

Inverse Neural Rendering for Explainable Multi-Object Tracking

Julian Ost*

Tanushree Banerjee*

Mario Bijelic

Felix Heide

Princeton University

Abstract

Today, most methods for image understanding tasks rely on feed-forward neural networks. While this approach has allowed for empirical accuracy, efficiency, and task adaptation via fine-tuning, it also comes with fundamental disadvantages. Existing networks often struggle to generalize across different datasets, even on the same task. By design, these networks ultimately reason about high-dimensional scene features, which are challenging to analyze. This is true especially when attempting to predict 3D information based on 2D images. We propose to recast 3D multi-object tracking from RGB cameras as an Inverse Rendering (IR) problem, by optimizing via a differentiable rendering pipeline over the latent space of pre-trained 3D object representations and retrieve the latents that best represent object instances in a given input image. To this end, we optimize an image loss over generative latent spaces that inherently disentangle shape and appearance properties. We investigate not only an alternate take on tracking but our method also enables examining the generated objects, reasoning about failure situations, and resolving ambiguous cases. We validate the generalization and scaling capabilities of our method by learning the generative prior exclusively from synthetic data and assessing camera-based 3D tracking on the nuScenes and Waymo datasets. Both these datasets are completely unseen to our method and do not require fine-tuning. We present result videos at <https://light.princeton.edu/publication/inr-mot/>

1. Introduction

The most successful image understanding methods today employ feed-forward neural networks for performing vision tasks, including segmentation [11, 37, 41], object detection [17, 33, 39, 57–59, 91], object tracking [9, 30, 54, 61, 72, 85, 89] and pose estimation [69, 80]. Typically, these approaches learn network weights using large labeled

datasets. At inference time, the trained network layers sequentially process a given 2D image to make a prediction. Despite being a successful approach across disciplines, from robotics to health, and effective in operating at real-time rates, this approach also comes with several limitations: (i) Networks trained on data captured with a specific camera/geography *generalize poorly*, (ii) they typically rely on high-dimensional internal feature representations which are *often not interpretable*, making it hard to identify and reason about failure cases, and, (iii) it is challenging to enforce 3D geometrical constraints and priors during inference.

We focus on multi-object tracking as a task that must tackle all these challenges. Accurate multi-object tracking is essential for safe robotic planning. While approaches using LiDAR point clouds (and camera image input) are successful as a result of the explicitly measured depth [13, 30, 40, 54, 73, 81, 85], camera-based approaches to 3D multi-object tracking have only been studied recently [9, 18, 23, 44, 47, 53, 70, 76, 82, 89]. Monocular tracking methods, typically consisting of independent detection, 3D dynamic models, and matching modules, often struggle as the errors in the distinct modules tend to accumulate. Moreover, wrong poses in the detections can lead to ID switches in the matching process.

We propose an alternative approach that recasts visual inference problems as inverse rendering (IR) tasks, jointly solving them at test time by optimizing over the latent space of a generative object representation. Specifically, we combine object retrieval through the inversion of a rendering pipeline and a learned object model with a 3D object tracking pipeline. This approach allows us to simultaneously reason about an object’s 3D shape, appearance, and three-dimensional trajectory from monocular image input only. The location, pose, shape, and appearance parameters corresponding to the anchor objects are then iteratively refined via test-time optimization to minimize the distance between their corresponding generated objects and the given input image. Rather than directly predicting scene and object attributes, we optimize over a latent object representation to synthesize image regions that best explain the ob-

*equal contribution

served image. We match the inverse-rendered objects then be matched by comparing their optimized latents.

Our method hinges on an efficient rendering pipeline and generative object representation at its core. While the approach is not tied to a specific object representation, we adopt GET3D [16] as the generative object prior, that *is only trained on synthetic data* to synthesize textured meshes and corresponding images with an efficient differentiable rendering pipeline. Note that popular implicit shape/object representations do either not support class-specific priors [45, 55], or require expensive volume sampling [62].

The proposed method builds on the inductive geometry priors embedded in our rendering forward model, solving *different several tasks simultaneously*. Our method refines object pose as a byproduct, merely by learning to represent objects of a given class. Recovering object attributes as a result of inverse rendering also provides *interpretability “for free”*: once our proposed method detects an object at test time, it can extract the parameters of the corresponding representation alongside the rendered input view. This ability allows for reasoning about failure cases.

We validate that the method naturally exploits 3D geometry priors and *generalizes across unseen domains and unseen datasets*. After training solely on simulated data, we test on nuScenes [7] and Waymo [67] datasets, and although untrained, we find that our method *outperforms both existing dataset-agnostic multi-object tracking approaches and dataset-specific learned approaches* [89] when operating on the same detection inputs. In summary, we make the following contributions.

- We introduce an inverse rendering method for 3D-grounded monocular multi-object tracking. Instead of formulating tracking as a feed-forward prediction problem, we propose to solve an inverse image fitting problem optimizing over the latent embedding space of generative scene representations.
- We analyze the single-shot capabilities and the interpretability of our method using the generated image produced by our method during test-time optimization.
- Trained only on synthetic data, we validate the generalization capabilities of our method by evaluating on unseen automotive datasets, where the method compares favorably to existing methods when provided the same detection inputs.

1.0.1 Scope and Limitations

While facilitating inverse rendering, the iterative optimization in our method makes it slower than classical object-tracking methods based on feed-forward networks. We hope to address this limitation in the future by accelerating the forward and backward passes with adaptive level-

of-detail rendering techniques.

2. Related Work

Object Tracking is a challenging visual inference task that requires the detection and association of multiple objects. Specific challenges include highly dynamic scenes with partial or full occlusions, changes in appearance, and varying illumination conditions [66, 77, 84]. In this section, we first review classical tracking methods and deep detection and association methods. Following, we review 3D scene representations and inverse rendering.

3D Object Tracking. An extensively investigated line of work proposes tracking by detection, i.e., to solve the task by first detecting scene objects and then learning to find the associations between the detected objects over multiple frames [3, 5, 6, 8, 25, 74, 75]. In addition to association, 3D tracking requires the estimation of object pose. Since directly predicting 3D object pose is challenging [24], most existing 3D tracking methods rely on some explicit depth measurements in the form of Lidar point clouds [1, 15, 85], hybrid camera-lidar measurements [24] or stereo information[18, 51]. Weng *et al.* [72] proposed a generic tracking method that combines a 3D Kalman filter and the Hungarian algorithm for matching on an arbitrary object detector.

Only recent work [9, 23, 44, 76, 89] tackles monocular 3D tracking. Hu *et al.* [23] relies on similarity across different viewpoints to learn rich features for tracking. DEFT [9] jointly trains the feature extractor for detection and tracking using the features to match objects between frames. In contrast, Marinello *et al.* [44] use an off-the-shelf tracker and enhance image features with 3D motion and bounding box information. Zhou *et al.* [89] rely on a minimal input of two frames and predicted heatmaps to perform simultaneous detection and tracking. Some 3D tracking methods rely on motion models [12, 46, 60] such as the Kalman Filter [26]. Recent methods also make use of optical flow predictions [42], learned motion models metrics [82], long short-term memory modules (LSTM) [9, 23, 44] and more recently transformer modules [53, 70]. All the above methods rely on a feed-forward image encoder backbone to predict object features. Departing from this approach, we propose a multi-object tracking method that directly optimizes a consistent three-dimensional reconstruction of objects and 3D motion via an inverted graphics pipeline.

3D Scene Representations, Generation and Neural Rendering. A growing body of work addresses joint 3D reconstruction and detection from monocular cameras. Existing methods have exploited different geometrical priors [43] for this task, including meshes [2], points [34], wire frames [21], voxels [79] CAD models or implicit func-

tions [52] signed distance functions (SDFs) [87]. Early approaches in neural rendering represent the scene explicitly by, e.g., encoding texture or radiance on the estimated scene geometry [68] or using volumetric pixels (Voxels) [65]. Other methods represent 3D scenes *implicitly*. This includes the successful NeRF method [45] and variants that have been extended to dynamic scenes [52, 56, 86]. To allow the handling of semi-transparent objects, these representation models refrain from explicitly representing object surfaces. Signed distance fields represent surfaces of watertight objects as a zero level-set [14, 29, 55] modeling a Signed Distance Function (SDF). Adding textures to surface models allows for disentangling object shape from appearance [32, 78]. In recent years ideas from generative imaging models, such as GANs [27, 28], VAEs and diffusion models [22, 48] have been applied to the 3D domain [14, 16, 20, 62]. Generative models can either be used for pure generation [62] or provide prior knowledge for downstream tasks. Starting from a good prior can drastically improve the efficiency of inverse tasks, such as IR. While Gina3D [62] provides a prior on in-the-wild objects its volumetric rendering pipeline adds another layer of complexity sampling the full volume. We therefore rely on GET3D [16] generating a mesh as a prior object model and renders through rasterization, profiting from graphic pipelines optimized over decades.

Inverse Rendering. Inverse rendering methods conceptually “invert” the graphics rendering pipeline, which generates images from 3D scene descriptions, and instead estimates the 3D scene properties, i.e., geometry, lighting, depth, and object poses based on input images. Recent work [38, 71, 83] successfully achieved joint optimization of a volumetric model and unknown camera poses from a set of images merely by back-propagating through a rendering pipeline. Another area of inverse rendering focuses on material and lighting properties [19, 49, 50], to find a representation that best models the observed image.

To the best of our knowledge, we present the first method that employs an inverse rendering approach for multi-object 3D tracking, *without any feed-forward prediction* of object features – only given 2D image input.

3. Tracking by Inverse Rendering

We cast object tracking as a test-time inverse rendering problem that fits generated multi-object scenes to the observed image frames. First, we discuss the proposed scene representation we fit. Next, we devise our rendering-based test-time optimization at the heart of the proposed tracking approach. We employ an object-centric scene representation. We model the underlying 3D scene for a frame observation as a composition of all object instances without the background scene.

3.1. Scene Generation

Object Prior. To represent a large, diverse set of instances per class, we define each object instance o as a sample from a distribution \mathbf{O} over all objects in a class, that is

$$\mathbf{O} \sim f(\mathbf{o}), \quad (1)$$

where f is a learned function over a known prior object distribution. Here, the prior distribution is modeled by a differentiable generative 3D object model $o_p = G(z_p)$, that maps a latent embedding z_p to an object instance o_p , the object p . In particular, the latent space comprises two disentangled spaces \mathbf{z}_S and \mathbf{z}_T for shape S and texture T .

Given an object-centric camera projection $\mathbf{P}_c = \mathbf{K}_c \mathbf{T}$, where \mathbf{K}_c is the camera intrinsic matrix, a transformation $\mathbf{T} = [\mathbf{R}|\mathbf{t}]$ to camera c that is composed of rotation \mathbf{R} and translation \mathbf{t} , a differentiable rendering method $R(o_p, c)$, such as rasterization for meshes or volumetric rendering for neural fields, this renders an image $I_{c,p}$, a 2D observation of the 3D object o_p . While our method is general, implementation details of the generator and rendering method are provided in the implementation section.

