
Generative Lifting of Multiview to 3D from Unknown Pose: Wrapping NeRF inside Diffusion

Xin Yuan
University of Chicago

Rana Hanocka
University of Chicago

Michael Maire
University of Chicago

{yuanx, ranahanocka, mmair}@uchicago.edu

Abstract

We cast multiview reconstruction from unknown pose as a generative modeling problem. From a collection of unannotated 2D images of a scene, our approach simultaneously learns both a network to predict camera pose from 2D image input, as well as the parameters of a Neural Radiance Field (NeRF) for the 3D scene. To drive learning, we wrap both the pose prediction network and NeRF inside a Denoising Diffusion Probabilistic Model (DDPM) and train the system via the standard denoising objective. Our framework requires the system accomplish the task of denoising an input 2D image by predicting its pose and rendering the NeRF from that pose. Learning to denoise thus forces the system to concurrently learn the underlying 3D NeRF representation and a mapping from images to camera extrinsic parameters. To facilitate the latter, we design a custom network architecture to represent pose as a distribution, granting implicit capacity for discovering view correspondences when trained end-to-end for denoising alone. This technique allows our system to successfully build NeRFs, without pose knowledge, for challenging scenes where competing methods fail. At the conclusion of training, our learned NeRF can be extracted and used as a 3D scene model; our full system can be used to sample novel camera poses and generate novel-view images.

1 Introduction

Structure from motion is a well-studied problem in computer vision, with a substantial history of research focusing on the specific task of reconstructing a 3D scene from a collection of 2D images captured from different viewpoints. When the 3D pose (camera extrinsics) for each 2D view is unknown, classic approaches [26, 1] explicitly estimate correspondence between 2D views (*e.g.*, by matching local feature descriptors) prior to optimizing a shared 3D geometry whose reprojections are consistent with those views. Neural Radiance Fields (NeRFs) [19, 2, 16, 3] have led a revolution toward widespread use of differentiable 3D scene representations [19, 8, 12] that are compatible with deep learning techniques. However, the problem of jointly solving for both the 3D reconstruction and the pose, when neither is known a priori, remains an open problem. Recent attempts to connect learning of camera pose with NeRFs operate under simplifying assumptions, such as coarse pose initialization (only learning adjustments) [14] or front-facing (as opposed to arbitrary 360°) views of the scene [29].

In parallel with the development of differentiable 3D representations, progress across a variety of paradigms for generative models [9, 13, 11], has transformed the landscape for designing and training systems using deep learning. Learning to synthesize data provides an unsupervised training objective and scaling compute, parameters, and datasets is a path toward foundation models [4] whose feature representations can subsequently be repurposed to specific downstream tasks. However, large-scale foundation models are not the only setting in which generative learning is appropriate. Nor is

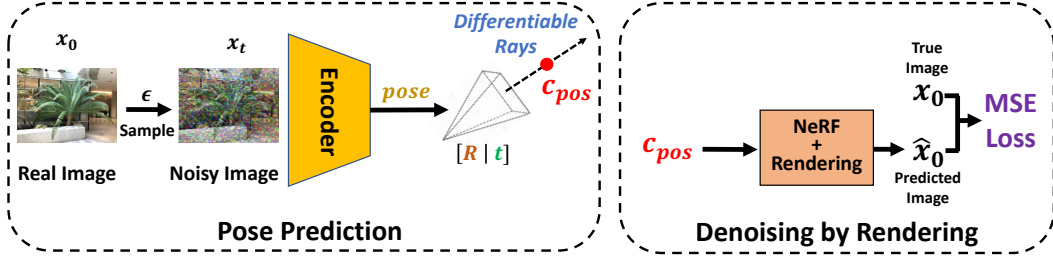


Figure 1: **Wrapping NeRF inside Diffusion.** We learn a 3D scene reconstruction by training a denoising diffusion model (DDPM) on a dataset of 2D views of the scene. The architecture of our DDPM consists of two components. **Left:** An *Encoder* predicts the pose of a single noisy 2D input image. **Right:** A *NeRF* is rendered from the predicted camera pose to create a 2D output image that is treated as the predicted denoising of the input view. The system must learn parameters of both the *Encoder* and *NeRF* so that any 2D view can be denoised by predicting a camera and rendering the scene. The NeRF rendering process is differentiable with respect to rays shot from the camera, which themselves depend on the camera-to-world transformation matrix produced by the encoder. All modules are end-to-end trainable, and the system is optimized by the simple MSE loss on denoising.

manipulation of pre-trained models (e.g., extracting features, fine-tuning, or prompting) the only strategy for applying generative learning to solve downstream tasks.

Yuan and Maire [31] demonstrate an alternative strategy that utilizes a generative model and relies solely on a generative learning objective, yet directly solves a downstream task (image segmentation) as a byproduct of training the generative model. Their strategy is to constrain the architecture of the generative model such that it must synthesize an image by first predicting a segmentation and then generating the corresponding image regions in parallel. Trained as a Denoising Diffusion Probabilistic Model (DDPM) [11], segmentation emerges as the bottleneck representation in a network that first partitions a noisy input into regions and then denoises each region in parallel.

We port this general concept to the problem of multiview 3D reconstruction from unknown pose, where we devise an internal pose prediction network and a NeRF comprising the task-specific architecture encapsulated within our DDPM; see Figure 1. We solve a small-scale generative modeling problem: learning to generate images in the collection of 2D views of a single scene. Training examples are noised 2D images (views), the DDPM output is a predicted denoised image, and the loss is the denoising objective. Inside our generative wrapper, the model architecture dictates that we map a noisy image to a predicted camera pose and then render the NeRF from that pose to synthesize the clean output image. Successfully performing denoising in this manner requires that: (a) the NeRF stores a 3D scene representation consistent with all of the 2D views, and (b) the pose prediction network implicitly solves the 2D view correspondence problem by mapping each 2D input image to camera coordinates from which it is reconstructed by rendering the NeRF.

Figure 2 illustrates how our trained system jointly solves pose prediction and 3D reconstruction. Our system enables predicting the 3D pose for new (unseen) images, and re-rendering the learned scene from different camera poses which can be generated from noise or provided explicitly. As Figure 3 shows and Section 3 describes in detail, we significantly expand the complexity of multiview reconstruction problems our system can solve by replacing our simple pose prediction network with a more expressive version. This alternative maintains a representation of uncertainty over a distribution of multiple possible poses, which gives our system, trained from scratch by gradient descent, the implicit capacity to explore more view correspondence configurations.

