# Entity-NeRF: Detecting and Removing Moving Entities in Urban Scenes

Takashi Otonari[1]    Satoshi Ikehata[2,3,1]    Kiyoharu Aizawa[1]

[1]The University of Tokyo    [2]National Institute of Informatics (NII)    [3]Tokyo Institute of Technology

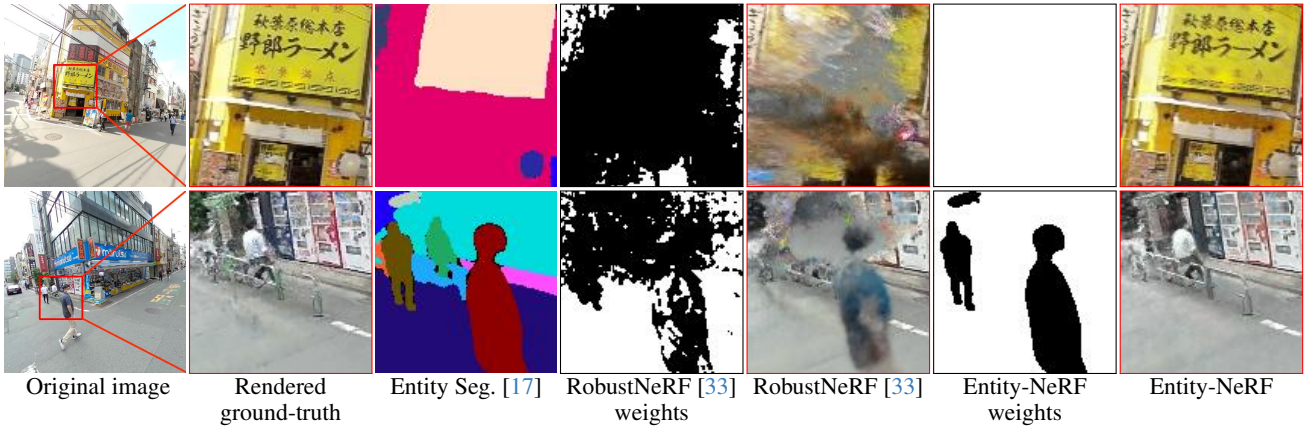{otonari,aizawa}@hal.t.u-tokyo.ac.jp    sikehata@nii.ac.jp

Figure 1. In urban scenes, statistical approach [33] mistakes complex backgrounds for moving objects (top) and fails to remove small moving objects (bottom). On the other hand, Entity-NeRF can reconstruct complex backgrounds and remove small moving objects.

## Abstract

*Recent advancements in the study of Neural Radiance Fields (NeRF) for dynamic scenes often involve explicit modeling of scene dynamics. However, this approach faces challenges in modeling scene dynamics in urban environments, where moving objects of various categories and scales are present. In such settings, it becomes crucial to effectively eliminate moving objects to accurately reconstruct static backgrounds. Our research introduces an innovative method, termed here as Entity-NeRF, which combines the strengths of knowledge-based and statistical strategies. This approach utilizes entity-wise statistics, leveraging entity segmentation and stationary entity classification through thing/stuff segmentation. To assess our methodology, we created an urban scene dataset masked with moving objects. Our comprehensive experiments demonstrate that Entity-NeRF notably outperforms existing techniques in removing moving objects and reconstructing static urban backgrounds, both quantitatively and qualitatively. [1]*

## 1. Introduction

Novel view synthesis is rapidly evolving, which enables the creation of new visual content such as the immersive views found in Google Maps and the free-viewpoint visualizations in sports broadcasts. However, a key innovation in this field, Neural Radiance Fields (NeRF) [19], faces challenges when dealing with urban scenes whose complexity is inherently high due to a large number of dynamic elements present, such as moving vehicles, pedestrians, changing lighting conditions, and varying shadows. The ability to accurately render and reconstruct such scenes is crucial for several applications including autonomous navigation, surveillance, and virtual urban exploration among others.

The challenge of handling dynamic scenes has been a notable point of extension within the domain of Neural Radiance Fields, and two major approaches prevail. The first involves explicit modeling of scene dynamics, which concurrently encodes both static and dynamic information, exemplified by methods such as D-NeRF [27], HyperNeRF [24], and RoDynRF [16]. The second approach adopts a more statistical perspective, treating scene dynamics as outliers in relation to the static elements [33]. The elimination of dynamic elements in the scene contributes to a reduction in clutter, enhances the comprehension of the scene, and improves visual fidelity.

Despite the progress in research on novel-view synthesis of dynamic scenes, to our knowledge, there is no effective method for unbounded scenes like urban environments, where a multitude of moving objects of various categories and scales, such as people, cars, and bicycles, coexist. For

---

[1]Our project page is available at https://otonari726.github.io/entitynerf/

instance, current methods in the former category often target specific objects, deal with a minority of moving objects within a scene, and are restricted to narrow scene boundaries. Moreover, while the latter approach can handle multiple objects simultaneously, it solely relies on the statistics of reconstruction errors for outlier separation and does not function effectively when dynamic objects vary in scale or when the background is complex, as shown in Fig. 1.

In this work, we address the task of learning static NeRFs of dynamic urban scenes. To identify a multitude of moving objects of various categories and scales, we propose a hybrid method that integrates the strengths of both knowledge-based and statistical approaches with Entity-wise Average of Residual Ranks (EARR) and stationary entity classification. EARR identifies distractors by entity-wise statistics based on entity segmentation [17]. In addition, the stationary entity classification using the thing/stuff segmentation [53, 60] enables more efficient learning by incorporating complex backgrounds such as building from the early stages of learning.

To evaluate our proposed method, we annotated moving objects in three videos captured in urban scenes with challenging settings and rendered images using a static NeRF which removed the masked moving objects. Using the overall Peak Signal-to-Noise Ratio (PSNR) of the image to measure whether moving objects, which constitute only a small portion of the image, have been appropriately removed is challenging. Therefore, we evaluate the moving object removal (foreground PSNR) and the static background reconstruction (background PSNR) separately. Our experiments demonstrate that our method effectively removes moving objects and reconstructs static backgrounds in urban scenes, while still maintaining accuracy on existing datasets.

## 2. Related Works

### 2.1. NeRF on Dynamic Scenes

NeRF [19] represents coordinate-based neural networks that predict the radiance from a specific view and opacity at any given 3-D coordinate. To render novel views of a scene via ray-tracing, NeRF is trained by minimizing the difference between each pixel's rendered and observed colors given calibrated multi-view images of the scene.

In the original NeRF, handling dynamic scenes is challenging due to the inherent assumption that the entire scene remains static. To address this issue, subsequent research has proposed methods that either explicitly learn scene dynamics by category-specific methods [6, 8, 14, 15, 25, 26, 30, 34, 37, 46, 50, 57], detection [10, 22], deformation [18, 23, 24, 27, 41, 44, 48, 54, 56], flow [4, 5, 7, 12, 55], multiple synchronized videos [11, 47, 58], depth-based approaches [52], or treat moving objects as outliers in a robust approach [33].

**Detection-based Approach:** Neural Scene Graphs [22] and Panoptic Neural Fields [10] provide a structured approach to explicitly detecting and modeling individual dynamic objects within dynamic scenes. While these methods facilitate object-level manipulation, moving object detection is hindered by occlusions, diverse object types, and scales in the urban environment.

**Deformation-based Approach:** Without object detection, some methods such as D-NeRF [27], Nerfies [23] and HyperNeRF [24] represent scenes using a deformation field, mapping observations to neighboring frames or a canonical scene. However, they are limited to small-motion, object-centric scenes due to challenges in representing entire sequences with a single canonical voxel.

In recent efforts, $D^2$NeRF [51] separates moving objects, static backgrounds, and shadows into three fields using regularization. DynIBaR [13] aggregates multi-view image features in a motion-adjusted ray space, while Ro-DynRF [16] uses a time-dependent MLP and single-view depth priors. FSDNeRF [45] uses data-driven optical flow for backward deformation computation to handle rapid motion.