Scene Composition. We model a multi-object scene as a scene graph composed of transformations in the edges and object instances in the leaf nodes, similar to Ost *et al.* [52]. Object poses are described by the homogeneous transformation matrix $\mathbf{T}_p \in \mathbb{R}^{4x4}$ with the translation \mathbf{t}_p and orientation \mathbf{R}_p in the reference coordinate system. The camera pose $\mathbf{T}_c \in \mathbb{R}^{4x4}$ is described in the same reference coordinate system. The relative transformation of the camera c and each object instance p can be computed through edge traversal in the scene graph as

$$\mathbf{T}_{c,p} = \text{diag}\left(\frac{1}{s_p}\right) \mathbf{T}_p \mathbf{T}_c^{-1}, \quad (2)$$

where the factor s_p is a scaling factor along all axes to allow a shared object representation of a unified scale. This canonical object scale is necessary to represent objects of various sizes, independent of the learned prior on shape and texture. The object centric projection $\mathbf{P}_{c,p} = \mathbf{K}_c \mathbf{T}_{c,p}$ is used to render the RGB image $I_{c,p} \in \mathcal{R}^{H \times W \times 3}$ and mask $M_{c,p} \in [0, 1]^{H \times W}$ for each object/camera pair.

Individual rendered RGB images are ordered by object distance $\|\mathbf{t}_{c,p}\|$, such that $p = 1$ is the shortest distance to c . Using the Hadamard product of the non-occluded mask γ_p all N_o object images are composed into a single image

$$\begin{aligned} \hat{I}_c &= \sum_{k=1}^{N_o} R(G(\mathbf{z}_{S,p}, \mathbf{z}_{T,p}), \mathbf{P}_{c,p}) \circ \gamma_p, \text{ where} \\ \gamma_p &= \max\left(\left(\mathbf{M}_{c,p} - \sum_{q=1}^p \mathbf{M}_{c,q}\right), \mathbf{0}^{H \times W}\right), \end{aligned} \quad (3)$$

where instance masks are generated in the same fashion.

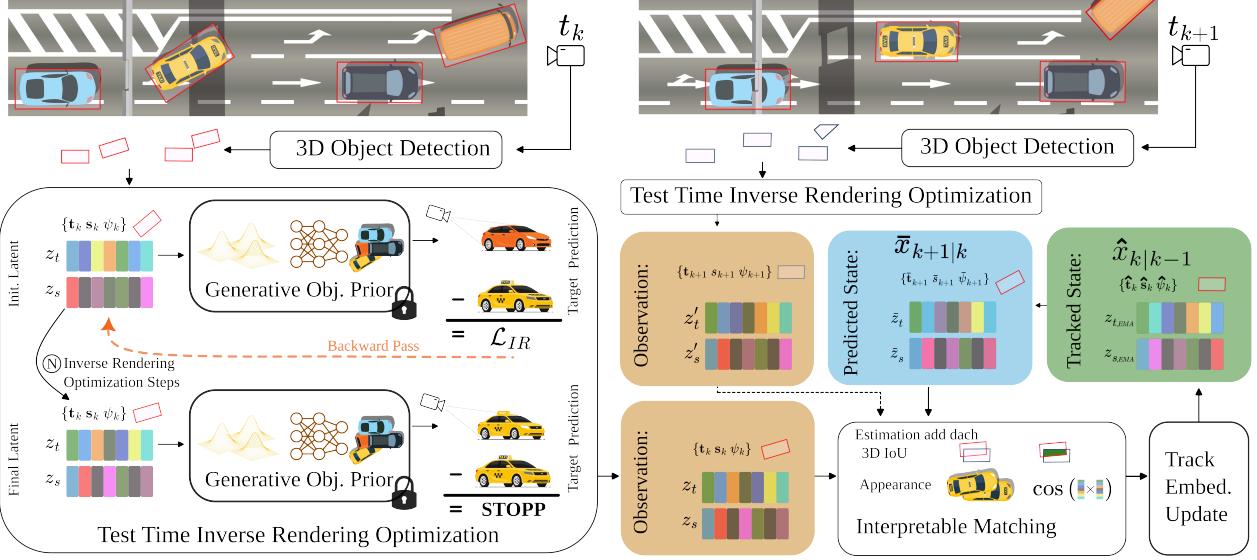


Figure 1. Inverse Rendering for Monocular Multi-Object Tracking. For each detection, we initialize the embedding codes of an object generator \mathbf{z}_S for shape and \mathbf{z}_T for texture. The generative object prior model is frozen and only embedding codes, pose, and size of each object instance are optimized through inverse rendering to best fit the image observation. Inverse-rendered texture and shape embeddings and refined object locations are provided to the matching stage to match predicted states of tracked objects of the past and new observations. Matched and new tracklets are updated, and unmatched tracklets are ultimately discarded before predicting states in the next step.

3.2. Inverse Multi-Object Scene Rendering.

We invert the differentiable rendering model defined in Eq. 3 by optimizing the set of all object representations in a given image I_c with gradient-based optimization. We assume that, initially, each object o_p is placed at a pose $\hat{\mathbf{T}}_{c,p}$ and scaled with \hat{s}_p near its underlying location. We represent object orientations in their respective Lie algebraic form $\mathfrak{so}(3)$. We further sample an object embedding $\hat{\mathbf{z}}_{S,p}$ and $\hat{\mathbf{z}}_{T,p}$ in the respective latent embedding space.

For in-the-wild images, I_c is not just composed of sampled object instances but other objects and the scene background. Since our goal for tracking is the reconstruction of all object instances of specific object classes, a naïve ℓ_2 image matching objective of the form $\|I_c - \hat{I}_c\|_2$ is noisy and challenging to solve with vanilla stochastic gradient descent methods. To tackle this issue, we optimize visual similarity in the generated object regions instead of the full image. We optimize only on rendered RGB pixels and minimize

$$\mathcal{L}_{RGB} = \| (I_c - \hat{I}_c) \circ \hat{M}_{I_c} \|_2, \text{ with } \hat{M}_{I_c} = \min \left(\sum M_{c,p}, 1 \right)$$

The mask of all foreground/object pixels \hat{M}_{I_c} is computed as the sum over all object masks $M_{c,p}$ in the frame rendered by camera c . We employ a learned perceptual similarity metric [88] (LPIPS) on object-centered image patches, that is

$$\mathcal{L}_{perceptual} = \text{LPIPS}_{patch} (I_c, \hat{I}_{c,p}). \quad (4)$$

The combined loss function of our method is

$$\mathcal{L}_{IR} = L_{RGB} + \lambda \mathcal{L}_{perceptual}, \quad (5)$$

which we optimize by refining the latent codes of shape and appearance, position, rotation, and scale, leading to

$$\hat{\mathbf{z}}_{S,p}, \hat{\mathbf{z}}_{T,p}, \hat{s}_p \hat{\mathbf{t}}_p, \hat{\mathbf{R}}_p = \arg \min (\mathcal{L}_{IR}). \quad (6)$$

Instead of using vanilla stochastic gradient descent methods, we propose an alternating optimization schedule of distinct properties that includes aligning \mathbf{z}_S before \mathbf{z}_T , to reduce the number of total optimization steps. A detailed implementation and validation of all design choices of the optimization are presented in the Supplementary Material.

Optimization. To solve Eq. 5, we propose an optimization schedule, that first optimizes a coarse color, and then jointly optimizes the shape and the positional state of each object. As the backbone of the learned perceptual loss, we use a pre-trained VGG16 [64] and utilize individual output feature map similarities at different points of the optimization. We find that color and other low-dimensional features are represented in the initial feature maps and those are better guidance for texture than high-dimensional features as outputs of the later blocks. These features have a more informative signal for shape and object pose. We use the average of the first and second blocks in the optimization for \mathbf{z}_T , while the combined perceptual similarity loss guides the optimization of \mathbf{z}_T and the pose.

We initialize all object embeddings with the same fixed values inside the embedding space, take two optimization steps solely on color utilizing the described loss, and then freeze the color for the joint optimization of the shape and



Figure 2. Tracking via Inverse Neural Rendering on nuScenes [7]. From left to right, we show (i) observed images from diverse scenes at timestep $k = 0$; (ii) an overlay of the optimized generated object and its 3D bounding boxes at timestep $k = 0, 1, 2$ and 3 . The color of the bounding boxes for each object corresponds to the predicted tracklet ID. We see that even in such diverse scenarios, our method does not lose any tracks and performs robustly across all scenarios, although the dataset is unseen.

pose. We regularize out-of-distribution generations with

$$\mathcal{L}_{embed} = \alpha_T \mathbf{z}_T + (1 - \alpha_T) \mathbf{z}_T^{avg} + \alpha_S \mathbf{z}_S + (1 - \alpha_S) \mathbf{z}_S^{avg} \quad (7)$$

that minimizes a weighted distance in each dimension with respect to the average embedding \mathbf{z}_S or \mathbf{z}_T respectively. For optimization, we use the ADAM optimizer [31]. The final loss function combines the RGB, perceptual cost 5 and the regularization 7 with $\lambda = 10$, $\alpha_T = 0.7$ and $\alpha_S = 0.7$. We freeze color after two steps of optimization and optimize the shape and scale for three more steps, adding translation and rotation only in the last two steps.

3.3. 3D Multi-Object Tracking via Inverse Rendering

Next, we describe the proposed method for tracking multiple dynamic objects with the inverse rendering approach from above. The approach tracks objects in the proposed representation across video frames and is illustrated in Fig. 1. For readability, we omit p and the split of \mathbf{z} into \mathbf{z}_S and \mathbf{z}_T in the following.

Initial Object and Pose Estimation. Common to tracking methods, we initialize with a given initial 3D detection on image $I_{c,k}$, and we set object location $\mathbf{t}_k = [x, y, z]_k$, scale $s_k = \max(w_k, h_k, l_k)$ using the detected bounding box dimensions and heading ψ_k in frame k . We then find an optimal representation z_k , and a refined location and rotation of each object o via the previously introduced inverse rendering pipeline for multi-object scenes. The resulting location,

rotation, and scale lead to the observation vector

$$\mathbf{y}_k = [\mathbf{t}_k, s_k, \psi_k]. \quad (8)$$

Prediction. While not confined to a specific dynamics model, we use a linear state-transition model \mathbf{A} , for the objects state $\mathbf{x}_k = [x, y, z, s, \psi, w, h, l, x', y', z']_k$, and a forward prediction using a Kalman Filter [26], a vanilla approach in 3D object tracking [72]. An instantiated object in $k-1$ can be predicted in frame k as

$$\hat{\mathbf{x}}_{k|k-1} = \mathbf{A}\hat{\mathbf{x}}_{k-1|k-1} \text{ and } \mathbf{P}_{k|k-1} = \mathbf{A}\mathbf{P}_{k-1|k-1}\mathbf{A}^T + \mathbf{Q}, \quad (9)$$

with the predicted *a priori* covariance matrix modeling the uncertainty in the predicted state.

Interpretable Latent Matching. In the matching stage, all optimal object representations o_p in frame k are matched with *tracked* and *lost* objects from $k-1$. Objects are matched based on appearance and location with a weighted affinity score

$$A = w_{IoU} A_{IoU,3D} + w_z A_z + w_c D_{centroid}, \quad (10)$$

where $A_{IoU,3D}$ is the IoU of the 3D boxes computed over the predictions of tracked object predictions $\mathbf{x}_{k|k-1}$ and refined observations. Here, the object affinity A_z is computed as the cosine distance of tracked object latent embeddings \mathbf{z} . In addition to that the Euclidean distance between the center $D_{centroid}$ adds additional guidance. We add no score for unreasonable distant tracked objects and detections.

We compute the best combination of tracked and detected objects using the Hungarian algorithm [35], again



Figure 3. Without changing the model or training on the dataset, our proposed method can generalize well to the Waymo Open Driving Dataset [67]. Similar to Fig 2, from left to right, we show (i) observed images from diverse scenes from the dataset at timestep $k = 0$; (ii) an overlay of the closest generated object and predicted 3D bounding boxes at timestep $k = 0, 1, 2$ and 3 . The color of the bounding boxes for each object corresponds to the predicted tracklet ID. Our method does not lose any tracks even on a different unseen dataset in diverse scenes, validating that the approach generalizes.

a conventional choice in existing tracking algorithms. Matched tracklet and object pairs are kept in the set of *tracked* objects and the representation of the corresponding detections is discarded, while unmatched detections are added as new objects. Unmatched tracklets are set to *lost* with a lost frame counter of one. Objects that were not detected in previous frames are set to *tracked* and their counter is reset to 0. Objects with a lost frame count higher than lifespan N_{life} , or outside of the visible field, are removed.

