# Non-line-of-sight Imaging via Neural Transient Fields

Siyuan Shen[†], Zi Wang[†], Ping Liu, Zhengqing Pan, Ruiqian Li, Tian Gao,
Shiying Li, and Jingyi Yu, *Fellow, IEEE*

**Abstract**—We present a neural modeling framework for non-line-of-sight (NLOS) imaging. Previous solutions have sought to explicitly recover the 3D geometry (e.g., as point clouds) or voxel density (e.g., within a pre-defined volume) of the hidden scene. In contrast, inspired by the recent Neural Radiance Field (NeRF) approach, we use a multi-layer perceptron (MLP) to represent the neural transient field or NeTF. However, NeTF measures the transient over spherical wavefronts rather than the radiance along lines. We therefore formulate a spherical volume NeTF reconstruction pipeline, applicable to both confocal and non-confocal setups. Compared with NeRF, NeTF samples a much sparser set of viewpoints (scanning spots) and the sampling is highly uneven. We thus introduce a Monte Carlo technique to improve the robustness in the reconstruction. Experiments on synthetic and real datasets demonstrate NeTF achieves state-of-the-art performance and can provide reliable reconstructions even under semi-occlusions and on non-Lambertian materials.

**Index Terms**—Computational photography, non-line-of-sight imaging, neural radiance field, neural rendering

✦

## 1 INTRODUCTION

NON-line-of-sight (NLOS) imaging employs time-resolved measurements for recovering hidden scenes beyond the direct line of sight from a sensor [17], [37]. Applications are numerous, ranging from remote sensing to autonomous driving and to rescue missions in hazardous environments. Most existing NLOS setups orient an ultra-fast pulsed laser beam towards a relay wall in the line of sight where the wall diffuses the laser into spherical wavefronts towards the hidden scene. As the wavefront hits the scene and bounces back onto the wall, a time-of-flight

- [†] *indicates equal contribution.*

- *S. Shen is with the School of Information Science and Technology, ShanghaiTech University, Shanghai 201210, China, and also with DGene, Inc., Shanghai, China. E-mail: shensy@shanghaitech.edu.cn.*

- *Z. Wang is with the Shanghai Institute of Technical Physics, Chinese Academy of Sciences, Shanghai 200083, China, and the School of Information Science and Technology, ShanghaiTech University, Shanghai 201210, China, and also with the University of Chinese Academy of Sciences, Beijing 100049, China. E-mail: wangzi@shanghaitech.edu.cn.*

- *P. Liu and R. Li are with the School of Information Science and Technology, ShanghaiTech University, Shanghai 201210, China. E-mail: {liuping, lirq1}@shanghaitech.edu.cn.*

- *Z. Pan is 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 200050, China, and also with the University of Chinese Academy of Sciences, Beijing 100049, China. E-mail: panzhq@shanghaitech.edu.cn.*

- *T. Gao is with the School of Physical Science and Technology, ShanghaiTech University, Shanghai 201210, China. E-mail: gaotian@shanghaitech.edu.cn.*

- *S. Li, J. Yu are with Shanghai Engineering Research Center of Intelligent Vision and Imaging, and also with the School of Information Science and Technology, ShanghaiTech University, Shanghai 201210, China. E-mail: {lishy1, yujingyi}@shanghaitech.edu.cn. (Corresponding authors: Shiying Li; Jingyi Yu).*
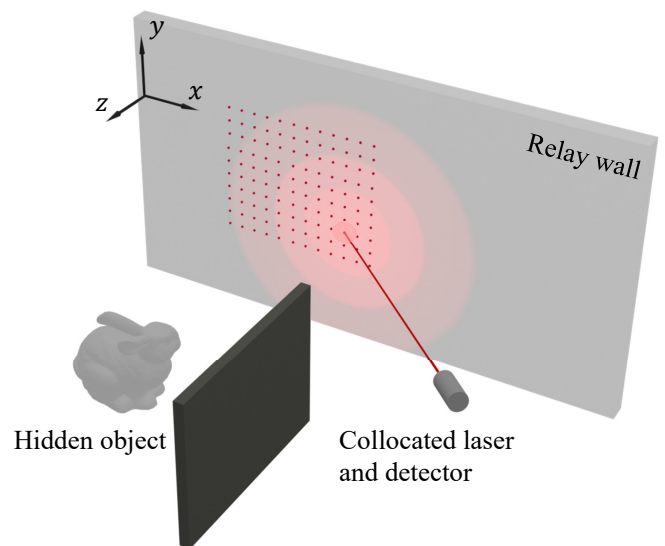
Fig. 1. A typical confocal NLOS imaging system aims a laser towards a diffuse wall that serves as a virtual reflector. The hidden scene is indirectly illuminated as spherical waves that intersect with the scene and are reflected back onto the wall. A SPAD sensor measures at different spots on the wall to form transient images.

(ToF) detector with picosecond resolution (such as the streak camera [1], [2], [3] or recently the more affordable single-photon avalanche diodes (SPADs) [9], [15], [16], [31], [32], [34]) can be used to record the arrival time and the number of the returning photons [5], [6], [7]. SPAD sensors in a time correlated single photon counting (TCSPC) mode can thus produce transients, of which a single pixel corresponds to a specific pair of illumination and detection spots on the wall and a histogram of the photon counts versus time bins.

The measured transients contain rich geometric information of a hidden scene, potentially usable for scene recov-

ery. In reality, the process corresponds to a typical inverse imaging problem that generally incurs high computational cost, especially because the transients are high dimensional signals. To make the problem tractable, the pioneering back-projection (BP) technique and its variations assume smooth objects so that scene recovery can be modeled as deconvolution [1], [2], [8], [31]. Alternatively, the light-cone transform (LCT) based methods collocate the illumination and sensing spots on the relay wall so that the forward imaging model can be simplified as 3D convolution [4], [15], [34], [35] where advanced signal processing techniques such as Wiener filters [34], [35], [38] can further reduce noise. Assuming the scene is near diffuse, analysis-by-synthesis algorithms can improve reconstruction [18], [36]. The seminal work of the Fermat path based approaches [16], [32] can handle highly specular objects by simultaneously recovering the position and normal of Fermat points on the surface.

We present a novel volumetric NLOS imaging framework by modeling the transient field via deep networks. Our Neural Transient Field (NeTF) technique is inspired by the recent neural radiance field (NeRF) that conducts 3D reconstruction and view synthesis from a set of input images. Different from existing multi-view stereo (MVS) techniques, NeRF assumes a volume rendering model and sets out to use multi-layer perception (MLP) to recover per-voxel scene density and per-direction color. We observe that NLOS resembles MVS in that each scanning point on the wall resembles a virtual camera and therefore a similar deep learning technique may be potentially used for scene recovery. Different from NeRF, though, NeTF measures the transient over spherical wavefronts rather than the radiance along lines. We therefore first formulate volumetric transient fields under the spherical coordinate and devise an MLP that trains on the measurements to predict per-voxel density and view-dependent albedo. Our NeTF formulation is applicable to both confocal and non-confocal setups.

Compared with NeRF, NeTF captures a much sparser set of viewpoints (scanning spots) and the distribution of scene points on the spherical wavefronts can be highly uneven. We therefore develop a Markov chain Monte Carlo (MCMC) technique based on importance sampling for matching the actual scene distribution. We conduct comprehensive experiments on existing synthetic and real datasets. We demonstrate that NeTF achieves state-of-the-art reconstruction quality under both confocal and non-confocal settings. In particular, the trained MLP provides a continuous 5D representation of the hidden scene without requiring digitizing the NLOS volume or optimizing surface parameters. The 5D representation with directional encoding can handle view-dependent albedo of non-Lambertian surface reflectance and strong self-occlusions, under both confocal and non-confocal setups. All our codes and data are available at https://github.com/zeromakerplus/NeTF_public.

## 2 RELATED WORK

As an emerging computational imaging technique, NLOS imaging has found broad applications in computer vision and computer graphics, ranging from recovering 3D shape of hidden objects [2], [14], [15], [22], [31], [32], [34], [35] to tracking hidden moving objects [2], [4], [10], [13]. Existing

solutions employ time-resolved optical detectors such as streak cameras [1], [2], SPADs [5], [31] and interferometry [12], [16] or non-optical acoustic [11] and thermal [10] sensors to indirectly measure the hidden scene and then apply inverse imaging techniques for recovery. We refer readers to recent surveys [17], [21], [37] for a comprehensive overview.

**Confocal vs. Non-Confocal.** Kirmani et al. [19], [20] designed and implemented the first prototype non-confocal NLOS system and derived a linear time-invariant model amenable to multi-path light transport analysis. In reality, varying both the laser beam and the measuring spot yield a high-dimensional transient field analogous to the light field. Many efforts have since been focused on imposing priors and constraints to accelerate data processing. Velten et al. [2] proposed a back-projection technique with ellipsoidal constraints: the observing point and the laser projection point on the wall correspond to the foci of a set of ellipsoids, each corresponding to a specific transient. The hidden scene can then be reconstructed by intersecting the ellipsoids. To further improve reconstruction quality and speed, subsequent work has applied filtering techniques such as sharpening and thresholding [1], [8], [31]. Alternatively, one can directly model the scene using parametric surfaces and then optimize the parameters over the observations [18], [36], [38]. Ahn et al. [38] model parameter fitting as a linear least-squares problem using a convolutional Gram operator. It is also possible to adopt wave optics for NLOS imaging [14], [15], [23], [24], by characterizing the problem as specific properties of a temporally evolving wave field in the Fourier domain.

To reduce data dimensionality, several recent approaches adopt a confocal setting [34], [35] where the laser and the detector (e.g., a SPAD) collocate, e.g., via a beam splitter. Consequently, the ellipsoidal constraints degenerate to be spherical, simplifying the inverse problem with a 3D deconvolution and system calibration. The seminal work of light-cone transform (LCT) [34] casts the NLOS reconstruction problem as Wiener filtering in the Fourier domain and can achieve a low computational complexity of $O(N^3 \log N)$ for $N^3$ voxels, compared to $O(N^5)$ in the traditional BP methods. Yong et al. [35] formulate the albedo and normal recovery based on directional LCT (DLCT) as a vector deconvolution problem. The confocal setting results in an overwhelming contribution of direct light first-bounce off the diffuse wall and subsequently produces useful geometric constraints. The seminal work by Xin et al. [16] exploits the Fermat flow induced by the transients for estimating surface normals. Lindell et al. [15] adapt an F-K migration in seismology to convert the surface reconstruction problem to a boundary value problem. The F-K migration method enables faster reconstruction and supports planar or non-planar diffuse walls.

