

# MP-MAT: A 3D-AND-INSTANCE-AWARE HUMAN MATTING AND EDITING FRAMEWORK WITH MULTIPLANE REPRESENTATION

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Human instance matting aims to estimate an alpha matte for each human instance in an image, which is challenging as it easily fails in complex cases requiring disentangling mingled pixels belonging to multiple instances along hairy and thin boundary structures. In this work, we address this by introducing a novel 3D-and-instance-aware matting framework with multiplane representation, where the multiplane concept is designed from two different perspectives: scene geometry level and instance level. Specifically, we first build feature-level multiplane representations to split the scene into multiple planes based on depth differences. This approach makes the scene representation 3D-aware, and can serve as an effective clue for splitting instances in different 3D positions, thereby improving interpretability and boundary handling ability especially in occlusion areas. Then, we introduce another multiplane representation that splits the scene in an instance-level perspective, and represents each instance with both matte and color. We also treat background as a special instance, which is often overlooked by existing methods. Such an instance-level representation facilitates both foreground and background content awareness, and is useful for other down-stream tasks like image editing. Once built, the representation can be reused to realize controllable instance-level image editing with high efficiency. Extensive experiments validate the clear advantage of MP-Mat in matting task. We also demonstrate its superiority in image editing tasks, an area under-explored by existing matting-focused methods, where our approach under zero-shot inference even outperforms trained specialized image editing techniques by large margins. Code will be released to inspire relevant fields.

## 1 INTRODUCTION

Human matting is one of the foundation tasks in computer vision that can widely serve for applications such as image editing, image compositing, and film post-production (Zhu et al., 2017; Chen et al., 2018; Sengupta et al., 2020; Lin et al., 2021; 2023). Despite the development of effective algorithms, most methods focus on human matting under single-instance scenarios, which cannot fully align with real-world applications, where multiple instances could exist in a complex scene. In recent years, InstMatte (Sun et al., 2022) formally introduced the multi-instance matting formulation, which separates the image into a combination of multiple instance layers and background layer:

$$I = \sum_{i=1}^n \alpha_i C_i + \left(1 - \sum_{i=1}^n \alpha_i\right) B \quad (1)$$

where  $\alpha_i \in [0, 1]$  denotes the opacity (alpha matte) of the  $i$ -th foreground, whose value is the ultimate task goal. The task is actually an ill-defined problem since the foreground color  $C_i$ , background color  $B$  and the alpha value  $\alpha_i$  are left unknown. Compared with single-instance assumption, multi-instance matting poses additional challenges. Specifically, the algorithm should be instance-aware (i.e., can localize and distinguish different human instances), and also needs to preserve complex and fine instance edges. Maintaining the integrity of each instance without blurring the edges is particularly challenging in cases where instances are in contact or occluded (Ke et al., 2021).

A core motivation of this work is to establish layered representation to facilitate instance matting in complex scenarios, where we split the scene into different layers based on 2 different perspectives:

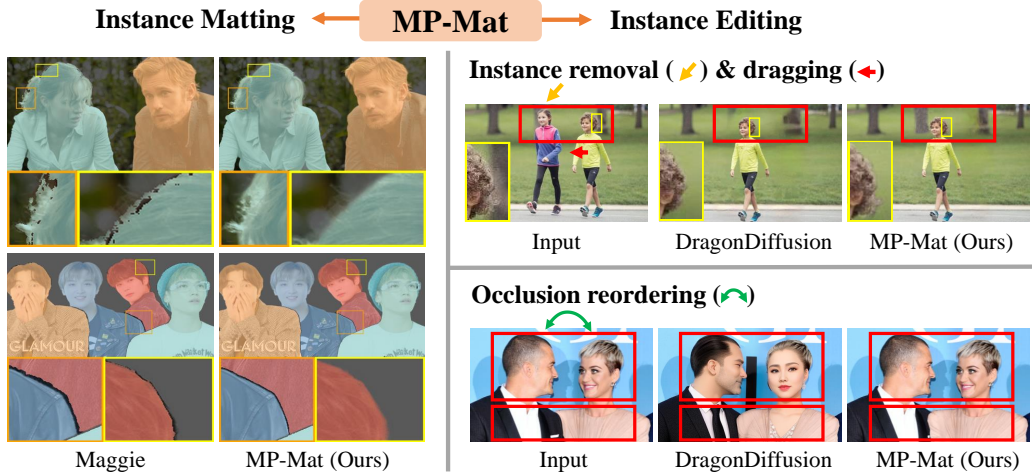


Figure 1: The proposed MP-Mat can perform well in both instance matting and editing tasks, outperforming existing state-of-the-art specialist models designed for individual tasks. Distinguished areas are highlighted with bounding boxes, where MP-Mat preserves finer details and better retains regions that should remain semantically unchanged.

depth-oriented and instance-oriented. Our design of such layered concept is partially inspired by the idea of Multiplane Image (MPI) (Tucker & Snavely, 2020), a learnable 3D representation that is proposed for novel view synthesis. MPI aims to learn a 3D scene representation with multiple RGBA planes from source image, which can then be used to synthesize different novel views of the scene. The learned multiple RGBA planes can split the scene based on the depth difference. We argue that such a depth-oriented plane-splitting idea could be potentially helpful for matting in multi-instance scenarios. Specifically, different instance may lie in different depth positions, so they can be effectively distinguished by grouping them into different planes that represent a distinct depth level, which could be potentially helpful for handling occlusions under complex scenes. However, we argue that naively using MPI for instance matting is not feasible, as it is build on low pixel-level based on depth differences for only pixel-level novel view synthesis and without instance-awareness.

Based on the aforementioned analysis, we propose MP-Mat, a 3D-and-instance-aware matting framework that is built on meticulously designed multiplane representations. Specifically, our formulation mainly consists of 2 parts: scene geometry-level multiplane representation (SG-MP) and instance-level multiplane representation (Inst-MP). The SG-MP is built to split the scene into multiple planes based on the depth differences, making the scene representation 3D-aware, and can serve as an effective clue for splitting instances in different 3D positions. Different from the existing MPI representations, the proposed SG-MP is built on feature-level, rather than the low pixel-level. The benefits lie in 2 aspects: (1) Compared with low-level RGBA, building MP features on high-level deep features can contain more semantic information to better represent the scene geometry as well as fine-grained texture context, which can be useful for the subsequent instance-level analysis based on it; (2) When optimized together with subsequent instance-level Inst-MP for instance-level perception tasks (e.g., instance localization and matting), the SG-MP feature will also receive relevant gradient, making the plane division and scene representation become more instance-aware.

