

Towards Realistic Example-based Modeling via 3D Gaussian Stitching

XINYU GAO*, State Key Lab of CAD&CG, Zhejiang University, China
ZIYI YANG*, State Key Lab of CAD&CG, Zhejiang University, China
BINGCHEN GONG*, The Chinese University of Hong Kong, Hong Kong
XIAOGUANG HAN, The Chinese University of Hong Kong, Shenzhen, China
SIPENG YANG, State Key Lab of CAD and CG, Zhejiang University, China
XIAOGANG JIN, State Key Lab of CAD and CG, Zhejiang University, China



Fig. 1. Our method can seamlessly stitch multiple 3D Gaussian fields together [Kerbl et al. 2023] interactively, resulting in new, highly detailed, and realistic objects. All of the geometric parts or models are derived from the BlendedMVS [Yao et al. 2020] and Mip360 [Barron et al. 2022] datasets.

Using parts of existing models to rebuild new models, commonly termed as example-based modeling, is a classical methodology in the realm of computer graphics. Previous works mostly focus on shape composition, making them very hard to use for realistic composition of 3D objects captured from real-world scenes. This leads to combining multiple NeRFs into a single 3D scene to achieve seamless appearance blending. However, the current SeamlessNeRF method struggles to achieve interactive editing and harmonious stitching for real-world scenes due to its gradient-based strategy and grid-based representation. To this end, we present an example-based modeling method that combines multiple Gaussian fields in a point-based representation using sample-guided synthesis. Specifically, as for composition, we create a GUI to segment and transform multiple fields in real time, easily obtaining a semantically meaningful composition of models represented by

*equal contribution

Authors' addresses: Xinyu Gao, State Key Lab of CAD&CG, Zhejiang University, China; Ziyi Yang, 14ziyiyang@gmail.com, State Key Lab of CAD&CG, Zhejiang University, China; Bingchen Gong, gongbingchen@gmail.com, The Chinese University of Hong Kong, Hong Kong; Xiaoguang Han, hanxiaoguang@cuhk.edu.cn, The Chinese University of Hong Kong, Shenzhen, China; Sipeng Yang, 12121024@zju.edu.cn, State Key Lab of CAD and CG, Zhejiang University, China; Xiaogang Jin, jin@cad.zju.edu.cn, State Key Lab of CAD and CG, Zhejiang University, China.

3D Gaussian Splatting (3DGS). For texture blending, due to the discrete and irregular nature of 3DGS, straightforwardly applying gradient propagation as SeamlessNeRF is not supported. Thus, a novel sampling-based cloning method is proposed to harmonize the blending while preserving the original rich texture and content. Our workflow consists of three steps: 1) real-time segmentation and transformation of a Gaussian model using a well-tailored GUI, 2) KNN analysis to identify boundary points in the intersecting area between the source and target models, and 3) two-phase optimization of the target model using sampling-based cloning and gradient constraints. Extensive experimental results validate that our approach significantly outperforms previous works in terms of realistic synthesis, demonstrating its practicality.

CCS Concepts: • **Computing methodologies** → **Image-based rendering**.

Additional Key Words and Phrases: Neural Rendering, 3D Model Synthesis, Composition

1 INTRODUCTION

As we all know, 3D scenes typically contain multiple 3D objects composed of various parts. Example-based modeling [Funkhouser

et al. 2004] is a technique that involves combining different parts from different objects to create new ones. This is a common tool in Computer Graphics (CG) modeling, where objects are designed in a non-realistic CG fashion. In this paper, we consider realistic example-based modeling, where all parts are captured from the real world, as shown in Fig. 1. This task becomes prominent with the emergence of Neural Radiance Fields, which enables photorealistic 3D reconstruction and rendering.

Among the various approaches designed for 3D modeling from multiple neural fields, a portion of the research [Gao et al. 2023; Liu et al. 2023a] is devoted to the inverse rendering process to achieve consistent lighting and shadowing. But these methods rarely consider a situation where the harmonious and seamless effect is required for merging or unifying two or more neural fields. SeamlessNeRF [Gong et al. 2023] is the first work to tackle seamless merging, attempting to address the consistency problem by propagating gradients on synthesis cases. Nonetheless, due to its implicit grid-based representation, SeamlessNeRF can neither achieve fine-grained editing (e.g. the face in the *Santa* case in Fig. 2) under real-world cases nor provide an interactive workflow in real-time. Additionally, its gradient-based strategy can produce significant artifacts (see Fig. 9) and fails to propagate structural characteristics when the condition becomes more complex (e.g., the *bottle* in the left-upper corner in Fig. 1). Therefore, achieving a harmonious and photorealistic stitching result on real-world data remains an unsolved challenge that needs further exploration.

To address the limitations mentioned above, we propose a new method for interactive editing and stitching multiple parts using explicit shape representation in 3D Gaussian Splatting. Our method has two significant advantages. First, its point-based representation enables fine-grained editing, allowing for detailed appearance optimization and the removal of artifacts. Second, its rasterizer pipeline provides a real-time interactive editing environment. Due to the discrete and irregular nature of 3D-GS, it is not feasible to conduct gradient propagation as SeamlessNeRF. Thus, we introduce a novel sampling-based optimization strategy that can seamlessly propagate not only color tones but also structural characteristics. Our evaluation benchmarks are primarily derived from real-world scenes, demonstrating our superior ability to handle complex cases.

More specifically, our pipeline takes multiple scenes as input, containing source and target objects represented by 3DGS. We then carefully segment these objects and apply rigid transformations in order to create a semantically meaningful composite in 3D space. An intersection boundary region between the objects is also identified before blending. The next is the key step in our process which aims to optimize the appearance of the target objects so that their texture and color match those of the source object. We achieve this by using a two-phase optimization scheme: the first phase involves sampling-based cloning (S-phase), and the second phase involves clustering-based tuning (T-phase). During the S-phase, the target field is optimized using a heuristic sampling strategy that considers the structural characteristics at the boundary. Additionally, an efficient 2D gradient constraint is applied to preserve the original texture content of the target field. However, optimizing solely with S-phase may lead to the appearance of artifacts or unintended color features that do not fit with the overall composite. Therefore, we

address this issue with T-phase, where we utilize a pre-calculated feature palette derived from the source field through aggregation and clustering. Subsequently, this palette is applied to tune the target field. It is important to note that the two-phase optimization is a joint procedure, where losses from the S-phase are always maintained while losses from the T-phase are added later during optimization.

In summary, our method makes the following contributions:

- The first work to use 3D-GS for realistic and seamless part compositing, enabling real-world example-based modeling.
- A novel sampling-based optimization strategy is proposed, with which not only the texture color but also the structural characteristics can be propagated seamlessly.
- A user-friendly GUI is carefully designed to support an interactive workflow of the modeling process in real time.

