

# Convolutional Neural Opacity Radiance Fields

Haimin Luo, Anpei Chen, Qixuan Zhang, Bai Pang, Minye Wu, Lan Xu, and Jingyi Yu, *Fellow, IEEE*

**Abstract**—Photo-realistic modeling and rendering of fuzzy objects with complex opacity are critical for numerous immersive VR/AR applications, but it suffers from strong view-dependent brightness, color. In this paper, we propose a novel scheme to generate opacity radiance fields with a convolutional neural renderer for fuzzy objects, which is the first to combine both explicit opacity supervision and convolutional mechanism into the neural radiance field framework so as to enable high-quality appearance and global consistent alpha mattes generation in arbitrary novel views. More specifically, we propose an efficient sampling strategy along with both the camera rays and image plane, which enables efficient radiance field sampling and learning in a patch-wise manner, as well as a novel volumetric feature integration scheme that generates per-patch hybrid feature embeddings to reconstruct the view-consistent fine-detailed appearance and opacity output. We further adopt a patch-wise adversarial training scheme to preserve both high-frequency appearance and opacity details in a self-supervised framework. We also introduce an effective multi-view image capture system to capture high-quality color and alpha maps for challenging fuzzy objects. Extensive experiments on existing and our new challenging fuzzy object dataset demonstrate that our method achieves photo-realistic, globally consistent, and fine detailed appearance and opacity free-viewpoint rendering for various fuzzy objects.

**Index Terms**—Computational Photography, Neural Rendering, Opacity Modelling, View Synthesis.



## 1 INTRODUCTION

THE past ten years have witnessed a rapid development of 3D reconstruction technologies for complex scenes with the popularity of commercial passive and active image sensors, which enables numerous immersive experience and virtual and augmented reality (VR and AR) applications and has recently attracted substantive attention. However, the photo-realistic modeling of fuzzy objects with complex opacity such as hair, fur and feathers remains unsolved, which suffers from strong view-dependent brightness, color changes, leading to difficulties in both geometry and appearance reconstruction.

For high-quality fuzzy object modeling, early solutions [1], [2], [3], [4] require costly capture devices and systems, coded lighting, or even manually effort to achieve high-fidelity hair strand geometry reconstruction, which is difficult to be deployed for daily usage. To avoid the heavy reliance on precise geometry modeling, researchers adopt image-based rendering (IBR) [5], [6], [7] to reconstruct the appearance of furry objects by interpolating new views from the captured ones. Specifically, to handle opacity objects, the traditional approach [6] utilizes multi-view images and alpha mattes to compute the angular opacity maps in novel views, which suffers from severe ghosting effect caused by insufficient view samples and inaccurate geometry proxies.

Moreover, obtaining accurate geometry of furry objects is intractable since fur and hair contain tens of thousands of thin fibers and their mutual occlusions are the fundamental causes of translucency.

Only recently, the neural rendering techniques [8], [9], [10], [11], [12] bring huge potential for photo-realistic novel view synthesis from only images input, with various data representations such as point-clouds [11], [13], voxels [9], [14], meshes [10], [15] or implicit representation [8], [12], [16], [17]. However, the literature on fuzzy object neural rendering remains sparse. The recent approach [11] utilizes a neural point renderer to generate texture maps and alpha mattes in novel views explicitly. However, extracting features from a coarse point cloud leads to insufficient sampling and severe artifacts when zooming in and out. Besides, researchers [12], [18], [19], [20] combine implicit representation with volume rendering to inherently model the density of the continuous 3D space, achieving state-of-the-art appearance rendering results, even for opacity objects. Though these methods achieve view-consistent fuzzy object modeling, their method still lacks fine details in both texture and alpha. Specifically, only using a Multi-Layer Perceptron (MLP) network is adopted to fit the continuous 3D space, leading to uncanny high-frequency appearance and opacity details. Most recently, Positional encoding [12] and Fourier feature mapping [21] schemes enable the MLP to handle richer texture details, but it still fails to recover fuzzy surface which contains extremely high-frequency variation.

In this paper, we attack the above challenges and propose a novel scheme to generate convolutional neural opacity radiance fields for fuzzy objects, which is the first to combine explicit opacity supervision with neural radiance field technique (See Fig. 1 for an overview). Our novel pipeline enables high-quality appearance and global consistent alpha mattes generation in arbitrary novel views for more immersive VR/AR applications.

More specifically, to provide explicit opacity supervi-

- H. Luo, Q. Zhang, B. Pang, L. Xu and J. Yu are with the School of Information Science and Technology, ShanghaiTech University, Shanghai 201210, China. E-mail: {luohm, zhangqx1, pangbai, xulan1, yu-jingyi}@shanghaitech.edu.cn.
- A. Chen and M. Wu are with the School of Information Science and Technology, ShanghaiTech University, Shanghai 201210, China, and the Shanghai Institute of Microsystem and Information Technology, Chinese Academy of Sciences, Shanghai 200031, China, and also with the University of Chinese Academy of Sciences, Beijing 100049, China. E-mail: {chenap, wumy}@shanghaitech.edu.cn.
- J. Yu is with the Shanghai Engineering Research Center of Intelligent Vision and Imaging, School of Information Science and Technology, ShanghaiTech University, Shanghai 201210, China. E-mail: yu-jingyi@shanghaitech.edu.cn.

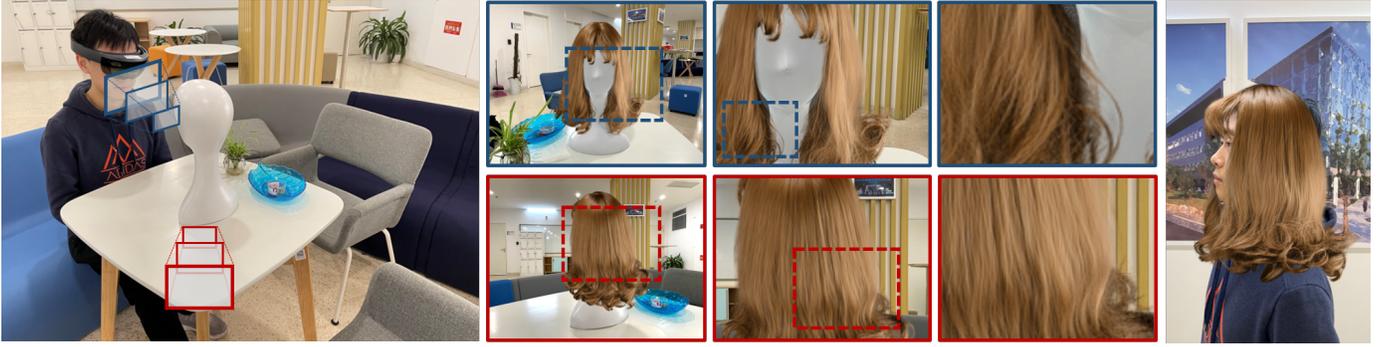


Fig. 1. VR/AR experience with photo-realistic opacity radiance field rendering (first column), our method is able to reconstruct fuzzy appearance at different scales (2 - 5 columns).

sion, from the system side we introduce an effective multi-view capture system equipped with a step turntable and specific transparent lighting design as well as a corresponding multi-color keying algorithm. Our novel system automatically captures both high-quality RGB images and corresponding opacity maps for challenging fuzzy objects in the input capture views. Based on such hybrid input, from the algorithm side, we introduce the convolutional mechanism in the image plane into the neural radiance field framework [19] so that enable photo-realistic appearance and global consistent opacity generation in arbitrary novel views. To this end, we first propose an efficient sampling strategy that utilizes the inherent silhouette prior along the rays and encodes the spatial information across the image plane, which enables efficient radiance field sampling and learning in a patch-wise manner. Then, we perform a novel volumetric integration scheme to generate a per-patch hybrid appearance and opacity feature maps, followed by a light-weight convolutional U-Net to reconstruct the view-consistent fine-detailed appearance and opacity output. Moreover, a patch-wise adversarial training scheme is proposed to preserve both high-frequency appearance and opacity details for photo-realistic rendering in a self-supervised framework. To summarize, our main contributions include:

- We present a novel convolutional neural radiance field generation scheme to reconstruct high frequency and global consistent appearance and opacity of fuzzy objects in novel views, achieving significant superiority to the existing state of the art.
- To enable convolutional mechanism, we propose an efficient sampling strategy, a hybrid feature integration as well as a self-supervised adversarial training scheme for patch-wise radiance field learning.
- We introduce an effective multi-view system to capture the color and alpha maps for challenging fuzzy objects, and our capture dataset will be made available to stimulate further research.

## 2 RELATED WORK

### 2.1 Neural 3D Shape Modeling

Recent work has made a significant process on 3D object modeling and realism free-viewpoint rendering with level

sets of deep networks that implicitly map spatial locations  $xyz$  to a geometric representation (i.g., distance field [16], occupancy Field [17], [22] etc.). The aforementioned explicit representations require discretization (e.g., in terms of the number of voxels, points or vertices), implicitly models shapes with a continuous function and naturally is able to handle complicated shape topologies. The implicit geometric modeling can be quickly learned from 3D point samples [23], [24], and the trained models can be used to reconstruct shapes from a single image or 3D part. However, these models are limited by their requirement of access to ground truth 3D geometry, typically obtained from synthetic 3D shape datasets such as ShapeNet [25]. Subsequent works relax this requirement by formulating differentiable rendering functions that allow neural implicit shape representations to be optimized using only 2D images [8], [26].

