# MODNet: Real-Time Trimap-Free Portrait Matting via Objective Decomposition

**Zhanghan Ke[1,2], Jiayu Sun[1], Kaican Li[2], Qiong Yan[2], Rynson W.H. Lau[1]**

[1]Department of Computer Science, City University of Hong Kong      [2]SenseTime Research
{zhanghake2-c,jiayusun5-c}@my.cityu.edu.hk, {likaican,yanqiong}@sensetime.com, rynson.lau@cityu.edu.hk

## Abstract

Existing portrait matting methods either require auxiliary inputs that are costly to obtain or involve multiple stages that are computationally expensive, making them less suitable for real-time applications. In this work, we present a light-weight matting objective decomposition network (MODNet) for portrait matting in real-time with a single input image. The key idea behind our efficient design is by optimizing a series of sub-objectives simultaneously via explicit constraints. In addition, MODNet includes two novel techniques for improving model efficiency and robustness. First, an Efficient Atrous Spatial Pyramid Pooling (e-ASPP) module is introduced to fuse multi-scale features for semantic estimation. Second, a self-supervised sub-objectives consistency (SOC) strategy is proposed to adapt MODNet to real-world data to address the domain shift problem common to trimap-free methods. MODNet is easy to be trained in an end-to-end manner. It is much faster than contemporaneous methods and runs at 67 frames per second on a 1080Ti GPU. Experiments show that MODNet outperforms prior trimap-free methods by a large margin on both Adobe Matting Dataset and a carefully designed photographic portrait matting (PPM-100) benchmark proposed by us. Further, MODNet achieves remarkable results on daily photos and videos. Our code and models are available at: https://github.com/ZHKKKe/MODNet, and the PPM-100 benchmark is released at: https://github.com/ZHKKKe/PPM.

## Introduction

Portrait matting aims to predict a precise alpha matte that can be used to extract the persons in a given image or video. It has a wide variety of applications, such as photo editing and movie re-creation. Most existing matting methods use a pre-defined trimap as a priori to produce an alpha matte (Cai et al. 2019; Hou and Liu 2019; Li and Lu 2020; Lu et al. 2019; Tang et al. 2019; Xu et al. 2017). However, trimaps are costly to annotate. Although a depth sensor (Foix, Alenyà, and Torras 2011) can ease the task, the resulting trimaps may suffer from low precision. Some recent trimap-free methods attempt to eliminate the model dependence on the trimap. For example, background matting (Sengupta et al. 2020) replaces the trimap by a separate background image. Other

methods (Chen et al. 2018; Liu et al. 2020; Shen et al. 2016) include multiple stages (*i.e.*, with several independent models that are optimized sequentially) to first generate a pseudo trimap or semantic mask, which is then used to serve as the priori for alpha matte prediction. Nonetheless, using the background image as input has to take and align two photos, while having multiple stages would significantly increase the inference time. These drawbacks make all aforementioned matting methods not suitable for real-time applications, such as camera preview. Besides, limited by insufficient amount of labeled training data, existing trimap-free methods often suffer from the domain shift problem (Sun, Feng, and Saenko 2016) in practice, *i.e.*, the models cannot generalize well to real-world data.

In this work, we present MODNet, a light-weight model for real-time trimap-free portrait matting. As illustrated in Fig. 1, unlike prior methods which require auxiliary inputs or consist of multiple stages, MODNet inputs a single RGB image and applies explicit constraints to solve matting sub-objectives simultaneously in one stage. There are two insights behind our design. First, applying explicit constraints to different parts of the model can address portrait matting effectively under a single input image. In contrast, to obtain comparable results, auxiliary inputs would be necessary for the model trained by only one matting constraint. Second, optimizing sub-objectives simultaneously can further exploit the model capability by sharing intermediate representations. In contrast, training multiple stages sequentially will accumulate the errors from the early stages and magnify them in subsequent stages. We further propose two novel techniques to improve MODNet's efficiency and robustness, including (1) an Efficient Atrous Spatial Pyramid Pooling (e-ASPP) module for fast multi-scale feature fusion in portrait semantic estimation, and (2) a self-supervised strategy based on sub-objective consistency (SOC) to alleviate the domain shift problem in real-world data.

MODNet has several advantages over previous trimap-free methods. First, MODNet is much faster, running at 67 frames per second ($fps$) on a GTX 1080Ti GPU with an input size of $512 \times 512$ (including data loading). Second, MODNet achieves state-of-the-art results on both open source Adobe Matting benchmark and our newly proposed PPM-100 benchmark. Third, MODNet can be easily optimized end-to-end as it is a single model instead of a complex
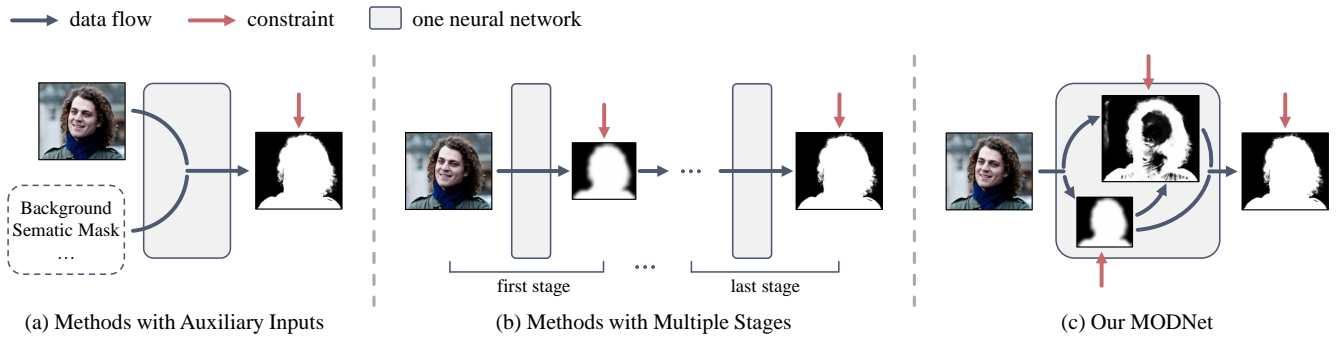
Figure 1: **Different Trimap-free Matting Approaches.** Existing trimap-free matting methods either (a) require auxiliary inputs to address the complex matting objective directly or (b) consist of multiple stages to address the matting sub-objectives sequentially. Both of them are less suitable for real-time applications. Instead, (c) our MODNet solves the matting sub-objectives simultaneously with only a single input image, which is more efficient and effective.

pipeline. Finally, MODNet has better generalization ability, due to our SOC strategy.

Since open-source portrait matting datasets (Shen et al. 2016; Xu et al. 2017) are typically small and have limited precision, prior works train and validate their models on private datasets of diverse quality and difficulty levels. As a result, it is difficult to conduct a fair evaluation. In this work, we evaluate existing trimap-free methods under the same environment, *i.e.*, all models are trained on the same dataset and validated on the portrait images from Adobe Matting Dataset (Xu et al. 2017) and our newly proposed benchmark. Our benchmark is labelled in high quality, and is more diverse than those used in previous works.

In summary, we present a light-weight network architecture, named MODNet, for real-time trimap-free portrait matting. MODNet includes two novel techniques, an e-ASPP module for efficient semantic feature fusion and a self-supervised SOC strategy to generalize MODNet to new data domain. In addition, we have also designed a new validation benchmark for portrait matting. Our code, pre-trained model, and benchmark are released at *https://github.com/ZHKKKe/MODNet* and *https://github.com/ZHKKKe/PPM*.

