

# Unsupervised Simultaneous Learning for Camera Re-Localization and Depth Estimation from Video

Shun Taguchi<sup>1</sup> and Noriaki Hirose<sup>1</sup>

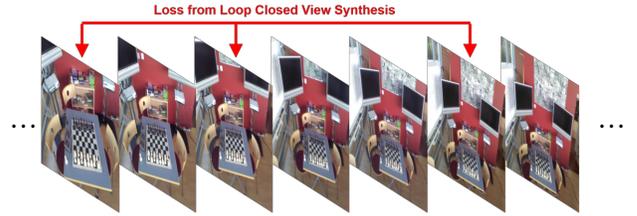
**Abstract**— We present an unsupervised simultaneous learning framework for the task of monocular camera re-localization and depth estimation from unlabeled video sequences. Monocular camera re-localization refers to the task of estimating the absolute camera pose from an instance image in a known environment, which has been intensively studied for alternative localization in GPS-denied environments. In recent works, camera re-localization methods are trained via supervised learning from pairs of camera images and camera poses. In contrast to previous works, we propose a completely unsupervised learning framework for camera re-localization and depth estimation, requiring only monocular video sequences for training. In our framework, we train two networks that estimate the scene coordinates using directions and the depth map from each image which are then combined to estimate the camera pose. The networks can be trained through the minimization of loss functions based on our loop closed view synthesis. In experiments with the 7-scenes dataset, the proposed method outperformed the re-localization of the state-of-the-art visual SLAM, ORB-SLAM3. Our method also outperforms state-of-the-art monocular depth estimation in a trained environment.

## I. INTRODUCTION

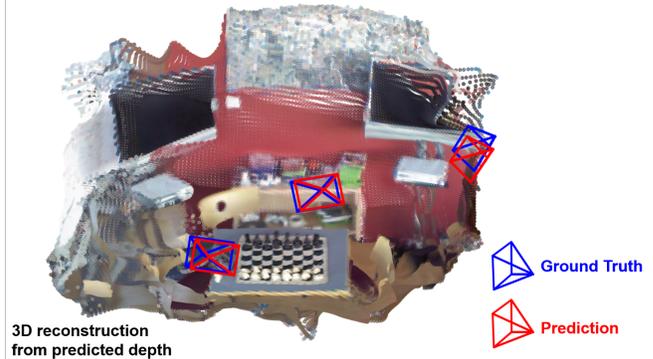
Relying on a single memory image, humans can easily recognize a previously visited environment. Humans can also comprehend the implicit map of any environment relying on only the memory of their first visit. Immediately they see a new image, they can recognize where it belongs on the map. How can humans construct an implicit map in their brains relying only on unlabeled image sequences? We consider that humans recognize relative pose and depth from the difference of appearance in the images. They also perform re-localization on the images to recognize a place as the same place when they see it for a second time. Inspired by this human ability, we propose a novel unsupervised learning framework to train neural networks to estimate camera pose from a single image.

To this end, we trained camera re-localization and depth estimation simultaneously from unlabeled monocular video sequences by utilizing view synthesis inspired by monocular depth estimation [1], [2]. By introducing re-localization networks, it is possible to construct implicit maps while determining that places with the same appearance are the same places. Our training led the networks to implicitly understand the environment captured in the training dataset,

<sup>\*</sup>This work was not supported by any organization  
<sup>1</sup>Shun Taguchi and Noriaki Hirose are with Toyota Central R&D Labs., Inc., 41-1 Yokomichi, Nagakute, Aichi, Japan s-taguchi@mosk.tytlabs.co.jp hirose@mosk.tytlabs.co.jp



(a) Training from unlabeled monocular video sequences.



(b) Test: camera re-localization (camera pose estimation from a single monocular camera image) with depth estimation.

**Fig. 1. Unsupervised simultaneous learning for camera re-localization and depth estimation.** We present an unsupervised simultaneous learning framework for camera re-localization and depth estimation from unlabeled monocular video sequences. Our models predict the camera pose and depth in a known environment at test time.

thereby, achieving accurate camera re-localization and depth estimation in inference (see Fig. 1).

Our method can train networks for localization and structure understanding of the environment via a single learning process. This information is essential for visual-based autonomous navigation of robots. Our method is smart and useful because it can simultaneously learn the networks required for such autonomous navigation.

This study makes three main contributions:

- We develop a novel unsupervised simultaneous learning framework for camera re-localization and depth estimation from unlabeled video sequences.
- We propose a novel directed scene coordinate network that enables accurate estimation of the camera pose from a single image by combining it with an estimated depth map.
- We propose a calculation process of the loss function with loop closed view synthesis, which improves the training of camera re-localization.

To the best of our knowledge, our study is among the first studies to propose unsupervised camera re-localization by training models from unlabeled video sequences and estimating camera pose from a single image in inference.

The proposed method outperformed the re-localization systems embedded in the state-of-the-art visual SLAM method, ORB-SLAM3 [3]. Our method also outperformed one of the unsupervised monocular depth estimation methods, monodepth2 [2] in depth estimation. Our method can suppress the deviation of scale in the test dataset on depth estimation because more consistent scales are learned via camera pose estimation. These evaluations are quantitatively and qualitatively conducted on 7-scenes dataset [4].

## II. RELATED WORK

In the following, we discuss camera re-localization and monocular depth estimation, which are most closely related to our study.

### A. Camera Re-Localization

In early studies of camera re-localization based on image retrieval [5], images were retrieved from a database by matching query images with global image descriptors [6], [7]. Because the re-localization accuracy of these retrieval approaches is limited by the sampling density of database images, various camera re-localization methods have been proposed to address this problem.

