

Weakly Supervised Learning of Rigid 3D Scene Flow

Zan Gojcic^{1,2}Or Litany^{2,3}Andreas Wieser¹Leonidas J. Guibas²Tolga Birdal²¹ETH Zurich²Stanford University³NVIDIA

3dsceneflow.github.io

Abstract

We propose a data-driven scene flow estimation algorithm exploiting the observation that many 3D scenes can be explained by a collection of agents moving as rigid bodies. At the core of our method lies a deep architecture able to reason at the **object-level** by considering 3D scene flow in conjunction with other 3D tasks. This object level abstraction, enables us to relax the requirement for dense scene flow supervision with simpler binary background segmentation mask and ego-motion annotations. Our mild supervision requirements make our method well suited for recently released massive data collections for autonomous driving, which do not contain dense scene flow annotations. As output, our model provides low-level cues like pointwise flow and higher-level cues such as holistic scene understanding at the level of rigid objects. We further propose a test-time optimization refining the predicted rigid scene flow. We showcase the effectiveness and generalization capacity of our method on four different autonomous driving datasets. We release our source code and pre-trained models under github.com/zgojcic/Rigid3DSceneFlow.

1. Introduction

Understanding dynamic 3D environments is a core challenge in computer vision and robotics. In particular, applications such as self driving and robot navigation rely upon a robust perception of dynamically changing 3D scenes. To equip autonomous agents with the ability to infer spatiotemporal geometric properties, there has recently been an increased interest in *3D scene flow* as a form of low-level dynamic scene representation [40, 72, 78, 55, 52, 58]. Scene flow is the 3D motion field of points in the scene [74] and is a generalization of 2D optical flow. In fact, optical flow [3, 26] can be understood as the projection of the scene flow onto a camera image plane [16]. Such dense motion fields can serve bottom-up approaches for high-level dynamic scene understanding tasks like semantic segmentation [41] or motion perception. However, representing dynamics via a free form velocity field has two major disadvantages. First, in most applications of interest, dynamics are attributed to rigid object motion [11, 45]. This notion

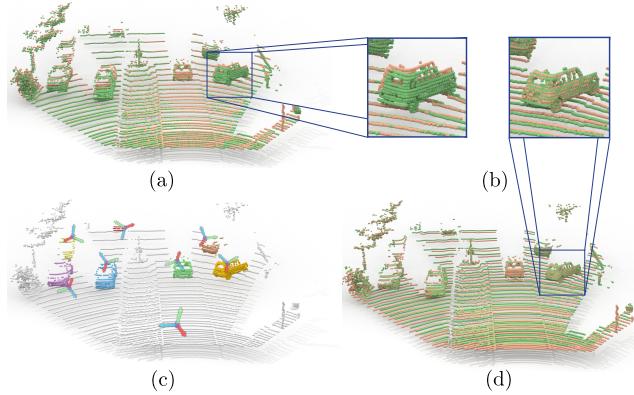


Figure 1: Our network takes two successive frames as input (a), and outputs a set of transformation parameters for each segmented rigid agent (c) which are used to recover per-point rigid scene flow. After applying the predicted flow to the first point cloud, the two frames are aligned (b, d).

has been extensively exploited in robotics [10, 15, 9, 8] and holds especially for vehicles in autonomous driving. Predicting unconstrained per-point flow may lead to non-viable results, e.g. parts of the same car might move in different directions. Second, accurately learning direct flow estimation necessitates dense supervision that is expensive to acquire and prone to annotation errors. As a result, many methods have resorted to training on simulated data [44, 78, 55], yet this comes at the price of a non-negligible domain gap. Other methods have attempted to solve the problem in a completely unsupervised manner [72, 80, 47], however they fail to provide competitive performance. In Fig. 2 we illustrate these two extremes of no- and full-supervision while spanning the many intermediate possibilities, sorted according to annotation effort¹.

Based on this observation, we seek a sweet spot between supervision effort and performance. To this end, we propose a scene abstraction approach that uses rigid objects as the basic components. More specifically, by splitting the scene into foreground (movable objects) and background (static objects), we explain the background (BG) flow as

¹Note that no prior work exists at different points on the spectrum given in Fig. 2. This leaves ample room for future exploration.

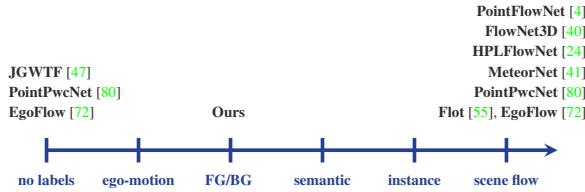


Figure 2: Recent scene flow methods either use full supervision (and suffer from domain gap) or no-supervision (and suffer from reduced performance). Instead, our method uses weak supervision and benefits from the best of both worlds.

the sensor ego-motion and the foreground (FG) flow as clusters of rigidly moving entities. As a result, we tackle both aforementioned challenges at the same time: (1) we enforce a rigidity constraint to get meaningful and more accurate (foreground and background) flow, (2) **we can relax the requirement for dense flow supervision with a much simpler binary mask annotation and ego-motion that can often be extracted directly from the agent’s IMU**. The result is a weakly supervised method for accurate flow estimation that, unlike completely unsupervised approaches, outperforms previous state of the art (SoTA) by a significant margin. For example, reducing the end-point-error on the *lidarKITTI* dataset by more than 30 cm relative to the SoTA. At the same time our result provides an interpretable and readily usable object-level scene representation. In brief, our contributions are:

- We exploit the geometry of the rigid scene flow problem to introduce an inductive bias into our network. This allows us to learn from weak supervision signals: background masks and ego-motion.
- Our data-driven method decomposes the scene into rigidly moving agents enabling us to reason on the level of objects rather than points. We use this notion to propose a new test-time optimization, further refining the flow predictions.
- Our method is backed by a novel, flexible, scene flow backbone, which can be adapted to solve various tasks.

As a result of these contributions, our method greatly outperforms the SoTA on several benchmarks: *FT3D* [44], *stereoKITTI* [45], *lidarKITTI* [22], and *semanticKITTI* [5], while generalizing to the *waymo-open* dataset [69] without additional fine-tuning.

2. Related Work

Data Driven 3D Scene Flow. While there is extensive literature on traditional 3D scene flow [74, 29, 79, 33, 70, 73, 51, 65, 66, 31, 19, 16, 11, 54], we focus our attention to recent data-driven methods that emerged based on advances in deep learning on unordered point sets [56, 88, 1].

Early methods for 3D scene flow estimation mimicked their 2D counterparts. SceneFlowNet [44] used 2D op-

tical flow and disparity maps to estimate the 3D scene flow. FlowNet3D [40] successfully adopted the ideas from FlowNet [30, 18]. FlowNet3D++ [78] extended FlowNet3D to incorporate additional geometric constraints. MeteorNet [41] used multiple temporally ordered frames to improve the accuracy of the inferred flow. Wang *et al.* [76] incorporated a continuous convolution into a 3D-FCN [71] to undo both the ego-motion and object-motion in two consecutive LiDAR frames. PointRNN [21] used recurrent neural networks to model temporal point sets, which yields 3D scene flow as a by-product. HPLFlowNet [24] ordered the points into a permutohedral lattice of SplatNet [68] to apply bilateral convolution layers. This allowed efficient and robust non-rigid 3D flow computation. Both OcCFlow [52] and CaSPR [58] introduced spatio-temporal representations to continuously and densely estimate the scene flow. Mustafa and Hilton used semantic coherence between multiple frames to improve 4D scene flow estimation, co-segmentation and reconstruction [49]. Mittal *et al.* [47] and PointPWCNet [80] proposed self-supervised losses to infer the scene flow in an end-to-end manner. Finally, FLOT [55] proposed a simple correspondence-based end-to-end scene flow network. While our backbone also estimates correspondences, decomposing the scene into rigid agents provides us further higher level scene understanding and enables test-time optimization, while requiring less supervision.

Local Rigidity and Multi-body Motion. Flow estimation has also been tackled by imposing physical priors such as multi-body rigidity. Initial attempts involved factorization [13, 37] to separate independently moving objects. Golyanik *et al.* [23] used rigidity constraints on over segmentation of RGBD-frames. [48] used detection & tracking to constrain the flow. Vogel *et al.* [75] modeled scene flow using piecewise rigidly moving planar patches. Deewan *et al.* [17] used 3D descriptors to enforce local geometric constancy on a factor graph. GraphFlow [2] and SphereFlow [27] considered large motions and relied on sparse keypoints that are not repeatable in 3D [61, 82]. Similar to us, Jamiez *et al.* [32] as well as recent Dynamic-SLAM pipelines [28, 62, 67] assumed that clustering would yield motion segmentation. In fact, MaskFusion [60] and EMFusion [67] explicitly used Mask-RCNN [25] to this end.

On a data driven front, considering the rigidity in the scenes [42], Ma *et al.* [43] made use of depth and flow estimates from a stereo RGB-D setup within an optimization framework to obtain the 3D motion of each instance. This method relied upon a given instance segmentation, a difficult problem to solve even for the SoTA approaches [77, 87, 83]. Based on VoxelNet [92], PointFlowNet [4] jointly predicted 3D scene flow, bounding boxes, and rigid motion of objects in the scene. Yi *et al.* [86] used a PointNet++ [57] based flow estimator for piecewise rigid 3D part induction.

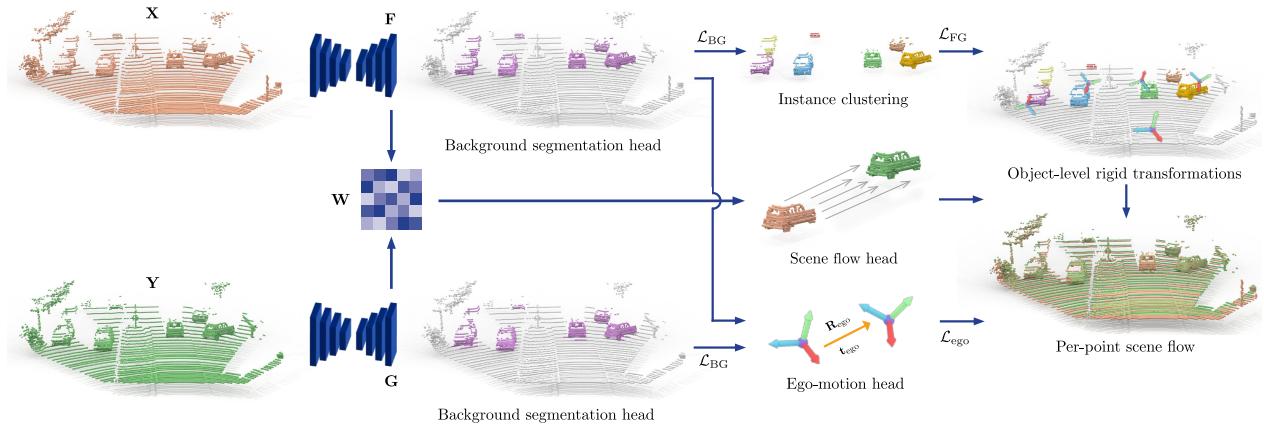


Figure 3: Architecture of our weakly-supervised scene flow estimation pipeline. Our module consumes point clouds \mathbf{X} and \mathbf{Y} of two consecutive frames and estimates per-object transformation parameters $\{\mathbf{T}_k\}_{k=1}^{K-1}$, ego-motion \mathbf{T}_{ego} , and object masks $\{\mathbf{z}_k\}_{k=1}^K$. These outputs can be combined into an **object-level scene abstraction** and **pointwise rigid scene flow**.

3. Method

Problem setting. Suppose that we observe a pair of 3D scenes \mathbf{X} and \mathbf{Y} acquired by a *single moving observer* in two consecutive instants t_0 and t_1 , respectively. Here, $\mathbf{X} \in \mathbb{R}^{3 \times N} = \{\mathbf{x}_i \in \mathbb{R}^3\}_i$ denotes a point cloud (so does \mathbf{Y}) and $\mathbf{V} \in \mathbb{R}^{3 \times N} = \{\mathbf{v}_i \in \mathbb{R}^3\}_i$ its corresponding vector field in 3D s.t. $\mathbf{X} + \mathbf{V} \approx \mathbf{Y}$ relates frame t_0 to t_1 . We further assume that \mathbf{X} , \mathbf{Y} , and also \mathbf{V} are *multi-body* i.e. composed of multiple objects. Hence, \mathbf{V} can be clustered into K objects $\mathcal{V} = \{\mathbf{V}_k \in \mathbb{R}^{3 \times N_k}\}_{k=1}^K$ each of which follows rigid dynamics, i.e. \mathbf{V} can be summarized by a set of K rigid transformations $\mathcal{T} \equiv \{\mathbf{T}_k \in SE(3)\}_{k=1}^K$ such that $\mathbf{V} \approx \{\mathbf{T}_k \circ \mathbf{X}_k - \mathbf{X}_k\}_{k=1}^K$ where $K \ll N$ and:

$$SE(3) = \left\{ \mathbf{T} \in \mathbb{R}^{4 \times 4} : \mathbf{T} = \begin{bmatrix} \mathbf{R} & \mathbf{t} \\ \mathbf{0}^\top & 1 \end{bmatrix} \right\}, \quad (1)$$

$\mathbf{R} \in SO(3)$, $\mathbf{t} \in \mathbb{R}^3$.² The motion of the immobile *background*, determines the *ego-motion* $\mathbf{T}_{ego} \subset \mathcal{T}$.

