
Binocular-Guided 3D Gaussian Splatting with View Consistency for Sparse View Synthesis

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Learning to synthesize novel views from sparse inputs is a vital yet challenging
2 task in 3D computer vision. Previous methods explore 3D Gaussian Splatting
3 with neural priors (e.g. depth priors) as the additional supervisions, demonstrating
4 promising quality and efficiency compared to the NeRF based methods. However,
5 the neural priors from 2D pretrained models are often noisy and blurred, which
6 struggle in precisely guiding radiance field learning. In this paper, we propose
7 a novel prior-free method for sparse view Gaussian Splatting by exploring the
8 self supervisions inherent in the binocular stereo consistency between each pair
9 of binocular images constructed with disparity-guided image warping. We addi-
10 tionally introduce a Gaussian Opacity constraint which regularizes the Gaussians
11 locations and avoids Gaussian redundancy for improving the robustness and ef-
12 ficiency of sparse-view Gaussian Splatting. Extensive experiments on the LLFF,
13 DTU, and Blender datasets demonstrate that our method significantly outperforms
14 the state-of-the-art methods.

15 **1 Introduction**

16 3D reconstruction technologies [22, 17] have demonstrated significant advances in synthesizing
17 realistic novel views given a set of dense input views. To explore the challenging task of sparse view
18 synthesis for harsh real-world situations where only sparse inputs are available, some researches train
19 NeRF[22] with specially designed constraints [15, 36, 35, 38] and regularizations [24, 41, 39] on the
20 view scarcity. However, NeRF-based methods often suffer from slow training and inference speeds,
21 leading to high computational costs that restrict their practical applications.

22 3D Gaussian Splatting (3DGS) [17] has achieved notable advantages in rendering quality and
23 efficiency. However, 3DGS faces significant challenges with sparse views, where the unstructured 3D
24 Gaussians with limited constraints are prone to overfitting to the input views, resulting in incorrectly
25 learned scene geometry. Some recent works [25, 46, 20, 40] on sparse view synthesis based on 3DGS
26 employ the commonly-used depth priors from pre-trained models as additional constraints on the
27 Gaussians geometries. However, the neural priors are often noisy and blurred, which struggle in
28 precisely guiding radiance field learning.

29 In this paper, we aim at designing a prior-free approach which directly exploring the self supervisions
30 from the limited input views for improving the quality and efficiency of 3DGS for sparse views.
31 We justify that the key factors in achieving these goals are 1) learning correct scene geometry of
32 Gaussians which leads to consistent views synthesis, and 2) avoiding redundant Gaussians near the
33 surface for better efficiency and filtering noises.

34 For learning correct scene geometry of Gaussians, we propose a novel prior-free method for sparse
35 view Gaussian Splatting. We explore the self supervisions inherent in the binocular stereo consistency

36 to constrain on the rendered depth of 3DGS solely based on existing input views and synthesized
37 novel views. The key insight of our method lies in the observation that binocular image pairs
38 implicitly involve the property of view consistency, as demonstrated in the binocular stereo vision
39 works [11, 10, 44]. Specifically, we first translate the camera of one input view slightly to the left
40 or right to obtain a translational view, from which we render the image and depth from 3DGS. The
41 rendered image and the input one thus form a left-right view pair as in binocular stereo vision. We
42 then leverage the rendered depth and the known camera intrinsics to compute the disparity of the
43 view pair. We conduct the supervisions by warping the rendered image of the translational view to
44 the viewpoint of the input image using the disparity, and constrains on the consistency between the
45 warped and input images.

46 To further reduce redundant Gaussians near the scene surface and enhance the quality and efficiency
47 of novel view synthesis, we propose a decay schema for the opacity of Gaussians. Specifically, we
48 simply apply a decay coefficient to the opacity property of the Gaussians, penalizing the opacity
49 during training. To this end, Gaussians with lower opacity gradients are pruned due to the continuous
50 reduction of their opacity, while those with higher opacity gradients are retained. This results in
51 cleaner and more robust Gaussians by filtering out the redundant ones and guiding the remaining
52 Gaussians to be closer to the scene surface. This opacity decay strategy significantly reduces artifacts
53 in novel views and decreases the number of Gaussians, improving both the rendering quality and
54 optimization efficiency of 3DGS.

55 Additionally, to achieve better geometry initialization for improving 3DGS quality when conducting
56 optimization on sparse views, we use pre-trained keypoints matching network [34] to generate a
57 dense initialization point cloud. The dense point cloud describes the geometry of the scene more
58 accurately, preventing Gaussians from appearing far from the scene surface, especially in low-texture
59 areas where the distribution of Gaussians is subject to limited constraints.

60 In summary, our main contributions are as follows:

- 61 • We propose a novel prior-free method for sparse view Gaussian Splatting. We explore the
62 self supervisions inherent in the binocular stereo consistency to constrain on the rendered
63 depth of 3DGS solely based on existing input views and synthesized views.
- 64 • We propose an opacity decay strategy which significantly regularizes the Gaussians and
65 reduces Gaussian redundancy, leading to better rendering quality and optimization efficiency
66 for sparse view Gaussian Splatting.
- 67 • Extensive experiments on widely-used forward-facing and 360-degree scene datasets demon-
68 strate that our method achieves state-of-the-art results compared to existing sparse novel
69 view synthesis methods.

70 2 Related Works

71 2.1 Neural Radiance Field

72 Detailed and realistic 3D scene representation has always been the research goal in the field of
73 computer vision, and the emergence of neural radiance fields (NeRFs) [22] has brought fundamental
74 innovation to this domain. NeRF can reconstruct high-quality 3D scenes from sparse 2D images and
75 generate realistic images from arbitrary viewpoints by representing scenes as continuous volume
76 radiance functions.

77 However, NeRF requires a large number of views as input during the training stage and exhibits
78 limitations in both training and inference speed. Consequently, the following researches have
79 primarily focused on addressing these bottlenecks by improving computational efficiency [9, 23, 32,
80 17, 14, 4, 5, 8] or reducing the number of input views [24, 15, 36, 35, 6, 38, 30, 19, 7, 31, 29, 33, 39],
81 while continuously striving for enhanced rendering quality [1, 2, 3].

82 Notably, recent approach 3D Gaussian Splatting [17] has shown promising results in achieving
83 real-time rendering capabilities without compromising rendering quality.

84 3D Gaussian Splatting [17] is an emerging method for novel view synthesis, which reconstructs scenes
85 rapidly and with high quality by utilizing a set of 3D Gaussians to represent the radiance field in the
86 scene. This method performs excellently when dealing with real-world scenes, particularly excelling

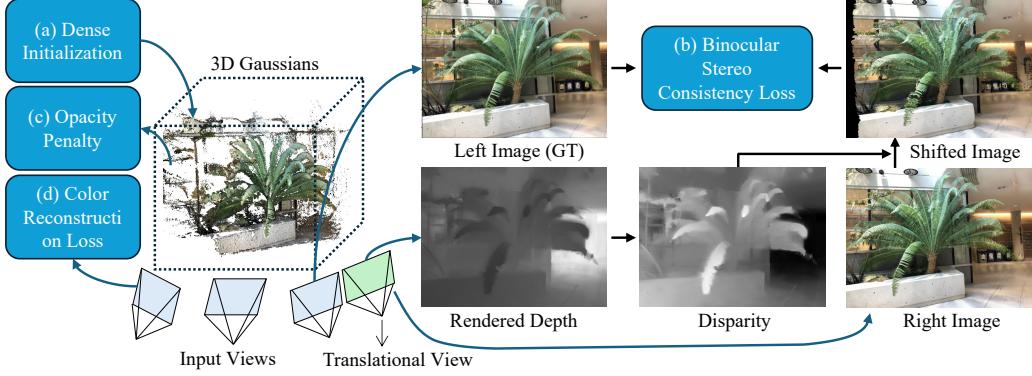


Figure 1: The overview of our method. (a) We leverage dense initialization for achieving Gaussian locations, and optimize the locations and Gaussian attributes with three constraints or strategies: (b) Binocular Stereo Consistency Loss. We construct a binocular view pair by translating an input view with camera positions, where we constrain on the view consistency of binocular view pairs in a self-supervised manner. (c) Opacity Penalty Strategy is designed to decay the Gaussian opacity during training for regularizing them. (d) The commonly-used Color Reconstruction Loss.

87 in handling high-frequency details. Moreover, 3D Gaussian Splatting demonstrates significant
88 advantages in inference speed and offers more intuitive editing and interpretability capabilities.

89 2.2 Novel View Synthesis from Sparse Views

90 Recent researches [24, 15, 36, 35, 6, 38, 30, 19, 7, 31, 29, 33, 39] in sparse novel view synthesis
91 have explored various approaches to generate novel views from sparse input images, with focus on
92 enhancing both rendering quality and efficiency. Methods such as NeRF[22] and 3DGS [17] have
93 been refined through various techniques, resulting in continuous improvement in the quality of novel
94 view synthesis.

