

Instance Neural Radiance Field

Benran Hu^{1*} Junkai Huang^{1*} Yichen Liu^{1*} Yu-Wing Tai^{1,2} Chi-Keung Tang¹

¹The Hong Kong University of Science and Technology ²Kuaishou Technology

Abstract

This paper presents one of the first learning-based NeRF 3D instance segmentation pipelines, dubbed as **Instance Neural Radiance Field**, or *Instance-NeRF*. Taking a NeRF pretrained from multi-view RGB images as input, *Instance-NeRF* can learn 3D instance segmentation of a given scene, represented as an instance field component of the NeRF model. To this end, we adopt a 3D proposal-based mask prediction network on the sampled volumetric features from NeRF, which generates discrete 3D instance masks. The coarse 3D mask prediction is then projected to image space to match 2D segmentation masks from different views generated by existing panoptic segmentation models, which are used to supervise the training of the instance field. Notably, beyond generating consistent 2D segmentation maps from novel views, *Instance-NeRF* can query instance information at any 3D point, which greatly enhances NeRF object segmentation and manipulation. Our method is also one of the first to achieve such results without ground-truth instance information during inference. Experimented on synthetic and real-world NeRF datasets with complex indoor scenes, *Instance-NeRF* surpasses previous NeRF segmentation works and competitive 2D segmentation methods in segmentation performance on unseen views.

1. Introduction

Neural Radiance Field (NeRF) [34] has become the mainstream approach to novel view synthesis nowadays. Given multi-view images with camera poses only, NeRF encodes the underlying scene in a multi-layer perceptron (MLP) by radiance propagation and generates very impressive results. Thus, subsequent to NeRF’s debut in [34] lot of works have made great progress in improving the quality [11, 46], efficiency [35, 54, 4] and generality [56, 50].

This excellent approach to associate 2D with 3D through radiance field leads us to rethink the *3D instance segmentation problem*. Unlike the 2D counterpart operating on plenty of training images, 3D instance segmentation is limited by

*Equal contribution. The order of authorship was determined alphabetically.

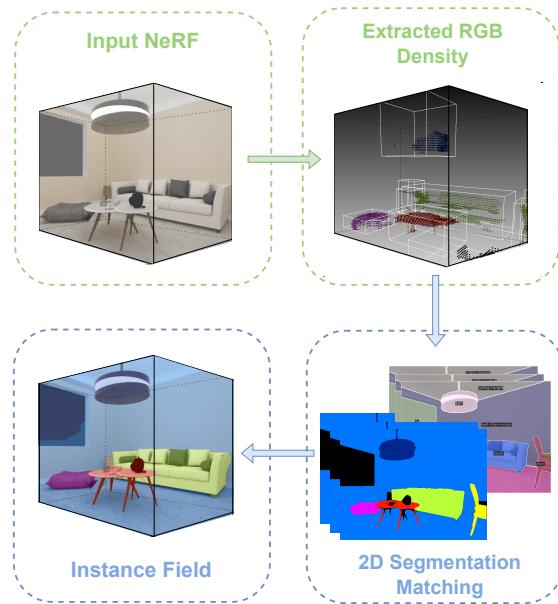


Figure 1: Pipeline of Instance-NeRF, demonstrated on the 3D-FRONT dataset. Instance-NeRF takes a pre-trained NeRF as input to detect objects within the scene and utilizes existing 2D panoptic segmentation to generate 2D segmentation maps, which are then matched and used to supervise the instance field training.

both the quantity and quality of the available data. Previously, 3D segmentation still relies on RGB-D images [23] or point clouds [38, 39, 47, 51, 43] captured by depth sensors or custom equipment as input, which are inconvenient to obtain and contain a variety of noises. To mitigate the dependence on explicit 3D geometry, there have been several investigations on embedding NeRF for addressing fundamental problems such as 3D semantic segmentation [57] and scene manipulation [28]. Some unsupervised methods [32, 55, 41] also involve 3D scene decomposition and instance segmentation, but it is hard to apply them on complex and large scenes akin to real world cases.

In this paper, inspired by NeRF-RPN [24] and Mask-RCNN [19], we propose *Instance-NeRF* for 3D instance segmentation in NeRF. More specifically, we incorporate NeRF-RPN with a mask head to predict 3D coarse seg-

mentation. After projecting the attained 3D masks back to 2D, Instance-NeRF leverages Mask2Former [7] and CascadePSP [8] to match the same instance in 2D segmentation from different views and refine the resultant masks. The refined 2D segmentation of multi-view images will be used to train an instance field which encode 3D instance information in a continuous manner as a neural field.

Our major contributions are:

- One of the first attempts to perform 3D instance segmentation in NeRF without using ground-truth segmentation information during inference.
- Propose the architecture and training approach of an *Neural Instance Field*, which can produce multi-view consistent 2D segmentation as well as continuous 3D segmentation using NeRF representation.
- Perform experiments and ablation studies on a synthetic indoor NeRF dataset to demonstrate the effectiveness of our method, which surpasses competitive 2D segmentation methods and previous works in NeRF segmentation.

2. Related Works

2.1. Neural Radiance Fields

Neural Radiance Field (NeRF) [34] is now the state-of-the-art for reconstructing new views of a given scene, by modeling the underlying 3D geometry and appearance using a continuous and implicit radiance field parameterized by a multilayer perceptron (MLP). Recent works include: Instant Neural Graphics Primitive [35] uses hash encoding to speed up training, PlenOctrees [54] uses an octree-based radiance field and spherical basis functions to improve rendering speed and appearance decoding, and TensoRF [4] encodes positional information by projecting 3D points onto three 2D planes. NeRF not only provides structural details of a 3D scene but is also conducive to 3D training, where only RGB images with camera parameters are required, thus making this alternative also suitable for 3D unsupervised object segmentation [32], where first encouraging results are demonstrated on real scenes using radiance propagation but with no geometry consideration.

2.2. R-CNN

In 2D, Region-based CNN (R-CNN) [17] focuses on a small number of regions [22], and uses convolutional networks to independently analyze each RoI. R-CNN was extended by incorporating RoIPool [21, 15] for attending to regions on feature maps, making it faster and more accurate. Faster R-CNN [40] incorporates Region Proposal Network (RPN) to learn the attention mechanism. Mask R-CNN [19] adds a branch to predict an object mask in parallel with the existing branch for bounding box recognition.

In 3D, NeRF-RPN [24] bridges RPN and NeRF and demonstrates great potential in direct 3D learning from NeRFs. This paper contributes to 3D object instance segmentation from NeRFs by capitalizing on the 3D boxes given by NeRF-RPN. In other words, analogous to Mask-RCNN, our work extends NeRF-RPN to enable 3D instance segmentation directly from a given NeRF.

2.3. Instance Segmentation

For 2D instance segmentation, two-stage methods [29, 19, 3, 5, 9] first detect candidate bounding boxes and then predict instance mask within each of them. On the other hand, [6, 27, 52, 2, 48, 49] release the network from proposal generation and achieve comparable results. Despite the great success of 2D instance segmentation, straightforward extension to 3D does not work in general. Specifically, applying 2D instance segmentation on each image capturing a single scene does not guarantee consistency across multi-view images of the underlying 3D scene.