Track and Embedding Update. In the update step, we refine each object embedding \mathbf{z} and motion model \mathbf{y}_k given the result of the matching step. Embeddings are updated through an exponential moving average

$$\mathbf{z}_{k,EMA} = \beta \mathbf{z}_k + (1-\beta) \mathbf{z}_{k-1,EMA} \text{ with } \beta = \frac{2}{T-1} \quad (11)$$

over all past observations of the object, where T is the number of observed time steps of the respective instance. The observation \mathbf{y}_k is used to update the Kalman filter. The optimal Kalman gain

$$\mathbf{K}_k = \mathbf{P}_{k|k-1} \mathbf{H}^T (\mathbf{H} \mathbf{P}_{k|k-1} \mathbf{H}^T + R)^{-1} \quad (12)$$

is updated to minimize the residual error of the predicted model and the observation. The observation \mathbf{y}_k is used to estimate the object state as

$$\hat{\mathbf{x}}_{k|k} = \hat{\mathbf{x}}_{k|k-1} + \mathbf{K}_k (\mathbf{y}_k - \mathbf{H} \hat{\mathbf{x}}_{k|k-1}) \quad (13)$$

and with

$$\mathbf{P}_k = \mathbf{P}_{k|k-1} - \mathbf{K}_k \mathbf{H} \mathbf{P}_{k|k-1} \quad (14)$$

the *a posteriori* of the covariance matrix is updated.

3.4. Implementation Details

Representation Model. We employ the GET3D [16] architecture as object model G . Following StyleGAN [27, 28]

embeddings z_T and z_S are mapped to intermediate style embeddings w_S and w_T in a learned W -space, which we optimize over instead of Z -space. Style embeddings condition a generator function that produces tri-planes representing object shapes as Signed Distance Fields (SDFs) and textures as texture fields. We deliberately *train our generator on synthetic data only*, see experiments below. Differentiable marching tetrahedra previously introduced in DMTet [63] extract a mesh representation and Images are rendered with a differentiable rasterizer [36].

Computational Cost. Each IR optimization step in our implementation takes ~ 0.3 seconds per frame. The generation and gradient computation through the generator determines the computational cost of the method. However, we note that the rendering pipeline, contributing the majority of the computational cost of the generator, has not been performance-optimized and can be naively parallelized when implemented in lower-level GL+CUDA graphics primitives.

4. Experiments

In the following, we assess the proposed method. Having trained our generative scene model solely on simulated data [10], we test the generalization capabilities on the nuScenes [7] and Waymo [67] dataset – both are unseen by the method. We analyze generative outputs of the test-time optimization and compare them against existing 3D multi-object trackers [23, 70, 72, 82, 89] on camera-only data.

4.1. Single-Shot Object Retrieval and Matching

Although trained only on general object-centric synthetic data, ShapeNet [10], our method is capable of fitting a sample from the generative prior to observed objects in real datasets that match the vehicle type, color, and overall ap-

Training Data Unseen	Method	AMOTA \uparrow	AMOTP (m) \downarrow	Recall \uparrow	MOTA \uparrow	Modality
\times	PF-Track [53]	0.622	0.916	0.719	0.558	Camera
\times	QTrack [82]	0.692	0.753	0.760	0.596	Camera
\times	QD-3DT [23]	0.425	1.258	0.563	0.358	Camera
\checkmark	QD-3DT [23] (trained on WOD)	0.000	1.893	0.226	0.000	Camera
\times (CP)	CenterTrack [89]	0.202	1.195	0.313	0.134	Camera
\checkmark (CP)	AB3DMOT [72]	<u>0.387</u>	1.158	0.506	0.284	Camera
\checkmark (CP)	Inverse Neural Rendering (ours)	0.413	1.189	0.536	0.321	Camera

Table 1. **Quantitative Evaluation for Camera-only Multi-Object Tracking.** Quantitative results on “cars” in the test split of the nuScenes tracking dataset [7]. Our IR-based tracker outperforms AB3DMOT [72] on all metrics and CenterTrack [89] on accuracy. All three methods use the same detection backbone for fair comparison, while only CenterTrack requires end-to-end training on the dataset. Additional results show on-par performance of our method with QD-3DT [23] trained on nuScenes [7]. QD-3DT trained on the Waymo Open Dataset (WOD) does not generalize to nuScenes and does not achieve competitive results. Only very recent transformer-based methods, such as PF-Track [53] and the metric learning approach of Q-Track [82] achieve a higher score. However, these methods require end-to-end training on each dataset. “CP” denotes here the vision-only version of CenterPoint [90] was used for object detection. **Bold** denotes best and underlined second best for methods that did not train on the dataset or use the same detection backbone.

pearance closely, effectively making our method dataset-agnostic. We analyze the generations during optimization in the following.

Optimization. Given an image observation and coarse detections, our method aims to find the best 3D representation, including pose and appearance, solely through inverse rendering. In Fig. 4 we analyze this iterative optimization process, following a scheduled optimization as described in Sec. 3.4. We observe that the object’s color is inferred in only two steps. Further, we can observe that even though the initial pose is incorrect, rotation and translation are optimized jointly through inverse rendering together with the shape and scale of the objects, recovering from sub-optimal initial guesses. The shape representation close to the observed object is reconstructed in just 5 steps.

4.2. Evaluation

To provide a fair comparison of 3D multi-object tracking methods using monocular inputs, we compare against existing methods by running all our evaluations with the method reference code. We only evaluate methods, that consider past frames, but have no knowledge about future frames, which is a different task. While our method does not store the full history length of all images, we allow such memory techniques for other methods. We only consider purely mono-camera-based tracking methods. We note that, in contrast to our method, most *baseline methods we compare to are finetuned on the respective training set*. For all two-stage detect-and-track methods, we use CenterPoint [85] as the detection method. We compare to CenterTrack [89] as an established learning-based baseline, and present results of the very recent PFTrack [53], a transformer-based tracking method, Qtrack [82] as a metric learning method, and QD-3DT [23] as an LSTM-based state tracker combined

with image feature matching. Of all learning-based methods, only CenterTrack [89] allows us to evaluate tracking performance with identical detections. Finally, we compare to AB3DMOT [72] that builds on an arbitrary 3D detection algorithm and combines it with a modified Kalman filter to track the state of each object. AB3DMOT [72] and the proposed method are the only methods that are data-agnostic in the sense that they have not seen the training dataset. For a fair evaluation of these generalization capabilities in learning-based methods, we include another version of QD-3DT solely trained on the Waymo Open Dataset [67] and evaluate on nuScenes [7]. We discuss the findings in the following.

Validation on nuScenes. Tab. 1 reports quantitative results on the test split of the nuScenes tracking dataset [7] on the car object class for all six cameras. We list results for the multi-object tracking accuracy (MOTA) [4] metric, the AMOTA [72] metric, average multi-object tracking precision (AMOTP) [72] and recall of all methods. First, we evaluate a version of QD-3DT [23] that has been trained on the Waymo Open Dataset [67] (WOD) but tested on nuScenes. This experiment is reported in row four of Tab. 1 and confirms that recent end-to-end detection and tracking methods do not perform well on unseen data (see qualitative results in the Supplementary Material). Moreover, perhaps surprisingly, even when using the same vision-only detection backbone as in our approach, the established end-to-end trained baseline CenterTrack [89], which has seen the dataset, performs worse than our method. Our IR-based method outperforms the general tracker AB3DMOT [72]. When other methods are given access to the dataset, recent learning-based methods such as the end-to-end LSTM-based method QD-3DT [23] perform on par. Only the most recent transformer-based methods such as PF-Track [53]



Figure 4. Optimization Process. From left to right, we show (i) the observed image, (ii) the rendering predicted by the initial starting point latent embeddings, (iii) the predicted rendered objects after the texture code is optimized (iv) the predicted rendered objects after the translation, scale, and rotation are optimized, and (v) the predicted rendered objects after the shape latent code is optimized. The ground truth images are faded to show our rendered objects clearly. Our method is capable of refining the predicted texture, pose, and shape over several optimization steps, even if initialized with poses or appearance far from the target – all corrected through inverse rendering.

and the QTrack [82], which employ a quality-based association model on a large set of learned metrics, such as heat maps and depth, achieve higher scores. Note again, that these methods, in contrast to the proposed method, have been trained on this dataset and cannot be evaluated independently of their detector performance.

We visualize the rendered objects predicted by our tracking method in Fig. 2. We show an observed image from a single camera at time step $k = 0$, followed by rendered objects overlaid over the observed image at time step $k = 0, 1, 2$ and 3 along with their respective bounding boxes, with color-coded tracklets. We see that our method does not lose any tracks in challenging scenarios in diverse scenes shown here, from dense urban areas to suburban traffic crossings, and handles occlusions and clutter effectively. By visualizing the rendered objects as well as analyzing the loss values, our method allows us to reason about and explain success and failure cases effectively, enabling explainable 3D object tracking. The rendered output images provide interpretable inference results that explain successful or failed matching due to shadows, appearance, shape, or pose. For example, the blue car in the IR inference in Fig. 5 top row was incorrectly matched due to an appearance mismatch in a shadow region. A rendering model including ambient illumination may resolve this ambiguity, see further discussion in the Supplementary Material.

Interpretation. Fig 5 shows the inverse rendered scene graphs in isolation and birds-eye-view tracking outputs on a layout level. Our method accurately recovers the object pose, instance type, appearance, and scale. As such, our approach directly outputs a 3D model of the full scene, i.e., layout and object instances, along with the temporal history of the scene recovered through tracking – a rich scene representation that can be directly ingested by downstream planning and control tasks, or simulation methods to train downstream tasks. As such, the method also allows us to reason about the scene by leveraging the 3D information provided by our predicted 3D representations. The 3D locations, object orientations, and sizes recovered from such

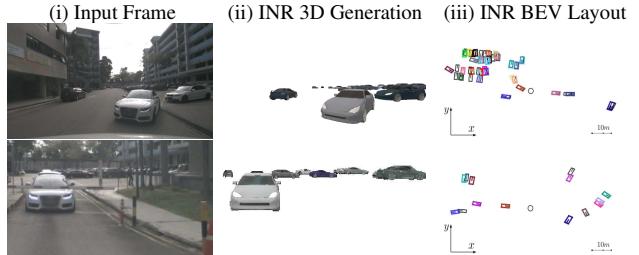


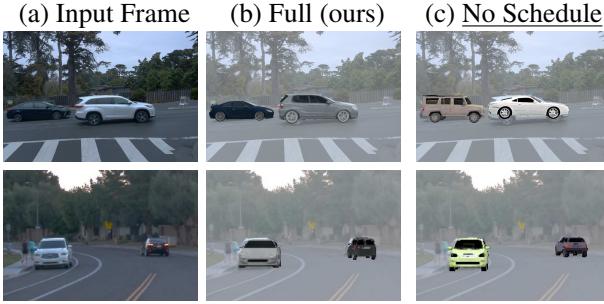
Figure 5. Layout Generation Through Inverse Rendering. From left to right, we show (i) observed image from a single camera for two scenes, (ii) test-time optimized inverse rendered (IR) objects of class “car”, and (iii) Bird’s Eye View (BEV) layout of the scene. In the BEV layout, black boxes represent ground truth, and the colored boxes represent predicted BEV boxes. The bottom shows a zoomed-in region at a 60 m distance (see BEV layout). Even in this setting, our method recovers the coarse appearance, shape of the objects, pose, and size,

visualizations can not only enable us to explain the predictions of our object tracking method, especially in the presence of occlusions or ID switches but also be used in other downstream tasks that require rich 3D understanding, such as planning.