95 Using depth priors obtained from pre-trained networks [42, 26, 27] to supervise neural radiance fields
96 is one of the effective techniques [31, 36]. Some methods [24, 41, 39, 7] introduce regularization terms
97 in NeRF to address the problem, such as frequency regularization [41] and ray entropy regularization
98 [18]. Additionally, there are also some methods [39, 38, 15] leverage pre-trained models such as
99 Diffusion model to enhance the rendering quality of novel views. However, most methods incur high
100 costs during training and inference. While some methods have improved inference efficiency through
101 generalizable models [5, 43] or by using voxel grids [33], they often sacrifice rendering quality.

102 Currently, with the advent of 3D Gaussian splatting, some methods use 3DGS for sparse view
103 synthesis, such as FSGS [46], SparseGS [40], DNGaussian [20], and CoherentGS [25]. These
104 methods utilize depth priors obtained from pre-trained models [26] as supervision. Depth constraints
105 enable the unstructured Gaussians to approximately distribute along the scene surface, thereby
106 enhancing the quality of novel view images. However, the depth priors from pre-trained models
107 often contain significant errors and cannot make the Gaussians distribute to optimal positions. In
108 contrast, our method employs view consistency constraints based on binocular vision, resulting in
109 more accurate depth information and thereby achieving a more optimal distribution of Gaussians.

110 3 Methods

111 The pipeline of our method is depicted in Figure 1. In this section, we first review the 3D information
112 representation method based on 3D Gaussians. Then, we explain how to construct stereo view pair
113 and utilize it to enforce view consistency in a self-supervised manner. Next, we introduce the opacity
114 penalty strategy and the dense initialization method for 3D Gaussians. Finally, we present the overall
115 loss function used for optimization.

116 **3.1 Review of 3D Gaussian Splatting**

117 Gaussian splatting [17] represents scene with a set of 3D Gaussians. Each 3D Gaussian is defined by
 118 a central location $\mu \in \mathbb{R}^3$, and a covariance matrix $\Sigma \in \mathbb{R}^{3 \times 3}$. Formally, it is defined as follows

$$G_i(x) = e^{-\frac{1}{2}(x-\mu_i)^T \Sigma^{-1}(x-\mu_i)}, \quad (1)$$

119 The covariance matrix Σ have physical meaning only when they are positive semi-definite. Therefore,
 120 it is decomposed into $\Sigma = RSS^T R^T$, $S \in \mathbb{R}^3$ is a diagonal scaling matrix with 3 parameters,
 121 $R \in \mathbb{R}^4$ is a rotation matrix analytically expressed with quaternions. In addition, for rendering the
 122 image, each Gaussian also stores an opacity value $\alpha \in \mathbb{R}$ and a color feature $f \in \mathbb{R}^k$.

123 The 3D Gaussian is projected into the 2D image space when rendering an image, the projected 2D
 124 Gaussian is sorted by its depth value, and then Alpha Blending is used to calculate the color of each
 125 pixel is calculated with α -blending.

$$c = \sum_{i=1}^n c_i \alpha'_i \prod_{j=1}^{i-1} (1 - \alpha'_j), \quad (2)$$

126 where c_i is the color computed from feature f and α' is the opacity of the 2D Gaussian, which is
 127 obtained by multiplying the covariance Σ' of the 2D Gaussian by the opacity α of the corresponding
 128 3D Gaussian. The 2D covariance matrix Σ' is calculated by $\Sigma' = JW\Sigma W^T J^T$, where J is the
 129 Jacobian of the affine approximation of the projective transformation. W is the view transformation
 130 matrix.

131 The optimization process of Gaussian splatting begins with creating a set of 3D Gaussians from a
 132 sparse SfM point clouds. Subsequently, the density of the Gaussian set is optimized and adaptively
 133 controlled. During optimization, a fast tile-based renderer is employed, allowing competitive training
 134 and inference times compared to the fast radiance field methods.

135 **3.2 Binocular Stereo Consistency Constraint**

136 The key factor in improving the rendering quality of 3DGS is to guide the Gaussians to be distributed
 137 close to the exact scene surfaces. To this end, we aim to design a constraint on the rendered depth of
 138 3DGS which directly represents the Gaussian geometry. Previous works commonly adopt the depth
 139 priors from pretrained models to guide the depth of 3DGS. However, the depth prior are often noisy
 140 and blurred, especially on complex scenes, which fails in providing accurate guidance.

141 In this paper, we propose a novel prior-free method for sparse view Gaussian Splatting by exploring
 142 the self supervisions inherent in the binocular stereo consistency. We constrain on the rendered depth of
 143 3DGS solely based on the existing input views and the synthesized novel views rendered
 144 from 3DGS itself. Specifically, we first obtain the translated image of the input view by translating
 145 the camera position corresponding to the input view, thus constructing a pair of binocular vision
 146 images. We then shift the novel view to the perspective of the input image using disparity. Finally,
 147 we constrain on the consistency between the input image and the warped image.

148 Given an image view I_l as input, we translate its corresponding camera position C_l to the right by
 149 a distance d_{cam} and obtain the translational camera position O_r . We then obtain a novel rendered
 150 image I_r by rendering 3DGS from O_r . The images I_l and I_r form a binocular stereo image pair,
 151 where I_l is the left image and I_r is the right image. Simultaneously, we can obtain the rendered depth
 152 D_r corresponding to the right image I_r by rendering 3DGS from O_r . According to the geometric
 153 theory of binocular stereo vision, the relationship between disparity d and depth D_r is given by
 154 $d = f \cdot d_{cam} / D_r$, where $d \in \mathbb{R}^{h \times w}$, each value d_{ij} in d indicates the horizontal shift required to
 155 align a pixel in the right image with its counterpart in the left image. f is the focal length of the
 156 camera. Then, we can use the disparity d to move each pixel of the right image I_r , obtaining the
 157 reconstructed left image $I_{shifted}$.

$$I_{shifted}[i, j] = I_r[i - d_i, j - d_j] \quad (3)$$

158 In practice, to make the reconstructed left image $I_{shifted}$ differentiable, we use a bilinear sampler to
 159 interpolate on the right image I_r to obtain each pixel of $I_{shifted}$. The sampling coordinates for each
 160 pixel are directly obtained from the disparity. Formally, we define the loss function as follows

$$L_{consis} = \frac{1}{N} \sum_{i,j} \left| I_{ij}^l - I_{ij}^{shifted} \right| \quad (4)$$

161 To achieve better optimization with the loss function, we do not use a fixed camera position shift in
162 practice. Instead, we randomly translate the input camera to the left or right by a distance within a
163 range $[-d_{max}, d_{max}]$.

164 Theoretically, it is also possible to implement the consistency constraint between input images by
165 only using the input images and the corresponding rendered depths. However, due to the large camera
166 pose variation among the sparse input views, the warped image can differ greatly from the target
167 image, making it difficult to achieve convergence during training.

168 3.3 Opacity Penalty Strategy

169 We further justify that relying solely on the depth constraints does not always lead to correct Gaussian
170 geometries that are closely aligned with the exact scene surfaces. The reason is that the rendered depth
171 varies with changes in the scale and opacity of the Gaussians, rather than being solely determined by
172 their positions. While 3DGS flexibly optimizes the scale and opacity during training, which leads to
173 deviations in the novel views. To address this issue, we design a simple strategy by applying a decay
174 coefficient λ to the opacity α of the Gaussians, penalizing the opacity during training.

$$\hat{\alpha} = \lambda\alpha, 0 < \lambda < 1 \quad (5)$$

175 We illustrate the opacity penalization strategy in Figure
176 2. The Gaussians close to the scene surface often con-
177 tain larger gradients due to constraints from various views,
178 while the Gaussians far from the surface are often with
179 much lower gradients due to weaker constraints. We aim
180 to improve and stabilize the optimization of 3DGS by
181 filtering out the far Gaussians which indicates incorrect
182 geometry while remain the ones close to the exact scene
183 surfaces. This is achieved by applying the opacity penalty
184 strategy. We justify that the strategy does not lead to a de-
185 crease in opacity for all Gaussians. As illustrated in Figure
186 2, the opacity of Gaussians with lower opacity gradients
187 gradually decreases until they are pruned. Conversely, the
188 increase in opacity for Gaussians near the scene surface
189 exceeds the penalty magnitude, ultimately achieving a bal-
190 ance between the opacity increase and the penalty, thereby
191 preserving Gaussians close to the surface.

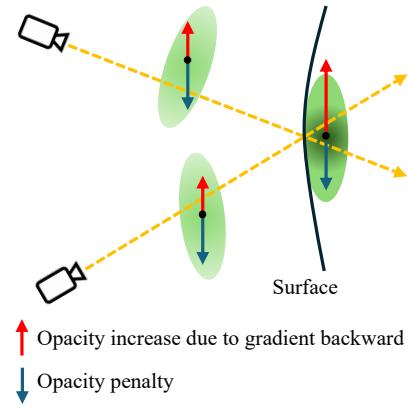


Figure 2: Illustration of the Gaussian opacity penalty strategy.

192 3.4 Initialization from dense point clouds

193 Previous 3DGS methods [17, 46, 40] usually utilize a sparse point cloud generated by Structure from
194 Motion (SfM) [28] to initialize 3D Gaussians. However, the point cloud produced by sparse views is
195 too sparse to adequately describe the scene to be reconstructed. Although the splitting strategy in
196 3DGS can replicate new Gaussians to cover the under-reconstructed area, they are subject to limited
197 geometric constraints and cannot adhere well to the scene surfaces, especially in low-texture areas
198 where the distribution of Gaussians may be arbitrary. We therefore seek a robust approach to achieve
199 better geometry initialization for improving 3DGS quality when optimizing from sparse views.

