

GS2Mesh: Surface Reconstruction from Gaussian Splatting via Novel Stereo Views

Yaniv Wolf*, Amit Bracha*, and Ron Kimmel

Technion - Israel Institute of Technology, Haifa, Israel
 {yaniv.wolf, amit.bracha, ron}@cs.technion.ac.il
<https://gs2mesh.github.io>

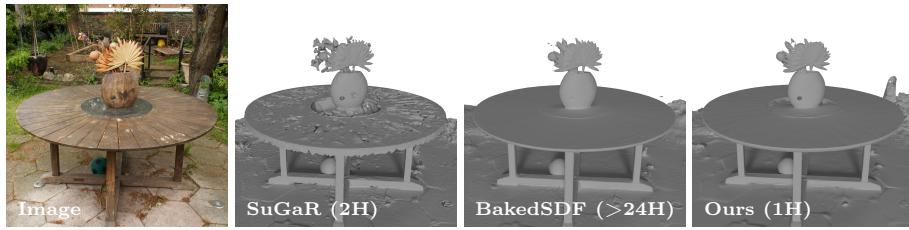


Fig. 1: Qualitative results on Mip-NeRF360 [1] dataset garden scene.

Abstract. Recently, 3D Gaussian Splatting (3DGS) has emerged as an efficient approach for accurately representing scenes. However, despite its superior novel view synthesis capabilities, extracting the geometry of the scene directly from the Gaussian properties remains a challenge, as those are optimized based on a photometric loss. While some concurrent models have tried adding geometric constraints during the Gaussian optimization process, they still produce noisy, unrealistic surfaces.

We propose a novel approach for bridging the gap between the noisy 3DGS representation and the smooth 3D mesh representation, by injecting real-world knowledge into the depth extraction process. Instead of extracting the geometry of the scene directly from the Gaussian properties, we instead extract the geometry through a pre-trained stereo-matching model. We render stereo-aligned pairs of images corresponding to the original training poses, feed the pairs into a stereo model to get a depth profile, and finally fuse all of the profiles together to get a single mesh. The resulting reconstruction is smoother, more accurate and shows more intricate details compared to other methods for surface reconstruction from Gaussian Splatting, while only requiring a small overhead on top of the fairly short 3DGS optimization process.

We performed extensive testing of the proposed method on in-the-wild scenes, obtained using a smartphone, showcasing its superior reconstruction abilities. Additionally, we tested the method on the Tanks and Temples and DTU benchmarks, achieving state-of-the-art results.

* Indicates equal contribution

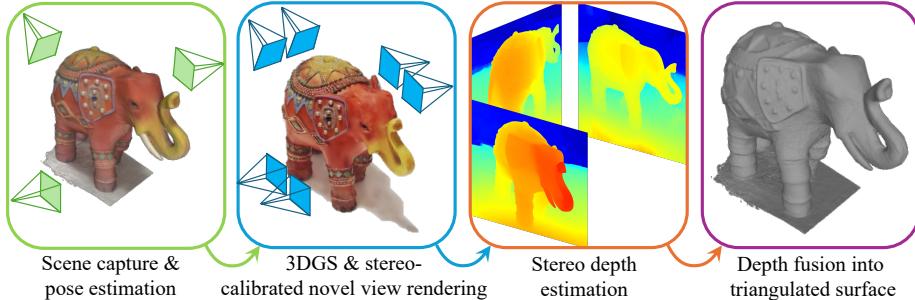


Fig. 2: The proposed pipeline for surface reconstruction. First, we represent the scene by applying a 3DGS model. We then use the 3DGS model to render stereo-aligned pairs of images corresponding to the original views. For each pair, using a shape from stereo algorithm, we reconstruct an RGB-D structure, which is then integrated from all views using TSDF [9] into a triangulated mesh of the scene.

1 Introduction

The Gaussian Splatting Model for radiance field rendering (3DGS) [19] has recently marked a significant leap forward in the realm of novel view synthesis, surpassing previous neural rendering methods in both speed and accuracy. By optimizing the distribution, size, color, and opacity of a cloud of Gaussian elements, and projecting, or splatting them onto virtual cameras, 3DGS is able to generate realistic images of complex scenes from novel viewing directions in real-time. However, direct reconstruction of surfaces from 3DGS involves significant challenges. The main problem is that the locations of the Gaussian elements in 3D space do not form a geometrically consistent surface, as those are optimized for best matching the input images when projected back onto their image planes. Consequently, reconstructing surfaces based on the centers of the Gaussians yields noisy and inaccurate results. Current state-of-the-art methods attempt to regularize the 3DGS optimization process by adding additional geometric constraints [14], flattening the Gaussian elements, [15], or extracting the geometry using opacity fields [51], but they still rely on the Gaussian locations and form noisy, unrealistic surfaces.

We propose an alternative approach for extracting depth from the optimized Gaussian point cloud, which does not rely on the noisy locations of the Gaussians. Instead, we take advantage of a powerful geometric regularizer, trained on real-world data - a pre-trained stereo matching model. Stereo matching models solve a correspondence problem on stereo-aligned pairs of images, from which accurate depth can be extracted. Our main observation is that through 3DGS rendering, we can artificially create stereo-aligned pairs of images corresponding to the original views, feed these pairs into a pre-trained stereo model, and fuse the resulting depths using the Truncated Signed Distance Function (TSDF)

algorithm [9]. The result is a smooth, geometrically consistent surface, that is extracted from the noisy 3DGS cloud using real-world regularization.

The proposed method reduces surface reconstruction time dramatically, taking only a small overhead on top of the 3DGS capturing of the scene, which is significantly faster compared to neural surface reconstruction methods. For instance, reconstruction of an in-the-wild scene taken by a standard smartphone camera requires less than five minutes of additional computation time after a 3DGS scene capture. Additionally, since we reconstruct the surface based on the 3DGS capture, it is straightforward to bind the mesh to the original model, as mentioned in [14, 37], for mesh-based manipulation of the Gaussian elements. Moreover, since our mesh is more accurate, it does not require any additional refinements [14].

We tested the proposed method on the Tanks and Temples (TnT) benchmark [21] as well as the DTU [16] benchmark, two commonly used 3D reconstruction datasets, and achieved state-of-the-art results. Additionally, we extensively tested our method on in-the-wild scenes captured with a smartphone, showing qualitative results of the proposed method’s reconstruction abilities. To summarize, our main contribution,

- We propose a novel method for fast and accurate in-the-wild surface reconstruction, by using a pre-trained stereo matching model as a geometric prior for extracting depth from a 3DGS model.

2 Related Efforts

2.1 Multi-View Stereo and Stereo Matching

Multi-View Stereo (MVS) is a fundamental geometry reconstruction method, where depth maps are extracted for each reference image based on correspondences with neighboring images. In the field of deep MVS methods, the pioneering work of MVSNet [46] introduced an end-to-end framework for MVS learning, which can be divided into three parts: 2D feature extraction, homography, and 3D cost volume with 3D convolutions. Latter methods presented an improvement to this scheme, by improving the 3D cost volume [26, 55], improving the architecture for 2D feature extraction [38], using a vision transformer (ViT) architecture for feature extraction [4], and improving 3D convolutions for more efficient computations using a coarse-to-fine method [13, 45]. To fuse the extracted depth maps into one point cloud or mesh there are two main methods: Fusibile [12], which has recently been generalized by [44], and TSDF [9]. MVS methods deeply rely on accurate camera poses for the calculation of epipolar lines, and on in-the-wild scenes they struggle to achieve high accuracy as in the controlled environment, since small errors in pose estimation result in a noisy reconstruction, as we show in our ablation study.

Deep stereo matching methods [5, 18, 22, 29, 34, 53] are related to deep MVS methods, however, since it is guaranteed that matching pixels between two images must lie in the same row, the cost volume layer works on the disparity

instead of the depth. Recent state-of-the-art stereo matching models, such as RAFT [25], IGEV [43], and DLNR [56], use iterative refinement using GRU or LSTM layers. Unlike MVS methods, stereo methods require only two images which are typically closer to each other compared to MVS, and share the same image plane, resulting in less occluded regions which are visible only in one of the views.

2.2 Neural Rendering for Novel View Synthesis and Surface Reconstruction

Novel view synthesis methods are trained on a set of images from a scene, and aim to render views of a scene from any given pose. The pioneering work of Neural Radiance Fields (NeRF) [27] presented a major leap forward in accuracy by incorporating importance sampling and positional encoding to enhance rendering quality. However, the use of relatively large Multi-Layer Perceptrons (MLP) to capture the scene resulted in long training times. Later, Mip-NeRF [1] improved the quality of the rendered view with a different sampling method, although training and rendering times remained long. InstantNGP [28] tackled the extended training times of previous efforts, by incorporating a hash grid and an occupancy grid with a small MLP.

