

Towards Non-Line-of-Sight Photography

Jiayong Peng¹, Fangzhou Mu², Ji Hyun Nam², Siddeshwar Raghavan²,
Yin Li², Andreas Velten², and Zhiwei Xiong¹

¹University of Science and Technology of China, ²University of Wisconsin Madison

jiayong@mail.ustc.edu.cn, {fmu2, jnam26, sraghavan7, yin.li, velten}@wisc.edu, zwxiong@ustc.edu.cn

Abstract

Non-line-of-sight (NLOS) imaging is based on capturing the multi-bounce indirect reflections from the hidden objects. Active NLOS imaging systems rely on the capture of the time of flight of light through the scene, and have shown great promise for the accurate and robust reconstruction of hidden scenes without the need for specialized scene setups and prior assumptions. Despite that existing methods can reconstruct 3D geometries of the hidden scene with excellent depth resolution, accurately recovering object textures and appearance with high lateral resolution remains an challenging problem. In this work, we propose a new problem formulation, called NLOS photography, to specifically address this deficiency. Rather than performing an intermediate estimate of the 3D scene geometry, our method follows a data-driven approach and directly reconstructs 2D images of a NLOS scene that closely resemble the pictures taken with a conventional camera from the location of the relay wall. This formulation largely simplifies the challenging reconstruction problem by bypassing the explicit modeling of 3D geometry, and enables the learning of a deep model with a relatively small training dataset. The results are NLOS reconstructions of unprecedented lateral resolution and image quality.

1. Introduction

Imaging methods mainly focus on recovering information in line-of-sight (LOS) scenarios, where there are no obstacles on the direct light path between the target and the camera. In non-line-of-sight (NLOS) imaging, the scene is hidden beyond the direct line of the cameras' sight and the light from the scene is scattered by a diffuse relay surface with dramatic loss in angular information. However, such complicated imaging processing puts much difficulties on the reconstruction methods.

Existing reconstruction methods can be roughly concluded into two categories: passive methods and active methods. Passive methods seek to perform a reconstruction

solely with ambient light and remain very challenging for general scenes. Active methods illuminate the scene with a controlled light source, usually a laser, and reconstruct from these active transient measurements. In addition, some high-level applications, such as object tracking and active recognition, have also been demonstrated using continuous or steady state light sources. While the most effective and robust methods are active transient methods that use pulsed illumination sources and fast detectors to measure the time of flight (ToF) measurements through the scenes. And all these reconstruction methods first try to reconstruct a 3D volume and then project it to a 2D coordinate, which not only requires massive computational demands but also relies on 3D cues.

Generating a 3D geometry that would project to a photorealistic 2D image is a challenging problem even in line of sight rendering. It requires accurate models of not only the scene geometry, but also scene BRDF and subsurface properties. A rendering approach would have to take into account multiple scattering in the scene that affects global illumination, and subsurface scattering. Rendering systems typically rely on approximate models of light transport to create the correct image appearance in reasonable time, but fail to incorporate the correct transient light transport seen in the NLOS measurement. To the best of our knowledge, rendering methods that can create photorealistic renderings of objects while also correctly modeling NLOS light transport do not exist, making the approach or reconstructing a 3D model with the hopes of projecting it down to an accurate 2D image of the scene extremely challenging.

While existing methods demonstrate exceptional depth resolution of few millimeters or better and therefore outperform most commercial line of sight 3D sensors, the lateral resolution of the NLOS reconstruction, is quite low with about 2 cm at a distance of 2 meters from the relay surface or an angular resolution of 0.6 degrees. Such a resolution is a little over an order of magnitude worse than the human eye. Further, prior methods also fail to correctly recover the texture and visual appearance of surfaces resulting in a reconstruction that appears very different from an actual

image of the hidden scene. Finally, existing data-driven methods are trained on simulated data using fast rendering algorithms, which approximate or omit potentially important aspects of real world light transport, such as multiple bounces, realistic BRDF, and background / sensor noise.

To bridge the gap, we present a new problem formulation, called NLOS photography. Our goal is to reconstruct intensity images resembling photographs captured by a LOS camera, based on the NLOS transient measurements using a data-driven deep model. The key difference between our work and previous methods is the focus on directly reconstructing high quality 2D images with only 2D supervisions, rather than recovering 3D volumes or surfaces and projecting them back to 2D images as previous methods do. This ensures several advantages for our model. The model is easy to train and requires less training samples; it is easier to obtain 2D ground-truth labels than that of 3D labels for training. Specifically, our model builds on a deep neural network previously developed for LOS ToF imaging [33], with several modifications tailored for NLOS imaging. Under our formulation and using our model, we demonstrate the first result on reconstructing intensity images of hidden objects using a ToF NLOS system, at a quality and resolution approaching that of a low end LOS camera. Our method thus provides significantly higher visual quality and lateral resolution than existing NLOS reconstruction methods.

To train and evaluate our model, we collect a new large-scale NLOS dataset of real-world scenarios with the help of a fast NLOS imaging system. The NLOS imaging system follows a non-confocal setup, and consists of an ultrafast pulsed laser, two single-photon avalanche diode (SPAD) arrays, and a conventional RGB camera embedded in the relay wall. The system captures synchronized pairs of NLOS measurement and LOS images in a fast manner (4 fps), making it possible to create a large-scale NLOS dataset. Our dataset consists of 400 transient measurements of distinct indoor scenes located around 1m away from a diffuse wall under different illumination conditions (with and without room light). To the best of our knowledge, it is the largest NLOS dataset with two orders of magnitude more samples than previous ones [17, 9, 22, 24]. Apart from the scale, our dataset provides a diverse set of scenes with free object shape, pose and appearance. We hope that our dataset can help to promote the development of data-driven methods in NLOS imaging.

Our contributions are summarized into four folds:

- *Problem formulation.* We introduce a new problem formulation of NLOS photography that is to recover intensity images resembling photographs captured by a LOS camera without explicit modeling of 3D geometry. The high quality reconstructions can be achieved with only 2D supervisions, which we identify as the

most significant deficits of the prior body of work.

- *High quality NLOS reconstructions.* We present the first result of high quality reconstructions of intensity images from NLOS imaging using a deep neural network. Our results demonstrate significantly higher visual quality and lateral resolution than existing methods, and approach the quality of conventional LOS cameras.
- *Learning from real-world data.* A key ingredient of our method is training using captured real-world data, rather than simulated synthetic data. This strategy rules out potential issues with the accuracy of the simulation, and allows us to generate training data faster than prior art.
- *Dataset.* We collect the first large-scale NLOS dataset of real-world scenarios. The dataset contains NLOS measurements coupled with synchronized intensity images. We will make the dataset publicly available to facilitate future research in NLOS imaging.

2. Related Work

Our work builds on a large body of relevant works on NLOS imaging and reconstruction methods. We now present a brief survey of these methods.