In addition, absolute pose regression methods [8]–[14] have been used to train neural networks to regress the camera pose using database images as the training set. However, in practice, absolute pose regression methods do not outperform the accuracy of image retrieval methods [15]. Relative pose regression methods [15], [16] train a neural network to predict the relative transformation between the query image and the database image most similar to that which is obtained by image retrieval. A recent study [17] suggested that relative pose regression can achieve accuracy comparable to that of the structure-based methods described below.

Structure-based approaches, such as scene coordinate regression [4], are some of the most successful approaches. Scene coordinate regression [4] directly predicts the 3D scene points corresponding to a given 2D pixel location, and the camera pose is calculated by solving the PnP problem. Originally, scene coordinate regression was proposed for RGB-D-based re-localization in indoor environments [4], [18]–[20]. Recently, scene coordinate regression has been shown to be effective for RGB-based re-localization [21], [22]. DSAC++ [23] presents the possibility of learning scene coordinate regression from RGB images and ground-truth poses only. In subsequent work, DSAC\* [24] makes several improvements to DSAC++ to achieve state-of-the-art performance.

Recently, some studies have extended the camera re-localization method to the time domain in order to address temporal re-localization [14], [25]–[28]. While these approaches are effective for some applications, we will focus

on one-shot camera re-localization because it is a more fundamental and widely used technique.

The reviewed literature referred to above has focused on supervised learning of camera re-localization. In contrast, we present an unsupervised framework for camera re-localization using unlabeled video sequences. We also estimate the scene coordinates based on the findings of these studies of camera re-localization. Unlike [23], [24], when solving the PnP problem, our method estimates the scene coordinate, which is a 6-dimensional coordinate with the gaze direction to handle the difference in appearance depending on the viewing direction. Our approach to estimating the camera pose is introduced in our unsupervised learning framework based on loop closed view synthesis, which is inspired by unsupervised monocular depth estimation.

### B. Monocular Depth Estimation

One of the earliest works in convolutional-based depth estimation was presented by Eigen et al [29]. Eigen et al used a multi-scale deep network trained on RGB-D sensor data to regress the depth directly from single images. Inspired by the two-view stereo disparity estimation [30] based on flow estimation [31], Umenhofer et al. [32] trained a depth and pose network simultaneously to predict depth and camera ego-motion between successive unconstrained image pairs. To address the difficulty of labeling the target depth, [33], [34] trained a monocular depth network with stereo cameras without requiring ground-truth depth labels. Godard et al. [34] used stereo images to geometrically transform the right image into a left image based on a predicted depth by leveraging spatial transformer networks [35]. The photometric re-projection loss between the synthesized and original left images can be used to train the depth network without the ground-truth depth.

Following [34] and [32], Zhou et al. [1] applied this unsupervised training to a purely monocular setting, where a depth and relative pose network are simultaneously learned from unlabeled monocular videos. Recent studies [36]–[43] have incorporated these methods, additional loss, and constraints. Monodepth2 [2] is a successful method that employs a ResNet encoder [44] and has achieved state-of-the-art performance. More recent methods employ improved network models with more parameters to achieve state-of-the-art performance [45]. Several methods [46]–[48] have focused on pose estimation based on an unsupervised monocular depth estimation framework, however, these approaches have focused only on relative pose estimation between two images, that is, visual odometry.

Instead of the visual odometry network of unsupervised monocular depth estimation, our method estimates the absolute camera pose in an implicit map via estimated scene coordinates, which is trained with our loop closed view synthesis. Note that the estimated camera pose in our method is not a relative pose between two images. To the best of our knowledge, our method is the first of its kind to train camera re-localization in an unsupervised setting.

### III. METHOD

#### A. Camera Re-Localization based on Directed Scene Coordinate and Depth

In this study, we propose a novel unsupervised simultaneous learning framework for camera re-localization and depth estimation, shown in Fig. 2. Our model consists of two networks for depth and directed scene coordinate. The depth network predicts the dense depth map  $D_t$  from a single RGB image  $I_t$  at time step  $t$ . The directed scene coordinate network predicts the direct scene coordinate  $S_t$ , which consists of subsampled scene coordinates with direction from a single RGB image  $I_t$ . The directed scene coordinate  $S_t = [\theta_t, \tau_t]^T$  consists of a 3-dimensional attitude  $\theta_t$  as an axis-angle representation and 3-dimensional position  $\tau_t$ . Here,  $\theta_t$  indicates the direction of the vector from the camera position to the position of corresponding scene coordinates.

The camera pose  $P_t$  can be calculated geometrically from a directed scene coordinate  $S_t$  and the depth  $D_t$  for each pixel. To calculate camera pose, we introduce the point cloud  $Q_t$ , which can be obtained by back projection from the depth  $D_t$  as follows:

$$Q_t(p) = D_t(p)K_t^{-1} \begin{bmatrix} p \\ 1 \end{bmatrix}, \quad (1)$$

where  $Q_t(p)$  is a vector to local point corresponding to pixel  $p = [u, v]^T$ , and  $K_t$  is the camera intrinsic parameter for  $I_t$ .

