

NeRFMeshing: Distilling Neural Radiance Fields into Geometrically-Accurate 3D Meshes

Marie-Julie Rakotosaona¹ Fabian Manhardt¹ Diego Martin Arroyo¹
Michael Niemeyer¹ Abhijit Kundu¹ Federico Tombari^{1,2}

¹ Google ² TU Munich

Abstract

With the introduction of Neural Radiance Fields (NeRFs), novel view synthesis has recently made a big leap forward. At the core, NeRF proposes that each 3D point can emit radiance, allowing to conduct view synthesis using differentiable volumetric rendering. While neural radiance fields can accurately represent 3D scenes for computing the image rendering, 3D meshes are still the main scene representation supported by most computer graphics and simulation pipelines, enabling tasks such as real time rendering and physics-based simulations. Obtaining 3D meshes from neural radiance fields still remains an open challenge since NeRFs are optimized for view synthesis, not enforcing an accurate underlying geometry on the radiance field. We thus propose a novel compact and flexible architecture that enables easy 3D surface reconstruction from any NeRF-driven approach. Upon having trained the radiance field, we distill the volumetric 3D representation into a Signed Surface Approximation Network, allowing easy extraction of the 3D mesh and appearance. Our final 3D mesh is physically accurate and can be rendered in real time on an array of devices.

1. Introduction

Accurate 3D scene and object reconstruction is a key problem in areas such as robotics, photogrammetry, AR/VR, where applications often rely on precise 3D geometry to perform physics-based simulations, real-time 3D visualizations, rendering and interactions. Moreover, the related field of novel view synthesis (NVS) has made tremendous advances in recent years. Recently Mildenhall *et al.* [13] proposed to perform NVS by means of neural radiance fields (NeRFs), a novel 3D representation where each 3D location in space can emit radiance, see Sec 3.1 for more details. Novel views are synthesized by means of differentiable volumetric rendering [13, 23]. Due to the impressive results and simplicity of the approach, most related work has focused on improving NeRF in terms of image qual-



Figure 1. Our method extracts meshes with accurate geometry and view dependent appearance given a collection of posed images. We show a composition of meshes extracted from the *Chair*, *Hot-dog*, *Lego* (Blender dataset) and *Garden* (MipNeRF 360 dataset) scenes using our method (Top left: scene rendered with colors, top right: geometry visualization). Our method enables physics based simulations. We show the results of a simulation of a cloth falling on the objects from the Blender Synthetic dataset (Bottom).

ity [1], robustness [9, 11, 16], as well as training speed [14] and rendering speed [3, 5]. Unfortunately, as these representations are commonly optimized for the NVS task and not explicitly for the underlying geometry [3, 6], it is yet unclear how to best obtain accurate 3D meshes from radiance fields. Indeed, while the volumetric representation of NeRF enables accurate renderings from new views, the underlying 3D geometry of each object is not uniquely defined as a level-set surface. NeRF methods often rely on layering and transparency effects to approximate complex appearance and geometry. The surface of objects is therefore approximated by dense regions of the volume instead of surfaces of zero thickness. Moreover, most related work still lacks the capability to be rendered in real-time [16, 21],

especially on commodity hardware. Finally, NeRFs cannot be directly integrated with most computer graphics (CG) pipelines, as they still rely on standard 3D meshes due to their compactness and physical properties.

Despite some recent work proposing alternative scene representations to re-enable real-time rendering even for NeRFs, they are again not designed to produce accurate 3D representations of the input objects or scenes to be used with standard CG pipelines [3, 6].

To deal with these limitations, we thus introduce *NeRFMeshing*, an end-to-end pipeline for efficiently extracting geometrically accurate meshes from trained NeRF-based networks, merely adding a very small overhead in time. Our method produces meshes with neural colors having accurate geometry that can be rendered in real time on commodity hardware. Introducing a novel signed surface approximation network (SSAN), we train a post-processing NeRF pipeline, defining the underlying surface and appearance. SSAN produces an accurate 3D triangle mesh of the scene that we render using a small appearance network to produce view-dependent colors. In contrast to other works that leverage distance fields, requiring significant modifications in the used NeRF architecture [22, 24], our method can be leveraged together with any NeRF, enabling to easily incorporate new advances, such as improved handling of unbounded scenes [1] or reflective objects [21]. Essentially, SSAN estimates a Truncated Signed Distance Field (TSDF) and a feature appearance field. Harnessing the NeRF approximated geometry as well as the used training views, we distill the trained NeRF into the SSAN model. We then extract the 3D mesh from the trained SSAN which can be rendered on embedded devices at high frame-rate using rasterization and the appearance network. Thanks to the flexibility of our method, we can generate these 3D meshes fast [14], are not tied to object-centric scenes [1], and can even model complex and non-lambertian surfaces [21].

To summarize, we propose *NeRFMeshing*, a novel method for capturing both accurate 3D meshes of the scene as well as enabling realistic view dependent rendering. The extracted meshes from our end-to-end pipeline can be integrated in graphics and simulation pipelines. Our model also preserves the high fidelity of neural radiance fields like view-dependent effects and reflections and can be used for real-time novel view synthesis.

2. Related Work

The base of our work is the neural radiance field (NeRF) formulation initially introduced in [13]. This work describes a trainable radiance field parameterized with a neural network. Subsequent works have addressed the main limitations of the original approach, such as the slow training speed [14], anti-aliasing effects or the ability to model unbounded scenes [1]. Another main caveat of the origi-

nal formulation is the lack of accurate underlying geometry, mostly caused by the fact that NeRFs are optimized exclusively for visual consistency. This gives way for the network to create occupancy regions to support the volumetric rendering process even when particular parts of the space are not occupied by an underlying surface. Additional supervision signals, such as depth [4, 19] are effective in regularizing the geometry, but require additional input modalities which might be difficult to source. In contrast, our work relies on NeRF networks trained from images.

