

Distractor-free Generalizable 3D Gaussian Splatting

Yanqi Bao*
Nanjing University
Jiangsu, Nanjing, China
yanqibao1997@gmail.com

Jing Liao†
City University of Hong Kong
Hong Kong, China
jingliao@cityu.edu.hk

Jing Huo†
Nanjing University
Jiangsu, Nanjing, China
huojing@nju.edu.cn

Yang Gao
Nanjing University
Jiangsu, Nanjing, China
gaoy@nju.edu.cn

Abstract

We present *DGGS*, a novel framework addressing the previously unexplored challenge of **Distractor-free Generalizable 3D Gaussian Splatting (3DGS)**. It accomplishes two key objectives: fortifying generalizable 3DGS against distractor-laden data during both training and inference phases, while successfully extending cross-scene adaptation capabilities to conventional distractor-free approaches. To achieve these objectives, *DGGS* introduces a scene-agnostic reference-based mask prediction and refinement methodology during training phase, coupled with a training view selection strategy, effectively improving distractor prediction accuracy and training stability. Moreover, to address distractor-induced voids and artifacts during inference stage, we propose a two-stage inference framework for better reference selection based on the predicted distractor masks, complemented by a distractor pruning module to eliminate residual distractor effects. Extensive generalization experiments demonstrate *DGGS*'s advantages under distractor-laden conditions. Additionally, experimental results show that our scene-agnostic mask inference achieves accuracy comparable to scene-specific trained methods. Homepage is <https://github.com/bbbby-99/DGGS>.

1. Introduction

The widespread availability of mobile devices presents unprecedented opportunities for 3D reconstruction, fostering demand for direct 3D synthesis capabilities from casually captured images or video sequences (referred to as references). Recent approaches introduce generalizable 3D

*This work was completed during a visit to City University of Hong Kong.

†Corresponding author

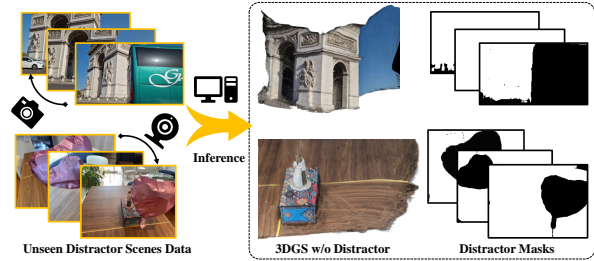


Figure 1. **Overview of Our Task.** DGGS enables direct 3DGS reconstruction from limited distractor-laden data while inferring distractor masks in a scene-agnostic manner.

representations to address this challenge, eliminating per-scene optimization requirements, with 3D Gaussian Splatting (3DGS) demonstrating particular promise due to its computational efficiency [3, 7, 17, 32]. In pursuit of scene-agnostic inference from references to 3DGS, these approaches simulate the complete pipeline from ‘references to 3DGS to novel query views’ within each training step, utilizing selected reference-query pairs while optimizing the process through query rendering losses. Following this paradigm, generalizable 3DGS requires both comprehensive training scenes and learned mechanisms for understanding geometric correlations between references to handle novel scenes. However, these essential components face fundamental challenges from distractors in unconstrained capture scenarios: (1) real-world scenes typically lack distractor-free training data, and (2) distractors disrupt 3D consistency among limited references.

To address these problems, a straightforward solution is to integrate distractor-free methods [5, 25] into generalizable 3DGS, enabling distractor mask prediction from residual loss. However, two fundamental limitations emerge in this approach: *First*, their loss-based masking strategies rely heavily on repeated optimization with sufficient single-

scene inputs and scene-specific hyperparameters. This approach faces significant challenges in scene-agnostic training settings, where residual loss uncertainty increases due to inter-iteration scene transitions and volatile reference-query pair selection mechanisms. This uncertainty undermines the core assumption that high-loss regions correspond to distractors, potentially misclassifying target objects as distractors and resulting in inadequate training supervision. *Second*, during reference-based inference paradigm, even when accurate masks are obtained, commonly occluded areas in references continue to affect spatial reconstruction and remain incomplete due to the limited number of references.

For the first challenge, we design a **Distractor-free Generalizable Training paradigm**, incorporating a **Reference-based Mask Prediction** and a **Mask Refinement** module to enhance training stability through precise distractor masking. Specifically, despite the absence of iteratively refined explicit scene representations when processing diverse scenes per iteration, our approach capitalizes on the stable reference renderings inherent in the ‘references to 3DGS’ paradigm. This facilitates the elimination of falsely identified distractor regions by utilizing the cross-view geometric consistency of static objects across references. After decoupling the filtered masks into distractor and disparity error components, we apply the Mask Refinement module, which incorporates pre-trained segmentation results to fill distractor regions and introduces reference-based auxiliary supervision in these areas for occlusion completion. Finally, to address the challenges posed by stochastic reference-query pairs, we introduce a proximity-driven **Training Views Selection** strategy based on translation and rotation matrices.

For the second challenge, despite accurate distractor region prediction, extensive occluded regions remain challenging to reconstruct with limited references. Therefore, we propose a two-stage **Distractor-free Generalizable Inference framework**. Specifically, in the first stage, we design a **Reference Scoring** mechanism based on predicted coarse 3DGS and distractor masks from pre-trained DGGS on initially sampled references. These scores guide the selection of minimally-distractor references for fine 3DGS reconstruction in the second stage. To further mitigate ghosting artifacts from residual distractors in this stage, we introduce a **Distractor Pruning** module that eliminates distractor-associated Gaussian primitives in 3D space.

Overall, we address a new task of *Distractor-free Generalizable 3DGS* as Fig. 1, and this is, to our knowledge, the first work to explore this problem. To tackle this challenge, we present **DGGS**, a framework designed to alleviate the adverse effects of distractors throughout the training and inference phases. Extensive experiments on distractor-rich datasets demonstrate that our approach successfully mitigates distractor-related challenges while improving generalization capability in conventional distractor-free mod-

els. Furthermore, our reference-based training paradigm achieves superior scene-agnostic mask prediction compared to existing scene-specific distractor-free methods.

2. Related Works

2.1. Generalizable 3D Reconstruction

Contemporary advances in generalizable 3D reconstruction seek to establish scene-agnostic representations, building upon early explorations in Neural Radiance Fields (NeRF) [20]. Benefiting from NeRF’s implicit representations, they treat the radiance field as an intermediary, effectively avoiding the need for explicit scene reconstruction and demonstrating the ability to infer novel viewpoints from only a few reference images, even in unseen scenes. The success of these works often relies on the sophisticated architectures such as Transformers [29, 30], Cost Volumes [4, 10], and Multi-Layer Perceptrons [1, 18]. However, the lack of explicit representations and rendering inefficiencies pose significant bottlenecks for them.

