

GD²-NeRF: Generative Detail Compensation via GAN and Diffusion for One-shot Generalizable Neural Radiance Fields

Xiao Pan, Zongxin Yang*, Shuai Bai, Yi Yang

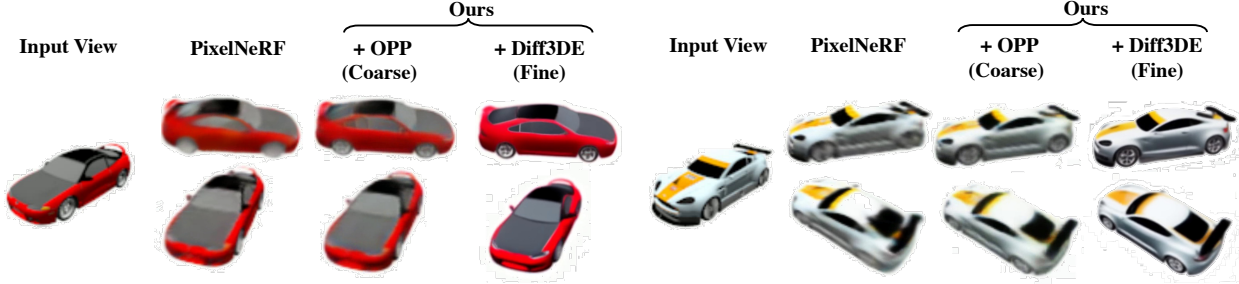


Fig. 1. Given a single reference image, our method GD²-NeRF synthesizes novel views with **vivid plausible details** in an **inference-time finetuning-free** manner. It is a coarse-to-fine generative detail compensation framework composed of OPP and Diff3DE. OPP first injects the GAN model into existing OG-NeRF pipelines, *e.g.*, PixelNeRF [1], for achieving in-distribution detail priors. Then, Diff3DE further incorporates the out-distribution detail priors from the pre-trained diffusion models [2], [3]. **We highly recommend readers to check our video demos for more intuitive comparisons.**

Abstract—In this paper, we focus on the One-shot Novel View Synthesis (O-NVS) task which targets synthesizing photo-realistic novel views given only one reference image per scene. Previous One-shot Generalizable Neural Radiance Fields (OG-NeRF) methods solve this task in an inference-time finetuning-free manner, yet suffer the blurry issue due to the encoder-only architecture that highly relies on the limited reference image. On the other hand, recent diffusion-based image-to-3d methods show vivid plausible results via distilling pre-trained 2D diffusion models into a 3D representation, yet require tedious per-scene optimization. Targeting these issues, we propose the GD²-NeRF, a Generative Detail compensation framework via GAN and Diffusion that is both inference-time finetuning-free and with vivid plausible details. In detail, following a coarse-to-fine strategy, GD²-NeRF is mainly composed of a One-stage Parallel Pipeline (OPP) and a 3D-consistent Detail Enhancer (Diff3DE). At the coarse stage, OPP first efficiently inserts the GAN model into the existing OG-NeRF pipeline for primarily relieving the blurry issue with in-distribution priors captured from the training dataset, achieving a good balance between sharpness (LPIPS, FID) and fidelity (PSNR, SSIM). Then, at the fine stage, Diff3DE further leverages the pre-trained image diffusion models to complement rich out-distribution details while maintaining decent 3D consistency. Extensive experiments on both the synthetic and real-world datasets show that GD²-NeRF noticeably improves the details while without per-scene finetuning.

Index Terms—One-shot novel view synthesis, generalizable neural radiance fields, 3D reconstruction, GAN, diffusion model.

I. INTRODUCTION

One-shot Novel View Synthesis (O-NVS) is a long-standing problem in computer vision and graphics which targets on synthesizing photo-realistic novel views of a scene given a single

reference image. An important technology solving this task is the One-shot Generalizable Neural Radiance Fields (OG-NeRF) which trains image-conditioned NeRF across scenes for learning general 3D priors and can generalize to a new scene by a single feed-forward pass, *i.e.*, *inference-time finetuning-free*.

However, 1) the existing OG-NeRF methods [1], [4], [5] mainly suffer the *blurry issue* since their encoder-only architectures highly rely on the reference images that contain limited information. For instance, [1], [4] first encodes the reference image into a 2D feature map and then indexes the condition features via pixel-wise projection. It works well when the target view is close to the reference view, yet tends to get blurry as the view difference becomes larger since the reference image can provide limited or even misleading scene information, *e.g.*, as shown by the upper row of Fig. 2, the query point requires the wheel information from the right view while the misleading body information from the back view is projected. 2) On the other hand, recent advances of diffusion-based image-to-3d methods [6], [7] show vivid plausible novel view results via distilling the 2D generative priors from pre-trained diffusion models [2], [3] into a 3D representation, yet requires tedious *per-scene optimization*.

Targeting these issues, we explore an O-NVS framework that is both *inference-time finetuning-free* and with *vivid plausible outputs*. To this end, we propose the GD²-NeRF, a coarse-to-fine generative detail compensation framework that hierarchically includes GAN and pre-trained diffusion models into OG-NeRF. GD²-NeRF is mainly composed of a One-stage Parallel Pipeline (OPP) that captures in-distribution priors via GAN model and a Diffusion-based 3D-consistent Enhancer (Diff3DE) that injects out-distribution priors from pre-trained diffusion models [2], [3], as illustrated by the second row of Fig. 2. In detail:

(i) *Coarse-stage OPP (GAN model)*. At the coarse stage,

Xiao Pan, Zongxin Yang, and Yi Yang are with the ReLER Lab, CCAI, Zhejiang University, Hangzhou, 310000, China, and Shuai Bai is with the Alibaba DAMO Academy, Hangzhou, 310000, China. (email: xiaopan@zju.edu.cn, yangzongxin@zju.edu.cn, baishuai.bs@alibaba-inc.com, yangyics@zju.edu.cn)

* Corresponding author.

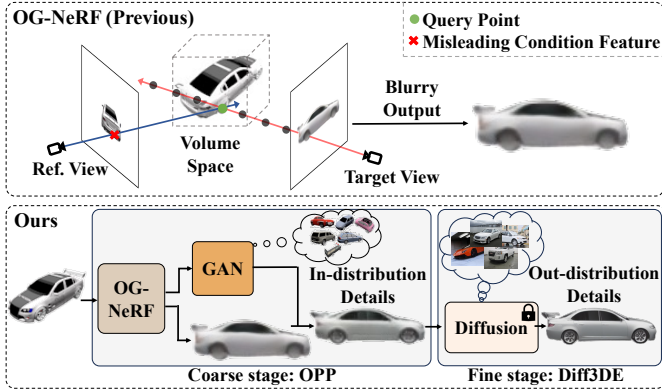


Fig. 2. Comparison between the existing encoder-only OG-NeRF and our generative detail compensation perspective (§ I). OG-NeRF suffers the blurry issue due to the projected misleading features while we propose to complement the object details via the prior learned by the generative model.

we intend to devise a pipeline that efficiently injects the GAN model into the existing OG-NeRF pipeline to primarily relieve the blurry issue using in-distribution detail priors captured from the training dataset. Targeting this, a naive solution is to directly build the GAN model on top of the OG-NeRF model tandemly, either in a two-stage or one-stage manner, as illustrated by the first two rows in Fig. 3 and will be detailed in § IV. However, though the *sharpness* (LPIPS, FID) is significantly improved, we empirically find it hard to maintain the *fidelity* (PSNR, SSIM), even with the more coherent one-stage tandem pipeline.

To address such contradiction between *fidelity* and *sharpness*, we further propose the One-stage Parallel Pipeline (OPP) that integrates the OG-NeRF and GAN model in a unified parallel framework, as shown by the bottom part of Fig. 3. It is built on the one-stage tandem pipeline with the proposed Dual-Paradigm Structure (DPS), Confidence Radiance Fields (CoRF), and Dual-Paradigm Fusion (DPF). With DPS, the OG-NeRF model and the GAN model can be optimized in parallel within a single framework. Then, CoRF takes the occlusion information as input and predicts a confidence map which adaptively gives the blurry part with low confidence. Finally, DPF aggregates the outputs from two paradigms via the learned confidence map.

(ii) *Fine-stage Diff3DE (diffusion model)*. Due to the limited size and quality of the training datasets, we find the in-distribution prior at the coarse stage is not enough for vivid outputs with rich plausible details. Therefore, at the fine stage, to break through such limitation, we turn to the large-scale diffusion models [2], [3] pre-trained on billions of high-quality images for more vivid out-distribution details.

However, naively using such an image diffusion model to process the rendered views from the coarse stage frame-by-frame leads to poor 3D consistency. Targeting this issue and inspired by the recent advances in zero-shot diffusion-based video editing methods [8], [9], we propose the Diffusion-based 3D Enhancer (Diff3DE). The main idea of Diff3DE is to first ensure the consistency between several keyframes dispersed around the dome. Then, given an arbitrary target view, the information from nearby keyframes is aggregated in the feature space based on view information, formulating a robust 3D-

consistent enhancer.

Extensive experiments on both the synthetic and real-world datasets show that, without any inference-time finetuning, 1) our OPP shows noticeable improvements over the baseline methods with balanced sharpness and fidelity while with little additional cost, and 2) Diff3DE can further compensate rich plausible details with decent 3D-consistency.

Our contributions are summarized as follows:

- We devise a coarse-to-fine generative detail compensation framework, GD²-NeRF, for O-NVS task that is both inference-time finetuning-free and with vivid plausible outputs.
- Our coarse-stage method OPP (§ IV) effectively inserts the GAN model into the existing OG-NeRF pipeline to primarily relieve the blurry issue with a good balance between fidelity and sharpness.
- To our best knowledge, our fine-stage method Diff3DE (§ V) makes the early attempt to directly use the pre-trained diffusion model as a 3D-consistent enhancer without any further finetuning.

II. RELATED WORKS

A. One-shot Novel View Synthesis

Recently, with the rapid development of the 3D computer vision community, there exist different technologies that can solve the one-shot novel view synthesis (O-NVS) task, though under different settings, including OG-NeRF [1], [4], [5], [10], Geometry-free Methods [11]–[13], 3D GAN [14]–[16], and Large-model-based Image-to-3D [6], [7], [17], [18].

Our work is motivated from the OG-NeRF perspective yet not limited by it. Specifically, we target relieving the blurry issue of existing OG-NeRF methods while maintaining its nice property of *inference-time finetuning-free*. However, in contrast to previous OG-NeRF methods that mainly focus on improving the in-distribution details from the limited dataset, we make the early attempt to also include the out-distribution details from the powerful diffusion models [2], [3] that pre-trained on billions of high-quality images for getting vivid plausible outputs. In the following paragraphs, we will distinguish between these technologies in detail, and we list comparisons with several representative methods in Tab. I.

