

# Anti-Aliased Neural Implicit Surfaces with Encoding Level of Detail

YIYU ZHUANG<sup>†</sup>, Njing University, China  
 QI ZHANG<sup>†</sup>, Tencent AI Lab, China  
 YING FENG, Tencent AI Lab, China  
 HAO ZHU, Njing University, China  
 YAO YAO, Njing University, China  
 XIAOYU LI, Tencent AI Lab, China  
 YAN-PEI CAO, Tencent AI Lab, China  
 YING SHAN, Tencent AI Lab, China  
 XUN CAO, Njing University, China

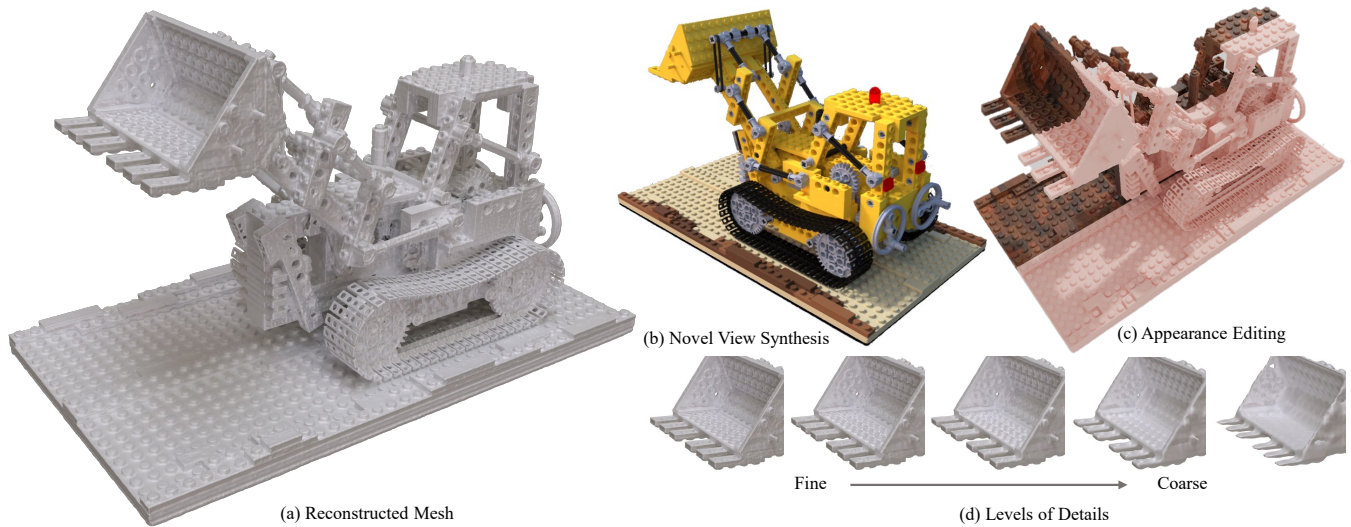


Fig. 1. Our method, called *LoD-NeuS*, adaptively encodes Level of Detail (LoD) features derived from the multi-scale and multi-convolved tri-plane representation. By optimizing a neural Signal Distance Field (SDF), our method is capable of reconstructing high-fidelity geometry (a). *LoD-NeuS* effectively captures varying levels of detail (d), resulting in anti-aliasing reconstruction, and thus, enabling photorealistic view synthesis (b) and appearance editing (c).

<sup>†</sup>Both authors contributed equally to this work. Zhuang did this work during the internship at Tencent AI Lab mentored by Zhang.

Authors' addresses: Yiyu Zhuang, Njing University, Nanjing, China, yiyu.zhuang@smail.nju.edu.cn; Qi Zhang, Tencent AI Lab, Shenzhen, China, nwpuzhang@gmail.com; Ying Feng, Tencent AI Lab, Shenzhen, China, yfeng.von@gmail.com; Hao Zhu, Njing University, Nanjing, China, zhuhaoese@nju.edu.cn; Yao Yao, Njing University, Nanjing, China, yyaoag@cse.ust.hk; Xiaoyu Li, Tencent AI Lab, Shenzhen, China, xliea@connect.ust.hk; Yan-Pei Cao, Tencent AI Lab, Shenzhen, China, caoyanpei@gmail.com; Ying Shan, Tencent AI Lab, Shenzhen, China, yingsshan@tencent.com; Xun Cao, Njing University, Nanjing, China, caoxun@nju.edu.cn.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](mailto:permissions@acm.org).

© 2023 Association for Computing Machinery.  
 0730-0301/2023/9-ART \$15.00  
<https://doi.org/10.1145/nnnnnnn.nnnnnnn>

We present *LoD-NeuS*, an efficient neural representation for high-frequency geometry detail recovery and anti-aliased novel view rendering. Drawing inspiration from voxel-based representations with the level of detail (LoD), we introduce a multi-scale tri-plane-based scene representation that is capable of capturing the LoD of the signed distance function (SDF) and the space radiance. Our representation aggregates space features from a multi-convolved featurization within a conical frustum along a ray and optimizes the LoD feature volume through differentiable rendering. Additionally, we propose an error-guided sampling strategy to guide the growth of the SDF during the optimization. Both qualitative and quantitative evaluations demonstrate that our method achieves superior surface reconstruction and photorealistic view synthesis compared to state-of-the-art approaches.

CCS Concepts: • **Computing methodologies** → **Volumetric models; Antialiasing.**

Additional Key Words and Phrases: Neural Implicit Surface, Signed Distance Function, Volume Rendering, Neural Radiance Fields, Anti-aliasing

## ACM Reference Format:

Yiyu Zhuang, Qi Zhang, Ying Feng, Hao Zhu, Yao Yao, Xiaoyu Li, Yan-Pei Cao, Ying Shan, and Xun Cao. 2023. Anti-Aliased Neural Implicit Surfaces

with Encoding Level of Detail. *ACM Trans. Graph.* 1, 1 (September 2023), 11 pages. <https://doi.org/10.1145/nnnnnnn.nnnnnnn>

## 1 INTRODUCTION

Recent advances in implicit representation and neural rendering (i.e., NeRF [Mildenhall et al. 2020] approaches) have provided a new alternative for geometric modeling and novel view rendering. However, applying the vanilla NeRF with the soft density representation to accurately reconstruct the geometry with fine-grained surface details remains challenging. In contrast, the neural implicit surface (NeuS) [Wang et al. 2021] was proposed to apply the signed distance function (SDF) rather than the soft density to model the object surface within the NeRF framework explicitly. The object surface is represented as the zero-level set of the SDF modeled by the multi-layer perceptron (MLP). NeuS and its variants have shown that SDF can flexibly represent the scene geometry with arbitrary topologies, and produce significantly better results in neural surface reconstruction than the vanilla NeRF approach.

One of the major challenges of neural surface reconstruction is the reconstruction of high-frequency surface details. While frequency position encoding has been employed in NeuS, it still struggles to capture fine-grained geometry details accurately, resulting in low-fidelity and over-smooth geometric approximations for intricate models. HF-NeuS [Wang et al. 2022] attempts to mitigate the issue by introducing a displacement network tailored specifically for learning high-frequency geometry details. However, the problem persists due to the inherent limitations of frequency position encoding, which lacks locality and fails to adaptively capture different level of detail (LoD) in surface geometry. Consequently, the undersampling and inadequate representation of high-frequency information inevitably results in aliasing artifacts during novel view rendering.

On the other hand, explicit voxel-based representations have long employed multi-scale prefiltering techniques, such as mipmaps and octrees, to enable fine-grained surface recovery and anti-aliasing in object rendering. Recent advancements in NeuS-based approaches have also explored the potential of hybrid implicit-explicit representations. These methods replace multi-layer perceptrons (MLPs) with discretized volumetric representations, such as voxel grids [Yu et al. 2022] and tri-planes [Wang et al. 2023], resulting in better geometric approximations. However, these methods possess inherent limitations (Sec. 2), which pose challenges when combining the anti-aliasing advantages of explicit methods with hybrid representations in surface reconstruction with continuous LoD.

In this paper, we present a neural implicit surface representation with the encoding level of detail (*LoD-NeuS*) for high-quality geometry reconstruction from multi-view images. The implicit surface is represented by a multi-scale tri-plane-based feature volume, which is optimized through differentiable cone sampling and volume rendering.

- (1) We present a tri-plane position encoding, optimizing multi-scale features, to effectively capture different levels of detail;
- (2) We design a multi-convolved featurization within a conical frustum, to approximate cone sampling along a ray, which enables the anti-aliasing recovery with finer 3D geometric details;

- (3) We develop a refinement strategy, involving error-guided sampling, to facilitate SDF growth for thin surfaces.

In experiments, our method outperforms state-of-the-art NeuS-based approaches at high-quality surface reconstruction and view synthesis, particularly for objects and scenes with high-frequency details and thin surfaces.