### 2.2 Free-Viewpoint Rendering

Free-viewpoint synthesis methods are generally model input/target images as a collection of rays and essentially aims to recover the plenoptic function [27] from the dense samples. Earlier Image-Based Rendering (IBR) work [28] used two planes ( $uvst$ ) parametrization or 2PP to represent rays and render new rays via a weighted blending of the ray samples, i.e., fusing nearby rays by considering view angle and camera spacial distance. They are able to achieve real-time interpolation but require much memory as they need to cache all rays. Following work [5] bring in proxy geometric to select suitable views and filter occluded rays by cross-projection to the image plane when ray fusion. However, those methods are still limited by the linear blending function, leading to severe ghosting and blurring artifact.

Most recently, seminal researches seek to implicitly represent the radiance field and render novel views with a neural network. Deep Surface Light Fields [29] use an MLP network to fix per-vertex radiance and learns to fill up the missing data across angles and vertices. Deferred neural rendering [10] presents a novel learnable neural texture to model rendering as image translation, which uses a coarse geometry for texture projection and offers flexible content editing. Recent works [14], [19], [30] present a learned representation that encodes the view-dependent appearance of a 3D scene without modeling its geometry explicitly. Another line of research extends the free-viewpoint to animation

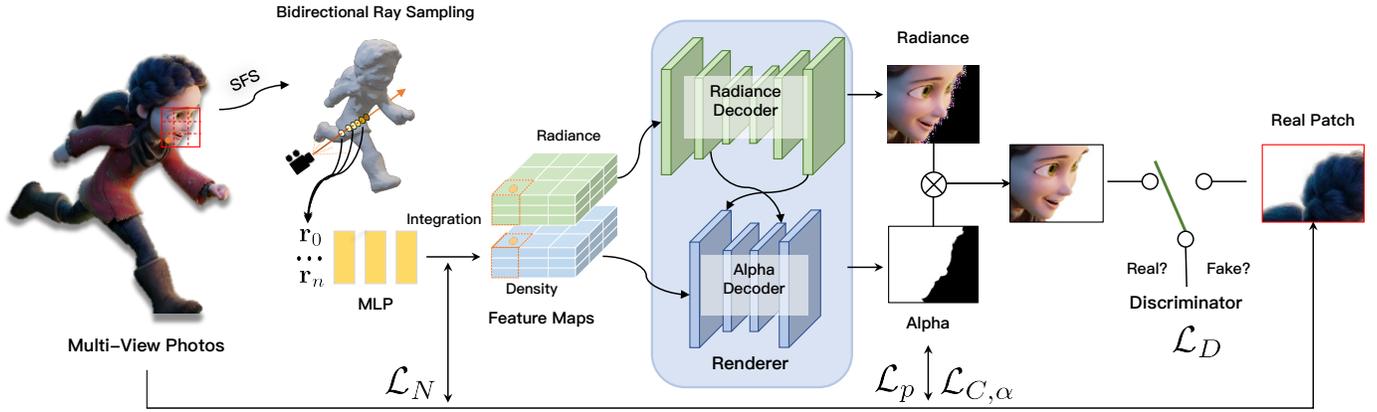


Fig. 2. Overview of our end-to-end ConvNeRF pipeline. Given multi-view RGBA images, we use an SFS to infer proxy geometric for Efficient Ray Sampling. For each sample point in the volume space, the position and direction are feeding to an MLP based feature prediction network to represent the object at a global level. We next concatenate nearby rays into local feature patches and decoded them into RGB and matte with the convolutional volume renderer. An adversarial training strategy is used on the final output to encourage fine surface details. In the reference period, we render the entire image at once rather than per patch rendering.

squeezing or scene relighting by modeling and rendering dynamic scenes through embedding spacial feature with sparse dynamic point cloud [31], using volumetric representation to reconstruct dynamic geometry and appearance variations jointly with only image-level supervision [9], or modeling image formation in terms of environment lighting, object intrinsic attributes and the light transport function [32]. A Notable exception is NeRF [12], which implicitly models the radiance field and the density of a volume with a neural network, then uses a direct volume rendering function to synthesize novel views. They also demonstrate a heretofore unprecedented level of fidelity on a range of challenging scenes. The following work NeRF-W [19] relaxes the NeRF’s strict consistency assumptions through modeling per-image appearance variations such as exposure, lighting, weather, and post-processing with a learned low-dimensional latent space. However, such pure ray-based rendering schemes are failing to recover high-frequency surface, such as fur.

### 2.3 Image Matting

Traditional natural image matting algorithms usually rely on user-defined trimap [33] or scribble [34] as additional input and can generally be divided into sampling-based methods and propagation-based methods. Sampling based methods [35], [36], [37], [38], the known foreground and background regions are sampled as candidates to find the best pixel pair by a carefully designed metric for estimating the alpha value of a query pixel in the unknown area. In propagation-based methods, the alpha values in unknown regions are propagated from the known foreground and background regions using a reformulated image matting model according to different affinities between pixels, using various propagation schemes [39], [40], [41], [42], [43].

Recently, learning-based methods have made an impressive process. Deep Image Matting [44] builds a large matting dataset by alpha composition and proposes an end-to-end deep learning framework to automatically compute alpha mattes using RGB image and trimap as input. Following this work, different methods have been proposed. Disentangled

Image Matting [45] adapts the input trimap as well as estimates alpha with a novel multi-task loss. HDMatting [46] estimate the alpha matte of high resolution images patch by patch through a Cross-Patch Contextual module guided by the given trimap. Recent methods not only estimate the alpha matte but the foreground image [47] and the background image [48] such that the popular perceptual loss [47], [49] or a fusion mechanism [48] can be adopted. To achieve trimap-free matting, recent works implicitly generate trimap using a late fusion model with a soft segmentation network [50] or require additional background photos as input [51] for self-supervised training.

Rather than recover alpha from a 2D image, we set out to generate free-viewpoint alpha matte directly following 3D constraints.

## 3 ALGORITHM DETAILS

In this section, we introduce the design of our convolutional neural opacity radiance fields (ConvNeRF in short), which enables photo-realistic global-consistent appearance and opacity rendering in novel views based on the RGBA input from our capture system, as illustrated in Fig. 2.

Our key insight is to encode opacity information explicitly with a spatial convolutional mechanism to enhance the neural radiance field approach NeRF [12] for high-frequency detail modeling. Inspired by NeRF, we adopt the similar implicit neural radiance field to represent a scene using a Multi-Layer Perceptron (MLP), as well as the volumetric integration of predicted density and color values along the casting rays. Please refer to NeRF [12] for more details.

Differently, our ConvNeRF further encodes opacity explicitly with spatial convolutional design to significantly improve the neural radiance field reconstruction. To this end, we first propose an efficient sampling strategy to not only utilize the inherent silhouette prior along the camera rays but also encode the spatial information across the image plane (Sec. 3.1). Then, a global geometric representor is adopted to map a 3D location to a high-level radiance feature and then a novel volumetric integration scheme is adopted to generate per-patch hybrid feature embeddings,

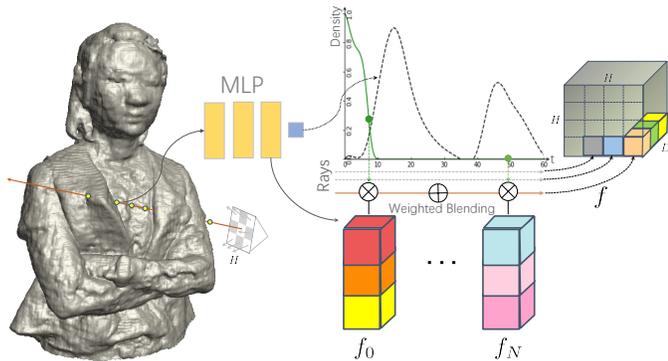


Fig. 3. Feature patches construction. To enable 2D convolution, we weighted blend the volume features of each ray as the ray feature with the “cumpruned” density of the sample points and concatenate nearby ray features into local patches for the following convolutional rendering.

which models the characteristics of appearance and opacity separately for more effective radiance field learning in a patch-wise manner (Sec. 3.2). Next, we choose a light-weight U-Net to decode the feature patches into view-consistent appearance and opacity output (Sec. 3.3). We further adopt a patch-wise adversarial training scheme to preserve both high-frequency appearance and opacity details in a self-supervised framework (Sec.3.4).

### 3.1 Efficient Ray Sampling

To speed up the training procedure, we introduce a patch-wise sampling scheme, which uses a coarse proxy to filter out redundant samples for more efficient radiance field sampling and rendering.

**Coarse Proxy Generation.** Note that the multi-view alpha mattes encodes the silhouette prior of the captured fuzzy object inherently, which provides a reliable proxy to guide the sampling process in the continuous 3D space. To this end, the input alpha mattes are binarized and dilated, and the Shape-from-Silhouette (SfS) [52] algorithm is applied to obtain a coarse 3D proxy of the fuzzy object.

