

Point'n Move: Interactive Scene Object Manipulation on Gaussian Splatting Radiance Fields

Jiajun Huang
Bournemouth University
jhuang@bournemouth.ac.uk

Hongchuan Yu
Bournemouth University
hyu@bournemouth.ac.uk

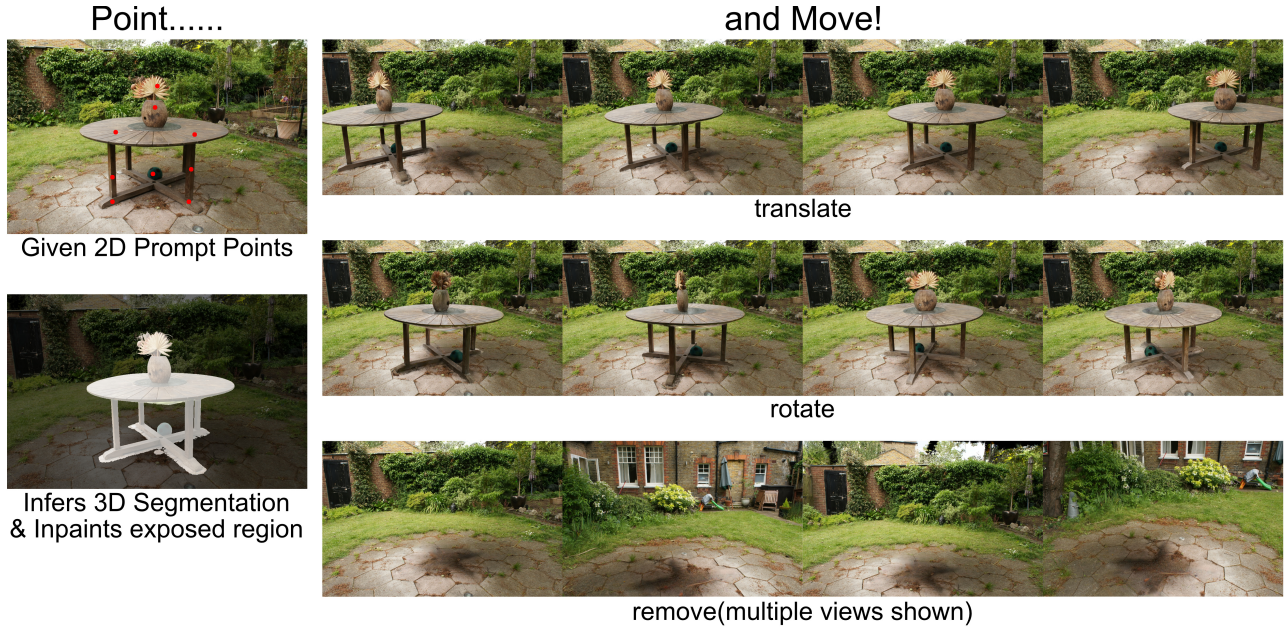


Figure 1. Capability highlight of our method. The user selects an object in the scene via 2D point prompts, which our method uses to infer a 3D segmentation mask. After segmentation and exposed region inpainting, the user can freely manipulate the selected object in the scene in real-time, with all exposed regions or holes inpainted. Please refer to the supplementary material for a video demo on real-time editing.

Abstract

We propose *Point'n Move*, a method that achieves interactive scene object manipulation with exposed region inpainting. Interactivity here further comes from intuitive object selection and real-time editing. To achieve this, we adopt *Gaussian Splatting Radiance Field* as the scene representation and fully leverage its explicit nature and speed advantage. Its explicit representation formulation allows us to devise a 2D prompt points to 3D mask dual-stage self-prompting segmentation algorithm, perform mask refinement and merging, minimize change as well as provide good initialization for scene inpainting and perform editing in real-time without per-editing training, all leads to superior quality and performance. We test our method by performing editing on both forward-facing and 360 scenes.

We also compare our method against existing scene object removal methods, showing superior quality despite being more capable and having a speed advantage.

1. Introduction

Radiance field-based scene representations, e.g., Neural Radiance Fields(NeRFs)[10], have achieved remarkable quality in capturing and representing real-life scenes. After capturing, there's a growing need to intuitively manipulate and rearrange objects within these scenes in a user-friendly manner. This enhanced editability holds immense promise for diverse applications, ranging from virtual home furnishing trials to effortless object removal from scenes and the production of VR/AR environments.

However, existing methods for editing are primarily con-

cerned with deformation and object-centric editing, rather than freely moving or rotating objects in scenes. Methods that produce directly editable representations, such as Control-NeRF[8], NeuralEditor[3] and NeuMesh[4] focus more on editing and does not provide ways to intuitively select an object from a captured scene. The real-time responsiveness poses a challenge, impeding interactive manipulation. Additionally, these methods do not inpaint the exposed surfaces or holes resulting from editing, yielding unrealistic rendering results.

In this paper, we introduce Gaussian Splatting Radiance Field(3DGS)[6] as the base scene representation, and present Point'n Move, a novel method for real-time interactive scene object manipulation, which includes intuitive object selection, 3D semantic segmentation, and inpainting (see Figure 1). Unlike existing neural implicit representations (e.g. NeRF[10]), 3DGS explicitly represents a radiance field using many 3D anisotropic balls, achieving rapid training and real-time rendering with impressive quality. Such point cloud-like formulation allows our method to perform segmentation and refinement from a point cloud perspective. It also serves as the foundation of our scene-content revealing pruning strategy and reprojection-based initialization, improving quality in exposed region inpainting. The speed of 3DGS also gives us a speed advantage throughout our method, especially in 3D scene inpainting. Finally, this approach enables us to directly manipulate the selected primitives in the scene rather than indirect fine-tuning, therefore achieving real-time editing.