Summary. We refrain from directly predicting unconstrained pointwise flow vectors and rather aim to estimate $\mathcal{T} \equiv \{\mathbf{T}_k\}_{k=1}^K$ for all rigid bodies from which the entire scene flow $\mathcal{V} = \{\mathbf{V}_k\}_{k=1}^K$ can be recovered. To this end, we propose to *learn* the task of rigid flow estimation by solving an optimization problem composed of a set of loss functions as illustrated in Fig. 3. To reason on the level of objects, we use both instance masking and motion. To obtain the masks, we use a FG / BG prediction module in conjunction with an FG clustering. To estimate ego-motion, we run a differentiable registration on the BGs extracted from both point sets. The motions of the individual objects in the FGs are obtained similarly under the assumption of local-rigidity. We

²We denote the action of \mathbf{T} as $\mathbf{X}' = \mathbf{T} \circ \mathbf{X}$ and $\hat{\mathbf{X}}' = \mathbf{T} \hat{\mathbf{X}}$ where $\hat{\mathbf{X}} \in \mathbb{R}^{4 \times N}$ is the *homogenized* \mathbf{X} .

avoid flow-level or instance-level supervision altogether and only assume the availability of the *binary FG/BG annotations* and *ego-motion information* as a much weaker supervision signal than dense scene flow. In the sequel, we first describe our formulation of the individual objectives (§ 3.1) before proceeding to our network architecture (§ 3.2). Additional details are available in our supplement.

3.1. Energy Formulation

Our solution to the 3D scene flow estimation is attained as the minimum of a non-convex energy composed of a BG segmentation loss \mathcal{L}_{BG} , an ego-motion loss \mathcal{L}_{ego} , and an FG loss \mathcal{L}_{FG} :

$$\boldsymbol{\Gamma}^* = \operatorname{argmin}_{\boldsymbol{\Gamma}} \mathcal{L}_{BG} + \mathcal{L}_{ego} + \mathcal{L}_{FG} \quad (2)$$

where the optimal rigid scene flow $(\mathcal{V}^*, \mathcal{T}^*)$ results from the output of a *deep neural network* φ with learnable parameters $\boldsymbol{\Gamma}$: $(\mathcal{V}^*, \mathcal{T}^*) = \varphi_{\boldsymbol{\Gamma}^*}(\mathbf{X}, \mathbf{Y})$. Next, we detail the individual loss terms; each involves an *unknown* or *latent variable* obtained as a network prediction as specified in § 3.2.

Background segmentation error (\mathcal{L}_{BG}). To decompose the scene into agents that move as rigid bodies, we follow a coarse-to-fine approach. In the first step we aim to split the background and foreground points, where the foreground represents all points belonging to the movable objects (e.g., cars, cyclists, people, ...). In order to learn this binary segmentation of a point cloud, we minimize the loss $\mathcal{L}_{BG} = \frac{1}{2}(\mathcal{L}_{BG}^{\mathbf{X}} + \mathcal{L}_{BG}^{\mathbf{Y}})$ where:

$$\mathcal{L}_{BG}^{\mathbf{X}} = \frac{1}{N} \sum_{i=1}^N \text{BCE}(h_i^{\mathbf{X}}, \bar{h}_i^{\mathbf{X}}). \quad (3)$$

$\bar{h}^{\mathbf{X}}, \bar{h}^{\mathbf{Y}}$ denote the GT binary masks of point clouds \mathbf{X} and \mathbf{Y} , respectively. $h^{\mathbf{X}} = \{h_i^{\mathbf{X}}\}_{i=1}^N$ and $h^{\mathbf{Y}} = \{h_i^{\mathbf{Y}}\}_{i=1}^N$ are

the inferred foreground probabilities of points in \mathbf{X} and \mathbf{Y} yet to be clarified in § 3.2, and $\text{BCE}(h_i, \bar{h}_i) = \bar{h}_i \log(h_i) + (1 - \bar{h}_i) \log(1 - h_i)$.

Ego-motion error (\mathcal{L}_{ego}). The scene flow of the physically static background can be fully explained by its transformation parameters, the ego-motion. We estimate these parameters by first extracting the background points of both the source and target point cloud³. To reduce the computational complexity, we then randomly sample $N^b = 1024$ points therefrom such that $\mathbf{X}^b \in \mathbb{R}^{3 \times N^b} \subset \mathbf{X}$ and $\mathbf{Y}^b \in \mathbb{R}^{3 \times N^b} \subset \mathbf{Y}$. The goal of ego-motion estimation is to compute the optimal $\mathbf{R}_{\text{ego}} \in \text{SO}(3)$ and $\mathbf{t}_{\text{ego}} \in \mathbb{R}^3$ in the weighted least-squares sense

$$\mathbf{R}_{\text{ego}}^*, \mathbf{t}_{\text{ego}}^* = \underset{\mathbf{R}_{\text{ego}}, \mathbf{t}_{\text{ego}}}{\operatorname{argmin}} \sum_{l=1}^{N^b} w_l \|\mathbf{R}_{\text{ego}} \mathbf{x}_l^b + \mathbf{t}_{\text{ego}} - \phi(\mathbf{x}_l^b, \mathbf{Y}^b)\|^2,$$

where $\phi(\mathbf{x}, \mathbf{Y})$ is a *soft assignment function* returning a point from \mathbf{Y} that corresponds to $\mathbf{x} \in \mathbf{X}$. The weights w_l will be clarified below.

We approximate the optimal soft assignment ϕ via the entropy-regularized Sinkhorn algorithm [63, 14, 84]. To this end, we momentarily assume the availability of an affinity matrix $\mathbf{M} \in \mathbb{R}^{N^b \times N^b}$ describing the similarity of the background points in \mathbf{X}^b and \mathbf{Y}^b . Prediction of the currently unknown \mathbf{M} will be made precise in § 3.2. Given \mathbf{M} , we perform an alternating row and column normalization on it for $k_S \triangleq 3$ iterations, which yields $\mathbf{A} \in \mathbb{R}_+^{N^b \times N^b}$, a *doubly stochastic* (DS) assignment matrix.

In practice, due to the occlusions and sampling pattern, not all background points will have correspondences. We therefore add a slack row and column to \mathbf{M} (hence to \mathbf{A}), which enable down-weighting the outliers, while still returning a DS-matrix. The soft correspondence function then reads $\phi(\mathbf{x}_i^b, \mathbf{Y}^b) = \mathbf{Y}^b \mathbf{a}_i / \|\mathbf{a}_i\|_1$ where \mathbf{a}_i is the i -th column of \mathbf{A} after removing the slack row. Given the correspondences, the ego-motion can be recovered in closed-form using a (differentiable) *weighted* Kabsch algorithm [36] where the weights w_i are obtained as the total contribution $w_i = \sum_{j=1}^{N^b} a_{ij}$. Our ego-motion penalty measures the l_1 -discrepancy between the points transformed with the estimated $(\mathbf{R}_{\text{ego}}, \mathbf{t}_{\text{ego}})$ and the GT parameters $(\bar{\mathbf{R}}_{\text{ego}}, \bar{\mathbf{t}}_{\text{ego}})$:

$$\mathcal{L}_{\text{trans}} = \frac{1}{B} \sum_{i=1}^B \|(\bar{\mathbf{R}}_{\text{ego}} \mathbf{x}_i^b + \bar{\mathbf{t}}_{\text{ego}}) - (\mathbf{R}_{\text{ego}} \mathbf{x}_i^b + \mathbf{t}_{\text{ego}})\|_1,$$

where B denotes the number of all background points in point cloud \mathbf{X} . To stabilize the training, we further add a regularizer loss that discourages the assignment of large

³During training we use the GT BG segmentation mask $1 - \bar{\mathbf{h}}$, while during inference we threshold the inferred FG probabilities $1 - \mathbf{h}$.

values to the slack rows and columns [84]:

$$\mathcal{L}_{\text{inlier}} = \frac{1}{N^b} \sum_{i=1}^{N^b} \left(1 - \sum_{j=1}^{N^b} a_{ij} \right) + \frac{1}{N^b} \sum_{j=1}^{N^b} \left(1 - \sum_{i=1}^{N^b} a_{ij} \right).$$

The overall ego-motion loss is then the weighted sum:

$$\mathcal{L}_{\text{ego}} = \mathcal{L}_{\text{trans}} + \lambda_{\text{inlier}} \mathcal{L}_{\text{inlier}} \quad (4)$$

where $\lambda_{\text{inlier}} := 0.005$ in all our experiments.

FG instance-level rigidity error (\mathcal{L}_{FG}). To define our *per-instance rigidity loss*, we assume the availability of the following entities: (i) \mathbf{X}^f , the points of the source frame belonging to the FG; (ii) \mathbf{V}^f , the flow vectors associated to \mathbf{X}^f ; and (iii) foreground clusters $\mathcal{C} = \{\mathbf{C}^k \in \mathbb{R}^{3 \times N_k} = \{\mathbf{c}_j^k \in \mathbb{R}^3\}_{j=1}^{N^c}\}_{k=1}^{N^c}$ aggregating the individual rigid entities. (i) is a by-product of BG segmentation, *i.e.* during training, the indices of these points are obtained from the GT mask and during inference by thresholding the *inferred* FG probabilities. The flow vectors in (ii) are the result of the *scene flow module* (see § 3.2). Finally, (iii) is computed by a simple DBSCAN clustering [20] of the 3D coordinates in \mathbf{X}^f , which is based on the hypothesis that the foreground objects scattered across the scene are naturally separated by *void space* [35]. We refrain from using a data driven instance segmentation module, because our simple approach alleviates the need for instance segmentation labels.

Our rigidity loss $\mathcal{L}_{\text{rigid}}$ encourages the predicted flow vectors \mathbf{V}_k^f of each cluster k to be *congruent*, *i.e.* \mathbf{V}_k^f can be well approximated by a rigid transformation \mathbf{T}_k composed of the rotation \mathbf{R}_k and translation \mathbf{t}_k :

$$\mathcal{L}_{\text{rigid}} = \frac{1}{N^c} \sum_{k=1}^{N^c} \frac{1}{N^k} \sum_{j=1}^{N^k} \|\mathbf{R}_k \mathbf{c}_j^k + \mathbf{t}_k - (\mathbf{c}_j^k + \mathbf{v}_j^k)\|_1 \quad (5)$$

The supervision signals \mathbf{R}_k and \mathbf{t}_k are computed on the fly, such that they best explain the underlying flow \mathbf{V}^f :

$$\mathbf{T}_k^* = \underset{\mathbf{T}_k}{\operatorname{argmin}} \|\mathbf{T}_k \circ \mathbf{C}_k - (\mathbf{C}_k + \mathbf{V}_k^f)\|. \quad (6)$$

We solve Eq (6) for each individual cluster once again using the Kabsch algorithm [36]. We additionally complement the per-cluster rigidity objective with a two way Chamfer distance (CD) computed across all the foreground points:

$$\mathcal{L}_{\text{CD}} = \sum_{\mathbf{x} \in \mathbf{X}_v^f} \min_{\mathbf{y} \in \mathbf{Y}^f} \|\mathbf{x} - \mathbf{y}\|_2 + \sum_{\mathbf{y} \in \mathbf{Y}^f} \min_{\mathbf{x} \in \mathbf{X}_v^f} \|\mathbf{x} - \mathbf{y}\|_2 \quad (7)$$

where $\mathbf{X}_v^f := \mathbf{X}^f + \mathbf{V}^f$. The overall FG loss is then a weighted sum of the above objectives:

$$\mathcal{L}_{\text{FG}} = \mathcal{L}_{\text{rigid}} + \lambda_{\text{CD}} \mathcal{L}_{\text{CD}} \quad (8)$$

where $\lambda_{\text{CD}} := 0.5$ in all our experiments.

3.2. Network implementation

Provided a large dataset with *FG-BG mask* annotations as well as *ego-motion*, we learn to minimize Eq (2) using a deep neural network φ_Γ as shown in Fig. 3. In the sequel, we describe the individual modules of our network and the full inference of rigid scene flow. We refer the reader to the supplement for more details.

Backbone. Our formulation (§ 3.1) involves the estimation of different entities, requiring our network to solve *multiple tasks* similar to [34]. While it would be possible to deploy a specialized network for each task, this would increase the memory footprint and would not encourage tasks to reinforce each other [90, 89]. Instead, we propose a flexible backbone suited for solving multiple tasks through specialized heads. Specifically, our backbone network is based on Minkowski-Net [12] and follows a U-Net [59]-like encoder-decoder architecture with skip connections. Its input is a *sparsely voxelized point cloud* $\mathbf{X}^v \in \mathbb{R}^{3 \times N^v}$ and its outputs are per-point latent features $\mathbf{F}^v \in \mathbb{R}^{64 \times N^v}$. The same backbone with shared weights is also applied to $\mathbf{Y}^v \in \mathbb{R}^{3 \times M^v}$ to obtain the latent features $\mathbf{G}^v \in \mathbb{R}^{64 \times M^v}$. In the following, we omit the superscript v for clarity and unless specified differently, use \mathbf{X} and \mathbf{Y} to refer to the voxelized point clouds and \mathbf{F} and \mathbf{G} to refer to their associated latent features.

Background segmentation head. Our background segmentation head consists of two sparse convolutional layers with instance normalization and the ReLU [50] activation function after the first one. It takes the latent features \mathbf{F} and \mathbf{G} as input and outputs per-point foreground probabilities $\mathbf{h}^\mathbf{X} \in \mathbb{R}^{N^v}$ and $\mathbf{h}^\mathbf{Y} \in \mathbb{R}^{M^v}$.

