# NeuSD: Surface Completion with Multi-View Text-to-Image Diffusion

Savva Ignatyev[1*]      Daniil Selikhanovych[1*]      Oleg Voynov[1,2]      Yiqun Wang[3]
Peter Wonka[4]              Stamatios Lefkimmiatis[5]              Evgeny Burnaev[1,2]

[1]Skoltech, Russia          [2]AIRI, Russia          [3]Chongqing University, China
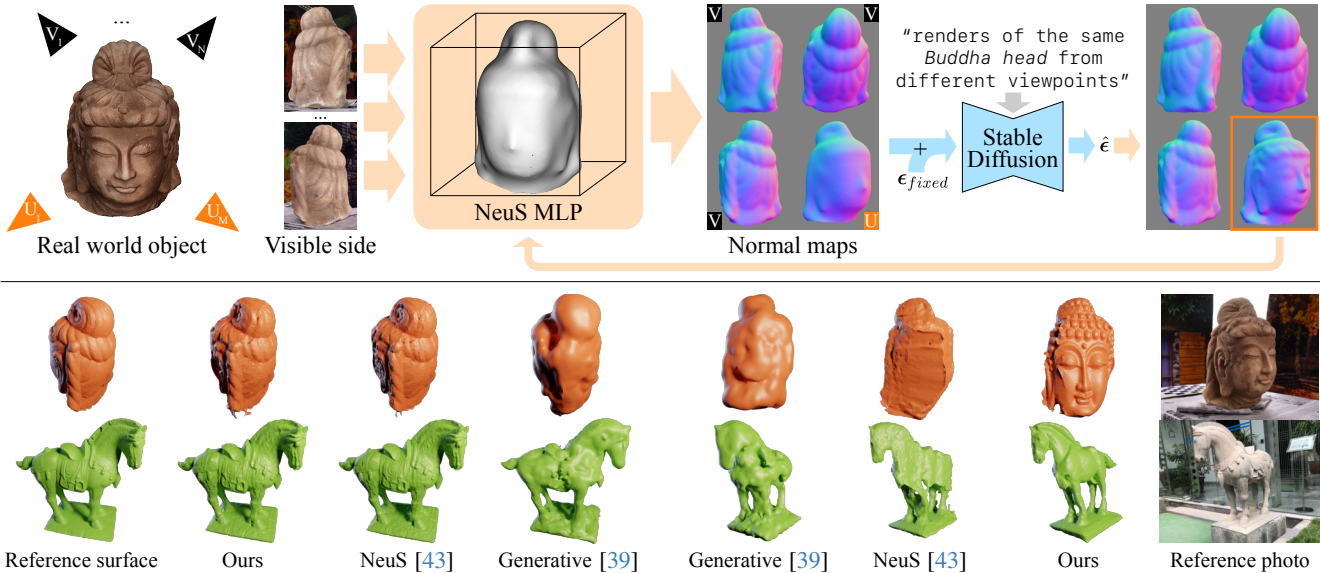[4]KAUST, Saudi Arabia          [5]AI Foundation and Algorithm Lab, Russia

Figure 1. We consider reconstruction of a partially observed surface. We train a neural implicit surface representation using the photos of the observed part while guiding the training using a 2D diffusion model. Our key insights are applying the diffusion model to normal maps instead of color renderings (Sec. 3.3), freezing the noise $\epsilon_{fixed}$ added to the input of the model (Sec. 3.4), and applying the model to a grid of normal maps obtained from visible (V) and unobserved (U) sides of the surface (Sec. 3.5). Our method accurately reconstructs the surface from the observed side (bottom left) and completes it from the unobserved side in a plausible manner (bottom right).

## Abstract

*We present a novel method for 3D surface reconstruction from multiple images where only a part of the object of interest is captured. Our approach builds on two recent developments: surface reconstruction using neural radiance fields for the reconstruction of the visible parts of the surface, and guidance of pre-trained 2D diffusion models in the form of Score Distillation Sampling (SDS) to complete the shape in unobserved regions in a plausible manner. We introduce three components. First, we suggest employing normal maps as a pure geometric representation for SDS instead of color renderings which are entangled with the appearance information. Second, we introduce the freezing of the SDS noise during training which results in more coherent gradients and better convergence. Third, we propose Multi-View SDS as a way to condition the generation of the non-observable part of the surface without fine-tuning or making changes to the underlying 2D Stable Diffusion model. We evaluate our approach on the BlendedMVS dataset demonstrating significant qualitative and quantitative improvements over competing methods.*

## 1. Introduction

We propose a novel framework called *NeuSD* that addresses the problem of 3D surface reconstruction from a set of the posed RGB input images. Specifically, we consider the case where the space of camera locations of the input images around an object is not evenly sampled. In practice, there are many realistic scenarios where a human observer can

---

[*]Equal contribution
For correspondence: o.voynov@skoltech.ru

1

nicely capture images from a set of viewpoints from one side of the object, but not from the other side. For example, imagine a statue in a museum placed right next to a wall or behind a fence or a car being parked close to a wall. This important problem has not been systematically studied in the literature and we set out to propose a solution to this challenge. We will first put this problem in the context of the existing literature.

The current state of the art in surface reconstruction uses neural implicit surfaces. Notably, [43] introduced NeuS, which represents a surface as a signed distance field. NeuS demonstrated reconstruction results on par with the classical Multi-View Stereo (MVS) methods, such as [36, 37], and more importantly, outputs watertight meshes that are ready to be used in other vision or graphics applications.

Both traditional and deep learning approaches to 3D reconstruction usually operate under the assumption of full observability. Nevertheless, there are two strands of literature that are related to the problem of partially-observed input shapes. One is shape completion. However, there are two reasons why this is typically viewed as a separate problem and solutions cannot easily be transferred to our problem setting. First, the majority of state-of-the-art approaches are trained on category-based datasets like ShapeNet [10], which limits their applicability. Second, these methods use point clouds as input instead of a set of images. We, therefore, do not consider shape completion methods as direct competitors. The other important line of work started with the seminal paper DreamFusion by [32]. They suggested a new mechanism, Score Distillation Sampling (SDS), which allowed us to employ pretrained 2D text-to-image diffusion models as "critics" that provide the gradient to learn a 3D neural field. Since then, a substantial amount of work on this topic emerged to advance generative methods for 3D shape synthesis working either with a text prompt [62] or a single input image (and a text prompt) [25, 28]. We extend these excellent methods with two important new capabilities. First, existing methods are not specialized in surface reconstruction and they use a volumetric instead of a surface representation to characterize the resulting 3D shape. We provide a solution that can output clean surfaces in the form of distance fields. Second, the single and few-shot methods do not generalize well to additional viewpoints and are often inherently limited in the viewpoints they can consider, requiring the view-dependent prompt or fixed w.r.t. each other camera positions. By contrast, we propose a method that generalizes better to many input images from relatively closeby viewpoints, and we demonstrate this capability in our experimental validation.

Our proposed solution starts from NeuS as our baseline and complements it with SDS guidance using the Stable Diffusion model [35]. We further develop our method with the following three contributions: (1) SDS Surface Shape Guidance: instead of using the RGB rendering for SDS guidance we use the normal map, which contains only the geometric information disentangled from the appearance; (2) Frozen SDS: for faster convergence and improved quality, we fix the noise map used in SDS; (3) Multi-View SDS: to condition the diffusion model on the visible part of the object we concatenate multiple normal maps in a grid for a joint SDS evaluation.

We test our approach on the BlendedMVS dataset [51] and show that our method is preferred by CLIP scores and a user study (users preferred our method $> 85\%$) compared to generative and reconstructive baselines.

## 2. Related work

**Implicit surface reconstruction.** Since its introduction, Neural Radiance Fields (NeRF) [29] achieved impressive photorealistic results in novel view synthesis. However, NeRF is a poor representation for multi-view surface reconstruction, since its density field is highly ambiguous and ill-regularized. To adapt NeRF to surface reconstruction, several works [43, 52, 53] proposed Neural Implicit Surfaces. NeuS [43] specifically, implicitly parameterizes the NeRF density as a rapidly decreasing function of the signed distance to the surface (SDF), and maintains the consistency of the SDF using an Eikonal penalty [19]. Subsequent works proposed constraining the implicit surface using the points obtained with Structure-from-Motion [17], improving multi-view consistency via a warping-based loss term [14], improving performance in featureless areas via a normal-based regularization [42, 44], or improving accuracy and convergence speed by using multi-resolution hash encoding [45] or hybrid representations of the surface [15]. While some works extended NeuS or NeRF to sparse image sampling [27, 54], still the vast majority of the work is concerned with the reconstruction of the visible part of the surface only, and is not applicable as is in our setting with a partially observed surface.

**2D diffusion for 3D generation.** Diffusion generative networks started as a fascinating concept and gradually made their way to the state of art through the recent years. Their idea for inference is to start from the Gaussian noise and gradually remove noise with the help of the *denoising network* step-by-step in the so-called *reverse diffusion process*. The training samples for the denoising network are drawn from the *forward diffusion process*, where real photos are taken and the noise is added to them. The denoising network can be conditioned on additional input such as 2D segmentation mask [58], depth map, or text prompt. Latent diffusion [35] was the first work that achieved generating very high-quality 2D images from a general-case text prompt-conditioned generative image model. Later, Stable Diffusion extended this work and went viral by producing near-photorealistic samples from user requests.

Despite its undeniable success and some signs that Stable Diffusion (SD) operates the internal 3D-like representation [11], obtaining the 3D model from SD is a non-trivial task. DreamFusion [32] introduced the Score Distillation Sampling (SDS) approach, which allows propagation of the learning signal from a pre-trained 2D diffusion model into the 3D implicit neural representation model. Score Jacobian Chaining [41] provided a theoretical framework for practically the same algorithm. Multiple works have further improved the text-conditioned 3D generation. Magic 3D [24] introduced a two-stage approach, first optimizing the low-resolution radiance field, and then switching to high-resolution mesh-texture representation. HIFA [62] suggested gradually decreasing SDS noise and making multiple consecutive denoising steps before the backpropagation.

Recently, several works adapted 2D diffusion models for 3D generation conditioned on a single image, *i.e.*, generative single-view 3D reconstruction. RealFusion [28] utilized textual inversion [18] to condition a pre-trained text-to-image diffusion model on the input image, and then train a NeRF using SDS. One-2-3-45 [25] used a special multi-view diffusion model conditioned on the input image [26] to synthesize images from different viewpoints and then train a neural surface representation on them. DreamGaussian [39] utilized a similar approach but trained a surface represented with Gaussian Splatting [22] and used SDS instead of the synthesized images directly. To the best of our knowledge, there are no works on generative *multi-view* 3D reconstruction that address the same problem as ours, so we compare with the three works described above. A different work, SparseFusion [61], utilizes a diffusion model for few-view 3D reconstruction, but it is only applicable to class-specific data.

**Shape completion** is a vast area of research. We only discuss recent shape completion methods that leverage generative modeling techniques, and refer the reader to a survey [16] for a broader review of the topic. One approach is to use 3D generative models conditioned on the partial shape, either using GANs [9, 47, 57], transformer [48, 49, 55], or diffusion [12, 13, 60]. Typically, these methods are trained on smaller datasets like ShapeNet [10] and on synthetically generated inputs. These two limitations preclude a comparison in our setting, as these methods cannot handle our inputs. A concurrent work [21] shows nice initial results using SDS. In an informal comparison, we observe that our method yields higher visual quality, but a systematic comparison is only feasible once the code becomes available.

## 3. Method

We start with the description of the core ideas of NeuS [43] and Score Distillation Sampling [32] that our method builds on and then describe the novel components of our method.

### 3.1. NeuS

We choose NeuS [43] as the starting point and the baseline for our work. NeuS is a volumetric differentiable rendering method that parametrizes the surface using two deep neural networks $f : \mathbb{R}^3 \to \mathbb{R}$ and $c : \mathbb{R}^3 \times \mathbb{S}^2 \to \mathbb{R}^3$ that represent the geometry as signed distance function (SDF) and the appearance as radiance field respectively. The surface is defined as the zero-level set of the SDF network. The color of a pixel rendered along the ray with direction $\boldsymbol{v}$ starting at the point $\boldsymbol{o}$ is given by

$$C(\boldsymbol{o}, \boldsymbol{v}) = \int_0^{+\infty} W(l|f)c(\boldsymbol{p}(l), \boldsymbol{v})\mathrm{d}l, \qquad (1)$$

where $\boldsymbol{p}(l) = \boldsymbol{o} + l\boldsymbol{v}$, $l > 0$ is a point on the ray, and $W(l|f)$ is a weight function that depends on the SDF network.

