Taming Video Diffusion Prior with Scene-Grounding Guidance for 3D Gaussian Splatting from Sparse Inputs

Supplementary Material

The project page is available at HERE, which contains demo videos for better visualization.

A. Implementation Details

During the denoising sampling process, we employ the DDIM sampler [10] combined with our proposed guidance, setting the number of sampling steps to 50. Regarding the trajectory initialization strategy, for each input view in its camera space, we sample views by changing the polar/azimuth angle to $[-30^\circ, -15^\circ, 0^\circ, 15^\circ, 30^\circ]$, and setting the radial distance to $[1, \frac{1}{3}, \frac{1}{10}]$ of the depth of the central distance to $[1, \frac{1}{3}, \frac{1}{10}]$ of the depth of the central distance to $[1, \frac{1}{3}, \frac{1}{10}]$ of the depth of the central distance to $[1, \frac{1}{3}, \frac{1}{10}]$ of the depth of the central distance to $[1, \frac{1}{3}, \frac{1}{10}]$ of the depth of the central distance to $[1, \frac{1}{3}, \frac{1}{10}]$ of the depth of the central distance to $[1, \frac{1}{3}, \frac{1}{10}]$ of the depth of the central distance to $[1, \frac{1}{3}, \frac{1}{10}]$ of the depth of the central distance to $[1, \frac{1}{3}, \frac{1}{10}]$ of the depth of the central distance to $[1, \frac{1}{3}, \frac{1}{10}]$ of the depth of the central distance to $[1, \frac{1}{3}, \frac{1}{10}]$ of the depth of the central distance to $[1, \frac{1}{3}, \frac{1}{10}]$ of the depth of the central distance to $[1, \frac{1}{3}, \frac{1}{10}]$ of the depth of the central distance to $[1, \frac{1}{3}, \frac{1}{10}]$ of the depth of th ter pixel (from the prediction of ViewCrafter [17]). Out of 75 sampled views, we discard those whose renderings exhibit holes larger than 10% of the image size (to filter out uncommon viewpoints), then select the top 6 views with the largest holes from the remaining. To obtain the point cloud used for initialization, we follow the standard pipeline provided on the DUSt3R [13] webpage. Since our focus is sparse-input radiance fields reconstruction, the groundtruth camera poses and intrinsics are provided. During DUSt3R optimization, we fix both the poses and intrinsics to their groundtruth values. In the main paper, we conduct experiments on a new benchmark that is created from two indoor datasets, the synthetic Replica [11] and the realistic Scan-Net++ [16] datasets. Please refer to [18] for more details about the benchmark.

B. More Results

Our method focuses on holistic modeling of an indoor scene of a moderate size, and we conduct the experiments in the main paper with 6 input views, since 6 input views are basically sufficient to cover the entire room. To validate the effectiveness of our method, we also test our method with different number of views following the common 3/6/9view settings of sparse-input modeling. Tab. A1 validates that, our method is effective given different number of input views, with consistent improvements over our baseline. InstantSplat [2] is a strong baseline of sparse-input pose-free modeling, leveraging DUSt3R [13] point cloud for 3DGS initialization. Our method also consistently outperforms InstantSplat as shown in Tab. A1.

To obtain a thorough understanding of the source of the performance improvement, we show some quantitative results regarding performances of observable and the unobservable regions respectively in Tab. A2. The results show that our method brings improvement in both observable and unobservable regions.

We further compare our method with two representative methods that leverage diffusion models for sparse-input modeling, ReconFusion [14] and CAT3D [3] on the datasets of RealEstate10K and LLFF. We adhere to their settings for fair comparisons and the results are shown in Tab. A3. On the LLFF dataset, our method is based on the strong baseline of binocular-guided 3DGS [4]. The results show that our method achieves comparable performance with both ReconFusion and CAT3D.

We provide per-scene comparisons in Table A4, demonstrating that our method consistently achieves superior performance across all scenes. Additional qualitative results are shown in Fig. A3. These results highlight the effectiveness of our approach in addressing issues such as extrapolation and occlusion, as seen in examples like the wall behind the chair (second row) and the ceiling (third row). Furthermore, our method preserves more intact structures with finer details, such as the edges in the fifth and sixth rows.

We present a comparison of the generated sequences from the video diffusion model with and without the proposed guidance in Fig. A2. The results clearly show that our proposed guidance enhances the plausibility of the generated sequences by maintaining consistent appearances and ensuring that only elements present in the scene are generated. Consistency in the generated video is crucial for effective 3DGS optimization. Using inconsistent sequences for 3DGS optimization often leads to artifacts, such as black shadows in the renderings, which significantly degrade visual quality, as demonstrated on the demo page.