## 2 RELATED WORK

**Multi-view 3D Reconstruction.** Reconstructing the surfaces of the scene from multi-view images is a fundamental problem that has been extensively studied throughout the development of computer vision and graphics [Han et al. 2019]. Multi-view 3D reconstruction has three categories: point-based reconstruction [Barnes et al. 2009; Campbell et al. 2008; Furukawa and Ponce 2009; Schönberger et al. 2016; Tola et al. 2012], surface reconstruction [Dai et al. 2017; Hoppe et al. 1992; Izadi et al. 2011; Kazhdan et al. 2006], and volumetric reconstruction [Broadhurst et al. 2001; De Bonet and Viola 1999; Kutulakos and Seitz 2000; Seitz and Dyer 1999]. Traditionally, point-based or surface-based methods first estimate the geometry information (e.g., depth and normal maps) of each pixel by matching the correspondences of multi-view images [Schonberger and Frahm 2016], and then fuse the geometry information [Merrell et al. 2007; Zach et al. 2007] followed by mesh surface reconstruction processes such as Delaunay triangulation [Labatut et al. 2007] and ball-pivoting [Bernardini et al. 1999]. The performance of surface reconstruction largely depends on the accuracy of correspondence matching. Recovering surfaces with minimal textures can be challenging, leading to significant artifacts and partially missing reconstructed content. To circumvent issues with insufficient geometry correspondence, volumetric methods [Nießner et al. 2013; Sitzmann et al. 2019] estimate occupancy and color within a voxel grid from multi-view images and evaluate color consistency at each grid. Nevertheless, these volumetric approaches involve explicitly breaking down a scene into a vast number of samples, which necessitates substantial storage capacity, thereby limiting grid resolution and impacting the overall reconstruction quality.

**Neural Implicit Surface.** Recent advances in implicit neural representations have showcased the potential to reconstruct highly detailed surfaces and render photorealistic views [Tewari et al. 2022]. Neural Radiance Fields (NeRF) [Mildenhall et al. 2020], a notable breakthrough in this domain, learns the radiance fields (density and view-dependent color) of a scene and renders novel views based on volumetric ray tracing. NeRF and its variations have been applied to a range of tasks, including novel view synthesis [Barron et al. 2022; Chen et al. 2022b; Zhu et al. 2023], generalizable models [Wu et al. 2023; Xin et al. 2023; Zhuang et al. 2022], imaging processing [Huang et al. 2023, 2022; Ma et al. 2022], and inverse rendering [Srinivasan et al. 2021; Verbin et al. 2022; Zhuang et al. 2023]. However, compared to the signed distance function (SDF) [Chabra et al. 2020; Genova et al. 2020] or occupancy field [Mescheder et al. 2019; Oechsle et al. 2021], recovering smooth and accurate surfaces using the density function is challenging, often produce noisy low-fidelity geometry approximation since it lacks sufficient constraints on its level sets. Specifically, VolSDF [Yariv et al. 2021] incorporates an SDF into the density function, ensuring that it satisfies a derived

error bound on the transparency function. NeuS [Wang et al. 2021] proposes an unbiased formulation with a logistic sigmoid function and introduces a learnable parameter to control the slope of the function during the rendering and sampling processes. Building upon NeuS, NeuralWarp [Darmon et al. 2022] and Geo-NeuS [Fu et al. 2022] leverage prior geometry information from MVS methods but may struggle in the regions with less texture. HF-NeuS [Wang et al. 2022] integrates additional displacement networks to fit the high-frequency details. However, the frequency position encoding used in these methods struggles to adaptively capture varying levels of detail (LoD) across different regions. Additionally, using ray sampling instead of cone sampling leads to undersampled or inaccurately represented high-frequency information, which finally results in aliasing artifacts.

**Anti-aliased Representation.** Surfaces employing traditional explicit representations (e.g., polygon mesh, voxel grids), can be efficiently reconstructed without encountering aliasing artifacts, thanks to the application of multi-scale prefilter techniques. These techniques, such as mipmaps and octrees, offer a robust solution for handling different levels of detail in surfaces while maintaining efficiency. Continuous implicit surface representations achieve higher performance but can only be anti-aliased through supersampling, which further slows down their already time-consuming reconstruction. The hybrid explicit-implicit representation emerges as a result. In particular, Takikawa et al. [2021] proposes a multi-scale representation based on sparse voxel octrees for implicit surface with a learned geometry prior. MonoSDF [Yu et al. 2022] employs a multi-scale voxel-based representation with monocular geometric cues for SDF reconstruction, which introduces the hash encoding [Müller et al. 2022] to optimize the grid feature. Although hash encoding can enhance both memory efficiency and performance, it potentially causes hash collisions and representations that are insufficiently explicit. PET-NeuS [Wang et al. 2023] adopts a self-attention convolution to generate the tri-plane-based representation [Chan et al. 2022; Chen et al. 2022a] for enhancing quality, but following positional encoding on both tri-plane features and position increase model parameters and computational complexity. Besides, the performance of PET-NeuS heavily depends on the effectiveness of self-attention. Inspired by these ideas, we introduce an efficient neural representation to aggregate LoD features, for the first time, that enables continuous LoD with cone sampling while achieving state-of-the-art geometry reconstruction quality.

### 3 PRELIMINARIES

This section overviews the base priors of NeRF [Mildenhall et al. 2020] for volume rendering, as well as geometric improvements extended by SDF and NeuS [Wang et al. 2021] for surface reconstruction and view synthesis.

**NeRF** represents the scene with a continuous volumetric radiance field, which utilizes MLPs to map the position  $\mathbf{x}$  and view direction  $\mathbf{r}$  to a density  $\sigma$  and color  $\mathbf{c}$ . To render a pixel's color, NeRF casts a single ray  $\mathbf{r}(t) = \mathbf{o} + t\mathbf{d}$  through the pixel and samples a set of points with different  $\{t_i\}$  along the ray. The evaluated  $\{(\sigma_i, \mathbf{c}_i)\}$  at the sampled points are accumulated into the color  $C(\mathbf{r})$  of the pixel

via volume rendering [Max 1995]:

$$C(\mathbf{r}) = \sum_i T_i \alpha_i \mathbf{c}_i, \text{ where } T_i = \exp\left(-\sum_{k=0}^{i-1} \sigma_k \delta_k\right), \quad (1)$$

and  $\alpha_i = 1 - \exp(-\sigma_i \delta_i)$  indicates the opacity of the sampled point. Accumulated transmittance  $T_i$  quantifies the probability of the ray traveling from  $t_0$  to  $t_i$  without encountering other particles, and  $\delta_i = t_i - t_{i-1}$  denotes the distance between adjacent samples.

**NeuS** extends the basic NeRF formulation by integrating an SDF into volume rendering. It represents the scene's geometry with a learnable function  $f$ , which returns the signed distance  $f(\mathbf{x})$  from each point to the surface. The underlying surface can be derived from the zero-level set,

$$S = \{\mathbf{x} \in \mathbb{R}^3 | f(\mathbf{x}) = 0\}. \quad (2)$$

Subsequently, NeuS defines a function to map the signed distance to density  $\sigma$ , which attains a locally maximal value at surface intersection points. Specifically, accumulated transmittance  $T(t)$  along the ray  $\mathbf{r}(t) = \mathbf{o} + t\mathbf{d}$  is formulated as a sigmoid function:  $T(t) = \Phi(f(t)) = (1 + e^{sf(t)})^{-1}$ , where  $s$  and  $f(t)$  refers to a learnable parameter and the SDF value of point at  $\mathbf{r}(t)$ , respectively. Discrete opacity values  $\alpha_i$  can then be derived as:

$$\alpha_i = \max\left(\frac{\Phi_s(f(t_i)) - \Phi_s(f(t_{i+1}))}{\Phi_s(f(t_i))}, 0\right). \quad (3)$$

NeuS employs volume rendering to recover the underlying SDF based on Eqs. (1) and (3). The SDF is optimized by minimizing the photometric loss between the renderings and ground-truth images.

## 4 METHOD

Drawing inspiration from anti-aliasing techniques for explicit voxel-based surface reconstruction, we aim to develop a hybrid representation that combines the advantages of both explicit and implicit representations, to achieve anti-aliasing and the recovery of delicate geometric details. In particular, we firstly present a novel position encoding based on multi-scale tri-planes to enable continuous levels of details (Sec. 4.1). To alleviate aliasing, we consider the size of cast cone rays (similar to [Barron et al. 2022]) and specifically design multi-convolved features to approximate the cone sampling (Sec. 4.2). Meanwhile, we observe that thin surface reconstruction using SDF is challenging, thus propose a refined solution involving an error-guided sampling strategy to facilitate SDF growth (Sec. 4.4).

