

Robust 3D Gaussian Splatting for Novel View Synthesis in Presence of Distractors

Paul UngermaNN, Armin Ettenhofer, Matthias Nießner, and Barbara Roessle

Technical University of Munich, Munich, Germany

{paul.ungermann, armin.ettenhofer, niessner, barbara.roessle}@tum.de

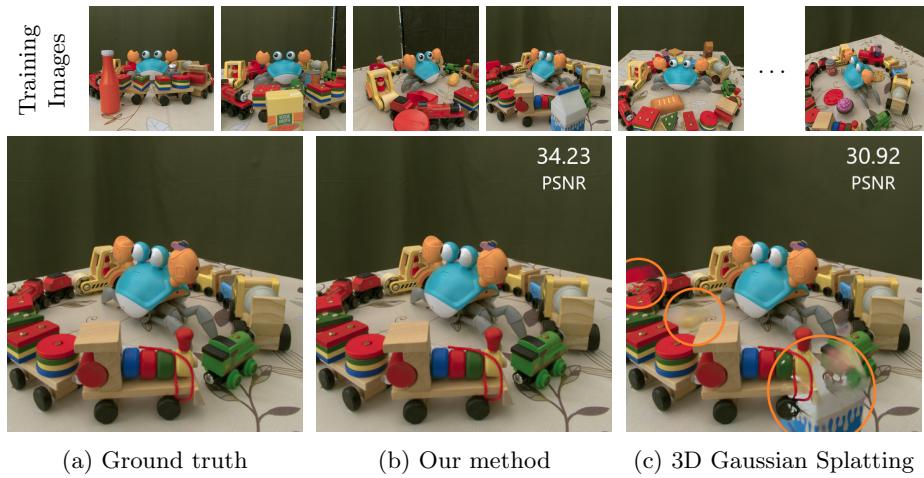


Fig. 1: Due to distractors in the scene 3D Gaussian Splatting creates floating artifacts in the image (highlighted with circles). Our method mitigates artifacts due to violations of the static scene assumption for Gaussian Splatting. As a key element to our approach, we optimize for semantic distractor masks simultaneous to the scene optimization, which allow us to effectively ignore distractors.

Abstract. 3D Gaussian Splatting has shown impressive novel view synthesis results; nonetheless, it is vulnerable to dynamic objects polluting the input data of an otherwise static scene, so called distractors. Distractors have severe impact on the rendering quality as they get represented as view-dependent effects or result in floating artifacts. Our goal is to identify and ignore such distractors during the 3D Gaussian optimization to obtain a clean reconstruction. To this end, we take a self-supervised approach that looks at the image residuals during the optimization to determine areas that have likely been falsified by a distractor. In addition, we leverage a pretrained segmentation network to provide object awareness, enabling more accurate exclusion of distractors. This way, we obtain segmentation masks of distractors to effectively ignore them in the loss formulation. We demonstrate that our approach is robust to various distractors and strongly improves rendering quality on distractor-polluted scenes, improving PSNR by 1.86dB compared to 3D Gaussian Splatting.

Keywords: 3D Gaussian Splatting · Robustness · Distractors.

1 Introduction

Neural Radiance Fields (NeRFs) [16] and 3D Gaussian Splatting [10] have shown remarkable improvements in novel view synthesis on complex scenes, enabling various tasks in the fields of virtual reality, autonomous systems, gaming or others. Given a set of input images along with camera poses, a 3D scene representation is optimized from which photo-realistic novel views can be rendered. NeRF represents the radiance and density distribution of a scene with a multi-layer perceptron (MLP) and employs differentiable volume rendering to synthesize novel views. In contrast, Gaussian Splatting represents the scene in an explicit manner, as a set of 3D Gaussians, defined by position, covariance, opacity, and spherical harmonic coefficients, which are efficiently rendered with a differentiable rasterizer, thereby enabling high-quality novel view synthesis in real-time.

3D Gaussian Splatting and NeRF are both optimized to minimize a re-rendering loss in RGB space. This procedure relies on a static scene assumption, i.e., the images must be photometrically consistent. In real-world scenarios, however, this assumption of a perfectly static scene is hardly ever fulfilled. Even with careful scene capture, the recordings often contain dynamics, such as lighting changes, moving shadows, or any unforeseeable moving objects or persons, e.g., tourists near a captured landmark. We refer to such undesired dynamic observations as distractors. Ignoring this problem severely degrades the optimized scene representation and the rendering quality, resulting in floating artifacts and blurriness (Fig. 2c). At the same time, the removal of distractors from a dataset as a post-processing step is non-trivial due to the variety of potential distractors. Clearly, manual pixel-wise annotation of distractors is impractical. Finally, for making Gaussian Splatting widely adopted and applicable to in-the-wild settings, the data capture has to remain simple with little effort. Therefore, it is highly desirable to increase the robustness of 3D Gaussian Splatting to distractors to be able to obtain clean reconstructions even from imperfect data. RobustNeRF [21] provides a solution for handling distractors in NeRF using a robust loss that computes distractor masks through iteratively reweighted least squares. Applying this approach to Gaussian Splatting, however, causes too aggressive masking and reduced performance (Section 4). Hence, we take a different approach and learn flexible neural decision boundaries to distinguish between distractors and static scene content.

Our work simultaneously optimizes for distractor masks to support the static scene reconstruction. To this end, we leverage image residuals in the training process and apply different transformations to obtain local smoothness and distractor contiguity. Using a neural classifier to classify distractor pixels, we compute a first semantic distractor mask. We refine these segmentation masks and establish object awareness using the object segmentation mask from SegmentAnything [11]. At the same time, our method remains independent from the type of distractors and can handle arbitrary distractors. We evaluate our approach on challenging distractor-polluted scenes [21] and obtain remarkable improvements over 3D Gaussian Splatting and RobustNeRF by 1.9dB and 4.3dB in PSNR, respectively.

In summary, we provide the following contributions:

- We introduce distractor masks by optimizing a neural decision boundary based on image residuals to effectively track and exclude distractors during 3D Gaussians optimization.
- We propose to leverage a pretrained segmentation network to enhance the distractor masks, making them object aware for more accurate exclusion of distractors.

2 Related Work

Neural Radiance Fields and 3D Gaussian Splatting Neural Radiance Fields (NeRF) [16] achieve outstanding novel view synthesis performance. The NeRF scene representation is realized as an implicit function, where a MLP maps a 3D position and viewing direction to radiance and density. Differentiable volume rendering combines radiance and density along target camera rays to produce pixel colors in the output view. The scene representation is optimized on a set of posed input images by minimizing a photometric loss. At inference time, novel views can be rendered from arbitrary view points. Follow-up works have extended NeRF in many directions, for instance towards alternative scene representations, e.g., voxel grids [6,14,24], decomposed tensors [2,22,3], or hash maps [17] to increase optimization and rendering speed. Various extensions also exist towards novel view synthesis on dynamic scenes [19,27,18,7], where typically a deformation field is optimized in addition to a canonical NeRF. A bit less explored, but still highly relevant is the application of NeRF to distractor-polluted scenes to which we dedicate an individual paragraph below.

