# Boundary-sensitive Network for Portrait Segmentation

Xianzhi Du[1], Xiaolong Wang[2], Dawei Li[2], Jingwen Zhu[2], Serafettin Tasci[2],
Cameron Upright[2], Stephen Walsh[2], Larry Davis[1]
[1]Computer Vision Lab, UMIACS, University of Maryland, College Park, 20742
[2]Samsung Research America, Mountain View, 94043
{xianzhi,lsd}@umiacs.umd.edu
{xiaolong.w,dawei.l,jingwen.z,s.tasci,c.upright,s1.walsh}@samsung.com

## Abstract

*Compared to the general semantic segmentation problem, portrait segmentation has higher precision requirement on boundary area. However, this problem has not been well studied in previous works. In this paper, we propose a boundary-sensitive deep neural network (BSN) for portrait segmentation. BSN introduces three novel techniques. First, an individual boundary-sensitive kernel is proposed by dilating the contour line and assigning the boundary pixels with multi-class labels. Second, a global boundary-sensitive kernel is employed as a position sensitive prior to further constrain the overall shape of the segmentation map. Third, we train a boundary-sensitive attribute classifier jointly with the segmentation network to reinforce the network with semantic boundary shape information. We have evaluated BSN on the current largest public portrait segmentation dataset,* i.e.*, the PFCN dataset, as well as the portrait images collected from other three popular image segmentation datasets: COCO, COCO-Stuff, and PASCAL VOC. Our method achieves the superior quantitative and qualitative performance over state-of-the-arts on all the datasets, especially on the boundary area.*

## 1. Introduction

Semantic segmentation is a fundamental problem in computer vision community which aims to classify pixels into semantic categories. In this paper, we target a special binary class semantic segmentation problem, namely portrait segmentation, which generates pixel-wise predictions as foreground (*i.e*., people) or background. Recently, it is becoming a hot topic and has been widely used in many real-world applications, such as augmented reality (AR), background replacement, portrait stylization, depth of field, advanced driver assistance systems [5], etc. Although numerous deep learning based approaches (*e.g*., [6]

[1] [19] [24] [16] [15] [14]) were proposed to solve the general semantic segmentation problem, direct adaptation of these methods cannot satisfy the high precision requirement in the portrait segmentation problem.

In portrait segmentation, precise segmentation around object boundaries is crucial but challenging. For applications like background replacement, accurate and smooth boundary segmentation (such as hair and clothes) is the key for better visual effects. However, this has long been one of the most challenging part of portrait segmentation, especially when using convolutional neural networks (CNN). Since the neighborhood of boundary pixels contains a mixture of both foreground and background labels, convolutional filters fuse information of different classes, which may confuse the network when segmenting boundary pixels. Previous CNN based semantic segmentation methods, which use either the conventional hard-label method or ignore the boundary pixels during training [6] [1] [10], fail to solve this problem. These methods aim to train a better model to separate foreground and background while sacrificing the accuracy when predicting the boundary pixels.

In this paper, we propose a new boundary-sensitive network (BSN) for more accurate portrait segmentation. In contrast to conventional semantic image segmentation systems, we dilate the contour line of the portrait foreground and label the boundary pixels as the third class with the proposed soft-label method. Two boundary-sensitive kernels are introduced into the loss function to help the network learn better representations for the boundary class as well as govern an overall shape of the portrait. The first boundary-sensitive kernel is designed for each training image such that a floating point vector is assigned as a soft label for each pixel in the boundary class. The second boundary-sensitive kernel is a global kernel where each location in the kernel indicates the probability of the current location belonging to the boundary class. Furthermore, a boundary-sensitive attribute classifier is trained jointly with the segmentation network to reinforce the training process. We evaluate our
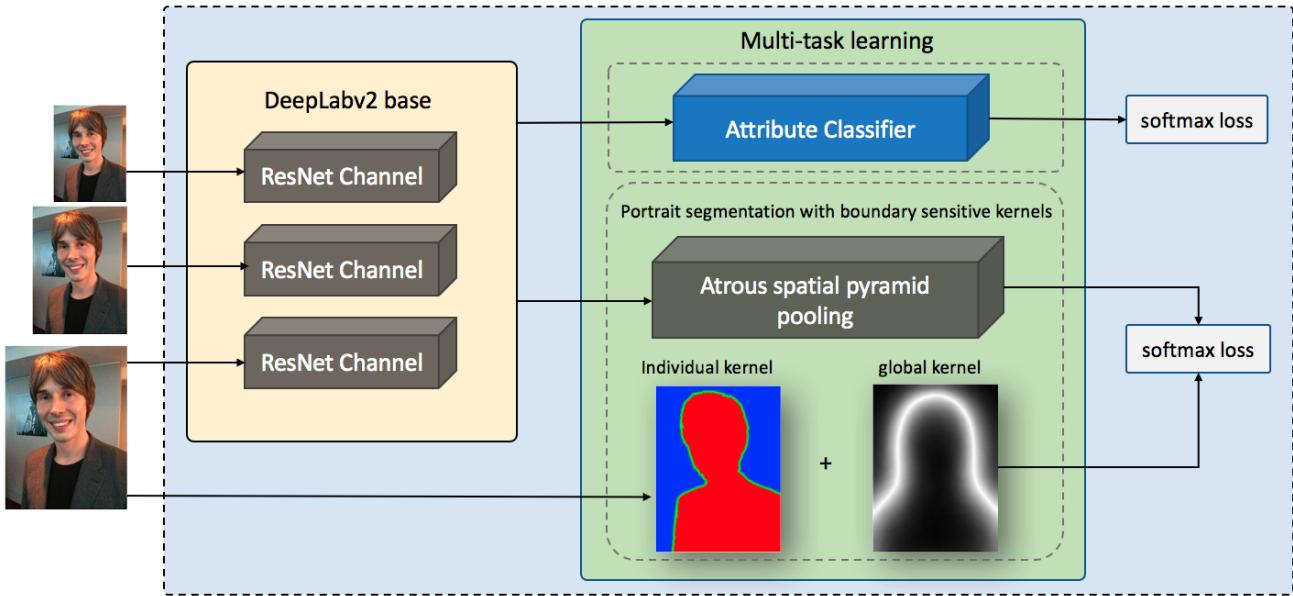
Figure 1: The whole architecture of our framework.