The camera pose  $P_t(p)$  corresponding to a pixel  $p$  also consists a 3-dimensional attitude  $\theta_{P_t}(p)$  as an axis-angle representation and a 3-dimensional position  $\tau_{P_t}(p)$  as

$$P_t(p) = \begin{bmatrix} \theta_{P_t}(p) \\ \tau_{P_t}(p) \end{bmatrix}, \quad (2)$$

and it is calculated from the directed scene coordinate  $S_t(p)$  and the point  $Q_t(p)$  (see. Fig. 3). The camera attitude  $\theta_{P_t}(p)$  is calculated by rotating the gaze direction  $\theta_t(p)$  using the direction of the pixel  $p$  corresponding to the directed scene coordinates.

$$R(\theta_{P_t}(p)) = R_p^{-1}R(\theta_t(p)) \quad (3)$$

where  $R_p$  is the rotation matrix of gaze direction to the pixel  $p$ ,  $R(\cdot)$  represents rotation matrix corresponding the 3-dimensional attitude. The camera position  $\tau_{P_t}(p)$  is calculated by translating the scene coordinates  $\tau_t(p)$  by the distance of the local point  $Q_t(p)$  in the gaze direction  $\theta_t(p)$ .

$$\tau_{P_t}(p) = \|Q_t(p)\| \cdot R(\theta_t(p))e_z + \tau_t(p), \quad (4)$$

where  $e_z$  is a unit vector to depth direction, and  $\|Q_t(p)\|$  is the distance of the point  $Q_t(p)$

In this procedure, the camera poses are calculated from all pixels of the directed scene coordinate; therefore, the final pose  $P_t$  is obtained as the average (median in testing) of the estimated camera poses from all pixels.

$$P_t = \frac{1}{n_S} \sum_p P_t(p), \quad (5)$$

where,  $n_S$  is the number of pixels of directed scene coordinate.

#### B. Unsupervised Learning based on Photometric Re-projection Loss

Here, we propose a framework for jointly training two networks that estimate directed scene coordinates and depths from unlabeled video sequences. Each model can be used independently during test-time inference, and camera re-localization can be performed by combining the two networks, as shown in Sec. III-A.

Our models were trained from image sequences obtained from moving cameras in the target environment. We assume that the corresponding environments are mostly static, that is, the scene appearance change across different frames is dominated by the camera motion. Similar to most unsupervised monocular depth estimations [1], [2], we also formulate our problem as the minimization of a photometric re-projection error at training time. Geometrically, a pixel mapping  $M_{s \rightarrow t}$  between  $I_t$  and  $I_s$  can be derived from the depth  $D_t$  and  $T_{t \rightarrow s}$  as follows [1], [2]:

$$M_{s \rightarrow t} = K_s T_{t \rightarrow s} D_t K_t^{-1}, \quad (6)$$

where  $K_t$  and  $K_s$  are the camera intrinsic matrices for  $I_t$  and  $I_s$ , respectively. The transformation matrix  $T_{t \rightarrow s}$  between two images can be calculated from estimated camera poses  $P_t$  and  $P_s$  from each image independently as

$$T_{t \rightarrow s} = T(P_s)^{-1}T(P_t), \quad (7)$$

where  $T(\cdot)$  denotes a transformation matrix corresponding to a 6-DoF pose.

By using the relative pose between two images, the base of our unsupervised learning framework for camera re-localization can be formulated. The synthesized image  $\hat{I}_t$  can be constructed by sampling from  $I_s$  by following  $M_{s \rightarrow t}$ .

Similar to previous works [2], we employ L1 and SSIM loss as a photometric re-projection loss, which is formulated as

$$L_p = \frac{1}{|V|} \sum_{p \in V} \alpha \frac{(1 - SSIM(I_t(p), \hat{I}_t(p)))}{2} + \frac{1}{V} \sum_{p \in V} (1 - \alpha) \|I_t(p) - \hat{I}_t(p)\|_1, \quad (8)$$

where  $V$  is a set of valid points  $p$  that are successfully transformed between  $I_t$  and  $I_s$ , and  $|V|$  denotes the number of points in  $V$ . The parameter  $\alpha = 0.85$  to follow previous works [2].

The photometric re-projection loss is effective for texture-rich scenes, however fragile for low-texture or homogeneous regions. Therefore, another smoothness loss  $L_s$  is used together with the photometric re-projection losses to solve this problem [2],

$$L_s = \frac{1}{n_D} \sum_p \left( e^{-\nabla I_t(p)} \cdot \nabla D_t(p) \right)^2, \quad (9)$$

where  $n_D$  is the number of pixels of  $D_t$ , and  $\nabla$  denotes the first derivative. This ensures that the smoothness of the depth map is constrained by the primary gradient of