C. Discussion

While our approach significantly improves overall quality by addressing extrapolation and occlusion challenges, we observe that it occasionally produces over-smoothed results. We hypothesize that this is due to the limited resolution supported by the video diffusion model during generation. On a 32GB V100 GPU, we are constrained to generating sequences at resolutions of 320×448 for the Replica dataset and 320×512 for the ScanNet++ dataset, which are subsequently upsampled to rendering resolutions of 480×640 and 480×720, respectively, for supervision during 3DGS optimization. This upsampling process introduces undersampling, which can smooth out certain regions and result in over-smoothed effects. Addressing the challenge of preserving high-frequency details during 3DGS optimization under resource-limited sequence generation remains an open problem and is a direction for future work.



Figure A1. Point clouds from DUSt3R [13] optimized with sparse input views on the Replica dataset. The yellow parts represent unobserved regions, e.g., regions that are outside the field of view or occluded. Note that the ceilings are removed for better visualization.



Figure A2. Generated frames from the video diffusion model with and without the proposed guidance. The numbers at the top indicate the frame IDs. The first frame corresponds to an image from the sparse input views, while other frames are generated. Without guidance, the generated sequences exhibit significant inconsistencies: (i) appearance inconsistencies, highlighted by the blue boxes; and (ii) hallucinated elements that do not exist in the scene, highlighted by the red boxes. In contrast, with the proposed guidance, the generated sequences are more plausible and consistent.

	3-view				6-view		9-view			
	PSNR↑	SSIM↑	LPIPS↓	PSNR↑	SSIM↑	LPIPS↓	PSNR↑	SSIM↑	LPIPS↓	
Replica										
Baseline 3DGS	19.87	0.794	0.178	22.80	0.818	0.179	24.81	0.863	0.124	
InstantSplat [2]	20.49	0.766	0.226	20.35	0.760	0.290	18.44	0.708	0.373	
Ours	23.98	0.848	0.136	26.35	0.872	0.122	27.42	0.891	0.111	

Table A1. Our method brings performance improvement over the baseline with different number of input views, and consistently outperforms another strong sparse-input modeling baseline InstantSplat [2].

Danling 6 view		Full Image	e	Obse	rvable Re	gions	Unobservable Regions			
Keplica 0-view	PSNR ↑	SSIM↑	$\text{LPIPS}{\downarrow}$	PSNR↑	SSIM↑	LPIPS↓	PSNR ↑	SSIM↑	LPIPS↓	
Baseline 3DGS	22.80	0.818	0.179	25.45	0.860	0.129	14.27	0.967	0.029	
w/ Vanilla Generation	23.69	0.840	0.160	25.00	0.870	0.119	17.11	0.977	0.025	
Ours	26.35	0.872	0.122	27.12	0.894	0.091	20.85	0.985	0.020	
Baseline 3DGS+LaMa [12] Baseline 3DGS+SDInpaint [9]*	24.56 25.15	0.833 0.853	0.167 0.141	25.45 26.13	0.860 0.878	0.129 0.104	17.80 19.25	0.981 0.982	0.021 0.022	

Table A2. Analysis of performance regarding observable and unobservable regions. * refers to incorporating our trajectory initialization strategy. The methods in the second block utilize inpainting models.

	PSNR ↑	3-view SSIM↑	LPIPS↓	 PSNR↑	6-view SSIM↑	LPIPS↓	PSNR ↑	9-view SSIM↑	LPIPS↓
RealEstate10K									
ReconFusion [14]	25.84	0.910	0.144	29.99	0.951	0.103	31.82	0.961	0.092
CAT3D [3]	26.78	0.917	0.132	31.07	0.954	0.092	32.20	0.963	0.082
Ours	25.03	0.871	0.136	30.62	0.944	0.069	32.45	0.955	0.062
LLFF									
ReconFusion [14]	21.34	0.724	0.203	24.25	0.815	0.152	25.21	0.848	0.134
CAT3D [3]	21.58	0.731	0.181	24.71	0.833	0.121	25.63	0.860	0.107
Ours	21.35	0.746	0.173	25.13	0.851	0.102	26.29	0.880	0.084

Table A3. Comparisons with ReconFusion [14] and CAT3D [3] on the RealEstate10K and LLFF datasets.