In urban settings, accurate optical flow estimation faces challenges due to numerous cluttered objects of varying scales, which often result in incorrect dynamics modeling. Moreover, deformation-based methods struggle with frames having large temporal steps, constraining their urban modeling use from a discrete set of multi-view images.

**Robust Approach:** Relatively little attention has been paid to removing non-static elements from discrete multi-view images rather than a continuous video stream. One straightforward approach is to segment and ignore pixels during training that are likely to be transient objects [32, 39], for example, by applying a data-driven segmentation model. However, removing objects based on object semantics is risky since semantic segmentation is far from perfect and runs the risk of erroneously removing static objects that are typically mobile (e.g., cars, pedestrians).

We can also use a robust estimator. RobustNeRF [33] has proposed a purely statistical approach to remove moving elements as outliers by analyzing the patch-wise statistics of reconstruction errors. By formulating training as a form of iterative reweighted least squares, this method can robustly separate inliers and outliers, which are not limited to specific predefined categories, from photo collections rather than videos. While effective, this method fixes the hyperparameters such as patch size, which becomes problematic when applied to urban scenes where there is a variety of moving objects of different types and scales.
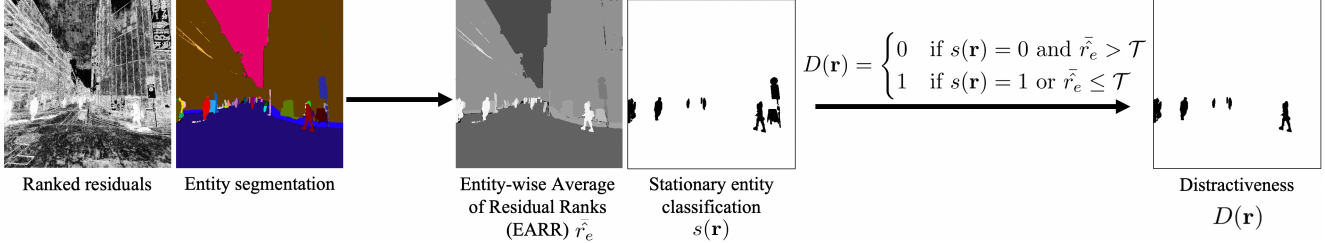
Figure 2. **Overview of our Entity-NeRF pipeline.** $D(\mathbf{r}) = 0$ if Entity-wise Average of Residual Ranks (4.1) of the entities labeled 'thing' in the stationary entity classification (4.2) is greater than a threshold value $\mathcal{T}$. The 'thing' label for the stationary entity classification is given as $s(\mathbf{r}) = 0$ and the 'stuff' label as $s(\mathbf{r}) = 1$.

## 2.2. NeRF on Unbounded Scenes

The original NeRF struggles with unbounded scenes due to sparse rays at greater distances. Adaptations like NeRF++[59], which introduced inverted sphere parameterization, F2-NeRF[49] with its perspective warping, and scene contraction approaches by MERF [31] and Mip-NeRF 360 [2], have been developed to address this. Recent developments, Zip-NeRF [3] and Nerfacto [40], further refine these methods for unbounded environments.

For large city-scale scenes, Block-NeRF [39], Mega-NeRF [42], and SUDS [43] segment scenes into blocks, applying NeRF within each. SUDS notably integrates additional data types like LiDAR for dynamic city-scale scenes.

## 2.3. Entity Segmentation

Entity segmentation is a new class of image segmentation tasks aiming to segment all entities in an image without predicting their semantic or instance labels [17, 28, 29, 36]. Eliminating the need for class labels is helpful for many practical applications, such as image manipulation and editing, where the quality of segmentation masks is crucial, but class labels are less important.

Recently, Qi *et al.* [17] presented a large-scale entity segmentation dataset and proposed CropFormer, a Transformer-based entity segmentation method. In our research, we utilize this result for training NeRF on urban dynamic scenes and demonstrate that Entity-wise Avarage of Residual Ranks (EARR) overcomes most issues found in patch-based counterparts in RobustNeRF [33].

## 3. Preliminaries

### 3.1. Problem Statement

In our context, there is no need to explicitly model moving objects. Therefore, our goal is to simply detect and segment out moving objects in the scene as distractors during the basic training pipeline of arbitrary NeRF models [1, 19]. More concretely, we want to label the distractiveness $D(\mathbf{r})$ to each ray $\mathbf{r}$ and reflect labels in photometric reconstruction

losses in training NeRF as

$$\mathcal{L}_{\mathbf{r}} = D(\mathbf{r}) \cdot \epsilon(\mathbf{r})$$
$$\epsilon(\mathbf{r}) = \|C_{\text{gt}}(\mathbf{r}) - C_{\text{pred}}(\mathbf{r})\|_2^2 \tag{1}$$

where

$$D(\mathbf{r}) = \begin{cases} 0 & \text{if } \mathbf{r} \text{ passes through a distractor,} \\ 1 & \text{otherwise.} \end{cases} \tag{2}$$

Here, $C_{\text{pred}}$ and $C_{\text{gt}}$ are rendered and observed pixel colors and $\epsilon(\mathbf{r})$ is the $\ell_2$ residual of them. While the task is formulated in a simple form, predicting $D(\mathbf{r})$ in urban scenes poses several fundamental challenges as detailed below.

### 3.2. Challenges in Urban Scenes

Firstly, prior knowledge of scene semantics is often incorporated to specify moving objects (*e.g.*, [32, 39]), as it provides accurate masks along the object's contours to a certain extent. However, relying solely on per-image scene semantics to specify distractors is insufficient. Urban environments feature a wide variety of moving objects, ranging from people, and vehicles, to minor elements like roadside cans. These objects are not always covered by standard semantic segmentation classes. Furthermore, even within common classes such as vehicles and pedestrians, they cannot always be identified as distractors based on class alone, as a parked car, for instance, should not be classified as a distractor.

In addition, a purely statistical approach (*e.g.*, [33]) may not always successfully identify distractors in urban scenes. For instance, RobustNeRF [33] assigns inlier/outlier labels to each non-overlapping $8 \times 8$ patch, considering only the statistics of reconstruction errors within $16 \times 16$ neighboring pixels. Due to this heuristic, RobustNeRF is only effective when the background is relatively simple when the reconstruction errors decrease rapidly, and the size of distractors is significantly larger than the predefined neighbor system. However, in most complex urban scenes, the varying distances to moving objects and the diverse coverage of
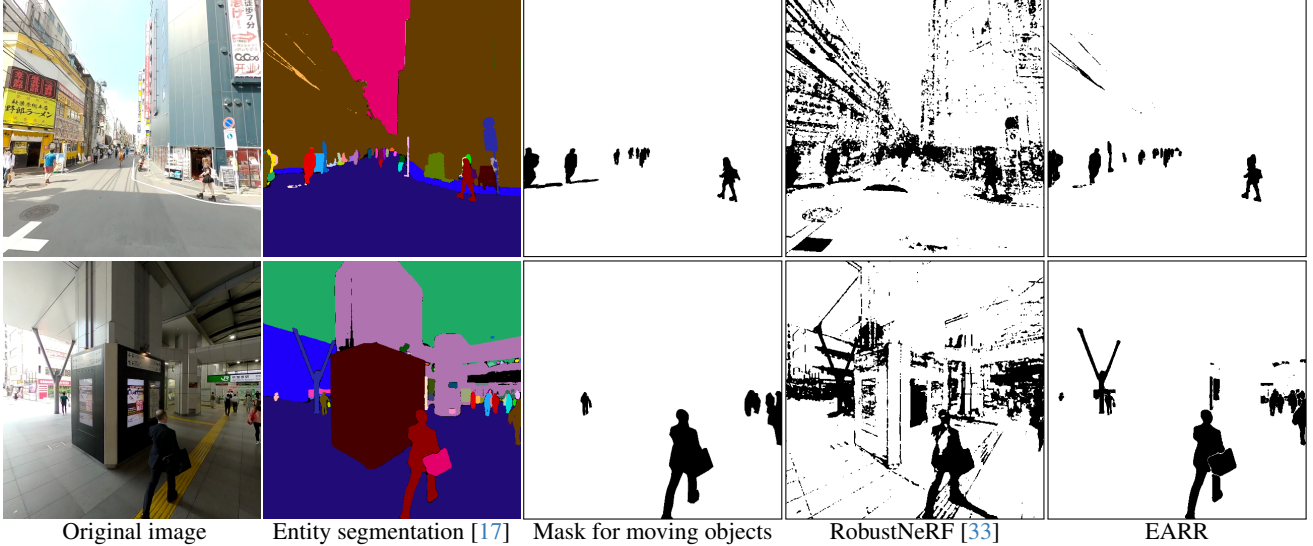
Figure 3. $D(\mathbf{r})$ of RobustNeRF [33] and our Entity-wise Average of Residual Ranks (EARR) at the end of training. Our EARR can more efficiently incorporate the background into learning.

