

Supplementary Material of MuSHRoom: Multi-Sensor Hybrid Room Dataset for Joint 3D Reconstruction and Novel View Synthesis

Xuqian Ren,¹ Wenjia Wang,² Dingding Cai,¹ Tuuli Tuominen,¹ Juho Kannala,³ Esa Rahtu¹
¹Tampere University, Finland ²The University of Hong Kong, China ³Aalto University, Finland
{xuqian.ren, dingding.cai, tuuli.tuominen, esa.rahtu}@tuni.fi wj2022@connect.hku.hk
Juho.Kannala@aalto.fi

1. Dataset details

1.1. Visualization System during capture

We leverage a visualization system developed by Spectacular AI SDK [1] to inspect the integrity of the point cloud reconstructed from the captured RGB-D images in real-time when using Kinect. In Figure 1, we show an example of the visualization interface.

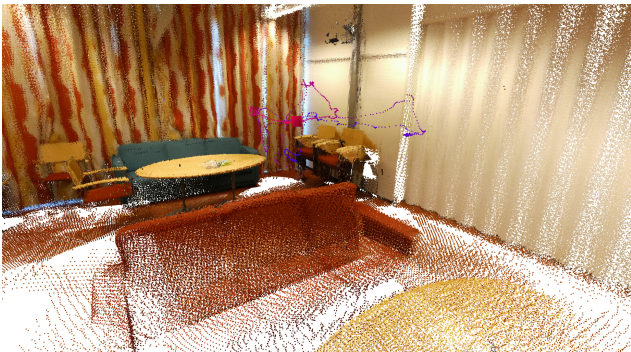


Figure 1. The visualization system. The colored area indicates successful capture, while the empty area is repeatedly scanned by checking the integrity of the point cloud.

1.2. The details of each room

In Figure 2, we show an example image of each room. The MuSHRoom is a room-scale dataset with various styles, colors, illumination, and objects, demonstrating real-world challenges. In Table 1, we show the details of the captured rooms, including the room names, scales, and camera settings.

When using COLMAP [8] to calculate the globally optimized pose for activity and koivu room captured with iPhone sequence, COLMAP failed to calculate accurate poses, so we walked around twice and captured a long sequence. We cut the long sequence from the middle and use the original Polycom of the first circle for training and the

Scene	Scale (m)	Exposure time (μ s)	White Balance (K)	Brightness	Gain
coffee room	$6.3 \times 5 \times 3.1$	41700	2830	128	130
computer room	$9.6 \times 6.1 \times 2.5$	33330	3100	128	255
classroom	$8.9 \times 7.2 \times 2.8$	33330	3300	128	88
honka	$6.1 \times 3.9 \times 2.3$	16670	3200	128	128
koivu	$10 \times 8 \times 2.5$	16670	4200	128	128
vr room	$5.1 \times 4.4 \times 2.8$	8300	3300	128	88
kokko	$6.7 \times 6.0 \times 2.5$	133330	4000	128	255
sauna	$9.9 \times 6.5 \times 2.4$	Auto	3300	Auto	Auto
activity	$12 \times 9 \times 2.5$	50000	3200	128	130
olohuone	$19 \times 6.4 \times 3$	Auto	3600	Auto	Auto

Table 1. The parameters of each room in our datasets. We introduce the room names, room scales and camera parameters.

frames in the second circle for testing.

2. Comparison methods

This section introduces the baseline methods we have compared with.

Volumetric Fusion. Volumetric Fusion [4] proposes to fuse depth from multiple views into the signed distance functions (SDF) and then extract the mesh model using marching cubes (MC) [6]. We use the implementation from Open3D [16], then further cluster the connected triangles and clean small clusters. Since the novel view images can only be synthesized from the textured mesh, the appearance has a large domain gap with the real images.

GO-Surf. Go-Surf [10] represents geometry and color features with the multi-resolution feature grid and decodes these representations into SDF and RGB with two shallow MLP networks. It improves speed and accuracy by simultaneously optimizing the feature volumes, decoders, and camera poses.

