

# Inductive Guided Filter: Real-time Deep Image Matting with Weakly Annotated Masks on Mobile Devices

Yaoyi Li<sup>1\*</sup>, Jianfu Zhang<sup>1</sup>, Weijie Zhao<sup>2</sup> and Hongtao Lu<sup>1</sup>

<sup>1</sup>Shanghai Jiao Tong University

<sup>2</sup>Versa

{dsamuel, c.sis}@sjtu.edu.cn, weijie.zhao@versa-ai.com, htlu@sjtu.edu.cn

## Abstract

Recently, significant progress has been achieved in deep image matting. Most of the classical image matting methods are time-consuming and require an ideal trimap which is difficult to attain in practice. A high efficient image matting method based on a weakly annotated mask is in demand for mobile applications. In this paper, we propose a novel method based on Deep Learning and Guided Filter, called Inductive Guided Filter, which can tackle the real-time general image matting task on mobile devices. We design a lightweight hourglass network to parameterize the original Guided Filter method that takes an image and a weakly annotated mask as input. Further, the use of Gabor loss is proposed for training networks for complicated textures in image matting. Moreover, we create an image matting dataset MAT-2793 with a variety of foreground objects. Experimental results demonstrate that our proposed method massively reduces running time with robust accuracy.

## 1 Introduction

Image matting is a fundamental problem in computer vision, which models an image as a linear combination of a foreground and a background image:

$$I_i = \alpha_i F_i + (1 - \alpha_i) B_i, \alpha_i \in [0, 1], \quad (1)$$

where  $\alpha_i$  is the linear coefficient at a pixel position  $i$ ,  $F_i$  for the foreground pixel at  $i$  and  $B_i$  for the corresponding background pixel. Image matting task is more than a high-accurate segmentation and proposed for the natural image decomposition, which takes transparency into consideration. The generated alpha matte can highly reduce the workload and special requirement of image or video editing for advertisement, design, Vlog, film and so on. With rapid growth of the users who are using mobile devices to edit images or videos, a fast and accurate model which can enhance user experience is in high demand.

However, most of the classical methods [Levin *et al.*, 2008; Chen *et al.*, 2013; Cho *et al.*, 2016; Xu *et al.*, 2017] are time-

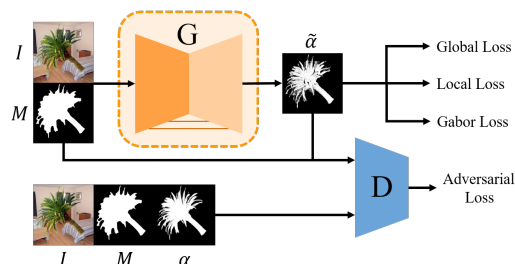


Figure 1: An illustration of our proposed method.

consuming which are not capable of running on mobile devices for real-time.

Another obstacle for image matting on mobile devices is that the classical methods are sensitive with the input mask. Most of the time we can only obtain weakly annotated masks on mobile devices due to the limitation of time latency and computing ability. We coin the term “weakly annotated mask” to describe a noisy or inexact mask which gives some inaccurate annotations for foreground and background. A weakly annotated mask can be an output binary image of a segmentation method, a thresholded depth map or an inexact annotated mask from user interaction. It contrasts conventional trimap which has an accurate annotation but consumes much time to compute. Reducing the quality of input masks or trimaps can massively degrade performance for the classical methods.

In this paper, we introduce a new deep learning framework that is specifically tailored for mobile devices by significantly reducing network parameters while retaining compatible accuracy. Compared with the classical methods that are highly dependent on the quality of trimap, our proposed model is robust with weakly annotated masks. We build our neural network for image matting as an image-to-image translation manner. With the help of GAN framework and three other different losses, we can generate highly detailed image mattes with a tiny network.

Our main contributions in this paper are three-fold and can be summarized as followings:

- We design a novel real-time framework for the weakly annotated image matting task, dubbed Inductive Guided Filter. We are the first to introduce the combination of the deep convolutional neural networks and Guided Filter into the image matting task.

\*Work done as an intern at Versa

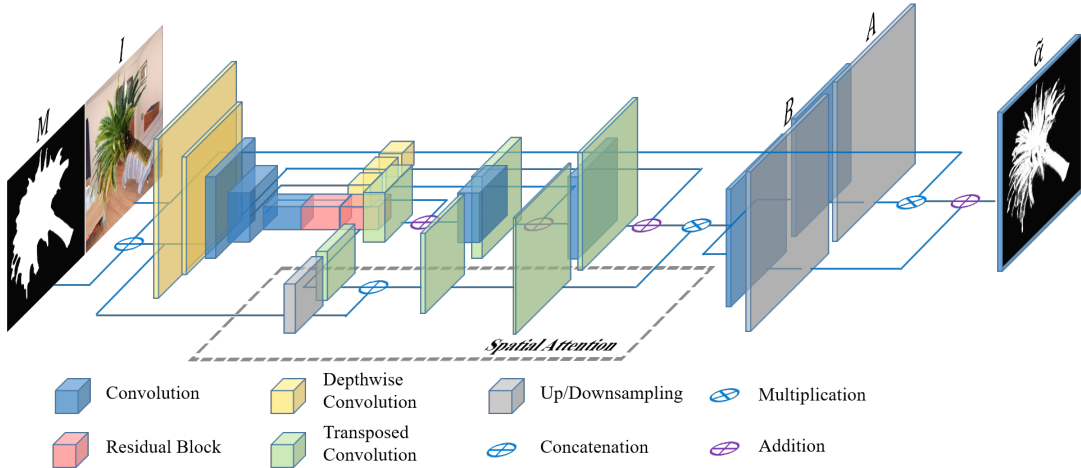


Figure 2: The overview of our generator. Linear coefficients  $A$  and  $B$  are generated from two branches sharing a same light weight Hourglass backbone. The right-most two upsampling layers both have a factor 4 which mimic Fast Guided Filter [He and Sun, 2015] for acceleration.

- We propose a Gabor loss based on a bundle of Gabor filters to extract more comprehensive high-frequency features. To the best of our knowledge, no prior works have introduced such a loss function.
- We further create an image matting dataset with 2793 foregrounds for the training of deep image matting called MAT-2793, which is the current biggest dataset for image matting to our knowledge. We evaluate our proposed method on MAT-2793 and Adobe Composition-1k testing dataset. Compared with the classical methods, our proposed method can achieve robust image matting effectively and efficiently.

## 2 Related Works

In this section, we review previous works on deep image matting and Guided Filter [He *et al.*, 2010] which are highly related to our method.

**Deep Image Matting** Many recent deep image matting approaches can be broadly categorized as general deep image matting approaches and ad hoc deep image matting approaches that are tailored for specific tasks.

General deep image matting approaches attempt to predict the alpha matte given any natural image and the ideal trimap. Cho *et al.* were the first to introduced deep learning into image matting task [2016]. Their DCNN matting aimed to learn convolutional neural networks to combine the output of different classical image matting approaches. The Adobe Deep Image Matting (Adobe DIM) proposed in [Xu *et al.*, 2017] was the first end-to-end deep image matting method, which significantly outperforms the conventional methods. Lutz *et al.* further employed the generative adversarial networks (GAN) in their proposed AlphaGAN [Lutz *et al.*, 2018].