An alternative to radiance fields is to learn a Signed Distance Function (SDF) [17, 22, 24]. A high-quality mesh can be extracted by means of marching cubes or similar approaches. [15] combines this approach with a tetrahedral grid that is optimized during training, followed by differentiable rasterizer to recover materials and lighting. However, this method is constrained by the choice of the grid resolution at training time. Compared to [15], our method does not rely on a fixed grid template during training, enabling us to reconstruct at any arbitrary resolution afterwards. Similarly, [25] uses an SDF with sphere tracing to determine the intersection with the surface. In this work, we are mainly interested in obtaining a unique and geometrically accurate 3D mesh surface from NeRF methods. We exploit the adaptive power of NeRFs to be able to robustly represent 3D scenes in a range of conditions and environments. In particular, we do not require any changes to the NeRF architecture in order to compute a surface. Also differentiable mesh rasterizers often suffer from optimization issues which are easily avoidable in differentiable volumetric rendering adopted by the NeRF methods.

Recent approaches advance speed and geometric accuracy of NeRFs. SNeRG [6] has achieved real-time rendering of radiance fields by restructuring the original NeRF architecture to precompute the predicted density, view-independent colors and feature vectors in a sparse voxel grid. A separate view-dependent network runs online to compute the final color. It results in a real-time visualization with a good trade-off in image quality, however, the underlying geometry is not well defined. MobileNeRF [3] builds on a similar principle, but relies on a triangle mesh that is optimized during training as the underlying geometry. However, the final geometry is far from being accurate. As the method relies on a single resolution grid at training time, the faces from the triangle soup cannot fit detailed regions. Moreover, multiple layers of faces can sometimes describe the same surface during optimization leading to inaccurate geometry to support the rendering loss. Voxel grids are also a viable alternative to efficiently store geometry [20, 26] and have efficient renderings. However, they are inherently constrained by the grid resolution and can suffer from discretization artifacts.

3. Method

In this section we present our approach for extracting accurate 3D meshes with neural features from NeRF for subsequent real-time rendering. We present an overview of the method in Fig. 2. We first briefly outline the general concept of NeRFs in Sec. 3.1. In Sec. 3.2 we present our method for approximating surface from NeRF. Finally in Sec. 3.3 we describe mesh extraction and real-time rendering.

3.1. Neural Radiance Fields

At the core, a neural radiance field is a continuous mapping from a 3D location $\mathbf{x} \in \mathbb{R}^3$ and a ray viewing direction $\mathbf{d} \in \mathbb{S}^2$ to an RGB color $\mathbf{c} \in [0, 1]^3$ and volume density $\sigma \in \mathbb{R}^+$. It can be formulated as,

$$[\mathbf{c}, \sigma] = F_\theta(\gamma_x(\mathbf{x}), \gamma_d(\mathbf{d})) \quad (1)$$

where F is modeled as an MLP with learnable parameters θ , and $\gamma : \mathbb{R}^3 \rightarrow \mathbb{R}^N$ is a positional encoding of the input \mathbf{x} required to capture high frequencies.

Given a camera pose $P_i = [R_i, t_i] \in SE(3)$, each pixel coordinate $\mathbf{p} \in \mathbb{R}^2$ determines a ray in the world coordinate system, whose origin is the camera center of projection $\mathbf{o}_i = \mathbf{t}_i$ and direction is defined as $\mathbf{d}_{i,p} = R_i K_i^{-1} \mathbf{p}$. We can express a 3D point along the viewing ray associated with \mathbf{p} at depth t as $\mathbf{r}_{i,p}(t) = \mathbf{o}_i + t\mathbf{d}_{i,p}$. To render the color $\hat{\mathbf{I}}_{i,p} \in [0, 1]^3$ at pixel \mathbf{p} , we sample M discrete depth values t_m along the ray within the near and far plane $[t_n, t_f]$, and query F_θ at the associated 3D points. The corresponding predicted color and volume density values $\{(\mathbf{c}_m, \sigma_m)\}_{m=1}^M$ are then composited as,

$$\hat{\mathbf{I}}_{i,p} = \hat{I}(\mathbf{p}; \theta, P_i) = \sum_{m=1}^M \alpha_m \mathbf{c}_m, \quad (2)$$

$$\text{where } \alpha_m = T_m (1 - \exp(-\sigma_m \delta_m)), \quad (3)$$

$$T_m = \exp\left(-\sum_{m'=1}^m \sigma_{m'} \delta_{m'}\right). \quad (4)$$

where T_m denotes the accumulated transmittance along the ray from t_n to t_m , and $\delta_m = t_{m+1} - t_m$ is the distance between adjacent samples. Finally we compute accumulated depth at percentile k as follows,

$$\hat{z}_{i,p,k} = \sum_{m=1}^{M_k} \alpha_m t_m. \quad (5)$$

Where M_k denotes the first index at which the accumulated transmittance $T_k > k/100$.

3.2. Surface approximation from NeRF

In this section we introduce our signed surface approximation network (SSAN) module that creates a truncated signed distance field (TSDF) from NeRFs.

Signed Surface Approximation Network (SSAN).

SSAN is optimized from pretrained NeRFs to enable the extraction of a unique and accurate surface, and appearance. Relying on the pre-trained NeRF representations brings two main advantages, which we harness when training the SSAN: a rough 3D approximation of the scene geometry learnt by the NeRF, and useful priors such as rendered depths via volumetric rendering. To this end, we feed a 3D coordinate $\mathbf{x} \in \mathbb{R}^3$ to SSAN $\phi : \mathbb{R}^3 \rightarrow (\mathbb{R}, \mathbb{S}^2, \mathbb{R}^8)$ to predict the truncated signed distance approximation $\hat{t} \in [-0.1, 0.1]$, the normal $\hat{\mathbf{n}} \in \mathbb{S}^2$, used for normal smoothness regularization as well as rendering, and 8-dimensional appearance feature $\hat{\mathbf{f}} \in [0, 1]^8$ with

$$[\hat{t}, \hat{\mathbf{n}}, \hat{\mathbf{f}}] = \phi(\mathbf{x}). \quad (6)$$

The outgoing features $\hat{\mathbf{f}}$ and normals $\hat{\mathbf{n}}$ are then further processed together with the viewing direction $\mathbf{d} \in \mathbb{S}^2$ by a small appearance network $\eta : (\mathbb{R}^8, \mathbb{S}^2, \mathbb{S}^2) \rightarrow \mathbb{R}^3$ to produce the view-dependant RGB color $\hat{\mathbf{c}} \in [0, 1]^3$ according to

$$\hat{\mathbf{c}} = \eta(\hat{\mathbf{f}}, \hat{\mathbf{n}}, \mathbf{d}). \quad (7)$$

We refer to Fig. 2 for an abstract overview over of SSAN.

