

# Hi-Gaussian: Hierarchical Gaussians under Normalized Spherical Projection for Single-View 3D Reconstruction

## Supplementary Material

### 1. More Theoretical Study

#### 1.1. Normalized Spherical Projection

In Sec. 3.2 of the main text, we propose normalized spherical projection, which introduces less spatial compression than the un-normalized one. In other words, the curvature of the upper edges in the images obtained from normalized spherical projection is flatter, which is shown in Fig. 8 of the main text. Now we provide a theoretical analysis of the phenomenon described above.

The formula for un-normalized spherical projection is given as follows:

$$\begin{bmatrix} \theta \\ \phi \end{bmatrix} = \begin{bmatrix} \pi - \arctan(x^{-1}) \\ \arccos(-y/r) \end{bmatrix} \quad (1)$$

where  $[x, y]^\top$  is coordinates in the camera coordinate system and  $r = \sqrt{x^2 + y^2 + 1}$ . We consider  $\theta$  and  $\phi$  as functions of  $x$  and  $y$ , denoting them as  $\theta = h_1(x, y)$  and  $\phi = h_2(x, y)$  respectively. From Eqn. 1, we can see that  $h_1(x, y) = \pi - \arctan(x^{-1})$  is solely dependent on  $x$ , while  $h_2(x, y) = \arccos(-y/r) = \arccos(-y/\sqrt{x^2 + y^2 + 1})$  is dependent on both  $x$  and  $y$ . Taking the partial derivative of  $h_2(x, y)$  with respect to  $x$ , we obtain

$$\frac{\partial h_2(x, y)}{\partial x} = -\frac{xy}{(x^2 + y^2 + 1)\sqrt{x^2 + y^2 + 1}} \quad (2)$$
$$\begin{cases} > 0 & \text{if } xy < 0 \\ = 0 & \text{if } xy = 0 \\ < 0 & \text{if } xy > 0 \end{cases}$$

Thus, for a given constant  $y_0 > 0$ ,  $h_2(x, y_0)$  is strictly monotonically increasing for  $x \in (-\infty, 0)$ , and strictly monotonically decreasing for  $x \in (0, +\infty)$ . Similarly,  $h_2(x, -y_0)$  is strictly monotonically decreasing for  $x \in (-\infty, 0)$ , and strictly monotonically increasing for  $x \in (0, +\infty)$ . Moreover, since  $h_2(x, y) = h_2(-x, y)$ , it is evident that  $h_2(x, y)$  is an even function with respect to  $x$ . Therefore, for any  $x \in \mathbb{R}$ , for any  $\hat{x} \in \{x' \mid |x'| < |x|\}$ , we have  $h_2(\hat{x}, y_0) > h_2(\pm x, y_0)$  and  $h_2(\hat{x}, -y_0) < h_2(\pm x, -y_0)$ .

Since normalization typically transforms  $x$  into a smaller  $\hat{x}$ , our normalized spherical projection maps Cartesian coordinates to a smaller range of spherical coordinates. This results in the upper and lower edges of the images obtained from normalized spherical projection appear relatively flatter than the un-normalized one.

### 2. More Experiments

Additional quantitative and qualitative results are presented to further elaborate on the superiority of our proposed method. We start by showcasing a more comprehensive comparison with other methods, then present a more detailed ablation study.

#### 2.1. Quantitative Comparison with other Methods

**Cross-dataset Novel View Synthesis.** We present quantitative results on cross-dataset generalization. Models are trained on BundleFusion [3] and are tested on all 8 sequences in NeRF-LLFF [7] dataset. Table 1 shows that our method achieves the best performance on most sequences in NeRF-LLFF dataset.

**Novel Depth Synthesis.** To better showcase the depth estimation capability of our method, we conduct evaluation on novel depth synthesis following the experimental setup of SceneRF [2]. Our approach outperforms other methods across all metrics, which is shown in Table 2 and Table 3.