3D Gaussian Splatting [10] is a recent method for novel view synthesis. In contrast to NeRF, a set of multivariate Gaussians with parameters such as position, covariance, opacity and appearance, is optimized as a scene representation. 3D Gaussian Splatting leverages a differentiable rasterizer that efficiently renders the Gaussians which allows real-time rendering, while at the same time outperforming NeRF in image quality metrics. Recent advances in Gaussian Splatting mainly focus on improving image quality in different distractor-free settings, namely dynamic scenes [30,35,34,13] and static scenes [36,15,33,9]. Other research directions tackle data efficiency using depth information [32,38,5]. Furthermore, Gaussian Splatting is also used to improve generative models like [26,4].

Handling Distractors. There are several approaches to handle distractors in general. In settings where distractors are known to be in specific classes we can employ a pretrained semantic segmentation model to remove distractors [20,25]. The main problem with this approach are distractors that do not belong to any known class such as shadows. Another approach is to exploit the time dependency of the images to classify static and dynamic (i.e., distractor) objects in the scene [31]. The problem with this method is that it requires time dependency, which typically is not available in the multi-view reconstruction setting.

Nonetheless, we take inspiration from the usage of pretrained semantic segmentation networks and leverage SegmentAnything [11] to provide object awareness in our distractor mask optimization, while at the same time maintaining flexibility to arbitrary distractor categories.

Robust Methods for NeRF. For traditional NeRFs, a promising method to mitigate problems with distractors, is to use a robust loss function [21]. RobustNeRF [21] is a robust technique for dealing with distractors. It computes a segmentation mask to ignore distractors in the loss during training. The masks are calculated in each iteration for each image through iteratively reweighted least squares. The idea is that the masks converge over time to the true distractor segmentation. Another approach is to use data-driven priors to remove artifacts from the image [29].

Up to now, there has been no work on distractor handling for 3D Gaussian Splatting.

3 Method

Given a set of input images from a scene with the corresponding camera poses, Robust 3D Gaussian Splatting optimizes for a clean 3D scene representation which ignores any distractors that may be present in the input. To mitigate the problems caused by distractors, we present an approach to identify and track distractors simultaneous to the scene optimization. First, we compute raw masks from image residuals and process them for better spatial smoothness and contiguous local support (Section 3.1). Then, we apply a logistic regression learning to distinguish distractor from non-distractor pixels (Section 3.2). Ultimately, we intersect the resulting mask with object segmentation masks from a pretrained network to enable object awareness (Section 3.3). The resulting distractor mask is used in the loss formulation, such that image parts containing distractors are effectively ignored in the optimization (Fig. 3).

3.1 Raw Mask Generation

Similar to [21], we build distractor masks during the optimization. However, instead of hard thresholds, we use a logistic regression to flexibly learn thresholds. Furthermore, we also compute a mask for each image channel, which increases performance.

We first center the residuals from the last iteration $\epsilon(\mathbf{R}) := |\mathbf{R}_{\text{GT}}^{(i)} - \mathbf{R}_{\text{render}}^{(i)}|$ for all pixels in the image \mathbf{R} of each channel $i \in C$ using the median:

$$\hat{\epsilon}(\mathbf{R})^{(i)} = \epsilon(\mathbf{R})^{(i)} - \text{median}\{\epsilon(\mathbf{R})^{(i)}\}. \quad (1)$$

For robustness, we use the median for centering since the residuals have high variance. Then, we apply a 3×3 box kernel to capture local smoothness and obtain a better local continuity:

$$\omega_1(\mathbf{R})^{(i)} = \hat{\epsilon}(\mathbf{R})^{(i)} \otimes \mathcal{B}_{3 \times 3}. \quad (2)$$

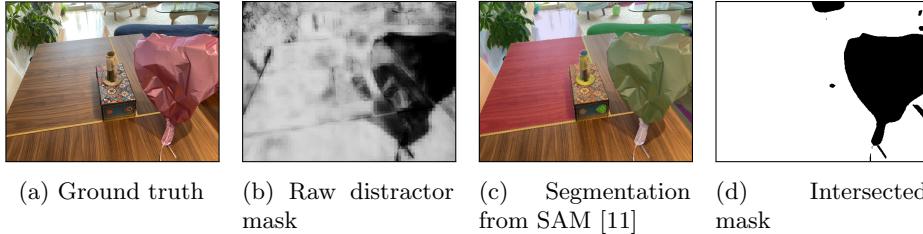


Fig. 2: Given the ground truth (Fig. 2a) and the rendered image we can calculate the raw distractor mask (see Eq. (5)). Next, we intersect the raw distractor mask with the object masks from SegmentAnything (Fig. 2c). The intersected mask (Fig. 2d) is then used in the loss.

Next, we compute the value of whole 8×8 patches in a 16×16 neighborhood. This allows us to capture a more contiguous behavior of distractors. We define

$$\omega_2(\mathcal{R}_8(\mathbf{R}))^{(i)} = \frac{1}{16^2} \sum_{\mathbf{s} \in \mathcal{R}_{16}(\mathbf{R})} \omega_1(\mathbf{s})^{(i)}, \quad (3)$$

where $\mathcal{R}_N(\mathbf{R})$ describes the $N \times N$ neighbourhood around \mathbf{R} . In the next step, we aggregate all information to obtain a better approximation

$$\omega_3(\mathbf{R})^{(i)} = \hat{\epsilon}(\mathbf{R})^{(i)} + \omega_1(\mathbf{R})^{(i)} + \omega_2(\mathbf{R})^{(i)}. \quad (4)$$

3.2 Neural Decision Boundary

Now, we learn the decision boundary using a logistic regression

$$\hat{\mathcal{W}}(\mathbf{R}) = \sigma(\mathbf{W}\omega_3(\mathbf{R}) + b), \quad (5)$$

where \mathbf{W} and b are learned parameters and $\sigma(\cdot)$ is the sigmoid function. Note that we apply the logistic regression pixel-wise and aggregate the channel after using the median.

The idea is that the higher the mask value of a pixel, the more likely it is to be a distractor. This is the case, because we calculate the masks using the residuals. Distractors cause artifacts in the renderings and pixels with artifacts have high residuals, because the artifacts do not match the target image. Another cause of high residuals is that the model has not been trained enough. However, since we are using a dynamic decision classifier that can change over time, we can mitigate this problem. For simplicity reasons, we define

$$\hat{\mathcal{W}}^c := 1 - \hat{\mathcal{W}}. \quad (6)$$

Note that in $\hat{\mathcal{W}}$, distractors are labeled as 0 so they can be directly ignored in the loss.