## Related Works

### Image Matting

The purpose of the image matting task is to extract the desired foreground $F$ from a given image $I$. This task predicts an alpha matte with a precise foreground probability value $\alpha$ for each pixel $i$ as:

$$I^i = \alpha^i F^i + (1 - \alpha^i) B^i, \qquad (1)$$

where $B$ is the background of $I$. This problem is ill-posed since all variables on the right hand side are unknown. Most existing matting methods take a pre-defined trimap as an auxiliary input, which is a mask containing three regions: absolute foreground ($\alpha = 1$), absolute background ($\alpha = 0$), and unknown area ($\alpha = 0.5$). They aim to estimate the foreground probability inside the unknown area only, based on the priori from the other two regions.

Traditional matting algorithms heavily rely on low-level features, *e.g.*, color cues, to determine the alpha matte through sampling (Chuang et al. 2001; Feng, Liang, and Zhang 2016; Gastal and Oliveira 2010; He et al. 2011; Johnson et al. 2016; Karacan, Erdem, and Erdem 2015; Ruzon and Tomasi 2000; Yang et al. 2018) or propagation (Aksoy, Aydin, and Pollefeys 2017; Aksoy et al. 2018; Bai and Sapiro 2007; Chen, Li, and Tang 2013; Grady et al. 2005; Levin, Lischinski, and Weiss 2007; Levin, Rav-Acha, and Lischinski 2008; Sun et al. 2004), which often fail in complex scenes. With the tremendous progress of deep learning, many methods based on CNNs have been proposed with notable successes. Cho *et al.* (Cho, Tai, and Kweon 2016) and Shen *et al.* (Shen et al. 2016) combined the classic algorithms with CNNs for alpha matte refinement. Xu *et al.* (Xu et al. 2017) proposed an auto-encoder architecture to predict an alpha matte from a RGB image and a trimap. Some works (Li and Lu 2020; Lu et al. 2019) used attention mechanisms to help improve matting performances. Cai *et al.* (Cai et al. 2019) suggested a trimap refinement process before matting and showed the advantages of using the refined trimaps.

Since creating trimaps requires users' additional efforts and the quality of the resulting mattes depends on users' experiences, some recent methods (including our MODNet) attempt to avoid them, as described below.

### Trimap-Free Portrait Matting

Image matting is extremely difficult if trimaps are not available, as a semantic estimation step will then be needed to locate the foreground, before an alpha matte can be predicted. To constrain the problem, most trimap-free methods focus on just one type of foreground objects, *e.g.*, portraits. Nevertheless, feeding RGB images to a single network often produce unsatisfactory alpha mattes. Sengupta *et al.* (Sengupta et al. 2020) proposed to capture a less expensive background image as a pseudo green screen to alleviate this problem. Other works designed their pipelines with multiple stages. For example, Shen *et al.* (Chen et al. 2018) assembled a trimap generation network before the matting network. Zhang *et al.* (Zhang et al. 2019) applied
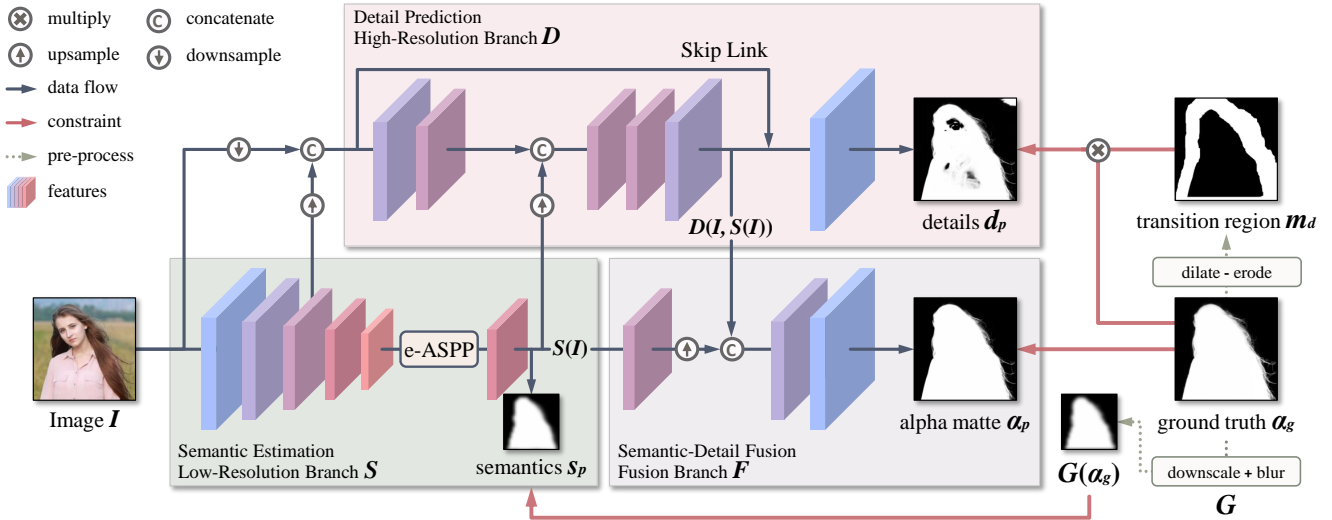
Figure 2: **Architecture of MODNet.** Given an input image $I$, MODNet predicts portrait semantics $s_p$, boundary details $d_p$, and final alpha matte $\alpha_p$ through three interdependent branches, $S$, $D$, and $F$, which are constrained by explicit supervisions generated from the ground truth matte $\alpha_g$. Since the decomposed sub-objectives are correlated and help strengthen each other, we can optimize MODNet end-to-end.

a fusion network to combine the predicted foreground and background. Liu *et al.* (Liu et al. 2020) concatenated three networks to utilize coarse labeled data in matting. The main problem of all these methods is that they cannot be used in interactive applications as: (1) the background images may change across frames, and (2) computationally expensive due to having multiple stages in the pipeline. Compared to the above methods, MODNet is light-weight in terms of pipeline complexity. It takes only one RGB image as input and uses a single model for real-time matting with better performances, by optimizing a series of sub-objectives simultaneously with explicit constraints.

### Other Techniques

We briefly discuss two techniques that are relevant to the design and optimization of our method.

**Atrous Spatial Pyramid Pooling (ASPP).** ASPP (Chen et al. 2017) has been widely explored and proved to boost the performance notably in segmentation-based tasks. Although ASPP has a huge number of parameters and a high computational overhead, some matting models (Qiao et al. 2020; Li and Lu 2020) still used it for better results. In MODNet, we design an efficient variant of ASPP that gives comparable performances with a much lower overhead.

**Consistency Constraint.** Consistency is an important assumptions behind many semi-/self-supervised (Schmarje et al. 2020) and domain adaptation (Wilson and Cook 2020) algorithms. For example, Ke *et al.* (Ke et al. 2020) designed a consistency-based framework that could be used for semi-supervised matting. Toldo *et al.* (Toldo et al. 2020) presented a consistency-based domain adaptation strategy for semantic segmentation. However, these methods consist of multiple models and constrain the consistency among their predictions. In contrast, MODNet imposes consistency among various sub-objectives within a model.

## MODNet

In MODNet, we divide the trimap-free matting objective into three parts: semantic estimation, detail prediction, and semantic-detail fusion. We optimize them simultaneously via three branches (Fig. 2). In the following subsections, we will delve into the design of each branch and the supervisions used to solve the sub-objectives.