2 RELATED WORK

2.1 Example-based Seamless Editing

Seamless editing, particularly in the context of example-based image and texture synthesis, is a well-studied editing technique in computer graphics and image processing. As for textures, example-based texture synthesis [Efros and Leung 1999; Wei et al. 2009] intends to seamlessly create textures at any size from exemplars, which has been widely employed in contemporary graphics pipelines and game engines. In 2D image synthesis, patch-based synthesis techniques have been widely researched to seamlessly combine visually inconsistent images [Darabi et al. 2012; Pérez et al. 2023]. Meanwhile, Kwatra et al. [2005] introduced “Texture Optimization,” which transfers photographic textures to a target image for example-based synthesis. To facilitate structural image editing tasks, “PatchMatch” [Barnes et al. 2009] found approximate nearest-neighbor correspondences between patches in images for seamless image region reshuffling. In terms of seamless editing in 3D objects, Rocchini et al. [1999] and Dessein et al. [2014] propose methods for stitching and blending textures on 3D objects, respectively, while Yu et al. [2004] use the Poisson equation to implicitly modify the original mesh geometry via gradient field manipulation. Additionally, example-based modeling can also generate novel models from parts of existing models [Funkhouser et al. 2004], allowing untrained users to create interesting and detailed 3D designs, such as city building [Merrell 2007], things arrangements [Fisher et al. 2012], mesh segmentation [Katz et al. 2005], and merging [Kreavoy et al. 2007]. Recently, deep learning methods have leveraged generative models to generate diverse instances from a single exemplar [Li et al. 2023a; Wu and Zheng 2022] or a cluster of examples [Zhang et al. 2023]. Definitely, the example-based methodology is a valuable tool for creating diverse and novel content, which can reduce the workload for the artists or can be leveraged by procedural content generation programs. In our work, we combine this valuable idea with the advanced technique of 3DGS to create content directly from the real world.

2.2 Neural Scene Composition

Neural scene composition primarily involves the synthesis of multiple neural objects represented by neural fields, such as free-viewport video [Lin et al. 2022; Wang et al. 2023; Zhang et al. 2021], autonomous driving [Fu et al. 2022; Kundu et al. 2022; Ost et al. 2021;

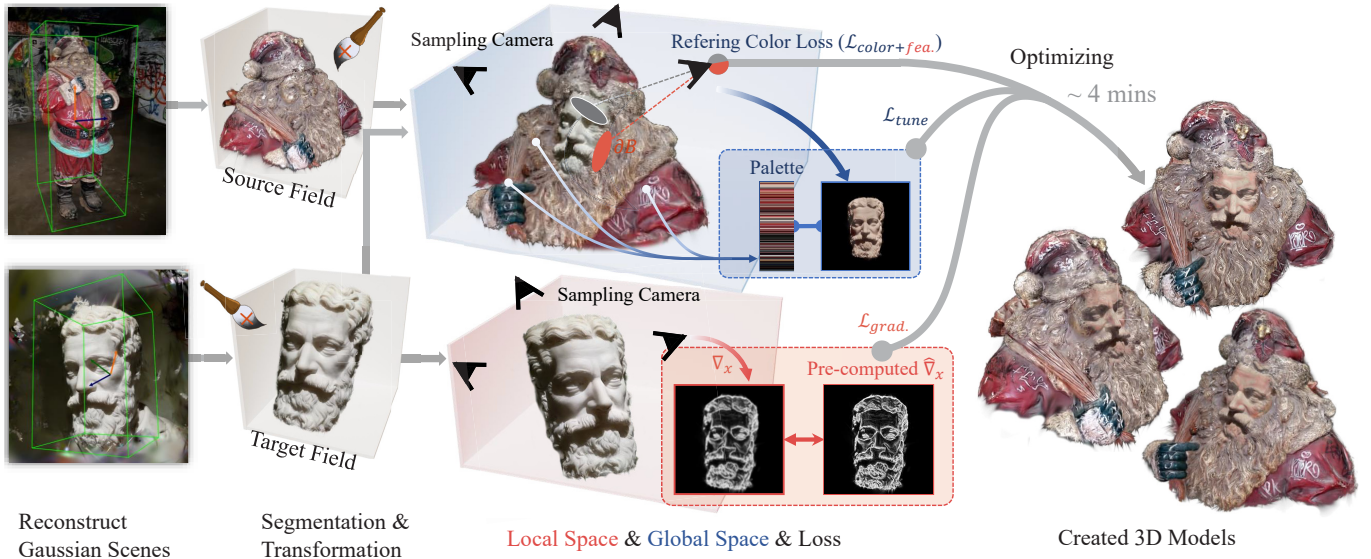


Fig. 2. Overview of our framework. Our novel pipeline provides an interactive editing experience and has real-time previewing capabilities to visualize the optimizing process, allowing for the seamless and interactive combination of multiple Gaussian fields.

Tancik et al. 2022; Yang et al. 2023; Zhou et al. 2023] and scene understanding [Kerr et al. 2023; Shuai et al. 2022; Wu et al. 2022; Yang et al. 2021]. And for those composition tasks with multiple pre-trained models, mesh scaffold [Yang et al. 2022; Yariv et al. 2023] or texture extraction [Chen et al. 2023b; Tang et al. 2023b] from the neural field are preferred to achieve higher render speed or rather fine-grained control. This type of work acts as a “bridge” between neural and traditional representations in order to improve performance using the classical graphics pipeline. A small portion of the work focuses on creating a mixed render pipeline for neural 3D scene composition tasks, combining traditional render techniques like ray tracing [Qiao et al. 2023], shadow mapping [Gao et al. 2024], and ambient occlusion [Gao et al. 2023]. There are also a few works that focus on creating a compositional scene with generative models like diffusion models [Po and Wetzstein 2023].

None of those works except Neural Imposter [Liu et al. 2023b] and SeamlessNeRF [Gong et al. 2023] focus on example-based modeling by stitching multiple part NeRFs. However, part objects in Neural Imposter are just placed together without any appearance blending, which cannot support a general case of 3D modeling. SeamlessNeRF achieved harmonious results on a small-scale synthesis dataset, making it the first work to discuss seamless example-based modeling with neural techniques today. However, SeamlessNeRF cannot handle real-world cases when the condition becomes more complex, nor can it perform interactive editing, which is commonly required in example-based modeling. On the contrary, our approach overcomes these limitations, performs well in real-world scenarios, and supports interactive editing using Gaussian fields.

2.3 3D Gaussians

3D Gaussian Splatting [Kerbl et al. 2023] is a point-based rendering method that has recently gained popularity [Chen et al. 2023a;

Huang et al. 2024; Li et al. 2024; Liang et al. 2023; Tang et al. 2023a; Yang et al. 2024a,b] due to its realistic rendering and significantly faster training time than NeRFs. Compared to the implicit representation of NeRF, 3DGS is more advantageous for editing tasks. The superior advance lies in the fact that, unlike previous work that embedded an object in a certain neural field (e.g., learnable grid or MLP network), once clusters of Gaussians are optimized, they can be easily fused together and fed into the rasterizer. The 3DGS pipeline was born with an intrinsic property suitable for composition.

3 SEAMLESS GAUSSIANS

Our approach starts with segmenting interesting parts from pre-trained Gaussian scenes. After acquiring target and source models represented by Gaussians, we carefully transform them to obtain a semantically meaningful composite. Then we optimize the target objects to achieve a harmonious composite through a two-phase (sampling-based cloning and clustering-based tuning) scheme. All these processes can be run interactively and previewed in real-time with our well-tailored GUI design.

3.1 Segmenting and Transforming Gaussians

Segmentation is the first step in example-based modeling, which involves picking out interesting parts as the components of the final artwork. Previous works have performed this task by providing guidance using 2D mask [Cen et al. 2023; Mirzaei et al. 2023] or injecting semantic label [Kerr et al. 2023] into a neural field. Now, benefiting from Gaussian representation (resembling point cloud), segmentation can become more practical at a finer-grained level. In our pipeline, we show that a combination of a simple bounding box and a user brush can work very well for a clean mask (see Fig. 14). For instance, we can mask the *sculpture* with a brush to match the shape of *Santa’s* face (see Fig. 2).