We demonstrate our method's effectiveness by selecting and editing objects in 360 and forward-facing scenes. We also benchmark our method against existing object removal methods, achieving competitive results in quality despite being able to perform full manipulation rather than just removal and a speed advantage in scene inpainting. The code of our method will be published upon paper acceptance.

In summary, our contributions are as follows:

- We propose Point'n Move, the first end-to-end method that achieves interactive scene object manipulation with exposed region inpainting on Gaussian Splatting Radiance Fields(3DGS). More concretely, our method leverages the rapid and explicit nature of 3DGS to enable intuitive selection, high-quality inpainting and real-time editing.
- For intuitive selection, we propose a dual-stage self-prompting mask propagation process that produces high-quality 3D semantic segmentation masks from 2D image prompt points.
- For high-quality exposed region inpainting, we propose a rapid inpainting procedure that minimizes unnecessary inpainting and a reprojection-based initialization scheme. Both contribute to high-quality results.
- Real-time editing is achieved by directly manipulating the

primitives in the scene representation, all without time-consuming fine-tuning.

2. Related Work

2.1. NeRF Editing

There has been a lot of literature on scene editing based on radiance field representations, including geometry editing[14][20][24] and object-centric editing[15][9][21]. They adopt many sophisticated techniques such as fine-tuning[15][9], rendering rays deformation[14][20][24], or editable representations[21].

However, most works prefer neural networks to represent radiance fields, and the implicit nature of this representation usually requires time-consuming fine-tuning for edit operations. As a result, these methods cannot support practical interactive editing use cases.

To tackle this challenge, another attempt is to represent radiance fields using explicit structures, including Control-NeRF[8] for volumetric grid, PAPR[26], NeuralEditor[3] and RIP-NeRF[17] for point clouds, NeuMesh[4] and Differentiable Blocks World[13] for meshes. The foundation of our method, Gaussian Splatting Radiance Field[6], also falls into this category.

The current body of work falls short in adequately supporting practical interactions, facing two challenges: intuitive 3D segmentation for object selection and inpainting of exposed regions or holes resulting from editing. Our method addresses these issues by leveraging the explicit representation of the Gaussian Splatting Radiance Field [6]. Building upon it, our method enables real-time, unrestricted editing.

2.2. NeRF Scene Object Removal

Scene object removal on NeRFs is another focus, aiming to select and remove objects from a trained neural radiance field. A pioneering work in object removal, SPIn-NeRF[12] achieves 3D segmentation by combining interactive segmentation methods, video segmentation methods, and NeRF-based semantic mask generation method that creates object removal masks. OR-NeRF[23] further simplifies this process via reprojecting the masks from the segmented 2D views onto new 2D views and then applying the Segment Anything Model[7] to perform segmentation on new views. Despite achieving good results, SPIn-NeRF and OR-NeRF both need to retrain NeRF after removal, which is not suitable for real-time unrestricted editing.

Inspired by SA3D[2], our approach incorporates a cross-view segmentation training process. Leveraging the explicit characteristics of the Gaussian Splatting Radiance Field[6], we realize explicit 3D segmentation through weighted point clouds, providing advantages for subsequent refinement and editing.

For inpainting exposed areas, both SPIn-NeRF and OR-NeRF introduced inpainting via 2D images. The scene NeRF associated with these images is updated by fine-tuning. A patch-wise perceptual loss in the masked regions is employed for multiview consistent effect from multiple inconsistent 2D inpaintings. Apart from these two approaches, [11] proposed an incremental scene inpainting and view-dependent effect estimation scheme for multiview-consistent effect. [19] also proposed a novel confidence-score-based scheme to filter out view-inconsistent inpainted images, and a re-training process to update the confidence scores. Despite achieving impressive results, the complexity of these methods leads to a long training time.

Our method follows the line of SPIn-NeRF and OR-NeRF. However, we employ a simpler perceptual loss that compares the entire area inside the bounding box that bounds the mask rather than small patches inside, reducing the number of views to back-propagate through. We also tap into the superior speed of Gaussian Splatting Radiance Field, enabled by a scene content revealing pruning scheme and reprojection-based initialization process to ensure good quality.

3. Method

Given a trained Gaussian Splatting Radiance Field[6] R , along with a set of cameras C in the scene, and some 2D point annotations P that selects an object in an image rendered from R at camera pose $C_0 \in C$, our method aims to achieve manipulation(translate/rotate/remove) of the object in real-time, with the newly exposed regions or holes properly inpainted.

Our proposed approach achieves this in three steps: segmentation, inpainting and recomposition(see Figure 2). Since our method heavily leverages its explicit formulation, we start by briefing Gaussian Splatting Radiance Fields.

