

R3D3: Dense 3D Reconstruction of Dynamic Scenes from Multiple Cameras

Aron Schmied^{1*} Tobias Fischer^{1*} Martin Danelljan¹ Marc Pollefeys^{1,2} Fisher Yu¹
¹ ETH Zürich ² Microsoft

<https://www.vis.xyz/pub/r3d3/>

Abstract

Dense 3D reconstruction and ego-motion estimation are key challenges in autonomous driving and robotics. Compared to the complex, multi-modal systems deployed today, multi-camera systems provide a simpler, low-cost alternative. However, camera-based 3D reconstruction of complex dynamic scenes has proven extremely difficult, as existing solutions often produce incomplete or incoherent results. We propose R3D3, a multi-camera system for dense 3D reconstruction and ego-motion estimation. Our approach iterates between geometric estimation that exploits spatial-temporal information from multiple cameras, and monocular depth refinement. We integrate multi-camera feature correlation and dense bundle adjustment operators that yield robust geometric depth and pose estimates. To improve reconstruction where geometric depth is unreliable, e.g. for moving objects or low-textured regions, we introduce learnable scene priors via a depth refinement network. We show that this design enables a dense, consistent 3D reconstruction of challenging, dynamic outdoor environments. Consequently, we achieve state-of-the-art dense depth prediction on the DDAD and NuScenes benchmarks.

1. Introduction

Translating sensory inputs into a dense 3D reconstruction of the environment and tracking the position of the observer is a cornerstone of robotics and fundamental to the development of autonomous vehicles (AVs). Contemporary systems rely on fusing many sensor modalities like camera, LiDAR, RADAR, IMU and more, making hardware and software stacks complex and expensive. In contrast, multi-camera systems provide a simpler, low-cost alternative already widely available in modern consumer vehicles. However, image-based dense 3D reconstruction and ego-motion estimation of large-scale, dynamic scenes is an open research problem as moving objects, uniform and repetitive textures, and optical degradations pose significant algorithmic challenges.

*Equal contribution.

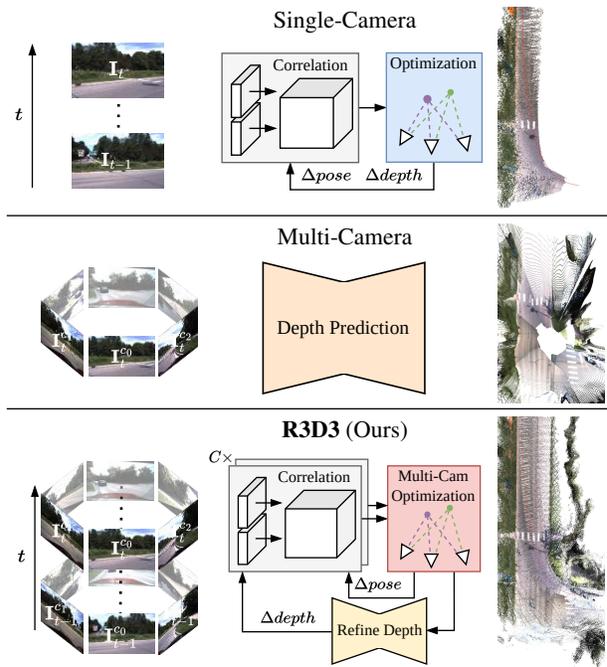


Figure 1. Many methods use temporal context but neglect inter-camera information, leading to incomplete results (**top**). Other works focus on exploiting inter-camera context but neglect temporal information, yielding incoherent predictions (**middle**). In contrast, our method achieves a consistent, dense 3D reconstruction by iteratively integrating geometric depth estimation from multiple cameras with monocular depth refinement (**bottom**).

Existing works approaching the aforementioned task can be divided into two lines of research. Many methods have focused on recovering 3D scene structure via structure-from-motion (SfM). In particular, simultaneous localization and mapping (SLAM) methods focus on accurate ego-motion estimation and usually only recover sparse 3D structure [9, 30, 35, 13, 12]. They typically treat dynamic objects or uniformly textured regions as outliers yielding an incomplete 3D reconstruction result, which makes them less suitable for AVs and robotics. In addition, only a few works have focused on multi-camera setups [25, 8, 34, 31, 27]. In contrast, multi-view stereo (MVS) methods [32, 33, 42, 29,

58, 28, 18] aim to estimate dense 3D geometry but focus on static scenes and highly overlapping sets of images with known poses.

The second line of research focuses on dense depth prediction from monocular cues, such as perspective object appearance and scene context [66, 16, 17, 60, 50, 62, 20]. However, due to the injective property of projecting 3D structures onto a 2D plane, reconstructing depth from a single image is an ill-posed problem, which limits the accuracy and generalization of these methods. Recent methods [52, 19, 22], inspired by MVS literature, combine monocular cues with temporal context but are focused on front-facing, single-camera setups. A handful of recent works extend monocular depth estimation to multi-camera setups [23, 53, 55]. These methods utilize the spatial context to improve accuracy and realize absolute scale depth learning. However, those works neglect the temporal domain which provides useful cues for depth estimation.

Motivated by this observation, we introduce R3D3, a system for dense 3D reconstruction and ego-motion estimation from multiple cameras of dynamic outdoor environments. Our approach combines monocular cues with geometric depth estimates from both spatial inter-camera context as well as inter- and intra-camera temporal context. We compute accurate geometric depth and pose estimates via iterative dense correspondence on frames in a co-visibility graph. For this, we extend the dense bundle adjustment (DBA) operator in [49] to multi-camera setups, increasing robustness and recovering absolute scene scale. To determine co-visible frames across cameras, we propose a simple yet effective multi-camera algorithm that balances performance and efficiency. A depth refinement network takes geometric depth and uncertainty as input and produces a refined depth that improves the reconstruction of, *e.g.*, moving objects and uniformly textured areas. We train this network on real-world driving data without requiring any LiDAR ground-truth. Finally, the refined depth estimates serve as the basis for the next iterations of geometric estimation, thus closing the loop between incremental geometric reconstruction and monocular depth estimation.

We summarize our **contributions** as follows. **1)** we propose R3D3, a system for dense 3D reconstruction and ego-motion estimation in dynamic scenes, **2)** we estimate geometric depth and poses with a novel multi-camera DBA formulation and a multi-camera co-visibility graph, **3)** we integrate prior geometric depth and uncertainty with monocular cues via a depth refinement network.

As a result, we achieve state-of-the-art performance across two widely used multi-camera depth estimation benchmarks, namely DDAD [20] and NuScenes [3]. Further, we show that our system exhibits superior accuracy and robustness compared to monocular SLAM methods [49, 5].

2. Related Work

Multi-view stereo. MVS methods aim to recover dense 3D scene structure from a set of images with known poses. While earlier works focused on classical optimization [32, 42, 4, 14, 15, 63, 33], more recent works capitalize on the success of convolutional neural networks (CNNs). They use CNNs to estimate features that are matched across multiple depth-hypothesis planes in a 3D cost volume [29, 58, 65, 28, 18, 59, 61]. While early approaches adopt multiple cost volumes across image pairs [65], recent approaches use a single cost volume across the whole image set [58]. These works assume a controlled setting with many, highly overlapping images and known poses to create a 3D cost volume. Instead, we aim to achieve robust, dense 3D reconstruction from an arbitrary multi-camera setup on a moving platform with an unknown trajectory.

Visual SLAM. Visual SLAM approaches focus on jointly mapping the environment and tracking the trajectory of the observer from visual inputs, *i.e.* one or multiple RGB cameras. Traditional SLAM systems are often split into different stages, where first images are processed into keypoint matches, which subsequently are used to estimate the 3D scene geometry and camera trajectory [9, 10, 6, 35, 36, 5, 40]. Another line of work focuses on directly optimizing 3D geometry and camera trajectory based on pixel intensities [13, 12, 57, 56]. A handful of works have focused on multi-camera SLAM systems [25, 8, 34, 31, 27]. Recent methods integrate CNN-based depth and pose predictions [46, 57, 56] into the SLAM pipeline. The common challenge these methods face is outliers in the pixel correspondences caused by the presence of low-textured areas, dynamic objects, or optical degradations. Hence, robust estimation techniques are used to filter these outliers, yielding an incomplete, sparse 3D reconstruction result.

In contrast, dense 3D reconstruction and ego-motion estimation has proven to be more challenging. Early works utilize active depth sensors [37, 65] to alleviate the above-mentioned challenges. Recently, several works on dense visual SLAM from RGB input have emerged [1, 45, 7, 47, 49]. While these methods are able to produce high-quality depth maps and camera trajectories, they inherit the limitations of their traditional counterparts, *i.e.* their predictions are subject to noise when there are outliers in the pixel correspondences. This can lead to artifacts in the depth maps, inaccurate trajectories, or even a complete failure of the system. On the contrary, we achieve robust, dense 3D reconstruction and ego-motion estimates by jointly leveraging multi-camera constraints as well as monocular depth cues.

Self-supervised depth estimation. The pioneering work of Zhou *et al.* [66] learns depth estimation as a proxy task while minimizing a view synthesis loss that uses geometric constraints to warp color information from a reference to a target view. Subsequent research has focused on improv-

ing network architectures, loss regularization, and training schemes [16, 17, 21, 60, 50, 62, 20]. Recent methods draw inspiration from MVS and propose to use 3D cost volumes in order to incorporate temporal information [52, 19, 54, 22]. While these methods achieve promising results, they still focus on single-camera, front-facing scenarios that do not reflect the real-world sensor setups of AVs [3, 44, 20]. Another recent line of work focuses on exploiting spatial information across overlapping cameras in a multi-camera setup [23, 53, 55]. These works propose to use static feature matches across cameras at a given time to guide depth learning and estimation, and to recover absolute scale. Our work aims to exploit both spatial and temporal cues in order to achieve a robust, dense 3D reconstruction from multiple cameras in dynamic outdoor environments.

3. Method

Problem formulation. At each timestep $t \in T$ we are given a set of images $\{\mathbf{I}_t^c\}_{c=1}^C$ from C cameras that are mounted on a moving platform with known camera intrinsics \mathbf{K}_c and extrinsics \mathbf{T}_c with respect to a common reference frame. We aim to estimate the depth maps $\mathbf{d}_t^c \in \mathbb{R}_+^{H \times W}$ and ego-pose $\mathbf{P}_t \in SE(3)$ at the current time step.

Overview. Our system is composed of three stages. First, given a new set of images $\{\mathbf{I}_t^c\}_{c=1}^C$, we extract deep features from each \mathbf{I}_t^c . We maintain a co-visibility graph $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ where a frame \mathbf{I}_t^c represents a node $v \in \mathcal{V}$ and co-visible frames are connected with edges. For each edge $(i, j) \in \mathcal{E}$ in this graph, we compute the feature correlation \mathbf{C}_{ij} of the two adjacent frames. Second, given initial estimates of the depth \mathbf{d}_i and the relative pose $\mathbf{G}_{ij} = (\mathbf{P}_{t_j} \mathbf{T}_{c_j})^{-1} \mathbf{P}_{t_i} \mathbf{T}_{c_i}$ between the two frames, we compute the induced flow \mathbf{f}_{ij}^* . We correct the induced flow with a recurrent update operator using the feature correlations \mathbf{C}_{ij} . We then globally align the updated flow estimates at each edge (i, j) with the current estimates of \mathbf{G}_{ij} and \mathbf{d}_i across the co-visibility graph \mathcal{G} with our proposed multi-camera DBA. Third, we introduce a depth refinement network that refines the geometric depth estimates with monocular depth cues to better handle scenarios that are challenging for geometric estimation methods. We illustrate our architecture in Fig. 2.

3.1. Feature Extraction and Correlation

Given a new observation $\{\mathbf{I}_t^c\}_{c=1}^C$, we first extract deep image features, update the co-visibility graph, and compute feature correlations. We describe these steps next.

Feature extraction. We follow [48, 49] and extract both correlation and context features from each image individually via deep correlation and context encoders g_ϕ and g_ψ .

Co-visibility graph. We store correlation and context features in a graph $\mathcal{G} = (\mathcal{V}, \mathcal{E})$. Each node corresponds to an image \mathbf{I}_t^c and an edge in the graph indicates that two

images are considered a pair in the feature correlation and bundle adjustment steps. Contrary to [49], we construct \mathcal{G} with three kinds of edges: temporal, spatial, and spatial-temporal. Since the number of possible edges is $|\mathcal{V}|^2$ and the system runtime scales linearly with $|\mathcal{E}|$, the degree of sparsity in the co-visibility graph is crucial. We thus carefully design a simple yet effective co-visibility graph construction algorithm in multi-camera setups (see Fig. 3).

For temporal edges (Fig. 3 in green), we examine the time window $t - \Delta t_{\text{intra}}$ and connect all frames of the same camera that are less than r_{intra} time steps apart. For spatial and spatial-temporal edges we exploit three priors. 1) we know camera adjacency from the given calibration, 2) we assume motion in the direction of the forward-facing camera and 3) the observer undergoes limited motion within the local time window. We establish spatial edges (Fig. 3 in red) between adjacent cameras within time step t . We further establish spatial-temporal edges (Fig. 3 in orange) given the static camera adjacency within $t - \Delta t_{\text{inter}}$. In particular, we connect camera c_i and c_j from any t' to $t' - r_{\text{inter}}$, if c_i and c_j connect within t and if c_j is closer to the forward-facing camera than c_i under forward motion assumption. Finally, we remove an edge (i, j) if node i or j is outside the local time windows $t - \Delta t_{\text{inter}}$ or $t - \Delta t_{\text{intra}}$. We remove unconnected nodes from the graph. We provide pseudo-code and additional discussion in the supplemental material.

Feature correlation. For each edge $(i, j) \in \mathcal{E}$ we compute the 4D feature correlation volume via the dot-product

$$\mathbf{C}_{ij} = \langle g_\phi(\mathbf{I}_i), g_\phi(\mathbf{I}_j) \rangle \in \mathbb{R}^{H \times W \times H \times W}. \quad (1)$$

As in [48, 49], we constrain the correlation search region to a radius r with a lookup operator

$$L_r : \mathbb{R}^{H \times W \times H \times W} \times \mathbb{R}^{H \times W \times 2} \rightarrow \mathbb{R}^{H \times W \times (r+1)^2}, \quad (2)$$

that samples a $H \times W$ grid of coordinates from the correlation volume using bilinear interpolation.

3.2. Depth and Pose Estimation

Given the correlation volume \mathbf{C}_{ij} , we here describe how to estimate the relative pose \mathbf{G}_{ij} and depth \mathbf{d}_i for each edge $(i, j) \in \mathcal{E}$ in the co-visibility graph.