200 To achieve this, we use a pre-trained keypoints matching network [34] to generate a dense initialization
201 point cloud. Specifically, we arbitrarily select two images from the input images, input them into
202 the matching network, and obtain matching points. We then leverage the triangulation method [13],
203 along with the camera parameters corresponding to these images, to project the matching points into
204 3D space. This forms a dense point cloud, providing a more robust initialization for the Gaussians.

205 Compared with the sparse point cloud, the dense point cloud describes the geometry of the scene
206 more accurately, preventing Gaussians from appearing far from the scene surface and ultimately
207 leading to improved quality in novel view synthesis.

208 3.5 Training Loss

209 The final loss function consists of two parts: the proposed binocular stereo consistency loss L_{consis}
210 as introduced in Eq.4 and the commonly-used color reconstruction loss L_{color} of 3DGS [17]. We

Table 1: Quantitative comparison on LLFF. We evaluate the NeRF-based and the 3DGS-based methods, our method achieves the best results in all metrics under different input-view settings.

Methods	PSNR↑			SSIM↑			LPIPS↓		
	3-view	6-view	9-view	3-view	6-view	9-view	3-view	6-view	9-view
DietNeRF [15]	14.94	21.75	24.28	0.370	0.717	0.801	0.496	0.248	0.183
RegNeRF [24]	19.08	23.10	24.86	0.587	0.760	0.820	0.336	0.206	0.161
FreeNeRF [41]	19.63	23.73	25.13	0.612	0.779	0.827	0.308	0.195	0.160
SparseNeRF [36]	19.86	23.26	24.27	0.714	0.741	0.781	0.243	0.235	0.228
ReconFusion [38]	21.34	24.25	25.21	0.724	0.815	0.848	0.203	0.152	0.134
3DGS [17]	15.52	19.45	21.13	0.405	0.627	0.715	0.408	0.268	0.214
FSGS [46]	16.49	19.67	20.45	0.483	0.621	0.662	0.367	0.286	0.267
DNGaussian [20]	19.12	22.18	23.17	0.591	0.755	0.788	0.294	0.198	0.180
Ours	21.44	24.87	26.17	0.751	0.845	0.877	0.168	0.106	0.090

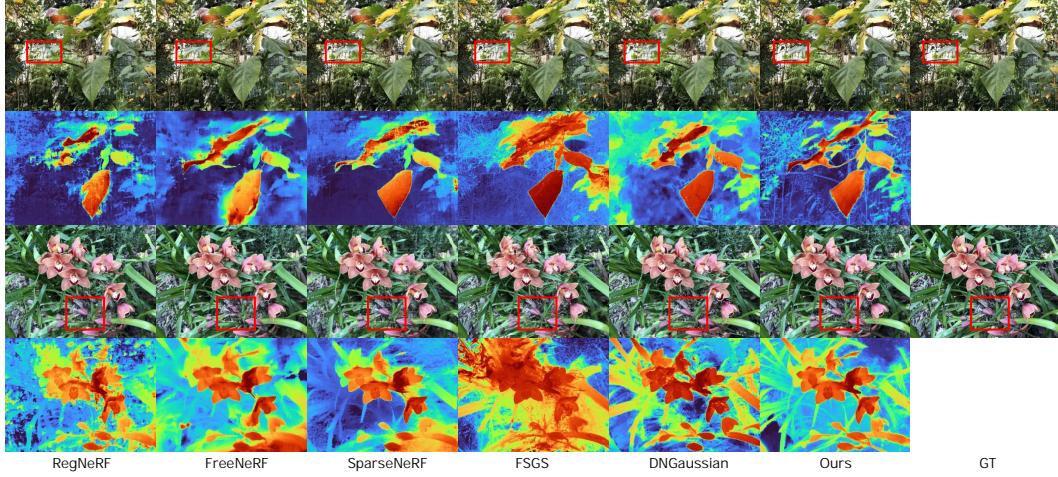


Figure 3: Visual comparison on LLFF dataset.

211 define the overall loss function as follows:

$$L = L_{consis} + L_{color}, \quad (6)$$

212 Where L_{color} is composed of an L_1 loss and a structural similarity loss L_{D-SSIM} , defined as

$$L_{color} = (1 - \beta)L_1 + \beta L_{D-SSIM} \quad (7)$$

213 4 Experiments

214 4.1 Datasets

215 We conducted experiments on three public datasets, including the LLFF dataset [21], the DTU dataset
216 [16], and the NeRF Blender Synthetic dataset (Blender) [22]. Following prior works [24, 41, 15], we
217 used 3, 6, and 9 views as training sets for the LLFF and DTU datasets, and 8 images for training on the
218 Blender dataset. The selection of test images remained consistent with previous works [15, 24, 41].
219 The downsampling rates for the LLFF, DTU, and Blender datasets were 8, 4, and 2, respectively.

220 4.2 Implementation details

221 Since The LLFF and DTU are datasets of forward-facing scenes, they cannot be constrained by views
222 from other directions during the optimization process. On the other hand, Blender is a dataset of
223 360-degree scenes, 8 images from different viewpoints are used as input during training, providing
224 stronger constraints. Therefore, for the LLFF and DTU datasets, we utilized a pre-trained matching
225 network PDC-Net+ [34] to obtain keypoints from input images, which were then used to generate

Table 2: Quantitative comparison on DTU. We evaluate the NeRF-based and the 3DGS-based methods, our method achieves the best results in most metrics under different input-view settings.

Methods	PSNR↑			SSIM↑			LPIPS↓		
	3-view	6-view	9-view	3-view	6-view	9-view	3-view	6-view	9-view
DietNeRF [15]	11.85	20.63	23.83	0.633	0.778	0.823	0.214	0.201	0.173
RegNeRF [24]	18.89	22.20	24.93	0.745	0.841	0.884	0.190	0.117	0.089
FreeNeRF [41]	19.52	23.25	25.38	0.787	0.844	0.888	0.173	0.131	0.102
SparseNeRF [36]	19.47	-	-	0.829	-	-	0.183	-	-
ReconFusion [38]	20.74	23.62	24.62	0.875	0.904	0.921	0.124	0.105	0.094
3DGS [17]	10.99	20.33	22.90	0.585	0.776	0.816	0.313	0.223	0.173
FSGS [46]	17.34	21.55	24.33	0.818	0.880	0.911	0.169	0.127	0.106
DNGaussian [20]	18.91	22.10	23.94	0.790	0.851	0.887	0.176	0.148	0.131
Ours	20.71	24.31	26.70	0.862	0.917	0.947	0.111	0.073	0.052

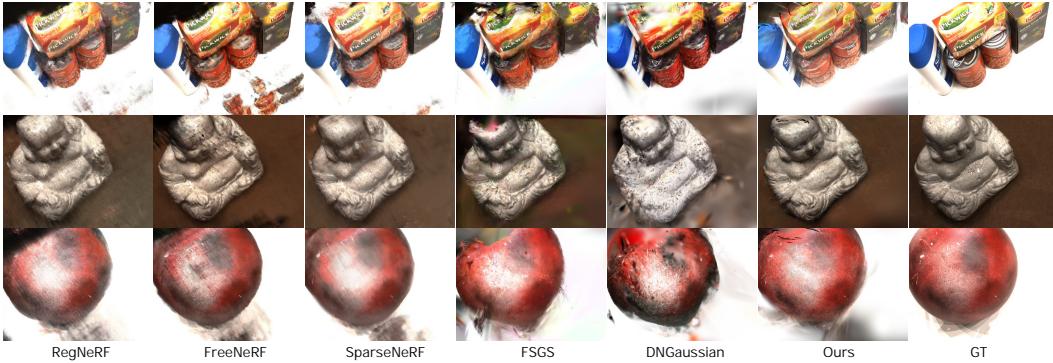


Figure 4: Visual comparison on DTU dataset.

226 dense initialization point clouds. For the Blender dataset, we adopted random initialization as in the
227 original 3DGS [17].

228 Moreover, based on the analysis above, we train the LLFF and DTU datasets for 30,000 iterations,
229 while the Blender dataset is trained for 7,000 iterations. Since the view consistency constraint based
230 on binocular stereo vision needs to be performed on the basis of the training views that can already
231 be rendered with high quality, we add the view consistency loss at 20,000 iterations for the LLFF
232 and DTU datasets, and at 4,000 iterations for the Blender dataset. The maximum distance d_{max} for
233 camera shift is set to 0.4, the opacity decay coefficient λ is set to 0.995, and the β in the loss function
234 is set to 0.2 as in the original 3DGS [17].

235 4.3 Baseline

236 We select some state-of-the-art NeRF-based and 3DGS-based sparse view synthesis methods to
237 contrast with our method. NeRF-based methods include DietNeRF [15], RegNeRF [24], FreeNeRF
238 [41], SparseNeRF [36], and ReconFusion [38]. 3DGS-based methods include DNGaussian [20] and
239 FSGS [46]. Additionally, we compared against the original 3DGS [17].