3.1. Background: Gaussian Splatting Radiance Field

Gaussian Splatting Radiance Field[6], also referred to as 3D Gaussian Splatting(3DGS), is an explicit radiance field-based scene representation that represents a radiance field using a large number of 3D anisotropic balls, each modelled using a 3D gaussian distribution(Equation (1)). More concretely, each anisotropic ball has mean $\mathcal{M} \in \mathbb{R}^3$, covariance Σ , opacity $\alpha \in \mathbb{R}$ and spherical harmonics parameters $\mathcal{C} \in \mathbb{R}^k$ (k is the degrees of freedom) for modelling view-dependent colour. For regularizing optimization, the covariance matrix is further decomposed into rotation matrix \mathbf{R} and scaling matrix \mathbf{S} by Equation (2). These matrices are further represented as quaternions $r \in \mathbb{R}^4$ and scaling factor $s \in \mathbb{R}^3$.

$$G(X) = e^{-\frac{1}{2}\mathcal{M}^T\Sigma^{-1}\mathcal{M}} \quad (1)$$

$$\Sigma = \mathbf{R}\mathbf{S}\mathbf{S}^T\mathbf{R}^T \quad (2)$$

For this scene representation, view rendering is performed via point splatting[22]. Specifically, all gaussian balls in the scene are first projected onto the 2D image plane, and their colour is computed from spherical harmonic parameters. Then, for every 16x16 pixel patch of the final image, the projected gaussians that intersect with the patch are sorted by depth. For every pixel in the patch, its colour is computed by alpha compositing the opacity and colour of all the gaussians covering this pixel by depth order, as in Equation (3).

$$C = \sum_{i \in N_{cov}} c_i \alpha_i \prod_{j=1}^{i-1} (1 - \alpha_j) \quad (3)$$

where, N_{cov} represents the splats that cover this pixel, α_i represents the opacity of this gaussian splat multiplied by the density of the projected 2D Gaussian distribution at the location of the pixel, and c_i represents the computed colour.

Gaussian Splatting Radiance Field achieves a significant advantage in training and inference speed due to its explicit approach and fast splatting rendering process. In addition to leveraging its superior speed, its point-cloud-like representation formulation enables us to perform processing from a point-cloud perspective, improving speed and quality.

3.2. Object Segmentation

While our method maintains the division of a scene into two parts, namely the object R_{obj} and the background R_{bg} , we devise a dual-stage segmentation process for implementation, as shown in Figure 3.

Coarse Stage Segmentation To achieve high-quality 3D segmentation, our method builds upon the cross-view self-prompting process proposed by SA3D[2]. We augment the scene representation for segmentation by adding and training a segmentation score attribute to all the gaussian balls. Such score could be used to render 2D segmentation maps by treating the score as colour and rendering via the splatting process. Following SA3D, each training step starts by rendering images and segmentation masks from a chosen camera pose using the augmented representation, passing them through a heuristic algorithm to extract prompt points, and feed the prompts as well as the rendered image to Segment Anything Model(SAM)[7] to produce a more accurate pseudo ground truth mask. The pseudo-ground truth is then compared to the rendered mask for loss. The loss is presented in Equation (4), where M_{sam} is the pseudo ground truth mask, M_r is the rendered mask. Herein L_m is namely the Equation 5 in [2].

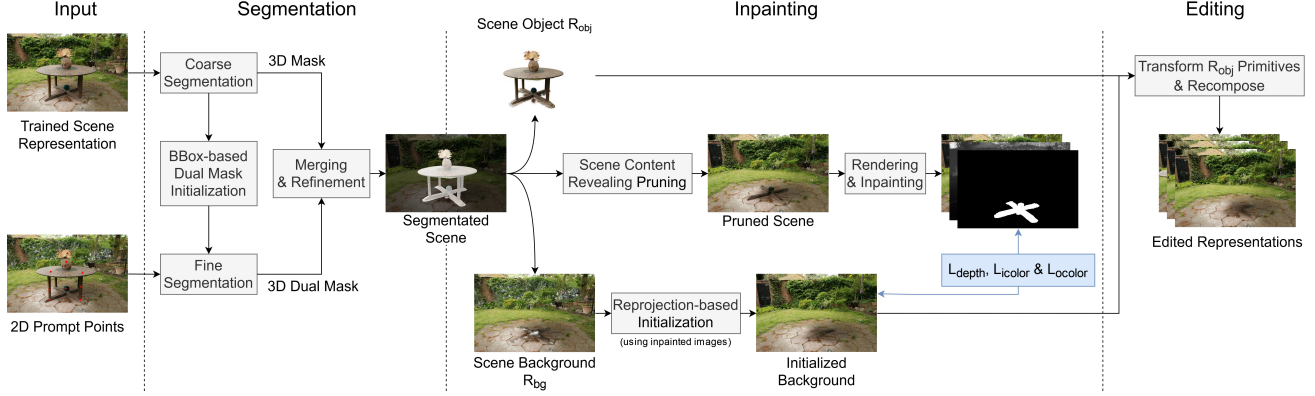


Figure 2. Overview of our method. Our method comprises of three stages: segmentation, inpainting and editing. We start with a dual-stage self-prompting cross-view mask propagation process to produce a 3D segmentation of the scene, splitting it into R_{obj} and R_{bg} . The segmentation is then used to derive 2D scene inpaintings with a scene-content revealing pruning strategy. After using 2D inpaintings to fine-tune R_{bg} , R_{obj} can be edited and recomposited into R_{bg} in real-time.