We train the logistic regression by calculating a custom mask loss using the Gaussian Splatting loss of a rendered image $\mathcal{L}_{\text{GS}}(\mathbf{R}, \mathbf{G})$, where \mathbf{G} is the ground truth image. We define the mask loss using the corresponding mask $\hat{\mathcal{W}}$ for the image \mathbf{R} as

$$\mathcal{L}_{\text{mask}}(\mathbf{R}, \mathbf{G}, \hat{\mathcal{W}}) = \mathcal{L}_{\text{GS}}(\mathbf{R} \circ \hat{\mathcal{W}}, \mathbf{G} \circ \hat{\mathcal{W}}) + \frac{\lambda}{mn} \sum_i^m \sum_j^n \hat{\mathcal{W}}_{i,j}^c \quad (7)$$

where λ is the regularization strength, m, n is the image height and width and \circ is the Hadamard product. The main challenge of this approach is that we do not have a ground truth mask. Because of that, we compute the mask's impact on the Gaussian Splatting loss. The trivial solution for the logistic regression is to classify every pixel as a distractor since we ignore the distractor pixel in the loss and thus obtain a loss of 0. Therefore, we regularize the loss by the proportion of distractor pixels in the mask. After the mask loss calculation, we round our mask

$$\bar{\mathcal{W}} := \text{round}(\hat{\mathcal{W}}). \quad (8)$$

The idea of the neural decision boundary is to dynamically find a threshold value to classify pixels as distractors. Instead of a fixed threshold we used a logistic regression because the logistic regression adapts itself during training and is able to shift the threshold value throughout the training process. In the following, we refer to the non-binary masks as raw masks.

3.3 Establishing Object Awareness

We can see in Fig. 2b that the raw masks capture the coarse structure of the distractor. However, we can also see that the distractor as a whole is not completely correctly classified, and some non-distractor parts also have high mask values. We propose to use only whole objects with an intersection ratio of more than 40% between the object and the mask. We define the object mask set for the image \mathbf{R} as $M(\mathbf{R})$. This set contains a segment mask for every object in the image, i.e., a matrix where every pixel is 1 if it belongs to the object. That means $|M(\mathbf{R})| = \#\text{objects}$ in the image. We obtain the object segmentations from SegmentAnything [11]. An example of the combined segments can be found in Fig. 2c. Now, we can define the intersection of each object with the mask as

$$\mathcal{I}(M(\mathbf{R})) = \{m \in M(\mathbf{R}) | J(m \circ \bar{\mathcal{W}}^c(\mathbf{R}), m) > 0.4\}. \quad (9)$$

where $J(\cdot, \cdot)$ is the Jaccard index and $\bar{\mathcal{W}}$ the rounded mask for image \mathbf{R} . Intuitively, we label a whole object as a distractor if enough pixels in this object are classified as a distractor in the mask computed from the residuals. The intersected mask as in Fig. 2d is defined as follows:

$$\mathcal{W}(\mathbf{R}) = \sum_{m \in \mathcal{I}(M(\mathbf{R}))} (1 - m). \quad (10)$$

Note that we switch between $\bar{\mathcal{W}}^c$ and $\bar{\mathcal{W}}$ to maintain the convention that distractor pixels are labeled as 0, hence the subtraction in Eq. (10).

Using the intersected masks, we can now ignore the distractors in the loss of the Gaussian Splatting optimization. The complete mask generation process is summarized in Fig. 3.

The computation of SAM during training does not significantly influence the runtime. Averaged over different scenes, we observe a runtime increase of less than 1%.

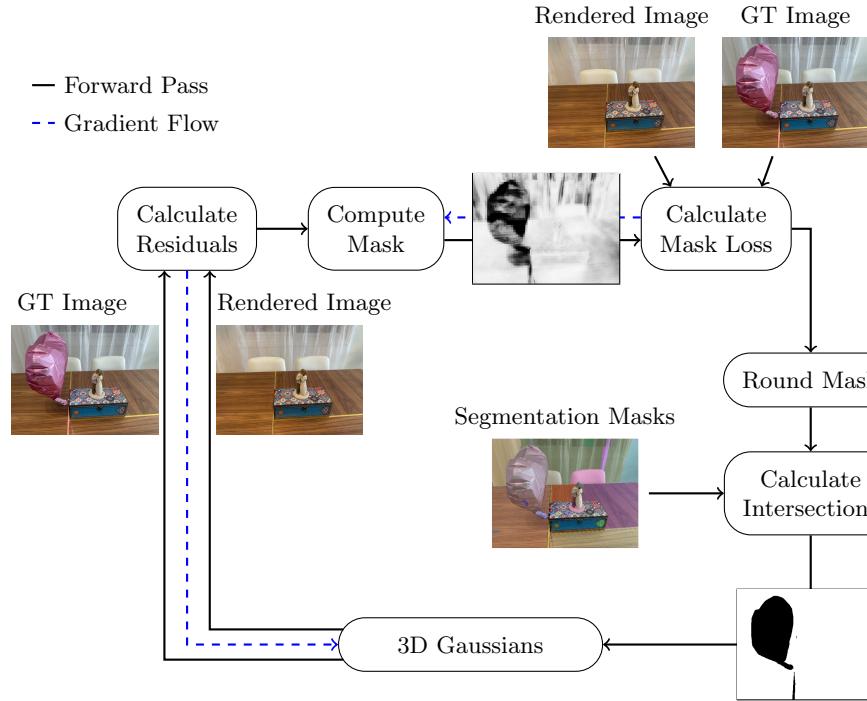


Fig. 3: The first step is to calculate the residuals using the ground truth image and the rendered image. Next, we compute the mask as described in Eq. (5) using the neural decision boundary. Then, we calculate the mask loss from Eq. (7) using the computed mask and propagate back to learn the logistic regression. After that, we intersect the mask with the segmentation masks from Segment Anything and round the masks. The ground truth and the rendered image are multiplied element-wise with the distractor mask and used in the Gaussian Splatting training.

4 Experiments

4.1 Dataset

Following RobustNeRF [21], we evaluate our method on the same four scenes of their released dataset. This dataset was captured to facilitate novel view synthesis on challenging scenes which contain various distractors. Each scene contains training images with distractors and test images that show the scene without any distractors from different camera poses. There is no temporal correlation between images in the dataset. As in the RobustNeRF paper, we also downsampled the images by a factor of 8 to ensure consistent comparisons. The scenes can be described as follows:

- Statue (225 train/19 test): A wooden statue on a painted box. Distractor: a red balloon.
- Yoda (63/202): A stuffed Baby Yoda and toy animals. Distractors: different household items.
- And-bot (122/19): Two Android figurines on a board game box. Distractors: wooden figures.
- Crab (72/72): A crab figure surrounded by a toy train. Distractors: different household items.

