# DITTO-NeRF: Diffusion-based Iterative Text To Omni-directional 3D Model

Hoigi Seo[1*]      Hayeon Kim[1*]      Gwanghyun Kim[1*]      Se Young Chun[1,2†]

[1]Dept. of Electrical and Computer Engineering, [2]INMC & IPAI

Seoul National University, Republic of Korea

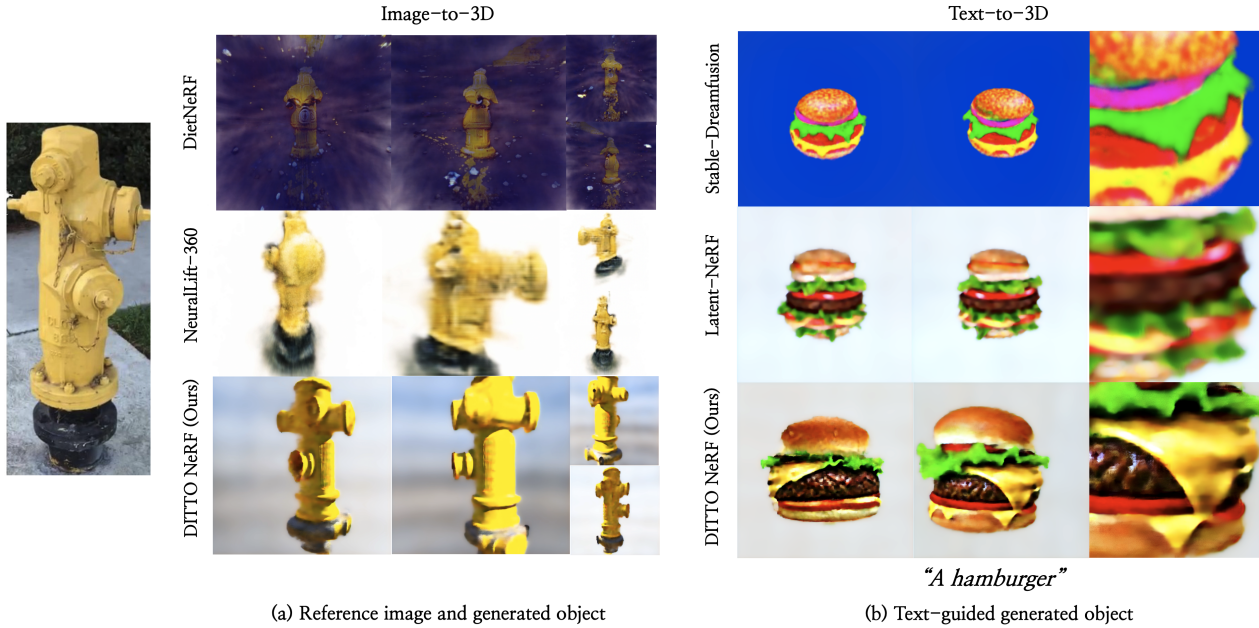{seohoiki3215, khy5630, gwang.kim, sychun}@snu.ac.kr

Image-to-3D

Text-to-3D

"A hamburger"

(a) Reference image and generated object

(b) Text-guided generated object

Figure 1: 3D NeRF models generated from (a) a reference image along with the text "*a yellow fire hydrant*" and (b) the text "*a hamburger*" using our DITTO-NeRF and prior arts (DietNeRF [25], NeuralLift-360 [68], Latent-NeRF [34], Stable-DreamFusion [54]). In (b), photos are zoomed in the rightmost column for comparisons. Our DITTO-NeRF yielded exquisite images in all views from text and/or image. See the supplementary videos at janeyeon.github.io/ditto-nerf.

## Abstract

*The increasing demand for high-quality 3D content creation has motivated the development of automated methods for creating 3D object models from a single image and/or from a text prompt. However, the reconstructed 3D objects using state-of-the-art image-to-3D methods still exhibit low correspondence to the given image and low multi-view consistency. Recent state-of-the-art text-to-3D methods are also limited, yielding 3D samples with low diversity per prompt with long synthesis time. To address these challenges, we propose DITTO-NeRF, a novel pipeline to generate a high-quality 3D NeRF model from a text prompt or a single image. Our DITTO-NeRF consists of constructing high-quality par-tial 3D object for limited in-boundary (IB) angles using the given or text-generated 2D image from the frontal view and then iteratively reconstructing the remaining 3D NeRF using inpainting latent diffusion model. We propose progressive 3D object reconstruction schemes in terms of scales (low to high resolution), angles (IB angles initially to outer-boundary (OB) later), and masks (object to background boundary) in our DITTO-NeRF so that high-quality information on IB can be propagated into OB. Our DITTO-NeRF outperforms state-of-the-art methods in terms of fidelity and diversity qualitatively and quantitatively with much faster training times than prior arts on image/text-to-3D such as DreamFusion, and NeuralLift-360.*

---

[*]Equal contribution. [†]Corresponding author.

# 1. Introduction

Recent advancements in virtual reality and augmented reality have led to a rapid increase in demand for 3D content. Nevertheless, traditionally the creation of high-quality 3D objects has been a time-consuming and costly process that requires human experts. The challenge of the high cost of making 3D objects has motivated to development of methods for synthesizing diverse 3D objects from simplified source inputs. The methods of generating 3D objects from a single image (image-to-3D) [10, 11, 12, 26, 30, 50, 57, 61, 66, 4, 73, 68] have been developed. The creation of a 3D object from a single image is a challenging task that is hindered by the insufficiency of available information. Recent image-to-3D studies such as DietNeRF [25] and NeuralLift-360 [68] have shown the impressive results leveraging the prior knowledge of pre-trained CLIP (Contrastive Language-Image Pre-training) [40] models or text-to-image diffusion models. Other approaches are the methods of generating 3D objects from a text prompt (text-to-3D) [39, 34, 32, 55] by leveraging text-to-image model to create 3D representation.

Recent image-to-3D methods still suffer low correspondence with the 2D input image, yielding unsatisfactory 3D output and modern text-to-3D methods also suffer low diversity in generated 3D objects from the same input text and high computation complexity. To mitigate these limitations in 3D object generation, we propose DITTO-NeRF, a diffusion-based iterative text to omni-directional 3d model, which utilizes the high diversity and high fidelity images generated by the latent diffusion model in response to a given text prompt. Specifically, our DITTO-NeRF incorporates a monocular depth estimation model [42] to predict the depth corresponding to the image and subsequently builds a high-quality partial 3D object for limited angles. NeRF is then trained using inpainting-SDS (Score Distillation Sampling) [47] loss for diffusion model to create an image corresponding to text prompt to fill the remaining part of the 3D representation. For better reconstruction of the 3D object in the early stage of training, we propose progressive global view sampling. Lastly, with the refinement stage, we could minimize discrepancies among generated parts. By utilizing these novel techniques, our method has the ability to construct 3D objects from text-generated images or given images.

We demonstrated the effectiveness of our method in both image-to-3D and text-to-3D tasks. In the image-to-3D task, our method outperformed those of the current state-of-the-art (SOTA) baselines in terms of both multi-view consistency and source image correspondence visually and in user studies. As such, DITTO-NeRF's effectiveness is extended beyond that of previous models. In the text-to-3D task, our method excelled those of the current SOTA baselines in terms of both output fidelity and diversity. Importantly,

these improvements were achieved while requiring reasonable training time and computation resources. Here is the summary of our contributions:

- Proposing a novel pipeline for generating a 3D object model from a single image or text prompt, called DITTO-NeRF, that iteratively propagates high-quality partial 3D model on in-boundary (IB) angles to the remaining 3D model on outer-boundary (OB) angles.

- Proposing progressive global view sampling (PGVS) from IB to OB, reliability-guidance masking for IB, and multi-scale consistency refinement for all IB and OB.

- Outperforming prior arts on text/image-to-3D [67, 25, 68], achieving remarkable results in terms of diversity/quality and speed/fidelity, respectively.

# 2. Related Works

## 2.1. Text-to-image guidance

In recent years, generative models have made significant strides in the field of image generation, an area that was once thought to be the exclusive domain of humans. Numerous studies have explored the use of generative models, including Variational AutoEncoder (VAE) [7, 46, 44], flow-based models [8, 9], and GAN [2, 18, 27]. However, the SOTA models are based on the denoising diffusion probabilistic model [6, 53, 20, 21, 19] (hereinafter referred to as the diffusion model). The diffusion model is comprised with forward diffusion steps that adds pre-defined noise based on time steps and reverse generative steps that denoise noisy images. These diffusion models have achieved exceptional image quality, but their applicability has been limited due to their high computational resource requirements. The latent diffusion model (LDM) [48] has emerged as a promising approach for addressing the computational resource requirements associated with the forward/reverse process in image generation. This model leverages the encoding process of the VAE to process the forward and reverse processes in the latent space, rather than the RGB space, thereby significantly reducing the amount of computation required. By passing the final denoised latent vector through the VAE decoder, high-quality images can be obtained with reduced computational effort. The development of LDM has led to the emergence of several novel applications in image generation, such as inpainting models that incorporate additional inputs, such as masks, to manipulate the generated images based on the surrounding image context.

In this study, we leveraged an inpainting model based on LDM to achieve the goal of creating a 3D object that matches the user-input or generated image.
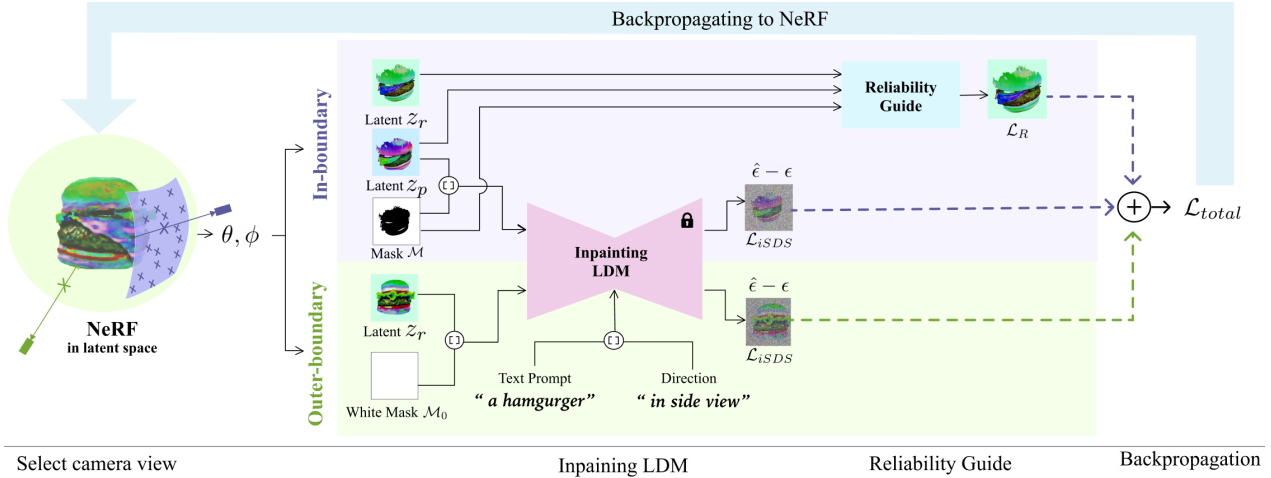
Figure 2: Our DITTO-NeRF pipeline for training from a 2D single image (given/generated). (Step 1) 2D latent $z_r$'s are sampled from latent NeRF at angles $\theta$ and $\phi$ using PGVS. (Step 2-1) Inpainting LDM with the given text and additional direction text yields the residual $\hat{\epsilon} - \epsilon$ from latent $z_r$ on OB or latent $z_p$ on IB. (Step 2-2) Latent $z_r$ on IB will be suppressed for outer-mask initially (so that latent $z_p$ will be relatively well-preserved for in-mask initially) and then progressively preserved for all area later with more reliable estimates.