Validation on Waymo. Next, we provide qualitative results from the 3D tracking on the validation set of WOD [67] in Fig. 3. The *only public results* on the provided test set are presented in QD-3DT [23], which may indicate it fails on this dataset. While the size of the dataset and its variety is of high interest for all autonomous driving tasks, Hu et al. [23] conclude that vision-only test set evaluation is not representative of a test set developed for surround view lidar data on partial unobserved camera images only. As such, we provide here qualitative results in Fig. 3, which validate that the method achieves tracking of similar quality on all datasets, providing a generalizing tracking approach. We show that our method does not lose tracks on Waymo scenes in diverse conditions.

Ablation Experiments As ablation experiments, we ana-

Table 2. **Ablation Experiments.**



Input frame is faded for visibility.

(a) **Effect of Optimization Schedule.** (a) observed image, (b) optimized generations using the proposed schedule in Sec. 3, (c) optimized generations using no schedule. This supports the quantitative to the left.

Method	AMOTA \uparrow	Recall \uparrow	MOTA \uparrow
\mathcal{L}_{IR} & \mathcal{L}_{embed} - Eq. 7	0.112	0.264	0.113
\mathcal{L}_{IR} - Eq. 5	0.103	0.236	0.112
$\mathcal{L}_{perceptual}$ - Eq. 4	0.100	0.251	0.101
\mathcal{L}_{RGB} - Eq. 4	N/A	N/A	N/A
No Schedule	0.102	0.224	0.110

(b) **Optimization Schedule and Loss Components.** Ablations were run on a small subset of the nuScenes [7] validation set. \mathcal{L}_{RGB} fails due to the optimizer fitting objects to the background instead, increasing the size of each object resulting in out of memory.

lyze the optimization schedule, the INR loss function components, and the weights of the tracker, applying them to a subset of scenes from the nuScenes validation set. We deliberately select this smaller validation set due to its increased difficulty.

The top row of Tab. 2b lists the quantitative results from our ablation study of the optimization scheduler. Our findings reveal a crucial insight: the strength of our method lies not in isolated loss components but in their synergistic integration. Specifically, the amalgamation of pixel-wise, perceptual, and embedding terms significantly enhances AMOTA, MOTA, and Recall metrics.

Moreover, the absence of an optimization schedule led to less robust matching as quantitative and qualitative results in Tab. 2a reveal. However, the core efficacy of our tracking method remained intact as indicated in the last row of Tab. 2b. This nuanced understanding underscores the importance of component interplay in our method.

5. Conclusion

We investigate inverse neural rendering as an alternative to existing feed-forward tracking methods. Specifically, we recast 3D multi-object tracking from RGB cameras as an inverse test-time optimization problem over the latent space of pre-trained 3D object representations that, when

rendered, best represent object instances in a given input image. We optimize an image loss over generative latent spaces that inherently disentangle shape and appearance properties. This approach to tracking also enables examining the reconstructed objects, reasoning about failure situations, and resolving ambiguous cases – rendering object layouts and loss function values provides interpretability “for free”. We validate that the method has high generalization capabilities, and without seeing a dataset, outperforms existing tracking methods’ generalization capabilities. In the future, we hope to investigate not only object detection with inverse rendering but broad, in-the-wild object class identification via conditional generation methods – unlocking analysis-by-synthesis in vision with generative neural rendering.

References

- [1] Álvarez-Aparicio, C., Guerrero-Higueras, Á.M., Rodríguez-Lera, F.J., Ginés Clavero, J., Martín Rico, F., Matellán, V.: People detection and tracking using lidar sensors. *Robotics* **8**(3), 75 (2019) 2
- [2] Beker, D., Kato, H., Morariu, M.A., Ando, T., Matsuoka, T., Kehl, W., Gaidon, A.: Monocular differentiable rendering for self-supervised 3d object detection. In: European Conference on Computer Vision. pp. 514–529. Springer (2020) 2
- [3] Bergmann, P., Meinhardt, T., Leal-Taixe, L.: Tracking without bells and whistles. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 941–951 (2019) 2
- [4] Bernardin, K., Elbs, A., Stiefelhagen, R.: Multiple object tracking performance metrics and evaluation in a smart room environment. In: Sixth IEEE International Workshop on Visual Surveillance, in conjunction with ECCV. vol. 90. Citeseer (2006) 7
- [5] Bewley, A., Ge, Z., Ott, L., Ramos, F., Upcroft, B.: Simple online and realtime tracking. In: 2016 IEEE international conference on image processing (ICIP). pp. 3464–3468. IEEE (2016) 2
- [6] Breitenstein, M.D., Reichlin, F., Leibe, B., Koller-Meier, E., Van Gool, L.: Robust tracking-by-detection using a detector confidence particle filter. In: 2009 IEEE 12th International Conference on Computer Vision. pp. 1515–1522. IEEE (2009) 2
- [7] Caesar, H., Bankiti, V., Lang, A.H., Vora, S., Lioung, V.E., Xu, Q., Krishnan, A., Pan, Y., Baldan, G., Beijbom, O.: nuscenes: A multimodal dataset for autonomous driving. In: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition. pp. 11621–11631 (2020) 2, 5, 6, 7, 9
- [8] Cao, J., Weng, X., Khirodkar, R., Pang, J., Kitani, K.: Observation-centric sort: Rethinking

- sort for robust multi-object tracking. arXiv preprint arXiv:2203.14360 (2022) 2
- [9] Chaabane, M., Zhang, P., Beveridge, J.R., O’Hara, S.: Deft: Detection embeddings for tracking. arXiv preprint arXiv:2102.02267 (2021) 1, 2
- [10] Chang, A.X., Funkhouser, T., Guibas, L., Hanrahan, P., Huang, Q., Li, Z., Savarese, S., Savva, M., Song, S., Su, H., Xiao, J., Yi, L., Yu, F.: ShapeNet: An Information-Rich 3D Model Repository. Tech. Rep. arXiv:1512.03012 [cs.GR], Stanford University — Princeton University — Toyota Technological Institute at Chicago (2015) 6
- [11] Chen, L.C., Papandreou, G., Kokkinos, I., Murphy, K., Yuille, A.L.: Semantic image segmentation with deep convolutional nets and fully connected crfs. arXiv preprint arXiv:1412.7062 (2014) 1
- [12] Chen, S.: Kalman filter for robot vision: a survey. IEEE Transactions on industrial electronics **59**(11), 4409–4420 (2011) 2
- [13] Chen, Y., Yu, Z., Chen, Y., Lan, S., Anandkumar, A., Jia, J., Alvarez, J.M.: Focalformer3d: Focusing on hard instance for 3d object detection (2023) 1
- [14] Chou, G., Chugunov, I., Heide, F.: Gensdf: Two-stage learning of generalizable signed distance functions. In: Proc. of Neural Information Processing Systems (NeurIPS) (2022) 3
- [15] Dewan, A., Caselitz, T., Tipaldi, G.D., Burgard, W.: Motion-based detection and tracking in 3d lidar scans. In: 2016 IEEE international conference on robotics and automation (ICRA). pp. 4508–4513. IEEE (2016) 2
- [16] Gao, J., Shen, T., Wang, Z., Chen, W., Yin, K., Li, D., Litany, O., Gojcic, Z., Fidler, S.: Get3d: A generative model of high quality 3d textured shapes learned from images. In: Advances In Neural Information Processing Systems (2022) 2, 3, 6
- [17] Girshick, R.: Fast r-cnn. In: Proceedings of the IEEE international conference on computer vision. pp. 1440–1448 (2015) 1
- [18] Gladkova, M., Korobov, N., Demmel, N., Ošep, A., Leal-Taixé, L., Cremers, D.: Directtracker: 3d multi-object tracking using direct image alignment and photometric bundle adjustment. arXiv preprint arXiv:2209.14965 (2022) 1, 2
- [19] Guo, Y.C., Kang, D., Bao, L., He, Y., Zhang, S.H.: Nerfren: Neural radiance fields with reflections. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 18409–18418 (2022) 3
- [20] Hao, Z., Mallya, A., Belongie, S., Liu, M.Y.: Gancraft: Unsupervised 3d neural rendering of minecraft worlds. arXiv preprint arXiv:2104.07659 (2021) 3
- [21] He, T., Soatto, S.: Mono3d++: Monocular 3d vehicle detection with two-scale 3d hypotheses and task priors. In: Proceedings of the AAAI Conference on Artificial Intelligence. vol. 33, pp. 8409–8416 (2019) 2
- [22] Ho, J., Jain, A., Abbeel, P.: Denoising diffusion probabilistic models. Advances in neural information processing systems **33**, 6840–6851 (2020) 3
- [23] Hu, H.N., Yang, Y.H., Fischer, T., Darrell, T., Yu, F., Sun, M.: Monocular quasi-dense 3d object tracking. IEEE Transactions on Pattern Analysis and Machine Intelligence (2022) 1, 2, 6, 7, 8
- [24] Huang, K., Hao, Q.: Joint multi-object detection and tracking with camera-lidar fusion for autonomous driving. In: 2021 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS). pp. 6983–6989. IEEE (2021) 2
- [25] Kalal, Z., Mikolajczyk, K., Matas, J.: Tracking-learning-detection. IEEE transactions on pattern analysis and machine intelligence **34**(7), 1409–1422 (2011) 2
- [26] Kalman, R.E.: A new approach to linear filtering and prediction problems (1960) 2, 5
- [27] Karras, T., Laine, S., Aila, T.: A style-based generator architecture for generative adversarial networks. In: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition. pp. 4401–4410 (2019) 3, 6
- [28] Karras, T., Laine, S., Aittala, M., Hellsten, J., Lehtinen, J., Aila, T.: Analyzing and improving the image quality of stylegan. In: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition. pp. 8110–8119 (2020) 3, 6
- [29] Kellnhofer, P., Jebe, L.C., Jones, A., Spicer, R., Pulli, K., Wetzstein, G.: Neural lumigraph rendering. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 4287–4297 (2021) 3
- [30] Kim, A., Osep, A., Leal-Taixé, L.: Eagermot: 3d multi-object tracking via sensor fusion. In: 2021 IEEE International Conference on Robotics and Automation (ICRA). pp. 11315–11321. IEEE (2021) 1
- [31] Kingma, D.P., Ba, J.: Adam: A method for stochastic optimization. CoRR **abs/1412.6980** (2015) 5
- [32] Koestler, L., Grittner, D., Moeller, M., Cremers, D., Lähner, Z.: Intrinsic neural fields: Learning functions on manifolds. arXiv preprint arXiv:2203.07967 2 (2022) 3
- [33] Ku, J., Mozifian, M., Lee, J., Harakeh, A., Waslander, S.L.: Joint 3d proposal generation and object detection from view aggregation. In: 2018 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS). pp. 1–8. IEEE (2018) 1

