

DeepMapping: Unsupervised Map Estimation From Multiple Point Clouds

Li Ding^{*†}

l.ding@rochester.edu

Chen Feng^{*‡§}

cfeng@nyu.edu

[†]University of Rochester [‡]NYU Tandon School of Engineering [§]Mitsubishi Electric Research Laboratories (MERL)

Abstract

We propose *DeepMapping*, a novel registration framework using deep neural networks (DNNs) as auxiliary functions to align multiple point clouds from scratch to a globally consistent frame. We use DNNs to model the highly non-convex mapping process that traditionally involves hand-crafted data association, sensor pose initialization, and global refinement. Our key novelty is that “training” these DNNs with properly defined unsupervised losses is equivalent to solving the underlying registration problem, but less sensitive to good initialization than ICP. Our framework contains two DNNs: a localization network that estimates the poses for input point clouds, and a map network that models the scene structure by estimating the occupancy status of global coordinates. This allows us to convert the registration problem to a binary occupancy classification, which can be solved efficiently using gradient-based optimization. We further show that *DeepMapping* can be readily extended to address the problem of Lidar SLAM by imposing geometric constraints between consecutive point clouds. Experiments are conducted on both simulated and real datasets. Qualitative and quantitative comparisons demonstrate that *DeepMapping* often enables more robust and accurate global registration of multiple point clouds than existing techniques. Our code is available at <https://ai4ce.github.io/DeepMapping/>.

1. Introduction

Advances in deep learning have led to many state-of-the-art methods for semantic computer vision tasks. Despite those successes, their compelling improvements on the geometric aspects of computer vision are yet to be fully demonstrated (especially for registration and mapping). This is perhaps because powerful deep semantic representations have limitations in accurately estimating and modeling of geometric attributes of the environment. This includes, but is not limited to, estimating camera motions from a sequence of images, or registering multiple point clouds into a complete model. As opposed to semantic attributes of

^{*}This work was partially done while the authors were with MERL. And Chen Feng is the corresponding author.

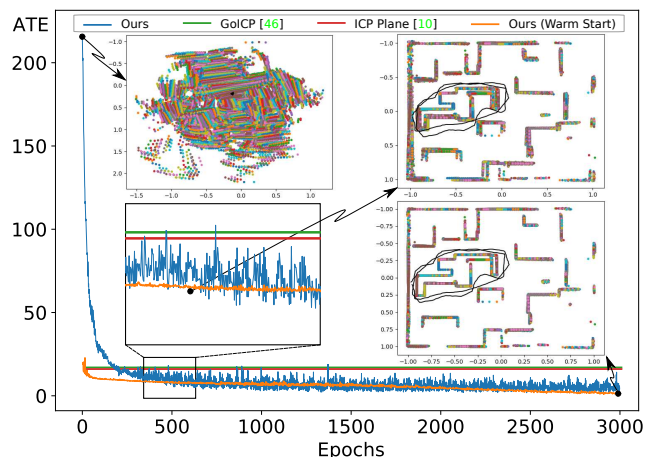


Figure 1. DeepMapping achieves better registration quality than other baselines on an example dataset. Best viewed in color.

objects/scenes that are often categorical and thus easily described by human language, those geometric attributes are more often continuous and can be better described numerically, such as poses and shapes. These spatial properties and relations between objects play as vital roles as semantic ones in robotics, augmented reality, medical, and other engineering applications.

Several works attempt to integrate deep learning methods into geometric vision problems [45, 52, 49, 20, 19, 29, 18]. Methods in [26, 9, 20] try to regress camera poses by training a DNN, inside which a map of the environment is implicitly represented. Methods in [45, 52] propose unsupervised approaches that exploit inherent relationships between depth and motion. Despite different tasks, most approaches follow the same train-and-test pipeline that neural networks are first learned from a set of training data (either supervised or unsupervised), and then evaluated on a testing set, expecting those DNNs to be able to generalize as much as possible to untrained situations.

The essence of our discussion is an open question: Will DNNs generalize well for geometric problems especially for registration and mapping? Semantic tasks can benefit from DNNs because those related problems are defined em-

pirically, and thus modeled and solved statistically. However many geometric problems are defined theoretically, and thus experiential solutions may not be adequate in terms of accuracy. Think of a simple scenario: given two images with overlapping field-of-views (FOV), without careful calculation, how accurate would a normal person be able to tell the Euclidean distance between the two camera centers? One may argue that reasonable accuracy can be achieved given adequate training. But if this means that it requires a large amount of data collection and training at each new location, the efficiency of this solution seems to be debatable.

Here we investigate another possibility of adopting powerful DNNs for the mapping/registration task. What we commonly agree from abundant empirical experiments and some theorems is that DNNs can model many arbitrarily complex mappings, and can be efficiently optimized through gradient-based methods, at least for categorical classification problems. This leads to our key idea in this paper: we convert the conventionally hand-engineered mapping/registration processes into DNNs, and solve them as if we are “training” them, although we do not necessarily expect the trained DNNs to generalize to other scenes. To make this meaningful, unlike the supervised training in [26, 9, 20], we need to properly define unsupervised loss functions that reflect the registration quality. Our exploration towards this line of thought shows promising results in our experiments, as shown in Figure 1. We summarize our contributions as follows:

- We propose DeepMapping to solve the point cloud mapping/registration problem as unsupervised end-to-end “training” of two DNNs, which is easier for parallel implementation compared to conventional methods requiring hand-crafted features and data associations.
- We convert this continuous regression problem to binary classification without sacrificing registration accuracy, using the DNNs and unsupervised losses.
- We demonstrate experimentally that DeepMapping is less sensitive to pose initialization compared with conventional baselines.

2. Related Work

Pairwise local registration: the methods for pairwise point cloud registration can be generally categorized into two groups: local vs. global methods. The local methods assume that a coarse initial alignment between two point clouds and iteratively update the transformation to refine the registration. The typical methods that fall into this category are the Iterative Closest Point (ICP) algorithms [6, 10, 34], probabilistic-based approaches [23, 31, 12] that model the point clouds as a probability distribution. The local methods are well-known for requiring a “warm start”, or a good initialization, due to limited convergence range.

Pairwise global registration: the global methods [46,

3, 30, 50, 28, 14] do not rely on the “warm start” and can be performed on point clouds with arbitrary initial poses. Most global methods extract feature descriptors from two point clouds. These descriptor are used to establish 3D-to-3D correspondences for relative pose estimation. Robust estimations, e.g., RANSAC [17], are typically applied to handle the mismatches. The feature descriptors are either hand-crafted such as FPFH [35], SHOT [43], 3D-SIFT [38], NARF [40], PFH [36], spin images [25], or learning-based such as 3DMatch [48], PPFNet [13], and 3DFeatNet [24].