### Semantic Estimation

Similar to existing multi-model approaches, the first step of MODNet is to locate the portrait in the input image $I$. The difference is that we extract high-level semantics only through an encoder, *i.e.*, the low-resolution branch $S$ of MODNet. This has two main advantages. First, semantic estimation becomes more efficient because a separate decoder with huge parameters is no longer required. Second, the high-level representation $S(I)$ is helpful for subsequent branches and joint optimization. An arbitrary CNN backbone can be applied to $S$. To facilitate real-time interaction, we adopt MobileNetV2 (Sandler et al. 2018), which is an ingenious model developed for mobile devices, as $S$.

To predict coarse semantic mask $s_p$, we feed $S(I)$ into a convolutional layer activated by the Sigmoid function to reduce its channel number to 1. We supervise $s_p$ by a thumbnail of the ground truth matte $\alpha_g$. Since $s_p$ is supposed to be smooth, we use the L2 loss as:

$$\mathcal{L}_s = \frac{1}{2} \left|\left| s_p - G(\alpha_g) \right|\right|_2 , \qquad (2)$$

where $G$ stands for $16\times$ downsampling followed by Gaussian blur. It removes the fine structures (such as hair) that are
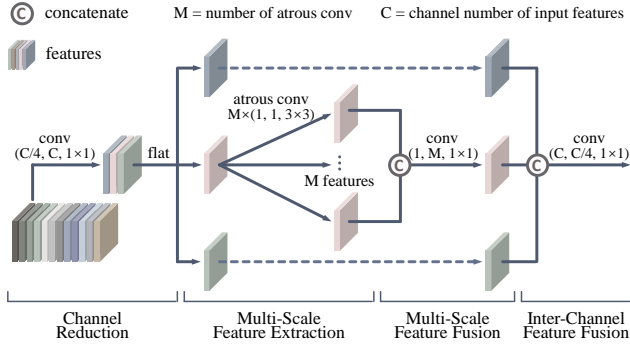
Figure 3: **Illustration of e-ASPP.** Our e-ASPP is efficient since it extracts and fuses multi-scale features depth-wise, followed by an inter-channel fusion. The tuple under convolution are (output channel, input channel, kernel size). The dotted lines indicate the same structure as the solid line in the center branch.

not essential to portrait semantics.

**Efficient ASPP (e-ASPP).** Semantic masks predicted by MobileNetV2 may have holes in some challenging foregrounds or backgrounds. Many prior works showed that ASPP was a feasible solution for improving such erroneous semantics. However, ASPP has a very high computational overhead. To balance between performance and efficiency, we design an efficient ASPP (e-ASPP) module to process $S(I)$, as illustrated in Fig. 3.

The standard ASPP utilizes atrous convolutions for multi-scale feature extraction and applies a standard convolution for multi-scale feature fusion. We modify it to e-ASPP via three steps. First, we split each atrous convolution into a depth-wise atrous convolution and a point-wise convolution. The depth-wise atrous convolution extracts multi-scale features independently on each channel, while the point-wise convolution is appended for inter-channel fusion at each scale. Second, we switch the order of inter-channel fusion and the multi-scale feature fusion. In this way, (1) only one inter-channel fusion is required, and (2) the multi-scale feature fusion is converted to a cheaper depth-wise operation. Third, we compress the number of input channels by $4\times$ for e-ASPP and recover it before the output.

Compared to the original ASPP, our proposed e-ASPP has only 1% of the parameters and 1% of the computational overhead [1]. In MODNet, our experiments show that e-ASPP can achieve performance comparable to ASPP.

### Detail Prediction

We process a transition region around the foreground portrait with a high-resolution branch $D$, which takes $I$, $S(I)$, and the low-level features from $S$ as inputs. The purpose of reusing the low-level features is to reduce the computational overhead of $D$. In addition, we further simplify $D$ in the following three aspects: (1) $D$ consists of fewer convolutional layers than $S$; (2) a small channel number is chosen for the

---
[1]Refer to Appendix A for more details of e-ASPP.

convolutional layers in $D$; (3) we do not maintain the original input resolution throughout $D$. In practice, $D$ consists of 12 convolutional layers, and its maximum channel number is 64. The resolution of the feature maps is reduced to $1/4$ of $I$ in the first layer and restored in the last two layers. The impact of the downsampling operation on $D$ is negligible since it contains a skip link.

We denote the outputs of $D$ as $D(I, S(I))$, which implies the dependency between sub-objectives — high-level portrait semantics $S(I)$ is a priori for detail prediction. We calculate the boundary detail matte $d_p$ from $D(I, S(I))$ and learn it through the L1 loss as:

$$\mathcal{L}_d = m_d \left|\left| d_p - \alpha_g \right|\right|_1, \qquad (3)$$

where $m_d$ is a binary mask to let $\mathcal{L}_d$ focus on the portrait boundaries. $m_d$ is generated through dilation and erosion on $\alpha_g$. Its values are 1 if the pixels are inside the transition region, and 0 otherwise.

### Semantic-Detail Fusion

The fusion branch $F$ in MODNet is a straightforward CNN module, combining semantics and details. We first upsample $S(I)$ to match its size with $D(I, S(I))$. We then concatenate $S(I)$ and $D(I, S(I))$ to predict the final alpha matte $\alpha_p$, constrained by:

$$\mathcal{L}_\alpha = \left|\left| \alpha_p - \alpha_g \right|\right|_1 + \mathcal{L}_c, \qquad (4)$$

where $\mathcal{L}_c$ is the compositional loss from (Xu et al. 2017). It measures the absolute difference between input image $I$ and the composited image obtained from $\alpha_p$, the ground truth foreground, and the ground truth background.

MODNet is trained end-to-end through the sum of $\mathcal{L}_s$, $\mathcal{L}_d$, and $\mathcal{L}_\alpha$, as:

$$\mathcal{L} = \lambda_s \mathcal{L}_s + \lambda_d \mathcal{L}_d + \lambda_\alpha \mathcal{L}_\alpha, \qquad (5)$$

where $\lambda_s$, $\lambda_d$, and $\lambda_\alpha$ are hyper-parameters balancing the three losses. The training process is robust to these hyper-parameters. We set $\lambda_s = \lambda_\alpha = 1$ and $\lambda_d = 10$.

## SOC for Real-World Data

The training data for portrait matting requires excellent labeling in the hair area, which is difficult to do for natural images with complex backgrounds. Currently, most annotated data comes from photography websites. Although these images have monochromatic or blurred backgrounds, the labeling process still needs to be completed by experienced annotators with considerable amount of time. As such, the labeled datasets for portrait matting are usually small. Xu *et al.* (Xu et al. 2017) suggested using background replacement as a data augmentation to enlarge the training set, and it has become a common setting in image matting. However, the training samples obtained in such a way exhibit different properties from those of the daily life images. Therefore, existing trimap-free models always tend to overfit the training set and perform poorly on real-world data.