Nerfacto. Nerfstudio [9] is an end-to-end workflow that encapsulates various state-of-the-art NeRF techniques, which is friendly to user-collected real-world data. Nerfacto is one of the NeRF pipelines assembled by Nerfstudio that combines components from practical novel methods to balance



Figure 2. The MuSHRoom dataset. Our dataset contains ten rooms with different shapes, colors, illumination, and objects. We show an image example of each room captured by Kinect.

speed and quality. With a proposal sampler [3], scene contraction [3], and density field, Nerfacto can achieve immersive novel view synthesis quality even with real-world noisy data. However, the density field is optimized exclusively for visual consistency, which means that this model sacrifices geometry accuracy and creates occupancy regions to support the volumetric rendering even in parts of the space that are not occupied by the underlying surface. The surface will be predicted with the help of the density rather than accurate zero thickness surface [7]. When the surface is not sharp enough, the consistency of the predicted depth and normal from multiple views cannot be guaranteed, leading the 3D model extracted from the density field to become ambiguity [12] when representing mesh with truncated signed distance function (TSDF). In our evaluation, we use Poisson surface reconstruction [5] to extract the mesh model.

NeuS-facto. SDFStudio [13] is a unified framework that focuses on 3D reconstruction based on Nerfstudio, combined with recent techniques designed from implicit surface reconstruction. Similar to Nerfacto, we chose NeuS-facto with components of proposal network, multi-resolution feature grid, SDF output, and background modeling [14]. SDF output can largely improve the geometry accuracy, but it constrains the occupancy predicting flexibility, which impedes the learning of details in appearance during volume rendering [12].

3. Per-room Evaluation

3.0.1 Implementation Details:

For GO-Surf, Nerfacto, NeuS-facto, and our method, we train each model with 10k, 40k, 60k, and 70k iterations on NVIDIA RTX-2080Ti separately. We train each model without camera optimization. For Go-Surf, when training with iPhone collected data, we set ADAM optimizer with a learning rate of 1×10^{-1} for MLP decoders, and the weights

for rgb, depth, sdf, fs loss are 10 times of the default one. When training Nerfacto, NeuS-facto and our method, we did not include camera pose optimization. Other settings are the same as the default setting reported in each paper. When synthesising pseudo images/depth for our data augmentation strategy, we set the interpolation number n to be 3 for kinect device and 4 for iPhone device. During data augmentation, we render pseudo RGB images from Nerfacto, and pseudo depths from mesh reconstructed by NeuS-facto. These two methods produce the relatively best synthesis and mesh results in our comparison.

3.1. Metrics

Metrics for comparing reconstruction We compare the mesh reconstruction ability from both the accuracy and completeness aspects. As introduced in [10], we measure accuracy (Acc), completion (Comp), Chamfer distance ($C-\ell_1$), normal consistency (NC), and F-score metrics when evaluating reconstruction results. Acc refers to what proportion of the predicted point cloud aligns with a reference point cloud with a certain threshold. Comp refers to how well a reconstructed mesh represents the full reference mesh. $C-\ell_1$ distance measures the similarity between the predicted point cloud and the ground truth point cloud. It computes the average distance from a point in one point cloud to the nearest point in the other point cloud, measuring how closely two point clouds are in space. NC refers to the alignment of normals between two surfaces, representing the influence of the orientation of the surface. F-score used to balance the precision P and recall R by

$$F_{score} = \frac{2PR}{P + R} \quad (1)$$

Precision comes from the percentage of Acc within a

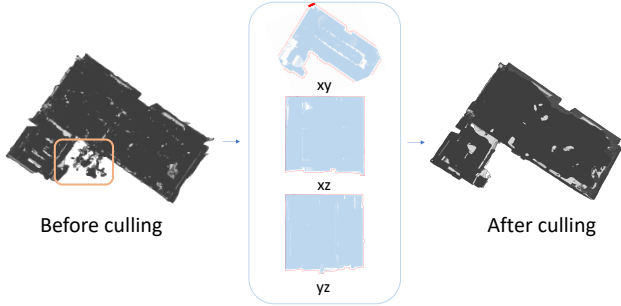


Figure 3. Effect of our culling protocol. On the left is a predicted mesh of the koivu room, reconstructed from a Kinect sequence and have culled using the prior culling protocol. Notably, there remains redundant mesh, as highlighted in the yellow square. However, when culled based on the contour of the reference mesh’s projections, the mesh is cleanly trimmed.

threshold:

$$P(t_i) = \frac{1}{n} \sum_{j=1}^n I(Acc_j \leq t_i) \quad (2)$$

Recall comes from the percentage of Comp within a threshold:

$$R(t_i) = \frac{1}{n} \sum_{j=1}^n I(Comp_j \leq t_i) \quad (3)$$

In our comparison, we set the threshold t_i to 5cm.

Metrics for comparing novel view synthesis

We use PSNR, SSIM [11], and LPIPS [15] to mesh the pixel and feature distances between synthesized images and real images.