Neural surface reconstruction methods [35, 39, 48, 49], in addition to accurately rendering novel views, are also capable of reconstructing the surface of the scene. IDR [49] trained an SDF represented by an MLP for both color and geometry reconstruction. Neus [39] reduced the geometric error by utilizing weighted volume rendering, and HF-Neus [40] enabled coarse-to-fine refinement for high-frequency detail reconstruction by decomposing the implicit SDF into a base function and a displacement function. RegSDF [54] used a point cloud obtained from shape-from-motion (SfM) as regularization, in addition to regularizing the curvature of the zero-level of the SDF function. NeuralWarp [10] suggested refining the geometry by regularizing image consistency between different views through warping based on implicit geometry. Neuralangelo [24], using a 3D hash encoded grid, enabled detailed reconstruction and achieved state-of-the-art results on leading benchmarks. However, reconstruction time of these methods can reach up to several days per scene. In the context of novel stereo views, a recent model [36] has managed to successfully perform unsupervised training of a stereo model using rendered stereo-aligned triplets from a neural scene reconstruction method, showcasing the possibility of novel view synthesis as a data factory.

2.3 Gaussian Splatting for Novel View Synthesis and Surface Reconstruction

Recently, a major leap forward was presented by 3DGS [19], a faster and more accurate method for scene capturing. 3DGS represents the scene as a point cloud of 3D Gaussians, where each Gaussian has the properties of opacity, rotation, scale, location, and spherical harmonics. The scaling of the Gaussians is anisotropic, which allows them to represent thin structures in the scene. The Gaussians are

initialized using an SfM algorithm [32, 33], which extracts the camera poses and provides an initial guess for the locations of the Gaussians. The capturing time is short compared to other methods based on MLPs, and it is capable of real-time rendering. Concurrent works on GS [6, 8, 11, 42, 50] have improved on the vanilla 3DGS in various ways, such as by reducing optimization time, increasing accuracy, reducing aliasing and removing the need for COLMAP [32, 33] poses.

As discussed earlier, the Gaussian locations in the vanilla 3DGS do not form a geometrically consistent surface. Recent methods try to manipulate the Gaussian elements to extract more accurate surfaces [7, 14, 15]. SuGaR [14], the pioneering method in surface reconstruction from Gaussian Splatting, added a regularization term for post-process optimization based on the opacity levels of the Gaussians, forcing the Gaussian element cloud to align with the surface. 2DGS [15] flattens the Gaussians into 2D elements, and GOF [51] extracts the surface by creating an opacity field from the Gaussians. However, since these methods utilize the location and opacity of the Gaussian elements, they reconstruct the surface with noisy undulations.

3 Method

We propose a novel pipeline for surface reconstruction from 3DGS, as illustrated in Fig. 2. In this section, we will explain in detail each step of the pipeline. We note that additionally, we can mask out specific objects by projecting segmentation masks from Segment Anything Model (SAM) [20] between consecutive images using depth maps. Additional information on masking is available in the supplementary material.

3.1 Scene Capture and Pose Estimation

We start with a video or images of a static scene as input. Following the vanilla 3DGS, we employ COLMAP [32, 33] for SfM to identify points of interest and deduce camera matrices from the provided images.

3.2 3DGS and Stereo-Aligned Novel View Rendering

The elements extracted from the previous stage are then fed into the 3DGS model to accurately represent the scene. For completeness, we will give a short formulation of the 3DGS process; In 3DGS, 3D Gaussian elements are defined in space by $G(\mathbf{x}) = \exp(-\frac{1}{2}(\mathbf{x} - \mathbf{x}_p)^\top \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \mathbf{x}_p))$, where \mathbf{x}_p is the center of the Gaussian, and $\boldsymbol{\Sigma}$ is its covariance matrix. During optimization, $\boldsymbol{\Sigma}$ is factorized into the rotation \mathbf{R} and scale \mathbf{S} matrices: $\boldsymbol{\Sigma} = \mathbf{RSS}^\top \mathbf{R}^\top$. When rendering, the Gaussians are projected onto the image plane: $\boldsymbol{\Sigma}' = \mathbf{JW}\boldsymbol{\Sigma}\mathbf{W}^\top\mathbf{J}^\top$, where \mathbf{W} is the view transformation, and \mathbf{J} is the Jacobian of the affine projective transformation onto the image plane. By removing the last row and column of $\boldsymbol{\Sigma}'$, we remain with 2D Gaussians in the image plane. To calculate the color of a pixel in the image plane, 3DGS employs alpha blending which applies weights

to the opacity from front to back, $C = \sum_{i \in N} \alpha_i c_i \prod_{j=1}^{i-1} (1 - \alpha_j)$, where α_i is the product of the the i^{th} 2D Gaussian with its opacity parameter, and c_i is the the directional appearance component. For more details, see the original 3DGS paper [19].

During the 3DGS process the Gaussians are optimized based on the photometric loss between the given source images and their corresponding rendered images. This creates a representation of the scene that allows rendering novel views which were not present in the original training data. It is important to note that the vanilla 3DGS relies on sufficient coverage of the scene, and in areas lacking sufficient coverage, noisy artifacts might appear, as seen in Fig. 8. Additionally, since the 3DGS is optimized based on the training images, staying close to a training image will result in a cleaner render.

Therefore, when generating novel stereo views of the scene, we input a sufficient amount of images that cover the region of interest. Additionally, we stay as close as possible to the original training poses, by choosing the left image of the stereo pair to be at the same pose R_L, T_L as a training image. Following this choice, the right pose with a horizontal baseline of b is formulated as follows: $R_R = R_L, T_R = T_L + (R_L \times [b, 0, 0])$. This ensures that the resulting left-right cameras are stereo-calibrated.

3.3 Stereo Depth Estimation

With the rendered stereo-aligned image pairs, we can essentially turn a scene captured from a single camera into a scene captured from a pair of stereo-calibrated cameras, using the novel view synthesis capabilities of 3DGS. We then apply a stereo matching algorithm to form depth profiles from every stereo pair. We have tested several stereo matching algorithms in the experimental section, and achieved the best qualitative and quantitative results with DLNR [56], a state-of-the-art neural stereo matching model, with the pre-trained Middlebury [31] weights. To further enhance the resulting reconstructions, we apply several masks to the output of the stereo model. The first mask is an occlusion mask, which is calculated by applying a threshold on the difference between the left-to-right and right-to-left disparities of the same pair of images. This masks out parts of the scene that were only visible in one of the cameras, and therefore the stereo model's output in these areas is unreliable. We justify the use of this mask by the fact that the occluded areas will be filled in from adjacent stereo views. An example of an occlusion mask can be seen in Fig. 3, and we added an experiment in the supplementary material which demonstrates the effectiveness of the occlusion mask.

The second mask is applied based on the depth of the stereo output. The relationship between stereo matching errors can be described as $\epsilon(Z) \approx \frac{\epsilon(d)}{f_x \cdot B} Z^2$ [3], where $\epsilon(d)$ represents the disparity output error, Z is the ground-truth depth, $\epsilon(Z)$ is the error of the depth estimate, f_x denotes the camera's horizontal focal length, and B is the baseline. Conversely, the disparity between matching pixels in two images of an object that is positioned at a short distance from the

cameras can exceed the maximum disparity limit produced by stereo matching algorithms. Thus, estimating the depth of an object that is too close to the camera can result in an error due to the limitation of the matching algorithms, and estimating the depth of an object that is distant results in a quadratic error. Therefore, we consider depth in the range $4B \leq Z \leq 20B$. This approach enhances the overall accuracy and reliability of the depth estimation process, ensuring more consistent geometric reconstructions. With the above considerations taken into account, we now have two contradicting factors when setting the horizontal baseline of the stereo pair; On the one hand, a larger baseline allows for a wider “sweet spot” for the stereo model. On the other hand, the 3DGS limits how far we can stray from the original training images without producing noisy renders. In the experimental section, we tested different baselines and found that a horizontal baseline of 7% of the scene radius, or 3.5% of the scene diameter, which allows for a “sweet spot” in the range of 14% to 70% of the scene diameter, provides the best results.

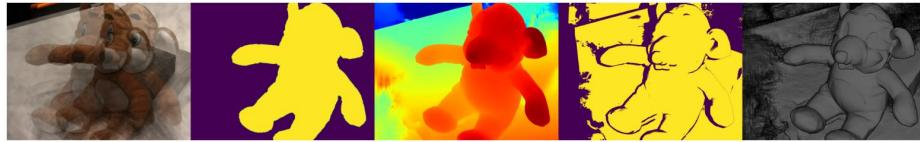


Fig. 3: Example of our method’s output on DTU [16] scan105. From left to right: The rendered left and right images, segmentation mask, left-right disparity, occlusion mask, and shading - depth gradient.

