# Multi-scale Information Assembly for Image Matting

Yu Qiao[†], Yuhao Liu[†], Qiang Zhu, Xin Yang[‡], Yuxin Wang, Qiang Zhang, and Xiaopeng Wei

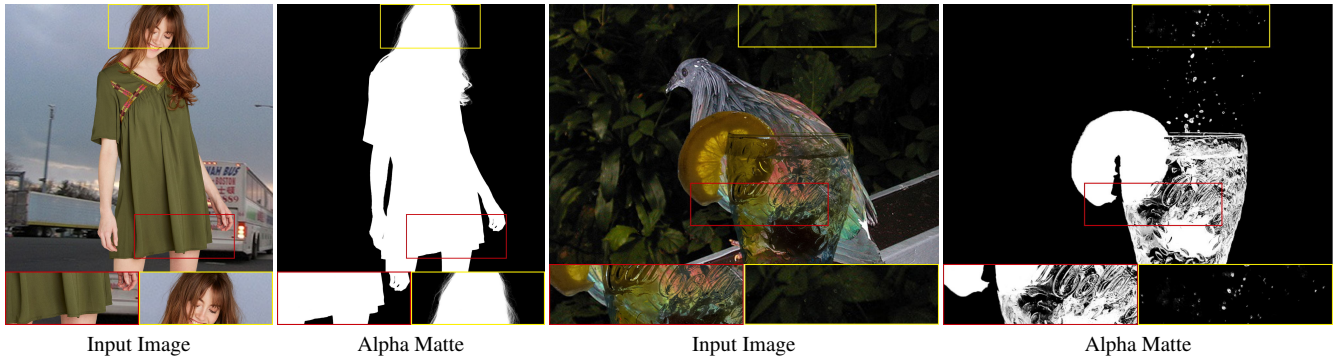College of Computer Science, Dalian University of Technology



**Figure 1:** *The alpha mattes generated by our proposed MSIA-matte. Some details are zoomed in and placed under the images.*

## Abstract

*Image matting is a long-standing problem in computer graphics and vision, mostly identified as the accurate estimation of the foreground in input images. We argue that the foreground objects can be represented by different-level information, including the central bodies, large-grained boundaries, refined details, etc. Based on this observation, in this paper, we propose a multi-scale information assembly framework (MSIA-matte) to pull out high-quality alpha mattes from single RGB images. Technically speaking, given an input image, we extract advanced semantics as our subject content and retain initial CNN features to encode different-level foreground expression, then combine them by our well-designed information assembly strategy. Extensive experiments can prove the effectiveness of the proposed MSIA-matte, and we can achieve state-of-the-art performance compared to most existing matting networks.*

**CCS Concepts**
*• Computing methodologies → Image segmentation; Image representations;*

## 1. Introduction

Natural images are composed of different kinds of objects, and we usually consider the regions of interest as foreground. Accurate prediction and separation of the foreground in the input image (noted as image matting) is a significant problem in industry and academia. Image matting has a wide range of applications, film production, live streaming, online image editing, etc. We can define image matting from the perspective of image synthesis as follows:

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

where $I$ represents the input image, $i$ refers to the pixel location of $I$, $F$ and $B$ refer to the foreground and background layers separately. $\alpha$ denotes the alpha matte, where $\alpha_i$ varies in the range of [0,1] to suggest the opacity of the foreground. We can summarize from Equation (1) that a relatively complete foreground object can be represented by different-level information. They could be subject content ($\alpha_i = 1$), clear borders or coarse textures ($\alpha_i$ are very large, even close to 1), semi-transparent areas or fine-grained details ($\alpha_i$ are relatively small).

Given only an input image, solving for other variables from Equation (1) can be intractable. Therefore, most existing matting methods import trimaps to confine the foreground and background. The trimap consists of three categories to indicate the foreground ($\alpha = 1$), background ($\alpha = 0$), and transition region ($\alpha \in (0, 1)$).

---

† Joint first authors. ‡Corresponding author and he led this project.

With explicit foreground and background as the reference, traditional matting methods resort to sampling or affinity to estimate the alpha values in the transition region. However, they mostly consider color distribution to predict alpha mattes and ignore the semantic information of the foreground object, which can result in poor performance when the foreground and background share similar colors.

Deep learning has contributed a lot to the development of computer vision, as well as image matting. [CTK16] et al. combined RGB images with the alpha mattes from [LLW07] and [CTK16] to predict refined results via deep convolutional neural networks. The concatenation of RGB images and trimaps as input was first introduced by [XPCH17], and the expression ability of neural networks for the foreground structure has been fully demonstrated. Many subsequent matting networks [HL19, CZF*19] follow this principle to optimize their results, using trimaps to confine the input images and designing comprehensive network architectures to present foreground information. However, generating a well-defined trimap from a natural image is fussy and difficult for novice users, limiting the application of the above methods in practice to some extent.

In recent years, some matting methods can also achieve alpha mattes without trimaps, and their motivations are mainly relying on semantic segmentation or saliency detection. [ZGF*19] et al. proposed a fusion network to blend the foreground and background weight maps, both of which are derived from a segmentation network. The attention mechanism proposed in [QLY*20] is more of a saliency-dominated model: the advanced semantics from the backbone network are applied to the improvement of the low-level features. Almost all trimap-free methods involve low-level features in their later fusion or refinement stage, which is the main reason that they can restore fine-grained boundary details in alpha mattes. However, many methods exploit low-level CNN features for image matting without considering their different-opacity foreground information.

In this paper, we argue that the foreground expression from low-level features still possesses different-level information. For a natural image, there is usually an explicit subject content of foreground object (it could be a human, a dog, or a glass), and there are also diverse boundaries or details (such as hands, hairs, branches, leaves), denoted as superficial traces in this paper. These superficial traces obviously require different alpha values indicated by Equation (1). Therefore, we propose a multi-scale information assembly network (*MSIA-matte*) to extract different-level foreground information for predicting alpha mattes. Specifically, we harness a backbone network to capture high-level semantics and design a superficial traces branch to extract different-level initial features. In the decoder part, we assemble multi-scale features by our proposed information assembly module to estimate alpha mattes. The feature maps from the superficial traces branch can complement the advanced semantics in multi-layered textures and details, and their integration can guarantee relatively complete and concrete foreground expression. Extensive experiments can prove the effectiveness of our proposed network, and we can achieve state-of-the-art performance compared to the existing matting methods with/without trimaps.

We list our contributions as follow:

- We propose a multi-scale information assembly network (MSIA-matte) to predict alpha mattes, which can combine different-level foreground expression to refine the textures and details in results.
- We design a novel extraction strategy, which can adapt to various opacity in the foreground and provide high-level semantics and superficial traces for information assembly.
- We conduct extensive experiments to validate our model, and the alpha mattes produced by our network can achieve state-of-the-art performance.