Multiple registration: in addition to pairwise registration, several methods have been proposed for multiple point clouds registration [42, 16, 22, 44, 11]. One approach is to incrementally add new a point cloud to the model registered from all previous ones. The drawback of the incremental registration is the accumulated registration error. This drift can be mitigated by minimizing a global cost function over a graph of all sensor poses [11, 42].

Deep learning approaches: recent works explore the idea of integrating learning-based approaches into mapping and localization problems. Methods in [45, 52] propose unsupervised approaches that exploit inherent relationships between depth and motion. This idea is further explored in [7, 49, 47, 29] using deep learning for visual odometry and SLAM problems. For instance, CodeSLAM [7] represents the dense geometry using a variational auto-encoder (VAE) for depth that is conditioned on the corresponding intensity image, which is later optimized during bundle adjustment. Differently, DeepMapping does not require any pre-training. [18] introduces a generative temporal model with a memory system that allows the agent to memorize the scene representation and to predict its pose in a partially observed environment. Although this method and DeepMapping are both able to determine the sensor pose from the observed data, [18] requires a supervised training stage for “loading” its memory, while DeepMapping is fully unsupervised and does not follow the aforementioned train-and-test pipeline. Methods in [20, 32] use the recurrent neural network (RNN) to model the environment through a sequence of images in a supervised setting. MapNet [20], for example, develops a RNN for RGB-D SLAM problem where the localization of camera sensor is performed using deep template matching on the discretized spatial domain that has a relatively small resolution. Unlike MapNet, the proposed DeepMapping does not require any partition of the space and is unsupervised.

Other related methods such as [26, 9] solve camera localization by training DNNs to regress camera poses and test the performance in the same environment as the training images. Related to that is DSAC [8] as a differentiable alternative to traditional RANSAC for its use in pose estimation DNNs. For place recognition, semantic scene completion is used in [37] as an auxiliary task for training an im-

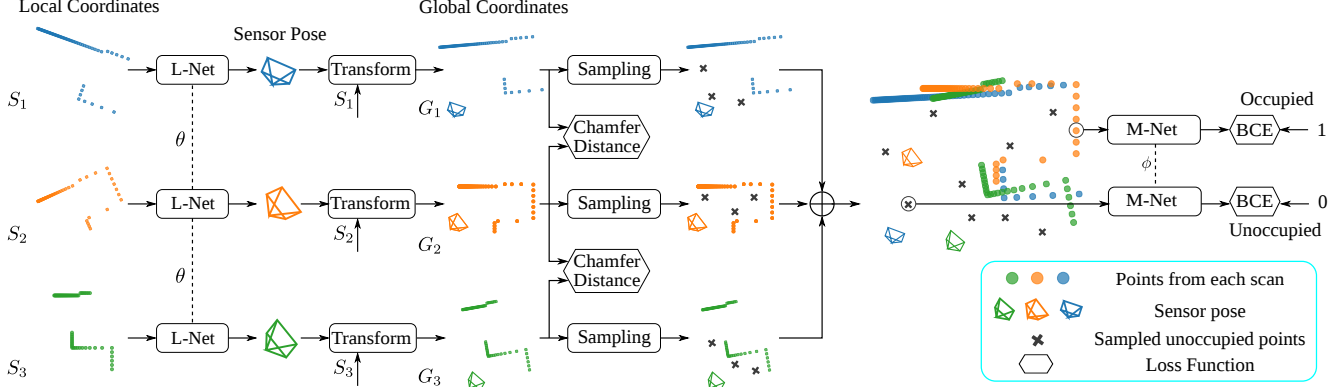


Figure 2. **DeepMapping pipeline.** Point clouds appear in different colors. Each input point cloud is fed into the shared L-Net to compute transformation parameters that map the input to the global coordinates where both occupied (colored solid circles) and unoccupied (gray cross marks) locations are sampled. The M-Net predicts the occupancy probabilities of the sampled locations. The global occupancy loss is the binary cross entropy (BCE) averaged over all sampled locations. DeepMapping is able to handle temporal information, if available, by integrating Chamfer distance loss between consecutive scans. Best viewed in color.

age VAE for long-term robustness. The method in [19] proposes an unsupervised approach with variational Bayesian convolutional auto-encoder to model structures from point clouds. In DeepMapping, we adopt this idea to model the scene structure but use DNN rather than Bayesian inference. Other prior works include: in [15] the generative query network (GQN) shows the ability to represent the scene from a given viewpoint and rendering it from an unobserved viewpoint in the simple synthetic environments. A neuroscience study [5] uses recurrent neural networks to predict mammalian spatial behavior.

As noted, most approaches follow a train-and-test pipeline. Our approach adopts DNNs but differs from the existing methods in the way that the process of “training” in DeepMapping is equivalent to solving the point clouds registration and that once trained, we do not expect the DNNs to generalize to other scenes.

3. Method

3.1. Overview

In this section, we describe the proposed DeepMapping that uses DNNs for registering multiple point clouds. Let $\mathbf{S} = \{S_i\}_{i=1}^K$ be the set of K input point clouds in the D -dimensional space that are captured by Lidar scanners, and the i^{th} point cloud S_i , represented as a $N_i \times D$ matrix, contains N_i points in sensor local frame. Given K point clouds, the goal is to register all point clouds in a common coordinate frame by estimating the sensor poses $\mathbf{T} = \{T_i\}_{i=1}^K$ for each point cloud S_i , where $T_i \in SE(D)$.

Conventional methods [50, 16] formulate this as an optimization problem that directly seeks the optimal sensor poses \mathbf{T} to minimize the loss function

$$\mathbf{T}^*(\mathbf{S}) = \arg \min_{\mathbf{T}} \mathcal{L}(\mathbf{T}, \mathbf{S}), \quad (1)$$

where $\mathcal{L}(\mathbf{T}, \mathbf{S})$ is the objective that scores the registration quality. As explained in Section 1, instead of directly optimizing \mathbf{T} , we propose to use a neural network, modeled as an auxiliary function $f_{\theta}(\mathbf{S})$, to estimate \mathbf{T} for the input point clouds \mathbf{S} where θ are the auxiliary variables to be optimized. The pose is then converted to a transformation matrix that maps S_i into its global version G_i .

Formally, we re-formulate this registration problem as finding the optimal network parameters that minimize a new objective function

$$(\theta^*, \phi^*) = \arg \min_{\theta, \phi} \mathcal{L}_{\phi}(f_{\theta}(\mathbf{S}), \mathbf{S}), \quad (2)$$

where \mathcal{L}_{ϕ} is an unsupervised learnable objective function which will be explained in Section 3.2 and 3.4.