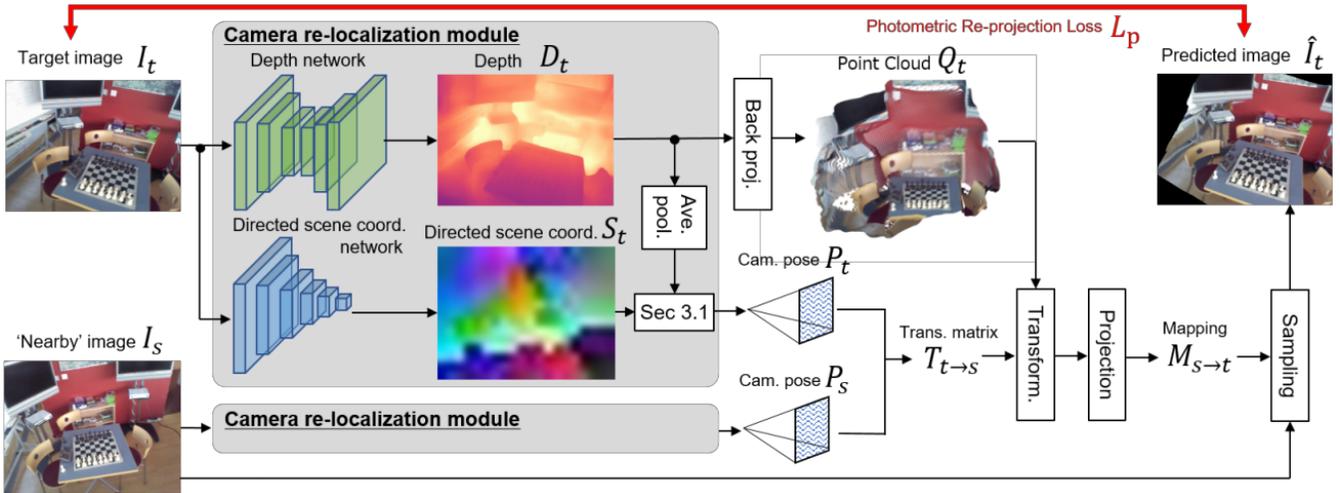


Fig. 2. **Overview of our unsupervised simultaneous learning framework of directed scene coordinate and depth networks.** Our model consists of two networks for depth and directed scene coordinate. The camera pose can be calculated geometrically from a directed scene coordinate and depth as shown in Sec.III-A, and the relative pose between an image  $I_t$  and a ‘nearby’ image  $I_s$  can be calculated from the predicted camera pose from each image independently. By using the relative pose and back-projected point cloud from depth, the corresponding map  $M_{s \rightarrow t}$  can be calculated. Our networks are trained by the photometric re-projection loss  $L_p$  between the original  $I_t$  and the view synthesis image  $\hat{I}_t$  from the ‘nearby’ image  $I_s$ .

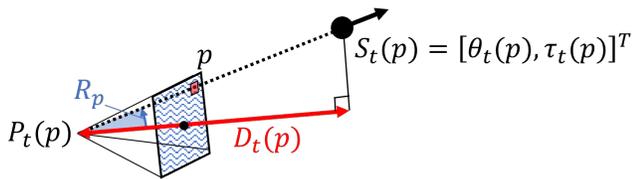


Fig. 3. **Camera pose calculation from directed scene coordinate and depth.** The camera pose  $P_t(p)$  can be calculated from a pixel of the directed scene coordinate  $S_t(p)$  and the corresponding depth  $D_t(p)$  by rotating the pose based on the direction of the pixel and translating the position based on the depth and direction of the scene coordinates. The final pose is obtained as the average of estimated camera poses from each pixel.

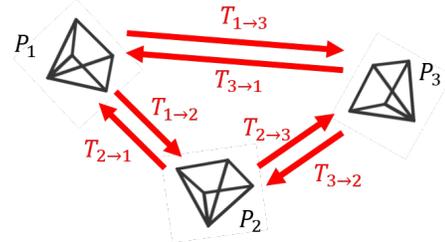


Fig. 4. **Loop closed view synthesis.** The monocular depth estimation [2] learn from view synthesis based on  $T_{1 \rightarrow 2}$  and  $T_{3 \rightarrow 2}$ , however, the two relative poses are insufficient to estimate the three absolute camera poses on an implicit map. To address this issue, we train our networks from view synthesis for all combinations of images.

the corresponding color image. Following [2],  $L_p$  and  $L_s$  are calculated for each estimated multi-scale depth at the resolution of each decoder’s layer.

In this study, we adopt an additional loss to these well-known losses in unsupervised depth estimation. The loss is the pose coordinate loss  $L_c$ , which is the error between the camera pose estimated from each pixel and the average final pose in (5).

$$L_c = \frac{1}{n_S} \sum_p \|P_t - P_t(p)\|_2. \quad (10)$$

Therefore, the final training loss  $L$  consists of photometric re-projection loss, smoothness loss, and pose coordinate loss, as follows:

$$L = L_p + w_s L_s + w_c L_c, \quad (11)$$

where  $w_s$  and  $w_c$  are the respective weight parameters of smoothness loss and pose coordinate loss.

### C. Loop Closed View Synthesis

In this study, to achieve effective training for camera pose, we introduced loop closed view synthesis. First, we picked three or more images from sequences. For each combination,

we bi-directionally predict the images by view synthesis, following Sec.III-B, and calculate the photometric re-projection loss. According to our loop closed view synthesis, we calculate six photometric re-projection losses and average them as  $L_p$  for three image cases (see Fig.4).

This is because the photometric re-projection loss from two images is insufficient to train networks for camera re-localization. Eq. (7) provides one constraint for two camera poses, which is insufficient to estimate absolute camera poses on an implicit map. Massive iteration in the learning process may provide sufficient constraints; however, the training process becomes unstable and does not converge smoothly. Our loop closed view synthesis attempts to solve this issue. Sufficient constraints to estimate the absolute camera pose stabilize the training process of unsupervised camera re-localization and leads to adequate results.

The final training loss  $L$  is averaged over the scale, batch, and view synthesis.

#### D. Network Architecture

Our depth network is the same as in [2]. The network is based on the general U-Net architecture [49], that is, an encoder-decoder network, with skip connections, enabling us to represent both deep abstract features as well as local information. We also used ResNet18 [44] as the depth encoder and started with weights pre-trained ImageNet [50]. The depth decoder has sigmoids at the output and we convert the sigmoid output  $\sigma$  to depth with  $D = 1/(a\sigma + b)$ , where  $a$  and  $b$  are chosen to constrain  $D$  between 0.1 and 100 units, following [2].

Our directed scene coordinate network also consists of an encoder-decoder network. We also use a pre-trained ResNet18 as the directed scene coordinate encoder. The decoder consists of three CNN layers with Rectified Linear Unit (ReLU) activations. The output of the decoder is the directed scene coordinates  $S_t$  with 6 channels of  $20 \times 15$  matrix. The output directed scene coordinate is subsampled by factor 32 from the original images ( $640 \times 480$ ). Therefore, we resize the estimated depth by  $32 \times 32$  average pooling layer to combine with directed scene coordinates for calculating the camera pose.

### IV. EXPERIMENTS

#### A. Dataset

In this study, an indoor camera dataset, 7-scenes dataset [4], was used for the evaluation. The 7-scenes dataset is an RGB-D indoor re-localization dataset of seven small indoor environments with challenging conditions such as motion blur, reflective surfaces, repetitive structures, and untextured areas. The dataset was collected using KinectFusion [51] to record the RGB image, the depth map, and the estimated camera pose. For each scene, several sequences with 500 or 1000 frames are available, splitting into training and test sets. We trained our network using only RGB images, essentially, we did not use depths and camera poses for training. For testing, the provided camera poses are used as the pseudo-ground-truth poses for evaluation.