Our contributions are:

- A new approach to 3D reconstruction from unknown pose based entirely on generative training. Denoising is a generative wrapper that “lifts” an architecture consisting of a forward model for pose prediction and differentiable rendering to learn view correspondence and 3D reconstruction. Compared to an autoencoder, this fully generative wrapper benefits learned reconstruction quality.
- A novel architecture for pose prediction that enables representing uncertainty during training, allowing us to learn 3D reconstructions from 2D views of arbitrary and unknown pose.

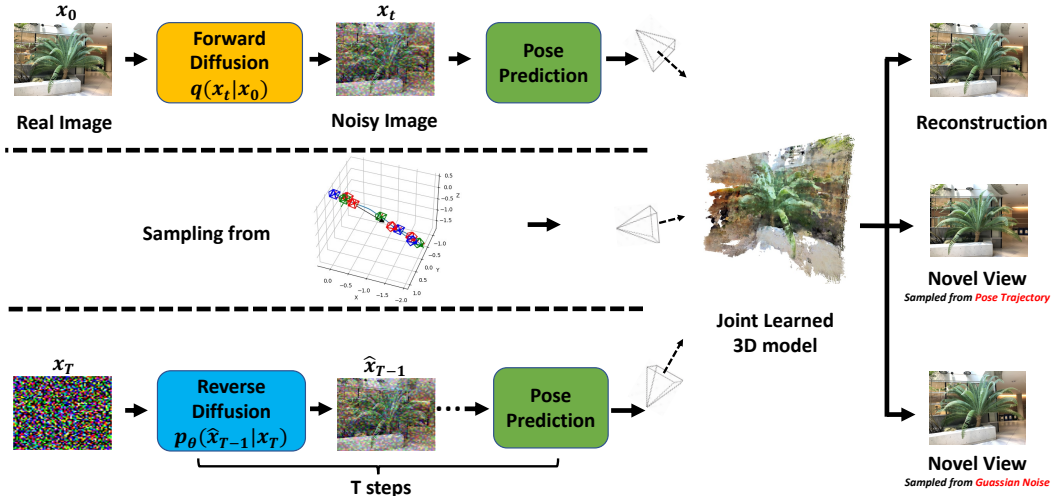


Figure 2: **Unifying pose prediction, 3D reconstruction, and novel-view image generation.** Our trained system (Figure 1) can be deployed for multiple tasks. *Pose prediction (top)*: We can predict the pose of a previously unseen real image by adding a small amount of noise (forward diffusion) and feeding it to our *Encoder* (Fig 1, left). Rendering our learned *NeRF* from that camera pose should reconstruct the real image. *Direct NeRF usage (middle)*: Our learned *NeRF* can be extracted and directly used to render the scene (e.g., along a manually specified camera path). *Sampling cameras and views (bottom)*: Performing sequential diffusion denoising from pure Gaussian noise input synthesizes a camera pose from which rendering the *NeRF* generates a novel view of the scene.

- New capabilities for 3D NeRF reconstruction which are demonstrated through experiments on arbitrary 360-degree poses. While Wang et al. [29] can reconstruct under certain assumptions about an unknown camera (e.g., forward-facing views of the scene), they fail on image collections from unconstrained pose (e.g., 360° views). Our method successfully reconstructs a NeRF and infers camera pose for these challenging datasets.

2 Related Work

Neural Radiance Fields (NeRFs) [19] have emerged as a powerful framework for 3D scene reconstruction and view synthesis, with multiple extensions and improvements [2, 16, 3]. NeRF++ [33] adds spatially-varying reflectance and auxiliary tasks for better training. PixelNeRF [30] extends NeRF to generate high-quality novel views from one or few input images, but still requires camera pose information. NeRF-- [29] jointly optimizes camera intrinsics and extrinsics as learnable parameters while training NeRFs. However, their proposed training scheme and camera parameterization cannot handle large camera rotation and is restricted to forward-facing views of the scene.

Generative models for 3D reconstruction aim to infer the underlying 3D structure of a scene from a set of 2D images. These models often learn a latent representation of the 3D scene and use it to generate novel views or perform other tasks such as object manipulation or scene editing. Generative Radiance Fields (GRAFs) [25] combine NeRFs with VAEs or GANs to generate novel 3D scenes without explicit 3D geometry. GRAFs learn a latent space encoding scene structure, with NeRF mapping points in this space to 3D radiance fields. DiffRF [20] leverages the diffusion prior to perform 3D completion in a two-stage manner, which is further improved by SSD-NeRF [7] with a single stage training scheme and an end-to-end objective that jointly optimizes a NeRF and diffusion. Multiple works combine NeRF with generative models for the purpose of 3D synthesis [5, 17, 10], including ones that place NeRF and diffusion models in series [21, 15, 27]. Our framework’s nested structure differs, as our aim is not to learn to generate novel 3D scenes; we aim to use generative training to solve the classic multiview 3D reconstruction problem.

Pose estimation is the challenging task of estimating object or camera position and orientation within a scene. COLMAP [23, 24] uses a Structure-from-Motion (SfM) [23] approach for pose estimation, can handle challenging scenes with varying lighting conditions and viewpoints, and is widely used in

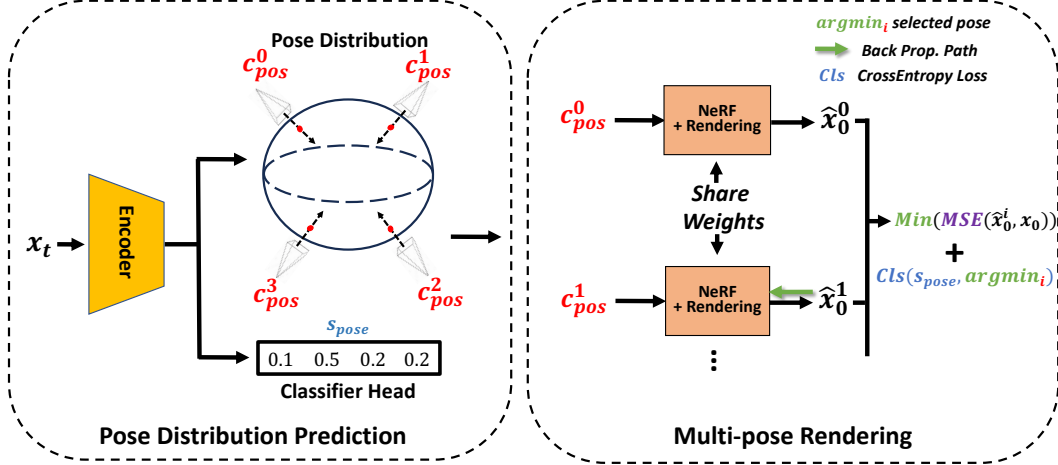


Figure 3: **Pose distribution representation and multi-pose rendering for 360° scenes.** In order to perform view denoising by learning a NeRF and predicting the pose from which to render it, our system (Figure 1) must implicitly solve multiview correspondence by mapping training images (of unknown pose) into consistent locations in the 3D environment. We enable training via gradient descent to discover such solutions for challenging multiview datasets (e.g., spanning 360°) by augmenting our architecture with the capacity to represent uncertainty over a pose distribution. **Left:** Our encoder, given a noisy image x_t , predicts parameters for multiple cameras and a corresponding probability distribution over cameras, s_{pose} . **Right:** During training, we render the NeRF from each predicted camera and use the best reconstruction to calculate the denoising loss; an auxiliary classification loss pushes the predicted camera distribution to upweight the selected output. At test time, we render using only the single camera predicted as most likely by the classifier.

