forms of user inputs such as a bounding box.

This study, proposes a principled approach to encoding information from extreme points and corrective clicks using a Soft Focus Generator (SFG) that produces a heat map, which is input to the model for generating a dense segmentation mask (Figure 1). Further, equipped with a dual-attention module, our approach generates high-quality segmentation masks on a variety of challenging datasets such PASCAL [9], COCO [20], SBD [11], GrabCut [31], and Berkeley [25]. Compared to several state-of-the-art approaches on object segmentation from user inputs, our method (FAIRS  Focus and Attention for Interactive Robust Segmentation) took fewer clicks to achieve superior evaluation metrics on all comparative experiments conducted in this study.

We demonstrate FAIRS's effectiveness in generating training data through an exacting experiment where we trained a previously-untrained version of our model using only the synthetic labels generated by a trained version of our model. This weakly-supervised version of FAIRS achieves results that are at-par with the state-of-the-art approaches in object segmentation from user inputs, and only lags behind the version of FAIRS trained using ground truth labels. Further, we also evaluate FAIRS's performance when presented with lower than 4 extreme points or when presented with corrective clicks for refinement during interactive segmentation. We find that FAIRS handles these diverse scenarios very well, while maintaining annotation quality.

## 2 Related Work

**Segmentation from User Inputs** In recent literature, segmentation from user-provided cues such as bounding box [35], patches [28,29], scribbles [18], and clicks [27,1,13], have been investigated. In state-of-the-art in segmentation from user inputs, Liew *et al.* [16] use an image and a distance transform map computed from user input to produce multiple segmentation masks that are fed to the computationally expensive post processing step of non-maximum suppression and graph cut. Jang *et al.* [14] encode the user inputs using distance map and input it along with the image to an FCN architecture; their method relies on 10 iterations of forward-backward propagation for refining the output.

Further, Majumder *et al.* [21] augment the RGB channels by adding 4 additional channels – 2 for superpixel based guidance, and 1 channel each for object-based guidance and distance transform. Lin *et al.* [19] propose a block annotation module for online annotation that asks users to annotate blocks iteratively until a satisfactory segmentation is achieved. Li *et al.* [15] compute positive and negative distance maps which are appended to activation maps extracted from an FCN pipeline fed to a segmentation model that produces multiple segmentation masks that are post-processed through a selection network for final output.

Lastly, extreme points (EPs) have also been used for segmentation from user input [22]. Maninis *et al.*. [22] presented deep extreme cut (DEXTR), where a 2D heat map is computed using EPs, with the purpose of guiding the network to the object-of-interest. Wang *et al.* built on the ideas presented in [22], by developing a post-processing module that combines image features extracted from CNN with a level set extraction method (DELSE) [34] for refinement. EPs were also used in a full-image segmentation study by Agustsson *et al.* [2], who proposed a method that takes four EPs per object and variable number of corrective scribbles in an image to produce an image-level segmentation mask, though in contrast to [34,22], they used a 6 pixel wide circle to represent the extreme points.

The studies cited above have attempted to incorporate rich spatial information as a guidance by relying on sophisticated pre-processing [21] and/or computationally expensive post-processing steps[34], which is not ideal for fast interactive segmentation. Further, while the methods using extreme points [34,22,2] work well, the approach of placing only small-width (=10 px) Gaussians or circles at EPs does not make any other distinction between foreground/background classes, since most of the background and foreground pixels ($\sim$99.85% at 512×512 resolution) are placed at nearly zero weight. This leads to issues with segmentation output for objects that may be occluded or have confusing textural properties, such as presence of texture-less high contrast patches (e.g. dalmatian dog), or suffer from annotation error. This avenue of improvement has inspired recent methods, though finding the best method to encode the background/foreground spatial information from user inputs remains an open problem.

**Contribution**: In this study, we set out to design a mechanism for improving guidance to the neural networks. We accomplish this in two ways. First, starting from the input end of the network, we propose a novel mechanism to incorporate a simple and scalable distance map by using the *n*-ellipse (also called

multi-foci ellipse) formulation. Second, at the output end of the network, we integrate a dual-attention module as proposed in [10] within our pipeline to encourage network's preferential attention to the foreground. While attention modules have been used to improve performance of a segmentation model where object classes are known, their efficacy on class-agnostic segmentation has so far not been reported. We further note that while the attention maps are part of the architecture and are learnt, $n$-ellipse heat map provide a computed (not learnt) diffused focus as an additional input to the network. Through a series of experiments we comprehensively demonstrate the effectiveness of our approach in producing high-quality segmentation masks as well as its ability to work with a combination of extreme points and corrective clicks. In the process, we achieve performance superior to several state-of-the-art methods on object segmentation from user inputs on multiple challenging datasets.

## 3  Methods

### 3.1  Soft Focus Generator (SFG)

The SFG comprises primarily three computation steps. First, we compute the potential field of an $n$-ellipse map. Subsequently, we post-process the potential field to generate the soft focus map and crop it with the bounding box defined by extreme points. When corrective clicks are added, they are compounded with the cropped soft focus map at this stage. Finally, we compute the Gaussians heat map as used by [34,22,2] and merge it with the soft focus map. These individual steps are described next.

**1. $n$-ellipse potential field**  $\pi$: To formulate $\pi(\cdot)$, we began with the goal of achieving the following properties in the soft focus map: (i) it should be simple and fast to compute, (ii) it should scale well with the number of points, (ii) it should encode spatial relationship between the extreme points. Multi-focal ellipses or $n$-ellipse, which generalize a simple ellipse (with two focal points) to higher number of focal points, fit all of these desired properties figure 2.