## 2. Related Work

In this paper, we propose a multi-scale integration model to generate alpha mattes without trimaps. In this section, we briefly review some matting approaches from three directions: traditional matting methods, trimap-based deep-learning networks, and trimap-free matting architectures.

**Traditional approaches.** Most traditional approaches resort to scribbles or trimaps to improve their sampling or propagation process. The trimap is generally a gray-scale image with three colors: white, black, and gray. The white indicates the foreground in the corresponding input RGB image, and the black and gray can suggest the background and transition region, respectively. Most of the foreground and background are recommended by trimap, and the core of the matting is to solve the transition region according to explicit parts. The theory of scribbles is similar, using simple scribbles to provide some explicit foreground and background information, which is user-friendly, while the produced results are inferior to trimap-based methods due to the limited available reference regions. According to the different ways of exploring explicit information, traditional methods can be divided into two directions: sampling-based methods and affinity-based methods.

Sampling methods [YCSS01, FLZ, KEE15, HRR*11, SRPC13a, WC07] collect sampling pixels from explicit parts and predict the pixels in transition regions by sampled ones. Affinity methods [AAP17, CLT13, GSAW05, LW11, LLW07, LRL08, SJTS04] propagate explicit pixels to the transition region to estimate the alpha values of the whole image, which is more robust than sampling methods. However, the above methods mostly exploit color or texture distribution to achieve the regression of transition regions, ignoring the high-level semantics of the foreground, which can represent the structure of the input images [XPCH17].

**Deep-learning matting with trimaps.** As pixel-wise dense annotations, trimaps can provide sufficient guidance for alpha perception. Therefore, many deep-learning-based matting methods utilize trimaps as additional inputs to improve the final results. The deep architecture in Cho et al. [CTK16] combined RGB images with the results from [LLW07] and [CLT13] as input and produce better alpha mattes. Xu et al. [XPCH17] proposed deep image matting (DIM), designing an encoder-decoder network to extract advanced semantics from RGB images and trimaps, and the importance and quality of the trimaps for DIM have been proved in [ZGF*19]. Lutz et al. [LAS18] adopted generative adversarial networks [GPAM*14] to enhance the matte prediction, which demonstrated the ability of discriminator for pixel-wise visual quality, and the later [QLY*20] also imported a discriminator to strengthen the visual effects. Sampling operation was integrated
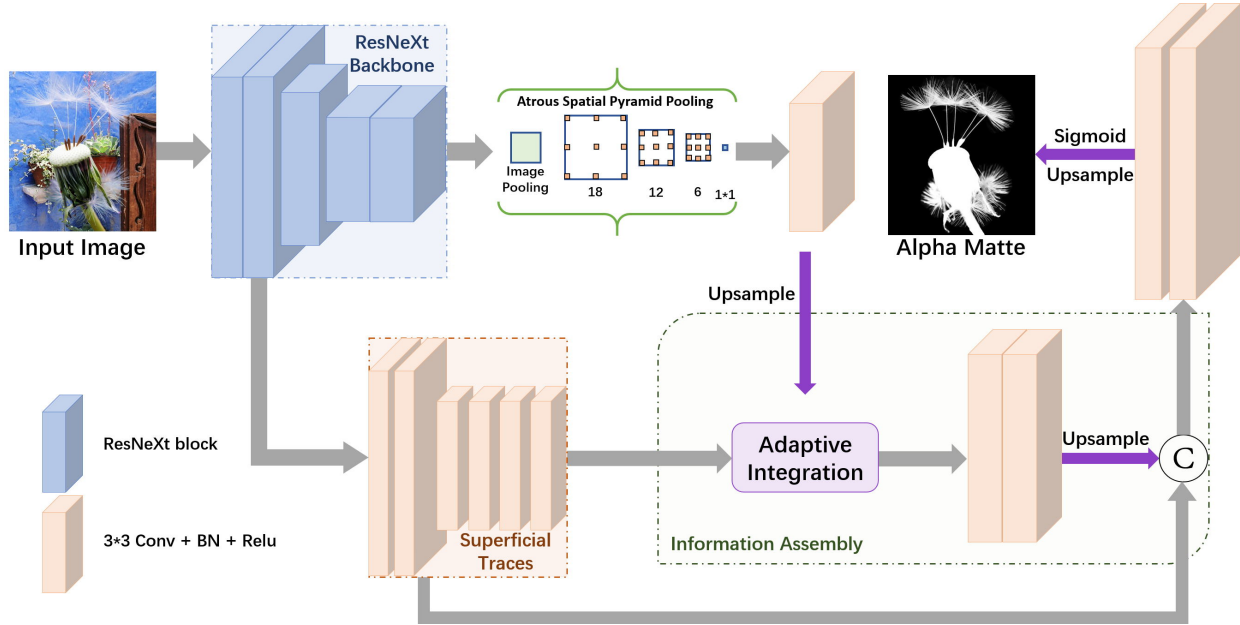
**Figure 2:** *The pipeline of the proposed MSIA-matte. We employ ResNeXt as our backbone to extract advanced semantics, which can represent the subject content of the foreground. Synchronously, we preserve initial features to infer different-level superficial traces, which contain foreground details and textures with distinguished opacity. In the information assembly decoder stage, we employ an adaptive strategy to assemble multi-scale foreground expression to predict an alpha matte.*

with neural networks in Tang et al. [TAO*19] for foreground opacity estimation. Cai et al. [CZF*19] and Hou et al. [HL19] both resorted to another branch to improve the alpha prediction, which exhibited the joint promotion of two branches. Hao et al. [LDSX19] embedded the index function into their encoder-decoder network to preserve more effective parameters. Li et al. [LL20] utilized guided contextual attention to improve the transmission of the opacity information. However, the above methods require carefully defined trimaps as assistant input, which partly restrict their practical values in real life.

**Trimap-free matting networks.** This solution is mostly the production of semantic segmentation or saliency detection, using intermediate results as the reference to confine the final alpha mattes. Some matting frameworks [CGX*18,CKTK17] employed segmentation models to produce coarse trimaps, and the results based on such trimaps can usually lose information in the foreground objects or edges. The reinforcement model proposed in [YXC*18] can promote alpha mattes progressively, while their feedback stage will trouble users a lot. The computation and time-consuming in [AOP*18] were very expensive, although they can produce comparable alpha mattes. Zhang et al. [ZGF*19] used a segmentation variant to represent the foreground and background of the input image, while subtle segmentation errors can discount their results. The attention mechanism in [QLY*20] can effectively aggregate the pyramidal features and appearance cues, while they use high-level semantics as primary guidance for attention, more like a saliency-dominated matting. In this paper, we investigate different-level superficial traces to complement the semantic extraction of the back-

bone, and the novel information assembly strategy can integrate multi-scale features effectively.