The advent of 3DGS [11], an explicit representation optimized for efficient rendering, has sparked renewed interest in the field. Existing works involve inferring Gaussian primitive attributes from references and rendering them from novel views. Analogous to NeRF-based approaches, 3DGS-related methods emphasize spatial comprehension from references, particularly focusing on depth estimation [3, 7, 15, 17, 32]. Subsequently, ReconX [16] and G3R [8] enhance reconstruction quality through the integration of additional video diffusion models and supplementary sensor inputs. The inherent reliance on high-quality references, however, makes generalizable reconstruction particularly susceptible to **distractors** - a persistent challenge in real-world applications. In this study, we examine Distractor-free Generalizable reconstruction, a topic that, to our knowledge, has not been addressed in existing literature.

2.2. Scene-specific Distractor-free Reconstruction

Scene-specific Distractor-free reconstruction focuses on accurately reconstructing one static scene while mitigating the impact of distractors [24] (or transient objects [25]). As a pioneering approach, NeRF-W [19] introduces additional embeddings to represent and eliminate transient objects under unstructured photo collections. Following a similar setting, subsequent extensive works focus on mitigating the impact of transient objects at the image level, which can generally be categorized into Knowledge-based methods, Heuristics-based methods and Hybrid methods [5, 22].

Knowledge-based methods predict transient objects using external knowledge sources, including pre-trained features or advanced segmentation models. Pre-trained features from ResNet [31, 33], Diffusion models [26], and DINO [13, 24] guide visibility map generation, effectively

weighting reconstruction loss. More recent works [5, 21, 22] directly employ state-of-the-art segmentation models like SAM [12] and Entity Segmentation [23] to establish clear distractors boundaries. While these approaches enhance earlier methods [6, 14, 19] with additional priors, they struggle to differentiate transient objects from complex static scene components, often serving mainly as auxiliary tools for mask prediction [5, 22].

Heuristics-based approaches employ handcrafted statistical metrics to detect distractors, predominantly emphasizing robustness and uncertainty analysis [9, 25, 28]. These methods exploit the observation that regions containing distractors typically manifest optimization inconsistencies. Therefore, they seek to predict outlier points based on loss residuals and mitigate their impact in loss functions. Regrettably, these approaches suffer from significant scene-specific data dependencies and frequently confound distractors with inherently challenging reconstruction regions, limiting their effectiveness in generalizable contexts.

Recently, there is growing advocacy for integrating the above-mentioned two methods [5, 22]. Entity-NeRF [22] integrates an existing Entity Segmentation model [23] and an extra entity classifier to determine distractors among each entity by analyzing the rank of loss residuals. Similarly, NeRF-HuGS [5] integrates pre-defined Colmap and Nerfacto [27] for capturing high and low-frequency features of static targets, while using SAM [12] to predict clear distractor boundaries. However, in our settings, acquiring additional entity classifiers or employing pre-defined knowledge such as Colmap and Nerfacto proves challenging, and loss residuals become unreliable compared to single-scene optimization due to the absence of iteratively refined explicit structures. Moreover, with limited references, despite obtaining accurate masks, Scene-specific Distractor-free methods struggle to handle commonly occluded regions and artifacts. Therefore, we present a novel **Distractor-free Generalizable** framework that jointly addresses distractor elimination in both training and inference phases.

3. Preliminaries

3.1. 3D Gaussian Splatting

3D Gaussian Splatting (3DGS) \mathcal{G} enables the representation of 3D scenes by splatting numerous anisotropic Gaussian primitives. Each Gaussian primitive is characterized by a set of attributes \mathbb{A} , including position \mathbf{p} , opacity α , covariance matrix Σ , and spherical harmonics coefficients for color $\hat{\mathbf{c}}$. To ensure positive semi-definiteness, the covariance matrix Σ is decomposed into a scaling matrix \mathbf{S} and a rotation matrix \mathbf{R} , such that $\Sigma = \mathbf{R}\mathbf{S}\mathbf{S}^\top\mathbf{R}^\top$. Conse-

quently, the color value after splatting on view \mathbf{P} is:

$$\hat{\mathbf{C}} = \mathcal{G}(\mathbf{P}) = \sum_{i \in M} \hat{\mathbf{c}}_i \alpha_i \prod_{j=1}^{i-1} (1 - \alpha_j), \quad (1)$$

where $\hat{\mathbf{c}}_i$ and α_i are derived from the covariance matrix Σ_i of the i -th projected 2D Gaussian, as well as the corresponding spherical harmonics coefficients and opacity values.

3.2. Generalizable 3DGS

Generalizable 3DGS presents a novel paradigm that directly infers Gaussian \mathcal{G}_{Ref} attributes from reference images, circumventing the computational overhead of scene-specific optimization. During the training phase, existing works optimize parameters θ (including En-Decoder, etc.) through randomly sampling paired references $\{\mathbf{I}_i\}_{i=1}^N$ and query image \mathbf{I}_T as inputs and ground truth under a sampled scene,

$$\mathcal{G}_{Ref} = \text{Decoder} \left(\mathcal{F} \left(\text{Encoder} \left(\{\mathbf{I}_i\}_{i=1}^N, \{\mathbf{P}_i\}_{i=1}^N \right) \right), \right. \quad (2)$$

$$\left. \arg \min_{\theta} \|\mathbf{I}_T - \mathcal{G}_{Ref}(\mathbf{P}_T)\|_2^2, \quad (3)$$

where $\{\mathbf{P}_i\}_{i=1}^N$ and \mathbf{P}_T are reference and query poses (views), and N denotes the number of references. Following Mvsplat [7], the \mathcal{F} denotes the process of feature warping, cost volumes $\{\mathbf{V}_i\}_{i=1}^N$ constructing, and depth estimation, etc.. After training across diverse training scenes, the model achieves scene-agnostic inference of 3DGS \mathcal{G}_{Ref} directly from given unseen scene references, as Eq. 2.

3.3. Robust Masks for 3D Reconstruction

Unlike conventional controlled environments, our research focuses on the challenges inherent in real-world, casually captured datasets. These in-the-wild scenarios contain not only static elements but also distractors [25] or transient objects [19], making it difficult to maintain 3D geometric consistency. Building upon prior research [25], we integrate a mask-based robust optimization process in our pipeline that can predict and filter out distractors. Eq. 2 is modified:

$$\arg \min_{\theta} \mathcal{M}_{Rob} \odot \|\mathbf{I}_T - \mathcal{G}_{Ref}(\mathbf{P}_T)\|_2^2. \quad (4)$$

Here, \mathcal{M}_{Rob} represents the predicted inlier/outlier mask on \mathbf{I}_T , where distractors are set to zero, which is typically associated with the residual loss and scene-specific thresholds.

$$\mathcal{M}_{Rob} = \mathbb{1} \left\{ \mathcal{C} \left(\mathbb{1} \left\{ \|\mathbf{I}_T - \mathcal{G}_{Ref}(\mathbf{P}_T)\|_2 > \rho_1 \right\} \right) > \rho_2 \right\}, \quad (5)$$

where \mathcal{C} represents the Convolution operator and ρ_1, ρ_2 are defined thresholds. Despite various proposed mask refinements in follow-up studies [5, 22], their heavy dependence on residual loss leads to extensive misclassification of static targets as distractor regions under the generalization setting, which will be addressed in subsequent sections.