#### B. Training

At the training time, we select three images to take our loop closed view synthesis. An image is picked at random from whole train sequences, and the other two ‘nearby’ images are picked within 20 steps before or after the target image. Each image is scaled in the range of  $1.0 \sim 1.1$  and cropped by its original size ( $640 \times 480$ ) as a data augmentation. We also augment its color via random brightness, contrast, saturation, and hue jitter with respective ranges of  $\pm 0.2$ ,  $\pm 0.2$ ,  $\pm 0.2$ , and  $\pm 0.1$ . Importantly, the color augmentations are only applied to the images that are fed to the networks, and not to those used to compute  $L_p$ . The weight parameters of the losses are set as  $w_s = 0.001$  and  $w_c = 0.03$ .

The training is performed through a training split for each scene of the 7-scenes dataset [4]. The models were trained within 300 epochs with a batch size of 6, and the ‘nearby’ images were selected at random from whole sequences for

50 % in the latter 100 epochs to learn consistency between images that are distant in time series. During training, we used batch normalization [52] for all the layers except for the output layers, and the Adam [53] optimizer with  $\beta_1 = 0.9$ ,  $\beta_2 = 0.999$ , and a learning rate of 0.0001. The training was executed on TITAN RTX GPUs. We implemented all our systems using the PyTorch [54] framework. In our setting, the training process takes about 50 sec per one training image.

#### C. Results

1) *Camera Re-Localization Performance*: First, we evaluate the camera re-localization performance using a 7-scenes test dataset [4]. At the test time, the final pose of our method is calculated by a median operator instead of the averaging in (5) to achieve robustness against outliers. For the comparison, we employ following baselines:

**ORB-SLAM3** [3] ORB-SLAM3 equips the camera re-localization system to detect the camera pose when tracking with VO is lost. The camera re-localization consists of image retrieval based on the bag-of-words (BOW) feature [57] and pose optimization via RANSAC using ORB features [58]. In this evaluation, we executed ORB-SLAM3 on the training dataset of each scene to create a map and evaluate the re-localization performance on the test dataset. The camera re-localization was not successful in all scenes; therefore, we output the pose of the most matched key-frame in the candidates retrieved by BOW as an estimated pose when the camera re-localization failed.

**DSAC\*** [24] (pose sup. [3], [55]) We trained DSAC\* with pose supervision via ORB-SLAM3 [3], and listed the results of DSAC\* trained with pose supervision via COLMAP [55] in [56] as a reference method of the supervised camera re-localization with SLAM and SfM from RGB sequences.

**monodepth2†** [2] We trained monodepth2 for unsupervised camera re-localization to estimate the camera pose from a single image instead of a relative camera pose in the original monodepth2. Following our method, we calculated the relative pose from estimated camera poses to predict images. This is also an ablation study to confirm the effectiveness of our unsupervised learning framework with the directed scene coordinates and loop closed view synthesis.

Unsupervised camera re-localization methods do not provide scale and alignment of trajectories; therefore, we applied Sim(3) alignments estimated on the training dataset using *evo* [59] for testing. The median re-localization error of position (m) and attitude ( $^\circ$ ) on the 7-scenes test dataset are shown in Table I, which also lists the median errors of conventional supervised camera re-localization methods.

Based on the results of unsupervised camera re-localization, our method outperforms all competitors in most scenes. ORB-SLAM3 camera re-localization cannot provide accurate results because the camera re-localization fails in most frames of each scene.

The monodepth2† is less effective than our method for all scenes. This result confirms the contribution of our learning

TABLE I

**Median re-localization error of position (m) and attitude ( $^{\circ}$ ) on the 7-scenes test dataset [4].** Unsupervised (Unsup.) methods cannot provide their scale and alignment; we applied them to the Sim(3) alignments estimated on the training dataset. Pose Sup. represents the pose supervision (pGT: pseudo-ground-truth poses, [3]; ORB-SLAM3, [55]; COLMAP) is used.