3.4 Depth Fusion into Triangulated Surface

To further enhance geometric consistency and smooth out any noise and errors which might have originated from the individual depth profiles, we aggregate all of the extracted depths using the Truncated Signed Distance Function (TSDF) algorithm [9], followed by the Marching-Cubes meshing algorithm [41].

4 Experiments and Results

We present experiments which demonstrate that our method is able to accurately reconstruct surfaces in a more geometrically consistent way than other 3DGS-based or MVS approaches, as well as achieve comparable performance to neural reconstruction methods while taking significantly less time to run. For quantitative results, we tested our method on the Tanks and Temples [21] and DTU [16] datasets, and compared our results to various neural and 3DGS-based reconstruction methods. We also compared between different versions of our model to justify our design choices. Additionally, we show qualitative reconstruction results from Mip-NeRF360 [1], demonstrating that our method achieves

comparable visual quality to neural reconstruction methods, and on in-the-wild videos taken from smartphones, we show our superiority in terms of geometric consistency and smoothness when compared to SuGaR [14]. Finally, we perform an ablation study on the MobileBrick [23] dataset, which validates the contribution of novel-view image generation and stereo, by replacing two different points in our pipeline with a deep MVS model. We note that in the MobileBrick dataset the camera poses are manually refined, and are shown to be more accurate than COLMAP [32,33] poses for reconstruction [23]. The comparison we present thus favors MVS models in that respect.

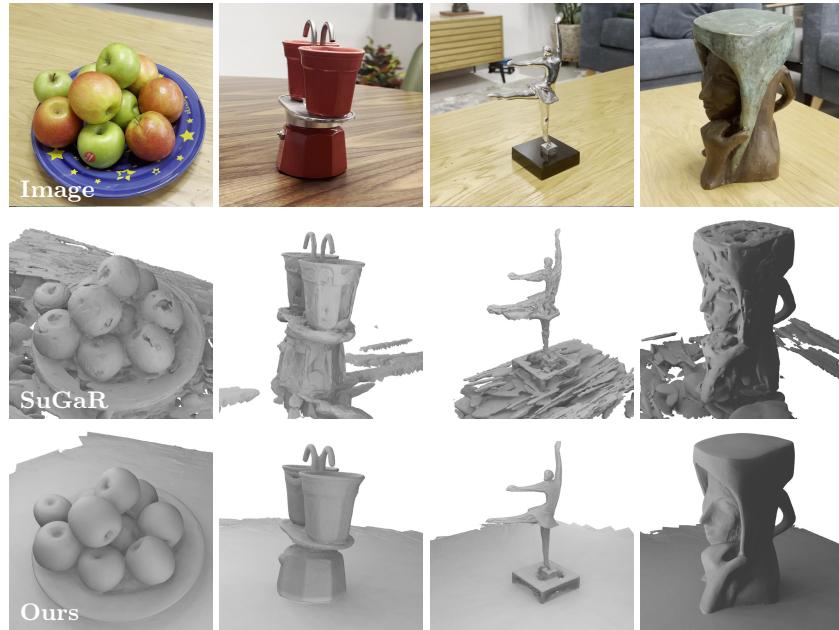


Fig. 4: Qualitative comparison of mesh reconstruction from in-the-wild videos between our method and SuGaR [14].

4.1 Datasets

DTU [16]. This dataset is an MVS dataset, containing scans of small objects, as well as accurate camera poses and 3D point clouds. We use the dataset and evaluation code from [15], which calculates the Chamfer Distance (CD) between the reconstructed and ground-truth point clouds.

Tanks and Temples (TnT) [21]. This dataset contains videos of large objects such as vehicles, buildings and statues. These objects are scanned with a laser scanner for an accurate ground-truth 3D point cloud. As the videos and the

laser scanned objects are difficult to align [30], for evaluation we use the official TnT [21] evaluation alignment method. It first aligns the point clouds using ICP [2], and then calculates the precision, recall, and F1 score.

Mip-NeRF360 [1]. This dataset contains scenes taken from a 360 degree view, with emphasis on minimizing photometric variations through controlled capture conditions. Since there is no ground-truth in terms of surface reconstruction, we leave this as a qualitative comparison only.

MobileBrick [23]. This dataset contains videos of LEGO models, with corresponding 3D ground-truth meshes created from the LEGO 3D model. The poses are manually refined and are more accurate than the COLMAP ones [32, 33]. This dataset is challenging since most of the videos in the test set are taken from a top view of the model, thus, creating occlusions and leaving areas in the model with little visibility. We use the official evaluation code.

In-the-wild videos. For reconstruction of in-the-wild objects, our method presents a favorable balance between accuracy and computation time. To validate this claim, we captured scenes containing various objects such as plants, sculptures, figures and everyday items, with intricate geometries and textures, and reconstructed their surface. Each video contains one or two cycles of moving around the object, depending on the object’s size, without any measures to maintain a persistent radius or pose of the camera, and without any control of the lighting in the environment. Since these objects are filmed only with a smartphone camera, there is no ground-truth reconstruction for these objects.

4.2 Baselines

Gaussian Splatting-based methods. For the DTU [16] dataset, We compare our method with the vanilla 3DGS [19] and SuGaR [14], as well as additional state-of-the-art methods, namely 2DGS [15] and Gaussian Opacity Fields (GOF) [51]. For the TNT [21] and Mip-NeRF360 [1] datasets, as well as for in-the-wild scenes, we compare with SuGaR [14].

Neural rendering methods. For the DTU [16] dataset, we compare our method with Neuralangelo [24], VolSDF [47], and NeuS [39]. For the TnT [21] dataset, we compare with Neuralangelo [24], NeuralWarp [10] and NeuS [39]. For the Mip-NeRF360 [1] dataset, we compare with BakedSDF [48].

Deep MVS. for in-the-wild scenes and the MobileBrick [23] dataset, we compare with MVSformer [4], a state-of-the art deep MVS network.

4.3 Results

DTU [16]. Tab. 1 presents quantitative results on the DTU [16] dataset. We ran the 3DGS step of our method for 30000 iterations, but as we show in the supplementary material, we can achieve nearly identical results with only 7000 iterations, reaching a mean Chamfer Distance of 0.70 with only $\sim 12m$ of total runtime per scan. We used the same TSDF as in 2DGS [15] and GOF [51], which is based on the Open3D implementation [9], for a fair comparison. We

Table 1: Quantitative results on the DTU [16] dataset, comparing our method with state-of-the-art neural and Gaussian Splatting-based methods. Chamfer distance - lower is better. Red-1st, Orange-2nd, Yellow-3rd. Table adapted from [51].

	24	37	40	55	63	65	69	83	97	105	106	110	114	118	122	Mean	Time
NeRF [27]	1.90	1.60	1.85	0.58	2.28	1.27	1.47	1.67	2.05	1.07	0.88	2.53	1.06	1.15	0.96	1.49	> 12h
VolSDF [47]	1.14	1.26	0.81	0.49	1.25	0.70	0.72	1.29	1.18	0.70	0.66	1.08	0.42	0.61	0.55	0.86	> 12h
NeuS [39]	1.00	1.37	0.93	0.43	1.10	0.65	0.57	1.48	1.09	0.83	0.52	1.20	0.35	0.49	0.54	0.84	> 12h
Neuralangelo [24]	0.37	0.72	0.35	0.35	0.87	0.54	0.53	1.29	0.97	0.73	0.47	0.74	0.32	0.41	0.43	0.61	> 12h
Gaussian Splatting																	
3DGSS [19]	2.14	1.53	2.08	1.68	3.49	2.21	1.43	2.07	2.22	1.75	1.79	2.55	1.53	1.52	1.50	1.96	11.2m
SuGaR [14]	1.47	1.33	1.13	0.61	2.25	1.71	1.15	1.63	1.62	1.07	0.79	2.45	0.98	0.88	0.79	1.33	~ 1h
2DGSS [15]	0.48	0.91	0.39	0.39	1.01	0.83	0.81	1.36	1.27	0.76	0.70	1.40	0.40	0.76	0.52	0.80	18.8m
GOF [51]	0.50	0.82	0.37	0.37	1.12	0.74	0.73	1.18	1.29	0.68	0.77	0.90	0.42	0.66	0.49	0.74	30m
Ours - DLNR Baseline 3.5%	0.61	0.85	0.64	0.39	0.96	1.25	0.80	1.52	1.10	0.68	0.59	0.93	0.45	0.60	0.54	0.79	~20m
Ours - DLNR Baseline 10.5%	0.69	0.81	0.95	0.51	0.82	1.06	0.72	1.18	0.93	0.61	0.54	0.66	0.37	0.54	0.50	0.73	~20m
Ours - RAFT Baseline 7%	0.59	0.81	0.68	0.40	0.83	1.15	0.73	1.35	1.05	0.62	0.53	0.80	0.390	0.55	0.49	0.73	~20m
Ours - DLNR Baseline 7%	0.59	0.79	0.70	0.38	0.78	1.00	0.69	1.25	0.96	0.59	0.68	0.37	0.50	0.46	0.68	~20m	

apply the mask supplied with the dataset before inputting the depths into the TSDF algorithm. The table is adapted from GOF [51] for consistency. Within the splatting-based methods, our method achieves the best score, while maintaining a similar runtime. Additionally, when compared to the neural methods, which take more than 12 hours to reconstruct a single scene, our method surpasses some of the methods, and is comparable with Neuralangelo [24], the state-of-the-art method. Additionally, we test our method with RAFT [25] as the stereo model with the RVC weights [17] and with DLNR [56] as the stereo model with the Middlebury weights [31], achieving better results with DLNR. Finally, we compare between three different horizontal baselines: 3.5%, 7% and 10.5% of the scene radius. We achieve the best results with 7%, noting that increasing or decreasing the horizontal baseline has a negative effect on the results. Fig. 3 shows an example of the intermediate representations of our method on one of DTU [16] scan105, and the full set of reconstructed meshes is available in the supplementary material.