### 4.1 Multi-scale Tri-plane Encoding

Recent progress [Müller et al. 2022; Yu et al. 2022] in the field of neural rendering has shown that incorporating learnable features extracted from multi-scale grids significantly enhances reconstruction quality and accelerates volume rendering. In contrast to voxel-based representations with heavy memory requirements and hash encoding with collision issues, tri-plane-based representations [Chan et al. 2022; Chen et al. 2022a] provide increased flexibility in handling complex geometry and effective spatial regularization. Inspired by these insights, we incorporate the multi-scale tri-plane representation into a NeuS-based framework for intricate surface reconstruction and high-quality rendering.

To address the challenges associated with reconstructing high-frequency details and achieve a more reasonable implicit surface

representation, we propose a learnable encoding based on multi-scale tri-planes. A tri-plane representation  $\mathbf{P}$  is a novel 3D data structure, which consists of three learnable feature planes  $\{P_{xy}, P_{xz}, P_{yz}\}$ . These planes are orthogonal to each other and form a 3D cube centered at the origin  $(0, 0, 0)$ . For each 3D point  $\mathbf{x} \in \mathbb{R}^3$ , we project it onto each of the three planes, gathering features  $F_{xy}, F_{xz}, F_{yz}$  using bilinear interpolation. The element-wise concatenation of these features yields the feature  $\mathbf{F}$  with dimensionality  $N$ .

Unlike previous methods [Chan et al. 2022; Chen et al. 2022a], we construct a set of tri-planes with different resolution  $\{R_l\}_{l=1}^L$ , where  $L$  indicates the number of levels. Each level is independent and stores feature at the vertices of tri-plane. Our position encoding function,  $\gamma(\mathbf{x})$ , concatenates the input  $\mathbf{x}$  with the feature  $\mathbf{F}_l$  from every level  $l$ , forming a multi-scale feature vector  $\vec{\mathbf{F}} = (\mathbf{x}, \mathbf{F}_1, \dots, \mathbf{F}_L)$ , whose length is  $3+L \times N$ . By replacing the traditional frequency position encoding with the multi-scale tri-plane feature vector, our method benefits from explicit representation while guarantees different levels of detail.

## 4.2 Anti-aliasing Rendering of Implicit Surfaces

Once we have acquired multi-scale tri-plane features, our goal is to estimate the SDF of samples along a ray for volume rendering. NeuS renders a pixel's color by casting a single ray through the pixel, without considering its size and shape. This approximation potentially leads to undersampling or ambiguous representation of high-frequency information, and results in aliasing artifacts. To alleviate this, we reformulate volume rendering by defining a ray as a cone, taking into account the pixel size. This enables the continuous LoD and recovers a high-quality SDF from undersampled images, leading to more accurately capture and reconstruction of the fine details of the scene.

A straightforward solution is to discretize a cone into a batch of rays, similar to super-sampling techniques [Cook 1986]. However, this approach increases the number of sampled rays and points for volume rendering, leading to prohibitively high computational costs and slow inference time.

Here we provide a more efficient solution, including a sampling strategy and featurization procedure, in which we cast a cone and integrate features within conical frustums, as shown in Fig. 2.

**Cone Discrete Sampling.** Assuming a casting cone ray through a camera pixel is divided into a series of conical frustums, we need to integrate the color and geometry information within each conical frustum. Drawing inspiration from Mip-NeRF [Barron et al. 2022], we attempt to integrate all features rather than just the network output. Note that, our target differs from Mip-NeRF, which focuses on rendering scenes at different resolutions rather than recovering scene details. Utilizing our tri-plane-based representation, we cast four additional rays through the pixel corners, thus takes the pixel size and shape into account. Each conical frustum along the cone is then represented by eight vertices. Given any 3D sampled position  $\mathbf{x}$  within a conical frustum, we blend the tri-plane features of each vertex  $\mathbf{x}_v$  using decreasing weights,

$$W(\mathbf{x}, \mathbf{x}_v) = \exp(-k|\mathbf{x}_v - \mathbf{x}|), \quad (4)$$

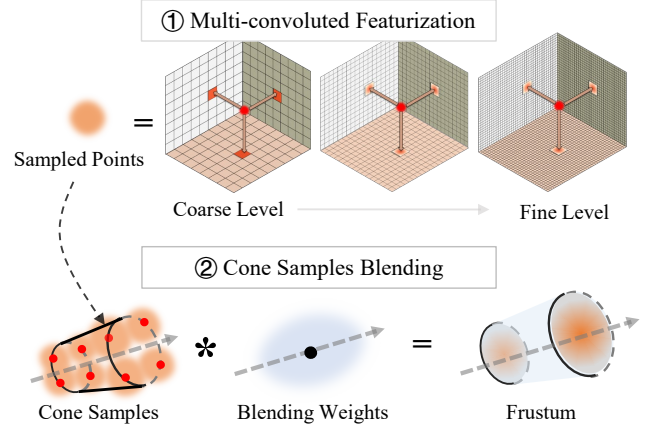


Fig. 2. Aggregation of LoD feature, including multi-convoluted featurization and cone discrete sampling. We obtain the feature of any sample within the conical frustum by blending the features of vertices. Additionally, considering the size of the sampled points, we introduce multi-convoluted features by Gaussian Kernel to efficiently represent ray sampling within a cone. Combining both of them, we aggregate the LoD feature of any sample in a continuous manner.

which decreases with the distance between the vertex  $\mathbf{x}_v$  and the sampled point  $\mathbf{x}$ .  $k$  is a learnable parameter that we initially set to 80 and update along with other parameters during training. It is important to note that the decreasing function should be aware of the size of the conical frustum. The smaller the conical frustum is, the more rapidly the function should decrease.

**Mult-convolved Featurization.** Though the multi-scale features of neighbor vertices along neighbor rays are applied for cone sampling, this approximation may be insufficient due to the sparse samples within the conical frustum. A straightforward way is to introduce more discretized samples, but this increases the computational cost and memory burden. Fortunately, the proposed explicit tri-plane-based representation makes it easy to integrate the features. Our approach utilizes the 2D Gaussian of each tri-plane to represent the region where the conical frustum should be integrated. In conjunction with our cone discrete sampling, we propose a multiple Gaussian convolved featurization to represent the features of neighbor vertices that approximate the sampled point and its corresponding conical frustum. Specifically, given a vertex  $\mathbf{x}_v$ , we project it onto the tri-planes and query the corresponding multi-scale feature vector  $\vec{\mathbf{F}}_v$ . Considering the grid resolution of the tri-plane, we apply multiple Gaussian convolutions with different kernel sizes. It represents the feature aggregation of samples within the conical frustum in a continuous manner. The multi-scale multi-convolved feature  $\mathbf{G}_v$  for each vertex of the conical frustum is defined as:

$$\mathbf{G}_v(\mathbf{x}_v) = \mathcal{G}(\vec{\mathbf{F}}_v, \{\tau_v\}_{l=1}^L) = \sqcup_{l=1}^L \mathcal{G}(\mathbf{F}_l, \tau_l), \quad (5)$$

where  $\mathcal{G}(\mathbf{F}, \tau)$  refers to our Gaussian convolution defined by covariance  $\tau$ . Through the 2D convolution, which aggregates the features of the compressed planes, we gather the features within a 3D sphere centered at  $\mathbf{x}_v$ . We choose different kernel sizes  $\{\tau_l\}$  for each level  $l$  to covers various frequency details. The convolved features are

combined using the concatenation operation, denoted as  $\sqcup$ . Consequently, similar to position encoding, our featurization is able to cover different frequency details of the scene.

According to Eqs. (4) and (5), for a sample at position  $\mathbf{x}$  within the corresponding conical frustum, its LoD feature with continuous levels of detail is defined as,

$$\mathbf{Z}(\mathbf{x}) = \sum_{v=1}^V W(\mathbf{x}, \mathbf{x}_v) \mathbf{G}_v(\mathbf{x}_v), \quad (6)$$

where  $V = 8$  is the number of vertices of a conical frustum. Compared to utilizing the 3D shape-adaptive Gaussian kernel to represent an ideal approximation of the cone discrete sampling, our formulation represents the solution spaces with a single learnable parameter  $k$  rather than dealing with the complexity of shape-adaptive Gaussian kernels.

### 4.3 Training and Loss

After obtaining LoD feature  $\mathbf{Z}$  of the samples along a ray, the colors and signed distance  $f$  can be predicted. We do this via a shallow 8-layer MLP:  $(f, \Theta) = \text{MLP}(\mathbf{Z})$ . In addition to  $f$ , it produces a feature vector  $\Theta \in \mathbb{R}^{256}$ , which is then passed to the color module. According to Eq. (3), we obtain the opacity  $\alpha$  of the sampled point. The color module is represented as a 3-layer MLP, which predicts the color  $\mathbf{c}$  from  $\Theta$  and view direction  $\mathbf{d}$  as  $\mathbf{c} = \text{MLP}_c(\Theta, \mathbf{d})$ . We finally follow Eq. (1) to render pixel color  $C_p$ .