To fit the model to a set of images, one optimizes the photometric term $\mathcal{L}_{\mathrm{c}} = \sum_k \|C(\boldsymbol{o}_k, \boldsymbol{v}_k) - C_k\|_1 / K$ over the set of observed pixel values $\{C_k\}$ and the respective camera positions $\boldsymbol{o}_k$ and view directions $\boldsymbol{v}_k$. The loss function of NeuS additionally includes the mask guidance term $\mathcal{L}_{\mathrm{m}}$, and the Eikonal penalty $\mathcal{L}_{\mathrm{eik}}$, intended to assure that the SDF network $f$ represents a valid signed distance function. We refer the reader to [43] for the definition of these terms.

### 3.2. Score Distillation Sampling

DreamFusion [32] introduced the concept of score distillation sampling (SDS) for training neural radiance fields (NeRF) using a pre-trained 2D diffusion model, without any real-world input images. The idea behind SDS is to pass the image rendered from the NeRF during training with added Gaussian noise through the denoising diffusion model conditioned on a textual description of the scene, and to use the output of the model to guide the NeRF towards a representation that better corresponds to the description.

For a 3D scene parameterized by $\theta$, the rendered image $\boldsymbol{x}(\theta)$, and the added Gaussian noise $\boldsymbol{\epsilon}$, the denoising model $\epsilon_\phi$ produces the output $\hat{\boldsymbol{\epsilon}} = \epsilon_\phi(\boldsymbol{y}, t, \alpha_t \boldsymbol{x} + \sigma_t \boldsymbol{\epsilon})$, that is then used to guide the NeRF with the gradient of the virtual SDS loss term

$$\nabla_\theta \mathcal{L}_{\mathrm{SDS}} = \mathbb{E}_{t, \boldsymbol{\epsilon}} \left[ w(t)(\hat{\boldsymbol{\epsilon}} - \boldsymbol{\epsilon}) \frac{\partial \boldsymbol{x}}{\partial \theta} \right], \qquad (2)$$

where $\boldsymbol{y}$ is the embedding of the text prompt, $t \sim \mathcal{U}(0, 1)$ is the diffusion timestep, and $\alpha_t$, $\sigma_t$, and $w(t)$ are the weighting factors that depend on the timestep. Please refer to [32] for their definition.

**SDS for surface completion.** We employ SDS to guide the formation of the unobserved part of the NeuS surface fitted to a set of images, the viewpoints for which we denote as *visible*. For this, we define a set of viewpoints directed to the unobserved side of the surface and apply the virtual SDS

loss term for these viewpoints, obtaining the rendered image $x$ through Eq. (1). The loss function for training NeuS with SDS guidance is given by

$$\mathcal{L}_{\text{NeuS+SDS}} = \mathcal{L}_{\text{c}} + \beta\mathcal{L}_{\text{m}} + \lambda\mathcal{L}_{\text{eik}} + \gamma\mathcal{L}_{\text{SDS}}. \quad (3)$$

For the SDS implementation, we follow [32] but use a latent diffusion model, Stable Diffusion 2.1, so in our formulation of the SDS loss term the rendered image $x$ in Eq. (2) is replaced with the respective latent code.

### 3.3. SDS Surface Shape Guidance

Typically SDS is used in conjunction with color renderings. In addition, it can be used to guide other types of renderings, *e.g.*, shaded albedo and textureless renderings to improve geometric details and avoid degenerate flat solutions. Based on our experiments, we argue that shading is a mere projection of the surface normals leading to information loss. Instead, we suggest using normal maps directly, which are easy to obtain from neural surfaces by rendering the SDF derivative similarly to Eq. (1), via

$$N(\boldsymbol{o}, \boldsymbol{v}) = \int_0^{+\infty} W(l|f)\nabla f(\boldsymbol{p}(l))\mathrm{d}l. \quad (4)$$

The normal map is then used instead of the rendered color image $x$ in Eq. (2) (or rather its respective latent code) to compute the normal-based SDS loss term $\mathcal{L}_{\text{SDS,N}}$.

A normal map with its components mapped to RGB, as in Fig. 1, could be viewed as a textureless render under special lighting conditions (not taking self-occlusion into account) with three distant colored light sources: red, green, and blue. Thus, normal maps should lie inside the generative manifold of the diffusion model with enough generalization capacity and be a valid input for the SDS algorithm.

The color renderings, despite their redundancy, may still contain important information that is not present in the normal maps, so we propose to use both the color- and normal-based SDS for training:

$$\begin{aligned}\mathcal{L}_{\text{+Normals}} = &\mathcal{L}_{\text{c}} + \beta\mathcal{L}_{\text{m}} + \lambda\mathcal{L}_{\text{eik}} + \\ &+ \gamma\mathcal{L}_{\text{SDS}} + \gamma_N\mathcal{L}_{\text{SDS,N}}.\end{aligned} \quad (5)$$

Since the diffusion model may have some form of bias and associate particular prompts with particular colors, *e.g.*, "photo camera" with black color, we randomly rotate the normals during training to ensure uniform color distribution of the color-mapped normals and avoid implicit bias.

### 3.4. Frozen SDS

During the initial phases of learning the 3D representation with SDS the surface is not yet formed and the gradient from the denoising network is highly inconsistent, which leads to artifacts and slow convergence. One way to address this issue [62] is to employ multi-step bootstrapping

for SDS to make it more consistent. Instead of computationally expensive multi-step bootstrapping, we propose to fix the additive Gaussian noise $\epsilon$ in Eq. (2) "frozen"

$$\nabla_\theta\mathcal{L}_{\text{SDS,frozen}} = \mathbb{E}_t\left[w(t)(\hat{\epsilon} - \epsilon_{\text{fixed}})\frac{\partial x}{\partial\theta}\right]. \quad (6)$$

Fixing the noise makes the gradient dependent only on the timestep, which drastically reduces the variation of the gradient, and not only leads to faster convergence but also improves the generation results.

In our method, we combine the idea of Eq. (6) with the loss function from Eq. (5).

### 3.5. Multi-view SDS

SDS is known to produce inconsistent updates from different viewpoints, leading to the well-known *Janus problem* [8]. While this particular issue is not central to our work, the observed and the unobserved parts of the surface still need to be reconstructed consistently. Previous works have addressed this issue in different ways. For example, RealFusion [28] fine-tunes the Stable Diffusion (SD) model using textual inversion [18], which is time-consuming and prone to overfitting. SparseFusion [61] trains a view-conditioned diffusion on a category-specific multi-view dataset. Unfortunately, such data is not widely available and the multi-view model lacks generalization capacity compared to the regular 2D diffusion model.

Instead, we propose to condition the diffusion model on multiple views without any tuning, architecture changes, or training of a different model from scratch. We observe that the publicly available SD model, given the prompt of the form *"renders of the same ⟨X⟩ from different viewpoints"*, produces a grid of samples that share distinctive stylistic and geometric similarities, although not perfectly aligned (see examples in the supplementary material). Given this, we propose to apply the SDS loss term to a grid of images compiled of one rendering from the unobserved viewpoint and several renderings from the visible viewpoints, as shown in Fig. 1. This provides the unobserved part of the surface with a guidance signal consistent with the observed parts. We note that we use the parts of the grid rendered from the visible viewpoints only for consistency guidance, and do not propagate the gradient from the SDS loss term to the surface representation through these parts.

We train our complete method using the loss function from Eq. (5) combined with the ideas of frozen and multi-view SDS.

## 4. Experiments

### 4.1. Evaluation setup

**Data.** We evaluate our method using the BlendedMVS [51] semi-synthetic dataset for multi-view stereo reconstruction,

| | | | ariadne | bear | boots | buddha | bull | camera | clock | david | dog | helen | horse | lion | man | plant | santa | shiba | xia | mean |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| CLIP similarity ← | | RealFusion [28] | 73.9 | 69.8 | 79.2 | 69.9 | 73.8 | 77.0 | 68.1 | 74.4 | 70.3 | 73.5 | 78.3 | 77.2 | 76.4 | 75.4 | 69.1 | 80.0 | 72.2 | 74.0 |
| | | One-2-3-45 [25] | 72.1 | 76.4 | 73.4 | 72.1 | 78.4 | 73.3 | 68.0 | 73.9 | 76.1 | 74.3 | 75.5 | 74.8 | 74.6 | 79.4 | 68.8 | 77.1 | 71.9 | 74.1 |
| | | DreamGaussian [39] | 73.1 | 79.9 | 78.6 | 75.7 | 82.1 | 74.1 | 67.2 | _84.2_ | 74.1 | 72.1 | 80.2 | 74.9 | 70.5 | 83.4 | 72.3 | 86.3 | 71.2 | 76.5 |
| | | NeuS [43] | _75.0_ | **87.0** | _84.3_ | _81.0_ | _89.2_ | _81.6_ | _76.4_ | 84.1 | _84.9_ | _81.4_ | _88.7_ | _86.5_ | _80.4_ | _89.4_ | _80.1_ | _87.3_ | _84.4_ | _83.6_ |
| | | Ours | **84.4** | **87.1** | **88.6** | **86.5** | **91.3** | **83.7** | **78.9** | **88.7** | **87.4** | **85.4** | **91.6** | **88.5** | **83.4** | **91.0** | **84.5** | **87.5** | **86.4** | **86.7** |
| Users ↑ | Img | Ours vs NeuS [43] | 100 | 63 | 100 | 98 | 93 | 100 | 87 | 98 | 100 | 100 | 100 | 94 | 16 | 95 | 100 | 84 | 99 | 90 |
| | | Ours vs DG [39] | 100 | 43 | 100 | 96 | 98 | 98 | 100 | 89 | 99 | 99 | 100 | 100 | 98 | 91 | 98 | 64 | 97 | 92 |
| | Surf | Ours vs NeuS [43] | 100 | 57 | 100 | 99 | 91 | 96 | 89 | 98 | 98 | 97 | 98 | 83 | 18 | 88 | 98 | 72 | 99 | 87 |
| | | Ours vs DG [39] | 100 | 35 | 100 | 100 | 97 | 98 | 100 | 94 | 100 | 99 | 100 | 100 | 100 | 91 | 98 | 61 | 99 | 93 |

Table 1. **Quantitative comparison.** Top: CLIP similarity, the **best** and <u>second best</u> results are highlighted. Bottom: results of the user study for the two questions described in the main text. The values represent the percentage of subjects that prefer our method over the respective competitor.

consisting of 117 scenes including small objects, statues, and architecture. Each scene usually contains a single salient object, the data for which consists of the reference 3D surface and the images captured from 30-200 different directions. We chose BlendedMVS because it provides viewpoints all around the object, required for evaluation in our setting, in contrast to other widely-used datasets for multi-view reconstruction like DTU [7, 20].

**Data preprocessing.** We picked 17 scenes from the dataset for our experiments. For each scene, we manually chose the observed side of the object, the viewpoints from the dataset that mostly capture only the observed side, and the guidance viewpoints that mostly capture only the unobserved side, 5-15 viewpoints in each set. To obtain the reference surface for evaluation, we only kept the salient object and manually removed the unrelated environment. We also used these surfaces to obtain image object masks required for the methods.

**Implementation details.** We took the implementation of the volumetric rendering from NeuS [43] and an open-source diffusion model Stable Diffusion 2.1 from Hugging-Face [34, 35]. We manually picked the text prompt for the diffusion model per scene and we refer to these prompts in the supplementary.

At each iteration, we randomly pick a viewpoint from the combined set of visible and unobserved views and compute the regular NeuS loss for the visible viewpoints, and the SDS-guided loss for the unobserved ones. To compute the SDS loss, we render dense images/normal maps at $64 \times 64$ resolution, to reduce the computational cost, and resize them to $512 \times 512$ resolution before feeding them to the diffusion model. We alternate between the single- and multi-view SDS to obtain more comprehensive guidance. For the multi-view SDS, we store the normal maps for the observed part of the object obtained at the last three iterations with visible viewpoints and combine them with the

normal map from the current unobserved viewpoint into a $2 \times 2$ grid. For every 4-th iteration of the single- or multi-view SDS, we perform the regular color-based SDS.