Flow correction. Given an initial estimate of \mathbf{G}_{ij} and \mathbf{d}_i , we first compute the induced flow \mathbf{f}_{ij}^* to sample the correlation volume \mathbf{C}_{ij} using (2). We then feed the sampled correlation features, the context features $g_\psi(\mathbf{I}_i)$, and the induced flow \mathbf{f}_{ij}^* into a convolutional GRU. The GRU predicts a flow residual \mathbf{r}_{ij} and confidence weights $\mathbf{w}_{ij} \in \mathbb{R}^{H \times W \times 2}$ as in [49]. A challenge of correspondence estimation in a multi-camera setup is that the flow magnitude varies greatly among edges in the co-visibility graph due to inter-camera edges with large relative rotation. To this end, we propose to subtract the rotational part of \mathbf{f}_{ij}^* in those edges via

$$\mathbf{f}'_{ij} = (\mathbf{x}_i + \mathbf{f}_{ij}^*) - \Pi_{c_j}(\mathbf{R}_{c_j}^{-1} \mathbf{R}_{c_i} \circ \Pi_{c_i}^{-1}(\mathbf{1})), \quad (3)$$

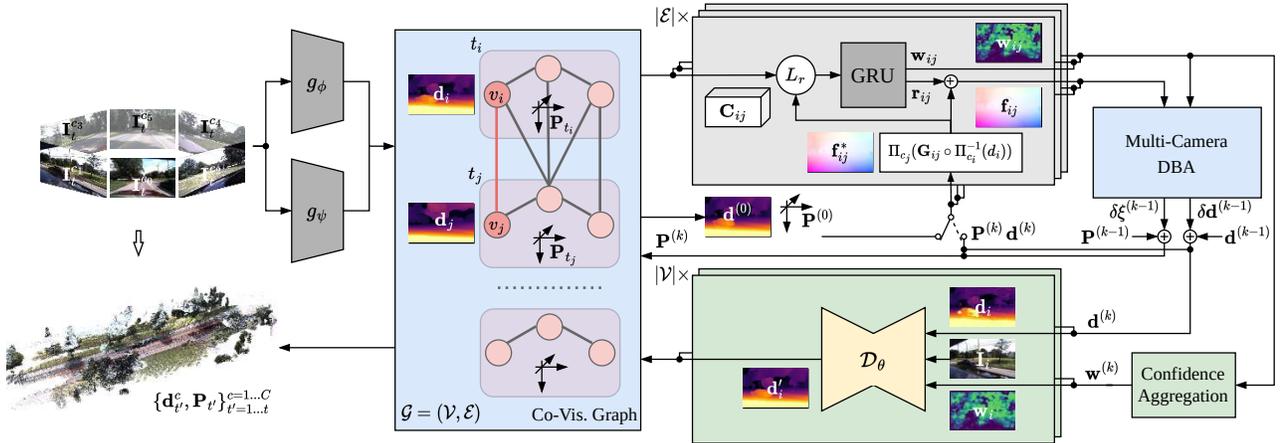


Figure 2. **Method overview.** First, frames $\{I_t^c\}_{c=1}^C$ from C cameras at time t are encoded and integrated into the co-visibility graph $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ with initial guesses of depth maps \mathbf{d}_i^c and ego-pose \mathbf{P}_t (Sec. 3.1). Second, for each edge $(i, j) \in \mathcal{E}$, we compute the induced flow \mathbf{f}_{ij}^* from \mathbf{d}_i and relative camera pose \mathbf{G}_{ij} derived from ego-poses \mathbf{P} and camera extrinsics \mathbf{T} . Given \mathbf{f}_{ij}^* , we pool feature correlations from \mathbf{C}_{ij} with operator L_r as input to a GRU that estimates flow update \mathbf{r}_{ij} and confidence \mathbf{w}_{ij} . We globally align depths \mathbf{d} and poses \mathbf{P} with the new flow estimates \mathbf{f} via our multi-camera DBA operator in k iterations (Sec. 3.2). Finally, for each node $i \in \mathcal{V}$, we refine the initial geometric estimate \mathbf{d}_i given the aggregated confidence \mathbf{w}_i via \mathcal{D}_θ (Sec. 3.3). We highlight our contributions in color.

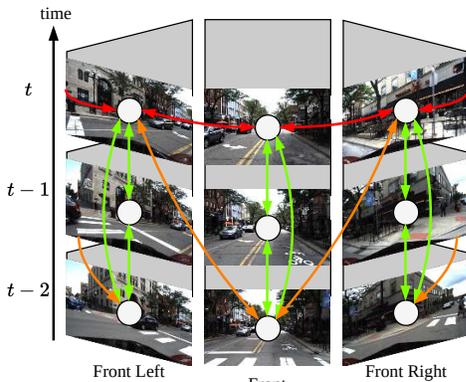


Figure 3. **Co-visibility graph.** We illustrate an example of the connectivity pattern of the co-visibility graph using our multi-camera graph construction algorithm. We use $\Delta t_{\text{intra}} = 3$, $\Delta t_{\text{inter}} = 2$, $r_{\text{intra}} = 2$, $r_{\text{inter}} = 2$ with six cameras. We depict **temporal**, **spatial**, and **spatial-temporal** edges.

where \mathbf{x}_i are pixel coordinates of frame i , Π_c denotes the projection operator under camera c and \mathbf{R} is the rotation of transformation $\mathbf{T} = (\mathbf{R} \mathbf{t})$. This aligns the flow magnitude with temporal, intra-camera edges. Note that, while we input \mathbf{f}_{ij}^l instead of \mathbf{f}_{ij}^* to the GRU, we still sample the correlation volume \mathbf{C}_{ij} based on \mathbf{f}_{ij}^* .

Pose and depth correction. We update our pose and depth estimates given the updated flow $\mathbf{f}_{ij} = \mathbf{f}_{ij}^* + \mathbf{r}_{ij}$ by minimizing the re-projection error, defined as

$$E = \sum_{(i,j) \in \mathcal{E}} \|\mathbf{x}_i + \mathbf{f}_{ij} - \Pi_{c_j}(\mathbf{G}_{ij} \circ \Pi_{c_i}^{-1}(\mathbf{d}_i))\|_{\Sigma_{ij}}^2, \quad (4)$$

where $\|\cdot\|_{\Sigma_{ij}}$ is the Mahalanobis norm and $\Sigma_{ij} =$

$\text{diag}(\mathbf{w}_{ij})$. Thus, only matched pixels where the confidence \mathbf{w}_{ij} is greater than zero contribute to the total cost. We extend the dense bundle adjustment (DBA) proposed in [49] to include the known extrinsics of the multi-camera setup. In particular, we decompose the relative poses between two frames \mathbf{G}_{ij} into the unknown, time-varying poses \mathbf{P} in the reference frame and the known static relative transformations \mathbf{T} between the cameras and the reference frame via

$$\mathbf{G}_{ij} = (\mathbf{P}_{t_j} \mathbf{T}_{c_j})^{-1} \mathbf{P}_{t_i} \mathbf{T}_{c_i}. \quad (5)$$

We linearize the residual of the energy function in Eq. 4 with a first-order Taylor expansion. We use Eq. 5 and treat \mathbf{T}_{c_j} and \mathbf{T}_{c_i} as constants when calculating the Jacobians, thus computing updates only for ego-pose \mathbf{P} . We apply the Gauss-Newton step to compute the updated $\mathbf{P}^{(k)} = \exp(\delta \xi^{(k-1)}) \mathbf{P}^{(k-1)}$ and depth $\mathbf{d}^{(k)} = \delta \mathbf{d}^{(k-1)} + \mathbf{d}^{(k-1)}$ with $\delta \xi^{(k-1)} \in \mathfrak{se}(3)$. We provide a detailed derivation in the supplementary material.

Our formulation has two advantages over DBA [49]. First, the optimization is more robust since C frames are connected through a single ego-pose \mathbf{P} . If one or multiple frames have a high outlier ratio in feature matching, *e.g.* through lens-flare, low textured regions, or dynamic objects, we can still reliably estimate \mathbf{P} with at least one undisturbed camera. Second, if there are not only pixels matched across the temporal context but also pixels matched across the static, spatial context, we can recover the absolute scale of the scene. This is because for static, spatial matches the relative transformation $\mathbf{G}_{ij} = \mathbf{T}_{c_j}^{-1} \mathbf{T}_{c_i}$ is known to scale so that E is minimized if and only if \mathbf{d} is in absolute scale.

Training. We train the networks g_ϕ and g_ψ as well as the flow correction GRU on dynamic scenes following the

procedure in [49]. As shown in [49], the geometric features learned from synthetic data can generalize to real-world scenes. We can therefore leverage existing synthetic driving datasets without relying on real-world ground-truth measurements from sensors like LiDAR or IMU to adjust our method to dynamic, multi-camera scenarios.

3.3. Depth Refinement

SfM relies on three assumptions: accurate pixel matches, sufficient camera movement, and a static scene. These assumptions do not always hold in the real world due to, *e.g.*, low-textured areas, a static ego vehicle, or many dynamic agents. Still, we would like the system to produce reasonable scene geometry to ensure safe operation.

On the contrary, monocular depth cues are inherently not affected by these issues. However, they usually lack the accurate geometric detail of SfM methods. Hence, for each node $i \in \mathcal{V}$, we complement the accurate, but sparse geometric depth estimates with monocular cues.

Network design. We use a CNN \mathcal{D}_θ parameterized by θ . We input the depth \mathbf{d}_i , confidence \mathbf{w}_i , and the corresponding image \mathbf{I}_i . The network predicts an improved, dense depth $\mathbf{d}'_i = \mathcal{D}_\theta(\mathbf{I}_i, \mathbf{d}_i, \mathbf{w}_i)$. We obtain the per-frame depth confidence \mathbf{w}_i for frame i by using the maximum across the per-edge confidence weights \mathbf{w}_{ij} . We compute $\mathbf{w}_i = \max_j \frac{1}{2}(\mathbf{w}_{ij}^x + \mathbf{w}_{ij}^y)$ where x and y are the flow directions. We use $\max(\cdot)$ because we observe that depth triangulation will be accurate if at least one pixel is matched with high confidence. We sparsify input depth and confidence weights by setting regions with confidence $\mathbf{w}_i < \beta$ to zero. We concatenate these with the image \mathbf{I}_i . We further concatenate depth and confidence with features at $1/8^{\text{th}}$ scale. As in [17], the output depth is predicted at four scales. To accommodate different focal lengths among cameras in the sensor setup, we do focal length scaling of the output [46].

Training. Contrary to geometric approaches, monocular depth estimators infer depth from semantic cues, which makes it hard for them to generalize across domains [24]. Hence, instead of relying on synthetic data, we train \mathcal{D}_θ on raw, real-world video in a self-supervised manner by minimizing a view synthesis loss [66]. We minimize the photometric error $\text{pe}(\mathbf{I}_t^c, \mathbf{I}_{t' \rightarrow t}^{c'})$ between a target image \mathbf{I}_t^c and a reference image $\mathbf{I}_{t'}^{c'}$ warped to the target viewpoint via

$$\mathbf{I}_{t' \rightarrow t}^{c'} = \Phi(\mathbf{I}_{t'}^{c'}; \Pi_{c'}(\mathbf{G}_{(t,c)(t',c')} \circ \Pi_c^{-1}(\mathbf{d}_t^c))), \quad (6)$$

with Φ the bi-linear sampling function, \mathbf{d}_t^c the predicted depth, $c' = c \pm 1$ and $t' = t \pm 1$. We compute photometric similarity with SSIM [51] and L_1 distances. We use $\mathbf{G}_{(t,c)(t',c')} = (\mathbf{P}_{t'} \mathbf{T}_{c'})^{-1} \mathbf{P}_t \mathbf{T}_c$ generated by the first two stages of our system (see Sec. 3.2). Further, self-supervised depth estimation is well-studied, and we follow the common practice of applying regularization techniques to filter the photometric error [17, 64, 21]. First, we mask areas

where $\mathbf{M}^{\text{st}} = [\text{pe}(\mathbf{I}_t^c, \mathbf{I}_{t' \rightarrow t}^{c'}) < \text{pe}(\mathbf{I}_t^c, \mathbf{I}_{t'}^{c'})]$, *i.e.* we filter regions where assuming a stationary scene would induce a lower loss than re-projection. Second, we compute an induced flow consistency mask from \mathbf{f}^* ,

$$\mathbf{M}^{\text{fc}} = \left[\left\| \mathbf{f}_{(t,c)(t',c')}^* + \Phi(\mathbf{f}_{(t',c')(t,c)}^*; \mathbf{x} + \mathbf{f}_{(t,c)(t',c')}^*) \right\|_2 < \gamma \right]. \quad (7)$$

This term warps pixel coordinates \mathbf{x} from target to reference view, looks up their value in $\mathbf{f}_{(t',c')(t,c)}^*$, and compares them to $\mathbf{f}_{(t,c)(t',c')}^*$. Therefore, pixels not following the epipolar constraint like dynamic objects are masked. Third, to handle the self-occlusion of the ego-vehicle, we manually draw a static mask \mathbf{M}^{oc} . We can use a single mask for all frames since the cameras are mounted rigidly. Further, to handle occlusions due to change in perspective, we take the minimum loss $\min_{t',c'} \text{pe}(\mathbf{I}_t^c, \mathbf{I}_{t' \rightarrow t}^{c'})$ across all reference views. The overall loss for \mathcal{D}_θ is thus

$$\mathcal{L} = \min_{t',c'} \left(\mathbf{M}^{\text{st}} \mathbf{M}^{\text{fc}} \mathbf{M}^{\text{oc}} \cdot \text{pe}(\mathbf{I}_t^c, \mathbf{I}_{t' \rightarrow t}^{c'}) \right) + \lambda \mathcal{L}_{\text{smooth}}, \quad (8)$$

with $\mathcal{L}_{\text{smooth}}$ a spatial smoothness regularization term [16].

3.4. Inference Procedure

Given an input stream of multi-camera image sets, we perform incremental 3D reconstruction in three phases.

Warmup phase. First, we add frames to the co-visibility graph if a single GRU step of a reference camera yields a large enough mean optical flow until we reach n_{warmup} frames. At this stage, we infer depth with our refinement network $\mathbf{d}_t^c = \mathcal{D}_\theta(\mathbf{I}_t^c, \mathbf{0}, \mathbf{0})$. We write the depth into the co-visibility graph so that it can be used during initialization.