Some deep image matting methods are specialized in practical application scenarios like portrait matting. Shen *et al.* proposed a portrait matting method [2016] with a deep network for trimap generation followed by a closed-form matting [Levin *et al.*, 2008], which can propagate gradients from closed-form matting to the neural network. Chen

*et al.* proposed a Semantic Human Matting [Chen *et al.*, 2018b] which incorporated person segmentation and matting in an end-to-end manner. More lightweight deep image matting methods were proposed for portrait and hair matting on mobile devices [Zhu *et al.*, 2017; Levinshtein *et al.*, 2018; Chen *et al.*, 2019].

**Guided Filter** Guided Filter was proposed in [He *et al.*, 2010] as an edge-preserving smoothing operator that had a theoretical connection with the matting Laplacian matrix. Deep Guided Filter [Wu *et al.*, 2018] applied the Guided Filter to the output image of an image-to-image translation network as a super-resolution block and propagates the gradient through Guided Filter to the low-resolution output. Zhu *et al.* proposed a fast portrait matting method with a feathering block inspired by Guided Filter [Zhu *et al.*, 2017], which will be elaborated in Section 3.1.

## 3 Method

We attempt to build an extremely high efficient image matting neural network which takes a weakly annotated mask as input. To this end, we employ the idea of linear model assumption in Guided Filter [He *et al.*, 2010], which is robust and efficient in feathering tasks.

Following [Lutz *et al.*, 2018], we adopt Generative Adversarial Network (GAN) [Goodfellow *et al.*, 2014] architecture to our model. The coarse architecture of our method is illustrated in Figure 1 and details of generator in Figure 2.

### 3.1 Inductive Guided Filter Formulation

In Guided Filter [He *et al.*, 2010], the basic assumption is that the output alpha is a linear transform of guidance image  $I$  in a small window  $\omega_k$  centered at pixel  $k$ :

$$\alpha_i = A_k I_i + B_k, \forall i \in \omega_k, \quad (2)$$

in which  $A_k$  and  $B_k$  are linear coefficients to be optimized. The optimization objective is to minimize the difference between output  $\alpha_i$  and the corresponding pixel  $M_i$  on input weakly annotated mask  $M$  with the regularization on  $A_k$ .

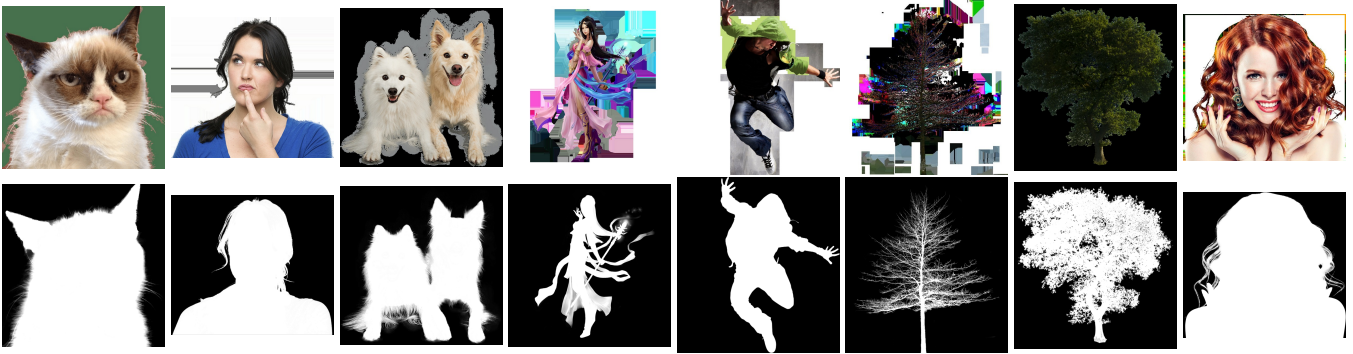


Figure 3: Some foreground object samples from the MAT-2793 dataset.

In the image matting setting, Guided Filter solves an optimization problem for each image and mask to generate a linear transformation from input image  $I$  to matte estimation  $\alpha$  which is as close to input mask  $M$  as possible.

Although Guided Filter is a fast and effective method for weakly annotated image matting task, it is limited by the constraint that the difference between optimal alpha matte and weakly annotated mask should be small enough. Empirically, the mask from a semantic segmentation method or user interaction will have a relatively large difference from the ground-truth alpha (see our testing set samples in Figure 9).

Different from Guided Filter, our method attempts to deal with a supervised learning task instead of an optimization problem. We abandon the objective function and remove the constraint on the difference between matte estimation and mask. An inductive model based on the linear transform assumption is built to leverage the ground truth information in an image matting dataset. We formula the Inductive Guided Filter as

$$\alpha = \phi_A(I, M) \circ I + \phi_B(I, M), \quad (3)$$

where  $\circ$  denotes Hadamard product and we parameterize  $A$  and  $B$  in Guided Filter by neural networks  $\phi_A(I, M)$  and  $\phi_B(I, M)$ . Networks  $\phi_A$  and  $\phi_B$  take image  $I$  and weakly annotated mask  $M$  as input and share backbone parameters (see last two layers in Figure 2 for details). The optimization objective of Inductive Guided Filter is to minimize the difference between the alpha matte prediction and ground truth.

For any image and mask,  $\phi_A$  and  $\phi_B$  can generate the specific coefficients  $A$  and  $B$ , to build a linear transform model for alpha matte.

A similar idea was also mentioned in [Zhu *et al.*, 2017]. The main difference between our method and theirs is that the authors of [Zhu *et al.*, 2017] formulated their feathering block based on the closed-form solution of Guided Filter:

$$\alpha_i = A_k M_i^F + B_k M_i^B + C_k, \forall i \in \omega_k, \quad (4)$$

where  $M^F$  and  $M^B$  are masks for foreground and background.  $A_k, B_k$  and  $C_k$  are coefficients that parameterized by a neural network like  $\phi(\cdot)$  in our method. From Equation (4) we can derive that the output of their feathering block will only preserve the edge and gradient of the mask instead of the input image. It can be seen as an attention map on the mask. Consequently, a weakly annotated mask may lead to performance degradation. On the contrary, Inductive Guided Filter can be regarded as an attention map on the original input

image, which is the same as the linear transform in Guided Filter. It is more robust to the noise in a mask.

### 3.2 Generator

The generator consists of a lightweight Hourglass backbone, spatial attention mechanism, and a linear transformation.

We build a lightweight Hourglass backbone following the structure of U-Net [Ronneberger *et al.*, 2015] and the Hourglass module in Stacked Hourglass Networks [Newell *et al.*, 2016] which prove to be effective to preserve low-level information from high-resolution features. Only two residual blocks are involved in the bottleneck. Moreover, depthwise convolution, which is widely used in lightweight deep neural networks [Chollet, 2017; Sandler *et al.*, 2018], is adopted in the first two convolution layers and the shortcut connections. We only introduce depthwise blocks to the layers that have high-resolution feature maps as [Nekrasov *et al.*, 2018] did. There is no  $1 \times 1$  convolution layer between the first two depthwise convolutions. We regard them as adaptive down-samplings. All efforts aim to reduce inference latency.

