

Information-Flow Matting[†]

Yağız Aksoy^{1,2}, Tunç Ozan Aydın², Marc Pollefeys¹

¹ ETH Zürich

² Disney Research Zürich

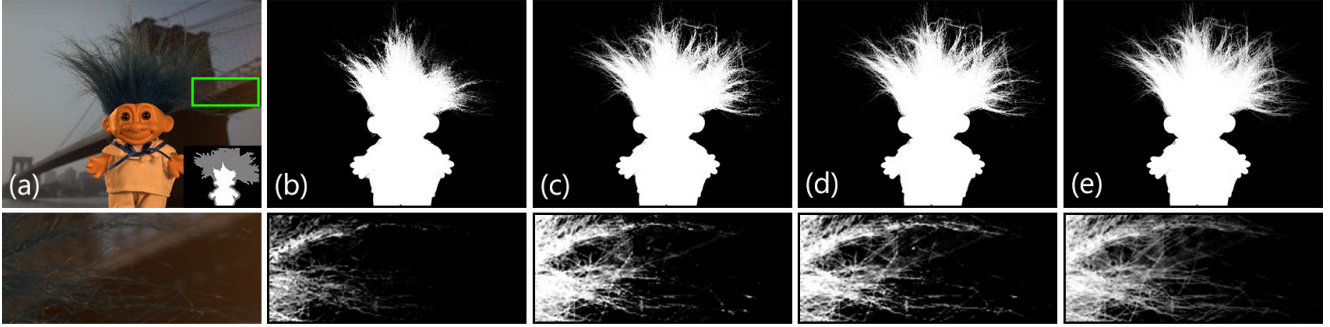


Figure 1: For an input image and a trimap (a), we construct our linear system by first using the color-mixture flow(b), then adding direct channels of information flow from known to unknown regions (c), and letting information be shared effectively inside the unknown region (d). We finally introduce local information flow to enforce spatial smoothness (e). Note that the intermediate results in this figure are solely for illustration. In practice, we construct a single energy function that accounts for all types of information flow and solve it once to obtain the end result.

Abstract

We present a novel, purely affinity-based natural image matting algorithm. Our method relies on carefully defined pixel-to-pixel connections that enable effective use of information available in the image. We control the information flow from the known-opacity regions into the unknown region, as well as within the unknown region itself, by utilizing multiple definitions of pixel affinities. Among other forms of information flow, we introduce color-mixture flow, which builds upon local linear embedding and effectively encapsulates the relation between different pixel opacities. Our resulting novel linear system formulation can be solved in closed-form and is robust against several fundamental challenges of natural matting such as holes and remote intricate structures. While our method is primarily designed as a standalone matting tool, we show that it can also be used for regularizing mattes obtained by sampling-based methods. The formulation is also extended to layer color estimation and we show that the use of multiple channels of flow increases the layer color quality. We also demonstrate our performance in green-screen keying and analyze the characteristics of the utilized affinities.

1. Introduction

Extracting the opacity information of foreground objects from an image is known as natural image matting. Natural image matting has received great interest from the research community through the last decade and can nowadays be considered as one of the classical research problems in visual computing. Mathematically, image matting requires expressing pixel colors in the transition regions from foreground to background as a convex combination of their underlying foreground and background colors. The weight, or the *opacity*, of the foreground color is referred to as the *alpha* value of that pixel. Since neither the foreground and background colors nor the opacities are known, estimating the opacity values is a highly ill-posed problem. To alleviate the difficulty of this problem, typically a *trimap* is provided in addition to the original image. The trimap is a rough segmentation of the input image into foreground, background, and regions with unknown opacity.

The main application of natural image matting is compositing, *i.e.* combining different scenes together to generate a new image. Image matting methods aim to provide accurate opacities such that when the foreground is overlaid onto a novel background, the transitions between them look natural. However, together with the matte, compositing requires the actual, *unmixed* layer colors for realistic composites. The layer colors appear as a mixture of foreground and background colors in the input image, and they are un-

[†]This document is an extended version of the 2017 CVPR publication titled *Designing effective inter-pixel information flow for natural image matting*.

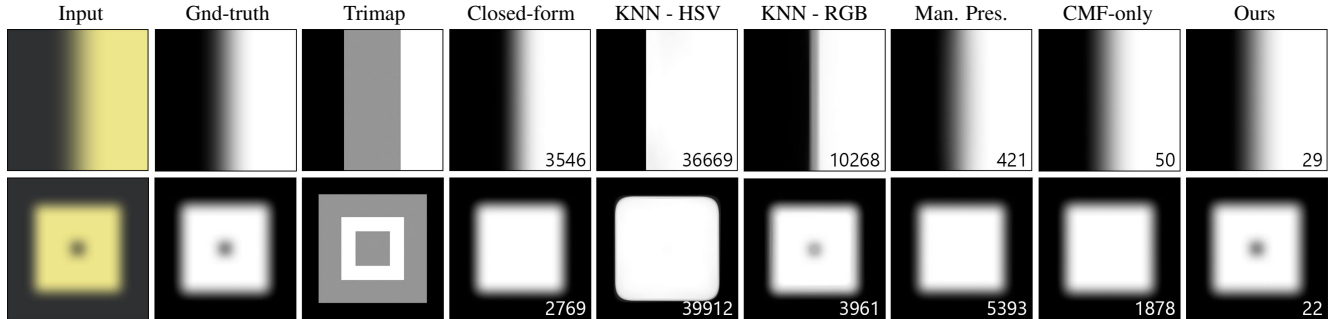


Figure 2: We created two duotone 500x500 images and blurred them to get soft transitions between regions. The numbers show the sum of absolute differences between the estimated alpha mattes and the ground truth. Closed-form matting [14] uses local information flow, KNN Matting [5] uses HSV- or RGB-based similarity measure, and manifold-preserving edit propagation [6] uses LLE weights [17]. We observe a performance improvement in large opacity gradients even when only the color-mixture flow (CMF) is used (Section 3.1). Notice also that both large gradients and holes are recovered with high performance using our final formulation. See text for further discussion.

derconstrained even with a given matte. Hence, accurate estimation of the layer colors is a critical component of a compositing pipeline and still an active research problem.

Affinity-based methods [14, 5, 6] constitute one of the prominent natural matting approaches in literature. These methods make use of pixel similarities to propagate the alpha values from the known-alpha regions to the unknown region. They provide a clear mathematical formulation, can be solved in closed-form, are easy to implement, and typically produce spatially consistent mattes. In addition, due to their formulation that can be modeled as a graph structure with each pixel as a node, affinity-based approaches can be generalized to related applications such as layer color estimation [14], edit propagation [6], and soft segmentation [15]. Studying affinity-based approaches for natural matting can open new directions for a larger set of applications in the image processing community.

In spite of these advantages, current affinity-based methods fail to effectively handle alpha gradients spanning large areas and spatially disconnected regions (i.e. *holes*) even in simple cases as demonstrated in Figure 2. This is because a straightforward formulation using the pixel-to-pixel affinity definitions can not effectively represent the complex structures that are commonly seen in real-life objects. In order to alleviate these shortcomings, we rely on a careful, case-by-case design of how alpha values should propagate inside the image. We conceptualize the affinities as *information flows* to help understanding and designing effective graph-based structures to propagate information in the image. We define several information flows, some of which target unknown-opacity regions that are remote and hence does not receive enough information in previous formulations. Other types of information flows address issues such as evenly distributing information inside the unknown region. We formulate this strategy through the use of a vari-

ety of affinity definitions including the *color-mixture flow*, which is based on local linear embedding and tailored for image matting. Step-by-step improvements on the matte quality as we gradually add new building blocks of our information flow strategy are illustrated in Figure 1. Our final linear system can be solved in closed-form and results in a significant quality improvement over the state-of-the-art. We demonstrate the matting quality improvement quantitatively, as well as through a visual inspection of challenging image regions. We also show that our energy function can be reformulated as a post-processing step for regularizing the spatially inconsistent mattes estimated by sampling-based natural matting algorithms.

This document is an extended version of our CVPR publication [1]. In this extended version, we additionally (i) propose a novel foreground color estimation formulation where we introduce a new form of local information flow, (ii) demonstrate that our method achieves state-of-the-art quality in green-screen keying, (iii) provide an in-depth spectral analysis of individual forms of information flow, and (iv) present a discussion on how our method relates to sampling-based matting methods, as well as new quantitative and qualitative results.

2. Related work

Natural Image Matting The numerous natural matting methods in the literature can be mainly categorized as sampling-based, learning-based or affinity-based. We briefly review the most relevant here and refer the reader to a comprehensive survey [24] for further information.

Sampling-based methods [10, 13, 18, 11] typically seek to gather numerous samples from the background and foreground regions defined by the trimap and select the best-fitting pair according to their individually defined criteria for representing an unknown pixel as a mixture of fore-

ground and background. While they perform well especially around remote and challenging structures, they require affinity-based regularization to produce spatially consistent mattes. Also, their methodology typically focuses solely on matting and they typically can not generalize any other applications unlike the affinity-based counterparts.

Machine learning has been used to aid in estimating the alpha in a semi-supervised setting [23], to estimate a trimap in constrained settings [19] or to combine results of other matting methods for a better matte [8]. Recently, a deep neural network architecture has been proposed [22] that generates high-quality mattes with the help of semantic knowledge that can be extracted from the image. In order to train such a network, Xu *et al.* [22] generated a dataset of 50k images with ground-truth mattes. Our method outperforms all current learning-based methods in the alpha matting benchmark [16] despite not taking advantage of a large dataset with labels. We hope that our formulation and the concepts presented in the paper will inspire next-generation learning-based matting methods.

Affinity-based matting methods mainly make use of pixel similarity metrics that rely on color similarity or spatial proximity and propagate the alpha values from regions with known opacity. Local affinity definitions, prominently the matting affinity [14], operate on a local patch around the pixel to determine the amount of information flow and propagate alpha values accordingly. The matting affinity is also adopted in a post-processing step in most sampling-based methods as proposed by Gastal and Oliveira [11].

Methods utilizing nonlocal affinities similarly use color similarity and spatial proximity for determining how the alpha values of different pixels should relate to each other. KNN matting [5] determines several neighbors for every unknown pixel and enforces them to have similar alpha values relative to their distance in a feature space. The manifold-preserving edit propagation algorithm [6] also determines a set of neighbors for every pixel but represents each pixel as a linear combination of its neighbors in their feature space.

Chen *et al.* [7] proposed a hybrid approach that uses the sampling-based robust matting [21] as a starting point and refines its outcome through a graph-based technique where they combine a nonlocal affinity [6] and the matting affinity. Cho *et al.* [8] combined the results of closed-form matting [14] and KNN matting [5], as well as the sampling-based method comprehensive sampling [18], by feeding them into a convolutional neural network.