Besides the built SG-MP, we also introduce an instance-level multiplane representation (Inst-MP) that splits the scene in an instance-level perspective, and represent each instance with both matte and color (as shown in Fig. 2). Our formulation has several benefits: (1) Different from most matting methods that only predict alpha matte, our proposed representation additionally estimates the color of each foreground, enabling better foreground content awareness, resulting in a higher matting accuracy; (2) Besides foreground instances, Inst-MP also explicitly models the background as a special instance. Such formulation is different from existing works that only focus on foreground matte modeling, and enables better handling of the boundary of instance and background, especially when occlusion occurs; (3) The proposed instance-level multiplane representation obeys the integration property for image rendering (i.e., the integral of the representation equals the whole RGB image), and thus can be easily used for downstream tasks like image editing where instance-level manipulation can be directly processed on separate feature planes. Once built, the representation can be reused to realize controllable instance-level image editing with very high efficiency.

We conduct extensive experiments to demonstrate the superiority of the proposed framework. Specifically, for instance matting, our MP-Mat outperforms existing methods by large margins (i.e., at least 2.76 SAD in HIM-100K dataset (Liu et al., 2024) and 4.16 SAD in SMPMat dataset (Jiao et al., 2024)). Besides, we also take a further step to explore the capacity of MP-Mat on image editing task that is actually under-explored by existing matting-focused methods. We found that our method can also perform well on this potential downstream task and even outperforms existing specialized image editing techniques by large margins and with high efficiency, even under zero-shot inference. A qualitative comparison is shown in Fig. 1, where distinguished regions are highlighted with bounding boxes. For both types of tasks, MP-Mat can actually preserve finer details and better retains regions that should remain semantically unchanged. We hope our work can inspire future research in related fields, including but not limited to image matting and image editing.

## 2 RELATED WORK

### 2.1 INSTANCE MATTING

Elder matting works (Xu et al., 2017; Tan et al., 2018; Lu et al., 2019; Zhang et al., 2019; Qiao et al., 2020; Li & Lu, 2020; Sun et al., 2021; Yu et al., 2021b; Ke et al., 2022; Li et al., 2022; Ma et al., 2023) assume a single-instance condition without multi-instance awareness, which remains a gap with many real-world scenarios. Human instance matting is a recently emerged task that differs from the traditional one, as it requires simultaneously localizing multiple instances and distinguishing their mattes in an instance-level manner. Such a task poses more challenges as it easily fails in complex cases requiring disentangling mingled pixels belonging to multiple instances along hairy and thin boundary structures. Mainstream methods (Hu & Clark, 2019; Sun et al., 2022; Huynh et al., 2024) for this task typically rely on instance-level masks (He et al., 2017; Wang et al., 2020) as input and gradually refine them to predict the matte. While these methods provide a certain level of instance-awareness, it is relatively weak as they largely depend on a pre-set instance segmentation module. Although a recent work (Liu et al., 2024) achieves independent instance-aware capability, its performance still falls short of the SoTA methods in the mainstream setting, primarily due to the lack of external guidance. Different from these works, we build layered representations to emphasize 3D and full instance awareness. Specifically, the build SG-MP decomposes the scene based on depth variation, thereby improving interpretability and boundary handling ability especially in occlusion areas. Our Inst-MP, besides representing the matte of foreground instances, also explicitly models the background and color of each instance, resulting a better context awareness for both foreground and background. Inst-MP also processes good properties for downstream image editing tasks that are under-explored by existing matting-focused methods.

### 2.2 INSTANCE CONTROLLABLE IMAGE EDITING

The task aims to impose instance-level editing (e.g., instance removal or dragging) on the image in a harmonious way. Recently, diffusion models (Yildirim et al. (2023); Shi et al. (2024); Ekin et al. (2024); Sheynin et al. (2024); Yang et al. (2024)) have brought significant breakthroughs in image editing tasks. Despite their effectiveness, they generally need more denoising steps to generated high quality result, which is usually more time costly. Existing methods are also of weak instance-aware capabilities and cannot fully guarantee that theoretically unchanged regions remain unaffected, resulting in poor performance on instance-level editing tasks. In this work, the intrinsic property of our proposed instance-level multiplane representation (Inst-MP) can also enable an easier and more direct way to achieve instance-level editing, where accurate instance-level feature manipulation can be done on separate feature planes without affecting other regions, thereby enhancing the consistency of the edited image. Experiments show that our approach, even under zero-shot inference, significantly outperforms trained specialized image editing techniques. Another distinguished advantage is that once Inst-MP is built, subsequent editing of the image takes negligible time. Besides, we also explore our potential in image editing tasks, which, to our knowledge, has not been well concerned by previous matting-focused works. We hope our work can inspire future research in related fields, including but not limited to image matting and image editing.

### 2.3 3D SCENE REPRESENTATION

3D scene representation has been widely used for tasks like 3D reconstruction (Mescheder et al., 2019; Wu et al., 2024) and view synthesis (Kong et al., 2024; Luiten et al., 2024). In this work, one of our motivations is to make instance matting become more 3D-aware. Our design is partially inspired by the idea of Multiplane Image (MPI) (Tucker & Snavely, 2020), a learnable 3D representation that

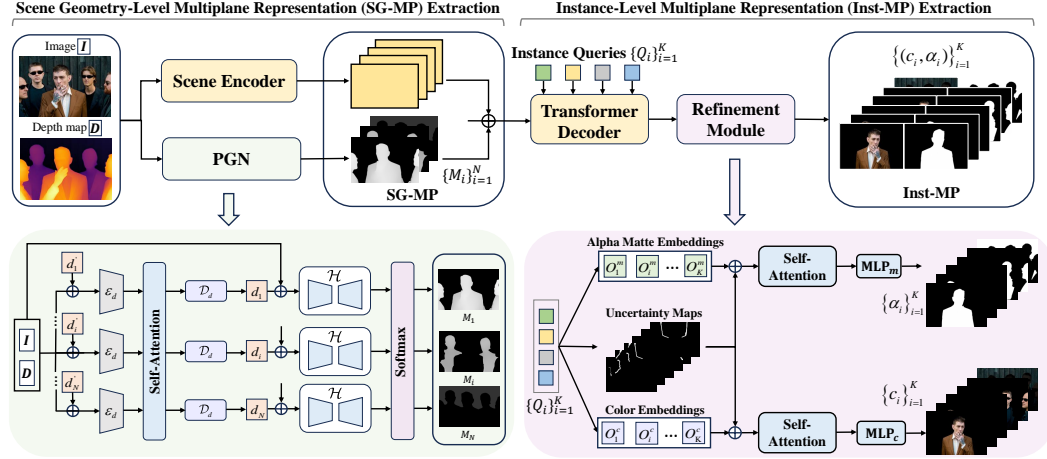


Figure 2: The overall framework of the proposed MP-Mat.  $\oplus$  indicates concatenation operation.