field-of-view may hinder the statistical method's ability to distinguish between inliers and outliers effectively.

To overcome these limitations, we propose a hybrid method that integrates the strengths of both knowledge-based and statistical approaches. Specifically, we leverage the knowledge-based entity segmentation method's ability to generate accurate object contours and semantic segmentation method to identify non-moving objects, while incorporating the capability of the statistical method in adaptively handling varied scene dynamics and object movements. This synergy aims to create a more robust and versatile system for identifying distractors in urban environments, addressing the limitations of each method when used in isolation.

## 4. Method

In this section, we present our hybrid approach combining knowledge-based and statistical methods to identify moving distractors of varying sizes in urban scenes. Concretely, we introduce the method of distractor labeling using Entity-wise Average of Residual Ranks (EARR) (4.1), utilizing both data-driven segmentation networks and entity-wise statistics of reconstruction losses. As the statistics of reconstruction losses become highly unstable in complex background areas, we incorporate knowledge of scene semantics to identify the non-moving stuff, such as buildings (4.2). The overall pipeline is illustrated in Fig. 2.

### 4.1. Entity-wise Average of Residual Ranks (EARR)

A pre-trained entity segmentation network [17] provides high-quality segmentation for objects in real-world scenes including urban scenes. This quality is maintained regard-less of the objects' semantics or sizes in the image. Although it's not possible to determine if an entity is moving based on the segmentation result alone, we can assume that there is consistency in the distractor label across pixels within the same entity. This observation leads to a departure from conventional methods. Instead of assigning moving distractor labels to rays that intersect each pixel [39] or each patch [33], our approach labels individual entities.

To determine whether each entity is moving or not, we utilize the statistics of reconstruction loss in each entity. Similar to RobustNeRF [33], we follow the principle that rays passing through distractors lead to a lack of consistency across multiple viewpoints, resulting in larger reconstruction loss.

To clarify, we denote $\epsilon(i)$ as the $\ell_2$ residual between rendered and observed colors of a ray passing through the $i$-th pixel. Considering $N$ pixels (*i.e.*, $N$ rays) in a batch, we define a rank function $R(\epsilon(i))$ that inputs $\epsilon(i)$ and outputs an ordered rank, with the largest residual assigned $N$ and the smallest assigned 1 among the $N$ pixels. Then, the normalized residual rank $\hat{r}(i)$ is calculated by normalizing these ranks to the [0, 1] range and used for the distractor labeling in the following steps. The use of normalized residual rank instead of raw residual values is justified because, during the initial stages of training, the residuals tend to be large. Relying solely on the raw residual values for decision-making can lead to excessive false detection of distractors. In contrast, by employing a rank function, it is possible to exclude only a specific proportion of samples with large residual values, while ensuring that all other samples are included in the training process.

The normalized residual rank for each ray tends to be

higher when the ray passes through a moving object, due to an increase in the residual. However, even in static scenes, the rank can become high if the ray passes through complex backgrounds or geometric shapes. Therefore, instead of using a single ray as the basis for decision-making, we gather statistics on the Rank of each entity to use as a basis for labeling.

However, since the shape and size of entities vary greatly, it is not practical to sample rays passing through all pixels of each entity during training. To address this, we sample a patch of size $k \times k$ pixels (*i.e.*, $k$ should be sufficiently large and we choose $64$ in our implementation), cluster the labels of entities within it, and then calculate statistical measures for each cluster. Specifically, for a set of pixels corresponding to an entity ID ($e$) within a patch, designated as $S(e)$, we calculate its average as follows.

$$\bar{\hat{r}}_e = \frac{\Sigma_{i \in S(e)} \hat{r}(i)}{|S(e)|} \ . \tag{3}$$

where, the number of elements in each entity is denoted as $|S(e)|$. This process is repeated for all entities in the patch. If the average exceeds a certain threshold $\mathcal{T}$, we assign $D(\mathbf{r}) = 0$ and otherwise $D(\mathbf{r}) = 1$ to all rays passing through pixels corresponding to that entity ID. The choice of $\mathcal{T}$ is crucial and its choice will be discussed later.

Despite its simplicity, our approach, which combines knowledge-based entity segmentation results and statistics of residual ranks, functions much more robustly than methods based solely on statistics. As depicted in Fig. 3, our proposed method has been shown to accurately detect nearly all moving objects without excessively excluding inliers. This is in contrast to RobustNeRF [33], a purely statistical approach, which is prone to excessively identify distractors.

### 4.2. Cooperative Stationary Entity Classification

In the early stages of training, all residuals are high, leading to low reliability of residual ranks and their statistical measures. This is especially true in urban scenes where backgrounds contain numerous elements, such as traffic signs and complex building structures. These elements make NeRF training difficult, resulting in many inliers being included in samples above the rank threshold. These inliers could potentially be excluded from training during large training steps. To address this issue, for entities of classes such as buildings, sky, and roads in urban scenes, which are certainly stationary, we attempt to assign a value of 1 to $D(\mathbf{r})$ to ensure their inclusion in the learning process regardless of their residual ranks.

To implement this, we train a stationary entity classification network, which is an MLP with three linear layers and one classification layer. This network identifies whether each entity belongs to a class of stationary objects. The input to this network is a feature vector calculated for each entity, and the output is a class label of the entity, either 'stuff' or 'thing'. 'stuff' and 'thing' are defined in ADE20K [60] as non-accumulative and accumulative objects, respectively. Specifically, all movable object classes are included in 'thing', and we can safely include 'stuff' entities in the training. Feature vectors for individual entities are computed by averaging pixel-wise features within each entity. Concretely, feature maps are extracted from an image by applying pre-trained SAM [9] and DINOv2 [21] encoders, individually.

For training this network, we also adopt a cooperative approach combining prior knowledge and statistics. Specifically, instead of training the stationary entity classification network entirely on ADE20K, we adapt the network for individual scenes. This adaptation is done by continuously fine-tuning the MLP using pixels classified as stationary ($D(\mathbf{r}) = 1$) during the training, based on the ranked residuals described above. Concretely, we first train SegFormer [53] on ADE20K to output 'stuff'/'thing' labels, then apply it to each scene in which NeRF is trained to generate initial pseudo ground truth labels for stationary entities specific to each scene. Note that the labels for each entity are determined based on the voting of pixel-wise labels for each entity. The four-layer MLP is then pre-trained based on these pseudo ground truth labels. Every 100 steps of NeRF training, the MLP is fine-tuned using the entities that have been determined as $D(\mathbf{r}) = 1$ based on the ranked residuals. Entities for which the stationary entity classification network assigns 'stuff' labels are consistently included in the NeRF training, whereas entities assigned 'thing' labels are trained solely based on ranked residuals.