- [34] Ku, J., Pon, A.D., Waslander, S.L.: Monocular 3d object detection leveraging accurate proposals and shape reconstruction. In: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition. pp. 11867–11876 (2019) [2](#)
- [35] Kuhn, H.W.: The hungarian method for the assignment problem. Naval research logistics quarterly **2**(1-2), 83–97 (1955) [5](#)
- [36] Laine, S., Hellsten, J., Karras, T., Seol, Y., Lehtinen, J., Aila, T.: Modular primitives for high-performance differentiable rendering. ACM Transactions on Graphics (TOG) **39**(6), 1–14 (2020) [6](#)
- [37] Li, Y., Qi, H., Dai, J., Ji, X., Wei, Y.: Fully convolutional instance-aware semantic segmentation. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 2359–2367 (2017) [1](#)
- [38] Lin, C.H., Ma, W.C., Torralba, A., Lucey, S.: Barf: Bundle-adjusting neural radiance fields. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 5741–5751 (2021) [3](#)
- [39] Liu, W., Anguelov, D., Erhan, D., Szegedy, C., Reed, S., Fu, C.Y., Berg, A.C.: Ssd: Single shot multibox detector. In: European conference on computer vision. pp. 21–37. Springer (2016) [1](#)
- [40] Liu, Z., Tang, H., Amini, A., Yang, X., Mao, H., Rus, D., Han, S.: Bevfusion: Multi-task multi-sensor fusion with unified bird’s-eye view representation. arXiv (2022) [1](#)
- [41] Long, J., Shelhamer, E., Darrell, T.: Fully convolutional networks for semantic segmentation. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 3431–3440 (2015) [1](#)
- [42] Luiten, J., Fischer, T., Leibe, B.: Track to reconstruct and reconstruct to track. IEEE Robotics and Automation Letters **5**(2), 1803–1810 (2020) [2](#)
- [43] Mao, J., Shi, S., Wang, X., Li, H.: 3d object detection for autonomous driving: A review and new outlooks. arXiv preprint arXiv:2206.09474 (2022) [2](#)
- [44] Marinello, N., Proesmans, M., Van Gool, L.: Triplet-track: 3d object tracking using triplet embeddings and lstm. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 4500–4510 (2022) [1, 2](#)
- [45] Mildenhall, B., Srinivasan, P.P., Tancik, M., Barron, J.T., Ramamoorthi, R., Ng, R.: Nerf: Representing scenes as neural radiance fields for view synthesis (2020) [2, 3](#)
- [46] Nguyen, H.T., Smeulders, A.W.: Fast occluded object tracking by a robust appearance filter. IEEE transactions on pattern analysis and machine intelligence **26**(8), 1099–1104 (2004) [2](#)
- [47] Nguyen, P., Quach, K.G., Duong, C.N., Le, N., Nguyen, X.B., Luu, K.: Multi-camera multiple 3d ob-
- ject tracking on the move for autonomous vehicles. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 2569–2578 (2022) [1](#)
- [48] Nichol, A.Q., Dhariwal, P.: Improved denoising diffusion probabilistic models. In: International Conference on Machine Learning. pp. 8162–8171. PMLR (2021) [3](#)
- [49] Nimier-David, M., Dong, Z., Jakob, W., Kaplanyan, A.: Material and Lighting Reconstruction for Complex Indoor Scenes with Texture-space Differentiable Rendering. In: Bousseau, A., McGuire, M. (eds.) Eurographics Symposium on Rendering - DL-only Track. The Eurographics Association (2021). <https://doi.org/10.2312/sr.20211292> [3](#)
- [50] Nimier-David, M., Vicini, D., Zeltner, T., Jakob, W.: Mitsuba 2: A retargetable forward and inverse renderer. Transactions on Graphics (Proceedings of SIGGRAPH Asia) **38**(6) (Dec 2019). <https://doi.org/10.1145/3355089.3356498> [3](#)
- [51] Osep, A., Mehner, W., Mathias, M., Leibe, B.: Combined image-and world-space tracking in traffic scenes. In: 2017 IEEE International Conference on Robotics and Automation (ICRA). pp. 1988–1995. IEEE (2017) [2](#)
- [52] Ost, J., Mannan, F., Thuerey, N., Knodt, J., Heide, F.: Neural scene graphs for dynamic scenes. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). pp. 2856–2865 (2021) [3](#)
- [53] Pang, Z., Li, J., Tokmakov, P., Chen, D., Zagoruyko, S., Wang, Y.X.: Standing between past and future: Spatio-temporal modeling for multi-camera 3d multi-object tracking. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (2023) [1, 2, 7](#)
- [54] Pang, Z., Li, Z., Wang, N.: Simpletrack: Understanding and rethinking 3d multi-object tracking. arXiv preprint arXiv:2111.09621 (2021) [1](#)
- [55] Park, J.J., Florence, P., Straub, J., Newcombe, R., Lovegrove, S.: Deepsdf: Learning continuous signed distance functions for shape representation. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (June 2019) [2, 3](#)
- [56] Park, K., Sinha, U., Barron, J.T., Bouaziz, S., Goldman, D.B., Seitz, S.M., Martin-Brualla, R.: Nerfies: Deformable neural radiance fields. Proceedings of the IEEE International Conference on Computer Vision (2021) [3](#)
- [57] Qi, C.R., Liu, W., Wu, C., Su, H., Guibas, L.J.: Frustum pointnets for 3d object detection from rgb-d data. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 918–927 (2018) [1](#)

- [58] Redmon, J., Divvala, S., Girshick, R., Farhadi, A.: You only look once: Unified, real-time object detection. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 779–788 (2016)
- [59] Ren, S., He, K., Girshick, R., Sun, J.: Faster r-cnn: Towards real-time object detection with region proposal networks. Advances in neural information processing systems **28** (2015) 1
- [60] Scheidegger, S., Benjaminsson, J., Rosenberg, E., Krishnan, A., Granström, K.: Mono-camera 3d multi-object tracking using deep learning detections and pmbm filtering. In: 2018 IEEE Intelligent Vehicles Symposium (IV). pp. 433–440. IEEE (2018) 2
- [61] Sharma, S., Ansari, J.A., Murthy, J.K., Krishna, K.M.: Beyond pixels: Leveraging geometry and shape cues for online multi-object tracking. In: 2018 IEEE International Conference on Robotics and Automation (ICRA). pp. 3508–3515. IEEE (2018) 1
- [62] Shen, B., Yan, X., Qi, C.R., Najibi, M., Deng, B., Guibas, L., Zhou, Y., Anguelov, D.: Gina-3d: Learning to generate implicit neural assets in the wild. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). pp. 4913–4926 (June 2023) 2, 3
- [63] Shen, T., Gao, J., Yin, K., Liu, M.Y., Fidler, S.: Deep marching tetrahedra: a hybrid representation for high-resolution 3d shape synthesis. Advances in Neural Information Processing Systems **34**, 6087–6101 (2021) 6
- [64] Simonyan, K., Zisserman, A.: Very deep convolutional networks for large-scale image recognition (2015) 4
- [65] Sitzmann, V., Thies, J., Heide, F., Nießner, M., Wetzstein, G., Zollhöfer, M.: Deepvoxels: Learning persistent 3d feature embeddings. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) (2019) 3
- [66] Smeulders, A.W., Chu, D.M., Cucchiara, R., Calderara, S., Dehghan, A., Shah, M.: Visual tracking: An experimental survey. IEEE transactions on pattern analysis and machine intelligence **36**(7), 1442–1468 (2013) 2
- [67] Sun, P., Kretzschmar, H., Dotiwalla, X., Chouard, A., Patnaik, V., Tsui, P., Guo, J., Zhou, Y., Chai, Y., Caine, B., et al.: Scalability in perception for autonomous driving: Waymo open dataset. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). pp. 2446–2454 (2020) 2, 6, 7, 8
- [68] Thies, J., Zollhöfer, M., Nießner, M.: Deferred neural rendering. ACM Transactions on Graphics **38**(4), 1–12 (Jul 2019).
- <https://doi.org/10.1145/3306346.3323035>, <http://dx.doi.org/10.1145/3306346.3323035> 3
- [69] Wang, C., Xu, D., Zhu, Y., Martín-Martín, R., Lu, C., Fei-Fei, L., Savarese, S.: Densefusion: 6d object pose estimation by iterative dense fusion. In: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition. pp. 3343–3352 (2019) 1
- [70] Wang, S., Liu, Y., Wang, T., Li, Y., Zhang, X.: Exploring object-centric temporal modeling for efficient multi-view 3d object detection. arXiv preprint arXiv:2303.11926 (2023) 1, 2, 6
- [71] Wang, Z., Wu, S., Xie, W., Chen, M., Prisacariu, V.A.: Nerf: Neural radiance fields without known camera parameters. arXiv preprint arXiv:2102.07064 (2021) 3
- [72] Weng, X., Wang, J., Held, D., Kitani, K.: 3d multi-object tracking: A baseline and new evaluation metrics. In: 2020 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS). pp. 10359–10366. IEEE (2020) 1, 2, 5, 6, 7
- [73] Weng, X., Wang, Y., Man, Y., Kitani, K.M.: Gnn3dmot: Graph neural network for 3d multi-object tracking with 2d-3d multi-feature learning. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 6499–6508 (2020) 1
- [74] Wojke, N., Bewley, A.: Deep cosine metric learning for person re-identification. In: 2018 IEEE Winter Conference on Applications of Computer Vision (WACV). pp. 748–756. IEEE (2018). <https://doi.org/10.1109/WACV.2018.00087> 2
- [75] Wojke, N., Bewley, A., Paulus, D.: Simple online and realtime tracking with a deep association metric. In: 2017 IEEE international conference on image processing (ICIP). pp. 3645–3649. IEEE (2017) 2
- [76] Wu, J., Cao, J., Song, L., Wang, Y., Yang, M., Yuan, J.: Track to detect and segment: An online multi-object tracker. In: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition. pp. 12352–12361 (2021) 1, 2
- [77] Wu, Y., Lim, J., Yang, M.H.: Online object tracking: A benchmark. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 2411–2418 (2013) 2
- [78] Xiang, F., Xu, Z., Hasan, M., Hold-Geoffroy, Y., Sunkavalli, K., Su, H.: Neutex: Neural texture mapping for volumetric neural rendering. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 7119–7128 (2021) 3
- [79] Xiang, Y., Choi, W., Lin, Y., Savarese, S.: Data-driven 3d voxel patterns for object category recognition. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 1903–1911 (2015) 2

- [80] Xiang, Y., Schmidt, T., Narayanan, V., Fox, D.: Posecnn: A convolutional neural network for 6d object pose estimation in cluttered scenes. arXiv preprint arXiv:1711.00199 (2017) 1
- [81] Xuyang Bai, Zeyu Hu, X.Z.Q.H.Y.C.H.F., Tai, C.L.: TransFusion: Robust Lidar-Camera Fusion for 3d Object Detection with Transformers. CVPR (2022) 1
- [82] Yang, J., Yu, E., Li, Z., Li, X., Tao, W.: Quality matters: Embracing quality clues for robust 3d multi-object tracking. arXiv preprint arXiv:2208.10976 (2022) 1, 2, 6, 7, 8
- [83] Yen-Chen, L., Florence, P., Barron, J.T., Rodriguez, A., Isola, P., Lin, T.Y.: inferf: Inverting neural radiance fields for pose estimation. In: 2021 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS). pp. 1323–1330. IEEE (2021) 3
- [84] Yilmaz, A., Javed, O., Shah, M.: Object tracking: A survey. Acm computing surveys (CSUR) **38**(4), 13–es (2006) 2
- [85] Yin, T., Zhou, X., Krahenbuhl, P.: Center-based 3d object detection and tracking. In: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition. pp. 11784–11793 (2021) 1, 2, 7
- [86] Yuan, W., Lv, Z., Schmidt, T., Lovegrove, S.: Star: Self-supervised tracking and reconstruction of rigid objects in motion with neural rendering. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 13144–13152 (2021) 3
- [87] Zakharov, S., Kehl, W., Bhargava, A., Gaidon, A.: Autolabeling 3d objects with differentiable rendering of sdf shape priors. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 12224–12233 (2020) 3
- [88] Zhang, R., Isola, P., Efros, A.A., Shechtman, E., Wang, O.: The unreasonable effectiveness of deep features as a perceptual metric. In: Proceedings of the IEEE conference on computer vision and pattern recognition (CVPR). pp. 586–595 (2018) 4
- [89] Zhou, X., Koltun, V., Krähenbühl, P.: Tracking objects as points. In: European Conference on Computer Vision. pp. 474–490. Springer (2020) 1, 2, 6, 7
- [90] Zhou, X., Wang, D., Krähenbühl, P.: Objects as points. In: arXiv preprint arXiv:1904.07850 (2019) 7
- [91] Zhou, Y., Tuzel, O.: Voxelnet: End-to-end learning for point cloud based 3d object detection. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 4490–4499 (2018) 1

Inverse Neural Rendering for Explainable Multi-Object Tracking

(Supplementary Information)

Julian Ost* Tanushree Banerjee* Mario Bijelic Felix Heide

Princeton University

This supplementary document provides additional information and experiments supporting the main manuscript. We list training and architecture details, further ablation experiments, and additional comparisons and analysis.

