

Learning to Predict Scene-Level Implicit 3D from Posed RGBD Data

Nilesh Kulkarni, Linyi Jin, Justin Johnson, David F. Fouhey
 University of Michigan

Abstract

We introduce a method that can learn to predict scene-level implicit functions for 3D reconstruction from posed RGBD data. At test time, our system maps a previously unseen RGB image to a 3D reconstruction of a scene via implicit functions. While implicit functions for 3D reconstruction have often been tied to meshes, we show that we can train one using only a set of posed RGBD images. This setting may help 3D reconstruction unlock the sea of accelerometer+RGBD data that is coming with new phones. Our system, D2-DRDF, can match and sometimes outperform current methods that use mesh supervision and shows better robustness to sparse data.

1. Introduction

Consider the image in Figure 1. From this ordinary RGB image, we can understand the complete 3D geometry of the scene, including the floor and walls behind the chairs. Our goal is to enable computers to recover this geometry from a single RGB image. To this end, we present a method that does so while learning *only* on posed RGBD images.

The task of reconstructing the full 3D of a scene from a single previously unseen RGB image has long been known to be challenging. Early work on full 3D relied on voxels [6, 17] or meshes [19], but these representations fail on scenes due to topology and memory challenges. Implicit functions (or learning to map each point in \mathbb{R}^3 to a value like the distance to the nearest surface) have shown substantial promise at overcoming these challenges. When conditioned on an image, these have led to several successful methods.

Unfortunately, the implicit function status quo mainly ties implicit function reconstruction methods to mesh supervision. This symbiosis has emerged since meshes give excellent direct supervision. However, methods are limited to training with an image-aligned mesh that is usually watertight (and often artist-created) [35, 38, 44] and occasionally non-watertight but professionally-scanned [27, 61].

We present a method, Depth to DRDF (D2-DRDF), that breaks the implicit-function/mesh connection and can train an effective single RGB image implicit 3D reconstruction

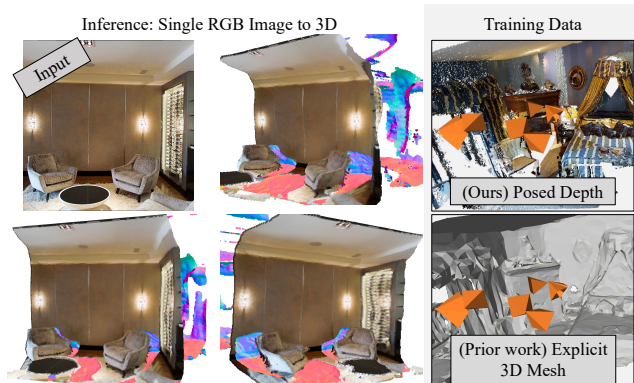


Figure 1. We propose a method, D2-DRDF, that reconstructs full 3D from a single RGB image. During inference (left), our method uses implicit functions to reconstruct the complete 3D including visible and occluded surfaces (visualized as surface normals) such as the **occluded wall** and **empty floor**. For training, our method uses *only* RGBD images and poses, unlike most prior works that need an explicit and often watertight 3D mesh.

system using a set of RGBD images with known pose. We envision that being able to entirely skip meshing will enable the use of vast quantities of lower-quality data from consumers (e.g., from increasingly common consumer phones with LiDAR scanners and accelerometers) as well as robots. In addition to permitting training, the bypassing of meshing may enable adaptation in a new environment on raw data without needing an expert to ensure mesh acquisition.

Our key insight is that we can use segments of observed free space in depth maps in other views to constrain distance functions. We show this using the Directed Ray Distance Function (DRDF) [27] that has recently shown good performance in 3D reconstruction using implicit functions and has the benefit of not needing watertight meshes for training. Given an input *reference* view, the DRDF breaks the problem of predicting the 3D surface into a set of independent 1D distance functions, each along a ray through a pixel in the reference view and accounting for only surfaces on the ray. Rather than use an ordinary unsigned distance function, the DRDF signs the distance using the location of the camera’s pinhole and the intersections along the ray. While [27] showed their method could be trained on *non*-watertight meshes, their method is still dependent

on meshes. In our paper, we show that the DRDF can be cleanly supervised using *auxiliary views* of RGBD observations and their poses. We derive constraints on the DRDF that power loss functions for training our system. While the losses on any one image are insufficient to produce a reconstruction, they provide a powerful signal when accumulated across thousands of training views.

We evaluate our method on realistic indoor scene datasets and compare it with methods that train with full mesh supervision. Our method is competitive and sometimes even better compared to the best-performing mesh-supervised method [28] with full, professional captures. As we degrade scan completeness, our method largely maintains its performance while mesh-based methods perform substantially worse. We conclude by showing that fine-tuning of our method on a handful of images enables a simple, effective method for fusing and completing scenes from a handful of RGBD images.

2. Related Work

Our approach infers a complete 3D scene from a single RGB image using implicit functions that are supervised via posed RGBD scans. This touches on several long-term goals of 3D computer vision that we discuss below.

Reconstructing Scenes from a Single Image. At test time, our system produces a *full* 3D reconstruction from a single RGB image, including occluded regions. This means that our desired output goes beyond 2.5D properties such as depth [4, 12], surface normals [12, 14, 55], or other intrinsic-image properties [23, 47]. Works on attempting to reconstruct complete 3D usually have focused on reconstructing objects from image by using voxels [6, 17] or meshes [18, 19], point clouds [13, 31] or CAD models [22]. These are usually trained on synthetic datasets like ShapeNet [2] or scene datasets [16, 43, 50] and do not generalize well to realistic scenes. Another line of works tries to learn holistic structures [37, 59] or planar surfaces [24, 32, 58] from realistic scanned mesh data [1, 8]. Creating realistic mesh based data for scenes requires post-processing using Poisson surface reconstruction [9, 25] which leads to deviation from raw captures. Manually aligning them is expensive hence 3D object-aligned datasets like [49, 52] are scarce and limited in diversity. Our method avoids these limitations by directly operating on the raw captured RGBD data.

Reconstructing Scenes from Posed Scans. There has been considerable work on using multiview RGBD data at inference time to produce 3D reconstructions, starting with analytic techniques [7, 9, 21, 29] and now using learning [20, 51, 54]. We use some similar tools to this community: for instance, the ray distances we use are known in this community as a projected distance functions [7]. However, their works solve a fundamentally different problem

by using posed RGBD data at test time: their goal is to produce a reconstruction from a set of posed RGBD images from one particular scene; our goal is to use a large dataset of posed RGBD data to train a neural network that can map a new single RGB image to a 3D reconstruction.

Our goal of learning to predict reconstructions from a single previously unseen image also distinguishes our work from NeRF [36] and other similar radiance field approaches [3]. Their goal is to learn a radiance field for a *particular* scene from a set of posed scans. There are methods that try to predict this radiance field from a single image [57]; we compare with a model using similar losses and find that an objective specialized for 3D works better.