TSDF. We aim to learn a TSDF approximation t that represents the underlying surface of a scene or object. t should have an accurate sign globally (positive outside the surface and negative inside). Moreover, t should be smooth and have near constant derivatives close to the surface to enable efficient marching cubes at a later stage. Note that we are not constrained to produce a well defined distance function. In this section, we introduce multiple loss functions to enforce these properties.

At training time, we exploit the NeRF occupancy aggregated along rays. While the weighted average of the occupancy is often used to render an approximation of depth, we instead use the distribution of the occupancy along the ray for more accuracy as in [18]. We call $\hat{z}_{i,p,50}$ the median value of depths along a ray as described in Sec. 3.1. We can estimate 3D points that are outside or inside the surface by constructing points before or after the median depth from the camera origin respectively. In practice, we choose to render the depth of the NeRF outside at the 16th percentile $\hat{z}_{i,p,16}$ and inside the object surface at the 84th percentile $\hat{z}_{i,p,84}$. Subsequently, we compute the corresponding 3D world coordinates and feed them to SSAN to estimate the respective SDFs with:

$$(\hat{\mathbf{t}}_{50,i,p}, \hat{\mathbf{n}}_{50,i}, \hat{\mathbf{f}}_{50,i}) = \phi(r_{i,p}(\hat{z}_{i,p,50})), \quad (8)$$

$$(\hat{\mathbf{t}}_{16,i,p}, \hat{\mathbf{n}}_{16,i}, \hat{\mathbf{f}}_{16,i}) = \phi(r_{i,p}(\hat{z}_{i,p,16})), \quad \text{and} \quad (9)$$

$$(\hat{\mathbf{t}}_{84,i,p}, \hat{\mathbf{n}}_{84,i}, \hat{\mathbf{f}}_{84,i}) = \phi(r_{i,p}(\hat{z}_{i,p,84})). \quad (10)$$

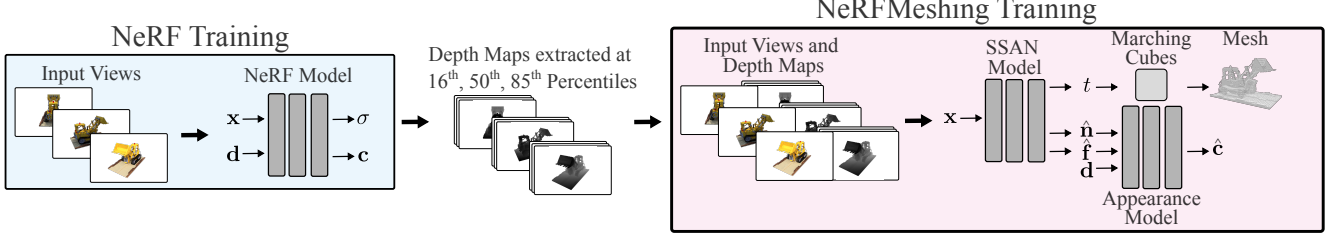


Figure 2. We exploit rendered depth distribution from NeRF to help supervise an approximated TSDF. We produce learned features that we feed to a small appearance network to predict RGB colors. We can extract the surface using marching cubes and store appearance features on the surface. Finally we render in real time using the mesh, appearance features and appearance network.

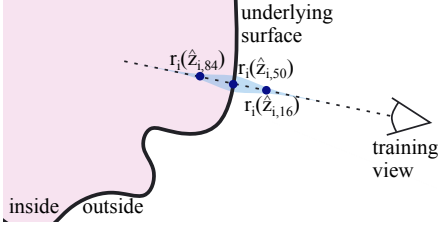


Figure 3. We exploit occupancy aggregated over training rays to estimate the outside ($r_i(\hat{z}_{i,16})$), the surface ($r_i(\hat{z}_{i,50})$) and the inside ($r_i(\hat{z}_{i,84})$), of the object.

For readability, we drop the term p in later references. During training, we enforce this construction by applying the loss L_i at sampled projections along training rays.

$$L_i = \sum_i \|\hat{t}_{16,i} - \epsilon\|^2 + \|\hat{t}_{50,i}\|^2 + \|\hat{t}_{84,i} + \epsilon\|^2, \quad (11)$$

We visualize these points in Fig. 3. We also provide an ablation in the supplementary material on the choice of the percentile values. We note that we choose to use values computed from the distribution and not a fixed constant to better account for the scale as well as uncertainties or different levels of details in one scene.

More importantly, we also enforce that the learned indicator function is locally smooth and has near constant derivatives close to the zero level set. Namely, the approximated TSDF function should have normals of constant norm n_c around the median depth. We recall that we are only approximating the behaviour of a SDF in order to run Marching Cubes. We enforce approximated signed distance values at positions computed from the ray distribution. However, the distribution can vary between scenes, ray direction or even regions of the same scenes. Therefore, we do not constrain the function to simulate exact euclidean distances.

$$L_n = \sum_i \|\|N(r_i(\delta))\| - n_c\|^2, \quad (12)$$

where $N : \mathbb{R}^3 \rightarrow \mathbb{S}^2$ is a function computing the normal at a given point x . δ is uniformly sampled value in

$[\hat{z}_{i,16}, \hat{z}_{i,84}]$. To increase the speed of our method, we compute the normals using finite difference of the SDF network at nearby 3D positions. While other methods such as [21] often compute normals using gradients from the network, their method is significantly more time consuming.