2.1. NLOS Imaging

Kirmanian *et al.* [16] propose the first transient imaging framework for inferring the geometry of hidden objects in NLOS settings. Pandharker *et al.* [30] first estimate the size and motion of the moving NLOS objects in cluttered environments by exploiting the relative times of arrival after reflection and analyzing the multipath. Velten *et al.* [39] first reconstruct the hidden objects from a transient NLOS imaging system with a ultra-fast laser and streak camera. Build on these influential works, the subsequent researches are focusing on designing various NLOS imaging systems as well as the reconstruction methods [25].

2.2. NLOS Imaging Systems and Datasets

The existing non-line-of-sight imaging systems can rely on the ambient light [4, 26, 3, 36, 44] or the active illuminations [16, 39, 5, 29, 13, 40, 21]. To name a few, Bouman *et al.* [4] demonstrate the possibility of using subtle spatial-temporal radiance variations to track the hidden objects behind the wall. Maeda *et al.* [26] propose a long-wave infrared NLOS imaging system, where the heat radiation is emitted from the human body simplifying the NLOS problem to a single bounce reflection. Batarseh *et al.* [3] capture the NLOS information with conventional RGB cameras and retrieve the geometric information of hidden objects with spatial coherence in the measurements. Different

from the systems that rely on the ambient light, the systems with active illuminations emit a short light pulse towards the hidden scenes and receive the back-reflected light after three bounces. For example, Kirmani *et al.* [16] build a NLOS imaging system with a femtosecond laser and an ultrafast photo detector array. After that, Velten *et al.* [39] replace the photo detector array with a streak camera, which achieves sub-millimeter depth precision. As a simplification, O’Toole *et al.* [29] propose a new imaging configuration, named confocal NLOS, where the detector and the light source share a common light path. Different from the existing NLOS techniques that densely scan the across the entire relay wall, Isogawa *et al.* [13] design an efficient scanning strategy, called circular confocal NLOS, that reduces both acquisition time and computational requirements. Wu *et al.* [40] introduce a long-range NLOS imaging system by increasing the imaging range from meters to kilometers. Instead of using light pulses, Lindell *et al.* [21] adopt the sound wave to sense the hidden objects. The existing active NLOS imaging systems generally require long acquisition time to capture one measurement, which hinders further applications in real-world scenarios. However, our imaging system can capture both transient measurements and the corresponding RGB images in a fast manner (4 fps), which makes building a large-scale dataset possible.

Very little work is proposed to provide a dataset for NLOS imaging. Jarabo *et al.* [14] introduce a framework for transient rendering, which has been adopted for NLOS data simulation in many literatures [22, 12, 19]. For example, Klein *et al.* [17] introduce a synthetic dataset for various NLOS reconstruction problems yet with only 16 data samples. Galindo *et al.* [9] build a synthetic NLOS transient dataset with varied complexities for benchmarking NLOS reconstructions. However, three major problems exist in the existing datasets: 1. the data are mainly simulated, which can be inaccurate for factors such as multiple bounces, realistic BRDF, and background / sensor noise; 2. the simulated process is time-consuming, usually taking several minutes to hours for one scene; 3. the amount of the reconstructible scenes are limited in a dozen. As a contrast, our dataset is captured by a real-world imaging system on 400 different scenes with both NLOS transient measurements and ground-truth intensities in a fast manner (4 fps), which facilitates deep-learning based NLOS reconstructions as well as many other tasks, *e.g.*, classification.

2.3. NLOS Reconstruction Methods

Plenty of methods in reconstructing the hidden objects from transient measurement have been developed [1, 39, 20, 12, 16, 18, 27, 32, 38, 41, 11, 29, 22, 24, 23, 45]. As a precursory work in this field, Velten *et al.* [39] propose a filtered back-projection method to recover the hidden ob-

jects from transient measurements. Arellano *et al.* [1] propose a new back-projection technique by optimizing voxelization of space-time manifolds in a GPU, which is up to a thousand times faster than previous methods. Heide *et al.* [11] formulate NLOS image formation as a constrained linear inverse problem and reconstruct the hidden objects by solving the inverse problem with an optimization framework. O’Toole *et al.* [29] propose a confocal NLOS imaging system where the image formation can be simplified to a convolution on the hidden objects. Thus the reconstruction can be expressed as a deconvolution process and solved efficiently. Heide *et al.* [10] model the occlusions in NLOS imaging and develop a factorization approach for inverse time-resolved light transport. Liu *et al.* [23] propose to recover the hidden objects in Fourier domain companied with Rayleigh Sommerfeld Diffraction (RSD) algorithm, which achieves high-fidelity NLOS reconstructions for challenging scenes. Xin *et al.* [41] present a novel theory of Fermat Path for NLOS transient measurements to help recover the surface hidden objects. While existing methods demonstrate superior performance in depth resolution, a major shortcoming of these methods is the limited lateral resolution of the reconstruction results.

Deep learning has demonstrated recent success in computational imaging [2, 7, 43, 42, 37]. Chopite *et al.* [8] propose to address the inverse problem with deep learning by first simulating plenty of training data according to confocal imaging system and then train a U-Net to recover the geometry of hidden objects. Chen *et al.* [6] propose a deep framework for both reconstruction and recognition tasks from NLOS measurements, which achieves the best performance so far. In real-world data, the improvements over other baselines are limited, which is mainly because of the differences between the simulated training data to the real-world datasets. To be specific, [6] adopts a simplified simulation process, which does not simulate multiple bounces, realistic BRDF, and background / sensor noise and jitter in real-world scenarios. In addition, even though [6] can recover 2D intensity as well, it does so by projecting down from a 3D feature space that has to be trained with supervisions of 3D cues. In contrast, our model is trained with only 2D supervisions and does not internally operate via a 3D feature space. This means that our model is much easier to train and requires less training samples than the existing ones [8, 6]. It also enables us to create training pairs of transient NLOS data and photorealistic hidden scene images. To the best of our knowledge, there is currently no way to capture or render 3D hidden scene data that would project to a photorealistic 2D image and pair that with the NLOS data of the same scene.

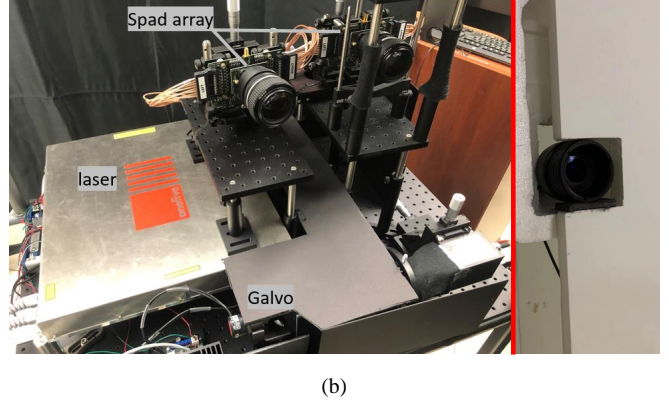
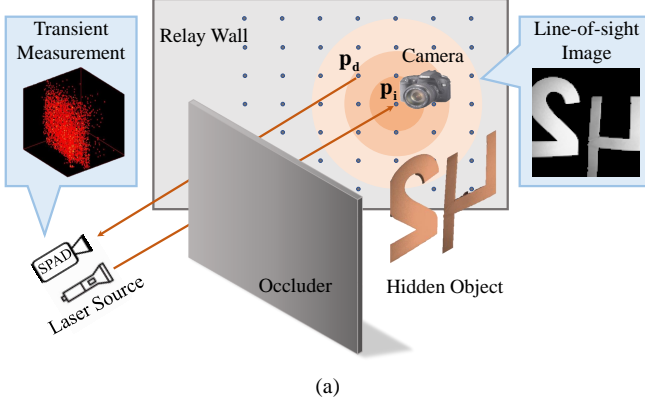