is proposed for novel view synthesis. MPI aims to learn multiple planes that decompose the scene according to depth variations. We argue that such a depth-oriented plane-splitting idea could be potentially helpful for matting in multi-instance scenarios. However, naively using MPI for instance matting is not feasible, as it is built on low pixel-level based on depth differences for only pixel-level novel view synthesis and without instance awareness. On the contrary, our formulation differs to it within 3 aspects: (1) MPI split the planes based on depth variations, while we design the multiplane representation from 2 different perspectives (scene geometry level and instance level) to be better aware of both depth and instance-level information; (2) MPI is built on low pixel-level, while our scene geometry level representation is built on high-level deep features, which is of better capacity for understanding scene geometry as well as fine-grained texture context. It also enables flexible end-to-end training with the subsequent instance-level task to enrich its instance-awareness ability; (3) MPI is generally used for low pixel-level tasks like novel view synthesis (Tucker & Snavely, 2020) and bokeh rendering (Peng et al., 2022; Rao, 2023) that lack of instance awareness, while our goal is to solve instance-level perception problems, with different formulation from original MPI.

### 3 METHOD

As shown in Fig. 2, our MP-Mat framework mainly consists of 2 parts: scene geometry-level multiplane representation and instance-level multiplane representation. In the following, we will first introduce MP-Mat for instance matting in detail. Then, we will also describe how our method can be applied for instance-controllable image editing tasks.

#### 3.1 SCENE GEOMETRY-LEVEL MULTIPLANE REPRESENTATION (SG-MP)

We first build multiplane representations to split the scene into multiple planes based on depth differences. This approach makes the scene representation 3D-aware, and can serve as an effective clue for splitting instances in different 3D positions, thereby alleviating occlusion effects. To be specific, we introduce plane-distinctive masks to represent plane information, which is obtained by a plane generation network. They will be concated with the deep feature encoded from RGBD input, to form the so-called scene geometry-level multiplane representation.

##### 3.1.1 PLANE GENERATION NETWORK (PGN)

This network aims to generate the plane-distinctive masks based on the input single view RGBD image, for plane information representation. To be more specific, we first define  $N$  planes that aim to split the scene based on distinction in depth and instance-level information, which we represent them using what we call plane-distinctive masks. The PGN will gradually generate  $N$  plane-distinctive masks  $\{M_i\}_{i=1}^N$  based on a set of predefined initial depths  $\{d'_i\}_{i=1}^N$  (uniformly sampled according to the scene depth map) and the actual RGBD image inputs, where the final refined masks will simultaneously possess scene depth-distinctive and instance-aware characteristics. Note that the

input depth map here can be easily estimated using off-the-shelf depth estimators, and the actual cost is comparable to pre-instance mask generation in existing mainstream mask-guided instance matting methods (Sun et al., 2022; Huynh et al., 2024).

As illustrated in Fig. 2, we first use a shared lightweight CNN  $\mathcal{E}_d$  to extract a global feature  $f'_i$  for plane  $i$  at the initial depth  $d'_i$ :

$$f'_i = \mathcal{E}_d(I, D, d'_i) \quad (2)$$

Then we apply the self-attention operation to  $\{f'_i\}_{i=1}^N$  to obtain  $\{f_i\}_{i=1}^N$ :

$$\{f_i\}_{i=1}^N = \text{Self-Attention}(\{f'_i\}_{i=1}^N) \quad (3)$$

The intuition here is to adjust the corresponding depth information of each plane at the feature level by exchanging the geometry and appearance information among  $\{f'_i\}_{i=1}^N$ . The adjusted feature  $f_i$  is then decoded to the adjusted depth  $d_i$  using a shared MLP  $\mathcal{D}_d$ :

$$d_i = \mathcal{D}_d(f_i) \quad (4)$$

After plane depth adjustment, we then use an interpreter to generate the plane-distinctive masks based on the adjusted plane depths and the original RGBD inputs:

$$\{M_i\}_{i=1}^N = \text{Softmax}(\{\mathcal{H}(I, D, d_i)\}_{i=1}^N) \quad (5)$$

where  $\mathcal{H}$  is a UNet-like architecture, and the spatial resolution of  $M_i$  aligns with the extracted deep feature from the scene encoder. The mask delicately allocates every visible pixel from the source viewpoint to each plane. We concat the plane-distinctive masks and the deep scene feature to form the so-called scene geometry-level multiplane representation.

### 3.2 INSTANCE-LEVEL MULTIPLANE REPRESENTATION (INST-MP)

Given the obtained scene geometry-level multiplane representation, we subsequently use it to derive instance-level multiplane representation  $\{(c_i, \alpha_i)\}_{i=0}^S$ , where  $i$  serves as the plane index, with  $i = 0$  corresponding to the background plane, and  $i$  ranging from 1 to  $S$  representing instance-level plane for each distinct human instances within the image.

Actually, the ultimate prediction of instance matting task (i.e., instance-level matte  $\{\alpha_i\}_{i=1}^S$ ) is a sub-set of our instance-level multiplane representation. Our formulation has several benefits: (1) Different from most matting methods that only predict alpha matte, our proposed representation additionally estimates the color of each foreground, enabling better foreground content awareness; (2) Different from existing works that only focus on foreground matte modeling, our representation also explicitly model the background by regarding it as a special foreground instance. Such formulation enables better handling of the boundary of instance and background, especially when occlusion occurs; (3) The proposed instance-level multiplane representation obeys the integration property for image rendering:

$$I = \sum_{i=0}^S c_i \alpha_i \quad (6)$$

Based on this property, we can reuse the built representation to realize controllable instance-level image editing with very high efficiency (please refer to Sec. 4 for details). In the following of this sub-section, we will introduce how to obtain the instance-level multiplane representation in detail.