**Scene Reconstruction.** We also evaluate mesh reconstruction following SceneRF [2] for a fair comparison. As demonstrated in Table 4, our method achieves the best reconstruction performance. It is worth noting that if we proposed an improved approach for mesh reconstruction, we could get better results.

**Ablation Study.** To further demonstrate the effectiveness of our method, we conduct more ablation studies. We compare the quantitative results of novel view synthesis at different distances from the input view. The results in Fig. 1 show that our Normalized Spherical Projection module and Hierarchical Gaussian Sampling strategy consistently enhance the performance of our model, regardless of the distance between the novel view and the input view.

#### 2.2. Qualitative Comparison with Other Methods

**Cross-dataset Novel View Synthesis.** To better demonstrate the generalization capability of our method, we conduct cross-dataset evaluations on novel view synthesis. Models are trained on BundleFusion [3] and are tested on NeRF-LLFF [7] dataset. The qualitative results in Fig. 2 indicate that our method renders the sharpest and clearest images in cross-dataset generalization.

Methods	Fern			Flower			Fortress			Horns		
	PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$	PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$	PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$	PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$
SceneRF	<u>13.64</u>	0.130	0.702	<u>13.31</u>	<b>0.263</b>	0.584	<u>15.46</u>	0.290	0.636	<u>14.55</u>	0.222	0.650
Splatter Image	11.89	<u>0.133</u>	<u>0.610</u>	12.36	0.107	<u>0.523</u>	8.95	0.218	<u>0.565</u>	14.48	<u>0.266</u>	<u>0.469</u>
Ours	<b>16.93</b>	<b>0.439</b>	<b>0.359</b>	<b>14.01</b>	<u>0.167</u>	<b>0.468</b>	<b>17.59</b>	<b>0.302</b>	<b>0.318</b>	<b>15.84</b>	<b>0.398</b>	<b>0.379</b>

Methods	Leaves			Orchids			Room			Trex		
	PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$	PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$	PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$	PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$
SceneRF	<b>13.77</b>	<b>0.161</b>	0.548	<u>11.31</u>	<u>0.128</u>	0.716	<u>12.74</u>	0.286	0.743	<u>10.65</u>	0.093	0.757
Splatter Image	12.66	0.120	<u>0.488</u>	9.88	0.067	<u>0.615</u>	9.62	0.207	<u>0.648</u>	9.22	0.074	<u>0.638</u>
Ours	<u>13.00</u>	<u>0.158</u>	<b>0.443</b>	<b>11.76</b>	<b>0.145</b>	<b>0.507</b>	<b>16.26</b>	<b>0.454</b>	<b>0.328</b>	<b>16.91</b>	<b>0.453</b>	<b>0.274</b>

Table 1. Quantitative evaluations on cross-dataset generalization from BundleFusion to NeRF-LLFF dataset.

Methods	SemanticKITTI						
	Abs Rel $\downarrow$	Sq Rel $\downarrow$	RMSE $\downarrow$	RMSE log $\downarrow$	$\delta_1 \uparrow$	$\delta_2 \uparrow$	$\delta_3 \uparrow$
PixelNeRF [9]	0.2364	2.080	6.449	0.3354	65.81	85.43	92.90
MINE [5]	0.2248	1.787	6.343	0.3283	65.87	85.52	93.30
VisionNeRF [6]	0.2054	1.490	5.841	0.3073	69.11	88.28	94.37
SceneRF [2]	<u>0.1681</u>	<u>1.291</u>	<u>5.781</u>	<u>0.2851</u>	<u>75.07</u>	<u>89.09</u>	<u>94.50</u>
SplatterImage [8]	0.2519	2.127	7.282	0.4205	58.41	79.30	89.02
Ours	<b>0.1165</b>	<b>0.812</b>	<b>4.702</b>	<b>0.2397</b>	<b>80.99</b>	<b>90.02</b>	<b>94.67</b>

Table 2. Novel depth synthesis on SemanticKITTI datasets.