Implicit Functions for 3D Reconstruction. We perform reconstruction with implicit functions. Implicit functions have been used for shape and scene modeling as level sets [34], signed distance functions [38], occupancy function [35, 40, 44], distance functions [5, 27, 48] and other modifications like [56]. Our work has two key distinctions. First, many approaches [5, 48, 56] *fit* an implicit function to one shape (i.e., there is no generalization to new shapes). In contrast, our work produces reconstructions from previously unseen RGB images. Second, the methods that predict an implicit function from a new image [27, 61] assume access to a mesh at training time. On the other hand, our work assumes access only to posed RGBD data for training.

The most similar work to ours is [27] that also aims to learn to reconstruct scenes from a single image by training on realistic data. While [27] shows how to reduce supervision requirements by enabling the use of non-watertight meshes, their method is still limited to mesh supervision. We show how to use posed RGBD data for supervision instead of using an image-aligned mesh. This substantially reduces the requirements for collecting training data.

3. Pixel-Aligned 3D Reconstruction & DRDF

We propose a method to predict full 3D scene structure from a previously unseen RGB image. At inference time, our method receives a single image with no depth information. We refer to this view as the *reference view*. As output, the method produces an estimate of a *distance function* at each point in a pre-defined set of 3D points. Given this volume of predicted distances, one can *decode* the predicted distance into a set of surfaces: e.g., if the predicted distance were the unsigned distance to the nearest surface, one could declare all points with sufficiently small predicted distance to be surfaces. At training time, our network is given supervision for predicting its distance function. Previously, supervision has been done via a mesh that provides oracle distance calculations. In §4, we show how to derive supervision from posed RGBD images at *auxiliary views* instead.

The approaches in the paper follow pixelNeRF [57] for predicting a distance at a 3D point x . The network accepts

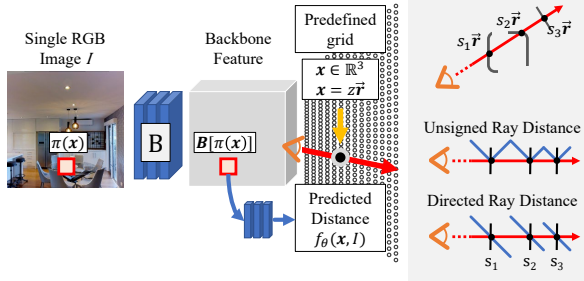


Figure 2. Model and DRDF overview. (Left) Given a single input RGB image I and query point \mathbf{x} on a ray $\vec{\mathbf{r}}$ in a pre-defined grid, all models, like [57] extract image features at the projection $\pi(\mathbf{x})$ of \mathbf{x} with a backbone and then predict the distance using a MLP on the backbone features and a positional encoding of \mathbf{x} . (Right) Given a ray $\vec{\mathbf{r}}$ through the scene, ray distance functions are defined as the distance to intersections of the ray with the scene. This permits analyzing the 3D scene as a set of 1D distances per ray. The Directed Ray Distance Function is further signed to be positive before and negative after an intersection.

an image I and 3D point \mathbf{x} and makes a prediction $f_\theta(\mathbf{x}, I)$. As shown in Fig. 2, the network extracts a convolutional feature from a backbone \mathbf{B} at the projection of \mathbf{x} , denoted $\pi(\mathbf{x})$. The image feature is concatenated to a positional encoding of \mathbf{x} and put into a MLP to produce the final prediction $f_\theta(\mathbf{x}, I)$. By repeating this inference for a set of 3D points, the network can infer the full 3D of the scene.

Training. In a conventional setup, the network $f_\theta(\mathbf{x}, I)$ is trained by obtaining samples from a mesh and directly training the network via regression. Given a mesh that is co-aligned with an image, one can obtain an n tuples consisting of an image I_i , 3D point \mathbf{x}_i and ground-truth distance d_i . Then the network is trained by minimizing the empirical risk $\frac{1}{n} \sum_{i=1}^n \mathcal{L}(f_\theta(\mathbf{x}_i, I_i), d_i)$ for some loss \mathcal{L} like L1.

What Distance Function Should We Predict? There are multiple distance functions in the literature, and the above formulation enables predicting any of them. To explain distinctions, assume we are given a point \mathbf{x} at which we want to evaluate a distance and a scene $\mathcal{S} \subseteq \mathbb{R}^3$ that we want to evaluate the distance to.

In a standard distance function, one defines a distance from a point \mathbf{x} to the *entire* scene \mathcal{S} . For instance, the *Unsigned Distance Function* (UDF) is defined as $\min_{\mathbf{x}' \in \mathcal{S}} \|\mathbf{x} - \mathbf{x}'\|$. In a *ray-based distance function* [7, 27], the distance is defined from \mathbf{x} only to the points in the scene that lie on the ray $\vec{\mathbf{r}}$ from the reference view pinhole towards \mathbf{x} . In other words, the ray distance function is defined as the distance to the intersections of the ray with the scene. For instance, the *Unsigned Ray Distance Function* (URDF) is defined as $\min_{\mathbf{x}' \in \mathcal{S}, \mathbf{x}' \text{ on } \vec{\mathbf{r}}} \|\mathbf{x} - \mathbf{x}'\|$. This definition means that when predicting the distance at a point \mathbf{x} , all of the points that define the ray distance function project to the same pixel, $\pi(\mathbf{x})$, which in turn means that neural networks do not need as large receptive fields [27].

Since a ray-based distance function is defined *only* by scene points that intersect a ray, we can reduce the 3D problem to a 1D problem, as illustrated in Fig. 2. Rather than represent each 3D point as a 3D point, we represent it as the distance along the ray $\vec{\mathbf{r}}$: the vector \mathbf{x} is represented via a scalar z such that $\mathbf{x} = z\vec{\mathbf{r}}$. Then, if the i^{th} intersection along the ray is $s_i\vec{\mathbf{r}}$ for $s_i \in \mathbb{R}$, the URDF is defined as the distance to the nearest intersection along that ray, or $d_{UR}(z) = \min_{i=1, \dots, k} |z - s_i|$.

Directed Ray Distance Function. The Directed Ray Distance Function (DRDF) is a ray distance function that incorporates a sign into the URDF. In particular, the DRDF is defined as $d_{DR}(z) = \text{dir}(z)d_{UR}(z)$, where the predicate $\text{dir}(z)$ is positive if the nearest intersection is ahead of z along the ray from the pinhole and negative otherwise. This is pictured in Fig. 2. The presence of both positive and negative components leads to better behavior under uncertainty [27]. *Critical Properties of DRDFs.* Two critical properties that we use to define supervision are that at z , the distance to the nearest intersection is $|d_{DR}(z)|$ and the nearest intersection on the ray is at $z + d_{DR}(z)$.

4. Depth-Supervised DRDF Reconstruction

Having introduced the setup in §3, we now explain how to train a network to predict an image-conditioned DRDF *without* access to an underlying mesh but instead access to posed RGBD images. Concretely, given a point $z\vec{\mathbf{r}}$, we aim to define a loss function \mathcal{L} that evaluates a network’s prediction of the DRDF $f_\theta((z\vec{\mathbf{r}}), I)$.