Current methods on 3D instance segmentation usually perform on RGB-D images or point clouds. Image-based approaches such as [23] use 2D convolution to extract RGB-D features and then project them back to infer 3D segmentation. For point-cloud-based methods, While some methods [51, 43, 18, 45, 30] voxelize the inputs and adopt 3D convolution, others [38, 39, 47, 53] reserve the irregularity of point clouds and utilize permutation-invariant function to extract features. In addition, some approaches [10, 25] take RGB images as extra inputs and fuse features from two types of data. However, all of them require some explicit 3D geometry e.g. obtained by LiDARs or other devices. 3D instance segmentation directly from multi-view images has not been explored in the context of NeRFs.

Segmentation and decomposition on NeRF were investigated. Semantic-NeRF [57] extends the NeRF by adding a semantic branch to produce view independent semantic labels. Semantic-NeRF predicts accurate semantic labels even when the ground truth labels are sparse and noisy, but this work mainly targets semantic map denoising and super-resolution instead of inference tasks. DM-NeRF [1], on the other hand, uses multi-view ground-truth instance maps to optimize an object field for NeRF object decomposition and manipulation. Panoptic NeRF [14] can render accurate semantic maps with sparse multi-view images and the corresponding noisy predicted masks given by a pre-trained model. However, it requires 3D bounding primitives to guide consistency and predict segmentation on instance level. Panoptic Neural Fields (PNF) [26] assumes objects are dynamic and tracks 3D object from images in panoptic images, which is not suitable for general scenes represented by NeRF. NeSF [44] utilizes a 3D UNet to generate semantic feature grid from density grid sampled from NeRF, and performs volume rendering over the feature grid

to produce 2D semantic maps. NeSF does not need ground truth semantics during inference, but it cannot perform instance segmentation.

There have been some investigation on unsupervised instance segmentation in NeRF. GIRAFFE [37] uses generative neural feature fields to decode shape and appearance, which decomposes foreground objects from background without explicit 3D segmentation learning. [41, 55] apply slot attention [33] on unsupervised object discovery, while others [12, 32] incorporate the advantages of NeRF to obtain multi-view 2D masks. But these methods only achieve qualified success in simple scenes with a limited number of objects. There is still a large room for quality improvement and generalization. Our work is one of the first learning-based approaches to generate 3D instance segmentation in NeRF, which can be applied to produce satisfactory 3D object masks on complex scenes during test time.

3. Method

Given a pre-trained NeRF, our method, Instance-NeRF, aims to detect all the objects within the underlying 3D scene and produces a bounding box, a continuous 3D mask, and a class label of each detected 3D object.

Instance-NeRF extends the given pre-trained NeRF model with an additional instance field, which can produce a view-independent instance label at any 3D position in the NeRF scene. To train the instance field component of Instance-NeRF, we propose a *NeRF-RCNN* for 3D sparse mask prediction and a *2D mask matching and refining* stage to produce multi-view consistent 2D masks based on the 3D masks for instance field supervision.

First, the *NeRF-RCNN*, which is extented from NeRF-RPN [24], takes the extracted radiance and density field of the pre-trained NeRF and output 3D bounding boxes, class labels, and discrete 3D masks for each detected object in the NeRF. Then, given a set of multi-view but possibly *inconsistent* 2D panoptic segmentation maps of the scene produced by existing methods, we project the 3D masks from the same camera poses to match the same instance across different views. The multi-view *consistent* 2D segmentation can then be used to train the instance field component and produce a continuous 3D segmentation using instance field representation.

3.1. Instance-NeRF

3.1.1 Model Architecture

Instance-NeRF optimizes a neural radiance field that maps a 3D position $\mathbf{x} = (x, y, z)$ and a viewing direction $\mathbf{d} = (\phi, \theta)$ to a volume density σ , RGB radiance $\mathbf{c} = (r, g, b)$, and an instance label distribution over L labels including the background label. As shown in Figure 3, Instance-NeRF is parameterized by three components - density net

Θ_σ , color branch Θ_c , and instance branch Θ_i . Θ_σ, Θ_c constitute the original NeRF, and are assumed to be pre-trained with posed RGB images. The pre-trained NeRF can be formulated as

$$\mathcal{F}_{\Theta_\sigma, \Theta_c}(\mathbf{x}, \mathbf{d}) = (\sigma, \mathbf{c}), \quad (1)$$

and the Instance Field can be formulated as:

$$\mathcal{G}_{\Theta_\sigma, \Theta_i}(\mathbf{x}) = \mathbf{i}, \quad (2)$$

where \mathbf{i} is a L -dimensional vector with its first dimension representing the background.

We follow the hierarchical stratified sampling method in [34] to render the RGB color, depth, and instance label for a single pixel. Specifically, for each pixel, we formulate a ray $\mathbf{r}(t) = \mathbf{o} + t\mathbf{d}$ emitted from the center of the projection camera to that pixel, and select K quadrature points $\{t_k\}_{k=1}^K$ on the ray \mathbf{r} within the traceable volume of Instance-NeRF. The numerical quadrature to accumulate the expected color $\hat{\mathbf{C}}(\mathbf{r})$, expected depth $\hat{\mathbf{D}}(\mathbf{r})$ and expected instance logits $\hat{\mathbf{I}}(\mathbf{r})$ for each pixel are respectively given by:

$$\hat{\mathbf{C}}(\mathbf{r}) = \sum_{k=1}^K \hat{T}(t_k) \alpha(\sigma(t_k) \delta_k) \mathbf{c}(t_k), \quad (3)$$

$$\hat{\mathbf{D}}(\mathbf{r}) = \sum_{k=1}^K \hat{T}(t_k) \alpha(\sigma(t_k) \delta_k) t_k, \quad (4)$$

$$\hat{\mathbf{I}}(\mathbf{r}) = \sum_{k=1}^K \hat{T}(t_k) \alpha(\sigma(t_k) \delta_k) \mathbf{i}(t_k), \quad (5)$$

where

$$\begin{aligned} \hat{T}(t_k) &= \exp\left(-\sum_{a=1}^{k-1} \sigma(t_a) \delta_a\right), \\ \alpha(x) &= 1 - \exp(-x), \\ \delta_k &= t_{k+1} - t_k. \end{aligned} \quad (6)$$

We discuss the NeRF pre-training and Instance Field training in 3.1.2 and 3.1.3, respectively.

3.1.2 Pre-training Density and Color

The density and radiance components of Instance-NeRF can be implemented and trained similarly as a regular NeRF, and the implementation is largely orthogonal to the Instance Field component. Posed RGB images are used to train the density net Θ_σ and color branch Θ_c with the appearance loss \mathcal{L}_p :

$$\mathcal{L}_p = \frac{1}{\mathcal{R}} \sum_{\mathbf{r} \in \mathcal{R}} \left\| \hat{\mathbf{C}}(\mathbf{r}) - \mathbf{C}(\mathbf{r}) \right\|_2^2, \quad (7)$$

where \mathcal{R} are the sampled rays within a training batch, $\mathbf{C}(\mathbf{r})$ and $\hat{\mathbf{C}}(\mathbf{r})$ are respectively the ground truth and predicted RGB color for ray \mathbf{r} . In this work, we assume a well-trained Θ_c and Θ_σ are given and focuses mainly on the Instance Field training.