The learnable parameters and networks are optimized by employing a loss function and the process of backward propagation. To be more specific, a batch of  $n$  pixels are randomly sampled, including their color  $\{C_p\}_{p=1}^n$  and optional masks  $\{M_p\}_{p=1}^n$ . We further sample  $m$  points along each ray, yielding the predicted color  $\{\hat{C}_p\}$ . Then the L1 loss is calculated to measure the reconstruction distance, which is defined as:

$$L_{rgb} = \frac{1}{n} \sum_p \|\hat{C}_p - C_p\|_1. \quad (7)$$

We also add an Eikonal term [Gropp et al. 2020] on all sampled points  $\{x_i\}_{i=1}^{nm}$  to regularize the SDF by:

$$L_{eikonal} = \frac{1}{nm} \sum_i (\|\nabla f(x_i)\|_2 - 1)^2. \quad (8)$$

The mask loss  $L_{mask}$  is optional and defined as:

$$L_{mask} = \frac{1}{n} \sum_p \text{BCE}(M_p, \hat{O}_p), \quad (9)$$

where  $\hat{O}_k = \sum_j^m T_j (1 - \exp(-\sigma_j \delta_j))$  is the opacity accumulated along the ray, and BCE is the binary cross entropy loss [Wang et al. 2021].

### 4.4 SDF Growth Refinement

We have observed that the SDF faces challenges when reconstructing thin objects, mainly for two reasons. First, representing a thin object necessitates a rapid flip in the SDF, which is difficult for the neural network [Yu et al. 2022]. Second, the image area corresponding to the thin object may have fewer samples compared to other areas, making it harder to learn. Based on this observation, a straightforward solution might be to increase the sampling frequency around this area. However, it's important to note that the optimization process of SDF is fundamentally different from NeRF.

---

#### ALGORITHM 1: SDF Growth Refinement

---

**Input:** SDF  $f$ , Rendered image  $I_r$ , Input image  $I$ , Step  $n$

**Output:** Refined SDF  $f_r$

$e \leftarrow \text{CalculateErrorMap}_{L1}(I_r, I)$ ;

$M_e \leftarrow \text{EstimateGrowthPoint}(I)$ ;

$M_s \leftarrow \text{ExpandRegion}(e)$ ;

**for** step  $\in \{1, \dots, n\}$  **do**

$M_e \leftarrow \text{ExpandRegion}(M_e) \cap M_s$ ;

$M_s \leftarrow M_s \setminus M_e$ ;

$\text{mask\_sequence.append}(M_e)$ ;

**end**

**for** step,  $\text{mask} \in \text{enumerate}(\text{mask\_sequence})$  **do**

$\text{sampld\_rays} \leftarrow \text{SampleRaysInsideMask}(\text{mask})$ ;

$f \leftarrow \text{OptimizeSDF}(f, \text{sampld\_rays})$ ;

**end**

$f_r \leftarrow f$ ;

---

In NeRF, the optimization runs in a spatially-independent manner, which means changes at one point does not affect another.

In contrast, the SDF represents the signed distance to the whole surface, which implies that when the implicit surface changes, the SDF values in a region are affected. Due to this interconnected nature, the optimization process of the SDF appears to deform the initial geometry to better fit the target surface. This property helps maintain connectivity and prevents floating artifacts. However, this also complicates the reconstruction of thin objects, as only the sampled rays located around the region are proven helpful for this reconstruction.

To utilize the property, we devise a strategy to refine the optimized SDF for better thin object reconstruction, as shown in Alg. 1 and Fig. 6. Our motivation is to guide the SDF growth from the spatial point where the missing thin segment meets the surface, utilizing the information from the 2D images. Specifically, we render the trained SDF at each training viewpoint, calculate the error map using the L1 distance against the inputs, sequentially binarize the map and dilate it to a candidate region  $M_e$ . To locate the beginning points of our growth method, we employ [Zhou et al. 2019] to detect the line endpoints and dilate them to our selected region  $M_s$ . We iterate through the following process: expand  $M_s$  and take the intersected set to form a new  $M_s$ , adding it into a mask list for training, and sequentially form a new  $M_e$  by removing the updated  $M_s$  region. After these preparations, the training is carried out one by one with the selected rays from the mask list.

## 5 EXPERIMENTS

### 5.1 Experimental settings

**Baselines.** We conduct a comparative analysis between our proposed method and prominent approaches, including NeuS [Wang et al. 2021], HF-NeuS [Wang et al. 2022], and NeRF [Mildenhall et al. 2020], which represents the state-of-the-art pipeline without any supplementary information. We have consciously excluded MoNoSDF [Yu et al. 2022], GeoNeuS [Fu et al. 2022], and NeuralWarp [Darmon et al. 2022] from the comparison, as these methods

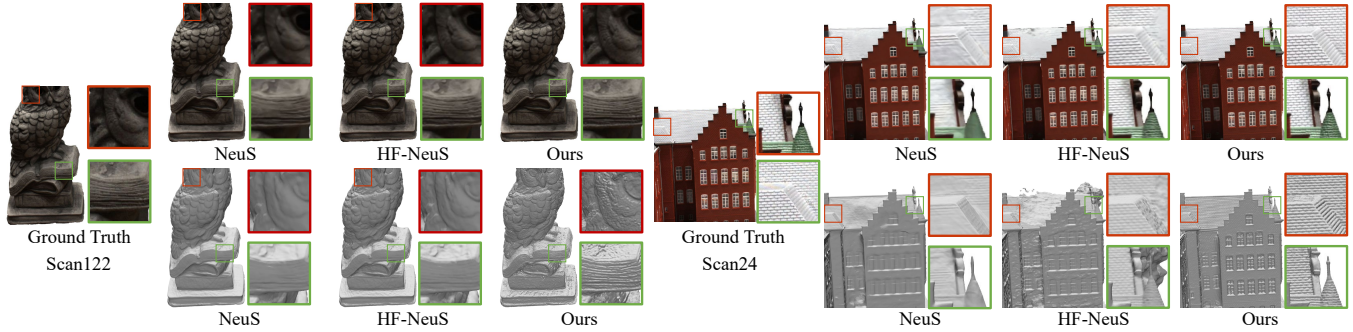


Fig. 3. Qualitative comparison with zoom-in details of our method against baselines on DTU dataset. Our method produces the most visually pleasing novel views and reconstructed geometry, especially on the intricate details and embossed patterns on the sculptures (left), along with the uneven surface created by the roof tiles (right), as shown in zoom-in details.

either introduce additional priors or employ constraints from multiple viewpoints, which are similarly applicable to our method. To qualitatively evaluate the performance, we adopt two criteria: the PSNR (Peak Signal-to-Noise Ratio) for gauging novel-view rendering quality, and the Chamfer distance for assessing the accuracy of the reconstructed mesh. We obtain the underlying mesh through the marching cubes algorithm on a grid with a resolution of 1500. For NeRF, following the approach used in NeuS, we extract the mesh using a threshold value of 25 for density.

**Dataset.** Following the setting of previous work, we report the metric on both the DTU dataset [Aanæs et al. 2016] and the NeRF-synthetic dataset. DTU is the multi-view stereo dataset. Each scene supplies 49 or 64 images with a resolution of  $1600 \times 1200$ , captured from various viewpoints. We adopt the foreground masks provided by IDR [Yariv et al. 2020] for these scenes. Additionally, we conduct further testing on 7 challenging scenes from the NeRF-synthetic dataset [Mildenhall et al. 2020], rendering 100 images each with resolution  $800 \times 800$  of black background, without the foreground mask. We designate every eighth image as the testing set, while the others as the training set.

**Implementation details.** In the following experiment, we select  $L = 5$ ,  $N = 6$ , including planes of resolution  $\{128, 256, 512, 1024, 2048\}$  and the Gaussian Kernel is of size  $\{1, 1, 1, 3, 5\}$ . We train our model for 300,000 iterations, with 512 rays randomly selected during each iteration. To ensure a fair and competitive comparison with NeuS, we maintain almost the same settings for our model.

## 5.2 Comparison.

Tab. 1 and Fig. 3 demonstrate quantitative and qualitative comparisons of our method against baseline methods, respectively. The results demonstrate that our method surpasses the baseline methods on the DTU dataset. Although the reported Chamfer distance does not exhibit a significant improvement, the qualitative comparison and PSNR clearly illustrate that our model can reproduce finer details. We believe this is attributed to the fact that the DTU ground truth point cloud is relatively coarse, and thus, enhancements in high-frequency details are not reflected in the metric. Additionally, the ground truth point cloud lacks some parts, which increases the distance between our reconstructed mesh and the ground truth. For further discussion, please refer to our supplementary materials.