NeRF training. However, this collection of techniques requires a large number of images for accurate pose estimation; pre-processing also restricts flexibility. PoseDiffusion [28] and Camera-as-Rays [32] use a diffusion model to denoise camera parameters and rays. Although sharing similar spirit in adopting diffusion, these methods require a supervised pertaining stage. More importantly, diffusion serves as a different role in our model: instead of denoising cameras to recover the pose distribution, we modulate a pose prediction system embedded inside the diffusion training process, yielding pose information as a latent representation.

3 Method

3.1 Unsupervised Pose Prediction from a Single Image

Our pose module (Figure 1, left) consists of several components designed to predict, from a 2D image, the position and orientation of a camera in the scene. We design the encoder based on a standard DDPM U-Net. We obtain input x_t , a noise version of x_0 , via forward diffusion:

$$\begin{aligned} q(x_t|x_0) &:= \mathcal{N}(x_t; \sqrt{\bar{\alpha}_t}x_0, (1 - \bar{\alpha}_t)I), \\ x_t &= \sqrt{\bar{\alpha}_t}x_0 + \sqrt{1 - \bar{\alpha}_t}\epsilon, \epsilon \sim \mathcal{N}(0, 1), \end{aligned} \quad (1)$$

where $\alpha_t = 1 - \beta_t$, $\bar{\alpha}_t = \prod_{s=1}^t \alpha_s$. We encode the pose information in the form of a camera-to-world transformation matrix, $T_{wc} = [Ro|ts]$, where $Ro \in SO(3)$ represents the camera’s rotation, and $ts \in \mathbb{R}^3$ represents its translation. Following [29], we adopt Rodrigues’ formula to form the rotation matrix Ro from axis-angle representation:

$$Ro = I + \sin(\phi)[\omega]_{\times} + (1 - \cos(\phi))[\omega]_{\times}^2 \quad (2)$$

where ϕ is the rotation angle, ω is a normalized rotation axis, and $[\omega]_{\times}$ is the skew-symmetric matrix of the rotation axis vector ω .

Different from [29], which formulates ω_i and translation ts_i as trainable parameters, for each input image x^i , our system maps the corresponding U-Net encoder features to two 3-dimensional input-

Algorithm 1 Generative Lifting to 3D: Single Camera Pose Prediction & NeRF

Input: Multiview image collection \mathcal{X}
Output: Encoder (pose predictor) & NeRF
Initialize: Model weights θ , Timesteps T
for iter = 1 **to** Iter_{total} **do**
 Sample $\mathbf{x}_0 \in \mathcal{X}$, $t \in [1, T]$.
 Sample \mathbf{x}_t using Eq. 1.
 Predict R_o using Eq. 2 and t_s .
 Compute $\hat{\mathbf{x}}_0$ by rendering the NeRF from the pose predicted for \mathbf{x}_t (see Figure 1).
 Backprop from loss in Eq. 3.
 Update model weights.
end for
return θ

Algorithm 2 Generative Lifting to 3D using Pose Distribution Modeling & Multi-pose NeRF Rendering

Input: Multiview image collection \mathcal{X}
Output: Encoder (multi-pose predictor & classifier) & NeRF
Initialize: Model weights θ , Timesteps T , and initial pose of candidate cameras
for iter = 1 **to** Iter_{total} **do**
 Sample $\mathbf{x}_0 \in \mathcal{X}$, $t \in [1, T]$.
 Sample \mathbf{x}_t using Eq. 1.
 Predict poses and corresponding \mathbf{s}_{pose} , as in Figure 3.
 Compute $\{\hat{\mathbf{x}}_0^i\}$ using multi-pose rendering (Figure 3).
 Backprop along the path of the best render via Eq. 4.
 Update model weights.
end for
return θ

dependent feature vectors. By doing so, we not only learn to fit the pose information during training, but also obtain a pose predictor network applicable to any input 2D image.

3.2 3D Optimization with Denoising Rendering

With the predicted camera pose from the U-Net Encoder, the system is able to perform denoising via differentiable rendering. Specifically, as shown in the right side of Figure 1, we sample the differentiable coordinates c_{pos} , which are fed into a NeRF MLP model to learn object density and opacity, generating a 2D image reconstruction $\hat{\mathbf{x}}_0$.

Benefiting from the compatibility between NeRF reconstruction loss and DDPM denoising objective, we train our model end-to-end by simply minimizing the pixel-wise distance from $\hat{\mathbf{x}}_0$ to \mathbf{x}_0 . Model weights of the camera predictor (U-Net encoder) and NeRF MLPs are optimized with loss:

$$L = \mathbb{E} \|\hat{\mathbf{x}}_0 - \mathbf{x}_0\|_2^2 \quad (3)$$

Algorithm 1 summarizes training.

3.3 Multi-pose Rendering for Scene Reconstruction from 360° Views

A failure to estimate poses accurately can occur when rotation perturbations exceed a certain threshold in Eq. 2, which prevents learning from 360° views of a scene [29]. Even with good reconstruction in 2D space, an overfitting issue can occur during optimization, where NeRF compensates by creating multiple disjoint copies of scene fragments instead of a unified 3D reconstruction. We address these issues via a higher capacity pose predictor capable of representing uncertainty (Figure 3).

Pose distribution prediction. A simple camera parameterization is to restrict position to a fixed-radius sphere with fixed intrinsics, and the constraint of always looking towards the origin at (0, 0, 0). Assuming the object in a 360° scene is centered, rotated, and scaled by some canonical alignment, such a parametrization has only two degrees of freedom. However, this simplistic approach restricts model capacity for capturing diverse viewpoints or extensive rotations.