Initialization phase. We initialize the co-visibility graph with the available frames and run $n_{\text{itr-wm}}$ GRU and bundle adjustment iterations. For the first $n_{\text{itr-wm}}/2$ iterations, we fix depths in the bundle adjustment and only optimize poses, which helps with inferring the absolute scale of the scene.

Active phase. We now add each incoming frame to the co-visibility graph. We initialize the new frames with $\mathbf{d}_t^c = \text{mean}(\mathbf{d}_{t-4:t-1}^c)$ and $\mathbf{P}_t = \mathbf{P}_{t-1}$. Then, we perform n_{iter1} GRU and bundle adjustment updates, and apply the refinement network to all updated frames. If the mean optical flow of the reference camera between t and $t-1$ is low, we remove all frames at $t-1$ from the co-visibility graph. This helps to handle situations with little to no ego-motion. In case no frame is removed, we do another n_{iter2} iterations. Finally, we write the updated depths and poses into the co-visibility graph.

4. Experiments

4.1. Experimental Setup

We perform extensive experiments on large-scale driving datasets that contain recordings from vehicles equipped

Table 1. **Method ablation.** We ablate our method components on the DDAD [20] dataset. We examine the influence of each component given the geometric estimation as in [49] with naive DBA as the baseline. Further, we show the influence of VKITTI [2] fine-tuning on the geometric estimation. We enumerate each variant for better reference. * median scaled depth, ** scaled trajectory.

No.	Geom. Est.	VKITTI	Multi-Cam DBA	Refinement	Abs Rel ↓	Sq Rel ↓	RMSE ↓	$\delta_{1.25}$ ↑	ATE [m] ↓	ATE** [m] ↓
1a*	✓	✓			0.438	26.97	19.109	0.636	-	2.134
1b	✓		✓		0.442	36.39	18.664	0.736	1.672	0.435
1c	✓	✓	✓		0.320	15.45	16.303	0.727	2.356	0.922
2a				✓	0.211	3.806	12.668	0.715	-	-
3a	✓	✓	✓	✓	0.162	3.019	11.408	0.811	1.235	0.433

Table 2. **Co-visibility graph ablation.** We compare our co-visibility graph construction algorithm to the original algorithm in [49] on the DDAD [20] dataset. We observe that our algorithm improves the overall runtime by an order of magnitude while maintaining the same level of performance.

Graph	Abs Rel ↓	Sq Rel ↓	RMSE ↓	$\delta_{1.25}$ ↑	ATE [m] ↓	t_{inf} [s] ↓
Original	0.162	2.968	11.156	0.816	0.982	3.23
Ours	0.162	3.019	11.408	0.811	1.235	0.35

Table 3. **Refinement network ablation.** We show on the DDAD [20] dataset that both adding confidence weights w_i and sparsifying the depth input via $w_i < \beta$ to filter outliers in the input of depth refinement network \mathcal{D}_θ is vital to depth accuracy.

w_i	$w_i < \beta$	Abs Rel ↓	Sq Rel ↓	RMSE ↓	$\delta_{1.25}$ ↑
		0.181	4.078	12.108	0.789
✓		0.173	3.497	11.924	0.792
✓	✓	0.162	3.019	11.408	0.811

with multiple cameras, LiDARs, and other sensors. Note that we only use camera data for training and inference, and use LiDAR data only as ground-truth for evaluation. We evaluate our method in inference mode (cf. Sec. 3.4), *i.e.* we obtain predictions in a fully online (causal) manner to mimic real-world deployment.

DDAD. This dataset consists of 150 training and 50 validation scenes from urban areas. Each sequence has a length of 50 or 100 time steps at a frame rate of 10Hz. The sensor setup includes six cameras in a surround-view setup with up to 20% overlap. We follow [20] and evaluate depth up to 200m averaged across all cameras. We use self-occlusion masks from [53] to filter the ego-vehicle in the images. Images have a resolution of 1216×1936 and are downsampled to 384×640 in our experiments.

NuScenes. This dataset contains 700 training, 150 validation, and 150 testing sequences of urban scenes with challenging conditions such as nighttime and rain. Each scene is composed of 40 keyframes at 2Hz synchronized across sensors. The six cameras have a sampling rate of 12Hz and are arranged in a surround-view setup with up to 10% overlap. We follow [23] and evaluate depth averaged across all cameras on the validation set. We use a single self-occlusion mask for all scenes. While the raw images have a resolution of 900×1600 , we use 768×448 .

Implementation details. We implement our system in PyTorch [38]. Our multi-camera DBA is implemented in

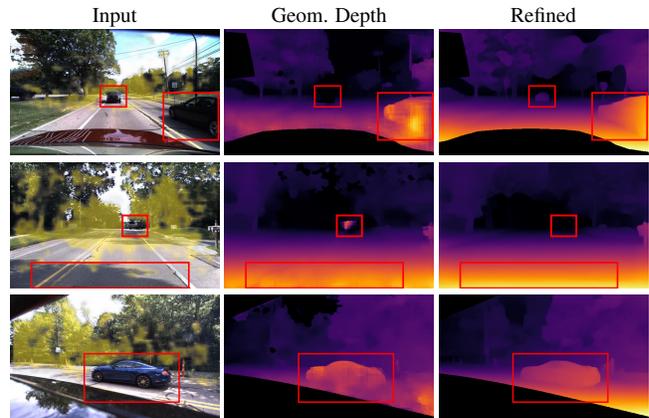


Figure 4. **Depth refinement.** We plot images overlaid with confidence w_i , the geometric depth estimate d , and the refined depth d' from \mathcal{D}_θ . We observe a smoother ground plane, corrected dynamic objects, and filtered outliers.

CUDA. We fine-tune g_ϕ , g_ψ and the flow correction GRU on VKITTI2 [2] for 10 epochs with a batch size of 1, sequence length of 7, and learning rate 10^{-4} with the pre-trained model in [49]. For \mathcal{D}_θ , we use a U-Net [39] with ResNet18 [26] encoder and initialize it with ImageNet [41] pre-trained weights. We use $\beta = 0.5$, $\gamma = 3$ and $\lambda = 10^{-3}$. We train for 20 epochs with batch size of 6 and the Adam optimizer with learning rate 10^{-4} , $\beta_1 = 0.9$, and $\beta_2 = 0.999$. We reduce the learning rate by a factor of 10 after 15 epochs. For inference, the runtime is measured for a set of six images on a single RTX3090 GPU. Each sequence is initialized with 3 warmup frames filtered at a mean flow threshold of 1.75. We build the co-visibility graph with $\Delta t_{intra} = 3$, $r_{intra} = 2$, $\Delta t_{inter} = 2$ and $r_{inter} = 2$.

4.2. Ablation Studies

Method ablation. We ablate our method components on the DDAD dataset in Tab. 1. The geometric estimation baseline (1a) consists of [49] with the naive DBA formulation applied to all cameras and the co-visibility graph as described in Sec. 3.1. It performs poorly in depth and pose estimation, with an Abs. Rel. score of 0.438 and an absolute trajectory error (ATE) [43] of 2.134m, even when adjusting for scale. Next, we add our multi-camera DBA (1b). We observe that, even without VKITTI fine-tuning, the pose esti-

Table 4. **Comparison to state-of-the-art on DDAD.** We compare favorably to existing methods on both scale-aware and scale-invariant depth prediction. * median scaled depth.

Method	Abs Rel ↓	Sq Rel ↓	RMSE ↓	$\delta_{1.25} \uparrow$
FSM [23]	0.201	-	-	-
SurroundDepth [53]	0.208	3.371	12.977	0.693
Ours	0.162	3.019	11.408	0.811
FSM* [23]	0.202	-	-	-
SurroundDepth* [53]	0.200	3.392	12.270	0.740
MCDP* [55]	0.193	3.111	12.264	0.811
Ours*	0.169	3.041	11.372	0.809

Table 5. **Per-camera evaluation on DDAD.** We show a per-camera breakdown of previous works, our \mathcal{D}_θ only (cf. 2a in Tab. 1), and our full method. Our full method performs the best with a particularly significant improvement on the side-view cameras while \mathcal{D}_θ performs similarly to previous works.

Method	Abs Rel ↓					
	Front	F.Left	F.Right	B.Left	B.Right	Back
FSM [23]	0.130	0.201	0.224	0.229	0.240	0.186
SurroundDepth [53]	0.152	0.207	0.230	0.220	0.239	0.200
\mathcal{D}_θ only	0.154	0.213	0.237	0.231	0.237	0.194
Ours	0.128	0.160	0.168	0.172	0.174	0.169

Table 6. **Comparison to multi-frame methods on DDAD.** We compare to methods that exploit temporal context from a single camera [52, 19]. We evaluate only the front camera since these methods were trained only on this camera.

Method	Front Camera			
	Abs Rel ↓	Sq Rel ↓	RMSE ↓	$\delta_{1.25} \uparrow$
ManyDepth [52]	0.146	3.258	14.098	0.822
DepthFormer [19]	0.135	2.953	12.477	0.836
Ours	0.128	3.168	13.214	0.868

Table 7. **Pose evaluation on DDAD.** We compare our method to state-of-the-art monocular SLAM methods for which we report the average ATE across the cameras that successfully tracked the camera path. ** scaled trajectory.

	DROID-SLAM [49]	ORB-SLAMv3 [5]	Ours
ATE** [m] ↓	7.500	6.179	0.433

mation accuracy improves dramatically over the naive DBA baseline. The relative scale ATE drops from 2.134m to only 0.435m. The absolute scale ATE is only about 1.2m higher, indicating the scene scale is recovered accurately.

Using VKITTI fine-tuning (1c), we achieve significant gains in depth accuracy, where Abs. Rel. drops to 0.32 and $\delta_{1.25}$ increases to 72.7%. However, note that VKITTI fine-tuning does not affect the $\delta_{1.25}$ score, *i.e.* it helps mostly with depth outliers. We attribute this to the absence of dynamic objects in the training data of [49] so that fine-tuning helps to adjust the confidence prediction to handle such outliers (cf. Sec. 3.2). We further test the depth refinement network (Sec. 3.3) without geometric estimation prior (2a). The depth obtained from \mathcal{D}_θ is less prone to outliers, as evidenced by its low Abs. Rel., Sq. Rel. and RMSE scores. However, note that its $\delta_{1.25}$ is only 71.5%, *i.e.* its

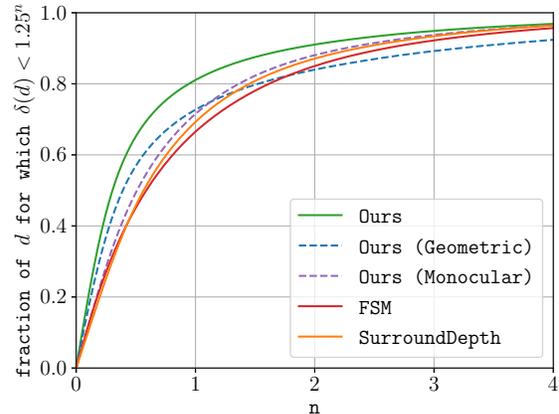


Figure 5. **Depth in $\delta_{1.25^n}$ range for different thresholds on DDAD.** We depict the ratio of depth estimates d for which $\delta(d) < 1.25^n$ where $\delta(d) = \max(\frac{d}{d^*}, \frac{d^*}{d})$ and d^* is ground-truth depth. Note that values at $n = 1, 2, 3$ represent $\delta_{1.25}, \delta_{1.25^2}, \delta_{1.25^3}$.

Table 8. **Comparison to state-of-the-art on NuScenes.** Evaluated up to 80m for scale-aware depth as in [23, 53] and 60m for median-scaled depth (denoted with *) as in [55].

Method	Abs Rel ↓	Sq Rel ↓	RMSE ↓	$\delta_{1.25} \uparrow$
FSM [23]	0.297	-	-	-
SurroundDepth [53]	0.280	4.401	7.467	0.661
Ours	0.253	4.759	7.150	0.729
PackNet*[20]	0.303	3.154	7.014	0.655
FSM* [23]	0.270	3.185	6.832	0.689
SurroundDepth* [53]	0.245	3.067	6.835	0.719
MCDP* [55]	0.237	3.030	6.822	0.719
Ours*	0.235	3.332	6.021	0.749

depth accuracy compared to the geometric depth estimates is actually lower when disregarding outliers. Finally, we test our full system (3a). Compared to geometric estimation (1d), we observe a significant improvement in outlier sensitive metrics like RMSE, but also in outlier robust metrics like $\delta_{1.25}$. Additionally, the pose estimation accuracy is increased and the difference between scaled and non-scaled ATE is smaller. Compared to depth obtained from \mathcal{D}_θ (2a), the increase in outlier sensitive metrics is smaller but still significant, *e.g.* 1.26m in RMSE. Further, we observe a sizeable gain in $\delta_{1.25}$ of 9.6 percentage points.

Further, in Fig. 5, we compare the $\delta_{1.25^n}$ scores at different accuracy thresholds n of our geometric depth estimation without refinement ('Ours (Geometric)'), our refinement network only ('Ours (Monocular)') and our full method ('Ours'). The results confirm our hypothesis that while geometric estimates are more accurate than monocular depth estimates, they are also noisier. Importantly, the results further illustrate that our full method can effectively combine the strengths of both approaches while not suffering from their weaknesses.

Co-visibility graph construction. In Tab. 2, we compare our co-visibility graph construction algorithm (Sec. 3.1) to a multi-camera adaption of the algorithm in [49]. Ours

Table 9. **Per-camera evaluation on NuScenes.** We compare our method with previous works on scale-aware depth estimation. Evaluated up to 80m as in Tab. 8. We observe the same trend as on DDAD, *i.e.* our method performs best across all cameras, with a particularly high improvement on the side views.

Method	Abs Rel ↓					
	Front	F.Left	F.Right	B.Left	B.Right	Back
FSM [23]	0.186	0.287	0.375	0.296	0.418	0.221
SurroundDepth [53]	0.179	0.260	0.340	0.282	0.403	0.212
Ours	0.174	0.230	0.302	0.249	0.360	0.201

Table 10. **Cross-dataset transfer.** We use models trained on DDAD to evaluate scale-aware depth prediction on the NuScenes dataset. We compare to state-of-the-art methods [23, 53] and observe that our method exhibits better generalization ability (cf. Tab. 8).

Method	Abs Rel ↓	Sq Rel ↓	RMSE ↓	$\delta_{1.25}$ ↑
FSM [23]	0.349	5.064	8.785	0.499
SurroundDepth [53]	0.364	5.476	8.447	0.525
Ours	0.292	4.800	7.677	0.660

yields similar performance in depth estimation, with the exact same Abs. Rel. score and only minor difference in other metrics. Further, the absolute scale ATE is only slightly (25cm) higher. In contrast, using our algorithm, we achieve a runtime reduction of nearly $10\times$ on the full system.

