

SDFDiff: Differentiable Rendering of Signed Distance Fields for 3D Shape Optimization

Yue Jiang, Dantong Ji, Zhizhong Han, Matthias Zwicker
 University of Maryland, College Park
 Maryland, USA

{yuejiang, dji, h312h, zwicker}@cs.umd.edu

Abstract

We propose *SDFDiff*, a novel approach for image-based shape optimization using differentiable rendering of 3D shapes represented by signed distance functions (SDF). Compared to other representations, SDFs have the advantage that they can represent shapes with arbitrary topology, and that they guarantee watertight surfaces. We apply our approach to the problem of multi-view 3D reconstruction, where we achieve high reconstruction quality and can capture complex topology of 3D objects. In addition, we employ a multi-resolution strategy to obtain a robust optimization algorithm. We further demonstrate that our SDF-based differentiable renderer can be integrated with deep learning models, which opens up options for learning approaches on 3D objects without 3D supervision. In particular, we apply our method to single-view 3D reconstruction and achieve state-of-the-art results.

1. Introduction

The “vision as inverse graphics” or “inverse rendering” strategy has long been attractive as a conceptual framework to solve inverse problems such as recovering shape or appearance models from images. In this analysis-by-synthesis approach, the goal is to reproduce given input images by synthesizing them using an image formation model, possibly including shape, appearance, illumination, and camera parameters. Solving this problem implies finding suitable model parameters (shape, appearance, illumination, camera) that describe the underlying scene. While conceptually simple, this approach can be challenging to use in practice, because it requires a suitable parameterization of a powerful image formation model, and effective numerical techniques to solve the resulting optimization problem. Recently, automatic differentiation has attracted renewed attention to implement differentiable renderers or image formation models, which can then be used in gradient-based optimization

techniques. In particular, it is attractive to combine differentiable rendering with neural networks to solve highly ill-posed inverse problems, such as single-view 3D reconstruction.

In this paper, we address the key issue of how to represent geometry in a differentiable image formation model. In particular, we advocate using signed distance fields (SDFs) because they have several advantages over other geometry representations for inverse rendering. In contrast to triangle meshes, the surface topology is not fixed in SDFs and can adapt to the actual scene topology during optimization. Point clouds can also represent arbitrary topologies but they do not provide continuous surface reconstructions, while SDFs inherently represent continuous, watertight surfaces. In addition, SDFs can easily be used in a multi-resolution framework, which is important to avoid undesired local minima during optimization.

The main contribution of our paper is *SDFDiff*, a differentiable renderer based on ray-casting SDFs. Our renderer is integrated with a popular deep learning framework such that it can be combined with neural networks to learn how to solve highly ill-posed inverse problems. Finally, we provide an effective multi-resolution strategy to improve the robustness of gradient-based optimization. We demonstrate the usefulness of our approach using several application studies, including multi-view 3D reconstruction and single-view 3D reconstruction using a neural network. Our results demonstrate the advantages of SDFs over other surface representations such as meshes or point clouds.

In summary, we make the following contributions:

- We introduce a differentiable renderer based on ray casting SDFs, and we describe an implementation that is integrated into a standard deep learning framework. Advantages of using SDFs for differentiable rendering and shape optimization include that we can adapt the topology freely, and that the resulting shapes are guaranteed to consist of watertight surfaces.
- We present results of a multi-view 3D reconstruction

tion approach using shape optimization via differentiable rendering. Using a multi-resolution approach, our gradient-descent optimization reliably converges to high quality solutions. Our approach is able to reconstruct geometry with high level of detail and complex topology, even with few input views.

- We leverage the SDF-based differentiable renderer to train deep neural networks to perform single-view 3D shape reconstruction without 3D supervision. We demonstrate the advantages of our approach by showing how it can recover accurate 3D shapes with arbitrary topology.

2. Related Work

Learning-based 3D Reconstruction. Reconstructing 3D models from 2D images is a classic problem in computer vision. Compared to traditional multi-view stereo algorithms and shading based approaches [43], learning-based 3D reconstruction techniques can achieve impressive performance even with very sparse input consisting of very few views. Deep learning models for 3D shape understanding have been proposed for different kinds of 3D representations, including multiple views [13, 14], point clouds [51, 18], triangle meshes [27, 28], voxel grids [50, 47], and signed distance functions (SDF) [39].

However, most learning-based methods [48, 46, 10, 1, 7, 19, 23] require 3D supervision. Even though some methods [4, 16] are less supervised, they are often limited by specific settings, such as restricted lighting conditions and required annotation of object orientation. In contrast, predicting 3D shapes with differentiable renderers has recently attracted increasing attention because it promises to enable 3D reconstruction without 3D supervision, that is, by optimizing neural networks only using images as training data.

Signed Distance Functions. A distance function is a level set method [37] that, at each point in 3D, stores the distance to the closest point on the surface of the object. Signed distance fields (SDFs) [8] store a signed distance function to distinguish between the inside of outside of an object. SDFs are often discretized using uniform voxel grids [34, 36]. Compared to meshes or parametric surfaces, the implicit surface representation [31, 6] of SDFs has the advantage that it can represent arbitrary topologies. In contrast to point clouds, SDFs always represent watertight surfaces.

SDFs just recently started attracting research interest for shape analysis via deep learning. DeepSDF [39] was proposed as a learned continuous signed distance function representation of a class of shapes. Similarly, deep level sets [32] were introduced as an end-to-end trainable model that directly predicts implicit surfaces of arbitrary topology. However, these methods require 3D supervision dur-

ing training, such as a set of pairs of 3D coordinates and their corresponding SDF values [39], or voxel occupancy maps [32].