OG-NeRF. The original NeRF technology [19] overfits the specific scene with tens or hundreds of posed input views and requires per-scene optimization. Targeting these issues, generalizable NeRF [1], [4], [5] is proposed which learns the general 3D prior across multiple scenes given very sparse reference images, which can naturally be applied to the O-NVS task.

(i) *Implicit condition*. Early works [20], [21] employ the implicit condition paradigm, which *implicitly* encode the scene information into the latent code based on the auto-decoder framework and requires tedious test-time optimization to find the latent code for new scenes.

(ii) *Explicit condition*. Another line of works [1], [4], [5] construct the condition *explicitly* with the help of an encoder module which extracts condition features from the reference

TABLE I
DISTINGUISH BETWEEN EXISTING REPRESENTATIVE TECHNOLOGIES
THAT CAN SOLVE THE O-NVS TASK UNDER DIFFERENT SETTINGS
(§ II).

Technology	Method	Inference-time Finetuning-free	GAN Model	Diffusion Model
3D GAN	EG3D-PTI [16]	✗	✓	✗
	Pix2NeRF [15]	✓	✓	✗
Geometry-free	3DiM [13]	✓	✗	✓
Image-to-3D	DietNeRF [18]	✗	✗	✗
	SinNeRF [17]	✗	✗	✗
	Zero-123-NVS [6]	✓	✗	✓
	Zero-123 [6]	✗	✗	✓
	Make-it-3D [7]	✗	✗	✓
OG-NeRF	NeRFDiff [10]	✗	✗	✓
	Ours	✓	✓	✓

images, and can generalize to a new scene by a single feed-forward pass. However, they inevitably suffer the blurry issues, especially when the target view is far from the source view, since they rely highly on the condition features from the limited reference image.

Targeting this issue, we propose GD²-NeRF, a coarse-to-fine framework to compensate for the details using generative models. At the coarse stage, we first inject the GAN model into the existing OG-NeRF pipeline to primarily relieve the blurry issue through learning in-distribution detail priors. Building on top of this, at the fine stage, we further exploit more vivid out-distribution priors from the pre-trained diffusion model [2], [3].

Recent work [10] also tries to relieve the blurry issue by adding a generative model. However, it co-trains a diffusion model with a conditional NeRF from scratch and needs tedious inference-time finetuning for each reference image. While our framework is finetuning-free and the pre-trained diffusion model [2], [3] on the large-scale high-resolution dataset is directly employed.

Geometry-free Methods. Several works [12], [13] also attempt to train a conditional model on posed images while without modeling the underlying geometry. For example, 3DiM [13] trains a pose-conditioned diffusion model on pairs of posed images, and inference in an auto-regression manner. However, though better at per-image quality (high FID), it fails to show smooth 3D consistency than the OG-NeRF methods due to the lack of geometry constraints.

3D GAN. In recent years, with the success of GANs in 2D image synthesizing [22]–[24] and the impressive performance of Neural Radiance Fields (NeRF) [19], a bunch of works [14], [16], [25], [26] include NeRF into GAN models as inductive bias for 3D-aware image synthesis. Typically, they can solve the O-NVS task via *GAN inversion* technology [16]. However, there mainly exist the following drawbacks:

(i) *Per-scene optimization.* It requires tedious optimization to find the corresponding latent code for each reference image. Though [15] further includes the encoder model into the existing framework [14] to obtain a conditional 3D GAN model, the performance is still unsatisfactory since it is trained on a collection of unposed images, *i.e.*, unsupervised.

(ii) *Specific category.* Similar to the traditional GAN methods, they usually only work on a specific category like cat,

car, etc.

In contrast to these methods [14], [16], [25], [26], our method does not require per-scene optimization given reference images and can generalize across different categories, even for real-world complex scenes.

Large-model-based Image-to-3D. As the recent advances on pre-trained large models [2], [3], [27], [28], a bunch of works [6], [7], [17], [18] explore using them to solve the O-NVS task in a per-scene optimization manner for getting plausible 3D representations.

(i) *DINO&CLIP-based.* Early attempts DietNeRF [18] and SinNeRF [17] use DINO [27] or CLIP [28] to constrain the distance between the rendered novel views and the reference view in the feature space. However, the feature space constraint struggles to provide fine-grained information, and [17] only works in nearby views.

(ii) *Diffusion-based.* Recently, several works [6], [7], [29], [30] attempt to lift the fine-grained 2D generative prior in pre-trained diffusion models [2], [3] to plausible 3D representations. For instance, Zero123 [6] first finetunes the latent diffusion [2] on a synthetic dataset [31] to inject the viewpoint condition (Zero123-NVS). Though Zero123-NVS can work in an inference-time finetuning-free manner, it achieves poor 3D consistency considering the undeterministic nature of diffusion models. Therefore, it further optimizes a 3D representation with the finetuned diffusion model using SJC [32] framework.

Different from all the methods above [6], [7], [17], [18], [29], [30], our framework requires NO per-scene optimization, and the pre-trained diffusion model is also fixed with no further finetuning in our fine-stage method 3DE.

B. Diffusion-based Video Editing

With the success of text-to-image diffusion models [2], [3], recent works [8], [9], [33]–[36] try to solve the video editing task via the pre-trained text-to-image diffusion models, either with finetuning [33] or in a zero-shot manner [8], [9], [34]–[36]. The main challenge for extending an image diffusion model to a video editing task is to ensure the consistency between video frames. Early works [33]–[35] mainly ensure the global appearance consistency by inflating the self-attention module to the temporal cross-attention module. Besides, [8] further includes the optical flow as the correspondence constraint, while [9] directly uses the correspondences calculated during DDIM inversion. However, all of them take consecutive video sequences with similar contents as inputs and can not process an arbitrary view in a 3D manner.

Differently, in this work, our Diff3DE makes the early attempts to extend the existing methods to a 3D-aware detail enhancer that maintains 3D consistency between different views that contain large content variance (*e.g.*, front and back views of a car) and supports the process of an arbitrary view.

III. PRELIMINARY

Overview. In this section, we first briefly introduce the Neural Radiance Fields (NeRF) [19] in § III-A and its one-shot generalizable variant [1] in § III-B. Then, we analysis the limitation of existing OG-NeRF in § III-C. Finally, we introduce the GAN model in § III-D and diffusion model in § III-E.

A. Neural Radiance Fields

NeRF parameterizes the 3D volume space as a continuous implicit function $f_{nerf}(\cdot)$ represented by a neural network, *e.g.*, multi-layer perceptron (MLP). Given a target view with pose \mathbf{P}_t , it queries the 3D volume space via marching a ray $\mathbf{r}(z) = \mathbf{o} + z\mathbf{d}$ for each pixel on its image plane, where $\mathbf{o} \in \mathbb{R}^3$ represents the camera center, $\mathbf{d} \in \mathbb{R}^3$ represents the ray unit vector, and $z \in \mathbb{R}^1$ is the depth between a pre-defined bounds $[z_n, z_f]$. Then, for each 3D point $\mathbf{x} \in \mathbb{R}^3$ on a marched ray, its density $\sigma \in \mathbb{R}^1$ and RGB color $\mathbf{c} \in \mathbb{R}^3$ is predicted by:

$$(\sigma, \mathbf{c}) = f_{nerf}(\mathbf{x}, \mathbf{d}). \quad (1)$$

After that, the predicted RGB color \mathbf{c} along a ray \mathbf{r} is accumulated via the differentiable volumetric rendering operation:

$$\hat{\mathbf{C}}(\mathbf{r}) = \int_{z_n}^{z_f} T(z) \sigma(z) \mathbf{c}(z) dz, \quad (2)$$

where $T(z) = \exp(-\int_{z_n}^z \sigma(s) ds)$ represents the probability that the ray travels from z to z_n . NeRF employs a hierarchical coarse-to-fine strategy for sampling discrete 3D points along rays and then approximates the integral using numerical quadrature [37]. Finally, for each pixel, the accumulated RGB color $\hat{\mathbf{C}}(\mathbf{r})$ is supervised by the ground truth color $\mathbf{C}(\mathbf{r})$ through the mean squared error:

$$\mathcal{L}_{NeRF} = \frac{1}{|\mathcal{R}(\mathbf{P}_t)|} \sum_{\mathbf{r} \in \mathcal{R}(\mathbf{P}_t)} \|\hat{\mathbf{C}}(\mathbf{r}) - \mathbf{C}(\mathbf{r})\|_2^2, \quad (3)$$

where $\mathcal{R}(\mathbf{P}_t)$ is the set of all marched rays from target pose \mathbf{P}_t .

B. One-shot Generalizable NeRF

NeRF is optimized per scene and requires tens or hundreds of posed input views and a time-consuming optimization process to memorize the scene. To relieve this issue and realize O-NVS, One-shot Generalizable Neural Radiance Fields (OG-NeRF) [1] is proposed to learn the general 3D prior across multiple scenes. At the core of OG-NeRF is to take an additional condition feature extracted from the reference image \mathbf{I}_s as input. A typically used method for getting the condition feature [1], [4] is to project the query point \mathbf{x} back to the feature map of the reference image, and fetch a feature via interpolation:

$$(\sigma, \mathbf{c}) = f_{ognerf}(\mathbf{x}, \mathbf{d}, W(\pi(\mathbf{x}))), \quad (4)$$

where $W = E(\mathbf{I}_s)$ is the extracted feature map via encoder E , and $\pi(\mathbf{x})$ indicates the projection of \mathbf{x} on the reference image plane. Note that, for better generalization ability, \mathbf{x} and \mathbf{d} here are under the coordinate of the reference view, *i.e.*, the relative one, instead of the world coordinate.

C. Limitations of Existing OG-NeRF Methods

The main drawback of such an encoder-only paradigm is the high reliance on the reference image, which may include misleading information, especially for the projection-based ones [1], [4] (see the upper part of Fig. 2). To address this

issue, we propose the coarse-to-fine generative detail compensation perspective. Specifically, at the coarse stage, we first propose the OPP that efficiently insert a GAN model into the existing OG-NeRF pipeline for capturing the in-distribution object details. Then, at the fine stage, with Diff3DE we further leverage the out-distribution detail priors from the pre-trained diffusion model [2], [3] in a 3D-consistency manner for more vivid outputs with plausible details.

D. GAN Model

We inject the GAN model to the OG-NeRF pipeline at the coarse stage for capturing primary in-distribution detail priors from the training dataset. For the training of the GAN model, we employ the commonly used non-saturating GAN objective [38] with R_1 gradient penalty [39]:

$$\mathcal{L}_{GAN} = \mathbb{E}_{\mathbf{I}_{in} \sim p_{in}} [D(G(\mathbf{I}_{in}))] + \mathbb{E}_{\mathbf{I}_{real} \sim p_{real}} [-D(\mathbf{I}_{real})], \quad (5)$$

where p_{in} represents the distribution of the generator input, p_{real} indicates the distribution of the training data, and the formulation of R_1 penalty is omitted here for simplicity.