#### 3.2.1 INSTANCE QUERY

Here we use learnable queries  $\{Q_i\}_{i=1}^K$  ( $K > S + 1$ ) to capture instance-level features, where  $S$  denotes the actual instance amount within image and we also regard background as a special instance, resulting in the total number as  $S + 1$ . The queries will collect the corresponding instance-level features from the obtained scene geometry-level multiplane representation, through a standard multi-layer transformer decoder that consists of multi-head self-attention and cross-attention with fully-connected feed-forward networks (FFNs). Then, we use different prediction heads (composed

of MLP) to map the query feature  $Q_i$  into distinct embeddings: opacity embeddings  $O_i^m \in \mathbb{R}^C$  that gauge transparency, and color embeddings  $O_i^c \in \mathbb{R}^C$  that capture hue information, where  $C$  is the feature dimension. We also derive an uncertainty map  $O_i^u \in \mathbb{R}^{C \times H \times W}$  ( $H \times W$  corresponds to the spatial resolution of input image) to estimate the uncertainty of the predictions, which will be used for subsequent feature refinement.

### 3.2.2 REFINEMENT MODULE

We design a refinement module that utilizes the uncertainty map  $\{O_i^u\}_{i=1}^K$  to enhance  $\{O_i^c\}_{i=1}^K$  and  $\{O_i^m\}_{i=1}^K$ , and finally output the instance-level multiplane representation. Specifically, given the acquired uncertainty map, the top  $T\%$  (i.e., 10% in our implementation) of pixels with the highest uncertainty values are selected, yielding filtered guidance masks  $\{R_i\}_{i=1}^K$  for refinement. The mask  $R_i$  is subsequently concatenated with  $O_i^m$  and  $O_i^c$  to indicate the region of high uncertainty for each instance-level representation. Then, self-attention is adopted among different instance-level features to refine (reconsider) the representations of each instance:

$$\begin{aligned} \{O_i^m\}_{i=1}^K &= \text{Self-Attention}(\{O_i^m, R_i\}_{i=1}^K) \\ \{O_i^c\}_{i=1}^K &= \text{Self-Attention}(\{O_i^c, R_i\}_{i=1}^K) \end{aligned} \quad (7)$$

Finally, we predict the alpha matte  $\alpha_i$  and the color  $c_i$  based on  $O_i^m$  and  $O_i^c$ :

$$\begin{aligned} \{\alpha_i\}_{i=1}^K &= \text{MattePredictor}(O_i^m) \\ \{c_i\}_{i=1}^K &= \text{ColorPredictor}(O_i^c) \end{aligned} \quad (8)$$

where both `MattePredictor` and `ColorPredictor` are two-layer MLPs that decode the alpha/color embeddings to predict the final alpha matte/color.

### 3.2.3 MODEL TRAINING

The whole MP-Mat framework is trained in an end-to-end manner. We employ a multi-task loss,

$$\begin{aligned} \mathcal{L} &= \lambda_1 \mathcal{L}_{Detect} + \lambda_2 \mathcal{L}_{Matting} \\ \mathcal{L}_{Matting} &= \mathcal{L}_{alpha} + \mathcal{L}_{lap} + \mathcal{L}_{comp} \end{aligned} \quad (9)$$

where  $\mathcal{L}_{Detect}$  is the loss for instance detection, including localization and classification loss. The detection predictions are derived from instance-level matte results (the bounding box that can tightly cover the pixel with  $m > 0$ ). We then perform bipartite matching (Carion et al., 2020) between the instance-level detection predictions and GTs to obtain the correspondence between predictions and GTs, and then calculate the standard matting loss  $\mathcal{L}_{Matting}$ , which encompasses alpha loss, pyramid Laplacian loss and composition loss following standard formulation (Sun et al., 2022).

## 4 MP-MAT FOR INSTANCE CONTROLLABLE IMAGE EDITING

Due to the good properties of Inst-MP as mentioned in Sec. 3.2, MP-Mat can achieve instance-level image editing in a direct, accurate, and nearly free manner. The overarching idea is that instance-level editing can be achieved by directly processing on separate instance-level planes within Inst-MP. Once the Inst-MP is built, subsequent image editing will only take negligible time, and the editing process is training-free. Here we will describe how it works for 3 different sub-tasks in detail.

**Instance removal** aims to delete specified foreground instance from the image without affecting its overall harmony. This problem can be treated as reassigning the alpha matte of the removed instance  $j$  to its backward instances (including background). For target instance  $j$ , We first define  $\Omega_j$  as the set of pixel  $(x, y)$  where  $\alpha_j(x, y) > 0$ . For every pixel in  $\Omega_j$ , we assign its alpha to the instance plane  $t$  that is closest behind  $j$  (can be estimated by the average depth of the plane):

$$\alpha_t(x, y) = \alpha_t(x, y) + \alpha_j(x, y). \quad (10)$$

Then, the edited image can be represented as:

$$I' = \sum_{i=0, i \neq j}^S c_i \alpha_i. \quad (11)$$

**Occlusion reordering** aims to switch the occlusion relationship between the source instance and target instance (e.g., changing from instance  $p$  occluding  $q$  to  $q$  occluding  $p$ ). We first consider a simplified case where no additional instances are positioned between  $p$  and  $q$  (indicated by depth position) for demonstration purposes. In this case, the task can be done by swapping the alpha values of the two instances within their intersection regions:

$$\alpha'_q(x, y) = \alpha_p(x, y), \quad (12)$$

$$\alpha'_p(x, y) = \alpha_q(x, y), \quad (13)$$

where  $(x, y)$  refers to the pixel position that satisfy  $(c_q(x, y) > 0) \wedge (c_p(x, y) > 0)$ . The edited image can be then synthesized with the updated  $\alpha_i$ :

$$I' = \sum_{i=0, i \neq p, i \neq q}^S c_i \alpha_i + c_p \alpha'_p + c_q \alpha'_q. \quad (14)$$

For general cases where other instances may exist between  $p$  and  $q$ , we first sort the plane based on their depth from the closest to the farthest to the camera plane, resulting in a sorted plane index set  $\{\dots, \mathbf{p}, p+1, p+2, \dots, q-1, \mathbf{q}, \dots\}$ . Then, we perform the aforementioned alpha swap between each pair of adjacent planes iteratively until the desired order is achieved (i.e.,  $\{\dots, \mathbf{q}, p+1, p+2, \dots, q-1, \mathbf{p}, \dots\}$ ). We provide more detailed description and pseudo code in the supplementary material.

