

PeRFception: Perception using Radiance Fields

Yoonwoo Jeong^{1†} Seungjoo Shin^{1†} Junha Lee^{1†} Christopher Choy²
 Animashree Anandkumar^{2,3} Minsu Cho¹ Jaesik Park¹
 POSTECH¹ NVIDIA² Caltech³

Abstract

The recent progress in implicit 3D representation, i.e., Neural Radiance Fields (NeRFs), has made accurate and photorealistic 3D reconstruction possible in a differentiable manner. This new representation can effectively convey the information of hundreds of high-resolution images in one compact format and allows photorealistic synthesis of novel views. In this work, using the variant of NeRF called Plenoxels, we create the first large-scale implicit representation datasets for perception tasks, called the **PeRFception dataset**, which consists of two parts that incorporate both object-centric and scene-centric scans for classification and segmentation. It shows a significant memory compression rate (96.4%) from the original dataset, while containing both 2D and 3D information in a unified form. We construct the classification and segmentation models that directly take as input this implicit format and also propose a novel augmentation technique to avoid overfitting on backgrounds of images. The code and data are publicly available in <https://postech-cvlab.github.io/PeRFception/>.

1 Introduction

Over the last few years, advances in implicit representations have demonstrated great accuracy, versatility, and robustness in representing 3D scenes by mapping low dimensional coordinates to the local properties of the scene, such as occupancy [1, 2], signed distance fields [3, 4], or radiance fields [5, 6, 7]. They offer several benefits that explicit representations (e.g., voxels, meshes, and point clouds) could not represent: smoother geometry, less memory space for storage, novel view synthesis with high visual fidelity, to name a few. Thus, implicit representations have been used for 3D reconstruction [1, 2, 8, 9], novel view synthesis [5, 6, 7, 10, 11, 12, 13, 14, 15], pose estimation [16, 17, 18, 19], image generation [20, 21], and many more.

In particular, Neural Radiance Fields [5] (NeRF) and many follow-up works [10, 11, 12, 14, 15, 22] have shown that implicit networks can capture accurate geometry and render photorealistic images by representing a static scene as an implicit 5D function which outputs view-dependent radiance fields. They use differentiable volumetric rendering, a scene geometry, and the view-dependent radiance that can be encoded into an implicit network using only image supervisions. These components allow the networks to capture high fidelity photometric features, such as reflection and refraction in a differentiable manner unlike the conventional explicit 3D representations.

Given the success of the implicit representations, it is only natural to consider the implicit representation as one of the standard data representations for 3D and for perception. However, these novel representations, which can capture a scene with high fidelity, have not yet been used for perception tasks such as classification and segmentation. One of the main reasons is that there is no large-scale

[†] Authors contributed equally to this work.

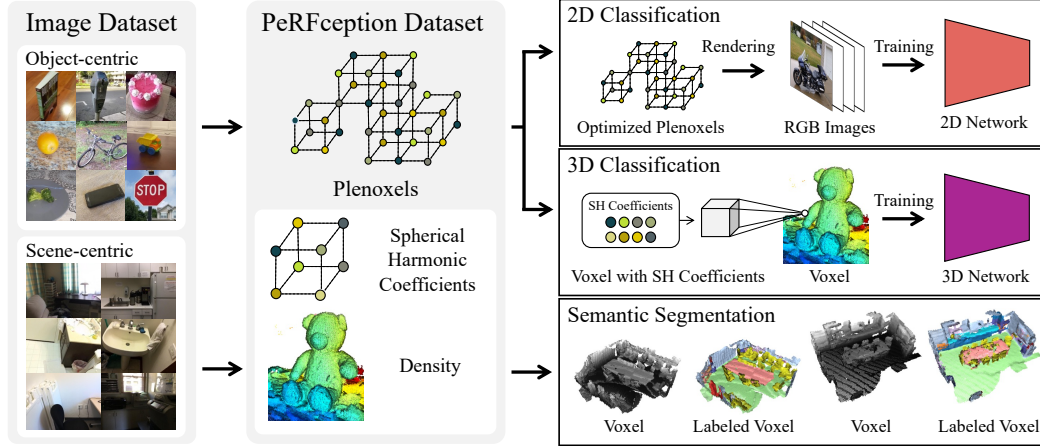


Figure 1: Overall illustration of PerFception dataset with its applications. Our PerFception dataset convey both visual (spherical harmonic coefficient) and geometric (density, sparse voxel grid) features in one compact format, it can be directly applied to various perception tasks, including 2D classification, 3D classification, and 3D segmentation.

dataset that allows training a perception system. Thus, in this work, we present the first large-scale implicit representation datasets to accelerate perception research.

Indeed, NeRFs have drawbacks that prevent the broad adoption of implicit representations as to the standard data format for 3D scenes and perception. First, training an implicit network is slow and can take up to days. Inference (volumetric rendering) also can take minutes, limiting the use of NeRFs in real-time applications. Second, the geometry and visual properties of a scene are implicitly encoded as weights in a neural network. These facts prevent an existing perception pipeline from processing the information directly. Third, implicit features or weights are scene-specific and are not transferrable between scenes. However, for perception, channels or features must have a consistent structure, such as RGB channels for images. For instance, if the order of channels is different from an image to an image, the image classification pipeline would not work properly.

Recent studies have resolved these limitations by adopting explicit sparse voxel grid geometry and basis functions for features. First, to tackle the slow speed, many works propose to use the explicit sparse voxel geometry, which reduces the number of samples along a ray by skipping empty space [10, 11, 12, 14, 15, 22]. Second, instead of using the implicit representations of weights of a network, directly optimizing features [14, 15, 22] assigned to explicit geometry reduces the time to extract features from a network. Lastly, for consistent features between scenes, which is crucial for perception or creating a scene with different objects in NeRF format, Yu et al. [14, 15] show that spherical harmonic coefficients can represent a scene as accurately as NeRFs while preserving consistent and structured features. In particular, Plenoxels [14] satisfy all criteria for data representation which supports fast learning and rendering while maintaining a consistent feature representation for perception and composition of scenes.

In this work, we adopt Plenoxels as the primary format for perception tasks and create both object-centric and scene-centric environments. We mainly use two image datasets and convert them into Plenoxels, the Common Object 3D dataset (CO3D) [23] and ScanNet [24], and name the converted datasets as PerFception-CO3D and PerFception-ScanNet, respectively. As the size of Plenoxels can be extremely large, we present a few techniques to compress the size of data and hyperparameters for each setup to maximize the accuracy while minimizing the data size.

We use the PerFception datasets to train networks for 2D image classification, 3D object classification, and 3D semantic segmentation. We successfully train networks for each perception task, indicating that our datasets effectively convey 2D and 3D information together in a unified format. Moreover, we show that our representation allows more convenient background augmentation and sophisticated camera-level manipulation.

We summarize our contributions as follows:

- We introduce the first large-scale implicit datasets that can be readily used in downstream perception tasks, including 2D image classification, 3D object classification, and 3D scene semantic segmentation.
- We conduct the first comprehensive study of visual perception tasks that directly process the implicit representation. The extensive experiments show that our datasets effectively convey the information for 2D and 3D perception tasks.
- We provide the ready-to-use pipeline to generate the implicit datasets with fully automatic processes. We expect this automatic process allows generating a very large scale 3D dataset in future.

2 Related Work

2.1 Neural Implicit Representations

Representing a scene using an explicit representation such as voxels, meshes, or point clouds has been the most widely used format, but these are discrete and introduce discretization errors. Neural implicit representations, on the other hand, use a neural network to approximate the geometry or properties of a scene continuously [1, 2, 3, 25]. Mildenhall et al. [5] showed that neural radiance representation can generate high fidelity renderings with view-dependent illumination effects using multi-layer perceptrons. Many of recent studies extend such implicit representation to dynamic scenes [26, 27, 28, 29], conditional generation [30, 31, 32, 33], pose-free [16, 17, 19], and many more. In particular, research on efficient rendering [10, 12, 15, 22] has been one of the major directions since volumetric rendering could take minutes. Hedman et al. [10] propose to create sparse voxel grids after training to accelerate rendering. Similarly, Plenocree [15] uses an octree data structure instead of sparse voxels for fast rendering. DVGO [22] and Plenoxels [14] also adopt the sparse voxel structure, and improve both inference and training time. INGP [12] proposes a multi-level hash encoding which enables the fast convergence. Recently, TensorRF [34] boosts both training and inference time by factorizing 3D radiance fields into lower dimensional vectors or matrices. In this work, we use Plenoxels for our data format since they have explicit geometry and consistent features in form of spherical harmonic coefficients.

2.2 3D Perception Datasets

Over the last decade, many public large-scale datasets of real objects for 3D perception have been published thanks to the advances in commodity sensors. In this section, we cover such large-scale 3D datasets for objects and scenes.

ShapeNet [35] and ModelNet [36] provide class and part annotations that are from synthetic CAD models. Early object-centric 3D datasets augment image datasets with 3D CAD model annotations. Pascal3D+ [37] and Objectron [38] contain 3D shapes that are matched with real-world 2D images containing objects; however, 3D models are chosen from approximately aligned 3D models, not precisely reconstructed from the corresponding 2D images. Redwood [39] is a large-scale object-centric RGB-D scan video dataset, where only a few categories include 3D models and camera poses. GSO [40] holds clear 3D models of real objects with textures, but missing physically rendered images. 3D-Future [41] provides synthetic CAD shapes with high-resolution informative textures developed by professional designers. CO3D [23] provides large-scale object-centric videos with camera poses and high-quality point cloud models. They assess quality of reconstructed 3D shapes using human-in-the-loop validation and marked 5,625 point clouds as successfully reconstructed. Recently, ABO [42] offers a dataset consisting of household object images and high-quality 3D models with 4K texture maps and full-view coverage. Professional artists manually designed its high-quality spatially-varying Bidirectional Reflectance Distribution Functions (BRDFs), indicating that the data generation processes were not fully automatic. We summarize the details of the aforementioned datasets in Table 1.

Many scene-centric 3D datasets use depth sensors to scan a section or an entire room and create dense annotations. SUN RGB-D [43] collected 13,355 RGB-D images, which are densely annotated with 2D polygons and 3D bounding boxes. However, it does not include camera parameters, which is essential information for surface reconstruction. NYUv2 [44] initially sparked interest for 3D scene understanding, with 464 indoor scans, 1,449 frames of which are annotated with 2D polygons for semantic segmentation. SUN3D[45] is comprised of 415 RGB-D indoor video sequences in 254

Table 1: Specs for publicly available 3D datasets. “Real” denotes whether the objects are from real-world images, “Full 3D” for the availability of 3D geometries for all the objects, and “Multi-view” for multi-view images and corresponding real-world catalog images. \blacktriangle is marked when the corresponding information is partially provided.

Dataset	# Classes	# Objects	Real	Full 3D	Multi-view
ShapeNet[35]	55	51K	\times	\checkmark	\checkmark
ModelNet[36]	40	128K	\times	\checkmark	\times
Pascal3D+[37]	12	36K	\checkmark	\times	\times
Redwood[39]	9	2K	\checkmark	\times	\checkmark
Objectron[38]	9	15K	\checkmark	\blacktriangle	\blacktriangle
GSO[40]	\times	2K	\checkmark	\checkmark	\times
3D-Future[41]	8	2K	\times	\checkmark	\times
ABO[42]	98	8K	\blacktriangle	\checkmark	\checkmark
CO3D[23]	51	19K	\checkmark	\blacktriangle	\checkmark
PeRFception-CO3D	51	19K	\checkmark	\checkmark	\checkmark

different spaces; only eight sequences are annotated. Each sequence was captured densely, with a large number of frames collected. 2D-3D-S [46, 47] is an instance-level annotated large-scale indoor scene dataset. It offers diverse modalities of six indoor scenes in RGB images, depth maps, surface normals, 3D meshes, and point clouds. Currently, ScanNet [24] is the most popular large-scale indoor scene dataset that collected instance-level annotated 1,513 scans of RGB-D images and 3D data. Matterport3D [48] contains large-scale RGB-D images annotated with surface and semantic information. In particular, it covers a wide area of 90 building-scale scenes by capturing panoramic views, but it does not provide annotations. Replica [49] is a small but high-quality surface-annotated indoor scene reconstruction database. In this paper, we use two datasets, CO3D and ScanNet, to cover both object- and scene-centric dataset respectively.

We select CO3D for object-centric dataset since it consists of camera-annotated real-world images and has sufficient number of classes. In addition, we use ScanNet since it is one of the most popular 3D indoor dataset providing adequate number of data with rich annotations.

2.3 3D Perception Models

Unlike perceptions in the 2D domain, where the image is the de facto standard representation, there is no canonical representation for spatial 3D data. Existing explicit representations, such as voxels, meshes, and point clouds, target different aspects of data and have pros and cons. We categorize methods into two groups based on the input representation. Point-based methods [50, 51] directly consume the continuous 3D coordinates of point clouds or meshes using MLP and continuous/graph convolutions. Recent studies have tried to define custom convolution layers [52, 53, 54, 55, 56] upon the continuous coordinate space, or non-local operations [57, 58, 59]. Overall, these methods exhibit fast and simple processing, but it often requires a large computational cost due to neighbor search.

On the other hand, voxel-based methods discretize input coordinates into voxels, which introduces small quantization errors but allows fast neighbor lookup using a data structure. Specifically, recent advances in spatial sparse convolutions [60, 61, 62] that operates on sparse voxels utilize an efficient GPU hash table and require small memory footprint for neighbor search. It has shown successful adoption in many perception tasks, including semantic segmentation [61, 62, 63], object detection [64], representation learning [65], and registration [66, 67, 68, 69]. We use the spatial sparse convnets to create the first perception network on our PeRFception datasets due to its scalability in terms of memory footprint and computational cost.

3 Preliminary

Yu et al. [14] proposed a novel scene representation called Plenoxels that combines a sparse voxel grid for coarse geometry with spherical harmonic coefficients for radiance fields. Unlike conventional MLP-based NeRFs [5, 6, 10, 11], which use a single neural network to represent an entire scene, Plenoxels optimize coefficients of spherical harmonic in each non-empty voxel independently, which

Table 2: Specs of CO3D and ScanNet, and our PeRFception-CO3D and PeRFception-ScanNet. **SH** denotes spherical harmonic coefficients, **D** for densities, and **C** for diffused color, "pcd" for point cloud. 3D-BG marks whether the 3D representation includes backgrounds of scenes. We note the number of frames in our proposed datasets as ∞ since our data representation is feasible to render frames from infinitely many camera intrinsics and extrinsics.