## 2.2. Image-to-3D

Various methods exist for representing 3D objects, including point clouds [1, 33, 36, 70, 75], meshes [15, 74], and voxels [14, 52, 62]. However, these methods require significant storage capacity to represent high-quality 3D objects. Neural Radiance Fields (NeRF) [35] offers a novel solution for it by learning a neural network from images of certain viewpoints. NeRF has been criticized for its long training times, rendering times, and dependency on a large number of camera positions and images. Consequently, much research has focused on optimizing training and rendering times [69, 13, 72, 37, 22, 5, 58, 45, 16, 3]. Instant-NGP [37] using a multi-resolution hash grid encoding shows promising results. Additionally, several studies have explored the potential for achieving high-quality 3D representations with a limited number of images [5, 3, 73, 66, 25, 68] by leveraging the prior knowledge of pre-trained models.

By combining the LDM approach with Instant-NGP, we were able to efficiently create a high-quality, multi-view consistent 3D representation from a single image.

## 2.3. Text-to-3D

In recent years, the diffusion model has shown outstanding performance in text-to-image tasks, and several attempts have been made to extend it to text-to-3D object tasks. DreamFusion has achieved remarkable results by leveraging that the noise residual value obtained from one step of the denoising generative process in the Diffusion Denoising Probabilistic Model [6] alters image to match the prompt.

The model adds Gaussian noise to the image rendered by NeRF in a forward process and backpropagates the noise residual obtained through the reverse process as a gradient to align the rendered image with the input text prompt. By performing this process iteratively for random views, the 3D object corresponding to the input text prompt can be obtained.

In this paper, we propose DITTO-NeRF by adopting the main idea of DreamFusion. However, since DreamFusion uses Imagen [49], a proprietary diffusion model, we utilize Latent NeRF [34] which was based on Stable diffusion [48] as the fundamental model.

## 3. Methods

This chapter delineates the pipeline of our DITTO-NeRF, which comprises four distinct components: 1) The process of generating partial 3D objects from an image using LDM through an input text or a given image by a user is executed and the resulting objects are referred to as in-boundary objects (IB 3D objects). 2) NeRF employs progressive global view sampling (PGVS) and reliability-guided ($\mathcal{L}_R$) to effectively learn the frontal view, using the IB 3D object created through the aforementioned process. 3) $\mathcal{L}_R$ and inpainting SDS loss ($\mathcal{L}_{iSDS}$) are utilized to enhance the fidelity of the remaining frontal view. 4) A refinement stage is introduced to enhance the overall quality of the representation. The complete pipeline of DITTO-NeRF is presented in Fig. 2.

## 3.1. IB partial 3D object and pre-rendering

The diffusion model described previously can generate diverse images that correspond to a given text prompt. To further enhance the capabilities of the model for learning 3D representations that match the quality and diversity of the image generation model, we introduced a new process of creating a 3D object with a single image, and the object is named with IB 3D object. This process involves the creation of a partial 3D object that is utilized to aid NeRF in learning a 3D representation that aligns with the generated or user-given image.

**IB 3D object.** A relative depth map was first extracted from the image using a monocular depth estimation model called MiDaS [42]. Using this depth map, a 3D object can be composed with the form of a point cloud. It has been observed that in the point cloud constructed in this step, extraneous points such as backgrounds or floors may be inadvertently included in the point cloud. This can result in unwanted points being incorporated into the generated mesh during the conversion of the point cloud to mesh. To circumvent this, the method identifies outlier vertices based on their proximity to other vertices and subsequently removes them using a predetermined standard deviation. The values of these thresholds were determined heuristically via experimentation. Following this, Poisson surface reconstruction [28] is employed to generate the mesh from the completed point cloud. After the Poisson surface reconstruction process, the 3D mesh of the desired object is acquired by removing the parts with less density than a certain quantile.
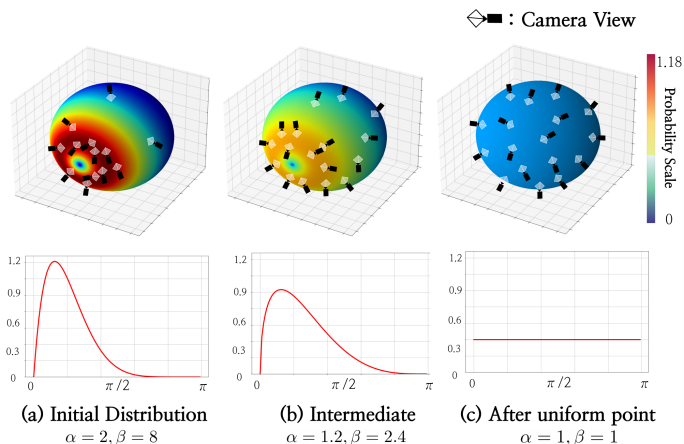


Figure 3: Examples of our PGVS over iterations based on varying Beta distribution. Initial camera view samples are heavily concentrated on IB (a) and then they are gradually spread to OB over iterations (b). Finally, camera view samples are uniformly distributed after uniform point.

**Setting IB and pre-rendering.** Rendering a 3D mesh in real-time for various viewpoints and using it for learning is computationally demanding. To mitigate this, the proposed method pre-renders the IB 3D object by sampling N views with uniform distribution within a limited angle range, which we call *in-boundary* (IB). Specifically, RGB images, depth, and latent vectors are rendered using the VAE encoder of the LDM.

## 3.2. Training NeRF using IB 3D object

**Progressive global view sampling (PGVS).** To enable the acquisition of the IB 3D object that was previously established and ensure coherence between the partial representation and subsequent portions generated through the diffusion model, it is imperative to initially sample numerous views of the in-boundary part. As the training progresses, views of the *outer-boundary* (OB) part need to be sampled to facilitate the generation of the overall 3D object. To accomplish this, we employ the beta distribution as the probability density function to sample the viewing position.

$$f(x;t) = \begin{cases} \frac{x^{\alpha(t)-1}(1-x)^{\beta(t)-1}}{B(\alpha(t),\beta(t))} & \text{if } t < t_u \\ U(0,1) & \text{if } t \geq t_u \end{cases} \quad (1)$$

with $\alpha(t) = \alpha_0 + 1 - \frac{\alpha_0}{t_u}t$ and $\beta(t) = \beta_0 + 1 - \frac{\beta_0}{t_u}t$.

The point which is named the *uniform point* $t_u$, is where $\alpha$ and $\beta$ decay to 1, after which uniform sampling is performed for all directions in the subsequent training steps which is illustrated in Fig. 3. Afterward, the process of obtaining the angle we want from the corresponding pdf is expressed by the following formula:

$$\theta, \phi = (\int f(x;t)dx)^{-1}(U_{\theta,\phi}(0,1)) \cdot range(\theta,\phi) \quad (2)$$

If the sampled camera position lies within the in-boundary region, the nearest pre-rendered position is selected instead of rendering a random view each time to avoid the unnecessary computational burden. Since the camera view is sampled progressively, we call this sampling Progressive global view sampling (PGVS).

## 3.3. Matching IB and OB 3D object

**Inpainting SDS loss.** Simultaneously applying the previously rendered prior images and the $\lambda_{iSDS}$ from the diffusion model can result in a conflict that depreciates the consistency of the IB 3D object and $\mathcal{L}_{iSDS}$ generated parts as the diffusion model may prioritize creating an object that corresponds to the built-in prior model. This can impede the goal of generating a 3D object that closely resembles the desired image. To address this, we employ a fine-tuned

diffusion model $\epsilon_\phi$ [48] for the inpainting task, rather than using a general diffusion model.

$$\mathcal{L}_{iSDS}(\phi, g(\theta)) = \mathbb{E}_{t,\epsilon}[||\epsilon_\phi(x_t; y, t, \mathcal{M}) - \epsilon_t||_2^2] \quad (3)$$

where $y$ is text embedding, $\mathcal{M}$ is the binary mask, $\epsilon_t$ is actual noise for time step $t$, and $x_t$ is noise image for each time step $t$. When the sampled viewpoint is in-boundary, the rendered latent vectors and pre-rendered latent images are optimized with $\mathcal{L}_R$. In addition, sparsity loss ($\mathcal{L}_{sp}$) was introduced, which was to obtain a cleaner representation by suppressing the lump caused by learning NeRF. For further details for $\mathcal{L}_{sp}$, see the supplementary material. The total loss is formulated as shown in the accompanying equation.

$$\mathcal{L}_{total} = \lambda_{iSDS}\mathcal{L}_{iSDS} + \delta_R\left(\theta, \phi\right)\mathcal{L}_R + \lambda_{sp}\mathcal{L}_{sp} \quad (4)$$

$\delta_R(\theta, \phi)$ here has a value of 1 within the IB area and a value of 0 otherwise.

**Reliability-guided loss.** When the camera position is sampled in the IB area, we optimize the loss between the rendered latent image in NeRF and the pre-rendered latent image. However, due to the difference in consistency and size between the IB 3D object and the generated object, we introduce a novel loss, which takes into account the differences between the part of the image corresponding to the desired object (foreground) and the remaining region (background), called the reliability-guided loss.

$$\mathcal{L}_R = [\zeta \cdot \mathcal{M} + \eta(t) \cdot (1 - \mathcal{M})] \odot ||z_r - z_p||_1 \quad (5)$$

with $\eta(t) = e^{-t/\lambda_\eta}$. $\lambda_\eta$ and $\zeta$ are constants that found with experiments. $z_r$ is rendered latent vector from NeRF and $z_p$ is pre-rendered latent vector.