In this work, we propose color-mixture flow and discuss its advantages over the affinity definition utilized by Chen *et al.* [6]. We also define three other forms of information flow, which we use to carefully distribute the alpha information inside the unknown region. Our approach differs from Chen *et al.* [7] in that our information flow strategy goes beyond combining various pixel affinities, as we dis-

cuss further in Section 3, while requiring much less memory to solve the final system. Instead of using the results of other affinity-based methods directly as done by Cho *et al.* [8], we formulate an elegant formulation that has a closed-form solution. To summarize, we present a novel, purely affinity-based matting algorithm that generates high-quality alpha mattes without making use of sampling or a learning step.

Layer Color Estimation For a given alpha matte, the corresponding foreground colors should also be estimated before compositing. Although the alpha matte is assumed to be given for the foreground color estimation, the problem is still underconstrained as there are 6 unknowns and 3 equations. Levin *et al.* [14] use the gradient of the alpha matte as a spatial smoothness measure and formulate the layer color estimation as a linear problem. Using only a smoothness measure limits their performance especially in remote regions of the foreground. Chen *et al.* [5] use the color-similarity measure they employ for matte estimation also for layer color estimation. Typically, using only a color-similarity metric results in incorrectly flat-colored regions and suppressed highlight colors in the foreground. In this work, we introduce a second spatial smoothness measure for the layer colors. We use in total 4 forms of information flow together for the layer estimation and show that our linear system improves the layer color quality especially in remote parts of the matte.

Green-Screen Keying A more constrained version of the natural image matting problem is referred as green-screen keying, where the background colors are homogeneous in a controlled setting. While this problem can be seen as a simpler version of natural matting, as green-screen keying is heavily utilized in professional production [2], the expected quality of the results is immense. In the movie post-production industry, multiple commercial software such as Keylight or Primatte are used by professional graphical artists to get high-quality keying results. These software typically use chroma-based or luma-based algorithms and provide many parameters that help the artist tweak the results. In their early work, Smith and Blinn [20] formulate the use of the compositing equation for a fixed background color. Recently, an unmixing-based green-screen keying method has been proposed [2] that uses a global color model of the scene and a per-pixel nonlinear energy function to extract the background color in high precision. In their paper, they compare their method to state-of-the-art natural matting methods and show that the current matting methods fail to give acceptable results in green-screen settings. In this paper, we show that our matting and color estimation methods outperform the natural matting methods and generate comparable results to that of specialized keying methods or commercial software without any parameter tweaking.

3. Method

Trimaps are typically given as user input in natural matting, and they consist of three regions: fully opaque (foreground), fully transparent (background) and of unknown opacity. \mathcal{F} , \mathcal{B} and \mathcal{U} will respectively denote these regions, and \mathcal{K} will represent the union of \mathcal{F} and \mathcal{B} . Affinity-based methods operate by propagating opacity information from \mathcal{K} into \mathcal{U} using a variety of affinity definitions. We define this flow of information in multiple ways so that each pixel in \mathcal{U} receives information effectively from different regions in the image.

The opacity transitions in a matte occur as a result of the original colors in the image getting mixed with each other due to transparency or intricate parts of an object. We make use of this fact by representing each pixel in \mathcal{U} as a mixture of similarly-colored pixels and defining a form of information flow that we call *color-mixture flow* (Section 3.1). We also add connections from every pixel in \mathcal{U} to both \mathcal{F} and \mathcal{B} to facilitate direct information flow from known-opacity regions to even the most remote opacity-transition regions in the image (Section 3.2). In order to distribute the information from the color-mixture and \mathcal{K} -to- \mathcal{U} flows, we define intra- \mathcal{U} flow of information, where pixels with similar colors inside \mathcal{U} share information on their opacity with each other (Section 3.3). Finally, we add local information flow, a pixel affecting the opacity of its immediate spatial neighbors, which ensures spatially coherent end results (Section 3.4). We formulate the individual forms of information flow as energy functions and aggregate them in a global optimization formulation (Section 3.5).

3.1. Color-mixture information flow

Due to transparent objects as well as fine structures and sharp edges of an object that cannot be fully captured due to the finite-resolution of the imaging sensors, certain pixels of an image inevitably contain a mixture of corresponding foreground and background colors. By investigating these color mixtures, we can derive an important clue on how to propagate alpha values between pixels. The amount of the original foreground color in a particular mixture determines the opacity of the pixel. Following this fact, if we represent the color of a pixel as a weighted combination of the colors of several others, those weights should correspond to the opacity relation between the pixels.

In order to make use of this relation, for every pixel in \mathcal{U} , we find $K_{CM} = 20$ similar pixels in a feature space by an approximate K nearest neighbors search in the whole image. We define the feature vector for this search as $[r, g, b, \tilde{x}, \tilde{y}]^T$, where \tilde{x} and \tilde{y} are the image coordinates normalized by image width and height, and the rest are the RGB values of the pixel. This set of neighbors, selected as similar-colored pixels that are also close-by, is denoted by \mathcal{N}_p^{CM} .

We then find the weights of the combination $w_{p,q}^{CM}$ that will determine the amount of information flow between the pixels p and $q \in \mathcal{N}_p^{CM}$. The weight of each neighbor is defined such that the weighted combination of their colors yields the color of the original pixel:

$$\arg \min_{w_{p,q}^{CM}} \left\| \mathbf{c}_p - \sum_{q \in \mathcal{N}_p^{CM}} w_{p,q}^{CM} \mathbf{c}_q \right\|^2, \quad (1)$$

where \mathbf{c}_p represents the 3x1 vector of RGB values. We minimize this energy using the method by Roweis and Saul [17]. Note that since we are only using RGB values, the neighborhood correlation matrix computed during the minimization has a high chance of being singular as there could easily be two neighbors with identical colors. So, we condition the neighborhood correlation matrix by adding $10^{-3} I_{K_{CM} \times K_{CM}}$ to it before inversion, where $I_{K_{CM} \times K_{CM}}$ is the identity matrix.

Note that while we use the method by Roweis and Saul [17] to minimize the energy in (1), we do not fully adopt their local linear embedding (LLE) method. LLE finds a set of neighbors in a feature space and uses all the variables in the feature space to compute the weights in order to reduce the dimensionality of input data. Manifold-preserving edit propagation [6] and LNSP matting [7] algorithms make use of the LLE weights directly in their formulation for image matting. However, since we are only interested in the weighted combination of colors and not the spatial coordinates, we exclude the spatial coordinates in the energy minimization step. This increases the validity of the estimated weights, effects of which can be observed even in the simplest cases such as in Figure 2, where manifold-preserving weight propagation and CMF-only results only differ in the weight computation step.

The energy term for the color-mixture flow is defined as:

$$E_{CM} = \sum_{p \in \mathcal{U}} \left(\alpha_p - \sum_{q \in \mathcal{N}_p^{CM}} w_{p,q}^{CM} \alpha_q \right)^2. \quad (2)$$

3.2. \mathcal{K} -to- \mathcal{U} information flow

The color-mixture flow already provides useful information on how the mixed-color pixels are formed. However, many pixels in \mathcal{U} receive information present in the trimap indirectly through their neighbors, all of which can possibly be in \mathcal{U} . This indirect information flow might not be enough especially for remote regions that are far away from \mathcal{K} .

In order to facilitate the flow of information from both \mathcal{F} and \mathcal{B} directly into every region in \mathcal{U} , we add connections from every pixel in \mathcal{U} to several pixels in \mathcal{K} . For each pixel in \mathcal{U} , we find $K_{\mathcal{KU}} = 7$ similar pixels in both \mathcal{F} and \mathcal{B} separately to form the sets of pixels $\mathcal{N}_p^{\mathcal{F}}$ and $\mathcal{N}_p^{\mathcal{B}}$.

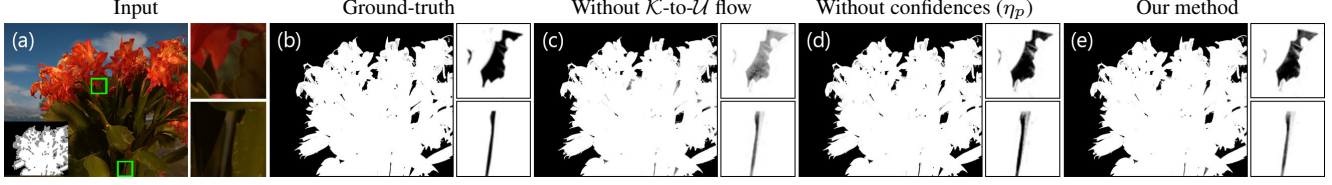


Figure 3: Direct information flow from both \mathcal{F} and \mathcal{B} to even the most remote regions in \mathcal{U} increases our performance around holes significantly (top inset). Using confidences further increases the performance, especially around regions where FG and BG colors are similar (bottom inset).

with K nearest neighbors search using the feature space $[r, g, b, 10 * \tilde{x}, 10 * \tilde{y}]^T$ to favor close-by pixels. We use the pixels in $\mathcal{N}_p^{\mathcal{F}}$ and $\mathcal{N}_p^{\mathcal{B}}$ together to represent the pixel color c_p by minimizing the energy in (1). Using the resulting weights $w_{p,q}^{\mathcal{F}}$ and $w_{p,q}^{\mathcal{B}}$, we define an energy function to represent the \mathcal{K} -to- \mathcal{U} flow:

$$E_{\mathcal{KU}} = \sum_{p \in \mathcal{U}} \left(\alpha_p - \sum_{q \in \mathcal{N}_p^{\mathcal{F}}} w_{p,q}^{\mathcal{F}} \alpha_q - \sum_{q \in \mathcal{N}_p^{\mathcal{B}}} w_{p,q}^{\mathcal{B}} \alpha_q \right)^2 \quad (3)$$

Note that $\alpha_q = 1$ for $q \in \mathcal{F}$ and $\alpha_q = 0$ for $q \in \mathcal{B}$. This fact allows us to define two combined weights, one connecting a pixel to \mathcal{F} and another to \mathcal{B} , as:

$$w_p^{\mathcal{F}} = \sum_{q \in \mathcal{N}_p^{\mathcal{F}}} w_{p,q}^{\mathcal{F}} \quad \text{and} \quad w_p^{\mathcal{B}} = \sum_{q \in \mathcal{N}_p^{\mathcal{B}}} w_{p,q}^{\mathcal{B}} \quad (4)$$

such that $w_p^{\mathcal{F}} + w_p^{\mathcal{B}} = 1$, and rewrite (3) as:

$$E_{\mathcal{KU}} = \sum_{p \in \mathcal{U}} (\alpha_p - w_p^{\mathcal{F}})^2. \quad (5)$$

The energy minimization in (1) gives us similar weights for all q when c_q are similar to each other. As a result, if $\mathcal{N}_p^{\mathcal{F}}$ and $\mathcal{N}_p^{\mathcal{B}}$ have pixels with similar colors, the estimated weights $w_p^{\mathcal{F}}$ and $w_p^{\mathcal{B}}$ become unreliable. We account for this fact by augmenting the energy function in (5) with confidence values.

