

Situational Perception Guided Image Matting

Bo Xu¹, Jiake Xie², Han Huang¹, Ziwen Li¹, Cheng Lu³, Yong Tang² and Yandong Guo^{1,*}
¹OPPO Research Institute, ²PicUp.AI, ³Xmotors

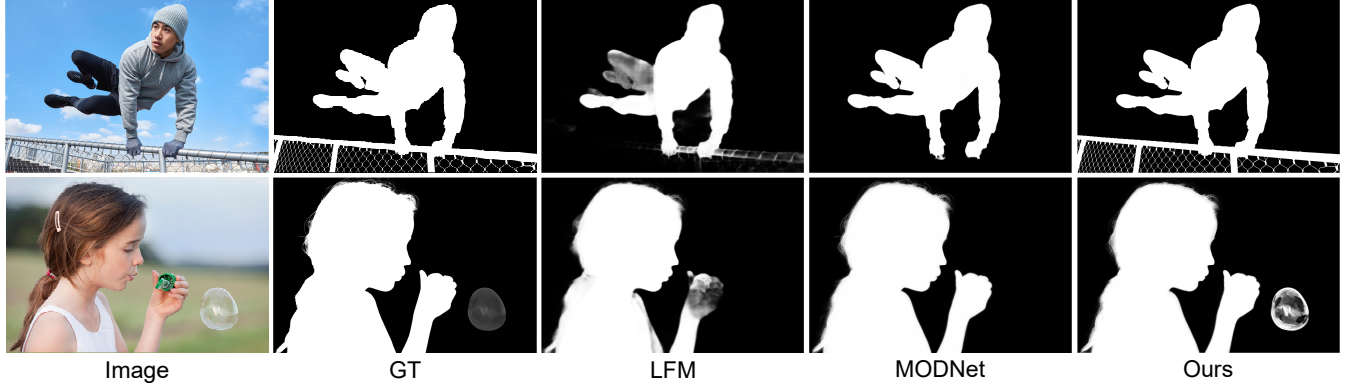


Figure 1: Visual comparisons of our SPG-IM with other trimap-free matting methods (LFM [57] and MODNet [28]) on the Real-World images.

ABSTRACT

Most automatic matting methods try to separate the salient foreground from the background. However, the insufficient quantity and subjective bias of the current existing matting datasets make it difficult to fully explore the semantic association between object-to-object and object-to-environment in a given image. In this paper, we propose a Situational Perception Guided Image Matting (SPG-IM) method that mitigates subjective bias of matting annotations and captures sufficient situational perception information for better global saliency distilled from the visual-to-textual task. SPG-IM can better associate inter-objects and object-to-environment saliency, and compensate the subjective nature of image matting and its expensive annotation. We also introduce a textual Semantic Transformation (TST) module that can effectively transform and integrate the semantic feature stream to guide the visual representations. In addition, an Adaptive Focal Transformation (AFT) Refinement Network is proposed to adaptively switch multi-scale receptive fields and focal points to enhance both global and local details. Extensive experiments demonstrate the effectiveness of situational perception guidance from the visual-to-textual tasks on image matting, and our model outperforms the state-of-the-art methods. We also analyze the significance of different components in our model. The code will be released soon.

KEYWORDS

image matting, trimap, visual-to-textual, cross modality, transformer

1 INTRODUCTION

Image matting is a fundamental computer vision task with great application value, which aims to separate the foreground from a single image or video stream and then composites it with a new

background. It has been practically applied in the background replacement scenarios of the multimedia such as entertainment image/video creation and special-effect film-making, without green screen backgrounds. Due to the rapid development of the deep neural networks in computer vision, automatic matting becomes increasingly matured [10, 30, 33, 36, 38, 43, 45, 52, 57].

However, there still remain two big challenges in the matting task. First, unlike some well-defined tasks such as object detection or segmentation, image matting is an ill-posed problem that requires a user input trimap or some interactive strokes/points to clearly separate foreground from background, which leads to the inconsistency matting annotation from the annotators' subjective understanding of the foreground. If we present the same image to a group of subjects, it's very natural that each individual will interpret the foreground of this image a bit differently based on his or her educational background, gender, ethnicity, or religious beliefs. Such variation impacts image matting annotations, especially in the human-object interactive and multi-object scenes. For example, as shown in Row 1 of Figure 1, annotator A is more willing to interpret the fence as the foreground object together with the man in this image, while annotator B justifiably prefers to only highlight human. One of the most straightforward ways to learn the underlining distribution of such variation is by introducing large-scale data-driven training. Unfortunately, it is unrealistically expensive to acquire such a representative dataset to cover the real data distribution due to the fine pixel-level granularity annotation of image matting. Recent existing image matting training datasets like AIM [52] and Distinction-646 [43] are examples of such data insufficiency.

Second, for most of the existing automatic matting models without extra inputs (*e.g.* trimap, interactive strokes/points, or known background), the learning of image saliency steams from object

detection or segmentation methodology but lacks the global situational awareness on multiple salient objects and their surrounding environment. For example in Row 2 of Figure 1, although no touch interaction with the girl, that bubble, which is critical for semantic completeness, is indisputably salient. Unfortunately, previous methods fail to extract such complete and meaningful saliency.

To address the above challenges, in this paper we propose a Situational Perception Guided Image Matting (SPG-IM) network, that aims to mitigate the subjective bias of matting annotations and capture sufficient situational perception information for global saliency learning. We seek semantic distilled information from the visual-to-textual task to guide the visual features of image matting, due to its large quantity but low-cost training samples. We believe that visual representations from the visual-to-textual task, *e.g.* image captioning, focus on more semantically comprehensive signals between a) object to object and b) object to the ambient environment to generate descriptions that can cover both the global info and local details. In addition, compared with the expensive pixel annotation of image matting, textual labels can be massively collected at a very low cost.

The SPG-IM network consists of two stages: Situational Perception Distillation (SPD) and Situational Perception Guided Matting (SPGM), both in an encoder-to-decoder fashion. In the first stage, we first pretrain the visual front-end and transformer decoder jointly to generate captions, and aim to learn visual representations including situational perception from visual-to-textual feature transformation. Then the visual front-end is spliced to a new back-end for saliency foreground mask prediction. In the second stage, the SPGM module takes both generated mask and the raw RGB image as inputs and outputs the estimated alpha matte. To leverage situational perception guidance, we propose a textual Semantic Transformation module that transforms and integrates the visual feature stream of the SPD module to guide the visual representations of the SPGM module at multi-scale levels. In addition, we propose an Adaptive Focal Transformation (AFT) Refinement Network that can adaptively select the size of the receptive field to process global context and local details separately, aiming to reach a good balance between complementary global information and local attributes when processing the fused situational perception guided visual features. To justify our solutions, we compare our algorithm, objectively and subjectively, with other methods. Also, we demonstrate by ablation study that the introduction of visual-to-textual transform as semantic guidance can mitigate subjective annotation bias and improve matting performance by leveraging inexpensive image captioning labeling.

Overall, the contributions of this paper are as follows:

- To the best of our knowledge, we are the first to underscore the subjective nature of foreground saliency in image/video matting and accordingly introduce situational perception guidance from the visual-to-textual transformation with low-cost labeling to semantically guide the visual features to compensate demographic bias and improve matting performance.
- We build a large-scale matting dataset consisting of 1000 images and corresponding alpha mattes for multi-foreground-object scenes. To the best of our knowledge, this is the first

large-scale and high-quality dataset for multi-foreground-object scenes.

- We propose a textual Semantic Transformation (TST) module to effectively transform and integrate more situational perception information that guides the matting network.
- We propose an Adaptive Focal Transformation (AFT) Refinement Network to adaptively select the size of receptive fields and focal regions to simultaneously improve global and local performance.
- Extensive experiments demonstrate the effectiveness of our situational perception-guided image matting method, outperforming the state-of-the-art (SOTA) approaches on both synthetic and real-world images.