**Volume vs. Surface.** Existing NLOS methods can also be categorized in terms of the form of the reconstruction results. Two most adopted forms are volume density and points/surfaces. Methods for recovering the former generally discretize the scene into voxels and compute the density, either using intersections of wavefronts under ellipsoidal [1], [2], [5], [8], [29], [31], [38] and spherical [34], [35] constraints, or via modeling the imaging process as convolution and recovering the volume via specially designed deconvo-
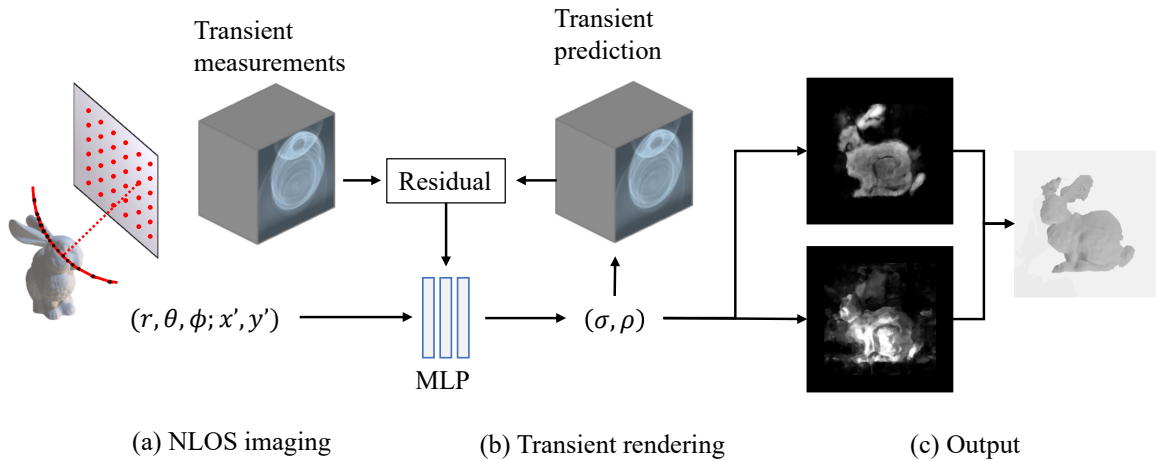
Fig. 2. Our neural transient field (NeTF) reconstruction pipeline. We parameterize every point on a spherical wavefront in terms of the origin on the wall, and the direction and radius of its corresponding spherical coordinates. We set out to recover the transient field under this parameterization via a multi-layer perception under spherical volume rendering.
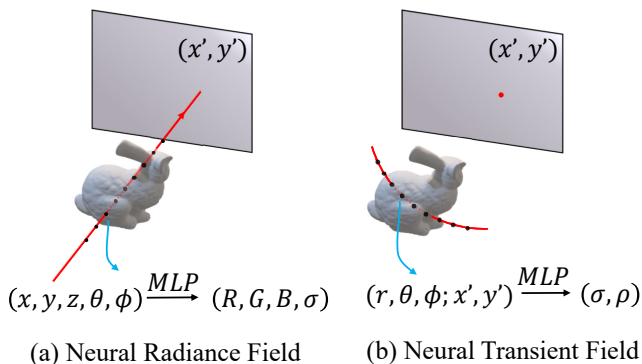


Fig. 3. NeRF vs. NeTF. In NeRF, volume density is accumulated along every line (ray) whereas in NeTF it is accumulated on a spherical wavefront.

lution filters. Methods for recovering the latter have relied on light transport physics [18], [36] for optimizing the shape and reflectance of the hidden scenes. Such methods are generally mathematically tractable but are computationally expensive particularly because higher order geometry such as the surface normal needs to integrated into the optimization process. Reconstruction results are either sparse as in Fermat [16] where only discontinuities in the transient were used, or rely heavily on the quality of the basis shape as in [18].

Our approach falls into the category of volume based technique. We are inspired by the recent multi-view reconstruction framework Neural Radiance Field (NeRF) that aims to recover the density and color at every point along every ray, implicitly providing a volumetric reconstruction. NeRF adopts a volume rendering model and sets out to optimize volume density that best matches the observation using a Multi-Layer Perception (MLP). It is also possible to modify NeRF to tackle photometric stereo (PS) problems where the camera is fixed but the lighting conditions vary. We observe the non-confocal NLOS imaging process greatly

resembles MVS/PS: fixing the laser beam but measuring the transient at different spots on the wall resembles MVS, whereas fixing the measuring spot but varying the laser beam resembles PS. In fact, the confocal setting is very similar to the NeRF AA setting [41] where the lighting and the camera move consistently. We therefore call our reconstruction scheme Neural Transient Field or NeTF.

Both NeTF and NeRF use MLP as an optimizer. However, there are several major differences between NeRF and NeTF. First, the volume rendering model used in NeRF is not directly applicable to NeTF. We therefore derive a novel volumetric image formation model under NLOS. Second, NLOS measures the transient rather than the radiance. In fact, the transient is measured from the sum of returning photons on a wavefront instead of a single ray. We hence formulate a spherical volume reconstruction pipeline. Finally, NeRF generally assumes dense ray samples, whereas the NLOS setting is much more sparse. We thus introduce a Monte Carlo technique to improve the robustness in the reconstruction.

## 3 NEURAL TRANSIENT FIELD

We recognize that the NLOS reconstruction problem resembles multi-view reconstruction in the line-of-sight (LOS) and adopt a neural reconstruction framework analogous to NeRF [39]. Each detection spot on the relay wall can be viewed as a virtual *camera*. These cameras capture the transients of the NLOS scene as if viewed from the wall. We adopt the plenoptic radiance field notion of NeRF and represent the NLOS scene as a continuous 5D function of transients, i.e., a plenoptic transient field. We then set out to infer scene density at every point along every spherical wavefront via deep network based optimization. It is important to note that, same as NeRF, our neural transient field (NeTF) representation chooses not to explicitly discretize the scene into volumes. Rather, we use multi-layer perception (MLP) to virtually represent the volume.

## 3.1 Scene Representation

The NeRF framework [39] uses the neural radiance field $L(x, y, z, \theta, \phi)$ as scene representation where $(x, y, z)$ corresponds to a point on a ray and $(\theta, \phi)$ the direction of the ray. Its trained network outputs both the density $\sigma$ at every position $(x, y, z)$ and the (view-dependent) color $c = (r, g, b)$ along direction $(\theta, \phi)$. The density can be further used for scene reconstruction and the color for image-based rendering. In our case, NeTF, instead of sampling on a single camera ray, samples a hemisphere of rays as light propagates as a spherical wave from the relay wall towards the hidden scene. We hence adopt a continuous 5D function of transients $L_{\text{NLOS}}$ under the spherical coordinates as:

$$L_{\text{NLOS}}(x', y', r, \theta, \phi) \rightarrow (\sigma, \rho) \tag{1}$$

where $P(x', y')$ is a detection spot on the wall that serves as the origin of the hemisphere. $Q(r, \theta, \phi)$ is a scene point parameterized using the spherical coordinate $(r, \theta, \phi)$ w.r.t. $P(x', y')$. Similar to NeRF, though, we set out to design a fully connected neural network, i.e., an MLP, to estimate $L_{\text{NLOS}}$. Different from NeRF, $L_{\text{NLOS}}$ in NeTF outputs a volume density $\sigma$ and a surface reflectance (albedo $\rho$) rather than color along the direction $(\theta, \phi)$.

Recall that NLOS needs to scan different spots on the relay wall, resulting in inconsistent spherical coordinates and casting challenges in network training and inference. We thus first transform the spherical coordinates $(x', y', r, \theta, \phi)$ to their corresponding Cartesian coordinates, i.e., $(x, y, z, \theta, \phi)$ as

$$R : \begin{cases} x = r \sin\theta \cos\phi + x' \\ y = r \sin\theta \sin\phi + y' \\ z = r \cos\theta \end{cases} \tag{2}$$

The transform $R$ ensures that the position of a 3D voxel is consistent when we scan over different detection spots. All subsequent training under MLP should be conducted under the Cartesian coordinate for density and view dependent albedo inferences.

$$L_{\text{NLOS}} : (x', y', r, \theta, \phi) \xrightarrow{R} (x, y, z, \theta, \phi) \xrightarrow{\text{MLP}} (\sigma, \rho) \tag{3}$$

Same as NeRF, a key benefit of NeTF is that we no longer need to discretize the scene into a fixed-resolution volume representation. Instead, the deep network representation can provide scene reconstructions at an arbitrary resolution, recovering fine details largely missing in prior art.

## 3.2 Forward Model

We first reformulate the NLOS reconstruction problem as a forward model under our NeTF representation. Under the confocal setting [15], [27], [34], the illumination and detection collocate at the same spot $P(x', y')$ on a relay wall, producing a spherical wave anchored at the spot. The transient $\tau_{\text{iso}}(x', y', t)$ recorded at each spot $P(x', y')$ is the summation of photons that are reflected back at a specific
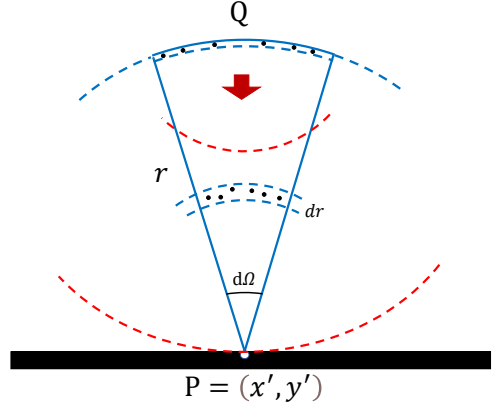


Fig. 4. The volume rendering model derived under the spherical coordinate system, suitable for processing in NeTF.

time instant $t$ from the NLOS scene in the 3D half-space $\Omega$ as:

$$\tau_{\text{iso}}(x', y', t) = \iiint_\Omega \frac{1}{r^4} \rho_{\text{iso}}(x, y, z) g(x', y', x, y, z) \tag{4}$$
$$\cdot \delta\left(2\sqrt{(x' - x)^2 + (y' - y)^2 + z^2} - tc\right) dx\, dy\, dz$$

where $c$ is the speed of light, $r$ is the distance between the wall and the NLOS scene as $r = \sqrt{(x' - x)^2 + (y' - y)^2 + z^2} = tc/2$, and $1/r^4$ is the light fall-off term. $\rho_{\text{iso}}(x, y, z)$ is the albedo of an NLOS point, where $x, y, z$ are the spatial coordinates of the point. Function $g$ models the time-independent effects, including the surface normal, bidirectional reflectance distribution functions (BRDFs), occlusions patterns, etc. The Dirac delta function relates the time of flight $t$ to the distance $r$.