Figure 2 illustrates the pipeline for DeepMapping. At the heart of DeepMapping are two networks, a localization network (L-Net) and an occupancy map network (M-Net), that estimates T_i and measures the registration quality, respectively. The L-Net is a function $f_{\theta} : S_i \mapsto T_i$ appeared in (2) that estimates the sensor pose of a corresponding point cloud, where the parameters θ in the L-Net are shared among all point clouds. The point cloud G_i in global coordinates is obtained using the estimated sensor pose. From the transformed point clouds, we first sample both occupied and unoccupied locations. Then the locations of these samples are fed into the M-Net, to evaluate the registration performance of the L-Net. The M-Net is a binary classification network that predicts probabilities of input locations being occupied. We denote M-Net as a function with learnable parameters ϕ . Those occupancy probabilities are used for computing the unsupervised loss \mathcal{L}_{ϕ} in (2) that measures the global occupancy consistency of the transformed point clouds and thus reflects the registration quality.

One may find that transforming from (1) to (2) could

increase the dimensionality/complexity of the problem. In fact, we provide a simple 1D version of this problem conversion and show that using DNNs as auxiliary functions and optimizing them in higher dimensional space using gradient-based methods could enable faster and better convergence than directly optimizing the original problem. Details are included in the supplementary material.

3.2. DeepMapping Loss

We use the M-Net m_ϕ to define the unsupervised loss function \mathcal{L}_ϕ that scores the registration quality. The M-Net is a continuous occupancy map $m_\phi : \mathbb{R}^D \rightarrow [0, 1]$ that maps a global coordinate to the corresponding occupancy probability, where ϕ are learnable parameters. If a coordinate in the global frame is associated with a binary occupancy label y indicating whether the location is occupied, then we can calculate the loss of a global coordinate as the binary cross entropy (BCE) B between the predicted occupancy probability p and the label y :

$$B[p, y] = -y \log(p) - (1 - y) \log(1 - p). \quad (3)$$

The question is how to determine the label y . We approach this question by considering the situation when all point clouds are already globally aligned precisely. In such a case, it is obvious that all observed/scanned points, in the global frame, should be marked as occupied with the label 1, due to the physical principle of the Lidar scanner.

It is also clear now that points lie between the scanner center and any observed points, i.e., on the line of sight, should be marked as unoccupied with label 0. Figure 3 depicts the mechanism to sample unoccupied points for Lidar scanners. The dash lines show the laser beams emitted from the scanner. We denote $s(G_i)$ as a set of sampled points from G_i that lie on these laser beams, illustrated as cross marks. These points are associated with label 0 indicating unoccupied locations.

By combining the binary cross entropy and the sampling function, the loss used in (2) is defined as an average of the binary cross entropy over all locations in all point clouds:

$$\mathcal{L}_{cls} = \frac{1}{K} \sum_{i=1}^K B[m_\phi(G_i), 1] + B[m_\phi(s(G_i)), 0], \quad (4)$$

where G_i is a function of L-Net parameters θ , and with a slight abuse of notation $B[m_\phi(G_i), 1]$ denotes the average BCE error for all points in point cloud G_i , and $B[m_\phi(s(G_i)), 0]$ means the average BCE error for correspondingly sampled unoccupied locations.

In Figure 3, we illustrate the intuition behind (4). The loss function can achieve smaller values if registrations are accurate, as shown in Figure 3 (a). Conversely, since misaligned point clouds will lead to self-contradictory occupancy status, the loss function will be larger due to the difficulties of the M-Net to model such noisy occupancy patterns, as shown in Figure 3 (b).

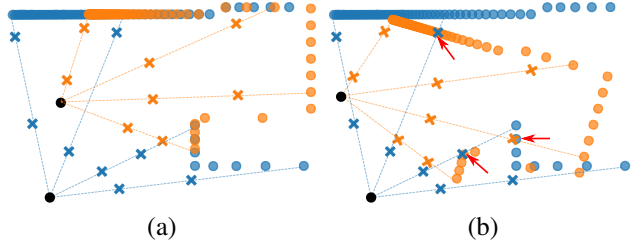


Figure 3. Illustration of sampling methods and self-contradictory occupancy status. Blue and orange circles are two point clouds in the global coordinates and cross marks represent the sampled unoccupied points. (a) and (b) show the correctly aligned and misaligned point clouds, respectively. The red arrows in (b) highlight the points with self-contradictory occupancy status. Best viewed in color.

Note that the loss in (4) is unsupervised because it relies only on the inherent occupancy status of point clouds rather than any externally labeled ground truth. For minimization, we adopt the gradient-based optimization because the loss function \mathcal{L}_{cls} is differentiable for both θ and ϕ . In each forward pass, unoccupied locations are randomly sampled.

It is also worth noting the potentials of M-Net in robotics. Unlike prior works that use discrete occupancy maps, we directly feed the M-Net with the floating number coordinates, resulting in a continuous occupancy map. Therefore, the M-Net offers a new way to represent the environment at an arbitrary scale and resolution.

3.3. Network Architecture

L-Net: the goal of the localization network, L-Net, is to estimate the sensor pose T_i in a global frame. The L-Net consists of a latent feature extraction module followed by a multi-layer perceptron (MLP) that outputs sensor pose T_i . This module depends on the format of the input point cloud S_i . If S_i is an *organized point cloud* that is generated from range images such as disparity maps or depth images, then S_i is organized as an array of points in which the spatial relationship between adjacent points is preserved. Therefore, we apply a convolutional neural network (CNN) to extract the feature vector of point cloud followed by a global pooling layer to aggregate local features to a global one. In the case where the inputs are a set of *unorganized point clouds* without any spatial order, we adopt PointNet [33] architecture for extracting features from the point cloud. We remove the input and feature transformations appeared in [33] and use a shared MLP that maps each D -dimensional point coordinate to a high dimensional feature space. A global pooling layer is applied across all points for feature aggregation. The extracted latent feature vector is processed with an MLP with the output channels corresponding to the degree of freedom of sensor movement.

M-Net: the occupancy map network, M-Net, is a binary classification network that takes as input a location coordinate.

dinate in the global space and predicts the corresponding occupancy probability for each input location. The M-Net is an MLP with a D -channel input and a 1-channel output from a sigmoid function, shared across all points.

3.4. Extension to Lidar SLAM

The loss function in (4) treats the input point clouds \mathbf{S} as an unordered set of scans instead of a temporally ordered sequence. In some applications, the temporal information may be available. For example, Lidar SLAM uses laser scanners to explore the unknown environment and captures an ordered sequence of point clouds at different time t .