Depth refinement network. We verify the choice of our inputs to \mathcal{D}_θ in Tab. 3. If we do not input the confidence w_i and do not sparsify the depth via $w_i < \beta$ before feeding it to \mathcal{D}_θ , we observe the lowest performance. If we input the confidence w_i but do not sparsify the depth, we see an improvement over not using them. We achieve the best performance with the inputs described in Sec. 3.3.

Moreover, we show qualitative examples of our depth refinement in Fig. 4. The confidence w_i (overlaid with the image) segments regions where geometric estimation is ambiguous, *i.e.* dynamic objects and uniformly textured areas like the sky and the road surface. The geometric depth maps exhibit severe artifacts in such regions. After refinement, these artifacts are corrected.

4.3. Comparison to State-of-the-Art

DDAD. In Tab. 4, we compare our method to existing self-supervised multi-camera depth estimation methods. We achieve significant improvements in scale-aware and scale-invariant depth prediction with only minor differences in absolute and relative scale results. In both settings, we achieve state-of-the-art results with the lowest Abs. Rel. scores of 0.162 and 0.169, respectively.

In Fig. 5 we compare the $\delta_{1.25^n}$ scores of our method to the state-of-the-art approaches on DDAD at different n . For small values of n , we observe a particularly pronounced performance gain over existing methods. This shows the advantage of combining highly accurate geometric estimates with monocular cues. Further, our method consistently out-

performs the baselines over all values of n for $\delta < 1.25^n$.

In Fig. 6, we show a qualitative comparison to existing works. On the left side, we illustrate depth maps across three cameras. Compared to FSM [23] and SurroundDepth [53], we produce sharper depth boundaries. On the right side, we illustrate point clouds accumulated over time. We produce a more consistent reconstruction, as can be seen when focusing on the vehicles or the fences beside the road.

Tab. 5 shows the per-camera breakdown of the scale-aware depth prediction comparison. The advantage of our method is particularly pronounced on the side cameras. Other methods struggle with the side views because the static spatial context is limited and they do not exploit any temporal or spatial-temporal context. The temporal and spatial-temporal contexts are conducive in the side views since sideward motion leads to better triangulation angles than forward motion. This is also evidenced by the results of \mathcal{D}_θ without geometric estimation, which performs similarly to the baseline methods, struggling with the side views. In contrast, our full system performs well across all cameras.

We further compare to self-supervised single-camera methods that leverage temporal information in Tab. 6. For a fair comparison, we evaluate only the front camera. We compare favorably to DepthFormer [19] and ManyDepth [52]. In particular, we substantially improve in the outlier robust metric $\delta_{1.25}$ and perform competitively in all others. Our Abs. Rel. score is the lowest at 0.128.

To evaluate ego-motion estimation, we compare our system to state-of-the-art monocular SLAM methods [49, 5] in Tab. 7. We run the baselines on each camera and report the average ATE across the cameras that successfully tracked the camera path. Note that, contrary to our approach, the competing methods frequently failed to produce a trajectory. We observe that state-of-the-art monocular SLAM methods exhibit a large error (ATE), while our method recovers accurate trajectories.

NuScenes. In Tab. 8, we compare our system to previous self-supervised multi-camera depth estimation methods. We outperform previous methods by a significant margin, corroborating our findings on DDAD. Notably, we observe substantial improvements in scale-aware depth estimation, while previous works usually struggle with recovering absolute scale on this dataset since the overlap between cameras is smaller than in DDAD [23]. Therefore, previous works like MCDP [55] report results on scale-invariant depth prediction. We outperform existing methods also in this setting, although the performance gap is smaller. Note that the gap between scale-aware and scale-invariant evaluation is only 2 percentage points in $\delta_{1.25}$ for our method, while for SurroundDepth [53] the gap is much larger with 5.8 percentage points. We show a per-camera breakdown of the comparison to the state-of-the-art in scale-aware depth estimation in Tab. 9. The results corroborate our findings

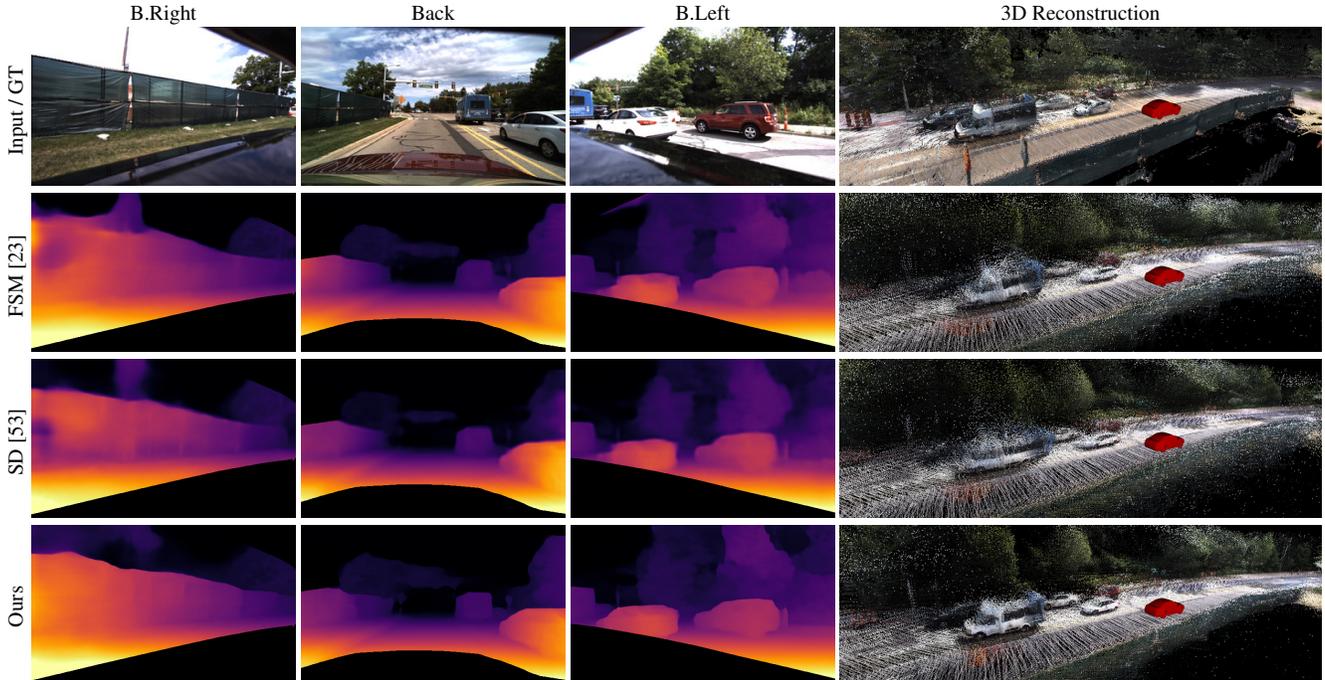


Figure 6. **Qualitative comparison on DDAD.** We show depth maps and point clouds accumulated over time along with input images, ground-truth LiDAR 3D Reconstruction, and the ego-vehicle in red. Our depth maps are sharper and more accurate. Our accumulated point clouds yield a more consistent 3D reconstruction. Note that we use our pose predictions for [23, 53] since these do not predict pose.

on DDAD, namely that our method outperforms previous works across all cameras with a particularly pronounced improvement on the side views.

Cross-dataset transfer. We finally test the cross-dataset generalization of our method versus FSM [23] and SurroundDepth [53]. In Tab. 10, we show *scale-aware* depth estimation results of models trained on DDAD, evaluated on NuScenes. We observe that the gap in Abs. Rel. widens compared to Tab. 8, with our method outperforming SurroundDepth by 0.072 in Abs. Rel. and FSM by 0.057. Further, our method maintains much higher levels of $\delta_{1.25}$ and Sq. Rel., while SurroundDepth drops significantly in both metrics. Note that we apply focal length scaling [46] to all methods to accommodate differences in camera intrinsics.

5. Conclusion

We introduced R3D3, a multi-camera system for dense 3D reconstruction of dynamic outdoor environments. The key ideas we presented are a multi-camera DBA operator that greatly improves geometric depth and pose estimation, a multi-camera co-visibility graph construction algorithm that reduces the runtime of our system by nearly $10\times$ without significant performance drop, and a depth refinement network that effectively fuses geometric depth estimates with monocular cues. We observe that our design choices enable dense 3D mapping in challenging scenes presenting many dynamic objects, uniform and repetitive textures, and

complicated camera phenomena like lens flare and auto-exposure. We achieve state-of-the-art performance in dense depth prediction across two multi-camera benchmarks.

References

- [1] Michael Bloesch, Jan Czarnowski, Ronald Clark, Stefan Leutenegger, and Andrew J Davison. Codeslam—learning a compact, optimisable representation for dense visual slam. In *CVPR*, 2018.
- [2] Johann Cabon, Naila Murray, and Martin Humenberger. Virtual KITTI 2. *CoRR*, 2020.
- [3] 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 multi-modal dataset for autonomous driving. In *CVPR*, 2020.
- [4] Neill DF Campbell, George Vogiatzis, Carlos Hernández, and Roberto Cipolla. Using multiple hypotheses to improve depth-maps for multi-view stereo. In *ECCV*, 2008.
- [5] Carlos Campos, Richard Elvira, Juan J Gómez Rodríguez, José MM Montiel, and Juan D Tardós. Orb-slam3: An accurate open-source library for visual, visual-inertial, and multi-map slam. *T-RO*, 2021.
- [6] Laura A Clemente, Andrew J Davison, Ian D Reid, José Neira, and Juan D Tardós. Mapping large loops with a single hand-held camera. In *RSS*, 2007.
- [7] Jan Czarnowski, Tristan Laidlow, Ronald Clark, and Andrew J Davison. Deepfactors: Real-time probabilistic dense monocular slam. *RA-L*, 2020.

- [8] Arun Das and Steven L Waslander. Entropy based keyframe selection for multi-camera visual slam. In *IROS*, 2015.
- [9] Andrew J Davison. Real-time simultaneous localisation and mapping with a single camera. In *ICCV*, 2003.
- [10] Andrew J Davison, Ian D Reid, Nicholas D Molton, and Olivier Stasse. Monoslam: Real-time single camera slam. *IEEE TPAMI*, 2007.
- [11] Ethan Eade. Lie groups for 2d and 3d transformations. *URL* <http://ethaneade.com/lie.pdf>, revised Dec, 2013.
- [12] Jakob Engel, Vladlen Koltun, and Daniel Cremers. Direct sparse odometry. *IEEE TPAMI*, 2017.
- [13] Jakob Engel, Thomas Schöps, and Daniel Cremers. Lsd-slam: Large-scale direct monocular slam. In *ECCV*, 2014.
- [14] Yasutaka Furukawa and Jean Ponce. Accurate, dense, and robust multiview stereopsis. *IEEE TPAMI*, 2009.
- [15] Silvano Galliani, Katrin Lasinger, and Konrad Schindler. Massively parallel multiview stereopsis by surface normal diffusion. In *ICCV*, 2015.
- [16] Clément Godard, Oisín Mac Aodha, and Gabriel J. Brostow. Unsupervised monocular depth estimation with left-right consistency. In *CVPR*, 2017.
- [17] Clément Godard, Oisín Mac Aodha, Michael Firman, and Gabriel J. Brostow. Digging into self-supervised monocular depth prediction. In *ICCV*, 2019.
- [18] Xiaodong Gu, Zhiwen Fan, Siyu Zhu, Zuozhuo Dai, Feitong Tan, and Ping Tan. Cascade cost volume for high-resolution multi-view stereo and stereo matching. In *CVPR*, 2020.
- [19] Vitor Guizilini, Rares Ambrus, Dian Chen, Sergey Zakharov, and Adrien Gaidon. Multi-frame self-supervised depth with transformers. In *CVPR*, 2022.
- [20] Vitor Guizilini, Rares Ambrus, Sudeep Pillai, Allan Raventos, and Adrien Gaidon. 3d packing for self-supervised monocular depth estimation. In *CVPR*, 2020.
- [21] Vitor Guizilini, Rui Hou, Jie Li, Rares Ambrus, and Adrien Gaidon. Semantically-guided representation learning for self-supervised monocular depth. In *ICLR*, 2020.
- [22] Vitor Guizilini, Kuan-Hui Lee, Rares Ambrus, and Adrien Gaidon. Learning optical flow, depth, and scene flow without real-world labels. In *RA-L*, 2022.
- [23] Vitor Guizilini, Igor Vasiljevic, Rares Ambrus, Greg Shakhnarovich, and Adrien Gaidon. Full surround monodepth from multiple cameras. In *RA-L*, 2022.
- [24] Xiaoyang Guo, Hongsheng Li, Shuai Yi, Jimmy Ren, and Xiaogang Wang. Learning monocular depth by distilling cross-domain stereo networks. In *ECCV*, 2018.
- [25] Adam Harmat, Inna Sharf, and Michael Trentini. Parallel tracking and mapping with multiple cameras on an unmanned aerial vehicle. In *ICRA*, 2012.
- [26] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, 2016.
- [27] Lionel Heng, Benjamin Choi, Zhaopeng Cui, Marcel Geppert, Sixing Hu, Benson Kuan, Peidong Liu, Rang Nguyen, Ye Chuan Yeo, Andreas Geiger, et al. Project autovision: Localization and 3d scene perception for an autonomous vehicle with a multi-camera system. In *ICRA*, 2019.
- [28] Po-Han Huang, Kevin Matzen, Johannes Kopf, Narendra Ahuja, and Jia-Bin Huang. Deepmvs: Learning multi-view stereopsis. In *CVPR*, 2018.
- [29] Abhishek Kar, Christian Häne, and Jitendra Malik. Learning a multi-view stereo machine. *NeurIPS*, 2017.
- [30] Georg Klein and David Murray. Parallel tracking and mapping for small ar workspaces. In *ISMAR*, 2007.
- [31] Juichung Kuo, Manasi Muglikar, Zichao Zhang, and Davide Scaramuzza. Redesigning slam for arbitrary multi-camera systems. In *ICRA*, 2020.
- [32] Kiriakos N Kutulakos and Steven M Seitz. A theory of shape by space carving. In *ICCV*, 1999.
- [33] Maxime Lhuillier and Long Quan. A quasi-dense approach to surface reconstruction from uncalibrated images. *IEEE TPAMI*, 2005.
- [34] Peidong Liu, Marcel Geppert, Lionel Heng, Torsten Sattler, Andreas Geiger, and Marc Pollefeys. Towards robust visual odometry with a multi-camera system. In *IROS*, 2018.
- [35] Raul Mur-Artal, Jose Maria Martinez Montiel, and Juan D Tardos. Orb-slam: a versatile and accurate monocular slam system. *T-RO*, 2015.
- [36] Raul Mur-Artal and Juan D Tardós. Orb-slam2: An open-source slam system for monocular, stereo, and rgb-d cameras. *T-RO*, 2017.
- [37] Richard A Newcombe, Steven J Lovegrove, and Andrew J Davison. Dtam: Dense tracking and mapping in real-time. In *ICCV*, 2011.
- [38] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative style, high-performance deep learning library. *NeurIPS*, 2019.
- [39] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *MICCAI*, 2015.
- [40] Antoni Rosinol, Marcus Abate, Yun Chang, and Luca Carlone. Kimera: an open-source library for real-time metric-semantic localization and mapping. In *ICRA*, 2020.
- [41] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael S. Bernstein, Alexander C. Berg, and Li Fei-Fei. Imagenet large scale visual recognition challenge. *IJCV*, 2015.
- [42] Johannes Lutz Schönberger, Enliang Zheng, Marc Pollefeys, and Jan-Michael Frahm. Pixelwise view selection for unstructured multi-view stereo. In *ECCV*, 2016.
- [43] Jürgen Sturm, Nikolas Engelhard, Felix Endres, Wolfram Burgard, and Daniel Cremers. A benchmark for the evaluation of rgb-d slam systems. In *RSS*, 2012.
- [44] Pei Sun, Henrik Kretschmar, 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 *CVPR*, 2020.
- [45] Chengzhou Tang and Ping Tan. Ba-net: Dense bundle adjustment network. In *ICLR*, 2019.