Spatial attention [Xu *et al.*, 2015; Chen *et al.*, 2017a] has shown effectiveness in various computer vision tasks. Some previous deep image matting methods adopted spatial attention in their structures, which conduces to a decent matting performance [Chen *et al.*, 2018b; Zhu *et al.*, 2017]. As for our attention mechanism, we fuse the feature from input and bottleneck to compute an attention map which is applied to the high-resolution features in the decoder.

Besides the adversarial loss, we impose three loss functions on the generator: global loss, local loss and Gabor loss. The real alpha matte and predicted alpha matte in triplet input are denoted by  $\alpha$  and  $\tilde{\alpha}$ .

**Global Loss** To supervise the alpha matte prediction, we leverage the global loss which is a L1 loss between the ground truth alpha and the estimated alpha,

$$\mathcal{L}_g = \|\alpha - \tilde{\alpha}\|_1, \quad (5)$$

**Local Loss** When we are training a network for image matting, in contrast to global loss, we would like objective function to focus more on the boundary of foreground object. Local loss is a weighted reconstruction based on a difference function  $\Delta(\alpha, M) = \delta(|\alpha - M| > \epsilon)$ . The difference function yields a binary boundary map, in which 1 for the same values in ground truth and mask and 0 for the other pixels.  $\delta$



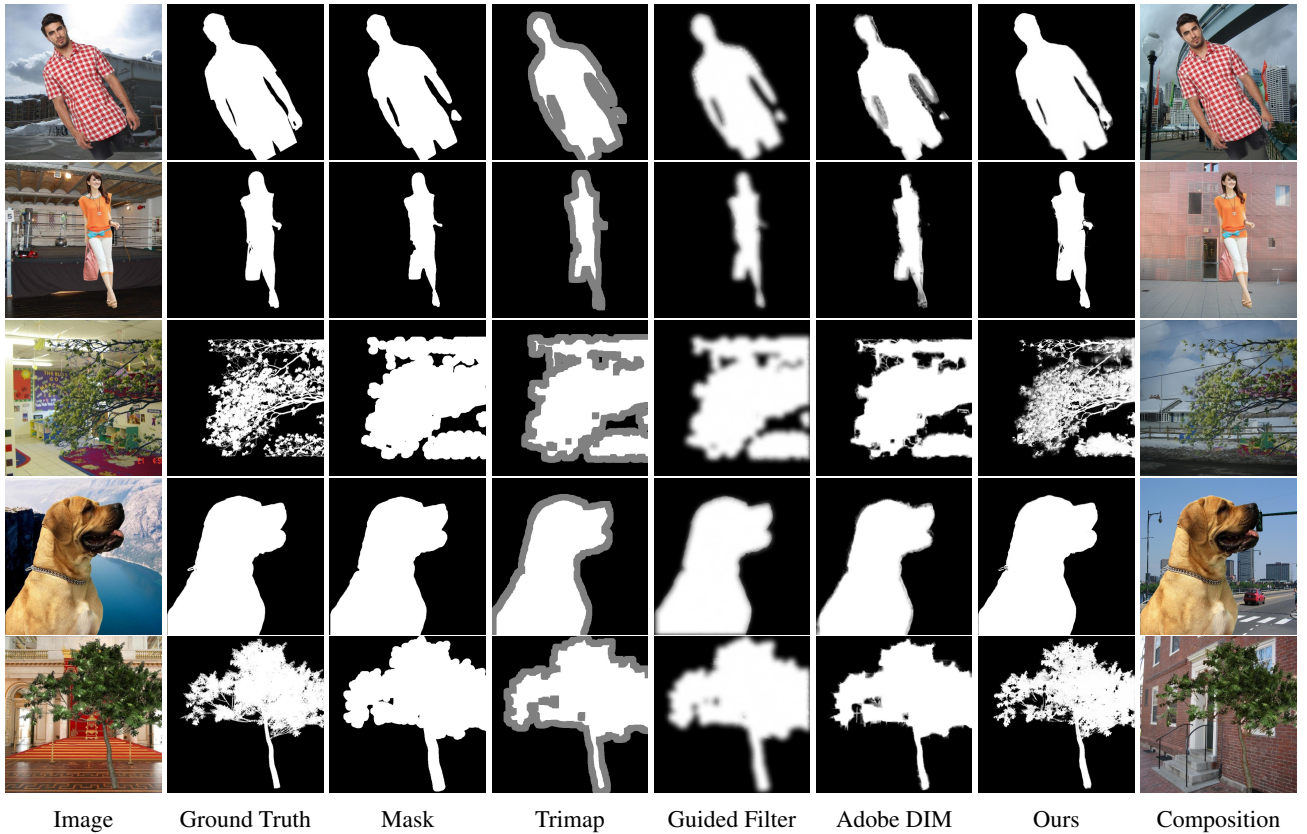


Figure 4: The visual comparison results on MAT-2793 testing set. The trimap is only for Adobe DIM. The composition is composed with our results and random backgrounds.

function enforces the small differences between below  $\epsilon$  are neglected. We use  $\epsilon = 0.01$  in the loss function. The local loss function is written as

$$\mathcal{L}_l = \|\Delta(\alpha, M) \circ (\alpha - \tilde{\alpha})\|_1, \quad (6)$$

and in practice we apply an additional morphological dilation with a  $7 \times 7$  kernel to the difference function for a larger boundary area.

**Gabor Loss** Perceptual loss proposed in [Johnson *et al.*, 2016] dramatically improves the visual quality of predictions in supervised image transformation tasks. It provides a series of supervisions from high-frequency features to semantic level features. Perceptual loss utilizes a VGG network pre-trained on RGB color images with specific classes. However, the alpha matte is a gray-scale image. Therefore, we design a Gabor loss to resemble the perceptual loss in our case. Gabor loss replaces the pre-trained multi-layer kernels in perceptual loss with a set of single-layer Gabor filters to extract high-frequency features.

Gabor filter was introduced into the neural network as a kernel or initialization in some previous works [Ouyang and Wang, 2013; Luan *et al.*, 2018] due to its comparability to the kernels from shallow layers. Thus, we define the Gabor loss by

$$\mathcal{L}_{gb} = \sum_{\phi_{gb} \in \Phi} \|\phi_{gb}(\alpha) - \phi_{gb}(\tilde{\alpha})\|_2^2, \quad (7)$$

where function  $\phi_{gb}(\cdot)$  denotes the convolution with Gabor filter,  $\Phi$  is the set of different Gabor filters. In our training, we design 16 different  $7 \times 7$  Gabor filters with 16 orientations in  $\Phi$ . All of the filters have wavelength  $\lambda = 5$ , spatial aspect ratio  $\gamma = 0.5$  and standard deviation  $\sigma = 0.5$ . We have also tried a larger set with different wavelengths and standard deviations, and no additional remarkable benefits were manifested in the training.

Some deep image matting method introduces gradient loss into their objective functions with a similar motivation [Levinshtein *et al.*, 2018; Chen *et al.*, 2019]. Gradient loss minimizes the difference between the image gradient of the input image and predicted alpha, globally or locally. Comparing with image gradient, Gabor loss is an extended version with the capability to extract more comprehensive features, especially for alpha matte which is rich in the high-frequency component.