|               | Method                    | Pose Sup.             | Scene                                  |  |  |  |  |  |   |
|---------------|---------------------------|-----------------------|--|--|--|--|--|--|---|
|               |                           |                       | Chess                                  | Fire                                   | Heads                                  | Office                                 | Pumpkin                                | Redkitchen                             | Stairs                                  |
| Supervised    | PoseNet2 [10]             | pGT                   | 0.13m, 4.5 $^{\circ}$                  | 0.27m, 11.3 $^{\circ}$                 | 0.17m, 13.0 $^{\circ}$                 | 0.19m, 5.6 $^{\circ}$                  | 0.26m, 4.8 $^{\circ}$                  | 0.23m, 5.4 $^{\circ}$                  | 0.35m, 12.4 $^{\circ}$                  |
|               | MapNet [12]               | pGT                   | 0.08m, 3.3 $^{\circ}$                  | 0.27m, 11.7 $^{\circ}$                 | 0.18m, 13.3 $^{\circ}$                 | 0.17m, 5.2 $^{\circ}$                  | 0.22m, 4.0 $^{\circ}$                  | 0.23m, 4.9 $^{\circ}$                  | 0.30m, 12.1 $^{\circ}$                  |
|               | NN-Net [13]               | pGT                   | 0.13m, 6.5 $^{\circ}$                  | 0.26m, 12.7 $^{\circ}$                 | 0.14m, 12.3 $^{\circ}$                 | 0.21m, 7.4 $^{\circ}$                  | 0.24m, 6.4 $^{\circ}$                  | 0.24m, 8.0 $^{\circ}$                  | 0.27m, 11.82 $^{\circ}$                 |
|               | ReLocNet [16]             | pGT                   | 0.12m, 4.1 $^{\circ}$                  | 0.26m, 10.4 $^{\circ}$                 | 0.14m, 10.5 $^{\circ}$                 | 0.18m, 5.3 $^{\circ}$                  | 0.26m, 4.2 $^{\circ}$                  | 0.23m, 5.1 $^{\circ}$                  | 0.28m, 7.5 $^{\circ}$                   |
|               | CamNet [17]               | pGT                   | 0.04m, 1.7 $^{\circ}$                  | 0.03m, 1.7 $^{\circ}$                  | 0.05m, 2.0 $^{\circ}$                  | 0.04m, 1.6 $^{\circ}$                  | 0.04m, 1.6 $^{\circ}$                  | 0.04m, 1.6 $^{\circ}$                  | 0.04m, 1.5 $^{\circ}$                   |
|               | DSAC* [24]                | pGT                   | 0.02m, 1.1 $^{\circ}$                  | 0.02m, 1.2 $^{\circ}$                  | 0.01m, 1.8 $^{\circ}$                  | 0.03m, 1.2 $^{\circ}$                  | 0.04m, 1.4 $^{\circ}$                  | 0.03m, 1.7 $^{\circ}$                  | 0.04m, 1.4 $^{\circ}$                   |
|               | - ORB-SLAM3 [3]           |                       | 1.23m, 38.6 $^{\circ}$                 | 1.28m, 59.1 $^{\circ}$                 | 0.42m, 21.1 $^{\circ}$                 | 1.42m, 48.4 $^{\circ}$                 | 1.29m, 36.6 $^{\circ}$                 | 1.72m, 37.7 $^{\circ}$                 | 0.63m, 157.9 $^{\circ}$                 |
| - COLMAP [56] | [55]                      | 0.04m, 1.4 $^{\circ}$ | 0.02m, 1.1 $^{\circ}$                  | 0.01m, 1.5 $^{\circ}$                  | 0.08m, 3.0 $^{\circ}$                  | 0.07m, 3.1 $^{\circ}$                  | 0.06m, 3.2 $^{\circ}$                  | 0.05m, 1.9 $^{\circ}$                  |   |
| Unsup.        | ORB-SLAM3 [3]             |                       | 0.91m, 27.3 $^{\circ}$                 | 0.92m, 29.3 $^{\circ}$                 | 0.17m, <b>5.0<math>^{\circ}</math></b> | <b>0.19m, 5.5<math>^{\circ}</math></b> | 0.43m, 10.0 $^{\circ}$                 | 0.87m, 25.4 $^{\circ}$                 | 0.63m, 175.3 $^{\circ}$                 |
|               | monodepth2 $\ddagger$ [2] |                       | 0.34m, 15.7 $^{\circ}$                 | 0.40m, 20.6 $^{\circ}$                 | 0.26m, 12.0 $^{\circ}$                 | 0.70m, 50.9 $^{\circ}$                 | 0.53m, 12.5 $^{\circ}$                 | 0.66m, 16.1 $^{\circ}$                 | 0.44m, 37.4 $^{\circ}$                  |
|               | <b>Ours (full)</b>        |                       | <b>0.08m, 3.1<math>^{\circ}</math></b> | <b>0.18m, 7.2<math>^{\circ}</math></b> | <b>0.12m, 7.8<math>^{\circ}</math></b> | 0.56m, 39.6 $^{\circ}$                 | <b>0.18m, 4.6<math>^{\circ}</math></b> | <b>0.21m, 7.0<math>^{\circ}</math></b> | 0.42m, 27.5 $^{\circ}$                  |
|               | - w/o DSC                 |                       | 0.24m, 9.4 $^{\circ}$                  | 0.55m, 36.3 $^{\circ}$                 | 0.20m, 9.5 $^{\circ}$                  | 0.29m, 9.8 $^{\circ}$                  | 0.38m, 10.4 $^{\circ}$                 | 0.53m, 12.5 $^{\circ}$                 | <b>0.41m, 38.6<math>^{\circ}</math></b> |
|               | - w/o LCVS                |                       | 0.14m, 5.9 $^{\circ}$                  | 0.30m, 15.1 $^{\circ}$                 | 0.15m, 9.2 $^{\circ}$                  | 0.56m, 32.5 $^{\circ}$                 | 0.28m, 8.0 $^{\circ}$                  | 0.37m, 14.7 $^{\circ}$                 | 0.43m, 31.5 $^{\circ}$                  |