Differentiable Rendering Voxel-based differentiable renderers [47, 11, 17, 20, 50, 12, 52, 49, 35] first drew research attention. However, these methods perform volumetric ray marching instead of computing ray-surface intersections, and they are limited by the low resolution of voxel grids. To render SDFs we use sphere tracing, which is also used in scene representation networks [45]. However, they focus more on novel view synthesis rather than 3D surface reconstruction, which is our goal. Recently, much work focused on mesh-based differentiable renderers [4, 16, 21, 29, 26, 5, 9, 41, 30, 2, 42, 38, 28]. Loper and Black [29] proposed the idea of differentiable rendering in OpenDR, a general framework of mesh-based differentiable renderer. With hand-designed gradients, Kato et al. [21] proposed a neural 3D mesh renderer. These two methods focus on the primary visibility gradients using triangle meshes as geometric representations. Paparazzi [27] employed analytically computed gradients to adjust the location of vertices. Similarly, SoftRas [28] assigned each pixel to all faces of a mesh in a probabilistic rasterization framework. Although these methods enable to learn 3D mesh reconstruction without 3D supervision, they are restricted to a fixed, usually spherical mesh topology. Many of these mesh-based rendering approaches are differentiable with respect to geometry [4, 21, 29, 26, 9, 42, 28], lighting models [5], textures [26, 41, 2], or materials [41, 30, 2]. Petersen et al. [40] applied their mesh-based differentiable renderer to real images.

Finally, differentiable rendering has also been applied to perform Monte Carlo ray tracing [25] and point cloud rendering [18, 19, 33]. Insafutdinov et al. [18] proposed a differentiable renderer for point clouds with visibility modeling by conducting orthogonal projection on voxelized 3D space holding the point cloud. Similarly, surface splatting [51] was employed to model the visibility in point cloud rendering. Although point clouds can be easily acquired using range sensing technology, such as the Microsoft Kinect and LIDAR systems, they do not explicitly represent topology and require post-processing to produce watertight surfaces.

3. Overview

We propose a novel approach for image-based 3D shape optimization by leveraging signed distance functions as the geometric representation to perform differentiable rendering. Given a set of parameters Θ representing the geometry description, lighting model, camera position, etc, a renderer R can be written as a forward operator that produces an

image I by computing $I = R(\Theta)$. In contrast, optimizing geometry and other scene parameters from images is a backward process that can be seen as an inverse rendering problem. Given a desired target image I , our goal is to get the set of parameters $\Theta = R^{-1}(I)$ that produces the target image. However, the rendering process itself is not invertible. Hence, instead of solving the inverse rendering problem directly, we can formulate it as an energy minimization problem,

$$\Theta^* = \arg \min_{\Theta} \mathcal{L}_{\text{img}}(R(\Theta), I) \quad (1)$$

where \mathcal{L}_{img} is a loss function measuring the distance between the target image and the rendered image from the 3D object. In practice, the loss is typically accumulated over a set of multiple target images. Getting the desired parameters Θ^* is equivalent to minimizing the loss \mathcal{L} . While all rendering parameters including geometry, illumination, camera pose, and surface appearance could in theory be recovered from images this way, we focus on shape optimization in this paper and assume the other parameters are known.

To enable practical, gradient-based optimization algorithms, a key issue is to obtain the gradient of $\mathcal{L}(R(\Theta), I)$ with respect to the parameters Θ . To achieve this, a differentiable renderer produces not only images from a description of the scene, but it can also provide the derivatives of pixel values with respect to scene parameters.

In this paper, we propose a novel differentiable renderer which uses signed distance functions (SDF) and camera pose as inputs and renders an image. Our SDF-based differentiable renderer leverages the ray casting algorithm and uses automatic differentiation to compute the derivatives.

4. Differentiable SDF Rendering via Sphere Tracing

We use an usual discrete SDF representation by sampling SDF values on a regular grid, and we apply a standard ray casting algorithm based on sphere tracing [15] to find the intersection points between rays and the object surface. For this purpose we employ trilinear interpolation to reconstruct a continuous SDF that can be evaluated at any desired location. This allows us to continuously represent the object surface, which is given by the zero level set of the interpolated SDF.

A key observation is that the derivatives of a given pixel with respect to rendering parameters only depend on a local neighborhood of eight SDF samples that define the value of the trilinearly interpolated SDF at the surface intersection point. In other words, the sphere tracing process itself does not need to be differentiable. Instead, only the local computations involving the local set of eight SDF samples around

the surface intersection need to be differentiable. Therefore, our approach proceeds in two stages: first, we apply sphere tracing to identify the eight samples nearest to the surface intersection. This step is not differentiable. Second, we locally compute the pixel color based on the local set of SDF samples. This step is implemented using an automatic differentiation framework to obtain the derivatives.

Note that differentiable ray marching of voxel grids has been used before [50, 52, 35]. However, in these approaches the voxels represent opacities, either as binary variables or as continuous occupancy probabilities. In these cases, ray marching through the entire volume needs to be differentiated because all voxels along a ray may influence the corresponding pixel. In addition, in these techniques the computation of surface normals and evaluation of shading models is not as straightforward as with SDFs.

4.1. Sphere Tracing

We perform ray casting via sphere tracing [15] in the first stage by starting from the ray origin, and evaluating the SDF using trilinear interpolation to find the minimum distance from the current point on the ray to the object. Then we move along the ray by that distance. Moving along the ray by the minimum distance to the object guarantees that we will never move across the boundary of the object, but it allows us to make a possibly large step towards the surface. We repeat this process until we reach the surface of the object, that is, until the SDF value at our current position on the ray is small enough, or until we leave the bounding box of the object. In practice, this stage is implemented in conventional CUDA code without support for automatic differentiation.