Notice that function $g$ makes the imaging process nonlinear. To make the problem tractable, previous linear approximation schemes such as [34] adopt $g = 1$ by assuming that the NLOS scene scatters isotropically and that no occlusions occur within the NLOS scene. Such assumptions, however, restrict NLOS scenes to being Lambertian and convex. In contrast, NeTF, by adopting a deep network to model the imaging process, can tackle non-linearity without imposing explicit constraints on $g$.

We intend to specify how much an NLOS point in the hemisphere contributes to the transient through photons propagation. Consider a detection spot $P(x', y')$ to an NLOS point $Q(r, \theta, \phi)$ in a hemisphere centered at $P$. Recall that the scattering equation [30] serves as the foundation for volume rendering and thereby NeRF. We first derive a photon version of the scattering equation, with more details provided in the supplementary materials.

In our NLOS setting, photons travel along spherical wavefronts. When they reach either the relay wall or the hidden surface, they are reflected and then continue to propagate along a hemisphere. We assume that the spot $P$ is a patch with radius $r_0$, and that the location $Q$ in the hidden scene with its neighbors that contribute to $P$ forms a spherical cross-section with radius $r$, thickness $dr$, and

a solid angle $d\Omega$. Fig. 4 shows that photons travel from $P$ to $Q$ and back from $Q$ to $P$. When $dr$ is sufficiently small, the inner and outer surface areas of the cross-section are $S = r^2 d\Omega$, and the volume is $S dr$. Recall that $\sigma$ is the density of particles in the cross-section, therefore the number of particles is $\sigma S dr$. Assuming a particle has a radius $a$, the projected area on the surface can then be computed as $A = \pi a^2$. We assume that light energy $E$ is attenuated due to particles' absorption and scattering and consequently the energy loss $dE$ can be computed as:

$$dE = -\frac{\pi a^2 \sigma r^2 d\Omega dr}{r^2 d\Omega}E = -A\sigma E dr \tag{5}$$

The attenuation coefficient can be computed as $e^{\int_0^r -A\sigma(r',\theta,\phi) dr'}$ along the radius $r$. Recall that the spot $P$ has a radius $r_0$ and emits radiant energy as a constant $E_P$. Taking integral of Eqn. 5, energy received at $Q$ is defined as:

$$E_Q = \exp\left(\int_0^r -A\sigma(r',\theta,\phi) dr'\right)E_P \frac{r^2 \cdot d\Omega}{r^2 \cdot 2\pi} \tag{6}$$

We now consider the reflection at $Q$. Assume that the cross-section is sufficiently thin, e.g., $dr = 2a$, the radiant energy at $Q$ attenuated due to absorption and reflection w.r.t. the reflectance $\rho$ can be defined as

$$E_Q'(r,\theta,\phi) = A \cdot \sigma(r,\theta,\phi) \cdot 2a \cdot \rho(r,\theta,\phi) \cdot E_Q \tag{7}$$

On the returning path from $Q$ to $P$, the wavefronts form hemispheres centered at $Q$ with radius $r$. The spot $P$ with area $\pi r_0^2$ receives the photons back to the relay wall, we thus have the energy at $P$ w.r.t. the solid angle $d\Omega$ as:

$$E_P'(r,\theta,\phi) = \exp\left(\int_0^r -A\sigma(r',\theta,\phi) dr'\right)E_Q' \frac{r_0^2 \cdot 2\pi}{r^2 \cdot 2\pi} \tag{8}$$

By taking the integral of Eqn. 8 w.r.t. the solid angle $d\Omega$, we obtain energy received at the detection $P(x',y')$ in the hemisphere at a time instant $t$ as:

$$\tau(x',y',t) = \iint_{H(x',y';\frac{ct}{2})} E_P'(r,\theta,\phi) d\Omega \tag{9}$$

Eqn. 9 serves our forward imaging model. It essentially maps an NLOS point $Q(r,\theta,\phi)$ to a transient $\tau$ detected at a spot $P(x',y')$ on a diffuse surface at a time instant $t$. For clarity, we abbreviate $\sigma(r,\theta,\phi;x',y')$ as $\sigma(r,\theta,\phi)$, and $\rho(r,\theta,\phi;x',y')$ as $\rho(r,\theta,\phi)$. Substituting Eqns. 6, 7 and 8, Eqn. 9 is rewritten as:

$$\tau(x',y',t) =$$
$$\Gamma_0 \iint_{H(x',y';\frac{ct}{2})} \frac{1}{r^2}\sigma(r,\theta,\phi)\rho(r,\theta,\phi)exp\left(2\int_0^r -A\sigma dr'\right) d\Omega \tag{10}$$

where constant $\Gamma_0 = Aar_0^2 E_P/\pi$ is determined by particle radius $a$, initial energy $E_P$, and patch radius $r_0$. The integration domain $H(x',y';\frac{ct}{2})$ is a hemisphere centered at $P(x',y')$ on a relay wall, with a radius of $r = ct/2$. $\theta$ and $\phi$ are the elevation and azimuth angles in the viewing

direction from $P(x',y')$ to an NLOS point, equivalent to those in the direction of reflection from the NLOS scene. $\rho(r,\theta,\phi;x',y')$ models view-varying BRDFs of the NLOS scene. $e^{2\int_0^r -A\sigma(r',\theta,\phi) dr'}$ is an exponential actuation coefficient and reveals the visibility of an NLOS point with respect to varying detection spots $P(x',y')$. Since $d\Omega = \sin\theta d\theta d\phi$, we have:

$$\tau(x',y',t) = \Gamma_0 \iint_{H(x',y';\frac{ct}{2})} \frac{\sin\theta}{r^2}\sigma(r,\theta,\phi)\rho(r,\theta,\phi)\cdot$$
$$exp\left(2\int_0^r -A\sigma dr'\right) d\theta d\phi \tag{11}$$

Recall that the forward plenotpic transient field model in Eqn. 11 is computationally expensive if we use the MLP for training. If we further assume that the NLOS scene is all opaque and does not exhibit self-occlusions, we can further simplify the formulation to:

$$\tau(x',y',t) = \Gamma_0 \iint_{H(x',y';\frac{ct}{2})} \frac{\sin\theta}{r^2}\sigma(r,\theta,\phi)\rho(r,\theta,\phi) d\theta d\phi \tag{12}$$

Such a formulation reduces computations and is used in examples of Figs. 6, 8, 9, 10, and 11 to accelerate processing. Its downside though is that, unlike Eqn. 11 , it cannot handle occlusions. For scenes that contain heavy occlusion, we use Eqn. 11 , e.g., in Fig. 14.

Finally, it is noting that although both NeRF and our NeTF derive the forward model based on volume rendering, NeRF models how a ray propagates along a line (i.e., with a cylinder between two points) whereas NeTF models spherical wavefront propagation (i.e., with a cone model that accounts for attenuation). In addition, the volume rendering model used in NeRF only considers one-way accumulation, i.e., how light travels through light emitting particles towards the camera sensor. In contrast, our NeTF adopts a two-way propagation model, i.e., how light illuminates the scene and how scenes illuminate the wall.

### 3.3 Differentiable Rendering

Our forward model is differentiable, we can therefore numerically compute the continuous integral Eqn. 12 using quadrature, as:

$$\tau(x',y',t) = \frac{\Delta\theta\Delta\phi}{r^2}\sum_{i,j}\sin(\theta_{ij})\sigma(r,\theta_{ij},\phi_{ij})\rho(r,\theta_{ij},\phi_{ij}) \tag{13}$$

where $Q(r,\theta_{ij},\phi_{ij})$ are scene points uniformly sampled along the hemispherical rays. We transform these points into their corresponding Cartesian coordinates as inputs to the MLP. The network outputs the density and reflectance at each point. We then sum all the outputs as neural transient fields from the transients.

We optimize our NeTF by minimizing the following $l_2$-norm loss function as the difference between the predicted $\tau(x',y',t)$ and measured $\tau_m(x',y',t)$ transients:

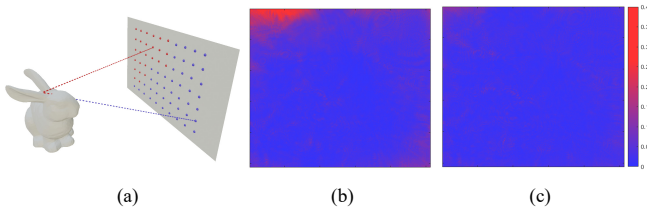$$L = \sum_{x',y',t} (\tau_m(x',y',t) - \tau(x',y',t))^2 \tag{14}$$

Fig. 5. We use the loss induced from the first stage in the training process to guide resampling. At each spot on the relay wall (a), we measure the loss (b) and apply our second stage for resampling. (c) shows the final loss after applying resampling.

Notice that the use of MLP allows minimizing arbitrary losses as long as they are differentiable with respect to our prediction $\tau(x', y', t)$, although $l_2$-norm is most commonly adopted same as in NeRF.