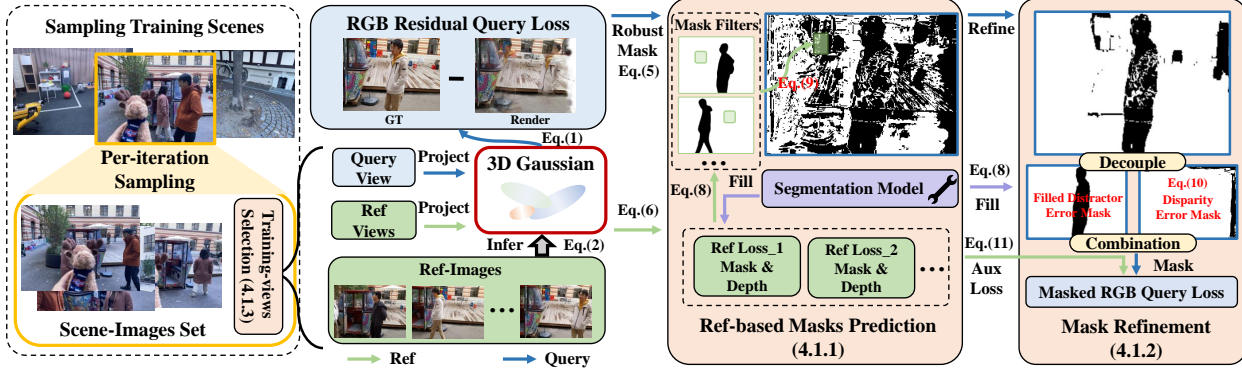


Figure 2. **Distractor-free Generalizable Training Paradigm.** DGGS first employs **Training Views Selection** for reference-query pair sampling and predict 3DGS attribute under sampled training scene. The **Reference-based Mask Prediction** module then generates filtered query robust masks, which are further refined through the **Mask Refinement** module to obtain final supervision for masked query loss.

4. Method

Given sufficient training reference-query pairs, the presence of distractors in either $\{\mathbf{I}_i\}_{i=1}^N$ or \mathbf{I}_T affects the 3D consistency relied upon by generalizable models, leading to training instability and artifacts during inference in the generalization paradigm. Therefore, we aim to design a **Distractor-free Generalizable Training** paradigm, Sec. 4.1 and a **Distractor-free Generalizable Inference** framework, Sec. 4.2 to mitigate these issues.

4.1. Distractor-free Generalizable Training

To mitigate the uncertainty in \mathcal{M}_{Rob} induced by scene transitions and stochastic reference-query pair sampling during each iteration, we propose a **Distractor-free Generalizable Training** paradigm, as illustrated in Fig. 2. Specifically, we introduce the **Reference-based Mask Prediction** (Sec. 4.1.1) and **Mask Refinement** (Sec. 4.1.2) modules to enhance per-iteration mask prediction accuracy scene-agnostically. Additionally, we design a **Training Views Selection** strategy (Sec. 4.1.3) to ensure stable views sampling.

4.1.1 Ref-based Masks Prediction

As discussed above, the excessive classification of target regions as distractor masks in Eq. 4 hinders geometric reconstruction of complex areas, as shown in Fig. 5. Therefore, we propose a scene-independent **Ref-based mask Prediction** method to maintain optimization focus across more non-distractor regions. Our inspiration stems from an intuitive observation: *3DGS inferred from references maintains stable rendering in non-distractor regions under reference views*. Therefore, we introduce a mask **Filter** that harnesses non-distractor regions from re-rendered references \mathcal{M}_{Ref_i} to identify and remove falsely labeled distractor regions in \mathcal{M}_{Rob} under query view based on the 3D consistency of static objects. Specifically, \mathcal{M}_{Ref_i} and \mathcal{M}_{Ref_i} -based query

view non-distractor regions \mathcal{M}_{Qry_i} as,

$$\{\mathcal{M}_{Ref_i} = \mathbb{1} \{ \mathcal{G}_{Ref}(\mathbf{P}_i) < \rho_{Ref} \} \}_{i=1}^N, \quad (6)$$

$$\{\mathcal{M}_{Qry_i} = \mathcal{W}(\mathcal{M}_{Ref_i}, \mathbf{D}_i, \mathbf{P}_i, \mathbf{P}_T, \mathbf{U})\}_{i=1}^N, \quad (7)$$

where \mathbf{U} represents the camera intrinsic matrix of image pairs, \mathbf{D}_i corresponds to the depth maps rendered from \mathbf{P}_i utilizing a modified rasterization library, \mathcal{W} defines the image warping operator that projects each \mathcal{M}_{Ref_i} from \mathbf{P}_i to \mathbf{P}_T using \mathbf{D}_i and \mathbf{U} , and ρ_{Ref} denotes the threshold parameter, experimentally determined as 0.001.

However, given the inherent inaccuracies in \mathbf{D}_i predictions and noise presence in \mathcal{M}_{Ref_i} , \mathcal{M}_{Qry_i} exhibits limited precision. Therefore, we incorporate a pre-trained segmentation model for mask filling and noise suppression, while designing a multi-reference masks fusion strategy to counteract warping-induced deviations. Following [5, 22], we incorporate a state-of-the-art Entity Segmentation Model [23] to improve \mathcal{M}_{Ref_i} into $\mathcal{M}_{Ref_i}^{En}$,

$$\mathcal{M}_{Ref_i}^{En} = \neg \left(\bigcup \mathcal{M}_i^{En_j} \right), \forall \frac{\mathcal{S}(\neg(\mathcal{M}_{Ref_i}) \cap \mathcal{M}_i^{En_j})}{\mathcal{S}(\mathcal{M}_i^{En_j})} \geq \rho_{En} \quad (8)$$

where \mathcal{S} represents the pixel-wise summation operator, \neg is the logical **NOT** operation, and $\mathcal{M}_i^{En_j}$ defines the j -th entity mask predicted from the segmentation model for \mathbf{I}_i . ρ_{En} is set to 0.8. After substituting \mathcal{M}_{Ref_i} with $\mathcal{M}_{Ref_i}^{En}$ in Eq. 7, we use an intersection operation to fuse multiple $\mathcal{M}_{Qry_i}^{En}$, then filter \mathcal{M}_{Rob} , obtaining Ref-based Mask \mathcal{M}_Q ,

$$\mathcal{M}_Q = \left\{ \bigcap \{ \mathcal{M}_{Qry_i}^{En} \}_{i=1}^N \right\} \bigcup \mathcal{M}_{Rob}. \quad (9)$$

The proposed approach ensures accurate distractor identification while filtering non-distractor regions, as shown in Fig. 5, which mitigates training instabilities induced by \mathbf{D}_i estimation errors. Excessively classified distractor regions undergo further refinement in the subsequent stage.