method on PFCN [21], the largest available portrait segmentation dataset. Our method achieves the best quantitative performance in mean IoU at 96.7%. In order to show the effectiveness and generalization capability of our method, we further test on the portrait images collected from COCO [17], COCO-Stuff [3], PASCAL VOC [11] and the experiment results demonstrate that our method significantly outperforms all other state-of-the-art methods.

The rest of this paper is organized as follows: Section 2 reviews the previous work on related problems. Section 3 describes the general framework and the three boundary-sensitive techniques in detail. Section 4 discusses and analyzes the experimental results. Section 5 draws conclusions and discusses further work.

## 2. Related work

**Semantic segmentation** systems can be categorized as unsupervised methods and supervised methods. Unsupervised methods solve the semantic segmentation problem with classic machine learning techniques include thresholding, hashing, K-means clustering, topic model, graph-cut [22] [8] [9], etc. On the other hand, conventional supervised methods treat the semantic segmentation problem as a pixel-wise classification problem which first build handcrafted features and then train classifiers such as Support Vector Machines [12], Random Forest [2], etc.

In recent years, convolutional neural network (CNN) based methods have been successfully applied to semantic segmentation. In 2014, Long *et al*. [18] introduced the end-to-end Fully Convolutional Networks (FCN) which takes a natural image as input and performs dense pixel-wise predictions to generate a segmentation map of the same size as the input image. Fully connected layers are removed from this network to preserve the spatial information and deconvolutional layers are proposed for up-sampling to recover the full image size. This paradigm popularized the CNN based method and was quickly adopted by subsequent approaches. In traditional CNN architectures, pooling layer was introduced to increase the receptive field as the network goes deeper. However, it also decreases the resolution of feature map. Yu *et al*. [23] proposed the dilated convolutional layer to replace the pooling layer, which allows for increasing the size of the receptive field without losing resolution in feature maps. Chen *et al*. [7] proposed the DeepLab system which passes multiple rescaled input images to different network branches in parallel and combines the features maps with max operation at the end.