**Mesh Visualization of Scene Reconstruction.** To offer a more intuitive representation of scene reconstruction, we display 3D meshes on the validation set of SemanticKITTI [1, 4] and BundleFusion [3]. These meshes are produced from the scene TSDF, which is obtained through the conversion of rendered images and depths. Fig. 4 and Fig. 5 demonstrate that our method reconstructs the clearest and sharpest meshes.

**Ablation study.** Due to the limitation of the main text’s length, we have only presented the qualitative results of ablation study on SemanticKITTI. Therefore, the results for BundleFusion are provided here, as shown in Fig. 3. Without Normalized Spherical Projection, our model fails to reconstruct outside the input FOV, leading to ghosting artifacts. Without Hierarchical Gaussian Sampling, there will be ripple artifacts and holes in the rendered novel views.

## References

- [1] Jens Behley, Martin Garbade, Andres Milioto, Jan Quenzel, Sven Behnke, Cyrill Stachniss, and Jurgen Gall. SemanticKITTI: A dataset for semantic scene understanding of lidar sequences. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 9297–9307, 2019. 2
- [2] Anh-Quan Cao and Raoul de Charette. Scenerf: Self-supervised monocular 3d scene reconstruction with radiance fields. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 9387–9398, 2023. 1, 2, 3
- [3] Angela Dai, Matthias Nießner, Michael Zollhöfer, Shahram Izadi, and Christian Theobalt. Bundlefusion: Real-time globally consistent 3d reconstruction using on-the-fly surface reintegration. *ACM Transactions on Graphics (ToG)*, 36(4):1, 2017. 1, 2
- [4] Andreas Geiger, Philip Lenz, and Raquel Urtasun. Are we ready for autonomous driving? the kitti vision benchmark suite. In *2012 IEEE conference on computer vision and pattern recognition*, pages 3354–3361. IEEE, 2012. 2
- [5] Jiaxin Li, Zijian Feng, Qi She, Henghui Ding, Changhu Wang, and Gim Hee Lee. Mine: Towards continuous depth mpi with nerf for novel view synthesis. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 12578–12588, 2021. 2, 3
- [6] Kai-En Lin, Yen-Chen Lin, Wei-Sheng Lai, Tsung-Yi Lin, Yi-Chang Shih, and Ravi Ramamoorthi. Vision transformer for nerf-based view synthesis from a single input image. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 806–815, 2023. 2, 3
- [7] Ben Mildenhall, Pratul P Srinivasan, Rodrigo Ortiz-Cayon, Nima Khademi Kalantari, Ravi Ramamoorthi, Ren Ng, and Abhishek Kar. Local light field fusion: Practical view synthesis with prescriptive sampling guidelines. *ACM Transactions on Graphics (ToG)*, 38(4):1–14, 2019. 1
- [8] Stanislaw Szymanowicz, Christian Rupprecht, and Andrea Vedaldi. Splatter image: Ultra-fast single-view 3d reconstruction. In *Proceedings of the IEEE/CVF Conference on Com-*

Methods	<b>BundleFusion</b>						
	Abs Rel↓	Sq Rel↓	RMSE↓	RMSE log↓	$\delta_1 \uparrow$	$\delta_2 \uparrow$	$\delta_3 \uparrow$
PixelNeRF [9]	0.6029	2.312	1.750	0.5904	46.34	72.38	83.89
MINE [5]	0.1839	0.098	0.386	0.2386	65.53	91.78	98.21
VisionNeRF [6]	0.5958	2.468	1.783	0.5586	55.47	79.29	86.68
SceneRF [2]	<u>0.1766</u>	<u>0.094</u>	<u>0.368</u>	<u>0.2100</u>	<u>72.71</u>	<u>94.89</u>	<u>99.23</u>
SplatterImage [8]	0.2407	0.142	0.454	0.2710	57.06	89.00	97.99
Ours	<b>0.0792</b>	<b>0.041</b>	<b>0.225</b>	<b>0.1101</b>	<b>95.28</b>	<b>99.30</b>	<b>99.75</b>

Table 3. Novel depth synthesis on BundleFusion datasets.