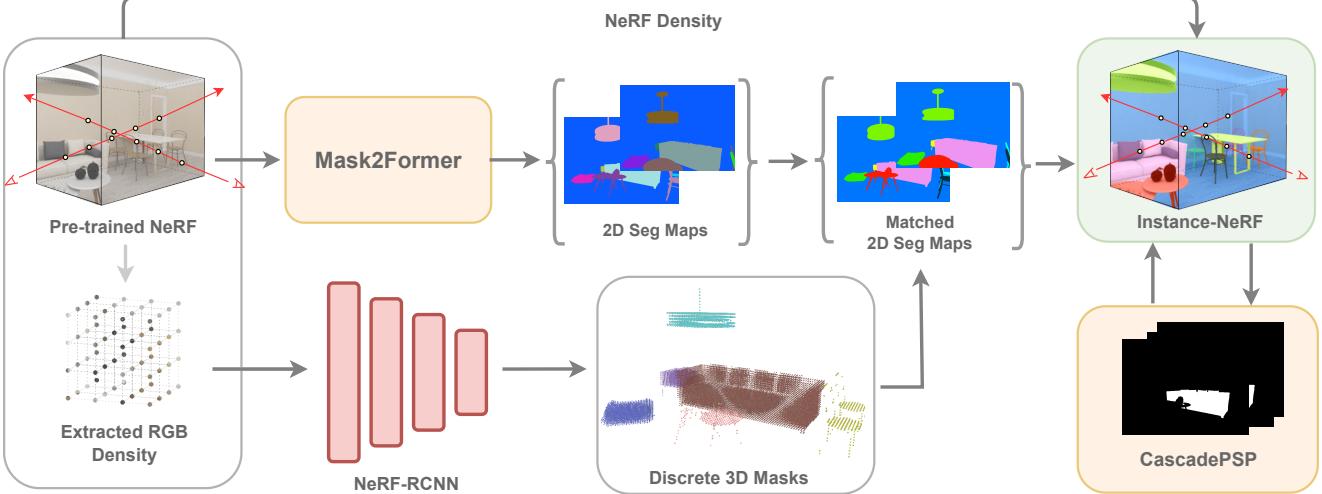


Figure 2: **Instance Field Training.** When training the instance field, NeRF-RCNN takes the extracted radiance and density field of the pre-trained NeRF and outputs discrete 3D masks for each detected object in the NeRF. Mask2Former generates 2D panoptic segmentation maps of images rendered from NeRF, which are *consistent* in terms of instance label across views. After projecting the 3D masks from the same camera poses to match the same instance across different views resulting in the multi-view *consistent* 2D segmentation maps, they can then be used to train the instance field component and produce a continuous 3D segmentation using instance field representation. In addition, we use CascadePSP to refine the preliminary instance segmentation results of Instance-NeRF, and use the refined instance masks to refine the Instance Field in turn.

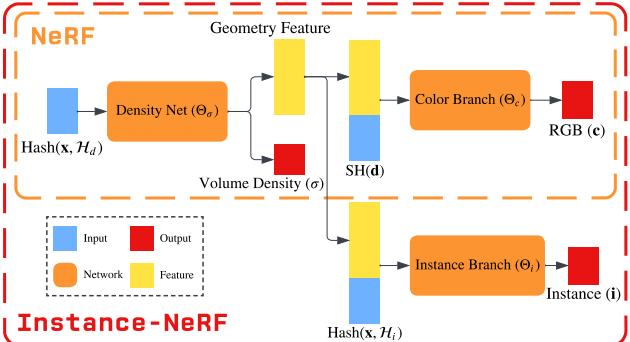


Figure 3: **Instance-NeRF Architecture.** Instance-NeRF consists of density net Θ_σ , color branch Θ_c , and instance branch Θ_i . Given the input spatial position $\mathbf{x} = (x, y, z)$ and a viewing direction $\mathbf{d} = (\phi, \theta)$, compared to the traditional NeRF, Instance-NeRF can predict a view-independent instance distribution \mathbf{i} with its additional instance branch. In this paper, we adopt the Multi-resolution Hash Encoding (Hash) in [35] for position encoding and Spherical Harmonics (SH) in [54] for viewing direction encoding. Note that the hash grid \mathcal{H}_i used for mask net is different from the one \mathcal{H}_σ used for density net.

3.1.3 Instance Field Training

The instance field is trained with multi-view *consistent* 2D instance segmentation maps, during which Θ_σ and Θ_c will be fixed and only Θ_i will be optimized. Given a set of multi-view consistent 2D instance segmentation maps, we train the instance branch with multi-class cross-entropy loss \mathcal{L}_i :

$$\mathcal{L}_i = -\frac{1}{\mathcal{R}L} \sum_{\mathbf{r} \in \mathcal{R}} \sum_{l=1}^L p^l(\mathbf{r}) \log \hat{p}^l(\mathbf{r}), \quad (8)$$

where $p^l(\mathbf{r})$ and $\hat{p}^l(\mathbf{r})$ are respectively the ground truth multi-class instance probability for class l and volume predictions for ray \mathbf{r} given by applying softmax to Eq. 5.

As the input 2D segmentation maps are generated by existing 2D segmentation methods, it cannot produce 3D geometry-consistent masks for different views. To compensate this issue and to encourage the smoothness of the instance label prediction, we take advantage from the prior that instance maps are generally smooth over local regions, and add an instance regularization loss \mathcal{L}_r adopted from [36]:

$$\begin{aligned} \mathcal{L}_r = \frac{1}{\mathcal{R}L} \sum_{\mathbf{r}_{i,j} \in \mathcal{R}} \sum_{l=1}^L & \left(\hat{\mathbf{I}}(\mathbf{r}_{i,j}) - \hat{\mathbf{I}}(\mathbf{r}_{i,j+1}) \right)^2 w_{i,j} + \\ & \left(\hat{\mathbf{I}}(\mathbf{r}_{i,j}) - \hat{\mathbf{I}}(\mathbf{r}_{i+1,j}) \right)^2 h_{i,j}, \end{aligned} \quad (9)$$

$$w_{i,j} = \frac{\delta_w(i, j)}{\sum_{\mathbf{r}_{i',j'} \in \mathcal{R}} \delta_w(i', j')}, \quad (10)$$

$$\delta_w(i, j) = \exp(-(\hat{\mathbf{D}}(\mathbf{r}_{i,j}) - \hat{\mathbf{D}}(\mathbf{r}_{i,j+1}))^2), \quad (11)$$

where $w_{i,j}$ is the weight determined by the similarity of NeRF depth $\hat{\mathbf{D}}$ between two sample rays, and likewise for

$h_{i,j}$. $\mathbf{r}_{i,j}$ is the ray passing through pixel (i, j) of the image. Recall that depth information is obtained from the pre-trained density branch. Similar depths of adjacent pixels imply that the pixels might be on the same surface in 3D space, and likely belong to the same object. Thus, this regularization loss encourages the model to predict the same instance label for them and generate smoother segmentation for each object.

The total loss \mathcal{L} for training the instance field is:

$$\mathcal{L} = \mathcal{L}_p + \lambda \mathcal{L}_i, \quad (12)$$

where λ is a hyperparameter.