2 RELATED WORKS

Currently, the matting is generally formulated as an image composite problem, which solves the 7 unknown variables per pixel from only 3 known values:

$$I_i = \alpha_i F_i + (1 - \alpha_i) B_i \quad (1)$$

where 3 dimensional RGB color I_i of pixel i , while foreground RGB color F_i , background RGB color B_i , and matte estimation α_i are unknown. In this section, we discuss the SOTA works trying to solve this under-determined equation.

2.1 Classic methods

Classic foreground matting methods can be generally categorized into two approaches: sampling-based and propagation-based. Sampling-based methods [2, 12, 15, 25, 27, 56] sample the known foreground and background color pixels, and then extend these samples to achieve matting in other parts. Various sampling-based algorithms are proposed, *e.g.* Bayesian matting [56], optimized color sampling [25], global sampling method [27], and comprehensive sampling [15]. Propagation-based methods [3, 10, 20, 29, 30, 47] reformulate the composite Eq. 1 to propagate the alpha values from the known foreground and background into the unknown region, achieving more reliable matting results. [26] provides a very comprehensive review on various matting algorithms.

2.2 Deep learning-based methods

Classic matting methods are carefully designed to solve the composite equation and its variant versions. However, these methods heavily rely on chromatic cues, which leads to bad quality when the color of the foreground and background show small or no noticeable difference.

Trimap-based methods. Automatic and intelligent matting algorithms are emerging, due to the rapid development of the deep neural network in computer vision. Initially, some attempts were made to combine deep learning networks with classic matting techniques, *e.g.* closed-form matting [30] and KNN matting [10]. Cho *et al.* [13] employ a deep neural network to improve the results of the closed-form matting and KNN matting. These attempts are not end-to-end, so not surprisingly the matting performance is limited by the convolution back-ends. Subsequently, full DL image matting algorithms appear [9, 19, 52]. Xu *et al.* [52] propose a two-stage deep neural network (Deep Image Matting) based on SegNet [4] for

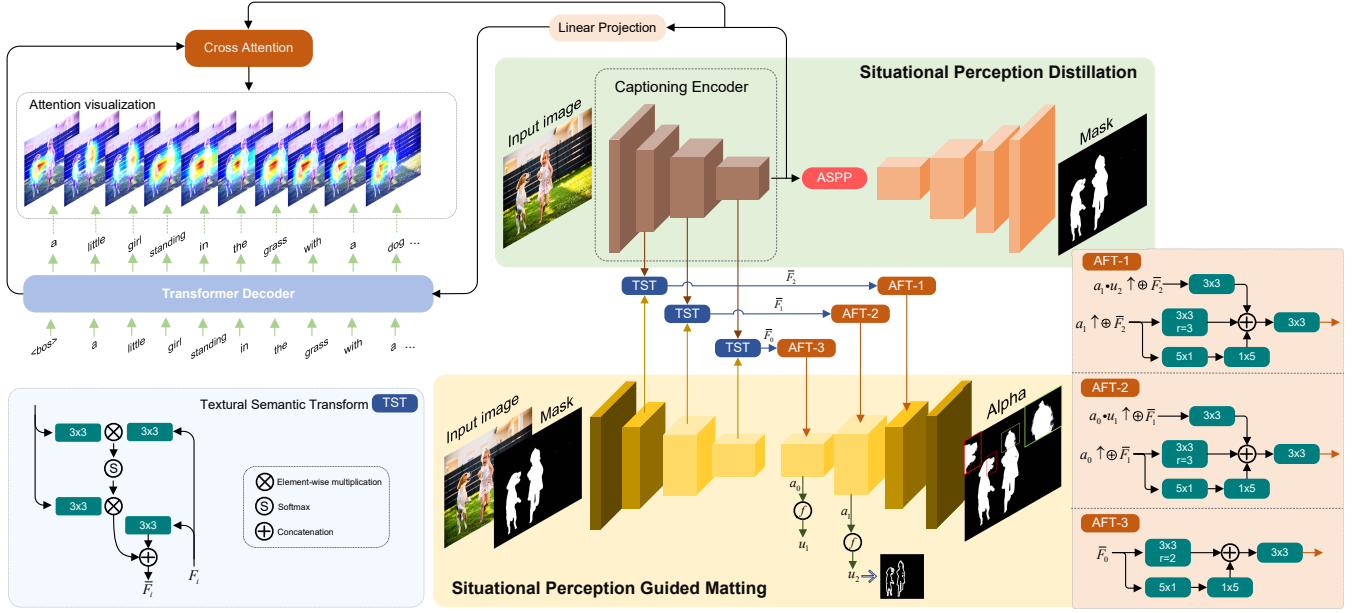


Figure 2: The architecture of the Situational Perception Guidance Image Matting (SPG-IM) network. The SPG-IM network consists of two branches: Situational Perception Distillation (SPD) and Situational Perception Guided Matting (SPGM). The SPD branch takes an RGB image as input and outputs the semantic distilled saliency mask. The SPGM then revisits the raw image input and combines it with the estimated saliency mask for alpha matte prediction, under the guidance of situational perception from the SPD. Both SPD and SPGM employ ResNet-50 as the encoder for visual representation extraction. We utilize ASPP [7] to extract and fuse multi-scale contextual information for semantic mask estimation.

alpha matte estimation and contribute a large-scale image matting dataset (Adobe dataset) with ground truth foreground (alpha) matte, which can be composited over a variety of backgrounds to produce training data. We also use this data for the first-step pre-training of our network. Lutz *et al.* [39] introduce a generative adversarial network (GAN) for natural image matting and improve the results of Deep Image Matting [52]. Cai *et al.* [5] investigate the bottleneck of the previous methods that directly estimate the alpha matte from a coarse trimap, and propose to divide the matting problem into trimap adaptation and alpha estimation tasks. Hou *et al.* [23] employs two encoder networks to extract essential information for matting, however it is not robust to faulty trimaps. Forte *et al.* [34] propose a low-cost upgrade to alpha matting networks to also predict the foreground and background colours. They study variations of the training regions and explore a wide range of existing and novel loss functions for optimal prediction. Liu *et al.* [36] propose a 3-branch encoder to accomplish comprehensive mining of the input RGB image and its corresponding trimap, and then develop a Tripartite Formation Integration (TI2) Module to transform and integrate the interconnections between the different branches.

Additional natural background. Qian *et al.* [42] compute a probability map to classify each pixel into the foreground or background by simple background subtraction. This algorithm is sensitive to the threshold and fails when the colors of the foreground and background are similar. Sengupta *et al.* [45] introduce a self-supervised adversarial approach - Background Matting (BGM),

achieving state-of-the-art results. However, as a prerequisite, the photographer needs to take a shoot of natural background first, which is not friendly to the intensive multi-scene shooting application. Liu *et al.* [33] propose the Background Matting V2 that employs two neural networks: a base network computes a low-resolution result which is refined by a second network operating at high-resolution on selective patches.

Trimap-free methods. Currently, a majority of deep image matting algorithms [5, 23, 39, 52] try to estimate a boundary that divides the foreground and background, with the aid of a user-generated trimap. Several trimap-free matting methods [9, 57] predict the trimap first, followed by alpha matting. Qiao *et al.* [43] employ spatial and channel-wise attention to integrating appearance cues and pyramidal features, they also introduce a hybrid loss function fusing Structural Similarity (SSIM), Mean Square Error (MSE), and Adversarial loss to guide the network to further improve the overall foreground structure in trimap-free matting. Lin *et al.* [34] propose a robust real-time matting method (RVM) training strategy that optimizes the network on both matting and segmentation tasks. Ke *et al.* [28] present a lightweight matting objective decomposition network (MODNet) by optimizing a series of sub-objectives simultaneously via explicit constraints. They also introduce an e-ASPP module to fuse multi-scale features, plus a self-supervised sub-objectives consistency (SOC) strategy to address the domain shift problem, which is common in trimap-free methods.