3.2. Mesh culling protocol

In the previous method [10], the mesh is culled based on several criteria: first subdivided to have the maximum edge length below 1.5cm and then culled by whether the parts are visible within the camera’s frustum, and if there’s valid depth in the corresponding region, and if they are occluded. However, we noticed that for some non-rectangular rooms without precise boundaries, not all redundant mesh parts are effectively culled. For instance, as depicted in Figure 3, the mesh takes on an “L” shape. The exterior mesh is not culled because it can be observed through a transparent window door. Therefore, meshes culled using the previous protocol can lead to imprecise comparison results. Here, we propose a new culling method that uses the boundary of the projection of the reference mesh to further cull the predicted mesh. In our culling protocol, after aligning the predicted mesh to the reference mesh, we project the reference mesh into the xy, xz, and yz planes. To avoid the boundaries being too close to the predicted mesh and causing some incorrect cuts, we first dilate the projections. Then we detect the

contours of three projections and cut the parts of the predicted mesh that are outside of the contours. For meshes reconstructed from Kinect sequences, we applied both the previous culling protocol and our cutting method. This was due to their non-rectangular geometry and unbounded areas. For meshes derived from iPhone sequences, we only employed our culling protocol. This is because the granularity of these meshes is too coarse to accurately determine if regions are within the camera’s frustum, occluded, or if they constitute valid mesh sections. We compare all regions culled by our protocol with the reference mesh.

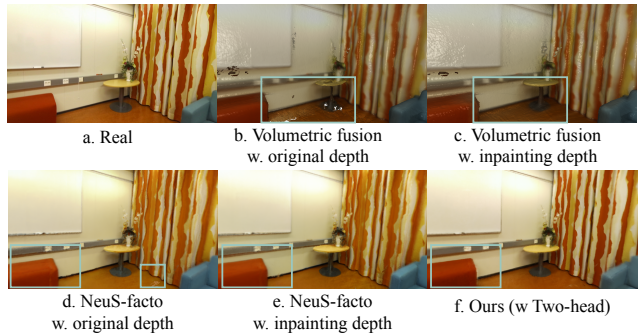


Figure 4. The ablation study of inpainting depth and two-head structure. We visualize the results of Volumetric fusion and NeuS-facto methods with original depth and with inpainting depth. We also present the results with or without the two-head structure. The two-head structure can help NeuS-facto fit the color better and avoid underfitting.

3.3. Per-room quantitative comparison result

In Table 3 and Table 4, we measure the reconstruction and rendering quality quantitatively for each room. Our method can obtain a good trade-off between reconstruction and rendering results. Note that we did not apply the data augmentation to the classroom, computer, and sauna room of Kinect sequences. Because the data augmentation requires accurate pseudo images and depths, the current NeuS-facto model still cannot render accurate pseudo depths and cannot further contribute to the final results.

We also try NeRF++ [14], Mip-NeRF [2] on the MuSH-Room dataset, but these pipelines cannot work on real-world dataset, which indicates the proposal sampling is very crucial for the noisy real-world data. The overall rendering and reconstruction quality of Kinect sequences are relatively better than the results of the iPhone. Potential estimation comes from Kinect can obtain more accurate depth map, which contribute to both the reconstruction and novel view synthesis. For methods that predict SDF for reconstruction, inaccurate depths are not only detrimental to the reconstruction, but also hinder the synthesis learning.

Method	Reconstruction quality					Rendering quality					
						Test within a single sequence			Test with a different sequence		
	Acc ↓	Comp ↓	C- ℓ_1 ↓	NC ↑	F-score ↑	PSNR ↑	SSIM ↑	LPIPS ↓	PSNR ↑	SSIM ↑	LPIPS ↓
Volumetric Fusion (original depth)	0.0178	0.0218	0.0198	0.8514	0.9217	14.61	0.6920	0.3774	12.90	0.6634	0.4150
Volumetric Fusion (inpainting depth)	0.0207	0.0212	0.0210	0.8407	0.9143	14.92	0.6873	0.3950	13.84	0.6556	0.4170
NeuS-facto (original depth)	0.0145	0.0183	0.0164	0.9121	0.9565	21.08	0.7658	0.2198	22.37	0.8483	0.1396
NeuS-facto (inpainting depth)	0.0136	0.0161	0.0149	0.9130	0.9655	21.21	0.7709	0.2132	21.98	0.8465	0.1427

Table 2. The ablation study of inpainting depth and two-head structure. Test within a single sequence means we uniform sample test frames from a single sequence and train on all left frames. Test with a different sequence means we train on one sequence and test on another individual sequence.