The background of pre-rendered images is initialized as white. However, learning the latent ground truth as it leads to suboptimal inpainting results in the IB region, as the $\mathcal{L}_{iSDS}$ is smaller than the $\mathcal{L}_R$. On the other hand, if the diffusion model generates the background and the camera position is sampled from the back of the object, the model may not recognize the IB 3D object, resulting in inconsistent objects with varying sizes. Therefore, to suppress the unreliable background area at the beginning of training, we also train the white part of the background at the beginning. After that, learning proceeds, and since the IB 3D object is sufficiently represented in NeRF, the part generated by the $\mathcal{L}_{iSDS}$ becomes reliable, and the weight of the part against the background is exponentially reduced to generate the rest part.

Specifically, during the calculation of the $\mathcal{L}_R$, the loss is divided into foreground and background parts, with the foreground part given a higher weight and the background part given a relatively lower weight to initially create a white background. Subsequently, during training, the background
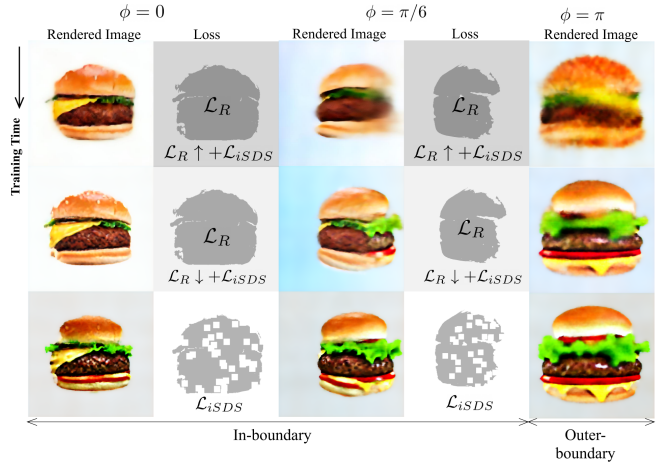


Figure 4: This figure illustrates how does the $\mathcal{L}_R$ and refinement works. At the beginning of the training process, $\mathcal{L}_R$ suppresses the unreliable part so that IB 3D object could be represented. As the training proceeds reliability of $\mathcal{L}_{iSDS}$ increases creating outer-boundary parts. In refinement procedure, random patches are applied to mask for enhancing overall quality.

weight is gradually decayed, allowing the $\mathcal{L}_{iSDS}$ for inpainting to have a more significant impact if the in-boundary region is learned after a certain amount of training. Through this process, we were able to train a NeRF model that created a 3D object that resembles an image prior with a constant size. This whole process is illustrated in Fig. 4 along with rendered training images for each step.

### 3.4. Refining details

**Random patch refinement.** Although the part generated by the $\mathcal{L}_{iSDS}$ and IB 3D object part continue semantically seamless, there is a noticeable discontinuity in the overall texture and color. This is attributed to the prior inside the diffusion model. To address this challenge and enable the created 3D object to exhibit continuity of texture and color while retaining the image prior of the front part as much as possible, a *refinement step* was introduced. This refinement step constitutes 10% of the entire training process, and if in-boundary is sampled during this step, the $\mathcal{L}_R$ is excluded, and a random patch is given to the foreground to refine the corresponding part with $\mathcal{L}_{iSDS}$. This process results in the color and texture of the generated 3D object becoming more similar to the part generated by the $\mathcal{L}_{iSDS}$, while maintaining the overall shape and content of the image prior. Through the refinement step, the desired objective of achieving continuity of texture and color while maintaining the image prior is attained. The corresponding method is shown in Fig. 4
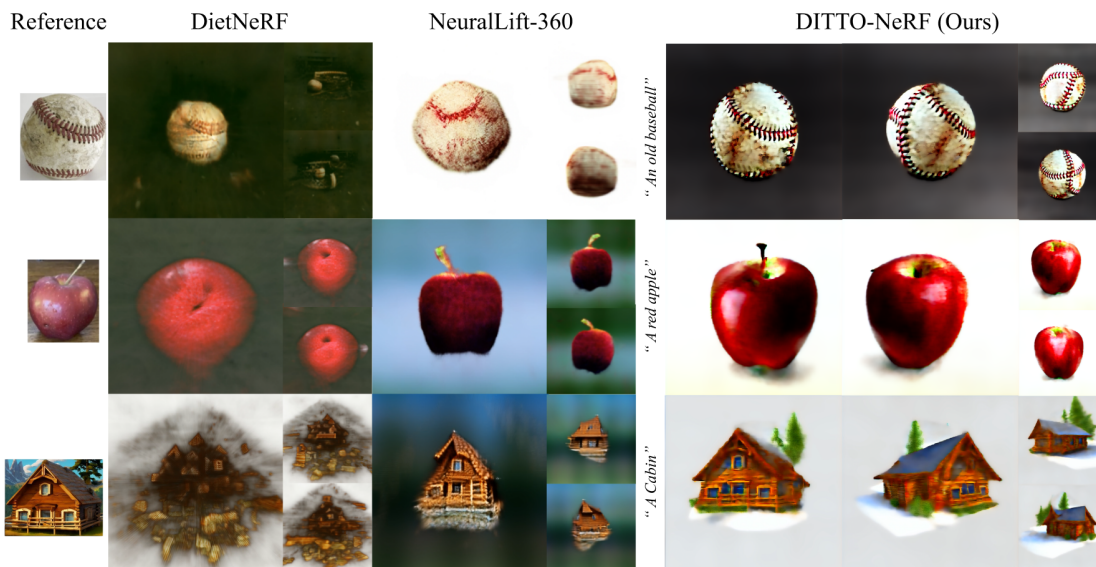
Figure 5: Qualitative comparison with other image-to-3D models [25, 68]. The outputs of the Neurallift-360 [68] were obtained and cropped from the Neurallift-360 website for fair comparison. Our model used the reference images and the corresponding texts on the left.
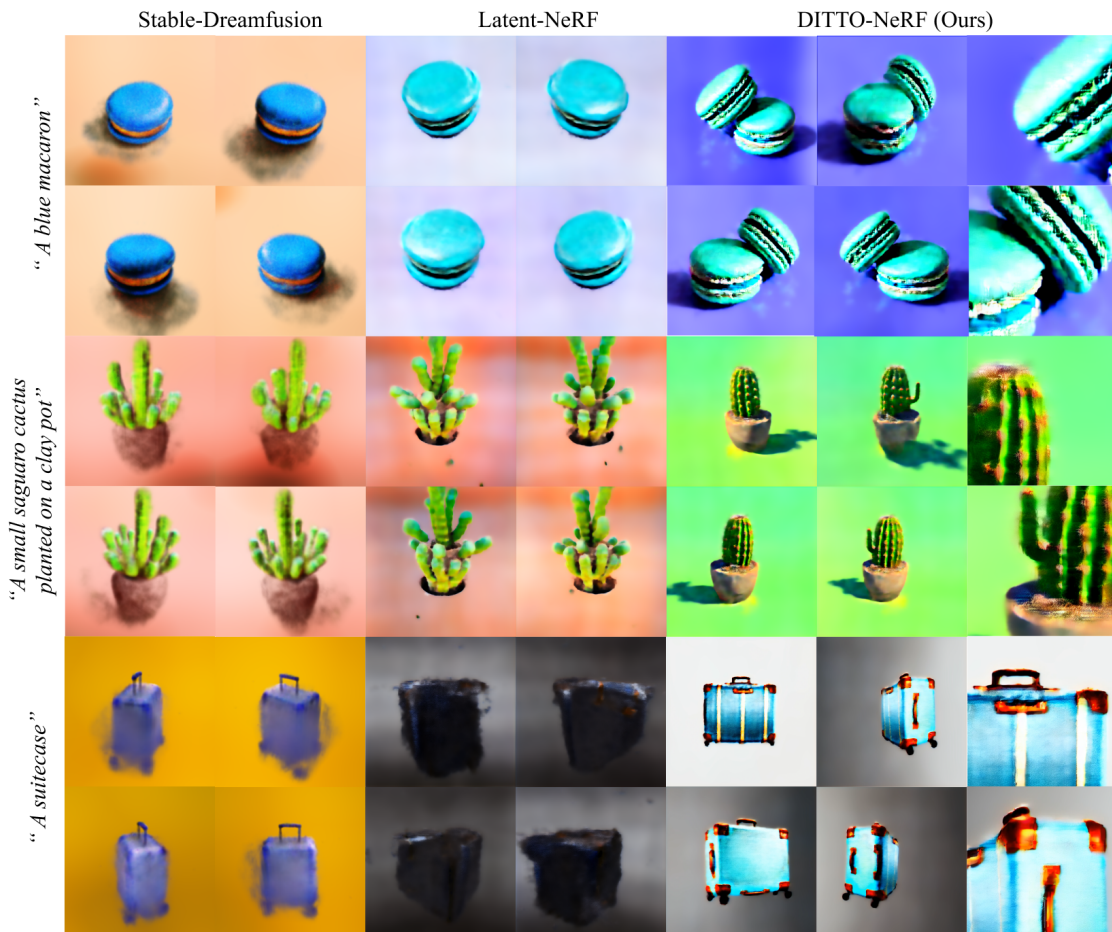


Figure 6: Qualitative comparison with other text-to-3D models [54, 34]. The last column is our model's zoomed-in results to show our excellent details. The number of iterations for the text-to-3D model to generate outputs was set to the default value in the baseline.

Figure 7: Qualitative comparison with other text-to-3D models [54, 34] in terms of diversity. All seeds were taken randomly.

| Method | Corr. ↑ | Fidelity ↑ |
|---|---|---|
| DietNeRF | 2.29 | 2.194 |
| SinNeRF | 2.084 | 2.031 |
| NeuralLift-360 | 3.221 | 3.24 |
| **DITTO-NeRF (Ours)** | **3.928** | **4.019** |

Table 1: Mean opinion score for evaluating image-to-3D models.

| Method | Diversity ↑ | Fidelity↑ | Corr.↑ | Time (m) ↓ |
|---|---|---|---|---|
| Stable DreamFusion | 3.704 | 2.6 | 2.959 | 60 |
| Latent-NeRF | 2.995 | 3.094 | 3.362 | **10** |
| **DITTO-NeRF (Ours)** | **3.966** | **4.158** | **4.242** | 25 |

Table 2: Mean opinion score and training time for evaluating text-to-3D models.