Besides, most current trimap-free methods focus only on human/portrait matting but ignore the objects that are interacting with or attached to people. In addition, they learn the saliency of the images by data-driven training, which lacks the situational perception between salient objects and the surrounding environment, leading to biased or incomplete foreground prediction, especially in multi-object scenes. This is the main reason why we propose the method of Situational Perception Guided Image Matting. In this paper, we quantitatively evaluate the performance of our model for alpha matting in human-object interactive and multi-object scenes.

2.3 Image Captioning.

The problem of generating natural language descriptions from visual data has long been studied in computer vision. Early methods use pre-defined templates, such as object detectors and attribute predictors to generate captions [46, 53]. With the rise of Deep Learning-based networks, RNNs [44, 51] are adopted as language models to decode corresponding visual features.

Due to the wide success of transformers in natural language processing and multi-media, image caption methods use transformers to either generate captions directly or fuse visual and language features. Herdade *et al.* [22] propose a object relation transformer and build image captions based on inter-object relations. Liu *et al.* [31] introduce an enTangled attention-based transformer that simultaneously exploits visual and semantic information. Huang *et al.* [24] propose an attention model that first generates an information vector and an attention gate, and then adds secondary attention using element-wise multiplication to aggregate the attended features.

Recent works have demonstrated that image captions can guide feature learning on various visual tasks. Karan Desai *et al.* [14] propose a pre-train approach using semantic dense captions to learn visual representations. We believe distilled semantic caption information can guide the feature learning of salient object detection tasks.

3 METHODOLOGY

The network architecture of the Situational Perception Guided Image matting (SPG-IM) is designed to automatically extract the accurate saliency foreground from an RGB image instead of using interactive strokes/points or extra inputs, *e.g.* trimap and known background. The SPG-IM network consists of two branches: Situational Perception Distillation (SPD) and Situational Perception Guided Matting (SPGM). We pretrain the front-end of SPD under an image caption generation framework before it is transferred to the downstream semantic distilled saliency mask estimation task. The SPGM then revisits the raw image input and combines it with the estimated saliency mask for alpha matte prediction, under the guidance of situational perception from the SPD. The overall architecture of the SPG-IM network is shown in Figure 2.

3.1 Situational Perception Distillation

The front-end of our Situational Perception Distillation (SPD) branch is pretrained joint with a transformer-based textual decoder [14] to generate the textual description of an image.

At the visual-to-textual pretraining stage, the Transformer-based textual decoder decodes the visual features output by the visual

front-end (captioning encoder) and generates a corresponding image caption $C = (c_0, c_1, \dots, c_T, c_{T+1})$. The start of the image caption is $c_0 = [SOS]$, while $c_{T+1} = [EOS]$ indicates the end of a caption sequence.

Visual front-end (captioning encoder). The visual encoder uses a convolutional network to compute the downsampled visual features. For the input image I , we use the ResNet-50 [21] as the visual encoder to extract grid feature $C \in \mathbb{R}^{2048 \times (7 \times 7)}$, followed by a linear projection layer before sending it to the textual decoder.

Textual decoder: The textual decoder receives a set of grid visual features and outputs the corresponding image caption. We predict the captions in both forward and backward manner and utilize Transformer [50] as the backbone of the textual decoder, which adopts a self-attention and cross attention mechanism to fuse visual features using textual queries.

The inputs of the textual decoder module are a set of image features from the visual encoder and a list of caption tokens. Grid visual features fed into the textual decoder are tokenized to a sequence of the patch features $G \in \mathbb{R}^{D_I \times N_I}$, where each $N_I = 7 \times 7$ patch has a feature vector with D_I -dimension.

The first token $c_0 = [SOS]$ indicates the start of the sentence. The transformer backbone iteratively predicts each word in the caption sentence. The prediction ends when transformer output $C_{T+1} = [EOS]$ label. We visualize the word-level attention maps in the cross attention module, where the highlighted regions illustrate that the visual representations from the visual-to-textual task can focus more on the global situational perception. More implementation details are described in the supplementary material.

Situational perception distillation. We transfer the visual features learned by the visual-to-textual transformation to the downstream dense prediction task. We adopt an encoder-to-decoder framework in the SPD network. An Atrous Spatial Pyramid Pooling (ASPP) [6] module is set between the visual front-end and decoder to enhance the fusion capabilities of multi-scale features for semantic situational perception. The semantic distilled saliency mask output by the SPD network is supervised by the Gaussian transformed thumbnail of the ground truth alpha matte at a L_2 loss:

$$L_{SPD} = \|M_s - S(a^*)\|_2 \quad (2)$$

where a^* is the ground truth of alpha matte, $S(a^*)$ is a Gaussian blur operation after the downsampling of a^* . We utilize L_2 loss to smooth the boundary and details of the estimated semantic distilled saliency mask M_s . Then M_s is fed into the SPGM branch along with the raw image for alpha prediction.

3.2 Situational Perception Guided Matting

The situational perception guided matting (SPGM) branch receives the RGB image I and the semantic distilled saliency mask M_s as inputs to transfer the high-level semantic generalizations into fine-grained foreground alpha mattes. The coarse foreground masks are experimentally proven to be effective as semantic priors for image matting in previous works [45, 55]. As for situational perception guidance, we propose a textual Semantic Transformation (TST) module that effectively transforms and integrates the visual feature stream of the SPD module, and guides the visual representations of the SPGM module at multi-scale levels.

Textual Semantic Transformation. The textual Semantic Transformation (TST) module is performed in a non-local fashion, its network architecture is shown in Figure 2. We first encode the visual presentation F'_i of SPD and the visual representation F_i of SPGM separately into the key (k'_i, k_i) maps and value (v'_i, v_i) maps at each feature scale. The fusion attention map f_i is computed by comparing the pixel-by-pixel similarity between k'_i and k_i :

$$f_i = \text{softmax}(k' k^\top) \quad (3)$$

Then we utilize f_i as indexes to retrieve the effective situational perception information and concatenate it with v_i to update the visual representation of SPGM:

$$\tilde{F}_i = v_i \oplus (v'_i{}^\top \odot f_i) \quad (4)$$

where \odot denotes an element-wise multiplication and \oplus denotes the concatenation. We conduct textual semantic transformation on the features from ResNet50's 2, 3, 4 - th res blocks. Consequently, the semantic distilled situational perception information can work as a guiding role at multiple feature levels.

Adaptive Focal Transformation Refinement Network. As described in previous dense prediction work [41], larger receptive fields establish dense connections between feature maps and per-pixel classifiers which improve the accuracy of internal regions, while smaller receptive fields benefit the localization focus on local fine-grained details near the object boundaries. It is impossible to find correct and accurate local details if the network is already confused on the global or regional level, which can only generate fine-grained but semantically incorrect results at best. However, it's equally impossible to get good matting if the network is heavily biased to global or regional features since that can only cause image blur.

That holds also true for the human visual system. When humans observe objects with their eyes, they usually use an adaptive focal strategy to first capture the main body and then narrow their vision on the particular details, such as hair or other small texture. After accumulating a few scrutinized details, the human in return re-gauge or re-evaluate their general perception on those objects. Inspired by this, we propose an Adaptive Focal Transformation (AFT) Refinement Network that can adaptively switch dimensions of receptive fields and focal regions, aiming to complement global information and local attributes when processing the fused situational perception guided features. Specifically, we first generate the focal region mask u_i at i th level from the output a_{i-1} of the previous feature level by the following formula:

$$u_i(x, y) = \begin{cases} 1 & \text{if } 0 < a_{i-1}(x, y) < 1 \\ 0 & \text{otherwise} \end{cases} \quad (5)$$

where we set the low confidence regions $0 < a_{i-1}(x, y) < 1$ which consists mainly of boundary details as shown in Figure 2 and need to be adaptively focused and progressively refined. Then we upsample a_{i-1} and $a_{i-1} \cdot u_i$, and concatenate them on the fused situational perception guided feature \tilde{F}_i respectively at i th level. As shown in the pink box of Figure 2, we apply multi-size kernels, larger sizes for $(a_{i-1} \uparrow \oplus \tilde{F}_i)$ to predict main body regions, and smaller sizes for $(a_{i-1} \cdot u_i \uparrow \oplus \tilde{F}_i)$ to refine boundary regions. Then additional convolutions are utilized to fuse of main body and refined boundary

details. To reduce the computation cost, we apply Atrous convolution [8] kernels along with spatially separable convolution kernels for large receptive field at multi-scale features.