3.4. Per-room qualitative comparison result

In this section we show more visualization comparison of each methods with both test within a single sequence and test with a different sequence methods. We show mesh comparison qualitatively in Figure 5 and Figure 6 for Kinect and iPhone sequences. Our method provides a relatively smoother and more completed mesh. The iPhone mesh is more coarse than Kinect mesh, except the mesh produced by Go-Surf [10] method, which shows this method is more robust to devices.

3.5. Ablation study

We further evaluate the effect of inpainting depth quantitatively and qualitatively and show the results in Table 2 and Figure 4. Volumetric fusion heavily relies on the completeness of the depth. Without inpainting the holes, the mesh will have parts missing where depth is invalid, as shown in Figure 4b. We also visualize the effect of the two-head structure in Figure 4. Without the two-head structure, the color of the object exhibits sub-optimal learning, as shown in the red sofa marked by the green square in Figure 4e and 4f, in which the color transitions from a rich red to a paler shade.

References

- [1] Spectacular ai sdk. <https://www.spectacularai.com>, 2021. 1
- [2] Jonathan T Barron, Ben Mildenhall, Matthew Tancik, Peter Hedman, Ricardo Martin-Brualla, and Pratul P Srinivasan. Mip-nerf: A multiscale representation for anti-aliasing neural radiance fields. In *ICCV*, pages 5855–5864, 2021. 3
- [3] Jonathan T Barron, Ben Mildenhall, Dor Verbin, Pratul P Srinivasan, and Peter Hedman. Mip-nerf 360: Unbounded anti-aliased neural radiance fields. In *CVPR*, pages 5470–5479, 2022. 2
- [4] Brian Curless and Marc Levoy. A volumetric method for building complex models from range images. In *Proceedings of the 23rd annual conference on Computer graphics and interactive techniques*, pages 303–312, 1996. 1, 8, 9, 10, 11
- [5] Michael Kazhdan, Matthew Bolitho, and Hugues Hoppe. Poisson surface reconstruction. *Symposium on Geometry Processing, Symposium on Geometry Processing*, Jun 2006. 2
- [6] William E Lorensen and Harvey E Cline. Marching cubes: A high resolution 3d surface construction algorithm. In *Seminal graphics: pioneering efforts that shaped the field*, pages 347–353. 1998. 1
- [7] Marie-Julie Rakotosaona, Fabian Manhardt, Diego Martin Arroyo, Michael Niemeyer, Abhijit Kundu, and Federico Tombari. Nerfmeshing: Distilling neural radiance fields into geometrically-accurate 3d meshes. *arXiv preprint arXiv:2303.09431*, 2023. 2
- [8] Johannes L. Schonberger and Jan-Michael Frahm. Structure-from-motion revisited. In *CVPR*, May 2016. 1
- [9] Matthew Tancik, Ethan Weber, Evonne Ng, Ruilong Li, Brent Yi, Justin Kerr, Terrance Wang, Alexander Kristoffersen, Jake Austin, Kamyar Salahi, et al. Nerfstudio: A modular framework for neural radiance field development. *arXiv preprint arXiv:2302.04264*, 2023. 1, 8, 9, 10, 11
- [10] Jingwen Wang, Tymoteusz Bleja, and Lourdes Agapito. Go-surf: Neural feature grid optimization for fast, high-fidelity rgb-d surface reconstruction. In *2022 International Conference on 3D Vision (3DV)*, pages 433–442. IEEE, 2022. 1, 2, 3, 4, 8, 9, 10, 11
- [11] Zhou Wang, A.C. Bovik, H.R. Sheikh, and E.P. Simoncelli. Image quality assessment: from error visibility to structural similarity. *IEEE Transactions on Image Processing*, 13(4):600–612, 2004. 3
- [12] Yuting Xiao, Yiqun Zhao, Yanyu Xu, and Shenghua Gao. Resnerf: Geometry-guided residual neural radiance field for indoor scene novel view synthesis. *arXiv preprint arXiv:2211.16211*, 2022. 2
- [13] Zehao Yu, Anpei Chen, Bozidar Antic, Songyou Peng Peng, Apratim Bhattacharyya, Michael Niemeyer, Siyu Tang, Torsten Sattler, and Andreas Geiger. Sdfstudio: A unified framework for surface reconstruction, 2022. 2, 8, 9, 10, 11
- [14] Kai Zhang, Gernot Riegler, Noah Snavely, and Vladlen Koltun. Nerf++: Analyzing and improving neural radiance fields. *arXiv: Computer Vision and Pattern Recognition, arXiv: Computer Vision and Pattern Recognition*, Oct 2020. 2, 3
- [15] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 586–595, 2018. 3