### 3.3 Discriminator

Lutz *et al.* first introduced GAN into their proposed AlphaGAN [2018]. In AlphaGAN, the discriminator takes a trimap and a newly composited image from predicted alpha matte as its input. Since image matting does not focus on semantic information, it is ambiguous to judge whether a composited image is real or not. A composited image may have a high fidelity when it is generated from an incorrect or a partial alpha matte.



Methods	MSE	SAD	Grad. ( $\times 10^3$ )	Conn. ( $\times 10^3$ )
Guided Filter	0.101	9.57	15.53	5.53
Adobe DIM	0.148	8.54	16.53	4.42
Our method	<b>0.028</b>	<b>2.51</b>	<b>6.45</b>	<b>1.51</b>

Table 1: The quantitative results on the MAT-2793 testing set

To overcome this, we feed the discriminator with a conditional triplet input which consists of an original image, a weakly annotated mask and an alpha matte, analogous to some methods with pair input [Chen *et al.*, 2017b; Hu *et al.*, 2018]. Given a triplet input, the discriminator can predict the self-consistency of an input. Concretely, the critic is designed to predict whether an estimated alpha matte is correct conditioned on the input image and mask.

**Adversarial Loss** We employ LSGAN [Mao *et al.*, 2017] with gradient penalty [Gulrajani *et al.*, 2017] as the adversarial loss. The loss function is defined as

$$\begin{aligned} \mathcal{L}_{adv} &= \mathcal{L}_D + \mathcal{L}_G + \lambda_{gp} \mathbb{E}_{\hat{\alpha}} [(\|\nabla_{\hat{\alpha}} D(\hat{\alpha}|I, M)\|_2 - 1)^2] \\ \mathcal{L}_D &= \mathbb{E}_{\alpha} [(D(\alpha|I, M) - 1)^2] + \mathbb{E}_{\tilde{\alpha}} [D(\tilde{\alpha}|I, M)^2] \\ \mathcal{L}_G &= -\mathbb{E}_{\tilde{\alpha}} [(D(\tilde{\alpha}|I, M) - 1)^2], \end{aligned} \quad (8)$$

where  $\hat{\alpha}$  is a convex combination of  $\alpha$  and  $\tilde{\alpha}$  with a random coefficient sampled from uniform distribution. We use  $\lambda_{gp} = 10$  in our training.

### 3.4 Full Loss and Implement Details

The full loss function in Inductive Guided Filter is

$$\mathcal{L} = \lambda_g \mathcal{L}_g + \lambda_l \mathcal{L}_l + \lambda_{gb} \mathcal{L}_{gb} + \lambda_{adv} \mathcal{L}_{adv}, \quad (9)$$

in which we use  $\lambda_g = 10$ ,  $\lambda_l = 1$ ,  $\lambda_{gb} = 200$  and  $\lambda_{adv} = 1$ .

We leverage PatchGAN [Isola *et al.*, 2017], which is capable of discriminating the fidelity of local patches, to drive the attention of critic to detailed textures. Spectral Normalization [Miyato *et al.*, 2018] has shown an appealing superiority in the training of GAN. We apply Spectral Normalization layers in our discriminator as well as the Batch Normalization [Ioffe and Szegedy, 2015].

We incorporate training tricks: learning rate warm-up and cosine learning rate decay, following [Xie *et al.*, 2018]. Adam optimizer [Kingma and Ba, 2014] is adopted with  $\beta_1 = 0.5$ ,  $\beta_2 = 0.999$  and initial learning rate 0.0001 for both generator and discriminator.

The size of input images, masks, and output alpha matte are  $512 \times 512$ . Out-channels of the first 5 convolution layers in the generator are 4, 4, 8, 16, 32, and all 5 convolution layers have a stride 2. The whole network has 460, 456 trainable parameters with 400, 990 in discriminator and 59, 466 in the generator.

## 4 Experiments

In this section, we evaluate our method on two datasets, our MAT-2793 and Adobe Composition-1k [Xu *et al.*, 2017]. We compare the proposed method with the state-of-the-art Adobe Deep Image Matting (Adobe DIM) [Xu *et al.*, 2017] as well

Methods	MSE	SAD	Grad. ( $\times 10^3$ )	Conn. ( $\times 10^3$ )
Guided Filter	0.130	176.2	150.5	111.2
Adobe DIM	0.243	105.2	55.0	56.1
Our method	<b>0.064</b>	<b>46.2</b>	<b>42.1</b>	<b>35.7</b>

Table 2: The quantitative results on the Adobe Composition-1k with weakly annotations.

as Guided Filter [He *et al.*, 2010] quantitatively and qualitatively. Extensive experiments are conducted to make an efficiency comparison with some ad hoc real-time deep matting methods.

### 4.1 Dataset MAT-2793

We create an image matting dataset MAT-2793, which contains 2793 foreground objects, to tackle our weakly annotated image matting task. Most of the foregrounds and corresponding alpha mattes are gathered from the Internet as transparent PNG images especially from some free clipart website. Therefore, small parts of the foreground objects are not real-world images. We split our dataset into a training set with 2504 samples and a testing set with 289 samples.

In our experiments, we select 5360 high-resolution background images from the SUN dataset [Xiao *et al.*, 2010] and ADE20k [Zhou *et al.*, 2017]. We composite each object onto random backgrounds with a flipping and 11 different rotation angles to synthesis 22 different images. To generate weakly annotated masks, we first treat the alpha larger than 0.5 as foreground. Then we apply dilation and erosion together to the foreground mask with random kernel size both from  $5 \times 5$  to  $30 \times 30$ . If the area of the foreground in a generated mask is less than half of the foreground in alpha matte, this training sample will be abandoned. We resize the shorter side of composited images, corresponding alpha mattes, and masks to 600 pixels. For data augmentation, we randomly crop the image and resize the crop to  $512 \times 512$ . Moreover, we randomly change the hue of images in the training phase. In testing, we make a  $600 \times 600$  center crop and resize it to  $512 \times 512$ .

### 4.2 Results on MAT-2793 testing set

For a fair comparison between our method, Guided Filter and Adobe DIM\*, we generate a trimap from each mask as the input of Adobe DIM. We apply a dilation along with an erosion on each mask both with a  $20 \times 20$  kernel to create an unknown area. The errors are only computed in the unknown area. We implement the fast version of Guided Filter [He and Sun, 2015] with downsampling factor 4 in all of our experiments.

We follow the image matting evaluation metrics suggested in [Rhemann *et al.*, 2009]. We report the quantitative results under SAD, MSE, Gradient errors and Connectivity errors in Table 5. Furthermore, we display qualitative visual results in Figure 5. Results show that our proposed method is robust to noises in input masks and capable of capturing detailed texture in composited images.