Furthermore, we compare our method on a real life scene, see supplementary material.

4.2 Baseline Methods and Ablations

We compare against 3D Gaussian Splatting [10], which our method builds upon, as well as to RobustNeRF [21], the most recent and strongest baseline for novel view synthesis in presence of distractors. We further provide results on the straight-forward combination of 3D Gaussian Splatting with the robust loss from RobustNeRF. To analyze the effectiveness of our added components, we conduct ablation experiments for the neural decision boundary optimization, as well as the object awareness using a pretrained segmentation network. For all methods we use the same training configuration as described in the supplementary material. All in all, we conduct comparisons on the following six approaches.

- **Gaussian Splatting**: The unmodified implementation of 3D Gaussian Splatting as introduced by [10].
- **RobustNeRF (GS)**: A version of RobustNeRF implemented for 3D Gaussian Splatting.
- **RobustNeRF (NeRF)**: The original version of RobustNeRF [21] using mip-NeRF 360 [1].
- **Raw Masks + Segmentation (w/o Neural)**: Sequentially adds the segmentation overlap improvements after calculating the raw masks but does not use a logistic regression.

- **Raw Masks + Neural Decision Boundary (w/o Segmentation):** Sequentially adds a trainable neural decision boundary (logistic regression) after calculating the raw mask but does not use segmentation.
- **Raw Masks + Neural Decision Boundary + Segmentation (Robust Gaussian Splatting):** Our complete method as described in Section 3 and shown in Fig. 3.

We evaluate in terms of three visual quality metrics: peak signal-to-noise ratio (PSNR), SSIM [28], and LPIPS [37].

Table 1: Quantitative comparison to baselines and ablated versions averaged over multiple scenes. The individual metrics are provided in the supplementary material. Values for the RobustNeRF (NeRF) implementation are taken from [21].

Metric	Gaussian Splatting	Robust NeRF (GS)	Robust NeRF (NeRF)	w/o Segmentation	w/o Neural	Robust Gaussian Splatting (Ours)
SSIM \uparrow	0.87253	0.8579	0.76	0.8842	0.8273	0.8875
PSNR \uparrow	26.53	24.09	26.06	27.95	23.28	28.39
LPIPS \downarrow	0.1568	0.1773	0.25	0.1363	0.1698	0.1321

4.3 Quantitative Comparisons

The quantitative result of the evaluation can be seen in Table 1. As the table shows, our method outperforms all other models. We averaged the performance of the models over all scenes. The detailed per scene performance is provided in the supplementary material.

Additionally, we performed an ablation test to investigate if using three color channels as input to the mask calculation improves results over using their norm. The test showed that they provide a small but consistent improvement.

We analyzed the behaviour of Robust Gaussian Splatting on clean data. We found that it performs fairly even to normal Gaussian Splatting on clean images. For detailed values, see supplementary material.

4.4 Qualitative Comparisons

In Fig. 4 we can see a baseline comparison of our method. As shown in the figure, 3D Gaussian Splatting captures the scene in high quality. However, it is filled to a large degree with artifacts of the distractors. When using RobustNeRF (GS), most of these distractors disappear. However, image quality significantly decreases, resulting in poor backgrounds and blurry subjects. When inspecting the masks generated by RobustNeRF (GS), it becomes clear that they are overly aggressive. While all distractors are filtered out, large parts of the remaining input image are masked out as well, effectively reducing the amount of

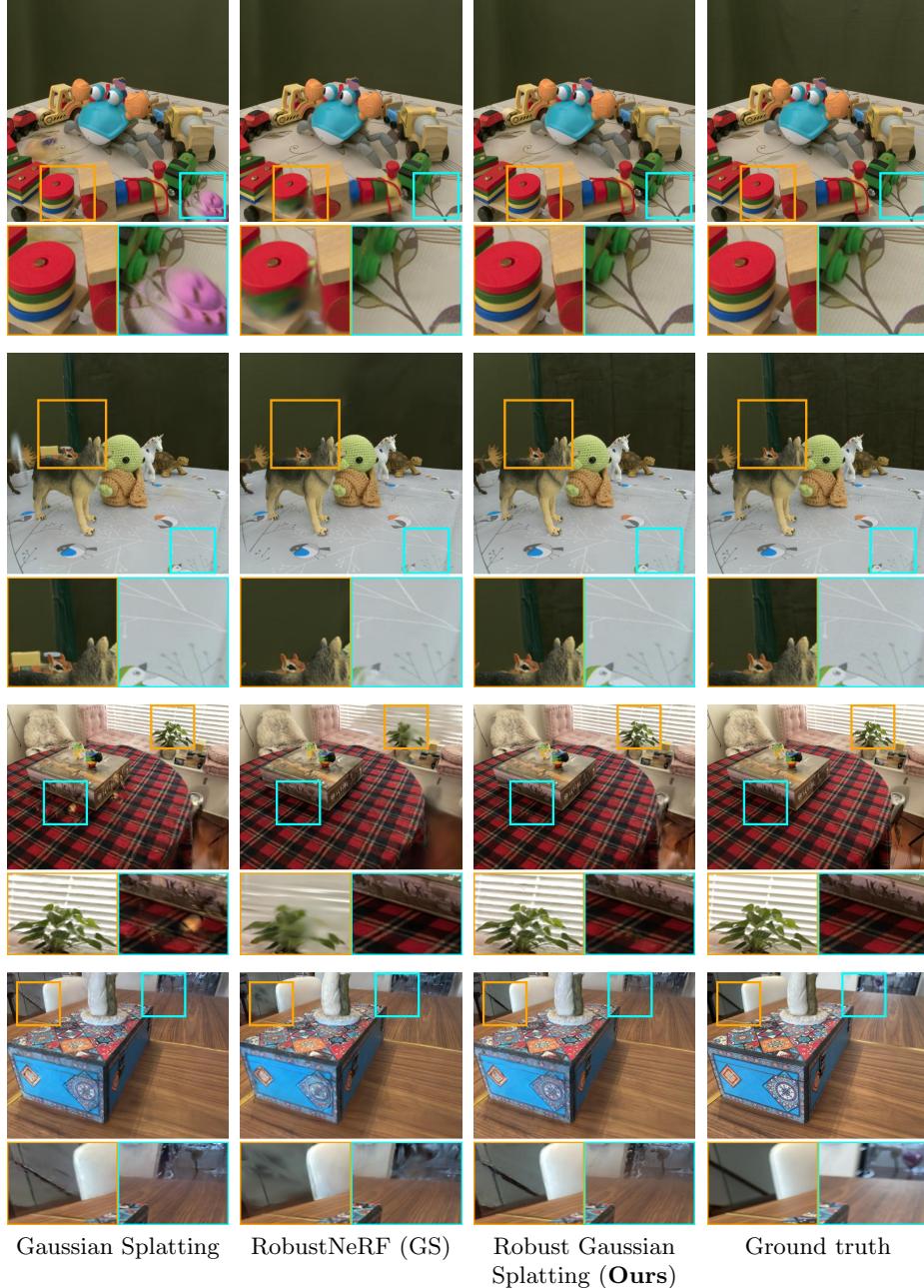


Fig. 4: Example comparison of qualitative results for all scenes from held-out test views. Robust Gaussian Splatting is most effective in ignoring distractors while maintaining a good background and general image quality. For more baseline comparisons see the supplementary material.