- [46] Keisuke Tateno, Federico Tombari, Iro Laina, and Nassir Navab. Cnn-slam: Real-time dense monocular slam with learned depth prediction. In *CVPR*, 2017.
- [47] Zachary Teed and Jia Deng. Deepv2d: Video to depth with differentiable structure from motion. In *ICLR*, 2020.
- [48] Zachary Teed and Jia Deng. Raft: Recurrent all-pairs field transforms for optical flow. In *ECCV*, 2020.
- [49] Zachary Teed and Jia Deng. DROID-SLAM: Deep Visual SLAM for Monocular, Stereo, and RGB-D Cameras. *NeurIPS*, 2021.
- [50] Chaoyang Wang, José Miguel Buenaposada, Rui Zhu, and Simon Lucey. Learning depth from monocular videos using direct methods. In *CVPR*, 2018.
- [51] Zhou Wang, Alan C Bovik, Hamid R Sheikh, and Eero P Simoncelli. Image quality assessment: from error visibility to structural similarity. *IEEE TIP*, 2004.
- [52] Jamie Watson, Oisín Mac Aodha, Victor Prisacariu, Gabriel Brostow, and Michael Firman. The Temporal Opportunist: Self-Supervised Multi-Frame Monocular Depth. In *CVPR*, 2021.
- [53] Yi Wei, Linqing Zhao, Wenzhao Zheng, Zheng Zhu, Yongming Rao, Guan Huang, Jiwen Lu, and Jie Zhou. Surround-depth: Entangling surrounding views for self-supervised multi-camera depth estimation. In *CoRL*, 2022.
- [54] Felix Wimbauer, Nan Yang, Lukas Von Stumberg, Niclas Zeller, and Daniel Cremers. Monorec: Semi-supervised dense reconstruction in dynamic environments from a single moving camera. In *CVPR*, 2021.
- [55] Jialei Xu, Xianming Liu, Yuanchao Bai, Junjun Jiang, Kaixuan Wang, Xiaozhi Chen, and Xiangyang Ji. Multi-camera collaborative depth prediction via consistent structure estimation. In *ACM MM*, 2022.
- [56] Nan Yang, Lukas von Stumberg, Rui Wang, and Daniel Cremers. D3vo: Deep depth, deep pose and deep uncertainty for monocular visual odometry. In *CVPR*, 2020.
- [57] Nan Yang, Rui Wang, Jorg Stuckler, and Daniel Cremers. Deep virtual stereo odometry: Leveraging deep depth prediction for monocular direct sparse odometry. In *ECCV*, 2018.
- [58] Yao Yao, Zixin Luo, Shiwei Li, Tian Fang, and Long Quan. Mvsnet: Depth inference for unstructured multi-view stereo. In *ECCV*, 2018.
- [59] Yao Yao, Zixin Luo, Shiwei Li, Tianwei Shen, Tian Fang, and Long Quan. Recurrent mvsnet for high-resolution multi-view stereo depth inference. In *CVPR*, 2019.
- [60] Zhichao Yin and Jianping Shi. Geonet: Unsupervised learning of dense depth, optical flow and camera pose. In *CVPR*, 2018.
- [61] Zehao Yu and Shenghua Gao. Fast-mvsnet: Sparse-to-dense multi-view stereo with learned propagation and gauss-newton refinement. In *CVPR*, 2020.
- [62] Huangying Zhan, Ravi Garg, Chamara Saroj Weerasekera, Kejie Li, Harsh Agarwal, and Ian Reid. Unsupervised learning of monocular depth estimation and visual odometry with deep feature reconstruction. In *CVPR*, 2018.
- [63] Enliang Zheng, Enrique Dunn, Vladimir Jovic, and Jan-Michael Frahm. Patchmatch based joint view selection and depthmap estimation. In *CVPR*, 2014.
- [64] Yiran Zhong, Pan Ji, Jianyuan Wang, Yuchao Dai, and Hongdong Li. Unsupervised deep epipolar flow for stationary or dynamic scenes. In *CVPR*, 2019.
- [65] Huizhong Zhou, Benjamin Ummenhofer, and Thomas Brox. Deeptam: Deep tracking and mapping. In *ECCV*, 2018.
- [66] Tinghui Zhou, Matthew Brown, Noah Snavely, and David G. Lowe. Unsupervised learning of depth and ego-motion from video. In *CVPR*, 2017.

Appendix

This supplementary material provides further details on our method, our experimental setup, and more quantitative and qualitative results and comparisons. In Sec. A, we provide further details and discussion on the geometric pose and depth estimation, the refinement network, our co-visibility graph, and training and inference procedures. In Sec. B, we provide additional ablation studies and more quantitative and qualitative comparisons and results. In Sec. C, we provide a detailed derivation of our multi-camera DBA. Note that we group large qualitative comparisons at the end of this document for better readability.

A. Implementation Details

A.1. Geometric Depth and Pose Estimation

We follow [49] for the context feature encoder g_ψ and correlation feature encoder g_ϕ and GRU architecture. In addition to the context features, g_ψ also outputs an initialization to the GRU hidden state $\mathbf{h}_{ij}^{(0)}$ from input \mathbf{I}_i . The hidden state after the last iteration is pooled among all outgoing edges from node i to predict a damping factor λ_i that serves as a parameter to the multi-camera DBA, which improves convergence when the current depth estimate is inaccurate. We perform feature matching and geometric depth updates at $1/8^{\text{th}}$ of the original image resolution. We predict upsampling masks from the pooled hidden state to upsample depth to the original resolution. The confidence maps \mathbf{w}_i are linearly interpolated. Furthermore, following RAFT [48], the correlation volume is built as a 4-level pyramid, and sampled features from all layers are concatenated. Further, each GRU update is followed by a multi-camera DBA with two Gauss-Newton steps.

A.2. Depth Refinement

We show the architecture of the depth refinement network in Tab. 11. Note that we input full-resolution frames and up-sampled depth and confidence. Further, we concatenate the depth and confidence predicted at $1/8^{\text{th}}$ scale with features after the third skip connection. We use the inverse of the input depth. Further, we obtain the refined depth prediction from the sigmoid output $\mathbf{o}_t^c \in [0, 1]^{H \times W}$ via

$$\frac{1}{\mathbf{d}_t^c} = \frac{f_c}{f_{\text{norm}}} \cdot \left(\frac{1}{d_{\text{max}}} + \left(\frac{1}{d_{\text{min}}} - \frac{1}{d_{\text{max}}} \right) \cdot \mathbf{o}_t^c \right), \quad (9)$$

where f_c is the focal length, f_{norm} a constant normalization factor and d_{min} and d_{max} are pre-defined minimum and maximum depth. For experiments on the DDAD dataset we set $d_{\text{min}} = 1$, $d_{\text{max}} = 200$, $f_{\text{norm}} = 715$ aligned with the focal length of the front-facing camera, for NuScenes we choose $d_{\text{min}} = 1$, $d_{\text{max}} = 80$, $f_{\text{norm}} = 500$.

Dataset generation. To train the refinement network, we first generate a dataset of samples that contain the prior geometric depth, confidence maps, and poses with the first two stages of our system. We filter scenes with inaccurate scene scale by measuring how many reliable feature matches there are across both temporal and spatial edges with the given confidences \mathbf{w}_{ij} . The fewer reliable matches, the weaker the constraint on the metric scale. Furthermore, based on the generated poses, we remove static scenes. This allows us to train an absolute scale monocular depth estimation model from the raw video data.

Training details. During the training of refinement network \mathcal{D}_θ , we randomly set input depth and confidence weights to zero to learn depth prediction with and without prior geometric estimates as input. We use color-jitter and random horizontal flipping as augmentations. Further, we follow a two-stage training paradigm as in [21]. After training \mathcal{D}_θ in the first stage, we remove training samples with outliers in the depth estimates of the current \mathcal{D}_θ . In particular, we apply RANSAC to determine the ground plane in the front view and calculate the height of each pixel in all views. We omit training samples with $\frac{1}{H \cdot W} \sum_{u,v} [h_{u,v} < -0.5m] > \epsilon$ where $h_{u,v}$ is the height of pixel (u, v) w.r.t. the ground plane and ϵ is set to 0.005 for the front and backward-facing cameras and 0.02 for the side views. This filters frames where a significant amount of pixels are below the ground plane. In the second stage, we re-train the network from scratch on the filtered dataset for 20 epochs with the same settings. On NuScenes, we train our refinement network with all available camera images (12Hz) instead of only keyframes (2Hz).

A.3. Co-visibility Graph

We detail our multi-camera co-visibility graph construction described in Sec. 3.1 of the main paper in Algorithm 1. Note that `GetAdjacentNodes` returns a different adjacency pattern than the spatial edges in A , as described in Sec. 3.1 of the main paper. In particular, we leverage the forward motion assumption in order to connect two frames (i, j) if camera c_j is closer to the forward-facing camera than camera c_i . For further clarification, please refer to Fig. 3 of the main paper.

Dynamic alternative. We further implement a dynamic algorithm without the assumptions stated in Sec. 3.1 of the main paper. The dynamic algorithm establishes edges based on camera frustum overlap, *i.e.* it measures the intersection over union (IoU) of the camera frustums of each frame across a local time window in world space given the current camera pose estimates \mathbf{G} . We order frame pairs by their IoU in descending order and choose the N highest overlapping pairs as co-visible frames. These establish the temporal and spatial-temporal edges.

We found empirically that the static algorithm performs

Table 11. **Depth refinement network architecture.** K describes the kernel size, S the stride. ResidualBlock consists of two convolutional layers and a skip-connection as proposed in [26]. We generate output at four scales in $[0, 1]$ which is normalized by focal length of the respective camera and scaled to $[1/d_{max}, 1/d_{min}]$.

No.	Input	Layer Description	K	S	Output Size
(#A, #B) UpBlock					
#i	(#A)	Conv2d → ELU → Up	3	1	
#ii	(#i, #B)	Concatenate → Conv → ELU	3	1	
Encoder					
#0		Input: Image + Inv. Geometric Depth + Confidence			$5 \times H \times W$
#1	(#0)	Conv → BN → ReLU	7	2	$64 \times H/2 \times W/2$
#2	(#1)	MaxPool → Skip	3	2	$64 \times H/4 \times W/4$
#3	(#2)	2xResidualBlock → Skip	3	1	$64 \times H/4 \times W/4$
#4	(#3)	2xResidualBlock → Skip	3	2	$128 \times H/8 \times W/8$
#5		Input: Inv. Geometric Depth + Confidence			$2 \times H/8 \times W/8$
#6	(#3, #5)	Concatenate			$130 \times H/8 \times W/8$
#7	(#6)	2xResidualBlock → Skip	3	2	$256 \times H/16 \times W/16$
#8	(#7)	2xResidualBlock	3	2	$512 \times H/32 \times W/32$
Decoder					
#9	(#8, #7)	UpBlock	3	1	$256 \times H/16 \times W/16$
#10	(#9, #4)	UpBlock	3	1	$128 \times H/8 \times W/8$
#11	(#10)	Conv2d → Sigmoid → Output	3	1	$1 \times H/8 \times W/8$
#12	(#10, #3)	UpBlock	3	1	$64 \times H/4 \times W/4$
#13	(#12)	Conv2d → Sigmoid → Output	3	1	$1 \times H/4 \times W/4$
#14	(#12, #2)	UpBlock	3	1	$32 \times H/2 \times W/2$
#15	(#14)	Conv2d → Sigmoid → Output	3	1	$1 \times H/2 \times W/2$
#16	(#15)	UpBlock	3	1	$16 \times H \times W$
#17	(#16)	Conv2d → Sigmoid → Output	3	1	$1 \times H \times W$

Algorithm 1 Co-visibility graph construction

Input: $\mathbf{K}, \mathbf{T}, \mathcal{G}, \{\mathbf{I}_t^c\}_{c=0}^C$

- 1: $A = \text{ComputeStaticAdjacency}(\mathbf{K}, \mathbf{T})$
- 2: $N = \text{GetNodes}(\mathcal{G})$
- 3: $M = \text{AddNodes}(\mathcal{G}, \mathbf{I}_t^c)$
- 4: **for** $i \in M$ # add temporal edges
- 5: **for** $j \in N$
- 6: **if** $\text{Radius}(i, j) < r_{\text{intra}}$
- 7: $\text{AddEdge}(\mathcal{G}, i, j)$
- 8: **end if**
- 9: **end for**
- 10: **end for**
- 11: **for** $(i, j) \in A$ # add spatial edges
- 12: $\text{AddEdge}(\mathcal{G}, i, j)$
- 13: **end for** # add spatial-temporal edges
- 14: $O = \text{GetNodesAtTime}(\mathcal{G}, t - r_{\text{intra}})$
- 15: $O' = \text{GetAdjacentNodes}(\mathcal{G}, A, M, O)$
- 16: **for** $(i, j) \in O'$
- 17: $\text{AddEdge}(\mathcal{G}, i, j)$
- 18: **end for**
- 19: # Remove out-of-context edges and nodes
- 20: $\text{RemoveTemporalEdges}(\mathcal{G}, t - \Delta_{t_{\text{intra}}})$
- 21: $\text{RemoveSpatialTemporalEdges}(\mathcal{G}, t - \Delta_{t_{\text{inter}}})$
- 22: $\text{RemoveUnconnectedNodes}(\mathcal{G})$

similarly to the dynamic alternative while being simpler and more efficient, so we use it in our experiments. However, for applications where the assumptions in Sec. 3.1 of the main paper do not hold, this algorithm provides a suitable alternative.