240 4.4 Comparisons

241 **LLFF** The Table 1 shows the quantitative results of the LLFF dataset with 3, 6, and 9 input views
242 respectively. Our method surpasses all baseline methods in terms of PSNR, SSIM [37], and LPIPS
243 [45] scores under different number of input views. When the number of input views increases to
244 9, it is almost adequate to provide sufficient color constraints. However, the evaluation scores of
245 DNGaussian [20] do not show a significant improvement, indicating that it does not perform well
246 with an increased number of input views. This is because errors in the depth prior negatively affect
247 the optimization process.

248 Figure 3 presents the visual comparison of
 249 novel view synthesis and depth rendering for the
 250 *leaves* and *orchids* scenes in the LLFF dataset.
 251 From the rendered results of novel views, sev-
 252 eral state-of-the-art methods are all perceptually
 253 acceptable. However, regarding the depth maps,
 254 RegNerf [24], FreeNerf [41], and FSGS [46] ex-
 255 hibit poor quality, indicating significant errors
 256 in their geometric estimation. SparseNeRF [36]
 257 and DNGaussian [20] indirectly utilize relative
 258 depth information from depth priors to supervise
 259 geometry estimation, hence achieving adequate
 260 rendering depth. Although FSGS [46] also ex-
 261 ploits depth priors, its attempt to directly employ
 262 depth prior information by using depth corre-
 263 lation loss. It does not yield accurate scene geometry, instead, inaccurate depth prior information
 264 exacerbates scene geometry. In contrast, our method utilizes self-supervised view consistency loss to
 265 obtain more precise geometry, thereby recovering more structural details in novel views.

266 **DTU** The Table 2 shows the quantitative re-
 267 sults of the DTU dataset under 3, 6, and 9 input
 268 views respectively. Figure 4 illustrates the visual
 269 comparison of novel view synthesis for three
 270 scenes in the DTU dataset. To distinctly discern
 271 the differences among various methods, we se-
 272 lect the synthesized images far from the training
 273 views for comparison in each scene. From the
 274 rendered results of novel views, it can be ob-
 275 served that NeRF-based methods produce blurry
 276 results, while FSGS [46] and DNGaussian [20]
 277 exhibit numerous artifacts due to insufficient
 278 constraints on Gaussians. Our method outper-
 279 forms previous state-of-the-art methods on most
 280 evaluation metrics and achieves better visual quality as well.

281 **Blender** We evaluate our method on 360-degree scenes using the Blender dataset. Table 3 presents
 282 the quantitative results under 8 input views. Our approach outperforms all baselines in the PSNR
 283 score, and achieves comparable SSIM and LPIPS. Figure 5 illustrates the visual comparison results of
 284 several 3DGS-based methods. It can be observed that FSGS [46] exhibits noticeable artifacts in the
 285 synthesized images due to insufficient Gaussian constraints, and the rendering results of DNGaussian
 286 [20] also has some distortion. In contrast, our method performs better in terms of detail preservation.

287 4.5 Ablation Study

288 To verify the effectiveness of the view consis-
 289 tency loss, opacity penalty strategy, and the
 290 dense initialization for Gaussians, we isolate
 291 each of these modules separately while keeping
 292 the other modules unchanged. We then evaluate
 293 the metrics and compared the visual results. As
 294 shown in Table 4, the performance decreases when any of our proposed modules is not used.

295 **Effectiveness of view consistency loss.** To verify the effectiveness of the view consistency loss,
 296 we compare the depth maps from novel views. Figure 6 presents the visual comparison of rendered
 297 depth maps. The comparison in the figure shows that it is evident that the view consistency loss
 298 significantly improves the estimation of scene geometry.

299 **Effectiveness of opacity penalty.** Figure 7 (b) and (c) present the comparison of novel view
 300 synthesis results and Gaussian point cloud visualizations for the *trex* scene in LLFF before and after

Table 3: Quantitative comparison on Blender for 8 input views. We evaluate the NeRF-based and the 3DGS-based methods, our method achieves comparable results to the state of the art methods.

Methods	PSNR↑	SSIM↑	LPIPS↓
DietNeRF [15]	22.50	0.823	0.124
RegNeRF [24]	23.86	0.852	0.105
FreeNeRF [41]	24.26	0.883	0.098
SparseNeRF [36]	22.41	0.861	0.199
3DGS [17]	22.23	0.858	0.114
FSGS [46]	22.76	0.829	0.157
DNGaussian [20]	24.31	0.886	0.088
Ours	24.71	0.872	0.101

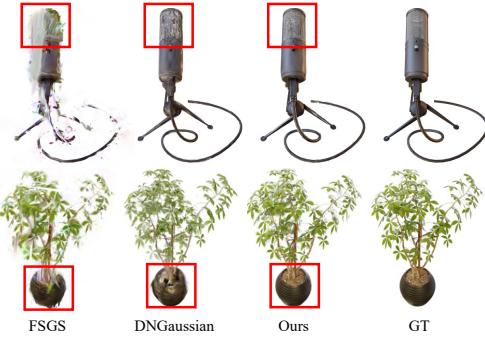


Figure 5: Visual comparison on Blender dataset.
 FSGS DNGaussian Ours GT

Table 4: Ablation studies on LLFF dataset with 3 input views.

Dense Init	L_{consis}	Opacity Penalty	PSNR ↑	SSIM ↑	LPIPS ↓
✓			16.67	0.494	0.384
✓	✓		19.30	0.651	0.238
	✓	✓	20.48	0.715	0.218
✓		✓	20.82	0.716	0.189
✓	✓	✓	21.13	0.738	0.180
✓	✓	✓	21.44	0.751	0.168

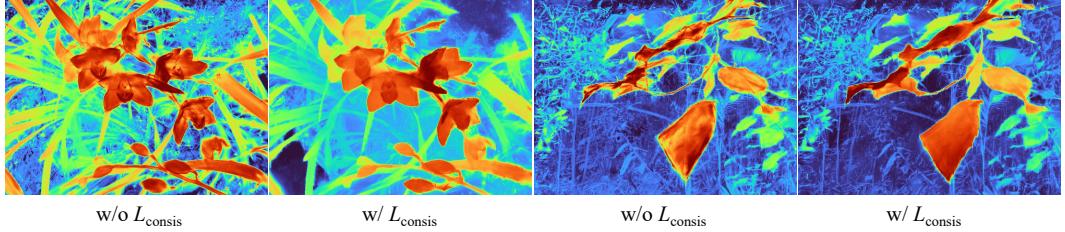


Figure 6: Visual comparison of depth maps before and after using view consistency loss.

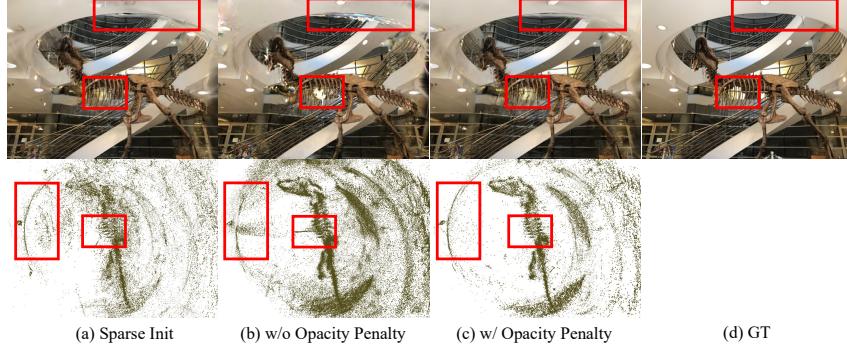


Figure 7: Visual comparison of novel view images and Gaussian point clouds.

301 applying the opacity penalty. It can be seen that without the opacity penalty, Gaussians appear far
 302 from the surfaces and there is noise near the surfaces. After applying the opacity penalty strategy,
 303 the Gaussians far from the scene surfaces are eliminated, and the noisy point clouds are significantly
 304 reduced, thereby further improving the rendering quality of novel views.

305 **Effectiveness of dense initialization.** Figure 7 (a) and (c) present the novel view images and the
 306 final Gaussian point clouds obtained using different initialization point cloud. It can be seen that
 307 without dense point cloud initialization, the spatial positions of the Gaussians become arbitrary in
 308 some low-texture regions or areas occluded in the training views due to insufficient constraints,
 309 resulting in reduced quality of the novel view images.

310 5 Conclusion

311 In this paper, we propose a novel method for the task of novel view synthesis from sparse views
 312 based on the 3DGS framework. We construct a self-supervised multi-view consistency constraint
 313 using the synthetic and input images, and introduce a Gaussian opacity penalty and a dense point
 314 cloud initialization strategy. These constraints ensure that the Gaussians are distributed as closely as
 315 possible to the scene surfaces and filter out those far from the surfaces. Our approach enables the
 316 unstructured Gaussians to accurately represent scene geometry even with sparse input views, resulting
 317 in high-quality novel rendering images. Extensive experiments on the LLFF, DTU, and Blender
 318 datasets demonstrate that our method outperforms existing state-of-the-art sparse view synthesis
 319 methods.