Dataset	# Scenes	# Frames	3D Shape	Features	3D-BG	Memory	Memory(Rel)
CO3D	18.6K	1.5M	pcd	C	\times	1.44TB	$\pm 0.00\%$
PeRFception-CO3D	18.6K	∞	voxel	SH,D	\checkmark	1.33TB	-6.94%
ScanNet	1.5K	2.5M	pcd	C	\times	966GB	$\pm 0.00\%$
PeRFception-ScanNet	1.5K	∞	voxel	SH,D	\checkmark	35GB	-96.4%

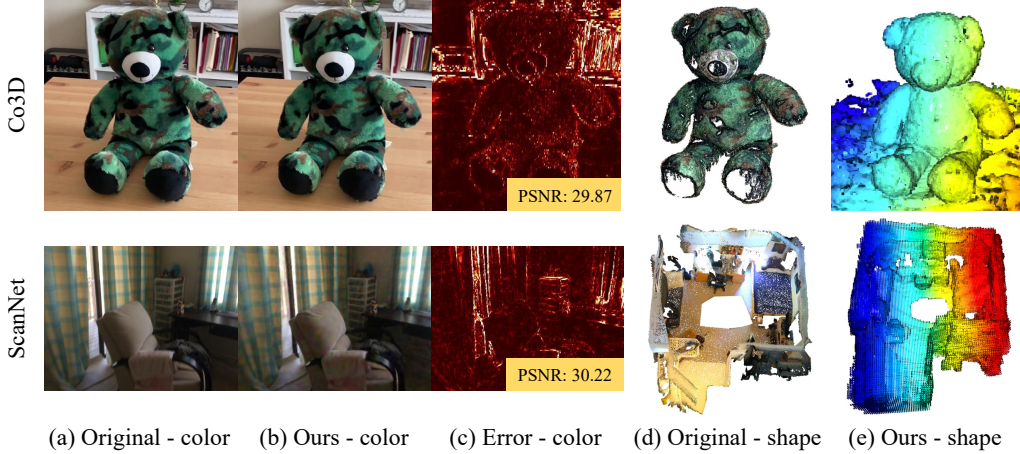


Figure 2: Visualization of a few example data of original datasets and our PeRFception datasets. From the source images and corresponding parameters, we successfully construct PeRFception datasets with both accurate geometry and photorealistic rendering. The color used in (e) is for visualization.

uses the same differentiable model for volumetric rendering described in NeRF [5] as follows:

$$\hat{C}(\mathbf{r}) = \sum_{i=1}^N T_i (1 - \exp(-\sigma_i \delta_i)) \mathbf{c}_i, \quad \text{where } T_i = \exp\left(-\sum_{j=1}^{i-1} \sigma_j \delta_j\right), \quad (1)$$

where T_i denotes light transmittance of a i -th sample, σ is opacity, \mathbf{c} is color, and δ is distance to the next sample on a ray \mathbf{r} . Plenoxels lookup the stored densities and spherical harmonic coefficients in the sparse voxel grids. For scenes with background, Plenoxels also use the lumisphere background representation to render the backgrounds.

We modified the official Plenoxels implementation to set the initial grid properly and hyperparameters. As the size of Plenoxels can be extremely large, we present a few techniques to compress the size of data and hyperparameters. More implementation details are in Section. 4.1, Section. 4.2, and the appendix.

4 Generating PeRFception dataset

We generate two datasets PeRFception-CO3D and PeRFception-ScanNet to train perception networks.

4.1 PeRFception-CO3D

CO3D [23] is a large-scale object-centric dataset that contains multi-view observations of objects. It contains 18,669 annotated videos with a total 1.5 million of camera-annotated frames and 50 classes from MS-COCO [70], and images crowd-sourced from Amazon Mechanical Turk (AMT). It also provides reconstructed pointclouds, generated by pretrained instance segmentation algorithm [71] and COLMAP [72]. Although reconstructing depth and pointclouds is automatic, its generation step

Table 3: Overall rendering qualities of PeRFception-CO3D and PeRFception-ScanNet on test set. For class-wise rendering scores are reported in the appendix.

Dataset	PSNR(\uparrow)	SSIM(\uparrow)	LPIPS ¹ (\downarrow)	Train Time (m)	PSNR > 15	PSNR > 20	PSNR > 25
PeRFception-CO3D	28.82	0.8564	0.3451	21.6	99.8%	98.2%	87.3%
PeRFception-ScanNet	23.42	0.7470	0.4637	11.3	99.7%	73.9%	40.4 %

still requires human-in-the-loop validation. When the amount of data increases, this human-in-the-loop validation becomes unsuitable. On the other hand, our 3D dataset generation does not require manual verification since image reconstruction qualities on unseen views are used as a proxy for reconstruction quality. We compare specs of the original CO3D and PeRFception-CO3D in Table 2.

Data Generation. We use the official implementation of Plenoxels [14] with a slight modification to the default configuration. We reduce the resolution of the background lumisphere from 1,024 to 512 and the number of background layers from 64 to 16. For sharper surface, we set the lambda sparsity value to 10^{-10} , 10 times larger than the default configuration. A voxel grid is initialized with 128^3 resolution and trained for 25,600 iterations. Then, it is upsampled once to 256^3 resolution and trained for further 51,200 iterations. Before saving the data, we quantize the trained parameters to unsigned 8-bit integers to minimize for storage except for density values. For each scene, we first filter out defective images and uniformly sample 10% of the images as the test set to assess the rendering quality. The quantitative and qualitative results of the rendering quality are reported in Table 3 and Figure 2. More details are in the appendix.

4.2 PeRFception-ScanNet

ScanNet is a 3D-scanned indoor dataset that captures more than 1.5K indoor scenes with the commercial RGB-D sensors. It provides 3D reconstructed point clouds of scenes with semantic labels containing 20 common object classes, as well as the raw RGB-D frames with corresponding semantic masks and camera parameters. In our experiment, we follow the official data split and report the numbers on the validation split since the test set annotations are not publicly available.

Data Generation. ScanNet videos are captured using handheld cameras where auto-exposure option is held. So, a fair number of frames contain motion blur which could lead to poor scene geometries. In practice, we generate the batches of rays before training and load them in CPU memory for efficient memory bandwidth utilization during training. Since the number of frames for each scene in ScanNet varies, we use uniformly-sampled 1,500 image frames at most. For the scenes with fewer than 1,500 images, we use them all after filtering out blurry images with a low variance of Laplacian [73]. Another characteristic of ScanNet is that, unlike object-centric datasets where cameras face inward, the images are captured from inside a room facing outward. These result in fewer images observing the same part of the space, which results in poor reconstruction of Plenoxel’s geometry on ScanNet dataset. Specifically, the Plenoxel reconstruction artificially creates an excessive number of voxels in the empty space (i.e. floaters) to minimize the image reconstruction loss.

Instead, to supply an additional geometric prior to Plenoxel training, we initialize the voxel grid using the unprojected depth maps provided in ScanNet rather than starting from the dense voxel grid. However, since the provided depth maps of ScanNet are contaminated with noisy observations, we incorporate the connected component analysis to filter out the disconnected outlier points in the unprojected point clouds. This leads to stable and more accurate reconstruction and does not excessively generate floaters to minimize the rendering loss. The resulting PeRFception-ScanNet dataset occupies only 35GB in disk whereas the original video streams of ScanNet requires about 966GB disk space. This is a significant compression rate (96.4%), which emphasizes the accessibility of our representation as a dataset. Detailed dataset specs of the original ScanNet and PeRFception-ScanNet are reported in Table 2. We report the rendering quality on Table 3 and visualize the qualitative novel view renderings on Figure 2.

¹The LPIPS metric sometimes generates “nan” although their visual qualities are great enough. Only 0.01% of scenes are noted as failure.

5 Experiment

We benchmark popular 2D image classification, 3D object classification, and 3D segmentation networks to demonstrate that our unified data format can be used for various perception tasks.

5.1 Classification on PeRFception-CO3D dataset

CO3D provides multi-view images of objects and 51 class labels for classification. We use the same class labels for classification of PeRFception-CO3D dataset. We adopt a few classification models for our dataset for both 2D [74] and 3D classification [62]. We split the dataset into the train, validation, and test set by scenes since the original CO3D does not provide such splits. We use 10% of the scenes for validation set and 10% for test set in each class. We use the same splits for 2D and 3D classification.

5.1.1 Implementation Details

All the 2D classification models are trained with the cross-entropy loss with the weight decay factor 10^{-4} . Following the recommendations from [75], we utilize the label smoothing with $\epsilon = 0.005$, remove bias decay, and initialize weights of the batch normalization layers on the residual connections to 0. We use the SGD optimizer with momentum 0.9 and trained for 500 epochs with a batch size 64. For 50 epochs, we linearly warmed up the learning rate from 0 to 0.1 and decayed it using the cosine annealing scheduler. It takes up to a day for training using a single RTX 3090 GPU. We have benchmarked variants of ResNet (ResNet18, ResNet34, ResNet50, ResNet101, ResNet152), ResNext [76] (ResNext50, ResNext101), and WideResNet [77] (WideResNet50, WideResNet101) networks.

For 3D classification, we train 3D version of ResNets [74] with varying depths that are implemented with spatially sparse convolutional layers [61, 62]. These networks directly take sparse voxels from Plenoxels as input. The Plenoxels consist of two components: coordinates of sparse voxels and their features (spherical harmonic coefficients and density values). To demonstrate the efficacy of such features in perception task, we train the networks by providing different input features. We use the SGD optimizer and set the initial learning rate as 0.1 for all experiments and decay it with the cosine annealing scheduler for 100K iterations with batch size 16 on a single RTX 3090 GPU. We augment the input data with both geometric augmentation (random rotation, coordinate dropout, random horizontal flip, coordinate uniform translation, and random scaling) and feature-level augmentation, random feature jittering. Further implementation details are in the appendix.

Background Augmentation. Plenoxels use both sparse voxels and lumispheres to render foregrounds and backgrounds respectively. In other words, we can render each of them separately or manipulate them to create various augmentations. Specifically, we create a novel augmentation that substitutes the background in a scene with backgrounds from other scenes while preserving the foreground object. We describe the composition of foregrounds and randomly selected backgrounds in the appendix. In addition, we visualize several background augmentation examples in the appendix.

5.1.2 2D Classification on PeRFception-CO3D

We train both the scratch and ImageNet [78] pretrained version of ResNet [74] variants on the original CO3D and PeRFception-CO3D to show that PeRFception-CO3D contains the same information as the original CO3D dataset. Using our random pose selection algorithm(described in the appendix), which discourages the selected pose to be extremely unobserved the train frames, we select 50 poses for each scene. As shown in Fig 3, the model trained with the original CO3D has a larger gap than the model trained with PeRFception-CO3D dataset. In addition, we observed using background augmentation is beneficial for improving generalizability of classification networks, especially performs the best when the augmentation is applied with probability $p = 0.5$. We conduct controlled experiments about the probability p in the appendix. In addition, using a popular vision analysis tool GradCAM++ [79], we demonstrate that using background augmentation helps the model not to memorize the backgrounds in the appendix.

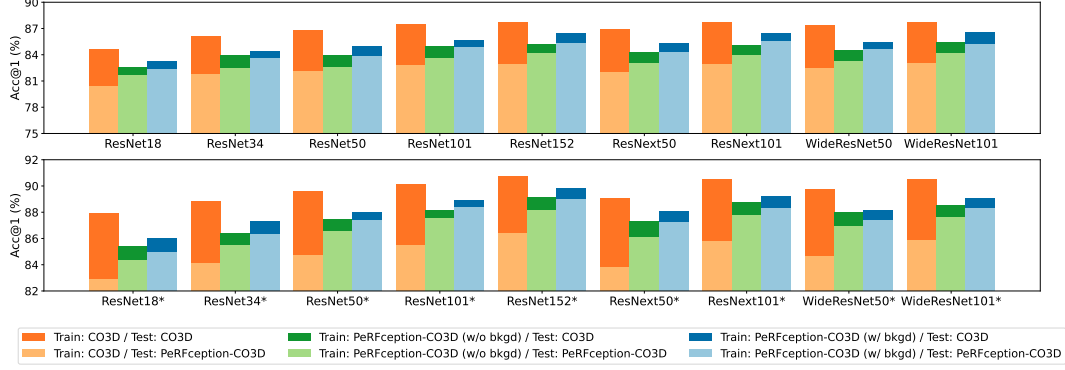


Figure 3: 2D classification accuracies (Acc@1) of the ResNet models trained either on CO3D or PeRFception-CO3D and evaluated either on CO3D or PeRFception-CO3D. * denotes the ImageNet[78] pretrained network. The models trained on PeRFception-CO3D dataset perform well on both CO3D test dataset and PeRFception-CO3D test dataset. Furthermore, the background augmentation of PeRFception-CO3D dataset improves the 2D classification performance. The score table is in the appendix.

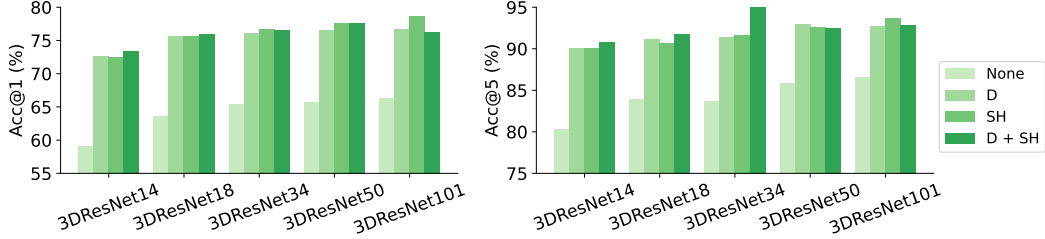


Figure 4: 3D classification performance of 3D ResNet [62] models on our PeRFception-CO3D. We visualize Acc@1 (Left) and Acc@5 (Right) score for each model and input features. "None" denotes the case where 3D classification **D** denotes the density and **SH** denotes spherical harmonic coefficients.

5.1.3 3D classification on PeRFception-CO3D

PeRFception-CO3D provides a novel 3D data representation which we can directly feed into a network without explicit rendering. We train spatially sparse 3D networks on PeRFception-CO3D and visualize the 3D classification accuracy on Figure 4. For each classification model, we utilize four types of input features that are in our Plenoxels representation: ones (None), densities (D), spherical harmonic coefficients (SH), and concatenation of spherical harmonic coefficients and densities (SH + D). One interesting observation is that using either density values or spherical harmonic coefficients as features improves performance much better. We conjecture this is because the density values provide information about where the model should focus more, and the spherical harmonic coefficients explicitly encode visual features.

5.2 Semantic Segmentation on PeRFception-ScanNet dataset

To further verify the fine-grained perception on the large-scale implicit data, we create and evaluate 3D semantic segmentation networks on our scene-centric PeRFception-ScanNet dataset. We assign semantic labels to each voxel by aligning the reconstructed PeRFception-ScanNet data with the provided ground truth point cloud data. Then, for each voxel of PeRFception-ScanNet data, we find the nearest point in the point cloud and when the distance is smaller than the predefined threshold (5cm for our experiments), we assign the class label of the nearest point to the voxel. Otherwise, we set the voxel label to IGNORE_CLASS.