**Portrait segmentation** is generally regarded as a sub-problem of semantic segmentation, and it is different from traditional segmentation in two aspects. First, the foreground object is limited to only people which provides additional prior information. Meanwhile, portrait segmentation has higher precision requirement on boundary area. Shen *et al*. [21] fine-tuned a portrait segmentation system from a pre-trained FCN network with portrait images. To provide more portrait-specific information to the network, two normalized $x$ and $y$ position channels and one mean mask shape channel are added to the input image. Shen *et al*. [20] proposed a joint correspondence and segmentation estimation method by using extra information provided by
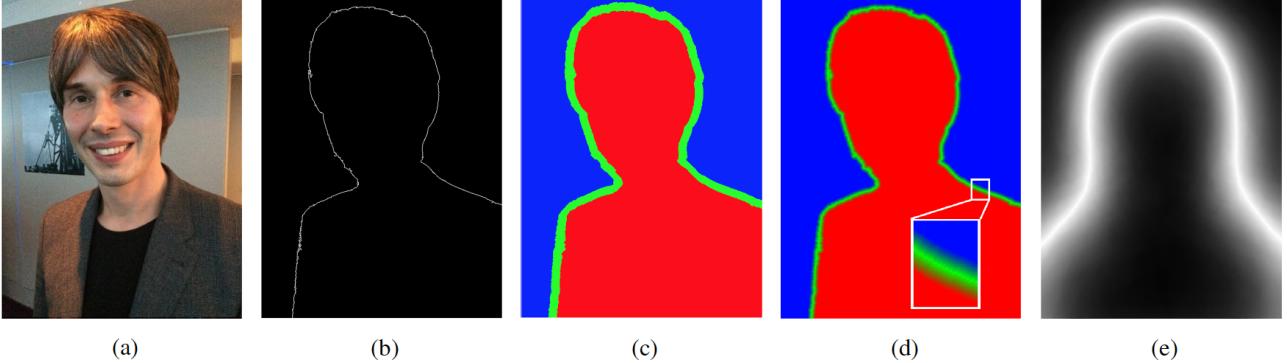
Figure 2: The kernel generating process in our method: (a) represents the original image; (b) represents the detected contour line; (c) shows the three class labels: foreground, background, and dilated boundary; (d) shows the individual boundary-sensitive kernel; (e) shows the global boundary-sensitive kernel.

dual-lens camera.

While most methods can easily generate a rough segmentation, they generally fail to provide precise segmentation near the object boundaries. For refining the predictions near the boundaries, the most commonly used solution is employing Conditional Random Fields (CRF) along with CNN. Deeplab[7] employs dense CRF after CNN as a post processing method to smooth out the predictions. However, CRF is generally used as a post-precessing step and may be quite time-consuming.

## 3. Boundary-sensitive portrait segmentation

The architecture of our framework is shown in Figure 1. We use DeepLabv2_ResNet101 model as the base segmentation network. DeepLabv2_ResNet101 consists of three ResNet101 branches at the base which process different scales of the input image. Then the three branches are followed by the atrous spatial pyramid pooling (ASPP) at different dilation rates and fused together at the end. For more details please refer to [7]. To make the model more sensitive to a portrait's boundary, during training, we label the training samples with three non-overlapping classes: foreground, boundary, and background, using the soft-label method described below. One individual boundary-sensitive kernel and one global boundary-sensitive kernel are introduced when updating the loss function, which affect both the forward pass and the back-propagation. The generation process of the two kernels are shown in Figure 2. Furthermore, an attribute classifier which shares the base layers with BSN is trained jointly with the segmentation task to reinforce the training process.

### 3.1. The individual boundary-sensitive kernel and the soft-label method

To better address the boundary prediction problem, we introduce the individual boundary-sensitive kernel. We label the boundary class as a third class to separate from foreground and background classes and assign soft-labels to pixels in the boundary class as follows. First, the portrait's contour line is identified in the ground truth segmentation map with the Canny edge detector [4]. The contour is then dilated to be P-pixels in width and that map is overlayed onto the ground truth segmentation map. We call the new label map the individual boundary-sensitive kernel. For each pixel in the kernel, a $1 \times 3$ floating-point vector $K^{indv} = [l^{fg}, l^{bdry}, l^{bg}]$ is assigned as the soft-label to represent how likely the current pixel belongs to each class. The $K^{indv}$ is computed as Equations (1) (2) (3).

$$
l_i^{bdry} = \begin{cases} \dfrac{\min\limits_{\forall C_j \in C} ||I_i - C_j||}{\sum\limits_{k} \min\limits_{\forall C_j \in C} ||I_k - C_j||} & , \text{if } i \in \text{boundary} \\ 0 & , \text{if } i \in \text{foreground} \\ 0 & , \text{if } i \in \text{background} \end{cases} \tag{1}
$$

$$
l_i^{fg} = \begin{cases} \mathbb{1}(M_i \in fg)(1 - l_i^{bdry}) & , \text{if } i \in \text{boundary} \\ 1 & , \text{if } i \in \text{foreground} \\ 0 & , \text{if } i \in \text{background} \end{cases} \tag{2}
$$

$$
l_i^{bg} = \begin{cases} \mathbb{1}(M_i \in bg)(1 - l_i^{bdry}) & , \text{if } i \in \text{boundary} \\ 0 & , \text{if } i \in \text{foreground} \\ 1 & , \text{if } i \in \text{background} \end{cases} \tag{3}
$$

where $\min_{\forall C_j \in C} ||I_i - C_j||$ represents the distance from the current pixel $I_i$ to the nearest point on the contour line $C$. $M_i$ represents the binary label of the current pixel in the original label map $M$. We can see that pixels in the foreground/background class are labeled as $[1, 0, 0]/[0, 0, 1]$ and pixels in the boundary class are labeled with a floating-point vector. The soft-label method computes $l^{bdry}$ as the normalized distance from the current pixel to the nearest point on the contour and sets $l^{fg}$ and $l^{bg}$ to either $(1 - l^{bdry})$ or 0 based on the class label of the current pixel in the ground truth segmentation map. During the forward pass for each pixel in one sample, the new formula for updating the loss function can be expressed as Equation (4):

$$\epsilon = -\sum_{j=1}^{c} K_j^{indv} \times log(\frac{e^{z_j}}{\sum_k e^{z_k}}) = -\sum_{j=1}^{c} K_j^{indv} \times log(y_j)$$