Loss function. For the supervision of the SPGM branch, we only utilize alpha loss to supervise the outputs of different levels in order to verify the validity of this pattern and to prevent bias caused by other losses:

$$\mathcal{L} = \sum_i \lambda_i \|a_i - a^*\|_1 \quad (6)$$

where λ_i is the loss weight assigning to the output alpha a_i of the i th level, a^* is the ground truth alpha.

4 EXPERIMENTS

We first describe the datasets used for training and testing. Subsequently, we compare our results with existing state-of-the-art (SOTA) foreground matting algorithms. Finally, we conduct ablation experiments to show the effectiveness of each branch and module. More implementation details are provided in the supplementary material.

4.1 Datasets

Composition-1K [52]. The training set consists of 431 foreground objects and each of them is composited over 100 random COCO [35] images to produce 43.1k composited training images. For the test set, we combine each foreground of Composition-1K with 20 random VOC [16] images to produce 1k composited testing images. And then we follow Eq. 2 in Section 3.1 to generate $S(a^*)$ for supervising the training of the Situational Perception Distillation (SPD) branch.

Distinction-646 [43]. It includes 431 and 50 foreground objects in training and test sets, respectively. We enforce the same rule and composited ratio with the Composition-1K.

Human-2K [36]. It provides 2100 foreground images (2000 for training and 100 for testing). The same rules and ratios as Composition-1K are used in the Human-2K to composite new images.

Multi-Object-1K. Although there are several typical datasets we can use for the matting task, most of them include only single-object foregrounds. To extend the image matting to the application of multi-foreground-object scenes, we propose our Multi-Object 1K which consists of both single-foreground-object and multi-foreground-object images. Consequently, this dataset can better evaluate the semantic situational perception ability of our method, especially in multi-object scenes.

Multi-Object-1K provides 1000+200 real-world images and high accuracy alpha mattes, where 70% of the datasets are multi-object scenes. We believe Multi-Object 1K can serve as a new challenging benchmark in the image matting area. We also apply the same rules and ratios as Composition-1K on our Multi-Object-1K for data composition.

4.2 Comparative study on composition datasets

We conduct comparative study on three composition benchmarks: Composition-1K [52], Distinction-646 [43], and Human-2K [36] datasets. We report mean square error (MSE), sum of the absolute difference (SAD), spatial-gradient (Grad), and connectivity (Conn)

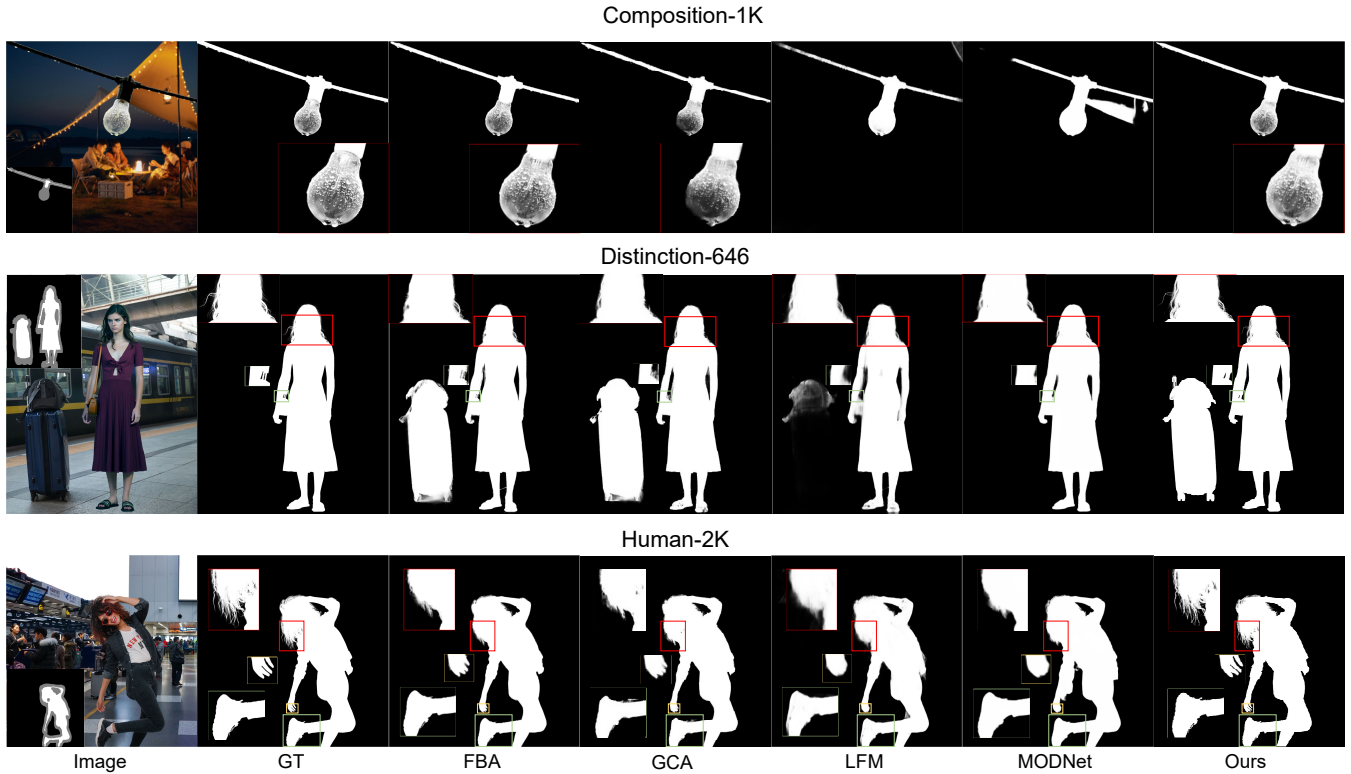


Figure 3: Visual comparison on public composition datasets. Trimap-based methods: FBA [17] and GCA [32]. Trimap-free methods: LFM [57], MODNet [17], and ours.

Methods	SAD↓	MSE↓	Grad↓	Conn↓
DIM [52]	50.4	0.014	31.0	50.8
IndexNet [38]	45.8	0.013	25.9	43.7
SampleNet [49]	40.4	0.010	-	-
CAM [23]	35.8	0.008	17.3	33.2
LFM [57]	49.0	0.020	34.3	50.6
HAttMatting [43]	44.0	0.007	29.3	46.4
MODNet [17]	43.7	0.012	29.1	42.0
GCA [32]	35.3	0.009	16.9	32.5
HDMatt [54]	33.5	0.007	14.5	29.9
SIM [48]	28.0	0.006	10.8	24.8
FBA [17]	25.8	0.005	10.6	20.8
TIMI-Net [36]	29.1	0.006	11.5	25.4
MG-Mask _{SPD} [55]	30.2	0.007	12.9	26.6
MODNet [17]	43.7	0.012	29.1	42.0
Ours	23.2	0.004	8.9	18.6

Table 1: The quantitative results on Composition-1K. *Mask_{SPD}* denotes the guidance input for MG matting [55].

Methods	SAD↓	MSE↓	Grad↓	Conn↓
DIM [52]	47.6	0.009	43.3	55.9
HAttMatting [36]	49.0	0.009	41.6	49.9
TIMI-Net [36]	22.3	0.011	14.4	20.5
MODNet	46.2	0.009	39.0	43.5
MG-Mask _{SPD} [55]	23.6	0.007	16.1	21.0
Ours	20.9	0.006	11.2	19.8

Table 2: The quantitative results on Distinction-646 [43].