Similar to 3D classification experiments, we use spatially sparse convolutional networks for prediction, but we use U-shaped convnets with varying depth and width for segmentation. For all networks,

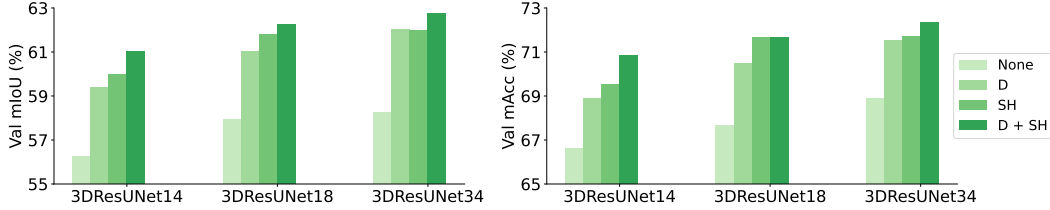


Figure 5: Evaluated semantic segmentation performance, mIoU (Left) and mAcc (Right), on PeRFception-ScanNet validation set with various input features. **D** denotes the density, **SH** denotes spherical harmonic coefficients.

we trained for 60K iterations, with batch size 8, SGD optimizer with initial learning rate 0.1, and cosine annealing scheduler. We train each network with different input features as same with the 3D classification in Sec 5.1 to analyze the effect of the plenoptic features to the 3D semantic segmentation task. We apply geometric augmentation (random rotation, random crop, random affine transform, coordinate dropout, random horizontal flip, random translation, elastic distortion), and feature-level augmentation (random feature jittering).

We use the standard experimental settings following [62], and report mean Intersection over Union (mIoU), mean per-point accuracy (mAcc) on the validation split in Figure 5, and scenewise statistics are in the appendix. We achieve up to 62.77% mIoU and 72.36% mAcc on the validation split of PeRFception-ScanNet, which shows that our PeRFception-ScanNet dataset has accurate geometry for networks to learn semantic of each object class. Consistent with 3D classification, the networks trained with spherical harmonic coefficients and density as input feature exhibit higher segmentation accuracy. This results indicate that the spherical harmonic coefficients and density values provide additional cues for fine-grained geometric perception as well. Additional quantitative and qualitative results are provided in the appendix.

6 Conclusion

In this work, we present the first perception networks for an implicit representation and conduct the comprehensive study of various visual perception tasks. To this end, we created two large-scale implicit datasets, namely PeRFception-CO3D and PeRFception-ScanNet, that cover object-centric and scene-centric environments, respectively. Extensive experiments with diverse perception scenarios, including 2D image classification, 3D object classification, and 3D scene semantic segmentation, show that our datasets effectively convey the same information for both 2D and 3D in a unified and compressed data format. This data format allows eliminating the need to separately store different data formats, 2D images and 3D shapes. Consequently, the required disk space for storage is reduced and the unified data format includes richer features. Furthermore, we propose a novel image augmentation method that was infeasible in image datasets. We expect our fully automatic pipeline should be a great candidate for establishing equally large datasets on 3D to tremendously large 2D image datasets, potentially enabling larger models to be trained.

Limitation. Plenoxels allow high-quality rendering for both indoor and outdoor scenes with fast training and rendering speed. However, the training step of Plenoxels strongly relies on calibrated camera information. The camera parameters would be inaccurately calibrated in the scenes when there are lots of symmetric or textureless patterns. Wrong camera information involves severe artifacts on rendered images, such as the occurrence of floater or geometrically deformed voxel shapes. We believe that jointly optimizing camera poses would be beneficial for improving the fidelity of our dataset. Our work opens up the potential of using radiance field representation in some conventional visual perception tasks and provides the first large-scale radiance field datasets that effectively convey both the 2D and 3D information. We expect future work relevant to more accurate and fast reconstruction could improve our work.

References

- [1] Lars Mescheder, Michael Oechsle, Michael Niemeyer, Sebastian Nowozin, and Andreas Geiger. Occupancy networks: Learning 3d reconstruction in function space. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4460–4470, 2019.
- [2] Songyou Peng, Michael Niemeyer, Lars Mescheder, Marc Pollefeys, and Andreas Geiger. Convolutional occupancy networks. In *European Conference on Computer Vision*, pages 523–540. Springer, 2020.
- [3] Jeong Joon Park, Peter Florence, Julian Straub, Richard Newcombe, and Steven Lovegrove. DeepSDF: Learning continuous signed distance functions for shape representation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 165–174, 2019.
- [4] Towaki Takikawa, Joey Litalien, Kangxue Yin, Karsten Kreis, Charles Loop, Derek Nowrouzezahrai, Alec Jacobson, Morgan McGuire, and Sanja Fidler. Neural geometric level of detail: Real-time rendering with implicit 3d shapes. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11358–11367, 2021.
- [5] 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 *European conference on computer vision*, pages 405–421. Springer, 2020.
- [6] 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. *ICCV*, 2021.
- [7] Kai Zhang, Gernot Riegler, Noah Snaveley, and Vladlen Koltun. Nerf++: Analyzing and improving neural radiance fields. *arXiv preprint arXiv:2010.07492*, 2020.
- [8] Zhiqin Chen and Hao Zhang. Learning implicit fields for generative shape modeling. *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019.
- [9] Shunsuke Saito, Zeng Huang, Ryota Natsume, Shigeo Morishima, Angjoo Kanazawa, and Hao Li. Pifu: Pixel-aligned implicit function for high-resolution clothed human digitization. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 2304–2314, 2019.
- [10] Peter Hedman, Pratul P Srinivasan, Ben Mildenhall, Jonathan T Barron, and Paul Debevec. Baking neural radiance fields for real-time view synthesis. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5875–5884, 2021.
- [11] Lingjie Liu, Jiatao Gu, Kyaw Zaw Lin, Tat-Seng Chua, and Christian Theobalt. Neural sparse voxel fields. *NeurIPS*, 2020.
- [12] Thomas Müller, Alex Evans, Christoph Schied, and Alexander Keller. Instant neural graphics primitives with a multiresolution hash encoding. *arXiv preprint arXiv:2201.05989*, 2022.
- [13] Vincent Sitzmann, Michael Zollhöfer, and Gordon Wetzstein. Scene representation networks: Continuous 3d-structure-aware neural scene representations. *Advances in Neural Information Processing Systems*, 32, 2019.
- [14] Alex Yu, Sara Fridovich-Keil, Matthew Tancik, Qinhong Chen, Benjamin Recht, and Angjoo Kanazawa. Plenoxels: Radiance fields without neural networks. *arXiv preprint arXiv:2112.05131*, 2021.
- [15] Alex Yu, Ruilong Li, Matthew Tancik, Hao Li, Ren Ng, and Angjoo Kanazawa. Plenotrees for real-time rendering of neural radiance fields. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5752–5761, 2021.
- [16] Yoonwoo Jeong, Seokjun Ahn, Christopher Choy, Anima Anandkumar, Minsu Cho, and Jaesik Park. Self-calibrating neural radiance fields. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5846–5854, 2021.
- [17] Zirui Wang, Shangzhe Wu, Weidi Xie, Min Chen, and Victor Adrian Prisacariu. NeRF-: Neural radiance fields without known camera parameters. <https://arxiv.org/abs/2102.07064>, 2021.
- [18] Lin Yen-Chen, Pete Florence, Jonathan T. Barron, Alberto Rodriguez, Phillip Isola, and Tsung-Yi Lin. iNeRF: Inverting neural radiance fields for pose estimation. <https://arxiv.org/abs/2012.05877>, 2020.

- [19] Chen-Hsuan Lin, Wei-Chiu Ma, Antonio Torralba, and Simon Lucey. Barf: Bundle-adjusting neural radiance fields. In *IEEE International Conference on Computer Vision (ICCV)*, 2021.
- [20] Michael Niemeyer and Andreas Geiger. Giraffe: Representing scenes as compositional generative neural feature fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11453–11464, 2021.
- [21] Katja Schwarz, Yiyi Liao, Michael Niemeyer, and Andreas Geiger. Graf: Generative radiance fields for 3D-aware image synthesis. In *Advances in Neural Information Processing Systems (NeurIPS)*, volume 33, 2020.
- [22] Cheng Sun, Min Sun, and Hwann-Tzong Chen. Direct voxel grid optimization: Super-fast convergence for radiance fields reconstruction. *arXiv preprint arXiv:2111.11215*, 2021.
- [23] Jeremy Reizenstein, Roman Shapovalov, Philipp Henzler, Luca Sbordone, Patrick Labatut, and David Novotny. Common objects in 3d: Large-scale learning and evaluation of real-life 3d category reconstruction. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 10901–10911, 2021.
- [24] Angela Dai, Angel X Chang, Manolis Savva, Maciej Halber, Thomas Funkhouser, and Matthias Nießner. Scannet: Richly-annotated 3d reconstructions of indoor scenes. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 5828–5839, 2017.
- [25] Michael Oechsle, Lars Mescheder, Michael Niemeyer, Thilo Strauss, and Andreas Geiger. Texture fields: Learning texture representations in function space. In *International Conference on Computer Vision (ICCV)*, 2019.
- [26] Keunhong Park, Utkarsh Sinha, Jonathan T. Barron, Sofien Bouaziz, Dan B Goldman, Steven M. Seitz, and Ricardo Martin-Brualla. Nerfies: Deformable neural radiance fields. *ICCV*, 2021.
- [27] Michael Niemeyer, Lars Mescheder, Michael Oechsle, and Andreas Geiger. Occupancy flow: 4d reconstruction by learning particle dynamics. In *International Conference on Computer Vision (ICCV)*, 2019.
- [28] Keunhong Park, Utkarsh Sinha, Peter Hedman, Jonathan T. Barron, Sofien Bouaziz, Dan B Goldman, Ricardo Martin-Brualla, and Steven M. Seitz. Hypernerf: A higher-dimensional representation for topologically varying neural radiance fields. *ACM Trans. Graph.*, 40(6), dec 2021.
- [29] Albert Pumarola, Enric Corona, Gerard Pons-Moll, and Francesc Moreno-Noguer. D-nerf: Neural radiance fields for dynamic scenes. *arXiv preprint arXiv:2011.13961*, 2020.
- [30] Alex Trevithick and Bo Yang. Grf: Learning a general radiance field for 3d scene representation and rendering. In *arXiv:2010.04595*, 2020.
- [31] Terrance DeVries, Miguel Angel Bautista, Nitish Srivastava, Graham W. Taylor, and Joshua M. Susskind. Unconstrained scene generation with locally conditioned radiance fields. *arXiv*, 2021.
- [32] Zekun Hao, Arun Mallya, Serge Belongie, and Ming-Yu Liu. GANcraft: Unsupervised 3D Neural Rendering of Minecraft Worlds. In *ICCV*, 2021.
- [33] Jiatao Gu, Lingjie Liu, Peng Wang, and Christian Theobalt. Stylenerf: A style-based 3d aware generator for high-resolution image synthesis. In *International Conference on Learning Representations*, 2022.
- [34] Anpei Chen, Zexiang Xu, Andreas Geiger, , Jingyi Yu, and Hao Su. Tensorf: Tensorial radiance fields, 2022.
- [35] Angel X Chang, Thomas Funkhouser, Leonidas Guibas, Pat Hanrahan, Qixing Huang, Zimo Li, Silvio Savarese, Manolis Savva, Shuran Song, Hao Su, et al. Shapenet: An information-rich 3d model repository. *arXiv preprint arXiv:1512.03012*, 2015.
- [36] Zhirong Wu, Shuran Song, Aditya Khosla, Fisher Yu, Linguang Zhang, Xiaoou Tang, and Jianxiong Xiao. 3d shapenets: A deep representation for volumetric shapes. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1912–1920, 2015.
- [37] Yu Xiang, Roozbeh Mottaghi, and Silvio Savarese. Beyond pascal: A benchmark for 3d object detection in the wild. In *IEEE winter conference on applications of computer vision*, pages 75–82. IEEE, 2014.

- [38] Yu Xiang, Wonhui Kim, Wei Chen, Jingwei Ji, Christopher Choy, Hao Su, Roozbeh Mottaghi, Leonidas Guibas, and Silvio Savarese. Objectnet3d: A large scale database for 3d object recognition. In *European conference on computer vision*, pages 160–176. Springer, 2016.
- [39] Sungjoon Choi, Qian-Yi Zhou, Stephen Miller, and Vladlen Koltun. A large dataset of object scans. *arXiv preprint arXiv:1602.02481*, 2016.
- [40] GoogleResearch. Google scanned objects, August.
- [41] Huan Fu, Rongfei Jia, Lin Gao, Mingming Gong, Binqiang Zhao, Steve Maybank, and Dacheng Tao. 3d-future: 3d furniture shape with texture. *International Journal of Computer Vision*, 129(12):3313–3337, 2021.
- [42] Jasmine Collins, Shubham Goel, Achleshwar Luthra, Leon Xu, Kenan Deng, Xi Zhang, Tomas F Yago Vicente, Himanshu Arora, Thomas Dideriksen, Matthieu Guillaumin, et al. Abo: Dataset and benchmarks for real-world 3d object understanding. *arXiv preprint arXiv:2110.06199*, 2021.
- [43] Shuran Song, Samuel P Lichtenberg, and Jianxiong Xiao. Sun rgb-d: A rgb-d scene understanding benchmark suite. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 567–576, 2015.
- [44] Nathan Silberman, Derek Hoiem, Pushmeet Kohli, and Rob Fergus. Indoor segmentation and support inference from rgb-d images. In *European conference on computer vision*, pages 746–760. Springer, 2012.
- [45] Jianxiong Xiao, Andrew Owens, and Antonio Torralba. Sun3d: A database of big spaces reconstructed using sfm and object labels. In *Proceedings of the IEEE international conference on computer vision*, pages 1625–1632, 2013.
- [46] Iro Armeni, Ozan Sener, Amir R Zamir, Helen Jiang, Ioannis Brilakis, Martin Fischer, and Silvio Savarese. 3d semantic parsing of large-scale indoor spaces. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1534–1543, 2016.
- [47] Iro Armeni, Sasha Sax, Amir R Zamir, and Silvio Savarese. Joint 2d-3d-semantic data for indoor scene understanding. *arXiv preprint arXiv:1702.01105*, 2017.
- [48] Angel Chang, Angela Dai, Thomas Funkhouser, Maciej Halber, Matthias Niessner, Manolis Savva, Shuran Song, Andy Zeng, and Yinda Zhang. Matterport3d: Learning from rgb-d data in indoor environments. *arXiv preprint arXiv:1709.06158*, 2017.
- [49] Julian Straub, Thomas Whelan, Lingni Ma, Yufan Chen, Erik Wijmans, Simon Green, Jakob J Engel, Raul Mur-Artal, Carl Ren, Shobhit Verma, et al. The replica dataset: A digital replica of indoor spaces. *arXiv preprint arXiv:1906.05797*, 2019.
- [50] Charles R Qi, Hao Su, Kaichun Mo, and Leonidas J Guibas. Pointnet: Deep learning on point sets for 3d classification and segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 652–660, 2017.
- [51] Charles R Qi, Li Yi, Hao Su, and Leonidas J Guibas. Pointnet++: Deep hierarchical feature learning on point sets in a metric space. *arXiv preprint arXiv:1706.02413*, 2017.
- [52] Jiageng Mao, Xiaogang Wang, and Hongsheng Li. Interpolated convolutional networks for 3d point cloud understanding. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 1578–1587, 2019.
- [53] Maxim Tatarchenko, Jaesik Park, Vladlen Koltun, and Qian-Yi Zhou. Tangent convolutions for dense prediction in 3d. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3887–3896, 2018.
- [54] Hugues Thomas, Charles R Qi, Jean-Emmanuel Deschaud, Beatriz Marcotegui, François Goulette, and Leonidas J Guibas. Kpconv: Flexible and deformable convolution for point clouds. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 6411–6420, 2019.
- [55] Wenxuan Wu, Zhongang Qi, and Li Fuxin. Pointconv: Deep convolutional networks on 3d point clouds. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9621–9630, 2019.
- [56] Mutian Xu, Runyu Ding, Hengshuang Zhao, and Xiaojuan Qi. Paconv: Position adaptive convolution with dynamic kernel assembling on point clouds. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3173–3182, 2021.