## 4 NEURAL TRANSIENT FIELDS OPTIMIZATION

Our NeTF forward model allows modeling the plenoptic transient field using an MLP. However, data acquired by NeTF are quite different from those in NeRF. In NeRF, a dense set of high resolution images is generally required to produce satisfactory density estimation and view interpolation. Under the dense viewpoint setting, the problem of occlusions is less significant as there will be a sufficient number of views capturing the occluded point to ensure reliable reconstruction. In NeTF, however, our SPAD only captures a sparse set of spots on the wall and an occluded point may be captured only from a very small number of viewpoints (spots). Consequently, occlusion can lead to strong reconstruction artifacts if not handled properly. We develop a two-stage training strategy along with a hierarchical sampling technique to address this sampling bias.

### 4.1 Two-stage Training

We observe that the sampling bias resembles the long-tailed classification problem in machine learning. A conventional solution is to resample the dataset to achieve a more balanced distribution by over-sampling the minority classes [28]. We therefore adopt a two-stage training strategy. We first conduct training using all samples to obtain an initial reconstruction. We then calculate the loss function between the predicted and measured transients at every measuring spot on the relay wall. We observe spots that correspond to a high loss imply undersampling and should incur more samples. We therefore normalize the calculated loss to form a probability density function (PDF). Next, we resample the detection spots using the PDF: a higher loss corresponds a higher PDF and should be more densely sampled. We thus use this sampling scheme to build a new training dataset and then retrain our network to refine reconstruction. The bunny scene (Fig. 5) shows a sample loss map (and therefore resampling density map) with 256x256 sampling spots on the relay wall, with red implying a higher loss (and thus high PDF in subsequent sampling) and blue a low one. Using two stage training, we manage to significantly reduce the loss near the upper left spots on the wall through which the ears of the bunny should be observed. Consequently, our



(a) One-stage    (b) Two-stage    (c) Two-stage with hierarchical sampling
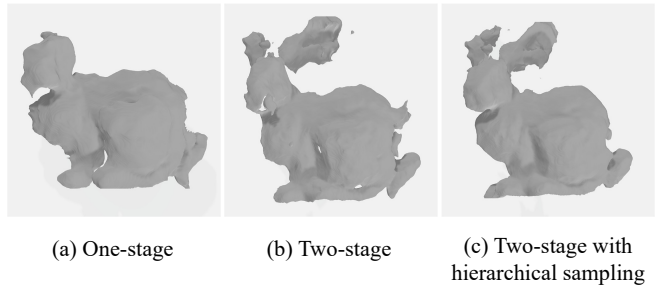
Fig. 6. From left to right: recovered results using our one-stage training, two-stage training without hierarchical sampling, and two-stage with hierarchical sampling. One-stage training, same as most prior art (Fig. 9), fails to recover the second ear of the bunny. Our two-stage schemes manage to recover both ears of the bunny. With hierarchical sampling, NeTF further improves reconstruction with more complete shape and more accurate silhouettes.

two-stage training manages to recover bunny ears largely missing in one-stage training and in prior art (Figs. 6 and 9).

The two-stage training process provides a viable solution to tackle imbalanced sampling for achieving more accurate reconstruction. Fig. 6 demonstrates the corresponding reconstruction results with and without using the second stage. Results with the second stage manage to recover many fine details largely missing from the first stage, e.g., the ear of Bunny and details of the abdomen regions.

### 4.2 Hierarchical Sampling

Denser samples produce higher quality reconstruction. At the same time, they lead to a much higher computation overhead. For example, by uniformly sampling $L$ hemispherical wavefronts at each detection spot and $N^2$ scene points on each $L$, the resulting training process requires a computational complexity of $O(N^2 L)$. We observe that under the confocal setting, spherical wavefronts only intersect with a very small portion of the NLOS scene. These wavefronts tend to converge at specific patches and contribute greatly to the final integral where the contributions from the rest are negligible.

Note that the hierarchical sampling scheme in NeTF is different from NeRF: NeRF calculates the integral along a ray, i.e., using 1D sampling, whereas NeTF on a hemisphere, i.e., using 2D sampling. To make this problem tractable, we develop a *coarse*-to-*fine* sampling scheme. Specifically, we first sample $N_c^2$ uniform scene points in the hemisphere and evaluate our *coarse* network with the estimated PDF $k(\theta, \phi)$. We then employ Metropolis-Hastings algorithm and conditional Gaussian distribution for state transition of Markov chain to produce a *fine* sampling of $N_f$ scene points as $K(\theta_{ij}^f, \phi_{ij}^f)$ along the hemispherical wavefronts that intersect with the NLOS scene. We finally combine the coarse and fine samples $N_c^2 + N_f$ to re-evaluate the reconstruction quality from our *fine* network.

Specifically,

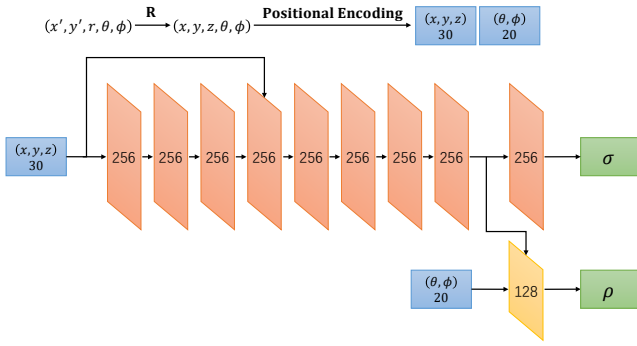$$\tau(x', y', t) = \frac{\tau_c(x', y', t) + \tau_f(x', y', t)}{2} \quad (15)$$

Fig. 7. NeTF network architecture: we adopt an MLP structure analogous to the one used in NeRF. The key differences are (1) NeTF uses ReLU vs. NeRF uses sigmoid and (2) the last four layers in NeRF are simplified to one layer.

where $\tau_c(x', y', t)$ is the integral with the coarse samples $N_c^2$, and $\tau_f(x', y', t)$ is estimated with samples $N_f$ from MCMC, as:

$$\tau_f(x', y', t) = \frac{1}{r^4} \sum_{i,j} \frac{\sigma(r, \theta_{ij}^f, \phi_{ij}^f)\rho(r, \theta_{ij}^f, \phi_{ij}^f)}{K(\theta_{ij}^f, \phi_{ij}^f)} \qquad (16)$$

It is important to note that our hierarchical sampling is intrinsically differentiable (3.3). Previous volume-based methods, e.g., [2], [15], [34], [35], in theory can also apply such a hierarchical sampling technique to refine their reconstruction. In reality, these methods use an explicit volumetric representation with a fixed resolution, making resampling on the hemisphere intractable.

## 5 EXPERIMENTAL RESULTS

In this section, we discuss our NeTF implementation and experimental validations.

### 5.1 MLP Settings

We train the NeTF using an MLP. Fig. 7 shows the structure of our MLP. Analogous to NeRF [39], we construct a fully connected network with nine 256-channel layers, and with one 128-channel layer. We use ReLU activations for all the layers. We transform the spherical coordinates of sampling points into their Cartesian coordinates $(x, y, z, \theta, \phi)$, and feed them into the MLP. We predict the volume density as a function of only position, and the view-dependent reflectance as a function of both position and direction.

We first normalize the spatial coordinates $(x, y, z)$ and the viewing direction $(\theta, \phi)$ to range [-1, 1]. Next, we apply the positional encoding (PE) technique and map each input from 1 dimension onto a 10-dimensional Fourier domain to represent high-frequency variation in geometry and reflectance. Our MLP then processes the coordinates $(x, y, z)$ as inputs with eight 256-channel layers and outputs a 256-dimensional feature vector. Note that we also concatenate $(x, y, z)$ with the fourth layer for skip connection. This feature vector is passed to an additional 256-channel layer and produces $\sigma$. Simultaneously, the feature vector is concatenated with the direction $(\theta, \phi)$ and passed to the 128-channel layer for reflectance $\rho$.

Under the NLOS setting, we consider a batch size of 1 to 4 transients and employ $32 \times 32$ or $64 \times 64$ samples for both uniform sampling $N_c^2$ and MCMC sampling $N_f$ on the hemisphere. We adopt the Adam optimizer [33] with hyperparameters $\beta_1 = 0.9$, and $\epsilon = 1 \times 10^{-7}$. In our experiments, we use a learning rate that begins at $1 \times 10^{-3}$ and decays exponentially to $1 \times 10^{-4}$ through the optimization. The training time of NeTF shares certain similarities to NeRF. In NeRF, the training cost depends on how densely we sample along each ray. In NeTF, it depends on two factors: how densely we sample the radius of the hemisphere (i.e., the number of layers) and how densely we sample each layer/hemisphere. For the Bunny scene, on a single GeForce RTX 3090 GPU the training takes 10 hours using 200 layers with $32 \times 32$ samples on each layer (5 epochs, batchsize 4). The training time quadruples with the same number of layers but at $64 \times 64$ samples.

### 5.2 Validations

We have validated our approach on two public NLOS datasets: a simulated ZNLOS dataset [25], [26], and a real Stanford dataset [15]. ZNLOS consists of multi-bounce transients of synthetic objects that are 0.5 m from the relay wall. The transients have a temporal resolution of 512 time bins with a width of 10 ps and a spatial resolution of $256 \times 256$ pixels. The Stanford dataset captures transients measured in real scenes that are 1.0 m away from the relay wall. The transients in this dataset have a temporal resolution of 512 time bins with a width of 32 ps and a spatial resolution of $512 \times 512$ or $64 \times 64$ pixels. We conduct quantitative and qualitative comparisons between NeTF and the state-of-the-art (SOTA) methods.

**Qualitative Comparisons.** On ZNLOS, we experiment on several simulated hidden objects: Bunny, Lucy, and Indonesian at a spatial resolution of $256 \times 256$ pixels that correspond to an area of size 1 m $\times$ 1 m on the relay wall. All three models are diffuse; Bunny does not contain a floor, but Lucy and Indonesian do. On the Stanford dataset, we have experimented on three real hidden objects with different materials: a diffuse Statue, a glossy Dragon, and a metal Bike. Their spatial resolution is of $512 \times 512$ spots and we down-sample them to $256 \times 256$, same as in [35] ($128 \times 128$ in [14], [15]).