4.2. Differentiable Shading

In the second stage, we compute the pixel color as a function of the local SDF samples that define the SDF at the intersection point, as determined by the first stage. In practice, these computations are implemented in a framework that supports automatic differentiation, which allows us to easily obtain the derivatives of the pixel. For each pixel, the input to this stage consists of the light and camera parameters, and the eight SDF samples closest to the ray-surface intersection point. In our current implementation the computations include: getting the intersection point and the surface normal at the intersection point as a function of the basis coefficients, and evaluating a simple shading model.

To take into account the dependence of the pixel value on the ray-surface intersection point, we express the intersection point as a function of the eight local SDF samples. More specifically, let us denote the local SDF values by d_0, \dots, d_7 . To express the intersection point as a function of d_0, \dots, d_7 , we use the same approximation as in the ray casting stage. In addition, denote the current po-

sition on the ray (obtained from the ray casting stage) s , and the unit ray direction v . Then, the approximate intersection is $p(d_0, \dots, d_7) = s + \text{trilinear}(d_0, \dots, d_7; s)v$, where $\text{trilinear}(d_0, \dots, d_7; s)$ is the trilinear interpolation of the SDF at location s and considered as a function of d_0, \dots, d_7 . This approximation is conservative in the sense that it is accurate only if the SDF represents a plane that is perpendicular to the ray direction v . Otherwise, $p(d_0, \dots, d_7)$ is guaranteed not to cross the true intersection along the ray.

As an alternative to our conservative approximation, one could express the intersection point exactly as the solution of the intersection of the ray $s + tv$ and the local trilinear interpolation of the SDF. That is, we could express the solution of $\text{trilinear}(d_0, \dots, d_7; s + tv) = 0$ with respect to t as a function of d_0, \dots, d_7 . However, this involves finding roots of a cubic polynomial, and we found that our much simpler approach works more robustly in practice.

To evaluate a shading model, we need the surface normal at the intersection point $p(d_0, \dots, d_7)$. Considering that the surface normal corresponds to the gradient of the SDF, we first compute gradients at the grid vertices using central finite differencing, and then trilinearly interpolate them at the intersection point $p(d_0, \dots, d_7; s)$. In summary, this leads to an expression of the normal at the intersection point as a function of SDF coefficients within an $4 \times 4 \times 4$ neighborhood around the intersection (because of central finite differencing). Finally, in our current implementation we evaluate a simple diffuse shading model.

4.3. Implementation

We implemented SDF ray casting using CUDA to leverage the computational power of GPUs. Differentiable shading is implemented in Python with the Pytorch library, which supports automatic differentiation and allows seamless integration of the renderer with deep learning and neural network training. Pytorch also leverages the GPU, and our Pytorch implementation directly accesses the output of the ray casting stage that is stored in GPU memory, avoiding any unnecessary memory transfers. We will open source our code upon publication of this work.

5. Multi-view 3D Reconstruction

In this section we describe how to perform multi-view 3D reconstruction using our differentiable renderer. This is a proof of concept, where we assume known camera poses, illumination, and surface appearance, and we only optimize over the 3D shape represented by the SDF.

5.1. Input

Our inputs for multi-view 3D reconstruction are multi-view images for the target objects with known camera poses. In addition, we initialize the SDF to a sphere. In

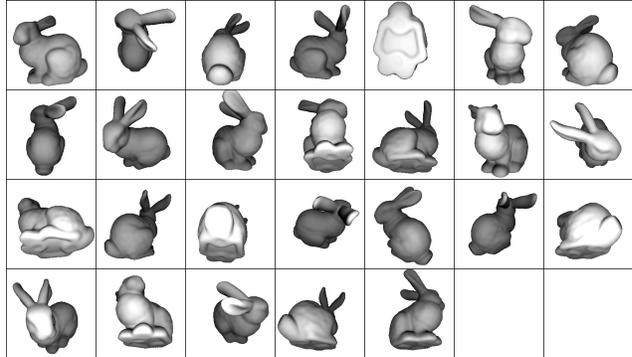


Figure 1. We use 26 input views in our multi-view reconstruction experiments as shown here for the bunny.

our current experiments, we set the camera poses to point from the center of each face and edge, and from each vertex of the bounding box towards its center, where the bounding box is a cube. Since the cube has 6 faces, 8 vertices, and 12 edges, we obtain 26 camera poses in total. Figure 1 shows the input images we used to reconstruct the bunny in our experiments.

5.2. Energy Function

As the image-based loss in our inverse rendering problem (Equation 1) we simply choose the L_2 distance between the rendered images and the target views, that is $\mathcal{L}_{\text{img}}(R(\Theta), I) = \|R(\Theta) - I\|^2$. The loss is summed over all target views. In this proof of concept scenario, the optimization parameters Θ include only the SDF values, as we assume the other rendering parameters are known.

In addition, we impose a regularization loss that ensures that the SDF values Θ represent a valid signed distance function, that is, its gradient should have unit magnitude. Writing the SDF represented by Θ as a function $f(x; \theta)$, where x is a point in 3D, the regularization loss is

$$\mathcal{L}_{\text{reg}} = \int \|1 - \|\nabla f(x; \theta)\|\|^2 dx \quad (2)$$

In practice, we obtain the gradients via finite differencing and we compute a discrete sum over the SDF grid vertices.