Methods	SAD↓	MSE↓	Grad↓	Conn↓
DIM [52]	7.5	0.008	6.4	6.7
TIMI-Net [36]	4.2	0.003	2.1	3.0
MODNet	7.8	0.008	7.2	7.4
MG-Mask _{SPD} [55]	4.4	0.004	2.5	3.2
Ours	4.0	0.002	2.0	2.8

Table 3: The quantitative results on Human-2K [36].

between predicted and ground truth alpha mattes. Lower values of these metrics indicate better estimated alpha matte.

We compare our method with state-of-the-art (SOTA) trimap-based methods: DIM [52], IndexNet [38], CAM [23], GCA [32],

Methods	Single-object				Multi-object			
	SAD↓	MSE↓	Grad↓	Conn↓	SAD↓	MSE↓	Grad↓	Conn↓
TIMI-Net [36]	26.4	0.012	15.7	30.0	28.3	0.013	26.1	42.4
MODNet	51.5	0.014	36.1	56.9	69.2	0.024	60.8	81.3
MG- $Mask_{SPD}$	32.9	0.017	26.4	47.0	33.5	0.019	32.7	55.6
Basic	40.2	0.023	43.2	64.5	47.2	0.024	45.0	68.1
Basic+PRN _{MG} [55]	29.9	0.015	20.4	41.8	32.4	0.017	31.2	53.5
Basic+AFT	26.9	0.012	16.5	30.4	27.3	0.015	19.8	32.9
SPG-IM w/o AFT	33.7	0.018	30.6	50.9	26.7	0.014	16.0	32.3
SPG-IM	22.6	0.008	12.5	24.9	22.7	0.008	13.1	27.7

Table 4: The quantitative results on our Multi-object-1K.

FBA [17], and SIM [48]; mask-based method: MG Matting [55]; trimap-free methods: LFM [57], HAttMatting [43], and MODNet [17]. To fairly compare them for fully automatic matting, we utilize our Situational Perception Distillation (SPD) branch to produce the semantic distilled foreground mask $Mask_{SPD}$ for MG matting. For trimap-based methods, we generate trimaps from the ground truth alpha mattes by thresholding and random dilation as discussed in [52]. For methods without publicly available codes, we follow their papers to reproduce the results with due diligence.

Table 1 to 4 show the quantitative results of our SPG-IM with other SOTA models on four datasets. Our SPG-IM outperforms all competing trimap-free methods (LFM [57], HAttMatting [43], and MODNet [17]) by a large margin. Meanwhile, our model also shows remarkable superiority over the state-of-the-art (SOTA) trimap-based and mask-guided methods in terms of all four metrics across the public datasets (*i.e.* Composition-1K, Distinction-646, and Human-2K), and our Multi-Object-1K benchmark.

To give an intuitive understanding of the significance of the situational perception guidance, we visualize sampled results from our SPG-IM and other SOTA models in Figure 3 and 5. It can be obviously observed that our method preserves fine details (*e.g.* hair tip sites, transparent textures, and boundaries) without the guidance of trimap. Moreover, compared to other competing trimap-free models, our SPG-IM can retain better global semantic completeness. For example in Row 2 of Figure 3, we composite a human foreground from Distinction-646 into a train station scene with luggage. Due to not directly touching the person, the saliency of the luggage is challenging to capture. Fortunately, with the help of the situational perception guidance, our SPG-IM explores the semantic association between object-to-object (*i.e.* the person and the luggage) and object-to-environment (*i.e.* the luggage and the station), and then effectively identifies the saliency of the luggage and considers it as part of the foreground, where other trimap-free methods fail.

4.3 Ablation Study

We validate the effectiveness of our key components on the Multi-Object-1K dataset, where more than 70% samples are multi-object foregrounds. As summarized in Table 4, we conduct quantitative comparisons in both single-object foreground and multi-object foreground scenes respectively. The model settings of the ablation

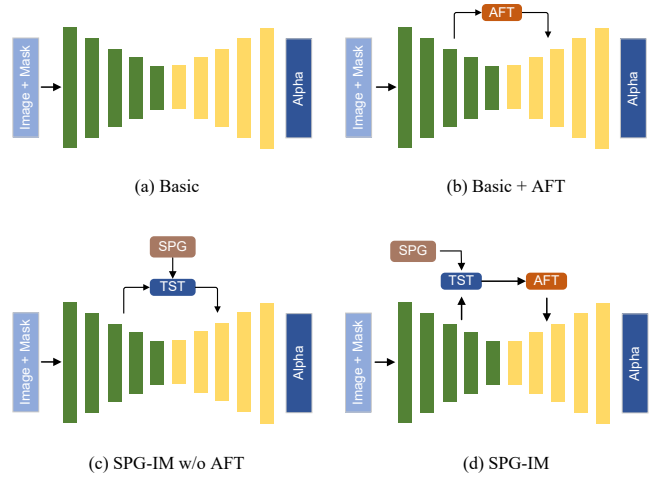


Figure 4: The settings of ablation study for our method.

study are illustrated in Figure 4, where Basic denotes the encoder-to-decoder structure of the situational perception guidance matting (SPGM) branch without situational perception guidance or Adaptive Focal Transformation (AFT) refinement module. Basic + AFT refers to performing adaptive focal transformation directly on multi-scale features of the Basic network. SPG-IM w/o AFT indicates that the fused situational perception guided feature \bar{F}_i is skip-connected respectively to the decoder of SPGM at i th feature level without adaptive focal transformation.

Adaptive Focal Transformation. We report the quantitative comparison results of our model with and without the Adaptive Focal Transformation (AFT) refinement module at dual settings, *i.e.* Basic vs. Basic+AFT and SPG-IM w/o AFT vs. SPG-IM. As shown in Table 4, both the baseline and its situational perception guided variants show performance gains by applying our AFT module, proving the necessity of the adaptive transformation of receptive fields for focal regions in image matting. We also apply the Progressive Refinement Module (PRM_{MG}) from the MG matting [55] on our Basic network for comparison. The quantitative results in