A.4. Inference Details

For both datasets, we observe that camera shutters are not well synchronized in both datasets. This poses a problem, especially at high speeds. Thus, instead of using constant camera extrinsics, we compute time-dependent relative camera extrinsics. For inference, we set $n_{\text{itr-wm}} = 16$, $n_{\text{iter1}} = 4$ and $n_{\text{iter2}} = 2$. For Nuscenets, use a different threshold $\beta = 0.8$.

B. Experiments

Evaluation metrics. We evaluate the proposed method in terms of depth accuracy and trajectory accuracy.

Depth. Given the estimated depths \mathbf{d}_t^c and ground truth depth \mathbf{d}_t^{*c} we use compare the following depth metrics

$$\text{Abs Rel: } \frac{1}{T \cdot C} \sum_{d,c} \frac{|\mathbf{d}_t^c - \mathbf{d}_t^{*c}|}{\mathbf{d}_t^{*c}} \quad (10)$$

$$\text{Sqr Rel: } \frac{1}{T \cdot C} \sum_{d,c} \frac{\|\mathbf{d}_t^c - \mathbf{d}_t^{*c}\|^2}{\mathbf{d}_t^{*c}} \quad (11)$$

$$\text{RMSE: } \frac{1}{T \cdot C} \sqrt{\sum_{d,c} \|\mathbf{d}_t^c - \mathbf{d}_t^{*c}\|^2} \quad (12)$$

$$\delta_{1.25}: \text{fraction of } d \in \mathbf{d} \text{ for which } \max\left(\frac{d}{d^*}, \frac{d^*}{d}\right) < 1.25 \quad (13)$$

For up-to-scale evaluation, we resort to the camera-wise metric scaling as described in [23] which results in scaling factor s described as

$$s = \frac{1}{C} \cdot \sum_c \frac{\text{median}(d^{*c})}{\text{median}(d^c)} \quad (14)$$

We then scale the predicted depth as $s \cdot \mathbf{d}$.

Trajectory. For trajectory evaluation we use the absolute trajectory error (ATE) score. Given an estimated trajectory $\mathbf{P}_1, \dots, \mathbf{P}_T \in SE(3)$ and ground truth trajectory $\mathbf{Q}_1, \dots, \mathbf{Q}_T$ the ATE is defined as

$$\mathbf{F}_i = \mathbf{Q}_i^{-1} \mathbf{S} \mathbf{P}_i$$

$$\text{ATE} = \text{RMSE}(\mathbf{F}_{1:T}) = \sqrt{\frac{1}{T} \sum_i \|\text{trans}(\mathbf{F}_i)\|^2} \quad (15)$$

where trans defines the translational part of \mathbf{F} and \mathbf{S} is identity \mathbf{I} for unscaled evaluation or $s_{\text{traj}} \cdot \mathbf{I}$ where s_{traj} is determined via least squares for scaled evaluation.

Ablation studies. In Tab. 13, we show a comparison of our depth refinement network trained with synthetic data only, with both synthetic data and real-world data, and real-world data only. As stated in Sec. 3.3 of the main paper, the results show that synthetic data cannot help the depth refinement network performance, even when fine-tuning on real-world data afterward. Instead, we observe the best performance when starting training from real-world data directly. This underlines the importance of our self-supervised training scheme for the refinement network, since contrary to the geometric parts of our system, here we cannot rely on synthetic data to provide us with ground-truth supervision.

For Tab. 14, we train two versions of our refinement network described in Sec. 3.3 of the main paper. We train \mathcal{D}_θ as described in the main paper with sparsified input depth and train \mathcal{D}_ω purely for monocular depth estimation without refining geometric estimates. We evaluate both networks on monocular depth estimation on DDAD. We observe that the networks perform similarly, while \mathcal{D}_θ can both refine geometric depth estimates and estimate depth without geometric depth input. This shows that the refinement network learns strong scene priors that can estimate depth even without any additional input. Further, we can conclude that the network generalizes well to both depth prediction and depth refinement.

In Fig. 7 we show an ablation of the covisibility graph density by varying the parameters described in Sec. 3.1 of

Table 12. **Masking scheme ablation on DDAD.** We compare the influence of the different masks used to train the refinement network (cf. Eq. 8 of the main paper) on the final performance.

M^{st}	M^{oc}	M^{fc}	Abs Rel ↓	Sq Rel ↓	RMSE ↓	$\delta_{1.25}$ ↑
			0.637	51.086	18.129	0.779
✓			0.288	10.304	12.982	0.786
✓	✓		0.162	3.304	11.638	0.811
✓	✓	✓	0.162	3.019	11.408	0.811

Table 13. **Refinement network training.** We show that training and pre-training the refinement network on the synthetic VKITTI [2] dataset does not yield improvement over self-supervised training with real-world data as proposed in Sec. 3.3 of the main paper.

VKITTI	DDAD	Abs Rel ↓	Sq Rel ↓	RMSE ↓	$\delta_{1.25}$ ↑
✓		0.282	5.025	15.281	0.572
✓	✓	0.163	3.313	11.580	0.809
	✓	0.162	3.019	11.408	0.811

Table 14. **Refinement network training comparison.** We train \mathcal{D}_θ as described in the main paper with sparsified input depth and train \mathcal{D}_ω purely for monocular depth estimation without prior geometric depth input. We evaluate both networks on monocular depth estimation *without geometric depth input* on DDAD.

$\mathcal{D}_\omega(\mathbf{I}_t^c)$	$\mathcal{D}_\theta(\mathbf{I}_t^c, \mathbf{0}, \mathbf{0})$	Abs Rel ↓	Sq Rel ↓	RMSE ↓	$\delta_{1.25}$ ↑
✓		0.204	3.583	12.652	0.723
	✓	0.211	3.806	12.668	0.715

our paper. An increase in the number of edges in the covisibility graph only yields marginal improvement, while resulting in a linear increase in runtime.

During training of \mathcal{D}_θ , we mask regions that do not provide useful information for depth learning when minimizing the view synthesis loss in Eq. 8 of the paper. We provide an ablation study of the three masks used in Eq. 8 in Tab. 12. The self-occlusion M^{oc} and static M^{st} masks are essential to depth learning by removing ego-vehicle and *e.g.* sky regions, respectively. The flow consistency mask M^{fc} further reduces outliers, caused by *e.g.* dynamic objects.

Qualitative comparisons. In Fig. 9, we qualitatively compare our method to existing works in terms of scale-aware depth estimation on DDAD. Further, we show our geometric depth estimate alongside the refined depth. We notice that, especially for the side views, existing works struggle to obtain accurate depth. On the other hand, the geometric depth of our method produces many accurate depth predictions, but is at the same time very noisy, especially in low-textured areas and for dynamic objects. Our full method demonstrates the best performance.

In Fig. 10, we show a comparison of our method to the state-of-the-art approach SurroundDepth [53] on the NuScenes dataset. We observe that our approach produces significantly sharper and more accurate depth maps for all three of the examples.

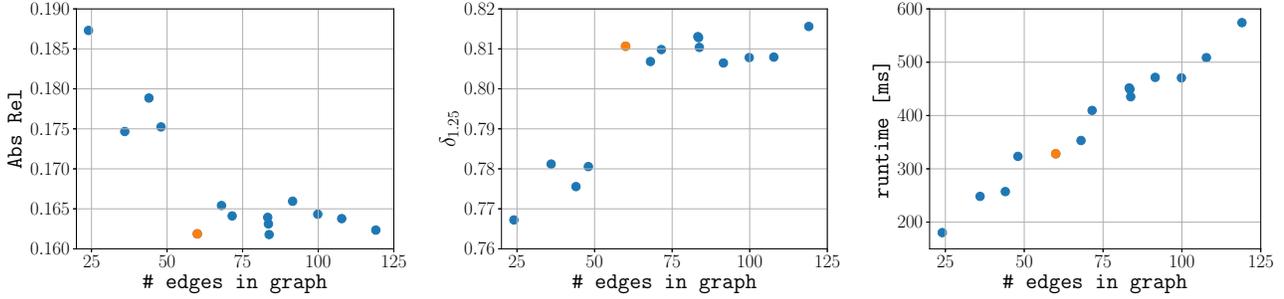


Figure 7. **Graph density analysis on DDAD.** We plot the AbsRel, $\delta_{1.25}$ scores and runtime w.r.t. the number of edges in the co-visibility graph. Our default setting is shown in orange.

Table 15. **Inference runtime analysis.** We list the runtime of each component of our system during inference.

Component	Complexity	Runtime [ms]	Percent [%]
Feature encoder	$\mathcal{O}(C)$	5	1.4
Context encoder	$\mathcal{O}(C)$	5	1.4
Create corr. volumes	$\mathcal{O}(\mathcal{E}_{new})$	13	3.7
Corr. volume sampling	$\mathcal{O}(n_{iter} \cdot \mathcal{E})$	31	8.8
GRU steps	$\mathcal{O}(n_{iter} \cdot \mathcal{E})$	196	55.4
Multi-cam. DBA steps	$\mathcal{O}(n_{iter})$	1	0.3
Completion	$\mathcal{O}(\mathcal{V})$	98	27.7
Others	-	5	1.4
Total		354	100

Finally, we show additional comparison on accumulated point clouds on DDAD in Fig. 11 and on NuScenes in Fig. 12. Our method produces significantly more consistent 3D reconstructions than competing methods, as can be seen by the more consistent reconstruction of road markings, sidewalks, and trucks in Fig. 11. In Fig. 12 we observe that our method approaches the 3D reconstruction accuracy of the LiDAR-based 3D reconstruction.

Runtime breakdown and memory consumption. In Tab. 15 we show a component-wise breakdown of time-complexity and measured runtime of our approach. Compared to [49], we tackle a more complex scenario with six images per timestep. This leads to many more possible edges in the co-visibility graph, increasing the computational burden. Thus, the runtimes are slower than reported in [49]. However, Tab. 15 shows that runtime is dominated by the GRU, which scales with the number of edges $|\mathcal{E}|$. The breakdown and our observed $10\times$ runtime improvement are both at test time. The peak GPU memory consumption in inference with our parameter setting is 6.08 GB (~ 61 MB per edge).

Limitations. The components of our system rely on deep neural networks with downsampling operations. This means a large chunk of the computation will happen at a lower resolution. While this provides a computational advantage, it also comes at the cost of losing high-frequency details that are important for thin structures like fences. In

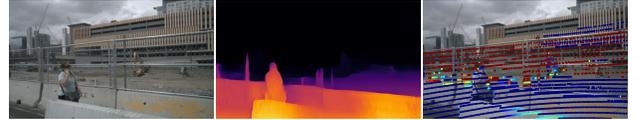


Figure 8. **Illustration of thin structures.** We depict an example frame, prediction, and error map from the NuScenes dataset. The fence in the picture poses a problem for our depth estimator since it is partially transparent. Further, since behind the fence, there is a large free space, we observe a large Sq. Rel. score for the areas that are predicted as background.

Fig. 8, we show an example of this phenomenon where our depth estimate results in a large error because there is an ambiguity between the background and the foreground fence. Similarly, other thin structures like poles could be missed, especially if they are far away.

C. Multi-Camera DBA

We provide a detailed derivation of our multi-camera DBA. The goal of the bundle adjustment is to minimize the energy function E , *i.e.* align edge-wise relative poses \mathbf{G}_{ij} and frame-wise depth \mathbf{d}_i whose reprojection minimizes the Mahalanobis distance to the estimated flow \mathbf{f}_{ij} over all edges \mathcal{E} in the co-visibility graph. In the following, let H and W be the depth map’s height and width and $T = |\mathcal{V}|/C$ the number of timesteps we consider, where C is the number of cameras.

$$E = \sum_{(i,j) \in \mathcal{E}} \|\mathbf{p}_{ij} - \Pi_{c_j}(\mathbf{G}_{ij} \circ \Pi_{c_i}^{-1}(\mathbf{d}_i))\|_{\Sigma_{ij}}^2 \quad (16)$$

With $\Sigma_{ij} = \text{diag}(\mathbf{w}_{ij})$, $\mathbf{p}_{ij} = \mathbf{x}_i + \mathbf{f}_{ij}$ and

$$\mathbf{G}_{ij} = (\mathbf{P}_{t_j} \mathbf{T}_{c_j})^{-1} \mathbf{P}_{t_i} \mathbf{T}_{c_i} \quad (17)$$

where \mathbf{T}_{c_i} and \mathbf{T}_{c_j} are known constants and \mathbf{P}_{t_i} and \mathbf{P}_{t_j} are optimization variables. This will lead to updates

$$\begin{aligned} \mathbf{P}^{(k)} &= \exp(\delta\xi) \mathbf{P}^{(k-1)} \\ \mathbf{d}^{(k)} &= \mathbf{d} + \mathbf{d}^{(k-1)} \end{aligned} \quad (18)$$

where $\delta\xi$ and $\delta\mathbf{d}$ are the solution to the normal equation

$$\mathbf{J}^\top \boldsymbol{\Sigma} \mathbf{J} \cdot \begin{bmatrix} \delta\xi \\ \delta\mathbf{d} \end{bmatrix} = -\mathbf{J}^\top \boldsymbol{\Sigma} \mathbf{r} \quad (19)$$

where $\mathbf{J} \in \mathbb{R}^{|\mathcal{E}| \cdot H \cdot W \cdot 2 \times (T \cdot 6 + |\mathcal{V}| \cdot H \cdot W)}$ is the Jacobian of the residuals w.r.t. all optimization variables and $\mathbf{r} \in \mathbb{R}^{|\mathcal{E}| \cdot H \cdot W \cdot 2}$ is the vector of all residuals

C.1. Pose - Depth Decomposition

We can make three observations. First, pose \mathbf{P}_k only appears in \mathbf{r}_{ij} if k is either i or j and $t_i \neq t_j$. Second, there are no loops, thus $i \neq j$. Third, depth \mathbf{d}_k only appears in \mathbf{r}_{ij} if $k = i$, therefore $\frac{\partial \mathbf{r}_{ij}}{\partial \mathbf{d}_k} = 0 \forall k \neq i$.