**Spatial Sampling.** To further enable patch-wise spatial sampling across the image plane, we first render two depth maps: the near-depth ( the first hit point) and far-depth (the last hit point) for all training views by projecting the proxy mesh to the image plane. Then, we uniformly divide the input images and depth maps into  $K \times K$  small patches without spacing and filter out the patches outside the valid near/far-depth (i.e., pure background patches).

**Ray-wise Sampling.** Given the above sampled ray patches, we perform a similar two-stage sampling scheme of the original NeRF [12] in a coarse-to-fine manner for ray-wise radiance fields sampling to sample  $N$  sampling points for each ray of the ray patches. Differently, note that we have obtained near-far depth priors for sampled rays, we perform the ray-wise sampling on the near-far region only in both coarse and fine stage. Thus it requires only 1/8 sample points of NeRF at least and achieves much more efficient radiance field training and rendering.

The sampled patches with tightly sampling points in valid regions are then utilized for the following training process of the neural radiance field in a patch-wise manner. In

practice, to balance the memory requirement and sampling effectiveness, during training  $K$  is set to be 32 and we use 12 patches as a batch with 64 samples per ray. It achieves about 4 times faster training and about 10 times faster inference compared to the original NeRF [12] in general.

### 3.2 Patch-wise integration

Unlike NeRF predicting pixel color by integrating in a low-dimensional color space, we instead perform a novel volumetric integration scheme to generate a per-patch hybrid appearance and opacity feature maps to enable more effective radiance field learning in a patch-wise manner.

In our ConvNeRF, for a given 3D location  $\mathbf{x}$  and viewing direction  $\mathbf{d}$ , we first extract its 3D features  $(\mathbf{f}, \sigma)$  with a global geometric representor  $\mathbf{E}_{\Theta}$ :

$$(\mathbf{f}, \sigma) = \mathbf{E}_{\Theta}(\mathbf{x}, \mathbf{d}). \quad (1)$$

Same as the original NeRF network architecture,  $\mathbf{E}_{\Theta}$  is achieved with an MLP block with optimizable weights  $\Theta$ , we remove the last layer of RGB color branch so that  $\mathbf{E}_{\Theta}$  plays as a global feature extractor and generates both radiance term  $\mathbf{f}$  and the geometric term  $\sigma$  output.

Then for each ray of the patches, we sample  $N$  points and extract their 3D features to predict the foreground and background probability (i.e., the alpha):

$$\alpha_i = T_i(1 - \exp(-\sigma_i \delta_i)), T_i = \exp(-\sum_{j=1}^{i-1} \sigma_j \delta_j), \quad (2)$$

where  $\delta$  is the distance between adjacent sampling points. The projected feature for each patch pixel is obtained by an integration scheme along the ray:

$$F_c = \sum_{i=1}^N \alpha_i \mathbf{f}_i, \quad (3)$$

As mentioned above, we decompose the free-viewpoint object rendering into radiance and alpha components; thus, we expect to decouple those two components when rendering. To this end, different from radiance feature extraction which fuses the radiance term of each sampling point along a ray by integration, we set out to encode the underlying depth prior provided by our efficient ray sampling scheme by concatenating  $\alpha$  of each sampling point estimated by Eqn.2 together, i.e.,  $F_d = [\alpha_1, \alpha_1, \dots, \alpha_N]$ .

After that, we can obtain view-dependent radiance feature map  $\mathcal{F}_c$  and density feature map  $\mathcal{F}_d$  for each patch for the following convolutional rendering. Our novel patch-wise volumetric integration scheme provides high-level appearance and opacity features that enable high-quality RGB and alpha image synthesis.

### 3.3 Convolutional Volume Renderer

Based on the above per-patch features which encode the characteristic of appearance and opacity, we introduce a convolutional volume renderer scheme for view-consistent appearance and alpha matte rendering by utilizing the spatial information, to address the issue that the original NeRF [19] fails to recover the high-frequency details of fuzzy

objects with a simple MLP. Specifically, let  $\mathbf{G}$  denote the network of our convolutional volume renderer with parameters  $\theta$ , which predicts the texture image  $\mathbf{F}$  with corresponding opacity map  $\alpha$  via  $(\mathbf{F}, \alpha) = \mathbf{G}(\mathcal{F}_c, \mathcal{F}_d; \theta)$ , as illustrated in Fig. 2. Different with NeRF which renders radiance pre-multiplied by alpha,  $\mathbf{F}$  is just foreground image to avoid retaining radiance from the background near boundary.

Note that we aim to achieve view consistent rendering, however, the local receptive field of convolutional neural network (CNN) itself may introduce view-inconsistency. Thus  $\mathbf{G}$  is designed to be light-weight to replace the last fully connect layer of the original NeRF network and decode the feature maps into radiance and opacity. In this way, the global implicit representation described in Sec. 3.2 plays a dominant role for view-consistency.

As illustrated in Fig. 2, our volume renderer  $\mathbf{G}$  consists of a radiance branch and an opacity branch with similar U-Net architectures. Note that in our radiance branch two downsample-upsample blocks are adopted to generate fine-detailed texture output, while only a single downsample-upsample block is utilized in our opacity branch. Such a sophisticated design is based on our observation that alpha mattes are sensitive to low-level features such as image gradients and are more suitable for a shallow network to preserve view-consistent output.

Besides, since the input feature maps  $\mathcal{F}_c$  and  $\mathcal{F}_d$  suffers from incomplete integration near the boundary regions, which is critical for fuzzy objects, we further adopt the gated convolution [53] in the U-Net architectures for both branches, so as to enhance the denoising and image completion capabilities of the network. For more detailed alpha matte prediction, the texture output from the radiance branch is further concatenated with density feature map  $\mathcal{F}_d$  to form a hybrid input of the opacity branch. Also, we accumulate the  $\mathcal{F}_d$  per-pixel first to form a coarse initial alpha matte, and the opacity branch predicts the opacity residuals. Finally, the coarse matte and the residual one are added to generate the final detailed alpha matte output.

To better encourage surface details (hairline, hair texture) and sharpness, we propose to apply a GAN-based discriminator loss to the final image patch (as shown in the right side of Fig. 2). To enable self-supervised adversarial learning, we randomly sample patches from all the input multi-view images as real samples of the discriminator, and the fake samples are generated using the above volumetric render and the estimated alpha matte, which is formulated as:

$$\mathbf{I}_{\text{fake}} = \alpha \mathbf{F} + (1 - \alpha) \mathbf{B}, \quad (4)$$

Where  $\mathbf{F}, \mathbf{B}$  denote the generated foreground texture and background image, respectively. We set  $B = [1.0, 1.0, 1.0]$  as a white background in all patches in our experiment.

### 3.4 Network Training

We train the convolutional volume renderer  $\mathbf{G}$  with 3 generation loss and one discriminator loss. Recall that our model targets on high quality free-viewpoint radiance rendering (RGB and its opacity) from given sparse views RGBA images, we decompose the prediction procedure into two layers (the RGB and alpha layers) and final compose them

together via the predicted alpha layer. We first propose to use  $L_2$  loss for the above two layers:

$$\mathcal{L}_{C,\alpha} = \sum_{i=1}^n \|\mathbf{I}_i - \tilde{\mathbf{I}}_i\|_2^2 + \|\alpha_i - \tilde{\alpha}_i\|_2^2 \quad (5)$$

where  $\tilde{\mathbf{I}}$  and  $\tilde{\alpha}$  indicate the ground truth image and alpha.

To encourage fine details and let outputs close to ground truth patches at the high level, We use a VGG19 perceptual loss [49] on the output feature maps of the  $l$ th layer of the VGG19 backbone, which is denoted as  $\phi^l$ , over both the compositional image and alpha map:

$$\mathcal{L}_P = \sum_{l \in \{3,8\}} \sum_{i=1}^n (\|\phi_i^l(\mathbf{I}) - \phi_i^l(\tilde{\mathbf{I}})\|_2^2 + \|\phi_i^l(\alpha) - \phi_i^l(\tilde{\alpha})\|_2^2) \quad (6)$$

Instead of computing losses on the final outputs only, we add an intermediate loss to the MLP output to encourage multi-view consistency by considering the regional patch-based CNN renderer has smaller perceptual fields and can not capture global geometric, while the MLP block taking 3D location and view direction as input and adding constraint on its output can better preserve global consistency:

$$\mathcal{L}_N = \sum_{i=1}^n (w_\alpha \|\alpha_i - \tilde{\alpha}_i\|_2^2 + w_i \|\mathbf{I}_i - \tilde{\mathbf{I}}_i\|_2^2 + \sum_{l \in \{3,8\}} \|\phi_i^l(\mathbf{I}) - \phi_i^l(\tilde{\mathbf{I}})\|_2^2) \quad (7)$$

where  $w_\alpha = 1$  if the scene is synthetic otherwise 0 as the ground truth alpha map of real objects is not view consistent, that is, we only implicitly supervise the alpha with RGB by considering stumpy gt alpha would result in jitter phenomenon on the boundary region. Note that,  $w_i = a - b\alpha_i$  where  $a - b = 1$ , this term penalities the false prediction on empty region. We set  $a = 2, b = 1$  in our experiment.