5.3. Iterative Optimization

We apply gradient descent optimization using ADAM [22] to iteratively optimize our SDF grid to match the target images. In addition, we accelerate convergence by greedily selecting a single view in each gradient descent step to compute the gradient, similar to active mini batch sampling. The intuition is that different views may incur image losses of varying magnitude, and we should focus on the views with large losses. Our approach first calculates the average loss for all the camera views from the result of the previous iteration. If a loss for a view is greater

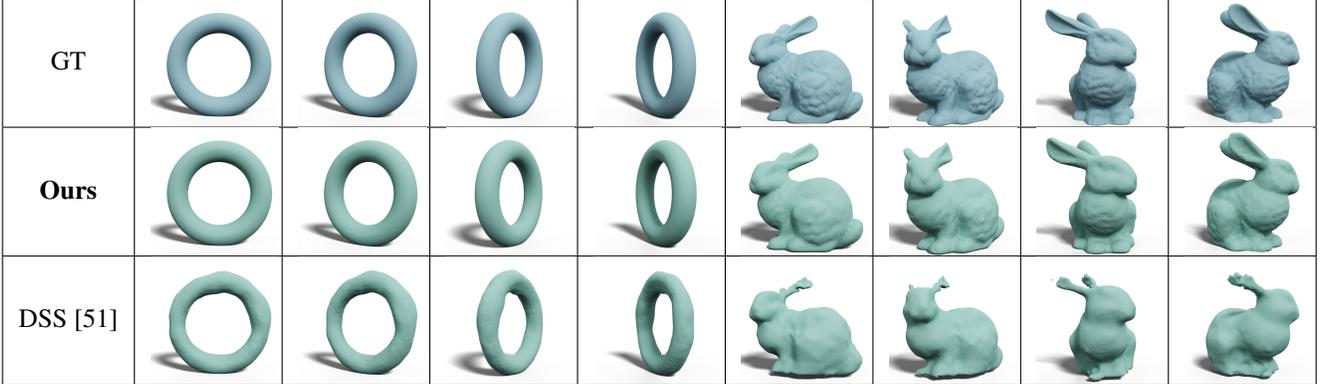


Figure 2. Multi-view reconstruction results comparing with the state-of-the-art point-based differentiable rendering approach DSS [51].

than the average loss, then during the current iteration, we update SDF until the loss for this view is less than the average (with 20 max updates). For the other views, we update the SDF five times. If one update increases the loss, then we switch to the next view directly. We stop our optimization process when the loss is smaller than a given tolerance or the step length is too small.

Reconstructing high-resolution 3D objects is challenging because gradient descent takes many iterations to eliminate low frequency errors. Therefore, we apply a coarse-to-fine multi-resolution approach. We start by initializing the SDF grid at a resolution of 8^3 to the SDF of a sphere. We then iterate between performing gradient descent optimization as described above, and increasing the grid resolution. We increase the resolution simply by performing trilinear interpolation. In our experiments, we stop at a SDF resolution of 512^3 .

To further improve the efficiency of the multiresolution scheme, we choose an appropriate image resolution for rendering corresponding to the SDF resolution at each resolution level. We determine the appropriate resolution by ensuring that a sphere with a radius equivalent to the grid spacing, and placed at the corner of the bounding box of the SDF furthest from the camera, has a projected footprint of at most the size of a 2×2 pixel block.

5.4. Experimental Results

Qualitative Results. We compare our results with DSS [51], which is a differentiable renderer for point clouds based on surface splatting [53]. We let both systems deform a sphere to fit the target object given as input. When running DSS, we adopt the same settings used in their original experiments: the system randomly selects 12 from a set of 26 views of the target in each optimization cycle, and optimizes for up to 16 cycles. We experimented with different numbers of 3D points and report the best result. For SDFDiff we use our optimization technique from Section 5.3 using the same set of 26 views. Figures 2 and 3 show the compar-

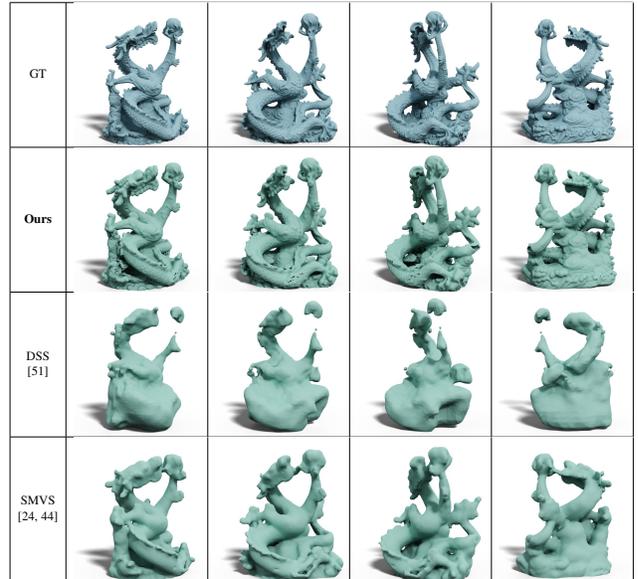


Figure 3. Multi-view reconstruction results from 26 views comparing with DSS [51] and SMVS [24, 44] on dragon.

ison between SDFDiff and DSS. DSS cannot recover the details of the mesh as accurately as SDFDiff, and is unable to reconstruct complex geometry like the Chinese dragon sculpture.

We also compare our result with SMVS [24, 44], which is a state-of-the-art shading-aware multi-view 3D reconstruction approach. We use the default settings, and provide 1000 randomly sampled views of the Chinese dragon rendered by our SDF renderer as input. Note that SMVS automatically recovers camera parameters from the input images and also estimates surface albedo and illumination, hence the comparison is not entirely fair. As shown in Figure 3, however, even with a large number of input views the SMVS output has holes and suffers from noise. In contrast, SDFDiff can obtain better results using only 26 views (with known camera poses, albedo, and illumination).

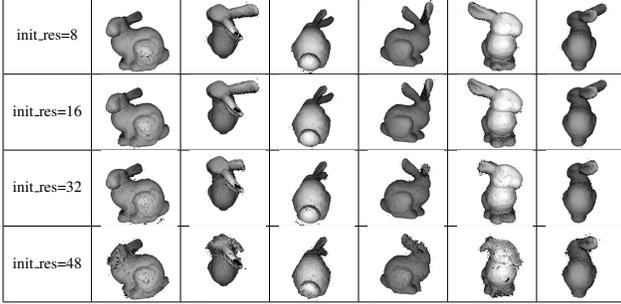


Figure 4. Given different initial resolutions, with 4 resolution stages, we can find that our 3D reconstruction results are better if the initial resolution is lower.