## 3. Methodology

In this section, we formally interpret our multi-scale information assembly network (MISA-matte). For a familiar natural image, there may be different kinds of objects, and we usually define the one of interest as the foreground. Existing matting methods mostly exploit complex models to extract the foreground semantics and adopt later refinement or fusion [XPCH17, CZF*19, ZGF*19] after the decoder to restore the details and edges in alpha mattes. However, there is various opacity in the foreground (like the hair, body, and hands in a portrait, the ears, legs, and hairs of a cat), and such relatively direct refinement or fusion is inapplicable to integrate different-level foreground expression, which can lead to some information absence when the opacity of the foreground or background shares a continuous variation. The motivation of our model is to capture multi-scale foreground information to represent manifold opacity, then integrate them to predict a relatively complete foreground.

### 3.1. Overview

The pipeline of the proposed MSIA-matte is unfolded in Figure 2. We adopt ResNeXt [XGD*17] as our backbone to extract advanced semantics, which can represent the subject content of the foreground objects. Then we use the Atrous Spatial Pyramidal Pooling module [CPK*18] to further capture different-level semantic

information. Similar to [ZGF*19, HL19, QLY*20], we first intend to combine low-level CNN features in the decoder parts. However, based on different-level opacity observations, we establish a superficial traces branch to encode low-level CNN features. The original textures, borders, or details can be converted to different-level foreground expression. Then we can integrate multi-scale expression by well-designed information assembly module, including upsampling, convolution, and concatenation operation in the decoder stage. The whole network architecture is optimized by a blended loss function.

### 3.2. Superficial traces branch

High-level semantic feature extraction has been widely developed by convolutional neural networks [SZ15,HZRS16], and such multi-layer semantics can be combined with upsampled features in the decoder stage [LSD15]. Inspired by this multi-scale semantics fusion strategy, we propose our superficial traces branch to further analyze low-level CNN features from ResNeXt block1. Low-level CNN features contain many complex image textures, which can be integrated into alpha perception to complement boundary details or semi-transparent regions [ZGF*19,QLY*20]. The motivation of our superficial branch is to extract different-level foreground traces. Thus we can restore various opacity in final alpha mattes.

We perform a feature transform module on the low-level features to extract superficial traces. Specifically, we first employ two convolutional layers with $3*3$ kernel size and 64 channels to process initial low-level features (the output features here are denoted as $\mathcal{F}_{ini}$). This operation can preliminarily filter scrambled image textures and preserve sufficient fine-grained foreground information. After this, we use a downsampled convolution on $\mathcal{F}_{ini}$ with stride=2, then three cascaded convolutional layers with $3*3$ kernel size and 256 channels are used to extract relatively deep superficial traces, recorded as $\mathcal{F}_{sed}$. We utilize such continuous convolutions to condense some coarse-grained borders or textures of the foreground, which belong to superficial traces and are defined in low-level CNN features but contain subtle semantics or clear transition (such as the hands, leaves, claws). Such filtration and concentration operations on superficial traces can mostly retain initial image textures and capture different-level foreground information, and their integration with high-level semantics can jointly regress alpha values.

### 3.3. Information assembly decoder

We assemble different-level superficial traces with high-level semantics during our upsampling integration stage to restore the foreground profiles in alpha mattes. According to our previous analysis, the secondary superficial traces $\mathcal{F}_{sed}$ from 4 cascaded convolutional layers contain foreground expression above the initial image textures. Nevertheless, such expression may be inapplicable for all kinds of matting images: the input image must include boundary details, but some examples may contain no middle-level superficial traces (hands, semi-transparent regions). Therefore, we first involve an information assembly to dynamically integrate secondary superficial traces $\mathcal{F}_{sed}$ with semantic features $\mathcal{F}_{aspp}$ from Atrous Spatial Pyramidal Pooling (ASPP) module. The assembled

foreground information $\mathcal{F}_{IA}$ is encoded by two $3*3$ convolutional layers and upsampled by a factor of 2, then integrated with preliminarily processed low-level features $\mathcal{F}_{ini}$ by a concatenation operation, the output of the concatenation is denoted as $\mathcal{F}_{cat}$. The final alpha mattes can be generated from the concatenated features $\mathcal{F}_{cat}$ through two $3*3$ convolutional layers, a sigmoid activation function, and the upsampling operation. The detailed upsampling process and convolutional layers in the decoder stage are all demonstrated in Figure 2.

Our information assembly module and concatenation operation can be formally defined as follows:

$$\mathcal{W}_{aspp} = \mathcal{W}_{aspp}/(\mathcal{W}_{aspp} + \mathcal{W}_{sed}) + \varepsilon, \qquad (2)$$

$$\mathcal{W}_{sed} = \mathcal{W}_{sed}/(\mathcal{W}_{aspp} + \mathcal{W}_{sed}) + \varepsilon, \qquad (3)$$

$$\mathcal{F}_{IA} = \mathcal{C}at(\mathcal{W}_{aspp} * \mathcal{F}_{aspp}, \mathcal{W}_{sed} * \mathcal{F}_{sed}), \qquad (4)$$

$$\mathcal{F}_{cat} = \mathcal{C}at[\mathcal{U}p(\mathcal{C}onv(\mathcal{C}onv(\mathcal{F}_{IA}))), \mathcal{F}_{ini}], \qquad (5)$$

where $\mathcal{W}_{aspp}$ and $\mathcal{W}_{sed}$ represent the assembly coefficient scalars which can balance dynamic weights between the ASPP features and the secondary superficial traces, the parameter $\varepsilon$ is to prevent the denominator from being zero. The feature maps $\mathcal{F}_{aspp}$ and $\mathcal{F}_{sed}$ can be adaptively integrated with this dynamic weights adjustment mode. The final concatenated features $\mathcal{W}_{cat}$ can be obtained by Equation (5), and here the $\mathcal{U}p$ and $\mathcal{C}onv$ denote the upsampling and convolution operation, respectively.

Such an information assembly strategy can effectively integrate multi-scale foreground expression to achieve high-quality alpha mattes. The high-level features from the ASPP module can produce foreground semantics, suggesting the subject content of the foreground. In contrast, the superficial traces from low-level CNN features contain different-level foreground details or textures which may share various opacity about the foreground boundaries. Their assembly can restore foreground information through dynamically adapted features integration. We conduct some experiments to evaluate the proposed MSIA-matte model, and the relevant ablation study can also prove the effectiveness of the proposed superficial traces branch and information assembly strategy.