	ScanNet++ [16]					Replica [11]						
	a2ccc	8a20d	94ee1	78318	avg	office2	office3	office4	room0	room1	room2	avg
Mip-NeRF [1]	18.28	23.48	16.93	19.63	19.58	17.43	19.04	19.08	17.46	16.57	19.16	18.12
	0.759	0.799	0.725	0.735	0.755	0.539	0.685	0.727	0.762	0.721	0.808	0.707
	0.351	0.321	0.431	0.451	0.389	0.486	0.421	0.393	0.342	0.386	0.317	0.391
InfoNeRF [7]	13.90	17.69	14.34	12.21	14.54	13.66	12.53	11.51	12.58	14.11	14.00	13.07
	0.662	0.691	0.627	0.605	0.646	0.463	0.545	0.592	0.618	0.689	0.678	0.598
	0.408	0.457	0.310	0.558	0.493	0.012	0.025	0.024	0.342	0.455	0.477	0.332
	20.67	23.00	15.34	20.02	19.76	19.12	19.35	18.97	19.84	17.18	19.46	18.99
DietNeRF [5]	0.751	0.776	0.627	0.725	0.719	0.612	0.695	0.419	0.783	0.749	0.797	0.676
	0.385	0.303	0.310	0.439	0.431	0.438	0.417	0.721	0.34	0.380	0.343	0.444
	19.93	22.37	19.42	18.94	20.17	20.89	21.06	20.25	22.55	19.69	21.43	20.99
FreeNeRF [15]	0.759	0.791	0.762	0.711	0.756	0.688	0.735	0.750	0.831	0.781	0.807	0.765
	0.307	0.299	0.417	0.449	0.308	0.559	0.340	0.304	0.234	0.323	0.321	0.324
S ³ NeRF [18]	21.81	25.60	20.05	21.36	22.21	22.79	23.83	23.08	24.01	19.66	21.87	22.54
	0.801	0.811	0.784	0.753	0.787	0.728	0.773	0.801	0.862	0.808	0.825	0.800
	0.324	0.330	0.357	0.444	0.364	0.326	0.309	0.301	0.213	0.277	0.293	0.287
*	20.65	23.49	20.38	21.11	21.41	25.03	23.60	22.14	20.32	22.68	23.07	22.80
3DGS↓ [6]	0.824	0.857	0.821	0.764	0.817	0.873	0.858	0.834	0.720	0.802	0.824	0.818
	0.193	0.136	0.218	0.298	0.211	0.141	0.147	0.180	0.204	0.203	0.196	0.179
	19.10	21.21	17.55	18.20	19.01	22.68	18.40	12.31	12.60	18.87	20.91	17.63
DNGaussian [8]	0.765	0.781	0.743	0.730	0.755	0.843	0.789	0.644	0.534	0.708	0.790	0.718
	0.343	0.292	0.382	0.450	0.367	0.233	0.291	0.628	0.722	0.397	0.338	0.435
•	20.47	23.73	18.90	19.61	20.68	25.31	23.34	21.83	20.33	22.59	22.88	22.71
DNGaussian [‡] [8]	0.805	0.842	0.784	0.722	0.788	0.890	0.853	0.837	0.729	0.800	0.820	0.821
	0.213	0.183	0.287	0.357	0.281	0.124	0.161	0.197	0.226	0.208	0.219	0.189
	19.19	18.98	15.77	17.87	17.95	20.70	20.26	21.62	19.65	19.23	19.89	20.22
FSGS [19]	0.760	0.735	0.719	0.708	0.730	0.802	0.790	0.825	0.654	0.712	0.779	0.760
	0.321	0.316	0.415	0.442	0.373	0.266	0.255	0.271	0.315	0.374	0.342	0.304
FSGS [‡] [19]	21.28	22.56	20.28	20.79	21.23	24.37	23.41	23.45	21.02	23.56	22.14	22.99
	0.826	0.844	0.815	0.767	0.813	0.873	0.856	0.862	0.759	0.823	0.822	0.833
	0.219	0.193	0.267	0.350	0.257	0.194	0.174	0.189	0.198	0.205	0.270	0.205
Ours	25.21	25.10	23.10	22.16	23.89	27.46	26.81	27.43	24.85	26.00	25.53	26.35
	0.857	0.882	0.860	0.803	0.850	0.916	0.902	0.897	0.796	0.851	0.872	0.872
	0.157	0.118	0.201	0.269	0.182	0.083	0.099	0.122	0.145	0.142	0.142	0.122

Table A4. Per-scene performance of various models on the ScanNet++ and Replica datasets. For each method, the three rows represent PSNR, SSIM, and LPIPS, respectively. *avg* indicates the average performance across all scenes in each dataset. Including our approach, 3DGS-based methods marked with \uparrow are initialized with the point cloud from DUSt3R [13].



Replica ScanNet++



Figure A3. Qualitative comparisons between other works on Replica and ScanNet++ datasets. All 3DGS-based methods are optimized using the initialized point cloud from DUSt3R [13].