We kept the implementation of NeuS and its hyperparameters unchanged. For the SDS loss, we picked the weight around $10^{-5}$, with the optimal exact value depending on the size of the unobserved part of the object. We used the classifier-free guidance scale of 100 and a timestep uniformly sampled from $[0, 0.5]$ with 1 being the possible maximum. For each scene, we trained the model for 300k iterations, which took 28 hours on a single A100 NVIDIA GPU, 13 of which were spent on Stable Diffusion inference.

**Competitors.** As a representative of purely reconstruction-based methods, we selected NeuS [43] because it constitutes the starting point for our work. Even though there are many follow-up works to NeuS, all of them are expected to behave similarly in reconstructing the unseen parts of the object. We used the official implementation and obtained the results for NeuS using the images and masks for the whole set of visible viewpoints.

As representatives of generative methods we selected RealFusion [28], One-2-3-45 [25], and DreamGaussian [39], that condition a diffusion model on a single input image and utilize this model to guide the surface reconstruction process. We used the official implementation of all methods and obtained the results for the input image that we manually selected from the visible viewpoints so that it provides the most context about the unobserved part of the object. We additionally provided the methods with the ground truth image object mask and made them aware of the ground truth input camera position w.r.t. the object, as we describe in the supplementary material. We initialized the prompt for RealFusion with the same prompt we used for our method.

**Metrics.** We evaluate the quality of the results using CLIP similarity and a perceptual user study. For CLIP similarity,

Figure 2. **Qualitative comparison.** The first row shows the image from the visible set used by the competing single-view methods as input. The next four rows show the results of the three generative single-view competitors: RealFusion (RF), One-2-3-45, and DreamGaussian (DG), and the reconstructive competitor (NeuS). The next two rows show the results of the multi-view baseline (NeuS + SDS) and of our full method. The last row depicts the reference surface. All meshes are rendered from the unobserved side.

we render the reconstructed and the reference surfaces using viewpoints distributed uniformly around the object, compute the cosine similarity between the respective CLIP [33] image embeddings, and report the average value across all views.

For the user study, we conducted an online survey in which each participant was asked to evaluate the results of our method, NeuS [43], and DreamGaussian [39] from the unobserved side using two tasks: picking the reconstruction that corresponds better to an input image, and picking the reconstruction that is more similar to the shown reference

surface. We report the responses of 130 participants.

We also compared the accuracy of the reconstructed surface. Our method accurately reconstructs the observed part of the object, as we show in Fig. 1. The accuracy computed for the whole surface is, as expected, similar for all methods. We report more on this in the supplementary material.

## 4.2. Quantitative and Qualitative Results

We show the quantitative comparison of the methods in Tab. 1 and the qualitative comparison in Fig. 2. We first discuss the quantitative results using CLIP and then the

qualitative results jointly with the user study.

Our method consistently outperforms both the generative competitors and the baseline NeuS w.r.t. CLIP similarity. This demonstrates that we can reconstruct the observed parts of the surface well, but at the same time achieve more plausible generative shape completions for the unobserved parts. We note that we measure CLIP similarity for both the observed and unobserved parts, so the shown numbers are a blend between reconstruction and generation quality. If we use CLIP similarity only for the generated parts, our method has an even bigger advantage compared to NeuS, as we show in the supplementary material.

The user study, evaluating the generation of the unobserved parts only, shows a very strong preference for our method on most scenes. We obtained close to 100 percent preference votes for the majority of the scenes. While a preference score of 100 percent is unusual, one can confirm in Figs. 1 and 2 that the reconstruction quality of the competitors is low for these scenes. There are multiple reasons why competitors struggle to generate plausible results for the unobserved parts. NeuS is not a generative method and generally produces overly smooth shape completions that cannot be judged to be realistic. The generative single-view competitors seem to be strongly dependent on an aligned frontal view. If the single input view does not contain the most salient object parts, the output degenerates. We additionally discuss the challenges of our setting for these methods in the supplementary material.

The scenes that show lower user scores for our method are *bear*, *man*, and *shiba*. The qualitative results shown in Fig. 3 confirm that our reconstructions for these scenes are indeed worse. Our conjecture for this problem is that our prompts "soft toy", "antique death mask", and "asian toy", respectively, are ill-suited for these scenes because the Stable Diffusion 2.1 interprets them differently than intended. For the *bear* it tends to generate the forward side of the regular teddy bear and for the *shiba* it associates the chosen prompt with some plastic toy. We generally opted for using simple and short prompts and chose not to spend a lot of time fine-tuning them for our method. While it could be possible to get better results when trying many different prompts, we also observed in some experiments that seemingly very fitting and elaborate prompts can lead to poor results.

For the *man*, which is the face of a statue for which we picked one half of the face as the observed part, we notice a conflict between reconstruction and generation with the prompt. While both sides of the face should be similar to the visible part, the reconstruction and the generation cannot agree, and the generated information overrides the information from the reconstructed part, rather than the other way around. We find that quite intriguing because for other examples, *e.g.*, *bull*, *helen*, or *santa*, we noticed that sym-
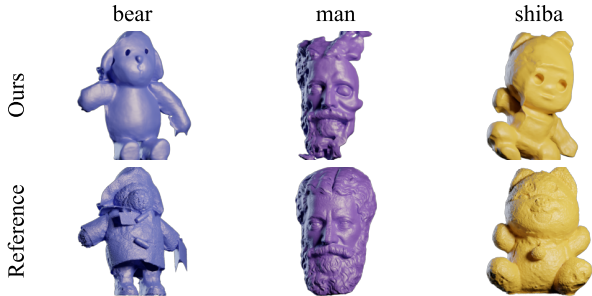


Figure 3. Scenes with lower user scores for our method.

metry is actually helpful for generation, and that the SDS loss helps to propagate information from the visible part to the unobserved part without creating a conflict.

Overall, we observe that single-view methods are limited by the fact that they cannot consider all views from the visible set as input. As there is no existing method available to tackle exactly the same problem as our method, we chose to implement a combination of NeuS with color-based SDS as an additional multi-view baseline. The advantage of this method is that it also serves as an ablation for our work as it is directly comparable. We compare with NeuS + SDS in the ablation study described in the next subsection. We also evaluated the code of multiple other single-view methods, attempting to extend them to the multi-view case. However, this requires a major engineering effort that exceeds a simple change. It would require an adaption of the prompting strategy and prompt generation, viewpoint generation, and searching for new hyperparameters.

### 4.3. Ablation study

For ablation, we start from the baseline NeuS method and add SDS as generative guidance to establish a multi-view baseline. Then, we gradually add components from Sec. 3, namely normal maps, frozen SDS, and the multi-view SDS. We provide some qualitative visual analysis of the ablation study at the bottom of Fig. 2 and in Fig. 4, and show the full comparison in the supplementary material.

We observe that the most significant component of our method is the introduction of normal maps in conjunction with SDS, which greatly improves the quality of the produced surfaces for all scenes compared to using color renderings only. The normal map supervision seems to work better with viewpoints sampled from only a side of the object compared to SDS on color images, and is therefore more suitable to mix reconstruction and generative losses.

The addition of the frozen SDS mainly speeds up the convergence of the method and in some cases slightly improves the results for large unobserved parts, as we show for the *santa*, *plant*, and *shiba*, making them more consistent with the visible part of the surface. We hypothesize that the static noise map decreases the variability in the learning
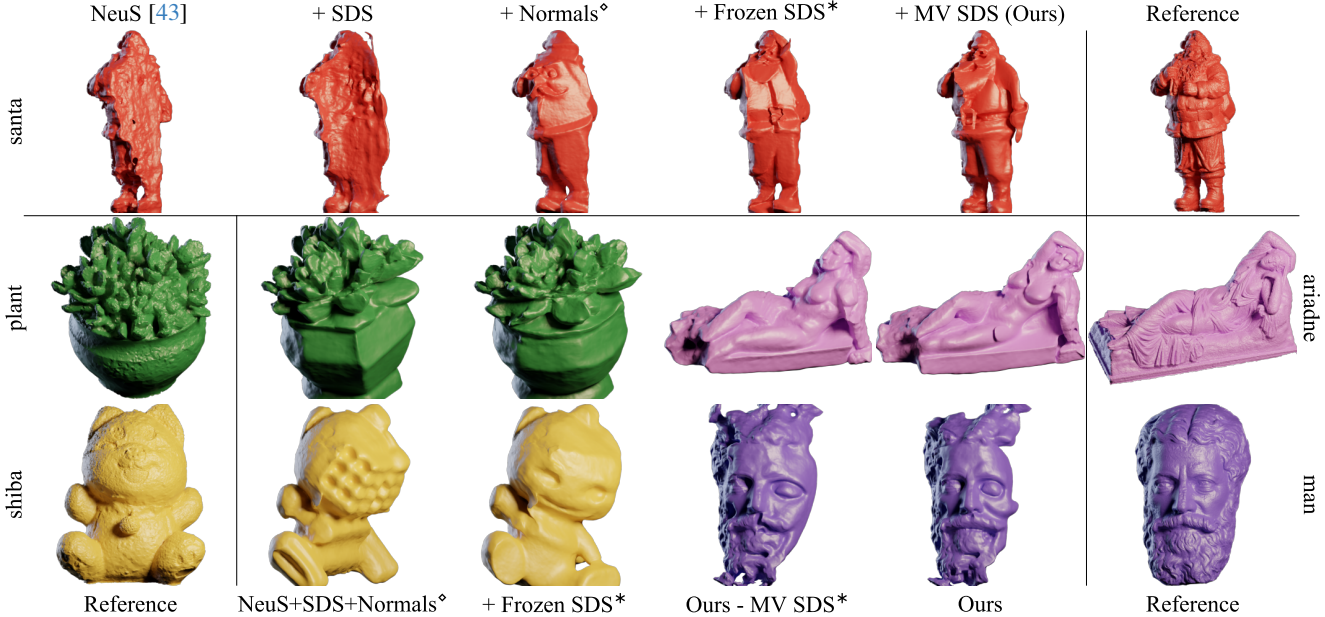
Figure 4. **Ablation study.** The first row shows how our method builds on NeuS by successively adding color image-based SDS guidance, normal maps, frozen SDS, and multi-view SDS guidance. The second and third rows show partial ablation studies for four scenes. We illustrate the effect of freezing the noise on the left, and the effect of multi-view SDS guidance on the right. The superscript markers ⋄ and ∗ mark the same versions of the method.

signal and thereby consistently guides the formation of the unseen part of the surface.

Finally, the multi-view SDS makes the generated part of the surface more consistent with the observed part over the number of scenes. Its effectiveness seems to largely depend on the stylistic similarity of the visible images with the unobserved images. For example, it makes the body pose of the *ariadne* more consistent with the visible images and more similar to the ground truth pose, makes the *man* more symmetric, and for the *santa* greatly improves the overall quality of the surface and makes the figure aligned with the overall pose.

We also tried to evaluate the differences between the different versions of the method in the ablation study using CLIP similarity. The differences are generally too subtle to be picked up by CLIP features and therefore the methods all show similar results within a narrow range. We show these results in the supplementary material.

## 5. Conclusions

We proposed a novel framework, called NeuSD. It mixes traditional NeuS-based 3D shape reconstruction with generative completion. NeuSD is especially useful for cases where part of a shape is observed by multiple images, and another part of the shape is not observed. It creates a plausible completion of the shape that is semantically meaningful and consistent with the observed parts. The main compo-

nents of our approach are: 1) surface diffusion guidance, 2) freezing noise, and 3) multi-view SDS. Our results outperform the current state of the art in qualitative and quantitative metrics. We also would like to discuss two limitations of our work. First, the computation times are still fairly high, compared to some of the fastest known methods that just use a single forward pass through a network to generate a shape. Second, we do not use a separately trained diffusion network to generate multiple-view images like [25]. This possibly limits the quality of the results. However, we do not think that current multi-view datasets are good enough for real-world reconstruction. In future work, we would like mainly to address the inference time problem to investigate faster approaches for mixing surface reconstruction and generative modeling. In addition, we would like to extend our framework to complete scenes containing multiple objects.

# NeuSD: Surface Completion with Multi-View Text-to-Image Diffusion

## Supplementary Material

In Sec. 6 we describe our test data in more detail. In Sec. 7 we provide additional details of selection of textual prompts for the text-to-image diffusion model used in our method, selection of the timestep for the diffusion model, and testing details for the competing methods. In Sec. 8 we describe our quantitative evaluation in detail. In Secs. 9 and 10 we provide the complete comparison of our method with the competing methods, and discuss the challenges of our setting for the competing methods. In Sec. 11 we show the complete set of qualitative and quantitative results of our ablation study. In Secs. 12 and 13 we additionally discuss the premise of multi-view SDS and theoretical aspects of frozen SDS.