Except for the GAN objective, we also employ the perceptual loss \mathcal{L}_{PER} [40] and MSE loss \mathcal{L}_{MSE} [26] for the reconstruction from $G(\mathbf{I}_{in})$ to its corresponding ground truth \mathbf{I}_{real} . Therefore, the overall objective \mathcal{L}_G for the GAN training is:

$$\mathcal{L}_G = \lambda_{GAN} \mathcal{L}_{GAN} + \lambda_{PER} \mathcal{L}_{PER} + \mathcal{L}_{MSE}, \quad (6)$$

where λ_{GAN} and λ_{PER} are the weights for GAN loss and perceptual loss, separately.

E. Diffusion Model

Diffusion probabilistic model has been broadly researched recently [2], [3], [41], [42] due to its strong generative power. It approximates the data distribution via progressively removing noises from a Gaussian *i.i.d* noised image. Stable Diffusion [2] makes the early attempt to operate on the latent space of a pre-trained auto-encoder to efficiently get high-resolution results. Building on this, ControlNet [3] enables more accurate and flexible controls via incorporating additional condition signals with a residual architecture, *e.g.*, edge, pose, etc.

DDIM Inversion. DDIM is a deterministic denoising-stage sampling algorithm proposed by [43]. It can be used in a reversed order to find the corresponding noises of an image in a *non-optimization manner*, which is commonly used in image/video editing tasks [8], [9], [44] for primarily maintaining the original contents of input images.

Inflated Self-attention. Usually, a U-Net ϵ_θ with attention blocks [2], [3] is employed for predicting the noise. To improve the temporal consistency when processing a sequence of frames, recent methods [8], [9], [33] inflate the original self-attention to incorporate information from other frames. In detail, given the sequence of projected query features $\{\mathbf{Q}_i\}_{i=1}^k$, key features $\{\mathbf{K}_i\}_{i=1}^k$, and value features $\{\mathbf{V}_i\}_{i=1}^k$, the Inflated Self-Attention (ISA) is calculated as follows:

$$\phi_i = \text{Softmax}\left(\frac{\mathbf{Q}_i[\mathbf{K}_1, \dots, \mathbf{K}_k]^T}{\sqrt{d}}\right)[\mathbf{V}_1, \dots, \mathbf{V}_k], \quad (7)$$

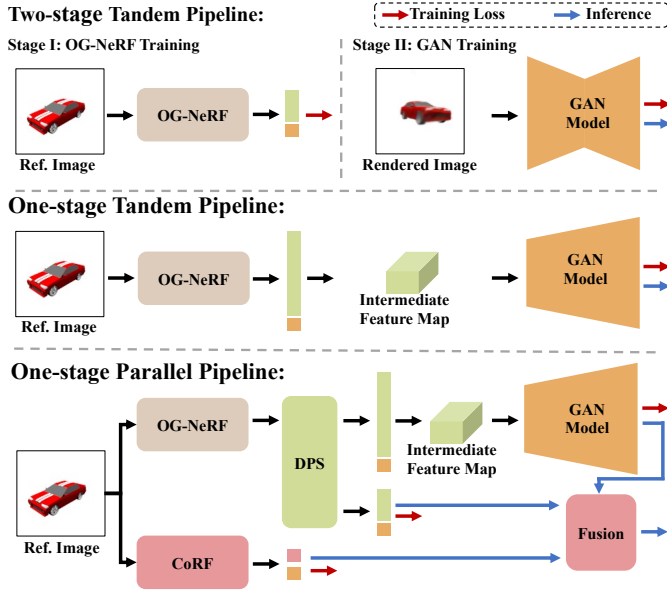


Fig. 3. Overview of two basic tandem pipelines (§ IV-A, § IV-B) and our proposed One-stage Parallel Pipeline (OPP, § IV-C) for integrating the GAN model into the OG-NeRF framework at the COARSE STAGE.

where ϕ_i is the ISA output tokens of the i -th frame.

IV. COARSE STAGE: FROM TANDEM TO PARALLEL PIPELINES

Overview. This section introduces the coarse stage method that injects the GAN model into the OG-NeRF pipeline. We first briefly analyze two basic tandem pipelines that directly build the GAN model on top of the OG-NeRF model (§ IV-A, § IV-B), as shown by Fig. 3. Then, in § IV-C, we further present a stronger One-stage Parallel Pipeline (OPP) that solves the contradiction between sharpness and fidelity efficiently. The pipeline of OPP is illustrated in Fig. 4.

A. Two-stage Tandem Pipeline

Different from OG-NeRF which supervises at the pixel level, the GAN model requires the input to be semantically complete at the image level, ideally, the full-sized image. However, rendering out the full-sized image during training is intractable considering the memory-intensive rendering process. A naive solution is to independently train the OG-NeRF at the first stage. Then, in the second stage, the full-sized novel view images are rendered with the previously trained OG-NeRF (fixed and inference only) pixel by pixel, and a GAN model is trained with these prepared full-sized images and their corresponding ground truth images, as illustrated by the first row in Fig. 3.

The main drawbacks of such a pipeline are two-fold: (i) Though the sharpness (*i.e.*, details) can be improved with such a pipeline, the fidelity of the generated details can hardly be guaranteed. The independently trained GAN model is easily biased towards the sharpness details without the joint optimization with the OG-NeRF model which constrains the fidelity. (ii) The two-stage training procedure is tedious.

B. One-stage Tandem Pipeline

To relieve the memory cost and facilitate the end-to-end training, the key is to reduce the rendered output size from $H \times W$ to $H_m \times W_m$. However, directly reducing the output size will inevitably cause the loss of 3D information from the volume space. To compensate for it, a feasible solution is to simultaneously improve the output dimension from $\mathbf{c} \in \mathbb{R}^3$ to a higher value $\mathbf{c}_m \in \mathbb{R}^{h_m}$, $h_m > 3$, as in [26] and illustrated by the second row of Fig. 3. Then, similar as Eq. 2, \mathbf{c}_m is accumulated along the ray \mathbf{r} as follows:

$$\hat{\mathbf{c}}_m(\mathbf{r}) = \int_{z_n}^{z_f} T(z)\sigma(z)\mathbf{c}_m(z)dz. \quad (8)$$

Notably, for the sampling of rays, different from the pixel-wise supervision method [1] which randomly samples rays among all target poses around the object, we employ the *grid sampling strategy* which first randomly samples a target pose, and then split its corresponding image plane of size $H \times W$ into $H_m \times W_m$ grids. After that, the rays tracing through the grid centers are sampled for accumulating with Eq. 8 and get the intermediate low-resolution yet high-dimensional feature map $\hat{\mathbf{I}}_m \in \mathbb{R}^{H_m \times W_m \times h_m}$ via reshaping. Finally, $\hat{\mathbf{I}}_m$ is sent to a light-weight up-sampling model $G(\cdot)$ for the full-sized output $\hat{\mathbf{I}}_t^G \in \mathbb{R}^{H \times W \times 3}$, which is finally supervised with a discriminator $D(\cdot)$.

Compared with the naive two-stage tandem pipeline, this end-to-end fashion makes the GAN model benefits more fidelity from the OG-NeRF model which extracts 3D information directly from the volume space. Therefore, it performs better than the two-stage tandem pipeline, especially at the fidelity metrics, *e.g.*, PSNR, SSIM.

C. One-stage Parallel Pipeline

After including the GAN model as a tandem pipeline, the *sharpness* (*e.g.*, LPIPS, FID) is significantly improved since the GAN model can capture the prior object details from the abundant training data. However, we find it hard to maintain the *fidelity* (*e.g.*, PSNR, SSIM) as well as the independently trained OG-NeRF model, even with the more coherent one-stage tandem pipeline.

To tackle such contradiction between two paradigms, we further present the simple yet effective One-stage Parallel Pipeline (OPP) that is built on the one-stage tandem pipeline with the proposed Dual-Paradigm Structure (DPS), Confidence Radiance Fields (CoRF), and Dual-Paradigm Fusion (DPF), as illustrated in Fig. 4. DPS makes it possible for these two paradigms to be optimized in parallel within a single framework, and CoRF further learns a confidence map that can adaptively give the blurry part with lower confidence. Finally, DPF integrates the outputs from two paradigms effectively with the learned confidence map.

Dual-Paradigm Structure. In the one-stage tandem pipeline, the output of the OG-NeRF MLP is represented as $[\mathbf{c}_m; \sigma] \in \mathbb{R}^{h_m+1}$, where $[\cdot]$ indicates the concatenation operation, $\mathbf{c}_m \in \mathbb{R}^{h_m}$ is the high-dimensional hidden color, and $\sigma \in \mathbb{R}^1$ is the density. Then, we send \mathbf{c}_m to an additional fully connected layer $FC(\cdot) : \mathbb{R}^{h_m} \rightarrow \mathbb{R}^3$ as follows:

$$\mathbf{c} = FC(\mathbf{c}_m), \quad (9)$$

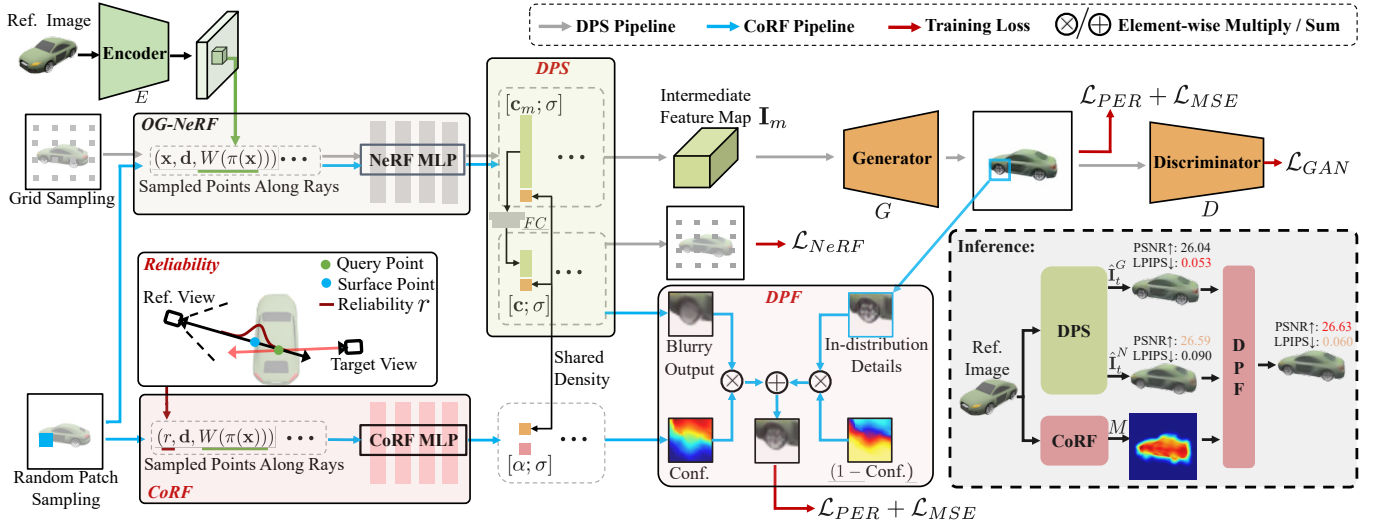


Fig. 4. **Overview of the COARSE-STAGE method OPP for including in-distribution details from the training data (§ IV-C).** It is built on the one-stage tandem pipeline (the first row) and efficiently integrates the GAN and OG-NeRF models in a unified parallel framework with DPS, CoRF, and DPF.

where $\mathbf{c} \in \mathbb{R}^3$ is the RGB color. After that, \mathbf{c} and \mathbf{c}_m are accumulated by Eq. 2 and Eq. 8 with the shared density σ and get $\hat{\mathbf{C}} \in \mathbb{R}^3$ and $\hat{\mathbf{C}}_m \in \mathbb{R}^{h_m}$, separately. Finally, $\hat{\mathbf{C}}$ is supervised with the ground truth pixel color with Eq. 3, and $\hat{\mathbf{C}}_m$ marched from the same target pose formulates the intermediate feature map $\hat{\mathbf{I}}_m$, which is sent to the up-sampling module for the GAN supervision.