Normal regularization. To ensure a smooth surface, we use the normal regularization used in RefNeRF [21] on the TSDF normals. We use the normal smoothness loss:

$$L_s = \sum_i \|N(r_i(\hat{z}_{i,50})) - \hat{n}_{50,i}\|^2 \quad (13)$$

Enforcing that the computed normals are close to the estimated normals produces smoother computed normals as the signals produced by the network are limited in frequency. We also use the normal orientation loss:

$$L_o = \sum_i \max(0, N(r_i(\hat{z}_{i,50})) \cdot \mathbf{d}_i). \quad (14)$$

Appearance. We train a small appearance network that takes as input shading features \hat{f}_i and predicted normals \hat{n}_i predicted by SSAN as well as input view directions \mathbf{d}_i to produce the surface color at a given coordinate from a given view similar to SNeRG [6]. We supervise this network with the ground truth colors from the training images. We notice that, the projected median depth can be an inaccurate estimation of the surface zero-level set, especially when the angle between the view direction and the surface is small. Therefore, we supervise the color at points projected on the zero-level set. We compute the projection using 4 steps of gradient descent along the viewing ray using the SDF approximated values, see the supplementary materials for more details. We refer to the depth of points $r_i(\hat{z}_{i,50})$ projected on the zero level set as l .

$$L_{color} = \sum_i \|\eta(\hat{f}_{l,i}, \hat{n}_{l,i}, \mathbf{d}_i) - \mathbf{c}_i\|^2, \quad (15)$$

where \mathbf{c}_i is the ground truth pixel color.

The different losses are combined as a weighted average during training into a final loss L .

$$L = \alpha L_i + \beta L_n + \gamma L_s + \delta L_o \quad (16)$$

Training stage. We explore training the SSAN module in two possible ways. The training process mostly relies on having available depth maps from NeRF. Depths maps can be computed directly during training of the NeRF based method or they can also be rendered from a pretrained NeRF based model. We choose to train our method from the rendered depth percentiles of a pretrained NeRF model. We provide more information about simultaneously training in the supplementary material.

3.3. Mesh extraction and Real-time rendering

Mesh extraction. The SSAN module converts the radiance field representation of the scene obtained from NeRF to distance field (TSDF) representation. Given the TSDF representation, we can easily reconstruct the surface by using a surface reconstruction algorithm [8, 10]. In this paper we use the PyMCubes¹ implementation of marching cubes. Moreover we use either vertex features, for instance on objects from the Blender Synthetic dataset or build a texture using per face parametrization for unbounded scenes from [1]. We fill the texture image with appearance features sampled from the SSAN appearance features at the interpolated face locations.

Rendering. The triangle mesh geometry extracted from the SSAN can be encoded in any common format like OBJ, glTF, and others. Thus output meshes from our pipeline can be directly integrated in traditional graphics pipeline. The neural view-dependent appearance is added by rasterizing the precomputed texture and feeding the results to the appearance network that produces the final RGB color values. This strategy ensures that the network only performs $w \cdot h$ evaluations and thus can run at a high frame rate on commodity hardware. We measure average FPS of 25 on a workstation and 30 on a MacBook on the Blender Synthetic dataset objects. This approach is similar to SNeRG [6], although our method only samples features at the surface boundary instead of aggregating them in a volumetric fashion, similar to MobileNeRF [3]. Contrary to these two methods however, our underlying geometry is by design significantly closer to the ground truth.

3.4. Implementation details

We use a backbone similar to Instant NGP [14] for the SSAN module. We divide the module in two separate branches with separate weights: i) the geometry branch that outputs TSDF approximation and the normal prediction. ii) The appearance branch which outputs color features. Each

¹<https://github.com/pmneila/PyMCubes>

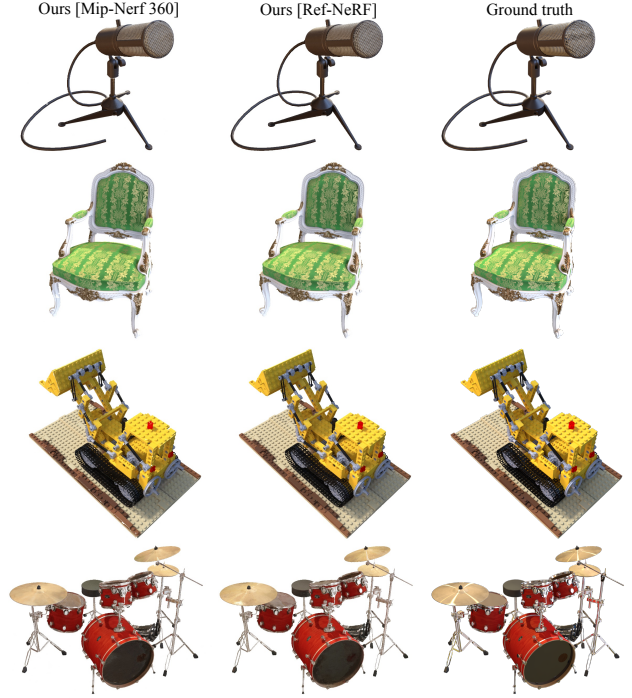


Figure 4. Rendering of the 3D meshes from test views on the Synthetic Blender dataset.

branch is a separate network with the same architecture. We use Instant NGP [14] architectures with hash table size 2^{19} , coarsest resolution of 16, highest resolution of 2048, 15 levels and a number of feature dimension per entry of 2. For the appearance network we use a network of 4 layers of MLP each with width of 32. We use the JAX framework for our implementation. We can train and extract a mesh end-to-end in less than an hour using 8 V100 NVIDIA GPUs. Finally we experimentally set the hyper-parameters $n_c = 10$ and $\epsilon = 0.1$.

4. Evaluation

We validate the effectiveness of our approach on synthetic blender scenes from [13] and real unbounded scenes from [1]. Our method can be easily integrated to NeRF pipelines and can help with regularizing the underlying 3D geometry to improve its accuracy. To demonstrate this we use the NeRF backbones of Mip-NeRF 360 [1] that we combine with Neural Graphics Primitives [14] to speed up training, and Ref-NeRF [21]. In addition to the results presented in this section, please see the supplementary materials for additional video results and demos.