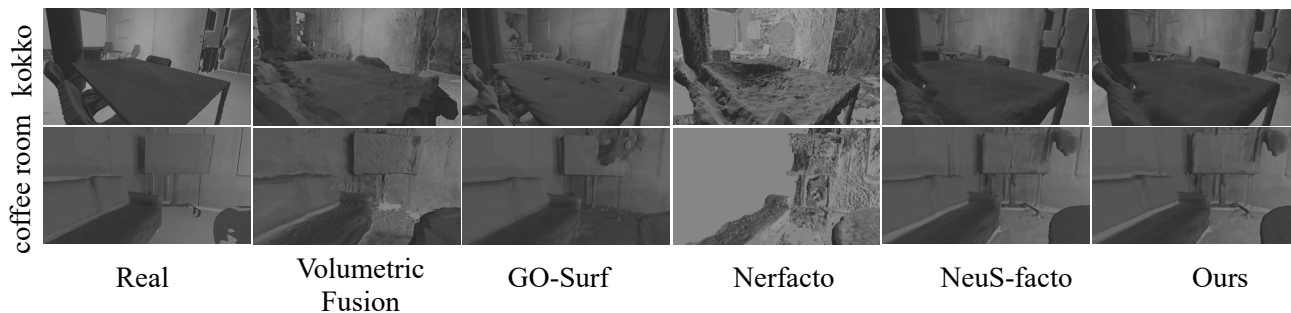


Figure 5. We compare the mesh reconstruction quality of Kinect sequences qualitatively. Please zoom in to see the details.

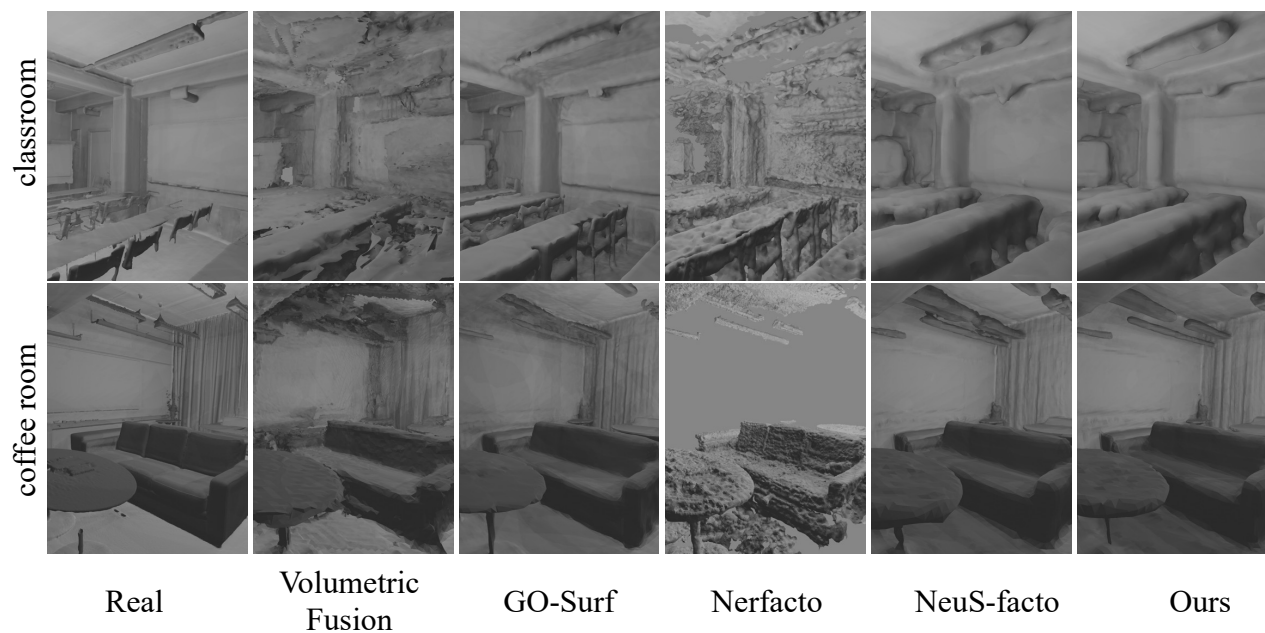


Figure 6. We compare the mesh reconstruction quality of iPhone sequences qualitatively. Please zoom in to see the details.