Provided with an image $I(x) \in \Omega; x \in \mathbb{R}^2$ and $n$ focal points of an object of interest $p_n = p_1, p_2, ..., p_i; \ p_i \in \mathbb{R}^2$, we compute the potential field $\pi(p_n, x)$ using: $\pi(p_n, x) = \sum_{i=1}^{n} \|x - u_i\|_2$. By definition, a 1-ellipse is a typical circle, and a 2-ellipse is an ellipse. This potential field is a smoothly varying distribution of weights over the image with nearly-convex isocontours. Using the $n$-ellipse formulation enables us to use variable number of extreme points, not necessary four, while preserving a consistent smooth assignment of weights to the foreground region.

**2. Post-processing** $\pi(p_n, x)$: Once the potential field is calculated, we transform it using:

$$\hat{\pi}(p_n, x) = N \left( \frac{1}{\pi(p_n, x)} \right)^\beta \circ B(x) \tag{1}$$

In this equation, $N(\cdot)$ normalizes the range of its argument to 01, $\beta$ is a hyper-parameter that controls the potential decay rate from the center of the $n$-ellipse,

and $B(x) : \mathbb{R}^2 \to \{0, 1\}$ is a mask that is 0 outside of the bounding box and 1 everywhere inside. Limiting $\hat{\pi}(p_n, x)$ to the extent of bounding box adds an implicit cue about background pixels. Figure 2 visualized potential field for a couple of images with 3 and 4 extreme points.

**3. Incorporating corrective clicks**: For interactive segmentation, corrective clicks can be incorporated over $\hat{\pi}(p_n, x)$ by composing Gaussians for false-positives-corrective (FPC) click as well as false-negative-corrective (FNC) click. The FPC clicks at $x^i_{fpc}; i = 1...n$ are encoded using a Gaussian heat map $g(x_{fpc})$, likewise, FNC clicks are encoded using $1\text{-}g(x_{fnc})$. The FPC and FNC maps are compounded with $\hat{\pi}(p_n, x)$ to produce $\tilde{\pi}(p_n, x)$ using the equations below:

$$\breve{\pi}(p_n, x) = \begin{cases} \hat{\pi}(p_n, x, x_{fnc}), & \text{if } \hat{\pi}(p_n, x) < 1 - g(x_{fnc}) \\ 1 - g(x_{fnc}), & \text{otherwise} \end{cases} \tag{2}$$

$$\tilde{\pi}(p_n, x, x_{fpc}, x_{fnc}) = \max(\breve{\pi}(p_n, x), g(x_{fpc})) \tag{3}$$

We note that when the corrective clicks are not provided, as is the case for segmentation using four extreme points, $\tilde{\pi}(p_n, x, x_{fpc}, x_{fnc}) = \hat{\pi}(p_n, x)$.

**4. Incorporating extreme points**:

The soft focus map $\psi(p_n, x, x_{fpc}, x_{fnc})$ is computed by compounding Gaussians placed at extreme points with the post-processed potential field using: $\psi(p_n, x, x_{fpc}, x_{fnc}) = \max(\tilde{\pi}(p_n, x, x_{fpc}, x_{fnc}), g(p_n))$.
The intermediate outputs of SFG are shown in 3.

### 3.2  Network Design

Similar to [22,5,6,7,10] we begin with a ResNet-101 backbone [12] equipped with dilated convolutions, and add a positional and spatial attention modules [10], to
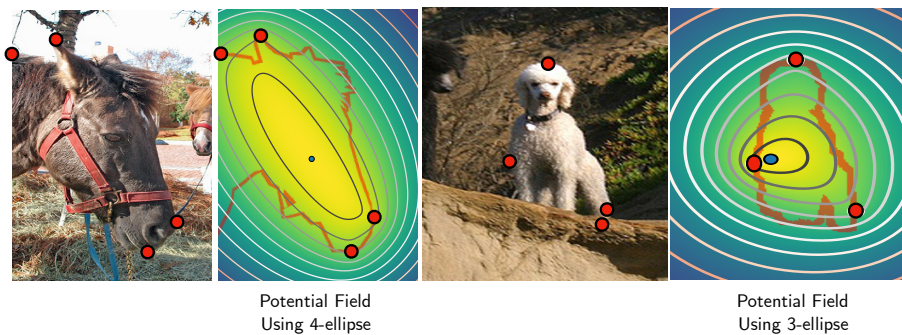


Potential Field
Using 4–ellipse

Potential Field
Using 3–ellipse

**Fig. 2.** An illustration of post-processed potential field calculated over the image using three and four extreme points. The contrast has been exaggerated for visualization. The potential field places a soft focus on the object that slowly decays away from towards the background.

model the semantic relations in spatial and channel dimensions. We remove the down-sampling operations and employ dilated convolutions in the last two blocks to preserve a reasonable spatial resolution (1/8 of the original image resolution).

Two parallel attention heads are applied to the output feature maps. Spatial attention head generates a map in which each pixel is a weighted sum of all pixels in the feature maps. Thus, it encodes global information and promotes the parts that are semantically related to the object blobs. Similarly, channel attention head models channel correlations by learning to promote channels that are relevant to the object segmentation task. The output of the two attention blocks is fused and used to produce the final segmentation mask.