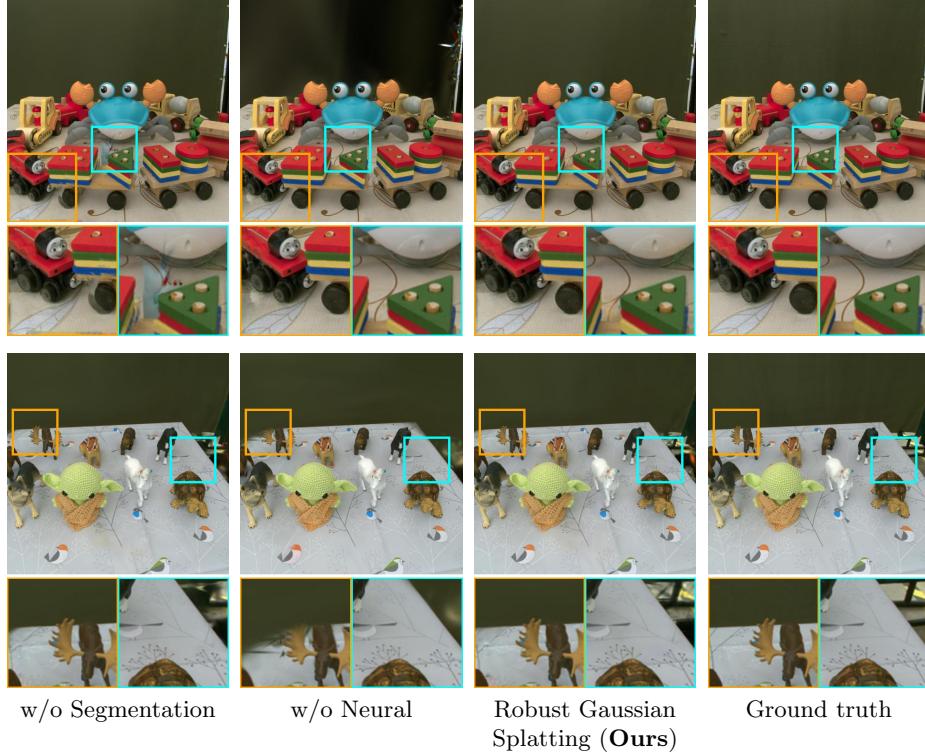


Fig. 5: Example comparison of all ablations. The background of w/o Neural is poor, but it filters the distractors efficiently. The w/o Segmentation version has a good background, but fails to remove all distractors and artifacts, resulting in blurred parts. We can see that our full version is most effective at ignoring distractors. Further ablation comparisons are provided in the supplementary material.

training data. Our method performs best and reduces artifacts to a minimum while maintaining good backgrounds and an overall high image quality.

Figure 5 shows an ablation study of our method. The version w/o Segmentation achieves good reconstruction quality in presence of distractors, however, some artifacts and blurry parts remain, compared to the full version of our method. The version w/o Neural removes many distractor artifacts but struggles with the background and is sensitive to the segmentation quality of SegmentAnything. Robust Gaussian Splatting (ours) yields the best results, achieving good image quality in the foreground and background and only containing a minimum amount of artifacts. Further qualitative comparisons are provided in the supplementary material. We found that our method performs well in various scenarios and show a scene where a person acts as distractor in the supplementary material.

Overall, RobustNeRF (GS) is too aggressive in masking and thus decreases image quality. By refining the masking process, we ensure that our method maintains its effectiveness in handling occlusions and complex scenes while preserving the fidelity of reconstructed images. Our improvements make the masks less aggressive without adding artifacts.

In Table 2 we show comparisons of our approach (Robust Gaussian Splatting) and 3D Gaussian Splatting [10] on clean image datasets without distractors. We can see that our model performs fairly comparable to the original Gaussian Splatting on distractor-free scenes. Checking the distractor masks in the clean settings shows that they are almost all blank. This demonstrates that our approach rarely misclassifies static scene content as distractor. Slightly reduced metrics of the robust method in the clean setting is known topic, that was equally reported in RobustNeRF [21]. Furthermore, Robust Gaussian Splatting is reliable in scenes with different amounts of distractors. For further detailed analysis see the supplementary material.

Table 2: Comparison between our approach and the original 3D Gaussian Splatting on clean image datasets. The scene Playroom is from the Deep Blending Dataset [8], Truck from Tanks&Temples [12] and Yoda (clean) from [21].

Model	Metric	Playroom	Truck	Yoda (clean)	Mean
Gaussian Splatting	SSIM \uparrow	0.9219	0.9244	0.9277	0.9247
	PSNR \uparrow	30.34	26.82	34.17	30.43
	LPIPS \downarrow	0.1403	0.0787	0.1596	0.1262
Robust Gaussian Splatting	SSIM	0.9174	0.9213	0.9366	0.8962
	PSNR	30.05	25.71	34.66	30.14
	LPIPS	0.1436	0.0727	0.1432	0.11198

4.5 Limitations

Our method reliably filters distractors and is able to render high quality novel views in scenes with different distractors. However, in some scenes SegmentAnything struggles to segment the correct objects. In Fig. 6 we can see that SegmentAnything segments each tile of the tablecloth as an object. Using different hyperparameters for SegmentAnything in different scenes can solve this problem and can be addressed in future research. Nevertheless, as we can see in the detailed per scene evaluations in the supplementary material our method performs well considering the incorrect segmentation of SAM.

Furthermore, it is possible to tune the regularization strength and the minimum intersection constant from Eq. (9) for individual scenes to get even better performance. However, we use the same parameter set across all tested scenes and found it generalizes well.



Fig. 6: SegmentAnything struggles with correct object segmentation in some scenes. In this scene, each tile of the tablecloth is segmented as an object. This leads to worse distractor masks and therefore worse results.