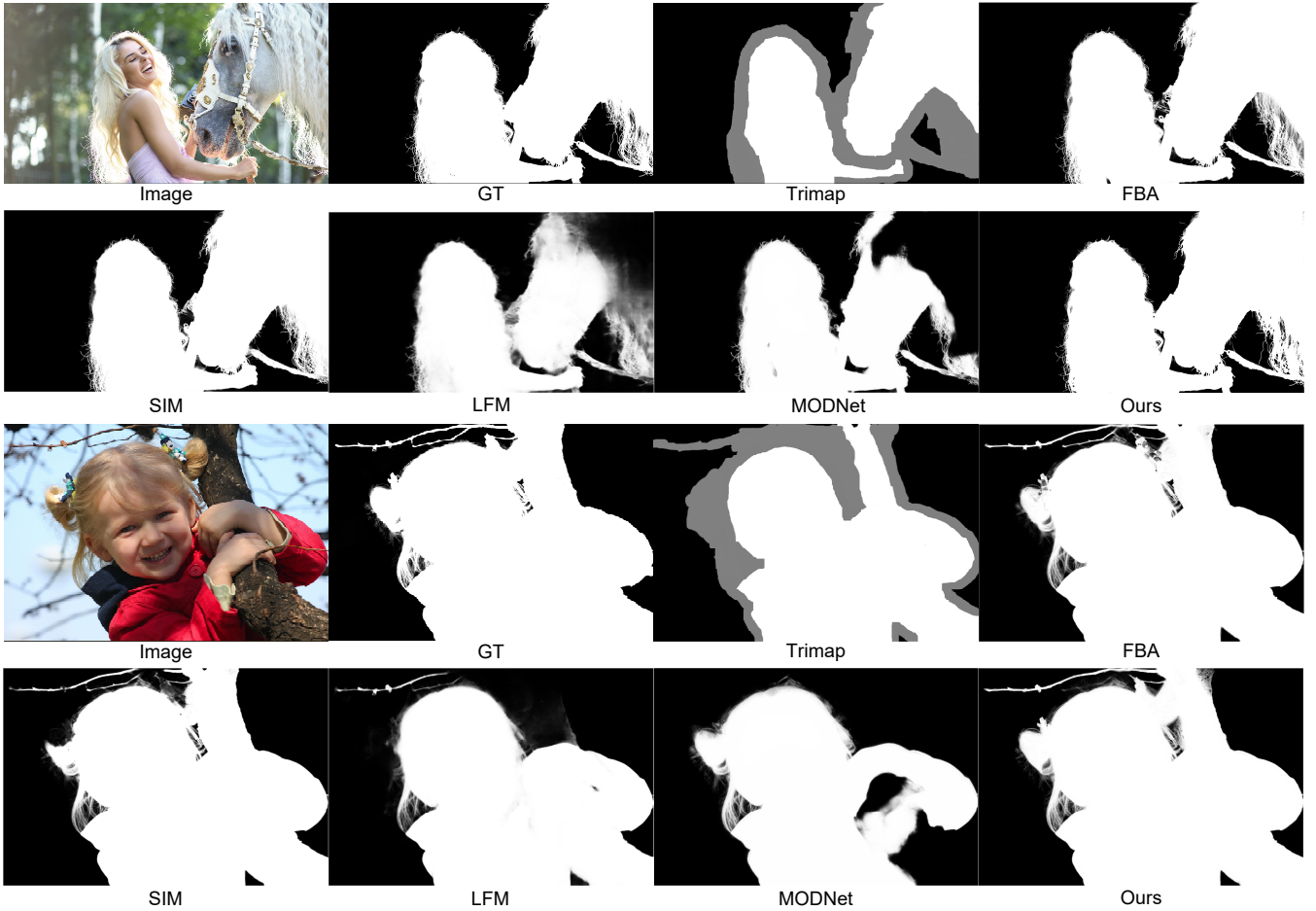


Figure 5: Visual comparisons on our Multi-Object-1K. Trimap-based methods: FBA [17] and SIM [48]. Trimap-free methods: LFM [57], MODNet [17], and ours

Table 4 illustrate the superiority of our AFT which can better complement the high-level semantic information completeness and the low-level subtle detail refinement.

Textual Semantic Transformation. We also evaluate our model under two ablation settings, *i.e.* Basic **vs.** SPG-IM w/o AFT and Basic + AFT **vs.** SPG-IM. The quantitative comparisons are shown in Table 4. The proposed TST module improves the performance of alpha estimation in both single-object and multi-object foreground scenes, particularly the improvement in the latter is expectedly significant. This is because the TST module can effectively transform and integrate more situational perception information to guide the matting network for better saliency association between inter-objects and object-to-environment. Additionally, we demonstrate the performance gain after combining textual semantic transformation and adaptive focal transformation. Some representative visualized comparisons of the original real-world samples from our Multi-Object-1K are provided in Figure 5, which also illustrate that our SPG-IM can enhance both global semantic awareness and local details.

5 CONCLUSION

In this paper, we present a situational perception guided matting technique (SPG-IM) that can capture situational perception information distilled from the visual-to-textual task for better global saliency, aiming to compensate the subjective nature of the image matting and mitigate the subjective bias of matting annotations without the expensive pixel annotation.

For the implementation of situational perception guidance, we introduce a Textual Semantic Transformation (TST) module to effectively transform and integrate more situational perception information that guides the matting network. Further, an Adaptive Focal Transformation (AFT) Refinement Network is proposed to adaptively select the size of receptive fields and focal regions to simultaneously improve global and local matting performance. Extensive experiments demonstrate that our model outperforms current state-of-the-art algorithms in both single-object and multi-object foreground scenes without extra inputs, *e.g.* trimap, known background, and interactive strokes/points. For future works, we will continue to focus on the study of multi-object foreground matting and conduct research on the real-time versions.

Parameter	Value
Optimizer	Adam
SPD	
Initial learning rate	1.0×10^{-2}
Input image size	512×512
Batch size of SPD	16
SPGM	
Initial learning rate	5.0×10^{-3}
Input image size	512×512
Batch Size of SPGM	4
Loss weight λ_1	1
Loss weight λ_2	2
Loss weight λ_3	3

Table 5: Implementation details and hyper-parameter setting.

A IMPLEMENTATION DETAILS

A.1 Implementation of the SPG-IM

The implementation of the Situational Perception Guided Image Matting (SPG-IM) framework is based on the public PyTorch [40] toolbox and trained on a single Tesla V100 GPU with 32GB memory. The training details and all hyper parameters are outlined in Table 5. The learning rate of SPD is decreased by a factor of 10 at the epoch of {20, 40}, {30, 60}, and {60, 80}, and {40, 60, 80} for Composition-1K, Distinction-646, and Human-2K, and our Multi-Object-1K, respectively. Meanwhile, the learning rate of SPGM decays at a rate of 0.1 in the epoch of {20, 30, 40}, {40, 60, 80}, and {80, 100, 120}, and {60, 80, 100} for Composition-1K, Distinction-646, and Human-2K, and our Multi-Object-1K, respectively.

A.2 Pretraining of the visual front-end in the SPD branch.

At the pretraining stage, we train the visual-to-textual network on the COCO Captions dataset [11] and use the SGD optimizer with momentum 0.9 and weight decay 10^{-4} . Follow [14], we utilize warmup [18] at the beginning iterations followed by cosine decay [37] to zero. The max learning rates for the visual front-end and textual decoder are set to 0.2 and 10^{-3} respectively.

A.3 Model size comparison

We compare our model size with other trimap-based (*e.g.* GCA [32] and FBA [17]), trimap-free (*e.g.* LFM [57]), and mask-guided (*e.g.* MG [55]) methods. As shown in Figure 6, the total model size of our method is smaller than GCA_{auto} , FBA_{auto} , and MG_{auto} . The SPD branch contains 65.7% of the parameters in our SPG-IM. Similarly, we observe that the automatic generation of higher-level semantic priors (such as trimap and mask) tends to be more computationally intensive for both trimap-based and mask-guided methods. For future work, we will optimize the front-end semantic distillation module to achieve the lightweight of the entire model.

Method	Parameters (M)	Size (MB)
LFM [57]	225.9	863.5
GCA [32]	25.3	96.6
GCA_{auto}	80.0	305.3
FBA [17]	34.7	138.8
FBA_{auto}	89.4	347.5
MG_{auto} [55]	84.4	322.7
SPD	40.2	153.7
SPG-IM	61.2	234.1

Table 6: Model size comparison. GCA_{auto} , FBA_{auto} , and MG_{auto} use DeepLabV3+ with Xception backbone for the segmentation prior (automatic trimap or mask generation). The implementation of LFM [57] network is based on the TensorFlow [1] library.

B MORE VISUALIZATION RESULTS

We display more representative visualizations on our Multi-Object-1K benchmark and real world images. Performance comparisons in Figure 6, 7, and 8 demonstrate the effectiveness and generalization of our situational perception guided image matting (SPG-IM), especially in the multi-object foreground scenes. Meanwhile, our SPG-IM can enhance both global semantic awareness and local details. The proposed Multi-Object-1K can further extend the image matting from the single-object foreground scenes to the complex multi-media situations.



Figure 6: More visual comparisons on our Multi-Object-1K. Trimap-based methods: GCA [32] and SIM [48]. Trimap-free methods: LFM [57], MODNet [17], and ours.