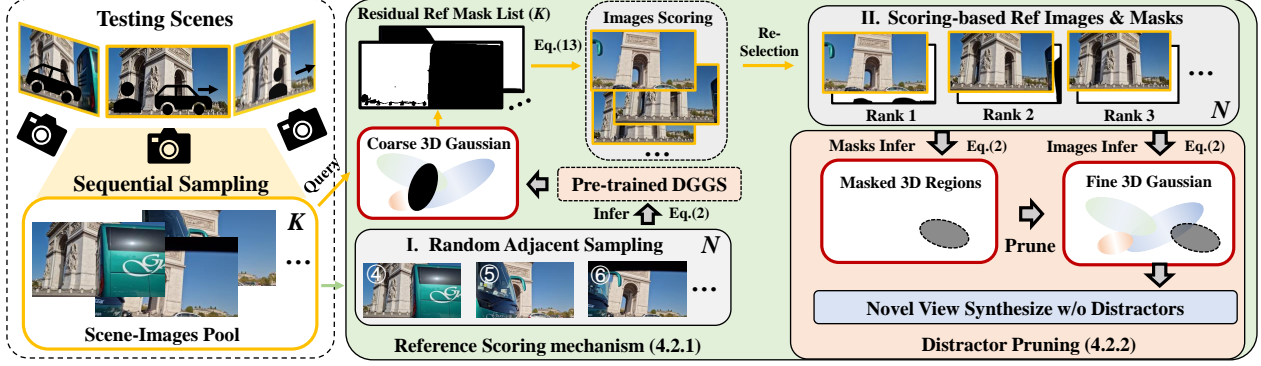


Figure 3. **Distractor-free Generalizable Inference Framework.** DGGs initially samples adjacent references from the scene-images pool and leverages trained DGGs for coarse 3DGS. Based on the **Reference Scoring mechanism**, quality scores and masks are computed for all pool images. These scores and masks subsequently guide reference selection and **Distractor Pruning** for fine 3DGS synthesis.

4.1.2 Mask Refinement

Given \mathcal{M}_Q , a straightforward approach is to utilize the segmentation results to remove excessive distractor regions and fill imprecise warping areas, as formulated in Eq. 8. In contrast to reference images, \mathcal{M}_Q contains both distractor regions and disparity-induced errors arising from reference-query view variations, with the latter being absent in references and primarily occurring at image margins. Thus, before introducing the segmentation model, regions decoupling is essential. The prediction of disparity-induced error mask follows a deterministic approach. Given N One Masks $\{\mathcal{M}_i^1\}_{i=1}^N$ corresponding to different poses \mathbf{P}_i , we warp them to \mathbf{P}_T as in Eq. 7. Then, the warped masks are merged using an union operation to ensure these regions are absent from all reference images.

$$\mathcal{M}_D = \bigcup \{ \mathcal{W}(\mathcal{M}_i^1, \mathbf{D}_i, \mathbf{P}_i, \mathbf{P}_T, \mathbf{U}) \}_{i=1}^N, \quad (10)$$

Finally, we decouple \mathcal{M}_D from \mathcal{M}_Q and recombine them after introducing the segmentation model [23] to refine the distractor error mask. The final refined mask, termed \mathcal{M} , substitutes \mathcal{M}_{Rob} in Eq. 4 to mitigate distractor effects during training. Note that all segmentation masks are pre-computed and cached to maintain training efficiency.

Additionally, in contrast to traditional distractor-free frameworks, reference images enable auxiliary supervision for masked regions under the query view, providing guidance for occluded area reconstruction. Thus, we re-warp \mathcal{M} to reference views and utilize $\mathcal{M}_{Ref_i}^{En}$ to determine the feasibility of occlusion completion. Specifically,

$$\mathcal{L}_A = \sum_{N=1}^n \mathcal{W}(\neg(\mathcal{M}), \mathbf{D}_T, \mathbf{P}_T, \mathbf{P}_i, \mathbf{U}) \odot \mathcal{M}_{Ref_i}^{En} \odot \|\mathbf{I}_i - \mathcal{G}_{Ref}(\mathbf{P}_i)\|_2^2. \quad (11)$$

The final form of Eq. 4 is modified to:

$$\arg \min_{\theta} \mathcal{M} \odot \|\mathbf{I}_T - \mathcal{G}_{Ref}(\mathbf{P}_T)\|_2^2 + \mathcal{L}_A. \quad (12)$$

4.1.3 Training Views Selection

As noted earlier, the selection strategy for references-query training pairs is critical. Intuitively, when query views are distant from references, suboptimal query rendering leads to significant residual losses in non-distractor regions and image margins. In contrast to prior approaches utilizing random sampling within a predefined range [3, 7], DGGs maintains minimal pose disparity between sampled reference and query views to enhance overall training stability.

In each training iteration, we randomly sample a scene and a corresponding query view, then choose references based on their translation and rotation matrix disparities relative to the query. Following the insights from work [2], we identify $2N$ views with minimal translation disparities, from which N views with the smallest rotation deviations are designated as reference views. Note that we must ensure the reference set do not include the query view.

4.2. Distractor-free Generalizable Inference

Despite improvements in training and mask prediction, DGGs's Inference faces two key limitations: (1) insufficient references compromise reliable reconstruction of commonly occluded regions, and (2) persistent distractors in references inevitably appear as artifacts in synthesized novel views. To address these challenges, we propose a two-stage **Distractor-free Generalizable Inference** framework, illustrated in Fig. 3. The first stage employs a **Reference Scoring mechanism** (Sec.4.2.1) to evaluate candidate references from the image pool, facilitating the selection of references with minimal distractor influence. The second stage implements a **Distractor Pruning** module (Sec.4.2.2) to suppress remaining distractor-induced artifacts.

4.2.1 Reference Scoring mechanism

Given a set of casually captured images or video frames containing distractors, a naive approach would be to select

Table 1. **Quantitative Experiments for distractor-free Generalizable 3DGS** under RobustNeRF Datasets. * denotes pre-trained models, + indicates baseline models augmented with existing mask prediction methods. More scenes are discussed in the supplementary materials.