With DPS, the OG-NeRF and GAN model can be optimized in parallel with a single training process. Intuitively, \mathbf{c}_m can be seen as the mapping of RGB color \mathbf{c} in a high-dimensional feature space. We inverse it back to the RGB space with a lightweight fully connected layer for the NeRF-style supervision. The shared density also profits the communications between two paradigms.

Confidence Radiance Fields. An ideal solution to integrate the two paradigms is to automatically detect the blurry part on the OG-NeRF output and then complement the corresponding details from the GAN output. Motivated by this, we propose the novel Confidence Radiance Fields (CoRF) that gives each pixel a confidence score reflecting the degree of clarity.

Specifically, for each sampled point \mathbf{x} , we assume that the reliability of the projected condition feature is determined by the distance between \mathbf{x} and the projected surface point \mathbf{s} since it reflects the occlusion information, as shown in Fig. 4. Therefore, we define reliability as $r = \mathcal{G}(z_{\mathbf{x}} - z_{\mathbf{s}})$, where \mathcal{G} is a gaussian function, and $z_{\mathbf{x}}$, $z_{\mathbf{s}}$ are the depth of \mathbf{x} and \mathbf{s} from the reference view, respectively. Notably, $z_{\mathbf{x}}$ can be achieved by simply using the pose information of the reference view, while for $z_{\mathbf{s}}$, we first render a coarse depth map (e.g., 16×16) from the reference view and resize it to the full size (e.g., 128×128), then index the depth via projection.

In a nutshell, the final CoRF is represented as:

$$\alpha = f_{corf}(r, \mathbf{d}, W(\pi(\mathbf{x}))), \quad (10)$$

where $\alpha \in [0, 1]$ represents the confidence score. Then, similar to Eq. 2, the confidence score along a ray \mathbf{r} is accumulated with the differentiable volumetric rendering using the same

density value as in DPS, and get the final confidence score $\hat{\alpha}(\mathbf{r}) \in \mathbb{R}^1$.

Dual-Paradigm Fusion. After that, for each ray \mathbf{r} , we fuse the RGB values from two paradigms with:

$$\hat{\mathbf{C}}'(\mathbf{r}) = \hat{\mathbf{C}}(\mathbf{r}) * \hat{\alpha}(\mathbf{r}) + \hat{\mathbf{I}}_t^G(\mathbf{r}) * (1 - \hat{\alpha}(\mathbf{r})), \quad (11)$$

where $\hat{\mathbf{C}}(\mathbf{r})$, $\hat{\mathbf{I}}_t^G(\mathbf{r}) \in \mathbb{R}^3$ are the RGB values from OG-NeRF and GAN model, respectively, and $\hat{\mathbf{C}}'(\mathbf{r})$ is the final fused one where we conduct CoRF training objectives on.

Training & Inference. During the training stage, since we expect the DPF output to have both high sharpness and fidelity, except for the *grid sampling* used for DPS training, we additionally employ a *random patch sampling strategy* for the CoRF learning, which can output a semantic patch that supports the perceptual supervision, as illustrated by Fig. 4. The overall loss for the training of OPP is:

$$\mathcal{L}_{OPP} = \mathcal{L}_G + \mathcal{L}_{NeRF} + \mathcal{L}_{CoRF}, \quad (12)$$

where $\mathcal{L}_{CoRF} = \mathcal{L}_{PER} + \mathcal{L}_{MSE}$. Considering the training stability, we empirically first train the two paradigms till convergence and then finetune the CoRF for several epochs with the two paradigms frozen.

During the inference stage of coarse-stage OPP, we have $\hat{\mathbf{I}}_t^G$ from the GAN model via grid sampling, $\hat{\mathbf{I}}_t^N$ from the OG-NeRF paradigm and the confidence map M from the CoRF with the full-sized sampling. Then, they are aggregated via DPF for the final output, as shown in Fig. 4. We emphasize that the rendering of grid sampling (e.g., 16×16) is quite efficient (over 40 times faster) compared with the full-sized rendering (e.g., 128×128), and the CoRF MLP is rather lightweight ($0.09M$), therefore the additional computational cost compared with the original OG-NeRF is rather small, which will be discussed in detail in the appendix.

V. FINE STAGE: DIFFUSION-BASED 3D ENHANCER

Motivation. With OPP, the coarse in-distribution details have been primarily compensated. However, considering the limited

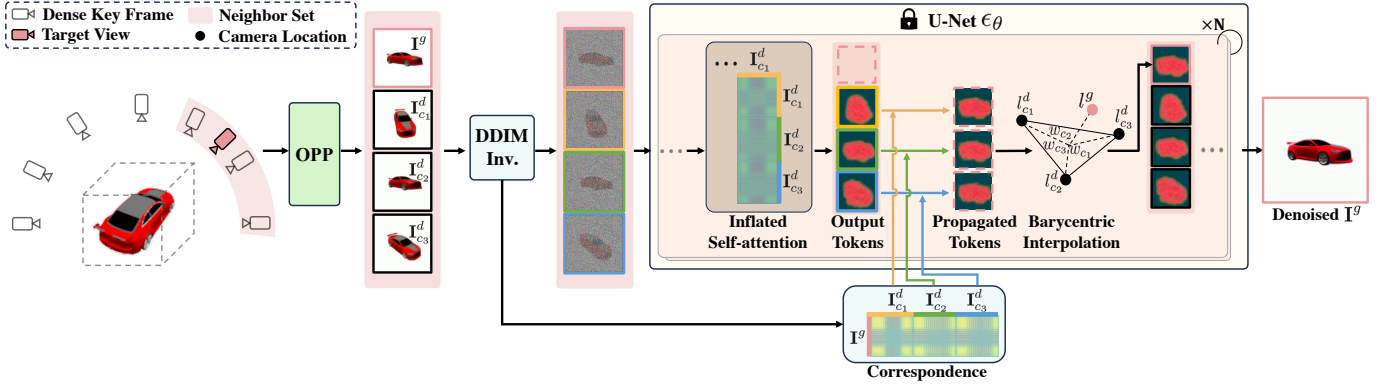


Fig. 5. **Overview of the FINE-STAGE method Diff3DE for including out-distribution details from the pre-trained diffusion model [2], [3] (§ V).** We first fix N_k dense keyframes around the dome. Then, for each target view, we select 3 neighbor keyframes based on the cosine similarity. For each diffusion time step and attention block, the output tokens of the target view are the barycentric interpolation of the propagated tokens from neighbor keyframes, using the correspondence calculated during DDIM inversion. The global 3D consistency is primarily achieved by the 3D-consistent constraint from OPP and further approximated by enforcing the local consistency for each neighbor area.

size and quality of the training dataset, the results are still unsatisfying in vivid detail. Thus, to break through such limitation, we turn to the rich out-distribution priors from the diffusion models [2], [3] pre-trained on billions of high-quality images.

A naive solution of using such an image diffusion model is to perform super-resolution [2], [3] on each rendered view from the OPP individually. Nevertheless, due to the lack of 3D constraints, this will lead to significant inconsistency between different views. Although the recent zero-shot video editing method [9] has explored maintaining the temporal consistency between edited video frames, it simply assumes the input video to contain very similar contents, and directly employing it in the 3D NVS task will lead to the following issues:

(i) *Dispersed sparse keyframes lead to blurry outputs.* At the core of [9] is to first maintain the consistency between several keyframes (e.g., 5) using Inflated Self-Attention (ISA) [33] and then propagate the U-Net self-attention output tokens of keyframes to nearby ones. A possible naive adaptation is to uniformly disperse the keyframes around the dome, however, due to the large variance of contents from dispersed sparse views, the ISA tends to get blurry outputs, as illustrated by Fig. 9. On the other hand, simply increasing the keyframe number will greatly increase memory usage, which is unaffordable.

(ii) *Unable to process arbitrary views.* A 3D enhancer should support the processing of an arbitrary given view. However, as a video processing method, [9] simply assumes the input videos are consecutive and calculates the propagation weights via the frame index, which is not suitable for the arbitrary view processing case.

Diffusion-based 3D Enhancer. Targeting the aforementioned issues, we propose the fine-stage method Diffusion-based 3D Enhancer (Diff3DE), a 3D extension of [9], as illustrated in Fig. 5. The main idea of Diff3DE is to relax the input of the original ISA from all the keyframes to *neighbor keyframe sets* selected based on *view distance*.

In detail, we first fix N_k dense keyframes $\{\mathbf{I}_i^d\}_{i=1}^{N_k}$ uniformly dispersed around the dome. Then, given a target view \mathbf{I}^g ,

we select its 3 neighbors $\{\mathbf{I}_{c_1}^d, \mathbf{I}_{c_2}^d, \mathbf{I}_{c_3}^d\}$ from $\{\mathbf{I}_i^d\}$ using the cosine similarity, which can be calculated by the provided camera poses. After that, the neighbor set $\{\mathbf{I}_{c_1}^d, \mathbf{I}_{c_2}^d, \mathbf{I}_{c_3}^d\}$ are taken as the input of the ISA to ensure the local consistency, and for each diffusion time step and U-Net attention block layer, their output tokens of the ISA module $\{\phi_{c_1}^d, \phi_{c_2}^d, \phi_{c_3}^d\}$ are further propagated to the target view using the correspondence calculated with original input frames during the DDIM inversion stage (see [9] for more details), formulating the propagated tokens $\{\phi_{c_1}^{d \rightarrow t}, \phi_{c_2}^{d \rightarrow t}, \phi_{c_3}^{d \rightarrow t}\}$. Finally, we perform a weighted sum on the propagated tokens using the barycentric interpolation:

$$w_{c_1}, w_{c_2}, w_{c_3} = \text{Bary}(l^g, l_{c_1}^d, l_{c_2}^d, l_{c_3}^d), \quad (13)$$

$$\phi^t = w_{c_1} * \phi_{c_1}^{d \rightarrow t} + w_{c_2} * \phi_{c_2}^{d \rightarrow t} + w_{c_3} * \phi_{c_3}^{d \rightarrow t},$$

where l indicates the camera location, $\text{Bary}(\cdot, \cdot, \cdot, \cdot)$ is the function for calculating barycentric weights, and ϕ^t is the final aggregated ISA output tokens for \mathbf{I}^g .