## 5. Results

### 5.1. Implementation Details

Our framework for labeling moving objects can be applied to all NeRF models that use photometric reconstruction loss similar to other robust approaches (*e.g.*, RobustNeRF [33]). In this paper, we incorporated our method into both Mip-NeRF 360 [2] and Nerfacto [40]. The implementation of each method is as follows. Note that we used the hyperparameter values from the original implementations. Mip-NeRF 360 trained on two Tesla A100 units and Nerfacto on one, using 16,384 samples during every iteration.

**Mip-NeRF 360**: Mip-NeRF 360 addresses challenges presented by unbounded scenes using a non-linear scene parameterization. We used the official implementation code from MultiNeRF [20], which contains an implementation of Mip-NeRF 360 [2] and RobustNeRF [33]. The models were trained for 250,000 iterations for each scene, taking approximately 24 hours.

**Nerfacto**: Nerfacto is implemented in NeRFStudio [40], which is a combination of various methods, rather than a

| Model | Loss | foreground PSNR↑ | background PSNR↑ | PSNR↑ | SSIM↑ | LPIPS↓ |
|---|---|---|---|---|---|---|
| Nerfacto [40] | Mean-squared error (MSE) | 12.10 | 25.07 | 24.96 | 0.87 | 0.10 |
| | RobustNeRF [33] | 17.63 | 21.74 | 23.19 | 0.84 | 0.12 |
| | Entity-NeRF (only EARR) | 19.48 | 23.68 | 24.63 | 0.84 | 0.13 |
| | **Entity-NeRF** | 19.82 | 24.00 | 24.93 | 0.85 | 0.12 |
| Mip-NeRF 360 [2] | Mean-squared error (MSE) | 11.40 | 27.36 | 24.22 | 0.88 | 0.13 |
| | RobustNeRF [33] | 20.15 | 22.52 | 22.87 | 0.83 | 0.18 |
| | Entity-NeRF (only EARR) | 20.20 | 25.49 | 25.21 | 0.85 | 0.14 |
| | **Entity-NeRF** | 20.74 | 25.50 | 25.23 | 0.84 | 0.15 |

Table 1. **Quantitative comparison with RobustNeRF [33] using Mip-NeRF 360 [2] and Nerfacto [40] on MovieMap Dataset.**



Figure 4. **MovieMap Dataset.** Only moving objects in the video are masked. Therefore, parked cars and stationary people are not masked.



Original image   D$^2$NeRF [51]   RoDynRF [16]   Entity-NeRF

Figure 5. **Qualitative comparison with dynamic NeRF methods (D$^2$NeRF [51] and RoDynRF [16]) on MovieMap Dataset.**

single published work, that has proven effective in real-world applications. The models were trained for 30,000 iterations for each scene, taking approximately 30 minutes.

The most important hyperparameter in our method is the threshold parameter for the averaged residual rank, $\mathcal{T}$. If this value is too high, it includes too many outliers as inliers, and if too low, it excludes inliers as outliers. To prevent the inclusion of outliers excessively as inliers during training, we set the threshold value to $\mathcal{T} = 0.8$. This decision is based on 78 manually annotated images from three scene images, with an average ratio of inliers being approximately 90.4%. The impact of this value is also evaluated in the next chapter.

### 5.2. Datasets

The proposed method was quantitatively evaluated using two real-world datasets. The first dataset is an urban scene dataset (MovieMap Dataset), generated from 360° videos from Movie Map [38], while the second comprises non-urban scenes published in [33] (RobustNeRF Dataset). Please see the supplementary for more details on RobustNeRF Dataset and quantitative comparisons.

**MovieMap Dataset**: Movie Map [38] offers an immersive interface for walking through cities in Japan using 360° videos. Movie Map currently offers the exploration of eight cities. From them, 360° videos of the Akihabara scene are used to create the dataset. Akihabara is characterized by an especially high number of pedestrians, vehicles, bicycles, and various moving objects, even by global standards. Additionally, the background comprises buildings of various shapes and colorful outdoor advertisements, embodying all the typical urban characteristics that make learning with Neural Radiance Fields (NeRF) challenging.

Our MovieMap Dataset comprises three different subsets, containing 12, 15, and 51 images each, sampled from the original 360° video in Movie Map. Manual annotation of moving objects was performed on each 360° image. As illustrated in Fig. 4, stationary cars and pedestrians are not labeled as distractors. For evaluation, we generated images without distractors by using Nerfacto [40] with annotated distractor labels to exclude moving objects from the training. This process allows us to render images from the same camera viewpoints to create background-only images.

In this study, we extract 14 perspective images from each 360° image, utilizing their degree of field (DoF) as input parameters for NeRF training. A full description of MovieMap Dataset is in the supplementary material.
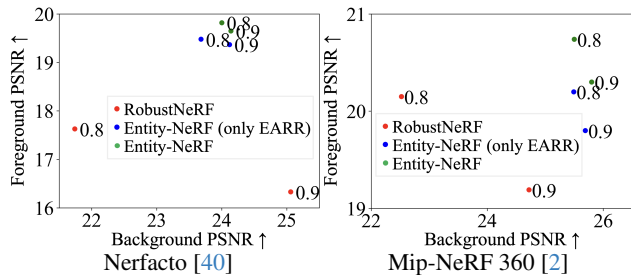
Figure 6. **The trade-off between the foreground/background PSNR.** The values in the figure indicate the threshold of inliers.
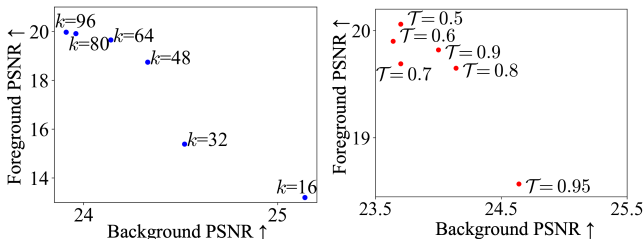


Figure 7. **Sensitivity to patch size ($k$) and threshold ($\mathcal{T}$).**



Figure 8. **Difference in the training curves.**



Figure 9. **Histogram of the average $D(\mathbf{r})$ per entity.**

## 5.3. Evaluation with dynamic NeRF models on MovieMap Dataset

To show urban environments' unsuitability for NeRF models encoding static and dynamic data, we tested $D^2$NeRF [51] and RoDynRF [16] as dynamic NeRF methods. For RoDynRF, designed solely for monocular video, we converted 360° images into 90° perspective images from the image center.

Fig. 5 demonstrates the qualitative results using $D^2$NeRF and RoDynRF in urban settings with motion. Both models struggle with moving objects and static background reconstruction, due to their limitations in handling excessive movement, complex motion, and scale variations. This makes developing dynamic NeRF models for urban scenes problematic. Subsequent experiments compare our work with RobustNeRF, which ignores moving objects.

## 5.4. Evaluation on MovieMap Dataset

**Qualitative comparison:** A comparison of Robust-NeRF [33] and our proposed method using Nerfacto [40] and Mip-NeRF 360 [2] on urban scenes is shown in Table 1. The mean-squared error (MSE) of incorporating all static backgrounds and moving objects into training enhances the PSNR of the backgrounds, which make up a larger percentage, leading to an increased overall PSNR. Therefore, it is crucial to evaluate foreground and background PSNR separately. Our proposed method achieved a background PSNR close to the mean-squared error while exceeding existing methods in foreground PSNR. The comparison with RobustNeRF showed consistent improvements.
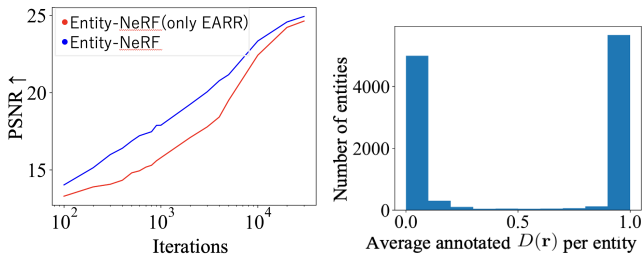
**Sensitivity of hyperparameters:** RobustNeRF and our proposed method present a trade-off between the foreground PSNR and background PSNR, influenced by a hyperparameter that determines the inlier to the residual ratio (denoted by $\mathcal{T}$ in our EARR). This trade-off is shown in the Fig. 6. Raising the inlier ratio improves background PSNR, but risks including moving objects in the learning, decreasing foreground PSNR. Similarly, lowering the inlier ratio worsens background PSNR, but removes many moving objects, boosting foreground PSNR. Entity-NeRF shows consistent improvements in foreground PSNR for Nerfacto and in background PSNR for Mip-NeRF 360. In addition, while RobustNeRF is biased toward improving one of the metrics, Entity-NeRF achieves more balanced results by increasing both metrics.