We can determine the colors contributing to the mixture estimated by (1) using the weights $w_{p,q}^{\mathcal{F}}$ and $w_{p,q}^{\mathcal{B}}$:

$$c_p^{\mathcal{F}} = \frac{\sum_{q \in \mathcal{N}_p^{\mathcal{F}}} w_{p,q}^{\mathcal{F}} c_q}{w_p^{\mathcal{F}}}, \quad c_p^{\mathcal{B}} = \frac{\sum_{q \in \mathcal{N}_p^{\mathcal{B}}} w_{p,q}^{\mathcal{B}} c_q}{w_p^{\mathcal{B}}}, \quad (6)$$

and define a confidence metric according to how similar the estimated foreground color $c_p^{\mathcal{F}}$ and background color $c_p^{\mathcal{B}}$ are:

$$\eta_p = \|c_p^{\mathcal{F}} - c_p^{\mathcal{B}}\|^2 / 3. \quad (7)$$

The division by 3 is to get the confidence values between $[0, 1]$. We update the new energy term to reflect our confidence in the estimation:

$$\tilde{E}_{\mathcal{KU}} = \sum_{p \in \mathcal{U}} \eta_p (\alpha_p - w_p^{\mathcal{F}})^2. \quad (8)$$

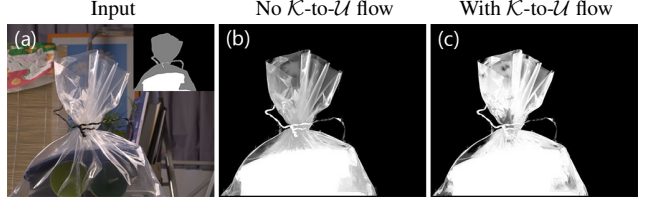


Figure 4: \mathcal{K} -to- \mathcal{U} flow does not perform well when the foreground object is highly-transparent. See text for discussion.

This update to the energy term increases the matting quality in regions with similar foreground and background colors, as seen in Figure 3.

It should be noted that the \mathcal{K} -to- \mathcal{U} flow is not reliable when the foreground is highly transparent, as seen in Figure 4. This is mainly due to the low representational power of $\mathcal{N}_p^{\mathcal{F}}$ and $\mathcal{N}_p^{\mathcal{B}}$ for c_p around large highly-transparent regions as the nearest neighbors search does not give us well-fitting pixels for $w_{p,q}^{\mathcal{F}}$ estimation. We construct our final linear system accordingly in Section 3.5.

3.2.1 Pre-processing the trimap

Prior to determining $\mathcal{N}_p^{\mathcal{F}}$ and $\mathcal{N}_p^{\mathcal{B}}$, we pre-process the input trimap in order to facilitate finding more reliable neighbors, which in turn increases the effectiveness of the \mathcal{K} -to- \mathcal{U} flow. Trimaps usually have regions marked as \mathcal{U} despite being fully opaque or transparent, as drawing a very detailed trimap is both cumbersome and prone to errors.

Several methods [10, 13] refine the trimap as a pre-processing step by expanding \mathcal{F} and \mathcal{B} starting from their boundaries with \mathcal{U} as proposed by Shahrian *et al.* [18]. Incorporating this technique improves our results as shown in Figure 5(d). We also apply this extended \mathcal{F} and \mathcal{B} regions after the matte estimation as a post-processing. Since this trimap trimming method propagates known regions only to nearby pixels, in addition to this edge-based trimming, we also make use of a patch-based trimming step.

To this end, we extend the transparent and opaque regions by relying on patch statistics. We fit a 3D RGB normal distribution N_p to the 3×3 window around each pixel p . In order to determine the most similar distribution

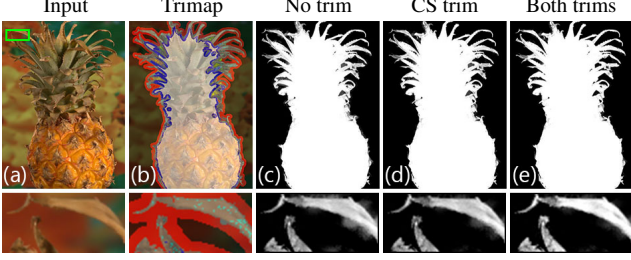


Figure 5: The trimap is shown overlaid on the original image (b) where the extended foreground regions are shown with blue (CS trimming [18]) and cyan (patch-search) and the extended background regions with red (CS trimming) and yellow (patch-search). CS trimming makes the fully opaque / transparent regions cleaner, while our trimming improves the results around remote structures.

in \mathcal{F} for a pixel $p \in \mathcal{U}$, we first find the 20 distributions with closest mean vectors. We define the foreground match score $b_p^{\mathcal{F}} = \min_{q \in \mathcal{F}} B(N_p, N_q)$, where $B(\cdot, \cdot)$ represents the Bhattacharyya distance between two distributions. We find the match score for background $b_p^{\mathcal{B}}$ the same way. We then select a region for pixel p according to the following rule:

$$p \in \begin{cases} \hat{\mathcal{F}} & \text{if } b_p^{\mathcal{F}} < \tau_c \text{ and } b_p^{\mathcal{B}} > \tau_f \\ \hat{\mathcal{B}} & \text{if } b_p^{\mathcal{B}} < \tau_c \text{ and } b_p^{\mathcal{F}} > \tau_f \\ \hat{\mathcal{U}} & \text{otherwise} \end{cases} \quad (9)$$

Simply put, an unknown pixel is marked as $\hat{\mathcal{F}}$, *i.e.* in foreground after trimming, if it has a strong match in \mathcal{F} and no match in \mathcal{B} , which is determined by constants $\tau_c = 0.25$ and $\tau_f = 0.9$. By inserting known-alpha pixels in regions far away from \mathcal{U} - \mathcal{K} boundaries, we further increase the matting performance in challenging remote regions (Figure 5(e)).

3.3. Intra- \mathcal{U} information flow

Each individual pixel in \mathcal{U} receives information through the color-mixture and \mathcal{K} -to- \mathcal{U} flows. In addition to these, we would like to distribute the information inside \mathcal{U} effectively. We achieve this by encouraging pixels with similar colors inside \mathcal{U} to have similar opacity.

For each pixel in \mathcal{U} , we find $K_{\mathcal{U}} = 5$ nearest neighbors only inside \mathcal{U} to determine $\hat{\mathcal{N}}_p^{\mathcal{U}}$ using the feature vector defined as $\mathbf{v} = [r, g, b, \tilde{x}/20, \tilde{y}/20]^T$. Notice that we scale the coordinate members of the feature vector we used in Section 3.1 to decrease their effect on the nearest neighbor selection. This lets $\hat{\mathcal{N}}_p^{\mathcal{U}}$ have pixels inside \mathcal{U} that are far away, so that the information moves more freely inside the unknown region. We use the neighborhood $\mathcal{N}_p^{\mathcal{U}} = \hat{\mathcal{N}}_p^{\mathcal{U}} \cup \{q \mid p \in \hat{\mathcal{N}}_q^{\mathcal{U}}\}$ to make sure that information flows both ways between p to $q \in \hat{\mathcal{N}}_p^{\mathcal{U}}$. We then determine the amount of information flow using the L^1 distance

between feature vectors:

$$w_{p,q}^{\mathcal{U}} = \max(1 - \|\mathbf{v}_p - \mathbf{v}_q\|_1, 0) \quad \forall q \in \mathcal{N}_p^{\mathcal{U}}. \quad (10)$$

The energy term for intra- \mathcal{U} flow then can be defined as:

$$E_{\mathcal{U}\mathcal{U}} = \sum_{p \in \mathcal{U}} \sum_{q \in \mathcal{N}_p^{\mathcal{U}}} w_{p,q}^{\mathcal{U}} (\alpha_p - \alpha_q)^2. \quad (11)$$

The information sharing between the unknown pixels increases the matte quality around intricate structures as demonstrated in Figure 1(d).

KNN matting [5] uses a similar affinity definition to make similar-color pixels have similar opacities. However, relying only on this form of information flow for the whole image creates some typical artifacts in the matte. Depending on the feature vector definition and the image colors, the matte may erroneously underrepresent the smooth transitions (KNN - HSV case in Figure 2) when the neighbors of the pixels in \mathcal{U} happen to be mostly in only \mathcal{F} or \mathcal{B} , or create flat alpha regions instead of subtle gradients (KNN - RGB case in Figure 2). Restricting information flow to be solely based on color similarity fails to represent the complex alpha transitions or wide regions with an alpha gradient.

3.4. Local information flow

Spatial connectivity is one of the main cues for information flow. We connect each pixel in \mathcal{U} to its 8 immediate neighbors denoted by \mathcal{N}_p^L to ensure spatially smooth mattes. The amount of local information flow should also adapt to strong edges in the image.

To determine the amount of local flow, we rely on the matting affinity definition proposed by Levin *et al.* [14]. The matting affinity utilizes the local patch statistics to determine the weights $w_{p,q}^L$, $q \in \mathcal{N}_p^L$. We define our related energy term as follows:

$$E_L = \sum_{p \in \mathcal{U}} \sum_{q \in \mathcal{N}_p^L} w_{p,q}^L (\alpha_p - \alpha_q)^2. \quad (12)$$

Despite representing local information flow well, matting affinity by itself fails to represent large transition regions (Figure 2 top), or isolated regions that have weak or no spatial connection to \mathcal{F} or \mathcal{B} (Figure 2 bottom).