Figure 1. (a) The diagram of our imaging system. The light emitted from the laser source reaches the relay wall at \mathbf{p}_i and then is scattered towards the hidden object. The back-reflected light from the hidden object is captured by the SPAD detector focusing on \mathbf{p}_d . A conventional RGB camera is mounted in the relay wall to capture line-of-sight images. (b) Our non-confocal NLOS imaging hardware. It consists of an ultrafast pulsed laser and two 16×1 SPAD arrays as well as a conventional RGB camera mounted in the relay wall (right part of the red line). The galvo drives laser to scan the visible relay wall at 4 fps, during when the transient measurements and line-of-sight images are captured simultaneously.

3. Non-line-of-sight Photography

3.1. Forward Model

Without loss of generality, we build our NLOS imaging system in a non-confocal manner, which contains a laser source, a SPAD detector and a hidden object, shown in Fig. 1 (a). The laser projects short periodic light pulses $\delta(t)$ towards the relay wall at position \mathbf{p}_i . From where, the light is diffusely scattered at time $t = 0$ and targets on the hidden object. After integrated with the object, a fraction of the light is reflected back to the relay wall at position \mathbf{p}_d after time interval t and then captured by the SPAD detector resulting in the transient measurement $\Phi(\mathbf{p}_i, \mathbf{p}_d, t)$.

The transient measurement $\Phi(\mathbf{p}_i, \mathbf{p}_d, t)$ is a function of illumination point \mathbf{p}_i , detection point \mathbf{p}_d , and the time interval t , which can be modeled as

$$\Phi(\mathbf{p}_i, \mathbf{p}_d, t) = \iiint_{\Omega} I(\mathbf{x}) \frac{\delta(\|\mathbf{x} - \mathbf{p}_i\| + \|\mathbf{x} - \mathbf{p}_d\| - tc)}{\|\mathbf{x} - \mathbf{p}_i\|^2 \|\mathbf{x} - \mathbf{p}_d\|^2} d\mathbf{x}, \quad (1)$$

where $I(\mathbf{x})$ denotes the intensity of the hidden object at point \mathbf{x} and c is the speed of light. The numerator describes the distance that light travels from the relay wall to the objects. While the denominator accounts for the attenuation of light intensity along with light traveling.

On top of an active NLOS imaging system, we further integrate a conventional RGB camera (see Fig. 1 (a)). The RGB camera is mounted in the center of the relay wall with its principal axis perpendicular to the wall. It is carefully calibrated and synchronized with the NLOS hardware to capture the intensity I of the hidden object.

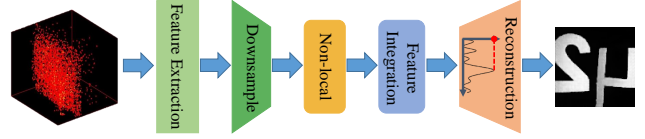


Figure 2. The flowchart of our network for NLOS photography. The network contains feature extraction block, downsampling operator, non-local block, feature integration block, and reconstruction block. It takes NLOS transient measurements as input and reconstructs 2D intensity images as output, which can be trained with only 2D supervisions. For more details about the network, please refer to the supplement.

3.2. Proposed Method

We focus on the NLOS photography which is to reconstruct intensity images resembling photographs captured by a LOS camera. In order to achieve it, we propose a deep learning-based network to recover the intensity from the NLOS transient measurements. Our network stems from a proven architecture previously developed for a different task [33] with several improvements for NLOS photography. Different from the existing networks [8, 6] that first recover a 3D volume then project it to a 2D coordinate to obtain depth maps or intensity images, our network can be trained directly with 2D supervisions, which is more suitable for real-world scenarios where 3D labels are difficult to obtain. The flowchart of our network is shown in Fig. 2.

Given a non-line-of-sight transient measurement, we first extract features with the *feature extraction* block, which consists of several 3D convolutions and dilated convolutions, then downscale the features along their spatial and temporal dimensions to improve the training efficiency with the *downsample* operation. After that, a *non-local* block is adopted to capture long-range spatial-temporal correlations within the measurement. Then *feature integration* block,

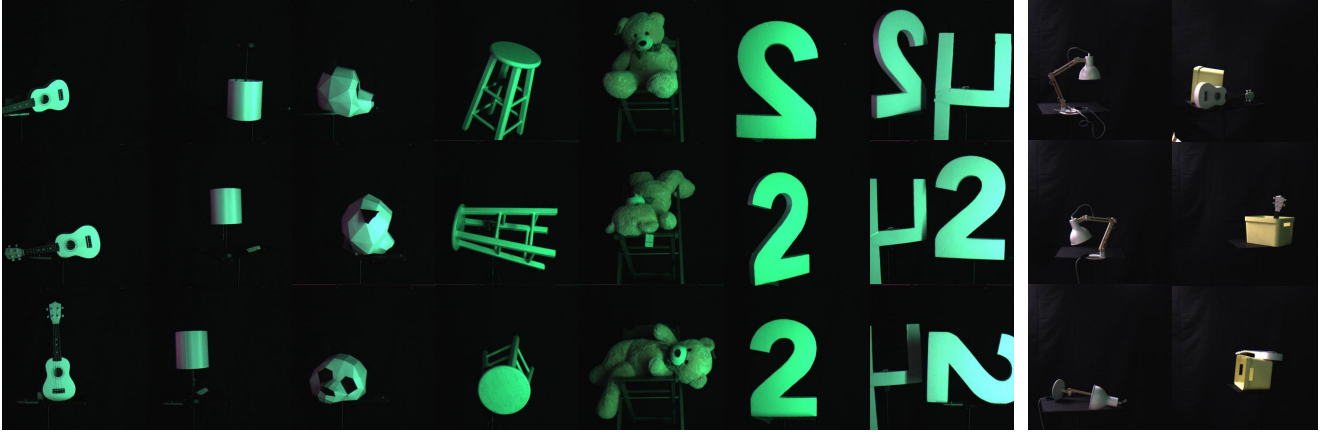


Figure 3. Thumbnails of several scenarios in our dataset with different orientations under various illumination conditions. Left part: without room light, right part: with room light. More scenarios can be found in the supplement.

containing several 3D dilated dense fusion sub-blocks [33], is performed to integrate the features across channels. The last *reconstruction* block first upsamples the input features in spatial dimension with 3D deconvolutions then generates 2D intensity images by reporting the maximum value along temporal dimension in each position. In order to train the network more efficiently, we adopt L2 norm computed between the ground-truth intensity images I , captured by the conventional RGB camera, and the reconstructed one \hat{I} from our network. Thus, the loss function can be written as

$$L_2 = \left\| I - \hat{I} \right\|_2. \quad (2)$$

L2 norm was widely used for image reconstruction [46]. We empirically verified that training using L2 loss leads to high quality reconstruction results for our model.