### 3.4. Loss Function

We use a blended loss function to optimize our MSIA-matte network model, which combines $\mathcal{L}_1$ and $\mathcal{L}_{SSIM}$ to supervise the network training. The $\mathcal{L}_1$ is defined as follows:

$$\mathcal{L}_1 = \sum_k^{\Omega} |\alpha_m^k - \alpha_t^k|, \quad \alpha_m^k, \alpha_t^k \in [0, 1], \qquad (6)$$

where $\Omega$ represents pixels set, and the value $k$ denotes the pixel index. $\alpha_m^k$ is the alpha value of pixel $k$ in the regressed alpha matte, while $\alpha_t^k$ represents the corresponding ground-truth value in the same pixel location. Compared to $\mathcal{L}_2$, $\mathcal{L}_1$ is more robust for noticeable alpha value differences (the foreground is 1, while the background is 0). With $\mathcal{L}_1$ in the optimization, the model can pay more

attention to absolute alpha values, which can promote the evacuation of different-level foreground information. The $\mathcal{L}_{SSIM}$ expression is listed as:

$$\mathcal{L}_{SSIM} = 1 - \frac{(2\mu_m\mu_t + c_1)(2\sigma_{mt} + c_2)}{(\mu_m^2 + \mu_t^2 + c_1)(\sigma_m^2 + \sigma_t^2 + c_2)}, \quad (7)$$

where $\mu_m$, $\mu_t$, and $\sigma_m$, $\sigma_t$ are the mean and standard deviations of $\alpha_m^k$ and $\alpha_t^k$, according to [ZBSS04]. The effectiveness of $\mathcal{L}_{SSIM}$ for image matting has been proved in [QLY*20]. The $\mathcal{L}_{SSIM}$ can improve the foreground integrity and enhance the correlation between local regions, while the $\mathcal{L}_1$ can provide pixel-wise accuracy. Thus our total loss function can be defined as:

$$\mathcal{L}_{total} = \lambda_1\mathcal{L}_1 + \lambda_2\mathcal{L}_{SSIM}, \quad (8)$$

where $\lambda_1$ and $\lambda_2$ are coefficients to balance these two losses. In the training procedure, $\lambda_1$ and $\lambda_2$ are set as 1 and 0.1 for the first epoch and modified to 1 and 0.025 in the following epochs.

## 4. Experiments

We conduct our experiments on two large-scale matting datasets: Composition-1K and Distinctions-646. The Composition-1K dataset is from [XPCH17], which contains 431 available alpha mattes for training and 50 alpha mattes for testing. The Distinctions-646, provided by [QLY*20], consists of 596 training alpha mattes and 50 testing ones. We follow the rules in [XPCH17] to produce 100 composite examples for each training alpha matte and 20 examples for each alpha matte in the testing set. We train our MSIA-matte on both datasets and evaluate the test and analysis of two trained models. We also perform an ablation study to verify different modules and finally demonstrate our results on natural images.

### 4.1. Implementation and comparison details

**Implementation environment and parameters.** We establish our MSIA-matte model by PyTorch and use 3 Tesla P100 graphics to train and test the model parallelly. For training, we use the pre-trained ResNeXt-101 network [XGD*17] to initialize the weights in our backbone and the other weights are assigned with a random Gaussian distribution. All input images during training are randomly cropped to 512 × 512, 640 × 640 or 800 × 800, then resized to a resolution of 512 × 512 and augmented by horizontal random flipping. We use the stochastic gradient descent (*SGD*) optimizer with a momentum of 0.9 and a weight decay of 0.0005 to optimize our network parameters. The learning rate is initialized to 0.01, adjusted by the "poly" policy [LRB15] with the power of 0.9 for 20 epochs. During training, the batch size is set as 4, which requires 20 hours to finish 20 epochs, and the test on 1000 images only needs several minutes.

**Experiment metrics for evaluation.** For quantitative comparisons with the state-of-the-art methods, we refer to the four standard metrics in [RRW*09]: the sum of absolute differences (SAD), mean square error (MSE), the gradient, and connectivity. To better display the visual quality, we zoom in some critical details of the alpha mattes and place them under the images to show the differences more concretely.

| Methods | SAD↓ | MSE↓ | Gradient↓ | Connectivity↓ |
|---|---|---|---|---|
| Shared Matting [GO10] | 125.37 | 0.029 | 144.28 | 123.53 |
| Learning Based [ZK09] | 95.04 | 0.018 | 76.63 | 98.92 |
| Comprehensive [SRPC13b] | None | None | None | None |
| Global Matting [HRR*11] | 156.88 | 0.042 | 112.28 | 155.08 |
| ClosedForm [LLW07] | 124.68 | 0.025 | 115.31 | 106.06 |
| KNN Matting [CLT13] | 126.24 | 0.025 | 117.17 | 131.05 |
| DCNN [CTK16] | 115.82 | 0.023 | 107.36 | 111.23 |
| Information Flow [AAP17] | 70.36 | 0.013 | 42.79 | 70.66 |
| DIM [XPCH17] | 48.87 | 0.008 | 31.04 | 50.36 |
| AlphaGAN [LAS18] | 90.94 | 0.018 | 93.92 | 95.29 |
| SampleNet [TAO*19] | 48.03 | 0.008 | 35.19 | 56.55 |
| Context Aware [HL19] | 38.73 | 0.004 | 26.13 | 35.89 |
| IndexNet [LDSX19] | 44.52 | 0.005 | 29.88 | 42.37 |
| Late Fusion [ZGF*19] | 58.34 | 0.011 | 41.63 | 59.74 |
| HAttMatting [QLY*20] | **44.01** | **0.007** | 29.26 | 46.41 |
| Ours | 47.86 | **0.007** | **28.61** | **43.39** |

**Table 1:** *The quantitative results on the Composition-1K testing set. The methods in gray (the Late Fusion [ZGF*19], HAttMatting [QLY*20], and our method) can produce alpha mattes with single RGB images, while others require trimaps to confine the transition regions.*

### 4.2. The Composition-1K testing dataset

Here we train the proposed MSIA-matte based on the training images in [XPCH17] and evaluate the optimized model on the Composition-1K testing set. We also compare our model with state-of-the-art methods, including Shared Matting [GO10], Learning-Based [ZK09], Comprehensive Sampling [SRPC13b], Global Matting [HRR*11], ClosedForm [LLW07], KNN Matting [CLT13], Information Flow [AAP17], DCNN [CTK16], DIM [XPCH17], AlphaGAN [LAS18], SampleNet [TAO*19], Context-aware [HL19], IndexNet [LDSX19], Late Fusion [ZGF*19] and HAttMatting [QLY*20]. Table 1 summarizes the results of all the methods in calculating the four metrics on the full alpha images. The results in Table 1 are from our re-implementations or relevant researches [QLY*20].