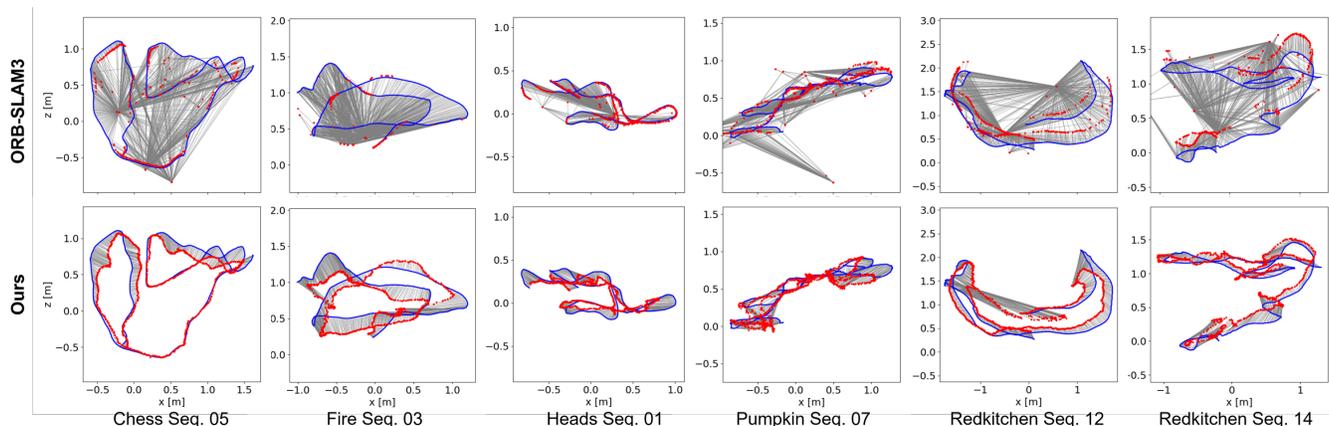


Fig. 5. **Qualitative results of camera re-localization on 7-scenes test dataset [4].** Blue lines represent pseudo-ground-truth camera trajectories, and red points represent predicted camera poses by our method and re-localization of ORB-SLAM3 [3]. Gray lines represent errors between predicted poses and corresponding pseudo-ground-truth poses. We output the pose of the most matched key-frame in the candidates retrieved by BOW as an estimated pose when the camera re-localization of ORB-SLAM3 fails. Therefore, many frames are matched to a key-frame pose which has much features.

framework for camera re-localization based on the directed scene coordinates and loop closed view synthesis.

Performance of our method degrades on the ‘Office’ and ‘Stairs’. This is because there is a similar view on different poses in the scenes. ‘Office’ has similar desks put in different positions in the room. ‘Stairs’ are very repetitive sequences because the camera goes up or down the stairs. Addressing these repetitive structures will be the scope of our future work.

The qualitative results of camera re-localization are shown in Fig. 5. As can be seen from the figures, our method can predict camera poses consistently in all frames. In contrast, ORB-SLAM3 cannot perform camera re-localization on many frames, the most matched candidates which are the alternative outputs cause the degradation in the performance. These results suggest the robustness of our camera re-localization.

Comparing the supervised and unsupervised methods, we can see how challenging unsupervised camera re-localization is. Nevertheless, our methods outperformed several previous

supervised camera re-localization methods, such as PoseNet2 [10], MapNet [12], NN-Net [13], and ReLocNet [16], in most scenes. Although it does not achieve the performance of state-of-the-art supervised camera re-localization, such as CamNet [17] and DSAC\* [24], the fact that our unsupervised learning can achieve an accuracy approaching that of supervised camera re-localization is a significant achievement.

From this result, the DSAC\* [24] can provide accurate results with COLMAP [55], however, it cannot be trained correctly with ORB-SLAM3 [3]. This denotes the performances of supervised methods depend on the accuracy of localization methods. Our method does not achieve the accuracy of DSAC\* based on COLMAP, however, it facilitates estimating the depth accurately online.

We also listed results of our ablation study:

- **w/o DSC** This represents our method without a directed scene coordinate network (DSC). The pose network is the same as in monodepth2 $\ddagger$ .
- **w/o LCVS** This represents our method without loop closed view synthesis (LCVS). The photometric re-projection losses

TABLE II

**Quantitative results for depth-estimation performance.** Comparison of our method to monodepth2 [2] on 7-scenes test dataset [4]. The predicted depth maps are scaled by a scalar that matches the median with the ground-truth for each frame. Standard deviations per median (std/med) of scale factors representing scale consistency are also shown. The evaluation was performed in the range of 0.1 m - 10 m, and averaged over scenes.

| Method                  | Scale factor<br>std / med | Error metric |              |              |              | Accuracy metric |                   |                   |
|-------------------------|---------------------------|--------------|--------------|--------------|--------------|-----------------|-------------------|-------------------|
|                         |                           | Abs Rel      | Sq Rel       | RMSE         | RMSE log     | $\delta < 1.25$ | $\delta < 1.25^2$ | $\delta < 1.25^3$ |
| monodepth2 [2]          | 0.080                     | 0.212        | 0.412        | 0.638        | 0.275        | 0.761           | 0.908             | 0.958             |
| monodepth2 <sup>†</sup> | 0.044                     | 0.220        | 0.445        | 0.663        | 0.270        | 0.748           | 0.913             | 0.961             |
| <b>Ours (full)</b>      | <b>0.017</b>              | <b>0.135</b> | <b>0.059</b> | <b>0.308</b> | <b>0.178</b> | <b>0.825</b>    | <b>0.962</b>      | <b>0.992</b>      |
| - w/o DSC               | 0.030                     | 0.171        | 0.164        | 0.432        | 0.218        | 0.771           | 0.942             | 0.982             |
| - w/o LCVS              | 0.022                     | 0.145        | 0.065        | 0.325        | 0.189        | 0.798           | 0.959             | 0.991             |

are calculated from two relative poses between three images as same as monodepth2<sup>†</sup> [2].

The benefits of DSC and LCVS are clearly shown in Table I, respectively. Focusing on each scene, we can see that the effectiveness of both methods is different for each scene. The effectiveness of DSC on ‘Fire’ and ‘Pumpkin’ is better than that of LCVS, however, DSC suffers from repetitive environments such as ‘Office’ and ‘Stairs.’ As mentioned above, addressing these repetitive structures will be the scope of future work. Although the best method is different for each scene, our full method achieved the best average performance by combining DSC and LCVS.

2) *Depth Estimation Performance:* Next, we evaluated the depth estimation performance using ground-truth depth in the 7-scenes dataset [4]. In this evaluation, we compared our method to a recent monocular depth estimation method, monodepth2 [2], trained on the 7-scenes dataset as follows:

**monodepth2 [2]** Monodepth2 evaluated here is trained by all train sequences for 7-scenes, because it is difficult to train by sequences from each scene respectively. The training was performed in 500 epochs. The monodepth2 uses three-frame image sequences for training; however, the sampling rate of the 7-scenes dataset is too high to train monodepth2. Therefore we generate a three-frame image sequence within 10 step intervals for training.