3.5. Linear system and energy minimization

Our final energy function is a combination of the four energies representing the individual information flows:

$$E_1 = E_{CM} + \sigma_{\mathcal{K}\mathcal{U}} E_{\mathcal{K}\mathcal{U}} + \sigma_{\mathcal{U}\mathcal{U}} E_{\mathcal{U}\mathcal{U}} + \sigma_L E_L + \lambda E_{\mathcal{T}}, \quad (13)$$

where $\sigma_{\mathcal{K}\mathcal{U}} = 0.05$, $\sigma_{\mathcal{U}\mathcal{U}} = 0.01$, $\sigma_L = 1$ and $\lambda = 100$ are algorithmic constants determining the strength of corresponding information flows, and

$$E_{\mathcal{T}} = \sum_{p \in \mathcal{F}} (\alpha_p - 1)^2 + \sum_{p \in \mathcal{B}} (\alpha_p - 0)^2$$

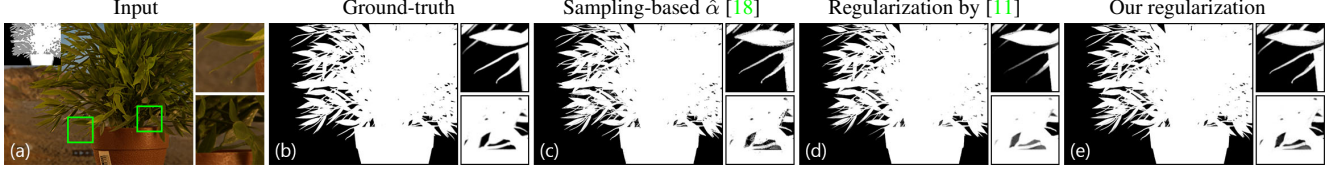


Figure 6: The matte regularization method by Gastal and Oliveira [11] loses remote details (top inset) or fills in holes (bottom inset) while our regularization method is able to preserve these details caught by the sampling-based method.

is the energy term to keep the known opacity values constant. For an image with N pixels, by defining $N \times N$ sparse matrices W_{CM} , W_{UU} and W_L that have non-zero elements for the pixel pairs with corresponding information flows and the vector w^F that has elements w_p^F for $p \in \mathcal{U}$, 1 for $p \in \mathcal{F}$ and 0 for $p \in \mathcal{B}$, we can write (13) in matrix form as:

$$E_1 = \alpha^T \mathcal{L}_{IFM} \alpha + (\alpha - w^F)^T \sigma_{KU} \mathcal{H} (\alpha - w^F) + (\alpha - \alpha_K)^T \lambda \mathcal{T} (\alpha - \alpha_K), \quad (14)$$

where \mathcal{T} is an $N \times N$ diagonal matrix with diagonal entry (p, p) 1 if $p \in \mathcal{K}$ and 0 otherwise, \mathcal{H} is a sparse matrix with diagonal entries η_p as defined in (7), α_K is a row vector with p^{th} entry being 1 if $p \in \mathcal{F}$ and 0 otherwise, α is a row-vector of the alpha values to be estimated, and \mathcal{L}_{IFM} is defined as:

$$\mathcal{L}_{IFM} = (D_{CM} - W_{CM})^T (D_{CM} - W_{CM}) + \sigma_{UU} (D_{UU} - W_{UU}) + \sigma_L (D_L - W_L), \quad (15)$$

where the diagonal matrix $D_{(\cdot)}(i, i) = \sum_j W_{(\cdot)}(i, j)$.

The energy in (14) can be minimized by solving

$$(\mathcal{L}_{IFM} + \lambda \mathcal{T} + \sigma_{KU} \mathcal{H}) \alpha = (\lambda \mathcal{T} + \sigma_{KU} \mathcal{H}) w^F. \quad (16)$$

We define a second energy function that excludes the \mathcal{K} -to- \mathcal{U} information flow:

$$E_2 = E_{CM} + \sigma_{UU} E_{UU} + \sigma_L E_L + \lambda E_{\mathcal{T}}, \quad (17)$$

which can be written in matrix form as:

$$E_2 = \alpha^T \mathcal{L}_{IFM} \alpha + (\alpha - \alpha_K)^T \lambda \mathcal{T} (\alpha - \alpha_K), \quad (18)$$

and can be minimized by solving:

$$(\mathcal{L}_{IFM} + \lambda \mathcal{T}) \alpha = \lambda \mathcal{T} \alpha_K. \quad (19)$$

We solve the linear systems of equations in (16) and (19) using the preconditioned conjugate gradients method [4].

As mentioned before, the \mathcal{K} -to- \mathcal{U} information flow is not effective for highly transparent objects. To determine whether to include the \mathcal{K} -to- \mathcal{U} information flow and solve for E_1 , or to exclude it and solve for E_2 for a given image, we use a simple histogram-based classifier to determine if we expect a highly transparent result.

If the matte is highly transparent, the pixels in \mathcal{U} are expected to mostly have colors that are a mixture of \mathcal{F} and \mathcal{B} colors. On the other hand, if the true alpha values are mostly 0 or 1 except for soft transitions, the histogram of \mathcal{U} will likely be a linear combination of the histograms of \mathcal{F} and \mathcal{B} as \mathcal{U} will mostly include very similar colors to that of \mathcal{K} . Following this observation, we attempt to express the histogram of the pixels in \mathcal{U} , $\mathcal{D}_{\mathcal{U}}$, as a linear combination of $\mathcal{D}_{\mathcal{F}}$ and $\mathcal{D}_{\mathcal{B}}$. The histograms are computed from the 20 pixel-wide region around \mathcal{U} in \mathcal{F} and \mathcal{B} , respectively. We define the error e , the metric of how well the linear combination represents the true histogram, as:

$$e = \min_{a,b} \|a \mathcal{D}_{\mathcal{F}} + b \mathcal{D}_{\mathcal{B}} - \mathcal{D}_{\mathcal{U}}\|^2. \quad (20)$$

Higher e values indicate a highly-transparent matte, in which case we prefer E_2 over E_1 .

4. Matte regularization for sampling-based matting methods

Sampling-based natural matting methods usually select samples for each pixel in \mathcal{U} either independently or by paying little attention to spatial coherency. In order to obtain a spatially coherent matte, the common practice is to combine their initial guesses for alpha values with a smoothness measure. Multiple methods [10, 11, 13, 18] adopt the post-processing method proposed by Gastal and Oliveira [11] which combines the matting affinity [14] with the sampling-based alpha values and corresponding confidences. This post-processing technique leads to improved mattes, but since it involves only local smoothness, the results can still be suboptimal as seen in Figure 6(d).

Our approach with multiple forms of information flow can also be used for post-processing in a way similar to that of Gastal and Oliveira [11]. Given the initial alpha values $\hat{\alpha}_p$ and confidences $\hat{\eta}_p$ found by a sampling-based method, we define the matte regularization energy:

$$E_R = E_2 + \sigma_R \sum_{p \in \mathcal{U}} \hat{\eta}_p (\alpha_p - \hat{\alpha}_p)^2, \quad (21)$$

where $\sigma_R = 0.05$ determines how much loyalty should be given to the initial values. This energy can be written in matrix form and solved as a linear system in the same way we did in Section 3.5.

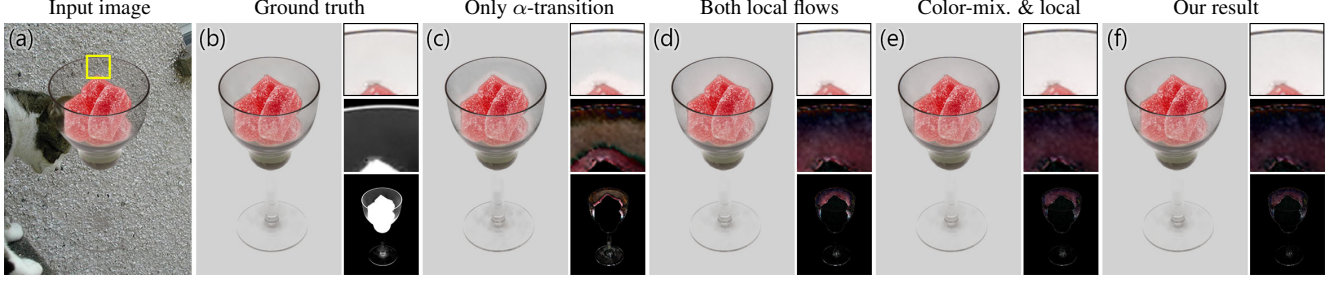


Figure 7: Color estimation results using a growing set of information flows using the ground truth matte. The bottom-right in each set shows per-pixel absolute difference between the estimation and ground truth multiplied by ten. See text for discussion.

Figure 6 shows that this non-local regularization of mattes is more effective especially around challenging foreground structures such as long leaves or holes as seen in the insets. In Section 6.2, we will numerically explore the improvement we achieve by replacing the matte regularization step with ours in several sampling-based methods.

5. Foreground color estimation

In addition to the alpha matte, we need the *unmixed* foreground colors [2] that got into the color mixture in transition pixels for seamlessly compositing the foreground onto a novel background. Similar to Levin *et al.* [14] and Chen *et al.* [5], we estimate the foreground colors for a given matte, after the matte estimation.

We propagate the layer colors from opaque and transparent regions in a similar way we propagate known alpha values in Section 3. We make use of the color-mixture and the intra- \mathcal{U} information flows by extending the search space and affinity computation to include the given alpha values together with spatial coordinates and pixel colors. We also use the spatial smoothness measure proposed by Levin *et al.* [14] in addition to a second spatial smoothness measure we introduce in this paper. Figure 7 shows how our color estimation result improves as we add more forms of information flow.

5.1. Information flow definitions

In the layer color estimation problem, the input is assumed to be the original image together with an alpha matte. This requires us to redefine the three regions using the matte instead of a trimap:

$$p \in \begin{cases} \tilde{\mathcal{F}} & \text{if } \tilde{\alpha}_p = 1 \\ \tilde{\mathcal{B}} & \text{if } \tilde{\alpha}_p = 0 \\ \tilde{\mathcal{U}} & \text{otherwise.} \end{cases} \quad (22)$$

$\tilde{\alpha}_p$ denote the alpha values that are given as input. The foreground and background colors to be estimated will be denoted by \mathbf{f} and \mathbf{b} . For a pixel p , the compositing equation

we would like to satisfy can be written as:

$$\mathbf{c}_p = \tilde{\alpha}_p \mathbf{f}_p + (1 - \tilde{\alpha}_p) \mathbf{b}_p \quad (23)$$

We will formulate the energy functions for a single color channel and solve for red, green and blue channels independently. The scalars f and b will denote the values for a single color channel.

5.1.1 Local information flows

Levin *et al.* [14] proposed the use of the gradient of the alpha channel as the amount of local information flow for the problem of layer color estimation. They solely rely on this form of information flow for propagating the colors. This local information flow basically enforces neighboring pixels to have similar colors if there is an alpha transition. This flow, which we refer to as *alpha-transition flow*, can be represented by the following energy:

$$E_{\nabla \tilde{\alpha}} = \sum_{\forall p} \sum_{q \in \mathcal{N}_p^L} |\nabla \tilde{\alpha}_{(p-q)}| \left((f_p - f_q)^2 + (b_p - b_q)^2 \right), \quad (24)$$

where $\nabla \tilde{\alpha}$ represents the alpha gradient. We compute the gradients in the image plane using the 3-tap separable filters of Farid and Simoncelli [9]. Note that the neighborhood is defined as the local 3×3 neighborhood similar to the local information flow in Section 3.4.

The transition flow helps around small regions with alpha gradient but does not propagate information in flat-alpha regions, such as pure foreground or background regions or regions with flat opacity. We propose a new smoothness measure to address this issue, which we call *no-transition flow*. The no-transition flow enforces spatial smoothness in regions with small color and alpha gradients:

$$E_{\nabla c \tilde{\alpha}} = \sum_{\forall p} \sum_{q \in \mathcal{N}_p^L} w_{p,q}^{\nabla c \tilde{\alpha}} \left((f_p - f_q)^2 + (b_p - b_q)^2 \right) \quad (25)$$

where $w_{p,q}^{\nabla c \tilde{\alpha}} = (1 - |\nabla \tilde{\alpha}_{(p-q)}|) (1 - \|\nabla \mathbf{c}_{(p-q)}\|)$ and $\|\nabla \mathbf{c}_{(p-q)}\|$ is the L_2 norm of the vector formed by gradients

of the individual color channels. This term increases the performance around slow alpha transitions and flat-alpha regions, as well as around sharp color edges in the image.