3.2. NeRF-RCNN

NeRF-RCNN takes as input a pre-trained NeRF and outputs 3D Axis-Aligned Bounding Boxes (AABB), class labels, and discrete 3D masks of the detected objects. This network extends NeRF-RPN [24] by appending two additional detection heads: a box prediction head (MLP) and a mask prediction head (3D CNN), much similar in spirit as [20]. As shown in Figure 2, we first uniformly sample the appearance and density on a grid that covers the full traceable volume of the pre-trained Instance-NeRF model, following the sampling method in [24]. The NeRF-RPN takes the sampled grid as input, and outputs proposals and the feature pyramid. We use 3D RoIAlign to obtain the feature within the ROI of each output proposal. The box prediction head takes ROI features as input and outputs class logits and AABB regression offset, while the mask prediction head takes ROI features as input and outputs 3D masks for each ROI.

The *NeRF-RCNN* can be trained in an end-to-end manner with sampled appearance and density grids as input, along with ground truth bounding boxes and 3D instance segmentation masks, by reducing the multi-class cross-entropy loss for class prediction, the smooth L1 loss for bounding box regression, and the binary cross-entropy loss for 3D mask prediction. In practice, we adopt the pre-trained NeRF-RPN, and optimize the classification head and mask head side by side. In our experiments, we train the *NeRF-RCNN* with a large NeRF segmentation dataset built on 3D-FRONT.

3.3. 2D Mask Matching and Refinement

It is infeasible to adopt a 2D instance segmentation model to perform segmentation on multi-view images and directly use the results to supervise instance field training. The main reason lies in the fact that 2D instance segmentation models generally do not guarantee any correspondence between the appearance of the same instance in different views, including the assigned instance IDs, predicted class labels, and the masks. The same object is likely to be as-

signed with different instance IDs in different images, making it difficult for direct use in instance field training.

To address this consistency issue, we utilize the 3D coarse mask produced by NeRF-RCNN to align masks in different images that correspond to the same object. We further propose an iterative refinement approach by feeding the intermediate results of Instance-NeRF to an existing mask refinement model, and use the refined masks to train the final instance field.

3.3.1 2D Mask Matching

To obtain the initial 2D segmentation, we use a Mask2Former [7] panoptic segmentation model pretrained on the COCO [31] dataset. We create a class mapping from COCO to the dataset we use, based on which we filter out masks that belong to background. To deal with prediction inconsistency across images, we match each Mask2Former predicted 2D mask with a 3D instance detected by NeRF-RCNN. To achieve this, we first project the 3D masks from NeRF-RCNN to the corresponding image spaces, and compute the Intersection over Union (IoU) between each pair of predicted 2D mask and projected 2D mask. Each predicted mask is then assigned to the instance with highest mask IoU. Those predicted masks that do not have an IoU greater than 0.05 with any projected masks are treated as unlabeled area, which do not participate in instance nerf training. After this process, we have consistently annotated 2D instance masks from multi-view images, which will be used to optimize the instance field.

3.3.2 2D Mask Refinement

Although 2D-to-3D matching emancipates the model from inconsistency, uniformly sampled images, which frequently contain objects under partial occlusion can easily mislead the prediction model pre-trained on general datasets, resulting in missing or erroneous mask prediction in different views. An instance field trained directly on the 2D masks therefore suffers from conflicting 2D masks of different views, and usually produces fragmented segmentation results. To alleviate this issue, we adopt a pre-trained 2D mask refinement model, CascadePSP [8], to refine the preliminary instance segmentation results of Instance-NeRF. Instead of directly feeding Mask2Former output into the refinement network, we use Instance-NeRF rendering results, because Instance-NeRF can enhance single-view results by fusing multi-view information, filling part of the area of incomplete or missing mask prediction in certain views. We split the instances in each preliminary Instance-NeRF mask and refine them independently, the results of which are superimposed to supervise the final instance field training.

Methods	mIoU \uparrow	PQ \uparrow
Mask2Former	31.6	22.5
Semantic-NeRF	35.0	-
DM-NeRF	11.3	5.4
Ours	43.0	32.4

Table 1: Comparison results of mean Intersection over Union (mIoU) and Panoptic Quality (PQ) on 3D-FRONT dataset.

\mathcal{L}_i	Mask Refinement	mIoU	PQ
		39.8	27.6
\checkmark		40.8	29.6
\checkmark	\checkmark	43.0	32.4

Table 2: Ablation over instance regularization loss and 2D mask refinement on 3D-FRONT dataset.

4. Experiments

4.1. Training and Testing

Training. In our model, only *NeRF-RCNN* requires a large number of scenes for training for generalization, while the instance field training are scene-specific and the models used in 2D mask matching and refinement are pre-trained. Following NeRF-RPN [24] strategy and configuration, NeRF-RCNN is trained on 3D-FRONT [13] NeRF segmentation dataset. We follow the dataset creation approach in [24] and extend the dataset to include 1016 scenes, 814 of which are used for training *NeRF-RCNN* and 101 scenes are used for validation and testing, respectively. We choose VGG19, which performs best on object detection, as our backbone. The weights of the three losses mentioned in Section 3.2 are all equal to 1.0. We apply random flipping and rotation by $\pi/2$ with probability 0.5 to the input as augmentation.

For the 2D panoptic segmentation model, we use the Mask2Former model pre-trained on COCO dataset, with Swin-L as the backbone. For the mask refinement, we adopt the publicly available CascadePSP model. No fine-tuning is performed for both methods on our dataset.

Testing We test our model on 8 scenes from 3D-FRONT in total. We use a non-maximum suppression (NMS) threshold of 0.3 for NeRF-RPN and 0.15 for NeRF-RCNN. We use a score threshold of 0.5 to filter the NeRF-RCNN detection. The threshold for Mask2Former segmentation is also 0.5. For the second training stage of Instance-NeRF, we implement our model based on [42]. The hyperparameter λ in Eq. (12) is 1.0. For each scene, we train the radiance and density field on a single NVIDIA 1080ti for 30k iterations. The instance field is then trained for 25k steps to obtain the intermediate segmentation results for refinement. We train the instance field for another 20k steps using the refined 2D masks. Figure 5 illustrates the high-quality 3D instance segmentation results.

4.2. Metrics

We evaluate the semantic segmentation performance of Instance-NeRF on novel views using mean Intersection over Union (mIoU) and instance segmentation performance using Panoptic Quality (PQ). The semantic segmentation is acquired by mapping the classes of the detected instances. We use PQ for evaluating instance segmentation as not all compared methods produce instance confidence scores, hence we cannot use common metrics such as mean Average Precision. Classes corresponding to background are filtered for PQ calculation. Note that we compute PQ in a conventional approach, where the consistency of instances across multiple frames is not taken into account. This allows us to compare with Mask2Former even it does not produce multi-view consistent segmentation. However, PQ also fails to reflect our method’s performance in producing consistent segmentation for different views.