To address this domain shift problem, we propose to utilize sub-objectives consistency (SOC) to adapt MODNet to unseen data distributions. The three sub-objectives
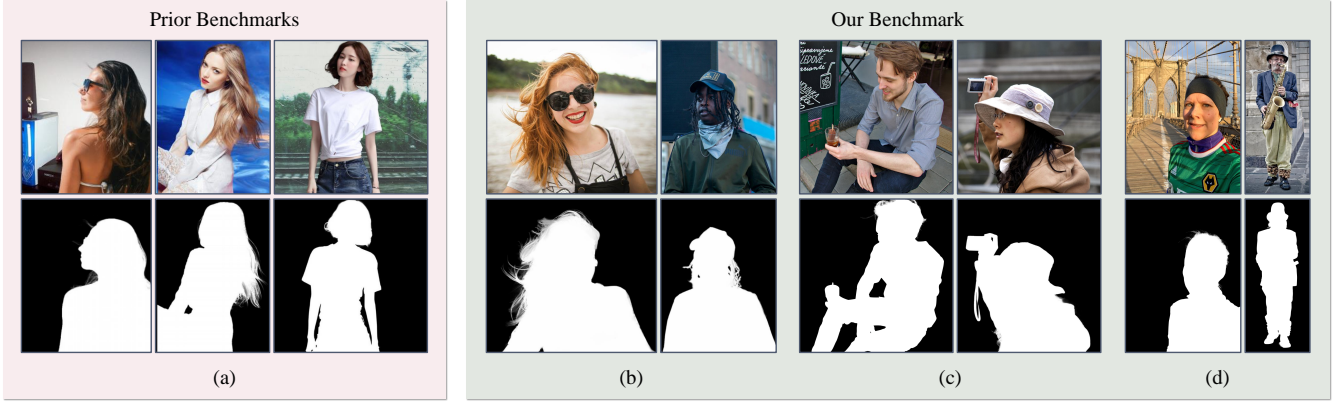
Figure 4: **Benchmark Comparison.** (a) Validation benchmarks used in (Chen et al. 2018; Liu et al. 2020; Zhang et al. 2019). Images are synthesized by replacing the original backgrounds with new ones. Instead, our PPM-100 contains original image backgrounds and has a higher diversity in the foregrounds, *e.g.*, (b) with fine hairs, (c) with additional objects, and (d) without bokeh or with full-body.

in MODNet should have consistent outputs in semantics or details. We take semantic consistency as an example, MODNet outputs portrait semantics $S(I)$ and alpha matte $F(S(I), D(S(I)))$ for input image $I$. As $S(I)$ is the prior of $F(S(I), D(S(I)))$, they should have consistent portrait semantics. In the labeled source domain, there is good consistency among the predictions of sub-objectives. However, inconsistent predictions occur in the unlabeled target domain, which may cause poor results. Motivated by this observation, our self-supervised SOC strategy imposes the consistency constraints among the predictions of the sub-objectives (Fig. 1(b)) to improve the performance of MOD-Net in the new domain, without ground truth labels.

Formally, we denote MODNet as $M$. As described in Sec. , $M$ has three outputs for an unlabeled image $\tilde{I}$:

$$\tilde{s}_p, \tilde{d}_p, \tilde{\alpha}_p = M(\tilde{I}) . \qquad (6)$$

We enforce the semantics in $\tilde{\alpha}_p$ to be consistent with $\tilde{s}_p$ and the details in $\tilde{\alpha}_p$ to be consistent with $\tilde{d}_p$ by:

$$\mathcal{L}_{cons} = \frac{1}{2} \left\| G(\tilde{\alpha}_p) - \tilde{s}_p \right\|_2 + \tilde{m}_d \left\| \tilde{\alpha}_p - \tilde{d}_p \right\|_1 , \qquad (7)$$

where $\tilde{m}_d$ indicates the transition region in $\tilde{\alpha}_p$, and $G$ has the same meaning as the one in Eq. 2. However, adding the L2 loss to blurred $G(\tilde{\alpha}_p)$ will smooth the boundaries in the optimized $\tilde{\alpha}_p$. As a result, the consistency between $\tilde{\alpha}_p$ and $\tilde{d}_p$ will remove the details predicted by the high-resolution branch. To prevent this problem, we duplicate $M$ to $M'$ and fix the weights of $M'$ before performing SOC. Since the fine boundaries are preserved in $\tilde{d}_p'$ output by $M'$, we append an extra regularization term to maintain the details in $M$ as:

$$\mathcal{L}_{dd} = \tilde{m}_d \left\| \tilde{d}_p' - \tilde{d}_p \right\|_1 . \qquad (8)$$

The sum of $\mathcal{L}_{cons}$ and $\mathcal{L}_{dd}$ is optimized during SOC.

## Experiments

In this section, we first introduce our PPM-100 benchmark for portrait matting. We then compare MODNet with existing matting methods on both Adobe Matting Dataset (AMD)

(Xu et al. 2017) and our PPM-100. We further conduct ablation experiments to evaluate various components of MOD-Net. Finally, we demonstrate the effectiveness of SOC in adapting MODNet to real-world data.

### Photographic Portrait Matting Benchmark

Existing works constructed their validation benchmarks from a small amount of labeled data through image synthesis. Their benchmarks are relatively easy due to unnatural fusion or mismatched semantics between the foreground and the background (Fig. 4(a)). Hence, trimap-free models may have comparable performances to the trimap-based models on these benchmarks, but unsatisfactory performances on natural images, *i.e.*, images without background replacement. This indicates that the performances of trimap-free methods have not been accurately assessed.

In contrast, we propose a Photographic Portrait Matting benchmark (PPM-100), which contains 100 finely annotated portrait images with various backgrounds. To guarantee sample diversity, we consider several factors in order to balance the sample types in PPM-100, including: (1) whether the whole portrait body is included; (2) whether the image background is blurred; and (3) whether the person is holding additional objects. We regard small objects held by a foreground person as a part of the foreground, which is more in line with practical applications. As shown in Fig. 4(b)(c)(d), the samples in PPM-100 have more natural backgrounds and richer postures. Hence, PPM-100 can be considered as a more comprehensive benchmark.

### Results on AMD and PPM-100[2]

We compare MODNet with trimap-free FDMPA (Zhu et al. 2017), LFM (Zhang et al. 2019), SHM (Chen et al. 2018), BSHM (Liu et al. 2020), and HAtt (Qiao et al. 2020). We use DIM (Xu et al. 2017) and IndexMatter (Lu et al. 2019)

---

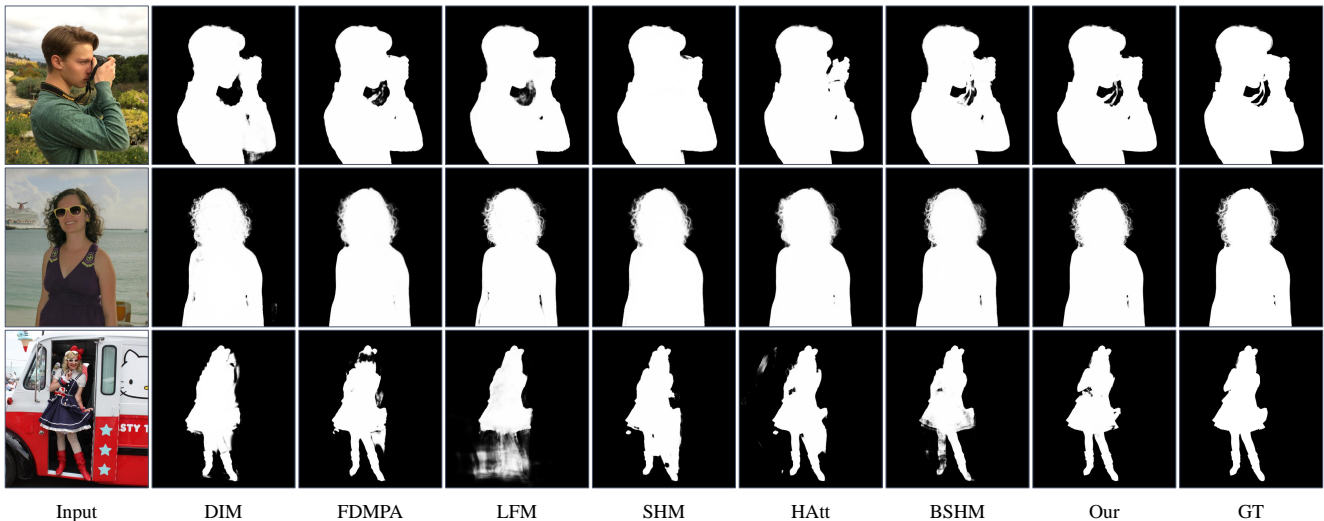[2]Refer to Appendix B for results on more benchmarks.