3.3. The Imaging System and The Dataset

To facilitate our network training and address the difficulties of current learning-based methods in generalization to real-world scenarios, we design a fast NLOS imaging system and build a large-scale real-world dataset.

The Imaging System. We start from a real-time NLOS imaging system [28] and embed an extra conventional RGB camera to capture both transient measurements and intensity images. Our system, shown in Fig. 1 (a), consists of an ultrafast pulsed laser (OneFive Katana HP) and two 16×1 SPAD arrays [34] as well as a conventional RGB camera. The 700 mW, 532 nm laser scans a 1.9×1.9 m area on the relay wall at 4 fps (exposure time per frame is 0.25 seconds). The effective temporal resolution of our system is 85 ps. Returning photons are captured by 28 SPAD pixels focused on a 1×9 cm area at the wall center. We apply the remapping operation in [28] to convert a raw measurement into a 3D data cube with resolution of $100 \times 100 \times 1440$.

The conventional RGB camera is positioned at the center of the relay wall with its imaging plane facing to the hidden

object. The field of view of the camera is carefully chosen to ensure sufficient coverage of the scene. During capturing, the camera is synchronized with the NLOS hardware and triggered at the middle of each scan. The exposure time of camera is kept short at 20 ms to approximate the desired illumination pattern. And the 2D intensity images are stored at a resolution of 700×700 . Thus, we can capture NLOS measurements and intensity images pairs in a fast manner (4 fps), which facilitates building a real-world dataset.

The Dataset. Our dataset consists of 400 transient measurements of distinct indoor scenes located 1m away from the diffuse relay wall under different illumination conditions (with and without room light, see examples in Fig. 3). One or more randomly oriented objects are placed at random locations within a $2 \times 2 \times 2$ m hidden volume. Each measurement is paired with a 2D intensity image of the scene. A backdrop is installed to insulate the hidden volume and forms the dark background in the intensity images. Our dataset captures a broad spectrum of real-world targets, ranging from rigid planar objects with uniform Lambertian reflectance to challenging deformable objects with complex shape and BRDF, for example, NUM2 in Fig. 3 has planar and uniform reflectance while BEAR has complicated shape. It is worth pointing that there are several practical issues for data capturing beyond resetting the objects, including the size, material, placement, and mounting of objects, which makes the collection of data non-trivial despite the fast capturing.

Our dataset offers several advantages, both as a source of training data and as a benchmark for algorithm evaluation. To the best of our knowledge, it is the largest dataset of *real* measurements, with two orders of magnitude more data than the existing ones [22, 24, 17, 9]. Apart from scale, our dataset represents a diverse set of scenes with *free object shape, pose and appearance*. This is in stark contrast to the current norm, which assumes large, bright, planar and cleverly positioned targets in favor of the algorithms. More-

over, our NLOS measurements come with high-resolution LOS images. We hope these additions enable the supervised training of learning-based models and ease the *quantitative* evaluation of reconstruction methods. Unlike the synthetic datasets [9, 6, 17], our data faithfully capture the entirety of light transport. We envision that a learning-based model could gain extra insight of the scene through the dissection of global light transport. Meanwhile, as current state-of-the-art typically do not account for multiple bounced light, our data would provide a means for assessing their robustness against the additional bounces present in real scenes.

Looking forward, a collection of NLOS measurements that probe complex scenes at full granularity is particularly important for the development of *high-fidelity* reconstruction methods. In this paper, we train our network on our dataset and demonstrate state-of-the-art reconstruction result at unprecedented quality and resolution approaching that of a line-of-sight camera (see Sec. 4). We believe that our dataset will pave the way for the supervised training of more sophisticated deep models and facilitate more future research.

4. Experiments

We conduct extensive experiments and compare with existing state-of-the-art methods to validate our NLOS photography that we would like to introduce as follows.

Dataset and metrics. We train our network on a subset of our dataset containing 10 different objects with 10 different views for each resulting in a total number of 100 distinct scenes. Note that, each scene has 50 data pairs (transient measurement – intensity image). In order to evaluate the efficiency of our network, we construct two test sets. The first set, denoted as *Unseen Poses* hereafter, contains *novel poses* of training objects. We randomly pick two additional views for five of the 10 objects in the training set. The more challenging *Unseen Objects* set assesses the model’s generalizability on *novel objects* completely missing in the training set. Note that the data in both test sets are absent from the training set. For training and evaluation, we downsample the intensity images to 100×100 using trilinear interpolation and normalize the input to $[0, 1]$. We adopt the green channel as a proxy to generate the ground-truth intensity map for network training. The evaluation metric is the generally PSNR and SSIM between recovered intensity and the ground truth.

Training details. We implement our model with PyTorch [31]. The weights in our model are randomly initialized and updated using Adam [15] with $\beta_1 = 0.9$ and $\beta_2 = 0.999$. The learning rate starts at 1×10^{-4} and decays by a factor of 0.95 after each epoch. We train the network on 4 NVIDIA 1080Ti GPUs with a batch size of 4. Training our model to full convergence takes about 20 hours.

Baselines. Since the existing NLOS reconstruction

methods are mainly focusing on recovering geometry instead of intensity from the transient measurements, it is non-trivial to make comparisons with them. We thus adopt two representative methods with some modifications to enable them to recover intensity images. The first baseline is Phasor Field RSD (PFRSD) [24], the state-of-the-art for physics-based non-confocal NLOS reconstruction. We adopt the fast implementation of Liu *et al.* [23]. The second baseline is a U-Net architecture originally proposed for image segmentation [35] and recently adopted for NLOS depth estimation [8]. This baseline is representative of learning-based reconstruction methods and is most similar to our training setup. We re-implement this model according to the network architecture in [8] and train the network with the same training data as our model.

4.1. Experimental Results

Results on Unseen Poses Set. We first evaluate our method on *Unseen Poses* set and make comparisons with PFRSD [23] and U-Net [8]. The quantitative results are listed in Table 1. As can be seen, our method achieves the best performance in terms of all scenarios and significantly surpasses the existing baselines in both PSNR and SSIM. Compared with U-Net [8], our method improve the reconstruction performance by a large margin (over 60% and 30% in PSNR and SSIM, respectively), which indicates the effectiveness of the method in intensity reconstruction for NLOS imaging.