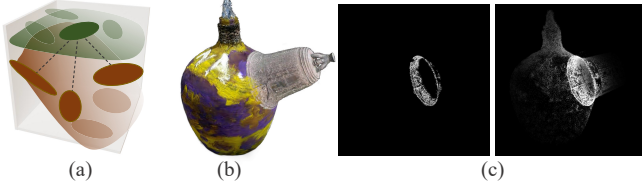


Fig. 3. For a Gaussian point in the target field, its (a) K-nearest neighbors in the source field can be leveraged to justify whether this point belongs to the intersection boundary region. We use the boundary of (b) as an example to demonstrate the effectiveness of this strategy, as shown in (c).

Transformation aims at placing multiple interesting parts \mathcal{G}_i represented by Gaussians to form a semantically meaningful composite \mathcal{M} , which can be denoted as:

$$\mathcal{G}_i^{\text{global}} = F(\mathcal{G}_i^{\text{local}} | \hat{\mathbf{q}}_i, \mathbf{t}_i, s_i), \quad \mathcal{G}_i \in \mathcal{M} \quad (1)$$

where F is the rigid transformation applied on one part of Gaussians with rotation $\hat{\mathbf{q}}_i$ (represented in quaternion), translation \mathbf{t}_i , and scale s_i , transforming the part from its local space to the global space. Specifically, the partial attributes of each \mathcal{G} should be modified, which includes position \mathbf{x} , scaling s , rotation \mathbf{q} (in quaternion), and feature \mathbf{f} (represented as spherical harmonics). The position and scaling can be performed trivially, while the transformed rotation \mathbf{q}' and feature \mathbf{f}' can be expressed as:

$$\begin{aligned} \mathbf{q}' &= \mathbf{q}\hat{\mathbf{q}}, \\ \mathbf{f}' &= M_{bands}(\mathbf{f} | \hat{\mathbf{q}}), \end{aligned} \quad (2)$$

where M_{bands} means we use a set of matrices to rotate each band of SH coefficients introduced by [Ivanic and Ruedenberg 1996].

3.2 Boundary Condition by KNN Analyzing

After transformation, certain points in one field approach another field (see Fig. 3), forming intersection boundary regions between all Gaussians. For the sake of simplicity, we will use two Gaussians, source field and target field, to demonstrate our approach.

Before optimization, the boundary points in the target field must be identified, as this is the critical and initial condition for harmonization. For each Gaussian point in target field \mathcal{T} , we search its K-nearest neighbors in source field \mathcal{S} , which can be denoted by:

$$\{b_i\}_K = \underset{\mathcal{S}}{\text{KNN}}(a), \quad a \in \mathcal{T}, b_i \in \mathcal{S} \quad (3)$$

where a is a point in the target field, and b_i is a point in the source field. Whether a point a belongs to boundary ∂B can be identified as $a \in \partial B$ iff:

$$\frac{1}{K} \sum_i^K |b_i - a| < \beta \quad \text{and} \quad o(a) > \tau, \quad (4)$$

where $o(a)$ is the opacity of that Gaussian point, $|b_i - a|$ is the Euclidean distance between b_i and a . τ and β are thresholds and we empirically set τ to 0.95, β to $0.05 \times L$. L is the size of the composite. (e.g. measured by the bounding box). An additional method for a better boundary condition on real-world data is that we discard outliers in both fields (e.g. some Gaussian points are far from the others, which may occur in some scenes).

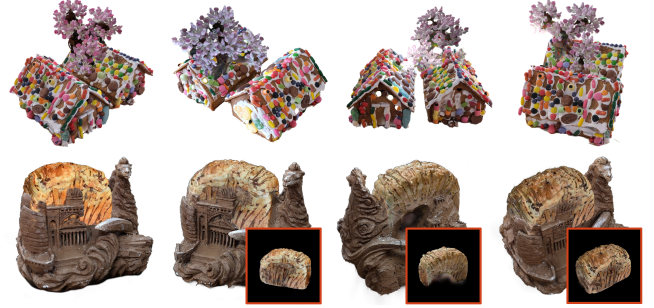


Fig. 4. ablation study on the color loss in the S-phase. Without color loss, the propagation is inefficient and will not begin. The cases shown above have been running for more than twice as long, but they are still trapped in insufficient propagation. It is because, without color loss, only a small number of points’ features need to be updated at first, as opposed to shared weights in an MLP applied to all points. That minor “forces” cannot drive the overall minimization of the gradient loss.

We calculate referenced features for these boundary points in order to confirm the boundary condition. For each $a \in \partial B$, its target feature is:

$$\hat{\mathbf{f}}(a) = \frac{1}{K} \sum_i^K \mathbf{f}'(b_i), \quad a \in \partial B, b_i \in \underset{\mathcal{S}}{\text{KNN}}(a) \quad (5)$$

where $\mathbf{f}'(b_i)$ is the feature of b_i after transformation. To achieve this boundary condition, we optimize boundary points toward their target features:

$$\mathcal{L}_{feature} = \sum_{a \in \partial B} \left\| \mathbf{f}'(a) - \hat{\mathbf{f}}(a) \right\|_2^2, \quad (6)$$

where $\mathbf{f}'(a)$ is the feature of a and we directly apply this loss on SH coefficients.

3.3 Sampling-based Cloning

We propose sampling-based cloning as our “S-phase” in optimization. The core idea is how to seamlessly propagate the style in boundary through the remaining points in the target field while preserving its rich content. In contrast to a regular grid suitable with a gradient-based strategy in SeamlessNeRF [Gong et al. 2023], Gaussian points are irregularly and discretely distributed in 3D space. As a result, alternative approaches need to be explored. A straightforward idea is that given a point in target field \mathcal{T} , one can calculate the feature difference between that point and its neighbors in \mathcal{T} , resembling “Laplacian coordinates”. Then, one can use that “difference” as the regularizer while minimizing $\mathcal{L}_{feature}$. However, this naive approach may fail even before propagation begins (see Fig. 4). Furthermore, the boundary’s structural characteristics (such as the *bottle-bell* intersection in the right-upper corner of Fig. 9) necessitate seamless cloning, which significantly improves the stitching quality.

Hence, we propose an effective sampling strategy to explicitly propagate features for each remaining point outside the boundary. The core idea lies in the way of searching several “driven points” for a candidate. The color of the candidate is driven by those points.

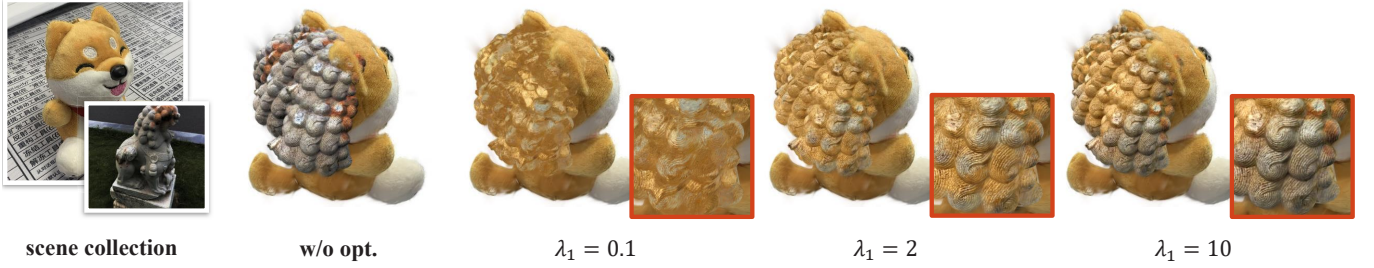


Fig. 5. ablation study on the effectiveness of gradient loss for different weights. Experiments show that higher weights can help to preserve more content while preventing harmonization.