(4)

where $l_j$ denotes the soft-label for class $j$ and $y_j = e^{z_j} / \sum_k e^{z_k}$ denotes the softmax probability for this class. $c$ represents all the three classes. The new back-propagation for this sample can be derived as in Equation (5):

$$
\begin{aligned}
\frac{\partial \epsilon}{\partial z_i} &= -\sum_{j=1}^{c} K_j^{indv} \times \frac{\partial log(y_j)}{\partial z_i} \\
&= -(\frac{K_i^{indv}}{y_i} \times \frac{\partial y_i}{\partial z_i} + \sum_{j \neq i}^{c} \times \frac{K_j^{indv}}{y_j} \frac{\partial y_j}{\partial z_i}) \\
&= -(\frac{K_i^{indv}}{y_i} \times y_i \times (1 - y_i) - \sum_{j \neq i}^{c} \frac{K_j^{indv}}{y_j} \times (y_j \times y_i)) \\
&= -(K_i^{indv} - y_i \times \sum_{j=1}^{c} l_j) = -(K_i^{indv} - y_i)
\end{aligned}
$$

(5)

The last step holds since the soft-label vector sums to one.

By using the soft-label method, we can see that boundary pixels contribute not only to the boundary class but also to the foreground/background class in a weighted manner based on how close it is to the contour line.

### 3.2. The global boundary-sensitive kernel

By the nature of aligned portrait images, it is likely that some locations in the image, such as the upper corner pixels, should belong to the background with very high probabilities while some other locations, such as the middle bottom pixels, should belong to the foreground with high probabilities. These pixels should be more easily classified, while pixels in between should be harder to classify. We estimate a position sensitive prior from the training data.

We design a global boundary-sensitive kernel to guide the network to learn a better shape prediction specifically for portrait images. The global kernel is designed as follows. First, a mean mask $\overline{M}$ is computed using the average of all ground truth segmentation maps from the training samples. This generates a probability map where the value at each location indicates how likely the current location belongs to foreground/background. Second, Equation (6) is employed to generate the global boundary-sensitive kernel. All the values are normalized to range $[a, b]$. A larger value close to $b$ in the global kernel indicates that the current location has a higher probability to be boundary. In other words, this location should be more difficult for the network to classify. To force the network to focus more on the possible boundary locations, we weight the locations with their corresponding kernel values when updating the loss function. When performing the forward pass for one pixel location in one sample, we update the loss function as equations (7)

$$K^{global} = b - (1 - \frac{|\overline{M} - 0.5|}{0.5}) \times (b - a) \qquad (6)$$

$$\epsilon \mathrel{-}= K_s^{global} \times \sum_j \mathbb{1}(j = c) \times log(y_j) \qquad (7)$$

where $K_s^{global}$ denotes the global kernel value at the pixel location $s$. $g$ denotes the ground truth class label for the current pixel location. During back-propagation, the new gradient is computed as Equation (8):

$$
\begin{aligned}
\frac{\partial \epsilon}{\partial z_i} &= -K_s^{global} \times \sum_j \mathbb{1}(j = g) \times \frac{\partial log(y_j)}{\partial z_i} \\
&= -K_s^{global} \times (\frac{1}{y_i} \times \frac{\partial y_i}{\partial z_i}) \\
&= -K_s^{global} \times (\frac{1}{y_i} \times y_i \times (\mathbb{1}(i = g) - y_i)) \\
&= -K_s^{global} \times (\mathbb{1}(i = g) - y_i)
\end{aligned}
$$

(8)

From the new forward pass and back-propagation functions we can see that the pixels that are more likely to be located in the boundary (*e.g.*, the pixels lying within the brighter region in Figure 2 (e)) are weighted higher so that they contribute more to the loss. This guides the network to be more sensitive to the difficult locations.

### 3.3. The boundary-sensitive attribute classifier

Portrait attributes such as long/short hair play an important role in determining a portrait's shape. Training a network which is capable of classifying boundary-sensitive attributes will give more prior information to the system, which further makes the system more accurate and efficient on boundary prediction. Motivated by this idea, we train an attribute classifier jointly with the portrait segmentation network for multi-task learning. An example of how the

hair style attribute changes the boundary shape is shown in Figure 3.