320 **Limitation.** Due to the utilization of view consistency constraints to estimate scene depth in our
 321 method, some areas with low texture in the scene cannot be accurately estimated the depth, thus
 322 failing to constrain the corresponding Gaussians, such as the white background areas in the DTU
 323 dataset. This results the white Gaussians to potentially occlude the object in the novel views. In
 324 contrast, DNGaussian [20] uses the depth prior estimated by the pre-trained network to constrain the
 325 Gaussians, preventing this scenario from happening. The visual comparisons are presented in the
 326 appendix.

327 **References**

- 328 [1] Jonathan T Barron, Ben Mildenhall, Matthew Tancik, Peter Hedman, Ricardo Martin-Brualla,
329 and Pratul P Srinivasan. Mip-nerf: A multiscale representation for anti-aliasing neural radiance
330 fields. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages
331 5855–5864, 2021.
- 332 [2] Jonathan T Barron, Ben Mildenhall, Dor Verbin, Pratul P Srinivasan, and Peter Hedman. Mip-
333 nerf 360: Unbounded anti-aliased neural radiance fields. In *Proceedings of the IEEE/CVF
334 Conference on Computer Vision and Pattern Recognition*, pages 5470–5479, 2022.
- 335 [3] Jonathan T Barron, Ben Mildenhall, Dor Verbin, Pratul P Srinivasan, and Peter Hedman.
336 Zip-nerf: Anti-aliased grid-based neural radiance fields. In *Proceedings of the IEEE/CVF
337 International Conference on Computer Vision*, pages 19697–19705, 2023.
- 338 [4] Anpei Chen, Zexiang Xu, Andreas Geiger, Jingyi Yu, and Hao Su. Tensorf: Tensorial radiance
339 fields. In *European Conference on Computer Vision*, pages 333–350. Springer, 2022.
- 340 [5] Anpei Chen, Zexiang Xu, Fuqiang Zhao, Xiaoshuai Zhang, Fanbo Xiang, Jingyi Yu, and
341 Hao Su. Mvsnerf: Fast generalizable radiance field reconstruction from multi-view stereo. In
342 *Proceedings of the IEEE/CVF international conference on computer vision*, pages 14124–14133,
343 2021.
- 344 [6] Di Chen, Yu Liu, Lianghua Huang, Bin Wang, and Pan Pan. Geoaug: Data augmentation for
345 few-shot nerf with geometry constraints. In *European Conference on Computer Vision*, pages
346 322–337. Springer, 2022.
- 347 [7] Kangle Deng, Andrew Liu, Jun-Yan Zhu, and Deva Ramanan. Depth-supervised nerf: Fewer
348 views and faster training for free. In *Proceedings of the IEEE/CVF Conference on Computer
349 Vision and Pattern Recognition*, pages 12882–12891, 2022.
- 350 [8] Sara Fridovich-Keil, Giacomo Meanti, Frederik Rahbæk Warburg, Benjamin Recht, and Angjoo
351 Kanazawa. K-planes: Explicit radiance fields in space, time, and appearance. In *Proceedings of
352 the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12479–12488,
353 2023.
- 354 [9] Sara Fridovich-Keil, Alex Yu, Matthew Tancik, Qinhong Chen, Benjamin Recht, and Angjoo
355 Kanazawa. Plenoxels: Radiance fields without neural networks. In *Proceedings of the IEEE/CVF
356 Conference on Computer Vision and Pattern Recognition*, pages 5501–5510, 2022.
- 357 [10] Clément Godard, Oisin Mac Aodha, and Gabriel J Brostow. Unsupervised monocular depth
358 estimation with left-right consistency. In *Proceedings of the IEEE conference on computer
359 vision and pattern recognition*, pages 270–279, 2017.
- 360 [11] Clément Godard, Oisin Mac Aodha, Michael Firman, and Gabriel J Brostow. Digging into
361 self-supervised monocular depth estimation. In *Proceedings of the IEEE/CVF international
362 conference on computer vision*, pages 3828–3838, 2019.
- 363 [12] Antoine Guédon and Vincent Lepetit. Sugar: Surface-aligned gaussian splatting for efficient 3d
364 mesh reconstruction and high-quality mesh rendering. *arXiv preprint arXiv:2311.12775*, 2023.
- 365 [13] Richard Hartley and Andrew Zisserman. *Multiple view geometry in computer vision*. Cambridge
366 university press, 2003.
- 367 [14] Wenbo Hu, Yuling Wang, Lin Ma, Bangbang Yang, Lin Gao, Xiao Liu, and Yuewen Ma. Tri-
368 miprf: Tri-mip representation for efficient anti-aliasing neural radiance fields. In *Proceedings of
369 the IEEE/CVF International Conference on Computer Vision*, pages 19774–19783, 2023.
- 370 [15] Ajay Jain, Matthew Tancik, and Pieter Abbeel. Putting nerf on a diet: Semantically consis-
371 tent few-shot view synthesis. In *Proceedings of the IEEE/CVF International Conference on
372 Computer Vision*, pages 5885–5894, 2021.
- 373 [16] Rasmus Jensen, Anders Dahl, George Vogiatzis, Engin Tola, and Henrik Aanæs. Large scale
374 multi-view stereopsis evaluation. In *Proceedings of the IEEE conference on computer vision
375 and pattern recognition*, pages 406–413, 2014.

- 376 [17] Bernhard Kerbl, Georgios Kopanas, Thomas Leimkühler, and George Drettakis. 3d gaussian
 377 splatting for real-time radiance field rendering. *ACM Transactions on Graphics*, 42(4):1–14,
 378 2023.
- 379 [18] Mijeong Kim, Seonguk Seo, and Bohyung Han. Infonerf: Ray entropy minimization for
 380 few-shot neural volume rendering. In *Proceedings of the IEEE/CVF Conference on Computer*
 381 *Vision and Pattern Recognition*, pages 12912–12921, 2022.
- 382 [19] Min-Seop Kwak, Jiuhan Song, and Seungryong Kim. Geconerf: Few-shot neural radiance fields
 383 via geometric consistency. *arXiv preprint arXiv:2301.10941*, 2023.
- 384 [20] Jiahe Li, Jiawei Zhang, Xiao Bai, Jin Zheng, Xin Ning, Jun Zhou, and Lin Gu. Dngaussian:
 385 Optimizing sparse-view 3d gaussian radiance fields with global-local depth normalization. *arXiv*
 386 *preprint arXiv:2403.06912*, 2024.
- 387 [21] Ben Mildenhall, Pratul P Srinivasan, Rodrigo Ortiz-Cayon, Nima Khademi Kalantari, Ravi
 388 Ramamoorthi, Ren Ng, and Abhishek Kar. Local light field fusion: Practical view synthesis
 389 with prescriptive sampling guidelines. *ACM Transactions on Graphics (TOG)*, 38(4):1–14,
 390 2019.
- 391 [22] Ben Mildenhall, Pratul P Srinivasan, Matthew Tancik, Jonathan T Barron, Ravi Ramamoorthi,
 392 and Ren Ng. Nerf: Representing scenes as neural radiance fields for view synthesis.
 393 *Communications of the ACM*, 65(1):99–106, 2021.
- 394 [23] Thomas Müller, Alex Evans, Christoph Schied, and Alexander Keller. Instant neural graphics
 395 primitives with a multiresolution hash encoding. *ACM transactions on graphics (TOG)*, 41(4):1–
 396 15, 2022.
- 397 [24] Michael Niemeyer, Jonathan T Barron, Ben Mildenhall, Mehdi SM Sajjadi, Andreas Geiger,
 398 and Noha Radwan. Regnerf: Regularizing neural radiance fields for view synthesis from
 399 sparse inputs. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern*
 400 *Recognition*, pages 5480–5490, 2022.
- 401 [25] Avinash Paliwal, Wei Ye, Jinhui Xiong, Dmytro Kotovenko, Rakesh Ranjan, Vikas Chandra,
 402 and Nima Khademi Kalantari. Coherentgs: Sparse novel view synthesis with coherent 3d
 403 gaussians. *arXiv preprint arXiv:2403.19495*, 2024.
- 404 [26] René Ranftl, Alexey Bochkovskiy, and Vladlen Koltun. Vision transformers for dense prediction.
 405 In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 12179–
 406 12188, 2021.
- 407 [27] René Ranftl, Katrin Lasinger, David Hafner, Konrad Schindler, and Vladlen Koltun. Towards
 408 robust monocular depth estimation: Mixing datasets for zero-shot cross-dataset transfer. *IEEE*
 409 *transactions on pattern analysis and machine intelligence*, 44(3):1623–1637, 2020.
- 410 [28] Johannes L Schonberger and Jan-Michael Frahm. Structure-from-motion revisited. In *Proceed-*
 411 *ings of the IEEE conference on computer vision and pattern recognition*, pages 4104–4113,
 412 2016.
- 413 [29] Seunghyeon Seo, Yeonjin Chang, and Nojun Kwak. Flipnerf: Flipped reflection rays for
 414 few-shot novel view synthesis. In *Proceedings of the IEEE/CVF International Conference on*
 415 *Computer Vision*, pages 22883–22893, 2023.
- 416 [30] Nagabhushan Somraj and Rajiv Soundararajan. Vip-nerf: Visibility prior for sparse input neural
 417 radiance fields. In *ACM SIGGRAPH 2023 Conference Proceedings*, pages 1–11, 2023.
- 418 [31] Jiuhan Song, Seonghoon Park, Honggyu An, Seokju Cho, Min-Seop Kwak, Sungjin Cho, and
 419 Seungryong Kim. Därf: Boosting radiance fields from sparse input views with monocular depth
 420 adaptation. *Advances in Neural Information Processing Systems*, 36, 2024.
- 421 [32] Cheng Sun, Min Sun, and Hwann-Tzong Chen. Direct voxel grid optimization: Super-fast
 422 convergence for radiance fields reconstruction. In *Proceedings of the IEEE/CVF Conference on*
 423 *Computer Vision and Pattern Recognition*, pages 5459–5469, 2022.