In summary, our optimizing objective function for generator  $\mathcal{L}_G$  is:

$$\mathcal{L}_G = \mathcal{L}_{C,\alpha} + \mathcal{L}_P + \mathcal{L}_N \quad (8)$$

For the discriminator  $D$ , we minimize the adversarial loss [54]:

$$\mathcal{L}_D = \|D(\mathbf{I})\|_2^2 + \|D(\tilde{\mathbf{I}}) - 1\|_2^2 \quad (9)$$

## 4 CAPTURE SYSTEM

Here, we describe our capture system, which produces high-quality multi-view RGBA (RGB and Alpha) images for explicit opacity modeling of challenging fuzzy objects. As illustrated in Fig. 4, our pipeline is equipped with an easy-to-use capture device and stable calibration and automatic matting methods.

**Device.** Opacity object data could be captured through well-designed acquisition systems proposed in NOPC [11]. However, there are some limitations, including disability in capturing objects without support (e.g., hair without a head model) or capturing without the calibration box. Thus we build a novel capture system as shown in Fig. 4(a), which consists of a precisely controllable step turntable placed in front of the green screen and six Nikon D750 DSLR Cameras facing towards the objects. For wig data, we use a transparent support covered with soft light fabric

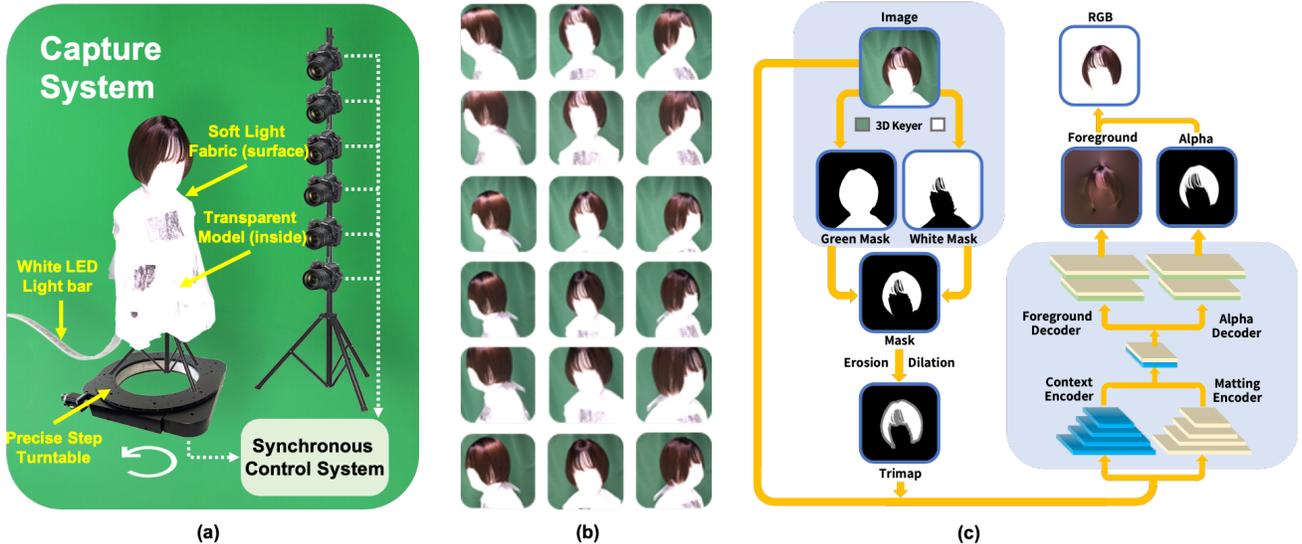


Fig. 4. (a) Our Capture system consists of a lit support, an array of calibrated cameras, a precise step turntable, and a green screen. The fuzzy object is placed on the lit support and the precise turntable is synchronized with the camera array. (b) Samples of images captured among 6 cameras in 3 different steps. (c) The pipeline of alpha matte generation.

to support the object and place them on the turntable. Insert an LED light bar into the support to make the soft light fabric lit. Before capturing, the cameras are adjusted to overexpose the soft fabric part, which will reduce the shadows of the captured image. When capturing images, the turntable moves forward in precise steps. The turntable will step forward 80 steps per lap. All six cameras capture an image per step, which will serve as the training data.

**Calibration.** Unlike general 3D object reconstruction, fuzzy objects are difficult to reconstruct using structure-from-motion(SFM) techniques, which means calibration via the reconstruction process would not work well. Previous work [11] use an auxiliary calibration camera with a pattern box to solve this, which takes much time to run an extra SFM pipeline for every capturing. In our system, the calibration process only needs one time for each camera setting. We know exactly how the camera’s external parameters are transformed from the initial step to each step via the precise step turntable. Let  $A_j$  denote the affine transformation from the initial step to the  $j$ th step. The calibration process is divided into two steps. Firstly, calibrate the intrinsic and extrinsic parameters of six cameras at the initial step via Zhang’s camera calibration [55], denote the  $i$ th camera’s intrinsic parameter as  $K_i$  and extrinsic parameter at its  $j$ th step as  $T_{i,j}$ . Then, calculate the extrinsic parameter of each view under the turntable’s coordinate

$$T_{i,j} = A_j T_{i,0} \quad (10)$$

where  $i$  is the index of cameras, and  $j$  is the number of steps.

**Opacity Decomposition.** Image matting is an ill-posed problem due to a lack of constraint in its formulation. To obtain alpha matte without loss of rich details from fuzzy objects, we apply both 3D keyer and deep learning-based matting algorithm (see Fig. 4(c)). Specifically, the keying benefits from our specially designed capture system. In addition to the traditional green screen, the overexposure of support reduces the shadow cast by the object, providing a clean white screen. We key out both green and white

from the image to extract a preliminary foreground mask. A trimap is then generated from the mask using erosion and dilation operations. We apply Context-aware Matting [47] to predict alpha matte and foreground simultaneously, which combines the benefits of local propagation and global contextual information and performs visually much better on intractable data like curly hair and air bangs.

## 5 EXPERIMENTS

We evaluate our ConvNeRF on various furry objects. Quantitative and qualitative evaluation in Sec. 5.1 show that our method can better preserve high fidelity appearance details than prior work and generate globally consistent alpha mattes in arbitrary novel views. We further perform extensive ablation studies to validate our design choices in Sec. 5.2. We urge the reader to view our supplementary video to better appreciate our method’s significant improvement over baseline methods when rendering novel views.

**Baselines.** We adopt the conventional IBR method *Image-based Opacity Hull (IBOH)* [6] and recent explicit neural rendering method *Neural Opacity Point Cloud (NOPC)* [11] and implicit volume rendering method *Neural Radiance Fields (NeRF)* [12] as baselines for comparisons.

**Training details.** We use half original image resolution for synthetic and real objects, 200 ~ 300 images for training. We use patches of size  $32 \times 32$  and 64 samples (32 in coarse and fine modules) per ray for ConvNeRF training. For the baseline models, we sample 300,000 points for each object from the same proxy mesh with us as the dense point cloud input for the *NOPC*. Also, we use 128 samples (64 in coarse and fine modules) for *NeRF* training. It takes 1 ~ 2 days per object to train our model with a single NVIDIA TITAN RTX GPU, while the *NeRF* takes about a week if trained on the same image number and resolution.

### 5.1 Comparison

We perform quantitative and qualitative comparisons on both RGB and alpha with the above baseline models. For



Fig. 5. Object gallery. Our method can generalize well to various fuzzy objects, including high frequency, view-dependency and translucency appearance, such as hairstyles, clothes, toys and animals etc.

qualitative evaluation, Fig. 6 shows several our novel viewpoint RGB rendering results together with visual comparison v.s. most recent neural scene representation methods. We can see that the *IBOH* suffers from ghosting and aliasing due to the input views are in-uniform sampled. The point cloud based *NOPC* can partially generate smooth texture details with sufficient training images but still suffers from color shifting and blur on the *Girl* object due to interpolation between the features projected from the discrete spacial point cloud. The NeRF can preserve low-frequency components but fails to recover high-frequency details due to the limited representation ability of the MLP network. In the fur region, there is noise caused by insufficient samples on both ray direction and the ray resolutions. In contrast, our ConvNeRF is able to reconstruct fine texture and geometry details and present sharply visual results with global consistency (refer to our supplemental video).

Fig. 7 shows the visual discrepancies of free-viewpoint alpha maps rendering results from given discrete view samples. Compared to other methods, our ConvNeRF is able to preserve the sharpness of the hair boundary. Note that we can easily obtain the perfect alpha map in synthesis scenes while not easy for real data; our method is not sensitive to the quality of input alpha maps and is able to partially recover missing part from nearby view samples (as shown in the first row of Fig. 7).

We quantitatively evaluate our method with PSNR, LPIPS (ALEX backbone) [56] and SSIM metric.<sup>1</sup> As shown

1. We computer quantity value only on foreground region (i.e., the object itself) by considering the background is not included in the final render.

in Tab. 1 and Tab. 2, our ConvNeRF achieves significant improvement on both RGB and alpha results. We replace the LPIPS metric with *Sum of Absolute Distance* (SAD) when evaluating the alpha result due to the domain gap between the alpha-like image and the training set of the network-based distance metric. Tab. 3 shows the average PSNR of all datasets on semi-translucent region (i.e.,  $0 < \alpha < 1$ ). Our method achieves state-of-the-art performance.