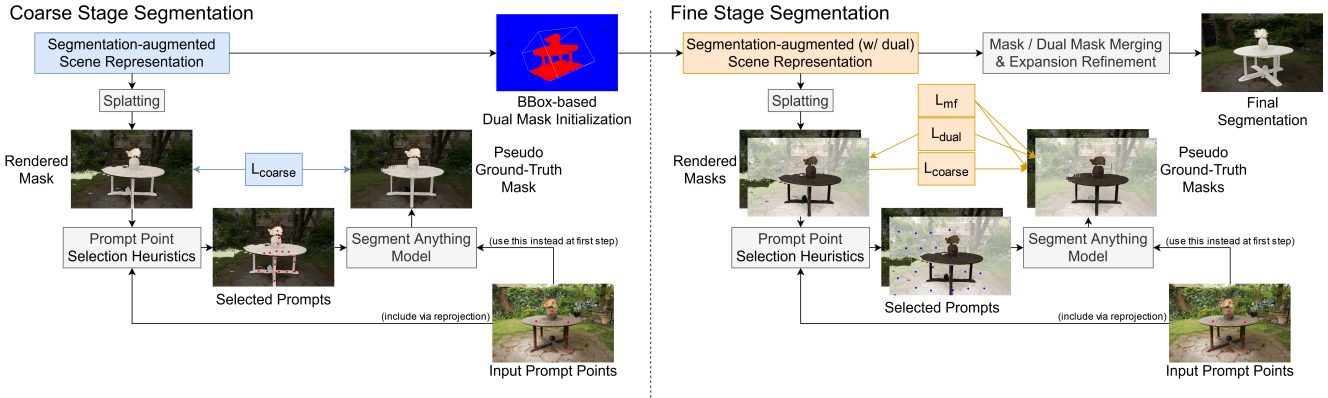


Figure 3. Detailed view of the segmentation stage. The cross-view transfer process involves selecting prompt points from rendered masks, which are used to produce pseudo-ground truths to supervise the rendering mask. The process is started using the input prompt points. After the coarse stage, a dual mask for content outside the selection and additional loss terms for training are added. Finally, the scores and dual scores are merged into the final segmentation.

$$L_{coarse} = L_m(M_{sam}, M_r) \quad (4)$$

Fine Stage Segmentation To improve segmentation quality, we expand the process by adding a fine stage and a merging process. After iterating over all the provided camera poses C , we enter the fine stage of segmentation. In the same way as adding segmentation scores, we further augment the representation with dual segmentation scores, aiming to capture scene contents that should not be selected. We then initialize the dual scores by computing a 3D bounding box containing all gaussians with high scores in the first stage. Every gaussian outside the bounding box has its dual score set close to one, while gaussians inside the bounding box have their dual score set to zero. We then continue the process for another full iteration, training the segmentation score and the newly added dual score, with the dual score

attribute following the same self-prompting process and a new loss. The loss terms for the fine stage are presented in Equation (5). L_{mf} (Equation (6)) is the loss for dual mask, where M_{sd} is the pseudo ground truth dual mask, M_{rd} is the rendered dual mask. We also add an additional term L_{dual} (Equation (7)) that encourages mask and dual mask to have no intersections.

$$L_{fine} = L_{coarse} + L_{mf} + L_{dual} \quad (5)$$

$$L_{mf} = L_m(M_{sd}, M_{rd}) \quad (6)$$

$$L_{dual} = L_m(-M_r, M_{rd}) + \lambda_{dd} L_m(-M_{rd}, M_r) \quad (7)$$

Merging and Refinement Finally, we merge the scores and dual scores to create the final segmentation. We com-

pute the 3D bounding box in the same way as in the initialization process. Then, in addition to the gaussians with a high enough segmentation score, we also include every Gaussian in the bounding box whose dual score is below a threshold into the final segmentation, effectively accepting those "rejected" by dual mask. This rejection inclusion scheme aims to cover all gaussians that are difficult to train, such as those hiding under the surface of the object. We also expand the selection by including gaussians whose mean is close enough to the already selected gaussians.

Remark Thanks to the explicit representation of the Gaussian Splatting Radiance Field, we can efficiently implement the fine-stage initialization scheme, merging process, and expansion refinement. The resultant 3D segmentation mask, after merging and expansion, serves as our final output, effectively delineating the gaussians and thereby segmenting the scene into R_{obj} and R_{bg} .

3.3. Exposed Surface Inpainting

After splitting, R_{bg} will likely contain newly exposed surface regions or holes. For better rendering results after reintegrating the manipulated R_{obj} , these defects should be inpainted. We achieve this by inpainting rendered images and depth maps and then fine-tuning the scene representation with them.

2D inpainting To minimize the amount of inpainting needed, we remove the masked gaussians far away from the not masked gaussians in R . This removes all selected gaussians except those in close contact with the rest of the scene, which are very likely where the exposed regions are. These regions can be identified by rendering a segmentation mask using the pruned representation. We also add the exposed holes to the rendered mask, which is computed by including new pixels with low total opacity or colour values close to the background colour after performing the removal above. Finally, we refine the acquired inpainting masks before inpainting. (Please refer to the supplementary material for details.)

This scene content revealing pruning strategy allows us to minimize the amount of inpainting and preserve the existing scene contents as much as possible, as presented in the Pruned Scene image of Figure 2.

We render RGB images and depth maps using this pruned representation and generate inpainting masks using the abovementioned method. The RGB images and depth maps are then inpainted by a 2D inpainter.

3D rapid Fine-tuning With the inpainted images I_i and depth maps D_i , we then perform fine-tuning to inpaint R_{bg} . We initialize by reprojecting the masked region of an inpainted depth map and its associated image back into the representation as new gaussians. As presented in the original 3DGS paper[6], good initialization is crucial for high-quality training results.