Methods	Statue (RobustNeRF)			Android (RobustNeRF)			Mean (RobustNeRF)			Train Data
	PSNR↑	SSIM↑	LPIPS↓	PSNR↑	SSIM↑	LPIPS↓	PSNR↑	SSIM↑	LPIPS↓	
Pixelsplat [3]* (2024 CVPR)	18.65	0.673	0.254	17.98	0.557	0.364	20.10	0.704	0.279	Pre-Train on Re10K
Mvsplat [7]* (2024 ECCV)	18.88	0.670	0.225	18.24	0.586	0.301	20.03	0.722	0.255	
Pixelsplat [3] (2024 CVPR)	15.49	0.378	0.531	16.34	0.331	0.492	16.02	0.422	0.511	Re-Train on Distractor-Datasets
Mvsplat [7] (2024 ECCV)	15.05	0.412	0.391	16.17	0.509	0.381	15.45	0.515	0.426	
+RobustNeRF [25] (2023 CVPR)	16.17	0.463	0.382	16.46	0.470	0.411	17.11	0.534	0.400	
+On-the-go [24] (2024 CVPR)	14.73	0.366	0.522	15.05	0.440	0.472	15.44	0.476	0.526	
+NeRF-HuGS [5] (2024 CVPR)	18.21	0.694	0.266	18.33	0.640	0.299	19.18	0.700	0.283	
+SLS [26] (Arxiv 2024)	18.11	0.695	0.270	18.84	0.662	0.282	19.29	0.709	0.286	
DGGS-TR (w/o Inference Part)	<u>19.68</u>	<u>0.700</u>	0.238	<u>19.58</u>	0.653	<u>0.286</u>	<u>21.02</u>	<u>0.738</u>	<u>0.242</u>	
DGGS (Our)	20.78	0.710	<u>0.233</u>	20.93	0.711	0.236	21.74	0.758	0.237	

Table 2. **Components Ablation** for DGGS-TR and DGGS.

Methods	Mean (RobustNeRF)		
	PSNR↑	SSIM↑	LPIPS↓
Ablation on Our Training Paradigm			
Baseline (Mvsplat)	15.45	0.515	0.426
+Robust Masks	17.11	0.534	0.400
+ Ref-based Masks Prediction	20.35	0.701	0.283
+ Mask Refinement (DGGS-TR)	21.02	0.738	0.242
w/o Training Views Selection	16.33	0.551	0.441
w/o Entity Segmantation	20.79	0.733	0.248
w/o Aux Loss	20.64	0.725	0.253
Ablation on Our Inference Framework			
DGGS-TR	21.02	0.738	0.242
+ Reference Scoring mechanism	21.47	0.749	0.242
+ Distractor Pruning (DGGS)	21.74	0.758	0.237

reference images with minimal distractor influence for inference. Therefore, we propose a Reference Scoring mechanism based on pre-trained DGGS as the first stage of our Inference framework. Specifically, it first involves random sampling of N adjacent references from the scene-images pool $\{\mathbf{I}_i\}_P$ defined as K consecutive images in the test scene - for coarse 3DGS inference via DGGS. We then designate unselected views from the image pool as query views for masks \mathcal{M} prediction, while the chosen reference views represent distractor masks by $\mathcal{M}_{ref_i}^{en}$. All masks $\{\mathcal{M}_i\}_{i=1}^K$ from the image pool are collected as the basis for scoring,

$$\{\mathbf{I}_i\}_{i=1}^N = \{\mathbf{I}_i\}_P \mid i \in \max_N \left\{ \mathcal{S} \left(\{\mathcal{M}_i\}_{i=1}^K \right) \right\}. \quad (13)$$

In practice, besides distractor ratios, the poses of images in the pool are also crucial scoring factors. However, thanks to the disparity-induced error mask discussed in \mathcal{M} , we can directly utilize the count of positive pixels in the \mathcal{M} as the primary criterion. In the second stage, we employ top-ranked images as references to achieving fine 3DGS, effectively reweighting the originally equal reference without modifying N .

While this approach successfully handles distractor-heavy reference images, it comes at the cost of decreased

rendering efficiency. Optionally, we can mitigate this by halving image resolution in the first phase.

4.2.2 Distractor Pruning

Although ‘cleaner’ references are selected, obtaining N distractor-free images in the wild is virtually impossible. These residual distractors propagate via the Gaussian encoding-decoding process in Eq. 2, manifesting as phantom splats in rendered query views, as shown in Fig. 7. Therefore, we propose a Distractor Pruning protocol, which is readily implementable given the distractor masks corresponding to references, as described in Sec. 4.2.1. Instead of direct masking on the references, we selectively prune Gaussian primitives within the 3D spatial regions corresponding to masked areas by removing decoded attributes in distractor regions while preserving the remaining components. More details are provided in supplementary.

5. Experiments

This section presents both qualitative and quantitative experimental results for our DGGS under real-world generalization scenarios on distractor-laden datasets. Experimental results validate the reliability of our proposed training and inference paradigm. Additionally, multi-scene experiments demonstrate that DGGS enables traditional distractor-free methods to achieve generalization capability, which originally lack cross-scene training and inference abilities.

5.1. Experimental Details

5.1.1 Datasets

In accordance with existing generalization frameworks, DGGS is trained on extensive scenes with distractor presence and evaluated on novel, unseen distractor scenes to simulate real-world scenarios. Specifically, we utilize two widely-used mobile-captured datasets: On-the-go [24] and RobustNeRF [25], containing 12 and 5 distractor-laden