In addition, we performed a detailed sensitivity analysis on hyperparameters (*i.e.*, patch size $k$ and threshold $\mathcal{T}$ in EARR). As shown in Fig. 7, increasing $k$ improved the foreground PSNR with only a minor background PSNR impact. The choice of $\mathcal{T}$ proved less sensitive than $k$, provided it remains below the typical inlier ratio in urban scenes (*e.g.*, 90.4% in MovieMap dataset).

**Effects of stationary entity classification:** We conducted an analysis comparing the training curves of Entity-NeRF with and without stationary entity classification. As shown in Fig. 8, stationary entity classification not only significantly boosts training efficiency but also enhances final PSNR.

**Validity of entity segmentation:** To confirm the stable performance of our entity segmentation across various images, we calculated the average of annotated labels for each segmented entity using all images in the MovieMap dataset. Then, we constructed a histogram representing these average values for all entities. As depicted in Fig. 9, the distribution is noticeably skewed towards either 0 or 1, which indicates that entities are clearly segmented into moving ($= 0$) or static ($= 1$) entities.

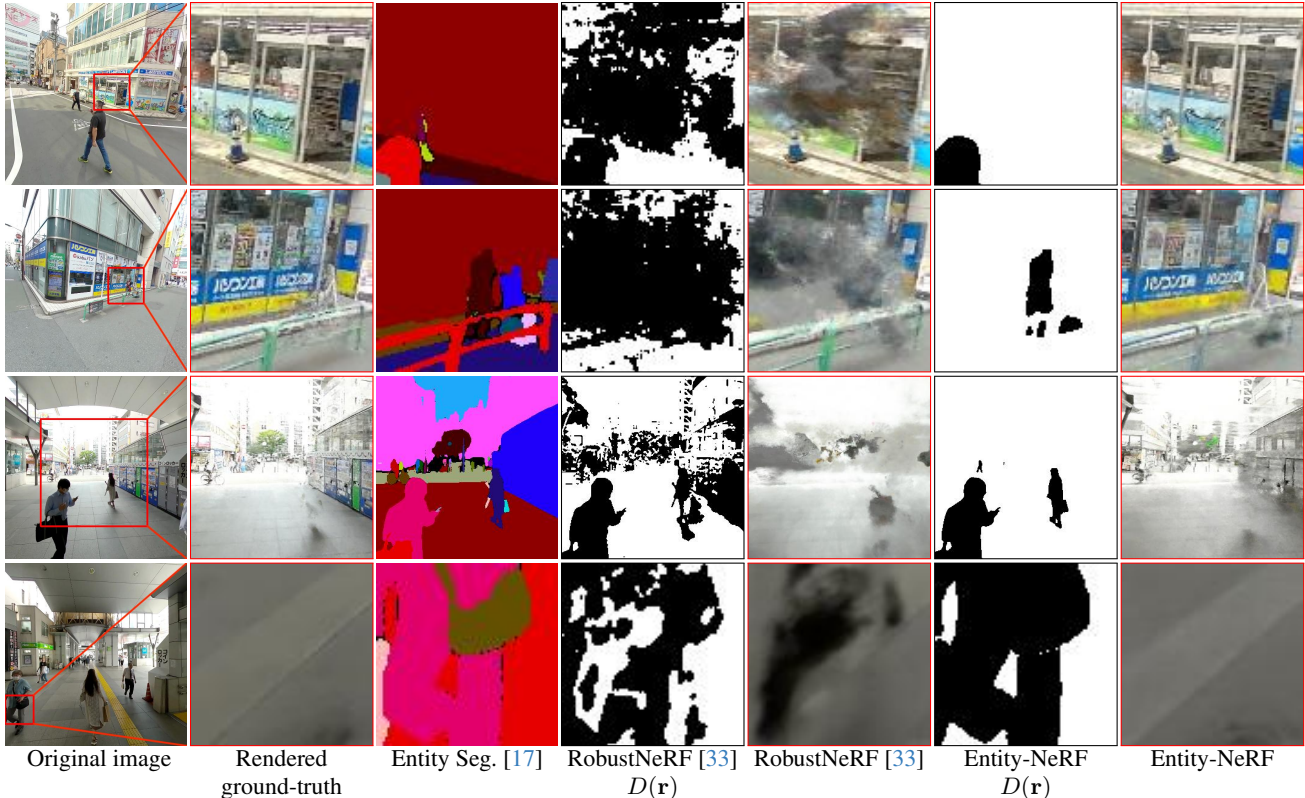**Qualitative comparison:** A qualitative comparison with

| Original image | Rendered ground-truth | Entity Seg. [17] | RobustNeRF [33] $D(\mathbf{r})$ | RobustNeRF [33] | Entity-NeRF $D(\mathbf{r})$ | Entity-NeRF |

Figure 10. **Qualitative comparison including Entity Seg. [17] on MovieMap Dataset.** $D(\mathbf{r})$ is calculated at the end of the training.



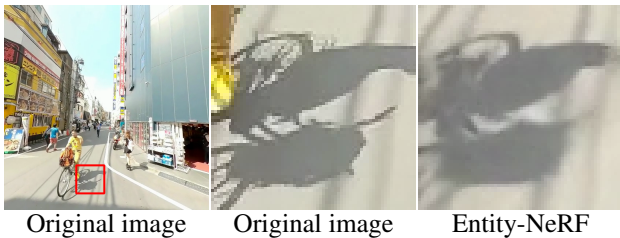| Original image | Original image | Entity-NeRF |

Figure 11. **Limitations.** Entity-NeRF cannot handle shadows.

RobustNeRF [33] is shown in Fig. 10. Our Entity-NeRF successfully reconstructed complex building walls that RobustNeRF mistakenly removed (Fig. 10-top three items). It also effectively removed moving objects that RobustNeRF failed to remove (Fig. 10-bottom two items). Thus, Entity-NeRF is clearly superior to RobustNeRF in both removing moving objects and reconstructing static backgrounds.

## 6. Conclusion

We address the problem of identifying and removing multiple moving objects of various categories and scales to build a NeRF for dynamic urban scenes. To solve this problem, we introduce Entity-wise Average of Residual Ranks designed to identify moving objects using entity-wise statistics and the stationary entity classification with thing/stuff segmentation to remove complex backgrounds in the early stages of NeRF training. Our evaluation using an urban scene dataset, where existing methods fail to model scene dynamics or remove moving objects, shows quantitatively and qualitatively that the proposed method works very well.

**Limitations**: While Entity-NeRF demonstrates outstanding performance in urban environments, it is subject to a few limitations. Firstly, if a large moving object dominates the image and thereby obscures the background from another perspective, there might be difficulty in successfully reconstructing the background hidden by the moving object. This issue, however, could potentially be mitigated by integrating existing inpainting techniques.

Moreover, as shown in Fig. 11, since shadows are not explicitly managed in the current framework, shadows cast by moving objects might be inadvertently incorporated into the training process. This issue may be resolved by using segmentation that includes shadows or by removing shadows in post-processing.