No-transition flow already improves the performance quite noticeably as seen in Figure 7(b). However, using only local information flows perform poorly in remote areas such as the end of long hair filaments (Figure 10(a)) or isolated areas (Figure 7, bottom inset). In order to increase the performance in these type of challenging areas, we make use of two types of non-local information flows.

5.1.2 Color-mixture information flow

The basic principle of color mixture as introduced in Section 3.1 also applies to the relationship between layer colors of pixels in the same neighborhood — if we represent the color and alpha of a pixel as a weighted combination of the colors and alpha of several others, those weights should also represent the layer color relation between the pixels. Since we have $\tilde{\alpha}$'s as additional information in the layer color estimation scenario, we extend the formulation of color-mixture flow to better fit the layer color estimation problem. Similar to its use in alpha estimation, it provides a well-connected graph and allows dense share of information. The performance improvement by the introduction of the color-mixture energy can be seen in Figure 7(c).

In the layer color estimation scenario, we optimize for both foreground and background colors in the same formulation. It should be emphasized that, as it is apparent from (23), the foreground and background colors are undefined for regions with $\tilde{\alpha} = 0$ and $\tilde{\alpha} = 1$, respectively. This requires us to avoid using color-mixture flow into \mathcal{U} from $\tilde{\mathcal{B}}$ for \mathbf{f} and from $\tilde{\mathcal{F}}$ for \mathbf{b} . We address this by defining two different neighborhoods and computing individual color-mixture flows for \mathbf{f} and \mathbf{b} .

For \mathbf{f} , we define the neighborhood $\mathcal{N}_p^{\tilde{\mathcal{U}}\tilde{\mathcal{F}}}$ by finding K_{CM} nearest neighbors in $(\tilde{\mathcal{U}} \cup \tilde{\mathcal{F}})$ using the feature vector $[r, g, b, \tilde{\alpha}, \tilde{x}, \tilde{y}]^T$. We then compute the weights $w_{p,q}^{C\tilde{\mathcal{F}}}$ as

$$\arg \min_{w_{p,q}^{C\tilde{\mathcal{F}}}} \left\| \begin{bmatrix} \mathbf{c}_p \\ \tilde{\alpha}_p \end{bmatrix} - \sum_{q \in \mathcal{N}_p^{CM}} w_{p,q}^{C\tilde{\mathcal{F}}} \begin{bmatrix} \mathbf{c}_q \\ \tilde{\alpha}_q \end{bmatrix} \right\|^2. \quad (26)$$

Notice that the search space and the weight computation includes $\tilde{\alpha}$ in addition to the color and location of pixels.

We compute the background conjugates of the neighborhood and weights, $\mathcal{N}_p^{\tilde{\mathcal{U}}\tilde{\mathcal{B}}}$ and $w_{p,q}^{C\tilde{\mathcal{B}}}$, in the same way, and define our color-mixture energy for layer color estimation:

$$E_{CM}^{fb} = \sum_{p \in \tilde{\mathcal{U}}} ((f_p - \sum_{q \in \mathcal{N}_p^{\tilde{\mathcal{U}}\tilde{\mathcal{F}}}} w_{p,q}^{C\tilde{\mathcal{F}}} f_q)^2 + (b_p - \sum_{q \in \mathcal{N}_p^{\tilde{\mathcal{U}}\tilde{\mathcal{B}}}} w_{p,q}^{C\tilde{\mathcal{B}}} b_q)^2).$$

5.1.3 Intra- $\tilde{\mathcal{U}}$ information flow

Intra- \mathcal{U} information flow, as detailed in Section 3.3, distributes the information between similar-colored pixels inside the unknown region without giving spatial proximity too much emphasis. Its behaviour is also very useful in the case of color estimation, as it makes the foreground colors more coherent throughout the image. For example, in Figure 7, bottom inset shows that the addition of intra- \mathcal{U} flow helps in getting a more realistic color to the isolated plastic region between the two black lines.

We make modifications to intra- \mathcal{U} flow similar to the modifications we made to color-mixture flow, in order to make use of the available information coming from $\tilde{\alpha}$'s.

We find $K_{\mathcal{U}}$ nearest neighbors only inside $\tilde{\mathcal{U}}$ to determine $\hat{\mathcal{N}}_p^{\tilde{\mathcal{U}}}$ using the feature vector defined as $\mathbf{v}^c = [r, g, b, \tilde{\alpha}, \tilde{x}/20, \tilde{y}/20]^T$. We then determine the amount of information flow between two non-local neighbors as:

$$w_{p,q}^{\tilde{\mathcal{U}}} = \max \left(1 - \|\mathbf{v}_p^c - \mathbf{v}_q^c\|_1, 0 \right) \quad \forall q \in \mathcal{N}_p^{\tilde{\mathcal{U}}}. \quad (27)$$

With the weights determined, we can define the energy function representing the intra- $\tilde{\mathcal{U}}$ flow:

$$E_{\tilde{\mathcal{U}}\tilde{\mathcal{U}}} = \sum_{p \in \tilde{\mathcal{U}}} \sum_{q \in \mathcal{N}_p^{\tilde{\mathcal{U}}}} w_{p,q}^{\tilde{\mathcal{U}}} \left((f_p - f_q)^2 + (b_p - b_q)^2 \right). \quad (28)$$

Note that in the color estimation formulation, we exclude the \mathcal{K} -to- \mathcal{U} information flow because we observed that the adaptation of the method in Section 3.2 to color estimation does not improve the quality of the final result.

5.2. Linear system and energy minimization

The final energy function for layer color estimation is the combination of the four types of information flow defined in Sections 5.1.1 to 5.1.3:

$$E_c = \sigma_L E_{\nabla \alpha} + \sigma_L E_{\nabla c \alpha} + E_{CM}^{fb} + \sigma_{\mathcal{U}\mathcal{U}} E_{\tilde{\mathcal{U}}\tilde{\mathcal{U}}} + \lambda E_{\text{COMP}}, \quad (29)$$

where σ_L , $\sigma_{\mathcal{U}\mathcal{U}}$ and λ are defined in Section 3.5 and E_{COMP} represents the deviation from the compositing equation constraint:

$$E_{\text{COMP}} = \sum_{\forall p} (c_p - \alpha_p^I f - (1 - \alpha_p^I) b)^2. \quad (30)$$

E_c is defined and minimized independently for each color channel.

Following the same strategy as we did in Section 3.5, we rewrite the energy function E_c in the matrix form, this time as a $2N \times 2N$ linear system, and solve it for foreground and background colors for 3 times, once for each color channel, using the preconditioned conjugate gradients method [4].

Table 1: Our scores in the alpha matting benchmark [16] together with the top-performing published methods at the time of submission. S , L and U denote the three trimap types, small, large and user, included in the benchmark. Bold and blue numbers represent the best scores in the benchmark.

	Average Rank				Troll			Doll			Donkey			Elephant			Plant			Pineapple			Plastic bag			Net		
	Overall	S	L	U	S	L	U	S	L	U	S	L	U	S	L	U	S	L	U	S	L	U	S	L	U	S	L	U
Sum of Absolute Differences																												
Ours	2.7	3.3	2.3	2.6	10.3	11.2	12.5	5.6	7.3	7.3	3.8	4.1	3	1.4	2.3	2.0	5.9	7.1	8.6	3.6	5.7	4.6	18.3	19.3	15.8	20.2	22.2	22.3
DIM [22]	2.9	3.6	2.3	2.8	10.7	11.2	11.0	4.8	5.8	5.6	2.8	2.9	2.9	1.1	1.1	2.0	6.0	7.1	8.9	2.7	3.2	3.9	19.2	19.6	18.7	21.8	23.9	24.1
DCNN [8]	4.0	5.4	2.3	4.3	12.0	14.1	14.5	5.3	6.4	6.8	3.9	4.5	3.4	1.6	2.5	2.2	6.0	6.9	9.1	4.0	6.0	5.3	19.9	19.2	19.1	19.4	20.0	21.2
CSC [10]	11	14.4	7.4	11.3	13.6	15.6	14.5	6.2	7.5	8.1	4.6	4.8	4.2	1.8	2.7	2.5	5.5	7.3	9.7	4.6	7.6	6.9	23.7	23.0	21.0	26.3	27.2	25.2
Mean Squared Error																												
Ours	4.0	5.4	2.8	3.8	0.3	0.4	0.5	0.3	0.4	0.5	0.3	0.3	0.2	0.1	0.1	0.1	0.4	0.4	0.6	0.2	0.3	0.3	1.3	1.2	0.8	0.8	0.8	0.9
DCNN [8]	4.3	5.3	2.5	5.0	0.4	0.5	0.7	0.2	0.3	0.4	0.2	0.3	0.2	0.1	0.1	0.1	0.4	0.4	0.8	0.2	0.4	0.3	1.3	1.2	1.0	0.7	0.7	0.9
DIM [22]	4.6	3.5	4.0	6.3	0.4	0.4	0.4	0.2	0.3	0.3	0.1	0.1	0.2	0	0	0.2	0.5	0.6	1	0.2	0.2	0.4	1.1	1.1	1.1	0.8	0.9	1
LNSP [7]	10.2	7.6	9.6	13.3	0.5	1.9	1.2	0.2	0.4	0.5	0.3	0.4	0.2	0.0	0.1	0.2	0.4	0.5	0.8	0.2	0.3	0.4	1.4	1.2	0.8	1.0	1.1	1.5

6. Results and discussion

We evaluate the proposed methods for matting, matte regularization, layer color estimation and green-screen keying with comparisons to the state-of-the-art of each application.

6.1. Matte estimation

We quantitatively evaluate the proposed algorithm using the public alpha matting benchmark [16] in Table 1. At the time of submission, our method ranks in the first place according to the sum-of-absolute-differences (SAD) and mean-squared error (MSE) metrics. Our proof-of-concept implementation in Matlab requires on average 50 seconds to process a benchmark image.

Our performance in the test set by Xu *et al.* [22] is shown in Table 2. This test set of 1000 images accompany a data-driven approach to matting. One advantage of using a deep network for this problem, such as DIM [22], is that the network can infer the matte even when there is no foreground region defined in the trimap due to heavy transparency, and their test set includes several such examples. Affinity-based and sampling-based approaches, however, assume both known regions are present when they are modeling the color models of affinities. While this can be seen as a shortcoming, the images without well-defined regions inadvertently skew the scores in this dataset. We perform better than competing methods except for DIM in this dataset, and our scores improve to be 76.5 (SAD) and 0.021 (MSE) when the images that violate our assumptions are removed.

Table 2: Matting performance on the test set of DIM [22].

	SAD	MSE
DIM [22]	50.4	0.014
Ours	100.6	0.038
DCNN [8]	161.4	0.087
CF [14]	168.1	0.091
KNN [5]	175.4	0.103

Table 3: Average percentage performance change with changing parameters using 27 images and 2 trimaps from the benchmark.