Our key insight is that we can *sometimes* see this point $z\vec{\mathbf{r}}$ in other views and these observations of $z\vec{\mathbf{r}}$ provide a constraints on what the value of the DRDF can be. Consider, for instance, the point 3m out from reference view camera pinhole along the red ray in Fig. 3. If we project this point into the auxiliary view, we know that its distance to the camera is closer to auxiliary view’s pinhole than the depthmap observed at the auxiliary view. Moreover, we can see that a segment of the ray is visible, starting at the point $s\vec{\mathbf{r}}$ and ending at the point $e\vec{\mathbf{r}}$ with the ending point $e\vec{\mathbf{r}}$ closer to $z\vec{\mathbf{r}}$ than $s\vec{\mathbf{r}}$ is.

We can derive supervision on $d_{DR}(z)$ from the segment of visible free space observable in the auxiliary view by using the fact that the nearest intersection at z is $|d_{DR}(z)|$ units away. Since e is nearer than s , and there are *no* intersections in the free space between s and e , we know that there cannot be an intersection within $(e - z)$ of z . This gives an inequality that $|d_{DR}(z)| \geq (e - z)$: otherwise there would be an intersection in the visible free space. If the ends of the segment are actual intersections where the ray visibly intersects an object, we obtain an equality constraint on $d_{DR}(z)$ since we actually see the nearest intersection. In other cases, the ray simply disoccludes or is occluded, and we can only have an inequality. We denote intersection start/end events as \mathcal{I}

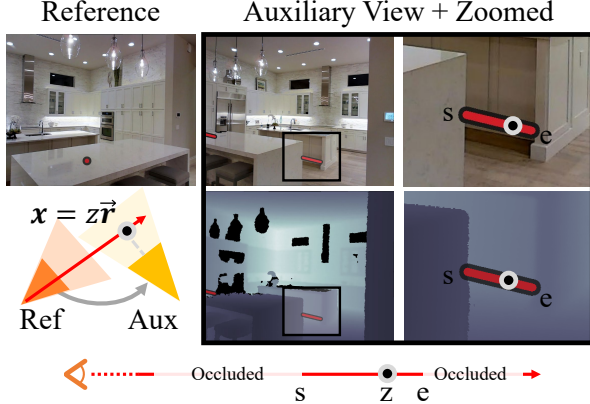


Figure 3. A red ray originates in the reference camera. Given a point x that is z units along the ray \vec{r} , we can test if $z\vec{r}$ is visible in an auxiliary view by comparing its projected depth with the RGBD image. The point z is visible, along with a segment of the ray from s to e units. Since the distance to the nearest intersection is $|d_{\text{DR}}(z)|$ and we know there are no intersections within $(e - z)$ units of z , we can constrain that $|d_{\text{DR}}(z)| \geq (e - z)$.

and occlusion start/end events as \mathcal{O} and distinguish them by checking if the projected z -value of the ray in auxiliary view is sufficiently close to the recorded auxiliary view’s depthmap value.

4.1. Losses

We now convert the concept of using freespace to constrain the DRDF into concrete losses. Recall that our goal is to evaluate a prediction $y = f_{\theta}(z\vec{r}, I)$ from our network and that the point at z along \vec{r} is in some segment of visible freespace from s to e along \vec{r} . Our goal is to penalize the prediction y . Since there are two start/end event classes (\mathcal{I} , \mathcal{O}), there are four segment types: \mathcal{II} , \mathcal{IO} , \mathcal{OI} , and \mathcal{OO} . We show several of these segments in Figure 5.

To assist the loss definitions, we define variables $l_s = s - z$, and $l_e = e - z$ which we plot in Fig. 4 as a function of z . l_s is the value of d_{DR} if the closest intersection is at s , and l_e is the value of d_{DR} if the closest intersection is at e . The values l_s, l_e define equalities for intersections and inequalities for occlusions.

$\mathcal{L}_{\mathcal{II}}$: \mathcal{II} segment. Given a segment bounded by two intersections, the nearest intersections are known exactly as l_s or l_e depending on whether $z < \frac{s+e}{2}$ or not. We penalize the ℓ_1 error between the prediction y and the known DRDF, or $\mathcal{L}_{\mathcal{II}}(y) = |y - l_s|$ if $z < \frac{s+e}{2}$ and $|y - l_e|$ otherwise. This penalty is zero only when y is equal to the known DRDF.

$\mathcal{L}_{\mathcal{OO}}$: \mathcal{OO} segment. Given a segment bounded by two occlusion events, the exact DRDF is not known, but the visible free space rules out potential values. Since $|d_{\text{DR}}(z)|$ is the distance to the nearest intersection, $d_{\text{DR}}(z)$ cannot lie in $[l_s, l_e]$ since such a value would imply an intersection in free space between s and e . We penalize incursions into $[l_s, l_e]$ with a ℓ_1 penalty: if we denote halfway be-

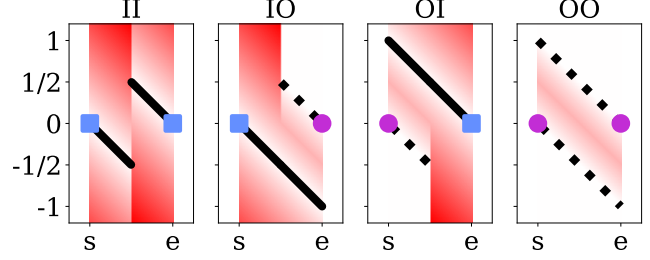


Figure 4. Loss functions for a starting s and ending e events spaced a unit apart. Intersection events (blue \blacksquare) define the DRDF precisely, and we penalize the network from deviating. Occlusion events (purple \bullet) provide bounds that we penalize the network for violating. We plot l_s as a line from the start event and l_e as a line from the end event. Loss: 0 1.7

tween l_s and l_e as h , then this can be done as $\mathcal{L}_{\mathcal{OO}}(y) = \max(0, l_e - h - |y - h|)$. The resulting penalty is zero if $y \leq l_s$ or if $y \geq l_e$.

$\mathcal{L}_{\mathcal{IO}}$: \mathcal{IO} segment. When the segment is bounded by an intersection event followed by an occlusion event, the situation is more complex and we define $\mathcal{L}_{\mathcal{IO}}$ piecewise. In the first half of the segment ($z < \frac{s+e}{2}$), d_{DR} is exactly known, and so we can use an ℓ_1 penalty like the \mathcal{II} case, so $\mathcal{L}_{\mathcal{IO}}(y) = |y - l_s|$. In the second half, there are two options. If the nearest intersection is s , then $d_{\text{DR}}(z) = l_s$. Otherwise, the nearest intersection is unknown but after e and so $d_{\text{DR}}(z) > l_e$ must hold (since $d_{\text{DR}}(z) \leq l_e$ would imply the existence of an intersection before e). We take the minimum of errors for the two cases: ℓ_1 distance to l_s and a ℓ_1 penalty function $\max(0, l_e - y)$, resulting in $\mathcal{L}_{\mathcal{IO}}(y) = \min(\max(0, l_e - y), |y - l_s|)$. This part of the penalty is zero if either $y = l_s$ or if $y > l_e$.

$\mathcal{L}_{\mathcal{OI}}$: \mathcal{OI} segment. The $\mathcal{L}_{\mathcal{OI}}$ loss is defined symmetrically to $\mathcal{L}_{\mathcal{IO}}$, simply by exchanging the role of s and e . In addition to occluded regions, this loss occurs in the reference view up to the depthmap, where a disocclusion into the reference camera’s view is followed by an intersection.