## 4   Experiments and Results

We extensively experimented using FAIRS on five public datasets: Berkeley segmentation dataset [25], PASCAL 2012 [9], GrabCut [31], COCO [20], and SBD. In this section, we first discuss the model implementation and training details, followed by details of datasets and the experiments and results. We discuss multiple experiments on using FAIRS for class-agnostic segmentation, including realistic evaluation on human-annotated extreme points, generalization across dataset, and generalization to seen and unseen object categories. Further, we discuss results from experiments where we use annotations generated by our model to train a weakly-supervised version of FAIRS, and demonstrate FAIRS's ability in generating high-quality annotations. Lastly, we discuss our results on using FAIRS for interactive segmentation and demonstrate the flexibility of our approach in achieving outstanding results with different number of clicks and corrective clicks in interactive mode.
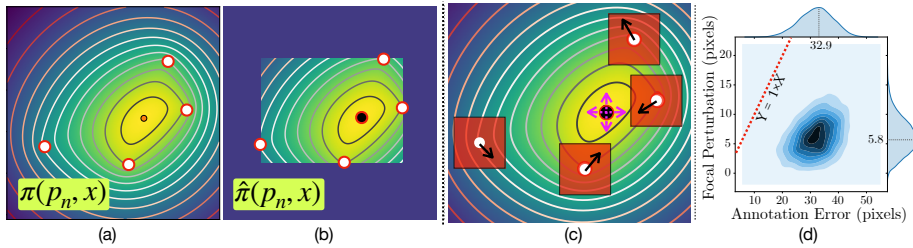


**Fig. 3.** Intermediate outputs of SFG are shown in (a) and (b). (c) illustrates an experiment done to assess the robustness of SFG to annotation error, the red boxes indicate the extent of simulated annotation error (-10–10 px), and the resulting induced perturbation on the focal point is measured. (d) show density plot of annotation error vs focal perturbation, we observe that focal point moves relatively little compared to the annotated extreme points.

### 4.1 Model Training and Datasets

**Model Training:** As is standard practice with segmentation architectures using ResNet backbone [12], we initialize our model's backbone using pre-trained ImageNet weights [32]. In order to work with pre-trained weights, we copy the $3^{rd}$ channel kernel-weights to the $4^{th}$ channel in the input layer. Attention module was initialized randomly. We used a multi-term loss with equal weights to train the model, where each term is a weighted cross entropy to alleviate the class imbalance. We use random scaling, rotations and horizontal flips for augmenting our dataset. We use SGD with momentum as the optimizer to train our model. We computed our heat map by setting $\beta = 5$, $\kappa = 2$ and Gaussian $\sigma = 10$. We compute our object masks on all images, including multiple objects within each image. Dataset-specific training details are mentioned next.

**COCO [20]:** We train using 82783 number of images from 2014 Coco train containing 80 object classes. We test on COCO 2017 validations set ($\sim$5k images, $\sim$36k objects). We train the model using a learning rate $1 \times 10^{-7}$, batch size of 48, and for 15 epochs due to the large number of images in the dataset.

**PASCAL [9,11]:** We train using PASCAL data in two different ways, one with PASCAL and SBD data (10582 images), and one with only PASCAL data (1464 images). We make the distinction between two versions of our model where necessary, otherwise, PASCAL, should be taken to mean PASCAL and SBD (10582 images) as this is the common practice when referring to this dataset. We train on PASCAL with an initial constant learning rate of $1 \times 10^{-7}$ for 100 epochs, and then reduce it to $5 \times 10^{-8}$ and train for another 50 epochs.

**Berkeley [23]:** We do not train using Berkeley train set. We test on 100 object masks extracted from 96 images, provided by [24,14].

**GrabCut: [31]** We evelute our model on GrabCut's test set (50 images).

**SBD [11]:** For reporting on SBD, we trained only on SBD data (8498 train images) for fair comparison with other methods. We test on SBD validation set (2820 images, all objects). We trained with an initial constant learning rate of $1 \times 10^{-7}$ for 100 epochs, and then reduce it to $5 \times 10^{-8}$ and train for 25 epochs.

**User Input:** We follow the approach used by [22], and infer extreme points by extracting them from the ground truth mask. To simulate noise in the extreme points, we add uniformly distributed noise of 10 px to their coordinates.

### 4.2 Ablation Study

We evaluated the relative gains by incorporating dual attention module and the output of soft focus generator (SFG). We adopt the DeepLab-V2 [5] with ResNet-101 backbone and PSP head [36,22] as the base. We do not use the PSP head with the attention module. We performed the ablation study using PASCAL-train and PASCAL+SBD-train data, and observed that, both attention module and our heat map improved the IoU scores as shown in table 1.

**Comparison with DEXTR**: To further assess the utility of SFG module, we replaced the Gaussian heat map in DEXTR pipeline with the output of SFG,

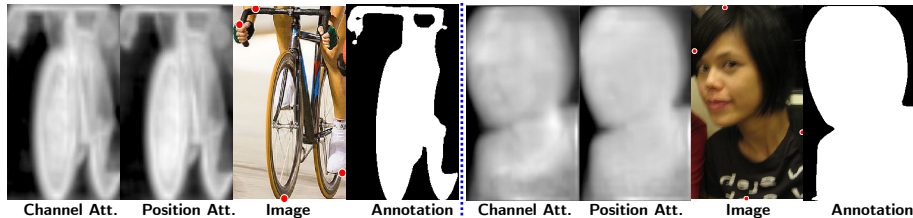| Channel Att. | Position Att. | Image | Annotation | Channel Att. | Position Att. | Image | Annotation |

**Fig. 4.** Attention maps generated for two different images are shown. On closer inspection we observe that channel attention module is attending to finer details, e.g. contours of the face, nose, rim of the cycle wheel etc, whereas the position attention module is attending preferentially to the foreground.

and trained the resulting model on PASCAL+SBD-train dataset. Compared to DEXTR's 91.50% IoU, DEXTR with SFG achieves 91.81% IoU.