Table 2: Quantitative results on the Tanks and Temples [21] benchmark. F1 score - higher is better.

	Barn	Caterpillar	Ignatius	Truck	Mean F1 ↑	Runtime
NeuralWarp [10]	0.22	0.18	0.02	0.35	0.19	
NeuS [39]	0.29	0.29	0.83	0.45	0.47	~16h-48h
Neuralangelo [24]	0.70	0.36	0.89	0.48	0.61	
SuGaR [14]	0.01 (0.08)	0.02 (0.09)	0.06 (0.34)	0.05 (0.17)	0.04 (0.17)	~2h
Ours	0.21 (0.22)	0.17 (0.12)	0.64 (0.68)	0.46 (0.40)	0.37 (0.36)	~1h

Tanks and Temples [21]. Tab. 2 presents a summary of the reconstruction results on the TnT [21] benchmark. Since SuGaR [14] yields a sparse mesh, and thus its recall drops significantly, we include a precision metric that is unaffected by mesh sparsity. However, it is important to note that this metric does not account for missing parts in the reconstruction. The results show that our

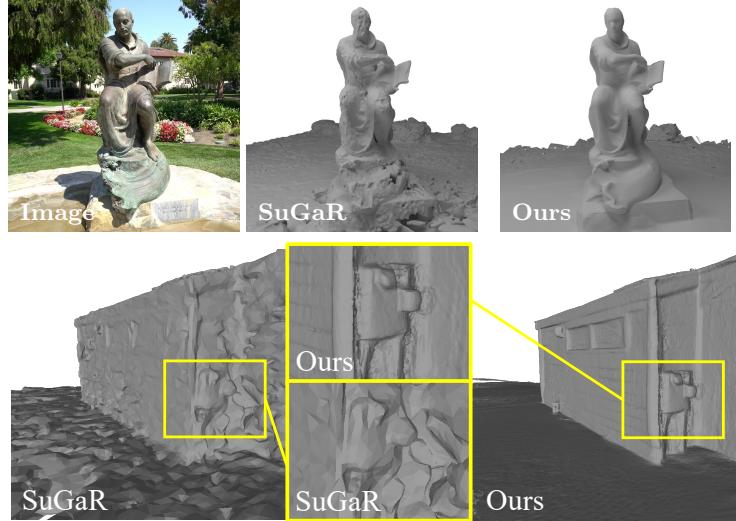


Fig. 5: Qualitative results on Tanks and Temples [21]. Top row: Ignatius scene, compared to SuGaR [14]. Bottom row: Barn scene, compared to SuGaR.

method outperforms SuGaR [14] in both F1 and precision. Additionally, it is evident from Fig. 5 that our method is able to reconstruct fine details such as in the Barn scene. Moreover, our method has a significant advantage in terms of processing time, requiring less than 60 minutes of total computation time per TnT [21] scene, compared to the 16-48 hours needed by neural reconstruction methods. The reason for our relatively longer computation times for TnT [21] is since each scene containing hundreds of frames, compared to a typical in-the-wild scene which contains less than 100 frames. It is important to note that the TnT [21] dataset predominantly features large scenes, whereas our method is based on 3DGS reconstruction that is designed for accurate reconstruction of smaller ones, and TSDF which is better suited for reconstruction of specific objects. This is particularly evident in the case of the Ignatius and Truck scenes, relatively small scenes, where our method performed on-par with the neural reconstruction methods.

Mip-NeRF360 [1]. As illustrated in Fig. 1 and Fig. 6, we present a qualitative analysis of scenes from the Mip-NeRF360 [1] dataset. This comparison reveals that our approach surpasses SuGaR [14] in terms of reconstruction quality and presents on-par results with BakedSDF [48]. Notably, our method excels in reconstructing fine details; for instance, even the small grooves in the garden scene’s table are evident in the reconstruction, and there are intricate details in the objects on the countertop scene. Furthermore, while BakedSDF [48] requires 48 hours for training, our method achieves comparable results in less than an hour. Compared to SuGaR [14], our method generates smoother and more realistic surfaces, especially in reflective areas; we note that our countertop is smooth



Fig. 6: Qualitative comparison on Mip-NeRF360 [1] dataset with BakedSDF [48] and SuGaR [14].

and flat, while SuGaR’s countertop has many bumps in areas with glare. This is likely due to our model’s use of a small baseline for stereo matching, where the reconstruction distortion is relatively small, and additionally, due to our model integrating the reconstructed patches from various viewing directions, which further reduces potential distortions.

In-the-wild comparison. Our comprehensive in-the-wild comparisons demonstrate the superior performance of our method across various scenes, as illustrated in Fig. 4, with additional results provided in the supplementary material. Our method surpasses SuGaR [14] in extracting accurate and noise-free meshes from 3DGS.

4.4 Ablation Study

Our main contribution is the use of a pre-trained stereo model to extract depth from a 3DGS scene using novel stereo views. To strengthen our claim of the benefit of using novel stereo views, we perform two ablations, which include replacing steps of our pipeline with deep MVS methods.

MVS on original images. Our method extracts depth from each original pose, by creating novel stereo-aligned views from that pose and applying a pre-trained stereo matching model. One obvious comparison would be to take each original pose and extract the depth using a pre-trained deep MVS model which will take as input the original training set of the scene. In the first ablation, we run a pre-trained deep MVS model on the original images, and fuse the resulting depths using TSDF [9].

MVS on rendered images. Applying 3DGS to the scene and re-rendering the images from the original poses can reduce distortion and camera noise, which may enhance the quality of the reconstruction regardless of the novel stereo views.

In the second ablation, we run a pre-trained deep MVS model on the rendered images from the original poses, which are the left image of each stereo-aligned novel view, and fuse the resulting depths using TSDF [52].

Evaluation. We evaluate on the MobileBrick [23] test set, and compare our method against MVSFormer [4], a state-of-the-art deep MVS model. To ensure a fair comparison, we use TSDF [9] as the fusion method both for our method and the deep MVS model.

Results. Tab. 3 shows the mean accuracy, recall, F1 and Chamfer Distance of Our method, compared to MVSFormer [4] with the original and rendered images as input. Fig. 7 shows a qualitative comparison on one of the scans in the MobileBrick [23] dataset, with the rest of the scans, as well as additional examples from in-the-wild scenes, available in the supplementary material. Qualitative comparison shows that applying deep MVS directly on the original images results in a reconstruction filled with holes. Applying MVS on the rendered images slightly improves the quality of the reconstruction, however, our method still produces a significantly smoother reconstruction. Quantitative comparison confirms that inputting rendered images to the same MVS model results in a smoother reconstructed surface, as evident by the higher recall, with a slight trade-off in accuracy. Overall, our method performs better, as evident by the higher recall and F1 and lower Chamfer distance, even though the manually refined poses given by the MobileBrick [23] dataset should give an advantage to MVSFormer [4].

Table 3: Ablation study results on MobileBrick [23] dataset. We compare our method against MVSFormer [4] with two types of inputs: the original images with the original refined poses, and the rendered images with the same poses.

Method	2.5mm Radius			5mm Radius			Chamfer Distance (mm) ↓
	Acc ↑	Recall ↑	F1 ↑	Acc ↑	Recall ↑	F1 ↑	
MVSFormer [4]	80.77	55.02	64.60	96.33	71.32	81.14	9.11
MVSFormer + Rendered	80.16	59.92	68.04	96.84	77.50	85.59	7.10
Ours	68.77	69.27	68.94	89.46	87.37	88.28	4.94

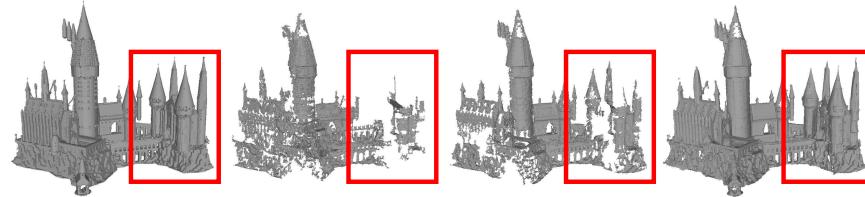


Fig. 7: Example from MobileBrick [23] dataset, on the castle scene. From left to right: the ground truth mesh, reconstruction of MVSFormer [4] with original images, reconstruction of MVSFormer with rendered images, and our reconstruction.