To assist the network, we add two auxiliary losses that are true statistically: \mathcal{L}_{sep} represents a prior that surfaces tend to be separated by distances and \mathcal{L}_{ent} captures a property of the DRDF that is true in the limit if our observations are randomly chosen.

\mathcal{L}_{sep} : Minimum Separation Loss. Since the cameras never sees the insides of objects, there is no incentive to predict inside objects. This prevents the generation of zero crossings, e.g., after the first intersection. To assist the network, we add a loss \mathcal{L}_{sep} that assumes that surfaces are separated by a minimum distance unless there is evidence otherwise. We continue the DRDF’s known value for $t = 0.2m$ before and after each intersection event to make a continued value c . We then penalize $\mathcal{L}_{\text{sep}}(y) = |y - c|$, so long as there is no conflicting free space evidence from another view.

\mathcal{L}_{ent} : Sign Entropy Loss. The occlusion-based constraints

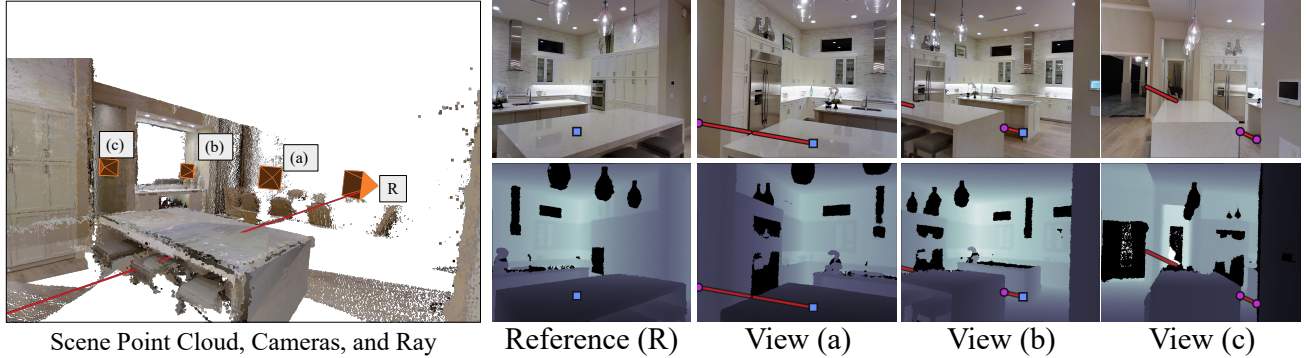


Figure 5. A red ray originates from the reference camera and intersects the table before entering the island. This ray is seen by three other auxiliary views (a, b, c). For each view we show just one of the many segments in that view for readability, showing intersections as blue ■s and occlusions as purple ●s. The segments in views (a) and (b) are of the form OI , where the O comes from the ray entering the frustum in (a) and disoccluding in (b). In view (c), we show an OO segment. We discuss the penalties for these segments in §4.1

can be satisfied by making y positive or negative. In theory, after the first intersection, d_{DR} is positive half of the time. Learning to produce the signs uniformly with gradient descent is difficult due to the large loss between the acceptable solutions. To encourage a uniform distribution of signs, we would like to maximize the entropy of the distribution of the signs of the predictions (achieved when the distribution is uniform). Since sign is non-differentiable, we optimize a differentiable surrogate. Suppose Y is the set of predictions at points that are occluded in the reference view and H is the binary entropy $H(p) = p \log(p) + (1-p) \log(1-p)$. Then our loss is $\mathcal{L}_{ent} = H(\sum_{y \in Y} \sigma(y/\tau) / |Y|)$ where σ functions like a soft-sign. As seen in the supplement, \mathcal{L}_{ent} can be minimized by distributing Y symmetric about 0. Using \mathcal{L}_{ent} improves the prediction and Scene F1 by 2 points.

4.2. Implementation Details

View Selection. For every candidate auxiliary view we compute the fraction of visible points in this view that are occluded in the reference view. Auxiliary views with large number of such points provide supervision for key occluded regions in the reference view. We sample up to 20 auxiliary views per reference image from this set of views.

Sampling Strategy. Given a set of fixed auxiliary views and a reference view, we sample over 200 rays per input image with 512 points per ray. We re-balance this set by sampling $20K$ points that are visible and $20K$ points that are occluded in reference view. We rebalance as most points are from the region between the camera and the first hit.

Combining segments from different views. We merge information from multiple posed RGBD images to produce a concise merged set of non-overlapping segments along the ray. This prevents double-counting losses (e.g., if a region of the ray is seen by multiple auxiliary views). We safely merge segments that provide the same information: e.g., if one depthmap provides an OO segment that is contained within another depthmap’s II segment, then the OO

segment can be safely dropped since $\mathcal{L}_{OO}(y) \leq \mathcal{L}_{II}(y)$. When segments disagree (e.g., due to inaccurate poses), we keep the segment with more auxiliary views in agreement. This approach handles merging \mathcal{L}_{sep} : \mathcal{L}_{sep} is seen by no auxiliary views, so any visible freespace overrules it.

Network Architecture. We follow [27] to facilitate fair comparison. Additionally, we clamp the outputs of our network to be $\in [-1, 1]$ by applying a tanh activation and adjust the loss to account for this clipping.

Training. We follow a two-stage training procedure. We first train with only the reference view followed by adding auxiliary views. In the first stage, we train for 100 epochs minimizing \mathcal{L}_{OI} and \mathcal{L}_{sep} . We then train for 100 epochs with auxiliary losses, minimizing a sum of the segment losses \mathcal{L}_{II} , \mathcal{L}_{IO} , \mathcal{L}_{OI} , \mathcal{L}_{OO} , and \mathcal{L}_{sep} as well as $\lambda \mathcal{L}_{ent}$ with λ set to 0.1 to balance loss scales. We minimize the loss with AdamW [26, 33] as the optimizer with learning rate warmup for 0.5% of the iterations followed by cosine learning rate decay with maximum value 3×10^{-4} . Our models are implemented using PyTorch [39] and visualizations in this paper use PyTorch3D [30, 42].

5. Experiments

We evaluate our method to address three experimental questions. First, we examine how training with RGBD data compares with mesh supervision. Next, we test how RGBD and mesh supervision compare when one has less complete scans. Finally, we show that our method can quickly adapt to multiple posed RGBD inputs.