Fig. 3 also showcases the incredible details that can be reproduced by our method, like the intricate details and embossed patterns on the sculptures (left), as well as the uneven surface created by the roof tiles (right), as illustrated in zoom-in details. This improvement can be demonstrated by the PSNR metric, in which our method outperforms the others, confirming the superior performance of our approach.

We also construct a comparison on 9 challenging models from NeRF-synthetic dataset [Mildenhall et al. 2020], which comprises objects exhibiting a higher level of high-frequency details. As depicted in Tab. 2, with the evaluation based on the ground truth meshes, the Chamfer distance reasonably demonstrates the superiority of our method, in contrast to the less accurate valuation on the DTU dataset. We provide a visual comparison in Fig. 4. Our method demonstrates superior performance in reproducing finer details, as evidenced by the zoom-in images of both the microphone’s grille and the Lego crawler belt, compared to other approaches. HF-NeuS and NeuS struggle with geometries featuring rapid changes, partly because they employ a fixed frequency of encoding and disregard the imaging model. Moreover, as the microphone’s electric wire is a relatively thin object, both HF-NeuS and NeuS fail to produce visually pleasing meshes. These methods introduce extra mesh sections connected to other parts and color them to match the background, resulting in accurate novel views but inherently inaccurate geometry. In contrast, our method not only excels in capturing rapid detail variation but also ensures smoothness on the microphone’s handle, thanks to our learnable multi-level encoding approach.

## 5.3 Ablation Study.

We developed our model building upon NeuS and retained most of its settings. To evaluate our proposed modules, we performed the following comparisons:

- TPE: Replacing Plane Encoding (PE) with Multi-scale Tri-plane Encoding;
- NGP: Replacing Plane Encoding (PE) with Multi-level Grid Compression via Hashing Encoding;
- TPE\_CS: Adding Cone Sampling to the model with Multi-scale Tri-plane Encoding;
- Ours: Combining Cone Sampling with Gaussian Convolution with Multi-scale Tri-plane Encoding.

	Chamfer Distance															Mean
	24	37	40	55	63	65	69	83	97	105	106	110	114	118	122	
NeuS	0.828	0.983	0.572	0.369	1.185	0.716	0.608	1.413	0.964	0.821	0.495	1.362	0.352	0.462	0.499	0.775
NeRF	1.418	1.611	1.665	0.799	1.856	1.288	1.203	1.603	1.645	1.113	0.947	2.101	0.977	1.027	0.918	1.345
HF-NeuS	1.113	1.276	0.609	0.465	0.973	0.682	0.619	1.344	0.914	0.728	0.534	1.816	0.378	0.536	0.510	0.833
Ours	0.652	0.913	0.373	0.482	1.049	0.869	0.821	1.216	0.954	0.693	0.564	1.301	0.416	0.584	0.569	0.764
	PSNR															Mean
	24	37	40	55	63	65	69	83	97	105	106	110	114	118	122	
NeuS	27.021	26.602	27.602	27.651	35.166	32.119	29.938	38.471	31.028	34.914	34.638	33.018	29.888	37.143	37.764	32.198
HF-NeuS	28.497	27.132	28.986	30.554	34.442	32.892	30.339	38.618	31.014	35.086	35.309	27.539	30.284	37.525	38.407	32.442
NeRF	29.564	26.608	28.351	29.537	35.838	32.853	29.941	38.576	31.225	35.389	36.324	33.504	30.379	37.332	38.154	32.905
Ours	30.489	27.325	30.052	31.387	36.111	32.348	29.985	39.189	31.824	36.318	36.519	34.370	31.089	38.251	39.235	33.633

Table 1. Quantitative results on the DTU dataset. Red and orange indicate the first and second best performing results.

	Chamfer Distance							PSNR						Mean
	Chair	Ficus	Lego	Materials	Mic	Ship	Mean	Chair	Ficus	Lego	Materials	Mic	Ship	
NeuS	1.350	0.121	0.143	0.103	0.364	0.758	0.473	28.590	25.234	29.348	29.197	29.998	26.412	28.130
HF-NeuS	0.531	0.074	0.070	0.113	1.971*	0.558	0.553	28.750	26.173	30.111	29.448	29.823	26.755	28.514
NeRF	1.501	3.660	1.299	1.069	1.396	4.272	2.199	28.196	25.545	28.474	30.882	27.074	26.644	27.803
Ours	0.502	0.063	0.038	0.039	0.094	0.368	0.184	30.468	26.713	31.488	30.710	34.107	28.360	30.308

Table 2. Quantitative results on the NeRF-synthetic dataset. The "Mic" case for HF-NeuS is marked with an asterisk (\*) as this metric is affected by the undesired reconstruction of mesh parts, as illustrated in Fig. 4.

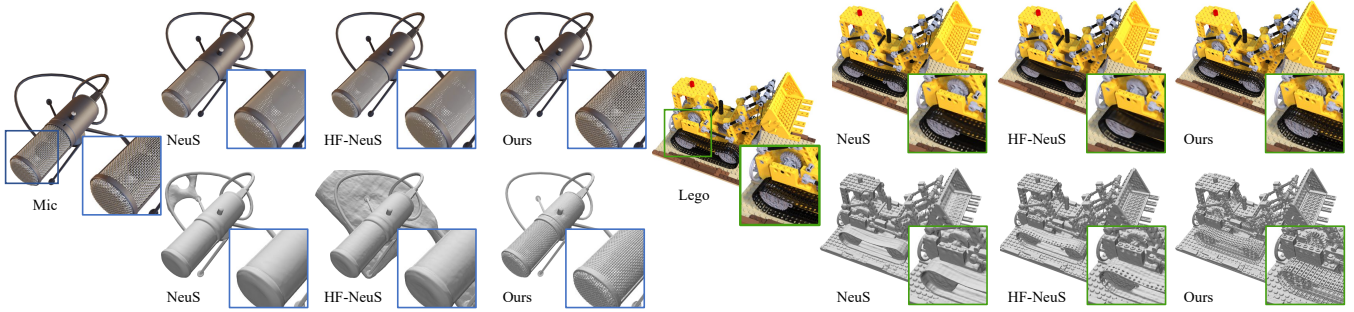


Fig. 4. Comparison of the novel-view synthesis and the reproduced mesh. We derive these detailed meshes with the marching cube of grid 1500. Our meshes show better details on both the microphone's grille and the crawler belt of Lego. Although HF-NeuS attempted to capture the high-frequency details, it struggles with geometries featuring rapid changes, such as the open hole on Lego's band. Simultaneously, the microphone's electric wire is a relatively thin object, both HF-NeuS and NeuS fail to reproduce visually pleasing meshes, as they introduce extra mesh sections connected to other parts and color them to match the background. As a result, while these methods produce appropriate novel views, the inherent geometry contains inaccuracies. Our method excels not only in reproducing details but also in maintaining smoothness on the microphone's handle, thanks to our learnable multi-level encoding method.

	NeuS	TPE	NGP	TPE_CS	Ours
Chamfer Distance	0.298	0.121	0.691	0.116	0.114
PSNR	28.375	29.945	27.977	29.951	30.023

Table 3. Comparison of Chamfer Distance and PSNR in Ablation Study.

	TPE	super-sampling+TPE	Ours
GPU Memory	12G	29G	13G
Chamfer Distance	0.121	0.117	0.114

Table 4. Comparison of GPU memory, and Chamfer distance.

In the experiment NGP, we set  $L = 16$ ,  $F = 2$ ,  $T = 19$ ,  $N_{min} = 16$  as specified in [Müller et al. 2022] and fix  $b = 2$  to define the level growth factor as 2 to provide a fair basis for assessment. We conducted the comparison on the last five cases in the NeRF-synthetic dataset. The quantitative mean values and qualitative results are reported in Tab. 3 and Fig. 5, respectively. In the "NGP" experiment, encoding changes led to performance degradation, likely due to hash collisions from NGP's hashing mapping. This works for NeRF's low-order volume but not for SDF, where spatial points record surface distance. Thus, artifacts emerge in areas with scarce textures or observations. Fig. 5 shows this on the Lego's bucket, as observed in other research [Liang et al. 2023] that employed extra regularization terms to fix it.

#### 5.4 Efficiency of "Cone Discrete Sampling".

Through multi-convolved featurization, our method approximates cone tracing on tri-planes. Compared to 4x super-sampling (evaluating multiple rays through each pixel and averaging them), our method only requires 1/4 of MLP queries, significantly reducing the computational burden. As shown in Tab 4, our method achieves computational efficiency and superior performance. Our method takes around 160s for 1600x1200 image inference and 9h for 300k iteration training on an A100, similar to NeuS but nearly half the time of HF-NeuS.