5 Limitations



Fig. 8: Examples of limitations of our method. On the left, we show a rendered image from the Caterpillar scene from TnT [21] dataset, highlighting an area with “floater” Gaussians. On the right, the Truck scene from TnT [21] dataset, highlighting the missing windshield.

Our pipeline consists of 3DGS, depth extraction via stereo, and TSDF fusion. Each of these steps exhibits limitations that can impact the final reconstruction: 3DGS can produce noisy “floater” Gaussians in areas which aren’t sufficiently covered in the original training images, as can be seen in the right side of Fig. 8. Additionally, stereo matching models are known to struggle with transparent surfaces, as can be seen in the left side of Fig. 8. Finally, TSDF fusion does not scale well for larger scenes, such as the Meetingroom and Courthouse scenes from TnT [21]. Swapping the 3DGS and stereo with future versions which will have improved accuracy and robustness, as well as adding fusion methods better suited for larger scenes, should help mitigate the effect of these limitations.

6 Conclusion

We introduce a novel approach for bridging the gap between noisy Gaussian point clouds and smooth surfaces in 3D. Instead of applying geometric optimizations directly on the Gaussians and extracting the depth using their locations, we use a pre-trained stereo model as a geometric prior with real-world knowledge to extract the depth. While this approach preserves the inherent properties of the 3DGS representation, it also enhances the accuracy and fidelity of the reconstructed surfaces. Our experimental results on DTU [16], Tanks and Temples [21], Mip-NeRF360 [1], MobileBrick [23] and real-world scenes captured using smartphones - demonstrate the superiority of our method over the current state-of-the-art methods for surface reconstruction from Gaussian splatting models, offering both improved accuracy and significantly shorter computation times compared to neural methods.

References

1. Barron, J.T., Mildenhall, B., Tancik, M., Hedman, P., Martin-Brualla, R., Srinivasan, P.P.: Mip-NeRF: A multiscale representation for anti-aliasing neural radiance fields. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 5855–5864 (2021)
2. Besl, P.J., McKay, N.D.: Method for registration of 3-D shapes. In: Sensor fusion IV: control paradigms and data structures. vol. 1611, pp. 586–606. Spie (1992)
3. Bracha, A., Rotstein, N., Bensaïd, D., Slossberg, R., Kimmel, R.: Depth refinement for improved stereo reconstruction. arXiv preprint arXiv:2112.08070 (2021)
4. Cao, C., Ren, X., Fu, Y.: Mvsformer: Learning robust image representations via transformers and temperature-based depth for multi-view stereo. arXiv preprint arXiv:2208.02541 (2022)
5. Chang, J.R., Chen, Y.S.: Pyramid stereo matching network. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 5410–5418 (2018)
6. Chen, G., Wang, W.: A survey on 3D Gaussian splatting. arXiv preprint arXiv:2401.03890 (2024)
7. Chen, H., Li, C., Lee, G.H.: Neusg: Neural implicit surface reconstruction with 3D Gaussian splatting guidance. arXiv preprint arXiv:2312.00846 (2023)
8. Cheng, K., Long, X., Yang, K., Yao, Y., Yin, W., Ma, Y., Wang, W., Chen, X.: GaussianPro: 3D Gaussian splatting with progressive propagation. In: Forty-first International Conference on Machine Learning (2024)
9. Curless, B., Levoy, M.: A volumetric method for building complex models from range images. In: Proceedings of the 23rd annual conference on Computer graphics and interactive techniques. pp. 303–312 (1996)
10. Darmon, F., Basclle, B., Devaux, J.C., Monasse, P., Aubry, M.: Improving neural implicit surfaces geometry with patch warping. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 6260–6269 (2022)
11. Fu, Y., Liu, S., Kulkarni, A., Kautz, J., Efros, A.A., Wang, X.: Colmap-free 3D Gaussian splatting. arXiv preprint arXiv:2312.07504 (2023)
12. Galliani, S., Lasinger, K., Schindler, K.: Massively parallel multiview stereopsis by surface normal diffusion (June 2015)
13. Gu, X., Fan, Z., Zhu, S., Dai, Z., Tan, F., Tan, P.: Cascade cost volume for high-resolution multi-view stereo and stereo matching. In: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition. pp. 2495–2504 (2020)
14. Guédon, A., Lepetit, V.: SuGaR: Surface-aligned Gaussian splatting for efficient 3D mesh reconstruction and high-quality mesh rendering. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 5354–5363 (2024)
15. Huang, B., Yu, Z., Chen, A., Geiger, A., Gao, S.: 2D Gaussian splatting for geometrically accurate radiance fields. In: SIGGRAPH 2024 Conference Papers. Association for Computing Machinery (2024)
16. Jensen, R., Dahl, A., Vogiatzis, G., Tola, E., Aanæs, H.: Large scale multi-view stereopsis evaluation. In: 2014 IEEE Conference on Computer Vision and Pattern Recognition. pp. 406–413. IEEE (2014)
17. Jiang, H., Xu, R., Jiang, W.: An improved RAFTstereo trained with a mixed dataset for the robust vision challenge 2022. arXiv preprint arXiv:2210.12785 (2022)

18. Kendall, A., Martirosyan, H., Dasgupta, S., Henry, P., Kennedy, R., Bachrach, A., Bry, A.: End-to-end learning of geometry and context for deep stereo regression. In: Proceedings of the IEEE international conference on computer vision. pp. 66–75 (2017)
19. Kerbl, B., Kopanas, G., Leimkühler, T., Drettakis, G.: 3D Gaussian splatting for real-time radiance field rendering. ACM Transactions on Graphics **42**(4) (2023)
20. Kirillov, A., Mintun, E., Ravi, N., Mao, H., Rolland, C., Gustafson, L., Xiao, T., Whitehead, S., Berg, A.C., Lo, W.Y., et al.: Segment anything. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 4015–4026 (2023)
21. Knapitsch, A., Park, J., Zhou, Q.Y., Koltun, V.: Tanks and Temples: Benchmarking large-scale scene reconstruction. ACM Transactions on Graphics **36**(4) (2017)
22. Laga, H., Jospin, L.V., Boussaid, F., Bennamoun, M.: A survey on deep learning techniques for stereo-based depth estimation. IEEE transactions on pattern analysis and machine intelligence **44**(4), 1738–1764 (2020)
23. Li, K., Bian, J.W., Castle, R., Torr, P.H., Prisacariu, V.A.: Mobilebrick: Building LEGO for 3D reconstruction on mobile devices. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 4892–4901 (2023)
24. Li, Z., Müller, T., Evans, A., Taylor, R.H., Underath, M., Liu, M.Y., Lin, C.H.: Neuralangelo: High-fidelity neural surface reconstruction. In: IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (2023)
25. Lipson, L., Teed, Z., Deng, J.: Raft-stereo: Multilevel recurrent field transforms for stereo matching. In: 2021 International Conference on 3D Vision (3DV). pp. 218–227. IEEE (2021)
26. Ma, Z., Teed, Z., Deng, J.: Multiview stereo with cascaded epipolar raft. In: European Conference on Computer Vision. pp. 734–750. Springer (2022)
27. Mildenhall, B., Srinivasan, P.P., Tancik, M., Barron, J.T., Ramamoorthi, R., Ng, R.: NeRF: Representing scenes as neural radiance fields for view synthesis. Communications of the ACM **65**(1), 99–106 (2021)
28. Müller, T., Evans, A., Schied, C., Keller, A.: Instant neural graphics primitives with a multiresolution hash encoding. ACM Trans. Graph. **41**(4), 102:1–102:15 (Jul 2022)
29. Poggi, M., Tosi, F., Batsos, K., Mordohai, P., Mattoccia, S.: On the synergies between machine learning and binocular stereo for depth estimation from images: a survey. IEEE Transactions on Pattern Analysis and Machine Intelligence **44**(9), 5314–5334 (2021)
30. Rotstein, N., Bracha, A., Kimmel, R.: Multimodal colored point cloud to image alignment. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 6656–6666 (2022)
31. Scharstein, D., Szeliski, R.: A taxonomy and evaluation of dense two-frame stereo correspondence algorithms. International journal of computer vision **47**, 7–42 (2002)
32. Schönberger, J.L., Frahm, J.M.: Structure-from-motion revisited. In: Conference on Computer Vision and Pattern Recognition (CVPR) (2016)
33. Schönberger, J.L., Zheng, E., Pollefeys, M., Frahm, J.M.: Pixelwise view selection for unstructured multi-view stereo. In: European Conference on Computer Vision (ECCV) (2016)
34. Shen, Z., Dai, Y., Rao, Z.: Cfnet: Cascade and fused cost volume for robust stereo matching. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 13906–13915 (2021)