It is noted that only the Late Fusion [ZGF*19], HAttMatting [QLY*20], and our MSIA-matte can achieve alpha mattes without trimaps in Table 1. The quality of trimap has a substantial impact on the final alpha mattes, which has been proved in [ZGF*19]. Compared with trimap-based methods, we can obtain more striking results than most of them. Only the latest Context-aware [HL19] and IndexNet [LDSX19] are better than our model. For trimap-free methods, we are second in the SAD metric to HAttMatting [QLY*20] and get the same results on the MSE metric. The proposed MSIA-matte can achieve the best performance on the gradient and connectivity metrics. According to [RRW*09], the SAD and MSE metrics can represent pixel-wise accuracy, while the gradient and connectivity can indicate visual quality. Therefore, the comparisons with trimap-free methods suggest that our model can produce a better judgment of visual quality. This promotion can contribute to our multi-scale information assembly, which can synthetically integrate different-level foreground expressions to achieve the estimation of alpha mattes. From another side, our multi-scale foreground expression may attenuate the subject se-
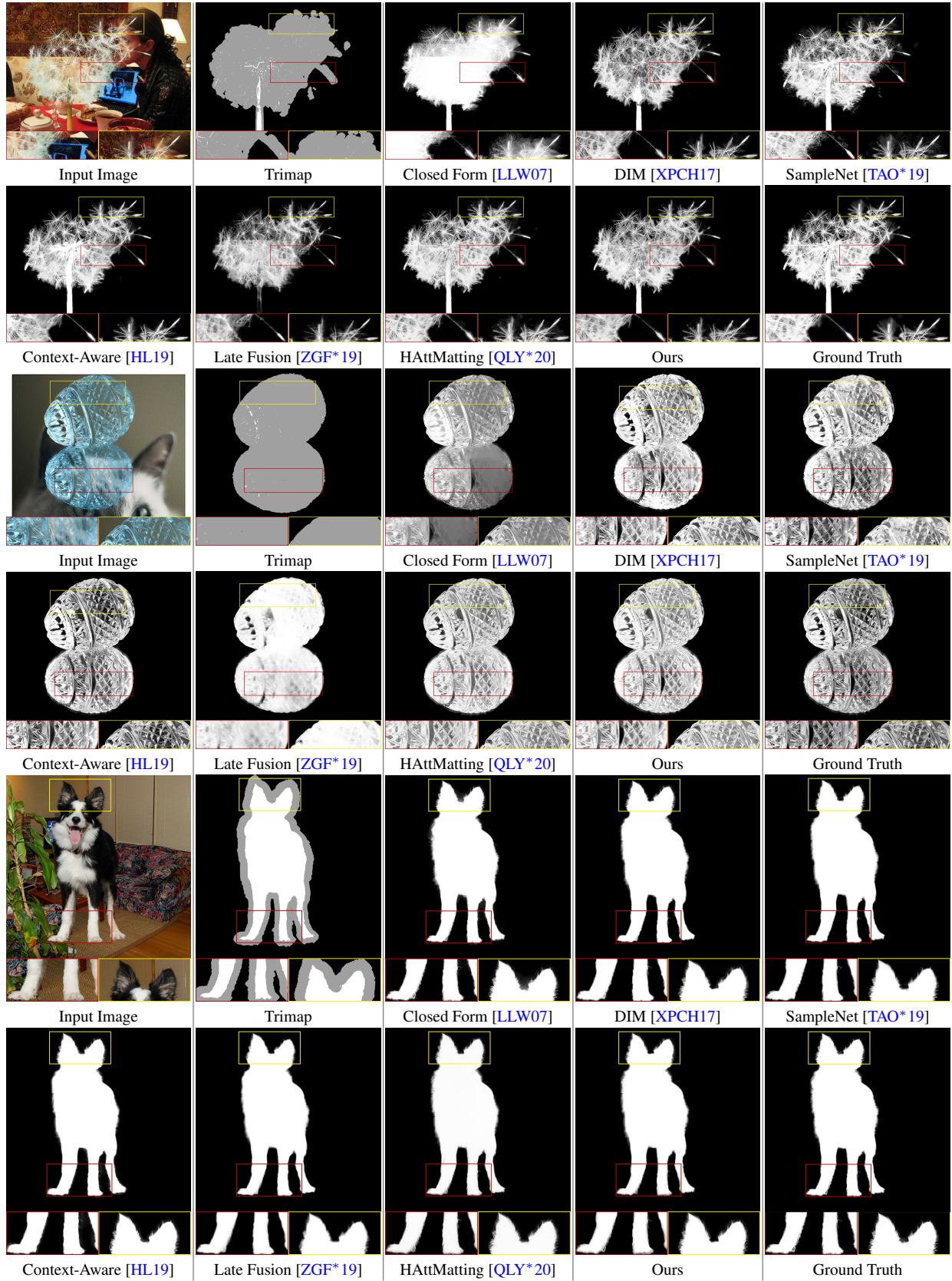
**Figure 3:** *Comparison of the alpha mattes on the Composition-1k testing set [XPCH17]. Some details and textures are zoomed in and placed under the images.*