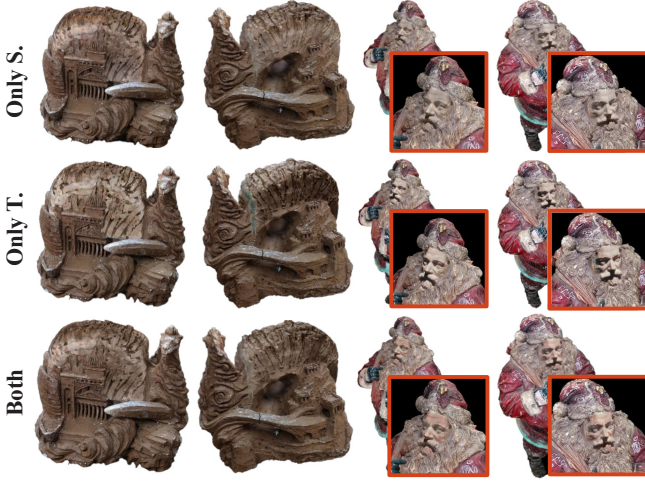


Fig. 6. ablation study on sampling-based cloning (S.) and clustering-based tuning (T.). Here, “Both” means the full scheme.

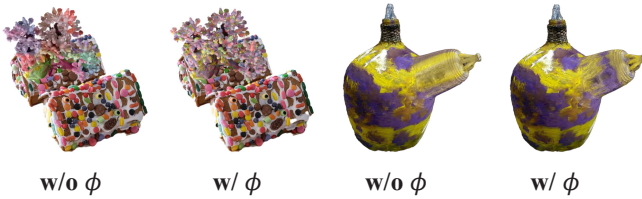


Fig. 7. Ablation study on the impact of mapping function ϕ in the S-phase. Random effects make composition more realistic.

For each $a \in \mathcal{T} - \partial B$, the optimizing target of its color in direction \mathbf{d}_a is:

$$\hat{\mathbf{f}}(a, \mathbf{d}_a) = \frac{1}{K} \sum_i^K \mathbf{f}'(b_i, \mathbf{d}_b), \quad a \in \mathcal{T} - \partial B, b_i \in KNN(\phi(a)) \quad (7)$$

where $\mathbf{f}(a, \mathbf{d}_a)$ means sampling SH color in view direction \mathbf{d}_a (from point a to camera), the same as $\mathbf{f}(b, \mathbf{d}_b)$. The camera centers are uniformly sampled from the surface of a sphere centered on the composite object’s origin. It is important to note that the sampling strategy $KNN(\phi(a))$ maps the locations of nearby candidate points a -s to the correlated neighboring “driven points” and inherits the

continuity of the textures from those “driven points”. We use $\phi(x) = x + \sin(\gamma \cdot \delta x)$ to add random effect by disturbing KNN searching (see Fig. 7), where x is the position of a , δx is the distance between a and its nearest b_i in boundary, and γ is empirically set to 10. A larger γ is suitable for higher structural frequencies. In this way, we can synthesize structurally aware stitching results. With Eq. (7), we add a color loss to the S-phase:

$$\mathcal{L}_{color} = \sum_{a \in \mathcal{T} - \partial B} \left\| \mathbf{f}'(a, \mathbf{d}_a) - \hat{\mathbf{f}}(a, \mathbf{d}_a) \right\|_2^2, \quad (8)$$

so that the color of those candidates can be optimized towards their target to achieve our explicit feature propagation.

To preserve the original rich content in \mathcal{T} , we present a more efficient gradient loss calculated in the local space of \mathcal{T} , leveraging the guidance in 2D screen space:

$$\mathcal{L}_{grad} = \sum_{x \in I} \left\| \nabla_x I^{\mathcal{T}}(p) - \hat{\nabla}_x I^{\mathcal{T}}(p) \right\|_2^2, \quad (9)$$

$$I^{\mathcal{T}}(p) = \mathcal{R}(\mathcal{G}_{\mathcal{T}}^{local}, p),$$

where p is the randomly sampled camera in the local space of target field \mathcal{T} , I is the rendered color image by rasterizer \mathcal{R} of 3DGS. We pre-calculate $\hat{\nabla}_x I$ for each camera with the Sobel operator [Sobel et al. 1968] before the optimization starts. We found that supervising gradients in screen space is more efficient than the straightforward one, as shown in Fig. 11.

3.4 Clustering-based Tuning

While S-phase optimization is effective in preserving local color consistency, relying solely on it may lead to misaligned global appearance, such as uneven brightness, hues, and saturation (See Fig. 6). Therefore, we propose using a clustering extracted color palette to perform global tuning, which we refer to as the “T-phase” in optimization. This approach enhances the overall harmony of the composite by performing dynamic matches to a palette. To implement the T-phase, we first aggregate and cluster the color of the source field from various angles:

$$\{\mathbf{c}_i\}_N, \{w_i\}_N \leftarrow \mathcal{A}(\mathcal{G}_S^{global}), \quad (10)$$

where \mathbf{c}_i is the color (cluster center) in palette, w_i is the sample percentage occupied by the center, and \mathcal{A} stands for our aggregation algorithm. Our approach, inspired by Li et al.’s work [Li et al. 2023b], uses a streaming method to accelerate color aggregation. We start



Fig. 8. Results for more real-world data from the BlendedMVS [Yao et al. 2020] and Mip360 [Barron et al. 2022] datasets, demonstrating that our method can produce realistic effects in real-world scenarios.

with three bins, collect color samples from a random view, and calculate the new color center for each bin by averaging the original center and new samples collected in it. The number of bins expands to accommodate far-off samples. Centers expire after 20 iterations with no sufficient votes. We repeat this process until all color centers are stable.

Once the aggregation process finishes, those color centers will form a palette (see Fig. 2). We employ the following loss in our T-phase as a pixel-wise summation:

$$\begin{aligned} \mathcal{L}_{tune} &= \sum_{c \in I'} w_{\chi_c} \|c - c_{\chi_c}\|_2^2, I' \leftarrow \{I_x^T(p) | \alpha(x) > 0.95\}, \\ \chi_c &= \arg \min_{1 \leq i \leq N} \{\|c - c_i\|_2 - w_i\}, \end{aligned} \quad (11)$$

where p is the randomly sampled camera in the global space, and α is the alpha mask corresponding to I^T . Both α and I^T are rendered by rasterizer \mathcal{R} . χ represents the target bin’s index, and it is determined by both the distance from color centers and the probability density of bins. Our final total loss function can then be expressed as:

$$\mathcal{L}_{total} = \mathcal{L}_{feature} + \mathcal{L}_{color} + \lambda_1 \mathcal{L}_{grad} + \lambda_2 \mathcal{L}_{tune}, \quad (12)$$

where both λ_1 and λ_2 are empirically set to 2 in our experiments.

4 EXPERIMENT

To test the effectiveness and generality of our approach, we conducted experiments on a variety of fascinating 3D objects. We interactively built 21 composite results, comprising a total of 39 part models: 17 from BlendedMVS [Yao et al. 2020], 4 from Mip360 [Barron et al. 2022], 16 from SeamlessNeRF datasets, and 2 created by ourselves in a graphics engine. For more results or the implementation details, please refer to our supplementary materials.