## Acknowledgments

# References

[1] Jonathan T. Barron, Ben Mildenhall, Matthew Tancik, Peter Hedman, Ricardo Martin-Brualla, and Pratul P. Srinivasan. Mip-nerf: A multiscale representation for anti-aliasing neural radiance fields. In *ICCV*, 2021. 3

[2] Jonathan T. Barron, Ben Mildenhall, Dor Verbin, Pratul P. Srinivasan, and Peter Hedman. Mip-nerf 360: Unbounded anti-aliased neural radiance fields. In *CVPR*, 2022. 3, 5, 6, 7, 2

[3] Jonathan T. Barron, Ben Mildenhall, Dor Verbin, Pratul P. Srinivasan, and Peter Hedman. Zip-nerf: Anti-aliased grid-based neural radiance fields. In *ICCV*, 2023. 3

[4] Quei-An Chen and Akihiro Tsukada. Flow supervised neural radiance fields for static-dynamic decomposition. In *ICRA*, 2022. 2

[5] Yilun Du, Yinan Zhang, Hong-Xing Yu, Joshua B Tenenbaum, and Jiajun Wu. Neural radiance flow for 4d view synthesis and video processing. In *ICCV*, 2021. 2

[6] Guy Gafni, Justus Thies, Michael Zollhöfer, and Matthias Nießner. Dynamic neural radiance fields for monocular 4d facial avatar reconstruction. In *CVPR*, 2021. 2

[7] Chen Gao, Ayush Saraf, Johannes Kopf, and Jia-Bin Huang. Dynamic view synthesis from dynamic monocular video. In *ICCV*, 2021. 2

[8] Shoukang Hu, Fangzhou Hong, Liang Pan, Haiyi Mei, Lei Yang, and Ziwei Liu. Sherf: Generalizable human nerf from a single image. In *ICCV*, 2023. 2

[9] Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C. Berg, Wan-Yen Lo, Piotr Dollar, and Ross Girshick. Segment anything. In *ICCV*, 2023. 5

[10] Abhijit Kundu, Kyle Genova, Xiaoqi Yin, Alireza Fathi, Caroline Pantofaru, Leonidas Guibas, Andrea Tagliasacchi, Frank Dellaert, and Thomas Funkhouser. Panoptic Neural Fields: A Semantic Object-Aware Neural Scene Representation. In *CVPR*, 2022. 2

[11] Tianye Li, Mira Slavcheva, Michael Zollhoefer, Simon Green, Christoph Lassner, Changil Kim, Tanner Schmidt, Steven Lovegrove, Michael Goesele, Richard Newcombe, et al. Neural 3d video synthesis from multi-view video. In *CVPR*, 2022. 2

[12] Zhengqi Li, Simon Niklaus, Noah Snavely, and Oliver Wang. Neural scene flow fields for space-time view synthesis of dynamic scenes. In *CVPR*, 2021. 2

[13] Zhengqi Li, Qianqian Wang, Forrester Cole, Richard Tucker, and Noah Snavely. Dynibar: Neural dynamic image-based rendering. In *CVPR*, 2023. 2

[14] Jia-Wei Liu, Yan-Pei Cao, Tianyuan Yang, Eric Zhongcong Xu, Jussi Keppo, Ying Shan, Xiaohu Qie, and Mike Zheng Shou. Hosnerf: Dynamic human-object-scene neural radiance fields from a single video. In *ICCV*, 2023. 2

[15] Lingjie Liu, Marc Habermann, Viktor Rudnev, Kripasindhu Sarkar, Jiatao Gu, and Christian Theobalt. Neural actor: Neural free-view synthesis of human actors with pose control. *ACM Trans. Graph.(ACM SIGGRAPH Asia)*, 2021. 2

[16] Yu-Lun Liu, Chen Gao, Andreas Meuleman, Hung-Yu Tseng, Ayush Saraf, Changil Kim, Yung-Yu Chuang, Johannes Kopf, and Jia-Bin Huang. Robust dynamic radiance fields. In *CVPR*, 2023. 1, 2, 6, 7

[17] Qi Lu, Jason Kuen, Shen Tiancheng, Gu Jiuxiang, Guo Weidong, Jia Jiaya, Lin Zhe, and Yang Ming-Hsuan. High-quality entity segmentation. In *ICCV*, 2023. 1, 2, 3, 4, 8

[18] Qi Ma, Danda Pani Paudel, Ajad Chhatkuli, and Luc Van Gool. Deformable neural radiance fields using rgb and event cameras. In *ICCV*, 2023. 2

[19] B. Mildenhall, P. P. Srinivasan, M. Tancik, J. T. Barron, R. Ramamoorthi, and R. Ng. Nerf: Representing scenes as neural radiance fields for view synthesis. In *ECCV*, 2020. 1, 2, 3

[20] Ben Mildenhall, Dor Verbin, Pratul P. Srinivasan, Peter Hedman, Ricardo Martin-Brualla, and Jonathan T. Barron. MultiNeRF: A Code Release for Mip-NeRF 360, Ref-NeRF, and RawNeRF, 2022. 5

[21] Maxime Oquab, Timothée Darcet, Theo Moutakanni, Huy V. Vo, Marc Szafraniec, Vasil Khalidov, Pierre Fernandez, Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, Russell Howes, Po-Yao Huang, Hu Xu, Vasu Sharma, Shang-Wen Li, Wojciech Galuba, Mike Rabbat, Mido Assran, Nicolas Ballas, Gabriel Synnaeve, Ishan Misra, Herve Jegou, Julien Mairal, Patrick Labatut, Armand Joulin, and Piotr Bojanowski. Dinov2: Learning robust visual features without supervision. *arXiv:2304.07193*, 2023. 5

[22] Julian Ost, Fahim Mannan, Nils Thuerey, Julian Knodt, and Felix Heide. Neural scene graphs for dynamic scenes. In *CVPR*, 2021. 2

[23] Keunhong Park, Utkarsh Sinha, Jonathan T Barron, Sofien Bouaziz, Dan B Goldman, Steven M Seitz, and Ricardo Martin-Brualla. Nerfies: Deformable neural radiance fields. In *ICCV*, 2021. 2

[24] 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. *arXiv:2106.13228*, 2021. 1, 2

[25] Sida Peng, Junting Dong, Qianqian Wang, Shangzhan Zhang, Qing Shuai, Xiaowei Zhou, and Hujun Bao. Animatable neural radiance fields for modeling dynamic human bodies. In *ICCV*, 2021. 2

[26] Sida Peng, Yuanqing Zhang, Yinghao Xu, Qianqian Wang, Qing Shuai, Hujun Bao, and Xiaowei Zhou. Neural body: Implicit neural representations with structured latent codes for novel view synthesis of dynamic humans. In *CVPR*, 2021. 2

[27] Albert Pumarola, Enric Corona, Gerard Pons-Moll, and Francesc Moreno-Noguer. D-NeRF: Neural Radiance Fields for Dynamic Scenes. In *CVPR*, 2020. 1, 2

[28] Lu Qi, Jason Kuen, Zhe Lin, Jiuxiang Gu, Fengyun Rao, Dian Li, Weidong Guo, Zhen Wen, Ming-Hsuan Yang, and Jiaya Jia. Ca-ssl: Class-agnostic semi-supervised learning for detection and segmentation. In *ECCV*, 2022. 3

[29] Lu Qi, Jason Kuen, Yi Wang, Jiuxiang Gu, Hengshuang Zhao, Philip Torr, Zhe Lin, and Jiaya Jia. Open world entity segmentation. *TPAMI*, 2022. 3

[30] Amit Raj, Michael Zollhofer, Tomas Simon, Jason Saragih, Shunsuke Saito, James Hays, and Stephen Lombardi. Pixel-aligned volumetric avatars. In *CVPR*, 2021. 2

[31] Christian Reiser, Richard Szeliski, Dor Verbin, Pratul P. Srinivasan, Ben Mildenhall, Andreas Geiger, Jonathan T. Barron, and Peter Hedman. Merf: Memory-efficient radiance fields for real-time view synthesis in unbounded scenes. *SIGGRAPH*, 2023. 3