We then train using the losses below:

outside mask color loss We supervise the color of regions outside the mask via a weighted sum of pixel L1 and SSIM[18]:

$$L_{ocolor} = (1 - \lambda_{ssim})L1(I_i(1 - M_i), I_r(1 - M_i)) + \lambda_{ssim}SSIM(I_i(1 - M_i), I_r(1 - M_i)) \quad (8)$$

where I_i and M_i are the inpainted image and mask, I_r is the rendered image of the representation, L1 stands for mean pixel L1 loss, λ_{ssim} is the weight for SSIM.

depth loss where for scene geometry, we employ depth map L1:

$$L_{idepth} = \lambda_{depth}L1(D_i, D_r) \quad (9)$$

where D_i is the inpainted depth map, D_r is the rendered depth map, L1 stands for mean pixel L1 loss, λ_{depth} is the weighting factor.

inside mask color loss where for color inside the masked region, we adopt a perceptual color loss:

$$L_{icolor} = \lambda_{lrips}LPIPS(\text{Box}(I_i, M_i), \text{Box}(I_r, M_i)) \quad (10)$$

where I_i and M_i are the inpainted image and mask, I_r is the rendered colour image, Box stands for a function that computes the bounding box of the mask and only keeps parts of the image in the bounding box. LPIPS is the LPIPS[25] perceptual metric.

Here, we employ a perceptual colour loss instead of directly comparing pixel values. This is because the 2D inpaintings are not guaranteed to be multiview consistent and a strict loss could harm quality. We also filter via the bounding box of the masked region instead of strictly in mask, this is to encourage better integration between the masked region and its surroundings.

3.4. Editing and Recomposition

With R_{bg} inpainted, R_{obj} can then be freely manipulated and composited back into R_{bg} in real-time for editing.

The selected object R_{obj} could be translated and rotated by transforming the position(or mean) \mathcal{M} and rotation r of the underlying gaussians. Recomposition is done by adding the gaussians in R_{obj} back into R_{bg} . Note that as these transformations could be done by multiplying transform matrices, no further training is needed. All possible surfaces and holes that could be exposed by editing have already been inpainted in the previous step.

3.5. Implementation Details

Both stages of the segmentation training are done using the SGD optimizer, with a learning rate of 1.0. λ_{dd} is 0.1. For

image inpainting, we employ state-of-the-art 2D image inpainting model Lama[16] for inpainting the color image and depth map. The scene inpainting training process also uses the SGD optimizer with learning rates equal to what is specified in the original 3D Gaussian Splatting paper[6]. λ_{ssim} is 0.2, λ_{depth} is 1.0, λ_{lips} is 1.0. All experiments are conducted using a single A5000.

4. Experiments

We demonstrate the effectiveness of our method by testing it on both 360 and forward-facing scene datasets. Our method is also compared to existing scene object removal methods, reporting competitive performance in segmentation and inpainting, despite being able to do more than just removal and has a speed advantage. Finally, we conclude with ablation studies for our key contributions.

4.1. Experiment Setups

Datasets For diversity in evaluation, we test our method on the MipNeRF360[1] dataset for 360 scenes and the SPIn-NeRF[12] dataset for forward-facing scenes. Both datasets consist of captured images of a scene with their associated camera parameters, which can be used to train the model to edit. As a dataset for evaluating object removal methods, the SPIn-NeRF dataset also contains scene images without the objects to remove and object segmentation masks for all captured images. For input prompt points, we make use of the point annotations from OR-NeRF[23].

Metrics For segmentation, we report the mean accuracy and IoU(Intersection over Union) score of the rendered 2D segmentation masks across all images. For the image quality of inpainted scenes, we report the mean PSNR and LPIPS[25] score between the rendered image of the inpainted scene representation and the ground truth scene image where the object is removed. We also calculate the Frechet inception distance (FID)[5] between all rendered images and ground truth images taken without the target object for image quality.

Baselines To evaluate the segmentation and inpainting ability of our method, we compare our approach against the state-of-the-art methods in 3D segmentation and scene object removal: SPIn-NeRF[12] and OR-NeRF[23]. More concretely, we compare against the TensorRF variant of OR-NeRF, which is the fastest and the best overall result quality variant presented in the paper.

4.2. Object Manipulation on 360 and Forward-facing Scenes

To demonstrate the effectiveness of our approach, we perform manipulation with our method on scenes from the forward-facing SPIn-NeRF dataset and the 360 MipNeRF360 dataset. The qualitative results are presented in Figure 4. It can be noted that our method can perform

high-quality selection and editing on various scenes, with exposed regions properly inpainted. Our method produces sensible results even for cases with extreme input, such as the scene at last row of Figure 4.

4.3. Comparing Against Object Removal Methods

We compare our model against the object removal baselines, which also support intuitive object segmentation. However, they can only remove the object, unable to move or rotate it. The experiments are conducted on the SPIn-NeRF dataset for the presence of ground-truth data.

Segmentation Quality For segmentation quality, the qualitative results are presented in Figure 5 and the quantitative results are presented in Table 1. The table shows that our method only achieves competitive performance against OR-NeRF. We attribute this to the fact that OR-NeRF directly operates on 2D images for segmentation and creates 2D masks instead of 3D segmentation. On the other hand, our method achieves true 3D segmentation and is rendered to 2D masks for comparison. The performance of our method is on par with SPIn-NeRF[12], a 3D segmentation method.

Scene Inpainting Quality For inpainting quality, the qualitative results are presented in Figure 6 and the quantitative results in Table 1. As can be seen, our methods produce more plausible results compared to existing methods. We attribute this to our adopted representation’s superior representation capability and our proposed initialisation and fine-tuning scheme. Our method also has a significant speed advantage in 3D scene inpainting. As shown in the last column of Table 1, our method is three times faster than the fastest baseline despite achieving superior result quality.