**Instance dragging** aims to drag a target instance to a new desired position, which can be divided into two categories: drag across images and drag within one image. The drag across image task can be divided into three steps: (1) Feed the reference image  $I_{ref}$  into MP-Mat to get its Inst-MP, and separate the plane  $t$  ( $c_t, \alpha_t$ ) that corresponds to the target instance. (2) Crop ( $c_t, \alpha_t$ ) to form ( $c'_t, \alpha'_t$ ) by extracting the rectangular region that tightly covers the pixels with positive alpha value. (3) Set up a new plane ( $c_{new}, \alpha_{new}$ ) with zero initialization (i.e., all pixel values are 0), and add the cropped ( $c'_t, \alpha'_t$ ) to it at the desired position on the target image, resulting in ( $c'_{new}, \alpha'_{new}$ ). Note that additional transformations, such as rescaling and rotation, can also be applied to ( $c'_t, \alpha'_t$ ) before adding. (4) Add the resulting new plane to the target image  $I$ :

$$I' = c'_{new} \alpha'_{new} + (1 - \alpha'_{new}) I. \quad (15)$$

For dragging within one image, it can be regarded as a combination of instance removal and dragging across images when the target image remains the same as source.

## 5 EXPERIMENTS

### 5.1 IMPLEMENTATION DETAILS

For both training and testing, we obtain the input depth map through an off-the-shelf monocular depth estimator (Yin et al., 2023). We train MP-Mat on 4 NVIDIA RTX 2080Ti GPUs with a total batch size of 4 (1 per GPU). The training employs the SGD optimizer with a momentum of 0.9 and a weight decay of 0.0005 for 50,000 iterations. The learning rate is initialized at 0.01 and is adjusted by multiplying with  $\left(1 - \frac{iter}{max-iter}\right)^{0.9}$ . Additional implementation specifics are detailed in the supplementary material. For the uncertainty map generation, the hyperparameter  $K$  is set to 5. In the loss function, the hyperparameters are  $\lambda_1 = 1$  and  $\lambda_2 = 5$ .

### 5.2 TASK, DATASET, AND METRICS

**Instance matting.** We first validate MP-Mat on the multi-instance matting task, where we conduct experiments on 2 datasets: the real image subset of the HIM-100k dataset (Liu et al., 2024), which contains 47,980 real images with ground truth manually annotated, and the SMPMat dataset, a recent synthetic matting dataset containing 40,000 synthetic images. We evaluate different models using two major quantitative metrics: sum of absolute differences (SAD) and mean square error (MSE, we report the  $10^2$  scaled value). Lower values for these metrics indicate better alpha matte results.

**Instance editing.** This is one of the potential downstream applications for matting, but existing matting-focused research has not explored their actual usability and performance on related tasks. To fill this gap, we conduct extensive experiments on 3 sub-tasks with compelling application value: (1)

Table 1: Quantitative comparison of instance matting.

Method	Dataset			
	HIM-100k		SMPMat	
	SAD	MSE	SAD	MSE
<b><i>Instance-agnostic</i></b>				
FBA(+Mask RCNN)	38.25	0.95	42.36	1.12
FBA(+SOLO)	38.18	0.94	41.23	0.96
FBA(+EVA)	37.76	0.91	39.82	0.95
MG(+Mask RCNN)	40.51	0.97	41.58	1.06
MG(+SOLO)	39.26	0.95	40.79	0.95
MG(+EVA)	38.19	0.94	39.73	0.95
<b><i>Instance-aware</i></b>				
InstMatte	37.34	0.93	40.19	0.96
E2E-HIM	32.22	0.84	38.55	0.93
Maggie	29.48	0.78	37.41	0.91
MP-Mat (Ours)	<b>26.75</b>	<b>0.49</b>	<b>33.25</b>	<b>0.81</b>

Table 2: Effects of SG-MP and Inst-MP in a global perspective.

SG-MP	Inst-MP	SAD	MSE
		35.85	0.87
	✓	33.78	0.86
✓		29.46	0.60
✓	✓	<b>26.75</b>	<b>0.49</b>

Table 3: Components effects within SG-MP extraction.

Depth	PGN	SAD	MSE
		33.78	0.86
✓		31.82	0.74
✓	✓	<b>26.75</b>	<b>0.49</b>

Table 4: Component effect in Inst-MP extraction.

Background Estimation	Color Estimation	Refinement	SAD	MSE
			29.46	0.60
✓			28.53	0.55
✓	✓		27.79	0.52
✓	✓	✓	<b>26.75</b>	<b>0.49</b>

Instance removal that aims to delete specified foreground instances from the image without affecting its overall harmony; (2) Occlusion reordering that aims to modify the occlusion relationships among foreground instance in a controlled and harmonious manner; (3) Instance dragging that aims to move any instance to the desired position harmoniously. conduct experiments on the GQA-inpaint dataset (Yildirim et al., 2023) for the instance removal task and use Mean L1 Loss, Mean L2 Loss, and PSNR metrics for evaluation, where lower values for Mean L2 Loss and Mean L1 Loss are preferable, while higher PSNR values indicate better quality. For occlusion reordering task, due to the lack of data with ground truth, we construct a synthetic dataset ORHuman that can theoretically ensure the correctness of the derived ground truth (see supplementary material for details). For instance dragging task, we only give qualitative comparisons due to the lack of evaluation benchmarks.

### 5.3 INSTANCE MATTING

**Compared methods.** Our method is compared with the methods designed for this tasks, including InstMatte (Sun et al., 2022), E2E-HIM (Liu et al., 2024) and Maggie (Huynh et al., 2024). We also compared with composite approaches following Liu et al. (2024). Specifically, we tailor the mask-guided matting model MG (Yu et al., 2021a) and FBA (Forte & Pitié, 2020) with off-the-shelf instance segmentation models (Mask-RCNN He et al. (2017), SOLO Wang et al. (2020), and EVA Fang et al. (2023)) to adapt it to the multi-instance matting task.