Figure 5: **Visual Comparison of Trimap-free Methods on PPM-100.** MODNet performs better in hollow structures (the $1st$ row) and hair details (the $2nd$ row). However, it may still make mistakes in challenging poses or costumes (the $3rd$ row). DIM (Xu et al. 2017) here does not take trimaps as input, but is pre-trained on the SPS (supervise.ly 2018) dataset.

| Method | Trimap | MSE ↓ | MAD ↓ |
|---|---|---|---|
| DIM (Xu et al. 2017) | ✓ | 0.0016 | 0.0067 |
| IndexMatter (Lu et al. 2019) | ✓ | 0.0015 | 0.0064 |
| MODNet (Our) | ✓ | **0.0013** | **0.0054** |
| DIM (Xu et al. 2017) | | 0.0221 | 0.0327 |
| DIM† (Xu et al. 2017) | | 0.0115 | 0.0178 |
| FDMPA† (Zhu et al. 2017) | | 0.0101 | 0.0160 |
| LFM† (Zhang et al. 2019) | | 0.0094 | 0.0158 |
| SHM† (Chen et al. 2018) | | 0.0072 | 0.0152 |
| HAtt† (Qiao et al. 2020) | | 0.0067 | 0.0137 |
| BSHM† (Liu et al. 2020) | | 0.0063 | 0.0114 |
| MODNet† (Our) | | **0.0044** | **0.0086** |

Table 1: **Quantitative Results on PPM-100.** A '†' indicates that the model is pre-trained on SPS.

| Method | Trimap | MSE ↓ | MAD ↓ |
|---|---|---|---|
| DIM (Xu et al. 2017) | ✓ | 0.0014 | 0.0069 |
| IndexMatter (Lu et al. 2019) | ✓ | 0.0013 | 0.0066 |
| MODNet (Our) | ✓ | **0.0011** | **0.0063** |
| DIM (Xu et al. 2017) | | 0.0075 | 0.0159 |
| DIM† (Xu et al. 2017) | | 0.0048 | 0.0116 |
| FDMPA† (Zhu et al. 2017) | | 0.0047 | 0.0115 |
| LFM† (Zhang et al. 2019) | | 0.0043 | 0.0101 |
| SHM† (Chen et al. 2018) | | 0.0031 | 0.0092 |
| HAtt† (Qiao et al. 2020) | | 0.0034 | 0.0094 |
| BSHM† (Liu et al. 2020) | | 0.0029 | 0.0088 |
| MODNet† (Our) | | **0.0023** | **0.0077** |

Table 2: **Quantitative Results on AMD.** We pick the portrait foregrounds from AMD for validation. A '†' indicates that the models pre-trained on SPS.

as the trimap-based baselines. For methods without publicly available codes, we follow their papers to reproduce them.

For a fair comparison, we train all models on the same dataset, which contains nearly $3,000$ annotated foregrounds. Background replacement (Xu et al. 2017) is applied to extend our training set. All images in our training set are collected from Flickr and are annotated by Photoshop. The training set contains $\sim 2,600$ half-body and $\sim 400$ full-body portraits. For each labeled foreground, we generate 5 samples by random cropping and 10 samples by compositing with the images from the OpenImage dataset (Kuznetsova et al. 2018) (as the background). We use MobileNetV2 pre-trained on the Supervisely Person Segmentation (SPS) (supervise.ly 2018) dataset as the backbone of all trimap-free models. For the compared methods, we explore the optimal hyper-parameters through grid search. For MODNet, we train it by SGD for 40 epochs. With a batch size of 16, the

initial learning rate is set to $0.01$ and is multiplied by $0.1$ after every 10 epochs. We use Mean Square Error (MSE) and Mean Absolute Difference (MAD) as quantitative metrics.

Table 1 shows the results on PPM-100. MODNet outperforms other trimap-free methods on both MSE and MAD. However, it is unable to outperform trimap-based methods, as PPM-100 contains samples with very challenging poses and costumes. When taking a trimap as input during both training and testing stages, *i.e.*, regarding MODNet as a trimap-based methods, it outperforms the compared trimap-based methods. This demonstrates the superiority of the proposed architecture. Fig. 5 shows visual comparison [3].

Table 2 shows the results on AMD (Xu et al. 2017). We pick the portrait foregrounds from AMD and composite 10

---

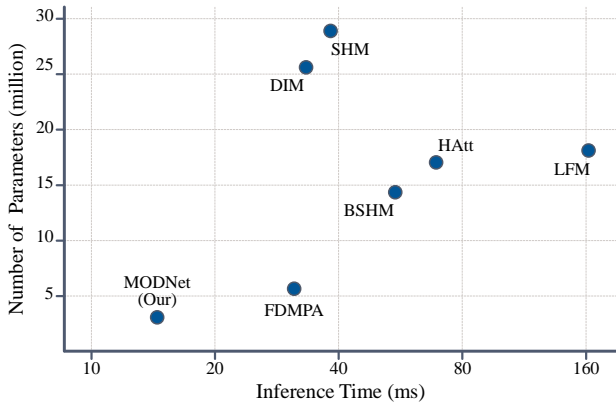[3]Refer to Appendix C for more visual results.

Figure 6: **Comparison on Model Size and Execution Efficiency.** $fps$ can be obtained by dividing $1,000$ with the inference time.

| $\mathcal{L}_s$ | $\mathcal{L}_d$ | e-ASPP | SPS | MSE $\downarrow$ | MAD $\downarrow$ |
|---|---|---|---|---|---|
| | | | | 0.0162 | 0.0235 |
| $\checkmark$ | | | | 0.0097 | 0.0158 |
| $\checkmark$ | $\checkmark$ | | | 0.0083 | 0.0142 |
| $\checkmark$ | $\checkmark$ | $\checkmark$ | | 0.0057 | 0.0109 |
| $\checkmark$ | $\checkmark$ | $\checkmark$ | $\checkmark$ | **0.0044** | **0.0086** |

Table 3: **Ablation of MODNet on PPM-100.** SPS indicates the model us pre-trained on SPS.

test samples for each foreground with diverse backgrounds. We validate all trained models on this synthetic benchmark. Unlike the results on PPM-100, the performance gap between trimap-free and trimap-based models is much smaller. The results show that trimap-free models can achieve results comparable to trimap-based models only on the synthetic benchmarks that have unnatural fusion or mismatched semantics between foreground and background.