Methods	<b>SemanticKITTI</b>			<b>BundleFusion</b>		
	IoU↑	Prec.↑	Rec.↑	IoU↑	Prec.↑	Rec.↑
SceneRF [2]	<u>13.84</u>	<u>17.28</u>	40.96	<u>20.16</u>	<u>25.82</u>	<u>47.92</u>
Splatter Image [8]	10.30	11.30	<u>53.93</u>	13.89	22.22	27.04
Ours	<b>15.56</b>	<b>17.39</b>	<b>59.72</b>	<b>40.42</b>	<b>48.91</b>	<b>69.96</b>

Table 4. Reconstruction evaluations on SemanticKITTI and BundleFusion datasets. We outperform all other methods across all metrics.

*puter Vision and Pattern Recognition*, pages 10208–10217, 2024. 2, 3

- [9] Alex Yu, Vickie Ye, Matthew Tancik, and Angjoo Kanazawa. pixelnerf: Neural radiance fields from one or few images. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 4578–4587, 2021. 2, 3

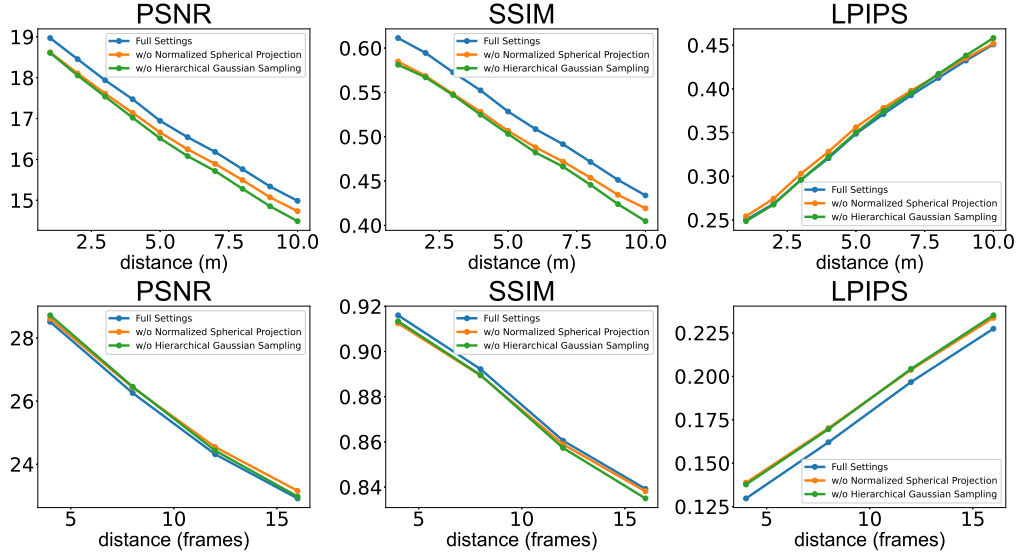


Figure 1. Performances at varying input view distances on SemanticKITTI (the first row) and BundleFusion (the second row) for ablation study. Our model with full settings usually achieves the best performance.



Figure 2. Qualitative evaluations on cross-dataset generalization from BundleFusion to NeRF-LLFF dataset.





Figure 3. Ablations on BundleFusion (val). Without Normalized Spherical Projection, our model fails to reconstruct outside the input FOV, leading to ghosting artifacts. Without Hierarchical Gaussian Sampling, there will be ripple artifacts and holes in the rendered novel views.



Figure 4. 3D meshes on SemanticKITTI (val.).