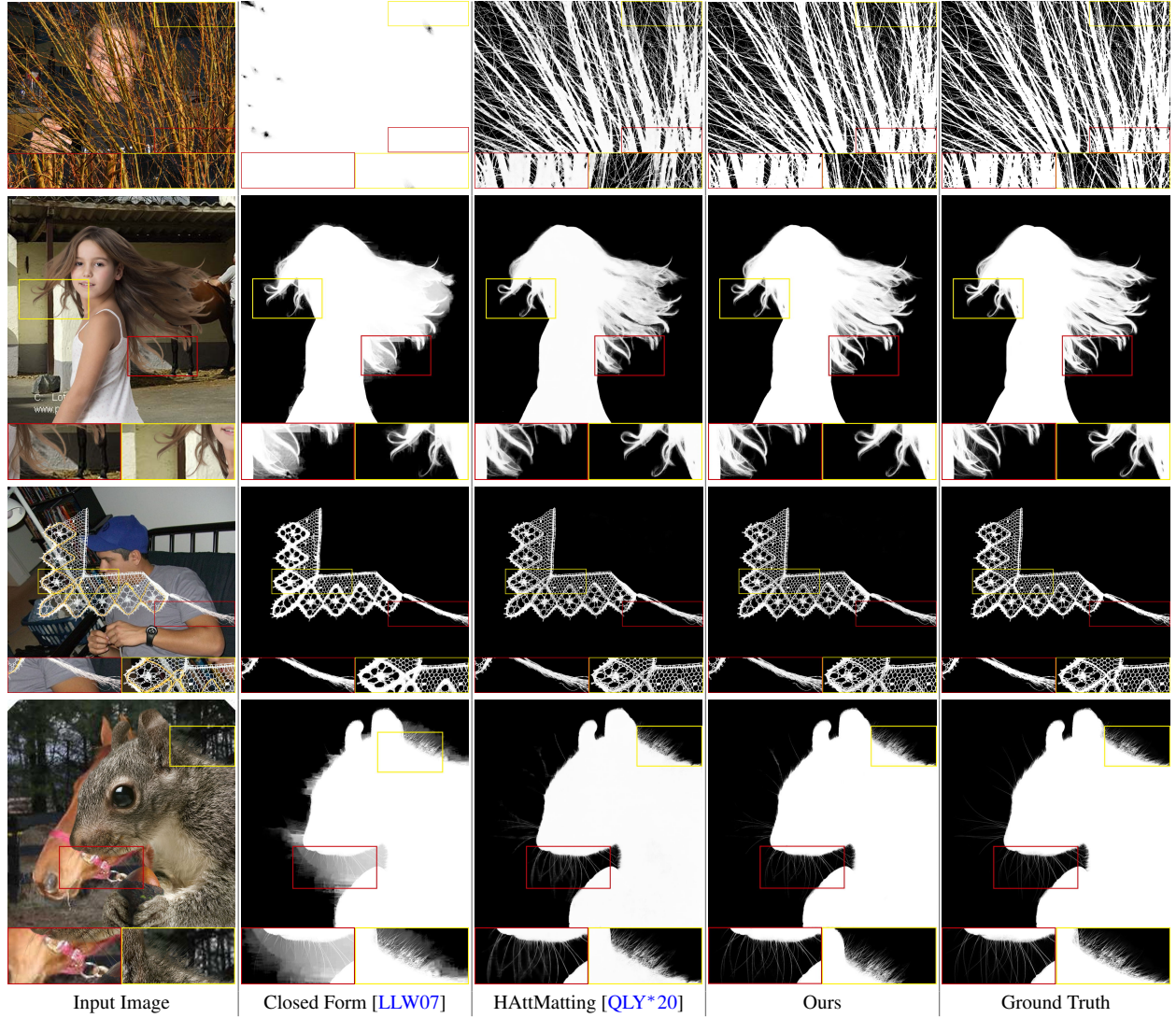|  |  |  |  |  |
|---|---|---|---|---|
| Input Image | Closed Form [LLW07] | HAttMatting [QLY*20] | Ours | Ground Truth |

**Figure 4:** *Comparison of the alpha mattes on the Distinctions-646 testing set [QLY*20]. Some details and texture are zoomed in and placed under the images.*

mantics of the foreground object to some extent and lead to some undersaturated pixels (e.g., some of the subject regions should be 255, but are only 252 or 253), which potentially reduces the pixel-wise accuracy. The visual results of some available methods are shown in Figure 3. From the first and second rows, we can see that the MSIA-matte can effectively capture subtle variations in opacity and correspond to visual quality promotion.

### 4.3. The Distinctions-646 dataset

In this experiment, we train and evaluate the MSIA-matte model on the Distinctions-646 dataset [QLY*20]. We also compare our model with traditional methods and some available matting networks, including Shared Matting [GO10], Learning-Based [ZK09], Global Matting [HRR*11], ClosedForm [LLW07], KNN Matting [CLT13], Information Flow [AAP17], DCNN [CTK16],

DIM [XPCH17], and HAttMatting [QLY*20]. Table 2 reports the results of all these methods.

Compared to Composition-1K [XPCH17] dataset, the Distinctions-646 [QLY*20] has more kinds of foreground objects. Correspondingly, all metrics have a slight rise, but the entire conclusion is similar to the results of the Composition-1K testing set. Compared to existing matting methods, we can achieve state-of-the-art performance. Our model has the best performance in gradient and connectivity, which suggests that the information assembly can improve visual quality. Moreover, the visual quality of our method can also be proved in Figure 4. We can adapt to continuous changes in the opacity and integrate multi-scale foreground information.

| Methods | SAD↓ | MSE↓ | Gradient↓ | Connectivity↓ |
|---|---|---|---|---|
| Shared Matting [GO10] | 119.56 | 0.026 | 129.61 | 114.37 |
| Learning Based [ZK09] | 105.04 | 0.021 | 94.16 | 110.41 |
| Global Matting [HRR*11] | 135.56 | 0.039 | 119.53 | 136.44 |
| ClosedForm [LLW07] | 105.73 | 0.023 | 91.76 | 114.55 |
| KNN Matting [CLT13] | 116.68 | 0.025 | 103.15 | 121.45 |
| DCNN [CTK16] | 103.81 | 0.020 | 82.45 | 99.96 |
| Information Flow [AAP17] | 78.89 | 0.016 | 58.72 | 80.47 |
| DIM [XPCH17] | 47.56 | 0.009 | 43.29 | 55.90 |
| HAttMatting [QLY*20] | **48.98** | 0.009 | 41.57 | 49.93 |
| Ours | 49.51 | **0.009** | **39.70** | **46.68** |

**Table 2:** *The quantitative results on the Distinctions-646 testing set.*

| IniST | SedST | AI | SAD↓ | MSE↓ | Gradient↓ | Connectivity↓ |
|---|---|---|---|---|---|---|
| | | | 69.02 | 0.013 | 39.91 | 67.62 |
| ✓ | | | 64.97 | 0.012 | 35.65 | 61.73 |
| ✓ | ✓ | | 52.35 | 0.010 | 34.75 | 53.40 |
| ✓ | ✓ | ✓ | **47.86** | **0.007** | **28.61** | **43.39** |
| | | | 64.73 | 0.016 | 71.34 | 72.48 |
| ✓ | | | 60.28 | 0.014 | 60.56 | 63.14 |
| ✓ | ✓ | | 57.47 | 0.012 | 45.34 | 60.03 |
| ✓ | ✓ | ✓ | **49.51** | **0.009** | **39.70** | **46.68** |

**Table 3:** *The quantitative results to compare the benefits of the proposed components. The "IniST", "SedST", and "AI" represent our initial and second phases of superficial traces branch and adaptive integration, respectively. The first four rows of results are evaluated on the Adobe Composition-1K testing dataset [XPCH17], and the last four results are based on the Distinctions-646 [QLY*20].*

## 4.4. Ablation study

We propose the superficial traces branch to approach low-level CNN features, and then combine different-level superficial traces with high-level semantics to integrate multi-scale foreground expression. In this section, we remove some modules to prove their benefits for the promotion of final alpha mattes. For concise description, here we denote the initial two convolutional layers of the superficial traces branch as IniST, and the later four cascaded convolutional layers are represented as SedST. For the baseline comparison, we remove IniST, SedST, and Adaptive Integration (AI), and use direct concatenation between low-level CNN features and upsampled ASPP features (the upsampling factor is 4) to produce alpha mattes. And then, we add IniST to the baseline and downsample the features from IniST as secondary superficial traces to perform following information assembly with direct concatenation. Finally, we add IniST and SedST to the baseline and only replace adaptive integration (AI) with direct concatenation to execute the decoder stage. The quantitative result of the above three models is summarized in Table 3.