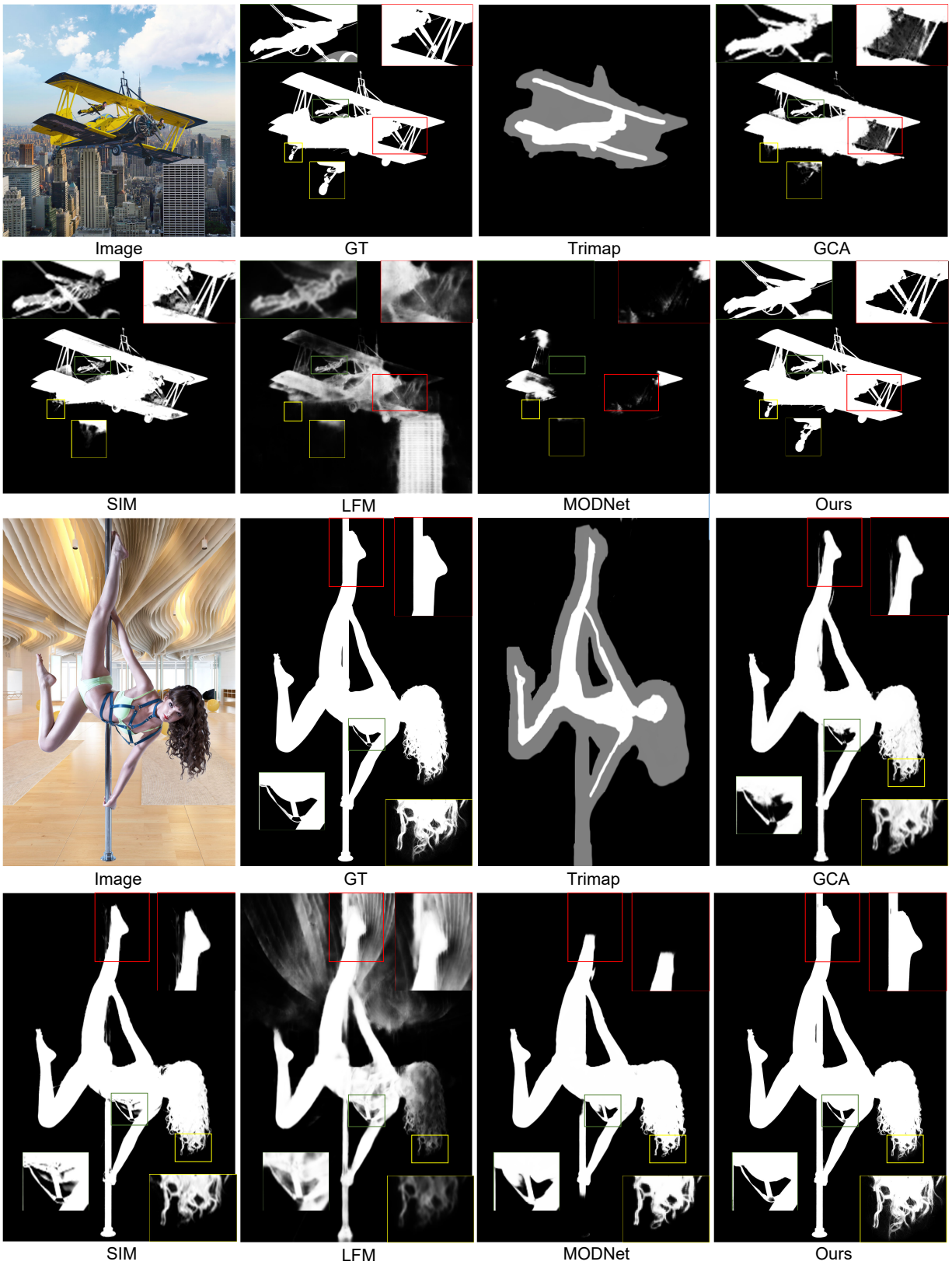


Figure 7: More visual comparisons on our Multi-Object-1K. Trimap-based methods: GCA [32] and SIM [48]. Trimap-free methods: LFM [57], MODNet [17], and ours.



Figure 8: More visual comparisons between MODNet [17] and our SPG-IM on real world images.