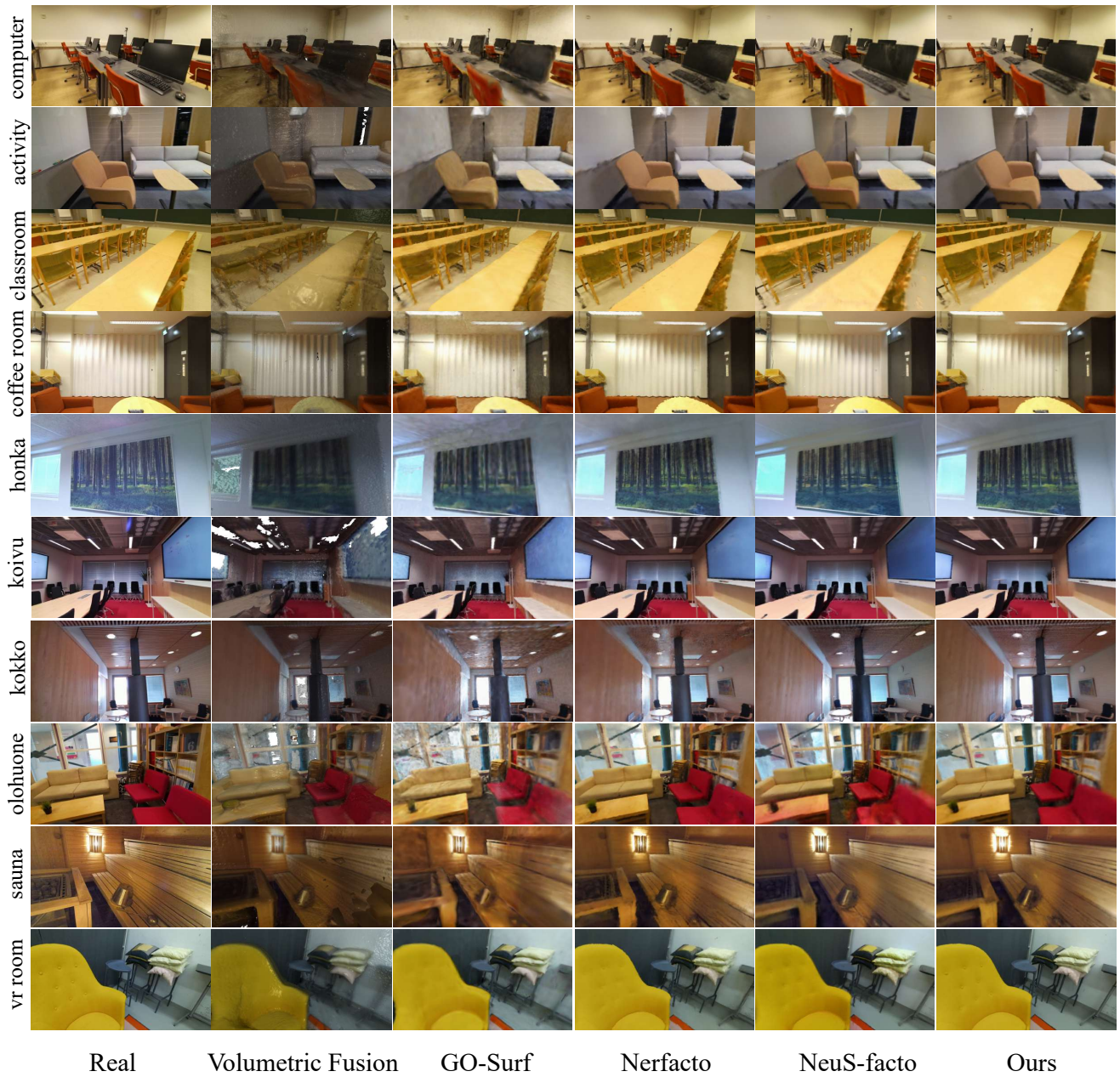


Figure 7. We compare the rendering quality of Kinect sequences with our test within a sequence method qualitatively. The color saturation level and fine-grained content of our method are comparable to the results of Nerfacto [9]. Volumetric Fusion [4] rendering results have a large content gap with the real images. GO-Surf [10] produces images lacking fine-grained details. Visualization results of NeuS-facto [13] still have some ripples, and colors are underfitting to some extent. Please zoom in to see the details.



Figure 8. We compare the rendering quality of Kinect sequences with our test with a different sequence method qualitatively. The color saturation level and fine-grained content of our method are comparable to the results of Nerfacto [9]. Volumetric Fusion [4] rendering results have a large content gap with the real images. GO-Surf [10] produces images lacking fine-grained details. Visualization results of NeuS-facto [13] still have some ripples, and colors are underfitting to some extent. Please zoom in to see the details.