Fig. 8 illustrates our results. Our NeTF outputs a volume density map $\sigma$ and a directional reflectance map $\rho$ of ZNLOS Bunny and Stanford Statue. From these two maps we can produce volumetric albedo, and reconstruct a 3D mesh of hidden objects. By sampling $256 \times 256$ transients, our NeTF produces high quality reconstructions of objects with complex textures (e.g., Lucy). The density and reflectance maps of both Statue and Lucy, and the volumetric albedo produce much less error. We can then apply the Marching Cubes algorithm to further convert the volume to surfaces.

Fig. 9 shows the comparisons with three most broadly adopted volume-based methods [14], [15], [35]. We compare the projected volumes to 2D maps of Indonesian, Lucy and Bunny. Note that the results from [14], [15] correspond to volumetric albedos whereas ours includes both the density and the albedo maps. We also show the normal volume using [35]. The recovered volume maps on these hidden
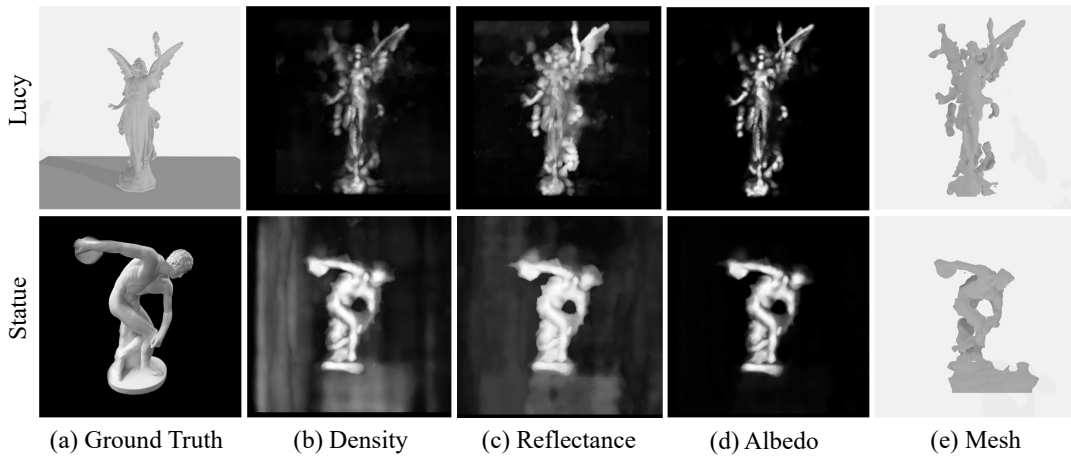
Fig. 8. From left to right: the ground truth, the recovered volume density, reflectance, albedo, and 3D mesh reconstruction using NeTF. Top shows the results on the Lucy model and bottom on the Statue.
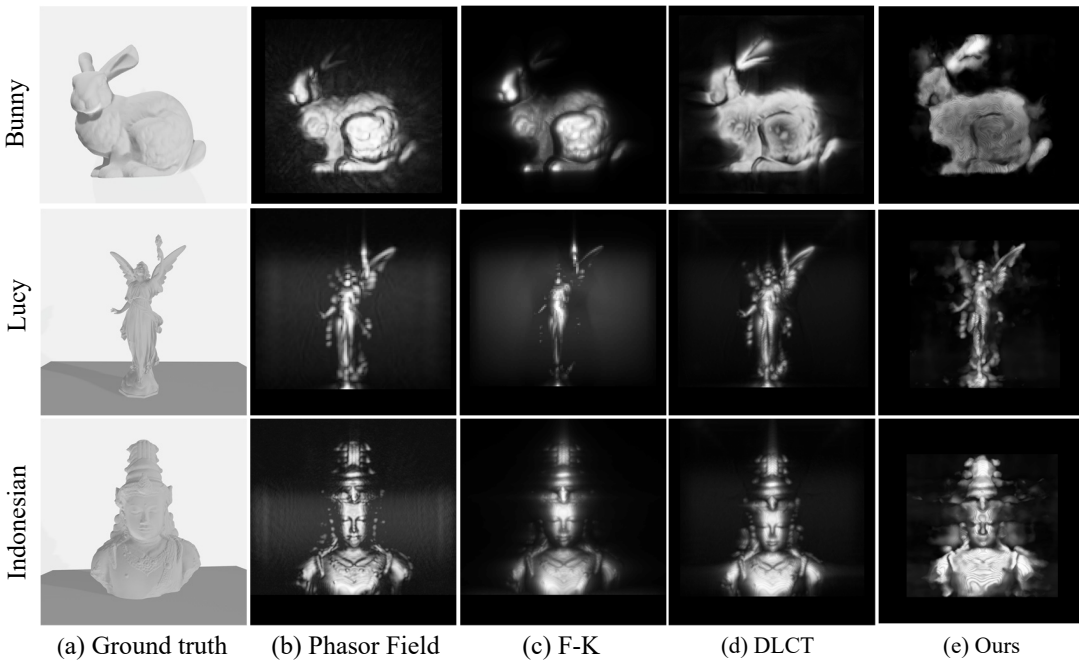


Fig. 9. Comparisons on simulated NLOS data. NeTF achieves comparable reconstructions as SOTA and further manages to recover challenging geometry such as the ear of Bunny, the wing of Lucy and the crown of Indonesian.

objects demonstrate NeTF achieves comparable reconstruction quality to SOTA. NeTF, however, can tackle challenging geometry, e.g., the ear of Bunny, the wing of Lucy and the head of Indonesian that are partially missing using prior art. The Phasor Field technique achieves the best performance on Indonesian but still misses the ear on bunny and wing on Lucy. This implies that such geometry may cast additional challenges to wave-based techniques but can be potentially recovered via volume reconstruction.

DLCT [35] produces comparable results to NeTF on Bunny and Indonesian. On Bunny, both methods manage to acquire the overall geometry yet DLCT misses one ear whereas NeTF captures both. In Fig. 9, and 10, DLCT further uses the mask (silhouettes) of the bunny to obtain the final mesh. The use of the mask can recover the shape (depth) of both ears but the geometry of the second ear is still incorrect. NeTF, in contrast, manages to recover both ears of Bunny. Similar reconstructions can be observed on Lucy. Fig. 10 shows in-depth comparisons between NeTF and DLCT for Bunny on the recovered albedo, normal (density in our results), and mesh reconstruction. Our method preserves fine details but is slightly more noisy, as shown in the depth error. A similar phenomenon is observed on NeRF for multi-view 3D reconstruction where the noise can be potentially filtered.

On the real Stanford dataset, Fig. 11 compares NeTF vs. SOTA for the glossy Dragon, diffuse Statue, and metal Bike. On Dragon and Statue where view-dependency is rel-

(a) Albedo     (b) Density     (c) Mesh     (d) Depth error
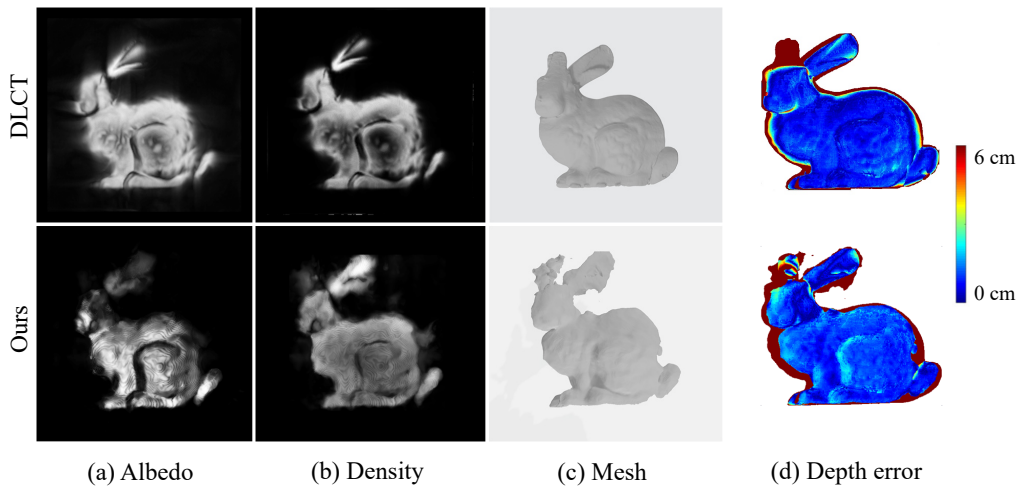
Fig. 10. Comparisons between NeTF and DLCT on the Bunny scene. Both methods manage to acquire the overall geometry yet DLCT misses one ear whereas NeTF captures both. In its own implementation [35], DLCT further uses the mask (silhouettes) of the bunny to further improve reconstruction to obtain the final mesh.