## 6. Data selection

We evaluated our method on 17 scenes from the Blended-MVS [51] dataset. For each scene, we used two sets of viewpoints: the *visible* viewpoints with the respective input images that capture the observed part of the surface, and the guidance viewpoints that our method uses to apply Score Distillation Sampling (SDS) [32]. We show the images for the visible viewpoints in Figs. 10 to 12. We selected both sets of viewpoints from all viewpoints in the dataset manually using the following algorithm.

First, we chose the unobserved side for each scene. In general, we chose the side of the surface that is intuitively harder to recover given only the information about the opposing side, such as the face of a statue or a figurine. To increase the diversity of the test data and to evaluate the ability of our method to exploit bilateral symmetry, we also chose a lateral unobserved side for some scenes, specifically, for *bull*, *camera*, *helen*, *horse*, *lion*, and *man*. Next, we picked 5–15 visible viewpoints that on the one hand capture the chosen unobserved side as little as possible, and on the other hand provide a sufficient parallax for multi-view stereo reconstruction of the visible side of the surface. Finally, we picked 5–15 guidance viewpoints based on similar considerations for swapped sides of the surface, *i.e.*, directed to the unobserved side.

We emphasize that we picked the data for our experiments solely based on our intuition about the difficulty of completion of the unobserved side, and not based on the performance of our method. We picked the viewpoints once and kept them fixed in all experiments. We note that the described algorithm could be formally implemented in software, similarly to view selection used for multi-view surface reconstruction (*e.g.*, [37] or [50], Section 4.1).

| Scene name | Prompt base | Project ID |
|---|---|---|
| ariadne | reclining statue | 59f87d0bfa6280566fb38c9a |
| bear | soft toy | 5bf3a82cd439231948877aed |
| boots | boots | 5c34529873a8df509ae57b58 |
| buddha | buddha head | 5ab85f1dac4291329b17cb50 |
| bull | running bull figurine | 5b22269758e2823a67a3bd03 |
| camera | photo camera | 5c34300a73a8df509add216d |
| clock | cuckoo clock | 5a969eea91dfc339a9a3ad2c |
| david | statue | 59ecfd02e225f6492d20fcc9 |
| dog | dog figurine | 5c1af2e2bee9a723c963d019 |
| helen | statue head | 59f363a8b45be22330016cad |
| horse | horse statue | 5b2c67b5e0878c381608b8d8 |
| lion | lion statue | 5b908d3dc6ab78485f3d24a9 |
| man | antique death mask | 5bf7d63575c26f32dbf7413b |
| plant | succulent pot | 5bfd0f32ec61ca1dd69dc77b |
| santa | santa claus | 5be47bf9b18881428d8fbc1d |
| shiba | asian toy | 5c1892f726173c3a09ea9aeb |
| xia | chinese statue | 5ab8713ba3799a1d138bd69a |

Table 2. **Textual prompts** that we used for the scenes in our experiments and their respective project IDs in the BlendedMVS dataset.

## 7. Implementation details

**Prompts.** To apply SDS guidance in our method we used Stable Diffusion 2.1 conditioned on a textual prompt. We manually picked simple and short prompts based on the input images for each scene: we show these prompts in Tab. 2. We used them as is for the NeuS+SDS baseline and as the initialization for textual inversion in RealFusion [28], and for our complete method with multi-view SDS we substituted these prompt bases with "renders of the same ( *prompt base* ) from different viewpoints".

**Timestep interval for SDS.** The effect of SDS depends on the timestep parameter $t$ of the underlying denoising diffusion model. This parameter controls the magnitude of the added noise and the scale of the guidance updates: greater values lead to lower-frequency updates and smaller values lead to higher-frequency updates [62]. Originally, Dream-Fusion [32] suggested to sample $t$ from the $[0, 1]$ interval. Further work [62] proposed to gradually decrease the upper threshold of the interval to improve the sharpness of the results. Unlike these two works which aim at generation from scratch, we focus on shape completion and want to keep a

large part of the surface intact. Therefore, we cut the upper threshold, finding that the interval of $[0, 0.5]$ is sufficient for our needs.

**Testing competitors.** For our comparisons we selected three generative methods, namely RealFusion [28], One-2-3-45 [25], and DreamGaussian [39]. All three methods take as input a single image, and in their official implementations assume some default camera model for this image and only require the elevation angle w.r.t. the reconstructed object to constrain the camera position. RealFusion additionally assumes the input image to be captured from the frontal side of the object, and uses this assumption to make view-dependent textual prompts for the text-to-image model.

For a fair comparison of these methods with NeuS [43] and our method that both use ground-truth camera models and positions, we modified the implementations of the generative methods to account for these camera parameters. First, we replaced the default camera models with the ground-truth ones. Next, we calculated the elevation angles for each scene from the ground-truth camera positions, setting the vertical axis of the object using the floor part of the reference mesh. Finally, for RealFusion we additionally calculated the horizontal position of the input image w.r.t. the object, setting the forward direction of the object manually.

We found that RealFusion tends to produce disconnected parts of the scene. Therefore, we additionally postprocessed the surfaces produced by all methods keeping only the largest connected component, which is also a common post-processing step in the literature on neural surface reconstruction.

## 8. Evaluation details

**CLIP similarity.** We used CLIP similarity to evaluate the perceptual quality of the whole surface, and additionally, to evaluate the unobserved part only.

To evaluate the whole surface, we rendered the reconstructed and the reference surfaces using viewpoints distributed uniformly around the object, computed the cosine similarity between the respective CLIP [33] image embeddings, and took the average value across all views. We used ViT-L/14 CLIP model [1].

To evaluate the unobserved part only, we discarded the viewpoints for which more than $1/3$ of the rendering of the reference surface corresponded to the visible part. To define the visible part of the reference surface formally, we tested the visibility of each triangle of the reference mesh from the visible viewpoints (shown in Figs. 10 to 12) taking self-occlusions into account. We marked a triangle visible if its center was traceable without occlusions from at least 3 viewpoints.

In Fig. 5 we show an example of a pair of renderings that we compared using CLIP. To obtain such renderings, we used an orthographic camera and a simple diffuse shader
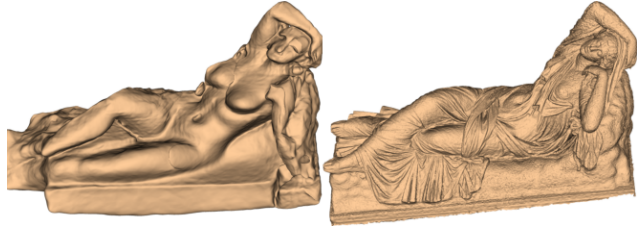


Figure 5. **Example of renderings used for CLIP similarity.**

without shadow mapping and with a single directional light source positioned behind the camera. A particular challenge of the CLIP similarity is its sensitivity and ability to distinguish differences in the geometry. Empirically, we observed that the sensitivity of CLIP similarity is higher if the surface is rendered in some color instead of grayscale. We also observed that there is no strong dependence on the actual color as long it is not gray, so we used the color shown in Fig. 5 in all computations with CLIP.

**Geometric surface quality.** To evaluate the geometric surface quality, that we report further in Sec. 9, we used the F-score computed for the whole surface, and additionally, the recall computed for the visible part only. We calculated these metrics similarly to how they are calculated in benchmarks on multi-view surface reconstruction [23, 38, 40], as we describe below.

First, to bring all scenes to the same scale, we transformed the reconstructed and the reference surfaces so that the reference surface is fitted into a unit sphere. Then, to prevent an uneven contribution of different parts of the surface to the metric, we resampled both surfaces uniformly with a sufficiently high resolution, specifically 0.003 (without changes to the geometry). After that, we calculated the precision of the reconstructed surface as the fraction of samples on this surface with the distance to the reference below a conservative threshold of 0.02 (which is an order of magnitude higher than what is usually used in benchmarks); we calculated the recall as the fraction of samples on the reference surface with the distance to the reconstructed surface below the same threshold. Finally, we computed the F-score for the whole surface as the harmonic mean of the precision and recall.

To evaluate only the visible part of the surface, we marked it as we formally described above, and calculated the recall only for the samples on the visible part. Since the calculation of the precision w.r.t. an incomplete visible part of the surface is ambiguous we only report the recall.

## 9. More quantitative and qualitative results

We show an additional qualitative comparison of the methods in Figs. 13 and 14, with the surfaces rendered from the unobserved side and from the observed side respectively;

| | | ariadne | bear | boots | buddha | bull | camera | clock | david | dog | helen | horse | lion | man | plant | santa | shiba | xia | mean |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| CLIP, Whole ↑ | RealFusion [28] | 73.9 | 69.8 | 79.2 | 69.9 | 73.8 | 77.0 | 68.1 | 74.4 | 70.3 | 73.5 | 78.3 | 77.2 | 76.4 | 75.4 | 69.1 | 80.0 | 72.2 | 74.0 |
| | One-2-3-45 [25] | 72.1 | 76.4 | 73.4 | 72.1 | 78.4 | 73.3 | 68.0 | 73.9 | 76.1 | 74.3 | 75.5 | 74.8 | 74.6 | 79.4 | 68.8 | 77.1 | 71.9 | 74.1 |
| | DreamGaussian [39] | 73.1 | 79.9 | 78.6 | 75.7 | 82.1 | 74.1 | 67.2 | 84.2 | 74.1 | 72.1 | 80.2 | 74.9 | 70.5 | 83.4 | 72.3 | 86.3 | 71.2 | 76.5 |
| | NeuS [43] | 75.0 | 87.0 | 84.3 | 81.0 | 89.2 | 81.6 | 76.4 | 84.1 | 84.9 | 81.4 | 88.7 | 86.5 | 80.4 | 89.4 | 80.1 | 87.3 | 84.4 | 83.6 |
| | Ours | 84.4 | 87.1 | 88.6 | 86.5 | 91.3 | 83.7 | 78.9 | 88.7 | 87.4 | 85.4 | 91.6 | 88.5 | 83.4 | 91.0 | 84.5 | 87.5 | 86.4 | 86.7 |
| | Ours – NeuS | 9.4 | 0.1 | 4.3 | 5.5 | 2.1 | 2.1 | 2.5 | 4.6 | 2.5 | 4.0 | 2.9 | 2.0 | 3.0 | 1.6 | 4.4 | 0.2 | 2.0 | 3.1 |
| CLIP, Unobserved ↑ | RealFusion [28] | 73.6 | 69.1 | 79.5 | 67.2 | 73.6 | 78.1 | 67.0 | 75.1 | 71.4 | 74.6 | 80.0 | 77.9 | 77.9 | 75.0 | 68.5 | 79.8 | 70.1 | 74.0 |
| | One-2-3-45 [25] | 71.9 | 75.8 | 74.4 | 69.9 | 78.8 | 73.7 | 66.5 | 74.0 | 77.5 | 75.4 | 76.9 | 75.6 | 77.3 | 79.7 | 69.6 | 78.5 | 70.5 | 74.5 |
| | DreamGaussian [39] | 72.9 | 79.2 | 77.5 | 73.1 | 82.3 | 74.6 | 66.4 | 84.2 | 75.1 | 73.0 | 80.9 | 75.7 | 71.6 | 83.1 | 72.6 | 85.3 | 68.8 | 76.2 |
| | NeuS [43] | 73.0 | 82.8 | 82.6 | 76.1 | 87.2 | 78.2 | 69.7 | 81.6 | 82.3 | 77.6 | 86.8 | 83.8 | 77.9 | 87.0 | 74.3 | 85.0 | 79.4 | 80.3 |
| | Ours | 83.2 | 83.5 | 86.8 | 87.1 | 90.9 | 82.8 | 75.2 | 87.9 | 86.6 | 85.3 | 90.2 | 87.3 | 81.5 | 89.9 | 81.7 | 86.8 | 84.3 | 85.4 |
| | Ours – NeuS | 10.2 | 0.7 | 4.2 | 11.0 | 3.7 | 4.6 | 5.5 | 6.3 | 4.3 | 7.7 | 3.4 | 3.5 | 3.6 | 2.9 | 7.4 | 1.8 | 4.9 | 5.1 |
| F-score, Whole ↑ | RealFusion [28] | 9.6 | 26.6 | 17.4 | 10.9 | 24.5 | 14.7 | 13.5 | 29.5 | 37.1 | 29.4 | 20.5 | 12.9 | 8.7 | 33.9 | 30.1 | 12.5 | 16.9 | 20.5 |
| | One-2-3-45 [25] | 11.5 | 17.9 | 8.9 | 16.3 | 26.8 | 15.8 | 13.7 | 7.5 | 14.0 | 21.9 | 20.4 | 11.4 | 2.1 | 18.6 | 14.1 | 16.4 | 16.4 | 14.9 |
| | DreamGaussian [39] | 12.5 | 29.9 | 23.7 | 18.8 | 27.4 | 17.4 | 14.0 | 41.8 | 18.5 | 28.2 | 27.7 | 22.0 | 18.5 | 29.3 | 36.2 | 39.5 | 27.8 | 25.5 |
| | NeuS [43] | 47.7 | 64.8 | 59.1 | 76.1 | 72.8 | 49.1 | 50.2 | 68.7 | 66.2 | 69.4 | 67.5 | 72.1 | 37.8 | 65.3 | 60.4 | 59.5 | 67.9 | 62.0 |
| | Ours | 56.6 | 63.4 | 49.9 | 62.8 | 71.0 | 53.2 | 49.9 | 74.0 | 64.6 | 73.2 | 69.9 | 68.9 | 54.1 | 66.2 | 67.1 | 57.7 | 64.5 | 62.8 |
| Recall, Visible ↑ | RealFusion [28] | 6.7 | 39.2 | 20.2 | 22.0 | 38.8 | 20.7 | 9.1 | 42.0 | 40.9 | 48.3 | 28.1 | 10.5 | 9.8 | 56.4 | 27.7 | 12.0 | 17.3 | 26.5 |
| | One-2-3-45 [25] | 12.9 | 13.6 | 9.5 | 14.8 | 33.7 | 14.4 | 22.1 | 7.4 | 13.0 | 17.5 | 19.6 | 5.5 | 4.0 | 22.2 | 12.4 | 17.6 | 15.6 | 15.0 |
| | DreamGaussian [39] | 17.9 | 37.5 | 40.0 | 24.0 | 31.5 | 38.3 | 19.4 | 46.6 | 36.6 | 32.1 | 36.1 | 37.4 | 40.5 | 44.4 | 47.7 | 44.3 | 30.8 | 35.6 |
| | NeuS [43] | 91.8 | 98.9 | 92.3 | 99.8 | 97.7 | 91.6 | 97.3 | 95.3 | 93.8 | 85.2 | 98.1 | 100.0 | 82.8 | 97.4 | 98.8 | 95.8 | 97.5 | 94.9 |
| | Ours | 95.5 | 97.3 | 83.2 | 98.9 | 96.7 | 83.0 | 89.5 | 97.0 | 89.1 | 86.5 | 99.1 | 99.3 | 92.2 | 96.3 | 98.0 | 95.4 | 95.9 | 93.7 |