**Dimension refinement.** In addition, it was observed that the pipeline resulted in the appearance of jagged edges in the generated object, known as the *jaggies* phenomenon. This issue arises due to the training process, which involves feeding a relatively high-resolution image into a relatively low-resolution latent space. To address this, we introduced a refinement step that linearly doubles the NeRF rendering resolution, resolving the issue of jagged edges. Through this process, we were able to generate an object with clear edges while maintaining the overall shape and content of the IB 3D object. The color and texture of the generated object were also found to be similar to those produced by the $\mathcal{L}_{iSDS}$. Check the supplementary for the details.

## 4. Experiments

In this section, we aim to evaluate the effectiveness of our proposed method and compare it with other existing models in the domains of image-to-3D and text-to-3D. For image-to-3D, we perform a comparative analysis between our method and existing models with respect to source-generated 3D object correspondency and fidelity. Similarly, for text-to-3D, we compare the performance of our method with other open-sourced models in terms of diversity, computational efficiency, and fidelity. Additionally, we investigate the impact of various factors that we introduced on method section.

The user study's respondents were asked to rate each item on a 5-point scale. 3150 questionnaires were collected from a total of 210 people. All subsequent experiments were done on a single NVIDIA A100 GPU.

### 4.1. Generating 3D objects from a single image

In the present section, a comparative analysis is carried out between our proposed method and existing single-image NeRF models. To evaluate the effectiveness of our approach, a questionnaire survey is conducted to obtain N responses on source-generated object correspondence and fidelity. The survey responses provide insight into the performance of our method relative to the SOTA models in terms of its ability to accurately reconstruct objects and maintain correspondence with the source.

**User study.** A survey was carried out to evaluate the effectiveness of our proposed method in comparison to existing image-to-3D generative models, namely SinNeRF [67], DietNeRF [25] and Neurallift-360 [68]. The survey was designed to gather user feedback on the fidelity of the generated object and source image-object correspondence. The results in Table 1 indicate that our method outperforms the existing models in both categories, establishing a new SOTA in image-to-3D generative models.
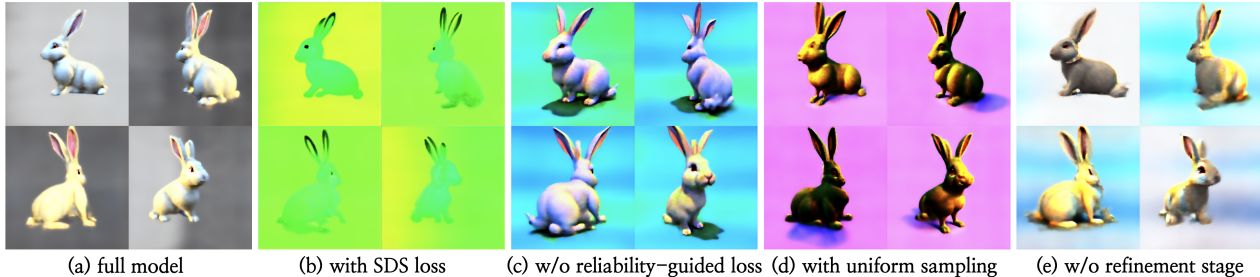
Figure 8: Ablation studies by generating 3D object with a text prompt *"a rabbit, animated film character, 3D rendered"* using full model (a), without i-SDS loss (b), without reliability-guided loss (c), with uniform sampling (not PGVS) (d), and without refinement stage (e).

## 4.2. Generating 3D objects from a text prompt

This section presents a comparative analysis of our proposed method with existing Stable diffusion-based text-to-3D generative models [54, 34]. The evaluation is conducted through a questionnaire survey that captures respondents' feedback on the diversity and fidelity of the generated 3D objects. The survey responses are scored on a 5-point scale for each item. Furthermore, to investigate the computational efficiency of the models, we report the time spent on the baseline's default setting.

**Diversity with a single prompt.**  Table 2 presents the results of our comparative analysis, indicating that our proposed method exhibits higher diversity in generated 3D objects as compared to the Stable diffusion-based text-to-3D generative models. An example of this diversity is illustrated in Fig. 7, where multiple outputs are generated for some given prompts.

**Fidelity and text correspondency.**  As observed in previous studies, the Stable diffusion-based text-to-3D generative models produce unrealistic, blurry, or unclear 3D objects. In contrast, our proposed method generates relatively realistic 3D objects, as evidenced by the results presented in Table 2, supported by users. Similarly, the correspondence between the input text and the generated object is observed to be more accurate and consistent in our method, owing to its superior performance in accurately capturing the semantic meaning of the input text.

**Computational efficiency.**  Stable-DreamFusion which is based on DreamFusion, suffered from the issue of requiring a long training time to achieve results with satisfactory quality. Latent-NeRF, on the other hand, addressed this issue by training the model on the latent space; however, it exhibited limitations in terms of diversity and fidelity. In contrast, our proposed method is capable of learning a 3D representation of satisfactory quality within a reasonable training time, as

illustrated in Table 2. While our method may have some limitations in fidelity when compared to models such as Magic3D or DreamFusion, it should be noted that these models utilize large and undisclosed models, limiting their applicability. In contrast, our proposed method offers an advantage in terms of computation.

### 4.3. Ablation study

In this section, we conduct an ablation study on the components of our method. We used inpainting SDS loss instead of conventional SDS loss to make the IB 3D object and diffusion-generated part continued smoothly. Otherwise, as shown in Fig. 8, it was obvious that the shape and color were lost. In the case where $\mathcal{L}_R$ does not exist, it was confirmed that the generated object was larger than IB 3D object or additional elements were added to it. There was a problem that the 3D prior could not be reconstructed when sampling viewpoints uniformly without using decaying beta distribution. Finally, if refinement was not performed, the semantic part continued, but discontinuity occurred in color and texture. Further detailed explanation is on the supplementary.

## 5. Limitation

Depth prediction is based on the shadows present in the image and the prior of the model. Therefore, if an image with minimal shadows or an image that is generated from outside of the distribution of the diffusion model is input, accurate depth estimation may not be possible. This can result in an inadequate IB 3D object being generated, leading to lower-quality 3D object outputs compared to those obtained in general. Additionally, if the quality of the image generated by the diffusion model is poor, learning can become more challenging as illustrated in Fig. 9. Also, in cases where a low-probability image is selected from the diffusion model, it may not follow the image prior, leading to a convergence of the 3D object within the diffusion model. This can also cause a discontinuity issue between the image prior and diffusion-generated parts.
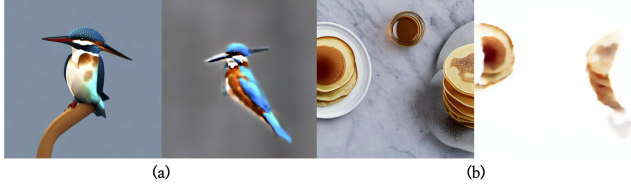
(a)                                  (b)

Figure 9: As the diffusion model failed to generate proper images, the generated 3D objects' quality are also limited.

## 6. Conclusion

We introduce DITTO-NeRF, a novel approach for obtaining a 3D representation from a single image or text input. DITTO-NeRF leverages an IB 3D object built from images obtained from user input or text, which is then used to train a NeRF model that generates 3D objects using a diffusion model. Based on user studies and qualitative comparisons, we draw the following conclusions: our proposed model achieves higher fidelity and quality in terms of image-to-3D correspondence than existing works. Moreover, in the context of text-to-3D generation, our method offers both higher fidelity and more computation-efficient diversity than existing Stable diffusion-based studies.

## Acknowledgements

# DITTO-NeRF: Diffusion-based Iterative Text To Omni-directional 3D Model (Supplementary Material)

## A. Additional Results

### A.1. Videos

We present a video showcasing a comparison of the results obtained using our proposed DITTO-NeRF model, which is a novel method for learning 3D representations from either image or text inputs, against other baselines in the same task. For video, please check the folllowing page: janeyeon.github.io/ditto-nerf.

## B. Details on Methods

### B.1. Algorithms

---

**Algorithm 1:** Overall pipeline

**Input:** $y$, $I$
**Output:** $\mathbf{V}_{out}$

   // Training & Refinement
1 **if** $\exists! \ I$ **then**
2    $y_{image} = CLIP((y, y_{bg}))$
3    $z_{image} \leftarrow$ inpainting-LDM$(y_{image})$
4    $I \leftarrow \mathcal{D}(z_{image})$
5 $\mathbb{Z}_p, \mathbb{M} \leftarrow$ Generate-IB-3D$(I)$
6 **for** $i = 1, 2, ..., t_{total}$ **do**
7    $\theta, \phi \leftarrow$ PGVS$(i)$
8    **if** $IB$ **then**
9       $z_p, \mathcal{M} \leftarrow$ Find-closest$(\mathbb{Z}_p, \mathbb{M}, \theta, \phi)$
10    $\mathbf{o}, \mathbf{d} \leftarrow$ Render-ray$(\theta, \phi)$
11    $z_r, \ \mathcal{P}_{ws} \leftarrow \mathcal{N}(\mathbf{o}, \mathbf{d})$
12    **if** $i \leq i_{refine}$ **then**
13       Training$(\mathcal{N}, z_p, \mathcal{M}, z_r, \mathcal{P}_{ws}, y)$
14    **else**
15       Refinement$(\mathcal{N}, z_p, \mathcal{M}, z_r, \mathcal{P}_{ws}, y, i)$

   // Rendering output
16 $\mathbf{V}_{out} \leftarrow \{\}$
17 **for** $\theta, \phi$ in rendering angles **do**
18    $\mathbf{o}, \mathbf{d} \leftarrow$ Render-ray$(\theta, \phi)$
19    $z_r \leftarrow \mathcal{N}(\mathbf{o}, \mathbf{d})$
20    $I_{out} \leftarrow \mathcal{D}(z_r)$
21    Append $I_{out}$ to $\mathbf{V}_{out}$

---

**Overall pipeline.** The overall algorithm of DITTO-NeRF is described in Algorithm 1. First, it requires text input $y$ and additional image input $I$ from the user. In the case of

image-to-3D synthesis, the algorithm receives a reference image $I$ and corresponding text value $y$ as input. Whereas in the case of text-to-3D synthesis, only text input is required. If no reference image is available, two text values $y, y_{bg}$ are concatenated to create a text embedding $y_{image}$ that is used to generate a latent vector $z_{image}$ in inpainting-LDM. $z_{image}$ is then passed through a decoder $\mathcal{D}$ to create a reference image $I$. $I$ is fed into the Generate-IB-3D function to obtain latent vector $z_p$ and mask $\mathcal{M}$ for randomly selected $N$ angles through uniform distribution within the in-boundary (IB) area.