Metrics. Throughout, we follow [27] and use two metrics that evaluate predictions against a ground-truth mesh. Following [45, 53], these metrics are based on: Accuracy/Acc (the fraction of predicted points that are within t to a ground-truth point), Completeness/Cmp (the fraction of ground-truth points that are within t to a predicted point), as well as F1 (the harmonic mean of Accuracy and Completeness). t for both Acc and Cmp is 0.5m. In

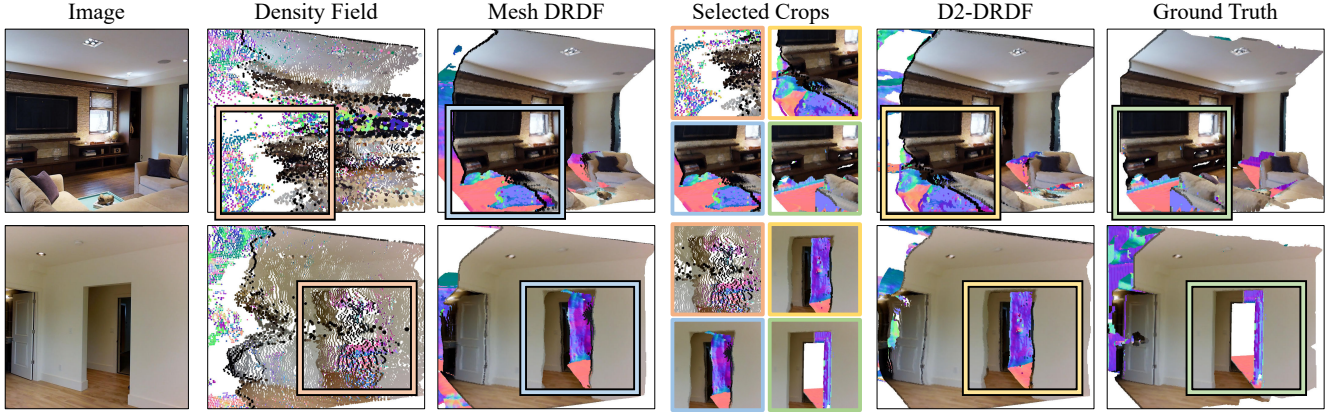



Figure 6. Comparison with baselines. 3D outputs generated by all methods trained on Matterport3D. We color the first intersection with image colors and occluded intersections with computed surface normals. We highlight regions of interest in the reconstructions in selected crops. D2-DRDF achieves results on par with Mesh DRDF while the density fields baselines fails to model the occluded parts faithfully. In row 1, our method recovers the back of sofa, and a hidden room behind the hallway in row 2. Surface Normal Map 

(Scene Acc/Cmp/F1), we evaluate the predicted mesh of the full scene against the ground-truth, using 10K samples per mesh. In (Ray Occ. Acc/Cmp/F1), We evaluate the performance on occluded points, evaluating per-ray and then averaging. We define occluded points for both the ground truth and prediction as any surface past the first intersection. Ray Occ. is a challenging metric as mistakes in one ray cannot be accounted for in another ray.

Datasets. We use Matterport3D [1] as our primary dataset following an identical setup to [28]. We choose Matterport3D because it has substantial occluded geometry (unlike ScanNet [8]) and was captured by a real scanner (unlike 3DFront [15], which is synthetic). We note that we use the *raw images* captured by the scanner rather than re-renders. We follow the split from [27], which splits train/val/test by house into 60/15/15 houses. After filtering and selecting images, there are 15K/1K/1K input images for each split.

5.1. Mesh Prediction Results

Baselines. Our primary point of comparison is (1) **Mesh DRDF** [27], which learns to predict the DRDF from direct mesh supervision. For context, we also report the baselines from [27] and summarize them: (2) **LDI** [46] predicts a set of 4 layered depthmaps, where the first represents the depth and the next three represent occluded intersections; (3) **UDF** [5] predicts an unsigned scene distance function; (4) **URDF** [5] predicts an unsigned **ray** distance function; (5) **ORF** predicts an occupancy function, or whether the surface is within a fixed distance. We note that for each of these approaches, there are a number of variants (e.g., of finding intersections in a URDF along a ray). We report the highest performance reported by [27], who document extensive and detailed tuning of these baselines.

Our final baseline, (6) **Density Field** [57] tests the value of predicting a DRDF compared to a density. Like our



Figure 7. Novels Views Comparison between D2-DRDF and Ground Truth (GT) from novel views. Rows 1, 2 are from unseen images on Matterport3D [1] and 3,4 from Omnidata [11]. Our method trained with *only* RGBD data recovers occluded **empty floors**, **kitchen cabinets** (row 1) and **sides of kitchen island**(row 3).

method, Density Field only uses posed RGBD data for supervision and does not depend on mesh supervision. We adapt pixelNeRF [57] to our setting, modifying the implementation from [41] to permit training from a single reference view as input and multiple auxiliary views for su-

Table 1. Matterport3D [1] Acc/Comp/F1Score. We separate methods that use ground-truth mesh from ones using posed RGBD by a horizontal line. D2-DRDF is comparable to the best Mesh supervised method DRDF, and is better than all other mesh based methods on Scene and Ray F1 scores.

Method	Scene			Ray		
	Acc	Cmp	F1	Acc	Cmp	F1
LDI [46]	66.2	72.4	67.4	13.9	42.8	19.3
UDF [5]	58.7	76.0	64.7	15.5	23.0	16.6
Mesh ORF	73.4	69.4	69.6	26.2	20.5	21.6
URDF [5]	74.5	67.1	68.7	24.9	20.6	20.7
DRDF [27]	75.4	72.0	71.9	28.4	30.0	27.3
D2-DRDF	73.7	73.5	72.1	28.2	22.6	25.1
Density Field [57]	45.8	80.2	57.5	24.8	14.0	17.9

pervision. In addition to the standard color loss, we supervised the network to match the ground-truth auxiliary depth RGBD depth as in [10]. At inference time, we integrate along the ray to decode a set of intersections as in [10].

Qualitative Results. We show qualitative results from our method and the baselines in Fig 6. Our methods outputs are comparable to the outputs of the mesh based DRDF that is trained with much stronger supervision. The density fields approach struggles to model the occluded hits and has a lots of floating blobs in the empty space. We show additional results from D2-DRDF in Fig. 7 where it recovers occlude kitchen cabinets, empty floor space behind kitchen island, sides of kitchen island, and hollow bed.

Quantitative Results. We report results in Table 1. Without access to ground-truth meshes, our approach slightly outperforms Mesh-based DRDF on Scene F1 and approaches its performance on Ray Occ. Our approach matches it in accuracy (i.e., the fraction of predicted points that are correct) but does worse on completeness. We hypothesize that this is due to mesh-based techniques obtaining supervision in adjacent rooms. Nonetheless, our approach outperforms all other baselines besides mesh-supervised DRDF by a large margin (25.1 vs 21.6 F1). D2-DRDF has substantially higher performance compared to the Density Field baseline. This shows that density fields representation has difficulties in learning from a supervision of less views with large occlusions.