- [57] Yue Wang, Yongbin Sun, Ziwei Liu, Sanjay E Sarma, Michael M Bronstein, and Justin M Solomon. Dynamic graph cnn for learning on point clouds. *Acm Transactions On Graphics (tog)*, 38(5):1–12, 2019.
- [58] Xu Yan, Chaoda Zheng, Zhen Li, Sheng Wang, and Shuguang Cui. Pointasnl: Robust point clouds processing using nonlocal neural networks with adaptive sampling. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5589–5598, 2020.
- [59] Hengshuang Zhao, Li Jiang, Jiaya Jia, Philip HS Torr, and Vladlen Koltun. Point transformer. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 16259–16268, 2021.
- [60] Benjamin Graham. Spatially-sparse convolutional neural networks. *arXiv preprint arXiv:1409.6070*, 2014.
- [61] Haotian* Tang, Zhijian* Liu, Shengyu Zhao, Yujun Lin, Ji Lin, Hanrui Wang, and Song Han. Searching efficient 3d architectures with sparse point-voxel convolution. In *European Conference on Computer Vision*, 2020.
- [62] Christopher Choy, JunYoung Gwak, and Silvio Savarese. 4d spatio-temporal convnets: Minkowski convolutional neural networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3075–3084, 2019.
- [63] Chunghyun Park, Yoonwoo Jeong, Minsu Cho, and Jaesik Park. Fast Point Transformer. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022.
- [64] JunYoung Gwak, Christopher Choy, and Silvio Savarese. Generative sparse detection networks for 3d single-shot object detection. In *Computer Vision—ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part IV 16*, pages 297–313. Springer, 2020.
- [65] Christopher Choy, Jaesik Park, and Vladlen Koltun. Fully convolutional geometric features. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 8958–8966, 2019.
- [66] Christopher Choy, Wei Dong, and Vladlen Koltun. Deep global registration. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 2514–2523, 2020.
- [67] Junha Lee, Seungwook Kim, Minsu Cho, and Jaesik Park. Deep hough voting for robust global registration. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, 2021.
- [68] Shengyu Huang, Zan Gojcic, Mikhail Usvyatsov, and Konrad Schindler Andreas Wieser. Predator: Registration of 3d point clouds with low overlap. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR*, 2021.
- [69] Christopher Choy, Junha Lee, René Ranftl, Jaesik Park, and Vladlen Koltun. High-dimensional convolutional networks for geometric pattern recognition. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 11227–11236, 2020.
- [70] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *European conference on computer vision*, pages 740–755. Springer, 2014.
- [71] Alexander Kirillov, Yuxin Wu, Kaiming He, and Ross Girshick. PointRend: Image segmentation as rendering. 2019.
- [72] Johannes L Schonberger and Jan-Michael Frahm. Structure-from-motion revisited. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4104–4113, 2016.
- [73] Said Pertuz, Domenec Puig, and Miguel Angel Garcia. Analysis of focus measure operators for shape-from-focus. *Pattern Recognition*, 46(5):1415–1432, 2013.
- [74] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.
- [75] Tong He, Zhi Zhang, Hang Zhang, Zhongyue Zhang, Junyuan Xie, and Mu Li. Bag of tricks for image classification with convolutional neural networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 558–567, 2019.

- [76] Saining Xie, Ross Girshick, Piotr Dollár, Zhuowen Tu, and Kaiming He. Aggregated residual transformations for deep neural networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1492–1500, 2017.
- [77] Sergey Zagoruyko and Nikos Komodakis. Wide residual networks. *arXiv preprint arXiv:1605.07146*, 2016.
- [78] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255. Ieee, 2009.
- [79] Aditya Chattopadhyay, Anirban Sarkar, Prantik Howlader, and Vineeth N Balasubramanian. Grad-cam++: Generalized gradient-based visual explanations for deep convolutional networks. In *2018 IEEE winter conference on applications of computer vision (WACV)*, pages 839–847. IEEE, 2018.
- [80] Arno Knapitsch, Jaesik Park, Qian-Yi Zhou, and Vladlen Koltun. Tanks and temples: Benchmarking large-scale scene reconstruction. *ACM Transactions on Graphics (ToG)*, 36(4):1–13, 2017.

Appendix

A.1 Effect of Background Augmentation

Reducing bias on the dataset is one of the fundamental problems in computer vision. A typical example of data bias in object-centric dataset is the backgrounds of image that contain locational context. For instance, couch images are mostly captured in the living room, and motorcycle images are taken outdoor hence the roadway or garage appears in the background most of the time. Such a strong correlation between background information and the class labels makes the classification networks biased toward background so that the networks make predictions by looking at the backgrounds rather than the foreground objects. Here, we demonstrate that our datasets, with additional augmentations, are more resistant to such bias. We use a popular analysis tool, Grad-CAM++ [79], to verify that the classification model trained with our dataset avoids overfitting on backgrounds. We use the official Grad-CAM++ implementation and feed test images. According to Figure a.1, the classification model trained on our dataset focuses on the object, whereas the classification model trained on the original CO3D focuses on the background for many outdoor scenes.

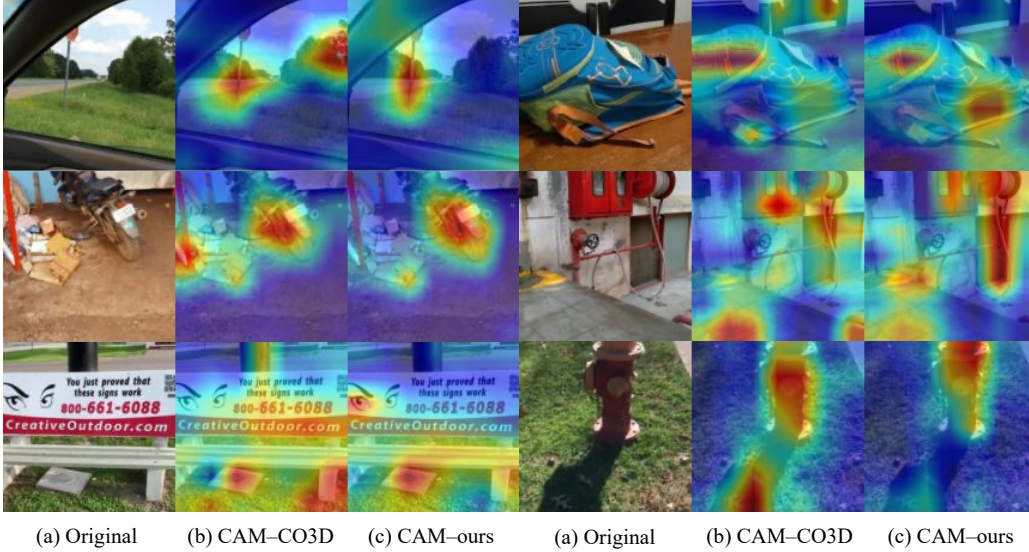


Figure a.1: Visualization of class activation maps generated by GradCAM++. The left and right activation map are generated from the ResNet18 trained with the original CO3D images and PerFception-CO3D images, respectively.

We also conduct control experiments about the strength of background augmentation. In detail, we compare evaluated performance by changing the probability, *i.e.*, 0%, 50%, and 100%, to apply the background augmentation for each iteration. In comparison to the results without background augmentation, the background augmentation with a probability of 50% increases the 2D classification accuracy; whereas the 100% probability of background augmentation does not improve the performance. In this ablation study, we determine the appropriate level of background augmentation assistance for improving 2D classification accuracy. Table a.1 reports the results of the control experiments.

We also visualize several examples of background augmented images on Figure a.4.

A.2 Implementation Details

Here, we describe the thorough details about the data generation with Plenoxels and the data augmentation methods proposed in our main paper. We also provide the architectural details of 3D semantic segmentation models trained on PerFception-ScanNet dataset.

Table a.1: 2D classification accuracies (Acc@1/Acc@5) of the ResNet model trained either on CO3D or PeRFception-CO3D and evaluated either on CO3D or PeRFception-CO3D. * denotes the ImageNet[78] pretrained network. PeRF-CO3D is an abbreviation of PeRFception-CO3D. p stands for the probability to apply background augmentation.

Train Dataset	CO3D		PeRF-CO3D (p=0.0)		PeRF-CO3D (p=0.5)		PeRF-CO3D (p=1.0)	
Test Dataset	CO3D	PeRF-CO3D	CO3D	PeRF-CO3D	CO3D	PeRF-CO3D	CO3D	PeRF-CO3D
Acc@1 ($\mu \pm \sigma$)								
ResNet18	84.74 \pm 0.04	80.39 \pm 0.19	82.56 \pm 0.24	81.43 \pm 0.14	83.15 \pm 0.10	82.05 \pm 0.24	81.93 \pm 0.10	81.31 \pm 0.12
ResNet34	86.18 \pm 0.05	81.62 \pm 0.10	83.98 \pm 0.07	82.71 \pm 0.17	84.62 \pm 0.14	83.61 \pm 0.04	83.40 \pm 0.10	82.54 \pm 0.22
ResNet50	86.83 \pm 0.02	82.19 \pm 0.09	84.12 \pm 0.28	82.80 \pm 0.37	84.83 \pm 0.14	83.77 \pm 0.08	83.60 \pm 0.11	82.73 \pm 0.04
ResNet101	87.42 \pm 0.04	82.86 \pm 0.06	84.97 \pm 0.13	83.71 \pm 0.13	85.95 \pm 0.22	85.11 \pm 0.23	85.01 \pm 0.08	83.82 \pm 0.05
ResNet152	87.75 \pm 0.11	83.04 \pm 0.11	85.67 \pm 0.31	84.35 \pm 0.19	86.40 \pm 0.04	85.28 \pm 0.02	85.11 \pm 0.09	84.06 \pm 0.09
ResNext50	86.84 \pm 0.29	81.99 \pm 0.14	84.21 \pm 0.08	83.04 \pm 0.04	85.25 \pm 0.16	84.32 \pm 0.18	84.06 \pm 0.38	83.23 \pm 0.28
ResNext101	87.73 \pm 0.03	82.90 \pm 0.07	85.22 \pm 0.09	84.17 \pm 0.15	86.37 \pm 0.10	85.48 \pm 0.06	85.46 \pm 0.16	84.46 \pm 0.11
WideResNet50	87.32 \pm 0.05	82.45 \pm 0.11	84.68 \pm 0.15	83.64 \pm 0.27	85.58 \pm 0.10	84.68 \pm 0.02	84.35 \pm 0.12	83.46 \pm 0.24
WideResNet101	87.80 \pm 0.12	83.04 \pm 0.10	85.36 \pm 0.26	84.15 \pm 0.19	86.32 \pm 0.21	85.30 \pm 0.11	85.45 \pm 0.13	84.23 \pm 0.23
ResNet18*	87.75 \pm 0.18	82.70 \pm 0.15	85.28 \pm 0.14	84.09 \pm 0.20	85.97 \pm 0.16	84.97 \pm 0.13	84.62 \pm 0.02	83.62 \pm 0.13
ResNet34*	88.83 \pm 0.06	84.01 \pm 0.12	86.60 \pm 0.27	85.43 \pm 0.22	87.19 \pm 0.31	86.25 \pm 0.19	86.09 \pm 0.11	85.01 \pm 0.09
ResNet50*	89.51 \pm 0.21	84.76 \pm 0.04	87.49 \pm 0.11	86.60 \pm 0.08	88.12 \pm 0.08	87.30 \pm 0.08	87.09 \pm 0.12	86.25 \pm 0.17
ResNet101*	90.21 \pm 0.09	85.60 \pm 0.16	88.39 \pm 0.22	87.46 \pm 0.17	89.00 \pm 0.07	88.32 \pm 0.13	88.28 \pm 0.07	87.26 \pm 0.01
ResNet152*	90.60 \pm 0.10	86.26 \pm 0.10	89.17 \pm 0.15	88.19 \pm 0.03	89.52 \pm 0.20	88.73 \pm 0.15	88.63 \pm 0.20	87.59 \pm 0.13
ResNext50*	89.21 \pm 0.14	83.90 \pm 0.28	87.28 \pm 0.06	86.28 \pm 0.14	87.82 \pm 0.19	87.30 \pm 0.06	86.81 \pm 0.27	86.09 \pm 0.25
ResNext101*	90.52 \pm 0.07	85.91 \pm 0.13	88.51 \pm 0.19	87.82 \pm 0.06	89.17 \pm 0.07	88.51 \pm 0.16	88.57 \pm 0.02	87.60 \pm 0.18
WideResNet50*	89.78 \pm 0.06	85.13 \pm 0.33	87.80 \pm 0.18	86.88 \pm 0.06	88.23 \pm 0.10	87.75 \pm 0.25	87.12 \pm 0.11	86.50 \pm 0.08
WideResNet101*	90.45 \pm 0.10	85.97 \pm 0.09	88.53 \pm 0.04	87.59 \pm 0.12	89.22 \pm 0.13	88.39 \pm 0.07	88.10 \pm 0.11	87.20 \pm 0.10
Acc@5 ($\mu \pm \sigma$)								
ResNet18	96.24 \pm 0.17	94.19 \pm 0.16	95.55 \pm 0.18	95.00 \pm 0.15	95.70 \pm 0.11	95.37 \pm 0.05	95.30 \pm 0.08	94.85 \pm 0.09
ResNet34	96.68 \pm 0.11	94.62 \pm 0.11	95.88 \pm 0.07	95.37 \pm 0.10	96.23 \pm 0.03	95.89 \pm 0.06	95.74 \pm 0.11	95.29 \pm 0.07
ResNet50	96.90 \pm 0.07	94.67 \pm 0.14	96.10 \pm 0.02	95.51 \pm 0.13	96.41 \pm 0.11	95.99 \pm 0.08	96.01 \pm 0.07	95.48 \pm 0.04
ResNet101	97.04 \pm 0.03	94.78 \pm 0.05	96.23 \pm 0.07	95.58 \pm 0.09	96.76 \pm 0.10	96.32 \pm 0.12	96.39 \pm 0.03	95.74 \pm 0.14
ResNet152	97.14 \pm 0.04	94.94 \pm 0.09	96.32 \pm 0.09	95.73 \pm 0.05	96.86 \pm 0.06	96.39 \pm 0.06	96.46 \pm 0.07	95.92 \pm 0.09
ResNext50	96.75 \pm 0.12	94.54 \pm 0.10	95.95 \pm 0.15	95.32 \pm 0.13	96.49 \pm 0.10	95.92 \pm 0.16	96.08 \pm 0.12	95.47 \pm 0.15
ResNext101	96.98 \pm 0.07	94.71 \pm 0.12	96.09 \pm 0.11	95.54 \pm 0.07	96.71 \pm 0.01	96.26 \pm 0.03	96.49 \pm 0.04	95.80 \pm 0.03
WideResNet50	96.84 \pm 0.13	94.60 \pm 0.10	96.05 \pm 0.04	95.52 \pm 0.05	96.44 \pm 0.11	96.03 \pm 0.03	96.10 \pm 0.16	95.55 \pm 0.01
WideResNet101	97.05 \pm 0.01	94.91 \pm 0.08	96.26 \pm 0.07	95.70 \pm 0.08	96.86 \pm 0.11	96.31 \pm 0.10	96.49 \pm 0.06	95.89 \pm 0.01
ResNet18*	97.13 \pm 0.10	95.04 \pm 0.16	96.39 \pm 0.08	95.89 \pm 0.03	96.73 \pm 0.04	96.24 \pm 0.09	96.19 \pm 0.08	95.67 \pm 0.10
ResNet34*	97.39 \pm 0.14	95.28 \pm 0.15	96.67 \pm 0.08	96.18 \pm 0.04	97.00 \pm 0.13	96.50 \pm 0.09	96.61 \pm 0.05	96.02 \pm 0.03
ResNet50*	97.60 \pm 0.04	95.46 \pm 0.02	96.84 \pm 0.04	96.39 \pm 0.07	97.11 \pm 0.05	96.69 \pm 0.08	96.95 \pm 0.06	96.37 \pm 0.13
ResNet101*	97.90 \pm 0.03	95.86 \pm 0.05	97.06 \pm 0.19	96.74 \pm 0.11	97.48 \pm 0.02	97.13 \pm 0.05	97.16 \pm 0.04	96.57 \pm 0.03
ResNet152*	97.97 \pm 0.04	95.97 \pm 0.04	97.32 \pm 0.17	96.87 \pm 0.13	97.66 \pm 0.07	97.24 \pm 0.08	97.22 \pm 0.17	96.69 \pm 0.09
ResNext50*	97.33 \pm 0.09	94.99 \pm 0.11	96.54 \pm 0.09	96.04 \pm 0.02	97.08 \pm 0.14	96.66 \pm 0.07	96.57 \pm 0.12	96.03 \pm 0.04
ResNext101*	97.67 \pm 0.03	95.51 \pm 0.06	96.90 \pm 0.02	96.61 \pm 0.03	97.32 \pm 0.05	96.93 \pm 0.07	97.07 \pm 0.04	96.52 \pm 0.05
WideResNet50*	97.60 \pm 0.07	95.53 \pm 0.17	96.75 \pm 0.10	96.39 \pm 0.08	97.13 \pm 0.02	96.84 \pm 0.09	96.72 \pm 0.07	96.22 \pm 0.05
WideResNet101*	97.76 \pm 0.07	95.67 \pm 0.04	97.01 \pm 0.05	96.62 \pm 0.07	97.48 \pm 0.07	97.10 \pm 0.06	96.93 \pm 0.09	96.48 \pm 0.06