4.4. Ablation Studies

Dual Stage Segmentation To validate the importance of dual-stage segmentation, we disable the fine stage in our segmentation process and compare by evaluating images rendered from R_{bg} . As shown in Figure 7, disabling the fine stage could leave floaters in the scene that would be fixed by the fine stage and merging process.

Content-Revealing Scene Pruning We validate the importance of scene content revealing pruning by checking the inpainted images of the garden scene that has or has not gone through the pruning process before rendering masks and images for inpainting. As presented in Figure 8, the pruning scheme reveals contents behind the table instead of relying on inpainting, thus improving output quality.

Inpainting Optimization Initialization We analyze the effect of initialization by disabling it and directly training the pruned scene. As shown in Figure 9, without proper initialization, the fine-tuning process would leave a cloud of tiny floating gaussian balls, instead of completely optimizing them away.



Figure 4. Input prompts, rendered segmentation mask and editing results of our method on 360 and forward-facing scenes. Note that our method gives sensible results even on scenes with extreme prompts(see last row). Please refer to the supplementary material for a video demo on real-time editing.



Figure 5. Qualitative comparison on segmentation against object removal methods on the SPIn-NeRF dataset. We compare by rendering our 3D segmentation mask into 2D mask images. As can be seen, our method has comparable performance against OR-NeRF(which is superior to SPIn-NeRF in segmentation), despite producing 3D segmentation instead of 2D masks on images.

5. Conclusion

This paper introduces Point’n Move, a methodology enabling interactive manipulation of scene objects coupled with exposed region inpainting. Leveraging the explicit, point cloud-like formulation and speed advantages offered by the Gaussian Splatting Radiance Field[6] as the fundamental framework of our design, our approach achieves intuitive object selection, high-fidelity exposed region inpainting, and real-time editing. These facets collectively deliver a user-friendly interactive editing experience characterized by high-quality results.

Limitations Currently, our method does not handle lighting or texture and focuses solely on geometry edit-

ing. Also, for some scenes, the inpaintings produced by our method are darker than expected, as presented in Figure 10. We attribute this to the precision of our segmentation as it does not include the shadow of the object, which misleads the 2D inpainter to inpaint with the shadow, making the output darker.

References

- [1] Jonathan T. Barron, Ben Mildenhall, Dor Verbin, Pratul P. Srinivasan, and Peter Hedman. Mip-nerf 360: Unbounded anti-aliased neural radiance fields. *CVPR*, 2022. 6
- [2] Jiazhong Cen, Zanwei Zhou, Jiemin Fang, Wei Shen, Lingxi Xie, Xiaopeng Zhang, and Qi Tian. Segment anything in 3d with nerfs. *arXiv:2304.12308*, 2023. 2, 3

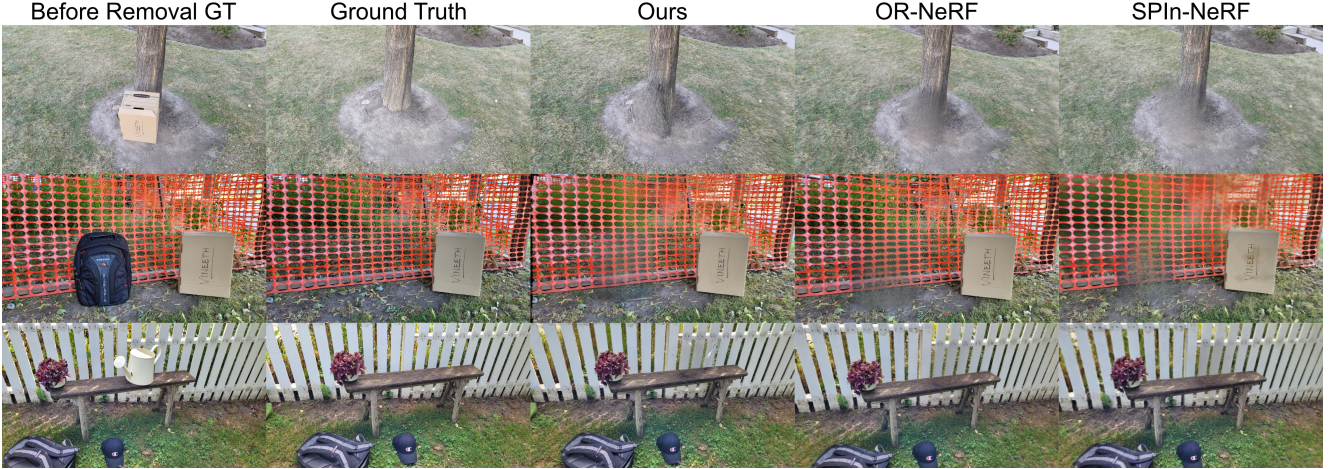


Figure 6. Qualitative comparison on exposed region inpainting on the SPIn-NeRF dataset. We present the ground truth with the object, without the object and the output of all the methods. As can be seen, our method is comparable, if not superior in terms of inpainting quality on the exposed regions.

method	Segmentation Metrics		Inpainting Metrics			
	Acc(\uparrow)	IoU(\uparrow)	PSNR(\uparrow)	FID(\downarrow)	LPIPS(\downarrow)	Inpainting Time(\downarrow)
Ours	99.51	92.71	19.83	40.33	0.2684	20.41 mins
OR-NeRF	99.71	95.42	13.94	49.91	0.6162	168.97 mins
SPIn-NeRF	98.91	91.66	14.83	67.26	0.6506	58.25 mins

Table 1. Quantitative results for comparing against object removal methods. The best value of each column is bolded for ease of reading. For segmentation, our method is slightly inferior compared to OR-NeRF and is on par with SPIn-NeRF. However, compared to OR-NeRF, we produce 3D segmentation rather than 2D segmentation. For scene inpainting, our method achieves superior performance both in terms of quality and speed. Please refer to the supplementary material for detailed results.