4.1. Synthetic blender scenes

In this section we focus on the Synthetic Blender dataset from [13]. We compare our method to state of the art baselines that focus on real time rendering of NeRF: SNeRG [6],

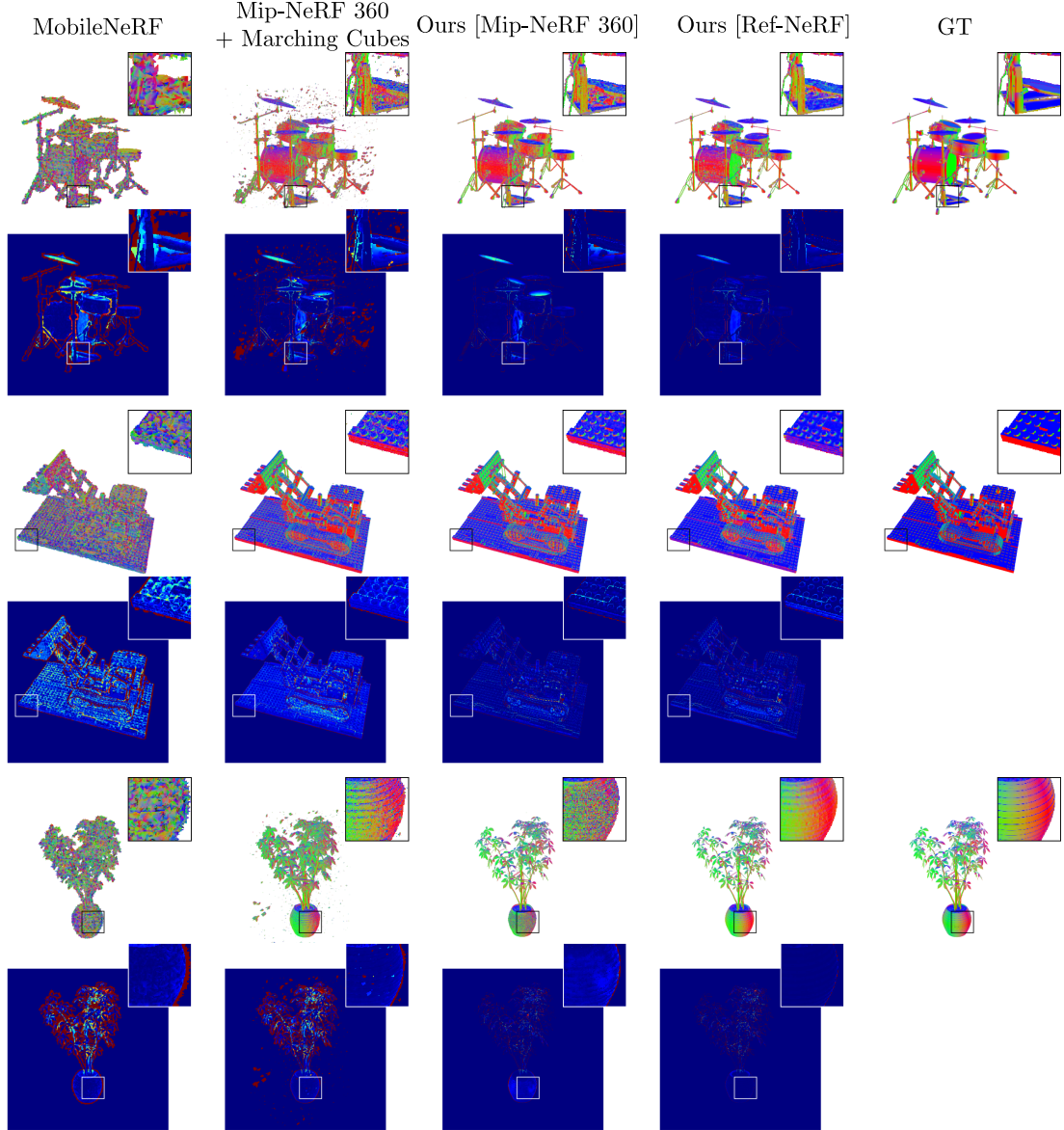


Figure 5. Geometry comparison on the Synthetic Blender scenes [13]. We visualize mesh normals and the depth absolute difference with the ground truth. We notice that meshes obtained using our method have smoother, more accurate and realistic geometry.

and MobileNeRF [3]. Additionally, we compare to Mip-NeRF 360 [1] from which we reconstruct the geometry with Marching Cubes (MC) at the same grid resolution as ours and using 0.5 as the zero-level set in the density grid. Similarly to our method we use Mip-NeRF 360 accelerated with Instant-NGP [14]. For our method, we use vertex based feature representation and a grid of size 1024 inside the bounding box obtained with the ground truth mesh for marching cubes. For the *drums* and *ficus*, we use a higher resolution of 2048 since these two scenes are more detailed. We compare the geometry and appearance produced by each method. We train our method together with the Mip-NeRF 360 backbone

and Ref-NeRF backbone.

Evaluation metrics. We measure Chamfer Distance (CD) that we compute using an observability mask similar to [7], as we want to evaluate only in observable regions to for fair comparisons. We construct an observability grid of resolution 256 constructed with rays from training views inside the mesh bounding box. Additionally, we compute the normal consistency (NC) as in [12] as the absolute value of the dot product between normals from a point cloud sampled on the ground truth mesh and the predicted mesh normals. Finally, we also evaluate PSNR to measure appear-

ance quality.

Comparison. We present quantitative results on the Blender synthetic dataset in Tab. 1 and qualitative results in Fig. 5. We observe that our method tends to produce better geometric results in terms of Chamfer Distance and normal consistency. In Fig. 5 we show the rendering of normal maps from different meshes (row 1, 3, 5) and the absolute difference with the ground truth depth rendered from the meshes (row 2, 4, 6). We notice that MobileNeRF produces a “triangle soup” that does not accurately reflect the geometry of the object. Computed meshes with marching cubes from pretrained NeRFs lead to more floaters and a surface that is not well aligned with the ground truth as we can notice a red region at the boundary of the depth maps as well as lighter color on the surface. Moreover, while we are expecting lower appearance quality than other methods, due to the constraints of producing accurate geometry and having real time rendering, we note that rendering quality remains very high and visibly close to the ground truth, as can be seen in Tab. 1 and Fig. 4.

Effect of the backbone We observe that using the appropriate NeRF based backbone can improve the quality of the final mesh significantly in Fig. 6. In particular, shiny objects such as the Materials scene highly benefit from the Ref-NeRF backbone as the method is specifically designed to deal with such objects (bottom of Fig. 6). On the other hand, the Ref-NeRF backbone can sometimes fail to approximate the right geometry due to the powerful appearance network while Mip-NeRF 360 produces more realistic geometry, for instance on the ship scene (top of Fig. 6). This effect is also visible in Tab. 1 as the Mip-NeRF 360 backbone produces more accurate geometry for the Ship scene and Ref-NeRF backbone produces more accurate geometry for the Materials scene. Finally we also observe the effect of the backbones in Fig. 5 as the ficus and drums scenes have shiny features that are better reconstructed when using the Ref-NeRF backbone.