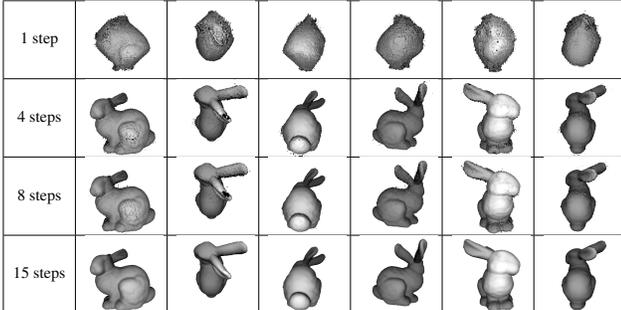


Figure 5. We fix the initial and target resolution to 8 and 64 respectively, but use different numbers of intermediate resolution stages. We find that more resolution stages can give us better results.

Quantitative Results. Table 1 compares the symmetric Hausdorff distance between ground truth and reconstructed meshes for torus, bunny and dragon. For a fair comparison, we report errors relative to the size of the bounding boxes. We observe that SDFDiff leads to smaller symmetric Hausdorff distances, which means our reconstruction results are closer to the ground truth than the other two approaches.

5.5. Parameter Study

Initial Resolution. Figure 4 shows the impact of the initial resolution in our multi-resolution scheme. We fix the number of multi-resolution steps and our target resolution being 64, and then set the initial resolution to be 8, 16, 32, and 48 respectively. We find that a lower initial resolution can reconstruct qualitatively better 3D shapes because it more robustly captures the large scale structure of the shape.

Number of Multi-Resolution Steps. In Figure 5, we show that given the fixed initial SDF resolution ($\text{init_res}=8$) and target resolution ($\text{target_res}=64$), adding more multi-resolution steps can give us better results. We demonstrate that single-resolution optimization (1 step) cannot reconstruct the object successfully, so multi-resolution scheme is necessary in our setup.

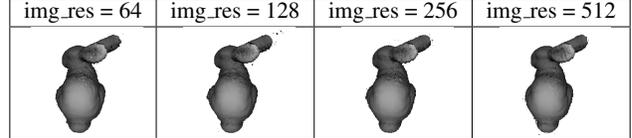


Figure 6. We show that the quality of reconstruction results are not affected much by the image resolution.

Object	Ours	DSS [51]	SMVS [24, 44]
Torus	0.015637	0.035398	N/A
Bunny	0.026654	0.109432	N/A
Dragon	0.074723	0.179456	0.097816

Table 1. Comparison of the symmetric Hausdorff distance between ground truth and reconstructed meshes for torus, bunny and dragon. SMVS could not reconstruct torus and bunny because camera pose estimation failed.

Image Resolution. As shown in Figure 6, image resolution does not significantly affect the quality of the final results, where we use the images with various resolutions in optimization.

6. Single-view Reconstruction using Deep Learning

In the following experiments, we leverage SDFDiff to train a neural network to perform single-view 3D reconstruction without 3D supervision. We use the same dataset as [21, 28], which includes 13 categories of objects from ShapeNet [3]. Each object had 24 rendered images from different views with image resolution of 64×64 . We use the same train/validate/test sets on the same dataset as in [21, 28, 50]. We use the standard reconstruction metric, *i.e.*, 3D intersection over union (IoU) [28] for quantitative comparisons.

Network. Our network contains two parts as shown in Figure 7. The first part is an Encoder-Decoder network which takes images as input and outputs coarse SDF results. The second part is a refiner network to further improve the quality of the 3D reconstruction results.

Loss Function. In addition to the energy function as shown in Section 5.2 containing the L_2 image-based loss \mathcal{L}_{img} and the SDF loss \mathcal{L}_{reg} ensuring the SDF values represent a valid signed distance function, we also add a geometry loss \mathcal{L}_{geo} that regularizes the finite difference Laplacian of the predicted SDFs. Furthermore, we use a narrow band technique to control the effects of the SDF and Laplacian losses. Since we care more about these two losses locally around the surfaces, we can use a distance-based binary mask \mathcal{M} to zero them out further away from the zero

Category	Airplane	Bench	Cabinet	Car	Chair	Display	Lamp	Speaker	Rifle	Sofa	Table	Phone	Vessel	Mean
NMR [21]	0.6172	0.4998	0.7143	0.7095	0.4990	0.5831	0.4126	0.6536	0.6322	0.6735	0.4829	0.7777	0.5645	0.6015
SoftRas (sil.) [28]	0.6419	0.5080	0.7116	0.7697	0.5270	0.6156	0.4628	0.6654	0.6811	0.6878	0.4487	0.7895	0.5953	0.6234
SoftRas (full) [28]	0.6670	0.5429	0.7382	0.7876	0.5470	0.6298	0.4580	0.6807	0.6702	0.7220	0.5325	0.8127	0.6145	0.6464
DIB-R [5]	0.570	0.498	0.763	0.788	0.527	0.588	0.403	0.726	0.561	0.677	0.508	0.743	0.609	0.612
Ours	0.6874	0.6860	0.7735	0.8002	0.6436	0.6584	0.5150	0.6534	0.5553	0.7654	0.6288	0.8278	0.6244	0.6674

Table 2. Comparison of IoU with the state-of-the-art approaches [21, 28, 5] on 13 categories of ShapeNet dataset.