Figure 3: An example of how boundary-sensitive attributes affect the portrait's shape: long hair vs. short hair.

To design the attribute classifier, the base layers from "conv1_1" to "pool5" are shared between the segmentation network and the classifier. Above this, for each channel, we add three more fully connected layers. The first two fully connected layers have 1024 neurons and are followed by a dropout layer and a ReLU layer. The last fully connected layer has two neurons for binary classification.

# 4. Experiments and results analysis

## 4.1. Training settings and evaluation settings

**Model details:** To train our portrait segmentation system, we fine-tune the DeepLabv2_ResNet101 model using the training set of the PFCN dataset. We will introduce this dataset in the next subsection. There are three ResNet branches in DeepLabv2. In each branch, 4 atrous convolution layers are added in parallel with dilation factors $[6, 12, 18, 24]$ and then summed together to produce the final feature map. Element-wise max operation is performed at the end over the three branches to produce the final prediction. To generate the individual kernel, we dilate the contour line to 10-pixels in width and label the dilated boundary using the soft-label method. We select the weight range in the global kernel as $[0.9, 1]$. Following PortraitFCN+, in addition to the three RGB channels, we add two normalized $x$ and $y$ position channels and one mean mask shape channel into the input. For more details please refer to [21]. At each iteration, a random patch of size $400 \times 400$ is cropped out from the original image and randomly flipped with probability $0.5$ for data augmentation. Then the input image is rescaled by factors of $[0.5, 0.75, 1.0]$ as the new input images to the three branches of the DeepLabv2 network. To train the attribute classifier, we label the training images into long/short hair classes. We use Stochastic Gradient Descent (SGD) with a learning rate of $2.5e^{-4}$ to train the model for 20K iterations without the attribute classifier. Then we decrease the learning rate by a factor of 10 and add

the attribute classifier to train the model for another 20K iterations. The whole network is built with the Caffe deep learning framework [13].

During testing, we ignore the boundary class and the attribute classifier. Only probabilities from foreground and background classes are used for segmentation.

**Mean IoU:** The standard mean Intersection-over-Union (IoU) metric is used to evaluate the segmentation performance. The mean IoU is computed as following.

$$\overline{\text{IoU}} = \frac{1}{N} \times \sum_{i}^{N} \frac{A_i^{seg} \cap A_i^{gt}}{A_i^{seg} \cup A_i^{gt}} \qquad (9)$$

where $A_i^{seg}$ and $A_i^{gt}$ represent the area of the segmentation results and the ground-truth label mask for the $i_{th}$ testing sample, respectively.

## 4.2. Results on the PFCN dataset

We evaluate the proposed method on the largest publicly available portrait segmentation dataset [21]. This dataset is collected from Flickr and manually labeled with variations in age, pose, appearance, background, lighting condition, hair style, accessory, etc. Most of the portrait images are captured by the frontal cameras of mobile phones. This dataset consists of 1800 portrait images which are split into 1500 training images and 300 testing images. All the images are scaled and cropped into size $800 \times 600$. In one portrait image, the pixels are labeled as either "foreground" or "background". We will refer to this dataset as PFCN dataset. Some sample images from the PFCN dataset are given in Figure 4.



Figure 4: Sample images from the PFCN dataset.

We compare with the state-of-the-art method reported on this dataset: PortraitFCN+ [21] and the

Figure 5: Sample images from COCO, COCO-Stuff, and Pascal VOC portrait datasets.

| Method | Mean IoU |
|---|---|
| Graph-cut | 80.0% |
| PortraitFCN+ | 95.9% |
| PortraitDeepLabv2 | 96.1% |
| BSN_AC (ours) | 96.2% |
| BSN_GK (ours) | 96.2% |
| BSN_IK (ours) | 96.5% |
| BSN (ours) | **96.7%** |

Table 1: Quantitative performance comparisons on the PFCN dataset.

DeepLabv2_ResNet101 fine-tuned model , which we will refer to as PortraitDeepLabv2. The PortraitDeepLabv2 model is fine-tuned using the same 6-channel training data as PortraitFCN+ and the same training settings as BSN. For ablation study, we report results of four models from our work: train with the attribute classifier only (BSN_AC), train with the global boundary-sensitive kernel only (BSN_GK), train with the individual boundary-sensitive kernel only (BSN_IK), and the all-in-one model (BSN). Our final model achieves the state-of-the-art mean IoU at 96.7%. The quantitative result comparison is given in Table 1. Result from graph-cut [22] is shown as the baseline.

| Dataset | Num. of Portrait |
|---|---|
| COCO portrait | 626 |
| COCO-Stuff portrait | 92 |
| PASCAL VOC portrait | 62 |

Table 2: Statistics of the three portrait datasets.