\*We use the implementation from <https://github.com/foamliu/Deep-Image-Matting>

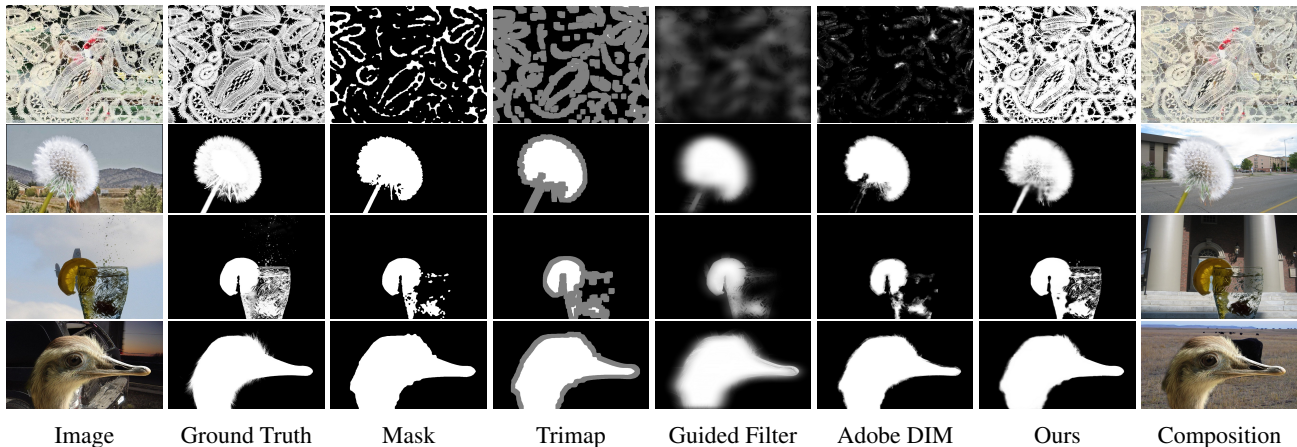


Figure 5: The visual comparison results on Adobe Composition-1k. The trimap is only for Adobe DIM. The composition is composed with our results and random backgrounds.

Methods	Device	Batch Size	Time (ms)
BANet-64 <sup>†</sup> (512x512)	1080 Ti	1	23.3
LDN+FB <sup>†</sup> (128x128)	TITAN X	1	13
Adobe DIM (512x512)	V100	1	51.1
Ours w/ I/O (512x512)	V100	1	3.48
		256	1.46
Ours w/o I/O (512x512)	V100	1	2.19
		256	0.18

Table 3: The results of speed evaluation on Nvidia GPU devices. We also display the speed of our method with or without GPU I/O.

### 4.3 Results on Adobe Composition-1k

Adobe Composition-1k is an image matting testing set with 1000 images and 50 unique foreground objects proposed in [Xu *et al.*, 2017].

In experiments on Adobe Composition-1k, we generate the weakly annotated mask in the same way as we did in Section 4.1. To generate weakly trimaps for DIM, considering the large image size in Adobe Composition-1k, we apply dilation and erosion on masks with a  $50 \times 50$  kernel. We then resize all the images to  $512 \times 512$  in our experiments.

We illustrate the visual comparison results in Figure 5 and the quantitative results in Table 2. All the errors are calculated in the unknown area of our generated weakly trimaps. The experiment results also show that our proposed method is capable of handling the weakly input.

### 4.4 Speed Evaluation

We evaluate the efficiency of the proposed method on different platforms including Nvidia GPU, computer CPU, and mobile devices. Some real-time deep image matting methods for portrait (LDN+FB [Zhu *et al.*, 2017], BANet-64 [Chen *et al.*, 2019] or hair (HairMatteNet [Levinshtein *et al.*, 2018]) are also included in this evaluation. We report the performance of these methods from the data in their original papers.

<sup>†</sup>Data from their original papers

Methods	Device	Time (ms)
HairMatteNet <sup>†</sup> (224x224)	iPad Pro GPU	30
LDN+FB <sup>†</sup> (128x128)	Core E5-2660	38
	Adreno 530	62
Guided Filter (512x512)	Core i7-7700	62.0
Adobe DIM (512x512)	Core i7-7700	4003.2
Ours (512x512)	Core i7-7700	13.0
	iPhone Xs	15.7
	iPhone X	22.3
	iPhone SE	25.9

Table 4: The results of speed evaluation on CPU and mobile devices.

We demonstrate the inference speed on Nvidia GPU devices in Table 3 and speed on CPU or difference mobile devices in Table 4. We deploy our model on different iPhone devices via the Apple CoreML framework. We can notice that taking a  $512 \times 512$  image as input, our method can run at 5000+ FPS on a single Nvidia V100 GPU with batch size 256 and achieve real-time performance on an iPhone SE in production in 2016.

## 5 Conclusion

In this paper, we propose a extremely efficient method for weakly annotated image matting on mobile devices, dubbed Inductive Guided Filter. A lightweight hourglass backbone and a novel Gabor loss are leveraged in the model. We also create a large image matting dataset MAT-2793. Evaluation on two testing datasets demonstrates that our proposed model is robust to the weakly annotated input mask and is competent to extract texture details in an image matting task.

## References

[Chen *et al.*, 2013] Qifeng Chen, Dingzeyu Li, and Chi-Keung Tang. Knn matting. *IEEE TPAMI*, 2013.