35. Sitzmann, V., Martel, J., Bergman, A., Lindell, D., Wetzstein, G.: Implicit neural representations with periodic activation functions. *Advances in neural information processing systems* **33**, 7462–7473 (2020)
36. Tosi, F., Tonioni, A., De Gregorio, D., Poggi, M.: Nerf-supervised deep stereo. In: *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. pp. 855–866 (2023)
37. Waczyńska, J., Borycki, P., Tadeja, S., Tabor, J., Spurek, P.: GaMeS: Mesh-based adapting and modification of Gaussian splatting. *arXiv preprint arXiv:2402.01459* (2024)
38. Wang, F., Galliani, S., Vogel, C., Speciale, P., Pollefeys, M.: Patchmatchnet: Learned multi-view patchmatch stereo. In: *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. pp. 14194–14203 (2021)
39. Wang, P., Liu, L., Liu, Y., Theobalt, C., Komura, T., Wang, W.: Neus: Learning neural implicit surfaces by volume rendering for multi-view reconstruction (2021)
40. Wang, Y., Skorokhodov, I., Wonka, P.: HF-NeuS: Improved surface reconstruction using high-frequency details. *Advances in Neural Information Processing Systems* **35**, 1966–1978 (2022)
41. WE, L.: Marching cubes: A high resolution 3D surface construction algorithm. *Computer graphics* **21**(1), 7–12 (1987)
42. Wu, T., Yuan, Y.J., Zhang, L.X., Yang, J., Cao, Y.P., Yan, L.Q., Gao, L.: Recent advances in 3D Gaussian splatting. *Computational Visual Media* pp. 1–30 (2024)
43. Xu, G., Wang, X., Ding, X., Yang, X.: Iterative geometry encoding volume for stereo matching. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. pp. 21919–21928 (2023)
44. Yan, J., Wei, Z., Yi, H., Ding, M., Zhang, R., Chen, Y., Wang, G., Tai, Y.W.: Dense hybrid recurrent multi-view stereo net with dynamic consistency checking. In: *European conference on computer vision*. pp. 674–689. Springer (2020)
45. Yang, J., Mao, W., Alvarez, J.M., Liu, M.: Cost volume pyramid based depth inference for multi-view stereo. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. pp. 4877–4886 (2020)
46. Yao, Y., Luo, Z., Li, S., Fang, T., Quan, L.: Mvsnet: Depth inference for unstructured multi-view stereo. In: *Proceedings of the European conference on computer vision (ECCV)*. pp. 767–783 (2018)
47. Yariv, L., Gu, J., Kasten, Y., Lipman, Y.: Volume rendering of neural implicit surfaces. *Advances in Neural Information Processing Systems* **34**, 4805–4815 (2021)
48. Yariv, L., Hedman, P., Reiser, C., Verbin, D., Srinivasan, P.P., Szeliski, R., Barron, J.T., Mildenhall, B.: BakedSDF: Meshing neural SDFs for real-time view synthesis. In: *ACM SIGGRAPH 2023 Conference Proceedings*. pp. 1–9 (2023)
49. Yariv, L., Kasten, Y., Moran, D., Galun, M., Atzmon, M., Ronen, B., Lipman, Y.: Multiview neural surface reconstruction by disentangling geometry and appearance. *Advances in Neural Information Processing Systems* **33**, 2492–2502 (2020)
50. Yu, Z., Chen, A., Huang, B., Sattler, T., Geiger, A.: Mip-splatting: Alias-free 3D Gaussian splatting. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. pp. 19447–19456 (2024)
51. Yu, Z., Sattler, T., Geiger, A.: Gaussian Opacity Fields: Efficient high-quality compact surface reconstruction in unbounded scenes. *arXiv preprint arXiv:2404.10772* (2024)
52. Zeng, A., Song, S., Nießner, M., Fisher, M., Xiao, J., Funkhouser, T.: 3dmatch: Learning local geometric descriptors from rgb-d reconstructions. In: *CVPR* (2017)

53. Zhang, F., Prisacariu, V., Yang, R., Torr, P.H.: Ga-net: Guided aggregation net for end-to-end stereo matching. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 185–194 (2019)
54. Zhang, J., Yao, Y., Li, S., Fang, T., McKinnon, D., Tsin, Y., Quan, L.: Critical regularizations for neural surface reconstruction in the wild. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 6270–6279 (2022)
55. Zhang, J., Yao, Y., Li, S., Luo, Z., Fang, T.: Visibility-aware multi-view stereo network. arXiv preprint arXiv:2008.07928 (2020)
56. Zhao, H., Zhou, H., Zhang, Y., Chen, J., Yang, Y., Zhao, Y.: High-frequency stereo matching network. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 1327–1336 (2023)

Surface Reconstruction from Gaussian Splatting via Novel Stereo Views

Supplementary Material

A Additional Ablations

Table A.1: Additional ablations on the DTU [16] dataset, regarding the number of 3DGS iterations and the use of occlusion masks. Chamfer distance - lower is better. Red-1st, Orange-2nd, Yellow-3rd.

	24	37	40	55	63	65	69	83	97	105	106	110	114	118	122	Mean	Time
Ours - 7000 w/o occlusion	0.75	0.81	0.67	0.40	0.91	0.98	0.82	1.46	1.03	0.76	0.79	0.89	0.47	0.58	0.46	0.78	~12m
Ours - 7000 w/ occlusion	0.62	0.78	0.70	0.37	0.80	1.01	0.71	1.28	0.98	0.60	0.55	0.76	0.35	0.53	0.46	0.70	~12m
Ours - 30000 w/o occlusion	0.71	0.82	0.66	0.42	0.88	0.98	0.78	1.46	1.03	0.75	0.72	0.82	0.48	0.57	0.46	0.77	~20m
Ours - 30000 w/ occlusion	0.59	0.79	0.70	0.38	0.78	1.00	0.69	1.25	0.96	0.59	0.50	0.68	0.37	0.50	0.46	0.68	~20m
Ours - 60000 w/o occlusion	0.72	0.82	0.66	0.43	0.86	0.96	0.77	1.42	1.02	0.75	0.74	0.82	0.49	0.56	0.45	0.76	~30m
Ours - 60000 w/ occlusion	0.60	0.80	0.71	0.38	0.75	1.00	0.68	1.25	0.97	0.60	0.53	0.68	0.37	0.50	0.46	0.69	~30m

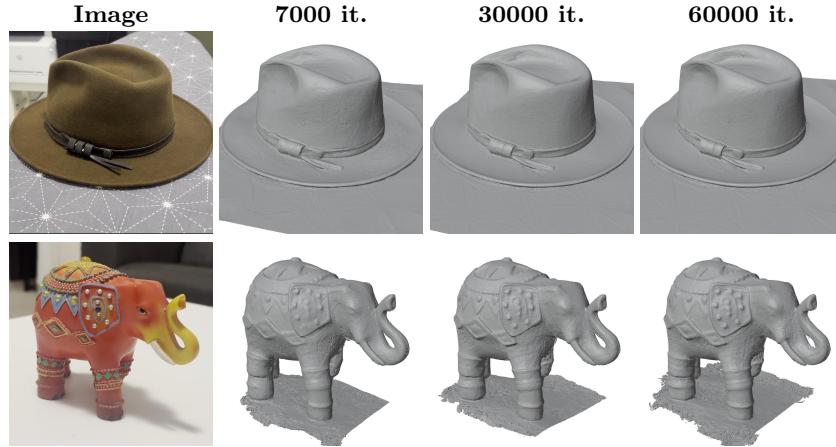


Fig. A.1: From left to right: One of the input images, and the reconstructed mesh using our method, after 7000, 30000 and 60000 3DGS iterations.

A.1 Number of GS Iterations

We claim that our method runs significantly faster compared to other reconstruction methods, and that the bottleneck of our runtime is the 3DGS optimization

time. For a typical in-the-wild scene containing a central object, with roughly 80 images, 3DGS optimization takes 5-30 minutes on an Nvidia L40 GPU, depending on the number of iterations: 5-10 minutes for 7000 iterations, 10-20 minutes for 30000 iterations, and 20-30 minutes for 60000 iterations. In Fig. A.1, we show that there are no noticeable differences in reconstructions of in-the-wild objects after 7000, 30000 and 60000 iterations. Additionally, in Tab. A.1, we evaluate our method on the DTU [16] dataset after 7000, 30000 and 60000 iterations, showing that even with only 7000 iterations we achieve nearly identical results, demonstrating the power of our real-world regularization using a pre-trained stereo model.