4.3. Comparison

To our best knowledge, we are one of the first to propose 3D instance segmentation for NeRF without utilizing ground-truth information at test time, which can be generalized to different kinds of scenes and requires pre-trained NeRF only during inference. Thus, it is hard to find comparable methods with the same configuration. Semantic-NeRF [58] also uses a 2D pre-trained model but the method cannot be straight-forwardly modified for instance segmentation. DM-NeRF [1] can achieve consistency in 3D with the supervision of 2D ground truth instance segmentation. In our experiments, we compare our method with Semantic-NeRF for semantic segmentation, and with DM-NeRF for instance segmentation. For fairness, we implement them on torch-npg [42] and use Mask2former [7] to predict 2D masks for input images, supervising the NeRF training of these two methods. As Mask2Former do not guarantee multi-view instance consistency, we use the majority class for each instance to create the mapping to class semantics. Figure 4 illustrates the comparison of the qualitative results and Table 1 tabulates the quantitative comparison.

4.4. Ablation Study

We perform ablations on the regularization loss and the 2D mask refinement of our method. Table 2 shows the quantitative results. For qualitative results, please refer to our Supplementary Material.

4.4.1 Instance Regularization Loss

Table 2 shows that the additional regularization loss \mathcal{L}_i consistently improves the mIoU and PQ of our methods. The loss helps smooth fragmented segmentation regions induced by inconsistent 2D segmentation maps and close

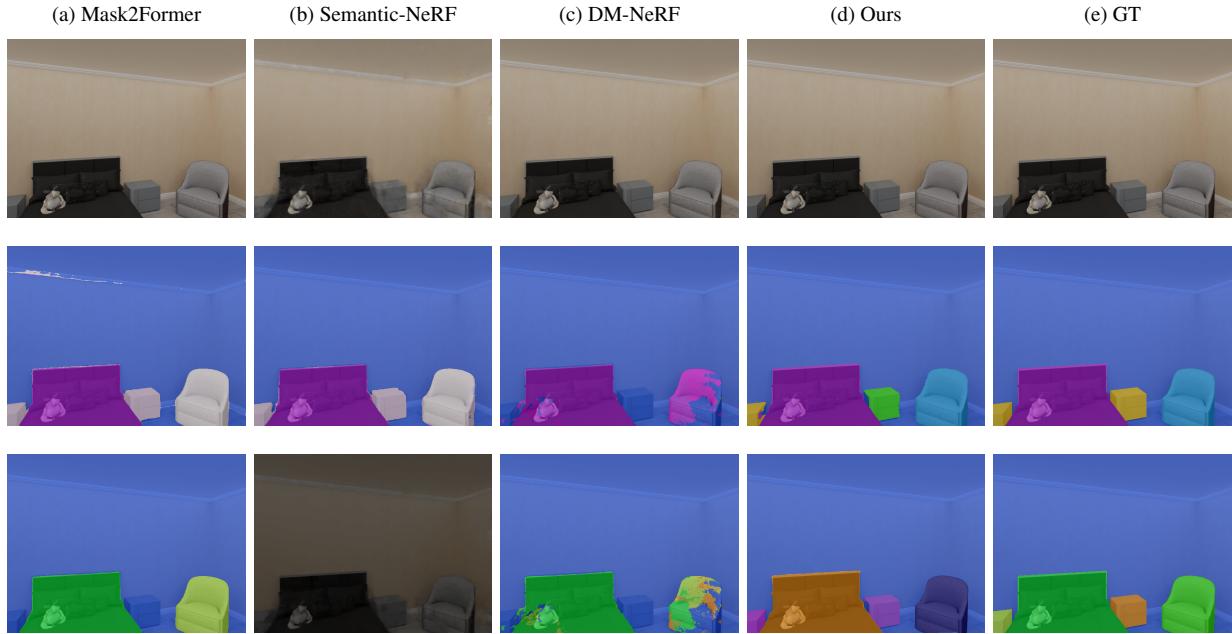


Figure 4: **Comparison.** This figure illustrates the comparison with other state-of-art methods. Rows from top to bottom are (a) ground truth RGB images or the rendered RGB image by the models (b) semantic segmentation and (c) instance segmentation.

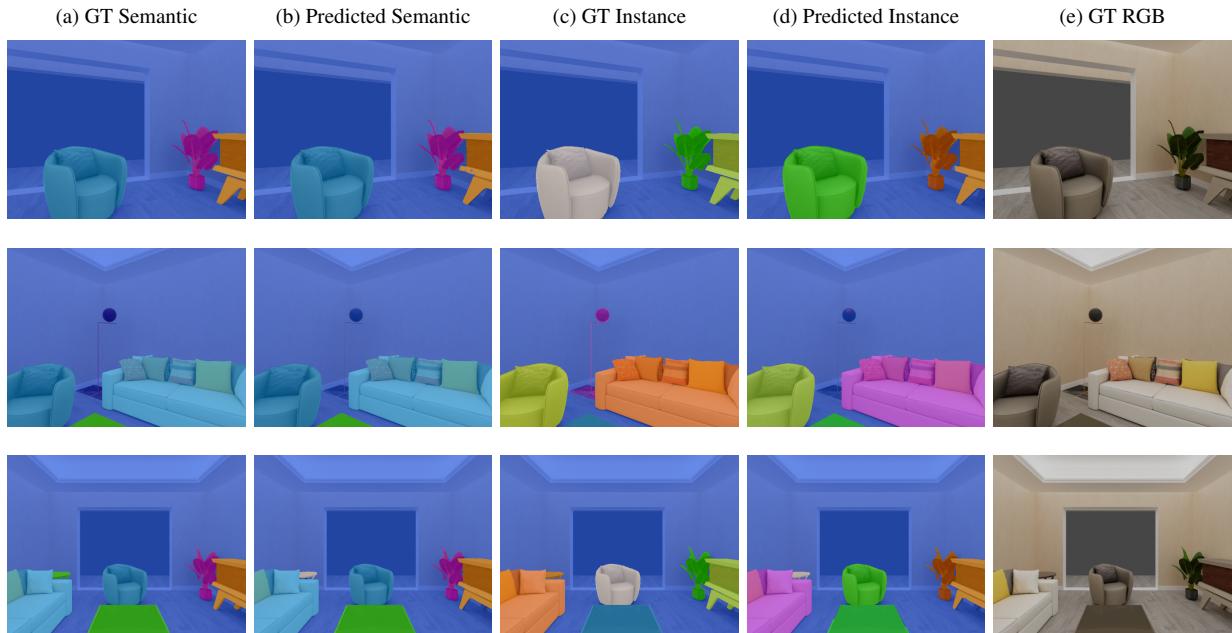


Figure 5: **Results.** This figure shows our multi-view results from a single scene. Our method achieves high consistency across images (Column (b) and Column (e)).

small holes on a continuous instance mask, which also suppresses noisy mask predictions floating around the objects for some scenes.

4.4.2 Mask Refinement

Using CascadePSP to refine the segmentation of the partially trained instance field helps improve the mask consis-

tency over different frames, especially for missing predictions which leads to wrong supervision for the instance field training and cannot be easily fixed with the instance regularization. Table 2 shows that such refinement significantly improves the mIoU and PQ of the model.