## 5.2 Evaluation

To explore each part’s contribution in our pipeline, we conduct ablation studies by choosing a set of controlled experiments, including randomized ray sampling v.s. efficient ray sampling, Pixel-wised v.s. Convolutional renderer, with/without GAN based discriminator loss, the range of the sampling Number and patch size.

**Randomize v.s. Efficient Ray Sampling** We first compare the proposed Efficient Ray Sampling (ERS) scheme described in Sec. 3.1, which include a baseline model (a) (i.e., *NeRF*) with 64 samples per ray same as our ConvNeRF and a model (b) with ERS. As the *Wolf* case shown in Fig. 8, model (a) fails to recover surface details while ERS brings significant improvements to the visual quality and quantitative quality shown by model (a) and (b) in Tab. 4.

**Pixel-wised v.s. Convolutional Volume Renderer** We train a generator model without the discriminator, as (c) denoted in Tab. 4 and Fig. 8, setting (c) can recover better texture details on the wolf’s forehead while (a) and (b) can only preserve low frequency information.

**W/o GAN Loss** We further evaluate the effect of the GAN discriminator adopted for fine details refinement. The sur-

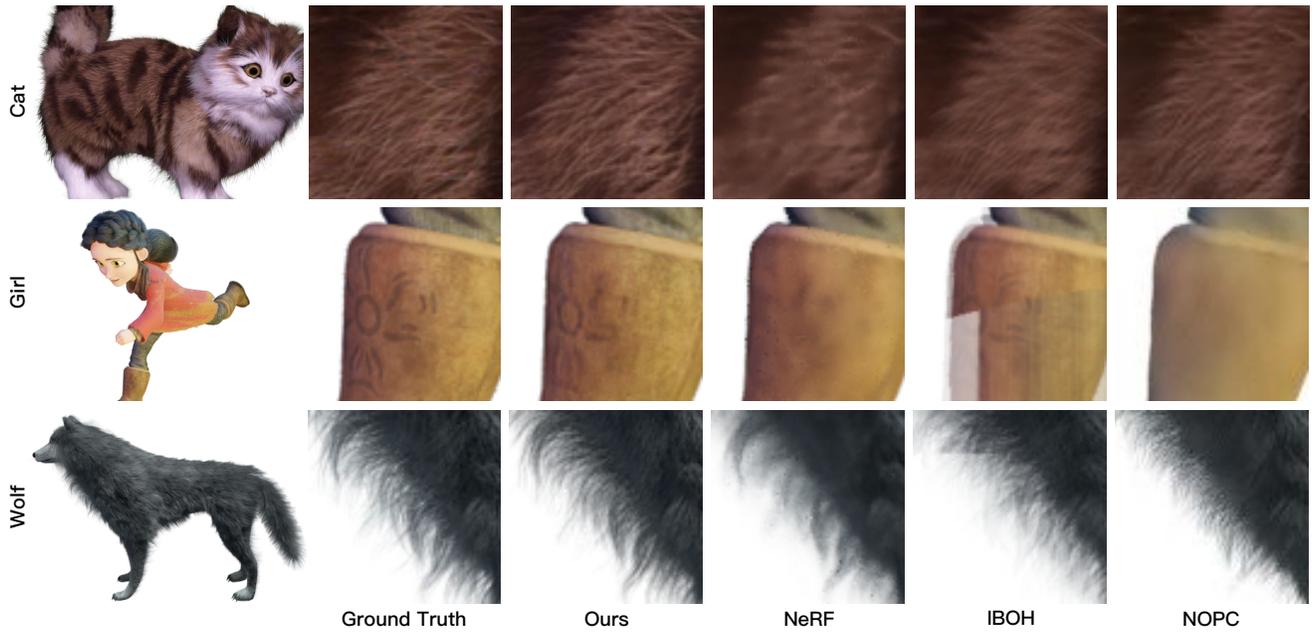


Fig. 6. Free Viewpoint RGB results comparison with IBOH [6], NOPC [11] and NeRF [12] on *Cat*, *Girl*, *Wolf* datasets. Our method is able to reconstruct fine details on both geometry and appearance while keeping global view-consistency, such as *Cat*'s fur texture, the pattern on *Girl*'s boots and the geometric details of wolf's hair. IBOH exhibits ghosting and aliasing, NOPC exhibits excessive blurs and loss of geometric details and NeRF exhibits excessive noise and blurs. See supplementary materials for more results.

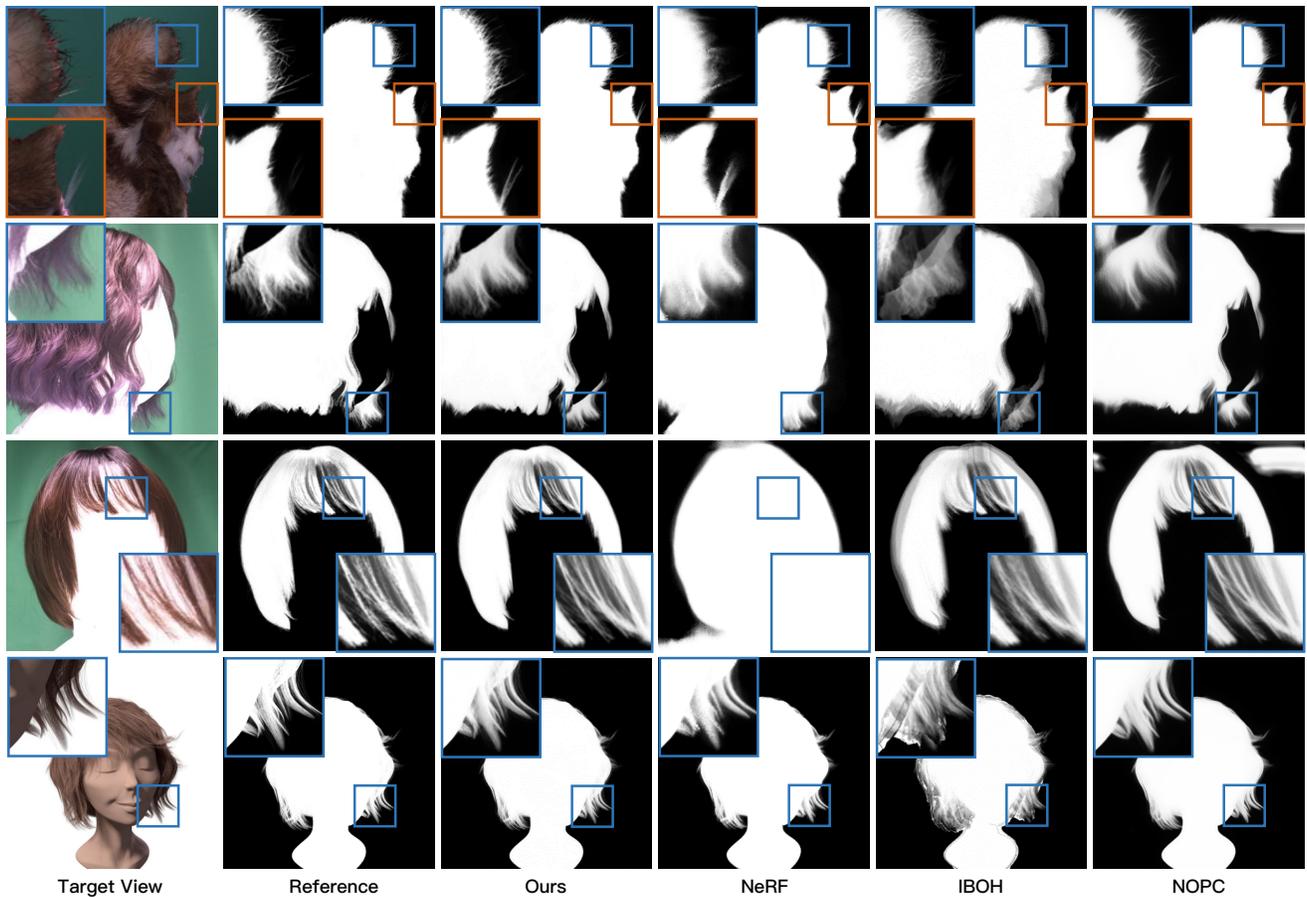


Fig. 7. Free Viewpoint Alpha results comparison with IBOH [6], NOPC [11] and NeRF [12] on *Cat*, *Hairstyle 2* datasets. Our method can recover missing part opacity such as *Cat*'s beard from view-inconsistent alpha mattes as shown in the first row, While IBOH fails and exhibits strong polygonal artifacts. Our method can produce sharper alpha mattes than NOPC which suffers from strong artifacts around the hair. NeRF fails on our challenging *Hairstyle 2* dataset. See supplementary materials for more results.