With Diff3DE, the processing of each view is determined by its pose instead of the frame index as in [9], making it work in a 3D manner. Notably, though the global consistency between all the keyframes is not explicitly constrained in Diff3DE, i.e., we do not send all the keyframes to ISA considering the computation cost, it is approximately maintained by: 1) the input frames are from the 3D-consistent OPP, which naturally provides certain 3D constraints, and 2) enforcing the local consistency for each neighbor area approximates the global consistency.

VI. EXPERIMENTAL RESULTS

A. Experimental Settings

ShapeNet Cars & Chairs. Following previous methods [1], [5], we evaluate on the synthetic large-scale ShapeNet benchmarks [45] with the Cars and Chairs categories following SRN [46], which contains 3514 cars and 6591 chairs with an image resolution of 128×128 . For each category, we train an individual model.

DTU Dataset. For the real-world DTU dataset, unless otherwise specified, we follow the split of [1] which includes 88

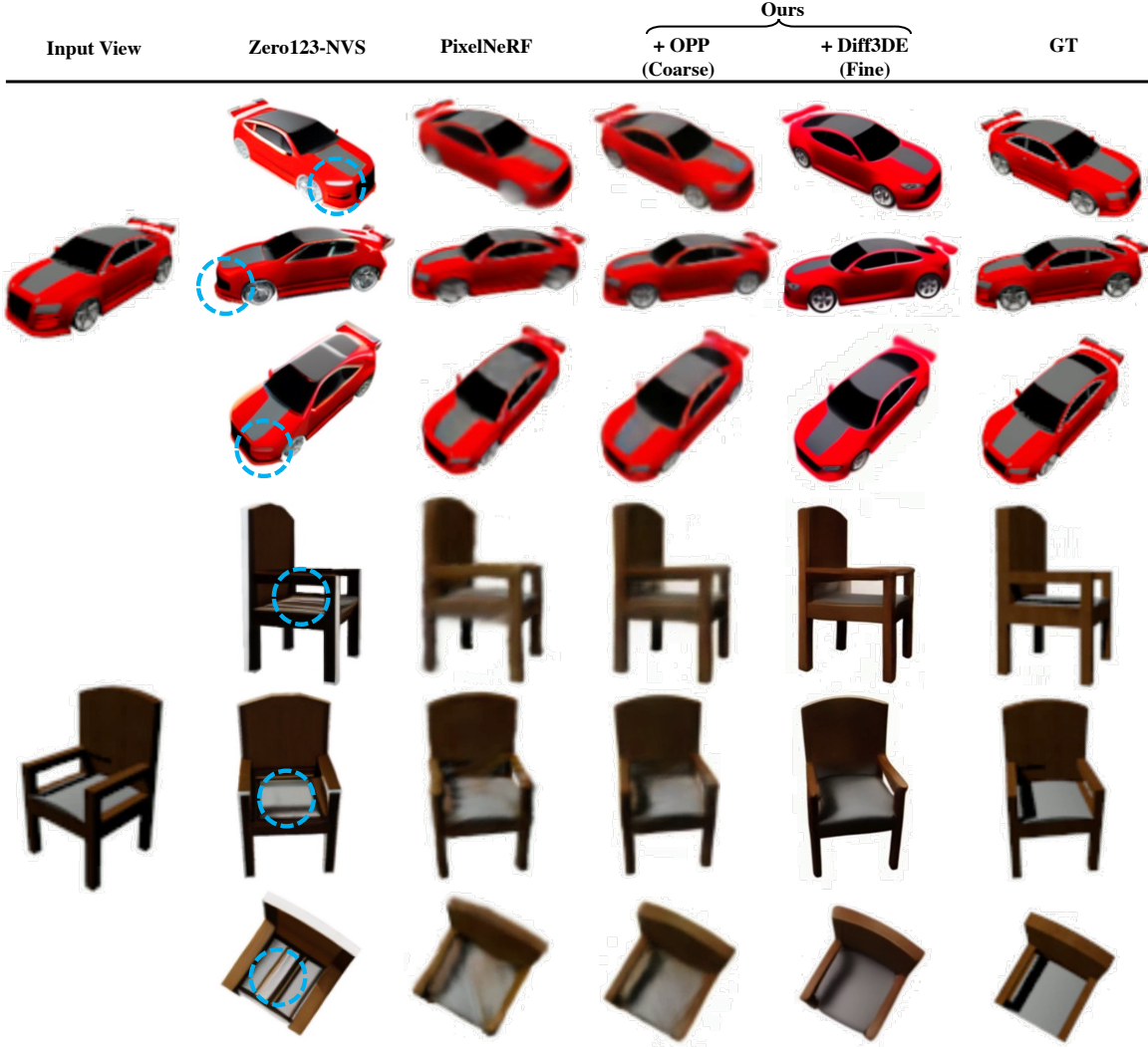


Fig. 6. **Qualitative comparisons with previous methods on ShapeNet Cars & Chairs (§ VI-C).** Zero123-NVS [6] shows significant inconsistency between different views (marked with blue circle), and PixelNeRF [1] gives blurry results. After including our GD²-NeRF framework, the details are greatly improved in a coarse-to-fine manner with decent 3D-consistency.

training scenes and 15 test scenes. We train at the 128×128 resolution and then resize to the original 300×400 resolution for comparison.

Out-distribution Metrics. For the final fine-stage out-distribution method Diff3DE, we evaluate the 3D consistency via rendering a consecutive video around the dome and calculating Pixel-MSE following [35], which is the averaged mean-squared pixel error between warped and original frames via optical flow [47]. Notably, since Diff3DE leverages the priors *outside the training dataset distribution*, and mainly focuses on achieving vivid plausible details with 3D-consistency, we do not calculate metrics with the ground truth from the dataset.

In-distribution Metrics. For the coarse-stage in-distribution method OPP, we report the commonly used PSNR and SSIM [48] as the fidelity metrics, while LPIPS [49] and FID [50] as the sharpness metrics.

B. Implementation Details

Architectures. (i) *OPP*. We directly take PixelNeRF [1] as the OG-NeRF part of our framework, *i.e.*, ResNet-34 as the encoder and the ResNet-like NeRF MLP. For the verification

of different pipelines, we employ the commonly used U-Net [51] as the generator for the two-stage pipeline, and the lightweight up-sampling module ($0.11M$) used in one-stage pipelines is in line with [26]. The discriminator in all the pipelines also shares the same architecture as in [26]. The CoRF is composed of three fully connected layers ($0.09M$).

(ii) *Diff3DE*. For Diff3DE, we employ the pre-trained ControlNet-Tile [3] as the diffusion model, which conditions on the text together with the input image to perform super-resolution.

Hyper Parameters. (i) *OPP*. For the intermediate feature map \mathbf{I}_m in one-stage pipelines, we set $H_m = W_m = 16$, $h_m = 128$, *i.e.*, 16×16 rays for grid sampling. The resolution of the random patch sampling is also set as 16×16 . For fair comparisons, we set the same training hyperparameters for all the pipelines. Specifically, for the generative learning objectives, we set $\lambda_{GAN} = 1e^{-3}$ and $\lambda_{PER} = 1e^{-2}$. To rigorously validate the effectiveness of our proposed method, we strictly keep the learning strategy of PixelNeRF such as the learning rate ($1e^{-4}$), optimizer (Adam solver [53]), and

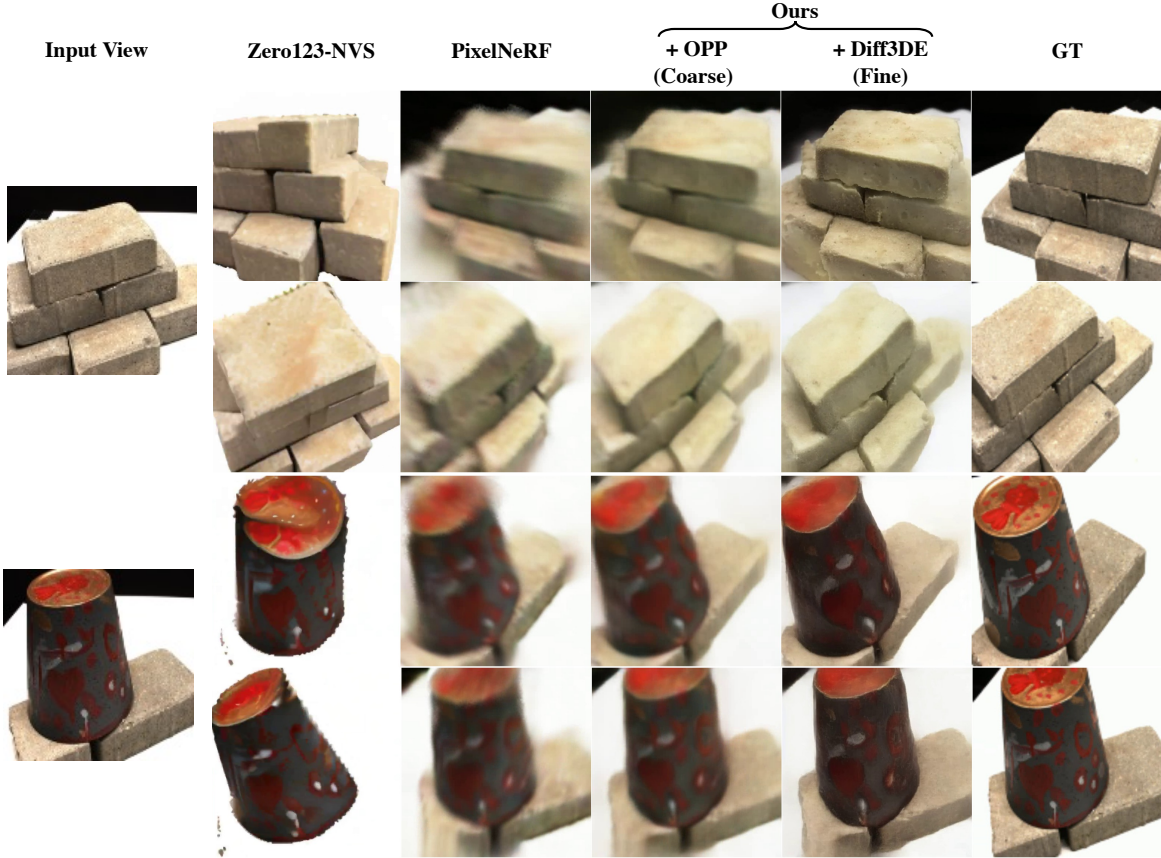


Fig. 7. **Qualitative comparisons with previous methods on the DTU dataset [52] (§ VI-C).** Notably, Zero123-NVS [6] only supports foreground objects, therefore, we mask out the foreground using their official code. Obviously, Zero123-NVS is sensitive to the estimated mask and shows significant inconsistency between different views with inaccurate geometry, while our method hierarchically includes the details in a coarse-to-fine manner with much better geometry and consistency.

sampling points per ray (64 coarse and 96 fine).