A.2 Necessity of Occlusion Mask

Tab. A.1 additionally validates the contribution of the occlusion mask that avoids using the depth information in areas which weren’t visible in both the left and right views, showing an improvement in metrics regardless of the number of 3DGS iterations.

B Object Segmentation using Depth and SAM

In some cases we want to extract only the surface of a given object in the scene. Other methods such as Segment-Any-Gaussians require additional training after the 3DGS, and are more suited for scenes with multiple objects rather than 360 object-centric scenes [?]. Thus, we choose to segment each image. The naive approach to object segmentation in a scene would involve independently segmenting every image, a method that can be labor-intensive. Instead, we employ a technique that leverages the power of Segment Anything Model (SAM) [20] in conjunction with depth information and geometric transformations. Initially, we annotate the first image of the scene using SAM to obtain a precise object mask. This initial segmentation acts as a foundation for tracking and segmenting the object across subsequent images. By performing this process after obtaining the depth of the identified object, we can project its mask onto the next image in the sequence using the extrinsic camera parameters. To accommodate for potential errors in SAM, we dilate the projected mask in the new image. Then, utilizing farthest point sampling, we select points that represent the extremities of the object within this dilated mask. These points are used as the seed for a new SAM annotation for the next image. This process is applied iteratively to each subsequent image in the series, allowing for dynamic and precise object segmentation throughout the scene. Fig. B.1 shows examples of our method’s segmentation abilities on in-the-wild videos taken by a smartphone.

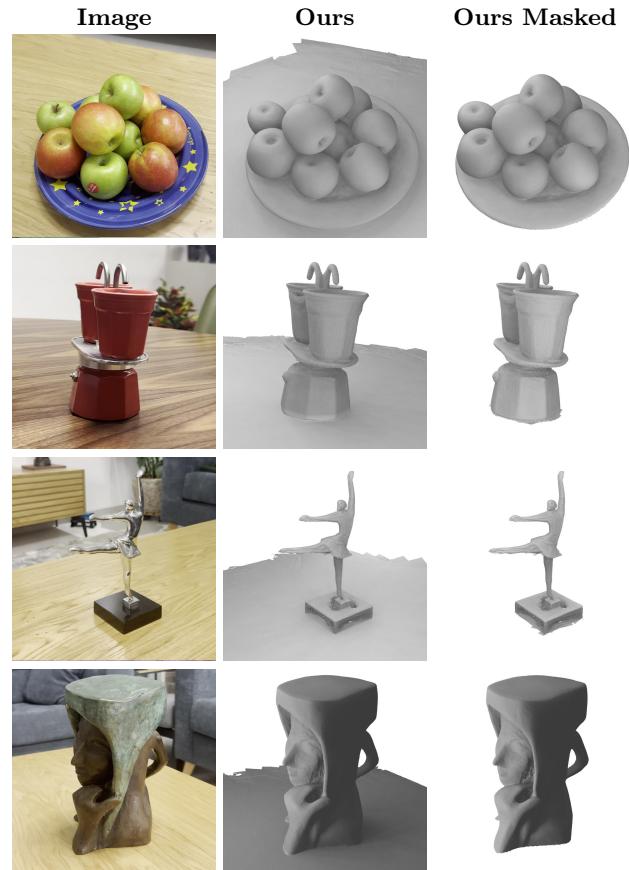


Fig. B.1: Examples of our method’s ability to segment objects using SAM [20] with depth projections, on in-the-wild videos taken by a smartphone.

C Additional Examples

C.1 Additional In-the-Wild Examples

In Fig. C.1 We present additional qualitative comparisons between reconstructions from our method and reconstructions from SuGaR [14], on in-the-wild videos taken by a smartphone in an uncontrolled environment. We run SuGaR according to the instructions from their official repository, with density regularization, no train/eval split, and 15000 refinement iterations.

C.2 Additional Examples from the DTU dataset

Fig. C.2 shows our method’s reconstruction of all of the DTU [16] dataset.

C.3 Additional Examples from The Ablation Study

Fig. C.3 presents results from the ablation study on the MobileBrick [23] test set, showing MVSFormer [4], MVSFormer with rendered images as input, and our method’s reconstructions, as well as the ground truth. Additional results on in-the-wild scenes taken by a smartphone are shown in Fig. C.4.



Fig. C.1: Additional qualitative comparisons between our method and SuGaR [14] on surface reconstruction from in-the-wild videos.

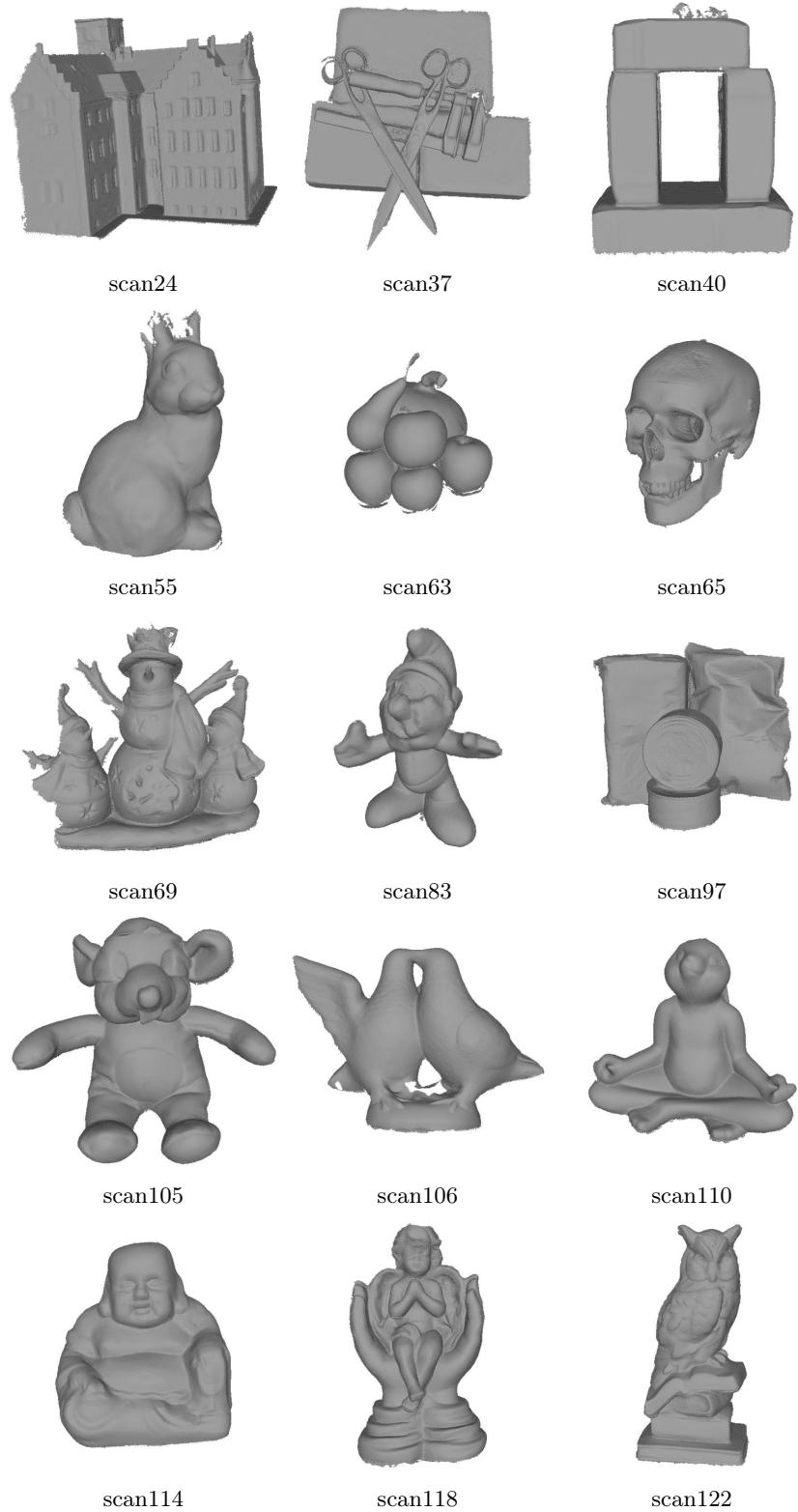
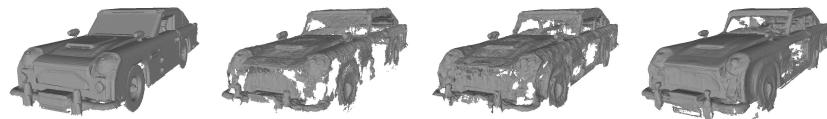
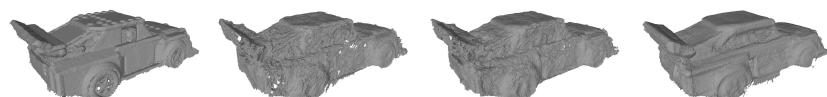


Fig. C.2: Our method’s reconstructions of the DTU [16] dataset.