- [Chen *et al.*, 2017a] Long Chen, Hanwang Zhang, Jun Xiao, Liqiang Nie, Jian Shao, Wei Liu, and Tat-Seng Chua. Sca-cnn: Spatial and channel-wise attention in convolutional networks for image captioning. In *CVPR*, 2017.
- [Chen *et al.*, 2017b] Mickaël Chen, Ludovic Denoyer, and Thierry Artières. Multi-view data generation without view supervision. *arXiv preprint arXiv:1711.00305*, 2017.
- [Chen *et al.*, 2018a] Liang-Chieh Chen, Yukun Zhu, George Papandreou, Florian Schroff, and Hartwig Adam. Encoder-decoder with atrous separable convolution for semantic image segmentation. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 801–818, 2018.
- [Chen *et al.*, 2018b] Quan Chen, Tiezheng Ge, Yanyu Xu, Zhiqiang Zhang, Xinxin Yang, and Kun Gai. Semantic human matting. In *ACM MM*, 2018.
- [Chen *et al.*, 2019] Xi Chen, Donglian Qi, and Jianxin Shen. Boundary-aware network for fast and high-accuracy portrait segmentation. *arXiv:1901.03814*, 2019.
- [Cho *et al.*, 2016] Donghyeon Cho, Yu-Wing Tai, and Inso Kweon. Natural image matting using deep convolutional neural networks. In *ECCV*, 2016.
- [Chollet, 2017] François Chollet. Xception: Deep learning with depthwise separable convolutions. In *CVPR*, 2017.
- [Goodfellow *et al.*, 2014] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In *NIPS*, 2014.
- [Gulrajani *et al.*, 2017] Ishaan Gulrajani, Faruk Ahmed, Martin Arjovsky, Vincent Dumoulin, and Aaron C Courville. Improved training of wasserstein gans. In *NIPS*, 2017.
- [He and Sun, 2015] Kaiming He and Jian Sun. Fast guided filter. *arXiv preprint arXiv:1505.00996*, 2015.
- [He *et al.*, 2010] Kaiming He, Jian Sun, and Xiaoou Tang. Guided image filtering. In *ECCV*, 2010.
- [Hu *et al.*, 2018] Yibo Hu, Xiang Wu, Bing Yu, Ran He, and Zhenan Sun. Pose-guided photorealistic face rotation. In *CVPR*, 2018.
- [Ioffe and Szegedy, 2015] Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. *arXiv preprint arXiv:1502.03167*, 2015.
- [Isola *et al.*, 2017] Phillip Isola, Jun-Yan Zhu, Tinghui Zhou, and Alexei A Efros. Image-to-image translation with conditional adversarial networks. In *CVPR*, 2017.
- [Johnson *et al.*, 2016] Justin Johnson, Alexandre Alahi, and Li Fei-Fei. Perceptual losses for real-time style transfer and super-resolution. In *ECCV*, 2016.
- [Kingma and Ba, 2014] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- [Levin *et al.*, 2008] Anat Levin, Dani Lischinski, and Yair Weiss. A closed-form solution to natural image matting. *IEEE TPAMI*, 2008.
- [Levinshtein *et al.*, 2018] Alex Levinshtein, Cheng Chang, Edmund Phung, Irina Kezele, Wenzhangzhi Guo, and Parham Aarabi. Real-time deep hair matting on mobile devices. In *CRV. IEEE*, 2018.
- [Luan *et al.*, 2018] Shangzhen Luan, Chen Chen, Baochang Zhang, Jungong Han, and Jianzhuang Liu. Gabor convolutional networks. *IEEE TIP*, 27(9):4357–4366, 2018.
- [Lutz *et al.*, 2018] Sebastian Lutz, Konstantinos Amliantis, and Aljosa Smolic. Alphagan: Generative adversarial networks for natural image matting. In *BMVC*, 2018.
- [Mao *et al.*, 2017] Xudong Mao, Qing Li, Haoran Xie, Raymond YK Lau, Zhen Wang, and Stephen Paul Smolley. Least squares generative adversarial networks. In *ICCV*, 2017.
- [Miyato *et al.*, 2018] Takeru Miyato, Toshiki Kataoka, Masanori Koyama, and Yuichi Yoshida. Spectral normalization for generative adversarial networks. *arXiv preprint arXiv:1802.05957*, 2018.
- [Nekrasov *et al.*, 2018] Vladimir Nekrasov, Thanuja Dharmasiri, Andrew Spek, Tom Drummond, Chunhua Shen, and Ian Reid. Real-time joint semantic segmentation and depth estimation using asymmetric annotations. *arXiv preprint arXiv:1809.04766*, 2018.
- [Newell *et al.*, 2016] Alejandro Newell, Kaiyu Yang, and Jia Deng. Stacked hourglass networks for human pose estimation. In *ECCV*, 2016.
- [Ouyang and Wang, 2013] Wanli Ouyang and Xiaogang Wang. Joint deep learning for pedestrian detection. In *ICCV*, 2013.
- [Rhemann *et al.*, 2009] Christoph Rhemann, Carsten Rother, Jue Wang, Margrit Gelautz, Pushmeet Kohli, and Pamela Rott. A perceptually motivated online benchmark for image matting. In *CVPR*, 2009.
- [Ronneberger *et al.*, 2015] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *MICCAI*, 2015.
- [Sandler *et al.*, 2018] Mark Sandler, Andrew Howard, Menglong Zhu, Andrey Zhmoginov, and Liang-Chieh Chen. Mobilenetv2: Inverted residuals and linear bottlenecks. In *CVPR*, 2018.
- [Shen *et al.*, 2016] Xiaoyong Shen, Xin Tao, Hongyun Gao, Chao Zhou, and Jiaya Jia. Deep automatic portrait matting. In *ECCV*, 2016.
- [Sup, ] Supervise.ly. <https://supervise.ly/>. Accessed: 2018.
- [Wu *et al.*, 2018] Huikai Wu, Shuai Zheng, Junge Zhang, and Kaiqi Huang. Fast end-to-end trainable guided filter. In *CVPR*, 2018.
- [Xiao *et al.*, 2010] Jianxiong Xiao, James Hays, Krista A Ehinger, Aude Oliva, and Antonio Torralba. Sun database: Large-scale scene recognition from abbey to zoo. In *CVPR*, 2010.
- [Xie *et al.*, 2018] Junyuan Xie, Tong He, Zhi Zhang, Hang Zhang, Zhongyue Zhang, and Mu Li. Bag of tricks for image classification with convolutional neural networks. *arXiv preprint arXiv:1812.01187*, 2018.
- [Xu *et al.*, 2015] Kelvin Xu, Jimmy Ba, Ryan Kiros, Kyunghyun Cho, Aaron Courville, Ruslan Salakhudinov, Rich Zemel, and Yoshua Bengio. Show, attend and tell: Neural image caption generation with visual attention. In *ICML*, 2015.
- [Xu *et al.*, 2017] Ning Xu, Brian Price, Scott Cohen, and Thomas Huang. Deep image matting. In *CVPR*, 2017.
- [Zhou *et al.*, 2017] Bolei Zhou, Hang Zhao, Xavier Puig, Sanja Fidler, Adela Barriuso, and Antonio Torralba. Scene parsing through ade20k dataset. In *CVPR*, 2017.
- [Zhu *et al.*, 2017] Bingke Zhu, Yingying Chen, Jinqiao Wang, Si Liu, Bo Zhang, and Ming Tang. Fast deep matting for portrait animation on mobile phone. In *ACM MM*, 2017.



# Appendices

## A Results of Real Images

We further evaluate the proposed method on real images. We employ DeepLab v3+ [Chen *et al.*, 2018a] as the segmentation method to demonstrate the performance of our method in practice. We only test on some images which have semantic objects that can be segmented by DeepLab v3+, and we adopt the segmentation as the input mask of our method. The results of potted plant images which are gathered from Google are shown in Figure 6. The results of images from Supervise.ly Person dataset [Sup, ] can be viewed in Figure 7. In addition, Figure 8 also demonstrate the results of samples from Adobe Composition-1k [Xu *et al.*, 2017].

## B Larger Unknown Area for Adobe DIM

We also evaluate the performance of different methods with a larger unknown area in trimap for Adobe DIM [Xu *et al.*, 2017]. We apply the dilation and erosion twice with a  $25 \times 25$  kernel to generate a larger unknown area. Since errors are only computed in the unknown area, the quantitative results of our method and Guided Filter [He *et al.*, 2010] are also changed. The experiment results are shown in Table 5 and Figure 9.

Methods	MSE	SAD	Grad. ( $\times 10^3$ )	Conn. ( $\times 10^3$ )
Guided Filter	0.112	24.49	21.33	14.25
Adobe DIM	0.146	21.36	20.54	14.20
Our method	<b>0.017</b>	<b>3.29</b>	<b>8.32</b>	<b>1.89</b>

Table 5: The quantitative results with larger unknown area in trimap on the MAT-2793 testing set