**Main results.** The main performance comparison on HIM-100K and SMPMat dataset is shown in Tab. 1. It can be observed that the proposed MP-Mat outperforms existing methods by large margins on both datasets (i.e., at least 2.76 SAD in HIM-100K dataset and 4.16 SAD in SMPMat dataset), demonstrating the superiority of the proposed method.

**Qualitative analysis.** We also give some qualitative results in Fig. 1 and Fig. 3. It can be observed that our method is superior in (1) distinguishing fine-grained boundaries such as hair regions, and (2) with better instance-aware ability, especially under occlusion areas caused by human interactions. More results can be found in the supplementary material.

### 5.4 ABLATION STUDIES

Here we conduct ablation studies on the HIM-100K dataset to analyze the effectiveness of the proposed components. Specifically, we first validate the individual effects of the proposed multiplane representations (i.e., SG-MP and Inst-MP) from a global perspective. Then, we go deeper to analyze

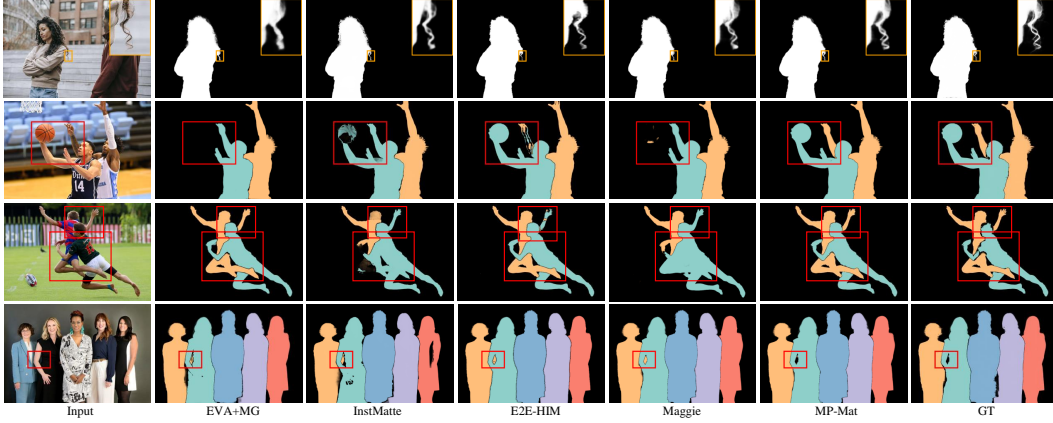


Figure 3: Qualitative comparisons among different methods. Distinguished areas are highlighted with bounding boxes.

the effect of each component within the proposed SG-MP and Inst-MP, demonstrating their effects while also revealing insights for further research. Besides, we also conduct parameter analysis used in the proposed refinement module.

**Effects of the proposed multiplane representations.** Here we analyze the individual effect of the proposed scene geometry-level multiplane representation (SG-MP) and instance-level multiplane representation (Inst-MP) from a global perspective. From Tab. 2, it can be seen that: (1) Without SG-MP, the performance drops drastically (SAD from 26.75 to 33.78). This validates the importance of explicit 3D scene representation for multi-instance matting task. Our designed SG-MP makes the scene representation 3D-aware, and can serve as an effective clue for splitting instances in different 3D positions, thereby alleviating occlusion effects; (2) Inst-MP can further boost SAD by a large margin (i.e., 2.71), verifying its effectiveness. Another benefits of such design is for its flexibility and high efficiency on down stream tasks like instance-level image editing, as discussed in Sec. 5.5.

**Component effects within SG-MP extraction.** From Tab. 3, it can be summarized that: (1) Depth information is useful, as it can bring auxiliary 3D information; (2) Only sending depth map as input is not sufficient enough. With our proposed Plane Generation Network (PGN) for multiplane representation at feature level, the performance further boosts by 5.07 in SAD, which is much larger than the performance gain solely from depth input (1.96 in SAD). This essentially demonstrates the effectiveness of our explicit multiplane feature representation extraction design.

**Component effects within Inst-MP extraction.** From Tab. 4, it can be summarized that: (1) Besides focusing on foreground instances, adding explicit background estimation is beneficial for matting accuracy, as it enables better handling of the boundary of instance and background, especially when occlusion occurs; (2) Besides matte estimation, adding color estimation for each instance can enable better content awareness and thus boost the performance; (3) The proposed uncertainty guided refinement can further facilitate performance, as it can adjust ambiguity at fine-grained boundaries.

## 5.5 INSTANCE EDITING

For all tasks mentioned in Sec. 5.2, we compare our method with SOTA image editing-focused methods (Yildirim et al., 2023; Shi et al., 2024). We finetune the image editing-focused methods on the target dataset for higher performance, while our MP-Mat adopts a zero-shot inference manner.

**Instance removal.** From Tab. 5, we can observe that: (1) MP-Mat significantly outperforms existing editing-based methods and with high efficiency (also refer to Fig. 4 (a) for qualitative comparison). (2) When combined with an off-the-shelf background inpainting model (Yu et al., 2018), the performance of MP-Mat can be further boosted by large margins, we attribute this to the inadequate background modeling capacity of our matting-based methods that is also partially caused by a lack of training data for such themes. (3) The editing time cost within the same image becomes negligible once Inst-MP is constructed in MP-Mat, further verifying the superiority of our design.

Table 5: Quantitative comparison on instance removal. BI refers to an off-the-shelf background inpainting method used to enhance the inpainting of the corresponding region. Bold text indicates the best performance, and underlined text represents the second-best performance.

Editing Method	Mean L1 Loss ( $\downarrow$ )	Mean L2 Loss ( $\downarrow$ )	PSNR ( $\uparrow$ )	Time per editing (s)	
Inst-inpaint	12.69%	2.58%	23.09db	0.1982	
Dragon Diffusion	12.13%	2.49%	22.17db	0.2139	
Matting Method	Mean L1 Loss ( $\downarrow$ )	Mean L1 Loss ( $\downarrow$ )	PSNR ( $\uparrow$ )	Time per image (s)	Time per editing (s)
Ours	<u>9.37%</u>	<u>1.96%</u>	<u>23.58db</u>	<b>0.1117</b>	<b>0.0003</b>
Ours (w/ BI)	<b>3.73%</b>	<b>0.84%</b>	<b>25.79db</b>	<b>0.1117</b>	<u>0.0169</u>

Table 6: Quantitative comparison on occlusion reordering.