Table 3. **Quantitative comparison.** Top: CLIP similarity for the whole surface and for the unobserved side only. "Ours – NeuS" is the difference between the respective values. Bottom: the F-score for the whole surface and the recall for the visible side only. The **best** and second best results are highlighted.

in Fig. 15 we show the renderings from the observed side for the scenes shown in Figure 2 in the main text. In Tab. 3 we show an additional quantitative comparison, using CLIP similarity for the whole surface and the unobserved side only, the F-score for the whole surface, and the precision for the visible side only.

Our method consistently outperforms all competitors w.r.t. CLIP similarity computed for both the whole surface and the unobserved side only, as we show at the top of Tab. 3. The advantage of our method compared to NeuS is even bigger for the unobserved side (the absolute values for both methods are lower since we exclude the more accurately reconstructed visible side from averaging).

In general, we noticed that CLIP similarity of shaded surface renderings is not always consistent with human perception. While being widely used as a quality measure for colored renderings, it on the one hand is sometimes not very sensitive to the differences between shaded renderings, and on the other hand, is sometimes unstable. For example, the CLIP similarity for the unobserved side of *david* is better for DreamGaussian than for NeuS, but the result is qualitatively

more similar to the reference for NeuS. We also experimented with other perceptual metrics, namely LPIPS [59] and SSIM [46], but they rather picked up local differences while we were interested in a more global semantic similarity. Therefore, we additionally evaluated the perceptual quality with a user study, as we describe in the main text.

As for the geometric surface quality, our method accurately reconstructs the visible part of the surface, on par with NeuS, as we show in Figs. 14 and 15, and confirm with the recall at the bottom of Tab. 3. The F-score computed for the whole surface shows mixed results since both methods reconstruct the unobserved part with an arbitrary geometric alignment to the reference. Our method generates a plausible completion, which does not have to be perfectly aligned, while NeuS smoothly connects the edges of the visible part.

The generative single-view methods struggle in our setting. Usually, they are only able to reconstruct the general shape of the object, as shown in Figs. 14 and 15. We discuss the challenges of our setting for these methods below.
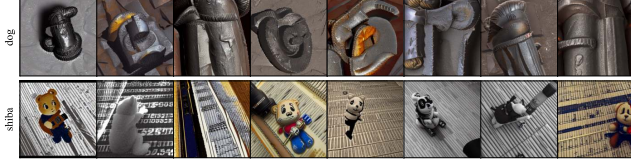
Figure 6. **Images generated with Stable Diffusion** 1.5 for textual embeddings estimated in experiments with RealFusion.



Figure 7. **Images generated with Zero123** in experiments with One-2-3-45.



Figure 8. **Adaptation of DreamGaussian to multi-view setting.** We show the color renderings from a fitted model for different numbers of input real-world images in the first row, the respective surfaces in the second row, and the density renderings from the model in the last row. In the first column, we show the reference data. For clarity, we show the results for a frontal viewpoint, not used in our main experiments.

## 10. Discussion of results of competitors

**RealFusion** trains an Instant NGP [30] representation of the scene that it supervises via SDS using text-to-image Stable Diffusion 1.5 model. It estimates the textual embedding, to condition the model on, from the input image using textual inversion [18].

In Fig. 6 we show some examples of images generated with Stable Diffusion 1.5 for the textual embeddings estimated by RealFusion for our test scenes. We observe that for some scenes the generated images are inconsistent with the input view, and hypothesize that in our setting, with limited information about the whole surface in the input view, the textual inversion produces incoherent textual embeddings, which leads to degraded results of RealFusion. Additionally, the Instant NGP representation of the surface lacks sufficient constraints on its level set, which leads to further degradation of the results.

**One-2-3-45** trains a SparseNeuS [27] representation of the scene that it supervises with synthetic images generated using Zero-123 [26] model, conditioned on the input real-world image and sampled camera poses for the synthetic images.

In Fig. 7 we show some examples of images produced by Zero-123 for our test scenes. We observe that for some scenes the generated images lack multi-view consistency, which inevitably leads to degradation of the 3D surface fitted to them. The authors of One-2-3-45 note that the quality of the results produced by their method is "often limited by the multi-view generation, since there are 3D inconsisten-
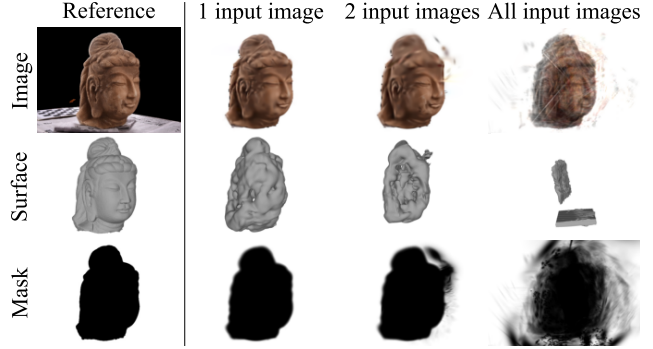
cies among the predicted views", and show a failure case for their method in a setting similar to ours, where the input view contains insufficient information about the whole surface. Additionally, we observe that for some scenes the generated images, while being rather consistent, represent an incorrect 3D shape.

**DreamGaussian,** similar to One-2-3-45, also uses Zero-123 to estimate multi-view information for the scene from a single input view. In contrast to One-2-3-45, it does not train the 3D representation of the scene on the generated images directly but instead uses SDS guidance. This presumably still leads to the same effects related to multi-view inconsistency of the estimation.

As the 3D representation, it uses Gaussian Splatting [22], with which we associate the "hollowness" of the extracted surfaces. We note that the results that we obtained with this method are qualitatively similar to the results previously obtained by others [2–5].

**Multi-view DreamGaussian.** For the single-view method that arguably produced the best results in our experiments, namely DreamGaussian, we tried to make an extension to our multi-view setting. For this, we followed a discussion on GitHub on a similar topic [6]. The modified multi-view version produces a worse surface than the single-view baseline, as we show in the first two rows of Fig. 8. We confirmed that the model fits the added real-world images one by one, but when the number of images increases the quality of the result degrades. We observe that the addition of real-world input images leads to the emergence of gaussian-splats with a nonzero density and the color of the background, floating around the object, as we show in the last row of Fig. 8. This is one of the possible reasons for the degradation of the surface. Overall, we hypothesize that the

values of the hyperparameters of the method picked for the single-view setting, such as the training parameters of Gaussian Splatting, the parameters of the mesh extraction, or the weights of loss terms, are not suitable for the multi-view setting.

Below, we compare to another multi-view baseline, namely a combination of NeuS with color-based SDS.

## 11. More ablation study

We show the qualitative results of the ablation study for all scenes in Figs. 16 and 17. In Tab. 4 we show the respective quantitative comparison, using CLIP similarity for the whole surface and the unobserved side only.

Our complete method consistently improves upon the baseline NeuS+SDS, which uses multiple real-world input images. The addition of the frozen SDS and the multi-view SDS improves the results in several cases qualitatively, but their effects are not picked up by CLIP similarity on average.

## 12. Premise of multi-view SDS

We obtain the SDS guidance signal for our method using an open-source diffusion model Stable Diffusion 2.1. We approach the problem of multi-view inconsistency of the SDS guidance via our multi-view SDS, that we describe in the main text. In Fig. 9 we show that Stable Diffusion, conditioned on the prompt "renders of the same ⟨object⟩ from different viewpoints", generates multi-view images of the object composed into a grid with high semantic and stylistic consistency. This demonstrates that Stable Diffusion can model complex interactions between different viewpoints, if they are combined in this way, and motivates our approach.

One may notice that the different views of the objects in Fig. 9 may be inconsistent *geometrically*. We note that in our method we do not use such images directly, but apply SDS guidance to normal maps rendered from the 3D representation of the scene, that are geometrically multi-view consistent by design.