Figure 7. R_{bg} part of segmentation, performed with(left) or without(right) the fine stage. As could be seen, disabling fine stage introduces floaters that should be included in R_{obj} .



Figure 9. Close-up view of the output of scene fine-tuning, performed with(left) or without(right) reprojection-based initialization. Please note the small residual grains present on the right.



Figure 8. 2D inpainting output, performed with(left) or without(right) content-revealing pruning. Note that pruning removes the blur caused by excessive inpainting in the red box region.



Figure 10. Artifacts caused by residual shadows misleading the inpainter. the shadows of the object(shadows in the red box on the left) could mislead the 2D inpainter, producing dark results.

- [3] Jun-Kun Chen, Jipeng Lyu, and Yu-Xiong Wang. NeuralEditor: Editing neural radiance fields via manipulating point clouds. In *CVPR*, 2023. 2
- [4] Chong Bao and Bangbang Yang, Zeng Junyi, Bao Hujun, Zhang Yinda, Cui Zhaopeng, and Zhang Guofeng. Neumesh: Learning disentangled neural mesh-based implicit field for geometry and texture editing. In *European Conference on Computer Vision (ECCV)*, 2022. 2
- [5] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. GANs trained by a two time-scale update rule converge to a local nash equilibrium. In *Proceedings of the 31st International Conference on Neural Information Processing Systems*, pages 6629–6640, 2017. 6
- [6] Bernhard Kerbl, Georgios Kopanas, Thomas Leimkühler, and George Drettakis. 3d gaussian splatting for real-time radiance field rendering. *ACM Transactions on Graphics*, 42(4), 2023. 2, 3, 5, 6, 7
- [7] 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. Segment Anything, 2023. 2, 3
- [8] Verica Lazova, Vladimir Guzov, Kyle Olszewski, S. Tulyakov, and Gerard Pons-Moll. Control-nerf: Editable feature volumes for scene rendering and manipulation. *2023 IEEE/CVF Winter Conference on Applications of Computer Vision (WACV)*, pages 4329–4339, 2022. 2
- [9] Steven Liu, Xiuming Zhang, Zhoutong Zhang, Richard Zhang, Jun-Yan Zhu, and Bryan Russell. Editing conditional radiance fields, 2021. 2
- [10] Ben Mildenhall, Pratul P. Srinivasan, Matthew Tancik, Jonathan T. Barron, Ravi Ramamoorthi, and Ren Ng. NeRF: Representing Scenes as Neural Radiance Fields for View Synthesis. In *Computer Vision – ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part I*, pages 405–421, 2020. 1, 2
- [11] Ashkan Mirzaei, Tristan Aumentado-Armstrong, Marcus A. Brubaker, Jonathan Kelly, Alex Levinshtein, Konstantinos G. Derpanis, and Igor Gilitschenski. Reference-guided controllable inpainting of neural radiance fields. In *ICCV*, 2023. 3
- [12] Ashkan Mirzaei, Tristan Aumentado-Armstrong, Konstantinos G. Derpanis, Jonathan Kelly, Marcus A. Brubaker, Igor Gilitschenski, and Alex Levinshtein. SPIn-NeRF: Multiview segmentation and perceptual inpainting with neural radiance fields. In *CVPR*, 2023. 2, 6, 1
- [13] Tom Monnier, Jake Austin, Angjoo Kanazawa, Alexei A. Efros, and Mathieu Aubry. Differentiable blocks world: Qualitative 3d decomposition by rendering primitives. *ArXiv*, abs/2307.05473, 2023. 2
- [14] Yicong Peng, Yichao Yan, Shengqi Liu, Yuhao Cheng, Shanyan Guan, Bowen Pan, Guangtao Zhai, and Xiaokang Yang. CageNeRF: Cage-based Neural Radiance Field for Generalized 3D Deformation and Animation. In *Advances in Neural Information Processing Systems*, pages 31402–31415, 2022. 2
- [15] Konstantinos Rematas, Ricardo Martin-Brualla, and Vittorio Ferrari. Sharf: Shape-conditioned Radiance Fields from a Single View. In *Proceedings of the 38th International Conference on Machine Learning*, pages 8948–8958, 2021. 2
- [16] Roman Suvorov, Elizaveta Logacheva, Anton Mashikhin, Anastasia Remizova, Arsenii Ashukha, Aleksei Silvestrov, Naejin Kong, Harshith Goka, Kiwoong Park, and Victor Lempitsky. Resolution-robust Large Mask Inpainting with Fourier Convolutions. In *2022 IEEE/CVF Winter Conference on Applications of Computer Vision (WACV)*, pages 3172–3182, 2022. 6
- [17] Yuze Wang, Junyi Wang, Yansong Qu, and Yue Qi. Rip-nerf: Learning rotation-invariant point-based neural radiance field for fine-grained editing and compositing. *Proceedings of the 2023 ACM International Conference on Multimedia Retrieval*, 2023. 2
- [18] Zhou Wang, Alan Conrad Bovik, Hamid R. Sheikh, and Eero P. Simoncelli. Image quality assessment: from error visibility to structural similarity. *IEEE Transactions on Image Processing*, 13:600–612, 2004. 5
- [19] Silvan Weder, Guillermo Garcia-Hernando, Áron Monszpart, Marc Pollefeys, Gabriel J. Brostow, Michael Firman, and Sara Vicente. Removing objects from neural radiance fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 16528–16538, 2023. 3
- [20] Tianhan Xu and Tatsuya Harada. Deforming Radiance Fields with Cages. In *Computer Vision – ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part XXXIII*, pages 159–175, 2022. 2
- [21] Bangbang Yang, Yinda Zhang, Yinghao Xu, Yijin Li, Han Zhou, Hujun Bao, Guofeng Zhang, and Zhaopeng Cui. Learning Object-Compositional Neural Radiance Field for Editable Scene Rendering. In *2021 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 13759–13768, 2021. 2
- [22] Wang Yifan, Felice Serena, Shihao Wu, Cengiz Öztireli, and Olga Sorkine-Hornung. Differentiable surface splatting for point-based geometry processing. *ACM Transactions on Graphics (TOG)*, 38(6):1–14, 2019. 3
- [23] Youtan Yin, Zhoujie Fu, Fan Yang, and Guosheng Lin. Ornerf: Object removing from 3d scenes guided by multiview segmentation with neural radiance fields, 2023. 2, 6, 1
- [24] Yu-Jie Yuan, Yang-Tian Sun, Yu-Kun Lai, Yüewen Ma, Rongfei Jia, and Lin Gao. NeRF-Editing: Geometry Editing of Neural Radiance Fields. In *2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 18332–18343, 2022. 2
- [25] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *CVPR*, 2018. 5, 6
- [26] Yanshu Zhang*, Shichong Peng*, Alireza Moazeni, and Ke Li. Papr: Proximity attention point rendering. *Advances in Neural Information Processing Systems*, 2023. 2