5.2. Training on Incomplete Data

One advantage to using posed RGBD data compared to meshes is that it opens the door to learning from more data that is of lower quality. One can use lots of lower-quality scene captures that would produce poor meshes due to substantial amounts of incomplete data. In the case of direct mesh supervision methods, having an incomplete mesh may lead to incorrect signals for a mesh-based system: for instance, the network may learn to predict that an intersec-

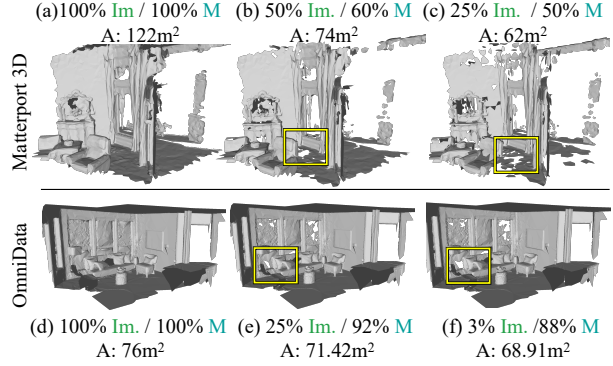


Figure 8. Mesh Degradation We show examples of how drop in images creates holes, reduces mesh area available for supervision. In (c), after retaining $\approx \frac{1}{4}$ of image data, we only lose $\approx \frac{1}{2}$ the mesh area. **Im.:** Image Coverage, **M:** Mesh Coverage, **A:** Area

tion *does not* exist simply because it was not scanned. In contrast, for D2-DRDF, missing data simply increases the fraction of points without supervision. We now test this hypothesis by reducing the number of views in datasets.

Optimistic Degradation Setup (ODS). We simulate the degradation of dataset collection by subsampling views. To avoid conflating errors in training with suboptimal meshing with incomplete data, we *optimistically* degrade the meshes to provide an upper bound on supervision. We assume that the mesh with fewer views is identical to the mesh from all views, minus triangles with no vertices in any view. We show ODS mesh examples in Fig. 8.

We degrade meshes by selecting $1/2^i$ views per dataset for an increasing i . While reducing the views linearly impacts the sample count for RGBD training setup, it has a non-linear impact on mesh completeness since a triangle is removed only if *all* of the views seeing it are removed (which is unlikely until most views are removed). For any given image retention (**Im.**) %, mesh coverage (**M**) % degrades less giving an edge to mesh based methods.

Usually, when dealing with limited data, we use a method called *Screened Poisson Reconstruction*(SPR) [25]. However, SPR does not perform well when there is not enough data available. To avoid conflating errors caused by poor quality inadequate meshing, we establish an upper limit on the performance of methods that rely on direct supervision. Our ODS strategy is much better than using SPR, but it cannot be used in real-world scenarios. We employ Open3D’s [60]’s SPR with hyper-parameters similar to [1] to reconstruct meshes.

Datasets. We evaluate on Matterport3D [1] and apply our method as is without any modifications on OmniData [11] which has a substantially different image view distribution, more rooms and more floors compared to Matterport3D.

Quantitative Results. We compare models trained on different amount of data available for supervision by using metrics defined in §5.1. In Tab. 2 we compare against the

Table 2. Robustness to Sparse Data Performance on partial data on Matterport3D[1]. We compare (SPR and ODS) trained DRDF [27] which uses mesh supervision and (Depth) Depth-based DRDF (ours), which uses posed RGBD supervision. In each row, we degrade the training data and report test performance of the trained model. At 100% data there is no M degradation for ODS or SPR. Our approach is more robust to drop in Im.: at 50% view sparsity (Im.), models using ODS or SPR suffer substantial performance drops. Scene F1 drop by 16.3 for SPR; 3.5 for ODS; 2.1 for Depth (ours).

Im. %	M %	Scene F1			Ray Occ F1		
		SPR	ODS	Depth	SPR	ODS	Depth
100	100	71.9	71.9	72.1	27.3	27.3	25.1
50	56	55.6	68.4	70.0	21.4	23.6	24.4
25	43	56.8	66.8	70.0	21.5	21.2	24.9

Table 3. Robustness to Sparse Data Performance of ODS (mesh) and Depth (ours) based DRDF on partial Omnidata [11] following the same setup as Table 2. RGBD-based training is substantially more robust to partial data.

Im. %	M %	Scene F1		Ray Occ F1	
		ODS	Depth	ODS	Depth
25	86	77.2	72.8	40.3	32.1
12.5	83	75.3 (-1.9)	70.9 (-1.9)	37.1. (-3.2)	29.3 (-2.8)
6.3	78	73.4 (-3.8)	71.8 (-1.0)	32.6 (-7.7)	28.1 (-3.0)
3	69	69.8 (-7.4)	70.4 (-2.4)	20.3 (-20.0)	26.7 (-5.4)

Mesh-DRDF trained ODS & SPR meshes. For any given view sampling level, the supervised M % area is high, resulting in stronger supervision for methods trained with ODS than RGBD. However, on Matterport3D, our method outperforms DRDF on all metrics at 25% Im. and outperforms all other baseline methods in Scene F1 at 100% Im..

In Tab. 3 we show robustness trends on OmniData. At 25% Im./ 86% M completion, mesh-based does better. We hypothesize this gain is due to better handling of estimates beyond the room. However, as Im. reduces mesh degrades, resulting steep fall for ODS DRDF: with 3% Im./ 69% M, MeshDRDF’s Ray Occ F1 drops by 20 points; ours is reduced by just 5.4. Moreover for SPR DRDF, at low Im. values, meshing performs dismally and the poor meshing performance translates into poor reconstruction performance: at 6% Im., training DRDF on SPR meshes produces a Ray F1 of just 5.7% and a scene F1 of 40.9%.

5.3. Adapting With Multiple Inputs

Since D2-DRDF can directly train on posed RGBD images, this enables test-time adaptation given a few auxiliary posed RGBD images. We start with the pre-trained model from §5.1 and then fine-tune for 500 iterations.

Dataset and Metrics. We generate 300 quadruplets of scenes consisting of a reference view as well as three auxiliary RGBD images with poses. These three auxiliary views

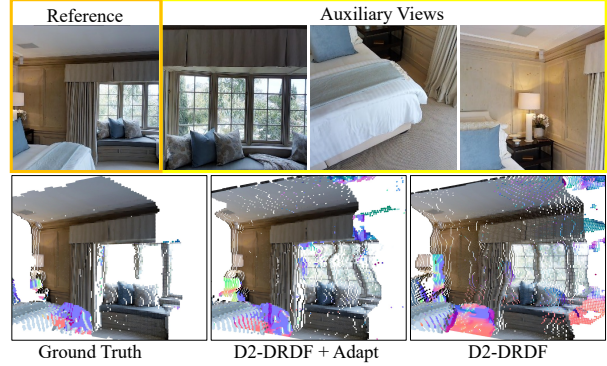


Figure 9. Adaptation to a few images We optimize our pre-trained model on the three auxiliary views (yellow) and a reference view (orange). The optimization lets the model fix reconstructions in occluded regions (e.g. bedside wall) as seen D2-DRDF + Adapt.

Table 4. Scene, Ray F1 Score for Adaptation Using D2-DRDF for adaptation produces better reconstructions that outperform the baselines by a large margin on Ray Occ F1 scores (by 8.9 points).

Method	Scene			Ray		
	Acc	Cmp	F1	Acc	Cmp	F1
D2-DRDF (Full)	79.2	76.4	76.0	36.2	33.6	34.9
D2-DRDF (Depth)	78.4	70.5	72.5	33.0	21.4	26.0
D2-DRDF (Scratch)	66.4	70.3	66.0	23.7	27.4	25.4
Density Field [57]	77.3	74.9	74.6	12.8	9.9	11.2

are randomly sampled from views that overlap with occluded parts of the reference view (see supp.). We evaluate inferred 3D using the metrics as §5.1.