For $N_{total}$ iterations, the algorithm performs the following procedures: the camera view position, $\theta$, and $\phi$ are determined using the PGVS function along the iteration. The IB flag is then used to determine whether the angle values are used directly or switched to pre-selected angles within the in-boundary area. Using the Find-closest function, $z_p$ and $\mathcal{M}$ values are obtained for each pre-selected angle. The Render-ray function is used to obtain the $\mathbf{o}$ and $\mathbf{d}$ values for each ray based on the $\theta$ and $\phi$. These values are used by the NeRF model ($\mathcal{N}$) to render the latent $z_r$ and weight values $\mathcal{P}_{ws}$, respectively, for each point of the ray. If the current iteration $i$ is less than $i_{refine}$, Training is performed, and if it is larger, Refinement is executed.

As the $\mathcal{N}$ is trained, the latent vector $z_r$ is extracted from the $\mathcal{N}$ along the angles corresponding to 'rendering angles'. At each rendering angle, the rendered image $\mathcal{D}(z_r)$ is obtained from the latent vector $z_r$ using the decoder $\mathcal{D}$. By collecting rendered images $I_{out}$, we can obtain the final results, $\mathbf{V}_{out}$, a video of a 3D rendered object. It should be noted that the output of the $\mathcal{N}$ is a 4-channel latent vector, unlike the RGB output in standard NeRF. So the $z_r$ should be passed through the $\mathcal{D}$ used by inpainting-LDM in the last rendering step to get an RGB image.

**Training procedure.** As described above, when the current iteration $i$ is less than $i_{refine}$, Training is performed. The entire algorithm of the Training stage is described in Algorithm 2. During this stage, the input text $y$ is concatenated with direction information corresponding to the angle and put into a CLIP to get $y_{train}$. The subsequent process differs depending on whether it is IB or not. In the case of IB, the $\mathcal{L}_{iSDS}$ is calculated by inputting the $z_p$ and $\mathcal{M}$ values into the inpainting-LDM. The difference between $z_r$ and $z_p$ is then corrected using the Reliability-Guide

**Algorithm 2:** Training

**Input:** $\mathcal{N}, z_p, \mathcal{M}, z_r, \mathcal{P}_{ws}, y$
**Output:** $\mathcal{N}$

1 **Function** Reliability-Guide($z_p, z_r, \mathcal{M}$):
2    **return** $[\zeta \cdot \mathcal{M} + \eta(t) \cdot (1 - \mathcal{M})] \odot ||z_r - z_p||_1$

3 $y_{train} \leftarrow CLIP((y, y_{dir}))$
4 **if** $IB$ **then**
5    $\mathcal{L}_{iSDS} \leftarrow$ inpainting-LDM($z_p, y_{train}, \mathcal{M}$)
6    $\mathcal{L}_R \leftarrow$ Reliability-Guide($z_p, z_r, \mathcal{M}$)
7 **else**
8    $\mathcal{M}_0 \leftarrow J_{H \cdot f_{scale}}$
9    $\mathcal{L}_{iSDS} \leftarrow$ inpainting-LDM($z_r, y_{train}, \mathcal{M}_0$)
10 $\mathcal{L}_{sp} \leftarrow$ Sparsity($\mathcal{P}_{ws}$)
11 $\mathcal{L}_{total} \leftarrow \lambda_{iSDS}\mathcal{L}_{iSDS} + \delta_R(\theta, \phi)\mathcal{L}_R + \lambda_{sp}\mathcal{L}_{sp}$
12 Update $\mathcal{N}$ with $\mathcal{L}_{total}$

---

**Algorithm 3:** Refinement

**Input:** $\mathcal{N}, z_p, \mathcal{M}, z_r, \mathcal{P}_{ws}, y, i$
**Output:** $\mathcal{N}$

1 **Function** Dimension-Refine($iter$):
2    $rate \leftarrow (iter - t_{total} \cdot f_{ref})/t_{total} \cdot (1 - f_{ref})$
3    $H, W \leftarrow 64 \times (1 + rate), 64 \times (1 + rate)$
4    **return** $H, W$

5 **Function** Add-Patch($\mathcal{M}$):
6    **for** $i = 0, 1, ..., N_{patch}$ **do**
7      $P_x = $ random($S_{mask}$)
8      $P_y = $ random($S_{mask}$)
9      $\mathcal{M}[P_x : P_x + S_{patch}, P_y : P_y + S_{patch}] \leftarrow 1$
10    **return** $\mathcal{M}$

11 $y_{refine} \leftarrow CLIP((y, y_{dir}))$
12 $H, W \leftarrow$ Dimension-Refine($i$)
13 **if** $IB$ **then**
14    $\mathcal{M} \leftarrow$ resize($H, W$)
15    $\mathcal{M} \leftarrow$ Add-Patch($\mathcal{M}$)
16    $z_p \leftarrow$ resize($H, W$)
17    $\mathcal{L}_{iSDS} \leftarrow$ inpainting-LDM($z_p, y_{refine}, \mathcal{M}$)
18 **else**
19    $\mathcal{M}_0 \leftarrow J_{H \cdot f_{scale}}$
20    $\mathcal{L}_{iSDS} \leftarrow$ inpainting-LDM($z_r, y_{refine}, \mathcal{M}_0$)
21 $\mathcal{L}_{sp} \leftarrow$ Sparsity($\mathcal{P}_{ws}$)
22 $\mathcal{L}_{total} \leftarrow \lambda_{iSDS}\mathcal{L}_{iSDS} + \lambda_{sp}\mathcal{L}_{sp}$
23 Update $\mathcal{N}$ with $\mathcal{L}_{total}$

---

function, which adjusts the background and foreground values using a $\mathcal{M}$ and expands the size of the reliable region as learning proceeds. When it is not IB, network training is performed using the inpainting-LDM with $z_r$ and $\mathcal{M}_0$ values. After obtaining the $\mathcal{L}_{sp}$ to eliminate the lump, back-propagation is performed by adding all of the losses into $\mathcal{L}_{total}$.

**Refinement procedure.** When the current iteration $i$ is greater than $i_{refine}$, a refinement process is performed as outlined in Algorithm 3. Unlike the previous training step, the Reliability-Guide process is omitted, and the Dimension-Refine and Add-Patch processes are added. The process of obtaining $y_{refine}$ is the same as the previous step, and the Dimension-Refine function calculates $H$ and $W$ based on the current iteration $i$. This function linearly increases the size of $H$ and $W$ from 64 to 128 as the refinement process progresses, for improvement of the resolution and quality in rendered views. To match the corresponding sizes of $H$ and $W$, the pre-calculated $\mathcal{M}$ and $z_p$ sizes are adjusted before entering. The Add-Patch process enters only the $\mathcal{M}$ in the IB area and maintains consistency of the front and back sides by randomly arranging $S_{patch}$ size of $N_{patch}$ patches in a mask of $S_{mask}$ size. Subsequently, $\mathcal{L}_{iSDS}$ and $\mathcal{L}_{sp}$ are added to $\mathcal{L}_{total}$, and the $\mathcal{N}$ is trained in the same way as in the training process.

### B.2. IB 3D object pre-rendering

To pre-render an IB 3D object, the angle of the in-boundary (IB) must be established first. In instances where this angle is too narrow, the matching between the IB and outer-boundary (OB) is insufficient. Whereas if the angle is too wide, additional views must be sampled, which can cause the bottleneck problem of Find-closest process in Algorithm 1. Through several experiments, we discovered

the setting $\phi \in [-30, 30]$ and $\theta \in [60, 120]$ provides optimal conditions. Moreover, if the number of pre-rendered images is too small, IB training cannot be performed sufficiently, while a large number of images increases the computational burden. Thus, based on our experimental findings, we determined that setting N=64 provides optimal conditions.

### B.3. Camera intrinsic parameters for point cloud

In order to generate a point cloud from an RGB-D image, it is essential to know the camera intrinsic parameters. For real-world cameras, these values are readily available; however, as our images are synthesized using a diffusion model, the actual parameters remain unknown. The required intrinsic parameters include skew rate, $c_x$, $c_y$, $f_x$, and $f_y$. Given that contemporary cameras exhibit minimal skew, we approximated the skew rate to be nearly zero. Additionally, considering the generated image's dimensions of $512^2$, we assigned $c_x = c_y = 256$. The most important value was the focal length, that is, $f_x$ and $f_y$. As these values are typically quite similar, we assumed equivalence and conducted multiple investigations to verify this assumption.

Fig. S1 illustrates the incorporation of an additional prompt for lens focal length alongside the three existing prompts: object, portrait, and landscape. The focal length
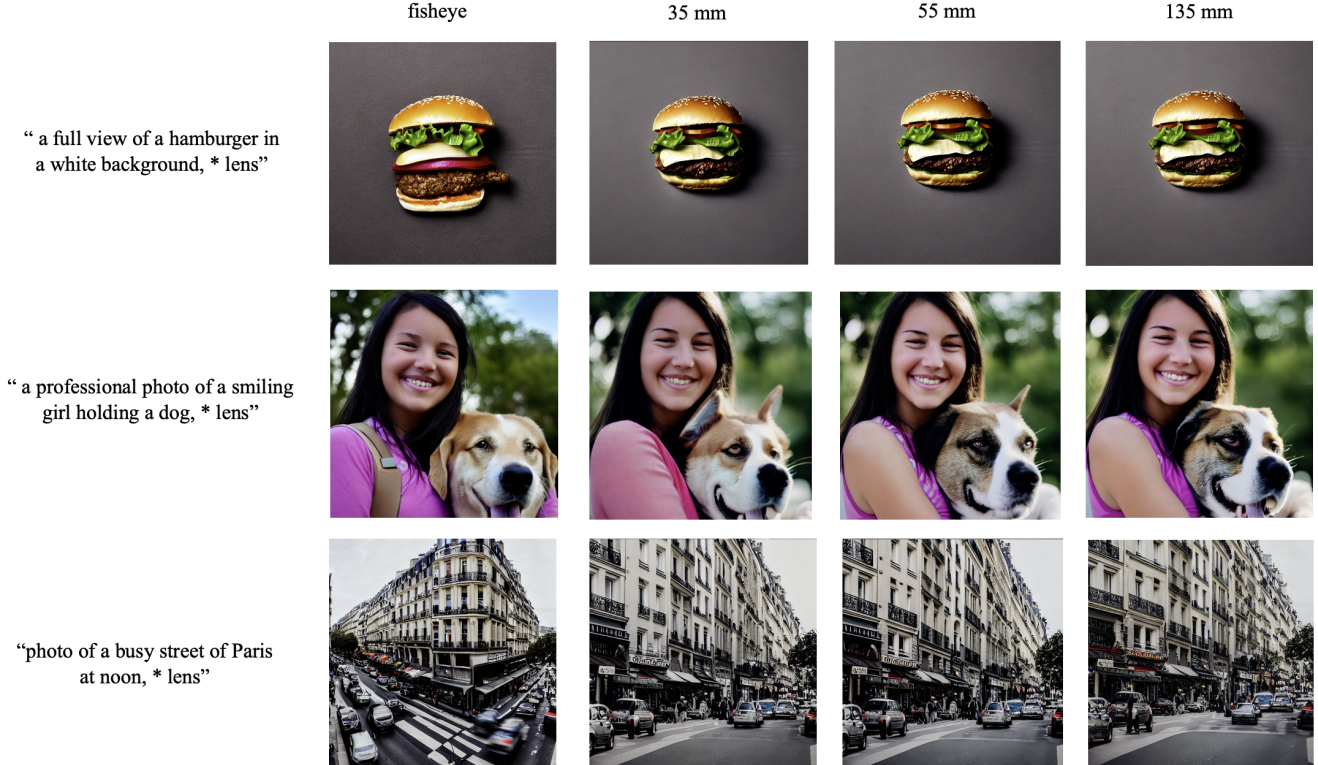
|  | fisheye | 35 mm | 55 mm | 135 mm |