4.2. Real unbounded scenes

In this section we show results on unbounded scenes from Mip-NeRF 360 [1]. We evaluate on the publicly available scenes. We use Mip-NeRF 360 [1] accelerated with NGP from [14] as our NeRF backbone architecture for experiments on this dataset. Note that there is no ground truth geometry available for these scenes so we do not evaluate the geometry numerically.

When evaluating on unbounded scenes, we extract the mesh using marching cubes separately for both the scene foreground and scene background. To limit the total size of the mesh, we use a foreground box of resolution 1024 to mesh the foreground and separately use a background

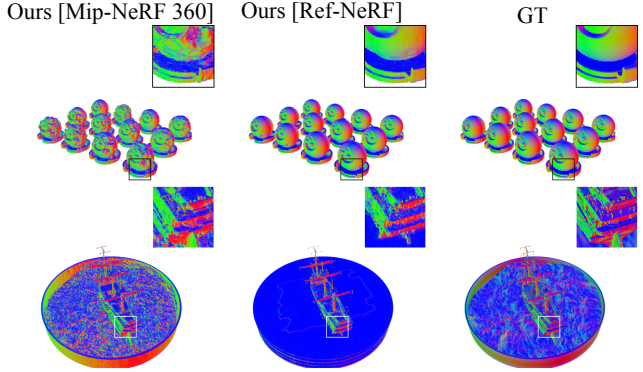


Figure 6. Effect of the NeRF backbone.

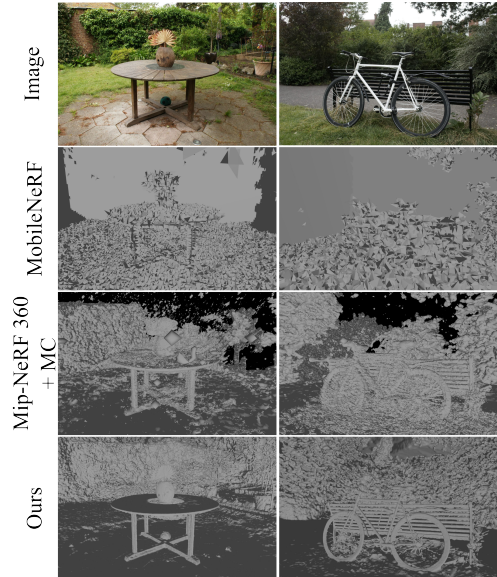


Figure 7. Results on scenes from the unbounded dataset [1] shown on the left. Rendered color images from our 3D meshes are shown in the right. Note that meshes from our method captures the scene geometry much better compared to MobileNeRF [3] and [1].

box of resolution 2048 to mesh the background. Due to the higher number of vertices we build a feature texture to retain image resolution. We construct a per face parametrization of the texture and map each face to a right-angled triangle of side 4 pixels in the texture image, we provide more details in the supplementary materials.

Fig. 7 provides a qualitative comparison of the underlying geometry obtained from our method with other methods. As shown in Fig. 7, our method produces more accurate geometry compared to MobileNeRF meshes and Mip-NeRF 360 meshes reconstructed via marching cubes. A key contributor of this advantage is the explicit regularization of the SDF via our SSAN. Rendered novel views from our method are shown in Fig. 8. We observe in Fig. 8 and Tab. 2 that our method remains of good visual quality even though it is constrained by accurate geometry.

Metric	Method	Chair	Drums	Ficus	Hotdog	Lego	Mats.	Mic	Ship	Avg
PSNR \uparrow	MobileNeRF [3]	34.09	25.02	30.20	35.46	34.18	26.72	32.48	29.06	30.90
	SNeRG [6]	33.24	24.57	29.32	34.33	33.82	27.21	32.60	27.97	30.38
	Ours [Mip-Nerf 360]	31.44	23.42	26.46	31.63	29.01	20.36	27.41	26.41	27.02
	Ours [Ref-NeRF]	31.93	23.49	25.95	32.38	28.89	23.30	27.83	24.70	27.31
CD \downarrow	Mip-Nerf 360 + Marching Cubes	0.016	0.047	0.032	0.027	0.024	0.019	0.022	0.044	0.029
	MobileNeRF	0.015	0.02	0.018	0.019	0.020	0.017	0.022	0.029	0.020
	Ours [Mip-Nerf 360]	0.012	0.018	0.012	0.026	0.018	0.023	0.016	0.019	0.018
	Ours [Ref-NeRF]	0.010	0.018	0.008	0.042	0.015	0.014	0.013	0.104	0.028
NC \uparrow	Mip-Nerf 360 + Marching Cubes	0.808	0.750	0.779	0.815	0.645	0.836	0.763	0.696	0.762
	MobileNeRF	0.615	0.550	0.655	0.605	0.546	0.49	0.519	0.535	0.565
	Ours [Mip-Nerf 360]	0.855	0.793	0.809	0.843	0.711	0.824	0.762	0.772	0.796
	Ours [Ref-NeRF]	0.857	0.803	0.865	0.854	0.714	0.908	0.773	0.611	0.798

Table 1. Quantitative results on the Blender Synthetic dataset [11]. We measure appearance with PSNR and geometry with Chamfer Distance (CD) and normal consistency (NC). Our method achieves better geometric reconstruction overall while producing reasonable PSNR.

Metric	Method	Bicycle	Garden	Stump	Counter	Room	Bonsai	Kitchen
PSNR \uparrow	NeRF [13]	21.76	23.11	21.73	25.67	28.56	26.81	26.31
	MobileNeRF [3]	21.70	23.54	23.95	-	-	-	-
	Ours	21.15	22.91	22.66	20.00	26.13	25.58	23.59

Table 2. Novel view evaluation on real unbounded scenes from [1].