In addition to quantitative comparisons, we also provide qualitative results shown in Fig. 4. Our results have the best reconstruction quality over all test scenes with visible improvements over the baselines. Specifically, PFRSD [23] fails to recover decent intensity, resulting in non-informative outputs. The learning-based U-Net [8] behaves slightly better than the PFRSD [23]. It recovers rough structures of the hidden objects but still with heavy blur. In contrast, our method can recover both main structures and fine details for the hidden objects and achieve decent results even in challenging cases. The distinct improvements over previous approaches on various objects across different poses further demonstrate the robustness of our method to viewpoints and orientations.

Results on Unseen Objects Set. In addition to the *Unseen Poses* set, we also evaluate our method on the *Unseen Objects* set, where the hidden objects have not been seen during training. The *Unseen Objects* set is adopted to evaluate the generalization capability of the method. The quantitative results are listed in Table 2. As can be seen, our method achieves the best performance in terms of all scenarios and significantly surpasses the existing baselines, which demonstrates the superiority of our method in generalization capability to unseen objects.

The qualitative comparisons are also provided in Fig. 5.

Table 1. The quantitative comparisons of intensity reconstructions from different methods on *Unseen Poses* set in terms of PSNR (upper part) and SSIM (lower part). “A” and “B” denote different poses of the objects. Note that, the quantitative results of PFRSD [23] can not be calculated, since it recovers depth instead of intensity, which does not match the groundtruth in relevant location and size.

Methods	BEAR		BOX		GUITAR		NUM2		NUM24		Ave.
	A	B	A	B	A	B	A	B	A	B	
U-Net [8]	23.28	22.82	22.26	22.69	20.25	18.29	17.74	16.15	14.35	14.59	19.24
Ours	29.63	30.74	34.02	34.62	31.21	30.48	31.04	28.61	30.35	31.90	31.26
U-Net [8]	.7916	.7653	.8300	.8536	.8630	.8126	.6543	.6055	.5231	.5878	.7287
Ours	.9586	.9591	.9674	.9827	.9677	.9571	.9201	.9172	.9488	.9521	.9531

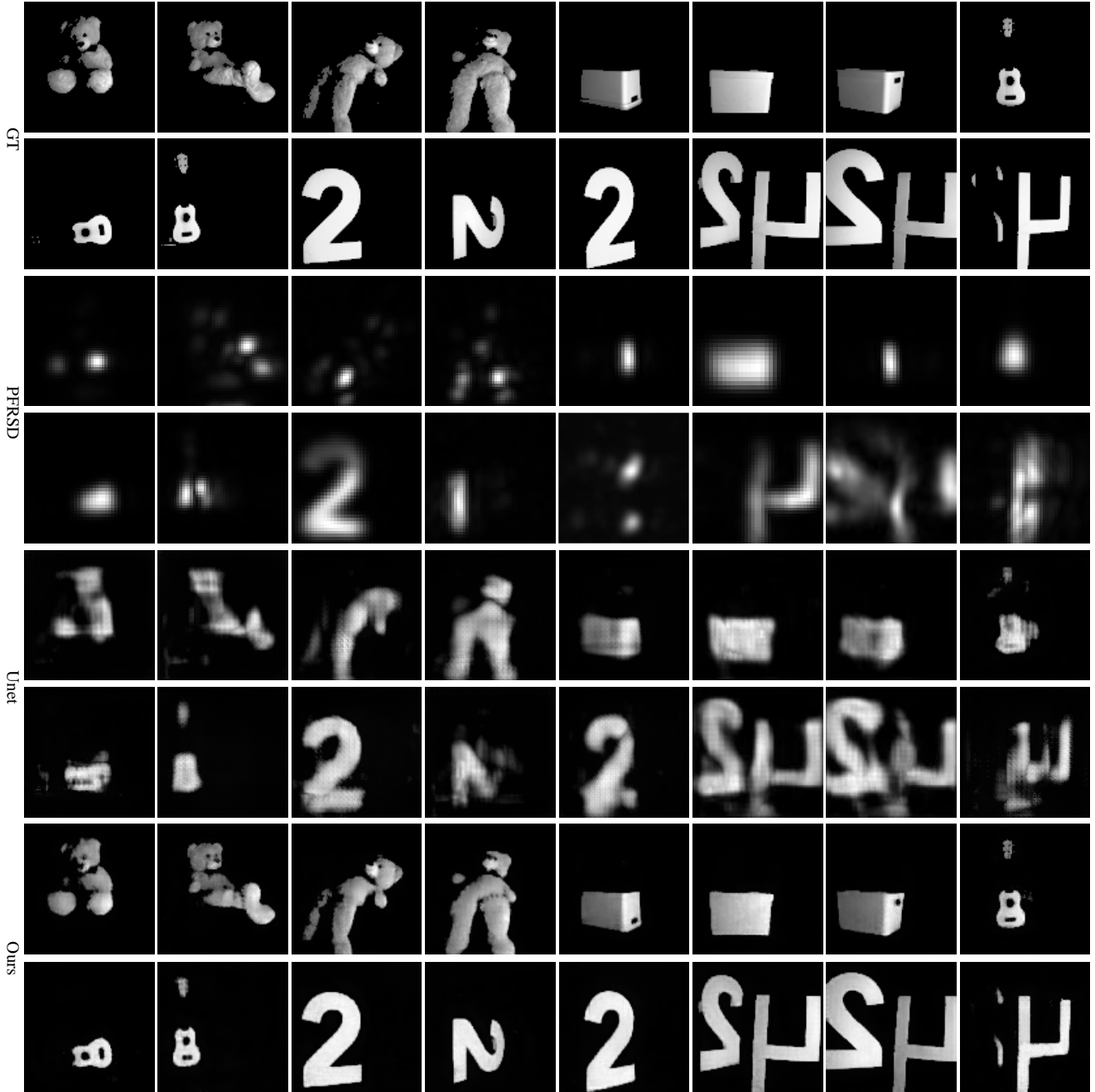


Figure 4. The qualitative comparisons of the reconstructed intensity from different methods on *Unseen Poses* set. GT denotes the captured intensity from our imaging system. See more results in supplementary material.

Table 2. The quantitative comparisons of intensity reconstructions from different methods on *Unseen Objects* set in terms of PSNR (upper part) and SSIM (lower part).

Methods \ Scenes	CHAIR		TRUCK		Ave.
	A	B	A	B	
U-Net [8]	18.49	20.27	19.89	19.16	19.45
Ours	21.11	22.60	27.93	24.37	24.00
U-Net [8]	.5719	.6327	.6176	.3345	.5392
Ours	.8508	.8453	.8901	.8662	.8631

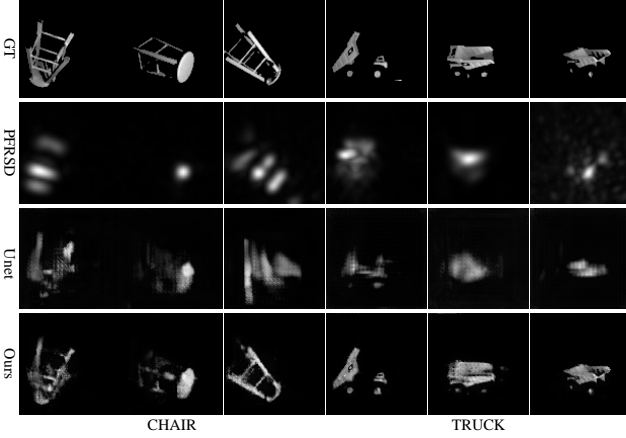


Figure 5. The qualitative comparisons of the reconstructed intensity from different methods on *Unseen Objects* set. GT denotes the captured intensity from our imaging system. See more results in supplementary material.