(ii) *Diff3DE*. The dense keyframe number N_k is set as 40, which uniformly dispersed around the dome. We resize the rendered images from the OPP module to 512×512 resolution as the Diff3DE inputs. The classifier-free guidance scale is set as 7.5 and the denoising step number is 25. For the ShapeNet Cars and Chairs categories, we use the simple text descriptions 'A car' and 'A chair', respectively; and for the DTU dataset, we use [54] to generate text descriptions.

C. Comparison with State-of-the-art

Out-distribution Baselines. Since the main goal of final fine stage is achieving vivid plausible details with 3D consistency in an inference-time finetuning-free manner via out-distribution priors, we choose recent Zero-123-NVS [6] as our main competitor, which is a finetuned latent diffusion model with viewpoint condition and can also synthesize novel views without per-scene finetuning using out-distribution priors.

In-distribution Baselines. The main in-distribution baseline of our method is PixelNeRF [1], which is taken as the OG-NeRF part of our GD²-NeRF. Additionally, we report the quantitative results of methods including 3D GAN methods [14]–[16], Geometry-free method [13], Large-model-based Image-to-3D methods [17], [18], and previous OG-NeRF methods [5], [20], [46]. Notably, VisionNeRF [4] and

NeRFDiff [10] are not listed here since they require much more computation than PixelNeRF [4], [10] (detailed in the appendix), and [10] requires tedious per-scene optimization using the co-trained diffusion model.

Qualitative Analysis. We perform qualitative comparisons with previous methods on ShapeNet Cars & Chairs in Fig. 6, and DTU dataset in Fig. 7. Obviously, Zero123-NVS shows significant inconsistency between different views together with inaccurate geometries. Also, it is sensitive to the foreground masking process, which is not suitable for relatively complex scenes as in the DTU dataset. In contrast, our method can gradually include the in- and out-distribution details while maintaining good 3D consistency and geometry on both synthetic and real-world complex datasets.

Quantitative Analysis. (i) *Out-distribution 3D-consistency*. We compare the 3D-consistency with our main out-distribution competitor, Zero123-NVS, quantitatively in Tab. II using 10 randomly picked videos from each dataset, where our method outperforms Zero123-NVS in Pixel-MSE score by almost twice. This can be credited to the primary consistency provided by the coarse-stage outputs together with the proposed global consistency approximation in fine-stage Diff3DE.

(ii) *In-distribution Image Quality*. We report the in-distribution metric comparisons in Tab. III. It is obvious that our in-distribution method OPP shows a good balance between

TABLE II

OUT-DISTRIBUTION COMPARISONS BETWEEN OUR FINE-STAGE METHOD DIFF3DE AND ZERO123-NVS ON SHAPENET AND DTU (§ VI-C). THE 3D-CONSISTENCY OF OUR DIFF3DE OUTPERFORMS ZERO123-NVS BY NEARLY TWO TIMES.

Methods	↓Pixel-MSE		
	ShapeNet Cars	ShapeNet Chairs	DTU
Zero123-NVS	939.67	642.67	4350.70
Diff3DE (Ours)	489.30	330.81	2352.69

TABLE III

IN-DISTRIBUTION COMPARISONS BETWEEN OUR COARSE-STAGE METHOD OPP AND PREVIOUS METHODS ON SHAPENET AND DTU (§ VI-C). WE ACHIEVE IMPROVEMENTS WITH BALANCED FIDELITY (PSNR, SSIM) AND SHARPNESS (LPIPS, FID). “†” MEANS PER-SCENE OPTIMIZATION/FINETUNING/AUTO-REGRESSION IS REQUIRED. “*” INDICATES FID IN 64×64 RESOLUTION.

Methods	Fidelity		Sharpness	
	↑ PSNR	↑ SSIM	↓ LPIPS	↓ FID*
<i>ShapeNet Chairs</i>				
π -GAN [14] †	-	-	-	15.47
Pix2NeRF [15]	18.14	0.84	-	14.31
SRN [46] †	22.89	0.89	0.104	-
CodeNeRF [20] †	22.39	0.87	0.166	-
FE-NVS [5]	23.21	0.92	0.077	-
PixelNeRF [1]	23.72	0.91	0.128	38.49
OPP (ours)	24.03	0.92	0.067	15.10
<i>ShapeNet Cars</i>				
EG3D-PTI [16] †	19.00	0.85	0.150	27.32
3DiM [13] †	21.01	0.57	-	8.99
SRN [46] †	22.25	0.89	0.129	-
CodeNeRF [20] †	22.73	0.89	0.128	-
FE-NVS [5]	22.83	0.91	0.099	-
PixelNeRF [1]	23.17	0.90	0.146	59.15
OPP (ours)	23.24	0.91	0.092	33.53
<i>DTU Dataset</i>				
DietNeRF [18] †	14.24	0.481	0.487	190.7
PixelNeRF [1]	15.55	0.537	0.535	-
OPP (ours)	16.51	0.659	0.399	146.56
<i>DTU Dataset (SinNeRF Split)</i>				
SinNeRF [17] †	11.18	0.424	0.571	283.86
OPP (ours)	17.27	0.730	0.354	146.30

sharpness and fidelity in general, outperforming the baseline methods even though many of them require per-scene optimization/finetuning/auto-regression. It is worth mentioning that though 3DiM [13] achieves the best FID, as a geometry-free method, it mainly focuses on the image quality of every single view while ignoring the 3D consistency.

D. Ablation Studies of Out-distribution Diff3DE

Effectiveness of Coarse-to-Fine. We verify the effectiveness of the coarse-to-fine strategy in Fig. 8. When directly adding Diff3DE on the blurry OG-NeRF outputs (“Ours w/o Coarse”), the results tend to lose rich details and even significantly wrong geometry, *e.g.*, the window of the green car. With the proposed coarse-to-fine method, the wrong geometry can be generally corrected and the details are gradually included, formulating vivid outputs.

Influence of Dense Keyframe Number N_k . The influence of dense keyframe number N_k is illustrated in Fig. 9. With small N_k , the input of ISA module tends to contain contents with large variance, therefore leads to blurry outputs. Since our Diff3DE takes neighbor keyframe sets as ISA inputs, we are able to increase the dense keyframe number N_k to relieve

TABLE IV

ABLATION OF DIFFERENT PIPELINES (§ VI-E). COMPARED WITH THE INDEPENDENTLY TRAINED OG-NeRF MODEL AND TWO BASIC TANDEM PIPELINES: TWO-STAGE TANDEM PIPELINE (TTP) AND ONE-STAGE TANDEM PIPELINE (OTP), OUR FINAL ONE-STAGE PARALLEL PIPELINE (OPP) ACHIEVES MORE BALANCED FIDELITY AND SHARPNESS.

Methods	Fidelity		Sharpness	
	↑ PSNR	↑ SSIM	↓ LPIPS	↓ FID
OG-NeRF	23.61	0.907	0.113	69.85
TTP	22.31	0.899	0.089	48.22
OTP	23.11	0.900	0.086	40.00
OPP (ours)	23.69	0.910	0.091	39.70

TABLE V

EFFECTIVENESS OF DPS (§ VI-E). COMPARED WITH THE INDEPENDENTLY TRAINED GENERATIVE AND OG-NeRF MODELS, THE JOINTLY OPTIMIZED ONES FROM DPS CAN ACHIEVE BETTER PERFORMANCE.

Methods	Fidelity		Sharpness	
	↑ PSNR	↑ SSIM	↓ LPIPS	↓ FID
Generative	23.61	0.907	0.113	69.85
Generative (DPS)	23.64	0.909	0.108	61.07
OG-NeRF	23.11	0.900	0.086	40.00
OG-NeRF (DPS)	23.16	0.902	0.085	41.35

the content variance of ISA inputs. Obviously, the blurry issue is relieved as N_k increases.

E. Ablation Studies of In-distribution OPP

Similar to [1], we perform in-distribution ablation studies with 10% random instances (fixed across all experiments) from the test set of the ShapeNet Cars category. We provide the efficiency analysis of OPP in the appendix.

Analysis of Pipelines. We report the performance of different pipelines in Table IV. Compared with the independently trained OG-NeRF, *i.e.*, PixelNeRF, the Two-stage Tandem Pipeline (TTP) improves the sharpness, *e.g.*, the FID score is decreased by 21.63, yet significantly suffers the loss of fidelity, *e.g.*, -1.30 in PSNR score. With the end-to-end training in the One-stage Tandem Pipeline (OTP), the fidelity is greatly improved, but there still exist large margins compared with the independently trained OG-NeRF, *e.g.*, 23.11 vs. 23.61 in PSNR score. In contrast, our final One-stage Parallel Pipeline (OPP) achieves both high fidelity and sharpness.

Effectiveness of DPS. We study the influence of jointly optimizing OG-NeRF and generative models with DPS in Table V. Compared with the independently trained OG-NeRF (*i.e.*, PixelNeRF) and generative models (*i.e.*, one-stage tandem pipeline), we find the DPS outputs give better performance, which proves that DPS is a lightweight yet effective structure that can facilitate the parallel optimization of OG-NeRF and generative models.

Effectiveness of CoRF & DPF. Table VI proves the effectiveness of fusing with the output confidence map of CoRF, where the fused one achieves both high fidelity and sharpness. We also illustrate the learned confidence map in Fig. 10. Obviously, CoRF successfully learns to give the blurry part with lower confidence, and the fused ones can achieve a good balance between $\hat{\mathbf{I}}_t^N$ and $\hat{\mathbf{I}}_t^G$.