5. Conclusion

In this paper, we propose the *Instance Neural Radiance Field* and the corresponding training approach and techniques to perform continuous 3D instance segmentation in NeRF without using ground-truth segmentation information during inference. Extensive experiments and ablation studies are performed on a synthetic indoor NeRF dataset to demonstrate the effectiveness of our method. We perform comparison to relevant 2D segmentation methods and prior works in NeRF segmentation. As one of the first successful attempts on 3D instance segmentation in NeRF, we believe that our Instance-NeRF will contribute significantly to fundamental research on detection, segmentation in NeRF, as well as down-stream applications leveraging NeRF representation.

References

- [1] Wang Bing, Lu Chen, and Bo Yang. Dm-nerf: 3d scene geometry decomposition and manipulation from 2d images. *arXiv preprint arXiv:2208.07227*, 2022.
- [2] Daniel Bolya, Chong Zhou, Fanyi Xiao, and Yong Jae Lee. Yolact: Real-time instance segmentation. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 9157–9166, 2019.
- [3] Zhaowei Cai and Nuno Vasconcelos. Cascade r-cnn: Delving into high quality object detection. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 6154–6162, 2018.
- [4] Anpei Chen, Zexiang Xu, Andreas Geiger, Jingyi Yu, and Hao Su. Tensorf: Tensorial radiance fields. In *ECCV*, 2022.
- [5] Kai Chen, Jiangmiao Pang, Jiaqi Wang, Yu Xiong, Xiaoxiao Li, Shuyang Sun, Wansen Feng, Ziwei Liu, Jianping Shi, Wanli Ouyang, et al. Hybrid task cascade for instance segmentation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 4974–4983, 2019.
- [6] Xinlei Chen, Ross Girshick, Kaiming He, and Piotr Dollár. Tensormask: A foundation for dense object segmentation. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 2061–2069, 2019.
- [7] Bowen Cheng, Ishan Misra, Alexander G. Schwing, Alexander Kirillov, and Rohit Girdhar. Masked-attention mask transformer for universal image segmentation. In *CVPR*, 2022.
- [8] Ho Kei Cheng, Jihoon Chung, Yu-Wing Tai, and Chi-Keung Tang. CascadePSP: Toward class-agnostic and very high-resolution segmentation via global and local refinement. In *CVPR*, 2020.
- [9] Tianheng Cheng, Xinggang Wang, Lichao Huang, and Wenyu Liu. Boundary-preserving mask r-cnn. In *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XIV 16*, pages 660–676. Springer, 2020.
- [10] Angela Dai and Matthias Niessner. 3dmv: Joint 3d-multi-view prediction for 3d semantic scene segmentation. In *Proceedings of the European Conference on Computer Vision (ECCV)*, September 2018.
- [11] Kangle Deng, Andrew Liu, Jun-Yan Zhu, and Deva Ramanan. Depth-supervised NeRF: Fewer views and faster training for free. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2022.
- [12] Zhiwen Fan, Peihao Wang, Xinyu Gong, Yifan Jiang, Dejia Xu, and Zhangyang Wang. Nerf-sos: Any-view self-supervised object segmentation from complex real-world scenes. *arXiv e-prints*, pages arXiv–2209, 2022.
- [13] Huan Fu, Rongfei Jia, Lin Gao, Mingming Gong, Binjiang Zhao, Steve Maybank, and Dacheng Tao. 3d-future: 3d furniture shape with texture. *International Journal of Computer Vision*, pages 1–25, 2021.
- [14] Xiao Fu, Shangzhan Zhang, Tianrun Chen, Yichong Lu, Lanyun Zhu, Xiaowei Zhou, Andreas Geiger, and Yiyi Liao. Panoptic nerf: 3d-to-2d label transfer for panoptic urban scene segmentation. In *International Conference on 3D Vision (3DV)*, 2022.
- [15] Ross Girshick. Fast r-cnn. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, December 2015.
- [16] Ross Girshick. Fast r-cnn. In *Proceedings of the IEEE international conference on computer vision*, pages 1440–1448, 2015.
- [17] Ross Girshick, Jeff Donahue, Trevor Darrell, and Jitendra Malik. Rich feature hierarchies for accurate object detection and semantic segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2014.
- [18] Lei Han, Tian Zheng, Lan Xu, and Lu Fang. Occuseg: Occupancy-aware 3d instance segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2020.
- [19] Kaiming He, Georgia Gkioxari, Piotr Dollar, and Ross Girshick. Mask r-cnn. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, Oct 2017.
- [20] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. Mask r-cnn. In *ICCV*, pages 2961–2969, 2017.
- [21] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Spatial pyramid pooling in deep convolutional networks for visual recognition. *ECCV*, abs/1406.4729, 2014.
- [22] Jan Hendrik Hosang, Rodrigo Benenson, Piotr Dollár, and Bernt Schiele. What makes for effective detection proposals? *PAMI*, abs/1502.05082, 2015.
- [23] Ji Hou, Angela Dai, and Matthias Nießner. 3d-sis: 3d semantic instance segmentation of rgb-d scans. In *Proc. Computer Vision and Pattern Recognition (CVPR), IEEE*, 2019.