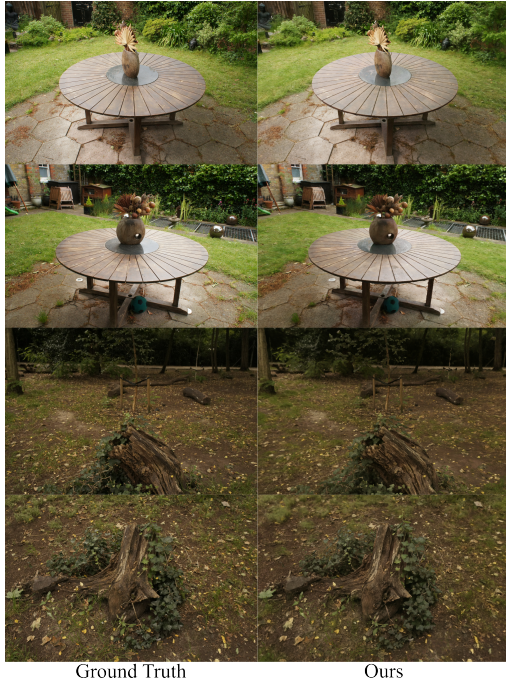


Figure 8. Rendered novel views of our method on real unbounded scenes from Mip-NeRF 360 [1]. Note that rendered images are of high fidelity and captures view dependent effects while being rendered real-time on commodity hardware.

4.3. Physics-based applications

Since our method produces accurate 3D meshes, several applications such as scene editing and physics simulation

can be easily conducted with traditional graphics and simulation pipelines. In Fig. 1, we show simple scene editing where we combine meshes extracted from *Chair*, *Hotdog*, *Lego* (Blender dataset), *Bonsai* and *Garden* (Mip-NeRF 360 dataset) sequences using our method. We further show a simulation of a cloth falling on objects from Blender Synthetic dataset. Please see the supplementary material for additional video results and demos.

5. Conclusion

In this work we propose a novel approach to extract geometrically accurate meshes from NeRF based architectures. Our SSAN model can be trained from any NeRF architecture without a significant penalty in training time. Our neural mesh representation can be rendered at high frame rates on commodity hardware. Thanks to their geometric accuracy, our extracted meshes can be quickly visualized, and also be used in physically accurate settings in simulations, to compute accurate occlusions and interactions with other objects. Our work is nevertheless limited in some aspects: rendering detailed surfaces causes us to generate meshes with a high number of faces and vertices, thus, our method would gain from using an adaptative mesh reconstruction strategy. Moreover, on large (unbounded) and detailed scenes, we are limited to lower resolutions to not hinder the overall size of the model. Finally, the appearance network that we are using is a non-standard component that requires customization of existing renderers. Being able to learn the underlying materials, as methods like [2] would further simplify the integration of meshes in other tools.

References

- [1] Jonathan T. Barron, Ben Mildenhall, Dor Verbin, Pratul P. Srinivasan, and Peter Hedman. Mip-nerf 360: Unbounded anti-aliased neural radiance fields. *CVPR*, 2022. 1, 2, 5, 6, 7, 8, 10
- [2] Mark Boss, Andreas Engelhardt, Abhishek Kar, Yuanzhen Li, Deqing Sun, Jonathan T. Barron, Hendrik P.A. Lensch, and Varun Jampani. SAMURAI: Shape And Material from Unconstrained Real-world Arbitrary Image collections. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2022. 8
- [3] Zhiqin Chen, Thomas Funkhouser, Peter Hedman, and Andrea Tagliasacchi. Mobilenerf: Exploiting the polygon rasterization pipeline for efficient neural field rendering on mobile architectures. *arXiv preprint arXiv:2208.00277*, 2022. 1, 2, 5, 6, 7, 8
- [4] Kangle Deng, Andrew Liu, Jun-Yan Zhu, and Deva Ramanan. Depth-supervised NeRF: Fewer views and faster training for free. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2022. 2
- [5] Peter Hedman, Pratul P. Srinivasan, Ben Mildenhall, Jonathan T. Barron, and Paul Debevec. Baking neural radiance fields for real-time view synthesis. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5875–5884, 2021. 1
- [6] Peter Hedman, Pratul P. Srinivasan, Ben Mildenhall, Jonathan T. Barron, and Paul Debevec. Baking neural radiance fields for real-time view synthesis. *ICCV*, 2021. 1, 2, 4, 5, 8
- [7] Rasmus Jensen, Anders Dahl, George Vogiatzis, Engin Tola, and Henrik Aanæs. Large scale multi-view stereopsis evaluation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 406–413, 2014. 6
- [8] Tao Ju, Frank Losasso, Scott Schaefer, and Joe Warren. Dual contouring of hermite data. In *Proceedings of the 29th annual conference on Computer graphics and interactive techniques*, pages 339–346, 2002. 5
- [9] Chen-Hsuan Lin, Wei-Chiu Ma, Antonio Torralba, and Simon Lucey. Barf: Bundle-adjusting neural radiance fields. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5741–5751, 2021. 1
- [10] William E Lorensen and Harvey E Cline. Marching cubes: A high resolution 3d surface construction algorithm. *ACM siggraph computer graphics*, 21(4):163–169, 1987. 5
- [11] Ricardo Martin-Brualla, Noha Radwan, Mehdi SM Sajjadi, Jonathan T Barron, Alexey Dosovitskiy, and Daniel Duckworth. Nerf in the wild: Neural radiance fields for unconstrained photo collections. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7210–7219, 2021. 1, 8
- [12] Lars Mescheder, Michael Oechsle, Michael Niemeyer, Sebastian Nowozin, and Andreas Geiger. Occupancy networks: Learning 3d reconstruction in function space. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2019. 6
- [13] 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. In *ECCV*, 2020. 1, 2, 5, 6, 8
- [14] Thomas Müller, Alex Evans, Christoph Schied, and Alexander Keller. Instant neural graphics primitives with a multiresolution hash encoding. *ACM Trans. Graph.*, 41(4):102:1–102:15, July 2022. 1, 2, 5, 6, 7, 10
- [15] Jacob Munkberg, Jon Hasselgren, Tianchang Shen, Jun Gao, Wenzheng Chen, Alex Evans, Thomas Mueller, and Sanja Fidler. Extracting Triangular 3D Models, Materials, and Lighting From Images. *arXiv:2111.12503*, 2021. 2
- [16] Michael Niemeyer, Jonathan T Barron, Ben Mildenhall, Mehdi SM Sajjadi, Andreas Geiger, and Noha Radwan. Regnerf: Regularizing neural radiance fields for view synthesis from sparse inputs. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5480–5490, 2022. 1
- [17] Michael Oechsle, Songyou Peng, and Andreas Geiger. Unisurf: Unifying neural implicit surfaces and radiance fields for multi-view reconstruction. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5589–5599, 2021. 2
- [18] Keunhong Park, Utkarsh Sinha, Jonathan T Barron, Sofien Bouaziz, Dan B Goldman, Steven M Seitz, and Ricardo Martin-Brualla. Nerfies: Deformable neural radiance fields. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5865–5874, 2021. 3
- [19] Barbara Roessle, Jonathan T. Barron, Ben Mildenhall, Pratul P. Srinivasan, and Matthias Nießner. Dense depth priors for neural radiance fields from sparse input views. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2022. 2
- [20] Sara Fridovich-Keil and Alex Yu, Matthew Tancik, Qinhong Chen, Benjamin Recht, and Angjoo Kanazawa. Plenoxels: Radiance fields without neural networks. In *CVPR*, 2022. 2
- [21] Dor Verbin, Peter Hedman, Ben Mildenhall, Todd Zickler, Jonathan T Barron, and Pratul P Srinivasan. Ref-nerf: Structured view-dependent appearance for neural radiance fields. In *2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 5481–5490. IEEE, 2022. 1, 2, 4, 5
- [22] Peng Wang, Lingjie Liu, Yuan Liu, Christian Theobalt, Taku Komura, and Wenping Wang. Neus: Learning neural implicit surfaces by volume rendering for multi-view reconstruction. *NeurIPS*, 2021. 2
- [23] Yiheng Xie, Towaki Takikawa, Shunsuke Saito, Or Litany, Shiqin Yan, Numair Khan, Federico Tombari, James Tompkin, Vincent Sitzmann, and Srinath Sridhar. Neural fields in visual computing and beyond. In *Computer Graphics Forum*, volume 41, pages 641–676. Wiley Online Library, 2022. 1
- [24] Lior Yariv, Jiatao Gu, Yoni Kasten, and Yaron Lipman. Volume rendering of neural implicit surfaces. In *Thirty-Fifth Conference on Neural Information Processing Systems*, 2021. 2
- [25] Lior Yariv, Yoni Kasten, Dror Moran, Meirav Galun, Matan Atzmon, Basri Ronen, and Yaron Lipman. Multiview neu-