4.1 Qualitative Comparison

We compare our method to SeamlessNeRF [Gong et al. 2023], the first and most recent work that approaches our goal. Fig. 9 depicts three comparison cases. In the first case (clay & bread), SeamlessNeRF failed to achieve high-level geometry editability and struggled with artifacts caused by implicit representation. In the second case (bottle and bell), SeamlessNeRF failed to maintain a harmonious seamless effect due to applying the gradient-based strategy on the complex boundary. In the third case, SeamlessNeRF failed to propagate sufficient color tones due to the complex gradients in the boundary. In addition, we show that the 2D-guided style-transfer method [Nguyen-Phuoc et al. 2022] cannot produce a seamless stitching effect, as shown in Fig. 10. On the contrary, ours can handle all of these situations while producing harmonious results.

4.2 Quantitative Comparison

Currently, there is neither a specialized dataset providing ground truth nor an established metric to assess the realism of a 3D model’s appearance, making it challenging to evaluate the effectiveness of our approach quantitatively. Nevertheless, we force an evaluation utilizing VQA (Video Quality Assessment) methods, as outlined by Wu et al. [2023], and explored the use of 2D projection in video display for assessment purposes. Our results, presented in Tab. 1, demonstrate that our average score surpasses the baseline. For a comprehensive understanding of the quantitative experiments, please refer to our supplementary materials.

4.3 Ablation Study

Effectiveness of 2D Gradient Loss. Fig. 5 depicts the effect of gradient loss at various weights. Higher weights can help to preserve

	ours	SeamlessNeRF
VQA average score \uparrow	0.784	0.753

Table 1. Quantitative comparison between ours and the baseline.

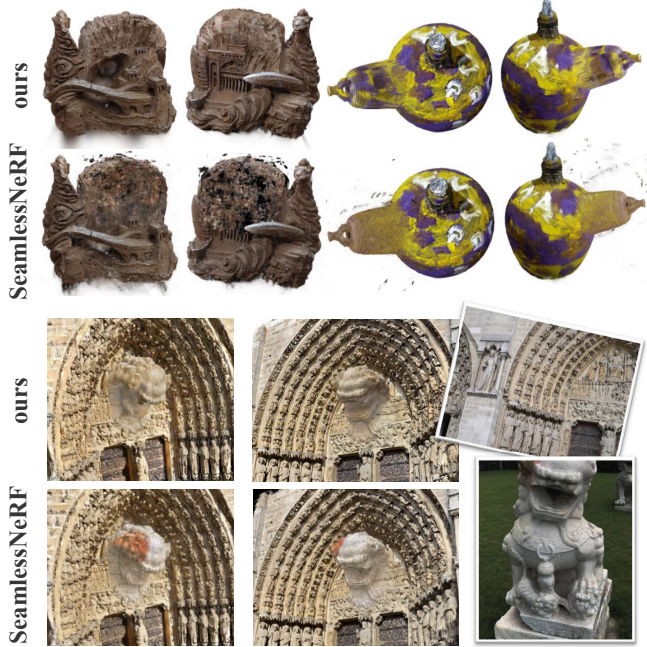


Fig. 9. Comparisons between our approach and the baseline methods [Gong et al. 2023]. SeamlessNeRF failed in all of these real-world scenarios.



Fig. 10. We show that these style-transfer results [Nguyen-Phuoc et al. 2022] fail to achieve our effect. Here, we re-implement SNeRF’s strategy [Nguyen-Phuoc et al. 2022] based on Gaussians to produce results above.

more content while obstructing harmonization. Fig. 11 demonstrates that 2D gradient loss with Sobel operator is significantly more effective than the simple one mentioned in Sec. 3.3.

Functionality of S-phase and T-phase. We demonstrate the efficacy of our two-phase scheme in Fig. 6. The S-phase aids in seamless boundary formation, while the T-phase aids in global harmonization when only the S-phase is present.

Effectiveness of Sampling Strategy for View-dependent Effects. We ablate the sampling strategy in the S-phase (see Fig. 12) to show that view-dependent effects can be properly propagated using this strategy instead of random sampling.

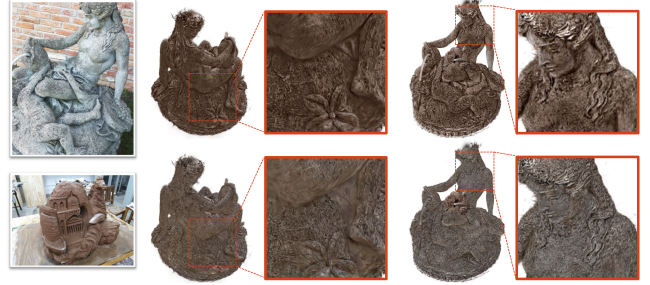


Fig. 11. ablation study on two kinds of gradient loss. The 2D gradient supervision (upper row) is more effective than the straightforward one since it focuses on the surface instead of the whole space.

Fig. 12. Ablation study on keeping view-dependent effects by sampling strategy in the S-phase. With that strategy used in Equ. (7), the upper-view color of paint on *bell* is properly propagated.

4.4 Editor and Application

To enable a practical and user-friendly workflow, we created an interactive GUI editor that can control and visualize any procedure in the entire process in real-time, including Gaussian segmentation and transformation, boundary identification, and optimization (see Fig. 16 and refer to the supplementary video for more details). Our framework can generate high-fidelity and seamless results across a wide range of real-world scenarios, providing distinct advantages in the direct creation of imaginative 3D models from reality.

5 CONCLUSIONS AND LIMITATIONS

We have developed a highly efficient and effective interactive framework for creating realistic 3D models. The method involves stitching Gaussian components seamlessly to create a harmonious 3D model that is an accurate representation of the real world. Our approach has been tested on real-world datasets and has proved to be capable of handling complex cases with a user-friendly interface. This presents a promising avenue for example-based modeling directly from the real world.

Limitations and Future Work. Currently, our work is unable to transform Gaussian models in a non-rigid manner, which may make it difficult to develop more imaginative cases. To enable a more flexible composition, we can use deformation methods such as ARAP [Igarashi et al. 2005] in the future. Furthermore, achieving a consistent lighting effect can help improve composition quality under intense lighting.