We propose a more flexible approach that allows for a wider range of camera positions and orientations. Given a 2D image captured from a specific viewpoint, instead of predicting a single transformation, we sample the camera’s position from a distribution of multiple cameras that cover a larger range of positions and orientations on a sphere.

As Figure 3 shows, different candidate cameras spread over the sphere, pointing to the origin at initialization. Each input view predicts parameters for all cameras in the distribution. An auxiliary classifier head predicts the probability of the input view corresponding to each camera in the distribution. Early in training, such a design facilitates searching over multiple pose hypotheses in order to discover a registration of all views into a consistent coordinate frame. Only one camera prediction per input view need be correct, as long as the system also learns which one.

Joint optimization with multi-pose rendering. We render the NeRF separately from each candidate camera to produce a set of 2D images $\{\hat{\mathbf{x}}_0^i\}$. During backpropagation, we only allow the gradient to

pass selectively to optimize the best match between the true image and the rendered reconstruction. The selected branch index serves as a pseudo-label to co-adapt the classifier head in a self-supervised bootstrapping manner. The total loss for joint training is:

$$L = \min_i \|\hat{\mathbf{x}}_0^i - \mathbf{x}_0\|_2^2 + \lambda \text{CrossEntropy}(\mathbf{s}_{pose}, \arg \min_i \|\hat{\mathbf{x}}_0^i - \mathbf{x}_0\|_2^2) \quad (4)$$

where λ is the trade-off parameter between view reconstruction and camera classification. We set λ as 0.1 in experiments and investigate its effect in an ablation study. Algorithm 2 summarizes training.

3.4 Novel View Generation

Figure 2 illustrates the different modalities in which our trained system can be used.

Pose prediction & reconstruction. Given input image \mathbf{x}_0 , we sample a noisy version \mathbf{x}_t through forward diffusion in Eq. 1. We then pass \mathbf{x}_t to the model. Our system estimates the camera pose of \mathbf{x}_t with respect to the scene and recovers a clean image reconstruction with one-step denoising.

Sampling from pose trajectory. Though not trained with any camera pose information, our system can generate novel views using a pre-defined camera trajectory, acting as a conventional NeRF model.

Sampling from Gaussian noise. A unique property of our system is its support for sampling cameras and scene views. Using reverse diffusion, our model can generate a realistic novel view and the corresponding camera pose, starting from a pure noise input $\mathbf{x}_T \sim \mathcal{N}(0, 1)$. We perform T steps of reverse diffusion (predict \mathbf{x}_{t-1} from \mathbf{x}_t) to progressively generate a novel view.

4 Experiments

We conduct the experiments on the face-forwarding dataset LLFF [18] with a resolution of 378×504 , using the single camera prediction system depicted in Alg. 1. To handle 360° scenes, we adopt the multi-pose rendering as in Alg. 2 on ShapeNet Car [6] (5 scenes), Lego and Drums [19] with a resolution of 128×128 . Instead of simultaneously optimizing two networks: one ‘coarse’ and one ‘fine’, we only use a single network to represent 360° scenes for all methods. We initialize 8 camera candidates spread over 8 quadrants of a sphere for the ShapeNet Car scene, and 12 candidates over 4 quadrants of a semi-sphere for Lego and Drums.

4.1 Implementation Details

For all experiments, we use a U-Net [22] encoder as the pose prediction module. The downsampling stack performs five steps of downsampling, each with 2 residual blocks. From highest to lowest resolution, U-Net stages use $[C, C, 2C, 2C, 4C]$ channels, respectively. We set $C = 64$ for all models. Figure 12 details the network architecture. We use 100 denoising steps for all models.

Our method involves two sets of trainable parameters: NeRF model weights and pose prediction network weights; we adopt separate Adam optimizers, with learning rates $1e^{-4}$ and $2e^{-5}$ for NeRF and pose prediction, respectively. We set (β_1, β_2) as $(0.9, 0.999)$ for both optimizers. We adopt the same batch size and learning rate scheduler used to train the corresponding baseline NeRF as in [19]. We train all models for 200k iterations. Code segments 1, 2 and 3 detail our camera parameterization.

4.2 Multi-view 3D Reconstruction

We measure the quality of reconstructions obtained by the *top* pipeline in Figure 2. We directly input previously unseen images from different views to generate reconstructions.

LLFF dataset. Figure 4 shows we obtain good-quality reconstructions and disparity predictions. Point cloud visualization plots the density and opacity output from our NeRF model at 3D coordinates.

360° scenes. Camera motions with large rotation perturbations cause failures in NeRF--; it cannot handle 360° scenes like Lego. Our method solves this challenging case from a single input image, without relative pose estimation between image pairs. To obtain reconstructions, we determine the camera pose for the input image based on the maximum score of the classification head in Figure 3 before rendering. As Figure 5 shows, we generate good reconstructions and point clouds.

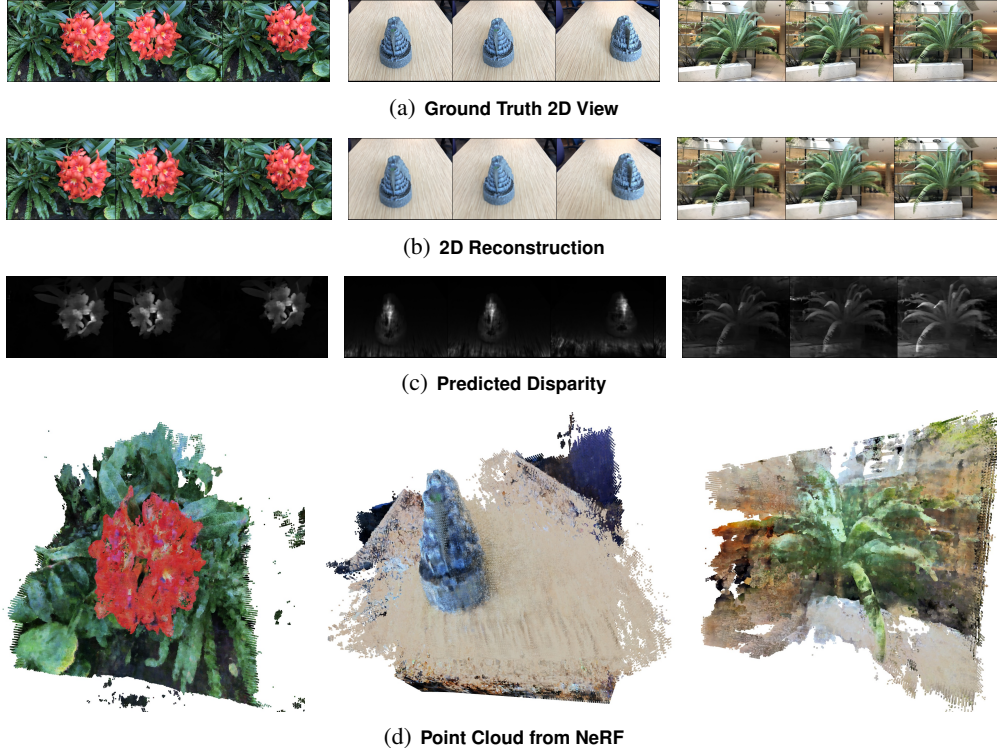


Figure 4: Reconstructions of images unseen during training on three scenes from LLFF [18].