## 13. Theoretical aspects of frozen SDS

The derivation of the Score Distillation Sampling [32] algorithm starts from the denoising diffusion training loss 7:

$$\mathcal{L}_{Diff}(\phi, \mathbf{x}) =$$
$$= \mathbb{E}_{t \sim \mathcal{U}(0,1), \epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})}[w(t)\|\epsilon_\phi(\alpha_t \mathbf{x} + \sigma_t \epsilon; t) - \epsilon\|_2^2] \quad (7)$$

The gradient of the denoising loss w.r.t. the generated datapoint $\mathbf{x}$ is taken, and by omitting the U-Net Jacobian loss we get 8:
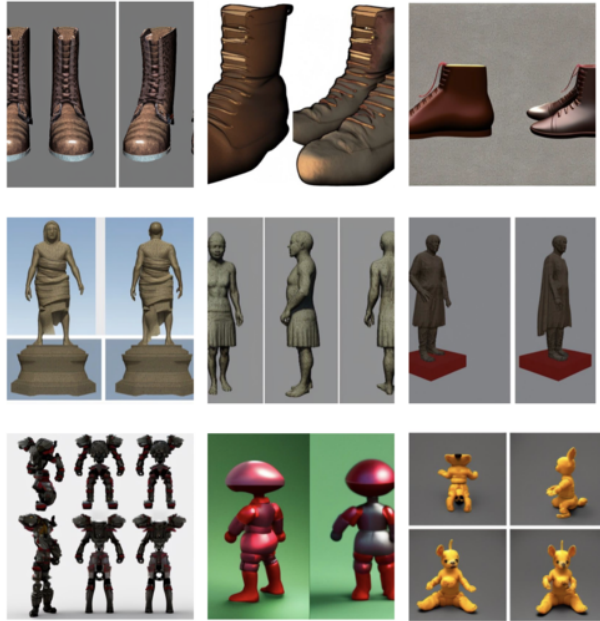


Figure 9. **Images generated with Stable Diffusion** for the prompt "renders of the same boots / statue / toy from different viewpoints".

$$\nabla\mathcal{L}_{Diff}(\phi, \mathbf{x} = g(\theta)) \triangleq$$
$$\triangleq \mathbb{E}_{t,\epsilon}\left[w(t)(\hat{\epsilon}(\mathbf{z}_t; y, t) - \epsilon)\frac{\partial \mathbf{x}}{\partial \theta}\right] \quad (8)$$

While this particular solution may appear to be ad-hoc the authors [32] show that the same gradient 8 is the gradient of the weighted probability density distillation loss [31]. Both lines of reasoning heavily rely on the assumption that $\epsilon \sim \mathcal{N}(0, \mathbf{I})$. Our results show, that fixing $\epsilon$ and getting rid of the assumption that it belongs to the Gaussian distribution changes very little in terms of outcome. This leads us to the following conclusions: a) the motivation, that $\nabla\mathcal{L}_{SDS}$ is the gradient of the denoising diffusion loss w.r.t. $\mathbf{x}$ is misleading because the denoising loss needs to be averaged across the distribution of $\epsilon \sim \mathcal{N}(0, \mathbf{I})$ b) the same equation seemingly produces the unbiased estimation of the gradient of the weighted probability density distillation loss even with the fixed $\epsilon$, with the only variable being $t$. Our observations lead us to the conclusion that in its nature the score distillation sampling seems to be much closer to the diffusion inference process than it was previously assumed. The main problem is that $\mathbf{x} = g(\theta)$ generally may lie out of the training distribution of the denoising U-net $\epsilon_\phi$. Thus we add the randomly picked at the beginning of the training, but **fixed** during training noise sample $\epsilon$, with the result now belonging to the training distribution $\mathbf{x}'_t = \alpha_t \mathbf{x} + \sigma_t \epsilon$. We use the "shifted" $\mathbf{x}'_t$ to calculate the update direction and than

apply the correction $\hat{\epsilon}((\mathbf{x}'_t; y, t)) - \epsilon)$, which mitigates the impact of the "shift" and decreases the variance [32]. Our intuition is supported by the fact that most recent works that employ the score distillation sampling to obtain impressive 3D generation results [62], use the decreasing schedule for $t$ resembling the denoising diffusion inference.
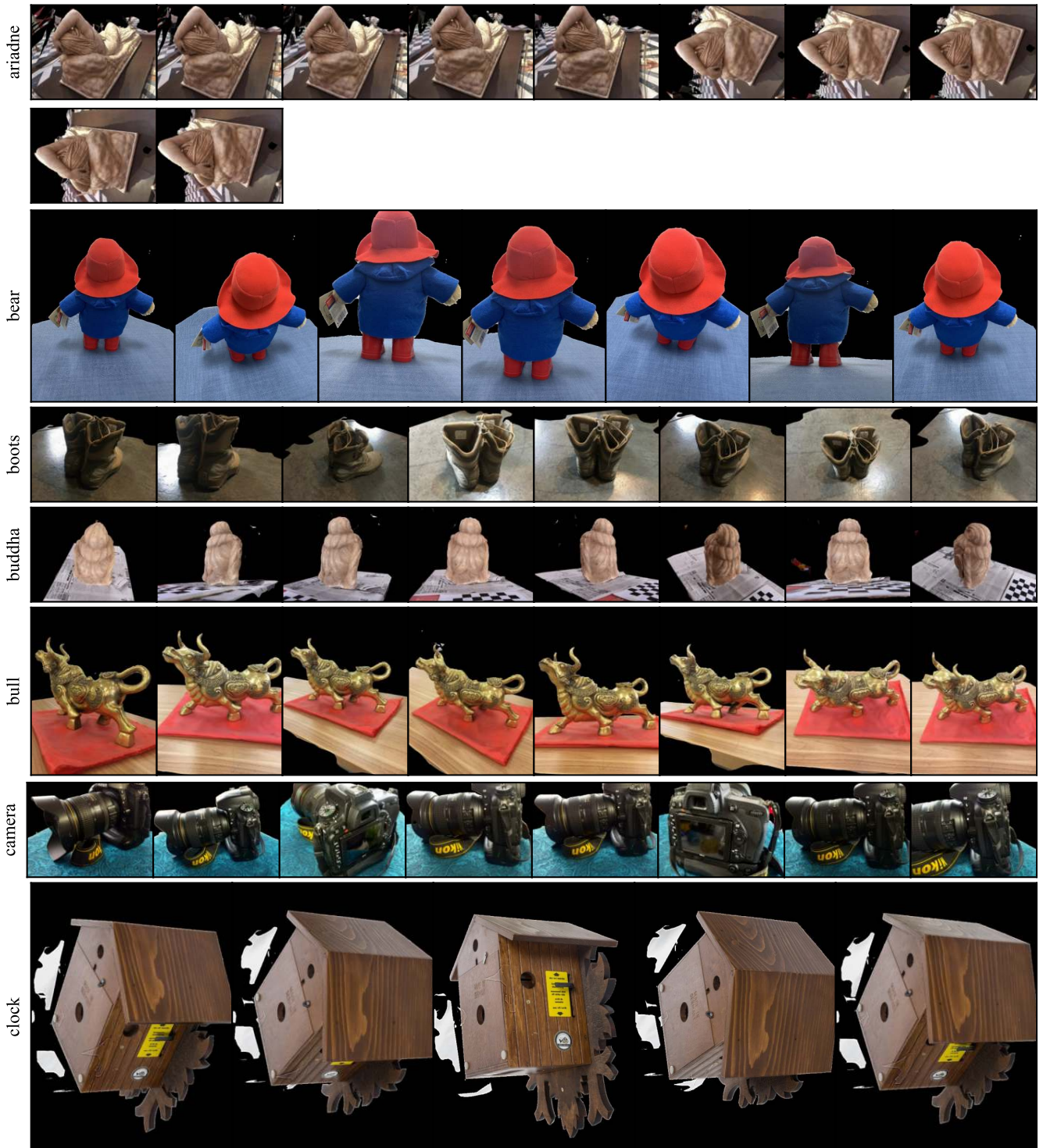
Figure 10. **Input images** that we used in our experiments.

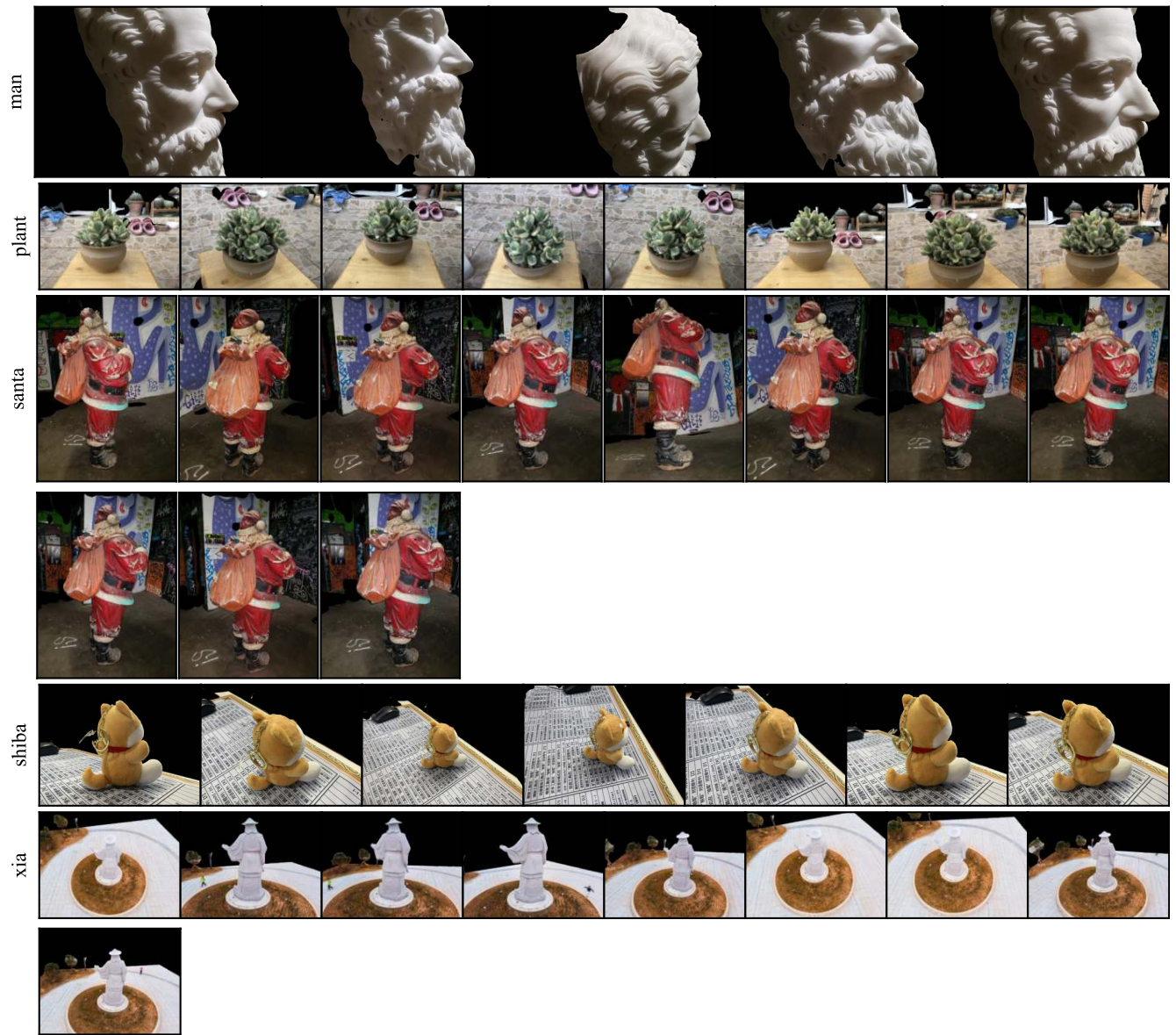Figure 11. **Input images** that we used in our experiments.

Figure 12. **Input images** that we used in our experiments.