Let us now consider a single edge $(i, j) \in \mathcal{E}$ with $\mathbf{r}_{ij} \in \mathbb{R}^{H \cdot W \cdot 2}$ the residuals and $\mathbf{J}_{ij} \in \mathbb{R}^{H \cdot W \cdot 2 \times (12 + H \cdot W)}$ the Jacobian w.r.t. all optimization variables. We can decompose the Jacobian into its components, i.e. $\mathbf{J}_{ij} = [\mathbf{J}_{\xi_i} + \mathbf{J}_{\xi_j} \quad \mathbf{J}_{\mathbf{d}_i}]$ where $\mathbf{J}_{\mathbf{d}_i}$ is diagonal. Now, when only considering the aforementioned edge, Eq. 19 can be written as

$$\begin{bmatrix} \mathbf{B}_{ii} & \mathbf{B}_{ij} & \mathbf{E}_{ii} \\ \mathbf{B}_{ji} & \mathbf{B}_{jj} & \mathbf{E}_{ji} \\ \mathbf{E}_{ii}^\top & \mathbf{E}_{ji}^\top & \mathbf{C}_i \end{bmatrix} \begin{bmatrix} \delta\xi_i \\ \delta\xi_j \\ \delta\mathbf{d}_i \end{bmatrix} = \begin{bmatrix} \mathbf{v}_i \\ \mathbf{v}_j \\ \mathbf{w}_i \end{bmatrix} \quad (20)$$

with

$$\begin{aligned} \mathbf{B}_{ii} &= \mathbf{J}_{\xi_i}^\top \boldsymbol{\Sigma}_{\mathbf{r}_{ij}} \mathbf{J}_{\xi_i} & \mathbf{E}_{ii} &= \mathbf{J}_{\xi_i}^\top \boldsymbol{\Sigma}_{\mathbf{r}_{ij}} \mathbf{J}_{\mathbf{d}_i} \\ \mathbf{B}_{ij} &= \mathbf{J}_{\xi_i}^\top \boldsymbol{\Sigma}_{\mathbf{r}_{ij}} \mathbf{J}_{\xi_j} & \mathbf{E}_{ji} &= \mathbf{J}_{\xi_j}^\top \boldsymbol{\Sigma}_{\mathbf{r}_{ij}} \mathbf{J}_{\mathbf{d}_i} \\ \mathbf{B}_{ji} &= \mathbf{J}_{\xi_j}^\top \boldsymbol{\Sigma}_{\mathbf{r}_{ij}} \mathbf{J}_{\xi_i} & \mathbf{v}_i &= -\mathbf{J}_{\xi_i}^\top \boldsymbol{\Sigma}_{\mathbf{r}_{ij}} \mathbf{r}_{ij} \\ \mathbf{B}_{jj} &= \mathbf{J}_{\xi_j}^\top \boldsymbol{\Sigma}_{\mathbf{r}_{ij}} \mathbf{J}_{\xi_j} & \mathbf{v}_j &= -\mathbf{J}_{\xi_j}^\top \boldsymbol{\Sigma}_{\mathbf{r}_{ij}} \mathbf{r}_{ij} \\ \mathbf{C}_i &= \mathbf{J}_{\mathbf{d}_i}^\top \boldsymbol{\Sigma}_{\mathbf{r}_{ij}} \mathbf{J}_{\mathbf{d}_i} & \mathbf{w}_i &= -\mathbf{J}_{\mathbf{d}_i}^\top \boldsymbol{\Sigma}_{\mathbf{r}_{ij}} \mathbf{r}_{ij} \end{aligned} \quad (21)$$

We now consider again all edges \mathcal{E} . Because the energy function in Eq. 16 is the sum of the energies for all edges, we can apply the sum rule for derivatives. Thus, we can combine the components of the normal equation. Since the block matrices from Eq. 21 are dependent on their respective residual \mathbf{r}_{ij} , we can combine the blocks via a scattered sum

$$\begin{aligned} \mathbf{B} &= \sum_{(i,j) \in \mathcal{E}} \Phi_{t_i t_i}^{T \times T} [\mathbf{B}_{ii}(\mathbf{r}_{ij})] + \Phi_{t_i t_j}^{T \times T} [\mathbf{B}_{ij}(\mathbf{r}_{ij})] \\ &\quad + \Phi_{t_j t_i}^{T \times T} [\mathbf{B}_{ji}(\mathbf{r}_{ij})] + \Phi_{t_j t_j}^{T \times T} [\mathbf{B}_{jj}(\mathbf{r}_{ij})] \\ \mathbf{E} &= \sum_{(i,j) \in \mathcal{E}} \Phi_{t_i}^{T \times |\mathcal{V}|} [\mathbf{E}_{ii}(\mathbf{r}_{ij})] + \Phi_{t_j}^{T \times |\mathcal{V}|} [\mathbf{E}_{ji}(\mathbf{r}_{ij})] \\ \mathbf{C} &= \sum_{(i,j) \in \mathcal{E}} \Phi_{ii}^{|\mathcal{V}| \times |\mathcal{V}|} [\mathbf{C}_i(\mathbf{r}_{ij})] \\ \mathbf{v} &= \sum_{(i,j) \in \mathcal{E}} \Phi_{t_i}^{T \times 1} [\mathbf{v}_i(\mathbf{r}_{ij})] + \Phi_{t_j}^{T \times 1} [\mathbf{v}_j(\mathbf{r}_{ij})] \\ \mathbf{w} &= \sum_{(i,j) \in \mathcal{E}} \Phi_i^{|\mathcal{V}| \times 1} [\mathbf{w}_i(\mathbf{r}_{ij})] \end{aligned} \quad (22)$$

where $\Phi_{mn}^{M \times N} : \mathbb{R}^{U \times V} \rightarrow \mathbb{R}^{M \cdot U \times N \cdot V}$ is the function which maps a block matrix to row $m \in \{1, \dots, M\}$ and column $n \in \{1, \dots, N\}$ while the other elements are set to zero. To improve convergence, the Levenberg-Marquardt method is used on \mathbf{C} , leading to

$$\begin{bmatrix} \mathbf{B} & \mathbf{E} \\ \mathbf{E}^\top & \mathbf{C} + \lambda \mathbf{I} \end{bmatrix} \begin{bmatrix} \delta\xi \\ \delta\mathbf{d} \end{bmatrix} = \begin{bmatrix} \mathbf{v} \\ \mathbf{w} \end{bmatrix} \quad (23)$$

With \mathbf{I} the identity matrix and pixel-wise damping factor λ is predicted by the GRU (see Sec. A) for each node i . We observe that \mathbf{C} is diagonal. We can use the Schur complement to solve Eq. 23 efficiently, due to $(\mathbf{C} + \lambda \mathbf{I})^{-1}$ being very easy to invert. Thus the updates are given by

$$\begin{aligned} \mathbf{S} &= \mathbf{B} - \mathbf{E} \mathbf{C}^{-1} \mathbf{E}^\top \\ \delta\xi &= \mathbf{S}^{-1} (\mathbf{v} - \mathbf{E} \mathbf{C}^{-1} \mathbf{w}) \\ \delta\mathbf{d} &= (\mathbf{C} + \lambda \mathbf{I})^{-1} (\mathbf{w} - \mathbf{E}^\top \delta\xi) \end{aligned} \quad (24)$$

Lastly, it can be shown that $\mathbf{S} \succ 0$, thus the Cholesky-decomposition can be used to efficiently solve for \mathbf{S}^{-1} .

C.2. Jacobians

Given the decomposition above, we now define the Jacobians \mathbf{J}_{ξ_i} , \mathbf{J}_{ξ_j} , and $\mathbf{J}_{\mathbf{d}_i}$. Let us consider a single depth d of node i at location (u, v) . The residual is given by

$$\mathbf{r}_{ij,uv} = \mathbf{p}_{ij,uv} - \hat{\mathbf{p}}_{ij,uv} \in \mathbb{R}^2 \quad (25)$$

where $\hat{\mathbf{p}}_{ij,uv} = \Pi_{c_j}(\mathbf{G}_{ij} \circ \Pi_{c_i}^{-1}(d))$. We further define the 3D point corresponding to pixel (u, v) as $\mathbf{X} = [X \ Y \ Z \ W]^\top$ which is given by $\mathbf{X} = \Pi_{c_i}^{-1}(d)$ and the transformed point $\mathbf{X}' = \mathbf{G}_{ij} \mathbf{X}$. We can thus define the projection and unprojection operators for pinhole cameras c_i and c_j

$$\Pi_{c_j}(\mathbf{X}') = \begin{bmatrix} f_{c_j} \frac{X'}{Z'} + c_{c_j}^x \\ f_{c_j} \frac{Y'}{Z'} + c_{c_j}^y \end{bmatrix} \quad (26)$$

$$\Pi_{c_i}^{-1}(d) = \begin{bmatrix} \frac{u - c_c^x}{f_c^x} \\ \frac{v - c_c^y}{f_c^y} \\ 1 \\ d \end{bmatrix} \quad (27)$$

where f_c^x, f_c^y are the camera c 's focal lengths in x and y direction, c_c^x, c_c^y are the respective principal points.

Depth The Jacobian \mathbf{J}_d w.r.t. depth d is defined as

$$\begin{aligned} \mathbf{J}_d &= \frac{\partial \mathbf{r}_{ij}}{\partial d} = -\frac{\partial \hat{\mathbf{p}}_{ij}}{\partial d} = -\frac{\partial \Pi_{c_j}(\mathbf{X}')}{\partial \mathbf{X}'} \frac{\partial \mathbf{X}'}{\partial d} \\ &= -\frac{\partial \Pi_{c_j}(\mathbf{X}')}{\partial \mathbf{X}'} \mathbf{G}_{ij} \frac{\partial \Pi_{c_i}^{-1}(d)}{\partial d} \end{aligned} \quad (28)$$

$$\frac{\partial \Pi_{c_j}(\mathbf{X}')}{\partial \mathbf{X}'} = \begin{bmatrix} f_{c_j}^x \frac{1}{Z'} & 0 & -f_{c_j}^x \frac{X'}{Z'^2} & 0 \\ 0 & f_{c_j}^y \frac{1}{Z'} & -f_{c_j}^y \frac{Y'}{Z'^2} & 0 \end{bmatrix} \quad (29)$$

$$\frac{\partial \Pi_{c_i}^{-1}(d)}{\partial d} = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \end{bmatrix} \quad (30)$$

Pose The Jacobian \mathbf{J}_ξ w.r.t. ξ where $\xi \in \mathfrak{se}(3)$ is either ξ_i or ξ_j .

$$\mathbf{J}_\xi = \frac{\partial \mathbf{r}_{ij}}{\partial \xi} = -\frac{\partial \hat{\mathbf{p}}_{ij}}{\partial \xi} = -\frac{\partial \Pi_{c_j}(\mathbf{X}')}{\partial \mathbf{X}'} \frac{\partial \mathbf{X}'}{\partial \xi} \quad (31)$$

The partial derivative $\frac{\partial \Pi_{c_j}(\mathbf{X}')}{\partial \mathbf{X}'}$ has been derived in Eq. 29. For $\frac{\partial \mathbf{G}_{ij}}{\partial \xi}$, we can again decompose \mathbf{G}_{ij} into the static parts \mathbf{T}_{c_i} and \mathbf{T}_{c_j} and the unknown, to be optimized parts, \mathbf{P}_{t_i} and \mathbf{P}_{t_j}

$$\mathbf{X}' = (\mathbf{P}_{t_j} \mathbf{T}_{c_j})^{-1} \mathbf{P}_{c_i} \mathbf{T}_{c_i} \mathbf{X} \quad (32)$$

First, similar to Eq. 39 to 44 in [11], for $\mathbf{A} = \mathbf{A}_1(\mathbf{A}_0)^{-1}$ we can write

$$\begin{aligned} \frac{\partial \mathbf{A}}{\partial \mathbf{A}_0} &= \frac{\partial \log(\mathbf{A}_1(\exp(\xi)\mathbf{A}_0)^{-1}(\mathbf{A}_1\mathbf{A}_0^{-1})^{-1})}{\partial \xi} \Big|_{\xi=0} \\ &= \frac{\partial}{\partial \xi} \Big|_{\xi=0} [\log(\mathbf{A}_1\mathbf{A}_0^{-1} \exp(-\xi)\mathbf{A}_0\mathbf{A}_1^{-1})] \\ &= \frac{\partial}{\partial \xi} \Big|_{\xi=0} [\log(\exp(-\text{Adj}_{\mathbf{A}_1\mathbf{A}_0^{-1}} \xi)\mathbf{A}_1\mathbf{A}_0^{-1}\mathbf{A}_0\mathbf{A}_1^{-1})] \\ &= \frac{\partial}{\partial \xi} \Big|_{\xi=0} [\log(\exp(-\text{Adj}_{\mathbf{A}_1\mathbf{A}_0^{-1}} \xi))] \\ &= \frac{\partial}{\partial \xi} \Big|_{\xi=0} [-\text{Adj}_{\mathbf{A}_1\mathbf{A}_0^{-1}} \xi] \\ &= -\text{Adj}_{\mathbf{A}_1\mathbf{A}_0^{-1}} \end{aligned} \quad (33)$$

In the following, we omit the explicit perturbation around $\xi = 0$. Let now $\mathfrak{G}_1, \dots, \mathfrak{G}_6 \in \mathbb{R}^{4 \times 4}$ be generators as defined in Eq. 65 in [11]. With the chain rule and Eq. 94 in [11] and Eq. 33 above, the derivative w.r.t. the pose of incoming node j is given by

$$\begin{aligned} \frac{\partial \mathbf{X}'}{\partial \xi_j} &= \frac{\partial}{\partial \xi_j} [(\exp(\xi_j)\mathbf{P}_{t_j} \mathbf{T}_{c_j})^{-1} \mathbf{P}_{t_i} \mathbf{T}_{c_i} \mathbf{X}] \\ &= \frac{\partial}{\partial \xi_j} \underbrace{[\mathbf{T}_{c_j}^{-1}(\exp(\xi_j)\mathbf{P}_{t_j})^{-1}]}_{\mathbf{A}} \underbrace{[\mathbf{P}_{t_i} \mathbf{T}_{c_i} \mathbf{X}]}_{\mathbf{X}^w} \\ &= \frac{\partial}{\partial \mathbf{A}} [\mathbf{A}\mathbf{X}^w] \cdot \frac{\partial}{\partial \xi_j} \underbrace{[\mathbf{T}_{c_j}^{-1}]}_{\mathbf{A}_1} \underbrace{[(\exp(\xi_j)\mathbf{P}_{t_j})^{-1}]}_{\mathbf{A}_0} \\ &= \frac{\partial}{\partial \mathbf{A}} [\mathbf{A}\mathbf{X}^w] \cdot \frac{\partial}{\partial \mathbf{A}_0} [\mathbf{A}_1\mathbf{A}_0^{-1}] \\ &= [\mathfrak{G}_1\mathbf{X}^w \quad \dots \quad \mathfrak{G}_6\mathbf{X}^w] \cdot (-\text{Adj}_{\mathbf{A}_1\mathbf{A}_0^{-1}}) \\ &= -[\mathfrak{G}_1\mathbf{X}^w \quad \dots \quad \mathfrak{G}_6\mathbf{X}^w] \cdot \text{Adj}_{\mathbf{T}_{c_j}^{-1}\mathbf{P}_{t_j}^{-1}} \end{aligned} \quad (34)$$