**Why does SFG work?** DEXTR places small Gaussians at extreme points (EP), this is problematic if the contrast around EP is low or if annotations are imperfect, as the texture cues around EP can be misleading. To this end, many studies [14,15,16,21], have used distance transforms as one of the cues. Both Gaussians and distance transforms have peaks (+ve/-ve) at the annotated points (e.g. [16] fig.2, [14] fig.1, [21] fig.1) and distribute density around these peaks, thus suffer from annotation error or low texture around EPs. Further, errors in annotation can aggravate the problem as all of these methods place multiple peaks in the heat map. We conjecture SFG overcomes these issues because it uses $n$-ellipse potential field which has unimodal peak of the density on object foreground, rather than EPs and is therefore robust to annotation error. To test the robustness of SFG, we did an experiment where we perturbed both the $x$ and $y$ index of all 4 EPs using a random number drawn from a uniform distribution to simulate annotation error ($r{\sim}U[-10, 10]$), and measured the perturbation induced on the focus of $n$-ellipse potential field. As indicated by the marginals in the density plot (fig. 3d, 10000 draws) of annotation error vs induced perturbation, for a mean annotation error $\approx 32.9$ px, the mean induced perturbation on the focal point is $\approx 5.8$ px. These results support our hypothesis. Qualitative results using our pipeline with SFG vs our pipeline with Gaussians (DEXTR's approach) have been shown in the supplementary material, where we also present example cases where FAIRS does not improve over DEXTR.

| Model | PASCAL | PASCAL+SBD |
|---|---|---|
| Base Model | 90.50% | 91.50% |
| + Dual Attention | 91.22% | 91.80% |
| + SFG | **91.56%** | **92.22%** |

**Table 1.** Results from ablation study. We observed a gain in IoU with dual attention module as well as $n$-ellipse heat map.

**Attention Module**: To qualitatively assess the output of the attention modules, we visualize their activations in figure 4. We find that although both the attention modules, have similar overall structure, on taking a closer look, we find that channel attention module seems to be focusing on finer details, whereas position attention module is attending to the coarser foreground. This nature of attending to different details helps improve performance on challenging cases, figure 5.

### 4.3   Class-Agnostic Segmentation

FAIRS can be used for class-agnostic segmentation, using extreme points as cues provided by the user. This object can be of any class, and it can be different than classes present in the training set. We perform a number of experiments to benchmark our method's performance on class-agnostic segmentation tasks, including generalization to unseen datasets and unseen classes.

**Human-Annotated Extreme Points:**  We use the extreme points provided by [27], covering a subset of PASCAL+SBD train-set images, for evaluating FAIRS under realistic conditions. Human-annotated extreme points collection was crowd-sourced on 5623 images by [27]. To be consistent with [22], we refer to this dataset as $PASCAL_{EXT}$. We predict segmentation masks using FAIRS trained on COCO and calculate the IoU. The results are shown in the Table 2, FAIRS outperforms all other methods on this evaluation.

| Method | IoU |
|---|---|
| Sharpmask from bounding box [27] | 69.3% |
| GrabCut using extreme points [27] | 73.6% |
| Sharpmask upper bound | 78.0% |
| DEXTR from extreme points[22] | 80.1% |
| FAIRS (**Ours**) from extreme points | **84.0%** |

**Table 2.** FAIRS (trained on COCO objects dataset) compared to other methods on class-agnostic segmentation from human-annotated extreme points on $PASCAL_{EXT}$.

| Dataset | DELSE | DEXTR | Ours |
|---|---|---|---|
| COCO | – | 87.8% | **90.6%** |
| PASCAL | 90.5% | 90.5% | **91.5%** |
| PASCAL + SBD | 91.3% | 91.5% | **92.2%** |

**Table 3.** Three different models were trained on large multi-class segmentation datasets using simulated extreme points. Resulting IoU scores on PASCAL 2012 validation set are shown.

Table 2 shows that the IoU using FAIRS is significantly better than DEXTR, GrabCut-based approach and sharpmask. This demonstrates our method's ability to generalize well to human-provided extreme points despite being trained on simulated extreme points.

**Segmentation From Simulated Extreme Points:** We experiment with large scale datasets by simulating extreme points as described previously. We trained three versions of FAIRS for this study. The three models were trained using
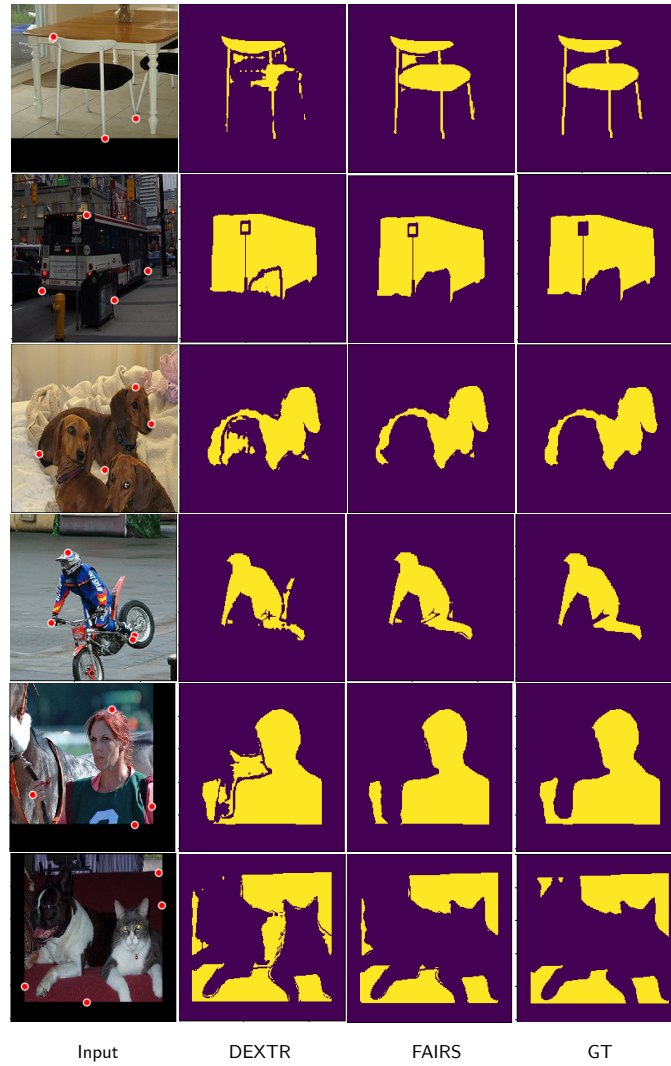
Input            DEXTR            FAIRS            GT