Param.	Def.	Val.	Perf.	Val.	Perf.	Val.	Perf.	Val.	Perf.
K_{CM}	20	10	1.07 %	15	0.44 %	25	-0.46 %	30	-0.62 %
K_{C-U}	7	1	-0.83 %	4	-0.41 %	10	0.12 %	13	0.22 %
K_{U-U}	5	1	-0.15 %	3	-0.1 %	7	0.08 %	9	0.11 %
σ_{K-U}	0.05	0.01	-6.44 %	0.025	-2.1 %	0.075	0.66 %	0.09	0.87 %
σ_{U-U}	0.01	0.001	-0.7 %	0.005	-0.1 %	0.02	-0.47 %	0.05	-3.12 %

We also compare our results qualitatively with the closely related methods in Figure 8. We use the results that are available on the matting benchmark for all except manifold-preserving matting [6] which we implemented ourselves. Figure 8(c,d,e) show that using only one form of information flow is not effective in a number of scenarios such as wide unknown regions or holes in the foreground. The strategy DCNN matting [8] follows is using the results of closed-form and KNN matting directly rather than formulating a combined energy using their affinity definitions. When both methods fail, the resulting combination also suffers from the errors as it is apparent in the pineapple and troll examples. The neural network they propose also seems to produce mattes that appear slightly blurred. LNSP matting [7], on the other hand, has issues around regions with holes (pineapple example) or when the foreground and background colors are similar (donkey and troll examples). It can also oversmooth some regions if the true foreground colors are missing in the trimap (plastic bag example). Our method performs well in these challenging scenarios mostly because of the intra-unknown and unknown-to-known connections which results in a more robust linear system.

We evaluate the sensitivity of our method against different parameter values on the training dataset of the matting benchmark [16]. Table 3 shows that different values for the parameters generally have only a small effect on the performance on average.

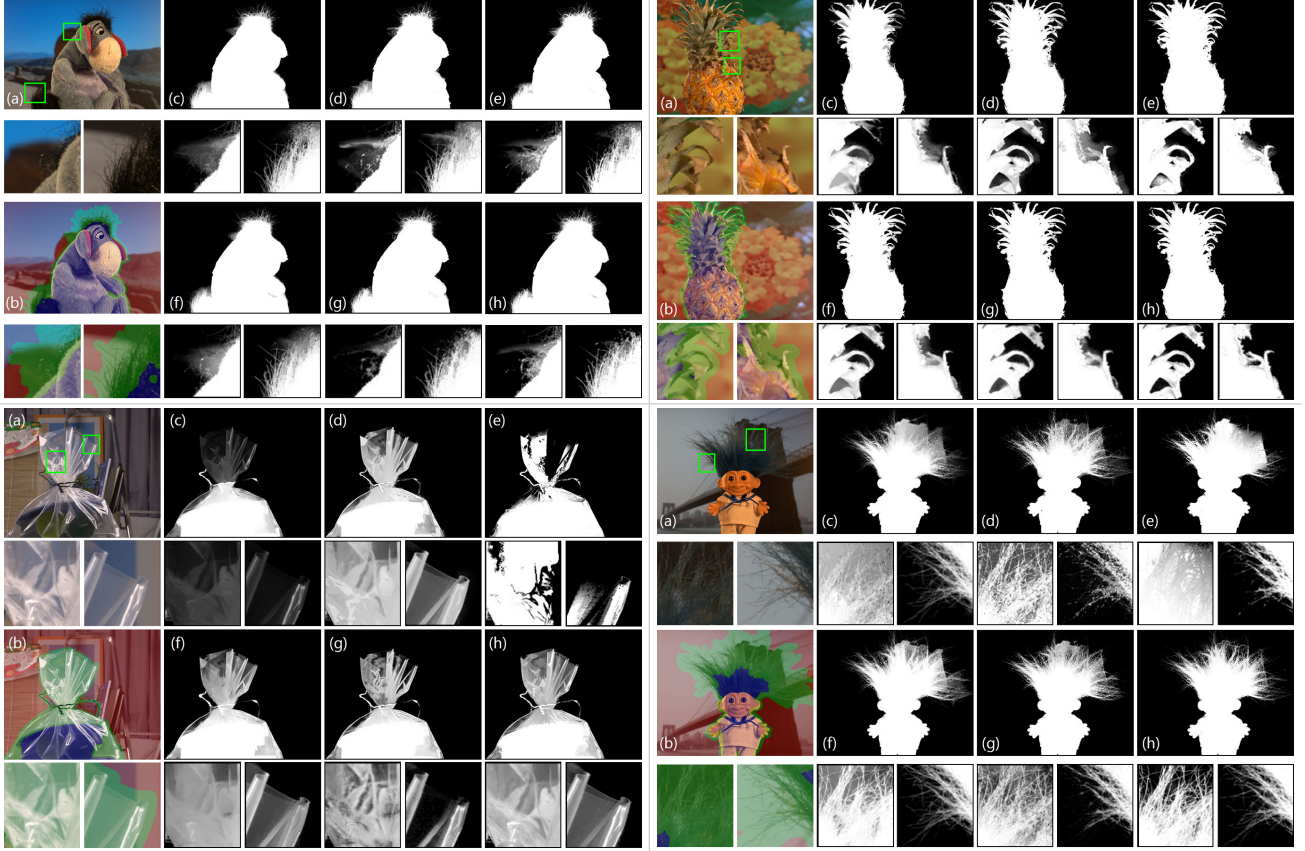


Figure 8: Several examples from the alpha matting benchmark [16] are shown (a) with trimaps overlaid onto the images (b). The mattes are computed by closed-form matting [14] (c), KNN matting [5] (d), manifold-preserving edit propagation [6] (e), LNSP matting [7] (f), DCNN matting [8] (g) and the proposed method (h). See text for discussion.

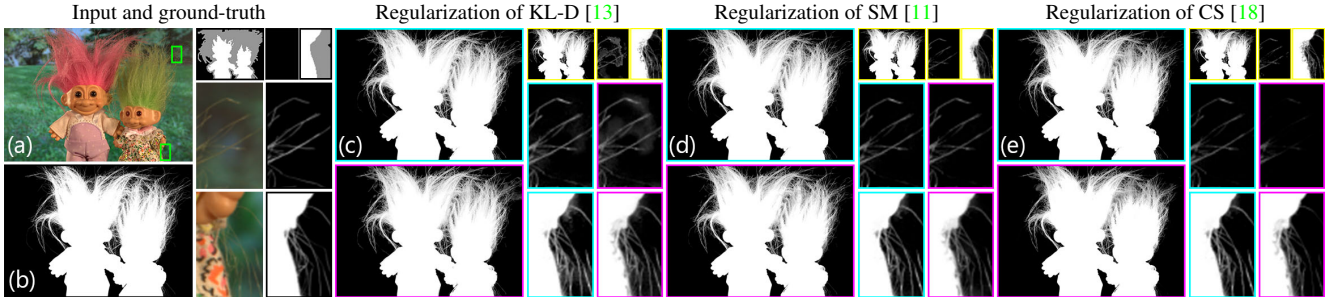


Figure 9: Matte regularization using the proposed method (cyan) or [11] (magenta) for three sampling-based methods (yellow). Our method is able to preserve remote details while producing a clean matte (top inset) and preserve sharpness even around textured areas (bottom).

6.2. Matte regularization

We also compare the proposed post-processing method detailed in Section 4 with the state-of-the-art method by Gastal and Oliveira [11] on the training dataset provided by Rhemann *et al.* [16]. We computed the non-smooth alpha values and confidences using the publicly available source code for comprehensive sampling [18], KL-divergence sampling [13] and shared matting [11]. Table 5

shows the percentage improvement we achieve over Gastal and Oliveira [11] for each algorithm using SAD and MSE as error measures. Figure 9 shows an example for regularizing all three sampling-based methods. As the information coming from alpha values and their confidences found by the sampling-based method is distributed more effectively by the proposed method, the challenging regions such as fine structures or holes detected by the sampling-based method are preserved when our method is used for post-processing.

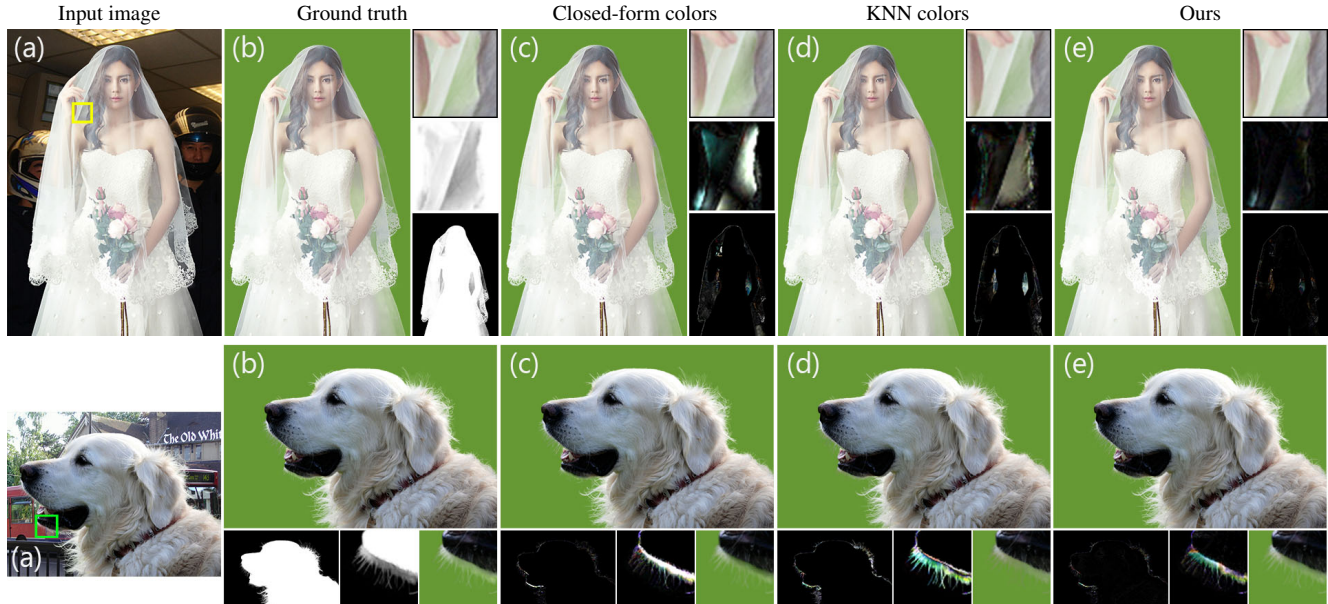


Figure 10: Color estimation results of three algorithms together with the ground truth colors and matte (b). The bottom-right in each set shows per-pixel absolute difference between the estimation and ground truth multiplied by ten. See text for discussion.