Now we extend DeepMapping to exploit such a temporal relationship. The underlying assumption is that the consecutive scans of point clouds are expected to have a reasonable overlapping with each other, which normally holds in the SLAM settings [41, 39]. We utilize the geometric constraints among point clouds that are temporally close to each other. Specifically, we adopt the Chamfer distance as the metric that measures the distance between two point clouds X and Y in the global coordinates

$$d(X, Y) = \frac{1}{|X|} \sum_{\mathbf{x} \in X} \min_{\mathbf{y} \in Y} \|\mathbf{x} - \mathbf{y}\|_2 + \frac{1}{|Y|} \sum_{\mathbf{y} \in Y} \min_{\mathbf{x} \in X} \|\mathbf{x} - \mathbf{y}\|_2. \quad (5)$$

The Chamfer distance measures the two-way average distance between each point in one point cloud to its nearest point in the other. The minimization of the Chamfer distance $d(G_i, G_j)$ results in a pairwise alignment between point clouds G_i and G_j . To integrate the Chamfer distance into DeepMapping, we modify the objective in (2) as

$$(\theta^*, \phi^*) = \arg \min_{\theta, \phi} \mathcal{L}_{cls} + \lambda \mathcal{L}_{ch}, \quad (6)$$

where λ is a hyperparameter to balance the two losses and \mathcal{L}_{ch} is defined as the average Chamfer distance between each point cloud G_i and its temporal neighbors G_j

$$\mathcal{L}_{ch} = \sum_{i=1}^K \sum_{j \in \mathcal{N}(i)} d(G_i, G_j). \quad (7)$$

3.5. Warm Start

Optimizing (6) with random initialization of network parameters (i.e., “cold start”) could take a long time to converge, which is undesirable for real-world applications. Fortunately, DeepMapping allows for seamless integration of a “warm start” to reduce the convergence time with improved performance. Instead of starting from scratch, we can first perform a coarse registration of all point clouds using any existing methods, such as incremental ICP, before further refinement by DeepMapping. Figure 1 shows an example

result of registration error versus optimization iterations. With such a warm start, DeepMapping converges much faster and more stable than the one starting from scratch and is more accurate and robust.

4. Experiments

We evaluate DeepMapping on two datasets: a simulated 2D Lidar point cloud dataset and a real 3D dataset called Active Vision Dataset (AVD) [4]. We implemented DeepMapping with PyTorch [2]. The network parameters are optimized using Adam optimizer [27] with a learning rate of 0.001 on an Nvidia TITAN XP.

For quantitative comparison, we use two metrics: the absolute trajectory error (ATE) and the point distance between a ground truth point cloud and a registered one. The ATE assesses the global consistency of the estimated trajectory, and the point distance is the average distance of corresponding points between the ground truth and the registered point cloud. Since the estimated position can be defined in an arbitrary coordinate, a rigid transformation is determined to map the points in estimated coordinates to the ground truth coordinates using the closed-form solution proposed in [21]. The metrics are calculated after such a registration.

4.1. Experiments on 2D Simulated Point Cloud

Dataset: to simulate the 2D Lidar point clouds captured by a virtual Lidar scanner, we create a large environment represented by a 1024×1024 binary image, as shown in Figure 4. The white pixels indicate the free space whereas the black ones correspond to obstacles. We assume a moving robot equipped with a horizontal Lidar scanner to explore this environment. We first sample the trajectory of the robot movement that consists of a sequence of poses. The rotation perturbation between two scans are randomly selected from -10° to 10° , and the translation perturbation on average is approximately 8.16 pixels. At each pose in the trajectory, the Lidar scanner spreads laser beams over the FOV and the laser beam is reflected back whenever it hits the obstacles within its path. Each observed point is computed as the intersection point between the laser ray and the boundary of obstacles. The scanning procedure yields a set of 2D points in the robot’s local coordinate frame.

In this experiment, we set the FOV of the Lidar scanner to 360° and assume an ideal scanner with an unlimited sensing range¹. If no object exists along the laser ray, we return the point where the ray hits the image boundary. We create 3 binary maps of the virtual environment and sample in total 75 trajectories. 50 trajectories have 128 poses and the remaining 25 trajectories contain 256 poses where point clouds are scanned. Each scan contains 256 points.

¹We tested the robustness of DeepMapping by imposing a limited sensing range, while our resulting performance does not change significantly.

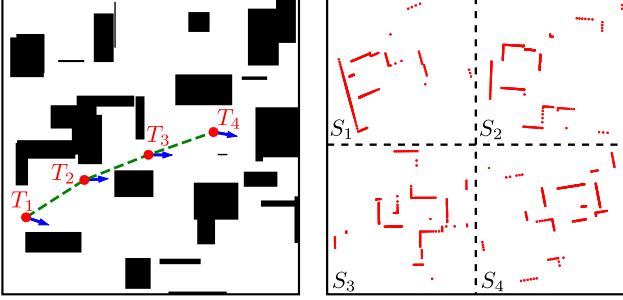


Figure 4. Illustration of 4 scans captured at time 20, 70, 120, 170. Left: the binary image shows the simulated environment with a size of 1024×1024 . The sensor poses are shown in red circles and blue arrows. The green dash line indicates the trajectory of the sensor. Right: observed point clouds captured at corresponding poses. Best viewed in color.

The transformation T_i is parameterized by a 2D translation vector (t_x, t_y) and a rotation angle α .

Baseline: we compare DeepMapping with the following baselines: incremental iterative closest point (ICP) with point-to-point distance metrics [6] and point-to-plane distance metrics [10], GoICP [46], and a particle swarm optimization PSO [1]. To justify the advantage of the neural network based optimization, we also test another baseline method, referred to as the direct optimization. Specifically, we remove the L-Net f_θ and perform the optimization of the same loss function with respect to the sensor poses \mathbf{T} rather than network parameters θ .

Implementation: the detailed architecture of the L-Net consists of C(64)-C(128)-C(1024)-M(1024)-FC(512)-FC(256)-FC(3), where C(n) denotes 1D atrous convolutions that have kernel size 3, dilation rate of 2 and n -channel outputs, M(n) denotes 1D global max-pooling over n channels, FC(n) denotes fully-connected layer with n -channel output. The M-Net can be described as FC(64)-FC(512)-FC(512)-FC(256)-FC(128)-FC(1). We did not optimize the network architectures. More ablation studies are included in the supplementary material. To sample unoccupied points, we randomly select 19 points on each laser ray. We run PSO for multiple rounds until they consume the same time as DeepMapping. To compare DeepMapping with the direct optimization, we use the same optimizer and the learning rate. To ensure the same initialization between DeepMapping and the direct optimization, we run DeepMapping with only on forward pass (no back-propagation is performed) and save the output sensor poses. These poses are used as the starting point for the direct optimization. The hyper-parameter λ is chosen based on the prior knowledge of the dataset. For the 2D dataset where sequential frames have reasonable overlaps, we assign more weights to Chamfer loss by setting λ to 10. The batch size is 128, and the optimization speed is approximately 100 epochs per minute on an Nvidia TITAN XP.