Ego-motion head. Given the latent features \mathbf{F}^b and \mathbf{G}^b of the background points \mathbf{X}^b and \mathbf{Y}^b , the ego-motion head computes the affinity matrix $\mathbf{M} \in \mathbb{R}^{N^b \times N^b}$ s.t.

$$M_{ij} = \exp(-\|\mathbf{f}_i^b - \mathbf{g}_j^b\|/\tau_{\text{ego}}), \quad (9)$$

where τ_{ego} controls the *softness* of the correspondences.

Scene flow head. Based on the notion that scene flow is tightly coupled to correspondences [55], we now devise our scene flow head. To assure differentiability, we estimate *soft* correspondences. To allow for large motion, we measure the similarity in the latent space of *features* \mathbf{F} and \mathbf{G} instead of in the physical space. Hence, we find the corresponding points in \mathbf{X} and \mathbf{Y} as:

$$\mathbf{X}_c := \mathbf{YD}, \quad d_{ij} := \text{softmax}\left(-\frac{1}{\tau_{\text{flow}}} \|\mathbf{f}_i - \mathbf{g}_j\|_2\right) \quad (10)$$

where τ_{flow} is again a learnable temperature that controls the *softness* of the correspondences in the same manner as τ_{ego} above. Soft correspondences can be used to compute the initial flow estimate as $\mathbf{V}^{\text{init}} := \mathbf{X}_c - \mathbf{X}$. However, this

initial flow vector field is likely to be noisy due to large motions, sampling, and imperfect latent features and thus still has to be refined [55]. Our *refinement* module takes \mathbf{V}^{init} as input and locally smoothes it by estimating a residual flow $\Delta \mathbf{V}^{\text{init}}$ thorough a series of sparse convolutional layers. The refined scene flow $\mathbf{V} \in \mathbb{R}^{3 \times N^v}$ is then obtained as: $\mathbf{V} = \mathbf{V}^{\text{init}} + \Delta \mathbf{V}^{\text{init}}$. The detailed architecture of the scene flow head is given in the supplement.

From transformations to per-point rigid scene flow. The output of our multi-task network comprises of: (i) transformation parameters of the ego-motion \mathbf{T}_{ego} and individual clusters $\{\mathbf{T}_k\}_{k=1}^{K-1}$; (ii) object level masks $\{\mathbf{z}_k\}_{k=1}^K$; and (iii) unconstrained pointwise scene flow estimates \mathbf{V} . The *pointwise rigid scene flow* $\mathbf{V}^{\text{rigid}}$ can then be recovered as $\mathbf{V}^{\text{rigid}} \approx \{\mathbf{T}_k \circ \mathbf{X}_k - \mathbf{X}_k\}_{k=1}^K$, where \mathbf{X}_k denotes the points of \mathbf{X} belonging to the cluster k according to the inferred object masks \mathbf{z}_k . For the points that are neither assigned to the background nor to any of the foreground rigid bodies, we use the unconstrained scene flow predictions \mathbf{V} .

Training and implementation details. Our method is implemented in PyTorch using the *MinkowskiEngine* [12]. Unless specified differently, we train our network in an end-to-end manner for 40 epochs (or until convergence), by minimizing Eq (2). We train on a single NVIDIA GTX2080Ti with batch size 8. We use the Adam [38] optimizer with an initial learning rate 10^{-3} , which is decayed every epoch according to an exponential schedule with $\gamma = 0.98$. The whole training takes about two and a half days. The detailed parameters of the Sinkhorn algorithm and DBSCAN clustering are available in the supplement.

Inference. The abstraction of the scene into a collection of rigid bodies enables us to run test-time optimization of their inferred transformation parameters. Specifically, we run an optimization scheme in which we iteratively minimize the closest point distance to the target points. For ego-motion we index the background points, and for individual clusters, all the foreground points of \mathbf{Y} . The indexing is performed using the inferred object masks $\mathbf{z}_k^\mathbf{X}$ and $\mathbf{z}_k^\mathbf{Y}$. Given the final transformation estimates $\{\mathbf{T}_k^*\}_{k=1}^K$ pointwise rigid scene flow can be recomputed as $\mathbf{V}^{\text{rigid}} \approx \{\mathbf{T}_k^* \circ \mathbf{X}_k - \mathbf{X}_k\}_{k=1}^K$. A detailed description of our optimization scheme including its run-time is available in the supplement.

In the following, we reintroduce the superscript v so as to denote the voxelized point \mathbf{X}^v . Note that $\mathbf{V}^{\text{rigid}}$ represents the flow for the voxel centers \mathbf{X}^v and during inference, still has to be *transferred* to the points \mathbf{X} . We perform this transfer by a simple inverse-distance weighted interpolation:

$$\mathbf{v}_i^* = \frac{\sum_{j: \mathbf{x}_j^v \in \mathcal{E}(\mathbf{x}_i)} \mathbf{v}_j^{\text{rigid}} \|\mathbf{x}_i - \mathbf{x}_j^v\|_2^{-1}}{\sum_{j: \mathbf{x}_j^v \in \mathcal{E}(\mathbf{x}_i)} \|\mathbf{x}_i - \mathbf{x}_j^v\|_2^{-1}}, \quad (11)$$

where $\mathcal{E}(\cdot)$ returns the set of k -NN in the Euclidean sense.

Dataset	Method	Supervision	EPE3D [m] ↓	Acc3DS ↑	Acc3DR ↑	Outliers ↓
FT3D	FlowNet3D [40]	Full	0.114	0.412	0.771	0.602
	HPLFlowNet [24]	Full	0.080	0.614	0.855	0.429
	PointPWC-Net [80]	Full	0.059	0.738	0.928	0.342
	FLOT [55]	Full	0.052	0.732	0.927	0.357
	EgoFlow [72]	Full	0.069	0.670	0.879	0.404
stereoKITTI	Ours	Full	0.052	0.746	0.936	0.361
	FlowNet3D [40]	Full	0.177	0.374	0.668	0.527
	HPLFlowNet [24]	Full	0.117	0.478	0.778	0.410
	PointPWC-Net [80]	Full	0.069	0.728	0.888	0.265
	FLOT [55]	Full	0.056	0.755	0.908	0.242
	EgoFlow [72]	Full	0.103	0.488	0.822	0.394
	Ours	Full	0.042	0.849	0.959	0.208

Table 1: Evaluation results in a fully supervised setting on *FT3D* and *stereoKITTI* datasets.

4. Experimental Evaluation

In this section, we first describe the datasets (§ 4.1) and evaluation metrics (§ 4.2) used in our experiments. We start the evaluation, by assessing the performance of our backbone flow estimation network under full supervision on point clouds lifted from stereo images (§ 4.3). We then proceed to evaluate our full pipeline in a weakly supervised setting on real LiDAR scans (§ 4.4). Finally, we showcase the generalization capability of our method (§ 4.5) and justify our design choices in an ablation study (§ 4.6).

4.1. Datasets

For all datasets, we follow a common preprocessing step [40, 24] and remove points whose depth or distance to the sensor is larger than 35 m. For training and evaluation, we randomly sample 8192 points from both frames independently. A detailed description of the datasets and preprocessing steps is available in the supplement.

FlyingThings3D (FT3D). [44] is a large-scale stereo dataset of synthetic man-made objects that are scattered in space and move randomly between the two frames. We generate the point clouds and GT scene flow in accordance with [24]. FT3D consists of 19640 training examples (from which we use 3928 for validation) and 3824 test examples. Note, FT3D is only used for training in the fully supervised evaluation of our backbone and scene flow head (§ 4.3).

stereoKITTI. [45, 46] is a real world scene flow dataset with 142 point cloud pairs, which are all used for testing. The point clouds and GT scene flow are obtained by lifting the annotated disparity maps and optical flow to 3D [24]. As a consequence, the points of the two frames are under direct correspondence. We remove the ground points by a naive thresholding of the height coordinate [40, 24].

lidarKITTI. [22] is a real world dataset acquired with a Velodyne 64-beam LiDAR. It consists of the same 142 pairs as *stereoKITTI*. GT is obtained by projecting the point clouds to the image plane and assigning them the annotated 3D flow vectors. In this dataset, the points of the two input

Dataset	Method	Supervision	EPE3D [m] ↓	Acc3DS ↑	Acc3DR ↑	Outliers ↓
<i>lidarKITTI</i> (w/o ground)	PointPWC-Net [80]	Full	0.390	0.387	0.550	0.653
	FLOT [55]	Full	0.653	0.155	0.313	0.837
	Ours (backbone)	Full	0.535	0.262	0.437	0.742
	Ours	Weak	0.150	0.521	0.744	0.450
<i>lidarKITTI</i> (with ground)	Ours+	Weak	0.110	0.745	0.844	0.353
	Ours++	Weak	0.094	0.784	0.885	0.314
	PointPWC-Net [80]	Full	0.710	0.114	0.219	0.932
	FLOT [55]	Full	0.773	0.084	0.177	0.943
	MeteorNet [41] ¹	Full	0.277	/	/	/
	Ours (backbone)	Full	0.820	0.102	0.190	0.934
	Ours	Weak	0.133	0.460	0.746	0.527
	Ours+	Weak	0.106	0.673	0.808	0.421
	Ours++	Weak	0.102	0.686	0.819	0.410

¹ MeteorNet uses three fine frames and was trained on 100 and evaluated on 42 scenes of this dataset.

Table 2: Evaluation results on *lidarKITTI*. Ours (backbone) denotes our model from § 4.3 trained with full supervision on *FT3D*. Ours are the direct estimates of our pipeline. Ours+ and Ours++ additionally denote test-time optimization of only ego-motion and all rigid bodies, respectively.

frames are not in direct correspondence and have a typical sampling pattern of a LiDAR sensor.

semanticKITTI. [5] provides per point semantic labels and accurate ego-motion for 21 LiDAR sequences of the KITTI odometry dataset. It is split into eleven (00-10) LiDAR sequences for training and eleven (11-21) for testing. We use sequences 03 and 05 for validation and the remaining nine for training. SemanticKITTI is used to train our method in a weakly supervised manner (§ 4.4 to § 4.6). Note that this dataset does not contain dense scene flow annotations.

4.2. Evaluation metrics

We use standard evaluation metrics to assess the performance of our approach and compare it with SoTA methods, *FlowNet3D* [40], *HPLFlowNet* [24], *PointPWCNet* [80], FLOT [55], and EgoFlow [72]. Our main evaluation metric is the 3D end-point-error (*EPE3D*), defined as the mean *l2* distance between the predicted and GT scene flow.

Additionally, we follow [40, 24] and also report: (i) strict accuracy (*Acc3DS*), defined as the percentage of points whose *EPE3D* < 0.05 m or relative error < 0.05, (ii) relaxed accuracy (*Acc3DR*) that denotes the ratio of points whose *EPE3D* < 0.10 m or relative error < 0.10, and (iii) *Outliers*, i.e. the ratio of points whose *EPE3D* > 0.30 m or relative error > 0.10.

For the experiments in a weakly supervised setting, we also report the relative angular error (*RAE*) and the relative translation error (*RTE*) of the estimated ego-motion.

4.3. Our backbone under full supervision

A core module of our proposed pipeline is the backbone network described in § 3.2. It is therefore valuable to first assess its performance in conjunction with the scene flow prediction head, before turning to evaluate the performance of our entire weakly-supervised pipeline. To this end, we follow the traditional setting used by our competitors and train in a fully supervised manner on *FT3D* by minimizing

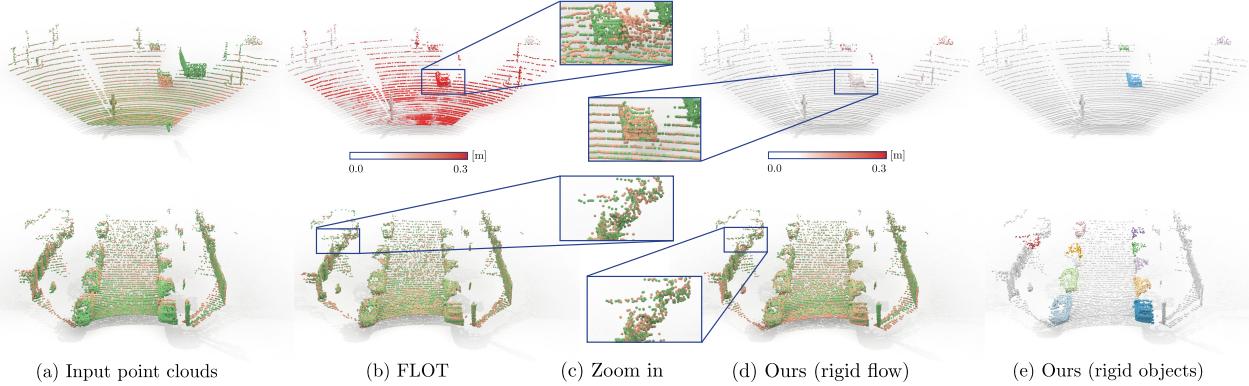


Figure 4: Qualitative results of our weakly supervised method on *lidarKITTI* (top) and *waymo open* (bottom). For improved visibility, the *EPE3D* (top row b,c) is clipped to the range between 0.0 m (white) at 0.3m (red). As a result of predicting an unconstrained pointwise sceneflow, the rigid objects (car) in the results of FLOT might get deformed (d).

the *l1* distance between the predicted and GT scene flow. We then evaluate our model on both *FT3D* and *stereoKITTI*.

When evaluated on *FT3D*, our method performs on par with FLOT [55] in terms of *EPE3D* and outperforms all methods in terms of *Acc3DS* and *Acc3DR* (Tab. 1). More importantly, it achieves superior generalization performance on *stereoKITTI*, where it consistently outperforms SoTA in all evaluation metrics, setting a new SoTA with 0.042 m *EPE3D* (\approx 1.5 cm better than the closest competitor). Based on these results, we conclude that sparse convolutions are an effective backbone for scene-flow estimation, and that our simple scene flow head can match the performance of SoTA, while enabling better generalization.