6.3. Layer color estimation

We evaluate our layer color estimation method against the closed-form color estimation [14] and KNN colors [5], on the test set of deep image matting [22] using the ground-truth alphas as input. Closed-form colors only use a single local affinity to propagate the colors from the foreground, and this creates artifacts around holes in the foreground (Figure 10, top) or incorrect colors being propagated to nearby regions (bottom). KNN colors, on the other hand, uses only the similarity affinity and it typically generates flat-colored regions, which results in erroneous values especially around hair and fur. Our multi-affinity approach is able to correctly estimate the colors even in the isolated regions or intricate structures. These properties are also reflected in the quantitative comparison, as shown in Table 4.

Table 4: Layer color estimation performance on the test set of DIM [22].

	SAD	MSE
Ours	3.8×10^3	6.9×10^{-4}
CF [14]	4.3×10^3	9.2×10^{-4}
KNN [5]	4.7×10^3	8.4×10^{-4}

6.4. Green-screen keying

Green-screen keying is a more constrained version of the natural image matting problem in which the background is mostly of single color. Despite the more constrained setup, it is challenging to get clean foregrounds for compositing. Aksoy *et al.* [2] show that common natural matting algorithms fail to get satisfactory results despite their perfor-

Table 5: Performance improvement achieved when our matte regularization method replaces [11] in the post-processing steps of 3 sampling-based methods. The training dataset in [16] was used for this experiment.

	Sum of Absolute Differences			Mean Squared Error		
	Overall	S	L	Overall	S	L
KL-D [13]	24.4 %	22.4 %	26.5 %	28.5 %	25.9 %	31.0 %
SM [11]	6.0 %	3.7 %	8.4 %	13.6 %	8.5 %	18.8 %
CS [18]	4.9 %	10.0 %	-0.1 %	18.7 %	25.5 %	11.8 %

mance on the matting benchmark.

We compare the performance of our method to that of the interactive green-screen keying method by Aksoy *et al.* [2] (GSK) and unmixing-based soft color segmentation [3] (SCS) as well as KNN matting [5] and comprehensive sampling [18] in Figure 11. GSK requires local color models, a subset of entries in their color model, and SCS requires a binary map to clean the noise in the background. The matting methods including ours require trimaps and we show results for two trimaps used for comparisons in [2]. We computed the foreground colors for our method and comprehensive sampling using our color estimation method, and KNN colors for KNN matting. We observed that the choice of color estimation method does not change the typical artifacts we see in KNN matting and comprehensive sampling. GSK and SCS compute foreground colors together with the alpha values.

Top example in Figure 11 shows that KNN matting over-



Figure 11: Green-screen keying results of GSK [2] with its input called *local color models* (a) and of SCS [3] with the mask needed for a clean result (b) together with the proposed method (c), comprehensive sampling [18] (d) and KNN matting [5] (e) using two trimaps, one narrow and one wide, for each example. See text for discussion.

estimates alpha values in critical areas and this results in a green halo around the foreground. In contrast, we see a reddish hue in the hair and around the glasses for comprehensive sampling. This is due to the underestimation of alpha values in those areas. The bottom example shows that both competing matting methods fail to get rid of the color spill, *i.e.* indirect illumination from the background. The proposed method successfully extracts the foreground matte and colors in both challenging cases and gives comparable results to the state-of-the-art in green-screen keying. It can also be seen that the effect of different trimaps is minimal in both cases. A successful matting approach requires less input than GSK (the local color models are conceptually similar to a multi-channel trimap and requires more time to generate than a trimap) and is robust against color spill unlike SCS, which makes our method a viable option for green-screen keying.

Although the images shown in Figure 11 have the resolution of 1080p, the average time our matte estimation was around 20 seconds, which is lower than our average for the matting benchmark. The reason is that the time required to construct and solve our linear system mostly depends on the number of unknown pixels in the image, rather than the image resolution. Hence, in a professional production setting where the unknown-opacity regions are typically narrower than the academic benchmarks, our algorithm has lower computational requirements.

7. Spectral analysis

The spectral clusters formed by Laplacians of affinity matrices can be effectively used to reveal characteristics of the constructed graph structure. For instance, Levin *et al.* [14] analyze the matting affinity by looking at eigenvectors corresponding to the smallest eigenvalues of the matting Laplacian. Spectral matting [15] uses the eigenvectors together with a sparsity prior to create a set of soft segments, or *alpha components*, that represent compact clusters of eigenvectors and add up to one for each pixel. The alpha components provide a more distilled and clear visualization to analyze the affinity matrix. In this section, we use the matting components computed using different subsets of information flows we defined for matte estimation to reveal the contribution of different flows at a higher level.

We compute the alpha components shown in Figure 12 using the public source code by Levin *et al.* [15]. We exclude the \mathcal{K} -to- \mathcal{U} flow, which is only defined for the unknown regions as it requires explicitly defined known regions. The resulting Laplacian matrix does not give meaningful spectral clustering because of the pixels with missing connections. We overcome this issue for intra- \mathcal{U} flow by defining it for the entire image instead of only the unknown region. In our matting formulation, we use the color-mixture flow to create the main source of information flow between close-by similarly-colored pixels. This approach creates densely connected graphs as both spatial and color distances are well accounted for in the neighborhood selec-

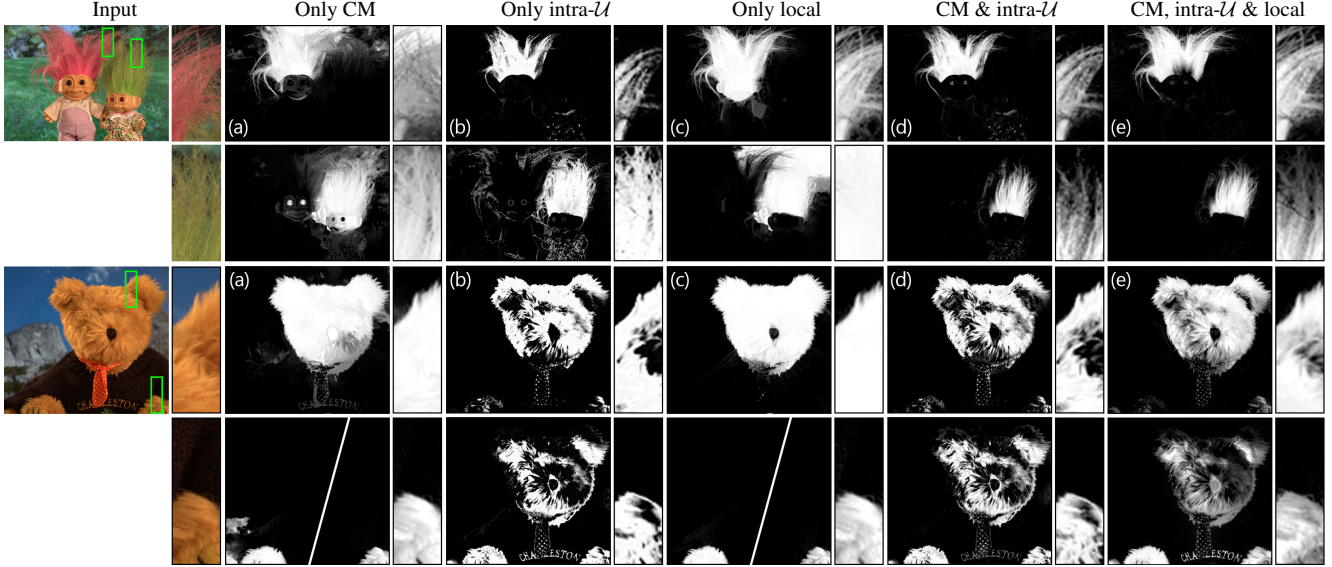


Figure 12: Selected matting components [15] computed from Laplacian matrices constructed using different subsets of information flow. Two components are included in the bottom examples for *only CM* and *only local* cases as the included parts appeared in separate components.

tion. We observed that spectral matting may fail to create as many components as requested (10 in our experiments) in some images, as many regions are heavily interconnected. Using the weighted average of neighboring colors for the flow creates soft transitions between regions.

The intra- \mathcal{U} flow connects pixels that have similar colors, with very little emphasis on the spatial distance. This creates a color-based segmentation of the pixels, but as we compute the weights based on the feature distances, it is not typically able to create soft transitions between regions. Rather, it creates components with alpha values at zero or one, or flat alpha regions with alpha values near 0.5.

The local information flow, used as the only form of flow in the original spectral matting, creates locally connected components with soft transitions.

We observed a harmonious combination of positive aspects of these affinity matrices as they are put together to create our graph structure. This provides a neat confirmation of our findings in the evaluation of our algorithm. We analyze the characteristics of each flow more in detail through visual examples in the remainder of this section.

The top example in Figure 12 shows an input image with the matting components that include the green and the pink hair. Color-mixture affinities give components that demonstrate the color similarity and soft transitions, but they typically bleed out of the confined regions of specific colors due to the densely connected nature of the graph formed by corresponding neighborhoods. We clearly see the emphasis on color similarity for intra- \mathcal{U} flow. While the color clusters are apparent, one can easily observe that unrelated pixels get mixed into the clusters especially around transition

regions between other colors. We see a significant improvement already when these two flows are combined. When the local information flow is added, which gives spatially confined clusters of many colors when used individually, we see smooth clusters of homogeneous colors. The intricate transitions that were missed in the lack of the local flow are successfully captured when all three flows are included in the Laplacian definition.

The spatial connectivity versus color similarity characteristics are even more clearly observable in the bottom example of Figure 12. We see that bright and dark brown of the fur is clearly separated by intra- \mathcal{U} flow in this example. In contrast, color-mixture and local flows separate the fur into three spatial clusters and the sweater into two separate clusters despite the uniform color. The combination, however, is able to successfully separate the dark and bright brown of the fur with smooth transitions.

The full Laplacian matrix we propose in this work blends the nonlocality of colors and spatial smoothness naturally. This is the key characteristic of the proposed matting method. When combined with \mathcal{K} -to- \mathcal{U} flow which addresses remote regions and holes inside the foreground, the proposed algorithm is able to achieve high performance in a variety of images as analyzed in Section 6.

8. Sampling-based methods and \mathcal{K} -to- \mathcal{U} flow

The \mathcal{K} -to- \mathcal{U} flow introduced in Section 3.2 connects every pixel in the unknown region directly to several pixels in both foreground and background. While the amount of flow from each neighbor is individually defined by the computed

Table 6: SAD scores of top sampling-based methods on the matting benchmark against the \mathcal{K} -to- \mathcal{U} flow as a sampling based method, regularized by [11]. Blue shows the best performance among the methods listed here for each image-trimap pair. Red marks the failure cases for the \mathcal{K} -to- \mathcal{U} flow.