Table 1: **Multiview reconstruction quality (PSNR, SSIM, & LPIPS).** Our system, without pose knowledge, reconstructs 3D scenes from challenging image collections (views spanning 360°) on which NeRF-- [29] fails. Supervised denotes standard NeRF training using ground-truth camera pose.

| Type | Scene | PSNR (\uparrow) | | | SSIM (\uparrow) | | | LPIPS (\downarrow) | | |
|----------------|----------|---------------------|----------|-------|---------------------|----------|------|------------------------|----------|------|
| | | Supervised | NeRF-- | Ours | Supervised | NeRF-- | Ours | Supervised | NeRF-- | Ours |
| Forward-Facing | Fern | 22.22 | 21.67 | 17.02 | 0.64 | 0.61 | 0.42 | 0.47 | 0.50 | 0.55 |
| | Flower | 25.25 | 25.34 | 22.42 | 0.71 | 0.71 | 0.58 | 0.36 | 0.37 | 0.43 |
| | Fortress | 27.60 | 26.20 | 22.02 | 0.73 | 0.63 | 0.50 | 0.38 | 0.49 | 0.51 |
| | Horns | 24.25 | 22.53 | 17.48 | 0.68 | 0.61 | 0.43 | 0.44 | 0.50 | 0.55 |
| | Leaves | 18.81 | 18.88 | 14.44 | 0.52 | 0.53 | 0.42 | 0.47 | 0.47 | 0.60 |
| | Orchids | 19.09 | 16.73 | 14.34 | 0.51 | 0.39 | 0.40 | 0.46 | 0.55 | 0.58 |
| | Room | 27.77 | 25.84 | 22.36 | 0.87 | 0.84 | 0.48 | 0.40 | 0.44 | 0.49 |
| | Trex | 23.19 | 22.67 | 19.96 | 0.74 | 0.72 | 0.62 | 0.41 | 0.44 | 0.51 |
| 360° | Car | 28.98 | \times | 26.43 | 0.95 | \times | 0.92 | 0.08 | \times | 0.08 |
| | Lego | 25.44 | \times | 21.38 | 0.92 | \times | 0.86 | 0.09 | \times | 0.12 |
| | Drums | 22.12 | \times | 18.65 | 0.89 | \times | 0.82 | 0.08 | \times | 0.16 |

To quantify the quality of our reconstructions, we compare the PSNR, SSIM, LPIPS [34] with supervised NeRF (using pre-processed pose information), and NeRF--. As Table 1 shows, for face-forwarding scenes, our method achieves reasonably high performance. We cannot beat NeRF-- because we aim to solve a more general pose prediction problem from a single input, instead of fitting camera parameters as trainable variables. For 360° scene views, which NeRF-- completely fails to handle, our method still yields good reconstructions. This validates the design choice of our multi-pose rendering system in tolerating large camera pose perturbations.

4.3 Visualization of Pose Optimization

As shown in Figure 6, we also demonstrate that our pose prediction system can generate reasonable pose estimates, though not in the same coordinate system, compared with ground-truth cameras.

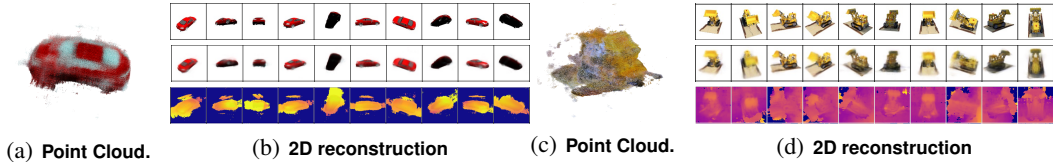


Figure 5: Reconstructions on 360° scenes.

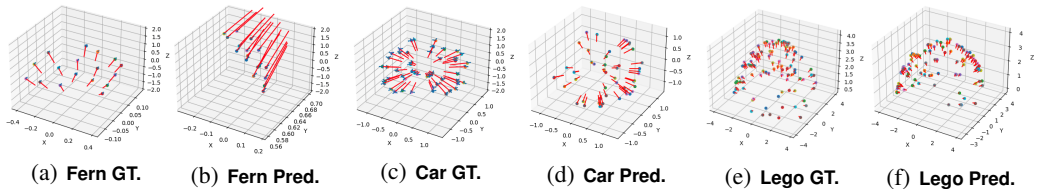


Figure 6: Visualization of camera poses for Fern, Car, and Lego.

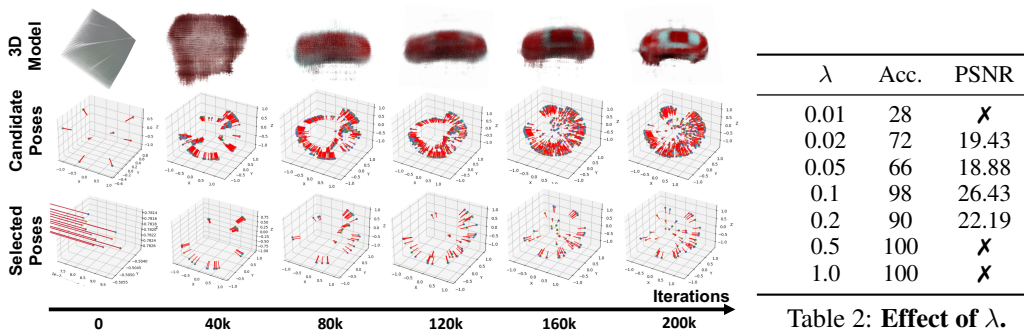


Table 2: Effect of λ .

Figure 7: Joint 3D and pose optimization for Car during training.

Camera distribution evolution during optimization To better demonstrate the pose prediction refinement during the optimization process, we visualize the camera poses at different training iterations for the Car scene. As shown in Figure 7, the candidate poses refers to all possible predictions over 8 quaternions while the selected poses represent those with maximum classification scores. During training, the learned candidate poses tend to cover the sphere uniformly, with the selected pose distributions gradually converge to that provided by the pre-processed dataset. We observe simultaneous refinements of both 3D model and pose prediction along the training process.