4.4. Our pipeline under weak supervision

Setting. Point clouds in both *FT3D* and *stereoKITTI* are obtained in the same manner: by lifting stereo images to 3D. Hence, their domain gap is relatively small. On the other hand, in LiDAR-based autonomous driving scenarios, point clouds are much sparser and assume a very different sampling pattern, resulting in a much more challenging setting for scene flow estimation.

We evaluate our entire weakly-supervised pipeline in this challenging setting by using the *lidarKITTI* dataset. Specifically, we consider two scenarios: 1) we remove ground points by naively thresholding the vertical coordinate, 2) we use the “raw” point clouds that also include the ground points for which the flow estimation is especially difficult.

We train a joint model for both scenarios in a weakly supervised manner using the point clouds from *semanticKITTI*. Unlike *semanticKITTI*, *lidarKITTI* only includes the points and annotations of the objects that are visible within the front camera images. We therefore process *semanticKITTI* in the same manner. Since there are no LiDAR datasets available with scene flow annotations, we use the models trained on *FT3D* for all the baselines.

Evaluation. Fig. 4 and Tab. 2 show that the domain gap between the *stereo* and *LiDAR* point clouds is too big for the traditional fully supervised methods to generalize effectively. Indeed, the performance of SoTA methods is up to 10 times worse when compared to the results on *stereoKITTI*. On the other hand, our weakly supervised model predicts accurate rigid scene flow with an *EPE3D* \approx 0.1 m for both scenarios (after test-time optimization), while also providing an object-level abstraction (Fig. 4). Since our fully supervised backbone model also fails to generalize, we conclude that the crucial advantage of our method does not lie in a stronger backbone, but rather in the ability to train on the same domain. Additional qualitative and quantitative results are available in the supplement.

4.5. Generalization to other datasets.

Waymo open [69] is a recently introduced large-scale autonomous driving dataset that would ideally be used for supervision of 3D scene flow methods. However, it does not provide dense flow annotations. While it does include all the annotations that our weakly supervised approach relies on, we are more interested in using it to evaluate the generalization capability of our method. To this end, we use the first three sequences⁴ of the *waymo open* validation set and quantitatively evaluate our weakly supervised model in terms of ego-motion estimation and background segmentation. We provide qualitative results of the rigid scene flow estimation and object-level scene abstraction in Fig. 4.

Remarkably, our model that was trained only on *semanticKITTI* can seamlessly generalize to *waymo open*. When evaluated on the task of ego-motion estimation, it achieves an *RRE* of 0.141° and *RTE* of 0.099 m. In the BG-segmentation task, its performance on the foreground points drops slightly to 0.960 precision and 0.689 recall, while on the background points it remains high with 0.957 and 0.996

⁴This results in more than 14k point cloud pairs, which is almost the same size as the whole *semanticKITTI* dataset

lidarKITTI						
\mathcal{L}_{ego}	\mathcal{L}_{CD}	$\mathcal{L}_{\text{rigid}}$	EPE [m] ↓	Acc3DS ↑	Acc3DR ↑	RRE [°] ↓
✓	✓		0.721	0.044	0.093	0.476
✓	✓		0.363	0.044	0.163	0.610
✓	✓		0.136	0.409	0.712	0.380
✓	✓	✓	0.134	0.460	0.746	0.320
						0.130

Table 3: Ablation study of the proposed training objective. All models are trained on *semanticKITTI* and evaluated without test-time optimization on *lidarKITTI (with ground)* dataset.

precision and recall, respectively. The drop in foreground performance can be accredited to the domain gap between the datasets [85]; *waymo open* includes many more foreground objects, especially pedestrians, than *semanticKITTI*.

4.6. Ablation Studies

Influence of different loss terms. We compare our model trained with the full loss function to multiple ablations in Tab. 3. Each model is trained on *semanticKITTI* and evaluated without test-time optimization on *lidarKITTI* with ground points. Tab. 3 shows that the terms \mathcal{L}_{ego} and \mathcal{L}_{CD} are crucial for the performance of our model. Adding $\mathcal{L}_{\text{rigid}}$ further regularizes the performance and leads to an improvement in all evaluation metrics. Note how the terms that are applied only on foreground points (e.g. $\mathcal{L}_{\text{rigid}}$) also improve the ego-motion estimation Tab. 3.

Task-specific networks. Instead of training a single network capable of solving multiple tasks, one could also devise a combination of task-specific networks. We ablate this design choice in Tab. 4 in which we compare our full model to BG segmentation and ego-motion specific networks. Both task specific networks comprise of our backbone with the corresponding head and are trained with the \mathcal{L}_{BG} and \mathcal{L}_{ego} objective function, respectively. Tab. 4 shows that the performance of our full model is slightly inferior to the task-specific one when compared on BG segmentation. This is expected, since in our full pipeline the BG segmentation head is also trained nearly in isolation, with a single loss function. On the other hand, our full model outperforms the task-specific ego-motion model, even when the task-specific model is combined with the GT background mask. This shows that individual tasks (e.g. flow and ego-motion estimation) can indeed reinforce each other, which leads to better downstream performance.

Pretraining the backbone with full supervision. We analyze the effect of initializing our weakly supervised model with pretrained backbone weights. To this end, we use the backbone weights from the model trained with full supervision in § 4.3. Initializing the backbone with the pretrained model leads to 1.4 cm and 2.3 cm improvement in terms of *EPE3D* on *lidarKITTI* without and with ground points, respectively. Further details and additional metrics of this ablation study are available in the supplement. Note that

Task BG seg. ego-motion	BG segmentation				ego-motion	
	prec. FG ↑	rec. FG ↑	prec. BG ↑	recall. BG ↑	RRE [°] ↓	RTE [m] ↓
<i>semanticKITTI (w/o ground)</i>						
✓	0.977	0.901	0.992	0.998	-	-
✓	✓	-	-	-	0.245	0.054
✓	✓	0.971	0.895	0.991	0.998	0.201
<i>semanticKITTI (with ground)</i>						
✓	0.970	0.911	0.996	0.999	-	-
✓	✓	-	-	-	0.307	0.071
✓	✓	0.966	0.904	0.996	0.999	0.249
						0.059

Table 4: Comparison of our full pipeline with specialized networks for BG segmentation and ego-motion estimation, respectively. Note, we provide GT background masks to the ego-motion specialized network also in the test phase.

in all evaluations presented in Sec. 4 we use the inferior model with randomly initialized weights trained only with weak supervision.

Run time. We now compare our method to FLOT [55] and PointPWC-Net [80] in terms of run-time and number of parameters⁵. We perform the evaluation on a standalone computer with Intel Xeon E5-1650, 32GB RAM, and a single NVIDIA Titan V. For FLOT [55] and PointPWC-Net [80] we use the official implementation provided by the authors. FLOT has the lowest number of trainable parameters (0.11 million) but, with 0.395 seconds on average, also the highest run time per point cloud pair. PointPWC-Net has a larger model with approximately 7.7 million parameters but performs one inference step in 0.147 seconds on average. Finally, our method contains about 8 million parameters and requires 0.154 seconds on average for a single point cloud pair. With added test-time optimization our run time increases to 0.234 seconds on average.

5. Conclusion

Scene flow is the lowest level in a hierarchy of dynamic scene perception. As such, while providing a useful cue to higher-level tasks, it is also the most demanding to supervise. Based on this observation, in this work, we have introduced a novel method that relaxes the dense supervision by integrating flow into a higher-level scene abstraction in the form of multi rigid-body motion. The result is a state-of-the-art flow estimation network that additionally outputs a concise dynamic scene representation. In particular, our mild supervision requirements are well suited for utilizing the annotation level of recently released massive data collections for autonomous driving. In future work, we plan to incorporate cues from multiple frames further seeking temporal consistency as well as increased accuracy.

Acknowledgements. This work is sponsored by Stanford-Ford Alliance, the Samsung GRO program, NSF grant IIS-1763268, the Vannevar Bush Faculty fellowship, and the NVIDIA GPU grant. We thank Barbara Verbost for her help with the visualizations and Davis Rempe for proofreading the article.

⁵The evaluation is performed with 8192 randomly sampled points on the *lidarKITTI* dataset.

References

- [1] Eman Ahmed, Alexandre Saint, Abd El Rahman Shabayek, Kseniya Cherenkova, Rig Das, Gleb Gusev, Djamilia Aouada, and Björn Ottersten. Deep learning advances on different 3d data representations: A survey. *arXiv preprint arXiv:1808.01462*, 1, 2018. 2
- [2] Hassan Abu Alhaija, Anita Sellent, Daniel Kondermann, and Carsten Rother. Graphflow-6d large displacement scene flow via graph matching. In *German Conference on Pattern Recognition*, pages 285–296. Springer, 2015. 2
- [3] Steven S. Beauchemin and John L. Barron. The computation of optical flow. *ACM computing surveys (CSUR)*, 27(3):433–466, 1995. 1
- [4] Aseem Behl, Despoina Paschalidou, Simon Donné, and Andreas Geiger. Pointflownet: Learning representations for rigid motion estimation from point clouds. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 7962–7971, 2019. 2
- [5] J. Behley, M. Garbade, A. Milioto, J. Quenzel, S. Behnke, C. Stachniss, and J. Gall. SemanticKITTI: A Dataset for Semantic Scene Understanding of LiDAR Sequences. In *IEEE International Conf. on Computer Vision*, 2019. 2, 6, 14
- [6] Paul J Besl and Neil D McKay. Method for registration of 3-d shapes. In *Sensor fusion IV: control paradigms and data structures*, volume 1611, pages 586–606. International Society for Optics and Photonics, 1992. 13
- [7] Tolga Birdal, Michael Arbel, Umut Simsekli, and Leonidas J Guibas. Synchronizing probability measures on rotations via optimal transport. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1569–1579, 2020. 13
- [8] Tolga Birdal, Umut Simsekli, Mustafa Onur Eken, and Slobodan Ilic. Bayesian pose graph optimization via bingham distributions and tempered geodesic mcmc. In S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 31, pages 308–319. Curran Associates, Inc., 2018. 1
- [9] Mai Bui, Tolga Birdal, Haowen Deng, Shadi Albarqouni, Leonidas Guibas, Slobodan Ilic, and Nassir Navab. 6d camera relocalization in ambiguous scenes via continuous multi-modal inference. In Andrea Vedaldi, Horst Bischof, Thomas Brox, and Jan-Michael Frahm, editors, *Computer Vision – ECCV 2020*, pages 139–157, Cham, 2020. Springer International Publishing. 1
- [10] Arunkumar Byravan and Dieter Fox. Se3-nets: Learning rigid body motion using deep neural networks. In *2017 IEEE International Conference on Robotics and Automation (ICRA)*, pages 173–180. IEEE, 2017. 1
- [11] Rodrigo L Carceroni and Kiriakos N Kutulakos. Multi-view scene capture by surfel sampling: From video streams to non-rigid 3d motion, shape and reflectance. *International Journal of Computer Vision*, 49(2-3):175–214, 2002. 1, 2
- [12] Christopher Choy, JunYoung Gwak, and Silvio Savarese. 4d spatio-temporal convnets: Minkowski convolutional neural networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019. 5, 12
- [13] João Paulo Costeira and Takeo Kanade. A multibody factorization method for independently moving objects. *International Journal of Computer Vision*, 29(3):159–179, 1998. 2
- [14] Marco Cuturi. Sinkhorn distances: Lightspeed computation of optimal transport. In *Advances in neural information processing systems*, pages 2292–2300, 2013. 4, 14
- [15] Haowen Deng, Mai Bui, Nassir Navab, Leonidas Guibas, Slobodan Ilic, and Tolga Birdal. Deep bingham networks: Dealing with uncertainty and ambiguity in pose estimation. *arXiv preprint arXiv:2012.11002*, 2020. 1
- [16] Frederic Devernay, Diana Mateus, and Matthieu Guilbert. Multi-camera scene flow by tracking 3-d points and surfels. In *2006 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR’06)*, volume 2, pages 2203–2212. IEEE, 2006. 1, 2
- [17] Ayush Dewan, Tim Caselitz, Gian Diego Tipaldi, and Wolfram Burgard. Rigid scene flow for 3d lidar scans. In *2016 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 1765–1770. IEEE, 2016. 2
- [18] Alexey Dosovitskiy, Philipp Fischer, Eddy Ilg, Philip Hausser, Caner Hazirbas, Vladimir Golkov, Patrick Van Der Smagt, Daniel Cremers, and Thomas Brox. Flownet: Learning optical flow with convolutional networks. In *Proceedings of the IEEE international conference on computer vision*, pages 2758–2766, 2015. 2
- [19] Mingsong Dou, Sameh Khamis, Yury Degtyarev, Philip Davidson, Sean Ryan Fanello, Adarsh Kowdle, Sergio Orts Escolano, Christoph Rhemann, David Kim, Jonathan Taylor, et al. Fusion4d: Real-time performance capture of challenging scenes. *ACM Transactions on Graphics (TOG)*, 35(4):1–13, 2016. 2
- [20] Martin Ester, Hans-Peter Kriegel, Jörg Sander, and Xiaowei Xu. A density-based algorithm for discovering clusters in large spatial databases with noise. In *KDD*, 1996. 4
- [21] Hehe Fan and Yi Yang. Pointrnn: Point recurrent neural network for moving point cloud processing. *arXiv preprint arXiv:1910.08287*, 2019. 2
- [22] Andreas Geiger, Philip Lenz, and Raquel Urtasun. Are we ready for autonomous driving? the kitti vision benchmark suite. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2012. 2, 6, 14
- [23] Vladislav Golyanik, Kihwan Kim, Robert Maier, Matthias Nießner, Didier Stricker, and Jan Kautz. Multiframe scene flow with piecewise rigid motion. In *2017 International Conference on 3D Vision (3DV)*, pages 273–281. IEEE, 2017. 2
- [24] Xiuye Gu, Yijie Wang, Chongruo Wu, Yong Jae Lee, and Panqu Wang. Hplflownet: Hierarchical permutohedral lattice flownet for scene flow estimation on large-scale point clouds. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 3254–3263, 2019. 2, 6, 14
- [25] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. Mask r-cnn. In *Proceedings of the IEEE international conference on computer vision*, pages 2961–2969, 2017. 2
- [26] Berthold KP Horn and Brian G Schunck. Determining optical flow. In *Techniques and Applications of Image Understanding*, volume 281, pages 319–331. International Society for Optics and Photonics, 1981. 1