REFERENCES

- Connelly Barnes, Eli Shechtman, Adam Finkelstein, and Dan B Goldman. 2009. Patch-Match: A Randomized Correspondence Algorithm for Structural Image Editing. *ACM Trans. Graph.* 28, 3, Article 24 (jul 2009), 11 pages. <https://doi.org/10.1145/1531326.1531330>
- Jonathan T Barron, Ben Mildenhall, Dor Verbin, Pratul P Srinivasan, and Peter Hedman. 2022. Mip-nerf 360: Unbounded anti-aliased neural radiance fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 5470–5479.
- Jiazhong Cen, Zanwei Zhou, Jiemin Fang, Chen Yang, Wei Shen, Lingxi Xie, Dongsheng Jiang, Xiaopeng Zhang, and Qi Tian. 2023. Segment Anything in 3D with NeRFs. In *NeurIPS*.
- Yiwen Chen, Zilong Chen, Chi Zhang, Feng Wang, Xiaofeng Yang, Yikai Wang, Zhong-gang Cai, Lei Yang, Huaping Liu, and Guosheng Lin. 2023a. GaussianEditor: Swift and Controllable 3D Editing with Gaussian Splatting. arXiv:2311.14521 [cs.CV]
- Zhiqin Chen, Thomas Funkhouser, Peter Hedman, and Andrea Tagliasacchi. 2023b. Mobilenerf: Exploiting the polygon rasterization pipeline for efficient neural field rendering on mobile architectures. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 16569–16578.
- Soheil Darabi, Eli Shechtman, Connelly Barnes, Dan B Goldman, and Pradeep Sen. 2012. Image melding: Combining inconsistent images using patch-based synthesis. *ACM Transactions on graphics (TOG)* 31, 4 (2012), 1–10.
- Arnaud Dessein, William AP Smith, Richard C Wilson, and Edwin R Hancock. 2014. Seamless texture stitching on a 3D mesh by poisson blending in patches. In *2014 IEEE International Conference on Image Processing (ICIP)*. IEEE, 2031–2035.
- Alexei A Efros and Thomas K Leung. 1999. Texture synthesis by non-parametric sampling. In *Proceedings of the seventh IEEE international conference on computer vision*, Vol. 2. IEEE, 1033–1038.
- Matthew Fisher, Daniel Ritchie, Manolis Savva, Thomas Funkhouser, and Pat Hanrahan. 2012. Example-based synthesis of 3D object arrangements. *ACM Transactions on Graphics (TOG)* 31, 6 (2012), 1–11.
- Xiao Fu, Shangzhan Zhang, Tianrun Chen, Yichong Lu, Lanyun Zhu, Xiaowei Zhou, Andreas Geiger, and Yiyi Liao. 2022. Panoptic nerf: 3d-to-2d label transfer for panoptic urban scene segmentation. In *2022 International Conference on 3D Vision (3DV)*. IEEE, 1–11.
- Thomas Funkhouser, Michael Kazhdan, Philip Shilane, Patrick Min, William Kiefer, Ayellet Tal, Szymon Rusinkiewicz, and David Dobkin. 2004. Modeling by example. *ACM transactions on graphics (TOG)* 23, 3 (2004), 652–663.
- Jian Gao, Chun Gu, Youtian Lin, Hao Zhu, Xun Cao, Li Zhang, and Yao Yao. 2023. Re-lightable 3D Gaussian: Real-time Point Cloud Relighting with BRDF Decomposition and Ray Tracing. arXiv:2311.16043 (2023).
- Xinyu Gao, Ziyi Yang, Yunlu Zhao, Yuxiang Sun, Xiaogang Jin, and Changqing Zou. 2024. A General Implicit Framework for Fast NeRF Composition and Rendering. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 38. 1833–1841.
- Bingchen Gong, Yuehao Wang, Xiaoguang Han, and Qi Dou. 2023. SeamlessNeRF: Stitching Part NeRFs with Gradient Propagation. In *SIGGRAPH Asia 2023 Conference Papers*. 1–10.
- Yi-Hua Huang, Yang-Tian Sun, Ziyi Yang, Xiaoyang Lyu, Yan-Pei Cao, and Xiaojuan Qi. 2024. Sc-gs: Sparse-controlled gaussian splatting for editable dynamic scenes. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 4220–4230.
- Takeo Igarashi, Tomer Moscovich, and John F Hughes. 2005. As-rigid-as-possible shape manipulation. *ACM transactions on Graphics (TOG)* 24, 3 (2005), 1134–1141.
- Joseph Ivanic and Klaus Ruedenberg. 1996. Rotation matrices for real spherical harmonics. Direct determination by recursion. *The Journal of Physical Chemistry* 100, 15 (1996), 6342–6347.
- Sagi Katz, George Leifman, and Ayellet Tal. 2005. Mesh segmentation using feature point and core extraction. *The Visual Computer* 21 (2005), 649–658.
- Bernhard Kerbl, Georgios Kopanas, Thomas Leimkühler, and George Drettakis. 2023. 3D Gaussian Splatting for Real-Time Radiance Field Rendering. *ACM Transactions on Graphics* 42, 4 (July 2023). <https://repo-sam.inria.fr/fungraph/3d-gaussian-splatting/>
- Justin Kerr, Chung Min Kim, Ken Goldberg, Angjoo Kanazawa, and Matthew Tancik. 2023. Lrf: Language embedded radiance fields. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 19729–19739.
- Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C. Berg, Wan-Yen Lo, Piotr Dollár, and Ross Girshick. 2023. Segment Anything. arXiv:2304.02643 (2023).
- Vladislav Kreevov, Dan Julius, and Alla Sheffer. 2007. Model composition from interchangeable components. In *15th Pacific Conference on Computer Graphics and Applications (PG'07)*. IEEE, 129–138.
- Abhijit Kundu, Kyle Genova, Xiaoqi Yin, Alireza Fathi, Caroline Pantofaru, Leonidas J Guibas, Andrea Tagliasacchi, Frank Dellaert, and Thomas Funkhouser. 2022. Panoptic neural fields: A semantic object-aware neural scene representation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 12871–12881.
- Vivek Kwatra, Irfan Essa, Aaron Bobick, and Nipun Kwatra. 2005. Texture Optimization for Example-Based Synthesis. *ACM Trans. Graph.* 24, 3 (jul 2005), 795–802. <https://doi.org/10.1145/1073204.1073263>
- Lingzhi Li, Zhen Shen, Zhongshu Wang, Li Shen, and Liefeng Bo. 2023b. Compressing volumetric radiance fields to 1 mb. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 4222–4231.
- Weiyeu Li, Xuelin Chen, Jue Wang, and Baoquan Chen. 2023a. Patch-based 3D Natural Scene Generation from a Single Example. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 16762–16772.
- Zhan Li, Zhang Chen, Zhong Li, and Yi Xu. 2024. Spacetime gaussian feature splatting for real-time dynamic view synthesis. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 8508–8520.
- Yixun Liang, Xin Yang, Jiantao Lin, Haodong Li, Xiaogang Xu, and Yingcong Chen. 2023. LucidDreamer: Towards High-Fidelity Text-to-3D Generation via Interval Score Matching. arXiv:2311.11284 [cs.CV]
- Haotong Lin, Sida Peng, Zhen Xu, Yunzhi Yan, Qing Shuai, Hujun Bao, and Xiaowei Zhou. 2022. Efficient Neural Radiance Fields for Interactive Free-viewpoint Video. In *SIGGRAPH Asia Conference Proceedings*.
- Ruiyang Liu, Jinxu Xiang, Bowen Zhao, Ran Zhang, Jingyi Yu, and Changxi Zheng. 2023b. Neural impostor: Editing neural radiance fields with explicit shape manipulation. In *Computer Graphics Forum*. Wiley Online Library, e14981.
- Yuan Liu, Peng Wang, Cheng Lin, Xiaoxiao Long, Jiepeng Wang, Lingjie Liu, Taku Komura, and Wenping Wang. 2023a. NeRO: Neural Geometry and BRDF Reconstruction of Reflective Objects from Multiview Images. *ACM Trans. Graph.* 42, 4, Article 114 (jul 2023), 22 pages. <https://doi.org/10.1145/3592134>
- Paul Merrell. 2007. Example-based model synthesis. In *Proceedings of the 2007 symposium on Interactive 3D graphics and games*. 105–112.
- Ashkan Mirzaei, Tristan Aumentado-Armstrong, Konstantinos G Derpanis, Jonathan Kelly, Marcus A Brubaker, Igor Gilitschenski, and Alex Levinshstein. 2023. SPIn-NeRF: Multiview segmentation and perceptual inpainting with neural radiance fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 20669–20679.
- Thu Nguyen-Phuoc, Feng Liu, and Lei Xiao. 2022. SNeRF: Stylized Neural Implicit Representations for 3D Scenes. 41, 4, Article 142 (jul 2022), 11 pages. <https://doi.org/10.1145/3528223.3530107>
- Julian Ost, Fahim Mannan, Nils Thuerey, Julian Knodt, and Felix Heide. 2021. Neural scene graphs for dynamic scenes. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2856–2865.
- Patrick Pérez, Michel Gangnet, and Andrew Blake. 2023. Poisson image editing. In *Seminal Graphics Papers: Pushing the Boundaries, Volume 2*. 577–582.
- Ryan Po and Gordon Wetzstein. 2023. Compositional 3d scene generation using locally conditioned diffusion. arXiv preprint arXiv:2303.12218 (2023).
- Yi-Ling Qiao, Alexander Gao, Yiran Xu, Yue Feng, Jia-Bin Huang, and Ming C Lin. 2023. Dynamic mesh-aware radiance fields. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 385–396.
- Claudio Rocchini, Paolo Cignoni, Claudio Montani, and Roberto Scopigno. 1999. Multiple textures stitching and blending on 3D objects. In *Rendering Techniques '99: Proceedings of the Eurographics Workshop in Granada, Spain, June 21–23, 1999*. Springer, 119–130.
- Qing Shuai, Chen Geng, Qi Fang, Sida Peng, Wenhao Shen, Xiaowei Zhou, and Hujun Bao. 2022. Novel view synthesis of human interactions from sparse multi-view videos. In *SIGGRAPH Conference Proceedings*. 1–10.
- Irwin Sobel, Gary Feldman, et al. 1968. A 3x3 isotropic gradient operator for image processing. *a talk at the Stanford Artificial Project in (1968)*. 271–272.
- Matthew Tancik, Vincent Casser, Xinchen Yan, Sabeek Pradhan, Ben Mildenhall, Pratul P Srinivasan, Jonathan T Barron, and Henrik Kretzschmar. 2022. Block-nerf: Scalable large scene neural view synthesis. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 8248–8258.
- Jiaxiang Tang, Jiawei Ren, Hang Zhou, Ziwei Liu, and Gang Zeng. 2023a. DreamGaussian: Generative Gaussian Splatting for Efficient 3D Content Creation. arXiv preprint arXiv:2309.16653 (2023).
- Jiaxiang Tang, Hang Zhou, Xiaokang Chen, Tianshu Hu, Errui Ding, Jingdong Wang, and Gang Zeng. 2023b. Delicate Textured Mesh Recovery from NeRF via Adaptive Surface Refinement. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*. 17739–17749.
- Liao Wang, Qiang Hu, Qihan He, Ziyu Wang, Jingyi Yu, Tinne Tuytelaars, Lan Xu, and Minye Wu. 2023. Neural Residual Radiance Fields for Streamably Free-Viewpoint Videos. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. 76–87.
- Li-Yi Wei, Sylvain Lefebvre, Vivek Kwatra, and Greg Turk. 2009. State of the art in example-based texture synthesis. *Eurographics 2009, State of the Art Report, EG-STAR (2009)*. 93–117.
- Haoning Wu, Erli Zhang, Liang Liao, Chaofeng Chen, Jingwen Hou Hou, Annan Wang, Wenxiu Sun Sun, Qiong Yan, and Weisi Lin. 2023. Exploring Video Quality Assessment on User Generated Contents from Aesthetic and Technical Perspectives. In *International Conference on Computer Vision (ICCV)*.
- Qianyi Wu, Xian Liu, Yuedong Chen, Kejie Li, Chuanxia Zheng, Jianfei Cai, and Jianmin Zheng. 2022. Object-compositional neural implicit surfaces. In *Computer Vision—ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings*.