(a) Ground Truth   (b) Phasor Field   (c) F-K   (d) DLCT   (e) Ours
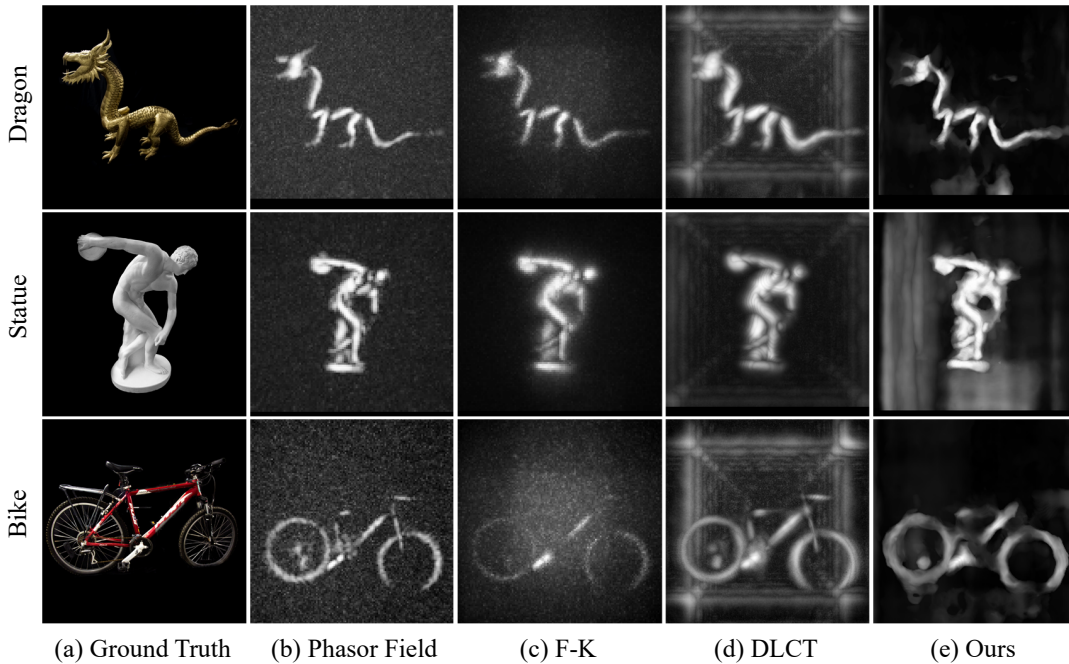
Fig. 11. Comparisons on real NLOS data. (a) The ground truth shows a photography of the hidden object. (b-c) show results using various techniques. The Bike dataset exhibit heterogeneous materials and complex topology and are particularly challenging. NeTF produces comparable reconstructions to SOTA. In particular, same as other neural modeling methods such as NeRF, NeTF reconstruction incurs much lower noise than SOTA.

atively small, NeTF and SOTA produce comparable results, although NeTF manages to better preserve high frequency features such as occluding edges. On the challenging Bike scene, NeTF achieves a similar performance to [14]. For DLCT, the reconstructed mesh exhibits adhesion between different parts, whereas reconstruction produced by NeTF manages to separate these parts.

We have further tested robustness and efficiency of NeTF under non-confocal settings as well as on low resolution transient inputs. Fig. 12 shows the density, albedo and mesh reconstruction. NeTF produces reasonable estimations to the

ground truth and significantly higher quality reconstructions compared with the results from BP and FBP [2]. Fig. 13 shows the NeTF results with down-sampled measurements from $256 \times 256$ spots to $32 \times 32$, $16 \times 16$, $8 \times 8$ and even $4 \times 4$. Even with very sparse sampling spots ($16 \times 16$ and $8 \times 8$), NeTF produces reasonable reconstructions without any prior. To test our approach under the non-confocal setting, we test on two additional objects from ZNLOS, i.e., the letter Z and the bunny, and their transients simulated under non-confocal setups.

To further test the robustness of NeTF vs. SOTA on
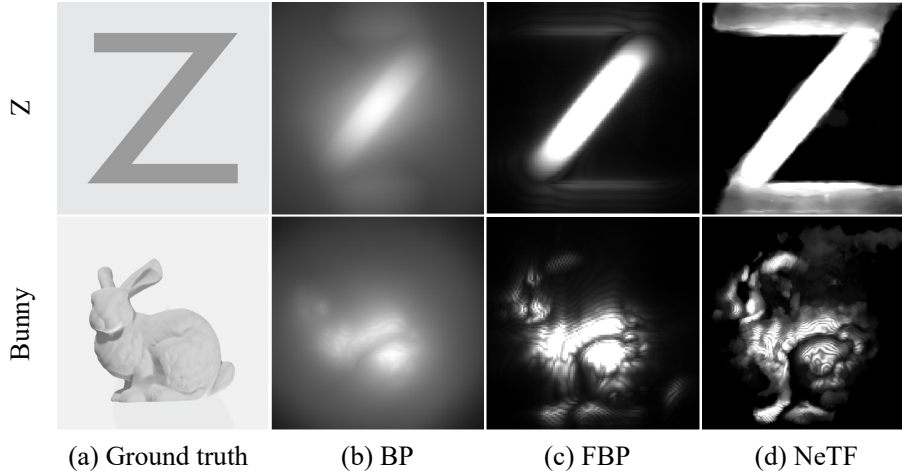
| (a) Ground truth | (b) BP | (c) FBP | (d) NeTF |

Fig. 12. Visual Comparisons of NLOS reconstructions by NeTF and SOTA under the non-confocal setting. NeTF manages to recover clearer silhouettes than SOTA.
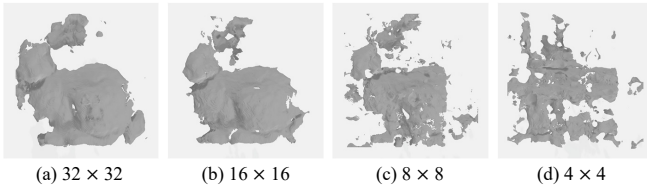


| (a) $32 \times 32$ | (b) $16 \times 16$ | (c) $8 \times 8$ | (d) $4 \times 4$ |

Fig. 13. Low Data Resolution: Mesh reconstructed from simulated transients of Bunny at $32 \times 32$, $16 \times 16$, $8 \times 8$, and $4 \times 4$ spots on the wall. Even at a very low resolution of $8 \times 8$, our NeTF produces reasonable reconstructions.

occlusions, we experiment on a semi-occluded scene from ZNLOS using Eqn. 11. Fig. 14 shows the frontal and top-viewed albedo maps of the reconstruction. Phasor Field is most sensitive to occlusions whereas DLCT and F-K can only recover one plane at a high accuracy. NeTF produces sharper edges of both front and back planes.

**Quantitative Comparisons.** Table 1 and 2 show that NeTF achieves comparable accuracy as the state-of-the-art (SOTA) in terms of Mean Absolute Error (MAE), demonstrating the feasibility and efficacy of deep neural networks in NLOS under both confocal and non-confocal settings. Under the MAE metric, the gain using NeTF does not seem significant. However, MAE does not fully reflect the reconstruction quality: our experiments have further revealed that NeTF can more robustly handle silhouettes and semi-occlusions, as shown in Fig. 12 and 14.

## 6 CONCLUSION AND FUTURE WORK

We have presented a novel neural modeling framework called the Neural Transient Field (NeTF) for NLOS imaging. Similar to the recent Neural Radiance Field that seeks to use a multi-layer perception (MLP) to represent the 5D radiance function, NeTF recovers the 5D transient function in both spatial location and direction. Different from NeRF, the input training data are parametrized on the spherical wavefronts in NeTF rather than along lines (rays) as in NeRF. We have hence formulated the NLOS process under spherical

TABLE 1
Reconstruction error using NeTF vs. SOTA on three confocal NLOS datasets measured by MAE. NeTF achieves comparable performance in MAE. Notice though MAE does not fully reflect the reconstruction quality: for example, Phasor Field produces the highest MAE on Indonesian, indicating lowest reconstruction quality; yet it manages to recover many fine details largely missing in F-K and DLCT, as shown in Fig. 9.

| MAE | Bunny | Lucy | Indonesian |
|---|---|---|---|
| Phasor Field | 2.89 cm | 1.36 cm | 1.69 cm |
| F-K | 2.43 cm | 2.05 cm | 0.61 cm |
| DLCT | 2.38 cm | 0.23 cm | 0.30 cm |
| NeTF | 2.65 cm | 1.05 cm | 0.31 cm |

TABLE 2
Reconstruction error using NeTF vs. SOTA on two non-confocal NLOS datasets measured by MAE. Same as in Table 1, we observe that low MAE does not sufficiently reflect reconstruction quality, e.g., on the Z letter scene, NeTF performs slightly worse than FBP in MAE but better preserves the silhouettes, as shown in Fig. 12.

| MAE | Bunny (non-confocal) | Z (non-confocal) |
|---|---|---|
| BP | 7.02 cm | 3.21 cm |
| FBP | 3.77 cm | 0.46 cm |
| NeTF | 7.45 cm | 0.60 cm |

coordinates, analogous to volume rendering under Cartesian coordinates. Another unique characteristic of our NeTF solution is the use of Markov chain Monte Carlo (MCMC) to account for sparse and unbalanced sampling in NeTF. MCMC enables more reliable volume density estimation and produces more accurate shape estimation by recovering details caused by occlusions and non-uniform albedo. Our experiments on both synthetic and real data demonstrate the benefits of NeTF over SOTA in both robustness and accuracy.

Same as NeRF, the final reconstruction of the hidden scene corresponds to a 3D density volume, implicitly represented by the MLP. Recovering the actual shape requires mapping the volume to surfaces, e.g., by thresholding followed by Marching Cubes. Such brute-force implementa-
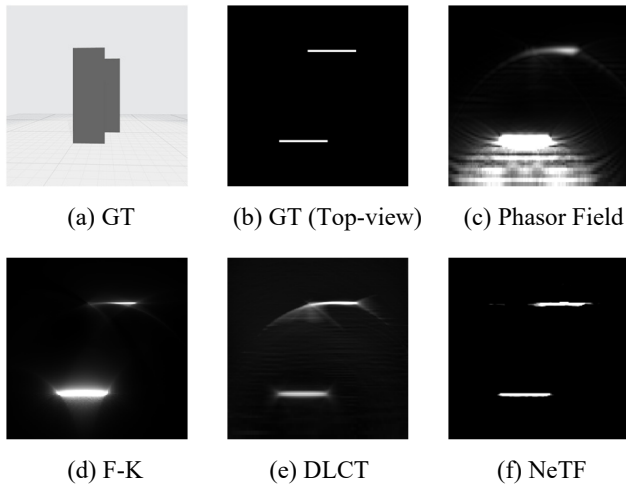
Fig. 14. Visual comparisons of NLOS reconstructions on the semi-occluded scene. (a) and (b) show the frontal and top-down views. Closest to NeTF is DLCT, which manages to recover the front plane but produces high errors on the back plane.

tions may lead to noise on smooth surfaces. There are a number of emerging neural modeling techniques that can potentially provide smooth reconstructions, by imposing shape priors [40]. In general, learning-based techniques (including NeRF and NeTF), in their current forms, are still substantially more computationally expensive than previous optimization techniques, although we observe a large number of emerging acceleration schemes. More importantly, NeTF demonstrates that deep learning provides an alternative and potentially feasible solution to a broader class of inverse imaging problems. There are also several acceleration schemes, e.g., using results from SOTA to initialize the network and then conduct training. It is our immediate future work to investigate how to integrate such approaches into our NeTF framework. Our current approach does not separately treat the confocal and non-confocal setups. Analogous to multi-view stereo vs. photometric stereo, it may be possible to tailor solutions such as [41] on top of NeTF to separately handle different settings.
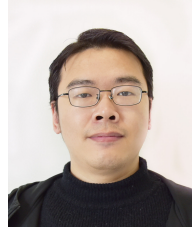
## ACKNOWLEDGMENTS

## REFERENCES

[1] Otkrist Gupta, Thomas Willwacher, Andreas Velten, Ashok Veeraraghavan, and Ramesh Raskar. Reconstruction of hidden 3d shapes using diffuse reflections. Optics express, 20(17):19096–19108, 2012.