Results: Figure 5 shows the qualitative comparison of 3 trajectories simulated from the dataset. As shown, the baseline method, direct optimization on sensor pose, fails to find the transformations, leading to globally inaccurate registration. While the incremental ICP algorithms are able to procedure a noisy registration, errors of the incremental ICP algorithms accumulate over frames. As opposed to the baselines, the proposed DeepMapping accurately aligns multiple point clouds by utilizing neural networks for the optimization. Once the registration is converged, we feed all points (both occupied and unoccupied) in the global coordinates into the M-Net to estimated the occupancy probability, which is then converted to an image of occupancy map shown in the third column in Figure 5. The unexplored regions are shown in gray. The estimated occupancy map agrees with registered point clouds.

The box plots in Figure 6 show the quantitative results of the ATE and the point distance. Note that the data is plotted with the logarithmic scale for the y-axis. DeepMapping has the best performance in both metrics, achieving a median ATE of 5.7 pixels and a median point distance of 6.5, which significantly outperforms the baseline methods.

Generalization of L-Net: while we do not expect the trained DeepMapping to generalize to other scenes, it is worth emphasizing that, to some extent, the “trained” L-Net is able to perform coarse re-localization for an unseen point cloud that is captured close to the “training” trajectory and has a similar orientation. For demonstration, we simulate local point clouds at all possible positions (e.g., all white pixels in Figure 4) and use the “trained” L-Net to estimate the sensor location. Figure 7 shows the re-localization errors of two L-Nets “trained” on different trajectories, thus having different generalization abilities. Only the point clouds captured close to the “training” trajectories have low errors. The ability of coarse re-localization without “re-training” has potential usages in robotic applications, such as the kidnapped robot problem.

4.2. Experiments on the Active Vision Dataset

Dataset: DeepMapping is also tested on a real 3D dataset: Active Vision Dataset [4] that provides RGB-D images for a variety of scans of the indoor environment. Unlike other RGB-D datasets [41, 39] that capture RGB-D videos around the scene, the AVD uses a robotic platform to visit a set of discrete points on a rectangular grid with a fixed width of 300mm. At each point, 12 views of images are captured by rotating the camera every 30° . This data collection procedure leads to a low frame rate and relatively small overlapping between two images.

The ground truth camera intrinsic and extrinsic parameters are provided. We use the camera intrinsic parameters to covert the depth map to point cloud representation and do not use any information provided by color images. We

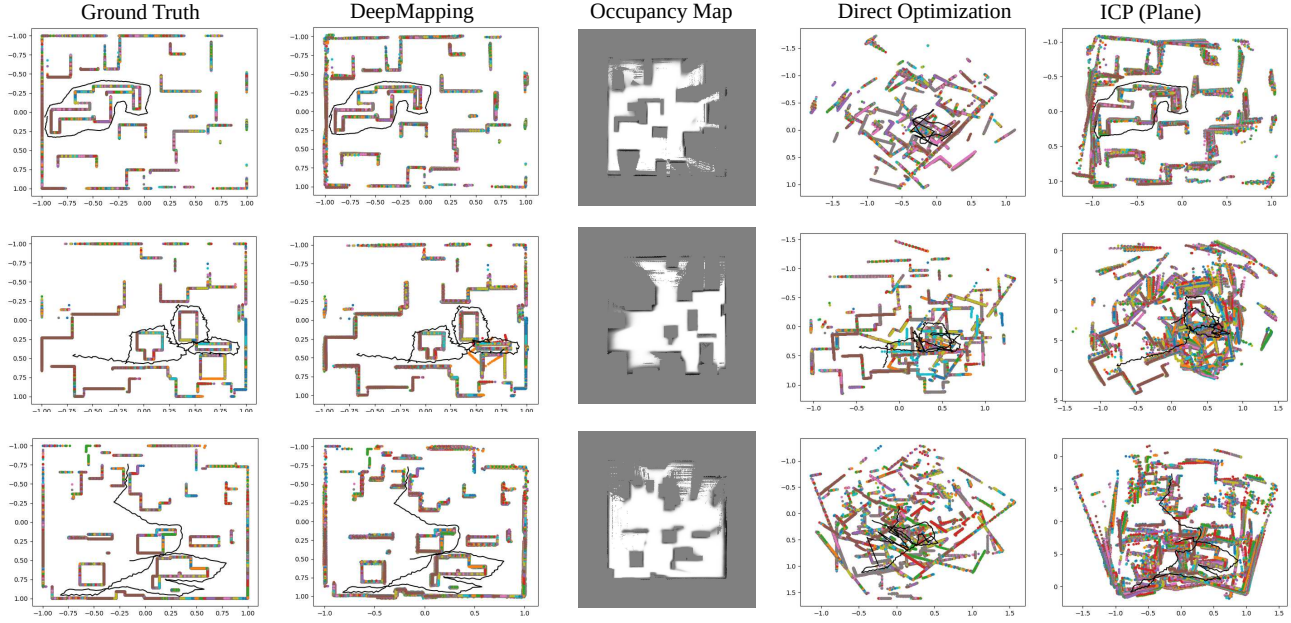


Figure 5. Qualitative results from the 2D simulated dataset. The black lines are the sensor trajectories. The third column shows the occupancy maps for each trajectory that are estimated by the M-Net. The black, white, and gray pixels show the occupied, unoccupied, and unexplored locations, respectively. While the results of each trajectory are defined in arbitrary coordinate systems, we aligned them with the ground truth coordinates for a clear comparison. More visualization results are included in supplementary material. Best viewed in color.

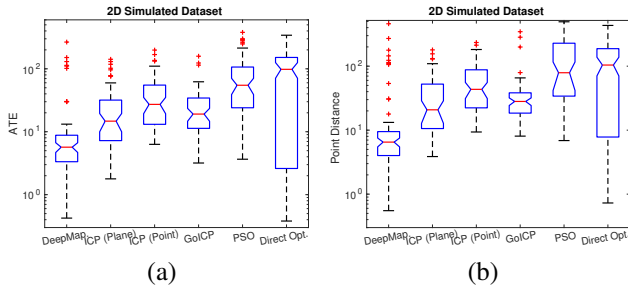


Figure 6. Box plots of the ATE and the point distance on the 2D simulated dataset. In each box, the red line indicates the median. The top and bottom blue edges of the box show the first (25th percentile) and the third (75th percentile) quartiles, respectively. Note the logarithmic scale for the y-axis.

randomly move or rotate the camera sensor in the space and collect 105 trajectories from the dataset. Each trajectory contains either 8 or 16 point clouds.