### 4.3. Evaluation on other datasets:

Since the performance on the PFCN dataset is pretty high and data for the boundary class is unbalanced compare to foreground/background, a good performance on boundary segmentation may only lead to marginal improvement in mean IoU on this dataset. Thus we further test our method on the portrait images collected from three more popular semantic segmentation datasets to evaluate the effectiveness of our boundary-sensitive techniques.

**COCO portrait:** We automatically collect all the portrait and portrait-like images from the COCO dataset. We run a face detector over the dataset and keep the images only containing one person where the face area covers at least 10% of the whole image. There are 626 images in total with ground truth segmentation maps. We will refer to this dataset as COCO portrait. COCO portrait is more challenging than the PFCN data in various ways such as large pose variations, large occlusions, unlabeled individuals appear on the background, large portion of background, different kinds of accessories, etc.

**COCO-Stuff portrait:** The COCO-Stuff dataset augments the COCO dataset with refined pixel-level stuff annotations on $10K$ images. We collect 92 portrait and portrait-like images from this dataset. The quality of images in this dataset are same as COCO portrait. We will refer to this dataset as COCO-Stuff portrait.

**Pascal VOC portrait:** We use the same method to collect portrait and portrait-like images from the Pascal VOC 2007, 2008, and 2012 datasets. Due to the lack of ground truth segmentation maps on this dataset, 62 images are collected. The images in this dataset are also challenging and unconstrained. We will refer to this dataset as PASCAL VOC portrait. Some sample images from the three datasets are illustrated in Figure 5 and the statistics are given in Table 2.

To test the generalization capability of our model, we directly test on these three datasets without fine-tuning. We achieve 77.7% mean IoU, 72.0% mean IoU, and 75.6% mean IoU on COCO portrait, COCO-stuff portrait, and PASCAL VOC portrait, respectively. We significantly outperform PortraitFCN+ on all the three datasets. The result comparisons are illustrated in Table 3. Since the DeepLabv2 model is trained on these dataset, we can not compare with it directly.

(a) Input image     (b) Ground truth     (c) PortraitFCN+     (d) PortraitDeepLabv2     (e) Our final model

Figure 6: Result visualizations of three challenging examples. The first row shows contains confusing objects in the background; the second row includes multiple people in the background; in the third row the background color is close to the foreground.

| Method | COCO | COCO-Stuff | PASCAL VOC |
|---|---|---|---|
| PortraitFCN+ | 68.6% | 60.8% | 59.5% |
| BSN (ours) | **77.7%** | **72.0%** | **75.6%** |

Table 3: Quantitative performance comparisons on COCO portrait, COCO-Stuff portrait and Pascal VOC portrait datasets.
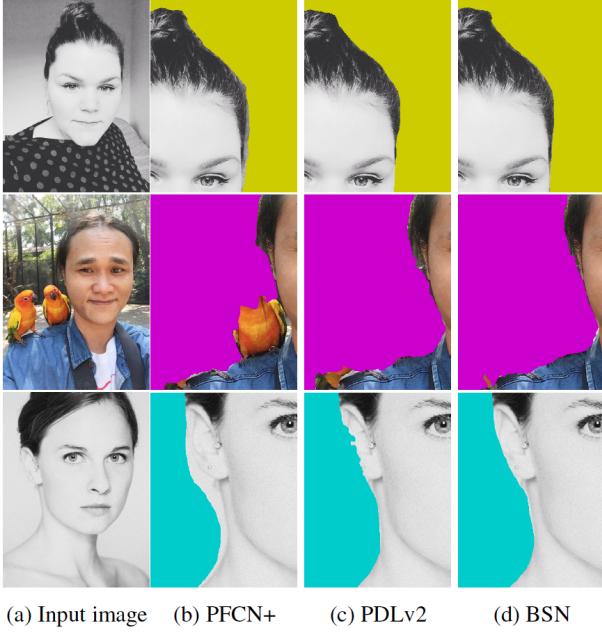
## 4.4. Result analysis

### 4.4.1 Results visualization on challenging scenarios

We visualize the overall performance of our BSN model compared to DeepLabv2 and PortraitFCN+ using three challenging scenarios: confusing objects in the background, multiple people appear in the image, and the background color theme is close to the foreground. Figure 6 shows that our model is more accurate and robust than other methods even under challenging conditions.