The first four rows of Table 3 are evaluated on the Adobe Composition-1K testing set, while the last four results are conducted on the Distinctions-646 dataset. Compared to the baseline network, the IniST can bring an obvious improvement for the alpha mattes. This comparison can prove the ability of the IniST to filter low-level CNN features, which are relatively primitive for informa-



**Figure 5:** *The results of natural images. The first row is the input images, and the second row is the alpha mattes produced by our method.*

tion assembly. Correspondingly, the SedST can further condense potential middle-levels semantics (such as hands, leaves) and preserve some primary textures and details, and the promotion of the results are also reflected on the four metrics (the SAD, MSE, Gradient and Connectivity dropped by 12.62, 0.002, 0.9 and 8.33, respectively). The subtle promotion on the MSE and Gradient two metrics are mainly due to the versatility of the SedST: almost all the images contain fine-grained textures and details, but not such middle-level foreground expression. Therefore, we import our adaptive integration to achieve information assembly effectively. The four metrics are further improved, as reported in the bold results of Table 3.

## 4.5. Results on natural images

The results of our model on natural images are shown in Fig. 5. Our MSIA-matte can effectively deal with different-level foreground information (the thick trunk and the thin branches of the first column, the various opacity of the net). Besides, we can even adapt to the background of continuous changes and predict a high-quality alpha matte (the boy in the second column).

## 4.6. Limitations analysis

There are also some failure cases of our method, as shown in Figure 6. We utilize different-level foreground expression to regress the alpha mattes, which possibly attenuate the significance of the subject semantics thus can sometimes result in subtle region missing. Besides, there are two additional possible limitations to our method: (1) just like most deep-learning-based matting networks, our model is fully trained on composite datasets, and the illumination, reflection, emptiness, etc. in natural images may have a huge influence on the alpha mattes. (2) the proposed MSIA-matte can predict high-quality alpha mattes without trimaps, but our model may lose efficacy if the input images are extremely complicated (the region of interest is not commonsensible foreground or there are multiple co-existing foreground objects), and the trimaps-based methods may work better in such situations.
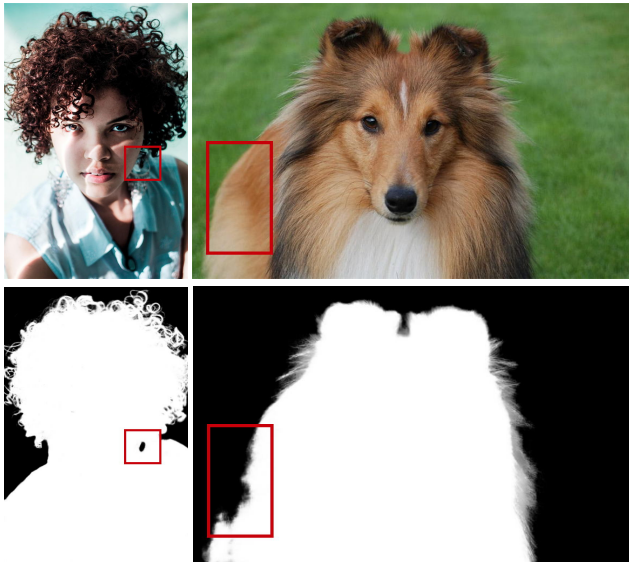
**Figure 6:** *The failure case of our model. The first row is the input images, and the second row is the alpha mattes produced by our method. Based on our analysis in section 4.2, the proposed multi-scale information assembly can possibly attenuate the impact of the advanced semantics and produce some subtle region missing in natural images (the red box).*

## 5. Conclusion

In this paper, we propose a multi-scale information assembly model (MSIA-matte) to extract and integrate different foreground expression from single RGB input images. The multi-scale information assembly strategy can effectively integrate the advanced semantics from ASPP and the superficial traces from low-level CNN features.

For future work, we intend to explore the possibility of the foreground information assembly in video matting. The temporary and spatial information may be considered to further enhance the correlation of multi-scale foreground expression between different frames.

## References

[AAP17] AKSOY Y., AYDIN T. O., POLLEFEYS M.: Designing effective inter-pixel information flow for natural image matting. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR'17)* (2017), pp. 228–236. 2, 5, 7, 8

[AOP*18] AKSOY Y., OH T.-H., PARIS S., POLLEFEYS M., MATUSIK W.: Semantic soft segmentation. *ACM Transactions on Graphics 37*, 4 (2018). 3

[CGX*18] CHEN Q., GE T., XU Y., ZHANG Z., YANG X., GAI K.: Semantic human matting. In *Proceedings of the ACM International Conference on Multimedia (MM'18)* (2018), p. 618–626. 3

[CKTK17] CHO D., KIM S., TAI Y. W., KWEON I. S.: Automatic trimap generation and consistent matting for light-field images. *IEEE Transactions on Pattern Analysis and Machine Intelligence 39*, 8 (2017), 1504–1517. 3

[CLT13] CHEN Q., LI D., TANG C.: Knn matting. *IEEE Transactions on Pattern Analysis and Machine Intelligence 35*, 9 (2013), 2175–2188. 2, 5, 7, 8

[CPK*18] CHEN L., PAPANDREOU G., KOKKINOS I., MURPHY K., YUILLE A. L.: Deeplab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected crfs. *IEEE Transactions on Pattern Analysis and Machine Intelligence 40*, 4 (2018), 834–848. 3

[CTK16] CHO D., TAI Y. W., KWEON I.: Natural image matting using deep convolutional neural networks. In *Proceedings of the European Conference on Computer Vision (ECCV'16)* (2016), pp. 626–643. 2, 5, 7, 8