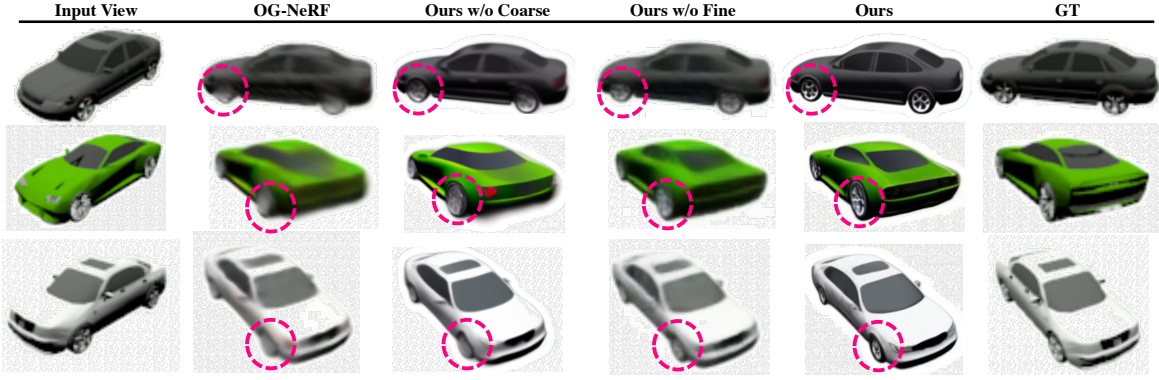


Fig. 8. **Effectiveness of coarse-to-fine (§ VI-D).** “w/o Coarse” indicates directly adding Diff3DE on the original blurry OG-NeRF, and “w/o Fine” means only using OPP. Obviously, our proposed coarse-to-fine scheme gradually includes the details and gives the best results.

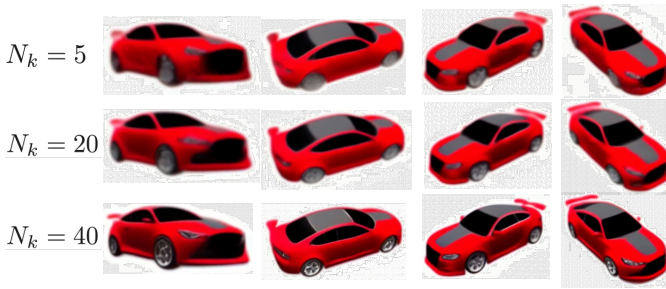


Fig. 9. **Influence of the dense key frame number N_k (§ VI-D).** Small N_k leads to more dispersed keyframes with large content variance, therefore generating blurry outputs. Increasing N_k relieves this issue via enabling more similar content in neighbor sets.

TABLE VI

EFFECTIVENESS OF FUSING VIA THE OUTPUT CONFIDENCE MAP OF CoRF (§ VI-E). THE FUSED ONE MINES THE BENEFITS OF BOTH THE OG-NeRF (PSNR, SSIM) AND GENERATIVE MODEL (LPIPS, FID).

Methods	Fidelity		Sharpness	
	↑ PSNR	↑ SSIM	↓ LPIPS	↓ FID
OG-NeRF (DPS)	23.64	0.909	0.108	61.07
Generative (DPS)	23.16	0.902	0.085	41.35
CoRF + DPF	23.69	0.910	0.091	39.70

VII. CONCLUSION

In this paper, targeting the blurry issue of existing OG-NeRF methods, we propose the GD²-NeRF, a generative detail compensation framework via GAN and pre-trained diffusion models for O-NVS task. It achieves vivid plausible outputs in an inference-time finetuning-free manner with decent 3D-consistency. We hope that our efforts will motivate more researchers in the future.

Limitations. First, same as most diffusion-based methods [8]–[10], the denoising process is inefficient. Second, as mentioned in [9], the decoder of the latent diffusion model introduces several high frequency flickering, which can be possibly addressed via an improved decoder or existing post-process deflickering method [9]. Third, our fine-stage method Diff3DE complements details mainly on the basis of the input contents, and cannot correct large geometry artifacts.

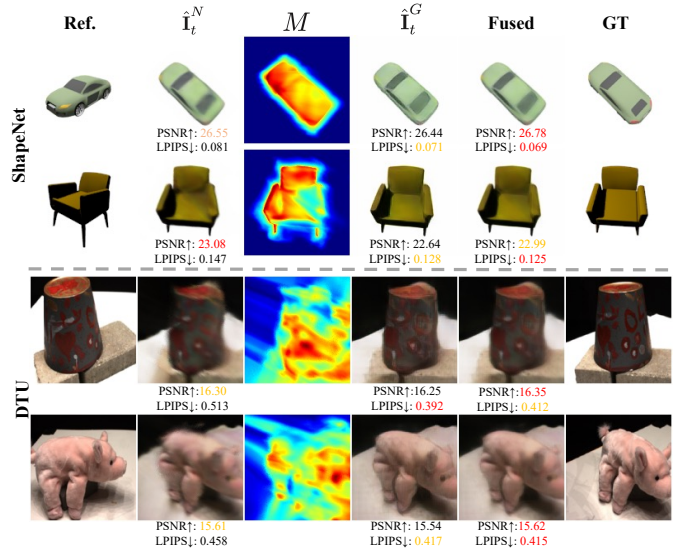


Fig. 10. **Visualization of the CoRF-predicted confidence map and the DPF results on ShapeNet and DTU (§ VI-E).** CoRF successfully learns to adaptively give the blurry part with low confidence, and the fused ones achieve good balance between two paradigms. Best viewed with zoom-in for details.

REFERENCES

- [1] A. Yu, V. Ye, M. Tancik, and A. Kanazawa, “pixelnerf: Neural radiance fields from one or few images,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, pp. 4578–4587.
- [2] R. Rombach, A. Blattmann, D. Lorenz, P. Esser, and B. Ommer, “High-resolution image synthesis with latent diffusion models,” in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2022, pp. 10 684–10 695.
- [3] L. Zhang, A. Rao, and M. Agrawala, “Adding conditional control to text-to-image diffusion models,” 2023.
- [4] K.-E. Lin, L. Yen-Chen, W.-S. Lai, T.-Y. Lin, Y.-C. Shih, and R. Ramamoorthi, “Vision transformer for nerf-based view synthesis from a single input image,” in *WACV*, 2023.
- [5] P. Guo, M. A. Bautista, A. Colburn, L. Yang, D. Ulbricht, J. M. Susskind, and Q. Shan, “Fast and explicit neural view synthesis,” in *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, 2022, pp. 3791–3800.
- [6] R. Liu, R. Wu, B. Van Hoorick, P. Tokmakov, S. Zakharov, and C. Von-drick, “Zero-1-to-3: Zero-shot one image to 3d object,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2023, pp. 9298–9309.