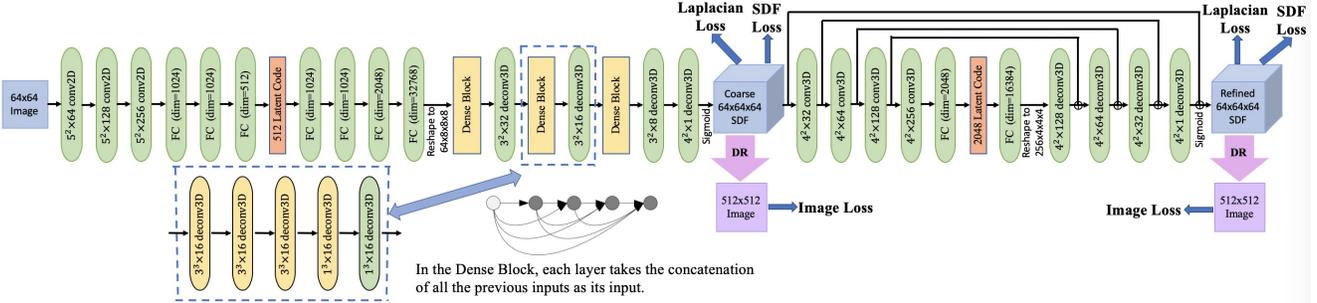


Figure 7. Network structure for single-view SDF reconstruction.

level-set. Hence the mask is defined as

$$\mathcal{M} = \|\text{SDF}\| \leq \mu \times \text{voxelSize}. \quad (3)$$

where μ is a hyperparameter to define the width of the narrow band. We currently set it to be 1.6. The final loss function is a weighted sum of the three losses,

$$\mathcal{L} = \mathcal{L}_{\text{img}} + \mathcal{M} \otimes (\lambda_1 \mathcal{L}_{\text{reg}} + \lambda_2 \mathcal{L}_{\text{geo}}). \quad (4)$$

Training Process. We first train the Encoder-Decoder part of the network alone based on the three loss terms. Then we fix the encoder and decoder and train the refiner network on the same three loss terms to get refined SDF shapes. In the end, we train all the three parts, *i.e.*, *encoder*, *decoder*, and *refiner* together to further improve the results. We do not use the multi-resolution approach in this application scenario.

Qualitative Evaluation. Figure 8 shows that our method can reconstruct detailed objects and accurately recover complicated topologies. In contrast, SoftRasterizer [28] relies on a template mesh with spherical topology and it cannot capture the complex topology of the chairs.

Quantitative Evaluation. We compare our method with the state-of-the-art approaches [21, 28, 5] in terms of 3D IoU scores in Table 2. Our method can reconstruct shapes with finer details in the 13 categories. In addition, the IoU show that our results achieve higher accuracy. Our scores have surpassed other approaches in most of the categories.

7. Discussion and Limitations

To our knowledge, our approach is the first to use SDFs as the geometric representation in a differentiable renderer.

We offer a higher level of freedom for topological changes compared to triangle meshes, which are typically restricted to the topology of a template mesh. In contrast to point clouds, SDFs inherently represent continuous watertight surfaces. We demonstrated applications of our approach in multi-view shape reconstruction, and single view 3D reconstruction using deep learning. Our experimental results showed that we can more robustly perform multi-view reconstruction than a state-of-the-art point-based differentiable renderer. In addition, we achieve state-of-the-art results on single view 3D reconstruction with deep learning models.

In our multi-view 3D reconstruction approach, our current simple shading model is not sufficient to perform inverse rendering from real images taken with a camera. For example, the current model does not include effects such as shadows, interreflections, texture, non-diffuse surfaces, or complex illumination. In contrast to rasterization-based differentiable renderers, our ray tracing-based renderer could be extended to include all such effects. A disadvantage of our deep learning approach is that we output a discrete SDF on a 3D grid. Instead, we could learn a continuous signed distance function represented by a deep network like in DeepSDF [39]. This would be more memory efficient, but it might be computationally too expensive for unsupervised 3D reconstruction with differentiable rendering, since it would require to evaluate the network for each ray marching step.

8. Conclusions

We proposed a novel approach to differentiable rendering using signed distance functions to represent watertight 3D geometry. Our rendering algorithm is based on sphere tracing, but we observe that only the local shading

Method	Input Image	Rendered Views				Input Image	Rendered Views			
GT										
Ours										
SoftRas[28]										
GT										
Ours										
SoftRas[28]										
GT										
Ours										
SoftRas[28]										
GT										
Ours										
SoftRas[28]										
GT										
Ours										
SoftRas[28]										

Figure 8. Single-view reconstruction results for airplanes, chairs, and benches.

computation needs to be differentiable in our framework, which makes the approach computationally more efficient

and allows for straightforward integration into deep learning frameworks. We demonstrate applications in multi-view

3D reconstruction and unsupervised single view 3D reconstruction using deep neural networks. Our experimental results illustrate the advantages over geometry representations such as point clouds and meshes. In particular, we report the state-of-the-art results in shape reconstruction.