[2] Andreas Velten, Thomas Willwacher, Otkrist Gupta, Ashok Veeraraghavan, Moungi G Bawendi, and Ramesh Raskar. Recovering three-dimensional shape around a corner using ultrafast time-of-flight imaging. Nature Communications 1747, 2012.

[3] Andreas Velten, Di Wu, Adrian Jarabo, Belen Masia, Christopher Barsi, Chinmaya Joshi, Everett Lawson, Moungi Bawendi, Diego Gutierrez, and Ramesh Raskar. Femto-photography: Capturing and Visualizing the Propagation of Light. ACM Transactions on Graphics, 32(4):44:1–44:8, 2013.

[4] Mariko Isogawa, Ye Yuan, Matthew O'Toole, and Kris Kitani. Optical Non-Line-of-Sight Physics-based 3D Human Pose Estimation, IEEE International Conference on Computer Vision and Pattern Recognition (CVPR), 2020.

[5] Genevieve Gariepy, Nikola Krstaji´c, Robert Henderson,Chunyong Li, Robert R Thomson, Gerald S Buller, Barmak Heshmat, Ramesh Raskar, Jonathan Leach, and Daniele Faccio. Single-photon sensitive light-in-fight imaging. Nature communications, 6:6021, 2015.

[6] Dongeek Shin, Feihu Xu, Dheera Venkatraman, Rudi Lussana, Federica Villa, Franco Zappa, Vivek K Goyal, Franco NC Wong, and Jeffrey H Shapiro. Photon-efficient imaging with a single-photon camera. Nature communications, 7:12046, 2016.

[7] Matthew O'Toole, Felix Heide, David B Lindell, Kai Zang, Steven Diamond, and Gordon Wetzstein. Reconstructing transient images from single-photon sensors. IEEE International Conference on Computer Vision and Pattern Recognition (CVPR), 2017.

[8] Victor Arellano, Diego Gutierrez, and Adrian Jarabo. Fast backprojection for non-line of sight reconstruction. Optics Express, 25(10):11574–11583, 2017.

[9] Feihu Xu, Gal Shulkind, Christos Thrampoulidis, Jeffrey H Shapiro, Antonio Torralba, Franco NC Wong, and Gregory W Wornell. Revealing hidden scenes by photonefficient occlusion-based opportunistic active imaging. Optics Express, 26(8):9945–9962, 2018.

[10] Tomohiro Maeda, Yiqin Wang, Ramesh Raskar, and Achuta Kadambi. Thermal Non-Line-of-Sight Imaging, IEEE International Conference on Computational Photography (ICCP), 2019.

[11] David B Lindell, Gordon Wetzstein, and Vladlen Koltun. Acoustic non-line-of-sight imaging. IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2019.

[12] Ioannis Gkioulekas, Anat Levin, Fr´edo Durand, and Todd Zickler. Micron-scale light transport decomposition using interferometry. ACM Transactions on Graphics (TOG), 34(4):37, 2015.

[13] Piergiorgio Caramazza1, Alessandro Boccolini, Daniel Buschek, Matthias Hullin, Catherine F. Higham , Robert Henderson, Roderick Murray-Smith and Daniele Faccio. Neural network identification of people hidden from view with a single-pixel, single-photon detector, Scientific Reports, 8:11945, 2018.

[14] Xiaochun Liu, Ibon Guillen, Marco La Manna, Ji Hyun Nam, Syed Azer Reza, Toan Huu Le, Adrian Jarabo, Diego Gutierrez, and Andreas Velten. Non-line-of-sight imaging using phasor-field virtual wave optics. Nature, 572:620–623, 2019.

[15] David B Lindell, Gordon Wetzstein, and Matthew O'Toole. Wavebased non-line-of-sight imaging using fast fk migration. ACM Transactions on Graphics (TOG), 38(4):116, 2019.

[16] Shumian Xin, Sotiris Nousias, Kiriakos N Kutulakos, Aswin C Sankaranarayanan, Srinivasa G Narasimhan, and Ioannis Gkioulekas. A theory of fermat paths for non-lineof-sight shape reconstruction. IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2019.

[17] Adrian Jarabo, Belen Masia, Julio Marco, and Diego Gutierrez. Recent advances in transient imaging: A computer graphics and vision perspective. Visual Informatics, 1(1):65–79, 2017.

[18] Chia-Yin Tsai, Aswin C Sankaranarayanan, and Ioannis Gkioulekas. Beyond volumetric albedo–a surface optimization framework for non-line-of-sight imaging. IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2019.

[19] Ahmed Kirmani, Tyler Hutchison, James Davis, and Ramesh Raskar. Looking around the corner using ultrafast transient imaging. International journal of computer vision, 95(1):13–28, 2011.

[20] Ahmed Kirmani, Tyler Hutchison, James Davis, and Ramesh Raskar. Looking around the corner using transient imaging. International Conference on Computer Vision, 2009.

[21] Yoann Altmann, Stephen McLaughlin, Miles J. Padgett, Vivek K. Goyal, Alfred O. Hero, and Daniele Faccio. Quantum-inspired computational imaging. Science, 361(6403): eaat2298, 2018.

[22] Felix Heide, Matthew O'Toole, Kai Zang, David B. Lindell, Steven Diamond, and Gordon Wetzstein. Non-line-of-sight imaging with partial occluders and surface normals. ACM Transactions on Graphics, 38(3):22:1–22:10, 2019.

[23] Syed Azer Reza, Marco La Manna, Sebastian Bauer, and Andreas Velten. Phasor field waves: A Huygens-like light transport model for non-line-of-sight imaging applications, Optics Express, 27(20): 29380-29400, 2019.

[24] Xiaochun Liu, Sebastian Bauer, Andreas Velten. Phasor field

diffraction based reconstruction for fast non-line-of-sight imaging systems, Nature communications, 11(1): 1-13, 2020.

[25] Adrian Jarabo, Julio Marco, Adolfo Muñoz, Raul Buisan, Wojciech Jarosz, and Diego Gutierrez. A framework for transient rendering. ACM Transactions on Graphics. 33(6): 177, 2014.

[26] Miguel Galindo, Julio Marco, Matthew O'Toole, Gordon Wetzstein, Diego Gutierrez, and Adrian Jarabo. A dataset for benchmarking time-resolved non-line-of-sight imaging, IEEE International Conference on Computational Photography (ICCP), 2019.

[27] Mariko Isogawa, Dorian Chan, Ye Yuan, Kris Kitani, and Matthew O'Toole. Efficient Non-Line-of-Sight Imaging from Transient Sinograms, European Conference on Computer Vision (ECCV), 2020.

[28] Bingyi Kang, Saining Xie, Marcus Rohrbach, Zhicheng Yan, Albert Gordo, Jiashi Feng, Yannis Kalantidis. Decoupling Representation and Classifier for Long-Tailed Recognition, International Conference on Learning Representations (ICLR), 2020.

[29] Marco La Manna, Fiona Kine, Eric Breitbach, Jonathan Jackson, Talha Sultan, and Andreas Velten. Error backprojection algorithms for non-line-of-sight imaging. IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI), 41(7): 1615-1626, 2019.

[30] James T. Kajiya, Brian P Von Herzen. Ray tracing volume densities, ACM Transactions on Graphics (TOG), 18(3):165-174, 1984.

[31] Mauro Buttafava, Jessica Zeman, Alberto Tosi, Kevin Eliceiri, and Andreas Velten. Non-line-of-sight imaging using a time-gated single photon avalanche diode. Optics express, 23(16):20997–21011, 2015.

[32] Chia-Yin Tsai, Kiriakos N Kutulakos, Srinivasa G Narasimhan, and Aswin C Sankaranarayanan. The geometry of first-returning photons for non-line-of-sight imaging. IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2017.

[33] Diederik P. Kingma, Jimmy Ba. Adam: A method for stochastic optimization. International Conference on Learning Representations (ICLR), 2020.

[34] Matthew O'Toole, David B. Lindell, and Gordon Wetzstein. Confocal non-line-of-sight imaging based on the light-cone transform. Nature, 555(7696):338–341, 2018.

[35] Sean I. Young, David B. Lindell, Bernd Girod, David Taubman, Gordon Wetzstein. Non-Line-of-Sight surface reconstruction Using the directional light-Cone transform, IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2020.

[36] Julian Iseringhausen and Matthias B Hullin. Non-line-of-sight reconstruction using efficient transient rendering. ACM Transactions on Graphics, 39(1): 8, 2020.

[37] Daniele Faccio, Andreas Velten and Gordon Wetzstein. Non-line-of-sight imaging, Nature Reviews Physics, 2:318-327, 2020.

[38] Byeongjoo Ahn, Akshat Dave, Ashok Veeraraghavan, Ioannis Gkioulekas, and Aswin C. Sankaranarayanan. Convolutional approximations to the general non-Line-of-sight imaging operator. IEEE International Conference on Computer Vision, 2019.

[39] Ben Mildenhall, Pratul P. Srinivasan, Matthew Tancik, Jonathan T. Barron, Ravi Ramamoorthi and Ren Ng. NeRF: Representing Scenes as Neural Radiance Fields for View Synthesis. European Conference on Computer Vision (ECCV), 2020.

[40] Lior Yariv, Yoni Kasten, Dror Moran, Meirav Galun, Matan Atzmon, Ronen Basri, Yaron Lipman. Multiview neural surface reconstruction by disentangling geometry and appearance. Neural Information Processing Systems (NeurIPS), 2020.