TABLE 1  
Quantitative comparisons of foreground images on different Objects

Method		IBOH	NOPC	NeRF	Ours
<b>Wolf</b>	PSNR	28.25	29.95	32.56	<b>37.23</b>
	SSIM	0.932	0.945	0.965	<b>0.975</b>
	LPIPS	0.073	0.058	0.061	<b>0.022</b>
<b>Hair</b>	PSNR	23.22	28.31	28.29	<b>33.56</b>
	SSIM	0.872	0.912	0.932	<b>0.952</b>
	LPIPS	0.104	0.058	0.070	<b>0.021</b>
<b>Girl</b>	PSNR	18.32	26.97	29.77	<b>32.17</b>
	SSIM	0.863	0.922	0.951	<b>0.958</b>
	LPIPS	0.150	0.099	0.052	<b>0.035</b>
<b>Dog</b>	PSNR	22.53	27.88	29.59	<b>31.87</b>
	SSIM	0.866	0.902	0.937	<b>0.950</b>
	LPIPS	0.127	0.098	0.053	<b>0.030</b>
<b>Koala</b>	PSNR	24.41	25.45	32.90	<b>33.48</b>
	SSIM	0.871	0.890	0.960	<b>0.965</b>
	LPIPS	0.103	0.117	0.053	<b>0.023</b>
<b>Cat</b>	PSNR	20.40	26.72	24.80	<b>28.95</b>
	SSIM	0.810	0.888	<b>0.896</b>	<b>0.894</b>
	LPIPS	0.198	0.101	0.187	<b>0.090</b>
<b>Hairstyle 1</b>	PSNR	25.32	27.34	17.84	<b>30.51</b>
	SSIM	0.905	0.921	0.904	<b>0.949</b>
	LPIPS	0.067	0.058	0.135	<b>0.037</b>
<b>Hairstyle 2</b>	PSNR	25.86	31.14	19.38	<b>37.21</b>
	SSIM	0.951	0.962	0.931	<b>0.983</b>
	LPIPS	0.104	0.033	0.131	<b>0.014</b>
<b>Hairstyle 3</b>	PSNR	14.26	29.38	22.69	<b>33.65</b>
	SSIM	0.877	0.942	0.944	<b>0.969</b>
	LPIPS	0.130	0.040	0.073	<b>0.016</b>

face details of (c) are a little over-smoothed while (d) present fine and sharp details on the fur texture and hairlines, which is almost close to the ground truth. As for quantitative metrics, (d), setting keeps high performance on foreground texture and alpha quality, as the (c) and (d) shown in Tab. 4.

**Number of Sample Points.** We further demonstrate the effect of different numbers of per ray sample points and image numbers. We evaluate on  $N \in \{16, 32, 48, 64\}$  and image number  $S \in \{50, 100, 150, 200, 250, 300\}$  respectively, as shown in the first and second row of Fig. 9. More sample points can lead to better rendering quality and higher alpha prediction accuracy such as hairlines. Note that our method is not sensitive to training data size and can produce satisfactory results even with a few samples per ray.

**Patch Size.** As we introduce the convolutional mechanism in the image plane, the patch size plays a significant role for synthesis quality. Larger patch size helps preserve high-order features through perceptual loss while leads to fewer patches for training, which reduces data diversity for adversarial training. To evaluate the effect of different patch sizes, we set the patch size  $K \in \{8, 16, 32, 64\}$  and conduct quantitative(Fig. 9) and qualitative(Fig. 10) comparisons. Notice that although the model with  $K = 8$  achieves similar quantitative performance with our proposed model, it exhibits severe visual artifacts (the first row of Fig. 10).

**Textured Background.** As we mask out the background for real data, it may affect the rendering quality of both NeRF and our model. Thus we render the synthetic scene

TABLE 2  
Quantitative comparisons of alpha mattes on different Objects

Method		IBOH	NOPC	NeRF	Ours
<b>Wolf</b>	SAD	18.46	7.977	4.777	<b>2.564</b>
	PSNR	25.10	29.03	32.14	<b>38.63</b>
	SSIM	0.956	0.983	0.990	<b>0.996</b>
<b>Hair</b>	SAD	18.93	6.117	5.078	<b>2.077</b>
	PSNR	21.87	29.54	28.25	<b>36.35</b>
	SSIM	0.930	0.981	0.981	<b>0.994</b>
<b>Girl</b>	SAD	48.61	12.66	5.290	<b>2.877</b>
	PSNR	19.59	27.83	30.58	<b>36.39</b>
	SSIM	0.937	0.986	0.991	<b>0.996</b>
<b>Dog</b>	SAD	30.77	11.32	6.828	<b>3.661</b>
	PSNR	22.77	27.78	30.41	<b>36.52</b>
	SSIM	0.945	0.983	0.988	<b>0.995</b>
<b>Koala</b>	SAD	26.04	71.32	190.15	<b>18.94</b>
	PSNR	22.04	21.76	14.56	<b>30.26</b>
	SSIM	0.939	0.962	0.904	<b>0.983</b>
<b>Cat</b>	SAD	50.85	21.46	<b>20.81</b>	23.79
	PSNR	19.17	<b>25.67</b>	24.25	<b>25.52</b>
	SSIM	0.913	0.958	<b>0.960</b>	0.944
<b>Hairstyle 1</b>	SAD	15.51	9.761	30.66	<b>5.445</b>
	PSNR	25.06	28.96	18.16	<b>32.48</b>
	SSIM	0.961	0.977	0.947	<b>0.987</b>
<b>Hairstyle 2</b>	SAD	7.677	5.948	22.69	<b>2.612</b>
	PSNR	28.80	32.01	19.46	<b>37.54</b>
	SSIM	0.980	0.988	0.957	<b>0.994</b>
<b>Hairstyle 3</b>	SAD	46.09	10.06	12.58	<b>3.037</b>
	PSNR	24.13	26.20	22.77	<b>34.62</b>
	SSIM	0.960	0.978	0.973	<b>0.993</b>

TABLE 3  
Average PSNR on semi-transparent region (U). U+ and U- represent the region after dilation and erosion.

Method	$\alpha F$			Alpha		
	U-	U	U+	U-	U	U+
IBOH	15.78	16.05	16.46	12.93	12.99	13.59
NOPC	16.24	17.49	18.75	14.15	15.14	16.55
NeRF	12.95	13.09	14.83	12.17	12.48	14.29
Ours	<b>25.23</b>	<b>26.21</b>	<b>27.06</b>	<b>20.73</b>	<b>22.45</b>	<b>24.21</b>

*Wolf* together with rich texture backgrounds and parallax to evaluate our approach against NeRF. Tab. 5 indicates that both models perform slightly better in radiance rendering, while using textured background seems to downgrade the alpha rendering quality of both our approach and NeRF. It follows the insight that utilizing explicit alpha supervision is more effective than relying on the network itself to extract such alpha cues implicitly.

### 5.3 Limitation and Discussion

Our model is able to reconstruct high-frequency appearance from given multi-view RGBA images, the combination of our patch-wise learning and the global implicit representation described in Sec. 3 encourage view-consistency, however, we note a bit of flicker in the around-view videos. It is mainly due to the patch-wise adversarial loss in Eqn. 9 which increases high-frequency details significantly with a slight trade-off of view consistency.

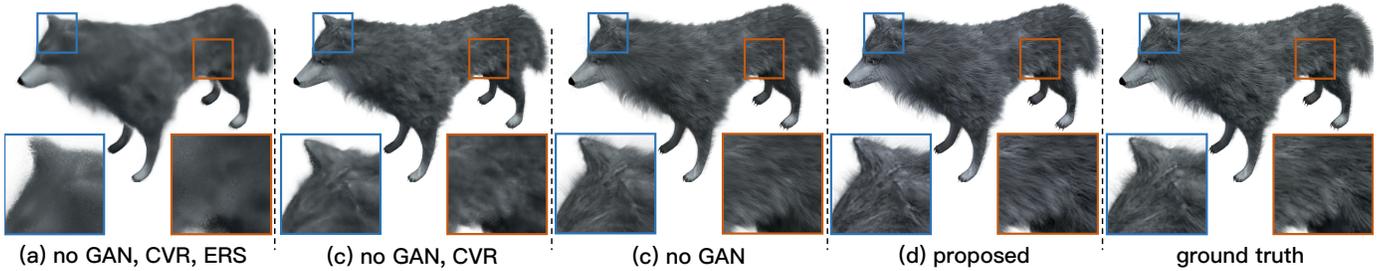


Fig. 8. Visualization of the ablation study on w/o the ERS, CVR and GAN loss modules (corresponding to the Tab 4).

TABLE 4  
Quantitative evaluation of the ablation studies.