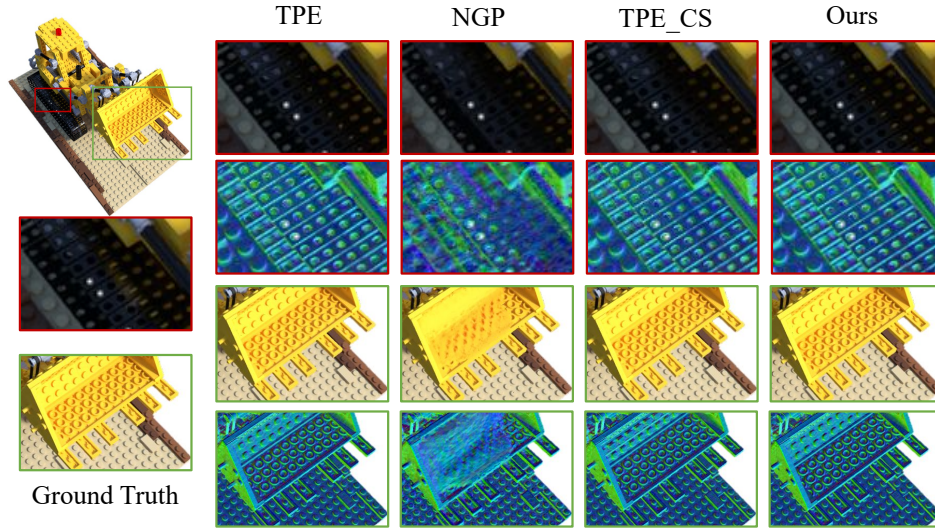


Fig. 5. We visualized the effects of different modules through novel-view synthesis, along with the corresponding surface normals. The 'NGP' encoding introduces undesired structures in regions with scarce textures or observations (e.g., the Lego's bucket). As we progress from 'TPE' and 'TPE\_CS' to our method, we obtain increasingly detailed and clearer results, as exemplified by the distinct hole in the Lego's crawler belt.

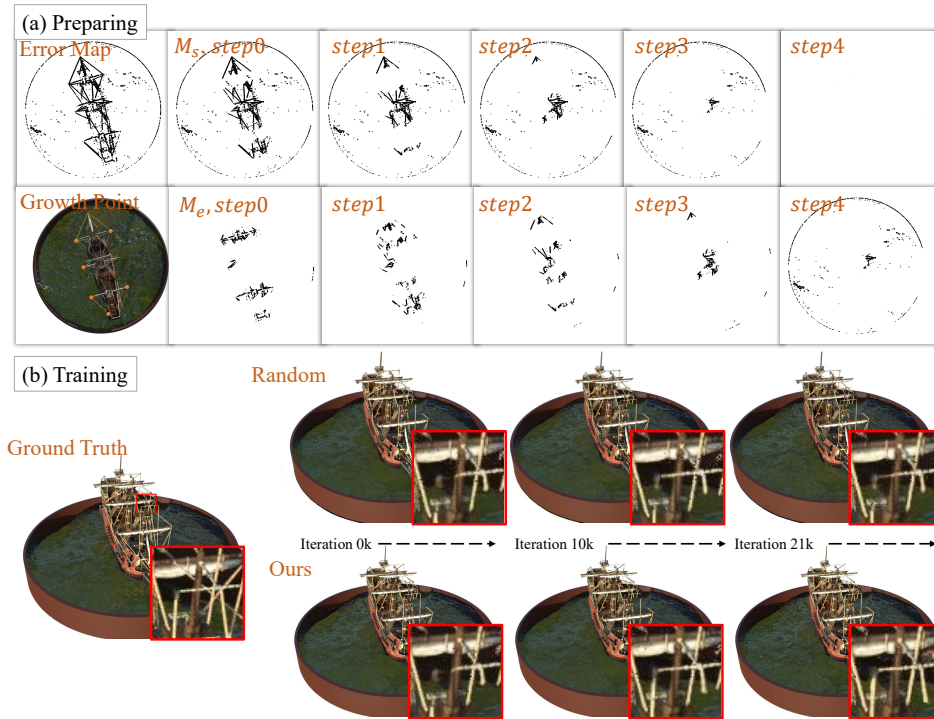


Fig. 6. We compare our SDF growth method ("Ours") with randomly selected rays around the error map ("Random"). In the preparation stage, we render SDF at each training view and compute the error map ("Error Map"), starting the region growth at the detected points ("Growth Point"). We apply the region growth as Alg. 1 with four steps for visualization. During training, the rays of  $M_s$  are selected as the training set. Comparing zoom-in results, our method achieves faster convergence relative to "Random" and helps minimize the effects on regions outside the error map during refinement.



## 5.5 Evaluation of SDF Growth Refinement.

We demonstrate our method that guides the SDF to converge through error-map-based guidance, as shown in Fig. 6. We set the number of steps to  $n = 14$  rained for 1500 iterations at each step, resulting in a total of 21,000 iterations of refinement for a trained model. Our method demonstrates faster convergence than the random sampling approach. Our SDF growth refinement enhances PSNR from 28.360dB to 28.427dB and reduces Chamfer distance from 0.368 to 0.359. This operation only adds a few seconds for initialization, making it highly efficient. Due to computational resource limitations and dataset constraints, we only tested this module on a single, challenging case. However, we believe that this innovative method can be explored and applied to other cases, potentially leading to further improvements in a broader range of scenarios.

## 6 CONCLUSION

In this paper, we present a method, *LoD-NeuS*, that encodes features with a continuous level of detail (LoD) from a novel tri-plane-based representation to adaptively reconstruct high-fidelity geometry. Specifically, we present a multi-scale tri-plane position encoding to capture different LoDs. To effectively represent the high-frequency sampling, we design a multi-convolved featurization to approximate the ray integral within a cone, and then aggregate LoD features from multi-convolved multi-scale features of vertices within a conical frustum along a ray. Besides, for thin surfaces, we develop an SDF growth refinement according to SDF sphere tracing for reconstruction improvement. The state-of-the-art results demonstrate the value of representing a continuous LoD to address aliasing concerns in advanced neural surface reconstruction.

## 7 ACKNOWLEDGE

This work was supported by the National Key Research and Development Program of China under Grant 2022YFF0902201, the National Natural Science Foundation of China under Grants 62001213, 62025108, and the Tencent Rhino-Bird Research Program. We thank the anonymous reviewers for their valuable feedback.

## A OVERVIEW

The supplementary material provides the implementation details (Section A.1) and the derivation of initialization of tri-plane encoding (Section A.2). We have also prepared a video and a reconstructed model for additional visualizations, please see the attachment.

### A.1 Experiment Setting.

**Network Architecture.** We adopt a network architecture similar to NeuS. The geometry network modeling SDF comprises 8 hidden layers with a hidden size of 256, and a skip connection concatenates the input with the output of the fourth layer. The geometry network output includes  $\{1, 256\}$ , representing the predicted SDF and the features for the color network. Sequentially, the color network takes this feature along with the spatial position, view direction and normal to predict the point's color. As in SAL [?], we employ weight normalization to the network to stabilize the training process.

**Training details.** We train our networks with ADAM optimizer. The learning rate is linearly warmed up from 0 to  $5 \times 10^{-4}$  in the

first 5k iterations and then controlled by a cosine decay schedule to reach a minimum learning rate  $2.5 \times 10^{-5}$ . With a ray batch size of 512, we train our network for 300k iterations, taking around 9 hours on a single Nvidia A100 GPU. Our module only introduces slight more computation compared to NeuS. We follow the setting of Hierarchical Sampling as [Wang et al. 2021]. We evaluate our SDF growth module on the "ship" case using a learning rate of  $5 \times 10^{-5}$  for 21k iterations.

### A.2 Geometric Initialization of Multi-scale Tri-plane.

As demonstrated by the previous work [?], a proper initialization is crucial to the training stabilization. However, the recent progress [Liang et al. 2023; Yu et al. 2022] of introducing explicit voxel from NeRF to SDF has rare discussion. The crude initialize the features of the grid using the normal distribution or uniform distribution, which may make the SDF with bad initialization and failed in some textureless area.

**Initialization of Grid representation.** We propose an initialization scheme for grid features  $G$  of dimension  $n$  at position  $\mathbf{x} = \{x, y, z\}$ . The initialized grid features consist of  $\{g_i \sim \mathcal{N}(0, \sigma^2)\}_i^n$ . In this case, the variance  $\sigma^2$  should be defined as  $\sigma^2 = x^2 + y^2 + z^2$ .