4.4 Novel View Synthesis

From camera trajectory: We show that our joint-learned 3D model, even learned with unknown pose, can generate novel views using a manually designed camera trajectory as in a supervised NeRF. In Figure 8, we generate novel views using a continuous spiral camera path (pointing to the origin) for three different challenging scenes.

From Gaussian noise: Given the nice property of denoising diffusion training, we have the flexibility to generate novel views from Gaussian noise progressively. In the DDPM Markov process, the model implicitly formulates a mapping between noise and data distributions. We validate this by visualizing the reconstructed \hat{x}_{t-1} and \hat{x}_0 along the sequential reversed diffusion process in Figure 9. We observe gradual refinement as denoising steps approach $t = 0$.

4.5 Ablation Studies

Training with clean images. Training NeRF with denoising diffusion is an essential part of our method. We show the effectiveness of this design by varying the input with clean images in our method, in which case our architecture downgrades to an autoencoder (AE). As shown in Figure 10, the baseline trained with clean images (denoted as AE) yields an incorrect 3D model for the Car scene. This suggests the failure of camera pose prediction, hence leading to a NeRF over-fitting issue.

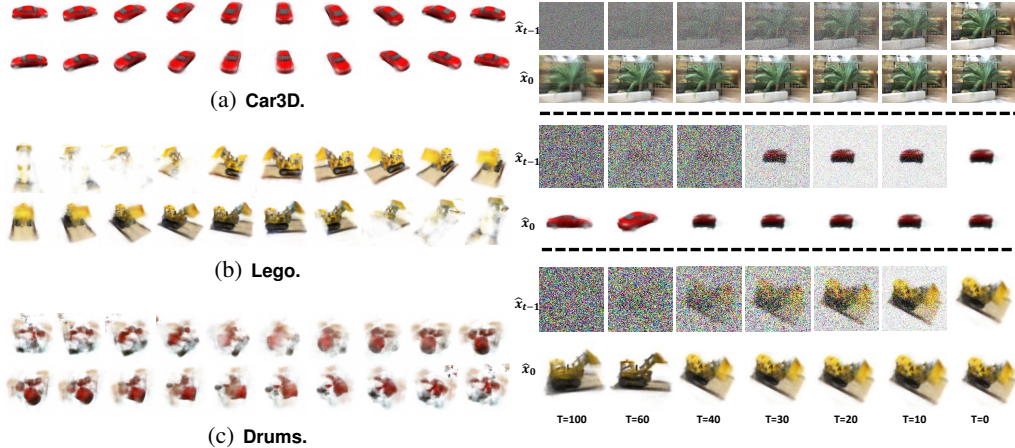


Figure 8: Novel view synthesis from circle trajectory.

Figure 9: Novel view synthesis from Gaussian noise.

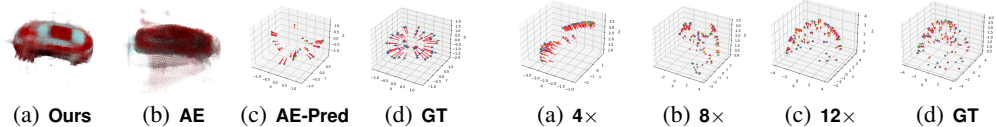


Figure 10: Learning w/o denoising diffusion.

Figure 11: Poses learned with different camera number.

Moreover, AE fails to perform novel-view synthesis given a pre-defined camera trajectory, as shown in Figure 13. Thus, denoising diffusion training not only provides a new way to perform novel view synthesis, but also significantly improves the learned 3D reconstruction.

Candidate camera numbers in multi-pose rendering. For Lego and Drums on a semi-sphere, we choose to use 12 candidate cameras instead of 4 at initialization for each view. The 4-way case poses an easier camera classification task to the system, but restricts flexibility for discovering view correspondence and thereby pushes the NeRF to overfit on incorrect pose predictions. Increasing capacity to 12 cameras during training addresses this issue and prevents the system from converging to a suboptimal solution. As shown in Figure 11, 4× and 8× variants fail to converge to the correct pose distributions, while the 12× succeeds.

Trade-off between classification and reconstruction. Due to the discrepancy between training and testing in our multi-pose system, the quality of rendering highly relies on the camera classification accuracy. To evaluate the performance of this self-supervised classifier, we use the camera index which produces the minimum reconstruction loss as ground-truth. We study the effect of the classification loss term by alternating $\lambda \in \{0.01, 0.02, 0.05, 0.1, 0.2, 0.5, 1.0\}$, generating trade-offs between the accuracy and PSNR for Car. As shown in Table 2, NeRF training yields a poor reconstruction and even an optimization failure when the classification accuracy is low. When we set λ as a large value, the model runs into a local minimum: the classifier obtains 100% accuracy at early training iterations and poses cannot be jointly optimized to discover correspondence.

5 Conclusion

We propose a novel technique that places NeRF inside a probabilistic diffusion framework to accurately predict camera poses and create detailed 3D scene reconstructions from collections of 2D images. Our approach enables training NeRF from images with unknown pose. Using a carefully constrained architecture and differentiable volume renderer, we learn a camera pose predictor and 3D representation jointly. Our experimental results and ablation studies confirm the effectiveness of this method, demonstrating its capability to produce high-quality reconstructions, localize previously unseen images, and sample novel-view images, all while trained in an entirely unsupervised manner.