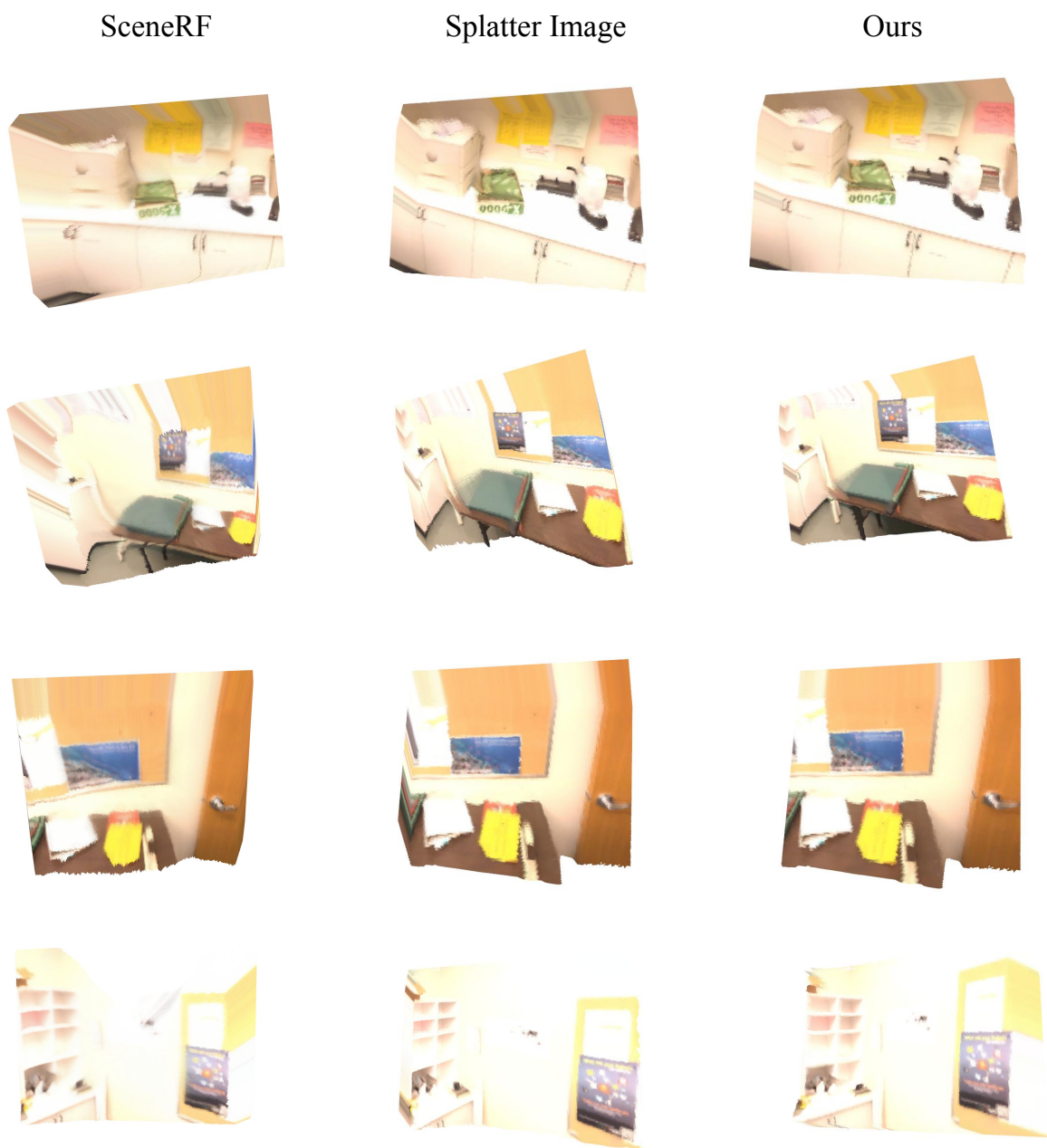


Figure 5. 3D meshes on BundleFusion (val.).