- 424 [33] Jiakai Sun, Zhanjie Zhang, Jiafu Chen, Guangyuan Li, Boyan Ji, Lei Zhao, Wei Xing, and
 425 Huaizhong Lin. Vgos: Voxel grid optimization for view synthesis from sparse inputs. *arXiv*
 426 preprint arXiv:2304.13386, 2023.
- 427 [34] Prune Truong, Martin Danelljan, Radu Timofte, and Luc Van Gool. Pdc-net+: Enhanced
 428 probabilistic dense correspondence network. *IEEE Transactions on Pattern Analysis and*
 429 *Machine Intelligence*, 2023.
- 430 [35] Prune Truong, Marie-Julie Rakotosaona, Fabian Manhardt, and Federico Tombari. Sparf: Neural
 431 radiance fields from sparse and noisy poses. In *Proceedings of the IEEE/CVF Conference on*
 432 *Computer Vision and Pattern Recognition*, pages 4190–4200, 2023.
- 433 [36] Guangcong Wang, Zhaoxi Chen, Chen Change Loy, and Ziwei Liu. Sparsenerf: Distilling
 434 depth ranking for few-shot novel view synthesis. In *Proceedings of the IEEE/CVF International*
 435 *Conference on Computer Vision*, pages 9065–9076, 2023.
- 436 [37] Zhou Wang, Alan C Bovik, Hamid R Sheikh, and Eero P Simoncelli. Image quality assessment:
 437 from error visibility to structural similarity. *IEEE transactions on image processing*, 13(4):600–
 438 612, 2004.
- 439 [38] Rundi Wu, Ben Mildenhall, Philipp Henzler, Keunhong Park, Ruiqi Gao, Daniel Watson,
 440 Pratul P Srinivasan, Dor Verbin, Jonathan T Barron, Ben Poole, et al. Reconfusion: 3d
 441 reconstruction with diffusion priors. *arXiv preprint arXiv:2312.02981*, 2023.
- 442 [39] Jamie Wynn and Daniyar Turmukhambetov. Diffusionerf: Regularizing neural radiance fields
 443 with denoising diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer*
 444 *Vision and Pattern Recognition*, pages 4180–4189, 2023.
- 445 [40] Haolin Xiong, Sairisheek Muttukuru, Rishi Upadhyay, Pradyumna Chari, and Achuta Kadambi.
 446 Sparsegs: Real-time 360 $\{\backslash\deg\}$ sparse view synthesis using gaussian splatting. *arXiv preprint*
 447 *arXiv:2312.00206*, 2023.
- 448 [41] Jiawei Yang, Marco Pavone, and Yue Wang. Freenerf: Improving few-shot neural rendering
 449 with free frequency regularization. In *Proceedings of the IEEE/CVF Conference on Computer*
 450 *Vision and Pattern Recognition*, pages 8254–8263, 2023.
- 451 [42] Lihe Yang, Bingyi Kang, Zilong Huang, Xiaogang Xu, Jiashi Feng, and Hengshuang
 452 Zhao. Depth anything: Unleashing the power of large-scale unlabeled data. *arXiv preprint*
 453 *arXiv:2401.10891*, 2024.
- 454 [43] Alex Yu, Vickie Ye, Matthew Tancik, and Angjoo Kanazawa. pixelnerf: Neural radiance fields
 455 from one or few images. In *Proceedings of the IEEE/CVF Conference on Computer Vision and*
 456 *Pattern Recognition*, pages 4578–4587, 2021.
- 457 [44] Ning Zhang, Francesco Nex, George Vosselman, and Norman Kerle. Lite-mono: A lightweight
 458 cnn and transformer architecture for self-supervised monocular depth estimation. In *Proceedings*
 459 *of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18537–18546,
 460 2023.
- 461 [45] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unrea-
 462 sonable effectiveness of deep features as a perceptual metric. In *Proceedings of the IEEE*
 463 *conference on computer vision and pattern recognition*, pages 586–595, 2018.
- 464 [46] Zehao Zhu, Zhiwen Fan, Yifan Jiang, and Zhangyang Wang. Fsgs: Real-time few-shot view
 465 synthesis using gaussian splatting. *arXiv preprint arXiv:2312.00451*, 2023.

466 **A More Visualizations**

467 **A.1 Additional Results on LLFF**

468 Figure 8 shows the visualization comparison of our method and DNGaussian [20] across all scenes in
 469 the LLFF [21] dataset. Each scene is trained using the same 3 input views. Although both methods
 470 achieve perceptually acceptable results, the novel view rendering results of our method are closer to
 471 the Ground Truth.

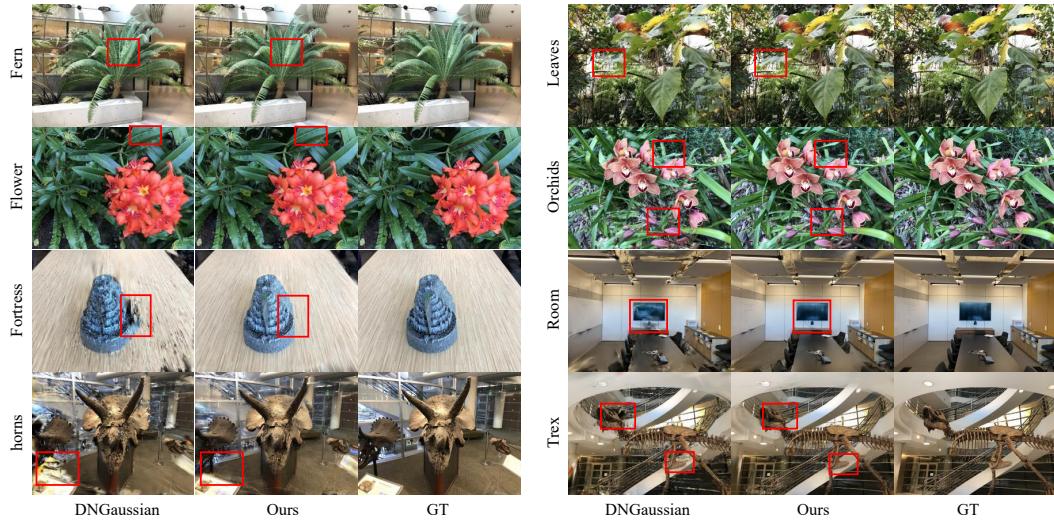


Figure 8: Visual comparison on LLFF dataset with 3 input views.

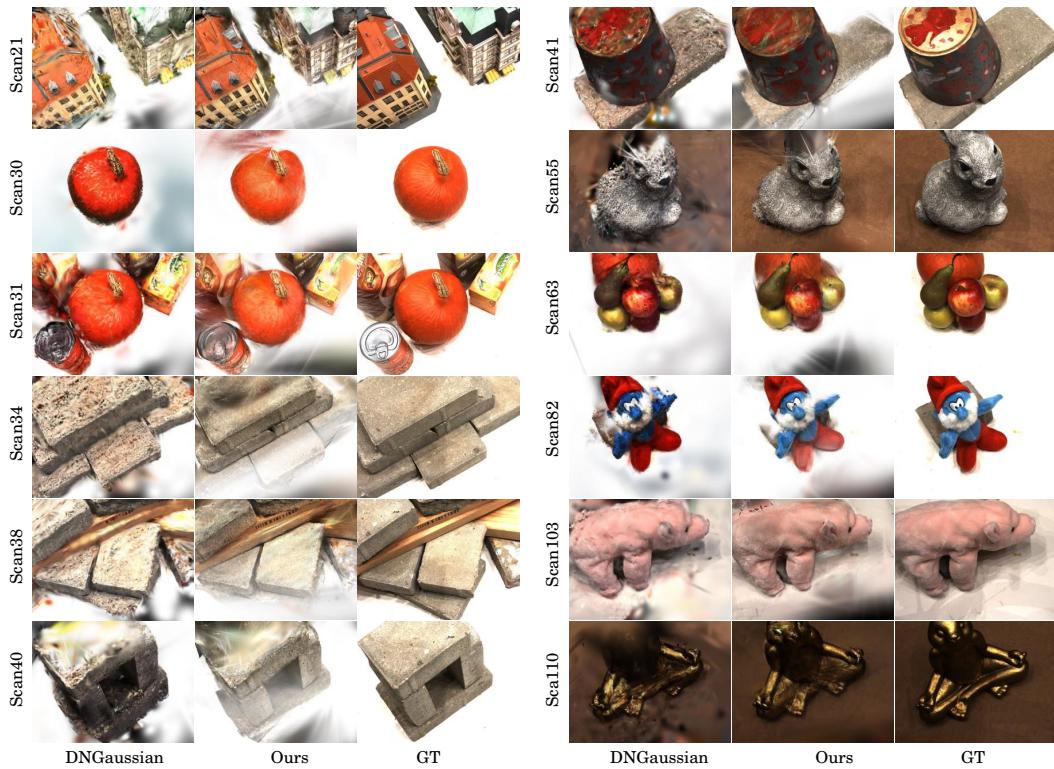


Figure 9: Visual comparison on DTU dataset with 3 input views.