**Proof:** According to SAL, an MLP  $f: \mathbb{R}^d \rightarrow \mathbb{R}$  with geometric initialization, can be considered as  $f(\mathbf{x}) \approx \|\mathbf{x}\| - r$ . That is,  $f$  is approximately the signed distance function to a  $d - 1$  sphere of radius  $r$  in  $\mathbb{R}^d$ . Assume the grid encoding as  $\gamma(\mathbf{x}) = \{\{g_i\}_i^n\}$ , with all features concatenated as results. We can further derive that

$$f(\gamma(\mathbf{x})) \approx \|\{g_i\}_i^n\| - r, \quad (10)$$

where  $\|\{g_i\}_i^n\| = \sqrt{\sum_i^n g_i^2}$ . According to the law of large numbers, we get  $\sum_i^n g_i^2 = n\mathbb{E}(g) = n\sigma^2$ . If we want to maintain a spatial sphere in  $\mathbf{x} \in \mathbb{R}^3$  after this encoding, we should set  $n\sigma^2 = \|\mathbf{x}\|$ , sequentially,  $\sigma^2 = \|\mathbf{x}\|/n$ .

**Initialization of tri-plane representation.** Extended this scheme to the case of tri-plane, which consist of three planes reflect the projection to  $\{G_{xy}, G_{yz}, G_{zx}\}$ . So we formulate  $g_i = g_i^{xy} + g_i^{yz} + g_i^{zx}$ , where  $g_i^{xy} \sim \mathcal{N}(0, \sigma_{xy}^2)$ ,  $g_i^{yz} \sim \mathcal{N}(0, \sigma_{yz}^2)$ ,  $g_i^{zx} \sim \mathcal{N}(0, \sigma_{zx}^2)$ , so  $g_i \sim \mathcal{N}(0, \sigma_{xy}^2 + \sigma_{yz}^2 + \sigma_{zx}^2)$ . Then the Equation 10 can be broke down as:

$$\begin{aligned} n(\sigma_{xy}^2 + \sigma_{yz}^2 + \sigma_{zx}^2) &= \|\mathbf{x}\| = x^2 + y^2 + z^2 \\ &= \frac{1}{2}((x^2 + y^2) + (y^2 + z^2) + (x^2 + z^2)). \end{aligned} \quad (11)$$

Therefore, we initialize the three plane as  $\sigma_{xy}^2 = \frac{1}{2n}(x^2 + y^2)$ ,  $\sigma_{yz}^2 = \frac{1}{2n}(y^2 + z^2)$ ,  $\sigma_{zx}^2 = \frac{1}{2n}(z^2 + x^2)$ .

### A.3 Comparison

We compared our method with the state-of-the-art methods (e.g. NeRF [Mildenhall et al. 2020], NeuS [Wang et al. 2021], HF-NeuS [Wang et al. 2022]) on the DTU dataset without mask supervision, as shown in Tab. 5. Our method achieves an average Chamfer distance of 0.72 v.s. 0.77 (from HF-NeuS's paper), demonstrating our method still exhibits superior performance. Both metrics improve compared to those with masks because the peripheries of meshes are masked before evaluation, following the post-processing in HF-NeuS.

	24	37	40	55	63	65	69	83	97	105	106	110	114	118	122	Mean
NeuS	1.37	1.21	0.73	0.40	1.20	0.70	0.72	1.01	1.16	0.82	0.66	1.69	0.39	0.49	0.51	0.87
NeRF	1.90	1.60	1.85	0.58	2.28	1.27	1.47	1.67	2.05	1.07	0.88	2.53	1.06	1.15	0.96	1.49
HF-NeuS	0.76	1.32	0.70	0.39	1.06	0.63	0.63	1.15	1.12	0.80	0.52	1.22	0.33	0.49	0.50	0.77
Ours	0.69	0.88	0.47	0.42	0.85	0.94	0.59	0.80	1.31	0.64	0.61	1.27	0.29	0.64	0.38	0.72

Table 5. Quantitative results on the DTU dataset. Red and orange indicate the first and second best-performing results.

## REFERENCES

- Henrik Aanæs, Rasmus Ramsbøl Jensen, George Vogiatzis, Engin Tola, and Anders Bjorholm Dahl. 2016. Large-scale data for multiple-view stereopsis. *International Journal of Computer Vision* 120 (2016), 153–168.
- Connelly Barnes, Eli Shechtman, Adam Finkelstein, and Dan B Goldman. 2009. Patch-Match: A randomized correspondence algorithm for structural image editing. *ACM Trans. Graph.* 28, 3 (2009), 24.
- Jonathan T Barron, Ben Mildenhall, Dor Verbin, Pratul P Srinivasan, and Peter Hedman. 2022. Mip-nerf 360: Unbounded anti-aliased neural radiance fields. In *CVPR*. 5470–5479.
- Fausto Bernardini, Joshua Mittleman, Holly Rushmeier, Cláudio Silva, and Gabriel Taubin. 1999. The ball-pivoting algorithm for surface reconstruction. *IEEE transactions on visualization and computer graphics* 5, 4 (1999), 349–359.
- Adrian Broadhurst, Tom W Drummond, and Roberto Cipolla. 2001. A probabilistic framework for space carving. In *Proceedings eighth IEEE international conference on computer vision. ICCV 2001*, Vol. 1. IEEE, 388–393.
- Neill DF Campbell, George Vogiatzis, Carlos Hernández, and Roberto Cipolla. 2008. Using multiple hypotheses to improve depth-maps for multi-view stereo. In *Computer Vision—ECCV 2008: 10th European Conference on Computer Vision, Marseille, France, October 12–18, 2008, Proceedings, Part I 10*. Springer, 766–779.
- Rohan Chabra, Jan E Lenssen, Eddy Ilg, Tanner Schmidt, Julian Straub, Steven Lovegrove, and Richard Newcombe. 2020. Deep local shapes: Learning local sdf priors for detailed 3d reconstruction. In *European conference on computer vision*. 608–625.
- Eric R Chan, Connor Z Lin, Matthew A Chan, Koki Nagano, Boxiao Pan, Shalini De Mello, Orazio Gallo, Leonidas J Guibas, Jonathan Tremblay, Sameh Khamis, et al. 2022. Efficient geometry-aware 3D generative adversarial networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 16123–16133.
- Anpei Chen, Zexiang Xu, Andreas Geiger, Jingyi Yu, and Hao Su. 2022a. Tensorf: Tensorial radiance fields. In *Computer Vision—ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part XXXII*. Springer, 333–350.
- Xingyu Chen, Qi Zhang, Xiaoyu Li, Yue Chen, Ying Feng, Xuan Wang, and Jue Wang. 2022b. Hallucinated neural radiance fields in the wild. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 12943–12952.
- Robert L Cook. 1986. Stochastic sampling in computer graphics. *ACM Transactions on Graphics (TOG)* 5, 1 (1986), 51–72.
- Angela Dai, Matthias Nießner, Michael Zollhöfer, Shahram Izadi, and Christian Theobalt. 2017. Bundlefusion: Real-time globally consistent 3d reconstruction using on-the-fly surface reintegration. *ACM Transactions on Graphics (ToG)* 36, 4 (2017), 1.
- François Darmon, Bénédicte Bascle, Jean-Clément Devaux, Pascal Monasse, and Mathieu Aubry. 2022. Improving neural implicit surfaces geometry with patch warping. In *CVPR*. 6260–6269.
- Jeremy S De Bonet and Paul Viola. 1999. Poxels: Probabilistic voxelized volume reconstruction. In *Proceedings of International Conference on Computer Vision (ICCV)*, Vol. 2. 3.
- Qiancheng Fu, Qingshan Xu, Yew Soon Ong, and Wenbing Tao. 2022. Geo-neus: Geometry-consistent neural implicit surfaces learning for multi-view reconstruction. *Advances in Neural Information Processing Systems* 35 (2022), 3403–3416.
- Yasutaka Furukawa and Jean Ponce. 2009. Accurate, dense, and robust multiview stereopsis. *IEEE transactions on pattern analysis and machine intelligence* 32, 8 (2009), 1362–1376.
- Kyle Genova, Forrester Cole, Avneesh Sud, Aaron Sarna, and Thomas Funkhouser. 2020. Local deep implicit functions for 3d shape. In *CVPR*. 4857–4866.
- Amos Gropp, Lior Yariv, Niv Haim, Matan Atzmon, and Yaron Lipman. 2020. Implicit geometric regularization for learning shapes. In *Proceedings of the 37th International Conference on Machine Learning*. 3789–3799.
- Xian-Feng Han, Hamid Laga, and Mohammed Bennamoun. 2019. Image-based 3D object reconstruction: State-of-the-art and trends in the deep learning era. *IEEE T-PAMI* 43, 5 (2019), 1578–1604.
- Hugues Hoppe, Tony DeRose, Tom Duchamp, John McDonald, and Werner Stuetzle. 1992. Surface reconstruction from unorganized points. In *Proceedings of the 19th annual conference on computer graphics and interactive techniques*. 71–78.
- Xin Huang, Qi Zhang, Ying Feng, Hongdong Li, and Qing Wang. 2023. Inverting the Imaging Process by Learning an Implicit Camera Model. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 21456–21465.
- Xin Huang, Qi Zhang, Ying Feng, Hongdong Li, Xuan Wang, and Qing Wang. 2022. Hdr-nerf: High dynamic range neural radiance fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 18398–18408.
- Shahram Izadi, David Kim, Otmar Hilliges, David Molyneux, Richard Newcombe, Pushmeet Kohli, Jamie Shotton, Steve Hodges, Dustin Freeman, Andrew Davison, et al. 2011. 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*. 559–568.
- Michael Kazhdan, Matthew Bolitho, and Hugues Hoppe. 2006. Poisson surface reconstruction. In *Proceedings of the fourth Eurographics symposium on Geometry processing*, Vol. 7. 0.
- Kiriakos N Kutulakos and Steven M Seitz. 2000. A theory of shape by space carving. *International journal of computer vision* 38 (2000), 199–218.
- Patrick Labatut, Jean-Philippe Pons, and Renaud Keriven. 2007. Efficient multi-view reconstruction of large-scale scenes using interest points, delaunay triangulation and graph cuts. In *2007 IEEE 11th international conference on computer vision*. IEEE, 1–8.
- Erich Liang, Kenan Deng, Xi Zhang, and Chun-Kai Wang. 2023. HR-NeuS: Recovering High-Frequency Surface Geometry via Neural Implicit Surfaces. arXiv:2302.06793 [cs.CV]
- Li Ma, Xiaoyu Li, Jing Liao, Qi Zhang, Xuan Wang, Jue Wang, and Pedro V Sander. 2022. Deblur-nerf: Neural radiance fields from blurry images. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 12861–12870.
- Nelson Max. 1995. Optical models for direct volume rendering. *IEEE Transactions on Visualization and Computer Graphics* 1, 2 (1995), 99–108.
- Paul Merrell, Amir Akbarzadeh, Liang Wang, Philippos Mordohai, Jan-Michael Frahm, Ruigang Yang, David Nistér, and Marc Pollefeys. 2007. Real-time visibility-based fusion of depth maps. In *2007 IEEE 11th International Conference on Computer Vision*. Ieee, 1–8.
- Lars Mescheder, Michael Oechsle, Michael Niemeyer, Sebastian Nowozin, and Andreas Geiger. 2019. Occupancy networks: Learning 3d reconstruction in function space. In *CVPR*. 4460–4470.
- B Mildenhall, PP Srinivasan, M Tancik, JT Barron, R Ramamoorthi, and R Ng. 2020. Nerf: Representing scenes as neural radiance fields for view synthesis. In *European conference on computer vision*.
- Thomas Müller, Alex Evans, Christoph Schied, and Alexander Keller. 2022. Instant neural graphics primitives with a multiresolution hash encoding. *ACM Transactions on Graphics (ToG)* 41, 4 (2022), 1–15.
- Matthias Nießner, Michael Zollhöfer, Shahram Izadi, and Marc Stamminger. 2013. Real-time 3D reconstruction at scale using voxel hashing. *ACM Transactions on Graphics (ToG)* 32, 6 (2013), 1–11.
- Michael Oechsle, Songyou Peng, and Andreas Geiger. 2021. Unisurf: Unifying neural implicit surfaces and radiance fields for multi-view reconstruction. In *ICCV*. 5589–5599.
- Johannes L Schonberger and Jan-Michael Frahm. 2016. Structure-from-motion revisited. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 4104–4113.
- Johannes L Schönberger, Enliang Zheng, Jan-Michael Frahm, and Marc Pollefeys. 2016. Pixelwise view selection for unstructured multi-view stereo. In *Computer Vision—ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11–14, 2016, Proceedings, Part III 14*. Springer, 501–518.
- Steven M Seitz and Charles R Dyer. 1999. Photorealistic scene reconstruction by voxel coloring. *International journal of computer vision* 35 (1999), 151–173.
- Vincent Sitzmann, Justus Thies, Felix Heide, Matthias Nießner, Gordon Wetzstein, and Michael Zollhofer. 2019. Deepvoxels: Learning persistent 3d feature embeddings. In *CVPR*. 2437–2446.
- Pratul P Srinivasan, Boyang Deng, Xiuming Zhang, Matthew Tancik, Ben Mildenhall, and Jonathan T Barron. 2021. Nerv: Neural reflectance and visibility fields for relighting and view synthesis. In *CVPR*. 7495–7504.
- Towaki Takikawa, Joey Litalien, Kangxue Yin, Karsten Kreis, Charles Loop, Derek Nowrouzezahrai, Alec Jacobson, Morgan McGuire, and Sanja Fidler. 2021. Neural geometric level of detail: Real-time rendering with implicit 3D shapes. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 11358–11367.