- [7] J. Tang, T. Wang, B. Zhang, T. Zhang, R. Yi, L. Ma, and D. Chen, "Make-it-3d: High-fidelity 3d creation from a single image with diffusion prior," *arXiv preprint arXiv:2303.14184*, 2023.
- [8] S. Yang, Y. Zhou, Z. Liu, and C. C. Loy, "Rerender a video: Zero-shot text-guided video-to-video translation," *arXiv preprint arXiv:2306.07954*, 2023.
- [9] M. Geyer, O. Bar-Tal, S. Bagon, and T. Dekel, "Tokenflow: Consistent diffusion features for consistent video editing," *arXiv preprint arXiv:2307.10373*, 2023.
- [10] J. Gu, A. Trevisan, K.-E. Lin, J. Susskind, C. Theobalt, L. Liu, and R. Ramamoorthi, "Nerfdiff: Single-image view synthesis with nerf-guided distillation from 3d-aware diffusion," in *International Conference on Machine Learning*, 2023.
- [11] V. Sitzmann, S. Rezkikov, B. Freeman, J. Tenenbaum, and F. Durand, "Light field networks: Neural scene representations with single-evaluation rendering," *Advances in Neural Information Processing Systems*, vol. 34, pp. 19 313–19 325, 2021.
- [12] M. S. M. Sajjadi, H. Meyer, E. Pot, U. Bergmann, K. Greff, N. Radwan, S. Vora, M. Lucic, D. Duckworth, A. Dosovitskiy, J. Uszkoreit, T. Funkhouser, and A. Tagliasacchi, "Scene Representation Transformer: Geometry-Free Novel View Synthesis Through Set-Latent Scene Representations," *CVPR*, 2022. [Online]. Available: <https://srt-paper.github.io/>
- [13] D. Watson, W. Chan, R. M. Brualla, J. Ho, A. Tagliasacchi, and M. Norouzi, "Novel view synthesis with diffusion models," in *The Eleventh International Conference on Learning Representations*, 2023. [Online]. Available: <https://openreview.net/forum?id=HtoA0oT30jC>
- [14] E. R. Chan, M. Monteiro, P. Kellnhofer, J. Wu, and G. Wetzstein, "pi-gan: Periodic implicit generative adversarial networks for 3d-aware image synthesis," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2021, pp. 5799–5809.
- [15] S. Cai, A. Obukhov, D. Dai, and L. Van Gool, "Pix2nerf: Unsupervised conditional p-gan for single image to neural radiance fields translation," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 3981–3990.
- [16] E. R. Chan, C. Z. Lin, M. A. Chan, K. Nagano, B. Pan, S. De Mello, O. Gallo, L. J. Guibas, J. Tremblay, S. Khamis *et al.*, "Efficient geometry-aware 3d generative adversarial networks," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 16 123–16 133.
- [17] D. Xu, Y. Jiang, P. Wang, Z. Fan, H. Shi, and Z. Wang, "Sinnerf: Training neural radiance fields on complex scenes from a single image," 2022.
- [18] A. Jain, M. Tancik, and P. Abbeel, "Putting nerf on a diet: Semantically consistent few-shot view synthesis," in *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, October 2021, pp. 5885–5894.
- [19] B. Mildenhall, P. P. Srinivasan, M. Tancik, J. T. Barron, R. Ramamoorthi, and R. Ng, "Nerf: Representing scenes as neural radiance fields for view synthesis," in *European conference on computer vision*. Springer, 2020, pp. 405–421.
- [20] W. Jang and L. Agapito, "Codenerf: Disentangled neural radiance fields for object categories," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 12 949–12 958.
- [21] K. Rematas, R. Martin-Brualla, and V. Ferrari, "Sharf: Shape-conditioned radiance fields from a single view," *arXiv preprint arXiv:2102.08860*, 2021.
- [22] A. Brock, J. Donahue, and K. Simonyan, "Large scale gan training for high fidelity natural image synthesis," *arXiv preprint arXiv:1809.11096*, 2018.
- [23] Y. Choi, M. Choi, M. Kim, J.-W. Ha, S. Kim, and J. Choo, "Stargan: Unified generative adversarial networks for multi-domain image-to-image translation," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 8789–8797.
- [24] T. Karras, S. Laine, and T. Aila, "A style-based generator architecture for generative adversarial networks," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2019, pp. 4401–4410.
- [25] K. Schwarz, Y. Liao, M. Niemeyer, and A. Geiger, "Graf: Generative radiance fields for 3d-aware image synthesis," *Advances in Neural Information Processing Systems*, vol. 33, pp. 20 154–20 166, 2020.
- [26] M. Niemeyer and A. Geiger, "Giraffe: Representing scenes as compositional generative neural feature fields," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, pp. 11 453–11 464.
- [27] M. Caron, H. Touvron, I. Misra, H. Jégou, J. Mairal, P. Bojanowski, and A. Joulin, "Emerging properties in self-supervised vision transformers," in *Proceedings of the IEEE/CVF international conference on computer vision*, 2021, pp. 9650–9660.
- [28] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell, P. Mishkin, J. Clark *et al.*, "Learning transferable visual models from natural language supervision," in *International conference on machine learning*. PMLR, 2021, pp. 8748–8763.
- [29] C. Deng, C. Jiang, C. R. Qi, X. Yan, Y. Zhou, L. Guibas, D. Anguelov *et al.*, "Nerdi: Single-view nerf synthesis with language-guided diffusion as general image priors," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 20 637–20 647.
- [30] D. Xu, Y. Jiang, P. Wang, Z. Fan, Y. Wang, and Z. Wang, "Neural-lift-360: Lifting an in-the-wild 2d photo to a 3d object with 360 views," *arXiv e-prints*, pp. arXiv–2211, 2022.
- [31] M. Deitke, D. Schwenk, J. Salvador, L. Weihs, O. Michel, E. VanderBilt, L. Schmidt, K. Ehsani, A. Kembhavi, and A. Farhadi, "Objaverse: A universe of annotated 3d objects," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 13 142–13 153.
- [32] H. Wang, X. Du, J. Li, R. A. Yeh, and G. Shakhnarovich, "Score jacobian chaining: Lifting pretrained 2d diffusion models for 3d generation," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 12 619–12 629.
- [33] J. Z. Wu, Y. Ge, X. Wang, S. W. Lei, Y. Gu, Y. Shi, W. Hsu, Y. Shan, X. Qie, and M. Z. Shou, "Tune-a-video: One-shot tuning of image diffusion models for text-to-video generation," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2023, pp. 7623–7633.
- [34] L. Khachatryan, A. Movsisyan, V. Tadevosyan, R. Henschel, Z. Wang, S. Navasardyan, and H. Shi, "Text2video-zero: Text-to-image diffusion models are zero-shot video generators," *arXiv preprint arXiv:2303.13439*, 2023.
- [35] D. Ceylan, C.-H. P. Huang, and N. J. Mitra, "Pix2video: Video editing using image diffusion," *arXiv preprint arXiv:2303.12688*, 2023.
- [36] C. Qi, X. Cun, Y. Zhang, C. Lei, X. Wang, Y. Shan, and Q. Chen, "Fatezero: Fusing attentions for zero-shot text-based video editing," *arXiv:2303.09535*, 2023.
- [37] N. Max, "Optical models for direct volume rendering," *IEEE Transactions on Visualization and Computer Graphics*, vol. 1, no. 2, pp. 99–108, 1995.
- [38] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio, "Generative adversarial nets," *Advances in neural information processing systems*, vol. 27, 2014.
- [39] L. Mescheder, A. Geiger, and S. Nowozin, "Which training methods for gans do actually converge?" in *International conference on machine learning*. PMLR, 2018, pp. 3481–3490.
- [40] S. Wu, C. Rupprecht, and A. Vedaldi, "Unsupervised learning of probably symmetric deformable 3d objects from images in the wild," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 1–10.
- [41] F.-A. Croitoru, V. Hondru, R. T. Ionescu, and M. Shah, "Diffusion models in vision: A survey," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2023.
- [42] P. Dhariwal and A. Nichol, "Diffusion models beat gans on image synthesis," *Advances in neural information processing systems*, vol. 34, pp. 8780–8794, 2021.
- [43] J. Song, C. Meng, and S. Ermon, "Denoising diffusion implicit models," *arXiv preprint arXiv:2010.02502*, 2020.
- [44] N. Tumanyan, M. Geyer, S. Bagon, and T. Dekel, "Plug-and-play diffusion features for text-driven image-to-image translation," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2023, pp. 1921–1930.
- [45] A. X. Chang, T. Funkhouser, L. Guibas, P. Hanrahan, Q. Huang, Z. Li, S. Savarese, M. Savva, S. Song, H. Su *et al.*, "Shapenet: An information-rich 3d model repository," *arXiv preprint arXiv:1512.03012*, 2015.
- [46] V. Sitzmann, M. Zollhöfer, and G. Wetzstein, "Scene representation networks: Continuous 3d-structure-aware neural scene representations," *Advances in Neural Information Processing Systems*, vol. 32, 2019.
- [47] Z. Teed and J. Deng, "Raft: Recurrent all-pairs field transforms for optical flow," in *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part II 16*. Springer, 2020, pp. 402–419.
- [48] Z. Wang, A. C. Bovik, H. R. Sheikh, and E. P. Simoncelli, "Image quality assessment: from error visibility to structural similarity," *IEEE transactions on image processing*, vol. 13, no. 4, pp. 600–612, 2004.
- [49] R. Zhang, P. Isola, A. A. Efros, E. Shechtman, and O. Wang, "The unreasonable effectiveness of deep features as a perceptual metric," in

Proceedings of the IEEE conference on computer vision and pattern recognition, 2018, pp. 586–595.

- [50] M. Heusel, H. Ramsauer, T. Unterthiner, B. Nessler, and S. Hochreiter, “Gans trained by a two time-scale update rule converge to a local nash equilibrium,” *Advances in neural information processing systems*, vol. 30, 2017.
- [51] O. Ronneberger, P. Fischer, and T. Brox, “U-net: Convolutional networks for biomedical image segmentation,” in *International Conference on Medical image computing and computer-assisted intervention*. Springer, 2015, pp. 234–241.
- [52] R. Jensen, A. Dahl, G. Vogiatzis, E. Tola, and H. Aanaes, “Large scale multi-view stereopsis evaluation,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2014, pp. 406–413.
- [53] D. P. Kingma and J. Ba, “Adam: A method for stochastic optimization,” *arXiv preprint arXiv:1412.6980*, 2014.
- [54] J. Li, D. Li, S. Savarese, and S. Hoi, “Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models,” *arXiv preprint arXiv:2301.12597*, 2023.

TABLE VII
INFLUENCE OF λ_{GAN} AND λ_{PER} FROM GENERATIVE LEARNING OBJECTIVES (EQ. 6).

Objectives	Weight	\uparrow PSNR	\uparrow SSIM	\downarrow LPIPS	\downarrow FID
λ_{GAN}	$1e^{-4}$	23.49	0.906	0.097	42.43
	$1e^{-3}$	23.69	0.910	0.091	39.70
	$1e^{-2}$	23.62	0.909	0.095	40.14
λ_{PER}	$1e^{-3}$	23.64	0.908	0.105	51.21
	$1e^{-2}$	23.69	0.910	0.091	39.70
	$1e^{-1}$	23.28	0.902	0.100	46.70

TABLE VIII
EFFICIENCY COMPARISONS WITH PIXELNERF. THE PARAMETER NUMBER OF OPP IS ONLY SLIGHTLY HIGHER THAN PIXELNERF, AND THE ADDITIONAL RENDERING OF $\hat{\mathbf{I}}^G$ AND M IS SIGNIFICANTLY FASTER THAN $\hat{\mathbf{I}}^N$. WE TEST FPS ON 1 V100 WITH A BATCH SIZE OF 1.

Methods	Param.	FPS ($\hat{\mathbf{I}}^N$)	FPS ($\hat{\mathbf{I}}^G$)	FPS (M)
PixelNeRF [1]	15.05M	0.54	-	-
OPP (ours)	15.46M	0.49	21.50	11.90

APPENDIX

A. Influence of λ_{GAN} and λ_{PER} in OPP

We study the influence of λ_{GAN} and λ_{PER} in Tab. VII. In line with Sec. VI-E, we report the performance of the 10% instances from the test set of the ShapeNet Cars category. We first fix λ_{PER} as $1e^{-2}$ and then vary λ_{GAN} as $\{1e^{-4}, 1e^{-3}, 1e^{-2}\}$. Obviously, $\lambda_{GAN} = 1e^{-3}$ gives the best performance. Then, we fix λ_{GAN} as $1e^{-3}$, and vary λ_{PER} as $\{1e^{-3}, 1e^{-2}, 1e^{-1}\}$, and find that $\lambda_{PER} = 1e^{-2}$ performs better. We notice that the weight used in this work is much smaller than the previous GAN methods, e.g., $\lambda_{GAN} = \lambda_{PER} = 5$ in [15]. We infer that this is due to the functional difference of the GAN model between us. Specifically, in our work, the intermediate feature map \mathbf{I}_m already contains the coarse information provided by the OG-NeRF model, and the GAN model merely needs to refine the coarse input with learned prior. Therefore, too large weight will lead to over-imagination. While in [15], the captured prior by the GAN model is the main component for final outputs, therefore they require a larger weight.

B. Efficiency Analysis of OPP

We report the efficiency comparisons with PixelNeRF in Table VIII. With only an additional lightweight up-sampling module and several modifications on fully connected layers, the parameter number of our OPP is only 0.41M higher than PixelNeRF. During the inference stage, we need to forward two times, i.e., one full-sized ($H \times W$) rendering for getting $\hat{\mathbf{I}}^N$ (same as PixelNeRF) and confidence map M , and one additional $H_m \times W_m$ rendering for $\hat{\mathbf{I}}^G$. However, the CoRF MLP and up-sampling module are rather lightweight (0.09M, 0.11M), and $\hat{\mathbf{I}}^G$ requires substantially fewer query points than $\hat{\mathbf{I}}^N$. Therefore, rendering $\hat{\mathbf{I}}^G$ and M are over 40 and 24 times faster than $\hat{\mathbf{I}}^N$ (21.50 vs. 0.49 FPS and 11.90 vs. 0.49 FPS). In a nutshell, the inference speed is hardly affected.

TABLE IX
COMPARISONS OF COMPUTATION COST WITH RECENT WORKS [4], [10]. NOTABLY, [10] REQUIRES ADDITIONAL PER-SCENE FINETUNING WITH UNKNOWN COST.

Methods	Param.	Traning Time
VisionNeRF [4]	122M	16*A100 5Days
NeRFDiff [10]	400M / 1B	8*A100 4Days + finetune
OPP (ours)	15M	8*V100 3Days

C. Comparisons of Computation Cost between OPP and Recent Works

Recent works [4], [10] are not listed as competitors in Tab. III since they employ much heavier architecture with more computation cost, as illustrated in Tab. IX. Note that, except for the larger parameter number and training cost, NeRFDiff [10] requires additional per-scene finetuning. In contrast, built on top of PixelNeRF, our OPP is quite lightweight, and requires no per-scene finetuning.