**Fig. 5.** We demonstrate FAIRS' ability to deal with a variety of challenging scenarios in comparison to DEXTR. The network is able to preferentially attend to foreground with the help of attention modules and soft focus map, which helps it perform under challenging conditions. For example, all extreme points of the chair are on the white patch, yet, FAIRS was able to produce a reliable segmentation with a small false +ve blob that can be further refined using corrective clicks. Similarly, in the second row, FAIRS is able to recognize that the newspaper dispenser is not a part of the bus, despite one of the extreme points being close to it. In the last row image, FAIRS is able to deal with presence of texture-less contrast in a better-controlled manner compared to DEXTR. Additional results in supplementary material.
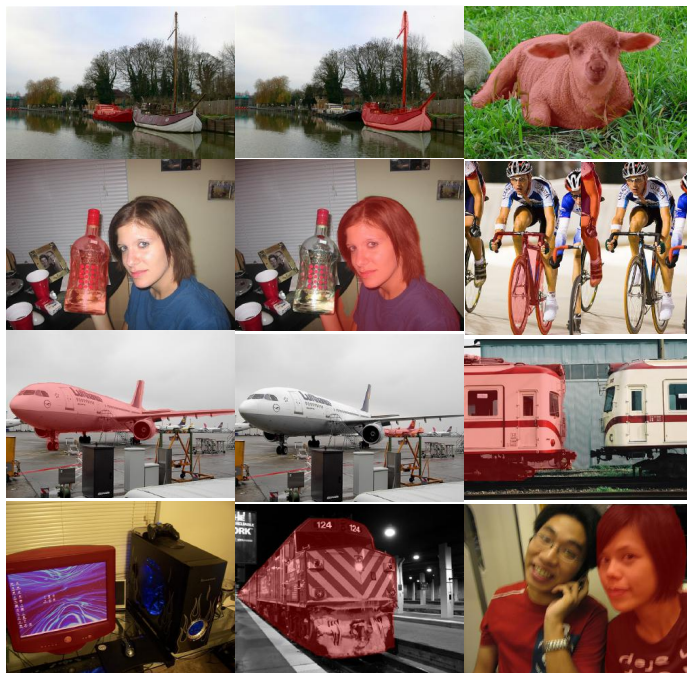
**Fig. 6.** A few example masks predicted using FAIRS shown overlaid on images from PASCAL dataset. We particularly highlight FAIRS's robustness to objects surrounded by clutter, e.g. top row middle image of the boat, and 2nd row right most image of the bicycle. Additional results in supplementary material.

the training dataset from COCO 2014 train set, PASCAL 2012 train set (1464 images), and PASCAL 2012 train set combined with SBD train set (10582 images). In table 3 below, we compare the IoU score on PASCAL 2012 validation set. Qualitative comparison of masks achieved with FAIRS and DEXTR are shown for a number of challenging cases in figure 5 and demonstrate FAIRS's robustness. Results in table 3 show that our method improves significantly over other state-of-the-art methods. Further in figure 6, we show FAIRS's ability to deal with cluttered and non-cluttered scenes for object segmentation using SFG output computed from four extreme points.

**Comparison to State-of-The-Art and Number of Clicks:** In order to compare with state-of-the-art methods that offer segmentation from user-inputs, we trained FAIRS, 3 and 4 extreme points, in addition to upto 2 corrective clicks. Our method of simulating 3 points, and corrective clicks is described as follows.

We simulate three extreme points by: (i) obtaining 4 extreme points, (ii) identifying a pair of extreme points that are closest to each other, (iii) dropping one of the points in the pair randomly during training. This closely simulates the actual use-case where an annotator would preferentially select extreme points

| Methods\Datasets | PASCAL 85% | Berkeley 90% | GrabCut 90% | SBD 85% |
|---|---|---|---|---|
| RIS-Net [17] | 5.7 | – | 6.0 | – |
| Latent Diversity [15] | – | – | 4.79 | 7.41 |
| DEXTR [34] | 4.0(91.5%) | 4+ (89.1 @4) | 4.0 | – |
| CAMLG [21] | 3.62 | 5.6 | 3.5 | – |
| FCTSFN [13] | 4.58 | 6.49 | 3.58 | – |
| MultiSeg [16] | 3.51 | 4.00 | – | – |
| BRS-DenseNet [14] | – | 5.08 | 3.60 | 6.59 |
| **FAIRS-WS (Ours)** | 4.0 (91.4%) | 4.0 (88.8%) | 4.0 (92.8%) | – |
| **FAIRS (Ours)** | **3.0 (88.9%)** | **4.0 (91.9%)** | **3.0 (91.9%)** | **4.0 (88%)** |

**Table 4.** We compare FAIRS's effectiveness with state-of-the-art methods on segmentation from user inputs on multiple datasets. Resulting IoU scores are shown. Most recent methods are in the last few rows of the table. We note that FAIRS-WS was trained using only the generated masks on COCO 2017 validation set using FAIRS trained on PASCAL. FAIRS-WS result demonstrate the effectiveness of generated labels in actual training for a realistic use-case.

such that the coverage of the objects is maximized. Further, corrective clicks were simulated by identifying largest false positive and false negative blobs, and sampling a point randomly from the largest blob across the false positives and negatives. We note that soft focus generator was able to handle all of these scenarios automatically for producing the soft focus map.