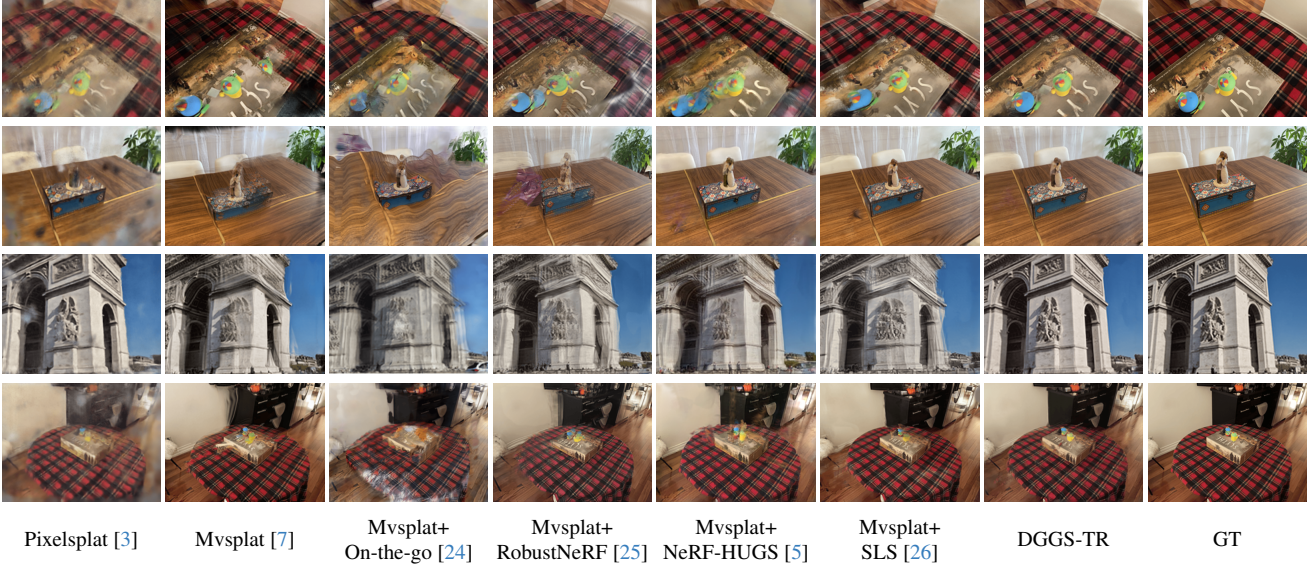


Figure 4. **Qualitative Comparison of Re-trained Existing Methods** across unseen scenes. More cases in the supplementary materials.

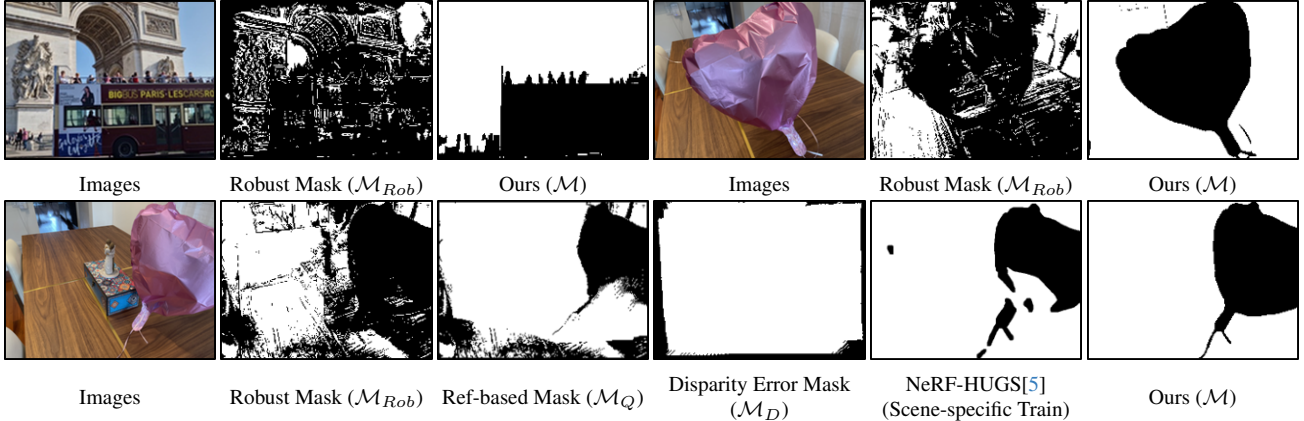


Figure 5. **The Ablation Visualize the Distractor Mask on Distractor Query and Comparison** with scene-specific trained method.

scenes respectively across outdoor and indoor environments. For fair comparison, we train all model on all On-the-go scenes except *Arcdetriomphe* and *Mountain*, which, along with the RobustNeRF dataset, serve as test scenes.

5.1.2 Training and Evaluation Setting

In all experiments, we set the number of references $N=4$ and the size of scene image pool $K=8$. During all re-training, query views are randomly selected and reference views are chosen following the **Training Views Selection** strategy, regardless of ‘clutter’ or ‘extra’ categorization. In the evaluation phase, we utilize all ‘extra’ images as query views for On-the-go scenes (*Arcdetriomphe* and *Mountain*), and for RobustNeRF scenes, query views are sampled from ‘clear’ images with a stride of eight. For evaluation metrics, we construct the scene-images pool using views closest to the query view, ensuring inclusion of both distractor-contaminated and distractor-free data to validate the effectiveness of Reference Scoring. Note that this setup is solely

for validation and evaluation purposes. In practical applications, the scene-images pool can be constructed using any adjacent views, independent of the query view and distractor presence. Finally, we compute scene-wide average PSNR, SSIM, and LPIPS metrics on the query render.

5.2. Comparative Experiments

5.2.1 Benchmark

Our Distractor-free Generalizable training and inference paradigms can be seamlessly integrated with existing generalizable 3DGS frameworks. We adopt Mvsplat [7] as our baseline model. Extensive comparisons are conducted against existing approaches **re-trained** under same settings on our distractor datasets, including: (1) original generalization methods [3, 7], and (2) Mvsplat [7] incorporating mask estimation from distractor-free approaches [5, 24–26]. We further evaluate **pre-trained** models (trained on clean datasets) on distractor-containing scenarios. Additional details are provided in the supplementary materials.

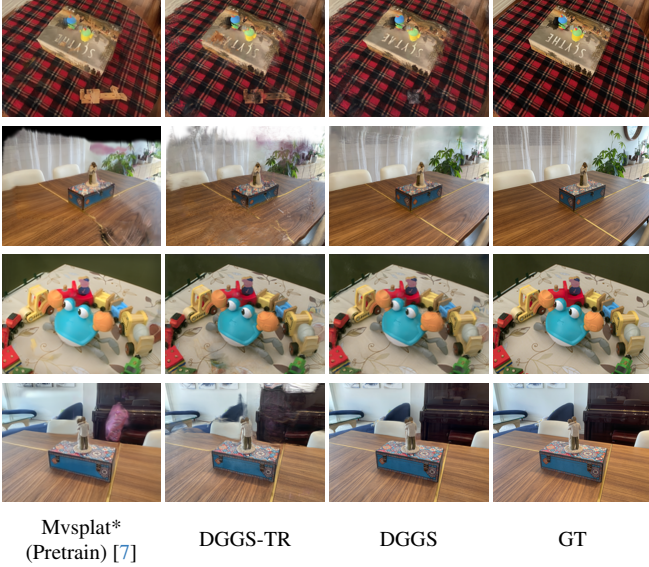


Figure 6. **Qualitative Comparison of Pre-trained Models** and our DGGs-TR as well as DGGs under unseen scenes.