Finally, with Eq. 94 and 97 in [11] and the derivative w.r.t. the pose of the outgoing node i

$$\begin{aligned} \frac{\partial \mathbf{X}'}{\partial \xi_i} &= \frac{\partial}{\partial \xi_i} \underbrace{[(\mathbf{P}_{t_j} \mathbf{T}_{c_j})^{-1}]}_{\mathbf{A}} \exp(\xi_i) \underbrace{[\mathbf{P}_{t_i} \mathbf{T}_{c_i} \mathbf{X}]}_{\mathbf{X}^{t_i}} \\ &= \frac{\partial}{\partial \mathbf{A}} [\mathbf{A}\mathbf{X}^{t_i}] \cdot \frac{\partial}{\partial \xi_i} [(\mathbf{P}_{t_j} \mathbf{T}_{c_j})^{-1} \exp(\xi_i)\mathbf{P}_{t_i}] \\ &= [\mathfrak{G}_1\mathbf{X}^{t_i} \quad \dots \quad \mathfrak{G}_6\mathbf{X}^{t_i}] \cdot \text{Adj}_{\mathbf{T}_{c_j}^{-1}\mathbf{P}_{t_j}^{-1}} \end{aligned} \quad (35)$$

Finally, we compose the component-wise Jacobians above into the Jacobian \mathbf{J} as stated in Sec. C.1. We can now use \mathbf{J} to compute the updates $\delta\xi$ and δd in Eq. 19.

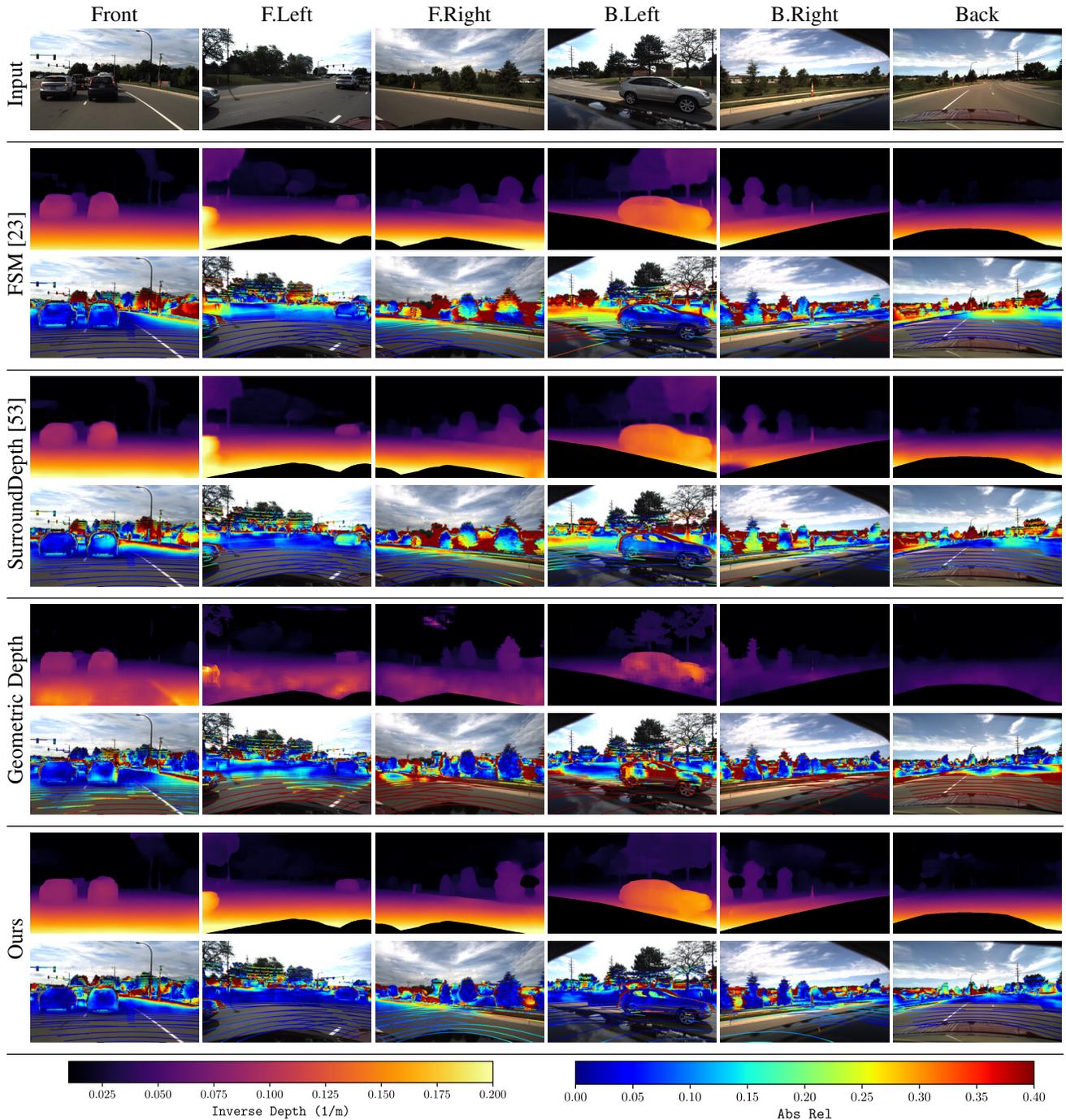


Figure 9. **Qualitative comparison on DDAD.** We compare existing works to both our geometric and refined depth estimates. Especially for the side views, existing works struggle to obtain accurate depth. The geometric depth produces many accurate depth predictions, but contains many noisy points, especially in low-textured areas and for dynamic objects. Our full method demonstrates the best performance.

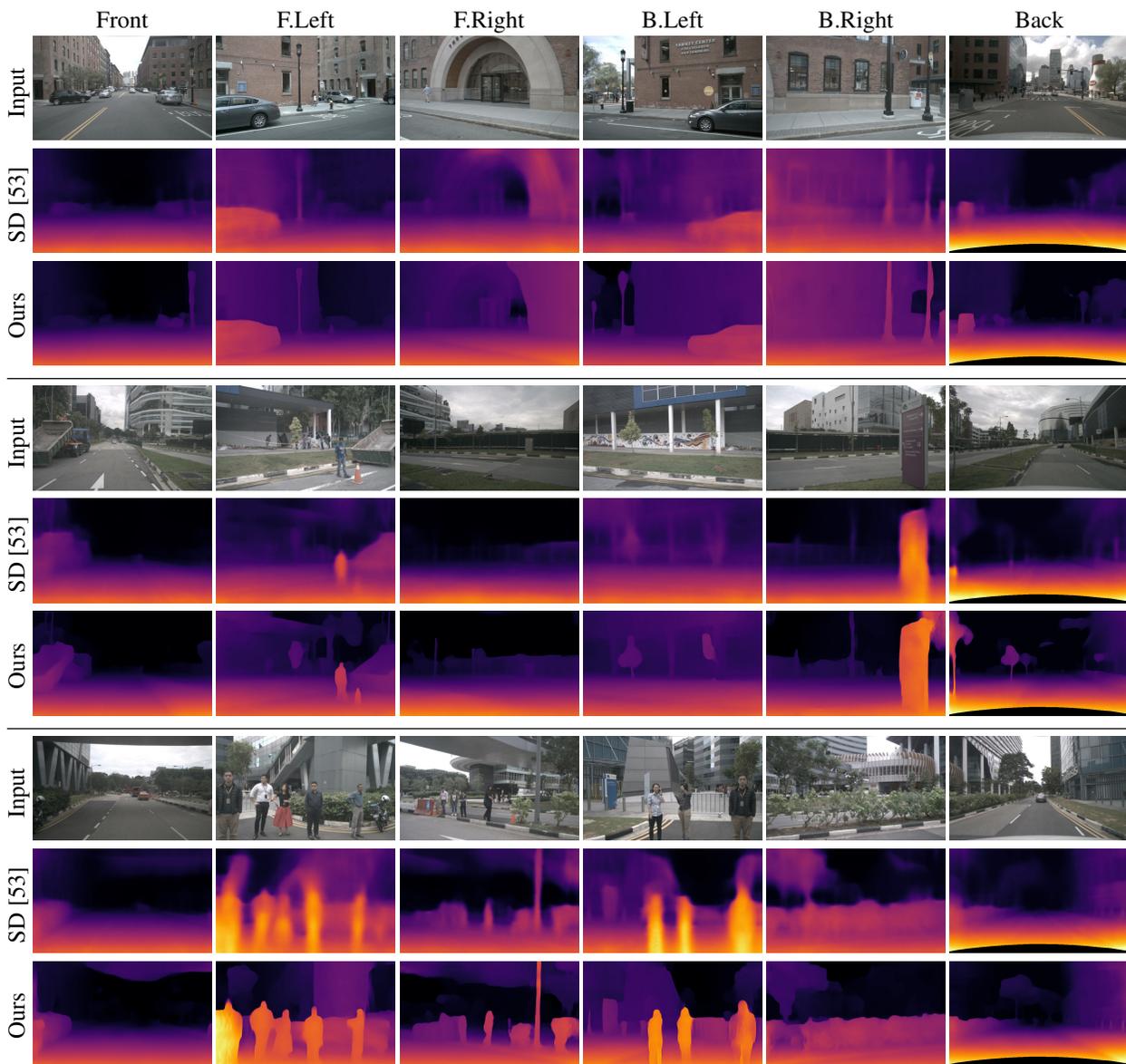


Figure 10. **Qualitative comparison on NuScenes.** We show a comparison of depth maps from our method to the depth maps of the state-of-the-art approach SurroundDepth [53]. We observe that our approach produces significantly sharper and more accurate depth predictions.

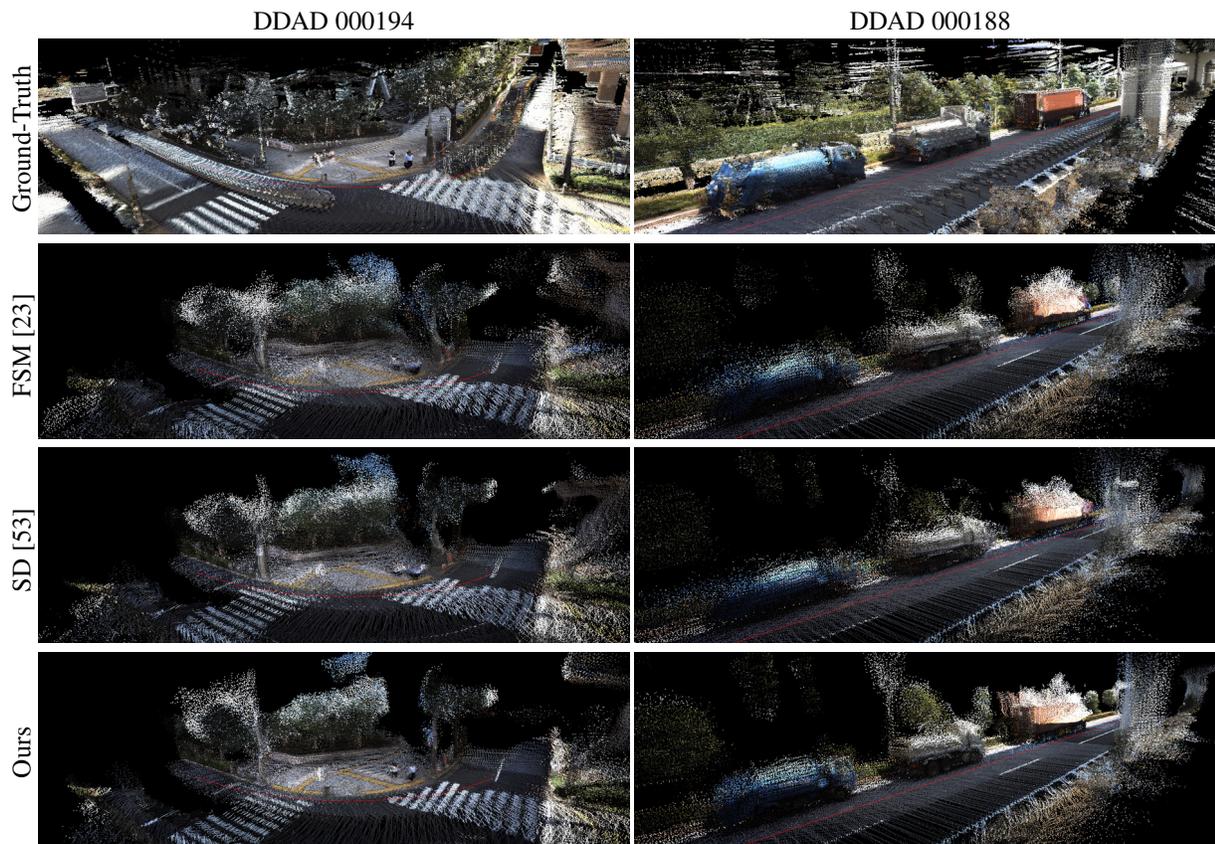


Figure 11. **Qualitative comparison of 3D reconstructions on DDAD.** We show the 3D reconstructions of our method compared to SurroundDepth [53] and FSM [23]. Additionally, we plot the ground-truth LiDAR 3D reconstruction at the top. The ego-vehicle trajectory is marked in red. We observe that our method produces significantly more consistent and accurate 3D reconstructions than competing methods, as can be seen when focusing on the trucks (right) and the street markings and pedestrians (left).



Figure 12. **Qualitative 3D reconstruction results on NuScenes.** We show our 3D reconstruction results alongside the LiDAR ground-truth 3D reconstruction. The ego-vehicle trajectory is marked in red. We observe that our method yields similarly consistent 3D reconstruction results as the LiDAR ground-truth while being denser.