Baselines. The baselines from §5.1 cannot operate in these settings, since they require meshes for training. Therefore we compare against a number of depth-map-based methods as well as ablations to give context to our results: (1) **Density Field** fine-tunes the density field baseline model from §5.1; (2) **Depth-Pretrained** fine-tunes D2-DRDF starting with the model from the first stage of training; (3) **Scratch Training** fine-tunes D2-DRDF from scratch.

Results. Our loss and penalty formulations lend themselves well to test-time adaptation as shown in Fig. 9. After adaptation, D2-DRDF can resolve uncertainties like the corner of the bed and wall. Table 4 shows D2-DRDF does the best.

6. Conclusion

We presented a method for learning to predict 3D from a single image using implicit functions while requiring only posed RGBD supervision. We believe our method can unlock new avenues with posed RGBD data becoming available from both consumers as well as robotic agents.

Acknowledgments. We thank Richard H., Sarah J., Dandan S., Shengyi Q., Karan D., Mohammed B., Chris R. for feedback. Toyota Research Institute (“TRI”) provided funds to assist the authors with their research but this article solely reflects the opinions and conclusions of its authors and not TRI or any Toyota entity.

References

- [1] Angel Chang, Angela Dai, Thomas Funkhouser, Maciej Halber, Matthias Niessner, Manolis Savva, Shuran Song, Andy Zeng, and Yinda Zhang. Matterport3d: Learning from rgb-d data in indoor environments. *3DV*, 2017. 2, 6, 7, 8
- [2] Angel X Chang, Thomas Funkhouser, Leonidas Guibas, Pat Hanrahan, Qixing Huang, Zimo Li, Silvio Savarese, Manolis Savva, Shuran Song, Hao Su, et al. Shapenet: An information-rich 3d model repository. *arXiv preprint arXiv:1512.03012*, 2015. 2
- [3] Anpei Chen, Zexiang Xu, Andreas Geiger, Jingyi Yu, and Hao Su. Tensorf: Tensorial radiance fields. *ECCV*, 2022. 2
- [4] Weifeng Chen, Zhao Fu, Dawei Yang, and Jia Deng. Single-image depth perception in the wild. *NeurIPS*, 2016. 2
- [5] Julian Chibane, Aymen Mir, and Gerard Pons-Moll. Neural unsigned distance fields for implicit function learning. In *NeurIPS*, 2020. 2, 6, 7
- [6] Christopher B Choy, Danfei Xu, JunYoung Gwak, Kevin Chen, and Silvio Savarese. 3d-r2n2: A unified approach for single and multi-view 3d object reconstruction. In *ECCV*, 2016. 1, 2
- [7] Brian Curless and Marc Levoy. A volumetric method for building complex models from range images. In *Proceedings of the 23rd annual conference on Computer graphics and interactive techniques*, pages 303–312, 1996. 2, 3
- [8] Angela Dai, Angel X Chang, Manolis Savva, Maciej Halber, Thomas Funkhouser, and Matthias Nießner. Scannet: Richly-annotated 3d reconstructions of indoor scenes. In *CVPR*, 2017. 2, 6
- [9] Angela Dai, Matthias Nießner, Michael Zollhöfer, Shahram Izadi, and Christian Theobalt. Bundlefusion: Real-time globally consistent 3d reconstruction using on-the-fly surface reintegration. *ACM Transactions on Graphics (ToG)*, 36(4):1, 2017. 2
- [10] Kangle Deng, Andrew Liu, Jun-Yan Zhu, and Deva Ramanan. Depth-supervised nerf: Fewer views and faster training for free. In *CVPR*, 2022. 7
- [11] Ainaz Eftekhari, Alexander Sax, Jitendra Malik, and Amir Zamir. Omnidata: A scalable pipeline for making multi-task mid-level vision datasets from 3d scans. In *ICCV*, 2021. 6, 7, 8
- [12] David Eigen and Rob Fergus. Predicting depth, surface normals and semantic labels with a common multi-scale convolutional architecture. *ICCV*, 2015. 2
- [13] Haoqiang Fan, Hao Su, and Leonidas J Guibas. A point set generation network for 3d object reconstruction from a single image. In *CVPR*, 2017. 2
- [14] David F Fouhey, Abhinav Gupta, and Martial Hebert. Data-driven 3d primitives for single image understanding. In *ICCV*, 2013. 2
- [15] Huan Fu, Bowen Cai, Lin Gao, Lingxiao Zhang, Cao Li, Qixun Zeng, Chengyue Sun, Yiyun Fei, Yu Zheng, Ying Li, Yi Liu, Peng Liu, Lin Ma, Le Weng, Xiaohang Hu, Xin Ma, Qian Qian, Rongfei Jia, Binqiang Zhao, and Hao Zhang. 3d-front: 3d furnished rooms with layouts and semantics. *ICCV*, 2021. 6
- [16] Huan Fu, Bowen Cai, Lin Gao, Ling-Xiao Zhang, Jiaming Wang, Cao Li, Qixun Zeng, Chengyue Sun, Rongfei Jia, Binqiang Zhao, et al. 3d-front: 3d furnished rooms with layouts and semantics. In *ICCV*, 2021. 2
- [17] Rohit Girdhar, David F Fouhey, Mikel Rodriguez, and Abhinav Gupta. Learning a predictable and generative vector representation for objects. In *ECCV*. Springer, 2016. 1, 2
- [18] Georgia Gkioxari, Jitendra Malik, and Justin Johnson. Mesh r-cnn. In *ICCV*, 2019. 2
- [19] Thibault Groueix, Matthew Fisher, Vladimir G Kim, Bryan C Russell, and Mathieu Aubry. A papier-mâché approach to learning 3d surface generation. In *CVPR*, 2018. 1, 2
- [20] Po-Han Huang, Kevin Matzen, Johannes Kopf, Narendra Ahuja, and Jia-Bin Huang. Deepmvs: Learning multi-view stereopsis. In *CVPR*, 2018. 2
- [21] Shahram Izadi, David Kim, Otmar Hilliges, David Molyneaux, Richard Newcombe, Pushmeet Kohli, Jamie Shotton, Steve Hodges, Dustin Freeman, Andrew Davison, et al. Kinectfusion: real-time 3d reconstruction and interaction using a moving depth camera. In *Proceedings of the 24th annual ACM symposium on User interface software and technology*, pages 559–568, 2011. 2
- [22] Hamid Izadinia, Qi Shan, and Steven M Seitz. Im2cad. In *CVPR*, 2017. 2
- [23] Michael Janner, Jiajun Wu, Tejas D Kulkarni, Ilker Yildirim, and Josh Tenenbaum. Self-supervised intrinsic image decomposition. *NeurIPS*, 2017. 2
- [24] Ziyu Jiang, Buyu Liu, Samuel Schulter, Zhangyang Wang, and Manmohan Chandraker. Peek-a-boo: Occlusion reasoning in indoor scenes with plane representations. In *CVPR*, 2020. 2
- [25] Michael Kazhdan, Matthew Bolitho, and Hugues Hoppe. Poisson surface reconstruction. In *Proceedings of the fourth Eurographics symposium on Geometry processing*, 2006. 2, 7
- [26] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *ICLR*, 2015. 5
- [27] Nilesh Kulkarni, Justin Johnson, and David F. Fouhey. Directed ray distance function for 3d scene reconstruction. In *ECCV*, 2022. 1, 2, 3, 5, 6, 7, 8
- [28] Nilesh Kulkarni, Ishan Misra, Shubham Tulsiani, and Abhinav Gupta. 3d-relnet: Joint object and relational network for 3d prediction. In *ICCV*, 2019. 2, 6
- [29] Kiriakos N Kutulakos and Steven M Seitz. A theory of shape by space carving. *IJCV*, 2000. 2
- [30] Christoph Lassner and Michael Zollhöfer. Pulsar: Efficient sphere-based neural rendering. *arXiv:2004.07484*, 2020. 5
- [31] Chen-Hsuan Lin, Chen Kong, and Simon Lucey. Learning efficient point cloud generation for dense 3d object reconstruction. In *AAAI*, volume 32, 2018. 2
- [32] Chen Liu, Kihwan Kim, Jinwei Gu, Yasutaka Furukawa, and Jan Kautz. Planercnn: 3d plane detection and reconstruction from a single image. In *CVPR*, 2019. 2
- [33] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. *ICLR*, 2019. 5
- [34] Ravi Malladi, James A Sethian, and Baba C Vemuri. Shape modeling with front propagation: A level set approach. *TPAMI*, 1995. 2
- [35] Lars Mescheder, Michael Oechsle, Michael Niemeyer, Sebastian Nowozin, and Andreas Geiger. Occupancy networks: Learning 3d reconstruction in function space. In *CVPR*, 2019. 1, 2