Supplementary Video. In addition, we provide a supplementary video demonstrating the performance of our proposed INR-based tracking method on a sample of diverse scenes from the nuScenes [1] dataset and the Waymo Open Dataset [12]. We overlay the observed image with the rendered objects through alpha blending with a weight of 0.4. Object renderings are defined by the averaged latent embeddings $\mathbf{z}_{k,EMA}$ and the tracked object state \mathbf{y}_k .

Contents

1. Additional Implementation Details	ii
1.1. Loss Terms	ii
1.2. Optimization Schedule	iii
1.3. Tracking Heuristics	iii
1.4. Details on Generative Object Model	iv
2. Interpretability of Failure Cases	iv
2.1. Visualization	iv
2.2. Analysis	iv
3. Additional Results	vi
3.1. Matching Ablations	vi
3.2. Comparison to QD-3DT	vi
3.3. Additional Results on nuScenes	vii
3.4. Additional Results on Waymo Open Dataset	vii

*equal contribution

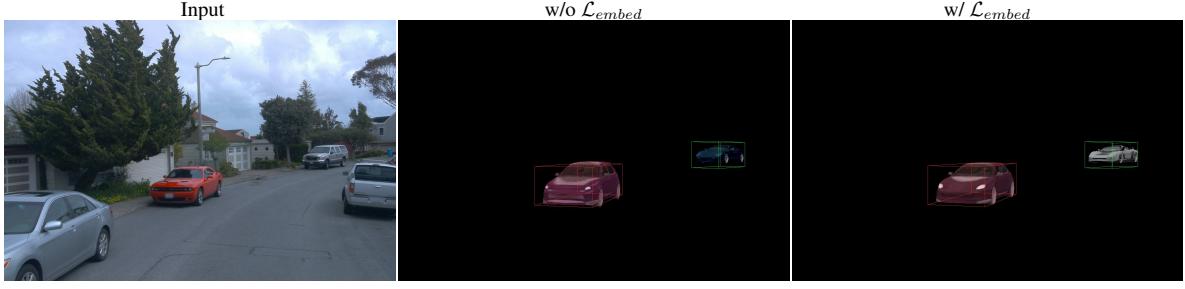


Figure 1. **Truncation Trick.** From left to right, (i) The input observed image, (ii) results without truncation regularizer applied and (iii) results with a truncation regularizer applied. For images (ii) and (iii), we also show bounding boxes for each object color-coded by the respective predicted object IDs. As shown here, applying the truncation regularizer helps us achieve more accurate textures and better shapes and colors for the predicted car surface by forcing the optimized embedding to be “well behaved”, i.e., close to the distribution of latent embeddings seen during the training by our representation model.

1. Additional Implementation Details

In the following, we describe the implementation of all design choices including the composition of the loss term, the proposed optimization schedule, heuristics applied in the matching stage of the multi-object tracker, and details on the generative object model.

1.1. Loss Terms

The test-time optimized inverse rendering of all objects in every scene and across datasets minimizes the loss term

$$\begin{aligned}
 \mathcal{L}_{IR+embed} &= \mathcal{L}_{RGB} + \lambda \mathcal{L}_{perceptual} + \mathcal{L}_{embed} \\
 &= \| (I_c - \hat{I}_c) \circ \hat{M}_{I_c} \|_2 \\
 &\quad + \lambda_1 \text{LPIPS}_{patch} (I_c, \hat{I}_{c,p}) \\
 &\quad + \lambda_2 (\alpha_S \mathbf{z}_S + (1 - \alpha_S) \mathbf{z}_S^{avg}) \\
 &\quad + \lambda_3 (\alpha_T \mathbf{z}_T + (1 - \alpha_T) \mathbf{z}_T^{avg})
 \end{aligned} \tag{1}$$

This combines RGB-MSE and a learned perceptual loss term with a regularization term, for which we describe implementation details as follows.

RGB Loss. The RGB loss is computed as the pixel-wise ℓ_2 -norm. Only pixels inside a mask are considered for which each pixel in the composed image at least one of the objects in the respective frame c also projects too. Only pixel values inside this mask \hat{M}_{I_c} are considered in the loss function. Please note, that in scenes with occluded objects only the object closest to the camera contributes to the rendered pixel in this non-volumetric rendering pipeline. We also assume this for the input images, given that considered objects are solid and mostly non-transparent.

Learned Perceptual Loss. The goal of the perceptual loss is to guide the inverse rendering to match the abstracted feature-level appearance of individual objects. We use the pre-trained LPIPS [14] loss with VGG16 [11] backbone for this, which operates on rectangular images with a minimum side length of 16 pixels. To consider objects individually we crop and resize patches from \hat{I}_c for each object p respectively. Mask $M_{c,p}(i, j)$ describes all pixels rendered from each object with their respective index i, j . $M_{c,p}(i, j) = 1$ if object p is projected into the respective pixel i, j from camera c . We then can describe an image patch by its top and bottom corner. The patch top corner is defined as $(u_{top}, v_{top}) = (i_{min, M_{c,p}}, j_{min, M_{c,p}})$, the upper and left corner of a tight axis-aligned rectangle around all rendered pixels. The patch bottom corner is defined as $(u_{bot}, v_{bot}) = (i_{max, M_{c,p}}, j_{max, M_{c,p}})$, the lower and right corner of the same axis-aligned rectangle.

Latent Embedding Regularization. Modern generative-adversarial (GAN) [3, 5–7, 9] methods, such as the used 3D object generator [3] first map an embedding sample from a multivariate Gaussian, called z-space or distribution, into a learned

Step \ Parameter (learning rate)	\mathbf{z}_S	\mathbf{z}_T	\mathbf{t}	Φ	s
1	-	3×10^{-1}	-	-	-
2	-	3×10^{-1}	-	-	-
3	6×10^{-2}	3×10^{-1}	3×10^{-2}	3×10^{-2}	1×10^{-6}
4	6×10^{-2}	-	-	-	-
5	6×10^{-2}	-	-	-	-
6	6×10^{-2}	-	-	-	-

Table 1. **Optimization Schedule.** Test time optimization of all object parameters, the shape and texture embeddings $\mathbf{z}_S, \mathbf{z}_T$, their location \mathbf{t} , rotation Φ in $\text{so}(3)$ and scale s follows this schedule. First, the texture is fitted in for two steps, followed by a pose adjustment in one step and inverse rendering of the shape, defined by the respective embedding code. Green denotes the optimization of the parameter in the respective step. The learning rates for all optimized parameters are noted in each field.

embedding space, called w-space, following a different distribution. The intuition behind this is that there are more optimal embedding distributions, which can be more easily mapped to the data distribution that the GAN is generating. High-quality samples are only generated from embeddings inside the high-dimensional embedding distribution. Therefore, we regularize the optimized embedding code through inverse rendering, with

$$\mathcal{L}_{\text{embed}} = \alpha_T \mathbf{z}_T + (1 - \alpha_T) \mathbf{z}_T^{\text{avg}} + \alpha_S \mathbf{z}_S + (1 - \alpha_S) \mathbf{z}_S^{\text{avg}}, \quad (2)$$

that is Eq.(14) in the main paper. Here, α_T and α_S are set to 0.7.

We employ the *truncation trick* widely used in GAN-based generators and first presented in StyleGAN [6], where z^{avg} represents an exponential-moving average of embedding codes generated from the Gaussian training during the training of the generator. This stabilizes the optimization through inverse rendering as Fig. 1 shows.

Weighting With empirical analysis of the validation set of both datasets, we find the weighting of loss terms $\lambda_1 = 0.4$, $\lambda_2 = 3$, and $\lambda_3 = 10$ for stable and truthfully generated objects via inverse rendering.

1.2. Optimization Schedule

For the loss function presented in Eq. (5) of the main paper, we found that the schedule presented in Tab. 1 solves this optimization problem effectively, while being stable across various scenes and datasets.

We first fit the texture embedding in only three steps during the test-time optimization of all object parameters. In step 3, we jointly solve for pose, scale, and shape, followed by three more steps on shape only. Details on the learning rate for the respective parameters are reported in Tab. 1.

An exhaustive search is impossible due to the number of hyper-parameters when including different loss functions and terms. We therefore performed empirical investigation on small, diverse subsets of scenes to find the parameter set used. The same setting works well on all datasets and *have not been changed* for the Waymo Open Dataset [12] and the nuScenes dataset [1].

1.3. Tracking Heuristics

The main paper describes the integration of test-time optimized object representations through inverse rendering into the tracking pipeline. Specifically, Eq. (6) in the main paper defines the following affinity

$$A = w_{IoU} A_{IoU} + w_z A_z + w_c D_{\text{centroid}}, \quad (3)$$

that describes the similarity of tracked and detected objects in the matching step. We follow [1, 13] and assign 0 to all object pairs whose center is more than 10m apart. In all other cases, we apply the weights $w_{IoU} = 0.7$, $w_z = 0.4$, and $w_c = 0.5$ between different affinity parameters. Here, A_{IoU} is the true 3D IoU of the bounding boxes, which is computed with the efficient PyTorch3D [8] implementation. The affinity A_z is computed as the pair-wise cosine distance between all tracklet-detection pairs. Centroid distances between pairs are computed in addition to the IoU to compensate for larger errors in line with the camera axis, common in vision-based object detectors. In such cases, IoU might be low, but object distances are still in a reasonable range which we empirically found at 5m for the object detector used. Finally, we define the distance-based affinity score as

$$D_{\text{centroid}} = \text{maximum} \left(\frac{-\|\mathbf{t}_{\text{tracklet}} - \mathbf{t}_{\text{detection}}\|_2}{d_{\text{max}}}, 0 \right), \text{ with } d_{\text{max}} = 5m. \quad (4)$$

Matches below the threshold of 0.48 for their affinity score are not considered matched and the respective tracklet is set to “lost” and the detection is added as a new tracklet in the next step. Tracklets are considered “dead” and removed after a maximum of 4 consecutive lost steps.

1.4. Details on Generative Object Model

The object generator, used as the prior for car representation, follows the architecture from GET3D [3]. The generator maps two sample codes z_{tex} and z_{geo} , drawn from a Gaussian distribution for the texture and shape respectively, to samples of a 3D SDF and a texture field. Details of this object model’s architecture, training, and usage are given below.