Figure 13. **Qualitative comparison, unobserved side.** The first row shows the image from the visible set used by the competing single-view methods as input. The next four rows show the results of the three generative single-view competitors: RealFusion (RF), One-2-3-45, and DreamGaussian (DG), and the reconstructive competitor (NeuS). The next two rows show the results of the multi-view baseline (NeuS + SDS) and of our full method. The last row depicts the reference surface. All meshes are rendered from the unobserved side.
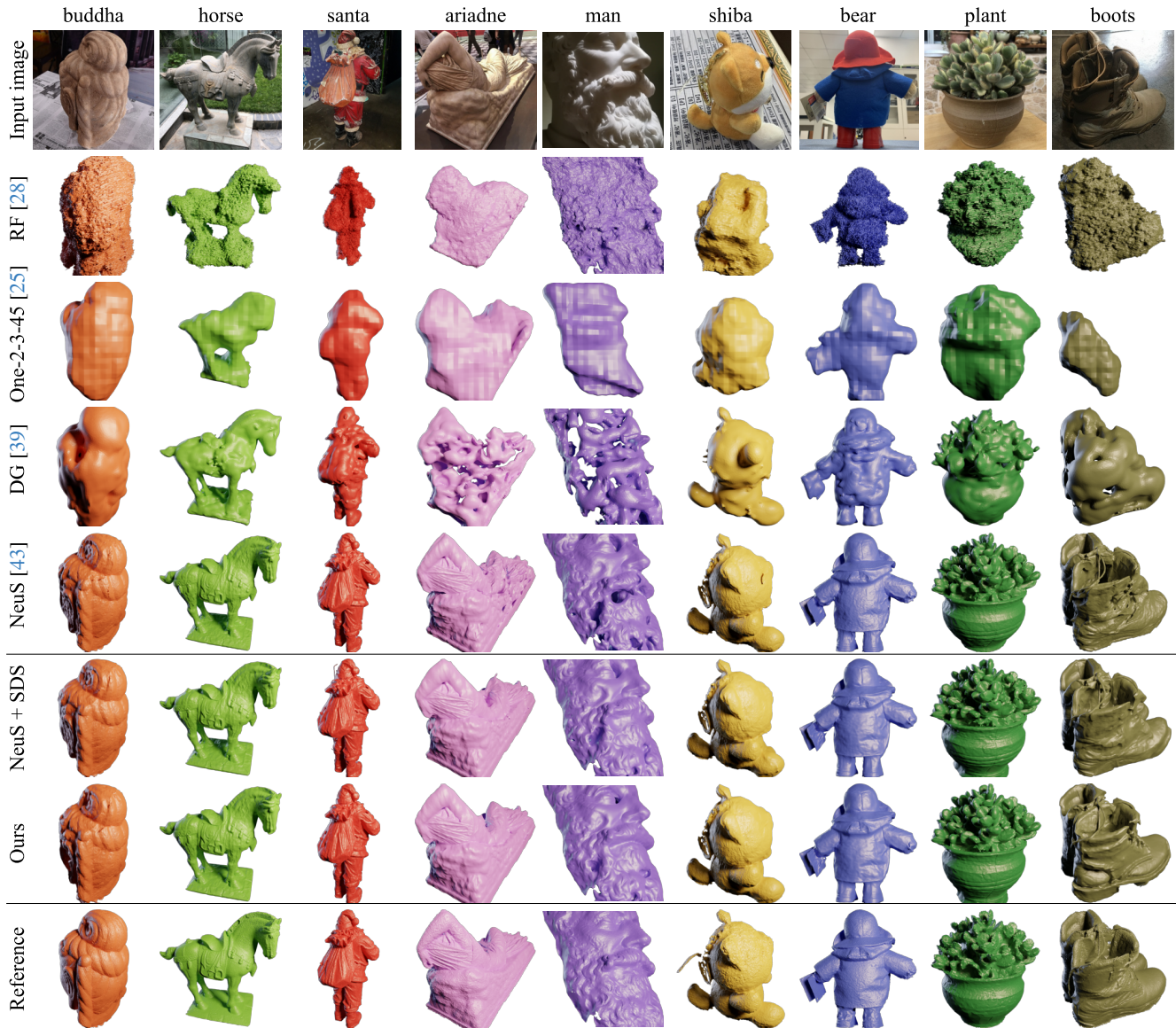
Figure 14. **Qualitative comparison, visible side.** The first row shows the image from the visible set used by the competing single-view methods as input. The next four rows show the results of the three generative single-view competitors: RealFusion (RF), One-2-3-45, and DreamGaussian (DG), and the reconstructive competitor (NeuS). The next two rows show the results of the multi-view baseline (NeuS + SDS) and of our full method. The last row depicts the reference surface. All meshes are rendered from the visible side.

Figure 15. **Qualitative comparison, visible side.** The first row shows the image from the visible set used by the competing single-view methods as input. The next four rows show the results of the three generative single-view competitors: RealFusion (RF), One-2-3-45, and DreamGaussian (DG), and the reconstructive competitor (NeuS). The next two rows show the results of the multi-view baseline (NeuS + SDS) and of our full method. The last row depicts the reference surface. All meshes are rendered from the visible side.

| | | ariadne | bear | boots | buddha | bull | camera | clock | david | dog | helen | horse | lion | man | plant | santa | shiba | xia | mean |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| CLIP, Unobserved ↑ | NeuS | 73.0 | 82.8 | 82.6 | 76.1 | 87.2 | 78.2 | 69.7 | 81.6 | 82.3 | 77.6 | 86.8 | 83.8 | 77.9 | 87.0 | 74.3 | 85.0 | 79.4 | 80.3 |
| | + SDS | 77.3 | 82.7 | 82.5 | 79.7 | 86.1 | 78.4 | 70.9 | 83.3 | 85.5 | 83.0 | 88.5 | 83.7 | 78.6 | 89.0 | 74.1 | 82.7 | 80.1 | 81.5 |
| | + Normals | 81.7 | 83.4 | 86.2 | **88.6** | 90.6 | **84.3** | 74.7 | 87.3 | **86.7** | 86.4 | 90.5 | _87.9_ | 79.4 | 89.3 | 80.0 | 83.8 | **85.5** | 85.1 |
| | + Fixed noise | _82.1_ | **84.3** | **87.6** | _87.6_ | **90.7** | 83.3 | **76.1** | _87.4_ | 86.4 | _86.0_ | 90.4 | **88.3** | _80.0_ | **90.3** | _80.9_ | _86.5_ | _84.7_ | **85.4** |
| | Ours | **83.2** | _83.5_ | _86.8_ | 87.1 | **90.9** | 82.8 | _75.2_ | **87.9** | _86.6_ | 85.3 | 90.2 | 87.3 | **81.5** | _89.9_ | **81.7** | **86.8** | 84.3 | **85.4** |
| CLIP, Whole ↑ | NeuS | 75.0 | 87.0 | 84.3 | 81.0 | 89.2 | 81.6 | 76.4 | 84.1 | 84.9 | 81.4 | 88.7 | 86.5 | 80.4 | 89.4 | 80.1 | 87.3 | 84.4 | 83.6 |
| | + SDS | 78.8 | 86.5 | 84.4 | 83.0 | 87.8 | 80.4 | 77.3 | 85.5 | 87.3 | 83.9 | 89.9 | 85.6 | 80.1 | 90.5 | 78.6 | 84.4 | 84.1 | 84.0 |
| | + Normals | 82.6 | 86.8 | 87.8 | **88.0** | 90.9 | **86.2** | _79.3_ | 88.4 | **88.0** | 86.8 | **91.7** | _89.1_ | 80.5 | 90.9 | 83.2 | 85.6 | **87.4** | _86.7_ |
| | + Fixed noise | _83.7_ | **87.7** | **89.1** | _87.3_ | _91.1_ | _84.9_ | **79.4** | _88.6_ | _87.9_ | 85.5 | **91.8** | **89.7** | _81.5_ | **91.4** | _83.6_ | _87.4_ | 86.4 | **86.9** |
| | Ours | **84.4** | _87.1_ | _88.6_ | 86.5 | **91.3** | 83.7 | 78.9 | **88.7** | 87.4 | 85.4 | 91.6 | 88.5 | **83.4** | _91.0_ | **84.5** | **87.5** | _86.4_ | _86.7_ |

Table 4. **Quantitative ablation study.** Top: CLIP similarity for the unobserved part of the surface; bottom: for the whole surface. The **best** and _second best_ results are highlighted.
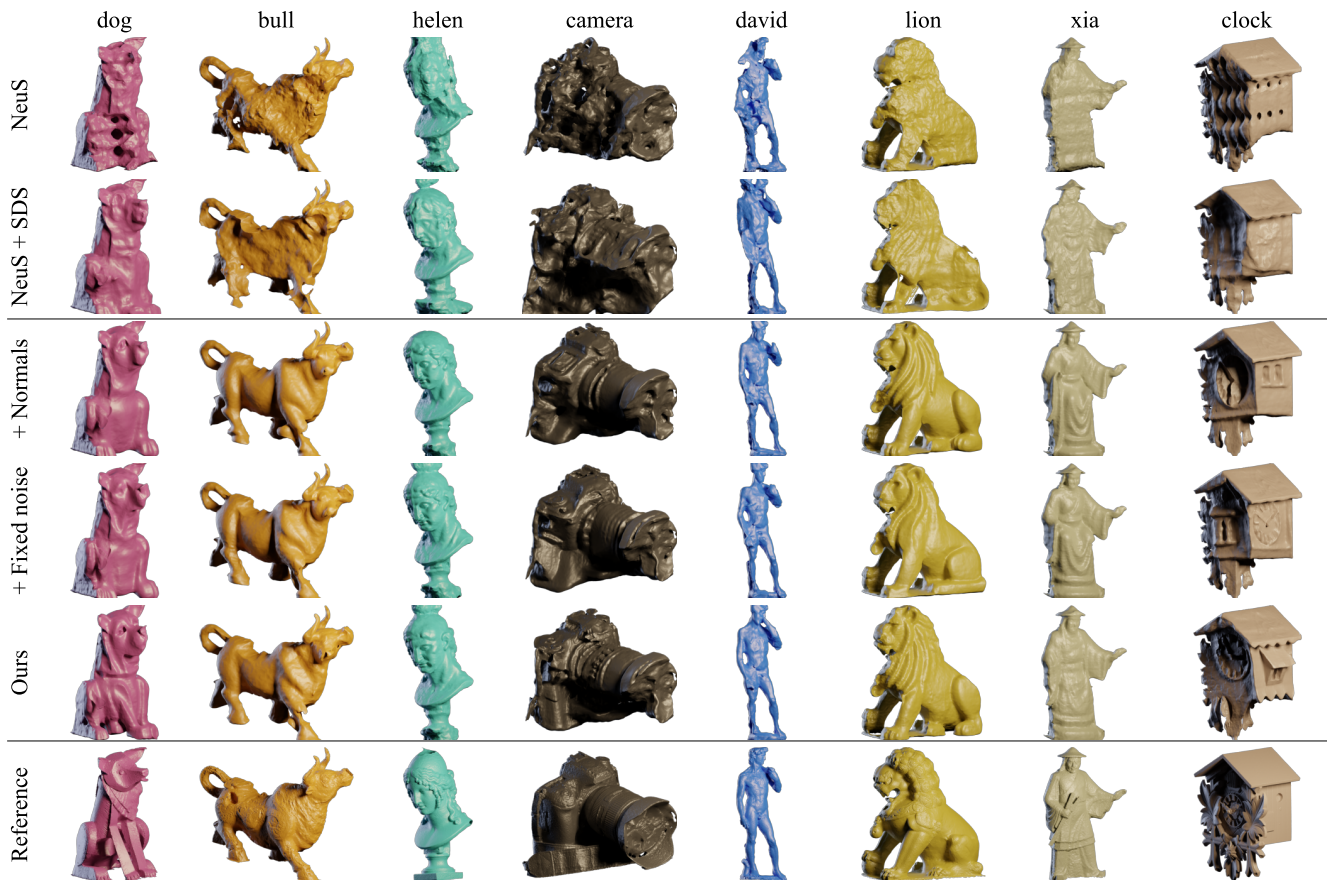


Figure 16. **Qualitative ablation study.** The first two rows show the results of NeuS and our multi-view baseline: NeuS with added color image-based SDS guidance. The next three rows show the results produced by successively adding normal maps, frozen SDS, and multi-view SDS guidance. The last row depicts the reference surface. All meshes are rendered from the unobserved side.

Figure 17. **Qualitative ablation study.** The first two rows show the results of NeuS and our multi-view baseline: NeuS with added color image-based SDS guidance. The next three rows show the results produced by successively adding normal maps, frozen SDS, and multi-view SDS guidance. The last row depicts the reference surface. All meshes are rendered from the unobserved side.

# References

[1] Clip ViT-L/14. https://huggingface.co/openai/clip-vit-large-patch14. 10

[2] A discussion of results of dreamgaussian. https://github.com/dreamgaussian/dreamgaussian/issues/15, 2023. 12

[3] A discussion of results of dreamgaussian. https://github.com/dreamgaussian/dreamgaussian/issues/33, 2023.

[4] A discussion of results of dreamgaussian. https://github.com/dreamgaussian/dreamgaussian/issues/43, 2023.