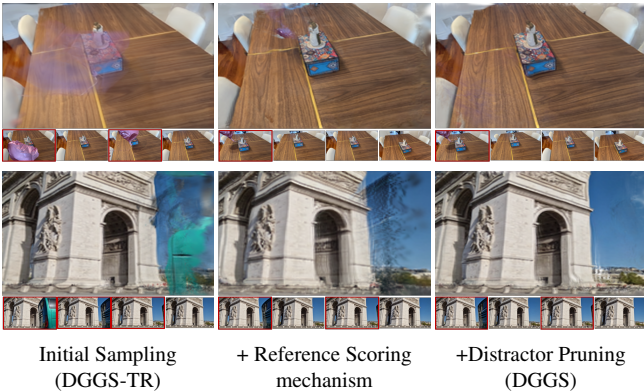


Figure 7. **Qualitative Ablation for Inference Strategy.**

5.2.2 Quantitative and Qualitative Experiments

Tab. 1, Fig. 4 and Fig. 6 quantitatively and qualitatively compares DGGs-TR (only TRaining) and DGGs with existing methods. The experimental results are analyzed from two aspects: re-training and pre-training models.

For Re-train Model: Evidence from Tab. 1 and Fig. 4 demonstrates that distractor data poses substantial challenges to our training paradigm. Although various single-scene distractor masking methods have been incorporated, they prove ineffective in generalizable multi-scene settings. As discussed above, overly aggressive distractor identification compromises reconstruction quality, particularly in regions containing fine details. Our DGGs addresses these challenges while enabling generalizability for scene-specific distractor-free methods.

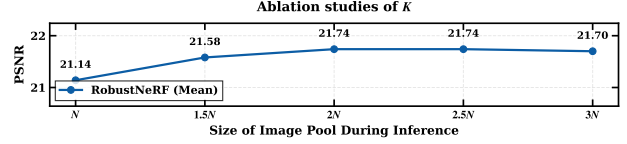


Figure 8. **Ablation of Image Pool Size K .** N is the refs number.

For Pre-train Model: Experimental results demonstrate that generalizable models, despite extensive dataset pre-training, suffer significant performance degradation in distractor-laden scenes in Tab. 1, primarily due to scene domain shifts and disrupted 3D consistency. DGGs-TR exhibits superior performance even with training limited to distractor scenes. Fig. 6 illustrates similar findings: although complete elimination of occlusion effects remains challenging, DGGs-TR effectively attenuates regions of 3D inconsistency. And then, DGGs achieves superior performance through references scoring and pruning strategies.

5.3. Ablation Studies

5.3.1 Ablation on Training Framework

The upper section of Tab. 2 and Fig. 5 present the impact of each component in the DGGs training paradigm. The Ref-based Masks Prediction combined with Mask Refinement mitigates the over-prediction of targets as distractors in the original Robust Masks, as shown in Fig. 5. Within the Mask Refinement module, the proposed Aux Loss demonstrates remarkable performance, with Entity Segmentation and Masks Decoupling providing substantial improvements. Also, Training Views Selection is essential during training. Our analysis reveals that DGGs achieves scene-agnostic mask inference capabilities, with direct inference results comparable to single-scene trained models (Fig. 5, second row). More cases are in supplementary.

5.3.2 Ablation on Inference Framework

The lower portion of Tab. 2 and Fig. 7 analyze the component effectiveness within the inference paradigm. Results indicate that although the Reference Scoring mechanism alleviates the impact of distractors in references by re-selection, certain artifacts remain unavoidable. Then, our Distractor Pruning strategy effectively mitigates these residual artifacts. We also analyze how the choice of K in Fig. 8, the scene image pool size, affects inference results. Generally, larger values of K yield better performance up to $2N$, beyond which performance plateaus, likely due to increased view disparity in the pool.

6. Conclusion

Distractor-free Generalizable 3D Gaussian Splatting presents a practical challenge, offering the potential to mitigate the limitations imposed by distractor scenes on

generalizable 3DGS while addressing the scene-specific training constraints of existing distractor-free methods. We propose novel training and inference paradigms that alleviate both training instability and inference artifacts from distractor data. Extensive experiments and discussions across diverse scenes validate our method’s effectiveness and demonstrate the potential of the refs-based paradigm in handling distractor data. We envision this work laying the foundation for future community discussions on Distractor-free Generalizable 3DGS and potentially extending to address 3D data challenges in broader applications.

7. Limitation

While our method enhances generalizability under distractor data during both training and inference, performance degradation under extensive mutual occlusions remains inevitable. Future work could potentially address this limitation by incorporating inpainting models based on predicted masks. Additionally, the increased inference time remains one of the challenges to be addressed in future work.