References

- Jonathan T Barron, Ben Mildenhall, Matthew Tancik, Peter Hedman, Ricardo Martin-Brualla, and Pratul P Srinivasan. Mip-nerf: A multiscale representation for anti-aliasing neural radiance fields. In *ICCV*, 2021. 3
- [2] Zhiwen Fan, Wenyan Cong, Kairun Wen, Kevin Wang, Jian Zhang, Xinghao Ding, Danfei Xu, Boris Ivanovic, Marco Pavone, Georgios Pavlakos, et al. Instantsplat: Unbounded sparse-view pose-free gaussian splatting in 40 seconds. arXiv preprint arXiv:2403.20309, 2024. 1, 2
- [3] Ruiqi Gao, Aleksander Holynski, Philipp Henzler, Arthur Brussee, Ricardo Martin-Brualla, Pratul Srinivasan, Jonathan T Barron, and Ben Poole. Cat3d: Create anything in 3d with multi-view diffusion models. In *NeurIPS*, 2024. 1, 3
- [4] Liang Han, Junsheng Zhou, Yu-Shen Liu, and Zhizhong Han. Binocular-guided 3d gaussian splatting with view consistency for sparse view synthesis. In *NeurIPS*, 2024. 1
- [5] Ajay Jain, Matthew Tancik, and Pieter Abbeel. Putting nerf on a diet: Semantically consistent few-shot view synthesis. In *ICCV*, 2021. 3
- [6] Bernhard Kerbl, Georgios Kopanas, Thomas Leimkühler, and George Drettakis. 3d gaussian splatting for real-time radiance field rendering. ACM Trans. Graph., 42(4):139–1, 2023. 3
- [7] Mijeong Kim, Seonguk Seo, and Bohyung Han. Infonerf: Ray entropy minimization for few-shot neural volume rendering. In CVPR, 2022. 3
- [8] Jiahe Li, Jiawei Zhang, Xiao Bai, Jin Zheng, Xin Ning, Jun Zhou, and Lin Gu. Dngaussian: Optimizing sparse-view 3d gaussian radiance fields with global-local depth normalization. In *CVPR*, 2024. 3
- [9] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *CVPR*, 2022. 3
- [10] Jiaming Song, Chenlin Meng, and Stefano Ermon.

Denoising diffusion implicit models. *arXiv preprint arXiv:2010.02502*, 2020. 1

- [11] Julian Straub, Thomas Whelan, Lingni Ma, Yufan Chen, Erik Wijmans, Simon Green, Jakob J Engel, Raul Mur-Artal, Carl Ren, Shobhit Verma, et al. The replica dataset: A digital replica of indoor spaces. arXiv preprint arXiv:1906.05797, 2019. 1, 3
- [12] Roman Suvorov, Elizaveta Logacheva, Anton Mashikhin, Anastasia Remizova, Arsenii Ashukha, Aleksei Silvestrov, Naejin Kong, Harshith Goka, Kiwoong Park, and Victor Lempitsky. Resolution-robust large mask inpainting with fourier convolutions. In WACV, 2022. 3
- [13] Shuzhe Wang, Vincent Leroy, Yohann Cabon, Boris Chidlovskii, and Jerome Revaud. Dust3r: Geometric 3d vision made easy. In *CVPR*, 2024. 1, 2, 3, 4
- [14] Rundi Wu, Ben Mildenhall, Philipp Henzler, Keunhong Park, Ruiqi Gao, Daniel Watson, Pratul P Srinivasan, Dor Verbin, Jonathan T Barron, Ben Poole, et al. Reconfusion: 3d reconstruction with diffusion priors. In CVPR, 2024. 1, 3
- [15] Jiawei Yang, Marco Pavone, and Yue Wang. Freenerf: Improving few-shot neural rendering with free frequency regularization. In *CVPR*, 2023. 3
- [16] Chandan Yeshwanth, Yueh-Cheng Liu, Matthias Nießner, and Angela Dai. Scannet++: A high-fidelity dataset of 3d indoor scenes. In *ICCV*, 2023. 1, 3
- [17] Wangbo Yu, Jinbo Xing, Li Yuan, Wenbo Hu, Xiaoyu Li, Zhipeng Huang, Xiangjun Gao, Tien-Tsin Wong, Ying Shan, and Yonghong Tian. Viewcrafter: Taming video diffusion models for high-fidelity novel view synthesis. arXiv preprint arXiv:2409.02048, 2024. 1
- [18] Yingji Zhong, Kaichen Zhou, Zhihao Li, Lanqing Hong, Zhenguo Li, and Dan Xu. Empowering sparse-input neural radiance fields with dual-level semantic guidance from dense novel views. arXiv preprint arXiv:2503.02230, 2025. 1, 3
- [19] Zehao Zhu, Zhiwen Fan, Yifan Jiang, and Zhangyang Wang. Fsgs: Real-time few-shot view synthesis using gaussian splatting. In ECCV, 2024. 3