A.2.1 Data Generation

We use the official implementation of Plenoxels [14], which is implemented in PyTorch with custom CUDA kernels. With a single RTX 3090 GPU, it takes up to 30 minutes per scene.

PeRFception-CO3D. Plenoxels uses a different learning rate for each parameter. Following the default configuration for rendering Tanks and Temples [80] dataset, we use the learning rate 30 for density values and 10^{-2} for spherical harmonic coefficients. We skip foreground rendering for 1,000 steps at the beginning of training for stable training. For total variation losses, we set the weight λ as 5.0×10^{-5} for foreground densities, 5.0×10^{-3} for foreground spherical harmonics, 1.0×10^{-3} for background colors, and 1.0×10^{-3} for background densities. The background brightness is set to 0.5. We set the lambda value to 10^{-5} for the beta loss and 10^{-10} for the sparsity loss, which is 10 times larger than the default configuration. We use the 9-dimensional spherical harmonic coefficients for each RGB channel. We prune voxels whose density values are below 1.28 during the upsampling step.

PeRFception-ScanNet. As we described in Section 4.2 in our main paper, we initialize the voxel grid with unprojected depth maps for PeRFception-ScanNet dataset generation. The unprojected depth maps are aggregated into a point cloud using the ground truth camera parameters. Since the ScanNet depth maps often contain noisy observations, the resulting point cloud is contaminated with outlier points. To filter out the outliers, we discretize the point cloud into coarse voxels with discretization size of 5cm and perform connected component analysis to detect disconnected voxels which are

likely to be outliers. The refined point cloud is then utilized as the initialization for Plenoxels. The voxel grid is initialized with 256^3 resolution and trained for 51,200 iterations and at 25,600th iteration we perform voxel pruning with the density threshold of 1.28. All the other hyperparameters are same with CO3D but removed background rendering since all the points of ScanNet datasets should be considered as foregrounds.

The input images are resized into 640*480 resolution. We exclude 16 scenes out of 1,513 training scenes from ScanNet that are severely affected by motion blurs. Detecting blurriness was done automatically by OpenCV blur detection algorithm with the threshold of 10. We consider a scene defective when the number of frames is too low after filtering blur images.

A.2.2 Random Pose Selection

Plenoxels show great rendering quality when the camera pose is close to the camera poses in the train set. Unfortunately, it fails to reliably render frames when attempting to render extremely unobserved parts, which are significantly out of train-view coverage. Thus, selecting appropriate poses is crucial for photorealistic image generation. We propose an operation that selects an intermediate pose from two input poses. The operator AA2R converts the axis-angle representation, i.e., (v, θ) where v stands for the axis and θ stands for the angle, to the rotation matrix. R2AA is exactly the inverse operation of AA2R. We compute the distance between the selected poses with distance functions D_R and D_t to choose a pair of poses that are close enough. We use the rotation threshold $\theta = \frac{1}{24}\pi$ and the translation threshold 0.5.

$$\text{ReduceAngle}(\mathbf{R}, s) = \text{AA2R}(\text{R2AA}(\mathbf{R})_v, s \cdot \text{R2AA}(\mathbf{R})_\theta)$$

$$\text{IntermediateRot}(\mathbf{R}_1, \mathbf{R}_2, s) = \text{ReduceAngle}(\mathbf{R}_2 \mathbf{R}_1^{-1}, s) \mathbf{R}_1$$

$$D_R(\mathbf{R}_1, \mathbf{R}_2) = \arccos(\text{Tr}(\mathbf{R}_2^{-1} \mathbf{R}_1)/2 - 1)$$

$$D_t(\mathbf{t}_1, \mathbf{t}_2) = \|\mathbf{t}_2 - \mathbf{t}_1\|_2^2$$

The main concept of the Algo 1 is to select an intermediate camera pose between randomly selected camera poses in the train set. For intrinsic parameters and image shape, we randomly select intrinsic matrix among intrinsic matrices in train set and use the corresponding image shape.

Algorithm 1 Random Pose Generation

Input
 $P = \{\{\mathbf{R}_i, \mathbf{t}_i\}\}_{i=1,2,\dots,N}$: The set of poses in train split
 θ : The rotation distance threshold
 d : The translation distance threshold
Output
 $\mathbf{p}_{\text{out}} \in \mathbb{R}^{4 \times 4}$: The output pose

s : A random value from $[0, 1]$
 $j, k = \text{RandomIndex}(N), \text{RandomIndex}(N)$
while $D_R(\mathbf{R}_j, \mathbf{R}_k) \geq \theta$ and $D_t(\mathbf{t}_j, \mathbf{t}_k) \geq d$ **do**
 $j, k = \text{RandomIndex}(N), \text{RandomIndex}(N)$
end while
 $\mathbf{R}_{jk} = \text{IntermediateRot}(\mathbf{R}_j, \mathbf{R}_k, s)$
 $\mathbf{p}_{\text{out}}[:3, :3] = \mathbf{R}_{jk}$
 $\mathbf{p}_{\text{out}}[:3, 3] = s\mathbf{t}_k + (1 - s)\mathbf{t}_j$
 $\mathbf{p}_{\text{out}}[3, 3] = 1$
return \mathbf{p}_{out}

A.2.3 Architectural Details

We train three variants of 3D ResUNets implemented with sparse convolutional layers [62] for semantic segmentation on PerFception-ScanNet dataset. The layer-wise architectural details are depicted in Table a.2.

Table a.2: Architectures of Res16UNet variants for semantic segmentation on PeRFception-ScanNet. We denote a convolution layer with its *kernel size*, *output channel size*, and *convolution stride size*. All convolution layers except for the last layer have a Batch Normalization and a ReLU layer after them. The layers with the tag "conv_tr" indicates the transposed convolution layers. We use a square bracket to denote a residual block, with the number of blocks stacked.

layer name	Res16UNet14A	Res16UNet18A	Res16UNet34C
init	$3^3, 32, \text{stride } 1$		
conv1	$2^3, 32, \text{stride } 2$		
block1	$\begin{bmatrix} 3^3, 32, \text{stride } 1 \\ 3^3, 32, \text{stride } 1 \end{bmatrix} \times 1$	$\begin{bmatrix} 3^3, 32, \text{stride } 1 \\ 3^3, 32, \text{stride } 1 \end{bmatrix} \times 2$	$\begin{bmatrix} 3^3, 32, \text{stride } 1 \\ 3^3, 32, \text{stride } 1 \end{bmatrix} \times 2$
conv2	$2^3, 32, \text{stride } 2$		
block2	$\begin{bmatrix} 3^3, 64, \text{stride } 1 \\ 3^3, 64, \text{stride } 1 \end{bmatrix} \times 1$	$\begin{bmatrix} 3^3, 64, \text{stride } 1 \\ 3^3, 64, \text{stride } 1 \end{bmatrix} \times 2$	$\begin{bmatrix} 3^3, 64, \text{stride } 1 \\ 3^3, 64, \text{stride } 1 \end{bmatrix} \times 3$
conv3	$2^3, 64, \text{stride } 2$		
block3	$\begin{bmatrix} 3^3, 128, \text{stride } 1 \\ 3^3, 128, \text{stride } 1 \end{bmatrix} \times 1$	$\begin{bmatrix} 3^3, 128, \text{stride } 1 \\ 3^3, 128, \text{stride } 1 \end{bmatrix} \times 2$	$\begin{bmatrix} 3^3, 128, \text{stride } 1 \\ 3^3, 128, \text{stride } 1 \end{bmatrix} \times 4$
conv4	$2^2, 128, \text{stride } 2$		
block4	$\begin{bmatrix} 3^3, 256, \text{stride } 1 \\ 3^3, 256, \text{stride } 1 \end{bmatrix} \times 1$	$\begin{bmatrix} 3^3, 256, \text{stride } 1 \\ 3^3, 256, \text{stride } 1 \end{bmatrix} \times 2$	$\begin{bmatrix} 3^3, 256, \text{stride } 1 \\ 3^3, 256, \text{stride } 1 \end{bmatrix} \times 6$
conv4_tr	$2^3, 128, \text{stride } 2$		$2^3, 256, \text{stride } 2$
block5	$\begin{bmatrix} 3^3, 128, \text{stride } 1 \\ 3^3, 128, \text{stride } 1 \end{bmatrix} \times 1$	$\begin{bmatrix} 3^3, 128, \text{stride } 1 \\ 3^3, 128, \text{stride } 1 \end{bmatrix} \times 2$	$\begin{bmatrix} 3^3, 256, \text{stride } 1 \\ 3^3, 256, \text{stride } 1 \end{bmatrix} \times 2$
conv5_tr	$2^2, 128, \text{stride } 2$		
block6	$\begin{bmatrix} 3^3, 128, \text{stride } 1 \\ 3^3, 128, \text{stride } 1 \end{bmatrix} \times 1$	$\begin{bmatrix} 3^3, 128, \text{stride } 1 \\ 3^3, 128, \text{stride } 1 \end{bmatrix} \times 2$	$\begin{bmatrix} 3^3, 128, \text{stride } 1 \\ 3^3, 128, \text{stride } 1 \end{bmatrix} \times 2$
conv6_tr	$2^3, 96, \text{stride } 2$		
block7	$\begin{bmatrix} 3^3, 96, \text{stride } 1 \\ 3^3, 96, \text{stride } 1 \end{bmatrix} \times 1$	$\begin{bmatrix} 3^3, 96, \text{stride } 1 \\ 3^3, 96, \text{stride } 1 \end{bmatrix} \times 2$	$\begin{bmatrix} 3^3, 96, \text{stride } 1 \\ 3^3, 96, \text{stride } 1 \end{bmatrix} \times 2$
conv7_tr	$2^3, 96, \text{stride } 2$		
block8	$\begin{bmatrix} 3^3, 96, \text{stride } 1 \\ 3^3, 96, \text{stride } 1 \end{bmatrix} \times 1$	$\begin{bmatrix} 3^3, 96, \text{stride } 1 \\ 3^3, 96, \text{stride } 1 \end{bmatrix} \times 2$	$\begin{bmatrix} 3^3, 96, \text{stride } 1 \\ 3^3, 96, \text{stride } 1 \end{bmatrix} \times 2$
final	$1^3, 20, \text{stride } 1$		

A.3 Dataset Statistics

PeRFception-CO3D dataset includes 18,618 object-centric scenes that cover all of the original CO3D dataset scenes, only except those with incorrect camera poses. It contains a fair quantity of 51 different object-class labeled scenes; see Figure a.2.

We also analyze the *resolution vs. rendering quality* trade-off on our PeRFception-ScanNet dataset. To measure the trade-off rates, we train Plenoxels with lower (128) and higher (512) resolutions than the default configuration (256) on a randomly selected subset of ScanNet scenes.

As reported in Table a.3 and Figure a.3, there is no direct correlation between the resolution and the average rendering quality on ScanNet reconstruction. This could be attributed to various non-trivial factors. We visualized the histogram of PSNR distribution on the ScanNet dataset in Figure a.3. The X-axis represents the PSNR score, and Y-axis represents the percentage of scenes. Note that the PSNR distribution of the higher resolution Plenoxel reconstructions exhibits fat-tailed distribution, whereas the lower resolution reconstructions show the long-tailed distribution. We speculate that this is due to the fact that higher resolution reconstruction results in each voxel learning spherical harmonics parameters from fewer rays. Thus, errors in camera parameters or motion blur would result in larger errors for smaller voxels as parameters are learned from fewer rays. Thus, the rendering

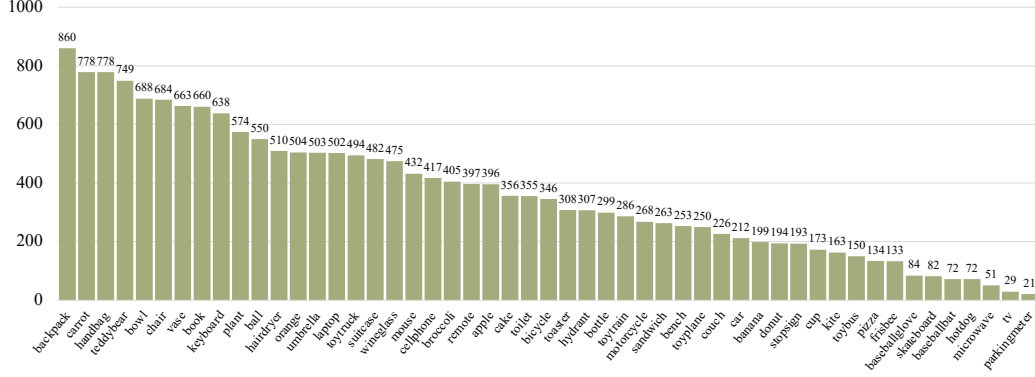


Figure a.2: Statistics of PerFception-CO3D dataset scene per object-class. The y-axis visualizes the number of scenes for each class.