- [24] Benran Hu, Junkai Huang, Yichen Liu, Yu-Wing Tai, and Chi-Keung Tang. Nerf-rpn: A general framework for object detection in nerfs. *arXiv preprint arXiv:2211.11646*, 2022.
- [25] Maximilian Jaritz, Jiayuan Gu, and Hao Su. Multi-view pointnet for 3d scene understanding. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV) Workshops*, Oct 2019.
- [26] Abhijit Kundu, Kyle Genova, Xiaoqi Yin, Alireza Fathi, Caroline Pantofaru, Leonidas J Guibas, Andrea Tagliasacchi, Frank Dellaert, and Thomas Funkhouser. Panoptic neural fields: A semantic object-aware neural scene representation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12871–12881, 2022.
- [27] Weicheng Kuo, Anelia Angelova, Jitendra Malik, and Tsung-Yi Lin. Shapemask: Learning to segment novel objects by refining shape priors. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 9207–9216, 2019.
- [28] Verica Lazova, Vladimir Guzov, Kyle Olszewski, Sergey Tulyakov, and Gerard Pons-Moll. Control-nerf: Editable feature volumes for scene rendering and manipulation. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 4340–4350, 2023.
- [29] Yi Li, Haozhi Qi, Jifeng Dai, Xiangyang Ji, and Yichen Wei. Fully convolutional instance-aware semantic segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2359–2367, 2017.
- [30] Zhidong Liang, Ming Yang, Hao Li, and Chunxiang Wang. 3d instance embedding learning with a structure-aware loss function for point cloud segmentation. *IEEE Robotics and Automation Letters*, 5(3):4915–4922, 2020.
- [31] 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, 2014.
- [32] Xinhang Liu, Jiaben Chen, Huai Yu, Yu-Wing Tai, and Chi-Keung Tang. Unsupervised multi-view object segmentation using radiance field propagation. In *Advances in Neural Information Processing Systems*, volume 35, 2022.
- [33] Francesco Locatello, Dirk Weissenborn, Thomas Unterthiner, Aravindh Mahendran, Georg Heigold, Jakob Uszkoreit, Alexey Dosovitskiy, and Thomas Kipf. Object-centric learning with slot attention. *Advances in Neural Information Processing Systems*, 33:11525–11538, 2020.
- [34] 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 *ECCV*, 2020.
- [35] Thomas Müller, Alex Evans, Christoph Schied, and Alexander Keller. Instant neural graphics primitives with a multiresolution hash encoding. *ACM TOG*, 41(4):102:1–102:15, July 2022.
- [36] Michael Niemeyer, Jonathan T. Barron, Ben Mildenhall, Mehdi S. M. Sajjadi, Andreas Geiger, and Noha Radwan. Regnerf: Regularizing neural radiance fields for view synthesis from sparse inputs. In *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2022.
- [37] Michael Niemeyer and Andreas Geiger. Giraffe: Representing scenes as compositional generative neural feature fields. In *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2021.
- [38] 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 (CVPR)*, July 2017.
- [39] Charles Ruizhongtai Qi, Li Yi, Hao Su, and Leonidas J Guibas. Pointnet++: Deep hierarchical feature learning on point sets in a metric space. In I. Guyon, U. Von Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc., 2017.
- [40] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster R-CNN: Towards real-time object detection with region proposal networks. In *Advances in Neural Information Processing Systems (NIPS)*, 2015.
- [41] Karl Stelzner, Kristian Kersting, and Adam R Kosiorek. Decomposing 3d scenes into objects via unsupervised volume segmentation. *arXiv:2104.01148*, 2021.
- [42] Jiaxiang Tang. Torch-npg: a pytorch implementation of instant-npg, 2022. <https://github.com/ashawkey/torch-npg>.
- [43] Lyne Tchapmi, Christopher Choy, Iro Armeni, JunYoung Gwak, and Silvio Savarese. Segcloud: Semantic segmentation of 3d point clouds. In *2017 international conference on 3D vision (3DV)*, pages 537–547. IEEE, 2017.
- [44] Suhaní Vora, Noha Radwan, Klaus Greff, Henning Meyer, Kyle Genova, Mehdi S. M. Sajjadi, Etienne Pot, Andrea Tagliasacchi, and Daniel Duckworth. Nesf: Neural semantic fields for generalizable semantic segmentation of 3d scenes, 2021.
- [45] Thang Vu, Kookhoi Kim, Tung M. Luu, Thanh Nguyen, and Chang D. Yoo. Softgroup for 3d instance segmentation on point clouds. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2708–2717, June 2022.
- [46] Chen Wang, Xian Wu, Yuan-Chen Guo, Song-Hai Zhang, Yu-Wing Tai, and Shi-Min Hu. Nerf-sr: High-quality neural radiance fields using super-sampling. *arXiv*, 2021.
- [47] Weiyue Wang, Ronald Yu, Qiangui Huang, and Ulrich Neumann. Sgnp: Similarity group proposal network for 3d point cloud instance segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2018.
- [48] Xinlong Wang, Tao Kong, Chunhua Shen, Yuning Jiang, and Lei Li. Solo: Segmenting objects by locations. In *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XVIII 16*, pages 649–665. Springer, 2020.
- [49] Xinlong Wang, Rufeng Zhang, Tao Kong, Lei Li, and Chunhua Shen. Solov2: Dynamic and fast instance segmentation. *Advances in Neural information processing systems*, 33:17721–17732, 2020.
- [50] Zirui Wang, Shangzhe Wu, Weidi Xie, Min Chen, and Victor Adrian Prisacariu. Nerf–: Neural radiance

- fields without known camera parameters. *arXiv preprint arXiv:2102.07064*, 2021.
- [51] Zhirong Wu, Shuran Song, Aditya Khosla, Fisher Yu, Ligang 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 (CVPR)*, June 2015.
- [52] Enze Xie, Peize Sun, Xiaoge Song, Wenhui Wang, Xuebo Liu, Ding Liang, Chunhua Shen, and Ping Luo. Polarmask: Single shot instance segmentation with polar representation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 12193–12202, 2020.
- [53] Bo Yang, Jianan Wang, Ronald Clark, Qingyong Hu, Sen Wang, Andrew Markham, and Niki Trigoni. Learning object bounding boxes for 3d instance segmentation on point clouds. *Advances in neural information processing systems*, 32, 2019.
- [54] Alex Yu, Rui long Li, Matthew Tancik, Hao Li, Ren Ng, and Angjoo Kanazawa. PlenOctrees for real-time rendering of neural radiance fields. In *ICCV*, 2021.
- [55] Hong-Xing Yu, Leonidas J. Guibas, and Jiajun Wu. Unsupervised discovery of object radiance fields. In *International Conference on Learning Representations*, 2022.
- [56] Kai Zhang, Gernot Riegler, Noah Snavely, and Vladlen Koltun. Nerf++: Analyzing and improving neural radiance fields. *arXiv preprint arXiv:2010.07492*, 2020.
- [57] Shuaifeng Zhi, Tristan Laidlow, Stefan Leutenegger, and Andrew Davison. In-place scene labelling and understanding with implicit scene representation. In *Proceedings of the International Conference on Computer Vision (ICCV)*, 2021.
- [58] Shuaifeng Zhi, Tristan Laidlow, Stefan Leutenegger, and Andrew J. Davison. In-place scene labelling and understanding with implicit scene representation. In *ICCV*, 2021.

A. NeRF 3D Instance Segmentation Dataset

Leveraging 3D-FRONT [13] and the data generating approach of [24], we produce a new benchmark for instance-level 3D scene understanding curated for NeRF. 3D-FRONT is a large-scale synthetic indoor scene dataset, from which NeRF-RPN renders RGB images and layout configuration and tailors it as a benchmark for object detection task in NeRF. As shown in Table A.1, apart from multi-view images with camera poses and ground truth 3D bounding boxes, 2D ground truth instance segmentation and 3D ground truth instance masks on grids with class labels are included in our new dataset, which can be used for 3D segmentation in NeRF and other research areas.

Dataset	NeRF-RPN	Ours
# scenes	152	1015
RGB images	✓	✓
Camera poses	✓	✓
3D bounding boxes	✓	✓
2D inst seg GT	-	✓
3D voxelized inst seg GT	-	✓

Table A.1: A comparison between the NeRF dataset in NeRF-RPN and ours.

B. NeRF-RCNN Architecture

In this section, we describe the architecture of NeRF-RCNN in detail. NeRF-RCNN is a proposal-based 3D mask prediction model that imitates the architecture of Mask-RCNN [20]. The input of NeRF-RCNN are the 3D radiance and density grid sampled from a pre-trained NeRF, and the Region of Interests (RoI) provided by NeRF-RPN [24]. For each RoI, we set the ground truth box with the largest intersection over union (IoU) as its regression target.