### 4.4.2 Accurate boundary segmentation

Our method also delivers more precise boundary predictions thanks to its novel boundary-sensitive segmentation

techniques. Figure 7 shows the comparison of our method with DeepLabv2 and PortraitFCN+ in three challenging scenarios: hair segmentation, accessory segmentation and ear segmentation. Results reveal that while other methods have difficulty in segmenting accessories and small body parts, our method can provide a smooth and accurate segmentation.



| (a) Input image | (b) PFCN+ | (c) PDLv2 | (d) BSN |

Figure 7: Boundary segmentation comparisons. The first column are the original images. The three subsequent columns represent the results from the PortraitFCN+ method, the fine-tuned DeepLabv2 model with the attribute classifier, and our final model (magnified for best viewing).

### 4.4.3 Generating trimap for image matting

Since our method can deliver an accurate boundary prediction, it is a natural extension to generate trimaps for image matting models. After performing segmentation, we use the same technique during training to dilate the boundary pixels to 10-pixels in width. Several examples are shown in Figure 8.

### 4.4.4 Applications of portrait segmentation

Portrait segmentation has been widely used in various image processing applications such as background replacement, depth of field, augmented reality, image cartoonization, etc. We show some applications in Figure 9.
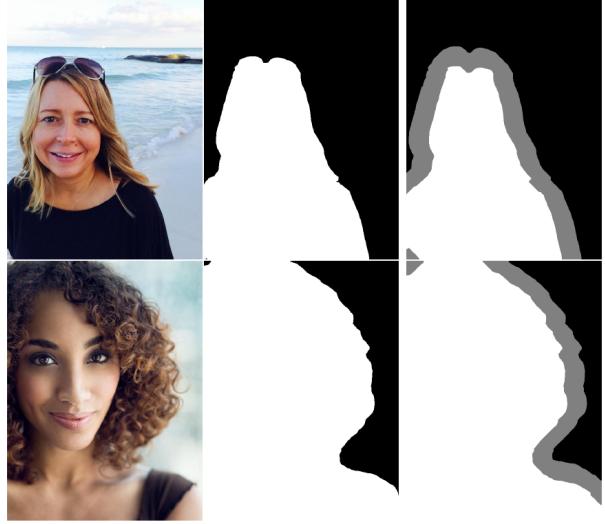


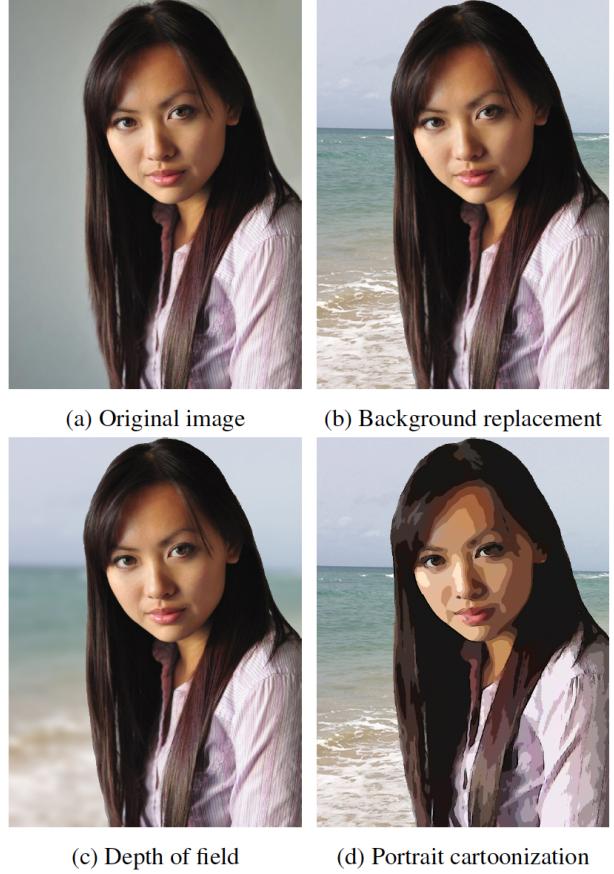Figure 8: Trimaps generated from our segmentation maps.



| (a) Original image | (b) Background replacement |
| (c) Depth of field | (d) Portrait cartoonization |

Figure 9: Some applications of portrait segmentation.

# 5. Conclusion and discussion

We present a boundary-sensitive portrait segmentation system. Two boundary-sensitive kernels are introduced into the loss function. One gives more boundary information for each individual image and one governs the overall shape of portrait prediction. An attribute classifier is trained jointly with the segmentation network to reinforce the training process. Experiments are conducted on the largest publicly available portrait segmentation dataset as well as portrait images collected from other three popular semantic segmentation datasets. We outperform the previous state-of-the-arts in both quantitative performance and visual performance.

For future work, we would like to extend our boundary-sensitive methods to general semantic segmentation problem and explore more semantic attributes to reinforce the training process.