References

- [1] Panos Achlioptas, Olga Diamanti, Ioannis Mitliagkas, and Leonidas J. Guibas. Representation learning and adversarial generation of 3d point clouds. *CoRR*, abs/1707.02392, 2017.
- [2] Anish Athalye, Logan Engstrom, Andrew Ilyas, and Kevin Kwok. Synthesizing robust adversarial examples. *CoRR*, abs/1707.07397, 2017.
- [3] Angel X. Chang, Thomas A. Funkhouser, Leonidas J. Guibas, Pat Hanrahan, Qi-Xing Huang, Zimo Li, Silvio Savarese, Manolis Savva, Shuran Song, Hao Su, Jianxiang Xiao, Li Yi, and Fisher Yu. Shapenet: An information-rich 3d model repository. *ArXiv*, abs/1512.03012, 2015.
- [4] Chengqian Che, Fujun Luan, Shuang Zhao, Kavita Bala, and Ioannis Gkioulekas. Inverse transport networks. *CoRR*, abs/1809.10820, 2018.
- [5] Wenzheng Chen, Jun Gao, Huan Ling, Edward J. Smith, Jaakko Lehtinen, Alec Jacobson, and Sanja Fidler. Learning to predict 3d objects with an interpolation-based differentiable renderer, 2019.
- [6] Zhiqin Chen and Hao Zhang. Learning implicit fields for generative shape modeling. *IEEE Conference on Computer Vision and Pattern Recognition*, 2019.
- [7] 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 *Proceedings of the European Conference on Computer Vision (ECCV)*, 2016.
- [8] 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 - SIGGRAPH '96*, 1996.
- [9] Amaël Delaunoy and Emmanuel Prados. Gradient flows for optimizing triangular mesh-based surfaces: Applications to 3d reconstruction problems dealing with visibility. *International Journal of Computer Vision*, 95(2):100–123, Nov 2011.
- [10] Haoqiang Fan, Hao Su, and Leonidas J. Guibas. A point set generation network for 3d object reconstruction from a single image. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, July 2017.
- [11] Matheus Gadelha, Subhransu Maji, and Rui Wang. 3d shape induction from 2d views of multiple objects. *2017 International Conference on 3D Vision (3DV)*, pages 402–411, 2016.
- [12] J. Gwak, C. B. Choy, M. Chandraker, A. Garg, and S. Savarese. Weakly supervised 3d reconstruction with adversarial constraint. In *2017 International Conference on 3D Vision (3DV)*, pages 263–272, Oct 2017.
- [13] Zhizhong Han, Mingyang Shang, Yu-Shen Liu, and Matthias Zwicker. View Inter-Prediction GAN: Unsupervised representation learning for 3D shapes by learning global shape memories to support local view predictions. In *AAAI*, pages 8376–8384, 2019.
- [14] Zhizhong Han, Mingyang Shang, Xiyang Wang, Yu-Shen Liu, and Matthias Zwicker. Y2Seq2Seq: Cross-modal representation learning for 3D shape and text by joint reconstruction and prediction of view and word sequences. In *AAAI*, pages 126–133, 2019.
- [15] John C. Hart. Sphere tracing: a geometric method for the antialiased ray tracing of implicit surfaces. *The Visual Computer*, 12(10):527–545, Dec 1996.
- [16] Paul Henderson and Vittorio Ferrari. Learning to generate and reconstruct 3d meshes with only 2d supervision. *CoRR*, abs/1807.09259, 2018.
- [17] Philipp Henzler, Niloy Mitra, and Tobias Ritschel. Escaping plato’s cave using adversarial training: 3d shape from unstructured 2d image collections. *arXiv preprint arXiv:1811.11606*, 2018.
- [18] Eldar Insafutdinov and Alexey Dosovitskiy. Unsupervised learning of shape and pose with differentiable point clouds. In *Advances in Neural Information Processing Systems*, pages 2807–2817, 2018.
- [19] Li Jiang, Shaoshuai Shi, Xiaojuan Qi, and Jiaya Jia. Gal: Geometric adversarial loss for single-view 3d-object reconstruction. In *The European Conference on Computer Vision (ECCV)*, September 2018.
- [20] Danilo Jimenez Rezende, S. M. Ali Eslami, Shakir Mohamed, Peter Battaglia, Max Jaderberg, and Nicolas Heess. Unsupervised learning of 3d structure from images. In D. D. Lee, M. Sugiyama, U. V. Luxburg, I. Guyon, and R. Garnett, editors, *Advances in Neural Information Processing Systems 29*, pages 4996–5004. Curran Associates, Inc., 2016.
- [21] Hiroharu Kato, Yoshitaka Ushiku, and Tatsuya Harada. Neural 3d mesh renderer. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018.
- [22] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization, 2014.
- [23] A. Kurenkov, J. Ji, A. Garg, V. Mehta, J. Gwak, C. Choy, and S. Savarese. Deformnet: Free-form deformation network for 3d shape reconstruction from a single image. In *2018 IEEE Winter Conference on Applications of Computer Vision (WACV)*, pages 858–866, March 2018.
- [24] F. Langguth, K. Sunkavalli, S. Hadap, and M. Goesele. Shading-aware multi-view stereo. In *Proceedings of the European Conference on Computer Vision (ECCV)*, 2016.
- [25] Tzu-Mao Li, Miika Aittala, Frédo Durand, and Jaakko Lehtinen. Differentiable monte carlo ray tracing through edge sampling. *ACM Trans. Graph. (Proc. SIGGRAPH Asia)*, 37(6):222:1–222:11, 2018.
- [26] Hsueh-Ti Derek Liu, Michael Tao, Chun-Liang Li, Derek Nowrouzezahrai, and Alec Jacobson. Adversarial geometry and lighting using a differentiable renderer. *CoRR*, abs/1808.02651, 2018.
- [27] Hsueh-Ti Derek Liu, Michael Tao, and Alec Jacobson. Paparazzi: Surface editing by way of multi-view image processing. *ACM Transactions on Graphics*, 2018.
- [28] Shichen Liu, Tianye Li, Weikai Chen, and Hao Li. Soft rasterizer: A differentiable renderer for image-based 3d reason-