- Part XXVII. Springer, 197–213.
- Rundi Wu and Changxi Zheng. 2022. Learning to Generate 3D Shapes from a Single Example. *ACM Trans. Graph.* 41, 6, Article 224 (nov 2022), 19 pages. <https://doi.org/10.1145/3550454.3555480>
- Bangbang Yang, Chong Bao, Junyi Zeng, Hujun Bao, Yinda Zhang, Zhaopeng Cui, and Guofeng Zhang. 2022. Neumesh: Learning disentangled neural mesh-based implicit field for geometry and texture editing. In *European Conference on Computer Vision*. Springer, 597–614.
- Bangbang Yang, Yinda Zhang, Yinghao Xu, Yijin Li, Han Zhou, Hujun Bao, Guofeng Zhang, and Zhaopeng Cui. 2021. Learning object-compositional neural radiance field for editable scene rendering. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 13779–13788.
- Ze Yang, Yun Chen, Jingkang Wang, Sivabalan Manivasagam, Wei-Chiu Ma, Anqi Joyce Yang, and Raquel Urtasun. 2023. UniSim: A Neural Closed-Loop Sensor Simulator. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 1389–1399.
- Ziyi Yang, Xinyu Gao, Wen Zhou, Shaohui Jiao, Yuqing Zhang, and Xiaogang Jin. 2024a. Deformable 3d gaussians for high-fidelity monocular dynamic scene reconstruction. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 20331–20341.
- Zeyu Yang, Hongye Yang, Zijie Pan, and Li Zhang. 2024b. Real-time Photorealistic Dynamic Scene Representation and Rendering with 4D Gaussian Splatting. In *International Conference on Learning Representations (ICLR)*.
- Yao Yao, Zixin Luo, Shiwei Li, Jingyang Zhang, Yufan Ren, Lei Zhou, Tian Fang, and Long Quan. 2020. BlendedMVS: A Large-Scale Dataset for Generalized Multi-View Stereo Networks. In *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. 1787–1796. <https://doi.org/10.1109/CVPR42600.2020.00186>
- Lior Yariv, Peter Hedman, Christian Reiser, Dor Verbin, Pratul P. Srinivasan, Richard Szeliski, Jonathan T. Barron, and Ben Mildenhall. 2023. BakedSDF: Meshing Neural SDFs for Real-Time View Synthesis. In *ACM SIGGRAPH 2023 Conference Proceedings (Los Angeles, CA, USA) (SIGGRAPH '23)*. Association for Computing Machinery, New York, NY, USA, Article 46, 9 pages. <https://doi.org/10.1145/3588432.3591536>
- Yizhou Yu, Kun Zhou, Dong Xu, Xiaohan Shi, Hujun Bao, Baining Guo, and Heung-Yeung Shum. 2004. Mesh Editing with Poisson-Based Gradient Field Manipulation. In *ACM SIGGRAPH 2004 Papers (Los Angeles, California) (SIGGRAPH '04)*. Association for Computing Machinery, New York, NY, USA, 644–651. <https://doi.org/10.1145/1186562.1015774>
- Jiakai Zhang, Xinhang Liu, Xinyi Ye, Fuqiang Zhao, Yanshun Zhang, Minye Wu, Yingliang Zhang, Lan Xu, and Jingyi Yu. 2021. Editable Free-Viewpoint Video Using a Layered Neural Representation. *ACM Trans. Graph.* 40, 4, Article 149 (jul 2021), 18 pages. <https://doi.org/10.1145/3450626.3459756>
- Yunzhi Zhang, Shangzhe Wu, Noah Snavely, and Jiajun Wu. 2023. Seeing a Rose in Five Thousand Ways. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. 962–971.
- Xiaoyu Zhou, Zhiwei Lin, Xiaojun Shan, Yongtao Wang, Deqing Sun, and Ming-Hsuan Yang. 2023. DrivingGaussian: Composite Gaussian Splatting for Surrounding Dynamic Autonomous Driving Scenes. arXiv:2312.07920 [cs.CV]



Fig. 13. Showcase in 3D with background. To demonstrate their natural appearance, we insert these composite models back into their unbounded backgrounds (the floaters are caused by the problem of 3DGS under unbounded scenes).