As we can see, our method achieves the best visualization quality over the baselines on the test scenes. Specifically, PFRSD [23] fail to reconstruct the hidden objects on such complicated scenarios. A U-Net [8] trained on our dataset can not recover the main structures of the scenes with informative intensity information missing and suffering from heavy blur. However, our method generates decent intensity reconstructions by reserving the main structures even in challenging scenarios. For example, our method can recover the cursory shape even for CHAIR, which contains very complicated structures, while both PFRSD [23] and U-Net [8] fail. The decent fidelity of our method on the *Unseen Objects* set demonstrates the superior generalization capability of our method over the existing baseline approaches, which is significant in real-world scenarios. Compared with the results in *Unseen Poses*, our network experiences much blurring and distortions in *Unseen Objects*. Note that it is much more challenging to recover unseen objects than seen objects with different poses, since the model can no longer rely on a strong shape prior of seen objects.

Inference Time. We further compare the inference time of our method with existing approaches. PFRSD [23] is tested on Intel Core i7-6700k @4GHz CPU, while U-Net [8] and our method are tested on a NVIDIA 1080Ti GPU. As shown in Table 3, our method is much faster than others, achieving nearly 17 and 46 times acceleration com-

Table 3. The inference time of different methods averaged on the test data with a resolution of $100 \times 100 \times 1440$.

Methods	PFRSD [23]	U-Net [8]	Ours
Time (s)	8.26	3.10	0.18

pared with U-Net [8] and PFRSD [23]. Benefited from the fast running time, we believe that our method can achieve video frame rate NLOS photography with proper optimization of the network architecture in near future, and thereby enabling many real-time applications in NLOS.

5. Ablations

We further investigate the performance of our method by reconstructing on data captured with various exposure time and recovering higher lateral resolution outputs.

5.1. Effects of Exposure Time

The exposure time of data capturing influences the performance of our method. Thus, we simulate training data with different exposure time by summing up 10, 20, 30, 40 and 50 frames (maximum number of frames), with corresponding exposure time of 2.5, 5, 7.5, 10, and 12.5 seconds, respectively. Different networks are trained on these data with different exposure time and evaluated on the *Unseen Poses* set. The quantitative results averaged over the corresponding test sets are shown in Fig. 6 (a) and (b).

Fig. 6 (a) and (b) shows improved reconstruction performance with increased exposure time (amount of adjacent frames in summation). Different from the PFRSD [23] that would saturate at a certain exposure time, the performance of our method continues to improve even with a long exposure time, for example 12.5 seconds (50 frames). The qualitative comparisons on one exemplar scene are also provided in Fig. 7. More details are included in our results with the increase of exposure time. Specifically, the reconstruction from short exposure time (e.g. 0.25 seconds) only contains main structures of the Teddy bear and is degraded with heavy blur. As the exposure time increases, the reconstructions start to contain more and more details, closely resembling the ground-truth intensity images. Our results show that the longer the exposure time is, the higher reconstruction performance our method can achieve.

5.2. Investigation on Higher Lateral Resolution Reconstruction

Finally, we demonstrate that our method can recover much higher lateral resolution results compared with existing reconstruction methods [24, 23, 8]. To this end, we first train networks to predict images with different lateral resolutions including 100×100 , 200×200 , 400×400 , and the full resolution 700×700 , then evaluate them on *Unseen Poses* set under lateral resolution of 700×700 . The quanti-

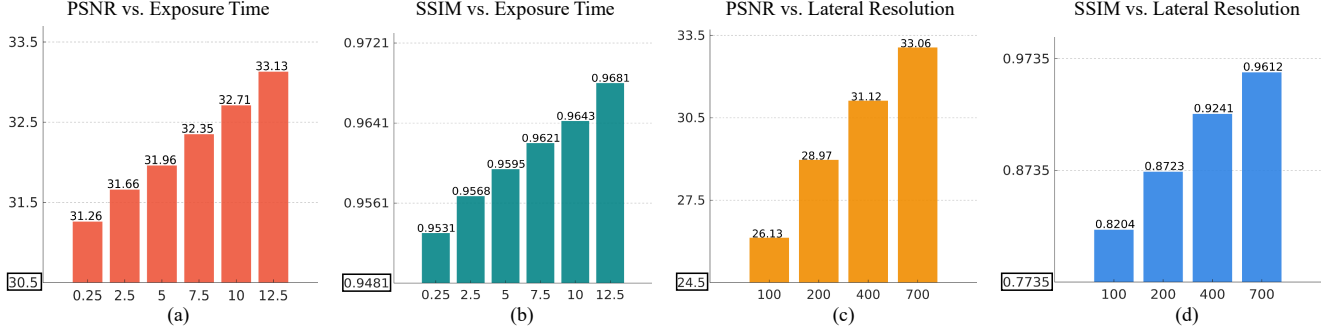


Figure 6. The quantitative comparisons of the reconstruction performance across different exposure time in terms of PSNR (a) and SSIM (b). The horizontal axis denotes the exposure time ranging from 0.25 to 12.5 seconds. The qualitative comparisons of the reconstruction performance on different lateral resolutions in terms of PSNR (c) and SSIM (d). Note that, here, we compute PSNR and SSIM on full resolution (700×700) by upsampling these low lateral resolution results to 700×700 . The numbers in horizontal axis denote the lateral resolutions of one side for a square image.