[CZF*19] CAI S., ZHANG X., FAN H., HUANG H., LIU J., LIU J., LIU J., WANG J., SUN J.: Disentangled image matting. In *Proceedings of the International Conference on Computer Vision (ICCV'19)* (2019), pp. 8818–8827. 2, 3

[FLZ] FENG X., LIANG X., ZHANG Z.: A cluster sampling method for image matting via sparse coding. In *Proceedings of the European Conference on Computer Vision (ECCV'16)*, pp. 204–219. 2

[GO10] GASTAL E. S. L., OLIVEIRA M. M.: Shared sampling for real-time alpha matting. *Computer Graphics Forum 29*, 2 (2010), 575–584. 5, 7, 8

[GPAM*14] GOODFELLOW I. J., POUGET-ABADIE J., MIRZA M., XU B., WARDE-FARLEY D., OZAIR S., COURVILLE A., BENGIO Y.: Generative adversarial nets. In *Proceedings of the International Conference on Neural Information Processing Systems (NeurIPS'14)* (2014), p. 2672–2680. 2

[GSAW05] GRADY L., SCHIWIETZ T., AHARON S., WESTERMANN R.: Random walks for interactive alpha-matting. In *Proceedings of the Visualization, Imaging, and Image Processing (VIIP'05)* (2005), pp. 423–429. 2

[HL19] HOU Q., LIU F.: Context-aware image matting for simultaneous foreground and alpha estimation. In *Proceedings of the International Conference on Computer Vision (ICCV'19)* (2019), pp. 4129–4138. 2, 3, 4, 5, 6

[HRR*11] HE K., RHEMANN C., ROTHER C., TANG X., SUN J.: A global sampling method for alpha matting. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR'11)* (2011), pp. 2049–2056. 2, 5, 7, 8

[HZRS16] HE K., ZHANG X., REN S., SUN J.: Deep residual learning for image recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR'16)* (2016), pp. 770–778. 4

[KEE15] KARACAN L., ERDEM A., ERDEM E.: Image matting with kl-divergence based sparse sampling. In *Proceedings of the International Conference on Computer Vision (ICCV'15)* (2015), pp. 424–432. 2

[LAS18] LUTZ S., AMPLIANITIS K., SMOLIC A.: Alphagan: Generative adversarial networks for natural image matting. 259. 2, 5

[LDSX19] LU H., DAI Y., SHEN C., XU S.: Indices matter: Learning to index for deep image matting. In *Proceedings of the International Conference on Computer Vision (ICCV'19)* (2019), pp. 3265–3274. 3, 5

[LL20] LI Y., LU H.: Natural image matting via guided contextual attention. In *AAAI* (2020). 3

[LLW07] LEVIN A., LISCHINSKI D., WEISS Y.: A closed-form solution to natural image matting. *IEEE Transactions on Pattern Analysis and Machine Intelligence 30*, 2 (2007), 228–242. 2, 5, 6, 7, 8

[LRB15] LIU W., RABINOVICH A., BERG A. C.: Parsenet: Looking wider to see better. *arXiv preprint arXiv:1506.04579* (2015). 5

[LRL08] LEVIN A., RAV-ACHA A., LISCHINSKI D.: Spectral matting. *IEEE Transactions on Pattern Analysis and Machine Intelligence 30*, 10 (2008), 1699–1712. 2

[LSD15] LONG J., SHELHAMER E., DARRELL T.: Fully convolutional networks for semantic segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR'15)* (2015), pp. 3431–3440. 4

[LW11] LEE P., WU Y.: Nonlocal matting. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR'11)* (2011), pp. 2193–2200. 2

[QLY*20] QIAO Y., LIU Y., YANG X., ZHOU D., XU M., ZHANG Q., WEI X.: Attention-guided hierarchical structure aggregation for image matting. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR'20)* (2020). 2, 3, 4, 5, 6, 7, 8

[RRW*09] RHEMANN C., ROTHER C., WANG J., GELAUTZ M., KOHLI P., ROTT P.: A perceptually motivated online benchmark for image matting. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR'09)* (2009), pp. 1826–1833. 5

[SJTS04] SUN J., JIA J., TANG C. K., SHUM H. Y.: Poisson matting. *ACM Transactions on Graphics 23*, 3 (2004), 315–321. 2

[SRPC13a] SHAHRIAN E., RAJAN D., PRICE B., COHEN S.: Improving image matting using comprehensive sampling sets. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR'13)* (2013), pp. 636–643. 2

[SRPC13b] SHAHRIAN E., RAJAN D., PRICE B., COHEN S.: Improving image matting using comprehensive sampling sets. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR'13)* (2013), pp. 636–643. 5

[SZ15] SIMONYAN K., ZISSERMAN A.: Very deep convolutional networks for large-scale image recognition. In *Proceedings of the International Conference on Learning Representations (ICLR'15)* (2015). 4

[TAO*19] TANG J., AKSOY Y., OZTIRELI C., GROSS M., AYDIN T. O.: Learning-based sampling for natural image matting. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR'19)* (2019), pp. 3050–3058. 3, 5, 6

[WC07] WANG J., COHEN M. F.: Optimized color sampling for robust matting. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR'07)* (2007), pp. 1–8. 2

[XGD*17] XIE S., GIRSHICK R., DOLLÁR P., TU Z., HE K.: Aggregated residual transformations for deep neural networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR'17)* (2017), pp. 5987–5995. 3, 5

[XPCH17] XU N., PRICE B., COHEN S., HUANG T.: Deep image matting. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR'17)* (2017), pp. 311–320. 2, 3, 5, 6, 7, 8

[YCSS01] YUNG-YU CHUANG, CURLESS B., SALESIN D. H., SZELISKI R.: A bayesian approach to digital matting. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR'01)* (2001), pp. II–II. 2

[YXC*18] YANG X., XU K., CHEN S., HE S., YIN B. Y., LAU R.: Active matting. In *Proceedings of the International Conference on Neural Information Processing Systems (NeurIPS'18)* (2018), pp. 4590–4600. 3

[ZBSS04] ZHOU WANG, BOVIK A. C., SHEIKH H. R., SIMONCELLI E. P.: Image quality assessment: from error visibility to structural similarity. *IEEE Transactions on Image Processing 13*, 4 (2004), 600–612. 5

[ZGF*19] ZHANG Y., GONG L., FAN L., REN P., HUANG Q., BAO H., XU W.: A late fusion cnn for digital matting. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR'19)* (2019), pp. 7461–7470. 2, 3, 4, 5, 6

[ZK09] ZHENG Y., KAMBHAMETTU C.: Learning based digital matting. In *Proceedings of the International Conference on Computer Vision (ICCV'09)* (2009), pp. 889–896. 5, 7, 8