Models	$\alpha F$		Alpha	
	PSNR $\uparrow$	LPIPS $\downarrow$	PSNR $\uparrow$	SAD $\downarrow$
(a) w/o GAN, CVR, ERS	20.27	0.175	20.18	23.51
(b) w/o GAN, CVR	31.38	0.074	34.46	3.808
(c) w/o GAN	<b>37.55</b>	0.026	37.93	2.642
(d) ConvNeRF	<b>37.23</b>	<b>0.022</b>	<b>38.63</b>	<b>2.564</b>

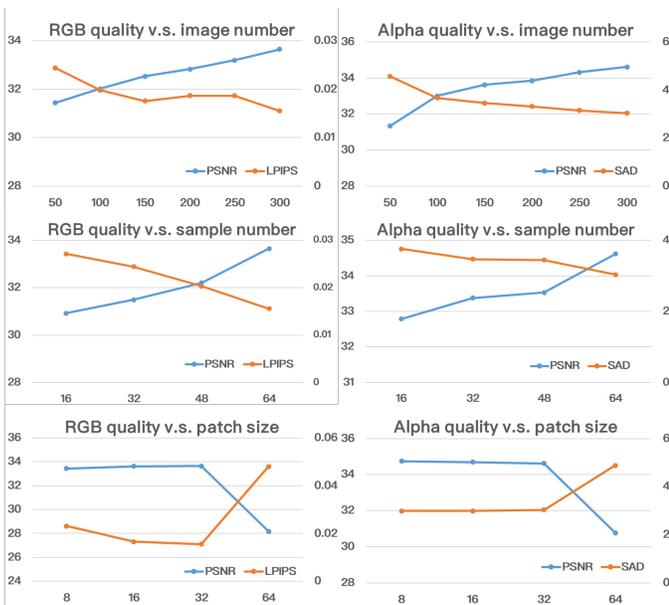


Fig. 9. From top to bottom: rendering quality v.s. training image number, samples per ray, patch size

On the other hand, our method relies on accurately calibrated camera pose; inaccurate camera poses for in the wild images will lead to loss of geometric details especially for fuzzy objects. Furthermore, since our method requires background images captured or separated by matting algorithms, it is almost impossible to manually capture the aligned background or obtain accurate alpha matte of images with a complex background, resulting in hard to extend it to in the wild data.

## 6 CONCLUSION

We have presented a novel neural rendering framework to combine both explicit opacity modeling and convolutional

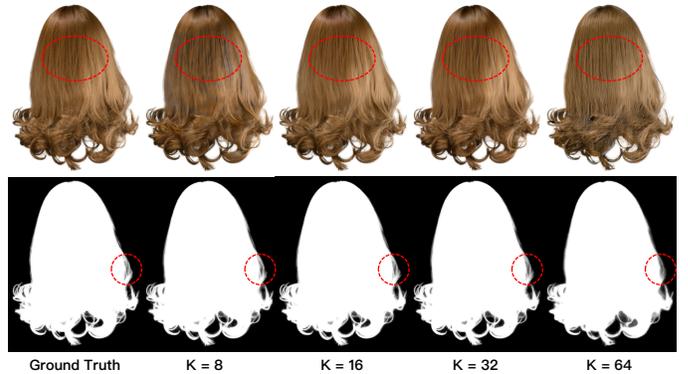


Fig. 10. Ablation study on the patch size of  $K = [8, 16, 32, 64]$ . Small and large sizes produce color cast and blurry artifacts.

TABLE 5  
Quantitative evaluation of textured background.

Background		$\alpha F$		Alpha	
		PSNR $\uparrow$	LPIPS $\downarrow$	PSNR $\uparrow$	SAD $\downarrow$
textured	NeRF	33.26	0.060	31.05	6.478
	Ours	<b>37.95</b>	0.023	35.51	4.309
white	NeRF	32.56	0.061	32.14	4.777
	Ours	<b>37.23</b>	<b>0.022</b>	<b>38.63</b>	<b>2.564</b>

mechanism into the neural radiance field, enabling high-quality, globally consistent and free-viewpoint appearance and opacity rendering for fuzzy objects. Our efficient sampling strategy enables efficient radiance field sampling and learning in a patch-wise manner, while our novel volumetric integration generates per-patch hybrid features to reconstruct the view-consistent fine-detailed appearance and opacity output. Our novel patch-wise adversarial training scheme further preserves the high-frequency appearance and opacity details for photo-realistic rendering in a self-supervised framework. Our experimental results demonstrate the effectiveness of the proposed convolutional neural opacity radiance field for high-quality appearance and opacity modeling. We believe that our approach is a significant step to enable photo-realistic modeling and rendering for challenging fuzzy objects, with many potential applications in VR/AR like gaming, entertainment and immersive telepresence.

## ACKNOWLEDGMENTS

This work was supported by NSFC programs (61976138, 61977047), the National Key Research and Development Program (2018YFB2100500), STCSM (2015F0203-000-06) and SHMEC (2019-01-07-00-01-E00003).

## REFERENCES

- [1] S. Paris, W. Chang, O. I. Kozhushnyan, W. Jarosz, W. Matusik, M. Zwicker, and F. Durand, "Hair photobooth: geometric and photometric acquisition of real hairstyles." *ACM Trans. Graph.*, vol. 27, no. 3, p. 30, 2008.
- [2] L. Luo, H. Li, and S. Rusinkiewicz, "Structure-aware hair capture," *ACM Transactions on Graphics (Proceedings SIGGRAPH 2013)*, vol. 32, no. 4, July 2013.
- [3] Z. Xu, H.-T. Wu, L. Wang, C. Zheng, X. Tong, and Y. Qi, "Dynamic hair capture using spacetime optimization," *ACM Trans. Graph.*, vol. 33, no. 6, Nov. 2014. [Online]. Available: <https://doi.org/10.1145/2661229.2661284>
- [4] G. Nam, C. Wu, M. H. Kim, and Y. Sheikh, "Strand-accurate multi-view hair capture," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2019.
- [5] C. Buehler, M. Bosse, L. Mcmillan, S. Gortler, and M. Cohen, "Unstructured lumigraph rendering," in *Conference on Computer Graphics and Interactive Techniques*, 2001, pp. 425–432.
- [6] W. Matusik, H. Pfister, A. Ngan, P. Beardsley, R. Ziegler, and L. McMillan, "Image-based 3d photography using opacity hulls," *ACM Transactions on Graphics (TOG)*, vol. 21, no. 3, pp. 427–437, 2002.
- [7] N. K. Kalantari, T.-C. Wang, and R. Ramamoorthi, "Learning-based view synthesis for light field cameras," vol. 35, no. 6, Nov. 2016. [Online]. Available: <https://doi.org/10.1145/2980179.2980251>
- [8] V. Sitzmann, M. Zollhöfer, and G. Wetzstein, "Scene representation networks: Continuous 3d-structure-aware neural scene representations," in *Advances in Neural Information Processing Systems*, 2019.
- [9] S. Lombardi, T. Simon, J. Saragih, G. Schwartz, A. Lehrmann, and Y. Sheikh, "Neural volumes: Learning dynamic renderable volumes from images," *arXiv preprint arXiv:1906.07751*, 2019.
- [10] J. Thies, M. Zollhöfer, and M. Nießner, "Deferred neural rendering: Image synthesis using neural textures," *ACM Transactions on Graphics (TOG)*, vol. 38, no. 4, pp. 1–12, 2019.
- [11] C. Wang, M. Wu, Z. Wang, L. Wang, H. Sheng, and J. Yu, "Neural opacity point cloud," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2020.
- [12] B. Mildenhall, P. P. Srinivasan, M. Tancik, J. T. Barron, R. Ramamoorthi, and R. Ng, "Nerf: Representing scenes as neural radiance fields for view synthesis," *arXiv preprint arXiv:2003.08934*, 2020.
- [13] K.-A. Aliev, A. Sevastopolsky, M. Kolos, D. Ulyanov, and V. Lempitsky, "Neural point-based graphics," *arXiv preprint arXiv:1906.08240*, 2019.
- [14] V. Sitzmann, J. Thies, F. Heide, M. Nießner, G. Wetzstein, and M. Zollhofer, "Deepvoxels: Learning persistent 3d feature embeddings," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 2437–2446.
- [15] H. Kato, Y. Ushiku, and T. Harada, "Neural 3d mesh renderer," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2018.
- [16] J. J. Park, P. Florence, J. Straub, R. Newcombe, and S. Lovegrove, "DeepSDF: Learning continuous signed distance functions for shape representation," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 165–174.
- [17] L. Mescheder, M. Oechsle, M. Niemeyer, S. Nowozin, and A. Geiger, "Occupancy networks: Learning 3d reconstruction in function space," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 4460–4470.
- [18] L. Liu, J. Gu, K. Z. Lin, T.-S. Chua, and C. Theobalt, "Neural sparse voxel fields," *NeurIPS*, 2020.
- [19] R. Martin-Brualla, N. Radwan, M. S. Sajjadi, J. T. Barron, A. Dosovitskiy, and D. Duckworth, "Nerf in the wild: Neural radiance fields for unconstrained photo collections," *arXiv preprint arXiv:2008.02268*, 2020.
- [20] K. Park, U. Sinha, J. T. Barron, S. Bouaziz, D. B. Goldman, S. M. Seitz, and R.-M. Brualla, "Deformable neural radiance fields," *arXiv preprint arXiv:2011.12948*, 2020.
- [21] M. Tancik, P. P. Srinivasan, B. Mildenhall, S. Fridovich-Keil, N. Raghavan, U. Singhal, R. Ramamoorthi, J. T. Barron, and R. Ng, "Fourier features let networks learn high frequency functions in low dimensional domains," 2020.
- [22] Z. Chen and H. Zhang, "Learning implicit fields for generative shape modeling," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 5939–5948.
- [23] S. Peng, M. Niemeyer, L. Mescheder, M. Pollefeys, and A. Geiger, "Convolutional occupancy networks," *arXiv preprint arXiv:2003.04618*, 2020.
- [24] S. Saito, Z. Huang, R. Natsume, S. Morishima, A. Kanazawa, and H. Li, "Pifu: Pixel-aligned implicit function for high-resolution clothed human digitization," in *Proceedings of the IEEE International Conference on Computer Vision*, 2019, pp. 2304–2314.
- [25] A. X. Chang, T. Funkhouser, L. Guibas, P. Hanrahan, Q. Huang, Z. Li, S. Savarese, M. Savva, S. Song, H. Su *et al.*, "Shapenet: An information-rich 3d model repository," *arXiv preprint arXiv:1512.03012*, 2015.
- [26] M. Niemeyer, L. Mescheder, M. Oechsle, and A. Geiger, "Differentiable volumetric rendering: Learning implicit 3d representations without 3d supervision," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 3504–3515.
- [27] P. Debevec, C. Bregler, M. Cohen, and L. Mcmillan, "Image-based modeling and rendering," 1998, p. 299.
- [28] M. Levoy, "Light field rendering," in *Conference on Computer Graphics and Interactive Techniques*, 1996, pp. 31–42.
- [29] A. Chen, M. Wu, Y. Zhang, N. Li, J. Lu, S. Gao, and J. Yu, "Deep surface light fields," *Proceedings of the ACM on Computer Graphics and Interactive Techniques*, vol. 1, no. 1, pp. 1–17, 2018.
- [30] K. Zhang, G. Riegler, N. Snavely, and V. Koltun, "Nerf++: Analyzing and improving neural radiance fields," *arXiv preprint arXiv:2010.07492*, 2020.
- [31] M. Wu, Y. Wang, Q. Hu, and J. Yu, "Multi-view neural human rendering," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2020.
- [32] Z. Chen, A. Chen, G. Zhang, C. Wang, Y. Ji, K. N. Kutulakos, and J. Yu, "A neural rendering framework for free-viewpoint relighting," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 5599–5610.
- [33] Y.-Y. Chuang, B. Curless, D. H. Salesin, and R. Szeliski, "A bayesian approach to digital matting," in *Proceedings of the 2001 IEEE Computer Society Conference on Computer Vision and Pattern Recognition*. CVPR 2001, vol. 2. IEEE, 2001, pp. II–II.
- [34] J. Wang and M. F. Cohen, "An iterative optimization approach for unified image segmentation and matting," in *Tenth IEEE International Conference on Computer Vision (ICCV'05) Volume 1*, vol. 2. IEEE, 2005, pp. 936–943.
- [35] E. S. Gastal and M. M. Oliveira, "Shared sampling for real-time alpha matting," in *Computer Graphics Forum*, vol. 29, no. 2. Wiley Online Library, 2010, pp. 575–584.
- [36] E. Shahrinan, D. Rajan, B. Price, and S. Cohen, "Improving image matting using comprehensive sampling sets," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2013, pp. 636–643.
- [37] L. Karacan, A. Erdem, and E. Erdem, "Image matting with kl-divergence based sparse sampling," in *Proceedings of the IEEE International Conference on Computer Vision*, 2015, pp. 424–432.
- [38] E. Elhamifar, G. Sapiro, and S. S. Sastry, "Dissimilarity-based sparse subset selection," *IEEE transactions on pattern analysis and machine intelligence*, vol. 38, no. 11, pp. 2182–2197, 2015.
- [39] J. Sun, J. Jia, C.-K. Tang, and H.-Y. Shum, "Poisson matting," in *ACM SIGGRAPH 2004 Papers*, 2004, pp. 315–321.
- [40] L. Grady, T. Schiwietz, S. Aharon, and R. Westermann, "Random walks for interactive alpha-matting," in *Proceedings of VIIP*, vol. 2005, 2005, pp. 423–429.
- [41] A. Levin, D. Lischinski, and Y. Weiss, "A closed-form solution to natural image matting," *IEEE transactions on pattern analysis and machine intelligence*, vol. 30, no. 2, pp. 228–242, 2007.
- [42] Q. Chen, D. Li, and C.-K. Tang, "Knn matting," *IEEE transactions on pattern analysis and machine intelligence*, vol. 35, no. 9, pp. 2175–2188, 2013.
- [43] X. Chen, D. Zou, S. Zhiying Zhou, Q. Zhao, and P. Tan, "Image matting with local and nonlocal smooth priors," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2013, pp. 1902–1907.