References

- [1] Sameer Agarwal, Yasutaka Furukawa, Noah Snavely, Ian Simon, Brian Curless, Steven M Seitz, and Richard Szeliski. Building Rome in a day. *Communications of the ACM*, 54(10):105–112, 2011.
- [2] Jonathan T. Barron, Ben Mildenhall, Dor Verbin, Pratul P. Srinivasan, and Peter Hedman. Mip-nerf 360: Unbounded anti-aliased neural radiance fields. *CoRR*, abs/2111.12077, 2021.
- [3] Jonathan T. Barron, Ben Mildenhall, Dor Verbin, Pratul P. Srinivasan, and Peter Hedman. Zip-nerf: Anti-aliased grid-based neural radiance fields. In *ICCV*, 2023.
- [4] Rishi Bommasani, Drew A Hudson, Ehsan Adeli, Russ Altman, Simran Arora, Sydney von Arx, Michael S Bernstein, Jeannette Bohg, Antoine Bosselut, Emma Brunskill, et al. On the opportunities and risks of foundation models. *arXiv preprint arXiv:2108.07258*, 2021.
- [5] Eric R Chan, Marco Monteiro, Petr Kellnhofer, Jiajun Wu, and Gordon Wetzstein. pi-gan: Periodic implicit generative adversarial networks for 3d-aware image synthesis. In *CVPR*, 2021.
- [6] Angel X. Chang, Thomas A. Funkhouser, Leonidas J. Guibas, Pat Hanrahan, Qi-Xing Huang, Zimo Li, Silvio Savarese, Manolis Savva, Shuran Song, Hao Su, Jianxiong Xiao, Li Yi, and Fisher Yu. Shapenet: An information-rich 3d model repository. 2015.
- [7] Hansheng Chen, Jiatao Gu, Anpei Chen, Wei Tian, Zhuowen Tu, Lingjie Liu, and Hao Su. Single-stage diffusion nerf: A unified approach to 3d generation and reconstruction. In *ICCV*, 2023.
- [8] 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.
- [9] Ian J. Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron C. Courville, and Yoshua Bengio. Generative adversarial nets. In *NeurIPS*, 2014.
- [10] Jiatao Gu, Lingjie Liu, Peng Wang, and Christian Theobalt. Stylenerf: A style-based 3d-aware generator for high-resolution image synthesis. *arXiv preprint arXiv:2110.08985*, 2021.
- [11] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. In *NeurIPS*, 2020.
- [12] Bernhard Kerbl, Georgios Kopanas, Thomas Leimkühler, and George Drettakis. 3D gaussian splatting for real-time radiance field rendering. *ACM Transactions on Graphics*, 42(4):1–14, 2023.
- [13] Diederik P. Kingma and Max Welling. Auto-encoding variational bayes. In Yoshua Bengio and Yann LeCun, editors, *ICLR*, 2014.
- [14] Chen-Hsuan Lin, Wei-Chiu Ma, Antonio Torralba, and Simon Lucey. Barf: Bundle-adjusting neural radiance fields. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5741–5751, 2021.
- [15] 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 *CVPR*, 2023.
- [16] Ricardo Martin-Brualla, Noha Radwan, Mehdi S. M. Sajjadi, Jonathan T. Barron, Alexey Dosovitskiy, and Daniel Duckworth. Nerf in the wild: Neural radiance fields for unconstrained photo collections. In *CVPR*, 2021.
- [17] Quan Meng, Anpei Chen, Haimin Luo, Minye Wu, Hao Su, Lan Xu, Xuming He, and Jingyi Yu. Gnerf: Gan-based neural radiance field without posed camera. In *ICCV*, pages 6351–6361, 2021.
- [18] Ben Mildenhall, Pratul P. Srinivasan, Rodrigo Ortiz-Cayon, Nima Khademi Kalantari, Ravi Ramamoorthi, Ren Ng, and Abhishek Kar. Local light field fusion: Practical view synthesis with prescriptive sampling guidelines. *ACM Transactions on Graphics (TOG)*, 2019.
- [19] Ben Mildenhall, Pratul P. Srinivasan, Matthew Tancik, Jonathan T. Barron, Ravi Ramamoorthi, and Ren Ng. Nerf: Representing scenes as neural radiance fields for view synthesis. In *ECCV*, 2020.
- [20] Norman Müller, Yawar Siddiqui, Lorenzo Porzi, Samuel Rota Buló, Peter Kotschieder, and Matthias Nießner. Diffirf: Rendering-guided 3d radiance field diffusion. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4328–4338, 2023.

- [21] Ben Poole, Ajay Jain, Jonathan T Barron, and Ben Mildenhall. Dreamfusion: Text-to-3d using 2d diffusion. *arXiv preprint arXiv:2209.14988*, 2022.
- [22] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-Net: Convolutional networks for biomedical image segmentation. In *MICCAI*, 2015.
- [23] Johannes Lutz Schönberger and Jan-Michael Frahm. Structure-from-motion revisited. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016.
- [24] 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.
- [25] Katja Schwarz, Yiyi Liao, Michael Niemeyer, and Andreas Geiger. GRAF: Generative radiance fields for 3D-aware image synthesis. *NeurIPS*, 2020.
- [26] Noah Snavely, Steven M Seitz, and Richard Szeliski. Photo tourism: Exploring photo collections in 3D. In *ACM siggraph 2006 papers*, pages 835–846. 2006.
- [27] Haochen Wang, Xiaodan Du, Jiahao Li, Raymond A Yeh, and Greg Shakhnarovich. Score jacobian chaining: Lifting pretrained 2d diffusion models for 3d generation. In *CVPR*, 2023.
- [28] Jianyuan Wang, Christian Rupprecht, and David Novotny. Posediffusion: Solving pose estimation via diffusion-aided bundle adjustment. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 9773–9783, 2023.
- [29] Zirui Wang, Shangzhe Wu, Weidi Xie, Min Chen, and Victor Adrian Prisacariu. NeRF—: Neural radiance fields without known camera parameters. *arXiv preprint arXiv:2102.07064*, 2021.
- [30] 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.
- [31] Xin Yuan and Michael Maire. Factorized diffusion architectures for unsupervised image generation and segmentation. *arXiv preprint arXiv:2309.15726*, 2023.
- [32] Jason Y Zhang, Amy Lin, Moneish Kumar, Tzu-Hsuan Yang, Deva Ramanan, and Shubham Tulsiani. Cameras as rays: Pose estimation via ray diffusion. *ICLR*, 2024.
- [33] Kai Zhang, Gernot Riegler, Noah Snavely, and Vladlen Koltun. Nerf++: Analyzing and improving neural radiance fields. *arXiv preprint arXiv:2010.07492*, 2020.
- [34] 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.