- [27] Michael Hornacek, Andrew Fitzgibbon, and Carsten Rother. Spherenet: 6 dof scene flow from rgb-d pairs. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3526–3533, 2014. 2
- [28] Jiahui Huang, Sheng Yang, Zishuo Zhao, Yu-Kun Lai, and Shi-Min Hu. Clusterslam: A slam backend for simultaneous rigid body clustering and motion estimation. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 5875–5884, 2019. 2
- [29] Frédéric Huguet and Frédéric Devernay. A variational method for scene flow estimation from stereo sequences. In *2007 IEEE 11th International Conference on Computer Vision*, pages 1–7. IEEE, 2007. 2
- [30] Eddy Ilg, Nikolaus Mayer, Tonmoy Saikia, Margret Keuper, Alexey Dosovitskiy, and Thomas Brox. Flownet 2.0: Evolution of optical flow estimation with deep networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2462–2470, 2017. 2
- [31] Matthias Innmann, Michael Zollhöfer, Matthias Nießner, Christian Theobalt, and Marc Stamminger. Volumedefor: Real-time volumetric non-rigid reconstruction. In *European Conference on Computer Vision*, pages 362–379. Springer, 2016. 2
- [32] Mariano Jaimez, Christian Kerl, Javier Gonzalez-Jimenez, and Daniel Cremers. Fast odometry and scene flow from rgbd cameras based on geometric clustering. In *2017 IEEE International Conference on Robotics and Automation (ICRA)*, pages 3992–3999. IEEE, 2017. 2
- [33] M. Jaimez, M. Souiai, J. Gonzalez-Jimenez, and D. Cremers. A primal-dual framework for real-time dense rgbd scene flow. In *Proc. of the IEEE Int. Conf. on Robotics and Automation (ICRA)*, 2015. 2
- [34] Huazhu Jiang, Deqing Sun, Varun Jampani, Zhaoyang Lv, Erik Learned-Miller, and Jan Kautz. Sense: A shared encoder network for scene-flow estimation. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 3195–3204, 2019. 5
- [35] Li Jiang, Hengshuang Zhao, Shaoshuai Shi, Shu Liu, Chi-Wing Fu, and Jiaya Jia. Pointgroup: Dual-set point grouping for 3d instance segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4867–4876, 2020. 4
- [36] Wolfgang Kabsch. A solution for the best rotation to relate two sets of vectors. *Acta Crystallographica Section A: Crystal Physics, Diffraction, Theoretical and General Crystallography*, 32(5):922–923, 1976. 4
- [37] Ken-ichi Kanatani. Motion segmentation by subspace separation and model selection. In *Proceedings Eighth IEEE International Conference on computer Vision. ICCV 2001*, volume 2, pages 586–591. IEEE, 2001. 2
- [38] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014. 5
- [39] Solomon Kullback and Richard A Leibler. On information and sufficiency. *The annals of mathematical statistics*, 22(1):79–86, 1951. 14
- [40] Xingyu Liu, Charles R Qi, and Leonidas J Guibas. Flownet3d: Learning scene flow in 3d point clouds. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 529–537, 2019. 1, 2, 6, 14
- [41] Xingyu Liu, Mengyuan Yan, and Jeannette Bohg. Meteor-net: Deep learning on dynamic 3d point cloud sequences. In *IEEE International Conference on Computer Vision*, pages 9246–9255, 2019. 1, 2, 6, 14
- [42] Zhaoyang Lv, Kihwan Kim, Alejandro Troccoli, Deqing Sun, James M Rehg, and Jan Kautz. Learning rigidity in dynamic scenes with a moving camera for 3d motion field estimation. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 468–484, 2018. 2
- [43] Wei-Chiu Ma, Shenlong Wang, Rui Hu, Yuwen Xiong, and Raquel Urtasun. Deep rigid instance scene flow. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3614–3622, 2019. 2
- [44] N. Mayer, E. Ilg, P. Häusser, P. Fischer, D. Cremers, A. Dosovitskiy, and T. Brox. A large dataset to train convolutional networks for disparity, optical flow, and scene flow estimation. In *IEEE International Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016. arXiv:1512.02134. 1, 2, 6, 14
- [45] Moritz Menze, Christian Heipke, and Andreas Geiger. Joint 3d estimation of vehicles and scene flow. *ISPRS Annals of Photogrammetry, Remote Sensing & Spatial Information Sciences*, 2, 2015. 1, 2, 6, 14
- [46] Moritz Menze, Christian Heipke, and Andreas Geiger. Object scene flow. *ISPRS Journal of Photogrammetry and Remote Sensing*, 140:60–76, 2018. 6, 14
- [47] Himangi Mittal, Brian Okorn, and David Held. Just go with the flow: Self-supervised scene flow estimation. *arXiv preprint arXiv:1912.00497*, 2019. 1, 2
- [48] Frank Moosmann and Christoph Stiller. Joint self-localization and tracking of generic objects in 3d range data. In *2013 IEEE International Conference on Robotics and Automation*, pages 1146–1152. IEEE, 2013. 2
- [49] Armin Mustafa and Adrian Hilton. Semantically coherent 4d scene flow of dynamic scenes. *International Journal of Computer Vision*, pages 1–17, 2019. 2
- [50] Vinod Nair and Geoffrey E Hinton. Rectified linear units improve restricted boltzmann machines. In *ICML*, 2010. 5, 13
- [51] Richard A Newcombe, Dieter Fox, and Steven M Seitz. Dynamicfusion: Reconstruction and tracking of non-rigid scenes in real-time. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 343–352, 2015. 2
- [52] Michael Niemeyer, Lars Mescheder, Michael Oechsle, and Andreas Geiger. Occupancy flow: 4d reconstruction by learning particle dynamics. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 5379–5389, 2019. 1, 2
- [53] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay. Scikit-learn: Machine

- learning in Python. *Journal of Machine Learning Research*, 12:2825–2830, 2011. 12
- [54] J-P Pons, Renaud Keriven, and Olivier Faugeras. Modelling dynamic scenes by registering multi-view image sequences. In *2005 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR'05)*, volume 2, pages 822–827. IEEE, 2005. 2
- [55] Gilles Puy, Alexandre Boulch, and Renaud Marlet. FLOT: Scene Flow on Point Clouds Guided by Optimal Transport. In *European Conference on Computer Vision*, 2020. 1, 2, 5, 6, 7, 8, 14
- [56] Charles R Qi, Hao Su, Kaichun Mo, and Leonidas J Guibas. Pointnet: Deep learning on point sets for 3d classification and segmentation. In *IEEE conference on computer vision and pattern recognition*, pages 652–660, 2017. 2
- [57] Charles Ruizhongtai Qi, Li Yi, Hao Su, and Leonidas J Guibas. Pointnet++: Deep hierarchical feature learning on point sets in a metric space. In *Advances in neural information processing systems*, pages 5099–5108, 2017. 2
- [58] Davis Rempe, Tolga Birdal, Yongheng Zhao, Zan Gojcic, Srinath Sridhar, and Leonidas J. Guibas. Caspr: Learning canonical spatiotemporal point cloud representations. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2020. 1, 2
- [59] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *International Conference on Medical image computing and computer-assisted intervention*, pages 234–241. Springer, 2015. 5
- [60] Martin Runz, Maud Buffier, and Lourdes Agapito. Maskfusion: Real-time recognition, tracking and reconstruction of multiple moving objects. In *2018 IEEE International Symposium on Mixed and Augmented Reality (ISMAR)*, pages 10–20. IEEE, 2018. 2
- [61] Samuele Salti, Federico Tombari, and Luigi Di Stefano. A performance evaluation of 3d keypoint detectors. In *2011 International Conference on 3D Imaging, Modeling, Processing, Visualization and Transmission*, pages 236–243. IEEE, 2011. 2
- [62] Muhamad Risqi U Saputra, Andrew Markham, and Niki Trigoni. Visual slam and structure from motion in dynamic environments: A survey. *ACM Computing Surveys (CSUR)*, 51(2):1–36, 2018. 2
- [63] Richard Sinkhorn. A relationship between arbitrary positive matrices and doubly stochastic matrices. *The annals of mathematical statistics*, 35(2):876–879, 1964. 4, 14
- [64] Richard Sinkhorn. Diagonal equivalence to matrices with prescribed row and column sums. *The American Mathematical Monthly*, 74(4):402–405, 1967. 14
- [65] Miroslava Slavcheva, Maximilian Baust, Daniel Cremers, and Slobodan Ilic. Killingfusion: Non-rigid 3d reconstruction without correspondences. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1386–1395, 2017. 2
- [66] Miroslava Slavcheva, Maximilian Baust, and Slobodan Ilic. Sobolevfusion: 3d reconstruction of scenes undergoing free non-rigid motion. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2646–2655, 2018. 2
- [67] Michael Strecke and Jorg Stuckler. Em-fusion: Dynamic object-level slam with probabilistic data association. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 5865–5874, 2019. 2
- [68] Hang Su, Varun Jampani, Deqing Sun, Subhransu Maji, Evangelos Kalogerakis, Ming-Hsuan Yang, and Jan Kautz. Splatnet: Sparse lattice networks for point cloud processing. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2530–2539, 2018. 2
- [69] 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: An open dataset benchmark. *arXiv preprint arXiv:1912.04838*, 2019. 2, 7, 15
- [70] Georg Tanzmeister, Julian Thomas, Dirk Wollherr, and Martin Buss. Grid-based mapping and tracking in dynamic environments using a uniform evidential environment representation. In *2014 IEEE International Conference on Robotics and Automation (ICRA)*, pages 6090–6095. IEEE, 2014. 2
- [71] Lyne Tchapmi, Christopher Choy, Iro Armeni, JunYoung Gwak, and Silvio Savarese. Segcloud: Semantic segmentation of 3d point clouds. In *2017 international conference on 3D vision (3DV)*, pages 537–547. IEEE, 2017. 2
- [72] Ivan Tishchenko, Sandro Lombardi, Martin R Oswald, and Marc Pollefeys. Self-supervised learning of non-rigid residual flow and ego-motion. *arXiv preprint arXiv:2009.10467*, 2020. 1, 2, 6
- [73] Arash K Ushani, Ryan W Wolcott, Jeffrey M Walls, and Ryan M Eustice. A learning approach for real-time temporal scene flow estimation from lidar data. In *2017 IEEE International Conference on Robotics and Automation (ICRA)*, pages 5666–5673. IEEE, 2017. 2
- [74] Sundar Vedula, Simon Baker, Peter Rander, Robert Collins, and Takeo Kanade. Three-dimensional scene flow. In *IEEE International Conference on Computer Vision*, pages 722–729. IEEE, 1999. 1, 2
- [75] Christoph Vogel, Konrad Schindler, and Stefan Roth. Piecewise rigid scene flow. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 1377–1384, 2013. 2
- [76] Shenlong Wang, Simon Suo, Wei-Chiu Ma, Andrei Pokrovsky, and Raquel Urtasun. Deep parametric continuous convolutional neural networks. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 2589–2597, 2018. 2
- [77] Weiyue Wang, Ronald Yu, Qiangui Huang, and Ulrich Neumann. Sgpn: Similarity group proposal network for 3d point cloud instance segmentation. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 2569–2578, 2018. 2
- [78] Zirui Wang, Shuda Li, Henry Howard-Jenkins, Victor Adrian Prisacariu, and Min Chen. Flownet3d++: Geometric losses for deep scene flow estimation. *arXiv preprint arXiv:1912.01438*, 2019. 1, 2