- [44] N. Xu, B. Price, S. Cohen, and T. Huang, "Deep image matting," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 2970–2979.
- [45] S. Cai, X. Zhang, H. Fan, H. Huang, J. Liu, J. Liu, J. Liu, J. Wang, and J. Sun, "Disentangled image matting," in *Proceedings of the IEEE International Conference on Computer Vision*, 2019, pp. 8819–8828.
- [46] H. Yu, N. Xu, Z. Huang, Y. Zhou, and H. Shi, "High-resolution deep image matting," *arXiv preprint arXiv:2009.06613*, 2020.
- [47] Q. Hou and F. Liu, "Context-aware image matting for simultaneous foreground and alpha estimation," in *Proceedings of the IEEE International Conference on Computer Vision*, 2019, pp. 4130–4139.
- [48] M. Forte and F. Pitié, "*f*, *b*, alpha matting," *arXiv preprint arXiv:2003.07711*, 2020.
- [49] J. Johnson, A. Alahi, and L. Fei-Fei, "Perceptual losses for real-time style transfer and super-resolution," in *European conference on computer vision*. Springer, 2016, pp. 694–711.
- [50] Y. Zhang, L. Gong, L. Fan, P. Ren, Q. Huang, H. Bao, and W. Xu, "A late fusion cnn for digital matting," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 7469–7478.
- [51] S. Sengupta, V. Jayaram, B. Curless, S. M. Seitz, and I. Kemelmacher-Shlizerman, "Background matting: The world is your green screen," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 2291–2300.
- [52] K. M. Cheung, S. Baker, and T. Kanade, "Shape-from-silhouette of articulated objects and its use for human body kinematics estimation and motion capture," in *2003 IEEE Computer Society Conference on Computer Vision and Pattern Recognition, 2003. Proceedings.*, vol. 1, 2003, pp. I–I.
- [53] J. Yu, Z. Lin, J. Yang, X. Shen, X. Lu, and T. S. Huang, "Free-form image inpainting with gated convolution," in *Proceedings of the IEEE International Conference on Computer Vision*, 2019, pp. 4471–4480.
- [54] P. Isola, J.-Y. Zhu, T. Zhou, and A. A. Efros, "Image-to-image translation with conditional adversarial networks," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 1125–1134.
- [55] Z. Zhang, "Flexible camera calibration by viewing a plane from unknown orientations," in *Proceedings of the seventh IEEE international conference on computer vision*, vol. 1. Ieee, 1999, pp. 666–673.
- [56] R. Zhang, P. Isola, A. A. Efros, E. Shechtman, and O. Wang, "The unreasonable effectiveness of deep features as a perceptual metric," in *CVPR*, 2018.



**Haimin Luo** received the B.S. degree from the School of Computer Science and Technology, Huazhong University of Science and Technology, Wuhan, China, in 2019. He is currently working toward the master's degree at ShanghaiTech University, Shanghai, China. His research interests include computer vision, computer graphics, and deep learning.



**Anpei Chen** received the B.S. degree from the School of Physics and Optoelectronic Engineering, Xidian University, Shanxi, China, in 2016. He is currently working toward the PhD degree at ShanghaiTech University, Shanghai, China. His research interests lie at the intersection of computer graphics and computer vision, including image synthesis/editing, geometric modeling, and realistic rendering.



**Qixuan Zhang** is currently working toward the B.S. degree at ShanghaiTech University, Shanghai, China. His research interests include computer vision, computational photography, and deep learning.



**Bai Pang** is currently working toward the B.S. degree at ShanghaiTech University, Shanghai, China. His research interests include computer vision, computational photography and machine learning.



**Minye Wu** received the B.S. degree from the School of Computer Engineering and Science, Shanghai University, Shanghai, China, in 2015. He is currently working toward the PhD degree at ShanghaiTech University, Shanghai, China. He is also with the Shanghai Institute of Microsystem and Information Technology, Chinese Academy of Sciences, China, and the University of Chinese Academy of Sciences, China. His research interests include computer vision, deep learning, and computational photography.



**Lan Xu** received the B.E. degree from Zhejiang University in 2015 and the Ph.D. degree from the Department of Electronic and Computer Engineering (ECE), The Hong Kong University of Science and Technology (HKUST), in 2020. He is currently an Assistant Professor with ShanghaiTech University. His research interests include computer vision, computer graphics, and machine learning.



**Jingyi Yu** received B.S. from Caltech in 2000 and PhD from MIT in 2005. He is currently the Vice Provost at the ShanghaiTech University. Before joining ShanghaiTech, he was a full professor in the Department of Computer and Information Sciences at University of Delaware. His research interests span a range of topics in computer vision and computer graphics, especially on computational photography and non-conventional optics and camera designs. He is a recipient of the NSF CAREER Award and the AFOSR YIP Award, and has served as an area chair of many international conferences including CVPR, ICCV, ECCV, IJCAI and NeurIPS. He is currently a program chair of CVPR 2021 and will be a program chair of ICCV 2025. He has been an associate editor of the IEEE Transactions on Pattern Analysis and Machine Intelligence, the IEEE Transactions on Image Processing, and the Elsevier Computer Vision and Image Understanding. He is a fellow of IEEE.