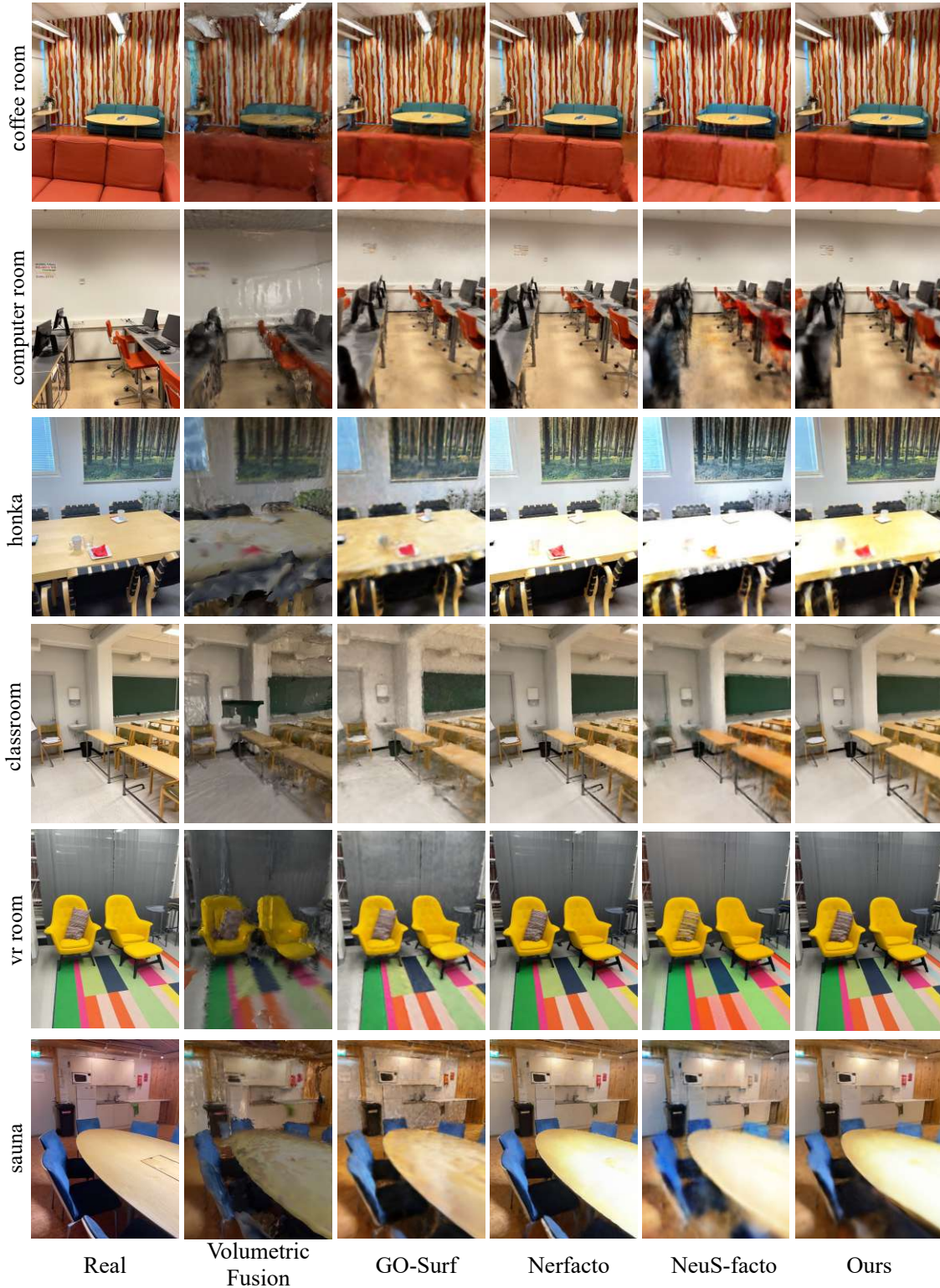


Figure 9. We compare the rendering quality of iPhone sequences with the test within a sequence method qualitatively. Nerfacto [9] method provides the most detailed and photorealistic results. Volumetric Fusion [4] rendering results have a large content gap with the real images. GO-Surf [10] produces images lacking fine-grained details. NeuS-facto [13] results are much more blurry. Our method improves the NeuS-facto from color and fine-grained details but still has a distance when compared with Nerfacto results. The blurry results also show the inaccurate depth in iPhone sequences can be detrimental to the synthesis quality. Please zoom in to see the details.