- [79] A. Wedel, C. Rabe, T. Vaudrey, T. Brox, U. Franke, and D. Cremers. Efficient dense scene flow from sparse or dense stereo data. In *ECCV*, Marseille, France, October 2008. 2
- [80] Wenxuan Wu, Zhiyuan Wang, Zhuwen Li, Wei Liu, and Li Fuxin. Pointpwc-net: A coarse-to-fine network for supervised and self-supervised scene flow estimation on 3d point clouds. *arXiv preprint arXiv:1911.12408*, 2019. 1, 2, 6, 8, 14
- [81] Saining Xie, Jiatao Gu, Demi Guo, Charles R Qi, Leonidas J Guibas, and Or Litany. Pointcontrast: Unsupervised pre-training for 3d point cloud understanding. *arXiv preprint arXiv:2007.10985*, 2020. 16
- [82] Jiaqi Yang, Yang Xiao, and Zhiguo Cao. Toward the repeatability and robustness of the local reference frame for 3d shape matching: An evaluation. *IEEE Transactions on Image Processing*, 27(8):3766–3781, 2018. 2
- [83] Jiancheng Yang, Qiang Zhang, Bingbing Ni, Linguo Li, Jinxian Liu, Mengdie Zhou, and Qi Tian. Modeling point clouds with self-attention and gumbel subset sampling. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 3323–3332, 2019. 2
- [84] Zi Jian Yew and Gim Hee Lee. Rpm-net: Robust point matching using learned features. In *CVPR*, 2020. 4
- [85] Li Yi, Boqing Gong, and Thomas Funkhouser. Complete & label: A domain adaptation approach to semantic segmentation of lidar point clouds, 2020. 8
- [86] Li Yi, Haibin Huang, Difan Liu, Evangelos Kalogerakis, Hao Su, and Leonidas Guibas. Deep part induction from articulated object pairs. *ACM Transactions on Graphics (TOG)*, 37(6):209, 2019. 2
- [87] Li Yi, Wang Zhao, He Wang, Minhyuk Sung, and Leonidas J Guibas. Gspn: Generative shape proposal network for 3d instance segmentation in point cloud. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 3947–3956, 2019. 2
- [88] Manzil Zaheer, Satwik Kottur, Siamak Ravanbakhsh, Barnabas Poczos, Ruslan Salakhutdinov, and Alexander Smola. Deep sets, 2018. 2
- [89] Amir R Zamir, Alexander Sax, Nikhil Cheerla, Rohan Suri, Zhangjie Cao, Jitendra Malik, and Leonidas J Guibas. Robust learning through cross-task consistency. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11197–11206, 2020. 5
- [90] Amir R Zamir, Alexander Sax, William Shen, Leonidas J Guibas, Jitendra Malik, and Silvio Savarese. Taskonomy: Disentangling task transfer learning. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3712–3722, 2018. 5
- [91] Qian-Yi Zhou, Jaesik Park, and Vladlen Koltun. Open3D: A modern library for 3D data processing. *arXiv:1801.09847*, 2018. 13
- [92] Yin Zhou and Oncel Tuzel. Voxelnet: End-to-end learning for point cloud based 3d object detection. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4490–4499, 2018. 2

A. Appendix

A.1. Implementation details

We now provide the implementation details of our network architecture, foreground clustering using DBSCAN, and test-time optimization scheme. Additionally, we summarize the Kabsch algorithm and the entropy regularized optimal transport. For a general overview please refer to Sec. 3 and Fig. 3.

A.1.1 Network architecture

Our whole network is built upon MinkowskiEngine [12], an auto-differentiation library for sparse tensors. In all experiments, we first randomly sample 8192 points from the source and target point cloud, respectively. We then voxelize the original point clouds into sparse tensors with a voxel size of 0.1 m. If more than 8192 voxels remain, we randomly select 8192 of them. For the evaluation, all the inferred quantities are transferred from the voxels to originally sampled points using Eq (11). The combined network has 8,078,149 learnable parameters. Fig. 5 shows the detailed architecture, where blue rectangles denote the Minkowski convolutional and transpose convolutional layers, violet rectangles the ResBlocks, and orange rectangles are ReLU activation functions.

Backbone network. The detailed architecture of the backbone network is depicted in Fig. 5 (a). We use absolute coordinates instead of the more common binary occupancy as the input feature of the first layer. All convolutional layers are followed by instance normalization and in ResBlocks additionally by the ReLU activation function. The output of the backbone network are 64-dimensional pointwise latent features. Due to the fixed voxel size, the number of points can vary between the point clouds.

Scene flow head. The scene flow head, depicted in Fig. 5 (b), takes the initial flow vector field \mathbf{V}^{init} as input and computes a residual flow $\Delta\mathbf{V}^{\text{init}}$. All layers except the last are followed by the ReLU activation function. Normalization functions are not used in the scene flow head.

Background segmentation head. The background segmentation head (Fig. 5 (c)) comprises two sparse convolutional layers, where the first one is followed by the instance normalization and ReLU activation function. The output of the second layer is passed through a sigmoid function to obtain the foreground probabilities \mathbf{h} .

A.1.2 Foreground clustering

We cluster the foreground points into objects using a simple DBSCAN clustering algorithm. Specifically, we first extract the foreground points \mathbf{X}^f of the source point cloud using the inferred foreground probabilities $\mathbf{h}^{\mathbf{X}}$. We then run a scikit-learn [53] implementation of DBSCAN with $\text{eps} = 0.75$ m and the minimum number of samples in the neighborhood equal to 5. After clustering, we only retain clusters with at least 10 points. Our clustering algorithm is based on the hypothesis that the objects scattered in the scene, are separated by void space. If two objects (e.g. cars) are too close to each other and get assigned to the same cluster, our method does not have the capacity to recover from this wrong assignment. However, we conduct an ablation study (§ A.3.2) and

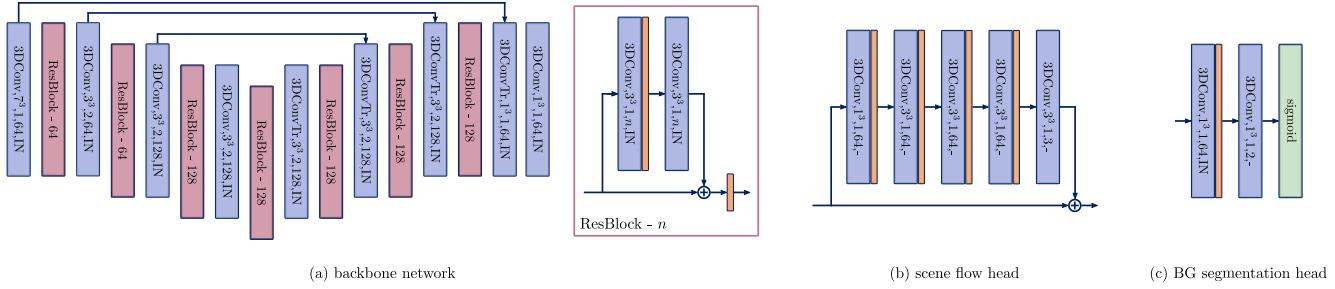


Figure 5: Network architecture of the backbone network (a), scene flow head (b), and background segmentation head (c). Blue blocks denote the convolutional (3DConv) and transpose convolutional (3DConvTr) layers, where $3DConv,a,b,c,d$ denotes a 3D convolutional layer with the kernel size a , stride b , output feature dimension c , and normalization function d . Orange rectangle denotes the ReLU activation function [50] and IN is the instance normalization.

empirically confirm the above-mentioned hypothesis for the autonomous driving datasets.

A.1.3 Test-time optimization

The object-level abstraction of the scene enables us to perform test-time optimization of per-object rigid body transformations. This optimization is performed independently for the background and each of the foreground objects. In the following we assume that our network outputs object-level masks $\{\mathbf{z}_k\}_{k=1}^K$ and transformation parameters $\{\mathbf{T}_k\}_{k=1}^K$.

Background transformation. Let \mathbf{z}_1 and \mathbf{T}_1 denote the inferred background mask and ego-motion transformation parameters, respectively. We use \mathbf{z}_1 to extract the background points \mathbf{X}^b and \mathbf{Y}^b of the source and target point cloud, transform \mathbf{X}^b with the inferred ego-motion \mathbf{T}_1 and use these point clouds as an input to the ICP [6] optimization of the transformation parameters. Specifically, we use the Open3D [91] implementation with the maximum correspondence distance of 0.15 m and a maximum number of iterations equal to 300. On *lidarKITTI* the optimization of the background transformation parameters takes 0.06 s on average for a point cloud pair, using the same computer as in the run time experiments in § 4.6.

Object level transformations. For the object level transformations, we follow a similar procedure. We start by extracting the points \mathbf{X}^k of the object k in the source point cloud using the object mask \mathbf{z}_k , and all the foreground points \mathbf{Y}^f of the target point cloud using the complementary mask of the background mask \mathbf{z}_1 . We then transform the points \mathbf{X}^k with the inferred transformation parameters \mathbf{T}_k and again use the point clouds as input to the ICP optimization of the transformation parameters. We use all the foreground points of the target point cloud as we do not have access to the instance level correspondences. The optimization of all foreground transformation parameters for a single *lidarKITTI* point cloud pair takes approximately 0.02 s on average. We perform at most 300 iterations of ICP per object with a maximum correspondences distance equal to 0.25 m.

A.1.4 Kabsch algorithm

In order to make the paper self contained, we now summarize the

weighted Kabsch algorithm, which represents the closed-form differentiable solution to the ego-motion estimation problem:

$$\mathbf{R}_{\text{ego}}^*, \mathbf{t}_{\text{ego}}^* = \underset{\mathbf{R}_{\text{ego}}, \mathbf{t}_{\text{ego}}}{\operatorname{argmin}} \sum_{l=1}^{N^b} w_l \|\mathbf{R}_{\text{ego}} \mathbf{x}_l^b + \mathbf{t}_{\text{ego}} - \phi(\mathbf{x}_l^b, \mathbf{Y}^b)\|^2.$$

In the following, we omit the superscript b and define $\mathbf{q}_l := \phi(\mathbf{x}_l, \mathbf{Y})$ and the resulting point cloud of correspondences as $\mathbf{Q} \in \mathbb{R}^{3 \times N} = \{\mathbf{q}_l \in \mathbb{R}^3\}_l$. Let $\bar{\mathbf{x}}$ and $\bar{\mathbf{q}}$:

$$\bar{\mathbf{x}} := \frac{\sum_{l=1}^N w_l \mathbf{x}_l}{\sum_{l=1}^N w_l}, \quad \bar{\mathbf{q}} := \frac{\sum_{l=1}^N w_l \mathbf{q}_l}{\sum_{l=1}^N w_l} \quad (12)$$

denote the weighted centroids of point clouds \mathbf{X} and \mathbf{Q} , respectively. The centered point coordinates can then be computed as $\tilde{\mathbf{x}}_l := \mathbf{x}_l - \bar{\mathbf{x}}$ and $\tilde{\mathbf{q}}_l := \mathbf{q}_l - \bar{\mathbf{q}}$. By arranging them back to the matrix form $\tilde{\mathbf{X}} \in \mathbb{R}^{N \times 3}$ and $\tilde{\mathbf{Q}} \in \mathbb{R}^{N \times 3}$, a weighted covariance matrix $\mathbf{S} \in \mathbb{R}^{3 \times 3}$ can be computed as

$$\mathbf{S} = \tilde{\mathbf{X}}^T \mathbf{W} \tilde{\mathbf{Q}} \quad (13)$$

where $\mathbf{W} = \text{diag}(w_1, \dots, w_N)$. Considering the singular value decomposition $\mathbf{S} = \mathbf{U} \mathbf{\Sigma} \mathbf{V}^T$, the optimal rotation matrix is given by

$$\mathbf{R}_{\text{ego}}^* = \mathbf{V} \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & \det(\mathbf{V} \mathbf{U}^T) \end{bmatrix} \mathbf{U}^T \quad (14)$$

where $\det(\cdot)$ denotes computing the determinant and is used to avoid generating a reflection matrix. Finally, the optimal translation vector is recovered as

$$\mathbf{t}_{\text{ego}}^* = \bar{\mathbf{q}} - \mathbf{R}_{\text{ego}}^* \bar{\mathbf{x}}. \quad (15)$$

Note, when estimating the per-object transformation parameters of the foreground points, we use an unweighted Kabsch algorithm, i.e. $\mathbf{W} = \mathbf{I}$.

A.1.5 Entropy-regularized optimal transport

We provide a brief overview of the Sinkhorn algorithm used to approximate the optimal transport following [7]. Consider a discrete probability measure μ on the simplex $\Sigma_d \triangleq \{\mathbf{x} \in \mathbb{R}_+^n : \mathbf{x}^\top \mathbf{1}_n = 1\}$ with weights $\mathbf{a} = \{a_i\}$ and locations $\{x_i\}$, $i = 1 \dots n$ as:

$$\mu = \sum_{i=1}^n a_i \delta_{x_i}, \quad a_i \geq 0 \wedge \sum_{i=1}^n a_i = 1 \quad (16)$$

where δ_x is a Dirac delta function at x .

For two probability measures μ and ν in the simplex, let $U(\mu, \nu)$ denote the polyhedral set of $n_\mu \times n_\nu$ matrices:

$$U(\mu, \nu) = \{ \mathbf{P} \in \mathbb{R}_+^{n_\mu \times n_\nu} : \mathbf{P}\mathbf{1}_{n_\nu} = \mu \wedge \mathbf{P}^\top \mathbf{1}_{n_\mu} = \nu \}$$

where $\mathbf{1}_d$ is a d -dimensional vector of ones. $U(\mu, \nu)$ is also referred to as the *transportation polytope*.

Let \mathbf{C} be a $n_\mu \times n_\nu$ cost matrix that is constructed from the *ground cost* function $c(x_i^\mu, x_j^\nu)$. Kantorovich’s optimal transport formulation seeks to find the transport plan optimally mapping μ to ν :

$$\mathcal{W}_c(\mu, \nu) = d_c(\mu, \nu) = \min_{\mathbf{P} \in U(\mu, \nu)} \langle \mathbf{P}, \mathbf{C} \rangle \quad (17)$$

with $\langle \cdot \rangle$ denoting the Frobenius dot-product. Note that this expression is also known as the *Wasserstein distance* (WD). Due to the computational difficulties in minimizing Eq (17) Cuturi [14] introduced an alternative convex set consisting of joint probabilities with small enough mutual information:

$$U_\alpha(\mu, \nu) = \{ \mathbf{P} \in U(\mu, \nu) : KL(\mathbf{P} \parallel \mu\nu^\top) \leq \alpha \} \quad (18)$$

where $KL(\cdot)$ refers to the Kullback Leibler divergence [39]. This restricted polytope leads to a tractable distance between discrete probability measures under the cost matrix \mathbf{C} constructed from the function c :

$$d_{c,\alpha}(\mu, \nu) = \min_{\mathbf{P} \in U_\alpha(\mu, \nu)} \langle \mathbf{P}, \mathbf{C} \rangle. \quad (19)$$

$d_{c,\alpha}$ can be computed by a modification of simple matrix scaling algorithms, such as Sinkhorn’s algorithm [63, 64]. Note that here the cost matrix is inversely related to the affinity matrix defined in the main paper.