- [36] 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. *arXiv preprint arXiv:2003.08934*, 2020. 2
- [37] Yinyu Nie, Xiaoguang Han, Shihui Guo, Yujian Zheng, Jian Chang, and Jian Jun Zhang. Total3dunderstanding: Joint layout, object pose and mesh reconstruction for indoor scenes from a single image. In *CVPR*, 2020. 2
- [38] Jeong Joon Park, Peter Florence, Julian Straub, Richard Newcombe, and Steven Lovegrove. DeepSDF: Learning continuous signed distance functions for shape representation. In *CVPR*, 2019. 1, 2
- [39] 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 *NeurIPS*. 2019. 5
- [40] Songyou Peng, Michael Niemeyer, Lars Mescheder, Marc Pollefeys, and Andreas Geiger. Convolutional occupancy networks. In *ECCV*, 2020. 2
- [41] Shengyi Qian, Alexander Kirillov, Nikhila Ravi, Deendra Singh Chaplot, Justin Johnson, David F Fouhey, and Georgia Gkioxari. Recognizing scenes from novel viewpoints. *arXiv preprint arXiv:2112.01520*, 2021. 6
- [42] Nikhila Ravi, Jeremy Reizenstein, David Novotny, Taylor Gordon, Wan-Yen Lo, Justin Johnson, and Georgia Gkioxari. Accelerating 3d deep learning with pytorch3d. *arXiv:2007.08501*, 2020. 5
- [43] Mike Roberts, Jason Ramapuram, Anurag Ranjan, Atulit Kumar, Miguel Angel Bautista, Nathan Paczan, Russ Webb, and Joshua M. Susskind. Hypersim: A photorealistic synthetic dataset for holistic indoor scene understanding. In *ICCV*, 2021. 2
- [44] Shunsuke Saito, Zeng Huang, Ryota Natsume, Shigeo Morishima, Angjoo Kanazawa, and Hao Li. Pifu: Pixel-aligned implicit function for high-resolution clothed human digitization. In *ICCV*, 2019. 1, 2
- [45] Steven M Seitz, Brian Curless, James Diebel, Daniel Scharstein, and Richard Szeliski. A comparison and evaluation of multi-view stereo reconstruction algorithms. In *CVPR*, 2006. 5
- [46] Jonathan Shade, Steven Gortler, Li-wei He, and Richard Szeliski. Layered depth images. In *Proceedings of the 25th annual conference on Computer graphics and interactive techniques*, pages 231–242, 1998. 6, 7
- [47] Jian Shi, Yue Dong, Hao Su, and Stella X Yu. Learning non-lambertian object intrinsics across shapenet categories. In *CVPR*, 2017. 2
- [48] Vincent Sitzmann, Julien Martel, Alexander Bergman, David Lindell, and Gordon Wetzstein. Implicit neural representations with periodic activation functions. *NeurIPS*, 2020. 2
- [49] Shuran Song, Samuel P Lichtenberg, and Jianxiong Xiao. Sun rgb-d: A rgb-d scene understanding benchmark suite. In *CVPR*, 2015. 2
- [50] Shuran Song, Fisher Yu, Andy Zeng, Angel X Chang, Manolis Savva, and Thomas Funkhouser. Semantic scene completion from a single depth image. *CVPR*, 2017. 2
- [51] Jiaming Sun, Yiming Xie, Linghao Chen, Xiaowei Zhou, and Hujun Bao. Neuralrecon: Real-time coherent 3d reconstruction from monocular video. In *CVPR*, 2021. 2
- [52] Xingyuan Sun, Jiajun Wu, Xiuming Zhang, Zhoutong Zhang, Chengkai Zhang, Tianfan Xue, Joshua B Tenenbaum, and William T Freeman. Pix3d: Dataset and methods for single-image 3d shape modeling. In *CVPR*, 2018. 2
- [53] Maxim Tatarchenko, Stephan R Richter, René Ranftl, Zhuwen Li, Vladlen Koltun, and Thomas Brox. What do single-view 3d reconstruction networks learn? In *CVPR*, 2019. 5
- [54] Kaixuan Wang and Shaojie Shen. Mvdepthnet: Real-time multiview depth estimation neural network. In *3DV*, 2018. 2
- [55] X. Wang, David F. Fouhey, and A. Gupta. Designing deep networks for surface normal estimation. In *CVPR*, 2015. 2
- [56] Jianglong Ye, Yuntao Chen, Naiyan Wang, and Xiao-long Wang. Gifs: Neural implicit function for general shape representation. In *CVPR*, 2022. 2
- [57] Alex Yu, Vickie Ye, Matthew Tancik, and Angjoo Kanazawa. pixelNeRF: Neural radiance fields from one or few images. In *CVPR*, 2021. 2, 3, 6, 7, 8
- [58] Zehao Yu, Jia Zheng, Dongze Lian, Zihan Zhou, and Shenghua Gao. Single-image piece-wise planar 3d reconstruction via associative embedding. In *CVPR*, 2019. 2
- [59] Cheng Zhang, Zhaopeng Cui, Yinda Zhang, Bing Zeng, Marc Pollefeys, and Shuaicheng Liu. Holistic 3d scene understanding from a single image with implicit representation. In *CVPR*, 2021. 2
- [60] Qian-Yi Zhou, Jaesik Park, and Vladlen Koltun. Open3D: A modern library for 3D data processing. *arXiv:1801.09847*, 2018. 7
- [61] Shizhan Zhu, Sayna Ebrahimi, Angjoo Kanazawa, and Trevor Darrell. Differentiable gradient sampling for learning implicit 3d scene reconstructions from a single image. In *ICLE*, 2022. 1, 2