Method	Mean L1 Loss ( $\downarrow$ )	Mean L2 Loss ( $\downarrow$ )	PSNR ( $\uparrow$ )	Speed (s)
Inst-inpaint	12.65%	3.42%	21.2db	0.2547
Dragon Diffusion	13.85%	3.66%	19.82db	0.2849
Ours	<b>3.26%</b>	<b>0.79%</b>	<b>28.32db</b>	<b>0.1739</b>

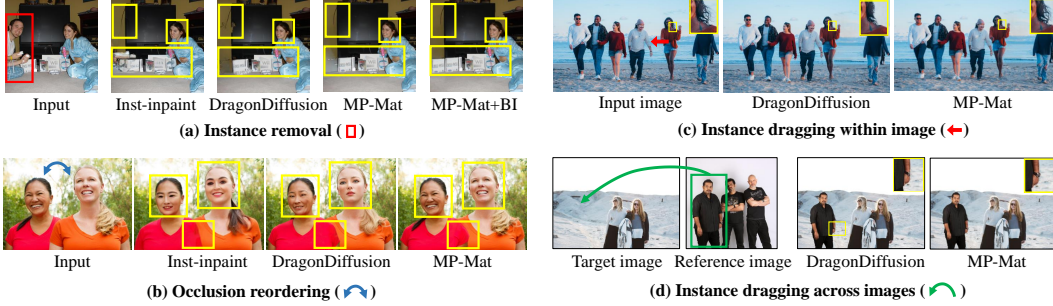


Figure 4: Qualitative comparisons for editing tasks. Yellow boxes highlight the distinguished areas.

**Occlusion reordering.** From Tab. 6, we can also observe the significant superiority of MP-Mat in both effectiveness and efficiency (i.e., more than 10% advantage on Mean L1 Loss, 8.5db on PSNR, and 68% on speed). From Fig. 4 (b), we can observe that when editing target semantics, existing SOTA methods will also unintentionally alter content that should remain unchanged, such as human faces. In contrast, our method better preserves these elements, supported by stronger theoretical guarantees derived from our mathematical transformations, as detailed in Sec. 4. This further demonstrates the superiority of our approach and highlights its potential in editing tasks.

**Instance dragging.** Here we only gave qualitative comparison due to the lack of benchmark datasets. As in Fig. 4 (c) and (d), MP-Mat can preserve finer details such as hair and watch after dragging, which highlights another advantage of our approach. Overall, the results highlight the potential of our matting-focused methods for image editing tasks.

## 6 CONCLUSIONS AND LIMITATIONS

In this work, we propose MP-Mat, a 3D-and-instance-aware matting framework that is built on meticulously designed multiplane representations. Specifically, we design layered representations from two perspectives: the scene geometry-level multiplane representation (SG-MP), which emphasizes scene decomposition based on depth differences, and the instance-level multiplane representation (Inst-MP), which focuses on instance-level modeling. These representations excel in handling occlusion effects and are better aware of both foreground and background content, leading to a significant performance improvement for instance matting. Additionally, our design demonstrates strong potential for instance-level image editing, a relatively underexplored area in existing matting-focused methods. Remarkably, our approach, even under zero-shot inference, outperforms specialized image editing techniques by large margins and with high efficiency. Despite the effectiveness, our work mainly focuses on human instances. This is partially due to the lack of data for multi-instance matting for other categories. We think that mixed training could be a potential workaround to alleviate this, and we left those aspects to our future works.

## REFERENCES

- Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. End-to-end object detection with transformers. In *European conference on computer vision*, pp. 213–229. Springer, 2020.
- Quan Chen, Tiezheng Ge, Yanyu Xu, Zhiqiang Zhang, Xinxin Yang, and Kun Gai. Semantic human matting. In *Proceedings of the 26th ACM international conference on Multimedia*, pp. 618–626, 2018.
- Yigit Ekin, Ahmet Burak Yildirim, Erdem Eren Caglar, Aykut Erdem, Erkut Erdem, and Aysegul Dundar. Clipaway: Harmonizing focused embeddings for removing objects via diffusion models. *arXiv preprint arXiv:2406.09368*, 2024.
- Yuxin Fang, Wen Wang, Binhui Xie, Quan Sun, Ledell Wu, Xinggang Wang, Tiejun Huang, Xinlong Wang, and Yue Cao. Eva: Exploring the limits of masked visual representation learning at scale. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 19358–19369, 2023.
- Marco Forte and François Pitié. *f, b, alpha matting*. *arXiv preprint arXiv:2003.07711*, 2020.
- Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. Mask r-cnn. In *Proceedings of the IEEE international conference on computer vision*, pp. 2961–2969, 2017.
- Guanqing Hu and James Clark. Instance segmentation based semantic matting for compositing applications. In *2019 16th Conference on Computer and Robot Vision (CRV)*, pp. 135–142. IEEE, 2019.
- Chuong Huynh, Seoung Wug Oh, Abhinav Shrivastava, and Joon-Young Lee. Maggie: Masked guided gradual human instance matting. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 3870–3879, 2024.
- Siyi Jiao, Wenzheng Zeng, Changxin Gao, and Nong Sang. Dfimat: Decoupled flexible interactive matting in multi-person scenarios. In *Proceedings of the Asian Conference on Computer Vision*, 2024.
- Lei Ke, Yu-Wing Tai, and Chi-Keung Tang. Deep occlusion-aware instance segmentation with overlapping bilayers. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 4019–4028, 2021.
- Zhanghan Ke, Jiayu Sun, Kaican Li, Qiong Yan, and Rynson WH Lau. Modnet: Real-time trimap-free portrait matting via objective decomposition. volume 36, pp. 1140–1147, 2022.
- Xin Kong, Shikun Liu, Xiaoyang Lyu, Marwan Taher, Xiaojuan Qi, and Andrew J Davison. Escher-net: A generative model for scalable view synthesis. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 9503–9513, 2024.
- Jizhizi Li, Jing Zhang, Stephen J Maybank, and Dacheng Tao. Bridging composite and real: towards end-to-end deep image matting. 130(2):246–266, 2022.
- Yaoyi Li and Hongtao Lu. Natural image matting via guided contextual attention. pp. 11450–11457, 2020.
- Chung-Ching Lin, Jiang Wang, Kun Luo, Kevin Lin, Linjie Li, Lijuan Wang, and Zicheng Liu. Adaptive human matting for dynamic videos. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 10229–10238, 2023.
- Shanchuan Lin, Andrey Ryabtsev, Soumyadip Sengupta, Brian L Curless, Steven M Seitz, and Ira Kemelmacher-Shlizerman. Real-time high-resolution background matting. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 8762–8771, 2021.
- Qinglin Liu, Shengping Zhang, Quanling Meng, Bineng Zhong, Peiqiang Liu, and Hongxun Yao. End-to-end human instance matting. *IEEE Transactions on Circuits and Systems for Video Technology*, 34(4):2633–2647, 2024. doi: 10.1109/TCSVT.2023.3306400.