To demonstrate the efficiency of our encoding method with lower number of extreme points, we report the number of clicks needed by each algorithm to reach certain IoU on Berkeley, PASCAL validation set, GrabCut, and SBD validation set. The results organized by dataset are shown in the table 4. FAIRS outperforms all state-of-the-art methods that we compared with on all of the datasets that we experimented with. [22]

**Generalization to Unseen Classes:** We evaluate FAIRS's ability to generalize to unseen classes by training the model on PASCAL training set, and evaluating on COCO 2017 validation set. In this experiment, we run FAIRS to compute masks for object classes that are not present in the PASCAL train set (COCO Unseen), this results in 60 object classes, with ~15k number of objects. We compute masks for all objects. The results are shown in table 5, which shows that FAIRS suffers from a negligible comparative performance drop when we use it to segment object classes that were absent from the training data.

**Generalization to Unseen Datasets:** In this experiment, we evaluate FAIRS's ability to generalize to new datasets. For fair comparison with DEXTR, we conducted this experiment consistent with DEXTR's approach. That is, we report results with our model trained on both COCO and PASCAL, with both of their validation sets, as shown in table 6. FAIRS achieves a higher IoU in these eval-

uations and reaches an IoU of 85% on COCO validation set when trained with COCO training set.

### 4.4   Assisted Annotation – Quality and Budget

A key application of a tool such as FAIRS is instance segmentation from user inputs. In order to evaluate the quality of masks produced by FAIRS we conducted the following experiment. First, we used our PASCAL-trained model to produce annotations for COCO validation set, which is our hypothetical new dataset to be annotated. We refer to these FAIRS-produced annotations as COCO-GT-Noisy (GT=Ground Truth). Second, we train a new version of our model (not trained previously on any other segmentation dataset), using COCO validation set images, but instead of actual ground truth, we use generated labels COCO-GT-Noisy. We refer to this version of our model as FAIRS-WS (WS: weakly supervised).

We, evaluate FAIRS-WS's performance on PASCAL, COCO, Berkeley, and GrabCut datasets. We note that FAIRS-WS was never trained on any segmentation ground truth data. Given the noise inherent in automatic labeling used to generate COCO-GT-Noisy, we highlight that extreme points input to the FAIRS-WS are noisy. Therefore, during training, FAIRS-WS has two difficulties to overcome: noisy annotations, and noisy user inputs (simulated extreme points from COCO-GT-Noisy). In this manner, we comprehensively test FAIRS's ability to create new annotations. We report the results with this model in table 4.3, and additional results using FAIRS-WS have been shown in table 4.
To elucidate these results, we highlight the following points:

**1.** Using only the annotations generated from extreme points on ~36k objects, and with no segmentation pre-training with any ground truth data, our model achieves an IoU of 91.4% on PASCAL validation set  on par with two state-of-the-art approaches [DELSE, DEXTR] that report similar result on PASCAL by training on ~25k objects with human-annotated ground truth annotations.

**2.** We demonstrate that two versions of our model, one trained on PASCAL ground truth (~10.5k images, ~25k objects), and the other trained on generated training data (5k images, ~36k objects), achieve nearly the same IoU. This demonstrates that annotations generated using FAIRS are effective for training fully-supervised segmentation models.

| Train | Test | | DEXTR | Ours |
|---|---|---|---|---|
| PASCAL | COCO | Seen | 80.3% | **81.8%** |
| PASCAL | COCO | Unseen | 79.9% | **81.7%** |

**Table 5.** Evaluation of FAIRS's ability to generalize on unseen classes of COCO 2017 Validation set.

**3.** Finally, we note that assuming 7.5 seconds as annotation time for extreme points [27], and ∼2 minutes as a very conservative estimate of full annotation time ([27] mention 55 seconds as median bounding box annotation time, and full annotation typically takes much longer), FAIRS-WS annotations (on COCO Validation set, ∼35k objects) could be obtained ∼11.4x faster than ground truth annotations (on PASCAL ∼25k objects) at the expense of marginally lower performance compared to FAIRS.

### 4.5   Interactive Object Segmentation

With FAIRS, user can start with 2 or more extreme points and add further positive or negative clicks. To demonstrate FAIRS's ability to encode corrective clicks for refinement, we trained our PASCAL model with additional positive and negative clicks. We simulated additional clicks by randomly sampling points within a distance of 15–60 pixels from the boundary of the mask to simulate refinement over false positive and false negative regions. We report results on PASCAL validation set for instances where IoU with 4 extreme points was less than 70%, representing a scenario where the user might add a corrective click. The corrective click at test time was sampled by randomly sampling a point from the largest blob that contains either the false positives or false negatives. With the $5^{th}$ click added by random sampling for these hard samples, IoU improved by 6.1% from 66.6%, for relative gain of 9.2%. We note that random sampling is not a fair representation of likely improvement, but represents a minimum gain achievable with the method. With an actual user click that is likely to be at a more conducive spot on the false positive or negative area, we expect the improvement to be greater. Lastly, in our experiments with only 3 extreme points, we observed that FAIRS was able to reach an IoU of 88.9% on PASCAL, and 91.9% on GrabCut. These results suggest that FAIRS can work well in an interactive segmentation mode, with a variety of click budgets.