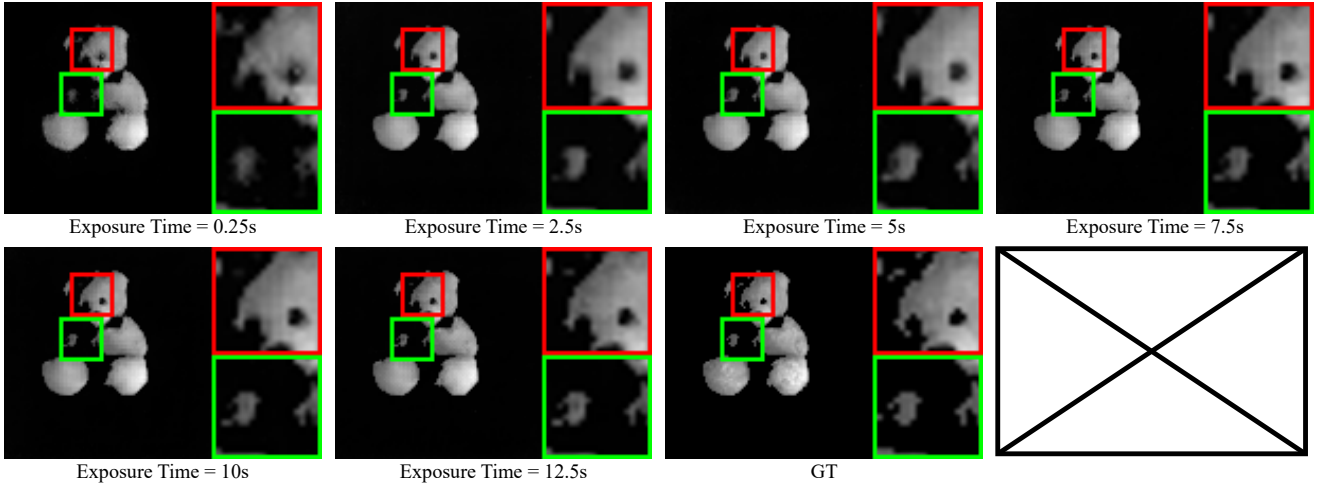


Figure 7. The qualitative comparisons of the reconstructed intensity images across different exposure time. The labels below each sub-figure denote the exposure time in seconds. GT denotes the ground-truth intensity image. Zoom in for a better visual experience.

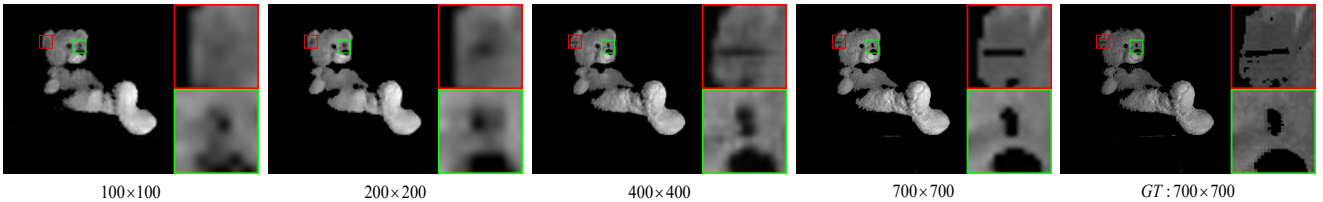


Figure 8. The qualitative comparisons of the reconstructed intensity images on different lateral resolutions. The labels below each sub-figure denote the lateral resolution of our outputs. GT denotes the ground-truth intensity image.

tative results averaged over the corresponding test sets are shown in Fig. 6 (c) and (d).

Fig. 6 (c) and (d) shows that our method can recover high lateral resolution reconstructions (up to 700×700) with superior quality. Quantitatively, our network achieves 26% improvements in terms of PSNR in 700×700 over that in 100×100 . The qualitative comparisons on one exemplar scene are provided in Fig. 8. As we can see, our method recovers more details of the object with increased lateral resolutions. Specifically, the reconstructions with resolution of 100×100 and 200×200 only contain the main structures of

the Teddy bear and suffer from heavy blur. The reconstruction with 400×400 is contaminated with noise. In contrast, the reconstruction on 700×700 contains fine details and is very close to the groundtruth. The significant improvements over various lateral resolutions demonstrate the high fidelity of our method for reconstructing intensity images.

6. Conclusion

In this paper, we proposed a new problem formulation of NLOS photography to recover intensity images with high lateral resolution and fine details from transient NLOS mea-

surements. In order to achieve this, we first collected a new NLOS dataset with the help of a fast NLOS imaging system. Our dataset contains 400 different real-world scenes, which is an order of magnitude larger than existing ones. Further, we designed a learning-based method and demonstrated the reconstruction of NLOS images with unprecedented quality. We hope that our work can draw the attention to the problem of NLOS photography, and our dataset would help to facilitate future research in this field. In particular, we anticipate our results of NLOS photography with high later resolution can be further combined with prior results of high depth resolution, leading to exciting avenues of high fidelity 3D NLOS reconstructions.

References

- [1] Victor Arellano, Diego Gutierrez, and Adrian Jarabo. Fast back-projection for non-line of sight reconstruction. *Optics express*, 25(10):11574–11583, 2017. 3
- [2] George Barbastathis, Aydogan Ozcan, and Guohai Situ. On the use of deep learning for computational imaging. *Optica*, 6(8):921–943, 2019. 3
- [3] Mufeed Batarseh, Sergey Sukhov, Zhiqin Shen, Heath Gemar, Reza Rezvani, and Aristide Dogariu. Passive sensing around the corner using spatial coherence. *Nature communications*, 9(1):1–6, 2018. 2
- [4] Katherine L Bouman, Vickie Ye, Adam B Yedidia, Frédo Durand, Gregory W Wornell, Antonio Torralba, and William T Freeman. Turning corners into cameras: Principles and methods. In *ICCV*, 2017. 2
- [5] 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. 2
- [6] Wenzheng Chen, Fangyin Wei, Kiriakos N Kutulakos, Szymon Rusinkiewicz, and Felix Heide. Learned feature embeddings for non-line-of-sight imaging and recognition. *ACM Transactions on Graphics (TOG)*, 39(6):1–18, 2020. 3, 4, 6
- [7] Zhen Cheng, Zhiwei Xiong, and Dong Liu. Light field super-resolution by jointly exploiting internal and external similarities. *IEEE Transactions on Circuits and Systems for Video Technology*, 2019. 3
- [8] Javier Grau Chopite, Matthias B Hullin, Michael Wand, and Julian Iseringhausen. Deep non-line-of-sight reconstruction. In *CVPR*, 2020. 3, 4, 6, 7, 8
- [9] 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. In *ACM SIGGRAPH 2019 Posters*, pages 1–2. 2019. 2, 3, 5, 6
- [10] 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 (ToG)*, 38(3):1–10, 2019. 3
- [11] Felix Heide, Lei Xiao, Wolfgang Heidrich, and Matthias B Hullin. Diffuse mirrors: 3d reconstruction from diffuse indirect illumination using inexpensive time-of-flight sensors. In *CVPR*, 2014. 3
- [12] Julian Iseringhausen and Matthias B Hullin. Non-line-of-sight reconstruction using efficient transient rendering. *ACM Transactions on Graphics (TOG)*, 39(1):1–14, 2020. 3
- [13] Mariko Isogawa, Dorian Chan, Ye Yuan, Kris Kitani, and Matthew O’Toole. Efficient non-line-of-sight imaging from transient sinograms. In *ECCV*, 2020. 2, 3
- [14] Adrian Jarabo, Julio Marco, Adolfo Munoz, Raul Buisan, Wojciech Jarosz, and Diego Gutierrez. A framework for transient rendering. *ACM Transactions on Graphics (ToG)*, 33(6):1–10, 2014. 3
- [15] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *ICLR*, 2015. 6
- [16] Ahmed Kirmani, Tyler Hutchison, James Davis, and Ramesh Raskar. Looking around the corner using transient imaging. In *ICCV*, 2009. 2, 3
- [17] Jonathan Klein, Martin Laurenzis, Dominik L. Michels, and Matthias B. Hullin. A quantitative platform for non-line-of-sight imaging problems. In *BMVC*, page 104, 2018. 2, 3, 5, 6
- [18] Jonathan Klein, Christoph Peters, Jaime Martín, Martin Laurenzis, and Matthias B Hullin. Tracking objects outside the line of sight using 2d intensity images. *Scientific reports*, 6(1):1–9, 2016. 3
- [19] Marco La Manna, Ji-Hyun Nam, Syed Azer Reza, and Andreas Velten. Non-line-of-sight-imaging using dynamic relay surfaces. *Optics express*, 28(4):5331–5339, 2020. 3
- [20] Martin Laurenzis and Andreas Velten. Non-line-of-sight active imaging of scattered photons. In *Electro-Optical Remote Sensing, Photonic Technologies, and Applications VII; and Military Applications in Hyperspectral Imaging and High Spatial Resolution Sensing*, volume 8897, page 889706, 2013. 3
- [21] David B Lindell, Gordon Wetzstein, and Vladlen Koltun. Acoustic non-line-of-sight imaging. In *CVPR*, 2019. 2, 3
- [22] David B Lindell, Gordon Wetzstein, and Matthew O’Toole. Wave-based non-line-of-sight imaging using fast fk migration. *ACM Transactions on Graphics (TOG)*, 38(4):1–13, 2019. 2, 3, 5
- [23] Xiaochun Liu, Sebastian Bauer, and Andreas Velten. Phasor field diffraction based reconstruction for fast non-line-of-sight imaging systems. *Nature communications*, 11(1):1–13, 2020. 3, 6, 7, 8
- [24] Xiaochun Liu, Ibón Guillén, 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(7771):620–623, 2019. 2, 3, 5, 6, 8
- [25] Tomohiro Maeda, Guy Satat, Tristan Swedish, Lagnojita Sinha, and Ramesh Raskar. Recent advances in imaging around corners. *arXiv preprint arXiv:1910.05613*, 2019. 2
- [26] Tomohiro Maeda, Yiqin Wang, Ramesh Raskar, and Achuta Kadambi. Thermal non-line-of-sight imaging. In *ICCP*, 2019. 2
- [27] Nikhil Naik, Shuang Zhao, Andreas Velten, Ramesh Raskar, and Kavita Bala. Single view reflectance capture using multiplexed scattering and time-of-flight imaging. In *Proceedings*