Figure S1: Focal length test using Stable diffusion [48]. We observed differences between generated images with additional lens information. There are no substantial differences in the result.

prompts employed in this instance include "fish eye," "35mm lens," "50mm lens," and "135mm lens". Upon examining the outcomes, we observed no substantial differences in the results, except the landscape category, particularly the "fish eye" prompt. This phenomenon is attributed to the majority of training data falling within the range of 35mm to 50mm focal lengths. Consequently, we generated point clouds assuming a focal length of approximately 45mm, yielding satisfactory results for the majority of tasks necessitating the generation of object-related images.

### B.4. Progressive global view sampling (PGVS)

The beta distribution for PGVS is initialized with $\alpha_0 = 2$ and $\beta_0 = 8$. Thereafter, both constants decrease linearly until the $t_u$ iteration, which is defined as 0.3 times the total number of iterations $t_{total}$. After $t_u$, both constants are fixed at a value of 1 for the remaining iterations. For the experiments reported in this paper, we set $t_{total}$ to 5000, resulting in $t_u$ being equal to 1500.

### B.5. Sparsity loss

The sparsity loss used in this study was introduced in the Dreamfields [24] and was also employed in Latent-NeRF. It has a similar form to the binary cross-entropy loss and is designed to encourage the density values to converge to either

0 or 1. This helps improve the quality of the generated 3D representation by ensuring that the rays either pass through completely or are blocked. The mathematical expression for the sparsity loss is shown below.

$$\mathcal{L}_{sp} = -\mathbb{E}[\alpha log(\alpha) + (1 - \alpha)log(1 - \alpha)] \quad \text{(S1)}$$

with $\alpha$ is the weighted sum of a ray which is clipped with $[\epsilon, 1 - \epsilon]$, where $\epsilon = 10^{-5}$ and $\lambda_{sp} = 5 \cdot 10^{-5}$.

### B.6. Reliability-guided loss

The equation for reliability-guided loss is on the following:

$$\mathcal{L}_R = [\zeta \cdot \mathcal{M} + \eta(t) \cdot (1 - \mathcal{M})] \odot ||z_r - z_p||_1 \quad \text{(S2)}$$

with $\eta(t) = e^{-t/\lambda_\eta}$. And all the experiments proceeded with $\lambda_\eta = 8, \zeta = 2$.

### B.7. Patch refinement

Patch refinement is a crucial method to ensure that the final 3D representation has uniformity in terms of overall tone and texture. This refinement applies $N_{patch}$ patches of size $k \times k$ to a mask with size $H \times W$. The choice of patch size $k$ is important. In the case of large $k$, it shows the unintended result such as a change of a semantic part of the existing IB.

When $k$ goes too small, the effect of patch refinement would disappear while downsizing the mask. Similarly, the number of patches $N_{patch}$ is also important, as an insufficient number of patches may not produce the appropriate refinement, while an excessive number of patches may not maintain the existing IB. Therefore, we selected the values of $k = 16$ and $N = 256$ for this experiment.

## C. Implementation Details

### C.1. Monocular depth estimation

The image/text-to-3D task at hand necessitates the ability to extract accurate relative depth maps from a set of images that consists of diverse objects. Hence, it was crucial to employ a robust and accurate depth estimation model that had been trained on a variety of datasets, and for this purpose, we adopted MiDaS [43]. MiDaS was trained using pre-existing models from 12 different datasets, including ReDWeb [63], DIML [29], Movies, MegaDepth [31], WSVD [56], TartanAir [60], HRWSI [64], ApolloScape [23], BlendedMVS [71], IRS [59], KITTI [17], and NYU Depth V2 [38]. The largest model provided by MiDaS 'DPT BeiT Large' [41] which offers the highest quality depth estimation among the available models, was employed. This model utilizes a transformer architecture, which enables more precise and detailed depth estimation as compared to convolutional structures.

### C.2. IB 3D construction based on depth map

An RGB-D image can be generated by combining a depth map obtained from monocular depth estimation with an image, followed by the creation of a point cloud using the Open3D [76] library. To eliminate outliers, points are removed if their distance from the five surrounding points exceeds one standard deviation. The normal vector of each vertex is then estimated from the resulting point cloud, and a mesh is created through Poisson surface reconstruction, with a depth value of 10. Vertices with a density below 0.1 quantiles are subsequently removed, based on the assumption that the object of interest is expected to exhibit a relatively high density. An example of the process's output is shown in Fig. S2

### C.3. Stable diffusion inpainting

We utilized Stable diffusion inpainting [48] for inpainting-LDM to train NeRF in OB. The model was fine-tuned with original Stable diffusion. The following steps shows how it was fine-tuned: To initialize the Stable diffusion inpainting model, the creater used the weights from the Stable diffusion-v-1-2 model. The training process consisted of two phases: regular training for 595k steps followed by 440k steps of training with inpainting task at a resolution of $512 \times 512$ using the 'LAION-Aesthetics v2 5+' dataset [51]. To improve the classifier-free guidance sampling, the text-conditioning



(a) Point cloud from RGB-D    (b) Outlier removed

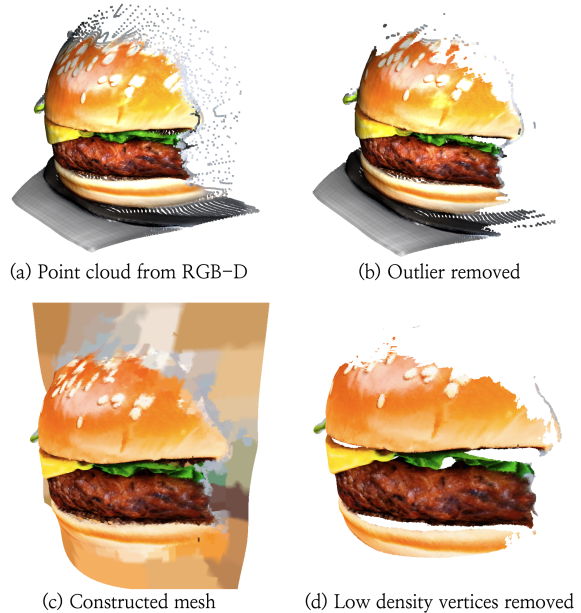(c) Constructed mesh    (d) Low density vertices removed

Figure S2: Outputs of each step in IB 3D construction. An example of established point cloud from RGB-D image is in (a). Representation of point cloud after removing outliers is shown in (b). By Poisson reconstruction, mesh (c) is created. And lastly, after removing low density vertices, we get final result like (d).

was dropped by 10%. For inpainting, the UNet was augmented with five additional input channels - four for the encoded masked-image and one for the mask itself. These additional channels were zero-initialized after restoring the non-inpainting checkpoint. During the training process, synthetic masks were generated and 25% of the images were masked in each iteration.

### C.4. Prompt conditioning

Monocular depth estimation is typically carried out by leveraging the shadow information present in an image, along with prior knowledge incorporated into the model. As a consequence, the depth map generated by such methods may be inadequate for images that lack substantial shading. Furthermore, it has been observed that the depth estimation process encounters difficulties when the background is either complex or positioned in close proximity to the subject. To overcome these issues, we propose a preprocessing step in which the initial image used for generating an IB 3D object is accompanied by the phrase "A whole photo of $\sim$" at the beginning of the prompt for non-cropped image and "$\sim$ in the white background taken with 50mm lens" at the end of the prompt. This ensures that the depth estimation process proceeds smoothly, resulting in an acceptable IB 3D object.

## D. Experimental Details

### D.1. Evaluation details

**Baselines.** In this study, we conducted a comparative analysis of existing image-to-3D models and text-to-3D models against our proposed DITTO-NeRF model. For the image-to-3D model, we evaluated our results against the current state-of-the-art model, NeuralLift-360 [68], for which the code has not been made publicly available. Therefore, we compared our results with the reference images provided on the NeuralLift-360 webpage [65]. Since the size of the NeuralLift-360 results was small, we cropped them to ensure a fair quality comparison. Regarding the text-to-3D model, we compared our results with two existing models based on Stable diffusion: Stable-Dreamfusion [54] and Latent-NeRF [34]. Stable-Dreamfusion replaces the generative diffusion model with Stable diffusion instead of Imagen which was used for the original Dreamfusion [39]. In contrast, Latent-NeRF learns the 3D representation directly in the latent space by learning the gradient through Stable diffusion in latent space. During the evaluation, the rendered latent vector is passed through the decoder of a Variational AutoEncoder (VAE) to extract the final RGB image, thereby achieving a fast learning time. All the evaluation of Stable-dreamfusion was conducted on 20000 iterations because the 10000 iterations' output quality was not enough for proper comparison and Latent-NeRF, it was 5000 iterations. Our work was conducted based on the code provided by Latent-NeRF.