[32] Konstantinos Rematas, Andrew Liu, Pratul P. Srinivasan, Jonathan T. Barron, Andrea Tagliasacchi, Tom Funkhouser, and Vittorio Ferrari. Urban radiance fields. In *CVPR*, 2022. 2, 3

[33] S. Sabour, S. Vora, D. Duckworth, D. J. Fleet I. Krasin, and A. Tagliasacchi. Robustnerf: Ignoring distractors with robust losses. In *CVPR*, 2023. 1, 2, 3, 4, 5, 6, 7, 8

[34] 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 *ICCV*, 2019. 2

[35] Schönberger, Johannes Lutz, Frahm, and Jan-Michael. Structure-from-Motion Revisited. In *CVPR*, 2016. 1

[36] Tiancheng Shen, Yuechen Zhang, Lu Qi, Jason Kuen, Xingyu Xie, Jianlong Wu, Zhe Lin, and Jiaya Jia. High quality segmentation for ultra high-resolution images. In *CVPR*, 2022. 3

[37] Shih-Yang Su, Frank Yu, Michael Zollhöfer, and Helge Rhodin. A-nerf: Articulated neural radiance fields for learning human shape, appearance, and pose. In *NeurIPS*, 2021. 2

[38] Naoki Sugimoto, Yoshihito Ebine, and Kiyoharu Aizawa. Building movie map - a tool for exploring areas in a city - and its evaluations. In *ACM International Conference on Multimedia*, 2020. 6, 1

[39] Matthew Tancik, Vincent Casser, Xinchen Yan, Sabeek Pradhan, Ben Mildenhall, Pratul P. Srinivasan, Jonathan T. Barron, and Henrik Kretzschmar. Block-nerf: Scalable large scene neural view synthesis. In *CVPR*, 2022. 2, 3, 4

[40] Matthew Tancik, Ethan Weber, Evonne Ng, Ruilong Li, Brent Yi, Justin Kerr, Terrance Wang, Alexander Kristoffersen, Jake Austin, Kamyar Salahi, Abhik Ahuja, David McAllister, and Angjoo Kanazawa. Nerfstudio: A modular framework for neural radiance field development. *arXiv:2302.04264*, 2023. 3, 5, 6, 7

[41] Edgar Tretschk, Ayush Tewari, Vladislav Golyanik, Michael Zollhöfer, Christoph Lassner, and Christian Theobalt. Non-rigid neural radiance fields: Reconstruction and novel view synthesis of a dynamic scene from monocular video. In *ICCV*, 2021. 2

[42] Haithem Turki, Deva Ramanan, and Mahadev Satyanarayanan. Mega-nerf: Scalable construction of large-scale nerfs for virtual fly-throughs. In *CVPR*, 2022. 3

[43] Haithem Turki, Jason Y Zhang, Francesco Ferroni, and Deva Ramanan. Suds: Scalable urban dynamic scenes. In *CVPR*, 2023. 3

[44] Chaoyang Wang, Ben Eckart, Simon Lucey, and Orazio Gallo. Neural trajectory fields for dynamic novel view synthesis. *arXiv:2105.05994*, 2021. 2

[45] Chaoyang Wang, Lachlan Ewen MacDonald, László A. Jeni, and Simon Lucey. Flow supervision for deformable nerf. In *CVPR*, 2023. 2

[46] Kangkan Wang, Guofeng Zhang, Suxu Cong, and Jian Yang. Clothed human performance capture with a double-layer neural radiance fields. In *CVPR*, 2023. 2

[47] Liao Wang, Jiakai Zhang, Xinhang Liu, Fuqiang Zhao, Yanshun Zhang, Yingliang Zhang, Minye Wu, Jingyi Yu, and Lan Xu. Fourier plenoctrees for dynamic radiance field rendering in real-time. In *CVPR*, 2022. 2

[48] Liao Wang, Qiang Hu, Qihan He, Ziyu Wang, Jingyi Yu, Tinne Tuytelaars, Lan Xu, and Minye Wu. Neural residual radiance fields for streamably free-viewpoint videos. In *CVPR*, 2023. 2

[49] Peng Wang, Yuan Liu, Zhaoxi Chen, Lingjie Liu, Ziwei Liu, Taku Komura, Christian Theobalt, and Wenping Wang. F2-nerf: Fast neural radiance field training with free camera trajectories. In *CVPR*, 2023. 3

[50] Chung-Yi Weng, Brian Curless, Pratul P Srinivasan, Jonathan T Barron, and Ira Kemelmacher-Shlizerman. Humannerf: Free-viewpoint rendering of moving people from monocular video. In *CVPR*, 2022. 2

[51] Tianhao Wu, Fangcheng Zhong, Andrea Tagliasacchi, Forrester Cole, and Cengiz Oztireli. $D^2$nerf: Self-supervised decoupling of dynamic and static objects from a monocular video. In *NeurIPS*, 2022. 2, 6, 7

[52] Wenqi Xian, Jia-Bin Huang, Johannes Kopf, and Changil Kim. Space-time neural irradiance fields for free-viewpoint video. In *CVPR*, 2021. 2

[53] Enze Xie, Wenhai Wang, Zhiding Yu, Anima Anandkumar, Jose M Alvarez, and Ping Luo. Segformer: Simple and efficient design for semantic segmentation with transformers. In *NeurIPS*, 2021. 2, 5

[54] Zhiwen Yan, Chen Li, and Gim Hee Lee. Nerf-ds: Neural radiance fields for dynamic specular objects. In *CVPR*, 2023. 2

[55] Jae Shin Yoon, Kihwan Kim, Orazio Gallo, Hyun Soo Park, and Jan Kautz. Novel view synthesis of dynamic scenes with globally coherent depths from a monocular camera. In *CVPR*, 2020. 2

[56] Heng Yu, Joel Julin, Zoltan A Milacski, Koichiro Niinuma, and Laszlo A Jeni. Dylin: Making light field networks dynamic. In *CVPR*, 2023. 2

[57] Zhengming Yu, Wei Cheng, xian Liu, Wayne Wu, and Kwan-Yee Lin. MonoHuman: Animatable human neural field from monocular video. In *CVPR*, 2023. 2

[58] Jiakai Zhang, Xinhang Liu, Xinyi Ye, Fuqiang Zhao, Yanshun Zhang, Minye Wu, Yingliang Zhang, Lan Xu, and Jingyi Yu. Editable free-viewpoint video using a layered neural representation. *ACM Transactions on Graphics (TOG)*, 40(4):1–18, 2021. 2

[59] Kai Zhang, Gernot Riegler, Noah Snavely, and Vladlen Koltun. Nerf++: Analyzing and improving neural radiance fields. *arXiv:2010.07492*, 2020. 3

[60] Bolei Zhou, Hang Zhao, Xavier Puig, Sanja Fidler, Adela Barriuso, and Antonio Torralba. Scene parsing through ade20k dataset. In *CVPR*, 2017. 2, 5

# Entity-NeRF: Detecting and Removing Moving Entities in Urban Scenes

## Supplementary Material

## 7. MovieMap Dataset Details

To evaluate our approach, we introduce three urban scenes from the MovieMap [38]. The MovieMap Dataset was created by sampling images from 360° videos of varying lengths. Specifically, we extracted 51 images from a 3-second video, 12 images from a 6-second video, and 15 images from a 7-second video. Each image has a resolution of 3840×1920. For each 360° image, we extracted 14 perspective projection images. An example of this extraction process from a single 360° image is illustrated in Fig. 12. A common challenge when capturing 360° images is the inclusion of the photographer in the frame. To address this, we created a mask to exclude the photographer from the images, which is demonstrated in Fig. 13. This masked area was subsequently omitted from both the training and evaluation phases.