The first part of NeRF-RCNN is a backbone identical to [24] for feature extraction. The second part takes the feature of each RoI as input and predicts the 3D bounding box, classification probability and discrete 3D mask. To obtain the feature of a single proposal on a feature map, we extend RoIAlign [20] with one more dimension, making all RoI features consistent. Aligned features are fed into two heads, namely *box head* and *mask head*. *Box head* first flattens the inputs for fully connected layer encoding and then separates into box branch and classification branch. The box branch further regresses a RoI to a more accurate bounding box for each class, while the classification branch predicts the classification scores. We follow similar network architecture in [20] by changing the 2D convolution and strided convolution layers to their corresponding 3D version. The loss function of *box head* consists of two parts:

$$\mathcal{L}_{cls} = \frac{1}{|\mathcal{N}|} \sum_{i \in \mathcal{N}} \mathcal{L}_{BCE}(\mathbf{p}_i, \mathbf{p}_i^*), \quad (13)$$

$$\mathcal{L}_{reg} = \frac{1}{|\mathcal{N}_p|} \sum_{i \in \mathcal{N}_p} \sum_{k=1}^L p_{i,k}^* \mathcal{L}_{smooth}(\mathbf{t}_{i,k}, \mathbf{t}_i^*), \quad (14)$$

where \mathbf{p}_i is the predicted classification score vector after sigmoid, $p_{i,k}$ is the k -th dimension of \mathbf{p}_i , $\mathbf{t}_{i,k}$ is the box offsets of class k , $\mathbf{p}_i^*, \mathbf{t}_i^*$ are ground-truth targets, \mathcal{N} is the set of sampled RoIs, \mathcal{N}_p is the set of positive samples, and L is the number of classes including background. \mathcal{L}_{BCE} and \mathcal{L}_{smooth} denote the binary cross entropy(BCE) loss and the smooth L1 loss in [16] respectively. Note that for an RoI associated with ground truth class c , only the c -th box regression BCE loss contributes to the total loss.

$\mathbf{t}_{i,k} = (t_{x,k}, t_{y,k}, t_{z,k}, t_{w,k}, t_{l,k}, t_{h,k})$ is the box head output. The relationship between $\mathbf{t}_{i,k}$ and bounding box parameters x, y, z, w, h, l is defined similarly to [24]:

$$\begin{aligned} t_{x,k} &= (x_k - x_a)/w_a, & t_{y,k} &= (y_k - y_a)/l_a, \\ t_{z,k} &= (z_k - z_a)/h_a, & t_{z,k} &= \log(w_k/w_a), \\ t_{l,k} &= \log(l_k/l_a), & t_{h,k} &= \log(h_k/h_a), \end{aligned} \quad (15)$$

where x_k, y_k, z_k are the center coordinate, w_k, l_k, h_k are the lengths of sides, and $x_a, y_a, z_a, w_a, l_a, h_a$ are the corresponding parameter of the RoI.

The mask head is a convolutional neural network which predicts L binary masks with size $m \times m \times m$ for each RoI. $m = 5$ is used for the box head, and $m = 10$ for the mask head. We also apply the sigmoid function as the activation. The loss for the mask head \mathcal{L}_M is defined as

$$\mathcal{L}_M = \frac{\lambda}{|\mathcal{N}_p|} \sum_{i \in \mathcal{N}_p} \sum_{k=1}^L p_{i,k}^* \mathcal{L}_p(\mathbf{m}_{i,k}, \mathbf{m}_i^*), \quad (16)$$

where \mathbf{m}_i^* is the ground truth mask and $\mathbf{m}_{i,k}$ is the predicted mask of class k . Similar to box regression, only the mask BCE loss corresponding to the ground truth label is included in \mathcal{L}_M .

The total loss of Instance-NeRF \mathcal{L} is

$$\mathcal{L} = \mathcal{L}_{cls} + \lambda_1 \mathcal{L}_{reg} + \lambda_2 \mathcal{L}_{mask}, \quad (17)$$

where λ_1, λ_2 are hyper-parameters.

C. Qualitative Results of Ablation

We present additional visualization and qualitative comparisons to demonstrate the effectiveness of our proposed instance label regularization loss, and the mask refinement stage. Figure C.6 shows the comparison between the instance segmentation results of Instance-NeRF optimized with and without instance label regularization. Notice how regularization helps suppress the noise in class labels resulting from inconsistent multi-view segmentation, and partially closes the holes in segmentation caused by false negative predictions from Mask2Former.

Although adding instance label regularization can help smooth the instance field, the segmentation quality of the preliminary results can still be unsatisfactory, especially on the silhouette of the objects. Besides, regularization can sometimes smooth out detailed structures in the segmentation, like thin chair legs or lamp stands. As illustrated in Figure C.7, performing 2D mask refinement using CascadePSP on the Instance-NeRF results and using it to guide further training can significantly improve the segmentation quality on the object boundaries.

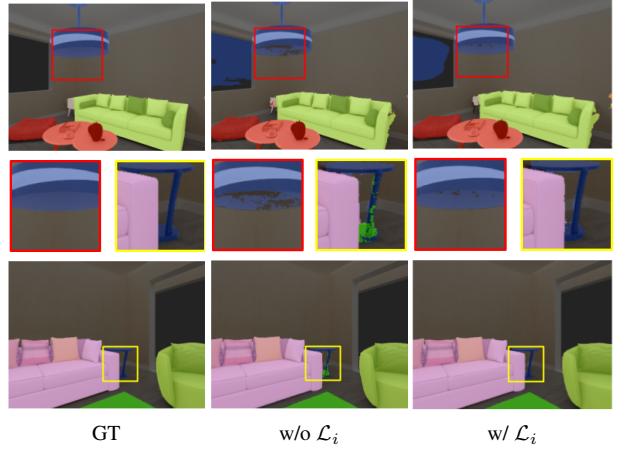


Figure C.6: **Ablation on instance label regularization.** Results are taken from the Instance-NeRF before going through the refinement process.

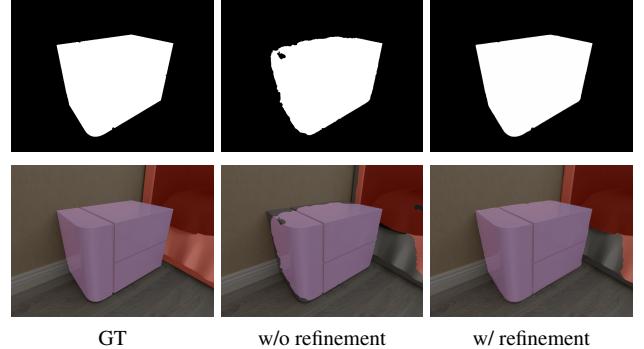


Figure C.7: **Ablation on 2D mask refinement.** The first row shows the separate mask for the nightstand, which is used as the input to CascadePSP. The separate masks after refinement are then composed into a single segmentation map to further optimize the instance field, the results of which are presented in the bottom row.

D. Additional Qualitative Comparison

We demonstrate extra qualitative comparisons between our method and other related methods as mentioned in the main paper. The results are given in Figure D.8. A video result display is also enclosed in the supplementary material.



Figure D.8: **Additional Comparison.** This figure illustrates the comparison between ours and other methods. For each group of comparison, rows from top to bottom are (a) ground truth RGB images or the rendered RGB images from the NeRF models, (b) semantic segmentation, and (c) instance segmentation. The instance segmentation results from Semantic-NeRF are left empty as it does not produce instance-level information.



Figure D.8: **Additional Comparison (cont.)** This figure illustrates the comparison between ours and other methods. For each group of comparison, rows from top to bottom are (a) ground truth RGB images or the rendered RGB images from the NeRF models, (b) semantic segmentation, and (c) instance segmentation. The instance segmentation results from Semantic-NeRF are left empty as it does not produce instance-level information.