	Troll			Doll			Donkey			Elephant			Plant			Pineapple			Plastic bag			Net		
	<i>S</i>	<i>L</i>	<i>U</i>	<i>S</i>	<i>L</i>	<i>U</i>	<i>S</i>	<i>L</i>	<i>U</i>	<i>S</i>	<i>L</i>	<i>U</i>	<i>S</i>	<i>L</i>	<i>U</i>	<i>S</i>	<i>L</i>	<i>U</i>	<i>S</i>	<i>L</i>	<i>U</i>	<i>S</i>	<i>L</i>	<i>U</i>
CSC [10]	13.6	15.6	14.5	6.2	7.5	8.1	4.6	4.8	4.2	1.8	2.7	2.5	5.5	7.3	9.7	4.6	7.6	6.9	23.7	23.0	21.0	26.3	27.2	25.2
Sparse coding [12]	12.6	20.5	14.8	5.7	7.3	6.4	4.5	5.3	3.7	1.4	3.3	2.3	6.3	7.9	11.1	4.2	8.3	6.4	28.7	31.3	27.1	23.6	25.1	27.3
KL-Div [13]	11.6	17.5	14.7	5.6	8.5	8.0	4.9	5.3	3.7	1.5	3.5	2.1	5.8	8.3	14.1	5.6	9.3	8.0	24.6	27.7	28.9	20.7	22.7	23.9
\mathcal{K} -to- \mathcal{U} inf. flow	12.0	13.1	14.6	7.5	9.1	8.9	3.9	4.3	3.8	1.4	2.0	2.0	5.3	5.9	8.0	2.7	3.6	3.3	37.2	39.1	35.8	47.2	56.0	41.9
Comp. Samp. [18]	11.2	18.5	14.8	6.5	9.5	8.9	4.5	4.9	4.1	1.7	3.1	2.3	5.4	9.8	13.4	5.5	11.5	7.4	23.9	22.0	22.8	23.8	28.0	28.1

color-mixture weights, we simplify the formulation and increase the sparsity of our linear system using some algebraic manipulations. These manipulations, in the end, give us the weights $w_p^{\mathcal{F}}$ that go into the final energy formulation.

These weights, which show the connection of the unknown pixel to the foreground, are essentially an early estimation of the matte. This estimation is done by individually selecting a set of neighbors for each pixel and computing an alpha based on the neighbor colors. While our approach is fundamentally defining affinities, it has parallels with sampling-based approaches in natural matting [18, 13, 10, 12], which also select samples from foreground and background and estimates alpha values based on sample colors. We compute confidence values for $w_p^{\mathcal{F}}$ that depends on the similarity of colors of neighbors from the foreground and background. Sampling-based approaches also define confidence values for their initial estimation, typically defined by the *compositing error*, $\|c - (\alpha f - (1 - \alpha)b)\|^2$.

Conceptually, there are several fundamental differences between our computation of \mathcal{K} -to- \mathcal{U} flow and common strategy followed by sampling-based methods. The major difference is how the samples are collected. Sampling-based methods first determine a set of samples collected from known-alpha regions and do a selection for unknown pixels from this predetermined set using a set of heuristics. We, on the other hand, select neighbors for each unknown pixel individually via a k nearest neighbors search in the whole known region. Using the samples, state-of-the-art methods typically use the compositing equation to estimate the alpha value from only one sample pair (a notable exception is CSC matting [10]), while we use 14 samples in total to estimate the alpha by solving the overconstrained system using the method by Roweis and Saul [17]. These differences also change the computation time. \mathcal{K} -to- \mathcal{U} flow can be computed in several seconds, while sampling-based algorithms typically take several minutes per image due to sampling and sample pair selection steps.

In order to compare the performance of \mathcal{K} -to- \mathcal{U} flow as a sampling-based method in a neutral setting, in this experiment, we post-process $w_p^{\mathcal{F}}$ and our confidence values

using the common regularization step [11] utilized by top-performing sampling-based methods in the benchmark. The quantitative results can be seen in Table 6.

As discussed in Section 3.2, \mathcal{K} -to- \mathcal{U} flow fails in the case of a highly-transparent matte (net and plastic bag examples). This is due to the failure to find representative neighbors using the k nearest neighbor search. Sampling-based methods are more successful in these cases due to their use of compositing error in the sample selection. However, in the other examples, \mathcal{K} -to- \mathcal{U} flow appears as the top-performing method among the sampling-based methods in 12 of 18 image-trimap pairs and gives comparable errors in the rest.

The performance of our affinity-inspired approach against the state-of-the-art [18, 13, 10, 12] gives us some pointers for a next-generation sampling-based matting method. While one can argue that the sampling algorithms have reached enough sophistication, selection of a single pair of samples for each unknown pixel seems to be a limiting factor. Methods that address the successful and efficient selection of *many* samples for each unknown pixel will be more likely to surpass state-of-the-art performance. Furthermore, determining the alpha values using more robust weight estimation formulations such as (1) instead of the more simple compositing equation (23) will likely improve the result quality.

9. Limitations

As discussed in corresponding sections, the \mathcal{K} -to- \mathcal{U} flow does not perform well in the case of highly-transparent mattes. We solve this issue via a simple classifier to detect highly-transparent mattes before alpha estimation. However, this does not solve the issue for foreground images that partially have transparent regions. For such cases, a locally changing set of parameters could be the solution.

The proposed matte estimation algorithm assumes dense trimaps as input. In the case of sparse trimaps, generally referred as scribble input, our method may fail to achieve its original performance, as seen in Figure 13. This performance drop is mainly due to the \mathcal{K} -to- \mathcal{U} flow, which fails to find good neighbors in limited known regions, and intra-



Figure 13: Our method fails gracefully in the case of sparse trimaps.

\mathcal{U} flow which propagates alpha information based solely on color to spatially far away pixels inside the unknown region.

10. Conclusion

We proposed a purely affinity-based natural image matting method. We introduced color-mixture flow, a specifically tailored form of LLE weights for natural image matting. By carefully designing flow of information from the known region to the unknown region, as well as distributing the information inside the unknown region, we addressed several challenges that are common in natural matting. We showed that the linear system we formulate outperforms the state-of-the-art in the alpha matting benchmark. The characteristic contributions of each form of information flow were discussed through spectral analysis. We extended our formulation to matte regularization and layer color estimation and demonstrate their performance improvements over the state-of-the-art. We demonstrated that the proposed matting and color estimation methods achieve state-of-the-art performance in green-screen keying. We also commented on several shortcomings of the state-of-the-art sampling-based methods by comparing them to our known-to-unknown information flow.

References

- [1] Y. Aksoy, T. O. Aydin, and M. Pollefeys. Designing effective inter-pixel information flow for natural image matting. In *Proc. CVPR*, 2017. 2
- [2] Y. Aksoy, T. O. Aydin, M. Pollefeys, and A. Smolić. Interactive high-quality green-screen keying via color unmixing. *ACM Trans. Graph.*, 35(5):152:1–152:12, 2016. 3, 8, 12, 13
- [3] Y. Aksoy, T. O. Aydin, A. Smolić, and M. Pollefeys. Unmixing-based soft color segmentation for image manipulation. *ACM Trans. Graph.*, 36(2):19:1–19:19, 2017. 12, 13
- [4] R. Barrett, M. Berry, T. Chan, J. Demmel, J. Donato, J. Dongarra, V. Eijkhout, R. Pozo, C. Romine, and H. van der Vorst. *Templates for the Solution of Linear Systems: Building Blocks for Iterative Methods*. SIAM, 1994. 7, 9
- [5] Q. Chen, D. Li, and C.-K. Tang. KNN matting. *IEEE Trans. Pattern Anal. Mach. Intell.*, 35(9):2175–2188, 2013. 2, 3, 6, 8, 10, 11, 12, 13
- [6] X. Chen, D. Zou, Q. Zhao, and P. Tan. Manifold preserving edit propagation. *ACM Trans. Graph.*, 31(6):132:1–132:7, 2012. 2, 3, 4, 10, 11
- [7] X. Chen, D. Zou, S. Zhou, Q. Zhao, and P. Tan. Image matting with local and nonlocal smooth priors. In *Proc. CVPR*, 2013. 3, 4, 10, 11
- [8] D. Cho, Y.-W. Tai, and I. S. Kweon. Natural image matting using deep convolutional neural networks. In *Proc. ECCV*, 2016. 3, 10, 11
- [9] H. Farid and E. P. Simoncelli. Differentiation of discrete multidimensional signals. *IEEE Trans. Image Process.*, 13(4):496–508, 2004. 8
- [10] X. Feng, X. Liang, and Z. Zhang. A cluster sampling method for image matting via sparse coding. In *Proc. ECCV*, 2016. 2, 5, 7, 10, 15
- [11] E. S. L. Gastal and M. M. Oliveira. Shared sampling for real-time alpha matting. *Comput. Graph. Forum*, 29(2):575–584, 2010. 2, 3, 7, 11, 12, 15
- [12] J. Johnson, E. S. Varnousfaderani, H. Cholakal, and D. Rajan. Sparse coding for alpha matting. *IEEE Trans. Image Process.*, 25(7):3032–3043, 2016. 15
- [13] L. Karacan, A. Erdem, and E. Erdem. Image matting with KL-divergence based sparse sampling. In *Proc. ICCV*, 2015. 2, 5, 7, 11, 12, 15
- [14] A. Levin, D. Lischinski, and Y. Weiss. A closed-form solution to natural image matting. *IEEE Trans. Pattern Anal. Mach. Intell.*, 30(2):228–242, 2008. 2, 3, 6, 7, 8, 10, 11, 12, 13
- [15] A. Levin, A. Rav-Acha, and D. Lischinski. Spectral matting. *IEEE Trans. Pattern Anal. Mach. Intell.*, 30(10):1699–1712, 2008. 2, 13, 14
- [16] C. Rhemann, C. Rother, J. Wang, M. Gelautz, P. Kohli, and P. Rott. A perceptually motivated online benchmark for image matting. In *Proc. CVPR*, 2009. 3, 10, 11, 12
- [17] S. T. Roweis and L. K. Saul. Nonlinear dimensionality reduction by locally linear embedding. *Science*, 290(5500):2323–2326, 2000. 2, 4, 15
- [18] E. Shahrian, D. Rajan, B. Price, and S. Cohen. Improving image matting using comprehensive sampling sets. In *Proc. CVPR*, 2013. 2, 3, 5, 6, 7, 11, 12, 13, 15
- [19] X. Shen, X. Tao, H. Gao, C. Zhou, and J. Jia. Deep automatic portrait matting. In *Proc. ECCV*, 2016. 3
- [20] A. R. Smith and J. F. Blinn. Blue screen matting. *ACM Trans. Graph.*, pages 259–268, 1996. 3
- [21] J. Wang and M. F. Cohen. Optimized color sampling for robust matting. In *Proc. CVPR*, 2007. 3
- [22] N. Xu, B. Price, S. Cohen, and T. Huang. Deep image matting. In *Proc. CVPR*, 2017. 3, 10, 12
- [23] Y. Zheng and C. Kambhamettu. Learning based digital matting. In *Proc. ICCV*, 2009. 3
- [24] Q. Zhu, L. Shao, X. Li, and L. Wang. Targeting accurate object extraction from an image: A comprehensive study of natural image matting. *IEEE Trans. Neural Netw. Learn. Syst.*, 26(2):185–207, 2015. 2