ral surface reconstruction by disentangling geometry and appearance. *Advances in Neural Information Processing Systems*, 33:2492–2502, 2020. 2

- [26] Alex Yu, Ruilong Li, Matthew Tancik, Hao Li, Ren Ng, and Angjoo Kanazawa. PlenOctrees for real-time rendering of neural radiance fields. In *ICCV*, 2021. 2

Appendix

A. Joint training

We choose to train our method from a pre-trained NeRF method. In this section we measure the performance of training our method simultaneously to an un-trained Mip-NeRF 360 using instant NGP [14]. We observe in Tab. 3 that joint training leads to better appearance. However, due to the increased simplicity we choose to evaluate our method as post-processing from a pre-trained NeRF method.

B. Ablation study

We evaluate the effect of some of our design choices for our method trained with a Mip-NeRF 360 with Instant NGP on the Blender Synthetic dataset in Tab. 4.

Effect of choice of percentiles. We train our method on the blender synthetic dataset with different choices of percentile values. We notice that this choice only has a minor effect (*c.f.* Tab. 4 row 2, 3, 5).

Effect of separate appearance network. We train our method with a single network to represent both the geometry and the appearance. We notice that separating the network leads to a small improvement over PSNR whilst geometric metrics remain comparable Tab. 4 row 1, 5.

Effect of the projection. We train our method without the projection defined in Equation 10. We notice that it leads to a small drop in PSNR while geometric metrics are not much affected (*c.f.* Tab. 4 row 4, 5).

C. Additional implementation details

C.1. Projection on the zero level set

The median depth computed from a pretrained NeRF can sometimes be inaccurate, especially when the angle between the view direction and the surface is small. To get a better approximation of the zero level-set of the surface we apply 4 steps of gradient descent along the viewing ray using the TSDF predicted value. We define one step applied to a 3D point p as :

$$p' = p + \alpha * \hat{t}, \text{ with } [\hat{t}, \hat{n}, \hat{f}] = \phi(p), \quad (17)$$



Figure 9. Representation of each face of the mesh in the texture image. Each face is mapped to a right angled triangle containing features predicted by SSAN.

where α is a small constant.

C.2. Per face parametrization

When evaluating on unbounded scenes from [1], we represent scene features using textures with per face parametrization. We map each face in the triangle mesh to a right angled triangle of side 4 in the texture image. We show the face shape in Fig. 9. Each pixel feature \hat{f} is computed using the SSAN at the corresponding 3D position.

C.3. Training times

Our method has 4.8 minutes training time on average on the Blender Synthetic Dataset.

Metric	Method	Chair	Drums	Ficus	Hotdog	Lego	Mats.	Mic	Ship	Avg
PSNR	Ours [post-processing]	31.44	23.42	26.46	31.63	29.01	20.36	27.41	26.41	27.02
	Ours [joint-training]	31.54	23.83	26.70	32.16	29.21	21.54	27.89	26.63	27.44
CD	Ours [post-processing]	0.012	0.018	0.012	0.026	0.018	0.023	0.016	0.019	0.018
	Ours [joint-training]	0.013	0.018	0.013	0.030	0.017	0.021	0.015	0.028	0.019
NC	Ours [post-processing]	0.855	0.793	0.809	0.843	0.711	0.824	0.762	0.772	0.796
	Ours [joint-training]	0.842	0.803	0.844	0.860	0.697	0.809	0.758	0.769	0.798

Table 3. Effect of training stages on the blender synthetic dataset.

	PSNR	CD	NC
w/o appearance network separation	26.74	0.018	0.802
w/ percentiles 5 & 95	27.01	0.018	0.796
w/ percentiles 25 & 75	27.02	0.018	0.796
w/o projection	26.82	0.018	0.798
Ours	27.02	0.018	0.796

Table 4. Ablation study on the blender synthetic dataset.