[5] A discussion of results of dreamgaussian. https://github.com/dreamgaussian/dreamgaussian/issues/65, 2023. 12

[6] Extending dreamgaussian to multi-view setting. https://github.com/dreamgaussian/dreamgaussian/issues/22, 2023. 12

[7] Henrik Aanæs, Rasmus Ramsbøl Jensen, George Vogiatzis, Engin Tola, and Anders Bjorholm Dahl. Large-scale data for multiple-view stereopsis. *International Journal of Computer Vision*, 120:153–168, 2016. 5

[8] Mohammadreza Armandpour, Huangjie Zheng, Ali Sadeghian, Amir Sadeghian, and Mingyuan Zhou. Re-imagine the negative prompt algorithm: Transform 2d diffusion into 3d, alleviate janus problem and beyond. *arXiv preprint arXiv:2304.04968*, 2023. 4

[9] Yingjie Cai, Kwan-Yee Lin, Chao Zhang, Qiang Wang, Xiaogang Wang, and Hongsheng Li. Learning a structured latent space for unsupervised point cloud completion. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5543–5553, 2022. 3

[10] Angel X Chang, Thomas Funkhouser, Leonidas Guibas, Pat Hanrahan, Qixing Huang, Zimo Li, Silvio Savarese, Manolis Savva, Shuran Song, Hao Su, et al. Shapenet: An information-rich 3d model repository. *arXiv preprint arXiv:1512.03012*, 2015. 2, 3

[11] Yida Chen, Fernanda Viégas, and Martin Wattenberg. Beyond surface statistics: Scene representations in a latent diffusion model. *arXiv preprint arXiv:2306.05720*, 2023. 3

[12] Yen-Chi Cheng, Hsin-Ying Lee, Sergey Tulyakov, Alexander G Schwing, and Liang-Yan Gui. Sdfusion: Multimodal 3d shape completion, reconstruction, and generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4456–4465, 2023. 3

[13] Ruihang Chu, Enze Xie, Shentong Mo, Zhenguo Li, Matthias Nießner, Chi-Wing Fu, and Jiaya Jia. Diffcomplete: Diffusion-based generative 3d shape completion. *arXiv preprint arXiv:2306.16329*, 2023. 3

[14] François Darmon, Bénédicte Bascle, Jean-Clément Devaux, Pascal Monasse, and Mathieu Aubry. Improving neural implicit surfaces geometry with patch warping. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6260–6269, 2022. 2

[15] Andreea Dogaru, Andrei-Timotei Ardelean, Savva Ignatyev, Egor Zakharov, and Evgeny Burnaev. Sphere-guided training of neural implicit surfaces. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 20844–20853, 2023. 2

[16] Ben Fei, Weidong Yang, Wen-Ming Chen, Zhijun Li, Yikang Li, Tao Ma, Xing Hu, and Lipeng Ma. Comprehensive review of deep learning-based 3d point cloud completion processing and analysis. *IEEE Transactions on Intelligent Transportation Systems*, 2022. 3

[17] Qiancheng Fu, Qingshan Xu, Yew Soon Ong, and Wenbing Tao. Geo-neus: Geometry-consistent neural implicit surfaces learning for multi-view reconstruction. *Advances in Neural Information Processing Systems*, 35:3403–3416, 2022. 2

[18] Rinon Gal, Yuval Alaluf, Yuval Atzmon, Or Patashnik, Amit H Bermano, Gal Chechik, and Daniel Cohen-Or. An image is worth one word: Personalizing text-to-image generation using textual inversion. *arXiv preprint arXiv:2208.01618*, 2022. 3, 4, 12

[19] Amos Gropp, Lior Yariv, Niv Haim, Matan Atzmon, and Yaron Lipman. Implicit geometric regularization for learning shapes. *arXiv preprint arXiv:2002.10099*, 2020. 2

[20] 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. 5

[21] Yoni Kasten, Ohad Rahamim, and Gal Chechik. Point-cloud completion with pretrained text-to-image diffusion models. *arXiv preprint arXiv:2306.10533*, 2023. 3

[22] Bernhard Kerbl, Georgios Kopanas, Thomas Leimkühler, and George Drettakis. 3d gaussian splatting for real-time radiance field rendering. *ACM Transactions on Graphics (ToG)*, 42(4):1–14, 2023. 3, 12

[23] Arno Knapitsch, Jaesik Park, Qian-Yi Zhou, and Vladlen Koltun. Tanks and temples: Benchmarking large-scale scene reconstruction. *ACM Transactions on Graphics*, 36(4), 2017. 10

[24] Chen-Hsuan Lin, Jun Gao, Luming Tang, Towaki Takikawa, Xiaohui Zeng, Xun Huang, Karsten Kreis, Sanja Fidler, Ming-Yu Liu, and Tsung-Yi Lin. Magic3d: High-resolution text-to-3d content creation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 300–309, 2023. 3

[25] Minghua Liu, Chao Xu, Haian Jin, Linghao Chen, Zexiang Xu, Hao Su, et al. One-2-3-45: Any single image to 3d mesh in 45 seconds without per-shape optimization. *arXiv preprint arXiv:2306.16928*, 2023. 2, 3, 5, 6, 8, 10, 11, 18, 19, 20

[26] Ruoshi Liu, Rundi Wu, Basile Van Hoorick, Pavel Tokmakov, Sergey Zakharov, and Carl Vondrick. Zero-1-to-3: Zero-shot one image to 3d object. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 9298–9309, 2023. 3, 12

[27] Xiaoxiao Long, Cheng Lin, Peng Wang, Taku Komura, and Wenping Wang. Sparseneus: Fast generalizable neural surface reconstruction from sparse views. In *European Conference on Computer Vision*, pages 210–227. Springer, 2022. 2, 12

[28] Luke Melas-Kyriazi, Iro Laina, Christian Rupprecht, and Andrea Vedaldi. Realfusion: 360deg reconstruction of any object from a single image. In *Proceedings of the IEEE/CVF*

*Conference on Computer Vision and Pattern Recognition*, pages 8446–8455, 2023. 2, 3, 4, 5, 6, 9, 10, 11, 18, 19, 20

[29] 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. *Communications of the ACM*, 65(1):99–106, 2021. 2

[30] Thomas Müller, Alex Evans, Christoph Schied, and Alexander Keller. Instant neural graphics primitives with a multiresolution hash encoding. *ACM Transactions on Graphics (ToG)*, 41(4):1–15, 2022. 12

[31] Aaron van den Oord, Yazhe Li, and Oriol Vinyals. Representation learning with contrastive predictive coding. *arXiv preprint arXiv:1807.03748*, 2018. 13

[32] Ben Poole, Ajay Jain, Jonathan T. Barron, and Ben Mildenhall. Dreamfusion: Text-to-3d using 2d diffusion. *arXiv*, 2022. 2, 3, 4, 9, 13, 14

[33] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PMLR, 2021. 6, 10

[34] Robin Rombach and Patrick Esser. Stable diffusion v2-1. https://huggingface.co/stabilityai/stable-diffusion-2-1. 5

[35] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10684–10695, 2022. 2, 5

[36] Johannes Lutz Schönberger and Jan-Michael Frahm. Structure-from-Motion Revisited. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016. 2

[37] Johannes Lutz Schönberger, Enliang Zheng, Marc Pollefeys, and Jan-Michael Frahm. Pixelwise View Selection for Unstructured Multi-View Stereo. In *European Conference on Computer Vision (ECCV)*, 2016. 2, 9

[38] Thomas Schöps, Johannes L. Schönberger, Silvano Galliani, Torsten Sattler, Konrad Schindler, Marc Pollefeys, and Andreas Geiger. A multi-view stereo benchmark with high-resolution images and multi-camera videos. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017. 10

[39] Jiaxiang Tang, Jiawei Ren, Hang Zhou, Ziwei Liu, and Gang Zeng. Dreamgaussian: Generative gaussian splatting for efficient 3d content creation. *arXiv preprint arXiv:2309.16653*, 2023. 1, 3, 5, 6, 10, 11, 18, 19, 20

[40] Oleg Voynov, Gleb Bobrovskikh, Pavel Karpyshev, Saveliy Galochkin, Andrei-Timotei Ardelean, Arseniy Bozhenko, Ekaterina Karmanova, Pavel Kopanev, Yaroslav Labutin-Rymsho, Ruslan Rakhimov, Aleksandr Safin, Valerii Serpiva, Alexey Artemov, Evgeny Burnaev, Dzmitry Tsetserukou, and Denis Zorin. Multi-Sensor Large-Scale dataset for Multi-View 3D reconstruction. In *2023 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 21392–21403. IEEE, 2023. 10

[41] Haochen Wang, Xiaodan Du, Jiahao Li, Raymond A Yeh, and Greg Shakhnarovich. Score jacobian chaining: Lifting pretrained 2d diffusion models for 3d generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12619–12629, 2023. 3

[42] Jiepeng Wang, Peng Wang, Xiaoxiao Long, Christian Theobalt, Taku Komura, Lingjie Liu, and Wenping Wang. Neuris: Neural reconstruction of indoor scenes using normal priors. In *European Conference on Computer Vision*, pages 139–155. Springer, 2022. 2

[43] 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. 1, 2, 3, 5, 6, 8, 10, 11, 18, 19, 20

[44] Yusen Wang, Zongcheng Li, Yu Jiang, Kaixuan Zhou, Tuo Cao, Yanping Fu, and Chunxia Xiao. Neuralroom: Geometry-constrained neural implicit surfaces for indoor scene reconstruction. *arXiv preprint arXiv:2210.06853*, 2022. 2

[45] Yiming Wang, Qin Han, Marc Habermann, Kostas Daniilidis, Christian Theobalt, and Lingjie Liu. Neus2: Fast learning of neural implicit surfaces for multi-view reconstruction. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 3295–3306, 2023. 2

[46] Zhou Wang, Alan C Bovik, Hamid R Sheikh, and Eero P Simoncelli. Image quality assessment: from error visibility to structural similarity. *IEEE transactions on image processing*, 13(4):600–612, 2004. 11

[47] Rundi Wu, Xuelin Chen, Yixin Zhuang, and Baoquan Chen. Multimodal shape completion via conditional generative adversarial networks. In *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part IV 16*, pages 281–296. Springer, 2020. 3

[48] Peng Xiang, Xin Wen, Yu-Shen Liu, Yan-Pei Cao, Pengfei Wan, Wen Zheng, and Zhizhong Han. Snowflakenet: Point cloud completion by snowflake point deconvolution with skip-transformer. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 5499–5509, 2021. 3

[49] Xingguang Yan, Liqiang Lin, Niloy J Mitra, Dani Lischinski, Daniel Cohen-Or, and Hui Huang. Shapeformer: Transformer-based shape completion via sparse representation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6239–6249, 2022. 3

[50] Yao Yao, Zixin Luo, Shiwei Li, Tian Fang, and Long Quan. Mvsnet: Depth inference for unstructured multi-view stereo. *European Conference on Computer Vision (ECCV)*, 2018. 9

[51] Yao Yao, Zixin Luo, Shiwei Li, Jingyang Zhang, Yufan Ren, Lei Zhou, Tian Fang, and Long Quan. Blendedmvs: A large-scale dataset for generalized multi-view stereo networks. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 1790–1799, 2020. 2, 4, 9

[52] Lior Yariv, Yoni Kasten, Dror Moran, Meirav Galun, Matan Atzmon, Basri Ronen, and Yaron Lipman. Multiview neural surface reconstruction by disentangling geometry and ap-

pearance. *Advances in Neural Information Processing Systems*, 33:2492–2502, 2020. 2

[53] Lior Yariv, Jiatao Gu, Yoni Kasten, and Yaron Lipman. Volume rendering of neural implicit surfaces. *Advances in Neural Information Processing Systems*, 34:4805–4815, 2021. 2

[54] Alex Yu, Vickie Ye, Matthew Tancik, and Angjoo Kanazawa. pixelNeRF: Neural radiance fields from one or few images. In *CVPR*, 2021. 2

[55] Xumin Yu, Yongming Rao, Ziyi Wang, Zuyan Liu, Jiwen Lu, and Jie Zhou. Pointr: Diverse point cloud completion with geometry-aware transformers. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 12498–12507, 2021. 3

[56] Igor Zacharov, Rinat Arslanov, Maksim Gunin, Daniil Stefonishin, Andrey Bykov, Sergey Pavlov, Oleg Panarin, Anton Maliutin, Sergey Rykovanov, and Maxim Fedorov. "zhores"-petaflops supercomputer for data-driven modeling, machine learning and artificial intelligence installed in skolkovo institute of science and technology. *Open Engineering*, 9(1): 512–520, 2019. 8

[57] Junzhe Zhang, Xinyi Chen, Zhongang Cai, Liang Pan, Haiyu Zhao, Shuai Yi, Chai Kiat Yeo, Bo Dai, and Chen Change Loy. Unsupervised 3d shape completion through gan inversion. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1768–1777, 2021. 3

[58] Lvmin Zhang, Anyi Rao, and Maneesh Agrawala. Adding conditional control to text-to-image diffusion models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 3836–3847, 2023. 2

[59] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *CVPR*, 2018. 11

[60] Linqi Zhou, Yilun Du, and Jiajun Wu. 3d shape generation and completion through point-voxel diffusion. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5826–5835, 2021. 3

[61] Zhizhuo Zhou and Shubham Tulsiani. Sparsefusion: Distilling view-conditioned diffusion for 3d reconstruction. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12588–12597, 2023. 3, 4

[62] Joseph Zhu and Peiye Zhuang. Hifa: High-fidelity text-to-3d with advanced diffusion guidance. *arXiv preprint arXiv:2305.18766*, 2023. 2, 3, 4, 9, 14