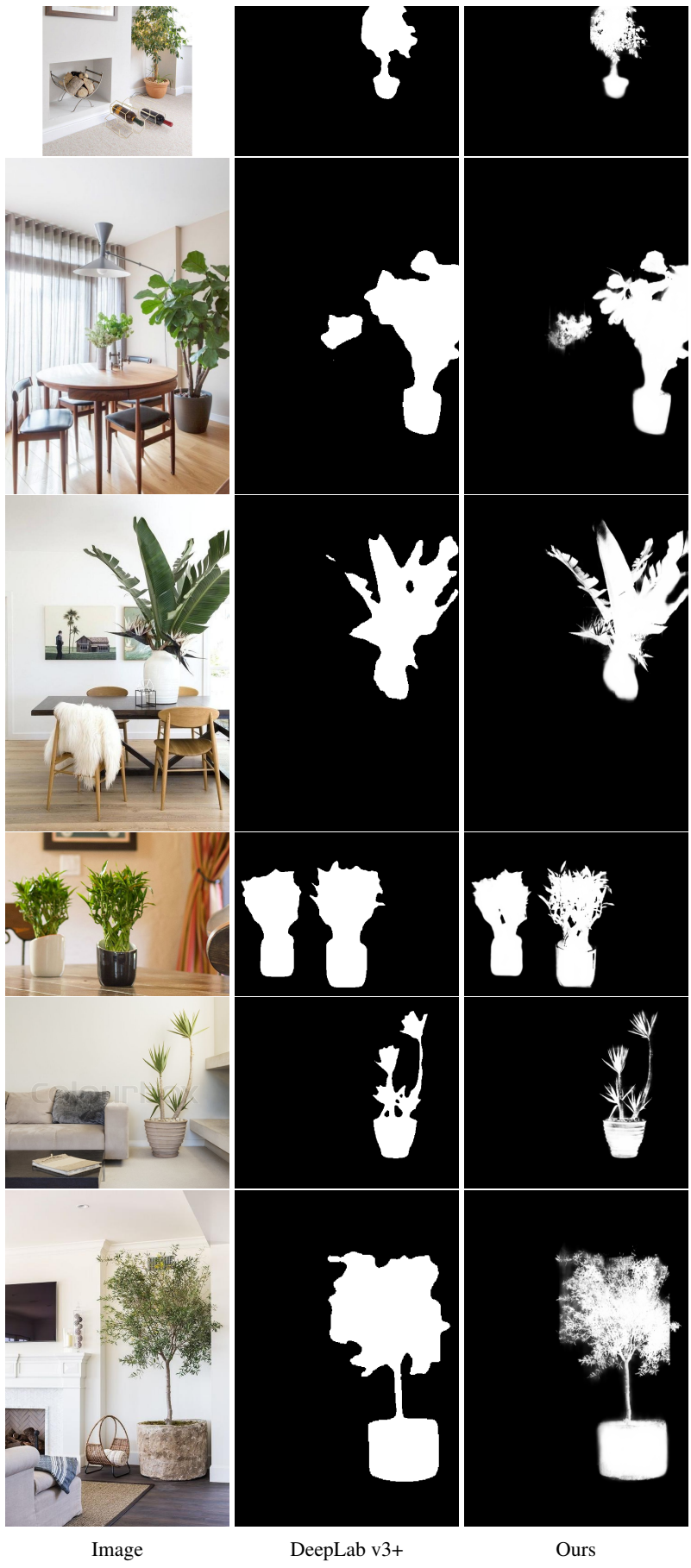


Figure 6: Results of some potted plant images. Masks are generated by DeepLab v3+.

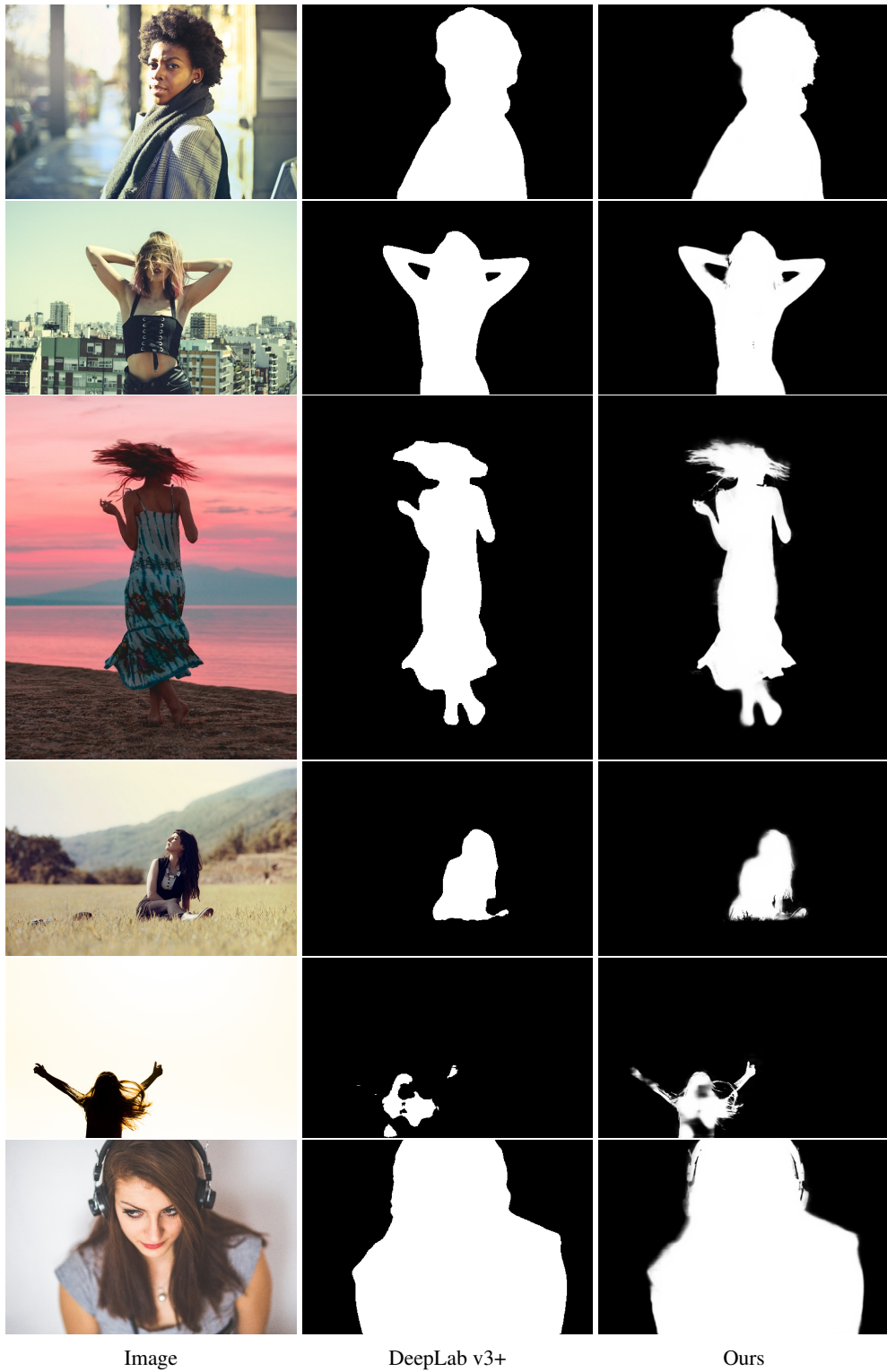


Figure 7: Results of samples from Supervisely Person dataset. Masks are generated by DeepLab v3+.