We further evaluate MODNet on model size and execution efficiency. A small model facilitates deployment on mobile/handheld devices, while high execution efficiency is necessary for real-time applications. We measure the model size by the total number of parameters, and we reflect the execution efficiency by the average inference time over PPM-100 on an NVIDIA GTX 1080Ti GPU (all input images are resized to $512 \times 512$). Note that fewer parameters do not imply faster inference speed due to large feature maps or time-consuming mechanisms, *e.g.*, attention, that the model may use. Fig. 6 summarizes the results. The inference time of MODNet is $14.9\,ms$ ($67\,fps$), which is twice the $fps$ of the fastest method, FDMPA ($31\,fps$). In addition, our MODNet has the smallest number of parameters among the trimap-free methods.

We have also conducted ablation experiments for MODNet on PPM-100, as shown in Table 3. Applying $\mathcal{L}_s$ and $\mathcal{L}_d$ to constrain portrait semantics and boundary details bring considerable performance improvements. The results
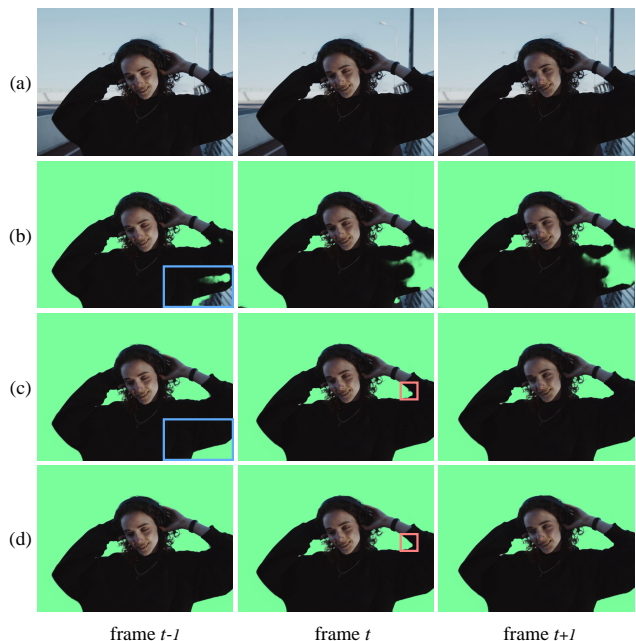
Figure 7: **Results on a Real-World Video.** We show three consecutive video frames from left to right. From top to bottom: (a) Input, (b) MODNet, (c) MODNet + SOC, and (d) MODNet + SOC + OFD. The blue region in frame $t-1$ shows the effectiveness of SOC, while the red region in frame $t$ highlights the flickers eliminated by OFD.

also show that the effectiveness of e-ASPP in fusing multi-level feature maps. Although SPS pre-training is optional to MODNet, it plays a vital role in other trimap-free methods. From Table 1, we can see that trimap-free DIM without pre-training performs far worse than the one with pre-training.

## Results on Real-World Data

To adapt MODNet to real-world data, we capture $\sim$400 video clips (divided into about 50,000 frames) as the unlabeled data for self-supervised SOC domain adaptation. In this stage, we freeze the BatchNorm layers within MODNet and finetune the convolutional layers by Adam at a learning rate of $0.0001$. The total number of fine-tuning epochs are 15. Here, we only provide visual results, as ground truth mattes are not available. In Fig. 7(b)(c), we composite the foreground over a green screen to emphasize that SOC is vital for generalizing MODNet to real-world data.

For video data, we also propose here a simple but effective One-Frame Delay (OFD) trick to reduce the flickers in the predicted alpha matte sequence. The idea behind OFD is that we can utilize the preceding and the following frames to fix the flickering pixels, because the corresponding pixels in adjacent frames are likely to be correct. Suppose that we have three consecutive frames, and their corresponding alpha mattes are $\alpha_{t-1}$, $\alpha_t$, and $\alpha_{t+1}$, where $t$ is the frame index. We regard $\alpha_t^i$ as a flickering pixel if the values of $\alpha_{t-1}^i$ and $\alpha_{t+1}^i$ are close, and $\alpha_t^i$ is very different from the values

of both $\alpha_{t-1}^i$ and $\alpha_{t+1}^i$. When $\alpha_t^i$ is a flickering pixel, we replace its value by averaging $\alpha_{t-1}^i$ and $\alpha_{t+1}^i$. As illustrated in Fig. 7(c)(d), OFD can further removes flickers along the boundaries.

## Conclusion

This paper has presented a simple, fast, and effective model, MODNet, for portrait matting. By taking only an RGB image as input, our method enables the prediction of a high-quality alpha matte in real time, which is benefited from objective decomposition and concurrent optimization with explicit supervisions. Besides, we have introduced (1) an e-ASPP module to speed up the multi-scale feature fusion process, and (2) a self-supervised sub-objectives consistency (SOC) strategy to allow MODNet to handle the domain shift problem. Extensive experiments show that MODNet outperforms existing trimap-free methods on the AMD benchmark, the proposed PPM-100 benchmark, and a variety of real-world data. Our method does have limitations. The main one is that it may fail to handle videos with strong motion blurs due to the lack of temporal information. One possible future work is to address the video matting problem under motion blurs through additional sub-objectives, such as optical flow estimation.

## Acknowledge

## Appendix A: Analysis of e-ASPP

Here we compare the proposed Efficient ASPP (e-ASPP) with the standard ASPP in terms of the number of parameters and computational overhead. For a convolutional layer, the number of its parameters $\mathcal{P}$ can be calculated by:

$$\mathcal{P} = C_{out} \times C_{in} \times K \times K, \qquad (9)$$

where $C_{out}$ is the number of output channels, $C_{in}$ is the number of input channels, and $K$ is the kernel size. We can use *FLOPs* to measure the computational overhead $\mathcal{O}$ of a convolutional layer as:

$$\mathcal{O} = C_{in} \times 2 \times K \times K \times H_{out} \times W_{out} \times C_{out}, \quad (10)$$

where $H_{out}$ and $W_{out}$ are the height and the width of output feature maps, respectively.

Following, we represent the size of the input feature maps by $(c, h, w)$, where $c$ is the number of channels, $h$ is the height of the input feature maps, and $w$ is the width of the input feature maps. We represent the number of atrous convolutional layers (with a kernel size of $k$) in both ASPP and e-ASPP by $m$.

**Standard ASPP (ASPP).** In ASPP, (1) all atrous convolutional layers are independently applied to the input features maps to extract multi-scale features. These multi-scale features are then (2) concatenated and processed by a pointwise convolutional layer (with a kernel size of 1). We have:

$$\begin{aligned} \mathcal{P}_{\mathcal{ASPP}} =& m \times (c \times c \times k \times k) \\ &+ c \times (m \times c) \times 1 \times 1 \qquad (11) \\ =& m \times c^2 \times (k^2 + 1), \end{aligned}$$

$$\begin{aligned} \mathcal{O}_{ASPP} =& m \times (c \times 2 \times k \times k \times h \times w \times c) \\ &+ (m \times c) \times 2 \times 1 \times 1 \times h \times w \times c \qquad (12) \\ =& ((2 \times k^2 + 2) \times m \times c) \times (h \times w \times c). \end{aligned}$$