Baseline: we compare DeepMapping with the baseline method of multiway registration [11]. The ICP algorithms perform poorly because of the low overlapping rate between consecutive point clouds and thus are not included in this section. The multiway registration aligns multiple point clouds via pose graph optimization where each node represents an input point cloud and the graph edges connect two point clouds. We follow the same procedure in Section 4.1 to test the performance of direct optimization with respect

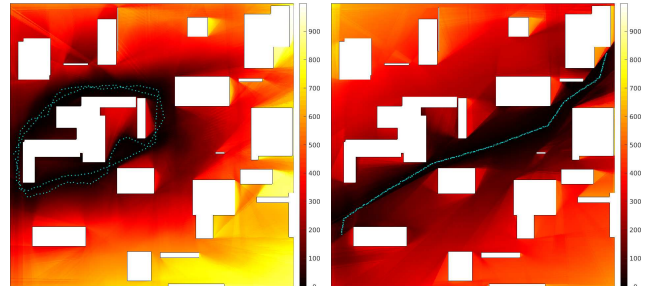


Figure 7. Re-localization errors of two L-Nets (black is better). Cyan points show the trajectories used in the “training” stage. Best viewed in color.

to sensor pose on the AVD. For DeepMapping and the direct optimization, we compare two loss functions: only the BCE loss \mathcal{L}_{bce} (by setting λ to 0), and the combination of the BCE loss \mathcal{L}_{bce} and the Chamfer distance \mathcal{L}_{ch} . For the latter, we set λ to 0.1 because of the small overlapping between two scans.

Implementation: the L-Net architecture consists of C(64)-C(128)-C(256)-C(1024)-AM(1)-FC(512)-FC(256)-FC(3), where AM(1) denotes 2D adaptive max-pooling layer. The M-Net has the same structure as that used in Section 4.1, except for the input layer that has 3 nodes rather than 2. We sample 35 points on each laser ray. The multiway registration [11] is implemented using Open3D library [51]. The batch size is 8, and the optimiza-

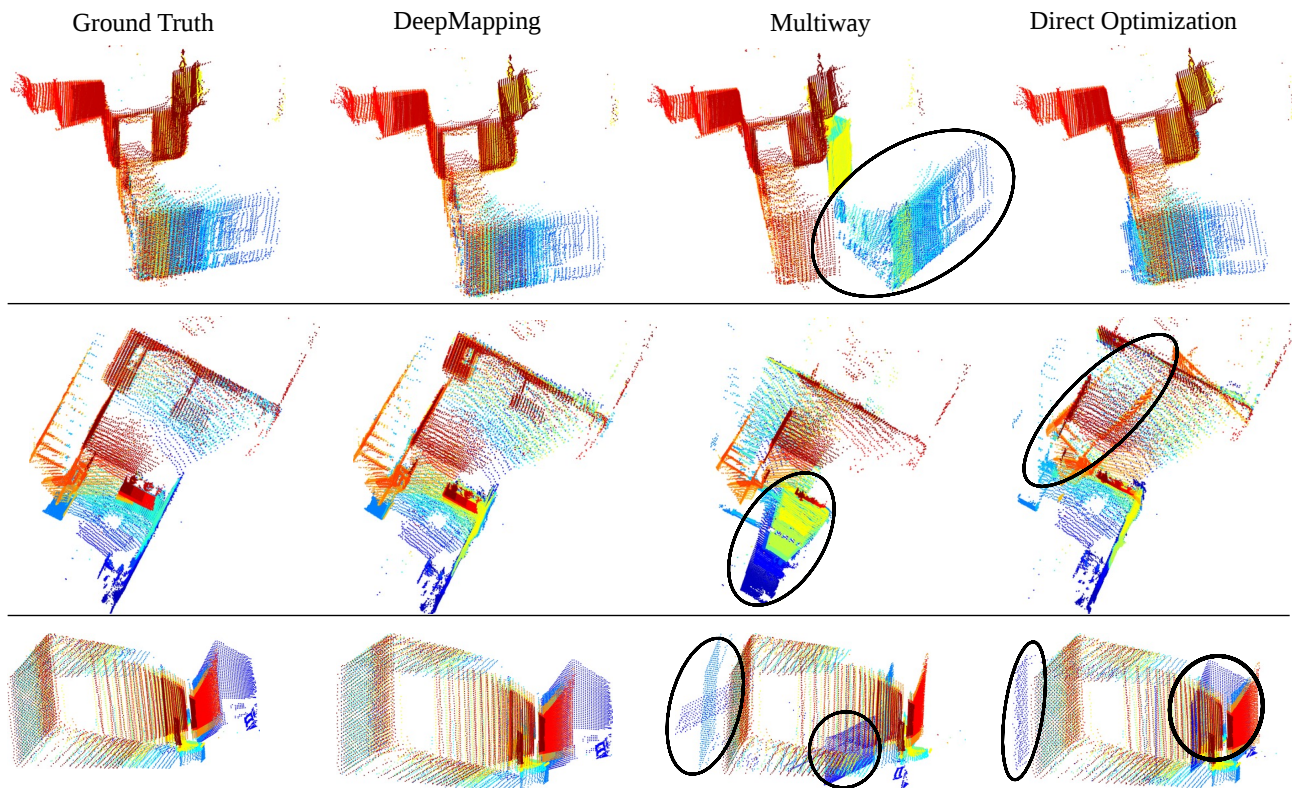


Figure 8. Qualitative results of multiple point clouds registration tested on the AVD [4] dataset. The black ellipses highlight the misaligned parts in baselines. Each color represents one point cloud. More visualizations are shown in supplementary material. Best viewed in color.

tion speed is approximately 125 epochs per minute on an Nvidia TITAN XP.

Results: The visual comparison of multiple point clouds registrations for different algorithms is shown in Figure 8. The misaligned point clouds are highlighted with black ellipses. The stairs from the direct optimization, for example, are not correctly registered, as shown in the second row in Figure 8. The multiway registration [11] and the direct optimization fails to align several planar structures that are shown in the last row in Figure 8. Failure cases of DeepMapping are shown in the supplementary material.

We show the box plots of the ATE and the point distance in Figure 9. The multiway registration [11] performs poorly that has large errors in both metrics. Comparing the direct optimization with DeepMapping, we show the advantage of using neural networks for optimization. In addition, integrating Chamfer distance into DeepMapping is also helpful to improve the registration accuracy.