Generator Architecture. In this work, we employ a variant described in the appendix of [3]. Following [6, 7], two Gaussian variables are mapped independently to intermediate style embedding w_S and w_T in a learned \mathbf{W} -space. Style embeddings are then used as an input to the CNN-based StyleGAN2 [7] generator. Both style embeddings for shape and texture jointly condition the generator in each block, allowing texture and shape to influence the other modality. Each intermediate feature map of the generator backbone is mapped to its respective modality through a mapping layer only conditioned on the respective style embedding. All feature maps are accumulated and reshaped into three feature planes in the last generator layer. These planes represent textures as texture fields, object shapes as Signed Distance Fields (SDFs), and vertex offsets of a corresponding mesh. This forms a feature volume representation of the textured 3D object on two tri-planes.

Rendering. We employ the differentiable marching tetrahedra [10] method and extract a mesh representation from the SDF and vertex offsets, allowing more efficient rendering compared to sampling-based volumetric SDF rendering. Differentiable rasterization for meshes efficiently renders a 2D image of the respective mesh. By querying the texture field only at visible surface points, the respective vertex color can be efficiently retrieved to render the RGB image output.

Training. The model is trained using adversarial losses defined on the 2D renderings of 3D objects from the ShapeNet version 1 dataset [2]. Specifically, we use a differentiable rasterizer to render RGB images together with the silhouette masks of the objects as in [3] with a training configuration that largely follows StyleGAN2 [7] including using a minibatch standard deviation in the discriminator, exponential moving average for the generator, non-saturating logistic loss, and R1 regularization.

2. Interpretability of Failure Cases

2.1. Visualization

Being able to visualize the reconstructed objects allows us to reason about failure cases. For example, in scene **(e)** in Fig. 2, we see that the initial object significantly overlaps with the background asphalt in the shadow region, causing the reconstructed car to erroneously generate a darker gray color. Thus, not only does our method allow us to reason about failure cases, but it also identifies ways in which our representation model can be modified to rectify such failures. For example, a generative object model with an additional component that can model different lighting conditions to account for shadows might allow us to identify and reconstruct cars in varied lighting conditions, including shadows. This can guide future work for perception tasks through inverse rendering.

2.2. Analysis

Next, we analyze common failure cases using the pre-trained generator as an object prior in the presented tracking method. The visualizations allow us to assess the types of cases where this pipeline fails to track objects. Common failure cases we observed are listed below, with visualizations of such failure cases shown in Figure 2. We describe the cases corresponding to the rows (a-d, see figure labels) as follows:

- (a) The apparent darker color of the car due to **shadows** often causes the predicted object color to be darker than the color a human would perceive the car as. In the presented case on the right, the white car is completely occupied by the shadow of the neighboring truck. While the human visual system perceives the color of the car as white, the numerical RGB value in the image is closer to grey/black. This causes the predicted embedding corresponding to the texture of the object to represent the darker color. This might cause the tracking to fail due to incorrect matching of corresponding objects in

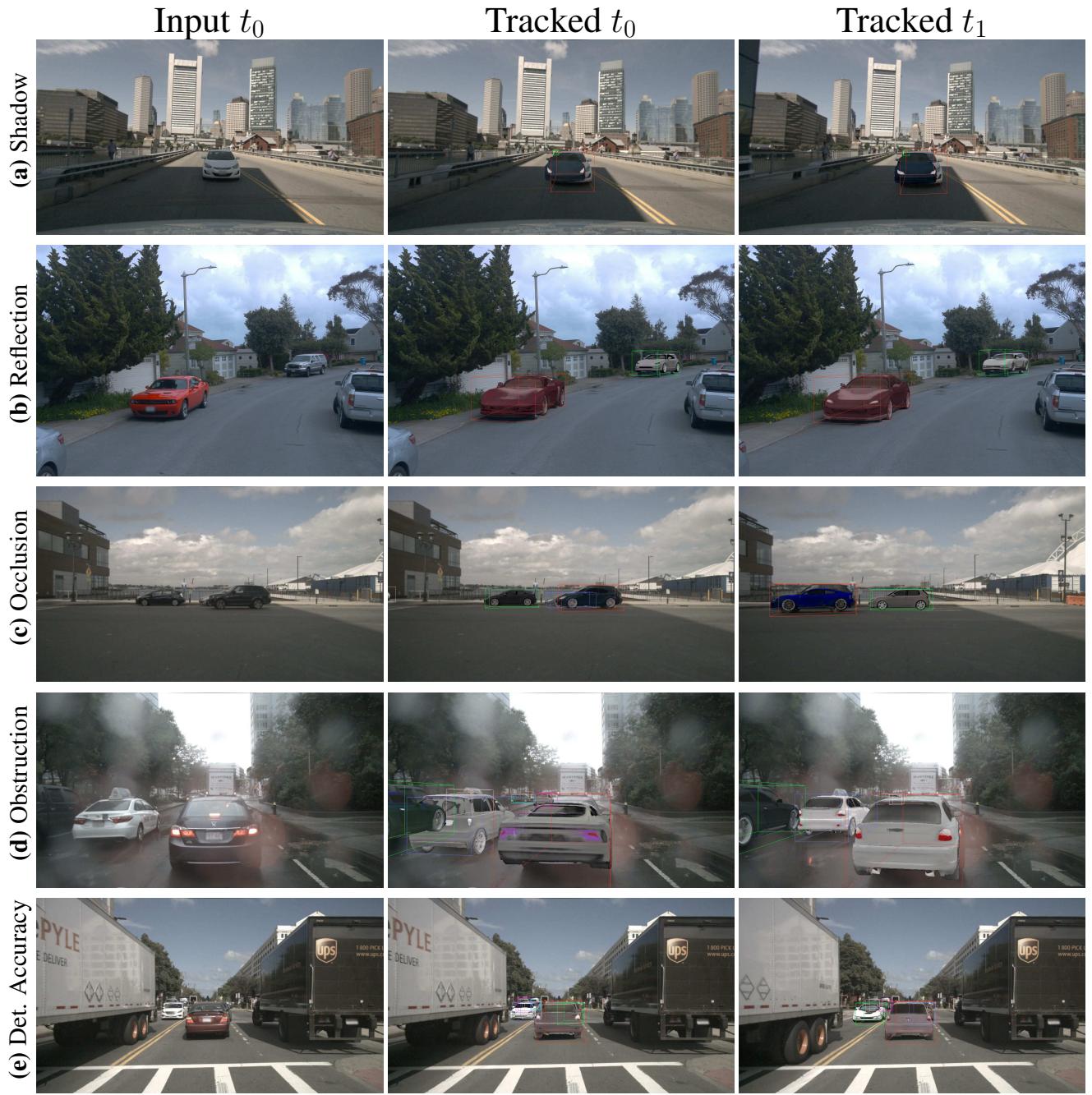


Figure 2. Examples of failure cases, such as lighting (shadows and reflections) or occluded objects, where the reconstructed object differs significantly from the observed object. These visualizations allow us to understand exactly why our model fails at reconstructing and tracking objects. This also allows us to identify ways the representation model and perception pipeline can be improved to incorporate effects that cause the method to fail.

consecutive frames with and without shadows.

- (b) Extreme **reflections** on the car due to the lighting conditions cause the model to try to model the RGB color of the reflection by erroneously modifying the texture of the generated object. Here, clouds in the sky are reflected as white spots on the hood and windshield of the red car, causing reflection and changes of these, influencing the generated texture as white spots. Future work that includes explicit models of BRDF would be beneficial in mitigating this class of failure

modes.

- (c) In addition to visualization and interpretation of the object prior more traditional aspects of the perception pipeline, such as ID-switches through **occlusion** in multi-object tracking scenarios can be observed. Here, the first object in the background (green box, t_0) is misidentified as the second object (green box, t_1). Such visualization allows further reasoning about the full pipeline.
- (d) Camera **obstructions** and extreme local **lighting**, such as in raindrops, lens flares, and bright lights, cause our method to erroneously predict the shape and texture of many cars to match the perceived shape of the car, causing matching to fail.
- (e) Inaccuracies in the predicted pose cause the predicted car in front to not perfectly align with the observed car patch, causing there to be an overlap between the predicted car and the immediate surroundings in the observed image, here the asphalt. Since only information on the detection is available, the color of the optimized object tends to be predicted gray (since the road is gray, and so the optimized texture embedding is closer to gray).

In general, our generative prior does not model specular textures and instead is restricted to diffuse reflectance. As such, it tends to reconstruct darker or lighter textures compensating for shadows from the environment and reflections of the sky. Modeling environmental lighting, complex material properties, and shadows may lead to a complex and less robust light simulation and will be restricted by the data available to train a generative prior model. Nevertheless, this is an exciting direction for future work.

3. Additional Results

In the following, we present additional tracking results on the real-world datasets on which we test our method (see the evaluation section in the main paper).

3.1. Matching Ablations

Adopting the same setting as AB3DMOT [13], we performed a hyper-parameter search for each matching weight and the detection confidence threshold, as denoted in Sec. 4. Our full method outperforms AB3DMOT, see Table 2, conducted on the validation split. Our best setting outperforms AB3DMOT, the only other method not trained on the dataset, by 3.9% AMOTA on the nuScenes test split.

3.2. Comparison to QD-3DT

The naive quantitative evaluation of multi-object tracking methods can easily be “unfair” in the sense that the tracker during training may be given access to future frames or rely on an improved detector backbone (making it challenging to evaluate the tracker in isolation). Evaluating generalization requires a nuanced setup to provide a fair evaluation. Therefore, we decide to focus the evaluation on methods that either build on the same detector backbone [15] and are not trained on the respective training set [13], achieving generalization by design, or end-to-end trained tracking methods for which we use a model trained on a different dataset. Especially we evaluate a model of QD-3DT [4] that has been trained on the Waymo Open Dataset [12] on nuscenes [1]. The authors of QD-3DT [4] were so kind to provide the respective checkpoints to us.

Method (split)	AMOTA \uparrow	Recall \uparrow	MOTA \uparrow
AB3DMOT (CP, test)	0.387	0.506	0.284
Best Hyper-param. (CP, test)	0.413	0.536	0.321
$w_{iou} = 1.4$ (val)	0.403	0.540	0.322
$w_{center} = 0.9$ (val)	0.417	0.514	0.332
$w_{embedd} = 0.4$ (val)	0.418	0.558	0.332
$\tau_{det} = 0.4$ (val)	0.397	0.567	0.326

Table 2. **Tracking Matching and Detection Confidence.** Parameters were optimized on the nuScenes [1] validation set. On the test split our best setting for w_{iou} , w_{center} , w_{embedd} , τ_{det} surpasses the performance of AB3DMOT [13], the only baseline not trained on the dataset.