**Efficient ASPP (e-ASPP).** As shown in Fig. 3 (in the paper), e-ASPP consists of four operations, including (1) Channel Reduction, (2) Multi-Scale Feature Extraction, (3) Multi-Scale Feature Fusion, and (4) Inter-Channel Feature Fusion. The total number of parameters and the total *FLOPs* are the sum of these four operations. We have:

$$\begin{aligned} \mathcal{P}_{e-ASPP} =& \frac{c}{4} \times c \times 1 \times 1 \\ &+ \frac{c}{4} \times m \times (1 \times 1 \times k \times k) \\ &+ \frac{c}{4} \times (1 \times m \times 1 \times 1) \qquad (13) \\ &+ c \times \frac{c}{4} \times 1 \times 1 \\ =& \frac{2 \times c^2 + (k^2 + 1) \times m \times c}{4}, \end{aligned}$$

$$\begin{aligned} \mathcal{O}_{e-ASPP} =& c \times 2 \times 1 \times 1 \times h \times w \times \frac{c}{4} \\ &+ \frac{c}{4} \times m \times (1 \times 2 \times k \times k \times h \times w \times 1) \\ &+ \frac{c}{4} \times (m \times 2 \times 1 \times 1 \times h \times w \times 1) \\ &+ \frac{c}{4} \times 2 \times 1 \times 1 \times h \times w \times c \\ =& (c + \frac{(k^2 + 1) \times m}{2}) \times (h \times w \times c). \end{aligned}$$
$$(14)$$

Following the standard ASPP, we set $k = 3$ and $m = 5$. Usually, $c \geq 256$ is applied in most networks. Therefore, we have:

$$\frac{\mathcal{P}_{e-ASPP}}{\mathcal{P}_{ASPP}} \approx 0.01, \qquad (15)$$

$$\frac{\mathcal{O}_{e-ASPP}}{\mathcal{O}_{ASPP}} \approx 0.01. \qquad (16)$$

It means that compared to the standard ASPP, our proposed e-ASPP has only 1% of the parameters and 1% of the computational overhead. In MODNet, our experiments show that e-ASPP can achieve performance comparable to ASPP. Note that when the Channel Reduction operation in e-ASPP is disabled, e-ASPP still has only 2% of the parameters and 2% of the computational overhead compared to ASPP.

## Appendix B: Results on CRGNN-R and D646

In Table 4, we provide the quantitative results on a video matting dataset proposed by (Wang et al. 2021) to show the effectiveness of the proposed SOC strategy. In Table 5, we compare MODNet with previous SOTA methods on the D646 dataset proposed by (Qiao et al. 2020).
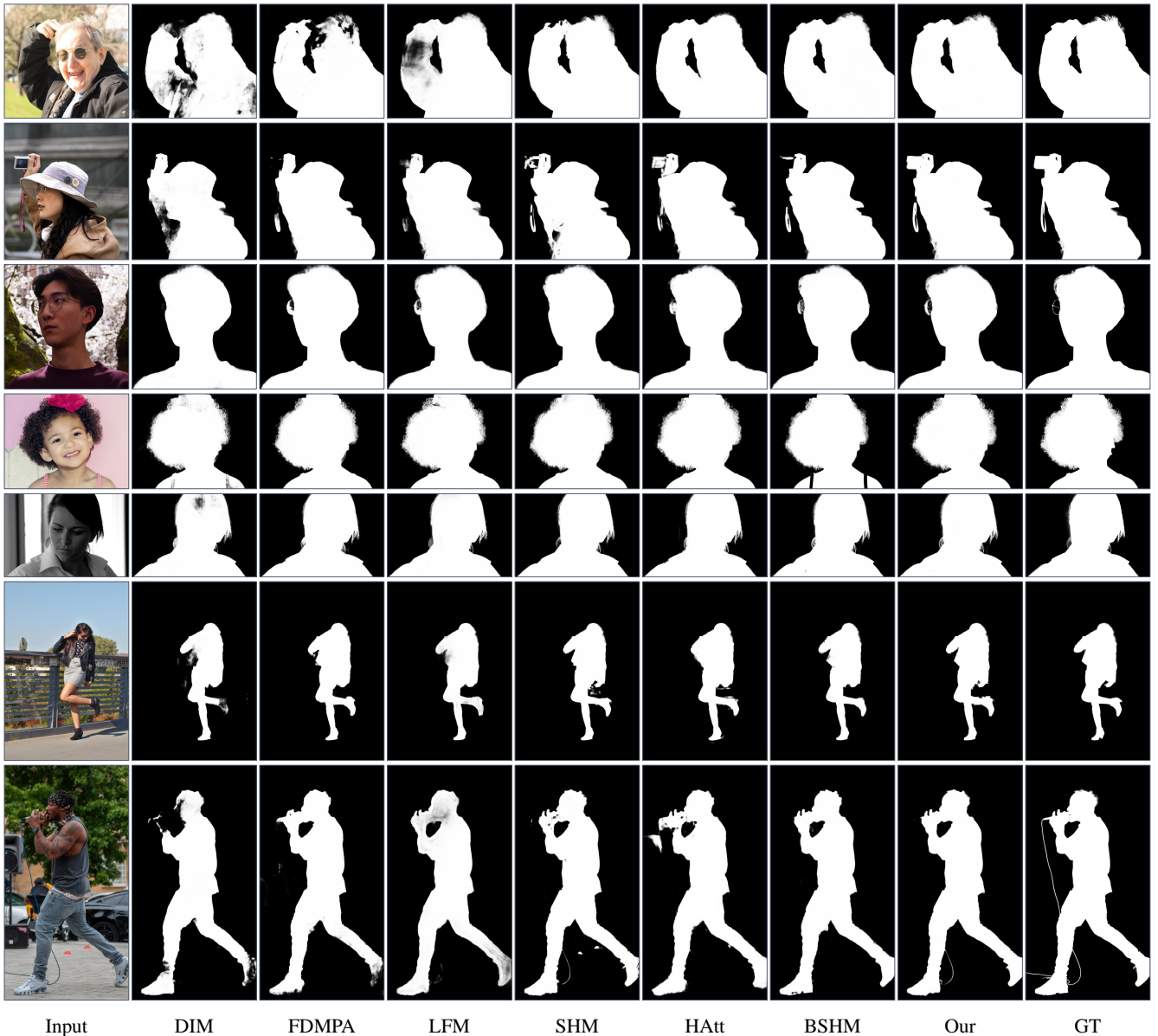
Figure 8: **More Visual Comparisons of Trimap-free Methods on PHM-100.** We compare our MODNet with DIM (Xu et al. 2017), FDMPA (Zhu et al. 2017), LFM (Zhang et al. 2019), SHM (Chen et al. 2018), HAtt (Qiao et al. 2020), and BSHM (Liu et al. 2020). Note that DIM here does not take trimaps as the input but is pre-trained on the SPS (supervise.ly 2018) dataset. Zoom in for the best visualization.

| Method | Trimap | MSE ↓ | MAD ↓ |
|---|---|---|---|
| CRGNN (Wang et al. 2021) | ✓ | 0.0010 | 0.0035 |
| MODNet (Our) | | 0.0082 | 0.0157 |
| MODNet + SOC (Our) | | **0.0033** | **0.0084** |

Table 4: Results on CRGNN-R (Wang et al. 2021).

| Method | Trimap | MSE ↓ | MAD ↓ |
|---|---|---|---|
| DIM (Xu et al. 2017) | ✓ | 0.0025 | 0.0081 |
| HAtt (Qiao et al. 2020) | | 0.0054 | 0.0126 |
| MODNet (Our) | | **0.0037** | **0.0098** |

Table 5: Results on D646 (Qiao et al. 2020).

## Appendix C: Visual Results on PHM-100

Fig. 8 provides more visual comparisons of MODNet and the existing trimap-free methods on PHM-100.