**User study.** In this study, a total of 3,150 responses were collected through administering a survey of 15 questions to 210 participants. The Image-to-3D evaluation was performed using four different models, namely, SinNeRF [67], DietNeRF [25], NeuralLift-360 [68], and our DITTO-NeRF about three different images, namely, "baseball" "apple" and "hydrant" all of which were sourced from the NeuralLift-360 webpage. The evaluation of the image-to-3D object correspondence and fidelity was carried out with respect to the models' outcome and scored on a five-point scale. Additionally, the text-to-3D task was evaluated for diversity using three prompts, namely, "a loaf of bread", "a small saguaro cactus planted on a clay pot" and "a single candle burning on an ornate silver candlestick". Outputs were generated by Stable-Dreamfusion, Latent-NeRF, and our DITTO-NeRF with three different seeds for each of the three prompts. The evaluation procedure was similar to that used in the image-to-3D evaluation and the diversity of the outputs was scored on a five-point scale. Lastly, the fidelity of the outputs in the text-to-3D task was evaluated using three different prompts, namely, "a hamburger" "an astronaut" and "a suitcase" with the same three models. The results of each evaluation item were averaged to create mean opinion scores.

## E. Additional Ablation Studies

### E.1. Effectiveness of Dimension refinement
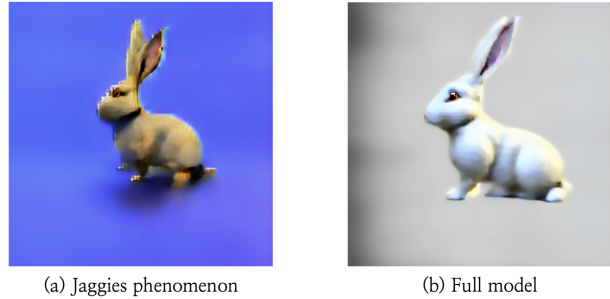


(a) Jaggies phenomenon          (b) Full model

Figure S3: An example of jaggies problem is shown in (a). With Dimension refinement, we can see there is no jaggies problem in (b).

If the dimension refinement stage does not exist, we can see 'Jaggies' phenomenon where the boundary looks like a saw, as shown in Fig. S3. As the cause of this, we identified that it is due to the projection of high-frequency information such as pre-rendered images into a low dimension such as latent space. We devised a method to increase the rendering dimension of NeRF to alleviate these artifacts. However, there was a problem that the training time increased rapidly when continuously using a large rendering dimension. Therefore, we devised a method of linearly increasing the rendering dimension during the refinement stage, finally doubling it, focusing on the fact that coarse shapes are generated even at low rendering dimensions. In actual implementation, the rendering dimension starting at $64$ starts to increase linearly in the refinement stage and finally reaches $128$ and learning ends.

### E.2. Effectiveness of PGVS

We tried other types of distributions for view sampling such as simple uniform sampling, 'moving beta distribution' and 'discrete accumulation distribution'. In this chapter, we would like to explain the two distributions mentioned above and their results.

**Moving beta distribution.** In our full model, the Probabilistic Gradient Vector Sampling (PGVS) method is employed to optimize the camera placement. In this method, the values of $\alpha$ and $\beta$ of the beta distribution decrease simultaneously and eventually converge to a value of 1. However, in the modified distribution used in this study, the two values gradually change to each other's first value. During the initial phase of training, many camera positions are sampled in the frontal part, and as the training progresses, a larger number of samples are taken in the back part. Similar to PGVS, the modified distribution also has a uniform point, after which the camera view is sampled uniformly. The results of this

(a) reference image

(b) full model

(c) frontal view of Mov. beta

(d) back view of Mov. beta

(e) frontal view of Discrete.
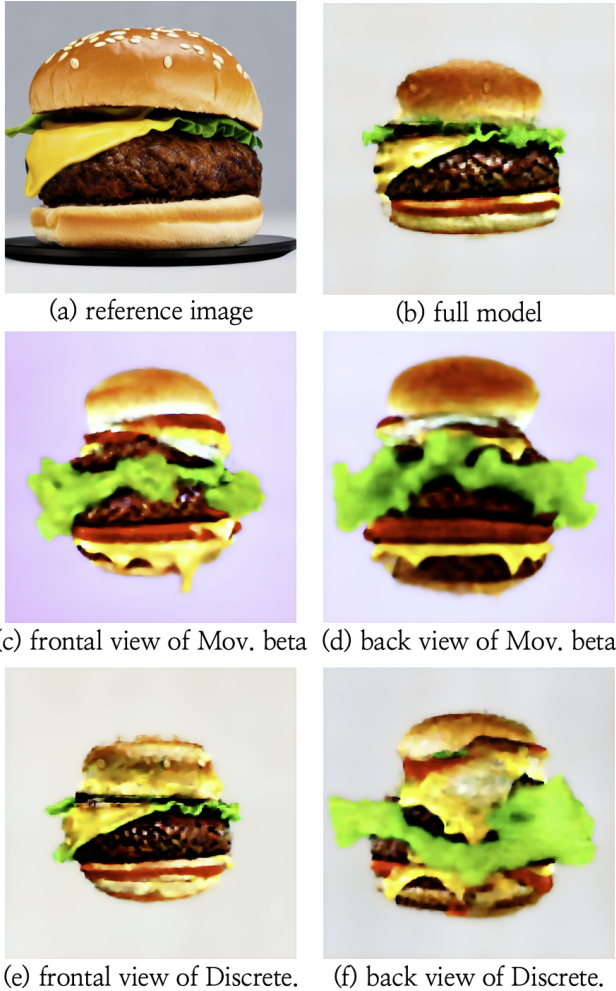
(f) back view of Discrete.

Figure S4: Comparison of results with different camera view sampler. (a) is a reference image generated by latent diffusion model with a prompt "a hamburger". (b) is the result of our full model, DITTO-NeRF. (c) is the frontal view of the result with Moving beta distribution sampling. (d) shows back view of the Moving beta distribution sampling. And (e) shows the frontal view of the result when Discrete accumulation sampling was utilized. (f) refers to the back view of the results with Discrete accumulation view sampling.

approach are presented in Fig. S4, where it can be observed that the frontal reconstruction did not work as effectively as in the full model.

**Discrete accumulation distribution.** The sampler used in this study is similar to the moving beta distribution described earlier, but it differs in that it has a discrete distribution. This sampler operates within a pre-defined interval, where the probability within a specific interval is set to $r$, and the probability in the remaining intervals is set to $1 - r$. At each specific iteration, the interval with the probability of $r$ is

passed on to the next iteration. Although the IB part was reconstructed effectively using this sampler, the back part did not converge to a normal shape. To address this issue, an experiment was conducted where $r$ was set to 0.65, as shown in Fig. S4.

## References

[1] Panos Achlioptas, Olga Diamanti, Ioannis Mitliagkas, and Leonidas Guibas. Learning representations and generative models for 3d point clouds. In *International conference on machine learning*, pages 40–49. PMLR, 2018.

[2] Andrew Brock, Jeff Donahue, and Karen Simonyan. Large scale gan training for high fidelity natural image synthesis. *arXiv preprint arXiv:1809.11096*, 2018.

[3] Anpei Chen, Zexiang Xu, Fuqiang Zhao, Xiaoshuai Zhang, Fanbo Xiang, Jingyi Yu, and Hao Su. Mvsnerf: Fast generalizable radiance field reconstruction from multi-view stereo. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 14124–14133, 2021.

[4] Kangle Deng, Andrew Liu, Jun-Yan Zhu, and Deva Ramanan. Depth-supervised NeRF: Fewer views and faster training for free. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2022.

[5] Kangle Deng, Andrew Liu, Jun-Yan Zhu, and Deva Ramanan. Depth-supervised nerf: Fewer views and faster training for free. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12882–12891, 2022.

[6] Prafulla Dhariwal and Alex Nichol. Diffusion models beat gans on image synthesis. *arXiv preprint arXiv:2105.05233*, 2021.

[7] P Kingma Diederik, Max Welling, et al. Auto-encoding variational bayes. In *Proceedings of the International Conference on Learning Representations (ICLR)*, volume 1, 2014.

[8] Laurent Dinh, David Krueger, and Yoshua Bengio. Nice: Nonlinear independent components estimation. *arXiv preprint arXiv:1410.8516*, 2014.

[9] Laurent Dinh, Jascha Sohl-Dickstein, and Samy Bengio. Density estimation using real nvp. *arXiv preprint arXiv:1605.08803*, 2016.

[10] Jean-Denis Durou, Maurizio Falcone, and Manuela Sagona. Numerical methods for shape-from-shading: A new survey with benchmarks. *Computer Vision and Image Understanding*, 109(1):22–43, 2008.

[11] Haoqiang Fan, Hao Su, and Leonidas J Guibas. A point set generation network for 3d object reconstruction from a single image. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 605–613, 2017.

[12] Paolo Favaro and Stefano Soatto. A geometric approach to shape from defocus. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 27(3):406–417, 2005.

[13] Sara Fridovich-Keil, Alex Yu, Matthew Tancik, Qinhong Chen, Benjamin Recht, and Angjoo Kanazawa. Plenoxels: Radiance fields without neural networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5501–5510, 2022.

[14] Matheus Gadelha, Subhransu Maji, and Rui Wang. 3d shape induction from 2d views of multiple objects. In *2017 International Conference on 3D Vision (3DV)*, pages 402–411. IEEE, 2017.

[15] Jun Gao, Tianchang Shen, Zian Wang, Wenzheng Chen, Kangxue Yin, Daiqing Li, Or Litany, Zan Gojcic, and Sanja Fidler. Get3d: A generative model of high quality 3d textured shapes learned from images. In *Advances In Neural Information Processing Systems*, 2022.

[16] Stephan J Garbin, Marek Kowalski, Matthew Johnson, Jamie Shotton, and Julien Valentin. Fastnerf: High-fidelity neural rendering at 200fps. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 14346–14355, 2021.

[17] Andreas Geiger, Philip Lenz, Christoph Stiller, and Raquel Urtasun. Vision meets robotics: The kitti dataset. *International Journal of Robotics Research (IJRR)*, 2013.

[18] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial networks. *Communications of the ACM*, 63(11):139–144, 2020.

[19] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *arXiv preprint arXiv:2006.11239*, 2020.

[20] Jonathan Ho, Chitwan Saharia, William Chan, David J Fleet, Mohammad Norouzi, and Tim Salimans. Cascaded diffusion models for high fidelity image generation. *J. Mach. Learn. Res.*, 23(47):1–33, 2022.

[21] Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance. *arXiv preprint arXiv:2207.12598*, 2022.

[22] Tao Hu, Shu Liu, Yilun Chen, Tiancheng Shen, and Jiaya Jia. Efficientnerf efficient neural radiance fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12902–12911, 2022.

[23] Xinyu Huang, Peng Wang, Xinjing Cheng, Dingfu Zhou, Qichuan Geng, and Ruigang Yang. The apolloscape open dataset for autonomous driving and its application. *IEEE transactions on pattern analysis and machine intelligence*, 42(10):2702–2719, 2019.