Table a.3: PSNR of PerFception-ScanNet dataset with varying resolutions. PCTL stands for *percentile*.

Reso	Mem (GB)	Avg. PSNR (dB)	50th PCTL	75th PCTL	90th PCTL	95th PCTL
128	28.7	23.20 \pm 3.69	23.59	26.14	27.94	29.12
256	43.8	23.01 \pm 3.96	23.38	26.10	28.02	29.20
512	113.8	22.94 \pm 4.26	23.22	26.35	28.34	29.49

quality increases for the scenes with accurate camera poses and less motion blur, while it decreases for noisy scenes. We conjecture that higher resolution reconstruction would yield better performance if we have high-resolution images with accurate camera poses.

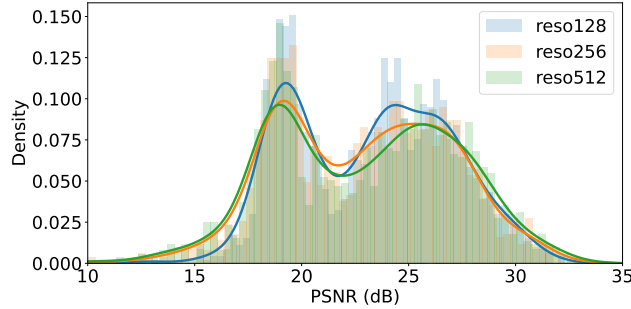


Figure a.3: Distributions of PSNR values of PerFception-ScanNet dataset with varying resolutions.

A.4 Camera Manipulation

Camera manipulation has been one of the inaccessible augmentation techniques on conventional image datasets. In contrast, our implicit data can be rendered from many continuous viewpoints and with arbitrary camera parameters (*e.g.*, intrinsics, extrinsics, and radial distortions), allowing us to generate images with non-standard but realistic camera-level augmentations. We have not adopted this augmentation skill while training 2D classification networks since most images do not have distortions. We perform camera intrinsics augmentation and camera distortion augmentation and Figure a.4 visualizes several examples of such camera manipulations.

A.5 Quantitative Results

In this section, we provide additional quantitative results that are not included in the main paper.

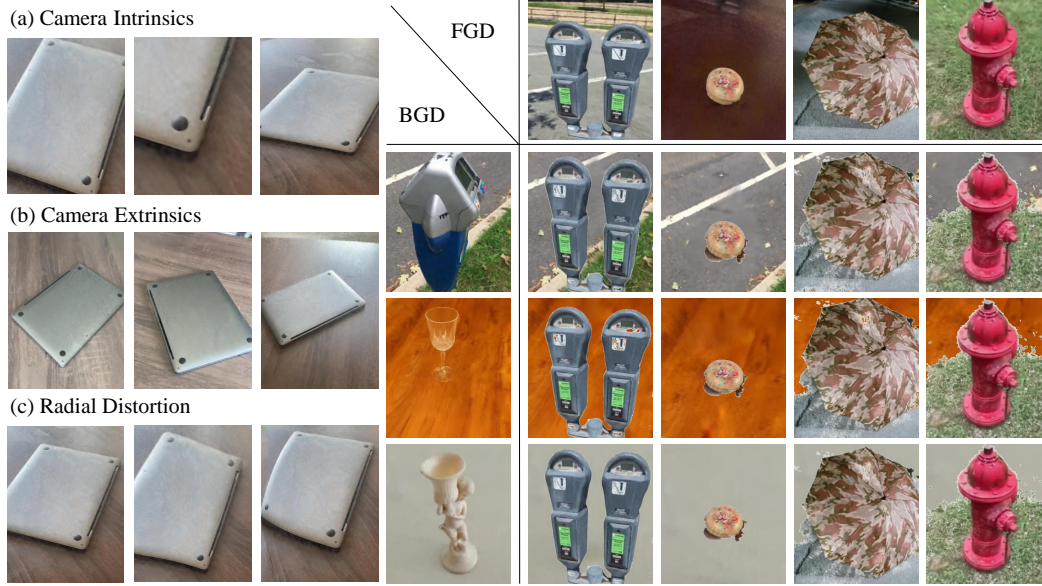


Figure a.4: Examples of camera manipulation and background augmentations. FGD and BGD are source images for foreground and background respectively.

Classwise Rendering Quality. We provide classwise PSNR, SSIM, and LPIPS scores of PeRFception-CO3D dataset in Table a.4.

2D and 3D Perception Tasks. Table a.1 reports the results of the 2D classification experiments. In Table a.5, we report the quantitative classification results on PeRFception-CO3D dataset. In Table a.7, we report the scenewise semantic segmentation results on PeRFception-ScanNet dataset. For all experiments, we perform three experiments with different random seed values and report the mean and standard deviation of evaluation metrics.

	apple	backp	ball	banana	bbb	bbg	bench	bicy	book	bottle	bowl	brocc	cake
PSNR	29.56	28.99	29.06	30.62	29.21	30.03	27.01	27.67	28.88	29.85	31.83	29.67	29.02
SSIM	0.8767	0.8702	0.8661	0.8951	0.8691	0.8898	0.7883	0.8246	0.8746	0.8671	0.8717	0.8584	0.8695
LPIPS	0.3263	0.3349	0.3614	0.2591	0.3566	0.3347	0.3471	0.3416	0.3453	0.3110	0.3135	0.3479	0.3452
	car	carrot	cellp	chair	couch	cup	donut	frisbee	hairdryer	handbag	hotdog	hydrant	kbd
PSNR	23.63	28.78	29.33	28.32	29.31	28.83	28.78	30.06	29.12	28.39	28.95	26.12	29.74
SSIM	0.7259	0.8573	0.8603	0.8524	0.8727	0.8628	0.8695	0.8720	0.8706	0.8622	0.8738	0.7428	0.8789
LPIPS	0.4003	0.3609	0.3511	0.3475	0.3358	0.3602	0.3484	0.3386	0.3351	0.3367	0.3639	0.3416	0.3301
	kite	laptop	micro	motor	mouse	orange	park	pizza	plant	remote	sandwch	skb	stop
PSNR	29.04	27.99	27.75	25.61	29.04	28.55	24.69	29.00	28.63	29.23	29.00	28.57	23.48
SSIM	0.8693	0.8535	0.8506	0.7772	0.8687	0.8521	0.7420	0.8687	0.8452	0.8667	0.8738	0.8458	0.7234
LPIPS	0.3293	0.3530	0.3650	0.3641	0.3486	0.3469	0.3835	0.3214	0.3436	0.3668	0.3459	0.3536	0.4121
	suit	teddy	toast	toil	tbus	tplane	ttrain	ttruck	tv	umb	vase	wine	overall
PSNR	28.87	29.39	28.05	32.61	28.15	29.30	28.42	28.63	27.69	28.62	28.47	28.13	28.82
SSIM	0.8609	0.8713	0.8564	0.9297	0.8549	0.8703	0.8438	0.8627	0.8709	0.8504	0.8536	0.8420	0.8564
LPIPS	0.3526	0.3345	0.3466	0.2841	0.3579	0.3537	0.3765	0.3546	0.3865	0.3601	0.3570	0.3820	0.3451

Table a.4: Classwise rendering quality of PeRFception-CO3D dataset.

A.6 Qualitative Results

We provide additional qualitative results introduced in the main paper. Figure a.5 illustrates several examples of PeRFception-CO3D with accurate geometric information and photorealistic rendering quality. Figure a.6, a.7 and a.8 compare visual reconstruction ability of CO3D and PeRFception-CO3D. Figure a.9, a.10, and a.11 visualize rendered novel views and their corresponding error maps. In Figure a.12, we provide the qualitative results of rendered novel views and their corresponding error maps on PeRFception-ScanNet dataset. Figure a.13 visualizes reconstructed sparse voxels of

Table a.5: Score tables for 3D classification networks. We repeat each experiment for 3 times and report the mean (μ) and standard deviation (σ). **D** denotes density and **SH** denotes spherical harmonic coefficients.

Acc@1 ($\mu \pm \sigma$)				
Input Feature	None	D	SH	SH + D
3DResNet14	59.36 \pm 0.30	72.44 \pm 0.29	71.87 \pm 0.61	72.92 \pm 0.42
3DResNet18	63.85 \pm 0.33	75.18 \pm 0.70	75.58 \pm 0.37	75.72 \pm 0.25
3DResNet34	64.55 \pm 0.84	76.38 \pm 0.34	76.51 \pm 0.61	76.50 \pm 0.03
3DResNet50	65.25 \pm 0.75	76.42 \pm 0.19	77.59 \pm 0.17	77.53 \pm 0.27
3DResNet101	66.21 \pm 0.93	77.27 \pm 0.61	78.04 \pm 0.58	77.19 \pm 0.89
Acc@5 ($\mu \pm \sigma$)				
Input Feature	None	D	SH	SH + D
3DResNet14	81.30 \pm 0.87	90.04 \pm 0.13	89.60 \pm 0.36	90.83 \pm 0.03
3DResNet18	84.47 \pm 0.54	91.10 \pm 0.24	90.98 \pm 0.40	91.54 \pm 0.21
3DResNet34	84.26 \pm 0.47	91.79 \pm 0.61	91.79 \pm 0.49	91.98 \pm 0.08
3DResNet50	85.71 \pm 0.64	92.91 \pm 0.24	92.73 \pm 0.21	92.93 \pm 0.54
3DResNet101	86.50 \pm 0.09	92.68 \pm 0.34	93.24 \pm 0.49	93.09 \pm 0.21

Table a.6: Evaluated semantic segmentation performance on PeRFception-ScanNet dataset without(left) and with(right) spherical harmonic coefficients as input feature. **SH** denotes spherical harmonic coefficients.

mIoU ($\mu \pm \sigma$)				
Input Feature	None	D	SH	SH + D
3D-Res16UNet14A [62]	57.6 \pm 2.0	60.4 \pm 0.1	60.2 \pm 0.2	60.3 \pm 0.1
3D-Res16UNet18A [62]	59.5 \pm 2.0	61.5 \pm 0.6	61.8 \pm 0.4	61.7 \pm 0.5
3D-Res16UNet34C [62]	60.2 \pm 2.3	62.2 \pm 0.9	62.5 \pm 0.7	62.5 \pm 0.7
mAcc ($\mu \pm \sigma$)				
Input Feature	None	D	SH	SH + D
3D-Res16UNet14A [62]	67.9 \pm 1.8	70.4 \pm 0.0	69.9 \pm 0.4	70.0 \pm 0.3
3D-Res16UNet18A [62]	69.7 \pm 1.7	71.7 \pm 0.4	71.7 \pm 0.3	71.4 \pm 0.4
3D-Res16UNet34C [62]	70.1 \pm 1.9	72.1 \pm 0.5	72.0 \pm 0.6	72.2 \pm 0.4

PeRFception-ScanNet with predicted semantic labels. We have attached an additional video to show our photorealistic rendering.

Table a.7: Scenewise statistics of semantic segmentation networks on the PeRFception-ScanNet validation split. All experiments are performed with three different random seed values and the mean and standard deviation are reported.