5 Conclusion

We address the problem of novel view synthesis in the presence of distractors using 3D Gaussian Splatting. We were consistently able to generate high-quality novel views from input data polluted by distractors. By introducing learnable neural decision boundaries and object awareness into the distractor tracking, we achieve considerably better quantitative and qualitative results than the RobustNeRF (GS and NeRF) method and 3D Gaussian Splatting.

In addition, Robust Gaussian Splatting maintains fairly comparable performance to Gaussian Splatting on clean scenes without distractors. Overall, we believe that our approach is an important step towards robust high-quality novel view synthesis that is easily applicable to in-the-wild data, where the efforts and costs of scene capture and data preprocessing can be kept minimal.

References

1. Barron, J.T., Mildenhall, B., Verbin, D., Srinivasan, P.P., Hedman, P.: Mipnerf 360: Unbounded anti-aliased neural radiance fields. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 5470–5479 (2022)
2. Chan, E.R., Lin, C.Z., Chan, M.A., Nagano, K., Pan, B., De Mello, S., Gallo, O., Guibas, L.J., Tremblay, J., Khamis, S., et al.: Efficient geometry-aware 3d generative adversarial networks. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 16123–16133 (2022)
3. Chen, A., Xu, Z., Geiger, A., Yu, J., Su, H.: Tensorrf: Tensorial radiance fields. In: European Conference on Computer Vision (ECCV) (2022)
4. Chen, Z., Wang, F., Liu, H.: Text-to-3d using gaussian splatting. arXiv preprint arXiv:2309.16585 (2023)
5. Chung, J., Oh, J., Lee, K.M.: Depth-regularized optimization for 3d gaussian splatting in few-shot images. arXiv preprint arXiv:2311.13398 (2023)
6. Fridovich-Keil, S., Yu, A., Tancik, M., Chen, Q., Recht, B., Kanazawa, A.: Plenoxels: Radiance fields without neural networks. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 5501–5510 (2022)

7. Gafni, G., Thies, J., Zollhofer, M., Nießner, M.: Dynamic neural radiance fields for monocular 4d facial avatar reconstruction. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 8649–8658 (2021)
8. Hedman, P., Philip, J., Price, T., Frahm, J.M., Drettakis, G., Brostow, G.: Deep blending for free-viewpoint image-based rendering. ACM Transactions on Graphics (ToG) **37**(6), 1–15 (2018)
9. Keetha, N., Karhade, J., Jatavallabhula, K.M., Yang, G., Scherer, S., Ramanan, D., Luiten, J.: Splatam: Splat, track & map 3d gaussians for dense rgb-d slam. arXiv preprint arXiv:2312.02126 (2023)
10. Kerbl, B., Kopanas, G., Leimkühler, T., Drettakis, G.: 3d gaussian splatting for real-time radiance field rendering. ACM Transactions on Graphics **42**(4) (2023)
11. 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. arXiv preprint arXiv:2304.02643 (2023)
12. Knapitsch, A., Park, J., Zhou, Q.Y., Koltun, V.: Tanks and temples: Benchmarking large-scale scene reconstruction. ACM Transactions on Graphics (ToG) **36**(4), 1–13 (2017)
13. Liang, Y., Khan, N., Li, Z., Nguyen-Phuoc, T., Lanman, D., Tompkin, J., Xiao, L.: Gaufre: Gaussian deformation fields for real-time dynamic novel view synthesis. arXiv preprint arXiv:2312.11458 (2023)
14. Liu, L., Gu, J., Zaw Lin, K., Chua, T.S., Theobalt, C.: Neural sparse voxel fields. Advances in Neural Information Processing Systems **33**, 15651–15663 (2020)
15. Matsuki, H., Murai, R., Kelly, P.H., Davison, A.J.: Gaussian splatting slam. arXiv preprint arXiv:2312.06741 (2023)
16. 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)
17. 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)
18. Park, K., Sinha, U., Barron, J.T., Bouaziz, S., Goldman, D.B., Seitz, S.M., Martin-Brualla, R.: Nerfies: Deformable neural radiance fields. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 5865–5874 (2021)
19. Pumarola, A., Corona, E., Pons-Moll, G., Moreno-Noguer, F.: D-nerf: Neural radiance fields for dynamic scenes. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 10318–10327 (2021)
20. Rematas, K., Liu, A., Srinivasan, P.P., Barron, J.T., Tagliasacchi, A., Funkhouser, T., Ferrari, V.: Urban radiance fields. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 12932–12942 (2022)
21. Sabour, S., Vora, S., Duckworth, D., Krasin, I., Fleet, D.J., Tagliasacchi, A.: Robustnerf: Ignoring distractors with robust losses. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 20626–20636 (2023)
22. Sara Fridovich-Keil and Giacomo Meanti, Warburg, F.R., Recht, B., Kanazawa, A.: K-planes: Explicit radiance fields in space, time, and appearance. In: CVPR (2023)
23. Schönberger, J.L., Frahm, J.M.: Structure-from-motion revisited. In: Conference on Computer Vision and Pattern Recognition (CVPR) (2016)
24. Sun, C., Sun, M., Chen, H.T.: Direct voxel grid optimization: Super-fast convergence for radiance fields reconstruction. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 5459–5469 (2022)