| Train | Test | DEXTR | Ours |
|---|---|---|---|
| Pascal | COCO 2017 Val | 80.1% | **81.76%** |
| COCO | COCO 2017 Val | 82.1% | **85%** |
| COCO | Pascal Val | 87.8% | **90.6%** |
| Pascal | Pascal Val | 91.5% | **92.2%** |

**Table 6.** Evaluation of FAIRS's performance on unseen datasets.

| Data | FAIRS-WS | DEXTR | DELSE |
|---|---|---|---|
| PASCAL | 91.4% | **91.5%** | 91.3% |
| COCO | **81.5%** | 80.3% | – |

**Table 7.** IoU results using FAIRS-COCO-Noisy on PASCAL and COCO datasets. We note that our method achieves IoUs on-par with DEXTR on both the datasets, despite never being trained on an actual ground truth.

## 5    Conclusion

In this study, we presented a novel scalable manner of incorporating cues from user-clicks, in a principled manner, in order to encode rich information for guiding a neural network towards the object of interest. Integrated with a dual attention module and a ResNet-101 backbone, we demonstrated through extensive experiments that FAIRS achieves its purpose of generating high quality data for fully supervised training, as evidenced by the results from FAIRS-WS. Finally, we demonstrated FAIRS's ability to handle ¡4 extreme points as well as corrective clicks in a single unified manner, enabled by soft focus generator. With these outcomes, we believe FAIRS can be an effective object segmentation tool.

# References

1. Acuna, D., Ling, H., Kar, A., Fidler, S.: Efficient Interactive Annotation of Segmentation Datasets with Polygon-RNN++. In: 2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 859–868. IEEE (6 2018). https://doi.org/10.1109/CVPR.2018.00096
2. Agustsson, E., Uijlings, J.R.R., Ferrari, V.: Interactive Full Image Segmentation by Considering All Regions Jointly. In: 2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (2019)
3. Arteta, C., Lempitsky, V., Zisserman, A.: Counting in the Wild. In: European Conference on Computer Vision, vol. 9911 LNCS, pp. 483–498. Springer Verlag (2016). https://doi.org/10.1007/978-3-319-46478-7_30, http://link.springer.com/10.1007/978-3-319-46478-7_30
4. Bhandari, A.K., Singh, V.K., Kumar, A., Singh, G.K.: Cuckoo search algorithm and wind driven optimization based study of satellite image segmentation for multi-level thresholding using Kapur's entropy. Expert Systems with Applications $41(7)$, 3538–3560 (6 2014). https://doi.org/10.1016/j.eswa.2013.10.059
5. Chen, L.C., Papandreou, G., Kokkinos, I., Murphy, K., Yuille, A.L.: Deeplab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected crfs. IEEE transactions on pattern analysis and machine intelligence $40(4)$, 834–848 (2017)
6. Chen, L.C., Papandreou, G., Schroff, F., Adam, H.: Rethinking atrous convolution for semantic image segmentation. arXiv preprint arXiv:1706.05587 (2017)
7. Chen, L.C., Zhu, Y., Papandreou, G., Schroff, F., Adam, H.: Encoder-decoder with atrous separable convolution for semantic image segmentation. In: Proceedings of the European conference on computer vision (ECCV). pp. 801–818 (2018)
8. Dai, J., He, K., Sun, J.: BoxSup: Exploiting Bounding Boxes to Supervise Convolutional Networks for Semantic Segmentation. In: 2015 IEEE International Conference on Computer Vision (ICCV). pp. 1635–1643. IEEE (12 2015). https://doi.org/10.1109/ICCV.2015.191
9. Everingham, M., Van Gool, L., Williams, C.K.I., Winn, J., Zisserman, A.: The PASCAL Visual Object Classes Challenge 2012 (VOC2012) Results, http://www.pascal-network.org/challenges/VOC/voc2012/workshop/index.html
10. Fu, J., Liu, J., Tian, H., Li, Y., Bao, Y., Fang, Z., Lu, H.: Dual attention network for scene segmentation. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 3146–3154 (2019)
11. Hariharan, B., Arbeláez, P., Bourdev, L., Maji, S., Malik, J.: Semantic contours from inverse detectors. In: 2011 International Conference on Computer Vision. pp. 991–998. IEEE (2011)
12. He, K., Zhang, X., Ren, S., Sun, J.: Deep Residual Learning for Image Recognition. In: 2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR). pp. 770–778. IEEE (6 2016). https://doi.org/10.1109/CVPR.2016.90
13. Hu, Y., Soltoggio, A., Lock, R., Carter, S.: A fully convolutional two-stream fusion network for interactive image segmentation. Neural Networks $109$, 31–42 (2019)
14. Jang, W.D., Kim, C.S.: Interactive Image Segmentation via Backpropagating Refinement Scheme. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 5297–5306 (2019)
15. Li, Z., Chen, Q., Koltun, V.: Interactive Image Segmentation with Latent Diversity. In: 2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 577–585. IEEE (6 2018). https://doi.org/10.1109/CVPR.2018.00067

16. Liew, J.H., Cohen, S., Price, B., Mai, L., Ong, S.H., Feng, J.: MultiSeg: Semantically Meaningful, Scale-Diverse Segmentations From Minimal User Input. In: Proceedings of the IEEE International Conference on Computer Vision. pp. 662–670 (2019)