IoU	D	SH	bath	bed	bksf	cab	chair	cntr	curt	desk	door	floor
14A			74.3±3.0	72.5±2.8	57.9±4.6	57.1±0.3	80.6±0.8	52.2±3.0	22.3±3.3	55.1±2.2	36.0±5.3	94.1±0.3
14A	✓		79.5±1.3	73.2±2.3	64.8±0.8	58.0±0.9	81.5±0.8	51.5±3.6	28.6±1.8	57.9±0.5	44.3±0.6	94.3±0.2
14A	✓	✓	77.0±0.9	74.5±0.7	64.6±1.0	57.6±0.7	81.5±0.2	51.8±2.9	23.7±1.7	57.9±2.5	43.3±0.3	94.5±0.0
14A	✓	✓	76.5±1.2	75.9±0.6	64.9±1.6	57.1±1.0	82.0±0.2	50.5±2.5	22.0±2.9	57.8±1.2	43.6±0.2	94.5±0.1
18A			77.3±2.4	73.5±2.5	59.6±3.8	59.8±1.1	81.4±1.2	58.2±2.0	21.8±2.1	58.0±2.7	39.7±6.1	94.2±0.3
18A	✓		77.3±2.2	73.0±2.8	64.9±0.4	58.6±1.4	82.6±0.5	57.3±2.1	31.4±5.6	60.0±2.1	47.7±0.9	94.4±0.2
18A	✓	✓	80.6±1.0	75.3±1.6	66.1±1.5	59.3±1.4	82.0±1.0	59.6±1.1	23.9±1.8	59.9±1.2	46.9±1.0	94.6±0.0
18A	✓	✓	80.3±3.1	75.9±0.8	66.0±1.0	59.6±1.0	82.5±0.7	55.3±3.7	26.2±2.7	60.8±0.6	46.5±1.3	94.6±0.1
34C			77.8±4.9	73.5±1.8	60.8±4.1	60.0±1.0	81.7±1.6	58.3±2.4	25.0±1.4	58.6±3.6	40.3±8.1	94.2±0.3
34C	✓		81.0±1.9	73.4±1.8	65.2±0.9	60.2±0.8	82.8±0.8	56.7±0.4	30.6±3.6	62.6±1.1	48.7±3.0	94.3±0.2
34C	✓	✓	79.9±3.1	76.3±0.3	66.2±0.3	59.6±1.3	83.7±0.4	55.9±2.7	25.8±0.4	62.2±1.5	48.6±2.9	94.6±0.1
34C	✓	✓	79.5±2.8	75.1±0.8	66.0±0.3	60.3±1.1	83.1±0.7	57.2±0.9	24.9±1.8	61.8±1.8	49.5±1.9	94.6±0.0
IoU	D	SH	othr	pic	ref	show	sink	sofa	tab	toil	wall	wind
14A			33.9±2.1	8.3±6.5	53.7±2.4	54.0±1.4	64.4±2.0	72.8±1.2	64.3±2.5	86.2±1.5	74.6±2.2	37.3±2.7
14A	✓		34.9±1.3	17.4±0.7	57.6±1.4	49.9±3.9	66.9±1.0	74.2±1.0	66.4±1.0	86.2±1.5	77.2±0.4	43.0±1.5
14A	✓	✓	35.9±0.6	17.3±1.3	58.5±1.2	52.0±2.1	66.5±0.4	74.2±0.6	66.8±0.8	88.4±1.0	77.0±0.5	39.9±0.8
14A	✓	✓	36.5±0.7	15.3±1.5	59.2±1.5	53.4±3.0	67.3±0.7	74.7±0.7	67.1±0.8	89.0±0.5	77.1±0.4	41.8±1.0
18A			35.0±1.1	9.9±5.8	57.6±1.4	54.7±2.0	65.9±1.9	75.7±0.9	65.8±1.6	86.6±1.4	75.0±1.8	41.3±2.5
18A	✓		34.8±1.2	18.7±0.9	57.1±0.9	51.2±5.9	66.7±1.0	77.2±0.5	68.0±0.4	85.6±2.1	77.5±0.1	45.4±0.9
18A	✓	✓	36.2±0.4	16.5±1.7	58.9±1.9	56.5±1.4	66.8±1.0	74.0±1.8	68.3±0.5	89.0±1.0	77.6±0.0	43.6±1.0
18A	✓	✓	36.4±0.3	17.2±1.6	56.7±2.9	55.5±1.6	67.1±0.9	75.5±1.7	67.5±0.4	88.7±0.8	77.6±0.2	43.1±1.4
34C			37.1±2.4	9.2±5.5	59.5±2.4	55.6±2.0	66.6±1.6	76.8±1.1	65.8±2.6	86.3±1.4	75.0±2.5	42.3±1.7
34C	✓		36.9±2.5	16.4±0.3	60.5±3.5	54.0±3.4	66.6±1.5	75.5±1.8	69.8±0.5	86.3±1.5	78.1±0.4	45.4±2.0
34C	✓	✓	37.3±2.7	17.2±1.3	61.2±2.6	57.9±1.4	67.2±1.1	77.6±0.4	69.1±0.6	88.8±0.4	78.1±0.6	43.5±0.6
34C	✓	✓	39.7±1.3	17.0±0.5	57.8±5.6	58.2±1.7	67.8±0.7	77.2±1.6	69.2±0.1	89.3±0.7	78.2±0.2	43.6±0.5
Acc	D	SH	bath	bed	bksf	cab	chair	cntr	curt	desk	door	floor
14A			81.6±2.3	80.8±2.0	74.1±4.7	72.3±1.2	87.7±0.7	66.4±2.3	26.7±3.3	74.0±1.6	45.5±5.6	98.1±0.2
14A	✓		86.9±1.8	80.6±1.4	79.3±1.1	73.3±0.6	88.2±0.4	64.3±2.6	33.8±3.0	75.9±0.9	55.5±1.6	98.2±0.1
14A	✓	✓	82.7±1.2	81.4±1.0	79.6±1.0	73.2±0.8	88.9±0.2	64.3±1.8	26.9±2.1	75.3±2.7	55.7±1.8	98.3±0.0
14A	✓	✓	82.8±1.2	82.2±0.3	80.0±0.8	73.2±0.7	89.1±0.7	62.7±1.3	25.7±3.0	75.1±1.9	56.0±1.9	98.3±0.0
18A			85.8±0.4	81.7±1.0	76.0±1.9	74.3±0.3	88.5±0.4	72.1±2.7	25.6±1.4	76.3±2.6	51.3±7.4	98.2±0.0
18A	✓		83.9±2.8	81.3±1.3	78.9±0.7	74.4±0.3	89.1±0.1	70.6±3.2	37.6±8.7	77.1±2.7	62.6±0.7	98.2±0.1
18A	✓	✓	87.1±0.7	81.5±1.4	79.6±1.0	74.7±0.4	88.6±0.7	74.6±1.7	27.3±2.3	77.9±2.7	59.9±2.2	98.2±0.0
18A	✓	✓	87.2±2.3	82.4±0.6	80.4±2.4	74.7±2.0	89.3±1.0	68.8±5.1	29.6±3.3	77.4±1.4	59.9±1.5	98.2±0.1
34C			86.1±3.7	81.4±1.6	75.6±2.3	74.6±1.7	89.3±0.6	70.3±3.2	28.7±1.7	78.1±3.1	52.9±7.6	98.1±0.2
34C	✓		87.7±1.5	81.2±0.8	78.2±2.0	74.9±1.4	89.6±0.4	71.4±1.9	36.1±5.7	81.0±1.1	60.0±3.1	98.2±0.1
34C	✓	✓	85.3±3.9	83.3±0.7	79.0±1.3	74.8±1.6	90.4±0.3	67.5±3.2	29.1±0.4	80.7±1.3	61.7±2.8	98.2±0.1
34C	✓	✓	85.3±3.1	80.8±1.9	80.6±1.5	75.6±0.9	89.8±0.5	71.4±1.9	28.0±2.0	80.9±1.8	61.4±1.4	98.2±0.1
Acc	D	SH	othr	pic	ref	show	sink	sofa	tab	toil	wall	wind
14A			41.1±1.6	10.7±8.7	65.6±1.5	64.4±0.8	76.6±0.8	86.9±0.6	75.0±2.6	92.3±0.1	91.7±1.0	45.9±4.5
14A	✓		41.9±1.0	22.2±0.7	67.9±0.4	59.2±4.3	75.5±1.1	87.5±0.3	78.8±1.4	92.9±0.4	92.3±0.5	54.2±1.5
14A	✓	✓	43.0±0.6	21.5±1.9	68.3±0.9	61.6±3.8	76.1±0.4	87.8±0.5	77.9±0.8	92.6±0.4	92.5±0.4	50.0±2.3
14A	✓	✓	43.3±1.0	19.3±2.6	68.2±1.6	63.3±1.9	75.9±0.4	88.5±1.3	78.5±0.5	93.1±0.5	92.3±0.5	52.6±1.4
18A			42.5±1.7	11.7±7.2	69.1±2.4	64.5±2.4	75.7±1.5	87.7±2.0	77.3±2.0	92.0±0.3	90.9±1.2	52.5±3.8
18A	✓		42.4±1.4	24.3±1.8	67.3±3.0	61.1±6.7	76.4±0.8	88.9±1.1	78.8±2.4	92.5±0.5	91.2±0.8	58.1±3.0
18A	✓	✓	43.5±1.2	20.3±1.8	69.5±0.8	65.4±2.5	76.0±1.2	88.2±0.6	80.6±1.8	92.5±0.5	92.1±0.3	55.7±2.6
18A	✓	✓	43.5±0.9	21.0±2.0	69.6±2.2	63.9±3.4	76.4±0.7	87.9±0.2	79.0±0.4	92.4±0.4	92.3±0.2	54.8±2.9
34C			44.6±1.9	11.4±7.4	72.2±4.3	63.1±2.5	76.8±0.8	86.9±2.7	75.9±2.0	91.8±0.5	90.8±1.5	53.8±2.9
34C	✓		43.4±2.8	22.5±1.8	70.9±3.0	61.2±3.7	77.1±1.2	86.7±3.0	80.2±1.1	92.8±0.5	92.8±0.2	55.6±2.8
34C	✓	✓	45.0±1.9	21.1±1.9	71.3±0.8	66.6±4.6	76.9±0.6	89.7±0.8	79.3±0.7	92.6±0.3	92.6±0.3	54.7±1.4
34C	✓	✓	46.6±1.4	21.4±0.2	72.0±1.9	66.5±0.9	76.7±1.0	88.6±1.3	78.9±0.3	92.7±0.3	92.7±0.3	56.0±1.8

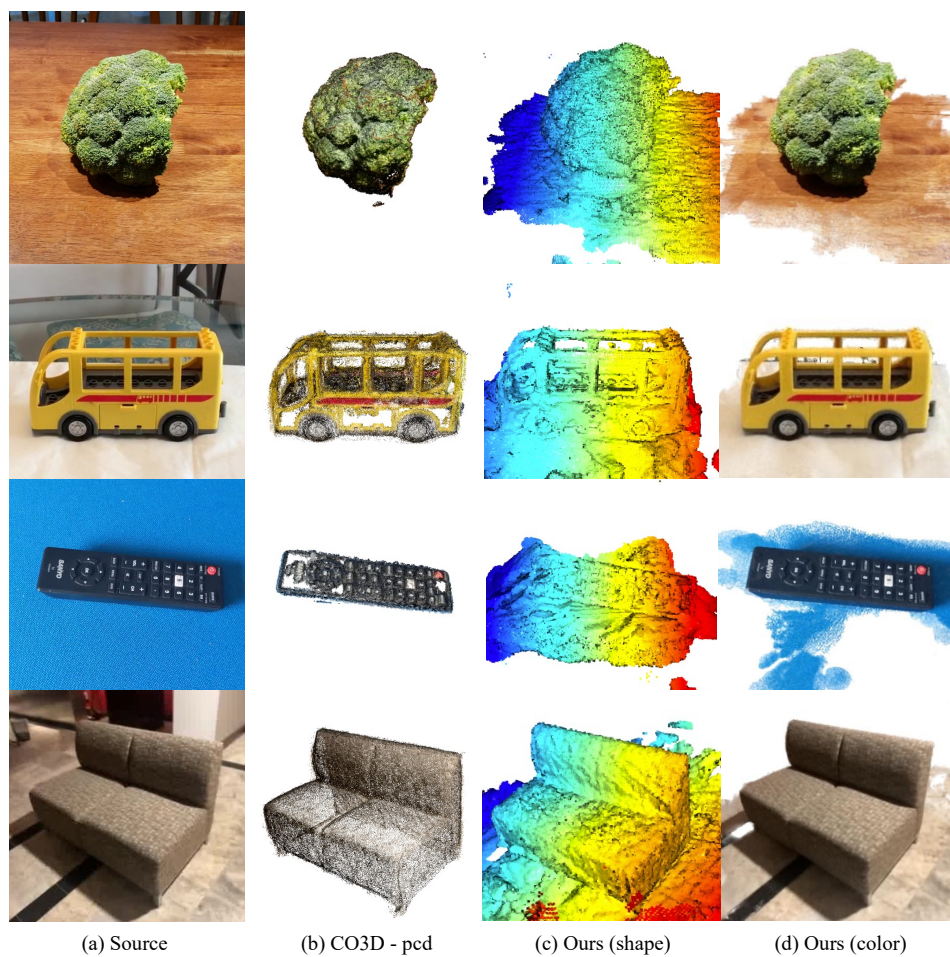


Figure a.5: Visualization of a few example data of original CO3D and PeRFception-CO3D.

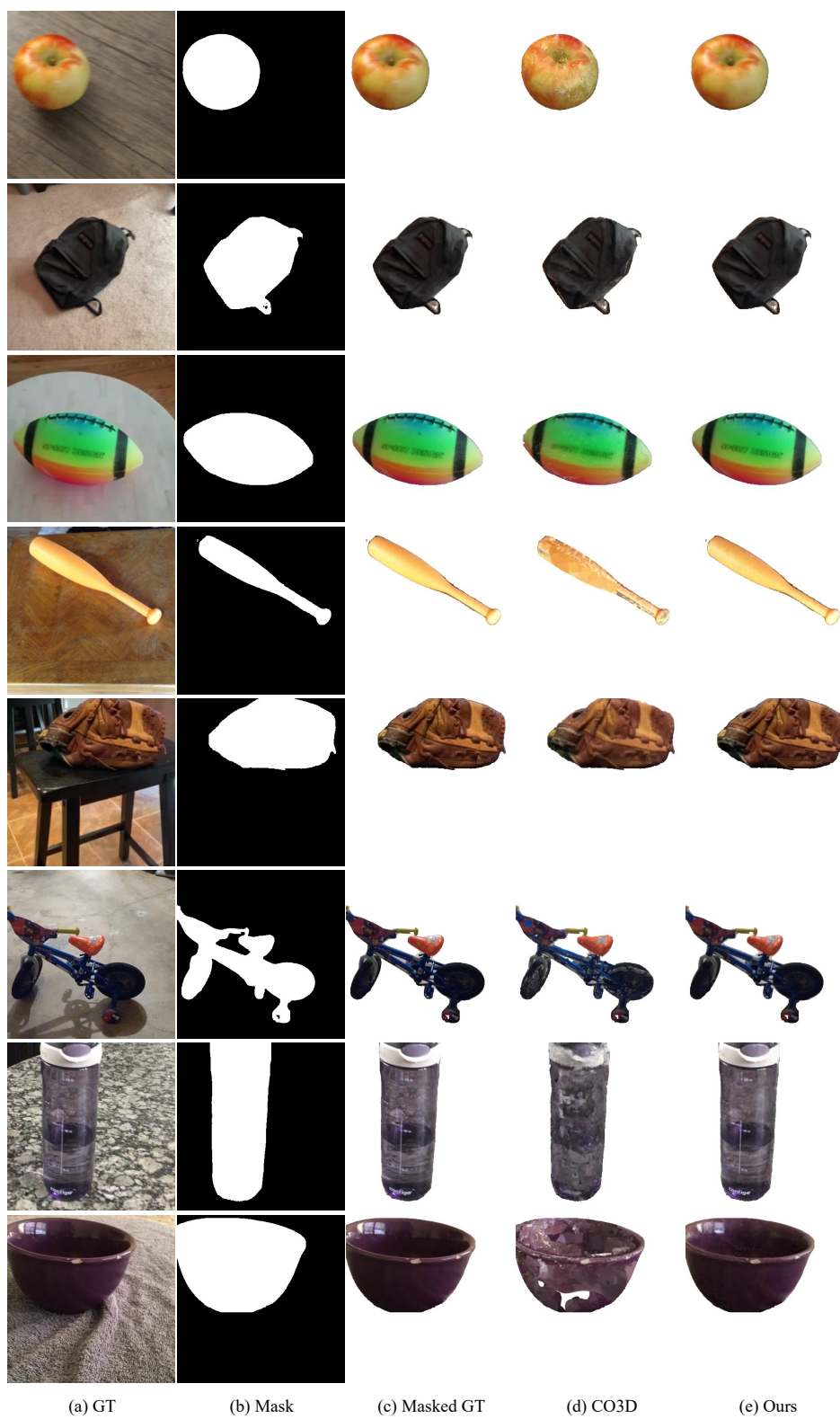


Figure a.6: Comparing visual reconstruction quality of original CO3D and PeRFception-CO3D.