25. Tancik, M., Casser, V., Yan, X., Pradhan, S., Mildenhall, B., Srinivasan, P.P., Barron, J.T., Kretzschmar, H.: Block-nerf: Scalable large scene neural view synthesis. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 8248–8258 (2022)
26. Tang, J., Ren, J., Zhou, H., Liu, Z., Zeng, G.: Dreamgaussian: Generative gaussian splatting for efficient 3d content creation. arXiv preprint arXiv:2309.16653 (2023)
27. Treitschke, E., Tewari, A., Golyanik, V., Zollhöfer, M., Lassner, C., Theobalt, C.: Non-rigid neural radiance fields: Reconstruction and novel view synthesis of a dynamic scene from monocular video. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 12959–12970 (2021)
28. Wang, Z., Bovik, A.C., Sheikh, H.R., Simoncelli, E.P.: Image quality assessment: from error visibility to structural similarity. IEEE transactions on image processing **13**(4), 600–612 (2004)
29. Warburg, F., Weber, E., Tancik, M., Holynski, A., Kanazawa, A.: Nerf-busters: Removing ghostly artifacts from casually captured nerfs. arXiv preprint arXiv:2304.10532 (2023)
30. Wu, G., Yi, T., Fang, J., Xie, L., Zhang, X., Wei, W., Liu, W., Tian, Q., Wang, X.: 4d gaussian splatting for real-time dynamic scene rendering. arXiv preprint arXiv:2310.08528 (2023)
31. Wu, T., Zhong, F., Tagliasacchi, A., Cole, F., Oztireli, C.: D[^]2nerf: Self-supervised decoupling of dynamic and static objects from a monocular video. Advances in Neural Information Processing Systems **35**, 32653–32666 (2022)
32. Xiong, H., Muttukuru, S., Upadhyay, R., Chari, P., Kadambi, A.: Sparsegs: Real-time 360 { \deg } sparse view synthesis using gaussian splatting. arXiv preprint arXiv:2312.00206 (2023)
33. Yan, C., Qu, D., Wang, D., Xu, D., Wang, Z., Zhao, B., Li, X.: Gs-slam: Dense visual slam with 3d gaussian splatting. arXiv preprint arXiv:2311.11700 (2023)
34. Yang, Z., Yang, H., Pan, Z., Zhu, X., Zhang, L.: Real-time photorealistic dynamic scene representation and rendering with 4d gaussian splatting. arXiv preprint arXiv:2310.10642 (2023)
35. Yu, H., Julin, J., Milacski, Z.Á., Niinuma, K., Jeni, L.A.: Cogs: Controllable gaussian splatting. arXiv preprint arXiv:2312.05664 (2023)
36. Yugay, V., Li, Y., Gevers, T., Oswald, M.R.: Gaussian-slam: Photo-realistic dense slam with gaussian splatting. arXiv preprint arXiv:2312.10070 (2023)
37. Zhang, R., Isola, P., Efros, A.A., Shechtman, E., Wang, O.: The unreasonable effectiveness of deep features as a perceptual metric. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 586–595 (2018)
38. Zhu, Z., Fan, Z., Jiang, Y., Wang, Z.: Fsgs: Real-time few-shot view synthesis using gaussian splatting. arXiv preprint arXiv:2312.00451 (2023)

A Training Details

We use a logistic regression as the neural decision boundary trained using SGD with a learning rate of 0.1. For the mask loss we have a regularization strength of 0.1. Furthermore, we use an intersection ratio of 40%. Other 3D Gaussian Splatting hyperparameters are set to default. All experiments are trained for 30k iterations. All scenes use the same hyperparameters.

The models on clean scenes were trained for 30k iterations as well. Images from the scene Playroom and Truck are downsized by the factor of 2 and Yoda (clean) is downsized by 8.

B Analysis of Fraction of Cluttered Images inside a Scene

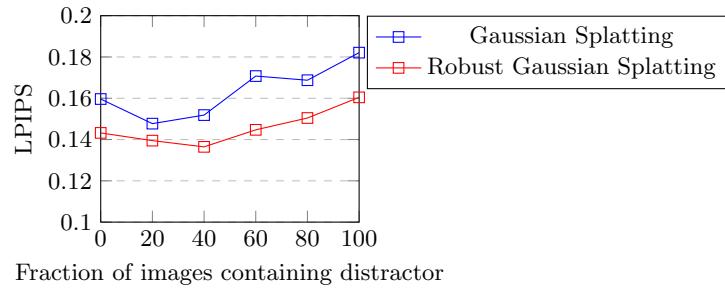


Fig. 7: Metric for the Yoda scene with different fractions of cluttered images. Our method performs better than Gaussian Splatting on all percentages of cluttered images.

Robust Gaussian Splatting is very reliable in scenes with different amounts of distractors. Fig. 7 presents a comparative analysis of our Robust Gaussian Splatting technique against 3D Gaussian Splatting in the Yoda scene under differing proportions of cluttered and clean images. A cluttered image is an image containing distractors. In particular, as the fraction of cluttered images increases, Gaussian Splatting struggles to maintain accuracy. Its performance declines, dropping from 0.159 on clean images to 0.182 on images containing only distractors. In contrast, our method exhibits remarkable stability, maintaining accuracy levels around 0.14-0.16. This robustness to the percentage of cluttered images highlights the effectiveness of our approach, particularly in challenging scenarios characterized by high levels of clutter.

C Detailed Quantitative Analysis

Table 3 shows that the w/o segmentation method yields less consistent improvements. This could be traced back to the high variance in the accuracy of the segments provided by SegmentAnything. For example, in the and-bot scene, each

square of the checkered tablecloth at its center is a separate segment, while the distractors are often not part of any segment. This could be prevented by using different hyperparameters for SegmentAnything, but those might not work for other scenes. The segmentation still considerably improves scores when applied on top of the neural decision boundary optimization, yielding improved scores for our full method over the w/o Segmentation ablation leading to the best mean improvement. Overall, our improvements perform considerably better.

Table 3: Quantitative comparison to baselines and ablated versions. The full version of our method performs better than Gaussian Splatting and RobustNeRF, as well as the ablations of our method. Values for the RobustNeRF (NeRF) implementation are taken from [21].

Model	Metric	Statue	Yoda	And-bot	Crab	Mean
Gaussian Splatting	SSIM \uparrow	0.8401	0.9111	0.8004	0.9385	0.87253
	PSNR \uparrow	21.56	29.99	23.63	30.92	26.53
	LPIPS \downarrow	0.1443	0.1821	0.15938	0.1414	0.1568
Robust NeRF (GS)	SSIM	0.8410	0.8821	0.7952	0.9132	0.8579
	PSNR	21.27	25.93	22.71	26.46	24.09
	LPIPS	0.2001	0.2177	0.1746	0.1891	0.1773
Robust NeRF (NeRF)	SSIM	0.75	0.83	0.65	0.81	0.76
	PSNR	20.89	30.87	21.72	30.75	26.06
	LPIPS	0.28	0.20	0.31	0.21	0.25
w/o Segmentation	SSIM	0.84164	0.9236	0.8235	0.9481	0.8842
	PSNR	21.52	32.41	24.40	33.46	27.95
	LPIPS	0.1376	0.1584	0.1308	0.1184	0.1363
w/o Neural	SSIM	0.8268	0.8994	0.7794	0.8034	0.8273
	PSNR	20.21	29.76	21.51	21.62	23.28
	LPIPS	0.1593	0.1654	0.1778	0.1765	0.1698
Robust Gaussian Splatting	SSIM	0.8514	0.9235	0.8240	0.9511	0.8875
	PSNR	22.21	32.65	24.48	34.23	28.39
	LPIPS	0.1214	0.1604	0.1314	0.1151	0.1321