A.2. Datasets

In this work, a total of five datasets are used to train and evaluate the performance of the proposed approach. In the following, we detail the generation of the training and evaluation data. For all datasets, we use the same coordinate system centered at the camera or LiDAR sensor. The z -axis of the coordinate system points in the viewing direction of the camera or to the front of the car, the positive y -axis points in the up direction, and x completes the right-handed coordinate system. All processed point clouds are available under <https://3dsценарий.github.io/>.

FlyingThings3D [44]. FlyingThings3D is a large-scale synthetic dataset of stereo RGB images proposed to train 2D based scene flow networks. We follow [24] and use the subset of the dataset in which the extremely hard examples are omitted. The point clouds are obtained by lifting the stereo images to 3D using the annotated disparity maps, optical flow, and camera parameters. We again follow [24] and remove the points whose disparity and optical flow are occluded (occlusion maps are part of the original dataset). Additionally, points with a depth larger than 35 m are removed [24]. The resulting point clouds are in direct correspondence, i.e. $\mathbf{Y} = \mathbf{X} + \mathbf{V}$, where \mathbf{X} and \mathbf{Y} are the source and target point cloud, and \mathbf{V} is the ground truth flow. For training and evaluation, 8192 points are randomly sampled from each point cloud independently.

Dataset	Method	Supervision	EPE3D [m] ↓	Acc3DS ↑	Acc3DR ↑	Outliers ↓
<i>stereoKITTI</i> (w/o ground)	HPLFlowNet [24]	Full	0.117	0.478	0.778	0.410
	PointPWC-Net [80]	Full	0.069	0.728	0.888	0.265
	FLOT [55]	Full	0.056	0.755	0.908	0.242
	Ours (backbone)	Full	0.042	0.849	0.959	0.208
	Ours	Weak	0.163	0.541	0.658	0.452
	Ours++	Weak	0.134	0.709	0.800	0.311
<i>stereoKITTI</i> (with ground)	HPLFlowNet [24]	Full	0.238	0.194	0.429	0.787
	PointPWC-Net [80]	Full	0.204	0.292	0.556	0.645
	FLOT [55]	Full	0.122	0.480	0.691	0.401
	Ours (backbone)	Full	0.143	0.392	0.660	0.533
	Ours	Weak	0.136	0.470	0.712	0.420
	Ours++	Weak	0.068	0.836	0.897	0.263

Table 5: Evaluation results on *stereoKITTI* dataset. Weakly supervised models are trained on *semanticKITTI* LiDAR point clouds and evaluated directly on stereo point clouds of *stereoKITTI*. Ours and Ours++ denote the output of the network before and after test-time optimization, respectively.

stereoKITTI [45, 46]. This is a dataset based on the KITTI Scene Flow benchmark, which is designed for the evaluation of stereo based scene flow methods. It consists of 200 training and 200 test scenes. Because the ground truth annotations of the test scenes are not available, only the training scenes are used for evaluation. Due to incomplete and noisy annotations, 58 scenes are further removed and the final evaluation set comprises 142 scenes [40, 24]. The stereo images are lifted to point clouds analogously to FlyingThings3D [24]. Again, occluded points and points with a depth larger than 35 m are removed. For evaluation, we sample 8192 points from each point cloud independently. Tab. 1 shows the result of the typical evaluation protocol in which the ground points are removed by a naive thresholding the y coordinate at -1.4 m [40, 24] (the cameras in the KITTI setup are mounted at ≈ 1.65 m height above ground). In Tab. 5 we additionally report the results without removing the ground points.

lidarKITTI [45, 46]. This dataset comprises the Velodyne 64-beam LiDAR point clouds corresponding to the same 142 scenes used in *stereoKITTI*. The ground truth scene flow vectors are obtained by projecting the points of the source point cloud to the image plane using the provided calibration parameters, and associating them with the scene flow vectors of the corresponding pixels from the *stereoKITTI* dataset [41]. Due to directly using the LiDAR point clouds the points of both frames are not in direct correspondence and exhibit the typical LiDAR sampling pattern with very uneven point density.

semanticKITTI [5]. This dataset provides pointwise semantic labels annotated in 3D and improved ego-motion information for all sequences of the KITTI odometry benchmark [22]. Different to *stereoKITTI* and *lidarKITTI*, which include only the points that map to the image plane of the front camera, *semanticKITTI* contains full 360° LiDAR sweeps. We therefore process *semanticKITTI* point clouds to make them consistent with *lidarKITTI*. Specifically, we first convert the 3D coordinates of the points to polar coordinates and consider only the points with an azimuth angle in the range $[-45^\circ, 45^\circ]$ and elevation angle in the range $[-24.9^\circ, 12.0^\circ]$. This maintains the points that would roughly map to the image plane of the front camera. We then additionally remove the points that are less than 1.5 m or more than 35 m away from the LiDAR sensor. The binary background mask is

Dataset	Method	Supervision	scene flow				BG - segmentation				ego motion	
			mean EPE3D [m] ↓	med. EPE3D [m] ↓	med. F-EPE3D [m] ↓	med. B-EPE3D [m] ↓	prec.	FG ↑	rec. FG ↑	prec. BG ↑	recall. BG ↑	RRE [°] ↓
<i>lidarKITTI</i> (w/o ground)	Ours	Weak	0.150	0.111	0.227	0.104	0.726	0.885	0.978	0.940	0.379	0.130
	Ours+	Weak	0.110	0.064	0.227	0.049	0.726	0.885	0.978	0.940	0.126	0.053
	Ours++	Weak	0.094	0.051	0.164	0.049	0.726	0.885	0.978	0.940	0.126	0.053
<i>lidarKITTI</i> (with ground)	Ours	Weak	0.133	0.109	0.186	0.109	0.734	0.855	0.991	0.980	0.327	0.130
	Ours+	Weak	0.106	0.083	0.186	0.083	0.734	0.855	0.991	0.980	0.142	0.091
	Ours++	Weak	0.103	0.083	0.131	0.083	0.734	0.855	0.991	0.980	0.142	0.091

Table 6: Detailed results of our weakly supervised method on *lidarKITTI* dataset. Ours denotes the direct output of the network. Ours+ and Ours++ are the results with only background and full test-time optimization, respectively

generated by combining the semantic labels into the background (class labels from 40 to 249) and foreground (other class labels)⁶. We split the point cloud pairs of *semanticKITTI* into 4350 validation (sequences 03 and 05) and 18840 training samples (remaining nine sequences of the training set). Because BG-FG labels of the test dataset are not available, we perform all evaluations on the validation dataset.

waymo open [69]. This is a large scale dataset collected by a fleet of waymo self-driving cars in various conditions. It contains 1950 sequences of 20 s duration each, collected with an acquisition rate of 10 Hz. In our experiments we use the training batches 0–16 ($\approx 50\%$ of the training data) for fine-tuning (§ A.3.2) and the validation batches 0–2 ($\approx 40\%$ of the validation data) for evaluating our approach⁷. Waymo cars are equipped with five LiDAR sensors altogether, one mid-range LiDAR on the top and four short range ones on the sides of the car. In our experiments, we only use the points acquired by the top, mid-range scanner. We then transform these points into a coordinate system centered at the location of the LiDAR sensor in the KITTI setup and follow the processing steps used in *semanticKITTI* dataset to extract only the points that would roughly project to the front camera. Along with the ego-motion information, *waymo-open* also provides the 3D bounding-boxes of vehicles, pedestrians, cyclists, and signs. We use the former three classes to extract the foreground and consider the remaining points as background.

A.3. Evaluation details and additional results

We start this section by providing additional details and results (§ A.3.1) supporting the evaluations presented in Sec. 4, before reporting additional evaluations (§ A.3.2) that were omitted from the main paper due to the space constraint.

A.3.1 Evaluation details

Additional results on *stereoKITTI*. In § 4.3, we evaluate the performance of our backbone under full supervision on *FlyingThings3D* and *stereoKITTI*. Tab. 5 supplements that section with the evaluation results on the *stereoKITTI* dataset without removing the ground points. Additionally, we report the generalization results of our weakly supervised model trained on the LiDAR point clouds of *semanticKITTI*.

As expected, the performance of all fully supervised methods drops significantly if the challenging ground points are not

⁶<https://github.com/PRBonn/semantic-kitti-api/blob/master/config/semantic-kitti-all.yaml>

⁷<https://waymo.com/open/download/>

Dataset	Method	Initialization	semanticKITTI (w/o ground)				semanticKITTI (with ground)			
			EPE3D [m] ↓	Acc3DS ↑	Acc3DR ↑	Outliers ↓	EPE3D [m] ↓	Acc3DS ↑	Acc3DR ↑	Outliers ↓
<i>lidarKITTI</i> (w/o ground)	Ours	Pretrained	0.102	0.706	0.833	0.357	0.942	0.996	0.998	0.177
	Ours++	Pretrained	0.080	0.834	0.912	0.279	0.942	0.996	0.998	0.116
	Ours	Random	0.150	0.521	0.744	0.450	0.942	0.996	0.998	0.092
<i>lidarKITTI</i> (with ground)	Ours++	Random	0.094	0.784	0.885	0.314	0.942	0.996	0.998	0.059
	Ours	Pretrained	0.091	0.601	0.788	0.445	0.942	0.996	0.998	0.044
	Ours++	Pretrained	0.080	0.742	0.850	0.369	0.942	0.996	0.998	0.032
<i>semanticKITTI</i>	Ours	Random	0.133	0.460	0.746	0.527	0.942	0.996	0.998	0.032
	Ours++	Random	0.103	0.686	0.819	0.410	0.942	0.996	0.998	0.022
	Ours	Pretrained	0.950	0.892	0.991	0.996	0.942	0.996	0.998	0.035
<i>semanticKITTI</i>	Ours++	Pretrained	0.950	0.892	0.991	0.996	0.127	0.031	0.996	0.998
	Ours	Random	0.971	0.895	0.991	0.998	0.942	0.996	0.998	0.047
	Ours++	Random	0.971	0.895	0.991	0.998	0.942	0.996	0.998	0.032

Table 7: Evaluation of our model trained with random initialization of the weights (Random) and with the weights pretrained on *FT3D* (Pretrained), on *lidarKITTI* dataset. Ours and Ours++ denote the direct output of the network and the result after test-time optimization, respectively.

Method	Initialization	BG segmentation				ego-motion			
		prec. FG ↑	rec. FG ↑	prec. BG ↑	recall. BG ↑	RRE [°] ↓	RTE [m] ↓		
semanticKITTI (w/o ground)									
Ours	Pretrained	0.950	0.892	0.991	0.996	0.145	0.035		
Ours++	Pretrained	0.950	0.892	0.991	0.996	0.127	0.031		
Ours	Random	0.971	0.895	0.991	0.998	0.201	0.047		
Ours++	Random	0.971	0.895	0.991	0.998	0.133	0.032		
semanticKITTI (with ground)									
Ours	Pretrained	0.942	0.909	0.996	0.998	0.177	0.044		
Ours++	Pretrained	0.942	0.909	0.996	0.998	0.116	0.029		
Ours	Random	0.966	0.904	0.996	0.999	0.249	0.059		
Ours++	Random	0.966	0.904	0.996	0.999	0.121	0.032		

Table 8: Evaluation of our model trained with random initialization of the weights (Random) and with the weights pretrained on *FT3D* (Pretrained), on *semanticKITTI* dataset. Ours and Ours++ denote the direct output of the network and the result after test-time optimization, respectively.

removed. Remarkably, our weakly supervised model generalizes from LiDAR to stereo point clouds, and when combined with the test-time optimization even outperforms all methods on *stereoKITTI* with ground points. Tab. 5 also hints that the generalization LiDAR \mapsto stereo is less challenging than the opposite stereo \mapsto LiDAR.

Additional results on *lidarKITTI*. Tab. 6 supplements Tab. 2 and provides detailed results of our method on the *lidarKITTI* dataset. Note, how the test-time optimization improves background (Ours+) as well as foreground (Ours++) scene flow estimates. The results of the BG-segmentation and FG/BG scene flow split should be interpreted with caution, due to the noisy annotation of the *lidarKITTI* dataset (see below). Further qualitative results are shown in Fig. 6 and failure cases in Fig. 7

Noisy annotations of *lidarKITTI*. During the evaluation, we have discovered that *lidarKITTI* annotations are noisy and contain outliers. These GT errors occur especially around the points

Method	Mode	BG segmentation				ego-motion	
		prec. FG \uparrow	rec. FG \uparrow	prec. BG \uparrow	recall. BG \uparrow	RRE [$^\circ$] \downarrow	RTE [m] \downarrow
Ours++	generalization	0.960	0.689	0.957	0.996	0.141	0.099
Ours++	fine-tuned	0.945	0.921	0.989	0.992	0.111	0.078

Table 9: Comparison of the model fine-tuned on *waymo open* with the model trained only on *semanticKITTI* (generalization), on the *waymo open* dataset. Fine-tuned model outperforms the directly generalized one in terms of FG precision and ego-motion error.

lying on or close to the object boundaries. Outliers are caused by the projection of the 3D point onto the 2D image plane, where two distant points in 3D can map to the same pixel, as there is no perception of depth. Fig. 8 shows two prominent examples: in some cases, the instance mask and motion of one object are also assigned to the other (top), and in other cases, the background points get assigned the instance label and scene flow of the object in the front (bottom). Because of training on *semanticKITTI* with accurate annotations, our method is still capable of predicting the correct object masks (Fig. 8 (b)). Wrong annotation however result in a lower BG-segmentation performance and cause apparent errors in our scene flow prediction (see Fig. 8 (c)). Unfortunately, this error is reflected in the quantitative evaluations not only for our method but also for many other scene flow algorithms out there. To reduce this effect, we additionally report median *EPE3D* values in Tab. 6, as well as the FG and BG *EPE3D* based on our predicted BG-mask, instead of the noisy GT mask.