5. Conclusion

In this paper, we propose DeepMapping to explore a possible direction for integrating deep learning methods into multiple point clouds registration that could also be informative for other related geometric vision problems. The

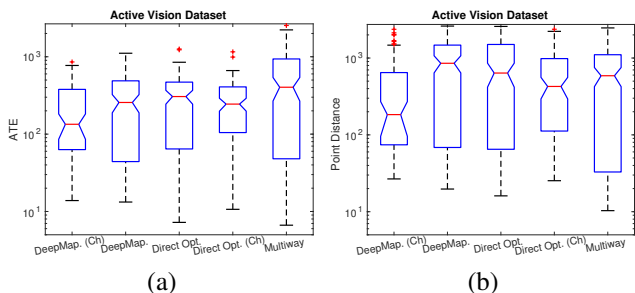


Figure 9. Box plots of the ATE and the point distance on the AVD [4]. The methods with “Ch” combine the Chamfer and the BCE losses. Legends are the same as those in Figure 5. Note the logarithmic scale for the y-axis.

novelty of our approach lies in the formulation that converts solving the registration problem into “training” some DNNs using properly defined unsupervised loss functions, with promising experimental performances.

Acknowledgment

The authors gratefully acknowledge the helpful comments and suggestions from Yuichi Taguchi, Dong Tian, Weiyang Liu, and Alan Sullivan.

References

- [1] Particle swarm optimization. https://github.com/iralabdisco/pso_registration. 6
- [2] PyTorch. <https://pytorch.org/>. 5
- [3] Dror Aiger, Niloy J Mitra, and Daniel Cohen-Or. 4-points congruent sets for robust pairwise surface registration. In *ACM Trans. Graphics*, volume 27, page 85, 2008. 2
- [4] Phil Ammirato, Patrick Poirson, Eunbyung Park, Jana Kosecka, and Alexander C. Berg. A dataset for developing and benchmarking active vision. In *Proc. the IEEE Intl. Conf. on Robotics and Auto.*, 2017. 5, 6, 8
- [5] Andrea Banino, Caswell Barry, Benigno Uribe, Charles Blundell, Timothy Lillicrap, Piotr Mirowski, Alexander Pritzel, Martin J Chadwick, Thomas Degris, Joseph Modayil, et al. Vector-based navigation using grid-like representations in artificial agents. *Nature*, 557(7705):429, 2018. 3
- [6] P. J. Besl and N. D. McKay. A method for registration of 3-D shapes. *IEEE Trans. Pattern Anal. Mach. Intel.*, 14(2):239–256, 1992. 2, 6
- [7] Michael Bloesch, Jan Czarowski, Ronald Clark, Stefan Leutenegger, and Andrew J. Davison. CodeSLAM - learning a compact, optimisable representation for dense visual SLAM. In *IEEE Intl. Conf. Comp. Vision and Pattern Recog.*, June 2018. 2
- [8] Eric Brachmann, Alexander Krull, Sebastian Nowozin, Jamie Shotton, Frank Michel, Stefan Gumhold, and Carsten Rother. DSAC - differentiable ransac for camera localization. In *IEEE Intl. Conf. Comp. Vision and Pattern Recog.*, July 2017. 2
- [9] Samarth Brahmabhatt, Jinwei Gu, Kihwan Kim, James Hays, and Jan Kautz. Geometry-aware learning of maps for camera localization. In *IEEE Intl. Conf. Comp. Vision and Pattern Recog.*, June 2018. 1, 2
- [10] Yang Chen and Gérard Medioni. Object modelling by registration of multiple range images. *Image and Vision Computing*, 10(3):145–155, 1992. 2, 6
- [11] Sungjoon Choi, Q. Zhou, and V. Koltun. Robust reconstruction of indoor scenes. In *IEEE Intl. Conf. Comp. Vision and Pattern Recog.*, June 2015. 2, 7, 8
- [12] Martin Danelljan, Giulia Meneghetti, Fahad Shahbaz Khan, and Michael Felsberg. A probabilistic framework for color-based point set registration. In *IEEE Intl. Conf. Comp. Vision and Pattern Recog.*, pages 1818–1826, 2016. 2
- [13] Haowen Deng, Tolga Birdal, and Slobodan Ilic. PPFNet: Global context aware local features for robust 3D point matching. In *IEEE Intl. Conf. Comp. Vision and Pattern Recog.*, June 2018. 2
- [14] Gil Elbaz, Tamar Avraham, and Anath Fischer. 3D point cloud registration for localization using a deep neural network auto-encoder. In *IEEE Intl. Conf. Comp. Vision and Pattern Recog.*, pages 2472–2481. IEEE, 2017. 2
- [15] SM Ali Eslami, Danilo Jimenez Rezende, Frederic Besse, Fabio Viola, Ari S Morcos, Marta Garnelo, Avraham Ruderman, Andrei A Rusu, Ivo Danihelka, Karol Gregor, et al. Neural scene representation and rendering. *Science*, 360(6394):1204–1210, 2018. 3
- [16] Georgios D Evangelidis, Dionyssos Kounades-Bastian, Radu Horaud, and Emmanouil Z Psarakis. A generative model for the joint registration of multiple point sets. In *Euro. Conf. on Comp. Vision*, pages 109–122. Springer, 2014. 2, 3
- [17] Martin A Fischler and Robert C Bolles. Random sample consensus: a paradigm for model fitting with applications to image analysis and automated cartography. *Commun. ACM*, 24(6):381–395, 1981. 2
- [18] Marco Fraccaro, Danilo Jimenez Rezende, Yori Zwols, Alexander Pritzel, SM Eslami, and Fabio Viola. Generative temporal models with spatial memory for partially observed environments. In *Intl. Conf. on Mach. Learning*, 2018. 1, 2
- [19] Vitor Guizilini and Fabio Ramos. Learning to reconstruct 3D structures for occupancy mapping from depth and color information. *Intl. J. of Robotics Research*, 2018. 1, 3
- [20] Joao F Henriques and Andrea Vedaldi. MapNet: An allocentric spatial memory for mapping environments. *IEEE Intl. Conf. Comp. Vision and Pattern Recog.*, 2018. 1, 2
- [21] Berthold KP Horn. Closed-form solution of absolute orientation using unit quaternions. *J. Opt. Soc. Am. A*, 4(4):629–642, 1987. 5
- [22] Shahram Izadi, David Kim, Otmar Hilliges, David Molyneaux, Richard Newcombe, Pushmeet Kohli, Jamie Shotton, Steve Hodges, Dustin Freeman, Andrew Davison, et al. Kinectfusion: real-time 3D reconstruction and interaction using a moving depth camera. In *ACM Symp. User Interface Software and Technology*, pages 559–568, 2011. 2
- [23] Bing Jian and Baba C Vemuri. A robust algorithm for point set registration using mixture of gaussians. In *IEEE Intl. Conf. Comp. Vision*, volume 2, pages 1246–1251, 2005. 2
- [24] Zi Jian Yew and Gim Hee Lee. 3DFeat-Net: Weakly supervised local 3D features for point cloud registration. In *Euro. Conf. on Comp. Vision*, September 2018. 2
- [25] Andrew E Johnson and Martial Hebert. Using spin images for efficient object recognition in cluttered 3D scenes. *IEEE Trans. Pattern Anal. Mach. Intel.*, (5):433–449, 1999. 2
- [26] Alex Kendall, Matthew Grimes, and Roberto Cipolla. PoseNet: A convolutional network for real-time 6-DOF camera relocalization. In *IEEE Intl. Conf. Comp. Vision*, December 2015. 1, 2
- [27] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *Intl. Conf. Learning Representations*, 2015. 5
- [28] Huan Lei, Guang Jiang, and Long Quan. Fast descriptors and correspondence propagation for robust global point cloud registration. *IEEE Trans. Image Proc.*, 26(8):3614–3623, 2017. 2
- [29] J. Li, H. Zhan, B. M. Chen, I. Reid, and G. H. Lee. Deep learning for 2D scan matching and loop closure. In *IEEE Intl. Conf. Intel. Robots and Sys.*, pages 763–768, Sept 2017. 1, 2
- [30] Nicolas Mellado, Dror Aiger, and Niloy J Mitra. Super 4PCS fast global pointcloud registration via smart indexing. In *Comp. Graphics Forum*, volume 33, pages 205–215, 2014. 2
- [31] A Myronenko and Xubo Song. Point set registration: Coherent point drift. *IEEE Trans. Pattern Anal. Mach. Intel.*, 32(12):2262–2275, 2010. 2