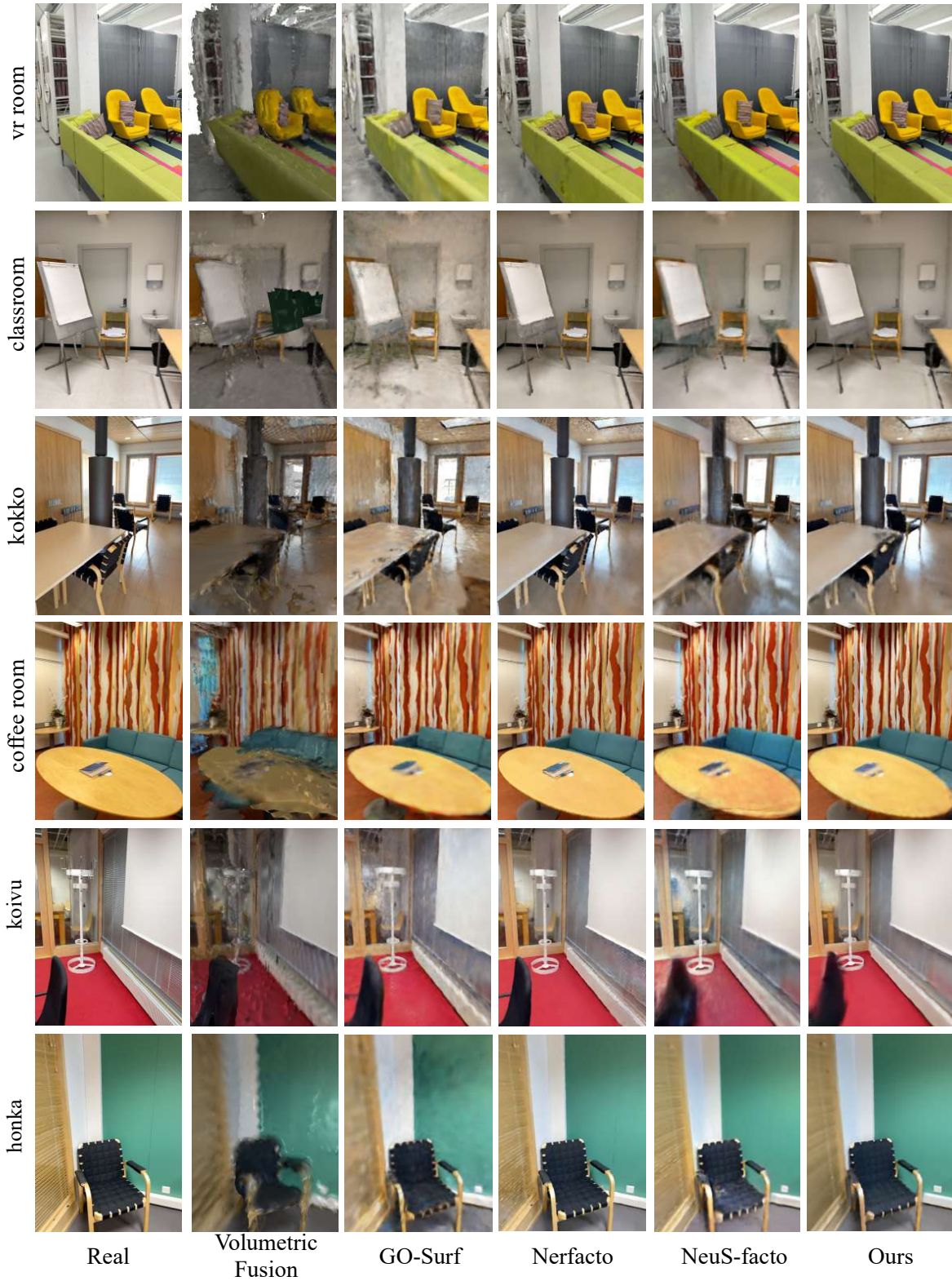


Figure 10. We compare the rendering quality of iPhone sequences with the test with a different sequence method qualitatively. Nerfacto [9] method provides the most detailed and photorealistic results. Volumetric Fusion [4] rendering results have a large content gap with the real images. GO-Surf [10] produces images lacking fine-grained details. NeuS-facto [13] results are much more blurry. Our method improves the NeuS-facto from color and fine-grained details but still has a distance when compared with Nerfacto results. The blurry results also show the inaccurate depth in iPhone sequences can be detrimental to the synthesis quality. Please zoom in to see the details.