Components	PSNR↑	SSIM↑	LPIPS↓
Opacity Entropy Reg.	15.75	0.480	0.361
Opacity Penalty	21.44	0.751	0.168

Table 5: Quantitative comparison of using different regularization components for opacity.



Figure 10: Visual comparison of using different regularization components for opacity.

472 A.2 Additional Results on DTU

473 Figure 9 shows additional visualization comparison of our method and DNGaussian [20] on the DTU
 474 [16] dataset. Each scene is trained using the same 3 input views. To clearly distinguish the differences
 475 between the two methods, we select a novel view that is far from the training views for comparison.
 476 As analyzed in the limitations section of the main text, our method cannot effectively constrain the
 477 white background in training views, resulting in object being occluded by white Gaussians in some
 478 scenes. Even so, the rendering quality of the novel views from our method is still better than the
 479 results from DNGaussian.

480 B More Experiments

481 B.1 Entropy regularization instead of Opacity Penalty

482 The Sugar [12] method applies entropy regularization to the opacity of Gaussians to bring Gaussians
 483 closer to the surface, similar to our proposed opacity penalty strategy. Therefore, we also evaluated
 484 the performance of opacity entropy regularization on the LLFF dataset. Following Sugar [12], the
 485 opacity threshold for Gaussians pruning is set to 0.5. Table 5 shows the metrics with opacity entropy
 486 regularization. It can be seen that our opacity penalty strategy has advantages. Figure 10 shows the
 487 novel view synthesis using different opacity regularization strategies. It can be seen that entropy
 488 regularization leads to noticeable overfitting with sparse view inputs, resulting in lower quality novel
 489 view rendering.

490 B.2 Training the DTU with Mask

491 To eliminate the undesirable effects of the solid color background on novel view rendering results
 492 in the DTU dataset, we perform additional experiments by using mask to filter out the background
 493 in the training views. The evaluation results are shown in Table 6. Figure 11 illustrates the visual
 494 comparison between our method and DNGaussian [20] when using masks in the training views. It
 495 can be seen that without the solid color background, our method shows a significant improvement,
 496 while the performance of DNGaussian decreases.

Methods	PSNR↑	SSIM↑	LPIPS↓
DNGaussian [20]	18.91	0.790	0.176
DNGaussian _{mask} [20]	14.94	0.699	0.237
Ours	20.71	0.862	0.111
Ours _{mask}	22.03	0.875	0.098

Table 6: Quantitative comparison of using background masks for input views on DTU dataset.

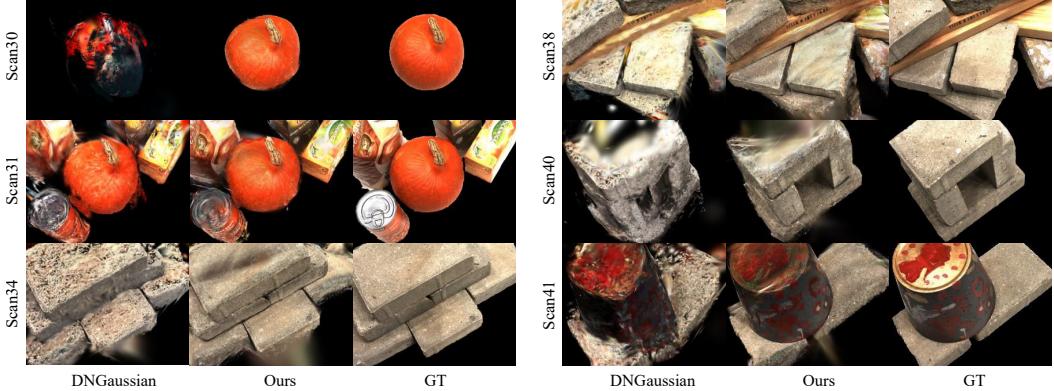


Figure 11: Visual comparison of using background masks for input views on DTU dataset.

497 C Implementations Details

498 Since our method uses an opacity penalty strategy, we do not use the opacity reset strategy from
499 the original 3DGS, and we also remove the scale threshold for Gaussians pruning from the original
500 3DGS.

501 We train our model on RTX 3090 GPU, running 30,000 iterations per scene on the LLFF and DTU
502 datasets, which takes approximately 9 minutes each. For the Blender dataset, we train for 7,000
503 iterations per scene, which takes approximately 3 minutes each. The storage space required for the
504 Gaussian point clouds is about one-third of that required by the original 3DGS.

505 D Dataset Details

506 **LLFF.** the LLFF [21] is a forward-facing dataset with 8 scenes. Following previous works [24, 41],
507 we take every 8th image as the novel views for evaluation. The input views are evenly sampled from
508 the remaining views. During training and evaluation, the images are downsampled by a factor of 8,
509 resulting in a resolution of 378×504 for each image.

510 **DTU.** The DTU [16] dataset consists of 124 scenes. We follow previous works [24, 41] to train and
511 evaluate our method on 15 test scenes. The test scene IDs are: 8, 21, 30, 31, 34, 38, 40, 41, 45, 55,
512 63, 82, 103, 110, and 114. In each scene, we use images with the following IDs as input views: 25,
513 22, 28, 40, 44, 48, 0, 8, 13. The first 3 and 6 image IDs correspond to the input views in 3-view and
514 6-view settings respectively. We use 25 images as novel views for evaluation, with the following IDs:
515 1, 2, 9, 10, 11, 12, 14, 15, 23, 24, 26, 27, 29, 30, 31, 32, 33, 34, 35, 41, 42, 43, 45, 46, 47. During
516 training and evaluation, the images are downsampled by a factor of 4, resulting in a resolution of 300
517 × 400 for each image.

518 **Blender.** The Blender [22] dataset consists of 8 synthetic scenes. We follow previous works [24, 41]
519 and use 8 images as input views for each scene, with the following IDs: 26, 86, 2, 55, 75, 93, 16, 73,
520 8. The 25 test views are sampled evenly from the testing images for evaluation. During training and
521 evaluation, the images are downsampled by a factor of 2, resulting in a resolution of 400 × 400 for
522 each image.

523 **NeurIPS Paper Checklist**

524 **1. Claims**

525 Question: Do the main claims made in the abstract and introduction accurately reflect the
526 paper's contributions and scope?

527 Answer: [Yes]

528 Justification: The abstract and introduction clearly state the contributions and scope of the
529 paper.

530 Guidelines:

- 531 • The answer NA means that the abstract and introduction do not include the claims
532 made in the paper.
- 533 • The abstract and/or introduction should clearly state the claims made, including the
534 contributions made in the paper and important assumptions and limitations. A No or
535 NA answer to this question will not be perceived well by the reviewers.
- 536 • The claims made should match theoretical and experimental results, and reflect how
537 much the results can be expected to generalize to other settings.
- 538 • It is fine to include aspirational goals as motivation as long as it is clear that these goals
539 are not attained by the paper.

540 **2. Limitations**

541 Question: Does the paper discuss the limitations of the work performed by the authors?

542 Answer: [Yes]

543 Justification: Our paper includes a dedicated "Limitation" paragraph where potential limita-
544 tions are discussed.

545 Guidelines:

- 546 • The answer NA means that the paper has no limitation while the answer No means that
547 the paper has limitations, but those are not discussed in the paper.
- 548 • The authors are encouraged to create a separate "Limitations" section in their paper.
- 549 • The paper should point out any strong assumptions and how robust the results are to
550 violations of these assumptions (e.g., independence assumptions, noiseless settings,
551 model well-specification, asymptotic approximations only holding locally). The authors
552 should reflect on how these assumptions might be violated in practice and what the
553 implications would be.
- 554 • The authors should reflect on the scope of the claims made, e.g., if the approach was
555 only tested on a few datasets or with a few runs. In general, empirical results often
556 depend on implicit assumptions, which should be articulated.
- 557 • The authors should reflect on the factors that influence the performance of the approach.
558 For example, a facial recognition algorithm may perform poorly when image resolution
559 is low or images are taken in low lighting. Or a speech-to-text system might not be
560 used reliably to provide closed captions for online lectures because it fails to handle
561 technical jargon.
- 562 • The authors should discuss the computational efficiency of the proposed algorithms
563 and how they scale with dataset size.
- 564 • If applicable, the authors should discuss possible limitations of their approach to
565 address problems of privacy and fairness.
- 566 • While the authors might fear that complete honesty about limitations might be used by
567 reviewers as grounds for rejection, a worse outcome might be that reviewers discover
568 limitations that aren't acknowledged in the paper. The authors should use their best
569 judgment and recognize that individual actions in favor of transparency play an impor-
570 tant role in developing norms that preserve the integrity of the community. Reviewers
571 will be specifically instructed to not penalize honesty concerning limitations.