[24] Ajay Jain, Ben Mildenhall, Jonathan T. Barron, Pieter Abbeel, and Ben Poole. Zero-shot text-guided object generation with dream fields. 2022.

[25] Ajay Jain, Matthew Tancik, and Pieter Abbeel. Putting nerf on a diet: Semantically consistent few-shot view synthesis. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5885–5894, 2021.

[26] Adrian Johnston, Ravi Garg, Gustavo Carneiro, Ian Reid, and Anton van den Hengel. Scaling cnns for high resolution volumetric reconstruction from a single image. In *Proceedings of the IEEE International Conference on Computer Vision Workshops*, pages 939–948, 2017.

[27] Tero Karras, Samuli Laine, and Timo Aila. A style-based generator architecture for generative adversarial networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4401–4410, 2019.

[28] Michael Kazhdan and Hugues Hoppe. Screened poisson surface reconstruction. *ACM Transactions on Graphics (ToG)*, 32(3):1–13, 2013.

[29] Youngjung Kim, Hyungjoo Jung, Dongbo Min, and Kwanghoon Sohn. Deep monocular depth estimation via integration of global and local predictions. *IEEE transactions on Image Processing*, 27(8):4131–4144, 2018.

[30] Zhengqi Li and Noah Snavely. Megadepth: Learning single-view depth prediction from internet photos. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2041–2050, 2018.

[31] Zhengqi Li and Noah Snavely. Megadepth: Learning single-view depth prediction from internet photos. In *Computer Vision and Pattern Recognition (CVPR)*, 2018.

[32] 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. *arXiv preprint arXiv:2211.10440*, 2022.

[33] Shitong Luo and Wei Hu. Diffusion probabilistic models for 3d point cloud generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2837–2845, 2021.

[34] Gal Metzer, Elad Richardson, Or Patashnik, Raja Giryes, and Daniel Cohen-Or. Latent-nerf for shape-guided generation of 3d shapes and textures. *arXiv preprint arXiv:2211.07600*, 2022.

[35] B Mildenhall, PP Srinivasan, M Tancik, JT Barron, R Ramamoorthi, and R Ng. Nerf: Representing scenes as neural radiance fields for view synthesis. In *European conference on computer vision*, 2020.

[36] Kaichun Mo, Paul Guerrero, Li Yi, Hao Su, Peter Wonka, Niloy Mitra, and Leonidas J Guibas. Structurenet: Hierarchical graph networks for 3d shape generation. *arXiv preprint arXiv:1908.00575*, 2019.

[37] 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.

[38] Pushmeet Kohli Nathan Silberman, Derek Hoiem and Rob Fergus. Indoor segmentation and support inference from rgbd images. In *ECCV*, 2012.

[39] Ben Poole, Ajay Jain, Jonathan T Barron, and Ben Mildenhall. Dreamfusion: Text-to-3d using 2d diffusion. *arXiv preprint arXiv:2209.14988*, 2022.

[40] 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. *arXiv preprint arXiv:2103.00020*, 2021.

[41] René Ranftl, Alexey Bochkovskiy, and Vladlen Koltun. Vision transformers for dense prediction. *ICCV*, 2021.

[42] René Ranftl, Katrin Lasinger, David Hafner, Konrad Schindler, and Vladlen Koltun. Towards robust monocular depth estimation: Mixing datasets for zero-shot cross-dataset transfer. *IEEE transactions on pattern analysis and machine intelligence*, 44(3):1623–1637, 2020.

[43] René Ranftl, Katrin Lasinger, David Hafner, Konrad Schindler, and Vladlen Koltun. Towards robust monocular depth estimation: Mixing datasets for zero-shot cross-dataset transfer. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 44(3), 2022.

7

[44] Ali Razavi, Aaron van den Oord, and Oriol Vinyals. Generating diverse high-fidelity images with vq-vae-2. In *Advances in neural information processing systems*, pages 14866–14876, 2019.

[45] Christian Reiser, Songyou Peng, Yiyi Liao, and Andreas Geiger. Kilonerf: Speeding up neural radiance fields with thousands of tiny mlps. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 14335–14345, 2021.

[46] Danilo Jimenez Rezende, Shakir Mohamed, and Daan Wierstra. Stochastic backpropagation and approximate inference in deep generative models. In *International conference on machine learning*, pages 1278–1286. PMLR, 2014.

[47] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models, 2021.

[48] 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.

[49] Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily Denton, Seyed Kamyar Seyed Ghasemipour, Burcu Karagol Ayan, S Sara Mahdavi, Rapha Gontijo Lopes, et al. Photorealistic text-to-image diffusion models with deep language understanding. *arXiv preprint arXiv:2205.11487*, 2022.

[50] Shunsuke Saito, Zeng Huang, Ryota Natsume, Shigeo Morishima, Angjoo Kanazawa, and Hao Li. Pifu: Pixel-aligned implicit function for high-resolution clothed human digitization. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 2304–2314, 2019.

[51] Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade Gordon, Ross Wightman, Mehdi Cherti, Theo Coombes, Aarush Katta, Clayton Mullis, Mitchell Wortsman, et al. Laion-5b: An open large-scale dataset for training next generation image-text models. *arXiv preprint arXiv:2210.08402*, 2022.

[52] Edward J Smith and David Meger. Improved adversarial systems for 3d object generation and reconstruction. In *Conference on Robot Learning*, pages 87–96. PMLR, 2017.

[53] Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. *arXiv preprint arXiv:2010.02502*, 2020.

[54] Jiaxiang Tang. Stable-dreamfusion: Text-to-3d with stable-diffusion, 2022. https://github.com/ashawkey/stable-dreamfusion.

[55] Can Wang, Menglei Chai, Mingming He, Dongdong Chen, and Jing Liao. Clip-nerf: Text-and-image driven manipulation of neural radiance fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3835–3844, 2022.

[56] Chaoyang Wang, Simon Lucey, Federico Perazzi, and Oliver Wang. Web stereo video supervision for depth prediction from dynamic scenes, 2019.

[57] Nanyang Wang, Yinda Zhang, Zhuwen Li, Yanwei Fu, Wei Liu, and Yu-Gang Jiang. Pixel2mesh: Generating 3d mesh models from single rgb images. In *Proceedings of the European conference on computer vision (ECCV)*, pages 52–67, 2018.

[58] Qianqian Wang, Zhicheng Wang, Kyle Genova, Pratul P Srinivasan, Howard Zhou, Jonathan T Barron, Ricardo Martin-Brualla, Noah Snavely, and Thomas Funkhouser. Ibrnet: Learning multi-view image-based rendering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4690–4699, 2021.

[59] Qiang Wang, Shizhen Zheng, Qingsong Yan, Fei Deng, Kaiyong Zhao, and Xiaowen Chu. Irs: A large naturalistic indoor robotics stereo dataset to train deep models for disparity and surface normal estimation. *arXiv preprint arXiv:1912.09678*, 2019.

[60] Wenshan Wang, Delong Zhu, Xiangwei Wang, Yaoyu Hu, Yuheng Qiu, Chen Wang, Yafei Hu, Ashish Kapoor, and Sebastian Scherer. Tartanair: A dataset to push the limits of visual slam. 2020.

[61] Jiajun Wu, Yifan Wang, Tianfan Xue, Xingyuan Sun, Bill Freeman, and Josh Tenenbaum. Marrnet: 3d shape reconstruction via 2.5 d sketches. *Advances in neural information processing systems*, 30, 2017.

[62] Jiajun Wu, Chengkai Zhang, Tianfan Xue, Bill Freeman, and Josh Tenenbaum. Learning a probabilistic latent space of object shapes via 3d generative-adversarial modeling. *Advances in neural information processing systems*, 29, 2016.

[63] Ke Xian, Chunhua Shen, Zhiguo Cao, Hao Lu, Yang Xiao, Ruibo Li, and Zhenbo Luo. Monocular relative depth perception with web stereo data supervision. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2018.

[64] Ke Xian, Jianming Zhang, Oliver Wang, Long Mai, Zhe Lin, and Zhiguo Cao. Structure-guided ranking loss for single image depth prediction. In *The IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2020.

[65] Dejia Xu. Neurallift-360: Lifting an in-the-wild 2d photo to a 3d object with 360° views. https://vita-group.github.io/NeuralLift-360/.

[66] Dejia Xu, Yifan Jiang, Peihao Wang, Zhiwen Fan, Humphrey Shi, and Zhangyang Wang. Sinnerf: Training neural radiance fields on complex scenes from a single image. 2022.

[67] Dejia Xu, Yifan Jiang, Peihao Wang, Zhiwen Fan, Humphrey Shi, and Zhangyang Wang. Sinnerf: Training neural radiance fields on complex scenes from a single image. In *Computer Vision–ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part XXII*, pages 736–753. Springer, 2022.

[68] Dejia Xu, Yifan Jiang, Peihao Wang, Zhiwen Fan, Yi Wang, and Zhangyang Wang. Neurallift-360: Lifting an in-the-wild 2d photo to a 3d object with 360° views. 2022.

[69] Qiangeng Xu, Zexiang Xu, Julien Philip, Sai Bi, Zhixin Shu, Kalyan Sunkavalli, and Ulrich Neumann. Point-nerf: Point-based neural radiance fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5438–5448, 2022.

[70] Guandao Yang, Xun Huang, Zekun Hao, Ming-Yu Liu, Serge Belongie, and Bharath Hariharan. Pointflow: 3d point cloud generation with continuous normalizing flows. In *Proceedings*

*of the IEEE/CVF international conference on computer vision*, pages 4541–4550, 2019.

[71] 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.

[72] Alex Yu, Ruilong Li, Matthew Tancik, Hao Li, Ren Ng, and Angjoo Kanazawa. Plenoctrees for real-time rendering of neural radiance fields. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5752–5761, 2021.

[73] Alex Yu, Vickie Ye, Matthew Tancik, and Angjoo Kanazawa. pixelnerf: Neural radiance fields from one or few images. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4578–4587, 2021.

[74] Yuxuan Zhang, Wenzheng Chen, Huan Ling, Jun Gao, Yinan Zhang, Antonio Torralba, and Sanja Fidler. Image gans meet differentiable rendering for inverse graphics and interpretable 3d neural rendering. *arXiv preprint arXiv:2010.09125*, 2020.

[75] 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.

[76] Qian-Yi Zhou, Jaesik Park, and Vladlen Koltun. Open3d: A modern library for 3d data processing. *arXiv preprint arXiv:1801.09847*, 2018.