**Ground
Truth****MVSFormer****MVSFormer
+ Rendered****Ours**

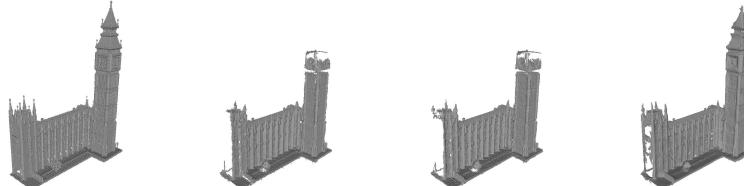
Aston



Audi



Beetles



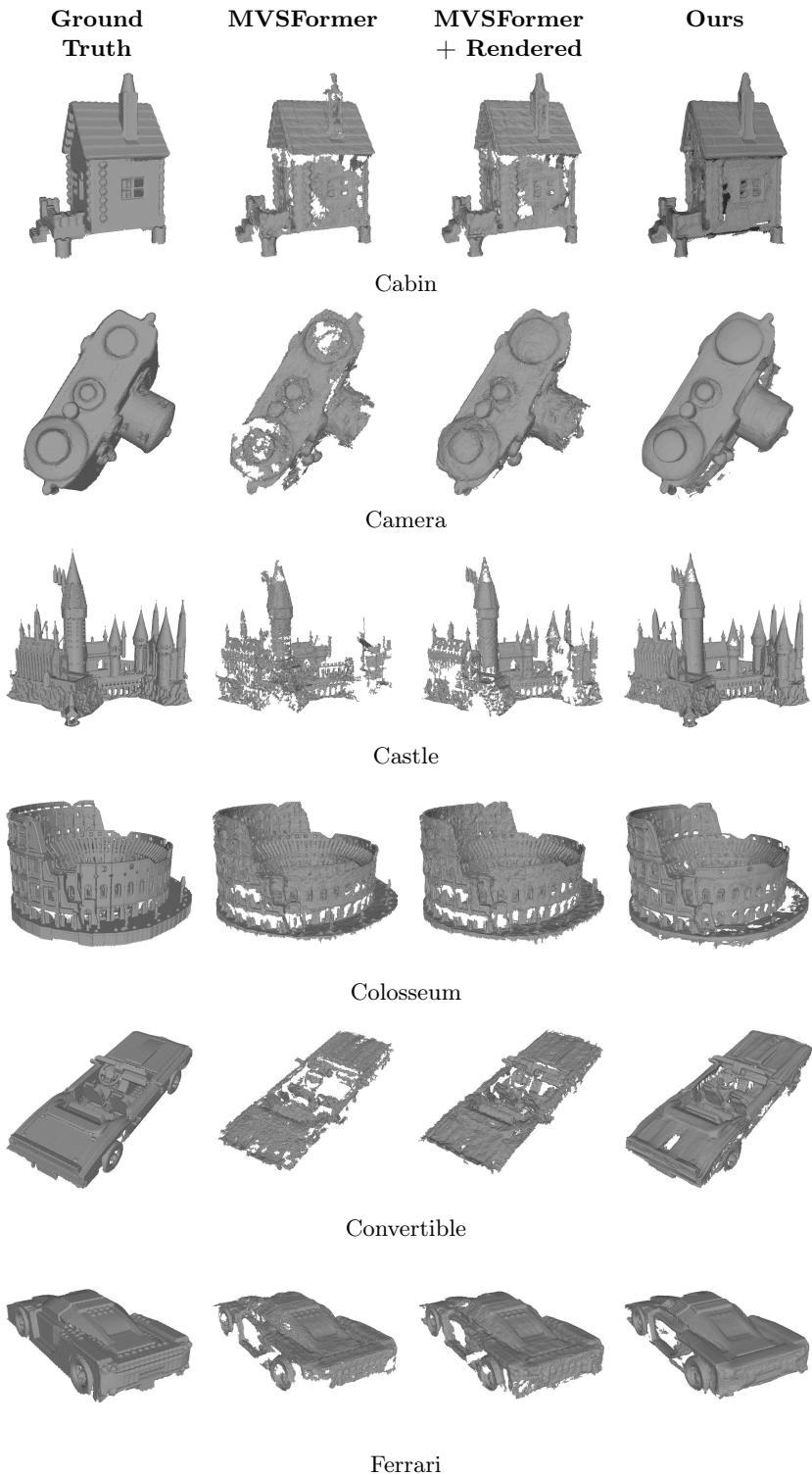
Big Ben



Boat



Bridge



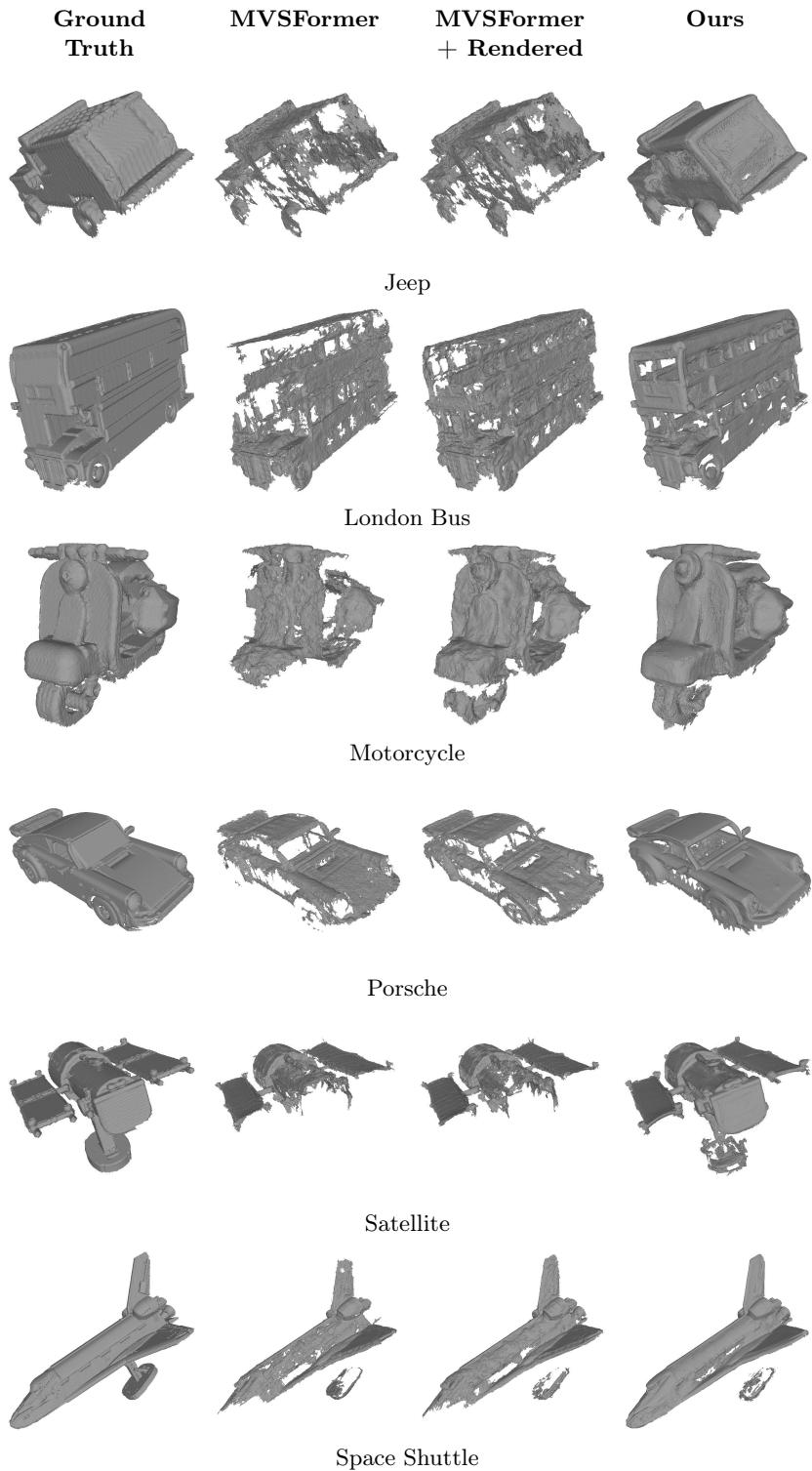


Fig. C.3: Left to right: ground truths of MobileBrick [23] dataset, reconstructions from MVSFormer [4] with original and rendered images, as well as our reconstruction.



Fig. C.4: Additional examples from the ablation study, showing reconstructions done by MVSFormer [4] on the original images and on the rendered images from the same poses, as well as our method’s reconstruction.