- Hao Lu, Yutong Dai, Chunhua Shen, and Songcen Xu. Indices matter: Learning to index for deep image matting. pp. 3266–3275, 2019.
- Jonathon Luiten, Georgios Kopanas, Bastian Leibe, and Deva Ramanan. Dynamic 3d gaussians: Tracking by persistent dynamic view synthesis. In *2024 International Conference on 3D Vision (3DV)*, pp. 800–809. IEEE, 2024.
- Sihan Ma, Jizhi Li, Jing Zhang, He Zhang, and Dacheng Tao. Rethinking portrait matting with privacy preserving. pp. 1–26, 2023.
- Lars Mescheder, Michael Oechsle, Michael Niemeyer, Sebastian Nowozin, and Andreas Geiger. Occupancy networks: Learning 3d reconstruction in function space. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 4460–4470, 2019.
- Juewen Peng, Jianming Zhang, Xianrui Luo, Hao Lu, Ke Xian, and Zhiguo Cao. Mpib: An mpi-based bokeh rendering framework for realistic partial occlusion effects. In *The Proceedings of the European Conference on Computer Vision (ECCV)*, 2022.
- Yu Qiao, Yuhao Liu, Xin Yang, Dongsheng Zhou, Mingliang Xu, Qiang Zhang, and Xiaopeng Wei. Attention-guided hierarchical structure aggregation for image matting. pp. 13676–13685, 2020.
- Zhefan Rao. Single portrait image matting and bokeh effect synthesis via multiplane images. Master’s thesis, HKUST, MonthOfPublication 2023.
- Soumyadip Sengupta, Vivek Jayaram, Brian Curless, Steven M Seitz, and Ira Kemelmacher-Shlizerman. Background matting: The world is your green screen. pp. 2291–2300, 2020.
- Shelly Sheynin, Adam Polyak, Uriel Singer, Yuval Kirstain, Amit Zohar, Oron Ashual, Devi Parikh, and Yaniv Taigman. Emu edit: Precise image editing via recognition and generation tasks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 8871–8879, 2024.
- Yujun Shi, Chuhui Xue, Jun Hao Liew, Jiachun Pan, Hanshu Yan, Wenqing Zhang, Vincent YF Tan, and Song Bai. Dragdiffusion: Harnessing diffusion models for interactive point-based image editing. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 8839–8849, 2024.
- Yanan Sun, Chi-Keung Tang, and Yu-Wing Tai. Semantic image matting. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 11120–11129, 2021.
- Yanan Sun, Chi-Keung Tang, and Yu-Wing Tai. Human instance matting via mutual guidance and multi-instance refinement. pp. 2647–2656, 2022.
- Fuwen Tan, Crispin Bernier, Benjamin Cohen, Vicente Ordonez, and Connelly Barnes. Where and who? automatic semantic-aware person composition. In *2018 IEEE Winter Conference on Applications of Computer Vision (WACV)*, pp. 1519–1528. IEEE, 2018.
- Richard Tucker and Noah Snavely. Single-view view synthesis with multiplane images. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 551–560, 2020.
- Xinlong Wang, Rufeng Zhang, Tao Kong, Lei Li, and Chunhua Shen. Solov2: Dynamic and fast instance segmentation. *Advances in Neural information processing systems*, 33:17721–17732, 2020.
- Rundi Wu, Ben Mildenhall, Philipp Henzler, Keunhong Park, Ruiqi Gao, Daniel Watson, Pratul P Srinivasan, Dor Verbin, Jonathan T Barron, Ben Poole, et al. Reconfusion: 3d reconstruction with diffusion priors. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 21551–21561, 2024.
- Ning Xu, Brian Price, Scott Cohen, and Thomas Huang. Deep image matting. pp. 2970–2979, 2017.
- Zhen Yang, Ganggui Ding, Wen Wang, Hao Chen, Bohan Zhuang, and Chunhua Shen. Object-aware inversion and reassembly for image editing. In *The Twelfth International Conference on Learning Representations*, 2024. URL <https://openreview.net/forum?id=dpcVXiMlcV>.

- Ahmet Burak Yildirim, Vedat Baday, Erkut Erdem, Aykut Erdem, and Aysegul Dundar. Inst-inpaint: Instructing to remove objects with diffusion models. *arXiv preprint arXiv:2304.03246*, 2023.
- Wei Yin, Chi Zhang, Hao Chen, Zhipeng Cai, Gang Yu, Kaixuan Wang, Xiaozhi Chen, and Chunhua Shen. Metric3d: Towards zero-shot metric 3d prediction from a single image. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 9043–9053, 2023.
- Jiahui Yu, Zhe Lin, Jimei Yang, Xiaohui Shen, Xin Lu, and Thomas S Huang. Generative image inpainting with contextual attention. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 5505–5514, 2018.
- Qihang Yu, Jianming Zhang, He Zhang, Yilin Wang, Zhe Lin, Ning Xu, Yutong Bai, and Alan Yuille. Mask guided matting via progressive refinement network. pp. 1154–1163, 2021a.
- Zijian Yu, Xuhui Li, Huijuan Huang, Wen Zheng, and Li Chen. Cascade image matting with deformable graph refinement. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 7167–7176, 2021b.
- Yunke Zhang, Lixue Gong, Lubin Fan, Peiran Ren, Qixing Huang, Hujun Bao, and Weiwei Xu. A late fusion cnn for digital matting. pp. 7469–7478, 2019.
- Bingke Zhu, Yingying Chen, Jinqiao Wang, Si Liu, Bo Zhang, and Ming Tang. Fast deep matting for portrait animation on mobile phone. In *Proceedings of the 25th ACM international conference on Multimedia*, pp. 297–305, 2017.