References

- [1] Yanqi Bao, Tianyu Ding, Jing Huo, Wenbin Li, Yuxin Li, and Yang Gao. Insertnerf: Instilling generalizability into nerf with hypernet modules. *arXiv preprint arXiv:2308.13897*, 2023. 2
- [2] Sibi Catley-Chandar, Richard Shaw, Gregory Slabaugh, and Eduardo Perez-Pellitero. Roguenerf: A robust geometry-consistent universal enhancer for nerf. *arXiv preprint arXiv:2403.11909*, 2024. 5
- [3] David Charatan, Sizhe Lester Li, Andrea Tagliasacchi, and Vincent Sitzmann. pixelsplat: 3d gaussian splats from image pairs for scalable generalizable 3d reconstruction. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 19457–19467, 2024. 1, 2, 5, 6, 7
- [4] Anpei Chen, Zexiang Xu, Fuqiang Zhao, Xiaoshuai Zhang, Fanbo Xiang, Jingyi Yu, and Hao Su. Mvsnr: Fast generalizable radiance field reconstruction from multi-view stereo. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 14124–14133, 2021. 2
- [5] Jiahao Chen, Yipeng Qin, Lingjie Liu, Jiangbo Lu, and Guanbin Li. Nerf-hugs: Improved neural radiance fields in non-static scenes using heuristics-guided segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 19436–19446, 2024. 1, 2, 3, 4, 6, 7
- [6] Xingyu Chen, Qi Zhang, Xiaoyu Li, Yue Chen, Ying Feng, Xuan Wang, and Jue Wang. Hallucinated neural radiance fields in the wild. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12943–12952, 2022. 3
- [7] Yuedong Chen, Haofei Xu, Chuanxia Zheng, Bohan Zhuang, Marc Pollefeys, Andreas Geiger, Tat-Jen Cham, and Jianfei Cai. Mvsplat: Efficient 3d gaussian splatting from sparse multi-view images. *arXiv preprint arXiv:2403.14627*, 2024. 1, 2, 3, 5, 6, 7, 8
- [8] Yun Chen, Jingkan Wang, Ze Yang, Sivabalan Manivasagam, and Raquel Urtasun. G3r: Gradient guided generalizable reconstruction. In *European Conference on Computer Vision*, pages 305–323. Springer, 2025. 2
- [9] Lily Goli, Cody Reading, Silvia Sellán, Alec Jacobson, and Andrea Tagliasacchi. Bayes’ rays: Uncertainty quantification for neural radiance fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 20061–20070, 2024. 3
- [10] Mohammad Mahdi Johari, Yann Lepoittevin, and François Fleuret. Geonerf: Generalizing nerf with geometry priors. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18365–18375, 2022. 2
- [11] 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. 2
- [12] Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C Berg, Wan-Yen Lo, et al. Segment anything. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 4015–4026, 2023. 3
- [13] Jonas Kulhanek, Songyou Peng, Zuzana Kukelova, Marc Pollefeys, and Torsten Sattler. Wildgaussians: 3d gaussian splatting in the wild. *arXiv preprint arXiv:2407.08447*, 2024. 2
- [14] Jaewon Lee, Injae Kim, Hwan Heo, and Hyunwoo J Kim. Semantic-aware occlusion filtering neural radiance fields in the wild. *arXiv preprint arXiv:2303.03966*, 2023. 3
- [15] Yiqing Liang, Numair Khan, Zhengqin Li, Thu Nguyen-Phuoc, Douglas Lanman, James Tompkin, and Lei Xiao. Gaufré: Gaussian deformation fields for real-time dynamic novel view synthesis. *arXiv preprint arXiv:2312.11458*, 2023. 2
- [16] Fangfu Liu, Wenqiang Sun, Hanyang Wang, Yikai Wang, Haowen Sun, Junliang Ye, Jun Zhang, and Yueqi Duan. Reconx: Reconstruct any scene from sparse views with video diffusion model. *arXiv preprint arXiv:2408.16767*, 2024. 2
- [17] Tianqi Liu, Guangcong Wang, Shoukang Hu, Liao Shen, Xinyi Ye, Yuhang Zang, Zhiguo Cao, Wei Li, and Ziwei Liu. Mvs gaussian: Fast generalizable gaussian splatting reconstruction from multi-view stereo. In *European Conference on Computer Vision*, pages 37–53. Springer, 2025. 1, 2
- [18] Yuan Liu, Sida Peng, Lingjie Liu, Qianqian Wang, Peng Wang, Christian Theobalt, Xiaowei Zhou, and Wenping Wang. Neural rays for occlusion-aware image-based rendering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7824–7833, 2022. 2
- [19] Ricardo Martin-Brualla, Noha Radwan, Mehdi SM Sajjadi, Jonathan T Barron, Alexey Dosovitskiy, and Daniel Duckworth. Nerf in the wild: Neural radiance fields for unconstrained photo collections. In *Proceedings of the IEEE/CVF*

- conference on computer vision and pattern recognition, pages 7210–7219, 2021. 2, 3
- [20] Ben Mildenhall, Pratul P Srinivasan, Matthew Tancik, Jonathan T Barron, Ravi Ramamoorthi, and Ren Ng. Nerf: Representing scenes as neural radiance fields for view synthesis. *Communications of the ACM*, 65(1):99–106, 2021. 2
- [21] Thang-Anh-Quan Nguyen, Luis Roldão, Nathan Piasco, Moussab Bennehar, and Dzmitry Tsishkou. Rodus: Robust decomposition of static and dynamic elements in urban scenes. *arXiv preprint arXiv:2403.09419*, 2024. 3
- [22] Takashi Otonari, Satoshi Ikehata, and Kiyoharu Aizawa. Entity-nerf: Detecting and removing moving entities in urban scenes. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 20892–20901, 2024. 2, 3, 4
- [23] Lu Qi, Jason Kuen, Weidong Guo, Tiancheng Shen, Jiuxiang Gu, Jiaya Jia, Zhe Lin, and Ming-Hsuan Yang. High-quality entity segmentation. *arXiv preprint arXiv:2211.05776*, 2022. 3, 4, 5
- [24] Weining Ren, Zihan Zhu, Boyang Sun, Jiaqi Chen, Marc Pollefeys, and Songyou Peng. Nerf on-the-go: Exploiting uncertainty for distractor-free nerfs in the wild. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8931–8940, 2024. 2, 6, 7
- [25] Sara Sabour, Suhani Vora, Daniel Duckworth, Ivan Krasin, David J Fleet, and Andrea Tagliasacchi. Robustnerf: Ignoring distractors with robust losses. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 20626–20636, 2023. 1, 2, 3, 6, 7
- [26] Sara Sabour, Lily Goli, George Kopanas, Mark Matthews, Dmitry Lagun, Leonidas Guibas, Alec Jacobson, David J Fleet, and Andrea Tagliasacchi. Spotlessplats: Ignoring distractors in 3d gaussian splatting. *arXiv preprint arXiv:2406.20055*, 2024. 2, 6, 7
- [27] Matthew Tancik, Ethan Weber, Evonne Ng, Ruilong Li, Brent Yi, Terrance Wang, Alexander Kristoffersen, Jake Austin, Kamyar Salahi, Abhik Ahuja, et al. Nerfstudio: A modular framework for neural radiance field development. In *ACM SIGGRAPH 2023 Conference Proceedings*, pages 1–12, 2023. 3
- [28] Paul Ungermann, Armin Ettenhofer, Matthias Nießner, and Barbara Roessle. Robust 3d gaussian splatting for novel view synthesis in presence of distractors. *arXiv preprint arXiv:2408.11697*, 2024. 3
- [29] Peihao Wang, Xuxi Chen, Tianlong Chen, Subhashini Venugopalan, Zhangyang Wang, et al. Is attention all that nerf needs? *arXiv preprint arXiv:2207.13298*, 2022. 2
- [30] Qianqian Wang, Zhicheng Wang, Kyle Genova, Pratul P Srinivasan, Howard Zhou, Jonathan T Barron, Ricardo Martin-Brualla, Noah Snavely, and Thomas Funkhouser. Ibrnet: Learning multi-view image-based rendering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4690–4699, 2021. 2
- [31] Jiacong Xu, Yiqun Mei, and Vishal M Patel. Wild-gs: Real-time novel view synthesis from unconstrained photo collections. *arXiv preprint arXiv:2406.10373*, 2024. 2
- [32] Chuanrui Zhang, Yingshuang Zou, Zhuoling Li, Minmin Yi, and Haoqian Wang. Transplat: Generalizable 3d gaussian splatting from sparse multi-view images with transformers. *arXiv preprint arXiv:2408.13770*, 2024. 1, 2
- [33] Dongbin Zhang, Chuming Wang, Weitao Wang, Peihao Li, Minghan Qin, and Haoqian Wang. Gaussian in the wild: 3d gaussian splatting for unconstrained image collections. *arXiv preprint arXiv:2403.15704*, 2024. 2