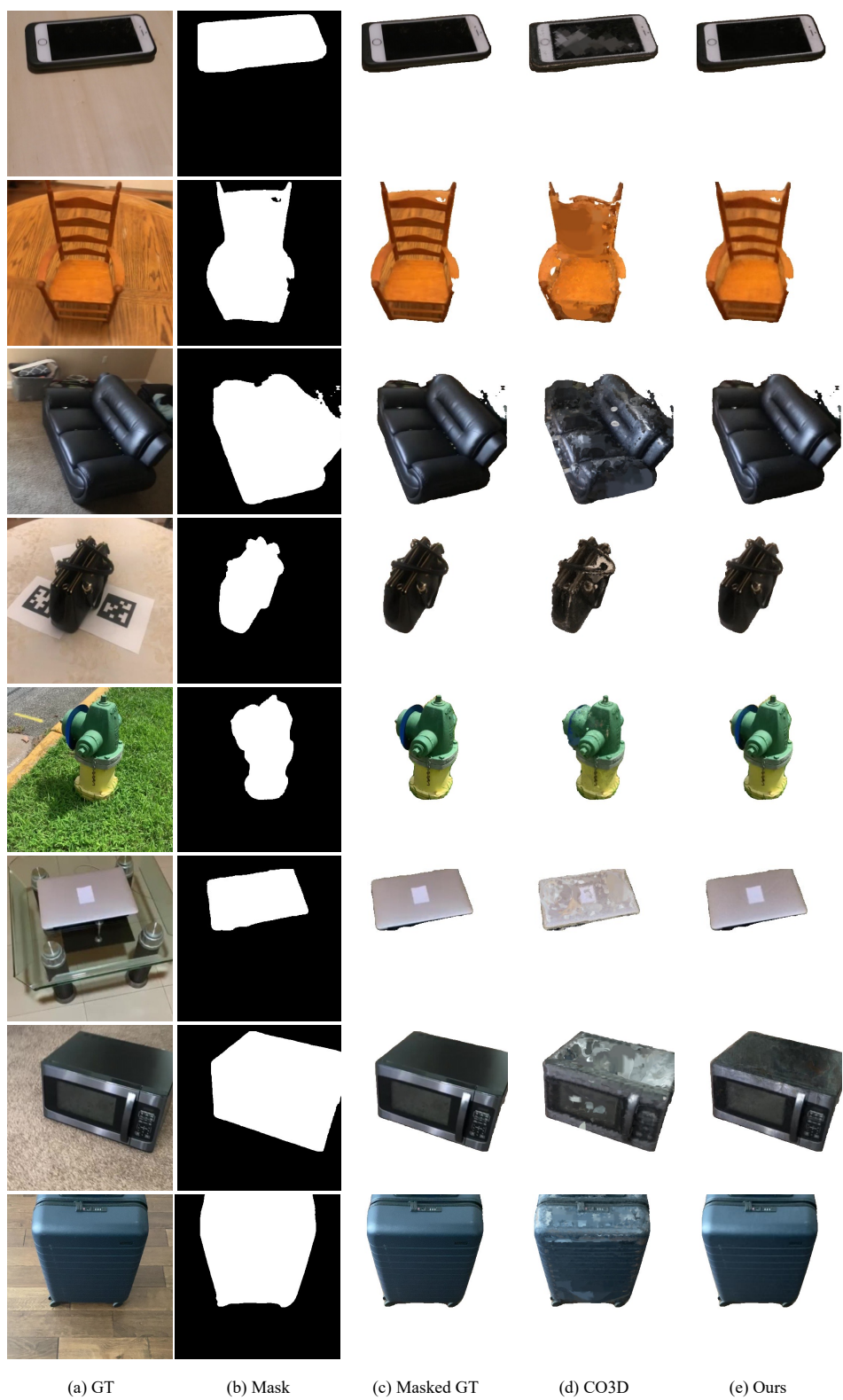


Figure a.7: Comparing visual reconstruction quality of original Co3D and PeRFception-Co3D.



Figure a.8: Comparing visual reconstruction quality of original Co3D and PeRFception-Co3D.

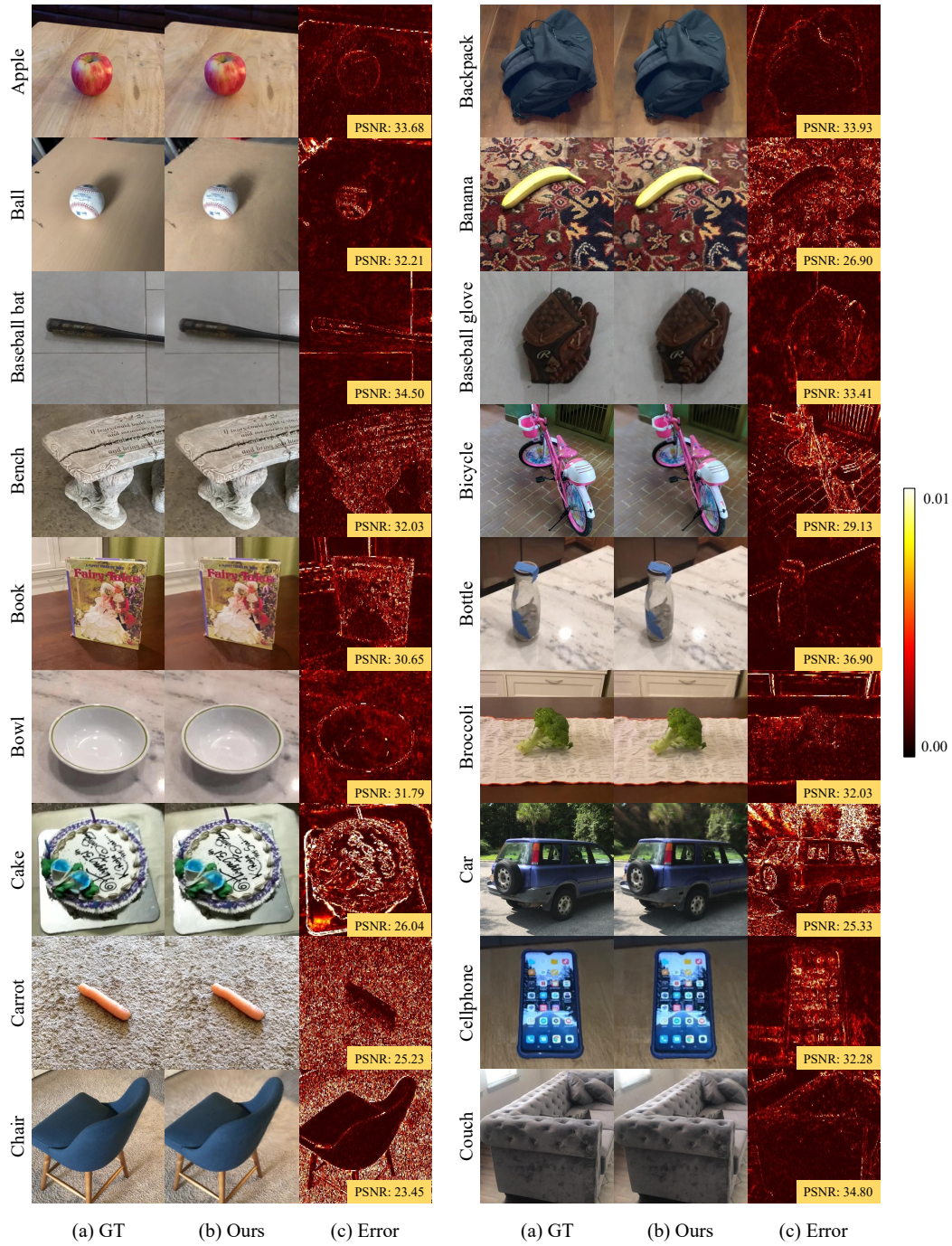


Figure a.9: Rendered class-wise novel views of PerFception-CO3D. The number in error maps denote the estimated PSNR.

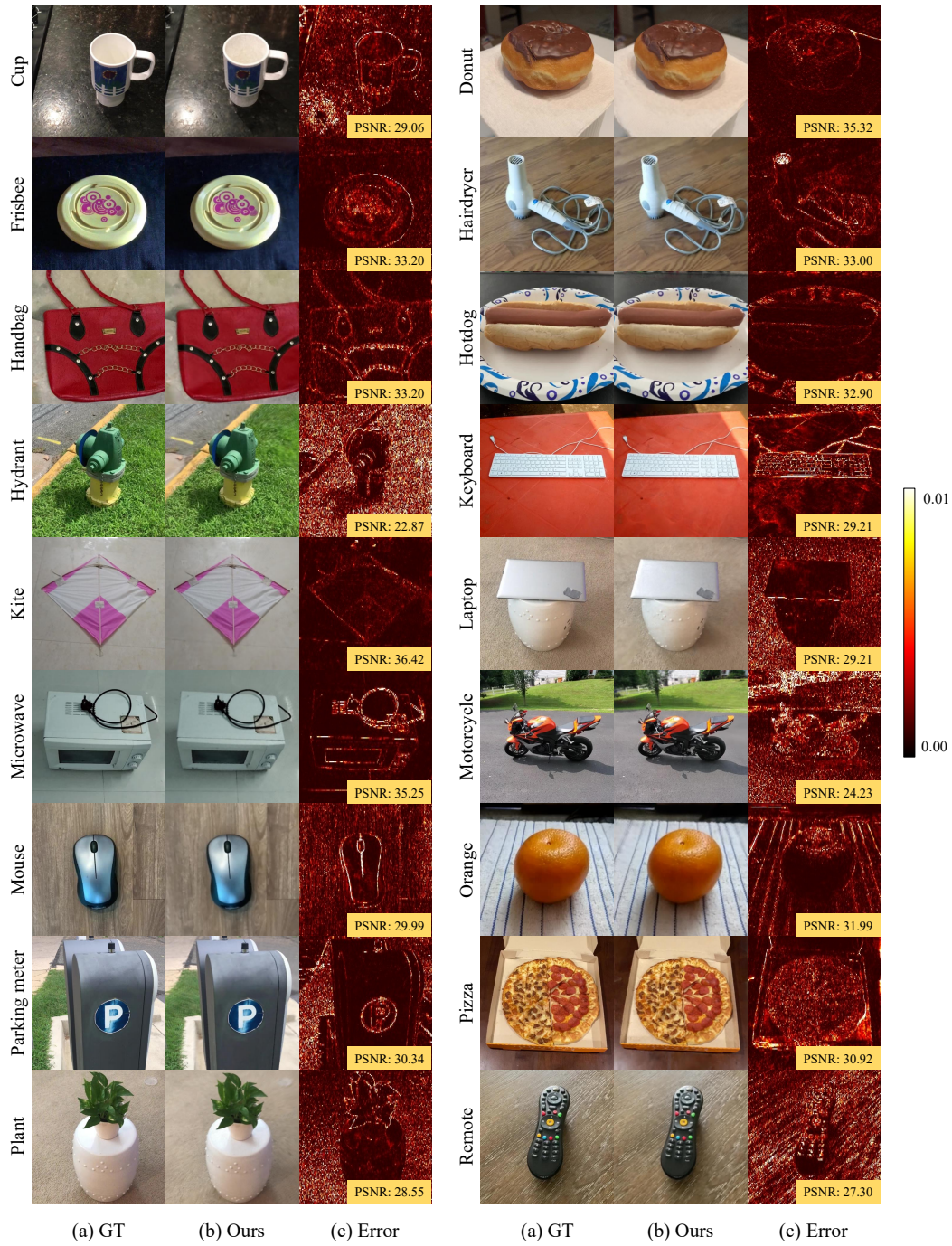


Figure a.10: Rendered class-wise novel views of PerFception-CO3D. The number in error maps denote the estimated PSNR.



Figure a.11: Rendered class-wise novel views of PerFception-CO3D. The number in error maps denote the estimated PSNR.

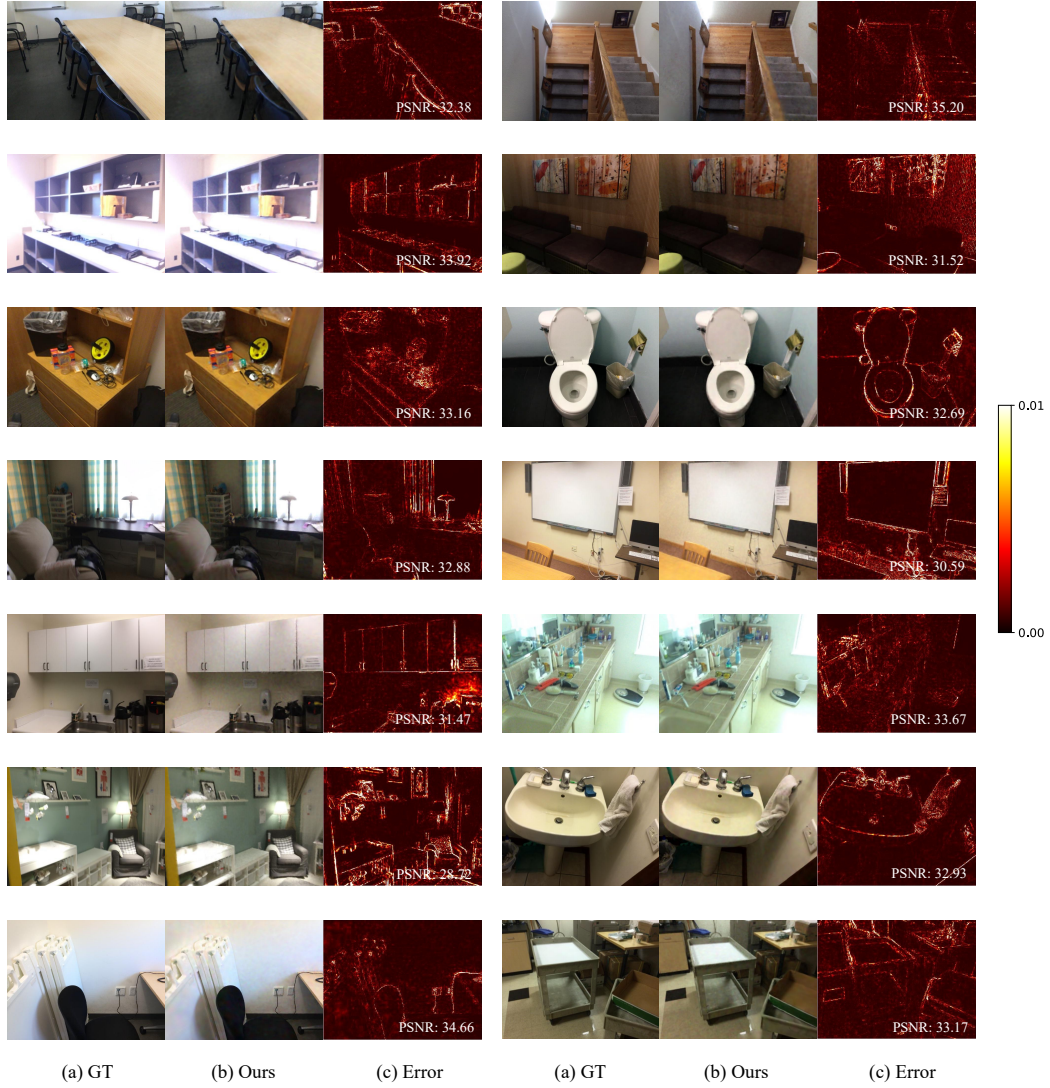


Figure a.12: Rendered class-wise novel views of PeRFception-ScanNet. The number in error maps denote the estimated PSNR



Figure a.13: Qualitative results of semantic segmentation on PeRFception-ScanNet dataset. (1st, 3rd columns) Ground truth point cloud with ground truth semantic labels, (2nd, 4th columns) Reconstructed sparse voxels with predicted semantic labels