## Appendix D: Comparison with BM

We compare MODNet against the background matting (BM) proposed by (Sengupta et al. 2020). Since BM does not support dynamic backgrounds, we conduct validations in the

**Input**   **BM**   **Our**

Figure 9: **MODNet versus BM with a fixed camera position.** MODNet outperforms BM (Sengupta et al. 2020) when a car is entering the background (red region).

fixed-camera scenes from (Sengupta et al. 2020). BM relies on a static background image, which implicitly assumes that all pixels whose value changes across frames belong to the foreground. As shown in Fig. 9, when a moving object suddenly appears in the background, the result of BM will be affected, but MODNet is robust to such disturbances.

# References

Aksoy, Y.; Aydin, T. O.; and Pollefeys, M. 2017. Designing effective inter-pixel information flow for natural image matting. In *CVPR*.

Aksoy, Y.; Oh, T.-H.; Paris, S.; Pollefeys, M.; and Matusik, W. 2018. Semantic soft segmentation. *TOG*.

Bai, X.; and Sapiro, G. 2007. A geodesic framework for fast interactive image and video segmentation and matting. In *ICCV*.

Cai, S.; Zhang, X.; Fan, H.; Huang, H.; Liu, J.; Liu, J.; Liu, J.; Wang, J.; and Sun, J. 2019. Disentangled Image Matting. In *ICCV*.

Chen, L.-C.; Papandreou, G.; Kokkinos, I.; Murphy, K.; and Yuille, A. L. 2017. Deeplab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected crfs. *PAMI*, 40(4): 834–848.

Chen, Q.; Ge, T.; Xu, Y.; Zhang, Z.; Yang, X.; and Gai, K. 2018. Semantic human matting. In *ACMMM*.

Chen, Q.; Li, D.; and Tang, C.-K. 2013. KNN Matting. *PAMI*.

Cho, D.; Tai, Y.-W.; and Kweon, I. 2016. Natural image matting using deep convolutional neural networks. In *ECCV*.

Chuang, Y.-Y.; Curless, B.; Salesin, D. H.; and Szeliski, R. 2001. A bayesian approach to digital matting. In *CVPR*.

Feng, X.; Liang, X.; and Zhang, Z. 2016. A cluster sampling method for image matting via sparse coding. In *ECCV*.

Foix, S.; Alenyà, G.; and Torras, C. 2011. Lock-in Time-of-Flight (ToF) cameras: A survey. *Sensors Journal*.

Gastal, E. S. L.; and Oliveira, M. M. 2010. Shared sampling for real-time alpha matting. In *Eurographics*.

Grady, L.; Schiwietz, T.; Aharon, S.; and Westermann, R. 2005. Random walks for interactive alpha-matting. In *VIIP*.

He, K.; Rhaemann, C.; Rother, C.; Tang, X.; and Sun, J. 2011. A global sampling method for alpha matting. In *CVPR*.

Hou, Q.; and Liu, F. 2019. Context-aware Image Matting for Simultaneous Foreground and Alpha Estimation. In *ICCV*.

Johnson, J.; Varnousfaderani, E. S.; Cholakkal, H.; and Rajan, D. 2016. Sparse coding for alpha matting. *TIP*.

Karacan, L.; Erdem, A.; and Erdem, E. 2015. Image matting with kl-divergence based sparse sampling. In *ICCV*.

Ke, Z.; Qiu, D.; Li, K.; Yan, Q.; and Lau, R. W. 2020. Guided Collaborative Training for Pixel-wise Semi-Supervised Learning. In *ECCV*.

Kuznetsova, A.; Rom, H.; Alldrin, N.; Uijlings, J. R. R.; Krasin, I.; Pont-Tuset, J.; Kamali, S.; Popov, S.; Malloci, M.; Duerig, T.; and Ferrari, V. 2018. The Open Images Dataset V4: Unified image classification, object detection, and visual relationship detection at scale. *IJCV*.

Levin, A.; Lischinski, D.; and Weiss, Y. 2007. A closed-form solution to natural image matting. *PAMI*.

Levin, A.; Rav-Acha, A.; and Lischinski, D. 2008. Spectral matting. *PAMI*.

Li, Y.; and Lu, H. 2020. Natural image matting via guided contextual attention. In *AAAI*.

Liu, J.; Yao, Y.; Hou, W.; Cui, M.; Xie, X.; Zhang, C.; and Hua, X.-S. 2020. Boosting Semantic Human Matting With Coarse Annotations. In *CVPR*.

Lu, H.; Dai, Y.; Shen, C.; and Xu, S. 2019. Indices Matter: Learning to Index for Deep Image Matting. In *ICCV*.

Qiao, Y.; Liu, Y.; Yang, X.; Zhou, D.; Xu, M.; Zhang, Q.; and Wei1, X. 2020. Attention-Guided Hierarchical Structure Aggregation for Image Matting. In *CVPR*.

Ruzon, M. A.; and Tomasi, C. 2000. Alpha estimation in natural images. In *CVPR*.

Sandler, M.; Howard, A.; Zhu, M.; Zhmoginov, A.; and Chen, L.-C. 2018. MobileNetV2: Inverted Residuals and Linear Bottlenecks. In *CVPR*.

Schmarje, L.; Santarossa, M.; Schröder, S.-M.; and Koch, R. 2020. A survey on Semi-, Self- and Unsupervised Learning for Image Classification. *ArXiv*, abs/2002.08721.

Sengupta, S.; Jayaram, V.; Curless, B.; Seitz, S.; and Kemelmacher-Shlizerman, I. 2020. Background Matting: The World is Your Green Screen. In *CVPR*.

Shen, X.; Tao, X.; Gao, H.; Zhou, C.; and Jia, J. 2016. Deep automatic portrait matting. In *ECCV*.

Sun, B.; Feng, J.; and Saenko, K. 2016. Return of Frustratingly Easy Domain Adaptation. In *AAAI*.

Sun, J.; Jia, J.; Tang, C.-K.; and Shum, H.-Y. 2004. Poisson matting. *TOG*.

supervise.ly. 2018. Supervisely Person Dataset. *supervise.ly*.

Tang, J.; Aksoy, Y.; Oztireli, C.; Gross, M.; and Aydin, T. O. 2019. Learning-based Sampling for Natural Image Matting. In *CVPR*.

Toldo, M.; Michieli, U.; Agresti, G.; and Zanuttigh, P. 2020. Unsupervised Domain Adaptation for Mobile Semantic Segmentation based on Cycle Consistency and Feature Alignment. *IMAVIS*.

Wang, T.; Liu, S.; Tian, Y.; Li, K.; and Yang, M.-H. 2021. Video Matting via Consistency-Regularized Graph Neural Networks. In *ICCV*.

Wilson, G.; and Cook, D. J. 2020. A Survey of Unsupervised Deep Domain Adaptation. *TIST*.

Xu, N.; Price, B.; Cohen, S.; and Huang, T. 2017. Deep Image Matting. In *CVPR*.

Yang, X.; Xu, K.; Chen, S.; He, S.; Yin, B. Y.; and Lau, R. 2018. Active matting. *Adv. Neural Inform. Process. Syst.*

Zhang, Y.; Gong, L.; Fan, L.; Ren, P.; Huang, Q.; Bao, H.; and Xu, W. 2019. A late fusion cnn for digital matting. In *CVPR*.

Zhu, B.; Chen, Y.; Wang, J.; Liu, S.; Zhang, B.; and Tang, M. 2017. Fast Deep Matting for Portrait Animation on Mobile Phone. In *ACMMM*.