572 **3. Theory Assumptions and Proofs**

573 Question: For each theoretical result, does the paper provide the full set of assumptions and
574 a complete (and correct) proof?

575 Answer: [NA]

576 Justification: Our paper does not include theoretical results.

577 Guidelines:

- 578 • The answer NA means that the paper does not include theoretical results.
- 579 • All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- 580 • All assumptions should be clearly stated or referenced in the statement of any theorems.
- 581 • The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- 582 • Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- 583 • Theorems and Lemmas that the proof relies upon should be properly referenced.

584 4. Experimental Result Reproducibility

589 Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

592 Answer: [Yes]

593 Justification: We clearly describe all the information needed to reproduce the experimental results in the paper, including our method, hyperparameter values, and datasets.

595 Guidelines:

- 596 • The answer NA means that the paper does not include experiments.
- 597 • If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- 599 • If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- 600 • Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- 601 • While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
 - 614 (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
 - 615 (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
 - 616 (c) If the contribution is a new model (e.g., a large language model), then there should 619 either be a way to access this model for reproducing the results or a way to reproduce 620 the model (e.g., with an open-source dataset or instructions for how to construct 621 the dataset).
 - 622 (d) We recognize that reproducibility may be tricky in some cases, in which case 623 authors are welcome to describe the particular way they provide for reproducibility. 624 In the case of closed-source models, it may be that access to the model is limited in 625 some way (e.g., to registered users), but it should be possible for other researchers 626 to have some path to reproducing or verifying the results.

627 5. Open access to data and code

628 Question: Does the paper provide open access to the data and code, with sufficient instruc-
629 tions to faithfully reproduce the main experimental results, as described in supplemental
630 material?

631 Answer: [Yes]

632 Justification: We use the publicly datasets and provided sufficient instructions on how to
633 use them for our method in the supplementary materials. The source code will be publicly
634 available.

635 Guidelines:

- The answer NA means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- While we encourage the release of code and data, we understand that this might not be possible, so "No" is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

655 6. Experimental Setting/Details

656 Question: Does the paper specify all the training and test details (e.g., data splits, hyper-
657 parameters, how they were chosen, type of optimizer, etc.) necessary to understand the
658 results?

659 Answer: [Yes]

660 Justification: We specify hyperparameters in the "Implementation Details" section of the
661 paper, with any unspecified hyperparameters being the same as the baselines. The supple-
662 mentary materials include detailed information on data splits.

663 Guidelines:

- The answer NA means that the paper does not include experiments.
- The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- The full details can be provided either with the code, in appendix, or as supplemental material.

669 7. Experiment Statistical Significance

670 Question: Does the paper report error bars suitably and correctly defined or other appropriate
671 information about the statistical significance of the experiments?

672 Answer: [No]

673 Justification: We report the average performance in the experimental results.

674 Guidelines:

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.

- 679 • The factors of variability that the error bars are capturing should be clearly stated (for
 680 example, train/test split, initialization, random drawing of some parameter, or overall
 681 run with given experimental conditions).
 682 • The method for calculating the error bars should be explained (closed form formula,
 683 call to a library function, bootstrap, etc.)
 684 • The assumptions made should be given (e.g., Normally distributed errors).
 685 • It should be clear whether the error bar is the standard deviation or the standard error
 686 of the mean.
 687 • It is OK to report 1-sigma error bars, but one should state it. The authors should
 688 preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis
 689 of Normality of errors is not verified.
 690 • For asymmetric distributions, the authors should be careful not to show in tables or
 691 figures symmetric error bars that would yield results that are out of range (e.g. negative
 692 error rates).
 693 • If error bars are reported in tables or plots, The authors should explain in the text how
 694 they were calculated and reference the corresponding figures or tables in the text.

695 **8. Experiments Compute Resources**

696 Question: For each experiment, does the paper provide sufficient information on the com-
 697 puter resources (type of compute workers, memory, time of execution) needed to reproduce
 698 the experiments?

699 Answer: [Yes]

700 Justification: We provide the information of GPU resource we used in the supplementary
 701 materials.

702 Guidelines:

- 703 • The answer NA means that the paper does not include experiments.
 704 • The paper should indicate the type of compute workers CPU or GPU, internal cluster,
 705 or cloud provider, including relevant memory and storage.
 706 • The paper should provide the amount of compute required for each of the individual
 707 experimental runs as well as estimate the total compute.
 708 • The paper should disclose whether the full research project required more compute
 709 than the experiments reported in the paper (e.g., preliminary or failed experiments that
 710 didn't make it into the paper).

711 **9. Code Of Ethics**

712 Question: Does the research conducted in the paper conform, in every respect, with the
 713 NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines>?

714 Answer: [Yes]

715 Justification: Our research conform with the NeurIPS Code of Ethics in every respect.

716 Guidelines:

- 717 • The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
 718 • If the authors answer No, they should explain the special circumstances that require a
 719 deviation from the Code of Ethics.
 720 • The authors should make sure to preserve anonymity (e.g., if there is a special consid-
 721 eration due to laws or regulations in their jurisdiction).

722 **10. Broader Impacts**

723 Question: Does the paper discuss both potential positive societal impacts and negative
 724 societal impacts of the work performed?

725 Answer: [NA]

726 Justification: Our paper proposes a method for synthesizing novel view images from sparse
 727 input views. The aim is to observe the objects in the input images from new perspectives.
 728 There is no generation of new or fake data involved.

729 Guidelines:

- The answer NA means that there is no societal impact of the work performed.
- If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.
- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

11. Safeguards

Question: Does the paper describe safeguards that have been put in place for responsible release of data or models that have a high risk for misuse (e.g., pretrained language models, image generators, or scraped datasets)?

Answer: [NA]

Justification: Our paper poses no such risks.

Guidelines:

- The answer NA means that the paper poses no such risks.
- Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

12. Licenses for existing assets

Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?

Answer: [Yes]

Justification: We cite the original papers of the code and datasets we used.

Guidelines:

- The answer NA means that the paper does not use existing assets.
- The authors should cite the original paper that produced the code package or dataset.
- The authors should state which version of the asset is used and, if possible, include a URL.
- The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.

- 783 • If assets are released, the license, copyright information, and terms of use in the
 784 package should be provided. For popular datasets, paperswithcode.com/datasets
 785 has curated licenses for some datasets. Their licensing guide can help determine the
 786 license of a dataset.
 787 • For existing datasets that are re-packaged, both the original license and the license of
 788 the derived asset (if it has changed) should be provided.
 789 • If this information is not available online, the authors are encouraged to reach out to
 790 the asset's creators.

791 **13. New Assets**

792 Question: Are new assets introduced in the paper well documented and is the documentation
 793 provided alongside the assets?

794 Answer: [NA]

795 Justification: Our paper does not release new assets.

796 Guidelines:

- 797 • The answer NA means that the paper does not release new assets.
- 798 • Researchers should communicate the details of the dataset/code/model as part of their
 799 submissions via structured templates. This includes details about training, license,
 800 limitations, etc.
- 801 • The paper should discuss whether and how consent was obtained from people whose
 802 asset is used.
- 803 • At submission time, remember to anonymize your assets (if applicable). You can either
 804 create an anonymized URL or include an anonymized zip file.

805 **14. Crowdsourcing and Research with Human Subjects**

806 Question: For crowdsourcing experiments and research with human subjects, does the paper
 807 include the full text of instructions given to participants and screenshots, if applicable, as
 808 well as details about compensation (if any)?

809 Answer: [NA]

810 Justification: Our paper does not involve crowdsourcing nor research with human subjects.

811 Guidelines:

- 812 • The answer NA means that the paper does not involve crowdsourcing nor research with
 813 human subjects.
- 814 • Including this information in the supplemental material is fine, but if the main contribu-
 815 tion of the paper involves human subjects, then as much detail as possible should be
 816 included in the main paper.
- 817 • According to the NeurIPS Code of Ethics, workers involved in data collection, curation,
 818 or other labor should be paid at least the minimum wage in the country of the data
 819 collector.

820 **15. Institutional Review Board (IRB) Approvals or Equivalent for Research with Human
 821 Subjects**

822 Question: Does the paper describe potential risks incurred by study participants, whether
 823 such risks were disclosed to the subjects, and whether Institutional Review Board (IRB)
 824 approvals (or an equivalent approval/review based on the requirements of your country or
 825 institution) were obtained?

826 Answer: [NA]

827 Justification: Our paper does not involve crowdsourcing nor research with human subjects.

828 Guidelines:

- 829 • The answer NA means that the paper does not involve crowdsourcing nor research with
 830 human subjects.
- 831 • Depending on the country in which research is conducted, IRB approval (or equivalent)
 832 may be required for any human subjects research. If you obtained IRB approval, you
 833 should clearly state this in the paper.

- 834 • We recognize that the procedures for this may vary significantly between institutions
835 and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the
836 guidelines for their institution.
837 • For initial submissions, do not include any information that would break anonymity (if
838 applicable), such as the institution conducting the review.