## 8. Additional Implementation Details

### 8.1. Rendered Background-only Images of MovieMap Dataset

When training static Neural Radiance Fields (NeRF) by removing masked objects, a significant challenge arises from errors near the edges of segmentation masks. These errors can disrupt the model's ability to accurately render static-only scenes, as they introduce inconsistencies at the boundaries of masked moving objects. To mitigate this, we dilated the mask area of moving objects. We implemented this by applying a convolution operation with a uniformly positive $3 \times 3$ kernel. Subsequently, in the output of this convolution, all positive values were converted to 1.

### 8.2. Robust Approaches for Entity Segmentation Errors

To prevent errors near the edges of entity segmentation, the area where the predicted entity-wise loss weights cover moving objects is increased through dilation. This is achieved by performing convolution with a uniform positive-valued $3 \times 3$ kernel and setting any positive values obtained in the result to 1.

Entity segmentation does not assign an entity to every pixel; some pixels are not assigned to any entity. Especially near the edges of objects, there are often pixels that were not assigned to any entity. We choose to include in training all pixels that are not classified as entities. However, we expect that the weight mask dilation process will exclude pixels near the edges of moving objects, which are not assigned to any entities, from the training process.


360° image


Perspective projection images

Figure 12. **Perspective projection images extracted from a single 360° image.**

## 9. More Results

### 9.1. Evaluation on RobustNeRF Dataset

**Dataset details**: Four natural scenes (i.e., Statue, Android, Crab, BabyYoda) from RobustNeRF [33]. Distractor objects are either moved or allowed to move between frames to simulate capture over extended periods. The number of unique distractors varies from 1 (Statue) to 150 (BabyYoda). Additional frames without distractors are provided to enable quantitative evaluation.

Note that we encountered issues with the provided camera parameters for the Statue and Android scenes, and the Crab scene does not provide camera parameters. Consequently, we calibrated the cameras using COLMAP [35]

1

| Loss | Statue | | | Android | | | Crab | | | BabyYoda | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | PSNR↑ | SSIM↑ | LPIPS↓ | PSNR↑ | SSIM↑ | LPIPS↓ | PSNR↑ | SSIM↑ | LPIPS↓ | PSNR↑ | SSIM↑ | LPIPS↓ |
| Mean-squared error (MSE) | 18.89 | 0.70 | 0.24 | 18.53 | 0.63 | 0.25 | 24.68 | 0.80 | 0.11 | 22.54 | **0.73** | **0.28** |
| RobustNeRF [33] | 21.14 | **0.74** | 0.19 | 19.47 | 0.65 | 0.21 | 30.32 | 0.83 | **0.10** | 25.16 | 0.69 | 0.33 |
| Entity-NeRF (only EARR) | 21.10 | **0.74** | **0.18** | 19.99 | **0.69** | **0.20** | 30.43 | 0.83 | **0.10** | 25.63 | 0.68 | 0.33 |
| Entity-NeRF | **21.20** | 0.73 | 0.19 | **20.23** | 0.67 | 0.21 | **30.65** | **0.84** | 0.11 | **25.65** | 0.68 | 0.33 |

Table 2. **Quantitative comparison with RobustNeRF [33] using Mip-NeRF 360 [2] on RobustNeRF Dataset.**



Figure 13. **Masks for photographers.**



Figure 14. **Visualization of $D(\mathbf{r})$ using stationary entity classification.** Compared to EARR, $D(\mathbf{r})$ in the early stages of training are improved.

| | IoU $D(\mathbf{r})=1$ ↑ | IoU $D(\mathbf{r})=0$ ↑ |
|---|---|---|
| RobustNeRF [33] | 0.84 | 0.14 |
| Entity-NeRF | **0.98** | **0.59** |

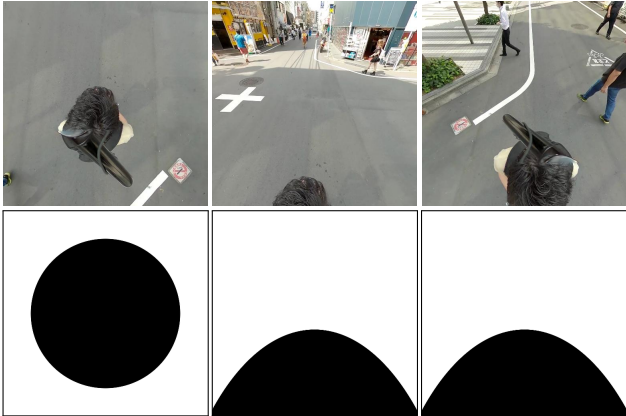Table 3. **Quantitative Comparison of Distractiveness with RobustNeRF [33] on MovieMap Dataset**

for these scenes and used the calibrated parameters for training in the three scenes (Statue, Android, and Crab). The BabyYoda scene was trained using the original camera parameters.

**Quantitative comparison**: A quantitative evaluation using Mip-NeRF 360 [2] on the RobustNeRF natural scenes (Statue, Android, Crab, and BabyYoda), which were shot with objects centered, is shown in Table 2. Although our proposed method is not intended to improve the performance of scenes shot with the object centered, it showed that the proposed method outperformed RobustNeRF in PSNR, and was equal or better in terms of SSIM and LPIPS. Even when a moving object is photographed at a large size, the same problem as in the urban scene may occur because the object appears at the edge of the patch, and EARR appeared to have solved this problem. In addition, the incorporation of knowledge by the stationary entity classification was also found to be effective in the indoor scenes.

## 9.2. Qualitative Comparison of Distractiveness using stationary entity classification

As shown in Fig. 14, our thing/stuff segmentation-based stationary entity classification provides more precise $D(\mathbf{r})$ assignments for each entity than EARR in initial learning stages, where predicting accurate diffuse $D(\mathbf{r})$ for all entities is challenging due to large residuals.
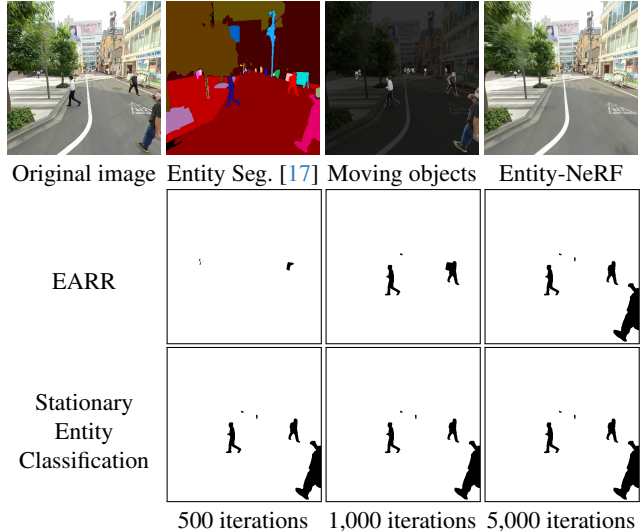
## 9.3. Quantitative Comparison of Distractiveness

In Table 3, the Intersection over Union (IoU) of Distractiveness $D(\mathbf{r})$ in Entity-NeRF at the end of training is compared with the IoU of Distractiveness $D(\mathbf{r})$ in RobustNeRF [33], using masks annotated on moving objects as ground-truth labels. Our proposed method achieves a better IoU for both $D(\mathbf{r}) = 0$ and $D(\mathbf{r}) = 1$, allowing for closer Distractiveness to the annotated mask to be given as a weight in the loss.

## 9.4. Novel View Synthesis

We conduct a qualitative comparison of our Entity-NeRF's performance in novel view synthesis. The novel view synthesis using the MovieMap Dataset is shown in Fig. 15. We created a circular trajectory around the straight-line path of

Reference image      RobustNeRF      Entity-NeRF

Figure 15. **Novel view synthesis.**

the original video and synthesized new views on the circular path. Entity-NeRF, although suffering from degradation due to the inability to learn correct geometry, shows less deterioration compared to RobustNeRF [33]. This is evident from the comparison of synthesized images from different viewpoints during training, as our approach avoids erroneously including moving objects in the training process and includes more static backgrounds into the training.