D Further Qualitative Comparisons

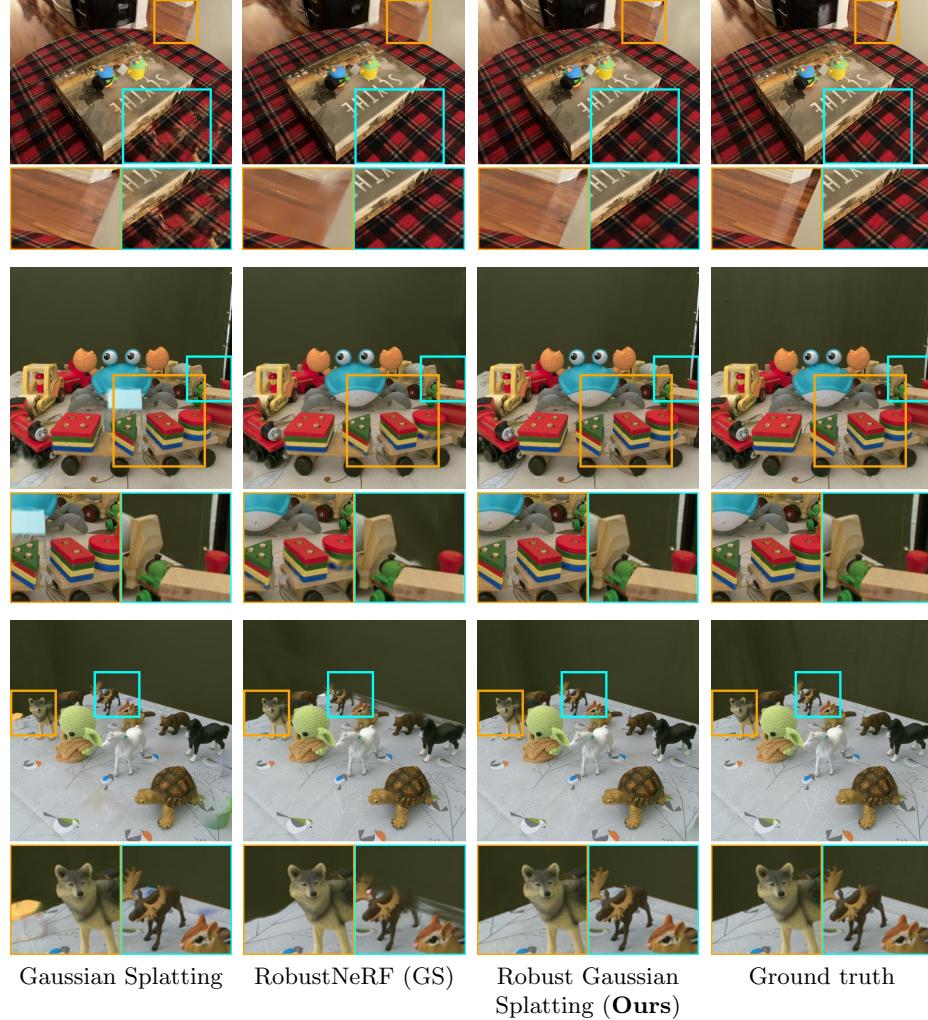


Fig. 8: Example comparison of qualitative results from held-out test views. We can see that our full version is most effective in ignoring distractors.

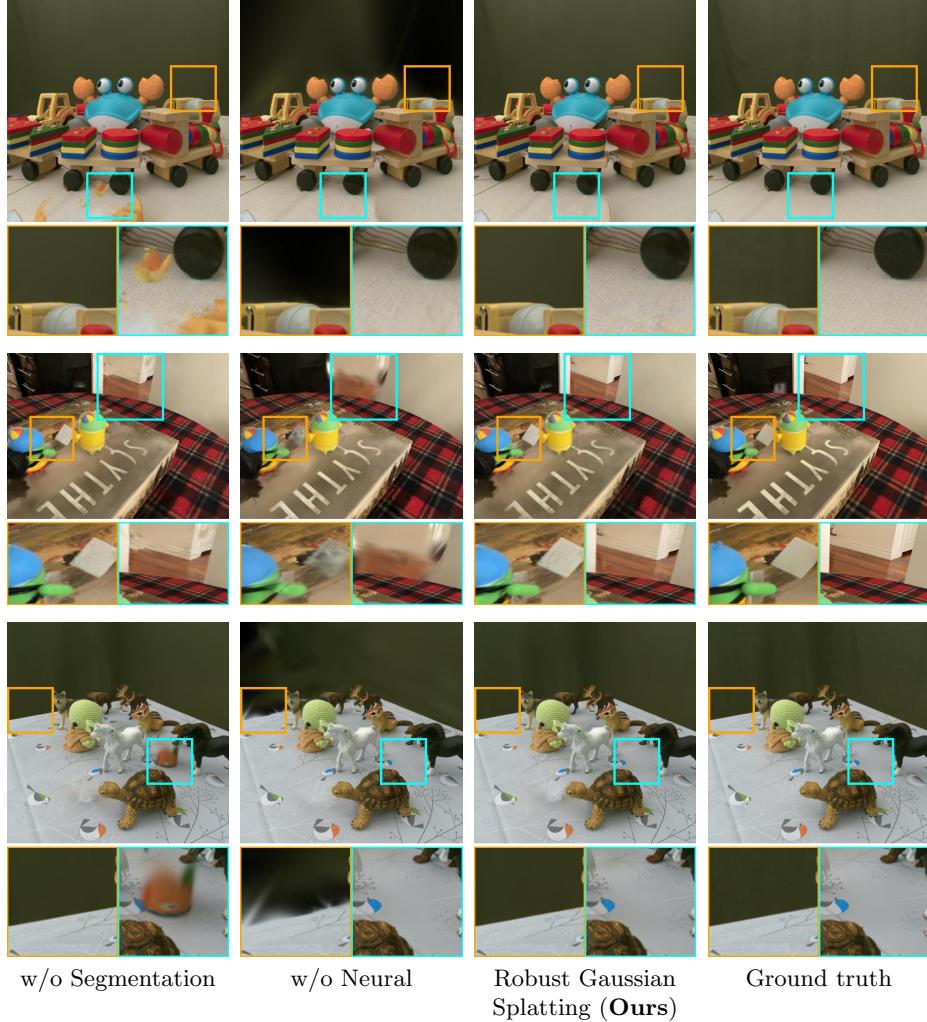


Fig. 9: Comparison of ablations using results from held-out test views. The w/o Neural version struggles to maintain good background but manages to minimize distractors. w/o Segmentation maintains good background but fails to filter out all distractors. We can see that our full version is most effective at ignoring distractors.

E Reallife Scene

Besides the scenes from RobustNeRF dataset, we test our method on the very common scenario, where a person acts as a distractor (Fig. 10). The images are recorded with a smartphone and COLMAP [23] is used to compute camera poses. 21% of the training images contain a person.

While Gaussian Splatting suffers from artifacts caused by the distracting person, our method successfully removes the distractor.



Fig. 10: Comparison of our method and Gaussian Splatting. We trained this scene with the same hyperparameters as from the main paper. 21% of the images contain distractors. Our method effectively filters out distractors and provides qualitative higher images.