REFERENCES

- [1] Martin Abadi, Ashish Agarwal, Paul Barham, Eugene Brevdo, Zhifeng Chen, Craig Citro, Greg S Corrado, Andy Davis, Jeffrey Dean, Matthieu Devin, et al. 2016. Tensorflow: Large-scale machine learning on heterogeneous distributed systems. *arXiv preprint arXiv:1603.04467* (2016).
- [2] Yagiz Aksoy, Tunc Ozan Aydin, and Marc Pollefeys. 2017. Designing effective inter-pixel information flow for natural image matting. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. 29–37.
- [3] Levin Anat, Rav-Acha Alex, and Lischinski Dani. 2008. Spectral matting. *IEEE transactions on pattern analysis and machine intelligence* 30, 10 (2008), 1699–1712.
- [4] Vijay Badrinarayanan, Alex Kendall, and Roberto Cipolla. 2017. Segnet: A deep convolutional encoder-decoder architecture for image segmentation. *IEEE transactions on pattern analysis and machine intelligence* 39, 12 (2017), 2481–2495.
- [5] Shaofan Cai, Xiaoshuai Zhang, Haoqiang Fan, Haibin Huang, Jiangyu Liu, Jiaming Liu, Jiaying Liu, Jue Wang, and Jian Sun. 2019. Disentangled image matting. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*. 8819–8828.
- [6] Liang-Chieh Chen, George Papandreou, and Iasonas Kokkinos et al. 2017. Deeplab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected crfs. *IEEE transactions on pattern analysis and machine intelligence* (TPAMI) 40, 4 (2017), 834–848.
- [7] Liang-Chieh Chen, George Papandreou, Florian Schroff, and Hartwig Adam. 2017. Rethinking atrous convolution for semantic image segmentation. *arXiv preprint arXiv:1706.05587* (2017).
- [8] Liang-Chieh Chen, Yukun Zhu, George Papandreou, Florian Schroff, and Hartwig Adam. 2018. Encoder-decoder with atrous separable convolution for semantic image segmentation. In *Proceedings of the European conference on computer vision (ECCV)*. 801–818.
- [9] Quan Chen, Tiezheng Ge, Yanyu Xu, Zhiqiang Zhang, Xinxin Yang, and Kun Gai. 2018. Semantic human matting. In *Proceedings of the 26th ACM international conference on Multimedia*. 618–626.
- [10] Qifeng Chen, Dingzeyu Li, and Chi-Keung Tang. 2013. KNN matting. *Proceedings of the IEEE transactions on pattern analysis and machine intelligence* 35, 9 (2013), 2175–2188.
- [11] Xinlei Chen, Hao Fang, Tsung-Yi Lin, Ramakrishna Vedantam, Saurabh Gupta, Piotr Dollár, and C Lawrence Zitnick. 2015. Microsoft coco captions: Data collection and evaluation server. *arXiv preprint arXiv:1504.00325* (2015).
- [12] Xiaowu Chen, Dongqing Zou, Steven Zhiying Zhou, Qimping Zhao, and Ping Tan. 2013. Image matting with local and nonlocal smooth priors. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. 1902–1907.
- [13] Donghyeon Cho, Yu-Wing Tai, and Inso Kweon. 2016. Natural image matting using deep convolutional neural networks. In *Proceedings of the European conference on computer vision (ECCV)*. Springer, 626–643.
- [14] Karan Desai and Justin Johnson. 2021. Virtex: Learning visual representations from textual annotations. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. 11162–11173.
- [15] Shahrian Ehsan, Rajan Deepu, Price Brian, and Cohen Scott. 2013. Improving image matting using comprehensive sampling sets. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. 636–643.
- [16] Mark Everingham, Luc Van Gool, Christopher KI Williams, John Winn, and Andrew Zisserman. 2010. The pascal visual object classes (voc) challenge. *International journal of computer vision* 88, 2 (2010), 303–338.
- [17] Marco Forte and François Pitié. 2020. F, B, Alpha Matting. *arXiv preprint arXiv:2003.07711* (2020).
- [18] Priya Goyal, Dhruv Mahajan, Abhinav Gupta, and Ishan Misra. 2019. Scaling and benchmarking self-supervised visual representation learning. In *Proceedings of the IEEE/CVF International Conference on computer vision*. 6391–6400.
- [19] Chen Guanying, Han Kai, and Wong Kwan-Yee K. 2018. TOM-Net: Learning transparent object matting from a single image. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. 9233–9241.
- [20] Kaiming He, Jian Sun, and Xiaoou Tang. 2010. Fast matting using large kernel matting laplacian matrices. In *Proceedings of the 2010 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, 2165–2172.
- [21] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition (CVPR)*. 770–778.
- [22] Simao Herdade, Armin Kappeler, Kofi Boakye, and Joao Soares. 2019. Image captioning: Transforming objects into words. *arXiv preprint arXiv:1906.05963* (2019).
- [23] Qiqi Hou and Feng Liu. 2019. Context-aware image matting for simultaneous foreground and alpha estimation. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*. 4130–4139.
- [24] Lun Huang, Wenmin Wang, Jie Chen, and Xiao-Yong Wei. 2019. Attention on attention for image captioning. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*. 4634–4643.
- [25] Wang Jue and Cohen Michael F. 2007. Optimized color sampling for robust matting. In *2007 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, 1–8.
- [26] Wang Jue and Cohen Michael F. 2008. *Image and video matting: a survey*. Now Publishers Inc.
- [27] He Kaiming, Rhemann Christoph, Rother Carsten, Tang Xiaoou, and Sun Jian. 2011. A global sampling method for alpha matting. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, 2049–2056.
- [28] Zhanghan Ke, Kaican Li, Yurou Zhou, Qiuhua Wu, Xiangyu Mao, Qiong Yan, and Rynson WH Lau. 2020. Is a green screen really necessary for real-time portrait matting? *arXiv preprint arXiv:2011.11961* (2020).
- [29] Philip Lee and Ying Wu. 2011. Nonlocal matting. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, 2193–2200.
- [30] Anat Levin, Dani Lischinski, and Yair Weiss. 2007. A closed-form solution to natural image matting. *IEEE transactions on pattern analysis and machine intelligence* 30, 2 (2007), 228–242.
- [31] Guang Li, Linchao Zhu, Ping Liu, and Yi Yang. 2019. Entangled transformer for image captioning. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*. 8928–8937.
- [32] Yaoyi Li and Hongtao Lu. 2020. Natural image matting via guided contextual attention. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 34. 11450–11457.
- [33] Shanchuan Lin, Andrey Ryabtsev, Soumyadip Sengupta, Brian L Curless, Steven M Seitz, and Ira Kemelmacher-Shlizerman. 2021. Real-time high-resolution background matting. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 8762–8771.
- [34] Shanchuan Lin, Linjie Yang, Imran Saleemi, and Soumyadip Sengupta. 2022. Robust High-Resolution Video Matting with Temporal Guidance. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*. 238–247.
- [35] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. 2014. Microsoft coco: Common objects in context. In *Proceedings of the European conference on computer vision (ECCV)*. Springer, 740–755.
- [36] Yuhao Liu, Jiake Xie, Xiao Shi, Yu Qiao, Yujie Huang, Yong Tang, and Xin Yang. 2021. Tripartite information mining and integration for image matting. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*. 7555–7564.
- [37] Ilya Loshchilov and Frank Hutter. 2016. Sgdr: Stochastic gradient descent with warm restarts. *arXiv preprint arXiv:1608.03983* (2016).
- [38] Hao Lu, Yutong Dai, Chunhua Shen, and Songcen Xu. 2019. Indices matter: Learning to index for deep image matting. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*. 3266–3275.
- [39] Sebastian Lutz, Konstantinos Amniamitis, and Aljosa Smolic. 2018. AlphaGAN: Generative adversarial networks for natural image matting. In *British Machine Vision Conference (BMVC)*. BMVA Press, 259.
- [40] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimeshein, Luca Antiga, et al. 2019. Pytorch: An imperative style, high-performance deep learning library. *Advances in neural information processing systems* 32 (2019).
- [41] Chao Peng, Xiangyu Zhang, Gang Yu, Guiming Luo, and Jian Sun. 2017. Large kernel matters—improve semantic segmentation by global convolutional network. In *Proceedings of the IEEE conference on computer vision and pattern recognition (CVPR)*. 4353–4361.
- [42] Richard J Qian and M Ibrahim Sezan. 1999. Video background replacement without a blue screen. In *Proceedings 1999 International Conference on Image Processing (Cat. 99CH36348)*, Vol. 4. IEEE, 143–146.
- [43] Yu Qiao, Yuhao Liu, Xin Yang, Dongsheng Zhou, Mingliang Xu, Qiang Zhang, and Xiaopeng Wei. 2020. Attention-Guided Hierarchical Structure Aggregation for Image Matting. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
- [44] Steven J Rennie, Etienne Marcheret, Youssef Mroueh, Jerret Ross, and Vaibhava Goel. 2017. Self-critical sequence training for image captioning. In *Proceedings of the IEEE conference on computer vision and pattern recognition (CVPR)*. 7008–7024.
- [45] Soumyadip Sengupta, Vivek Jayaram, Brian Curless, Steven M Seitz, and Ira Kemelmacher-Shlizerman. 2020. Background Matting: The World is Your Green Screen. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. 2291–2300.
- [46] Richard Socher and Li Fei-Fei. 2010. Connecting modalities: Semi-supervised segmentation and annotation of images using unaligned text corpora. In *2010 IEEE Computer Society Conference on Computer Vision and Pattern Recognition*. IEEE, 966–973.
- [47] Jian Sun, Jiaya Jia, Chi-Keung Tang, and Heung-Yeung Shum. 2004. Poisson matting. In *Proceedings of the ACM Special Interest Group on Computer Graphics (SIGGRAPH)*. 315–321.
- [48] Yanan Sun, Chi-Keung Tang, and Yu-Wing Tai. 2021. Semantic image matting. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 11120–11129.
- [49] Jingwei Tang, Yagiz Aksoy, Cengiz Oztireli, Markus Gross, and Tunc Ozan Aydin. 2019. Learning-based sampling for natural image matting. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 3055–3063.

- [50] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In *Advances in neural information processing systems*. 5998–6008.
- [51] Oriol Vinyals, Alexander Toshev, Samy Bengio, and Dumitru Erhan. 2016. Show and tell: Lessons learned from the 2015 mscoco image captioning challenge. *IEEE transactions on pattern analysis and machine intelligence (TPAMI)* 39, 4 (2016), 652–663.
- [52] Ning Xu, Brian Price, Scott Cohen, and Thomas Huang. 2017. Deep image matting. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. 2970–2979.
- [53] Benjamin Z Yao, Xiong Yang, Liang Lin, Mun Wai Lee, and Song-Chun Zhu. 2010. I2t: Image parsing to text description. *Proc. IEEE* 98, 8 (2010), 1485–1508.
- [54] Haichao Yu, Ning Xu, Zilong Huang, Yuqian Zhou, and Humphrey Shi. 2021. High-Resolution Deep Image Matting. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 35. 3217–3224.
- [55] Qihang Yu, Jianming Zhang, He Zhang, Yilin Wang, Zhe Lin, Ning Xu, Yutong Bai, and Alan Yuille. 2021. Mask guided matting via progressive refinement network. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 1154–1163.
- [56] Chuang Yung-Yu, Curless Brian, Salesin David H, and Szeliski Richard. 2001. A bayesian approach to digital matting. In *Proceedings of the 2001 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR)*, Vol. 2. IEEE, 264–271.
- [57] Yunke Zhang, Lixue Gong, Lubin Fan, Peiran Ren, Qixing Huang, Hujun Bao, and Weiwei Xu. 2019. A late fusion CNN for digital matting. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. 7469–7478.