A.3.2 Additional evaluations and ablation studies

Pretraining vs training from scratch. Evaluations reported in § 4.4 to § 4.6 are performed using the weakly supervised model trained on *semanticKITTI* from scratch. However, recent works show that general 3D backbone networks can benefit from pre-training on large (annotated) datasets [81].

To evaluate this in our setting, we consider a model whose backbone weights were initialized with weights trained (with full supervision) on *FlyingThings3D* rather than randomly. The evaluations on *lidarKITTI* (Tab. 7) and *semanticKITTI* (Tab. 8) show that, in line with the literature, our backbone network can indeed benefit from pre-training. The improvement is especially prominent in scene flow estimation (more than 1 cm lower *EPE3D* in Tab. 7) and in ego-motion estimation (lower rotation and translation errors in Tab. 8). In order to fully adhere to the weakly supervised setting, we use the randomly initialized model in all other evaluations.

Fine-tuning on *waymo open*. In the main paper, *waymo open* dataset is used only to evaluate the direct generalization performance of our model trained on *semanticKITTI*. However, as briefly discussed in the main paper, *waymo open* also provides all the annotations that our weakly supervised model relies on. We now use these annotations, to evaluate the gain obtained by fine-tuning our model. To this end, we initialize our model with the weights trained on *semanticKITTI* and fine-tune it for 22k iterations (less than 2 epochs) on *waymo open*. Tab. 9 shows that the fine-tuned model greatly outperforms the model trained only on *semanticKITTI*, especially in terms of foreground recall (gain of more than 20 percent) and relative translation error (im-

GT inst. mask	<i>lidarKITTI</i> (with ground)				
	train	test	EPE [m] \downarrow	F-EPE [m] \downarrow	B-EPE [m] \downarrow
✓			0.097	0.216	0.085
	✓		0.101	0.183	0.094
✓	✓		0.097	0.265	0.085
			0.102	0.195	0.094
					0.139
					0.091

Table 10: Performance evaluation of our simple foreground clustering algorithm compared to models using GT instance mask during training and/or testing on *lidarKITTI* dataset.

	<i>lidarKITTI</i> (with ground)		
	EPE [m] \downarrow	RRE [$^\circ$] \downarrow	RTE [m] \downarrow
w/o Sinkhorn	0.594	0.539	0.615
with Sinkhorn	0.133	0.327	0.130

Table 11: Ablation study of the Sinkhorn algorithm on *lidarKITTI* dataset (evaluation protocol follows Tab. 3 in the main paper).

provement of 2 cm). The relative improvement is qualitatively also depicted in Fig. 9

Training and testing with GT instance masks. To define the FG instance-level rigidity loss and to perform the test-time optimization of the foreground, we rely on a simple unsupervised clustering of the foreground into individual objects. Arguably, this part could be replaced by a more powerful instance segmentation head, which would however require instance annotations during training. It is of interest to see how much benefit a *perfect* clustering would bring to our method. We therefore ablate by replacing the output of our clustering with the GT instance labels, which are provided both in *semanticKITTI* and *lidarKITTI*. We assess the individual contributions from having GT labels during training, testing or both. Tab. 10 shows that our model with the simple clustering algorithm performs comparably to the models using GT instance masks during training (*semanticKITTI*) and/or testing (*lidarKITTI*). We conclude that our clustering algorithm, which does not require GT instance masks during training or testing, is the preferred option. Due to the noisy instance annotations of the *lidarKITTI* dataset, the performance of models using GT instance masks during testing should be interpreted with caution.

Sinkhorn algorithm. The benefit of using the entropy-regularized Sinkhorn algorithm for the ego-motion estimation is two fold: (i) the Sinkhorn distances should yield more accurate correspondences due to the optimality of the transport map, and (ii) in combination with the slack row and column it enables us to down-weight the outliers in a principled manner. We empirically confirm these benefits in an ablation study in which we directly use the affinity matrix \mathbf{M} to compute the soft correspondences $\phi(\mathbf{x}_i^b, \mathbf{Y}^b) = \mathbf{Y}^b \mathbf{m}_i / \|\mathbf{m}_i\|_1$ and estimate the ego-motion with an unweighted Kabsch algorithm (§ A.1.4).

Tab. 11 depicts a significant increase of the EPE when Sinkhorn is deactivated. This increase can be accredited to the inferior estimation of the ego-motion parameters, which results in a much higher RTE and RRE.

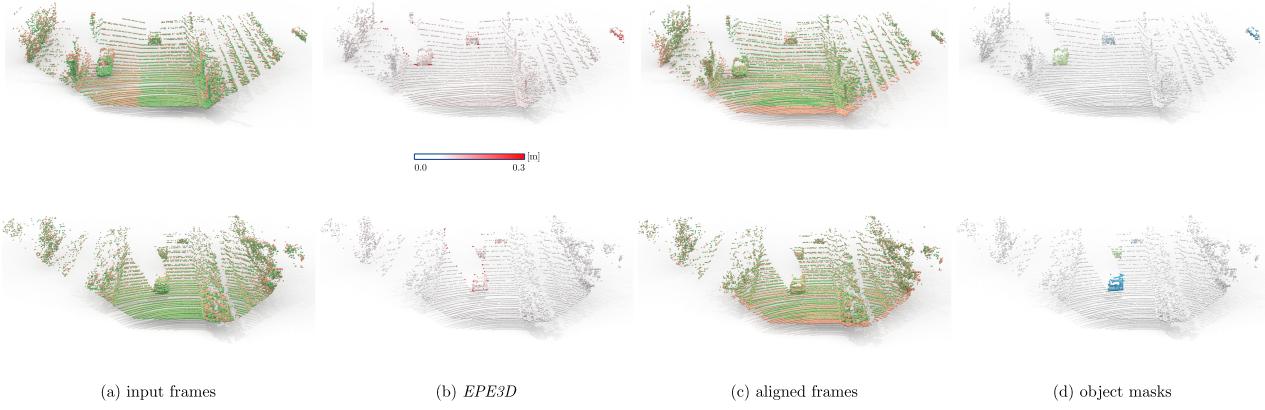


Figure 6: Successful cases of our method on the *lidarKITTI* dataset. By correctly splitting the scene into foreground and background (d), our method estimates the accurate scene flow vectors (b), which align the two frames (c).

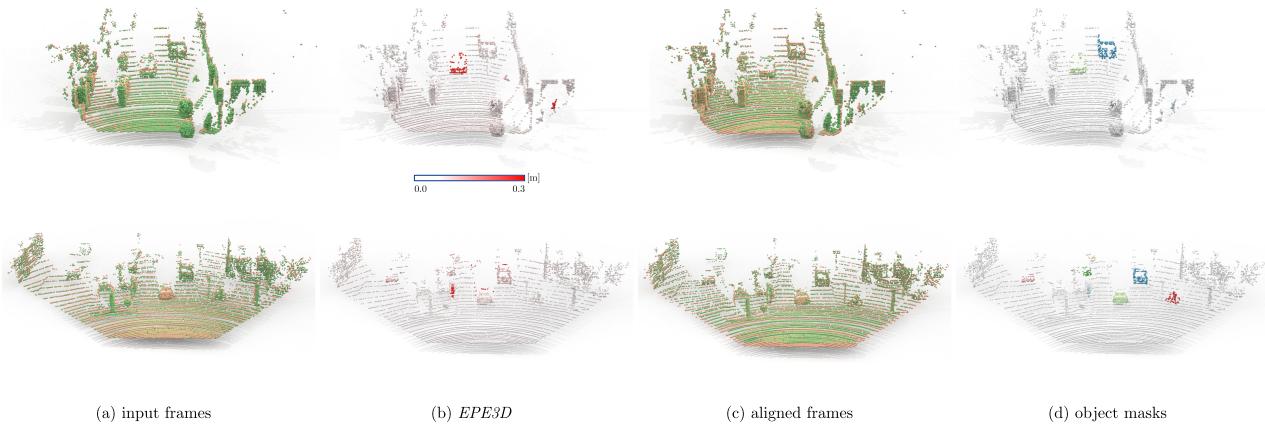


Figure 7: Failure cases of our method on the *lidarKITTI* dataset. **Top:** even though the car’s object mask (d) is correctly predicted, its predicted scene flow vectors yield large end-point-errors (b). **Bottom:** a pillar in the middle of the scene is wrongly predicted as foreground object (d), hence its scene flow does not agree with the background and GT (b).

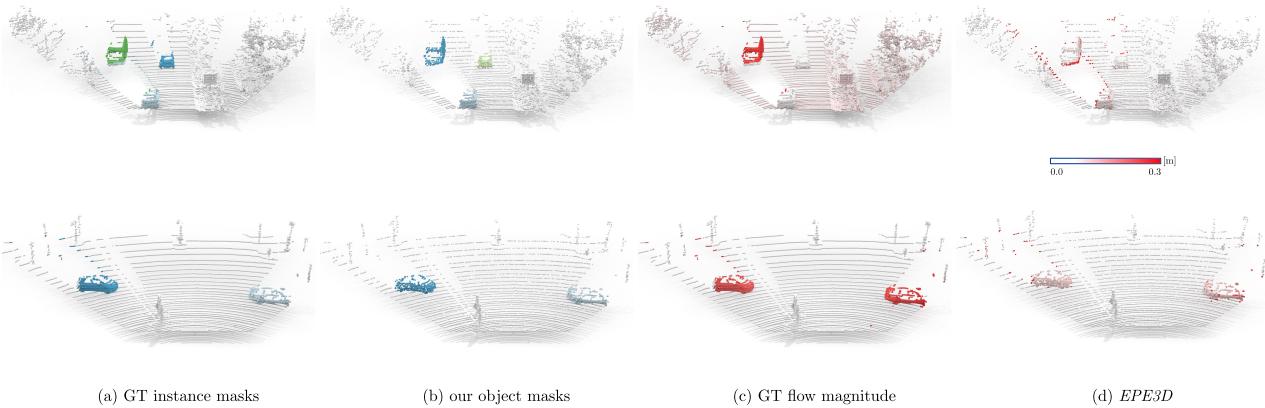


Figure 8: *lidarKITTI* annotations are obtained by projecting 3D points onto the image plane, which results in wrong instance (a) and scene flow (c) annotations for points with azimuth and elevation angles close to the object boundaries (e.g. green car is partially blue in (a) top). On the other hand, our method infers correct object masks (b) and scene flow (d), yet due to the wrong GT annotations, the scene flow appears to be erroneous.



Figure 9: Fine-tuning on *waymo open* improves performance and robustifies our model. **Top:** objects close to the sensor are not common in *semanticKITTI* and hence cannot be detected correctly by the generalized model (c). **Bottom:** a challenging example with 18 foreground objects (much larger than average number in *semanticKITTI*). Note, how more object masks are correctly inferred by our fine-tuned model compared to direct generalization (c-column).

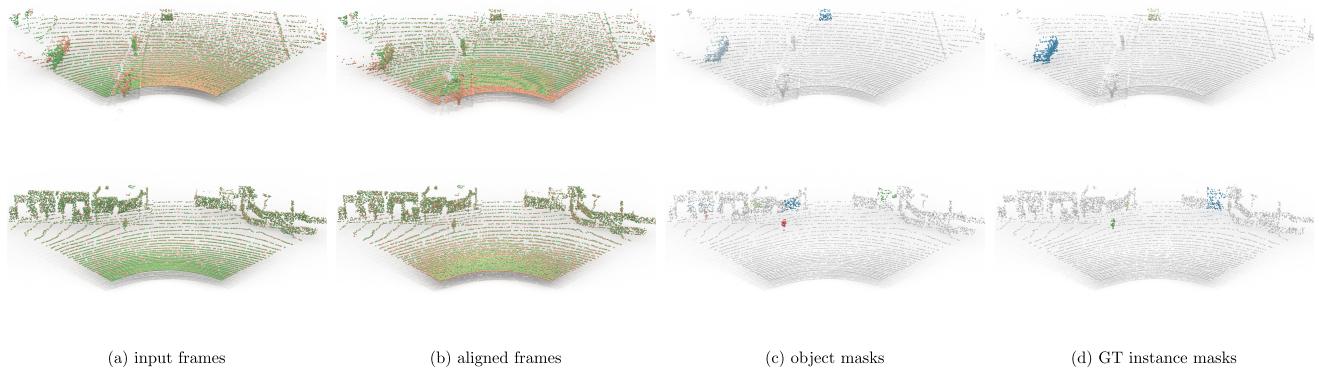


Figure 10: Failure cases on *waymo open* dataset. **Top:** our model is unable to estimate accurate ego-motion and scene flow (b) if the background points consists only of the ground points after foreground removal (c). **Bottom:** rare objects such as trucks (top right corner in c and d) appear ambiguous to our model and cause prediction of the wrong masks (c).