## References

[1] V. Badrinarayanan, A. Kendall, and R. Cipolla. Segnet: A deep convolutional encoder-decoder architecture for image segmentation. *CoRR*, abs/1511.00561, 2015.

[2] L. Breiman. Random forests. *Machine Learning*, 45(1):5–32, Oct 2001.

[3] H. Caesar, J. R. R. Uijlings, and V. Ferrari. Coco-stuff: Thing and stuff classes in context. *CoRR*, abs/1612.03716, 2016.

[4] J. Canny. A computational approach to edge detection. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, PAMI-8(6):679–698, Nov 1986.

[5] B. Chen, Z. Yang, S. Huang, X. Du, Z. Cui, J. Bhimani, and N. Mi. Cyber-physical system enabled nearby traffic flow modelling for autonomous vehicles. *36th IEEE International Performance Computing and Communications Conference, Special Session on Cyber Physical Systems: Security, Computing, and Performance (IPCCC-CPS). IEEE*, 2017.

[6] L. Chen, G. Papandreou, F. Schroff, and H. Adam. Rethinking atrous convolution for semantic image segmentation. *CoRR*, abs/1706.05587, 2017.

[7] L.-C. Chen, G. Papandreou, I. Kokkinos, K. Murphy, and A. L. Yuille. Deeplab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected crfs. *arXiv:1606.00915*, 2016.

[8] X. Du, W. Abdalmageed, and D. Doermann. Large-scale signature matching using multi-stage hashing. In *2013 12th International Conference on Document Analysis and Recognition*, pages 976–980, Aug 2013.

[9] X. Du, D. Doermann, and W. AbdAlmageed. A graphical model approach for matching partial signatures. In *2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1465–1472, June 2015.

[10] X. Du, M. El-Khamy, J. Lee, and L. Davis. Fused dnn: A deep neural network fusion approach to fast and robust pedestrian detection. In *2017 IEEE Winter Conference on Applications of Computer Vision (WACV)*, pages 953–961, March 2017.

[11] M. Everingham, L. Van Gool, C. K. I. Williams, J. Winn, and A. Zisserman. The pascal visual object classes (voc) challenge. *International Journal of Computer Vision*, 88(2):303–338, June 2010.

[12] M. A. Hearst, S. T. Dumais, E. Osuna, J. Platt, and B. Scholkopf. Support vector machines. *IEEE Intelligent Systems and their Applications*, 13(4):18–28, July 1998.

[13] Y. Jia, E. Shelhamer, J. Donahue, S. Karayev, J. Long, R. Girshick, S. Guadarrama, and T. Darrell. Caffe: Convolutional architecture for fast feature embedding. *arXiv preprint arXiv:1408.5093*, 2014.

[14] D. Li, T. Salonidis, N. V. Desai, and M. C. Chuah. Deepcham: Collaborative edge-mediated adaptive deep learning for mobile object recognition. In *2016 IEEE/ACM Symposium on Edge Computing (SEC)*, pages 64–76, Oct 2016.

[15] D. Li, X. Wang, and D. Kong. Deeprebirth: Accelerating deep neural network execution on mobile devices. *CoRR*, abs/1708.04728, 2017.

[16] G. Lin, A. Milan, C. Shen, and I. D. Reid. Refinenet: Multi-path refinement networks for high-resolution semantic segmentation. *CoRR*, abs/1611.06612, 2016.

[17] T. Lin, M. Maire, S. J. Belongie, L. D. Bourdev, R. B. Girshick, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick. Microsoft COCO: common objects in context. *CoRR*, abs/1405.0312, 2014.

[18] J. Long, E. Shelhamer, and T. Darrell. Fully convolutional networks for semantic segmentation. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2015.

[19] C. Peng, X. Zhang, G. Yu, G. Luo, and J. Sun. Large kernel matters - improve semantic segmentation by global convolutional network. *CoRR*, abs/1703.02719, 2017.

[20] X. Shen, H. Gao, X. Tao, C. Zhou, and J. Jia. High-quality correspondence and segmentation estimation for dual-lens smart-phone portraits. *arXiv preprint arXiv:1704.02205*, 2017.

[21] X. Shen, A. Hertzmann, J. Jia, , S. Paris, B. Price, E. Shechtman, and I. Sachs. Automatic portrait segmentation for image stylization. *Computer Graphics Forum*, 35(3):93–102, 2016.

[22] J. Shi and J. Malik. Normalized cuts and image segmentation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 22(8):888–905, Aug 2000.

[23] F. Yu and V. Koltun. Multi-scale context aggregation by dilated convolutions. In *ICLR*, 2016.

[24] H. Zhao, J. Shi, X. Qi, X. Wang, and J. Jia. Pyramid scene parsing network. *CoRR*, abs/1612.01105, 2016.