- ing. *The IEEE International Conference on Computer Vision*, 2019.
- [29] Matthew M. Loper and Michael J. Black. OpenDR: An approximate differentiable renderer. In *Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)*, 2014.
- [30] Abhimitra Meka, Maxim Maximov, Michael Zollhoefer, Avishek Chatterjee, Hans-Peter Seidel, Christian Richardt, and Christian Theobalt. Lime: Live intrinsic material estimation. In *Proceedings of Computer Vision and Pattern Recognition (CVPR)*, June 2018.
- [31] Lars Mescheder, Michael Oechsle, Michael Niemeyer, Sebastian Nowozin, and Andreas Geiger. Occupancy networks: Learning 3D reconstruction in function space. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2019.
- [32] Mateusz Michalkiewicz, Jhony K. Pontes, Dominic Jack, Mahsa Baktashmotlagh, and Anders P. Eriksson. Deep level sets: Implicit surface representations for 3D shape inference. *CoRR*, abs/1901.06802, 2019.
- [33] KL Navaneet, Priyanka Mandikal, Mayank Agarwal, and R Venkatesh Babu. Capnet: Continuous approximation projection for 3d point cloud reconstruction using 2d supervision. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pages 8819–8826, 2019.
- [34] Richard A. Newcombe, Shahram Izadi, Otmar Hilliges, David Molyneaux, David Kim, Andrew J. Davison, Pushmeet Kohli, Jamie Shotton, Steve Hodges, and Andrew Fitzgibbon. KinectFusion: Real-time dense surface mapping and tracking. In *2011 10th IEEE International Symposium on Mixed and Augmented Reality, ISMAR 2011*, 2011.
- [35] Thu Nguyen-Phuoc, Chuan Li, Stephen Balaban, and Yong-Liang Yang. RenderNet: A deep convolutional network for differentiable rendering from 3d shapes. In *Advances in Neural Information Processing Systems 31*, 2018.
- [36] Matthias Nießner, Michael Zollhöfer, Shahram Izadi, and Marc Stamminger. Real-time 3D reconstruction at scale using voxel hashing. *ACM Transactions on Graphics*, 2013.
- [37] S Osher, R Fedkiw, and K Piechor. Level Set Methods and Dynamic Implicit Surfaces. *Applied Mechanics Reviews*, 2004.
- [38] Andrea Palazzi, Luca Bergamini, Simone Calderara, and Rita Cucchiara. End-to-end 6-dof object pose estimation through differentiable rasterization. In *The European Conference on Computer Vision (ECCV) Workshops*, September 2018.
- [39] Jeong Joon Park, Peter Florence, Julian Straub, Richard Newcombe, and Steven Lovegrove. DeepSDF: Learning Continuous Signed Distance Functions for Shape Representation. *arXiv e-prints*, page arXiv:1901.05103, Jan 2019.
- [40] Felix Petersen, Amit H. Bermano, Oliver Deussen, and Daniel Cohen-Or. Pix2vex: Image-to-geometry reconstruction using a smooth differentiable renderer. *CoRR*, abs/1903.11149, 2019.
- [41] Ravi Ramamoorthi and Pat Hanrahan. A signal-processing framework for inverse rendering. In *Proceedings of the 28th Annual Conference on Computer Graphics and Interactive Techniques*, SIGGRAPH '01, pages 117–128, New York, NY, USA, 2001. ACM.
- [42] Elad Richardson, Matan Sela, Roy Or-El, and Ron Kimmel. Learning detailed face reconstruction from a single image. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 5553–5562, 07 2017.
- [43] S. M. Seitz, B. Curless, J. Diebel, D. Scharstein, and R. Szeliski. A comparison and evaluation of multi-view stereo reconstruction algorithms. In *2006 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR'06)*, volume 1, pages 519–528, June 2006.
- [44] Ben Semerjian. A new variational framework for multiview surface reconstruction. In *Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)*, 2014.
- [45] Vincent Sitzmann, Michael Zollhöfer, and Gordon Wetzstein. Scene representation networks: Continuous 3d-structure-aware neural scene representations. *CoRR*, abs/1906.01618, 2019.
- [46] S. Tulsiani, A. Kar, J. Carreira, and J. Malik. Learning category-specific deformable 3d models for object reconstruction. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 39(4):719–731, April 2017.
- [47] Shubham Tulsiani, Tinghui Zhou, Alexei A. Efros, and Jitendra Malik. Multi-view supervision for single-view reconstruction via differentiable ray consistency. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 209–217, 2017.
- [48] Qingxiang Wang, Cezary Kaliszyk, and Josef Urban. First experiments with neural translation of informal to formal mathematics. *CoRR*, abs/1805.06502, 2018.
- [49] Jiajun Wu, Yifan Wang, Tianfan Xue, Xingyuan Sun, William T Freeman, and Joshua B Tenenbaum. MarrNet: 3D Shape Reconstruction via 2.5D Sketches. In *Advances In Neural Information Processing Systems*, 2017.
- [50] Xinchen Yan, Jimei Yang, Ersin Yumer, Yijie Guo, and Honglak Lee. Perspective transformer nets: Learning single-view 3d object reconstruction without 3d supervision. In D. D. Lee, M. Sugiyama, U. V. Luxburg, I. Guyon, and R. Garnett, editors, *Advances in Neural Information Processing Systems 29*, pages 1696–1704. Curran Associates, Inc., 2016.
- [51] Wang Yifan, Felice Serena, Shihao Wu, Cengiz Öztireli, and Olga Sorkine-Hornung. Differentiable surface splatting for point-based geometry processing. *ACM Transactions on Graphics (proceedings of ACM SIGGRAPH ASIA)*, 38(6), 2019.
- [52] Jun-Yan Zhu, Zhoutong Zhang, Chengkai Zhang, Jiajun Wu, Antonio Torralba, Joshua B. Tenenbaum, and William T. Freeman. Visual object networks: Image generation with disentangled 3D representations. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2018.
- [53] Matthias Zwicker, Hanspeter Pfister, Jeroen van Baar, and Markus Gross. Surface splatting. In *Proceedings of the 28th Annual Conference on Computer Graphics and Interactive Techniques*, SIGGRAPH '01, pages 371–378, New York, NY, USA, 2001. ACM.