17. Liew, J., Wei, Y., Xiong, W., Ong, S.H., Feng, J.: Regional Interactive Image Segmentation Networks. In: 2017 IEEE International Conference on Computer Vision (ICCV). pp. 2746–2754. IEEE (10 2017). https://doi.org/10.1109/ICCV.2017.297

18. Lin, D., Dai, J., Jia, J., He, K., Sun, J.: ScribbleSup: Scribble-Supervised Convolutional Networks for Semantic Segmentation. In: 2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR). pp. 3159–3167. IEEE (6 2016). https://doi.org/10.1109/CVPR.2016.344

19. Lin, H., Upchurch, P., Bala, K.: Block Annotation: Better Image Annotation with Sub-Image Decomposition. In: Proceedings of the IEEE International Conference on Computer Vision. pp. 5290–5300 (2019)

20. Lin, T.Y., Maire, M., Belongie, S., Hays, J., Perona, P., Ramanan, D., Dollár, P., Zitnick, C.L.: Microsoft COCO: Common Objects in Context. pp. 740–755. Springer, Cham (2014). https://doi.org/10.1007/978-3-319-10602-1_48

21. Majumder, S., Yao, A.: Content-Aware Multi-Level Guidance for Interactive Instance Segmentation. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 11602–11611 (2019)

22. Maninis, K.K., Caelles, S., Pont-Tuset, J., Van Gool, L.: Deep Extreme Cut: From Extreme Points to Object Segmentation. In: 2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 616–625. IEEE (6 2018). https://doi.org/10.1109/CVPR.2018.00071

23. Martin, D., Fowlkes, C., Tal, D., Malik, J., others: A database of human segmented natural images and its application to evaluating segmentation algorithms and measuring ecological statistics. Iccv Vancouver: (2001)

24. McGuinness, K., OConnor, N.E.: A comparative evaluation of interactive segmentation algorithms. Pattern Recognition **43**(2), 434–444 (2 2010). https://doi.org/10.1016/j.patcog.2009.03.008, `https://linkinghub.elsevier.com/retrieve/pii/S0031320309000818`

25. McGuinness, K., OConnor, N.E.: Toward automated evaluation of interactive segmentation. Computer Vision and Image Understanding **115**(6), 868–884 (6 2011). https://doi.org/10.1016/j.cviu.2011.02.011, `https://linkinghub.elsevier.com/retrieve/pii/S1077314211000889`

26. Milioto, A., Stachniss, C.: Bonnet: An open-source training and deployment framework for semantic segmentation in robotics using CNNs. In: Proceedings - IEEE International Conference on Robotics and Automation. vol. 2019-May, pp. 7094–7100. Institute of Electrical and Electronics Engineers Inc. (5 2019). https://doi.org/10.1109/ICRA.2019.8793510

27. Papadopoulos, D.P., Uijlings, J.R.R., Keller, F., Ferrari, V.: Extreme Clicking for Efficient Object Annotation. In: 2017 IEEE International Conference on Computer Vision (ICCV). pp. 4940–4949. IEEE (10 2017). https://doi.org/10.1109/ICCV.2017.528

28. Pinheiro, P.O., Collobert, R., Dollár, P.: Learning to segment object candidates. In: Advances in Neural Information Processing Systems. pp. 1990–1998 (2015)

29. Pinheiro, P.O., Lin, T.Y., Collobert, R., Dollár, P.: Learning to refine object segments. In: European Conference on Computer Vision. pp. 75–91. Springer (2016)

30. Ronneberger, O., Fischer, P., Brox, T.: U-Net: Convolutional Networks for Biomedical Image Segmentation. In: Medical Image Computing and Computer Assisted

Interventions (MICCAI), pp. 234–241 (5 2015). https://doi.org/10.1007/978-3-319-24574-4_28

31. Rother, C., Kolmogorov, V., Blake, A.: "GrabCut" - Interactive foreground extraction using iterated graph cuts. In: ACM Transactions on Graphics. vol. 23, pp. 309–314 (8 2004). https://doi.org/10.1145/1015706.1015720, `http://portal.acm.org/citation.cfm?doid=1015706.1015720`

32. Russakovsky, O., Deng, J., Su, H., Krause, J., Satheesh, S., Ma, S., Huang, Z., Karpathy, A., Khosla, A., Bernstein, M., Berg, A.C., Fei-Fei, L.: ImageNet Large Scale Visual Recognition Challenge. International Journal of Computer Vision **115**(3), 211–252 (12 2015). https://doi.org/10.1007/s11263-015-0816-y

33. Snidaro, L., Micheloni, C., Chiavedale, C.: Video security for ambient intelligence. IEEE Transactions on Systems, Man, and Cybernetics Part A:Systems and Humans. **35**(1), 133–144 (1 2005). https://doi.org/10.1109/TSMCA.2004.838478

34. Wang, Z., Acuna, D., Ling, H., Kar, A., Fidler, S.: Object Instance Annotation with Deep Extreme Level Set Evolution. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 7500–7508 (2019). https://doi.org/CVPR 2019: 7500-7508, `https://github.com/fidler-lab/delse.`

35. Xu, N., Price, B., Cohen, S., Yang, J., Huang, T.: Deep GrabCut for Object Selection. In: Procedings of the British Machine Vision Conference 2017. British Machine Vision Association (2017). https://doi.org/10.5244/C.31.182

36. Zhao, H., Shi, J., Qi, X., Wang, X., Jia, J.: Pyramid Scene Parsing Network. In: 2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR). pp. 6230–6239. IEEE (7 2017). https://doi.org/10.1109/CVPR.2017.660