Point’n Move: Interactive Scene Object Manipulation on Gaussian Splatting Radiance Fields

Supplementary Material

6. 2D Mask Refinement Algorithm

Our 2D mask refinement algorithm(see Algorithm 1) is inspired by the refinement procedure present in the implementation of OR-NeRF[23].

In essence, our method dilates the mask, finds the contour with the largest area and takes the region inside it as the refined mask. This is to eliminate potential holes in the mask. All functions invoked in Algorithm 1 can be implemented using functions from OpenCV.

Algorithm 1 2D Mask Refinement Algorithm

Require: 2D Inpainting Mask m

Ensure: Refined Inpainting Mask ret

$m \leftarrow \text{dilate}(m, \text{kernel_size} = 3, \text{iteration} = 3)$

$\text{contours} \leftarrow \text{findContours}(m)$

$\text{contour} \leftarrow \text{contourWithMaxArea}(\text{contours})$

$ret \leftarrow \text{filledContour}(\text{contour})$

return ret

7. Additional Per-Scene Metrics

Quantitative results presented in Table 1 is the mean across all scenes in the SPIn-NeRF[12] dataset. We present the quantitative result for each individual scene here. Please see Table 2 for segmentation metrics and Table 3 for inpainting metrics.

metric	method	2	3	4	7	10	12	book	trash	mean
Acc(\uparrow)	Ours	99.80	99.79	99.74	99.59	99.84	98.72	99.01	99.56	99.51
	OR-NeRF	99.82	99.73	99.79	99.81	99.87	99.30	99.51	99.51	99.71
	SPIIn-NeRF	-	-	-	-	-	-	-	-	98.91
IoU(\uparrow)	Ours	96.16	98.09	98.15	94.68	94.67	86.51	83.45	89.99	92.71
	OR-NeRF	96.47	97.48	98.50	97.43	95.47	91.73	88.68	88.68	95.42
	SPIIn-NeRF	-	-	-	-	-	-	-	-	91.66

Table 2. Quantitative per-scene result on segmentation quality. Results for SPIIn-NeRF and OR-NeRF are taken from their original papers. SPIIn-NeRF did not provide per-scene metrics; hence, they are not shown. The best value of each column is bolded for ease of reading. Our method is slightly inferior to OR-NeRF but has an advantage over SPIIn-NeRF. However, we produce 3D segmentation, while OR-NeRF performs segmentation on 2D images, which is not usable for unrestricted manipulations

metric	method	2	3	4	7	10	12	book	trash	mean
PSNR(\uparrow)	Ours	18.48	18.04	20.88	21.4	19.75	16.62	22.28	21.18	19.83
	OR-NeRF	15.94	11.42	13.02	14.37	12.89	11.40	15.88	16.63	13.94
	SPIIn-NeRF	16.69	12.08	14.90	15.34	12.73	12.39	17.84	16.70	14.83
FID(\downarrow)	Ours	53.60	36.39	51.78	22.48	21.67	26.23	81.68	28.84	40.33
	OR-NeRF	72.10	34.72	74.04	38.66	43.89	38.02	64.91	32.96	49.91
	SPIIn-NeRF	71.75	68.35	61.10	43.95	91.73	50.52	102.71	47.98	67.26
LPIPS(\downarrow)	Ours	0.4544	0.2217	0.3229	0.2858	0.2264	0.3352	0.2147	0.2301	0.2864
	OR-NeRF	0.7909	0.4937	0.6684	0.6445	0.6165	0.7179	0.5094	0.4882	0.6162
	SPIIn-NeRF	0.8489	0.5472	0.6815	0.6552	0.7003	0.7518	0.4226	0.5972	0.6506

Table 3. Quantitative per-scene result on scene inpainting quality. Results for SPIIn-NeRF and OR-NeRF are taken from their original papers. The best value of each column is bolded for ease of reading. Please note our method’s overall best performance across most scenes.