- of the 2011 SIGGRAPH Asia Conference, pages 1–10, 2011. 3
- [28] Ji Hyun Nam, Eric Brandt, Sebastian Bauer, Xiaochun Liu, Eftychios Sifakis, and Andreas Velten. Real-time non-line-of-sight imaging of dynamic scenes. *arXiv preprint arXiv:2010.12737*, 2020. 5
- [29] 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. 2, 3
- [30] Rohit Pandharkar, Andreas Velten, Andrew Bardagjy, Everett Lawson, Mounsi Bawendi, and Ramesh Raskar. Estimating motion and size of moving non-line-of-sight objects in cluttered environments. In *CVPR*, 2011. 2
- [31] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. Pytorch: An imperative style, high-performance deep learning library. In *NIPS*. 2019. 6
- [32] Adithya Kumar Pediredla, Mauro Buttafava, Alberto Tosi, Oliver Cossairt, and Ashok Veeraraghavan. Reconstructing rooms using photon echoes: A plane based model and reconstruction algorithm for looking around the corner. In *ICCP*, pages 1–12, 2017. 3
- [33] Jiayong Peng, Zhiwei Xiong, Xin Huang, Zheng-Ping Li, Dong Liu, and Feihu Xu. Photon-efficient 3d imaging with a non-local neural network. In *ECCV*, 2020. 2, 4, 5
- [34] Marco Renna, Ji Hyun Nam, Mauro Buttafava, Federica Villa, Andreas Velten, and Alberto Tosi. Fast-gated 16 × 1 SPAD array for non-line-of-sight imaging applications. *Instruments*, 4(2):14, 2020. 5
- [35] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *International Conference on Medical image computing and computer-assisted intervention*, pages 234–241. Springer, 2015. 6
- [36] Charles Saunders, John Murray-Bruce, and Vivek K Goyal. Computational periscopy with an ordinary digital camera. *Nature*, 565(7740):472–475, 2019. 2
- [37] Zhan Shi, Chang Chen, Zhiwei Xiong, Dong Liu, and Feng Wu. Hscnn+: Advanced cnn-based hyperspectral recovery from rgb images. In *CVPR-W*, 2018. 3
- [38] Chia-Yin Tsai, Aswin C Sankaranarayanan, and Ioannis Gkioulekas. Beyond volumetric albedo—a surface optimization framework for non-line-of-sight imaging. In *CVPR*, 2019. 3
- [39] Andreas Velten, Thomas Willwacher, Otkrist Gupta, Ashok Veeraraghavan, Mounsi G Bawendi, and Ramesh Raskar. Recovering three-dimensional shape around a corner using ultrafast time-of-flight imaging. *Nature communications*, 3(1):1–8, 2012. 2, 3
- [40] Cheng Wu, Jianjiang Liu, Xin Huang, Zheng-Ping Li, Chao Yu, Jun-Tian Ye, Jun Zhang, Qiang Zhang, Xiankang Dou, Vivek K Goyal, et al. Non-line-of-sight imaging over 1.43 km. *Proceedings of the National Academy of Sciences*, 118(10), 2021. 2, 3
- [41] Shumian Xin, Sotiris Nousias, Kiriakos N Kutulakos, Aswin C Sankaranarayanan, Srinivasa G Narasimhan, and Ioannis Gkioulekas. A theory of fermat paths for non-line-of-sight shape reconstruction. In *CVPR*, 2019. 3
- [42] Zhiwei Xiong, Zhan Shi, Huiqun Li, Lizhi Wang, Dong Liu, and Feng Wu. Hscnn: Cnn-based hyperspectral image recovery from spectrally undersampled projections. In *ICCV-W*, 2017. 3
- [43] Mingde Yao, Zhiwei Xiong, Lizhi Wang, Dong Liu, and Xuejin Chen. Spectral-depth imaging with deep learning based reconstruction. *Optics Express*, 27(26):38312–38325, 2019. 3
- [44] Adam B Yedidia, Manel Baradad, Christos Thrampoulidis, William T Freeman, and Gregory W Wornell. Using unknown occluders to recover hidden scenes. In *CVPR*, 2019. 2
- [45] Sean I. Young, David B. Lindell, Bernd Girod, David Taubman, and Gordon Wetzstein. Non-line-of-sight surface reconstruction using the directional light-cone transform. In *Proc. CVPR*, 2020. 3
- [46] Hang Zhao, Orazio Gallo, Iuri Frosio, and Jan Kautz. Loss functions for image restoration with neural networks. *IEEE Transactions on computational imaging*, 3(1):47–57, 2016. 5