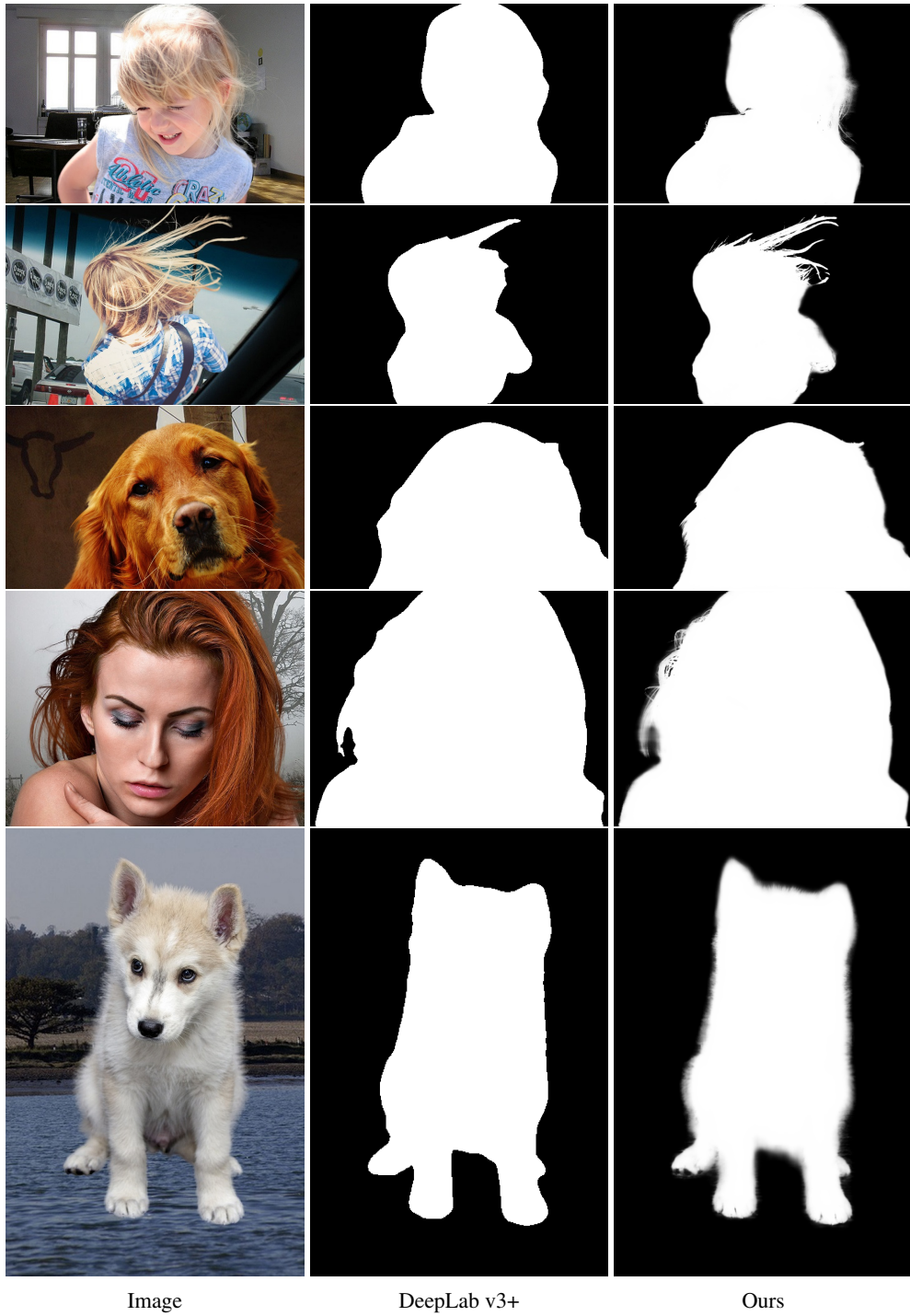


Figure 8: Results of samples from Adobe Composition-1k. Masks are generated by DeepLab v3+.

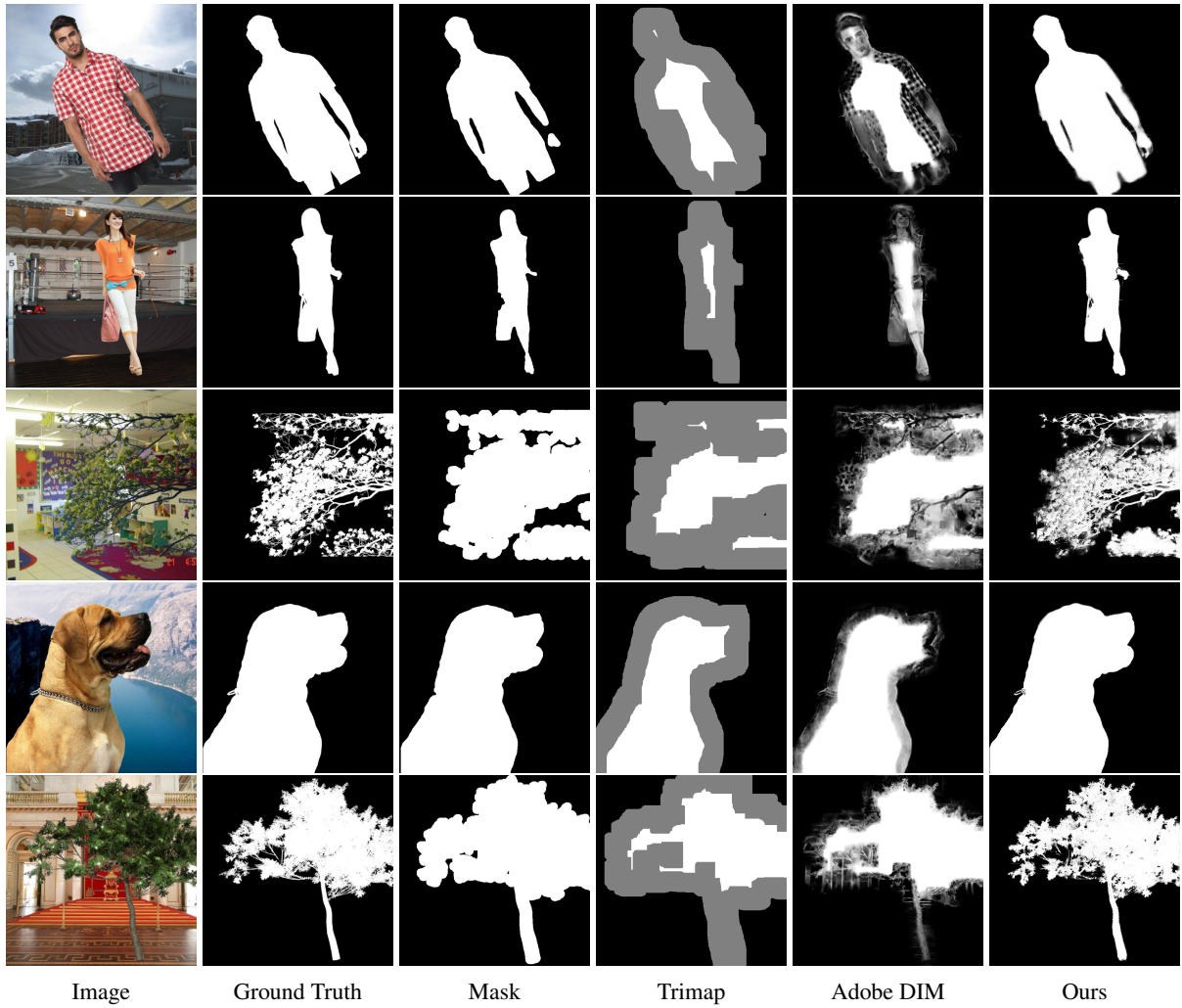


Figure 9: The visual comparison results with larger unknown area in trimap on MAT-2793 testing set. The trimap is only for Adobe DIM.