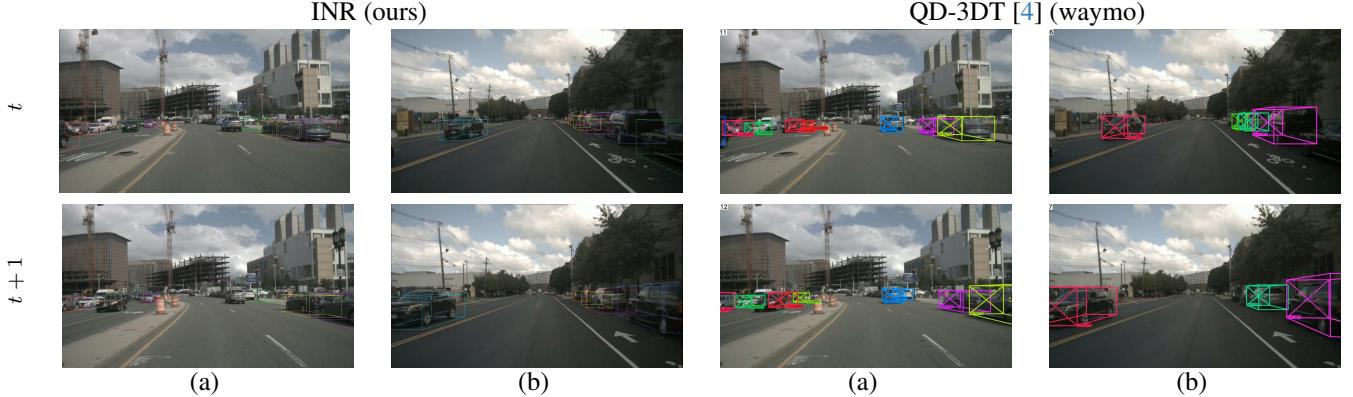


Figure 3. **Qualitative Comparison on Generalization.** We compare 3D bounding box outputs from QD-3DT [4] trained on waymo [12] and our inverse neural rendering based tracker (INR) overlayed on the respective input videos from nuscenes [1]. The same color in consecutive frames denote same tracklet. As we see in scene (a) on the left side QD-3DT is having ID switches especially in the far range and on the sides of a frame implying dataset specific performance caps, e.g. the training dataset has been tracking annotated in a rather short range. Another finding is exemplified in scene (b) on the right side were the tracker does not generalize well and is loosing a tracklet.

Fig. 3 shows tracking results from our method along with results from QD-3DT. A visual, qualitative inspection illustrates that the end-to-end trained tracking method QD-3DT [12] still performs well on objects in the center of the scene but does not generalize well on occluded or partially visible objects. This is reflected in the scores reported in Tab. 1 of the main manuscript.

3.3. Additional Results on nuScenes

Although the nuScenes [1] dataset consists of sensor data from 6 cameras, 5 radars, and 1 lidar, we tackle monocular camera-based 3D object tracking in this work. As such, we only use the data collected from the 6 cameras. The dataset comprises 1000 scenes, with each scene being 20s long. The test set contains 150 scenes. Each of these scenes is selected to be *interesting*, which include scenes with high traffic density (e.g., intersections, construction sites), rare classes (e.g. ambulances, animals), potentially dangerous traffic situations (e.g., jaywalkers, incorrect behavior), maneuvers (e.g., lane change, turning, stopping) and situations that may be difficult for an Autonomous Vehicle. Additional results of our method on the nuScenes dataset are listed in Fig. 4. We note that the colors of the cars are matched quite accurately. Moreover, the shapes of the cars get reconstructed as well. In turn, we can see that the tracking quality is high as visualized by bounding boxes. Note that the color of bounding boxes marks the same instance in consecutive frames.

3.4. Additional Results on Waymo Open Dataset

The Waymo Open Dataset [12] consists of 1150 scenes that each span 20s. Again, since we tackle monocular tracking, we only use the data from the 5 camera sensors. The dataset was collected by driving in Phoenix AZ, Mountain View CA, and San Francisco CA across daytime, nighttime, and dawn lighting conditions. Additional results of our method on the Waymo Open dataset are given in Figure 5. We find that our method generalizes effectively to this dataset. The colors of the cars are matched quite accurately, as shown in Figure 5. Moreover, the shapes of the cars get reconstructed as well. As such, again, tracking quality is high as visualized by bounding boxes. Note that the color of bounding boxes marks the same instance in consecutive frames.



Figure 4. Additional visualizations on nuScenes [1]. From left to right, we show (i) observed images from diverse scenes at timestep $k = 0$; (ii) an overlay of the optimized generated object and its 3D bounding boxes at timestep $k = 0, 1, 2$ and 3 . The color of the bounding boxes for each object corresponds to the predicted tracklet ID. We see that our method can accurately reconstruct objects in diverse scenarios.

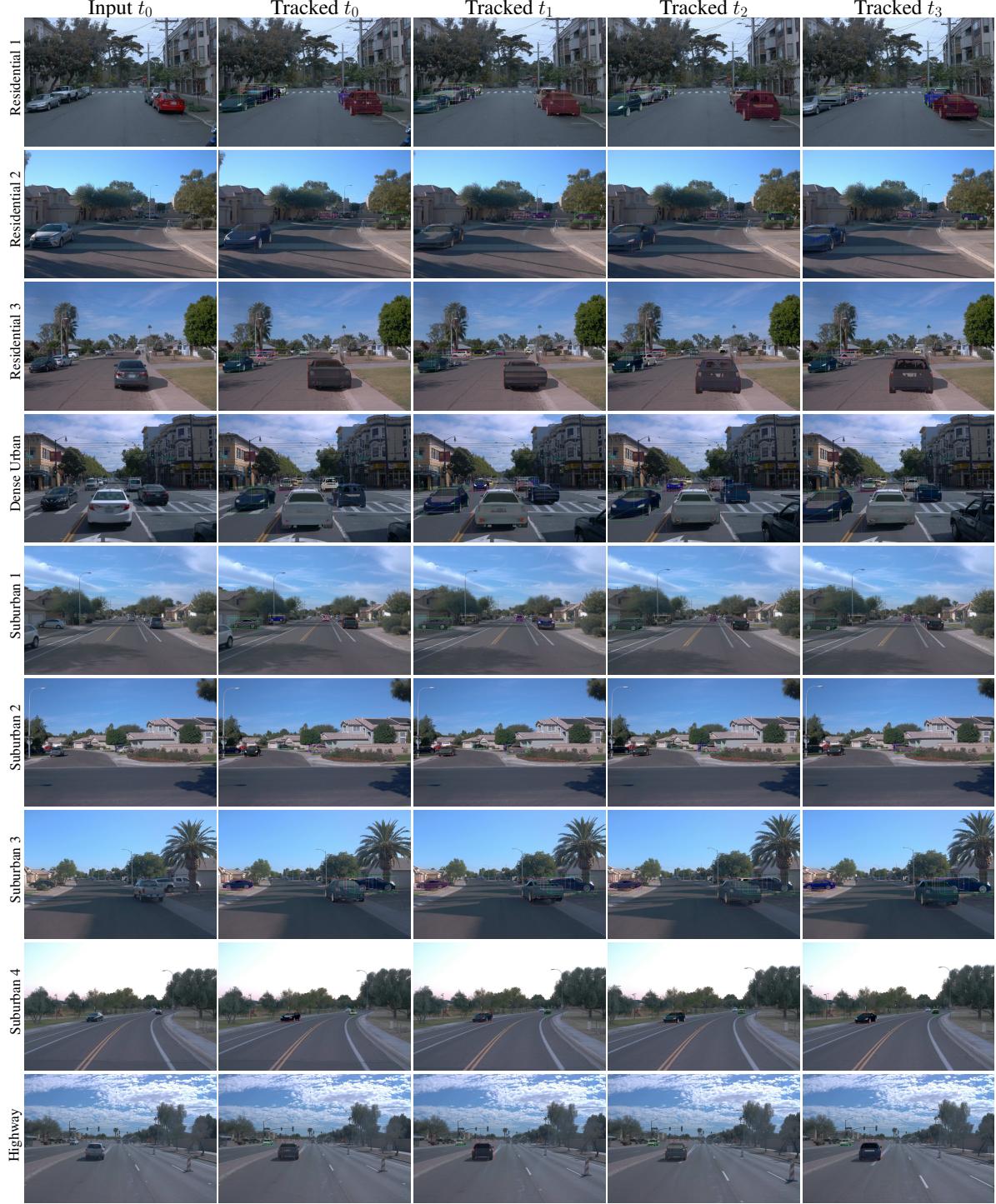


Figure 5. Additional visualizations on Waymo [12]. From left to right, we show (i) observed images from diverse scenes at timestep $k = 0$; (ii) an overlay of the optimized generated object and its 3D bounding boxes at timestep $k = 0, 1, 2$ and 3 . The color of the bounding boxes for each object corresponds to the predicted tracklet ID. We see that our method can accurately match and track tracklets in diverse scenarios in the Waymo dataset as well, confirming that the method is dataset-agnostic.

References

- [1] Holger Caesar, Varun Bankiti, Alex H Lang, Sourabh Vora, Venice Erin Liong, Qiang Xu, Anush Krishnan, Yu Pan, Giancarlo Baldan, and Oscar Beijbom. nuscenes: A multimodal dataset for autonomous driving. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 11621–11631, 2020. i, iii, vi, vii, viii
- [2] Angel X. Chang, Thomas Funkhouser, Leonidas Guibas, Pat Hanrahan, Qixing Huang, Zimo Li, Silvio Savarese, Manolis Savva, Shuran Song, Hao Su, Jianxiong Xiao, Li Yi, and Fisher Yu. ShapeNet: An Information-Rich 3D Model Repository. Technical Report arXiv:1512.03012 [cs.GR], Stanford University — Princeton University — Toyota Technological Institute at Chicago, 2015. iv
- [3] Jun Gao, Tianchang Shen, Zian Wang, Wenzheng Chen, Kangxue Yin, Daqing Li, Or Litany, Zan Gojcic, and Sanja Fidler. Get3d: A generative model of high quality 3d textured shapes learned from images. In *Advances In Neural Information Processing Systems*, 2022. ii, iv
- [4] Hou-Ning Hu, Yung-Hsu Yang, Tobias Fischer, Trevor Darrell, Fisher Yu, and Min Sun. Monocular quasi-dense 3d object tracking. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2022. vi, vii
- [5] Minguk Kang, Jun-Yan Zhu, Richard Zhang, Jaesik Park, Eli Shechtman, Sylvain Paris, and Taesung Park. Scaling up gans for text-to-image synthesis. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10124–10134, 2023. ii
- [6] Tero Karras, Samuli Laine, and Timo Aila. A style-based generator architecture for generative adversarial networks. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 4401–4410, 2019. iii, iv
- [7] Tero Karras, Samuli Laine, Miika Aittala, Janne Hellsten, Jaakko Lehtinen, and Timo Aila. Analyzing and improving the image quality of stylegan. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 8110–8119, 2020. ii, iv
- [8] Nikhila Ravi, Jeremy Reizenstein, David Novotny, Taylor Gordon, Wan-Yen Lo, Justin Johnson, and Georgia Gkioxari. Accelerating 3d deep learning with pytorch3d. *arXiv:2007.08501*, 2020. iii
- [9] Axel Sauer, Katja Schwarz, and Andreas Geiger. Stylegan-xl: Scaling stylegan to large diverse datasets. In *ACM SIGGRAPH 2022 conference proceedings*, pages 1–10, 2022. ii
- [10] Tianchang Shen, Jun Gao, Kangxue Yin, Ming-Yu Liu, and Sanja Fidler. Deep marching tetrahedra: a hybrid representation for high-resolution 3d shape synthesis. *Advances in Neural Information Processing Systems*, 34:6087–6101, 2021. iv
- [11] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition, 2015. ii
- [12] Pei Sun, Henrik Kretzschmar, Xerxes Dotiwalla, Aurelien Chouard, Vijaysai Patnaik, Paul Tsui, James Guo, Yin Zhou, Yuning Chai, Benjamin Caine, et al. Scalability in perception for autonomous driving: Waymo open dataset. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2446–2454, 2020. i, iii, vi, vii, ix
- [13] Xinshuo Weng, Jianren Wang, David Held, and Kris Kitani. 3d multi-object tracking: A baseline and new evaluation metrics. In *2020 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 10359–10366. IEEE, 2020. iii, vi
- [14] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *Proceedings of the IEEE conference on computer vision and pattern recognition (CVPR)*, pages 586–595, 2018. ii
- [15] Xingyi Zhou, Dequan Wang, and Philipp Krähenbühl. Objects as points. In *arXiv preprint arXiv:1904.07850*, 2019. vi