- [32] Emilio Parisotto, Devendra Singh Chaplot, Jian Zhang, and Ruslan Salakhutdinov. Global pose estimation with an attention-based recurrent network. In *IEEE Intl. Conf. Comp. Vision and Pattern Recog. Wksp.*, pages 237–246, 2018. 2
- [33] Charles R. Qi, Hao Su, Kaichun Mo, and Leonidas J. Guibas. PointNet: Deep learning on point sets for 3D classification and segmentation. In *IEEE Intl. Conf. Comp. Vision and Pattern Recog.*, July 2017. 4
- [34] Szymon Rusinkiewicz and Marc Levoy. Efficient variants of the ICP algorithm. In *3D Digital Imaging and Modeling*, pages 145–152, 2001. 2
- [35] Radu Bogdan Rusu, Nico Blodow, and Michael Beetz. Fast point feature histograms (FPFH) for 3D registration. In *Proc. the IEEE Intl. Conf. on Robotics and Auto.*, pages 3212–3217, 2009. 2
- [36] Radu Bogdan Rusu, Nico Blodow, Zoltan Csaba Marton, and Michael Beetz. Aligning point cloud views using persistent feature histograms. In *IEEE Intl. Conf. Intel. Robots and Sys.*, pages 3384–3391, 2008. 2
- [37] Johannes L Schönberger, Marc Pollefeys, Andreas Geiger, and Torsten Sattler. Semantic visual localization. *ISPRS J. Photographic and Remote Sensing*, 2018. 2
- [38] Paul Scovanner, Saad Ali, and Mubarak Shah. A 3-dimensional SIFT descriptor and its application to action recognition. In *ACM Intl. Conf. Multimedia*, pages 357–360, 2007. 2
- [39] Shuran Song, Samuel P Lichtenberg, and Jianxiong Xiao. Sun RGB-D: A RGB-D scene understanding benchmark suite. In *IEEE Intl. Conf. Comp. Vision and Pattern Recog.*, pages 567–576, 2015. 5, 6
- [40] Bastian Steder, Radu Bogdan Rusu, Kurt Konolige, and Wolfram Burgard. NARF: 3D range image features for object recognition. In *Wksp. on Defining and Solving Realistic Perception Problems in Personal Robotics at IEEE Intl. Conf. Intel. Robots and Sys.*, volume 44, 2010. 2
- [41] J. Sturm, N. Engelhard, F. Endres, W. Burgard, and D. Cremers. A benchmark for the evaluation of RGB-D SLAM systems. In *IEEE Intl. Conf. Intel. Robots and Sys.*, Oct. 2012. 5, 6
- [42] Pascal Willy Theiler, Jan Dirk Wegner, and Konrad Schindler. Globally consistent registration of terrestrial laser scans via graph optimization. *ISPRS J. Photographic and Remote Sensing*, 109:126–138, 2015. 2
- [43] Federico Tombari, Samuele Salti, and Luigi Di Stefano. Unique signatures of histograms for local surface description. In *Euro. Conf. on Comp. Vision*, pages 356–369, 2010. 2
- [44] Andrea Torsello, Emanuele Rodola, and Andrea Albarelli. Multiview registration via graph diffusion of dual quaternions. In *IEEE Intl. Conf. Comp. Vision and Pattern Recog.*, pages 2441–2448, 2011. 2
- [45] Benjamin Ummenhofer, Huizhong Zhou, Jonas Uhrig, Nikolaus Mayer, Eddy Ilg, Alexey Dosovitskiy, and Thomas Brox. DeMoN: Depth and motion network for learning monocular stereo. In *IEEE Intl. Conf. Comp. Vision and Pattern Recog.*, volume 5, page 6, 2017. 1, 2
- [46] J. Yang, H. Li, D. Campbell, and Y. Jia. Go-ICP: A globally optimal solution to 3D ICP point-set registration. *IEEE Trans. Pattern Anal. Mach. Intel.*, 38(11):2241–2254, Nov 2016. 2, 6
- [47] Nan Yang, Rui Wang, Jörg Stückler, and Daniel Cremers. Deep virtual stereo odometry: Leveraging deep depth prediction for monocular direct sparse odometry. In *Euro. Conf. on Comp. Vision*, pages 835–852, 2018. 2
- [48] Andy Zeng, Shuran Song, Matthias Nießner, Matthew Fisher, Jianxiong Xiao, and Thomas Funkhouser. 3Dmatch: Learning local geometric descriptors from RGB-D reconstructions. In *IEEE Intl. Conf. Comp. Vision and Pattern Recog.*, pages 199–208, 2017. 2
- [49] Huizhong Zhou, Benjamin Ummenhofer, and Thomas Brox. DeepTAM: Deep tracking and mapping. In *Euro. Conf. on Comp. Vision*, September 2018. 1, 2
- [50] Qian-Yi Zhou, Jaesik Park, and Vladlen Koltun. Fast global registration. In *Euro. Conf. on Comp. Vision*, pages 766–782, 2016. 2, 3
- [51] Qian-Yi Zhou, Jaesik Park, and Vladlen Koltun. Open3D: A modern library for 3D data processing. *arXiv:1801.09847*, 2018. 7
- [52] Tinghui Zhou, Matthew Brown, Noah Snavely, and David G Lowe. Unsupervised learning of depth and ego-motion from video. In *IEEE Intl. Conf. Comp. Vision and Pattern Recog.*, volume 2, page 7, 2017. 1, 2