- Ayush Tewari, Justus Thies, Ben Mildenhall, Pratul Srinivasan, Edgar Tretschk, Wang Yifan, Christoph Lassner, Vincent Sitzmann, Ricardo Martin-Brualla, Stephen Lombardi, et al. 2022. Advances in neural rendering. In *Computer Graphics Forum*, Vol. 41. Wiley Online Library, 703–735.
- Engin Tola, Christoph Strecha, and Pascal Fua. 2012. Efficient large-scale multi-view stereo for ultra high-resolution image sets. *Machine Vision and Applications* 23 (2012), 903–920.
- Dor Verbin, Peter Hedman, Ben Mildenhall, Todd Zickler, Jonathan T Barron, and Pratul P Srinivasan. 2022. Ref-nerf: Structured view-dependent appearance for neural radiance fields. In *CVPR*. IEEE, 5481–5490.
- Peng Wang, Lingjie Liu, Yuan Liu, Christian Theobalt, Taku Komura, and Wenping Wang. 2021. NeuS: Learning Neural Implicit Surfaces by Volume Rendering for Multi-view Reconstruction. *Advances in Neural Information Processing Systems* 34 (2021), 27171–27183.
- Yiqun Wang, Ivan Skorokhodov, and Peter Wonka. 2022. Hf-neus: Improved surface reconstruction using high-frequency details. *Advances in Neural Information Processing Systems* 35 (2022), 1966–1978.
- Yiqun Wang, Ivan Skorokhodov, and Peter Wonka. 2023. PET-NeuS: Positional Encoding Triplanes for Neural Surfaces. (2023).
- Menghua Wu, Hao Zhu, Linjia Huang, Yiyu Zhuang, Yuanxun Lu, and Xun Cao. 2023. High-fidelity 3D Face Generation from Natural Language Descriptions. In *Proceedings of IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Huang Xin, Zhang Qi, Feng Ying, Li Xiaoyu, Wang Xuan, and Wang Qing. 2023. Local Implicit Ray Function for Generalizable Radiance Field Representation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*.
- Lior Yariv, Jiatao Gu, Yoni Kasten, and Yaron Lipman. 2021. Volume rendering of neural implicit surfaces. *Advances in Neural Information Processing Systems* 34 (2021), 4805–4815.
- Lior Yariv, Yoni Kasten, Dror Moran, Meirav Galun, Matan Atzmon, Basri Ronen, and Yaron Lipman. 2020. Multiview neural surface reconstruction by disentangling geometry and appearance. *Advances in Neural Information Processing Systems* 33 (2020), 2492–2502.
- Zehao Yu, Songyou Peng, Michael Niemeyer, Torsten Sattler, and Andreas Geiger. 2022. Monosdf: Exploring monocular geometric cues for neural implicit surface reconstruction. *arXiv preprint arXiv:2206.00665* (2022).
- Christopher Zach, Thomas Pock, and Horst Bischof. 2007. A globally optimal algorithm for robust tv-l1 range image integration. In *2007 IEEE 11th International Conference on Computer Vision*. IEEE, 1–8.
- Yichao Zhou, Haozhi Qi, and Yi Ma. 2019. End-to-end wireframe parsing. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 962–971.
- Junyu Zhu, Hao Zhu, Qi Zhang, Fang Zhu, Zhan Ma, and Xun Cao. 2023. Pyramid NeRF: Frequency Guided Fast Radiance Field Optimization. *International Journal of Computer Vision* (2023), 1–16.
- Yiyu Zhuang, Qi Zhang, Xuan Wang, Hao Zhu, Ying Feng, Xiaoyu Li, Ying Shan, and Xun Cao. 2023. NeAI: A Pre-convoluted Representation for Plug-and-Play Neural Ambient Illumination. *arXiv preprint arXiv:2304.08757* (2023).
- Yiyu Zhuang, Hao Zhu, Xusen Sun, and Xun Cao. 2022. Mofanerf: Morphable facial neural radiance field. In *European conference on computer vision*.