A Appendix

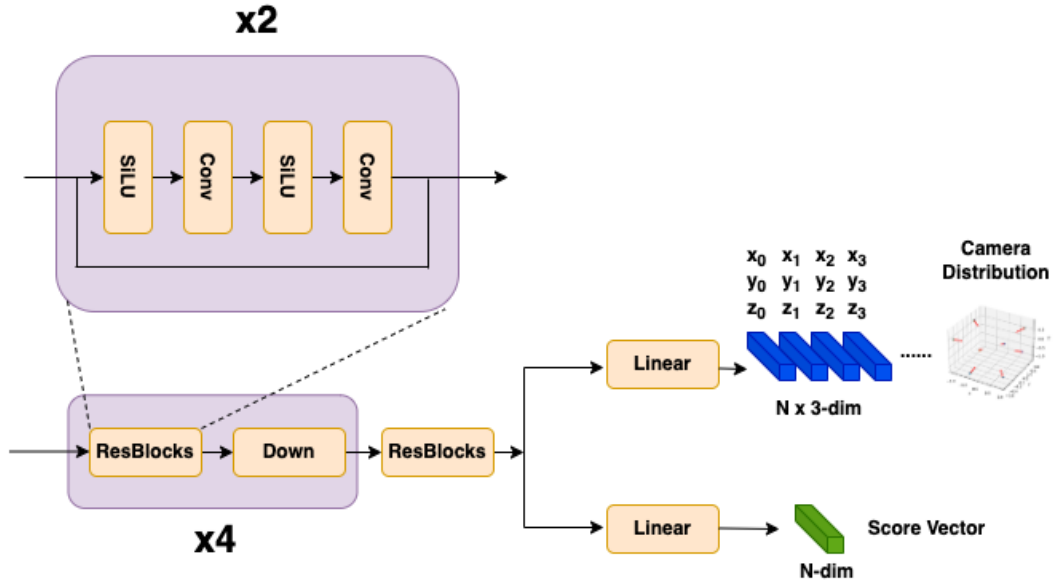


Figure 12: **Detailed architecture of pose prediction network.** A shared encoder trunk processes an input image and branches into heads for predicting a set of candidate camera poses, as well as a score vector indicating the probability the image was acquired from each of the predicted cameras.

```

1 import torch as th
2 import pytorch3d.transforms as tf
3 def gen_rotation_matrix_from_xyz(xyz, in_plane=th.from_numpy(np.array
4 ([0.0, 0.0, 0.0])).cuda().float()):
5     cam_from = xyz
6     cam_to = th.from_numpy(np.zeros(3)).float().cuda()
7     tmp = th.from_numpy(np.array([0.0, 1.0, 0.0])).float().cuda()
8
9     diff = cam_from - cam_to
10    forward = diff / th.linalg.norm(diff)
11    crossed = th.cross(tmp, forward)
12    right = crossed / th.linalg.norm(crossed)
13    up = th.cross(forward, right)
14
15    R = th.stack([right, up, forward])
16    R_in_plane = tf.rotation_conversions.euler_angles_to_matrix(in_plane,
17 "XYZ")
18    return R_in_plane @ R

```

Code 1: Obtain rotation matrix from camera position

```

1 import torch as th
2 def lkat(eye, target, up):
3     forward = normalize(target - eye)
4     side = normalize(th.cross(forward, up))
5     up = normalize(th.cross(side, forward))
6
7     zero = th.zeros(1).float().cuda()
8     one = th.ones(1).float().cuda()
9     trans_v0 = th.cat([side[0:1], up[0:1], -forward[0:1], zero]) # (3,
10 1)
11    trans_v1 = th.cat([side[1:2], up[1:2], -forward[1:2], zero])
12    trans_v2 = th.cat([side[2:3], up[2:3], -forward[2:3], zero])

```

```

12 trans_v3 = th.cat([-th.dot(side, eye)[None], -th.dot(up, eye)[None],
13 th.dot(forward, eye)[None], one])
14 c2w = th.stack([trans_v0, trans_v1, trans_v2,trans_v3], dim=0)
15 return c2w

```

Code 2: Camera transformation with pointing to the origin.

```

1 import torch as th
2 eye_candidates = th.Tensor([[1,1,1],[1,-1,1],[-1,1,1],[-1,-1,1],
3 [1,1,1],[1,-1,1],[-1,1,1],[-1,-1,1],
4 [1,1,1],[1,-1,1],[-1,1,1],[-1,-1,1]]).cuda()
5 r = th.tensor([4.0]).float().cuda()
6 target = th.from_numpy(np.zeros(3)).float().cuda()
7 up = th.from_numpy(np.array([0.0, 1.0, 0.0])).float().cuda()
8
9 # h: the output of the last residual block in pose prediction model.
10 h1 = linear1(h.squeeze(-1)).squeeze() # (12*3, )
11 h2 = linear2(h.squeeze(-1)).squeeze() # (12, ), score vector
12
13 zero = th.zeros(1).float().cuda()
14 init_cam_pos = th.cat([zero, zero, r] # (3, 1)
15 all_poses = []
16
17 for index in range(12):
18 h3= th.sigmoid(h1[3*index:index*3+3])
19 h3 = th.diag(eye_candidates[index])@h3
20 h3 = h3/th.linalg.norm(h3)
21 R1 = gen_rotation_matrix_from_xyz((h3)
22 eye = (R1 @ init_cam_pos)
23 look_at1 = lkat(eye, target, up)
24 pose1 = th.eye(4).float().cuda()
25 pose1[:3, :3] = look_at1[:3, :3]
26 pose1[:3, 3] = -look_at1[:3, :3] @ look_at1[3, :3]
27 all_poses.append(pose1[:3,:4])
28 all_poses = th.stack(all_poses, 0)
29 return all_poses, h2 #pose distribution (12,3,4), scores (12, 1)

```

Code 3: Pose distribution prediction with 12× camera candidates

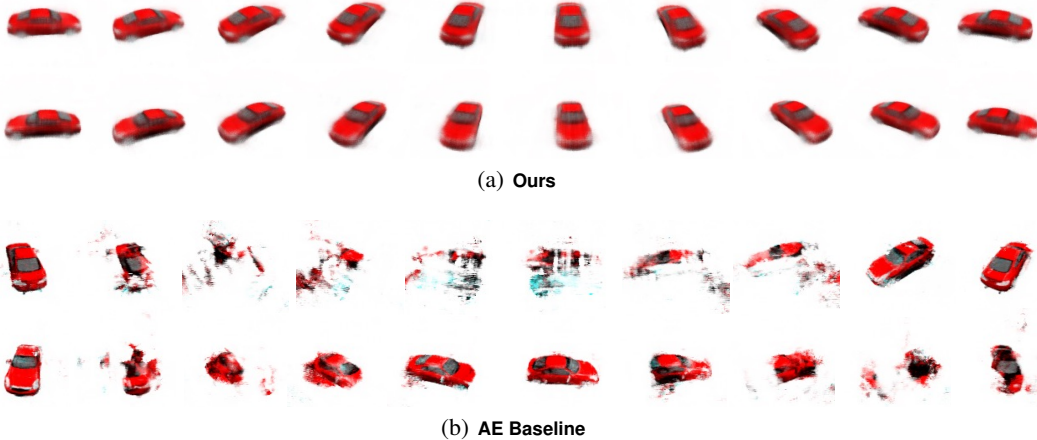


Figure 13: **Diffusion training benefits novel view synthesis on Car3D.** Our system, wrapped within a DDPM for training, significantly outperforms the same architecture trained as a simple autoencoder (AE). Training with the more challenging denoising task yields more robust generalization for the pose prediction network and NeRF scene representation.