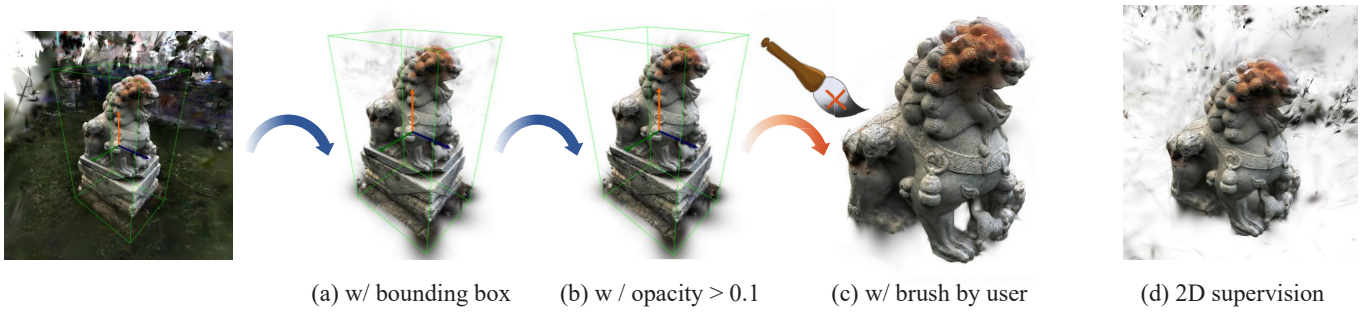


Fig. 14. We describe the segmentation workflow using our GUI and compare it to the result (d) from 2D mask supervision (for example, the Segment Anything Model (SAM) [Kirillov et al. 2023]). To segment with SAM, we re-implement the inverse-mask [Cen et al. 2023] strategy on 3DGS. A simple (a) bounding box with a (b) interactive (c) brush is demonstrated to be more practical in real-world scenes with numerous floaters. For more information, please refer to our supplementary video.

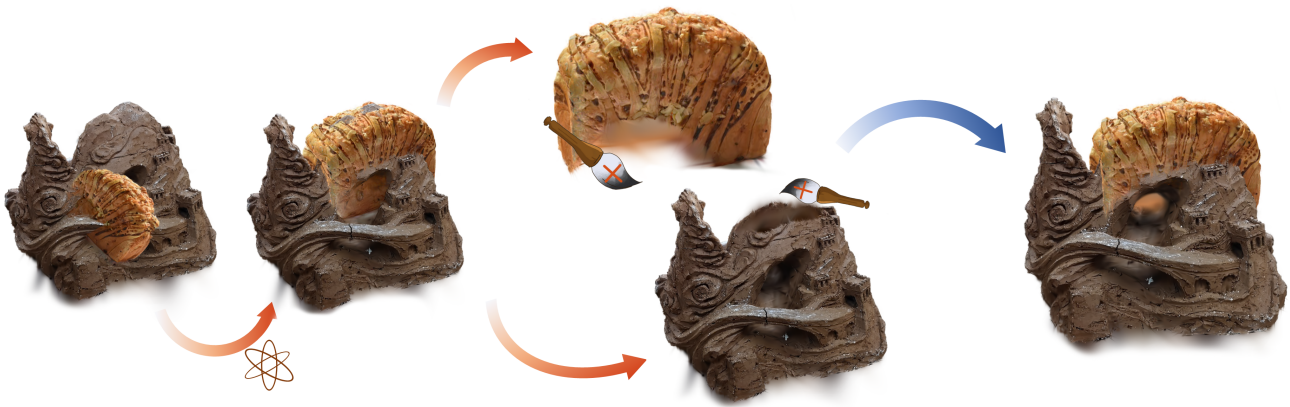


Fig. 15. We describe the transformation workflow using our GUI, as well as how to remove unwanted parts during composition. Users can adjust models to create a semantically meaningful composite, and then use a brush to remove unwanted parts, allowing for a more fine-grained composition. For more information, please refer to our supplementary video.

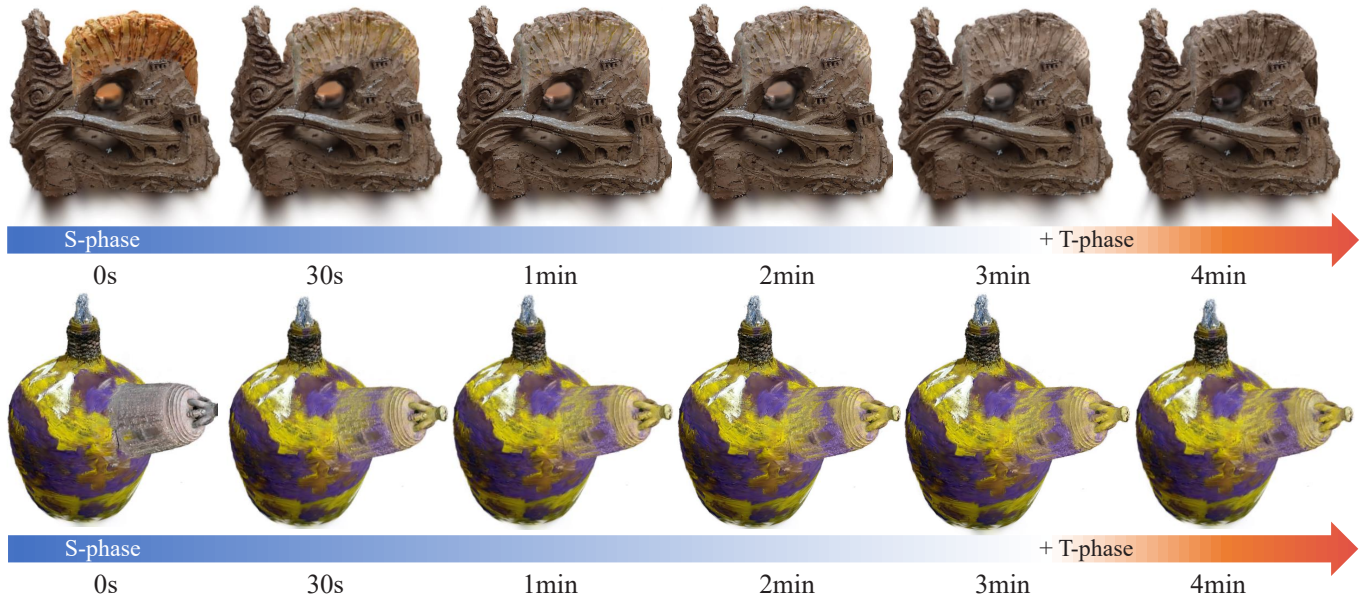


Fig. 16. Visualize how our optimization gradually and efficiently converges. In our comparison, SNeRF [Nguyen-Phuoc et al. 2022] takes over 10 hours, while SeamlessNeRF [Gong et al. 2023] takes more than an hour. For more information, please refer to our supplementary video.



Fig. 17. We demonstrate another compositing result using data from both the real world and the graphics engine. This additional case demonstrates our approach's versatility in dealing with both real and computer-generated models, validating its practical applicability. The two CG models were obtained from websites and rendered in Blender 3D on our own.