[41] Sai Bi, Zexiang Xu, Pratul Srinivasan, Ben Mildenhall, Kalyan Sunkavalli, Miloš Hašan, Yannick Hold-Geoffroy, David Kriegman, Ravi Ramamoorthi. Neural reflectance fields for appearance acquisition. arXiv preprint arXiv:2008.03824, 2020.

**Zi Wang** received the BS degree from the Beihang University, Beijing, China, in 2019. He is currently working toward the master's degree at ShanghaiTech University, Shanghai, China. His research interests include non-line-of-sight imaging and computational imaging.



**Ping Liu** received the BS degree from Central South University, Changsha, China, in 2020. He is working toward the master's degree in computer vision at ShanghaiTech University, Shanghai, China. His research interests include mainly in computational imaging and computer vision, especially on non-line-of-sight imaging.



**Zhengqing Pan** received the BS degree from Shanghaitech University, Shanghai, China, in 2019. He is currently working toward the master's degree at ShanghaiTech University, Shanghai, China. His research interests include single photon imaging and non-line-of-sight imaging.



**Ruiqian Li** is currently working toward the BS degree at ShanghaiTech University, Shanghai, China. His research interests include single photon imaging, noise model and denoising algorithm.
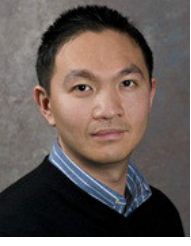


**Tian Gao** He is currently working toward the BS degree at ShanghaiTech University, Shanghai, China. His research interests include single photon imaging and computational imaging.



**Siyuan Shen** received the BS degree from Shanghaitech University, Shanghai, China, in 2019. He is currently working toward the master's degree at ShanghaiTech University, Shanghai, China. His research interests include non-line-of-sight imaging, computer vision and computational imaging.



**Shiying Li** received her MS and PhD degrees in computer science from the Nara Institute of Science and Technology (NAIST), Nara, Japan, in 2004 and 2007. She worked as a post-doctoral research fellow in the Tohoku University from 2007 to 2008, Sendai, Japan. She was an associate professor with the Hunan University from 2009 to 2017, Changsha, Hunan, and has been an associate researcher with the School of Information Science and Technology, ShanghaiTech University, Shanghai, China. Her research interests include computational imaging, computer vision and graphics.

**Jingyi Yu** received BS 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.

# APPENDIX A
## FORMULATING LCT VIA NeTF

We show LCT can be formulated as a simplified NeTF model. We first rewrite the forward model Eqn. 12 under triple integral with the Dirac delta function that correlates time of flight $t$ with distance $r$:

$$\tau(x',y',t) = \Gamma_0 \iiint_\Omega \frac{\sin\theta}{r^2} \sigma(r,\theta,\phi)\rho(r,\theta,\phi)\delta(r-\frac{ct}{2})\,\mathrm{d}r\,\mathrm{d}\theta\,\mathrm{d}\phi \quad (17)$$

where the integral domain $\Omega$ is defined under the spherical coordinates. Notice Eqn. 17 is consistent with the light-cone transform (LCT) model [34]: we can rewrite it under the Cartesian coordinates where $\mathrm{d}x\,\mathrm{d}y\,\mathrm{d}z = r^2\sin\theta\,\mathrm{d}r\,\mathrm{d}\theta\,\mathrm{d}\phi$ as:

$$\tau(x',y',t) = 2\Gamma_0 \iiint_\Omega \frac{1}{r^4}\sigma(x,y,z)\rho(x,y,z,\theta,\phi)\cdot$$
$$\delta(2\sqrt{(x-x')^2+(y-y')^2+z^2}-ct)\,\mathrm{d}x\,\mathrm{d}y\,\mathrm{d}z \quad (18)$$

If we further assume diffuse and isotropic albedo as $\rho_{\text{iso}}(x,y,z) = \sigma(x,y,z)\rho(x,y,z,\theta,\phi)$, Eqn. 18 degenerates to the LCT model (Eqn. 4 with $g=1$).

# APPENDIX B
## NON-CONFOCAL NeTF

Under the non-confocal setting, we set out to formulate the transient in terms of semi-ellipsoids with foci at the illumination and detection spots $P$ and $P'$ on the relay wall, as shown in Fig. 15. Given a scene point $Q$, assume $r_1$ and $r_2$ correspond to the distance from $P$ to $Q$ and $Q$ to $P'$, respectively. Following the same derivations of Eqns. 6, 7, and 8 under the confocal setting, we first compute the energy (transient) received at $P'$ from the location $Q$ as:

$$E_{P'} = \frac{\Gamma}{r_2^2}\sigma(Q)\rho(Q,P,P')\exp\left(-A\int_\Upsilon \sigma(s)\,\mathrm{d}s\right)\mathrm{d}\Omega \quad (19)$$

where $\Gamma = Aar_0^2 E_p/\pi$. $\exp\left(-A\int_\Upsilon \sigma(s)\,\mathrm{d}s\right)$ corresponds to the attenuation coefficient along optical path $\Upsilon : P \to Q \to P'$ with length $r_1 + r_2 = ct$.

To compute the complete transient received at $P'$ from $P$, recall $P'$ should be radiated by all points lying on a semi-ellipsoid $E$ with the foci $P, P'$, semi-major axis of length
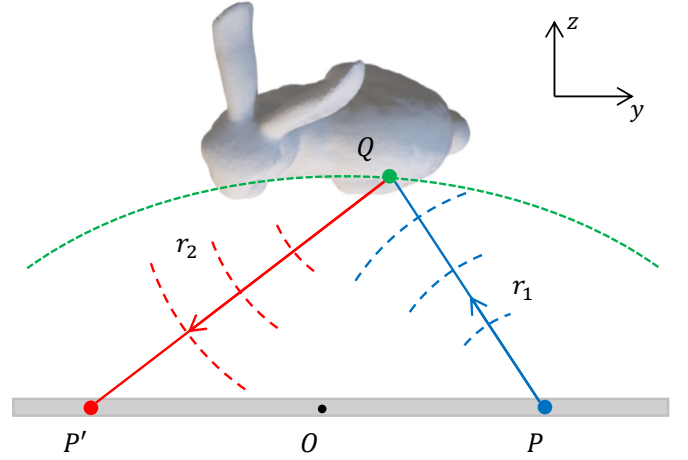


Fig. 15. Non-confocal NLOS imaging: The transient process from an illumination spot $P$ to an NLOS point $Q$ and from $Q$ to the detection spot $P'$ can be formulated under ellipsoidal coordinates.

$\alpha = ct/2$, focal length $\gamma = |\overrightarrow{OP} - \overrightarrow{OP'}|$, and the eccentricity $e = \gamma/\alpha$. For simplicity, we can set up the coordinate system so that $P$ and $P'$ are symmetric about origin $O$ and $\overline{PP'}$ parallel to $y$-axis. We can thus compute transient as:

$$\tau(P,P',t) = \iint_E E_{P'}\,\mathrm{d}\Omega \quad (20)$$

Since Eqn. 19 is integrated on the semi-ellipsoid $E$ but under spherical coordinates centered at $P$, we need to rewrite $E$ under ellipsoidal coordinates with foci $P$ and $P'$. Specifically, we first represent the ellipsoid in terms of $r_1$ and $\theta$ as:

$$r_1 = \frac{\alpha(1-e^2)}{1-e\cos\theta} \quad (21)$$

Eqn. 20 transforms to:

$$\tau(P,P',t) = \iiint_\Omega E_{P'}\delta\left(r_1 - \frac{\alpha(1-e^2)}{1-e\cos\theta}\right)\mathrm{d}r_1\,\mathrm{d}\Omega \quad (22)$$

Next, we transform spherical coordinates $(r_1,\theta,\phi)$ to ellipsoidal coordinates $(\mu,\nu,\varphi)$ as:

$$r_1\sin\theta\cos\phi = \gamma\sinh\mu\sin\nu\cos\varphi$$
$$r_1\sin\theta\sin\phi = \gamma\sinh\mu\sin\nu\sin\varphi \quad (23)$$
$$r_1\cos\theta = \gamma\cosh\mu\cos\nu$$

Recall that the Jacobian $J$ from the Cartesian to ellipsoidal coordinates are:

$$J = \frac{\mathrm{d}x\,\mathrm{d}y\,\mathrm{d}z}{\mathrm{d}\mu\,\mathrm{d}\nu\,\mathrm{d}\varphi} = \gamma^3\sinh\mu\sin\nu(\sinh^2\mu + \sin^2\nu) \quad (24)$$

We can map between spherical coordinates to ellipsoidal coordinates via $J$ as:

$$\mathrm{d}x\,\mathrm{d}y\,\mathrm{d}z = r_1^2\sin\theta\,\mathrm{d}r_1\,\mathrm{d}\theta\,\mathrm{d}\phi = r_1^2\,\mathrm{d}r_1\,\mathrm{d}\Omega = J\,\mathrm{d}\mu\,\mathrm{d}\nu\,\mathrm{d}\varphi \quad (25)$$

Substituting Eqn. 25 into Eqn. 22, we obtain the transient under the ellipsoidal coordinate system as:

$$\tau(P, P', t) = \iiint_{\Omega} \frac{1}{r_1^2} E_{P'} \delta\left(2\gamma \cosh\mu - ct\right) J \,\mathrm{d}\mu \,\mathrm{d}\nu \,\mathrm{d}\varphi$$

(26)

Notice that with a fixed $t$, we can find the corresponding $\mu$ for non-zero $\delta$ so that the triple integrate can be simplified to double integral in only $\nu$ and $\varphi$. In addition, if we further discard the attenuation term in $E_{P'}$, we can simplify the transient to:

$$\tau(P, P', t) = \Gamma_0 \iint_{E} \frac{J}{r_1^2 r_2^2} \sigma(\mu, \nu, \varphi) \rho(\mu, \nu, \varphi, P, P') \,\mathrm{d}\nu \,\mathrm{d}\varphi$$

(27)

where $\mu = arccosh(ct/2\gamma)$. In our non-confocal NeTF implementation (as in Sec. 3.3), we set out to solve Eqn. 27. It is important to note though that a downside of discarding attenuation ignores occlusions.