Table II shows the evaluation results of the depth estimation using the 7-scenes test dataset. Because neither model provides the depth scale, we multiply the predicted depth maps by a scalar that matches the median with the ground-truth for each frame. The relative standard deviations of the scale factors are listed in Table II.

As shown in Table II, our models outperformed monodepth2 on all metrics. It should be noted that with regards to the standard deviation of the scale factor, our model has a much smaller relative standard deviation of scale factors than the monodepth2 model, indicating that our models can predict depth with a more consistent scale across frames. Thus, our method can predict a better and more consistent depth than monodepth2 in known scenes.

We also performed an ablation study for the depth estimation. The benefits of DSC and LCVS for the depth estimation are also clearly shown in Table II, respectively. Our full method achieved the best performance by combining DSC

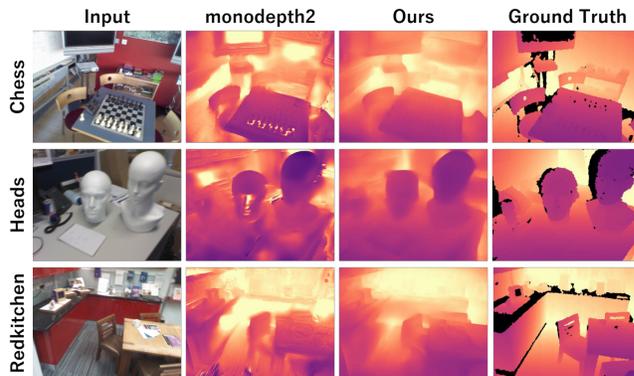


Fig. 6. **Qualitative results of depth on 7-scenes test dataset [4].** Our method has fewer artifacts than monodepth2 [2].

and LCVS.

The qualitative results are shown in Fig. 6, which shows our method has fewer artifacts than monodepth2 [2]. These results suggest that the proposed framework can train a depth network with better performance than a recent monocular depth estimation for indoor scenes. Note that the depth estimation of the proposed framework is trained jointly with camera re-localization; therefore, the depth model is also limited to the trained scene.

## V. CONCLUSION

In this study, we present an unsupervised simultaneous learning framework for camera re-localization and depth estimation from unlabeled video sequences. Our method simultaneously trained two networks for directed scene coordinates and depth map by view synthesis via a video sequence. Camera re-localization is performed by geometrically combining the directed scene coordinates and the depth map. To the best of our knowledge, this is the first study to propose a camera re-localization training method without supervised learning.

We evaluated our method on the 7-scenes dataset [4], and demonstrated its performance in camera re-localization and depth estimation. The results show that our method outperformed the camera re-localization embedded in ORB-SLAM3 [3]. Moreover, despite unsupervised learning, our method outperforms several previous supervised camera

re-localization methods. In addition, it also outperformed monodepth2 [2] on monocular depth estimation in known environments.

However, the current learning framework is limited to frequent and small indoor camera sequences, because the learning framework depends on the overlap of the viewed scene between three or more camera images. Thus, future research will aim at improving the method for applicability to wider scenes and data.

## REFERENCES

- [1] T. Zhou *et al.*, “Unsupervised learning of depth and ego-motion from video,” in *CVPR*, 2017, pp. 1851–1858.
- [2] C. Godard *et al.*, “Digging into self-supervised monocular depth estimation,” in *ICCV*, 2019, pp. 3828–3838.
- [3] C. Campos *et al.*, “Orb-slam3: An accurate open-source library for visual, visual-inertial, and multimap slam,” *IEEE Trans. on Robotics*, 2021.
- [4] J. Shotton *et al.*, “Scene coordinate regression forests for camera relocalization in rgb-d images,” in *CVPR*, 2013, pp. 2930–2937.
- [5] G. Schindler *et al.*, “City-scale location recognition,” in *CVPR*. IEEE, 2007, pp. 1–7.
- [6] A. Torii *et al.*, “24/7 place recognition by view synthesis,” in *CVPR*, 2015, pp. 1808–1817.
- [7] R. Arandjelovic *et al.*, “Netvlad: Cnn architecture for weakly supervised place recognition,” in *CVPR*, 2016, pp. 5297–5307.
- [8] F. Walch *et al.*, “Image-based localization using lstms for structured feature correlation,” in *ICCV*, 2017, pp. 627–637.
- [9] A. Kendall *et al.*, “Posenet: A convolutional network for real-time 6-dof camera relocalization,” in *ICCV*, 2015, pp. 2938–2946.
- [10] —, “Geometric loss functions for camera pose regression with deep learning,” in *CVPR*, 2017, pp. 5974–5983.
- [11] T. Naseer *et al.*, “Deep regression for monocular camera-based 6-dof global localization in outdoor environments,” in *IROS*. IEEE, 2017, pp. 1525–1530.
- [12] S. Brahmabhatt *et al.*, “Geometry-aware learning of maps for camera localization,” in *CVPR*, 2018, pp. 2616–2625.
- [13] Z. Laskar *et al.*, “Camera relocalization by computing pairwise relative poses using convolutional neural network,” in *ICCV Workshops*, 2017, pp. 929–938.
- [14] A. Valada *et al.*, “Deep auxiliary learning for visual localization and odometry,” in *ICRA*. IEEE, 2018, pp. 6939–6946.
- [15] T. Sattler *et al.*, “Understanding the limitations of cnn-based absolute camera pose regression,” in *CVPR*, 2019, pp. 3302–3312.
- [16] V. Balntas *et al.*, “Relocnet: Continuous metric learning relocalisation using neural nets,” in *ECCV*, 2018, pp. 751–767.
- [17] M. Ding *et al.*, “Camnet: Coarse-to-fine retrieval for camera relocalization,” in *ICCV*, 2019, pp. 2871–2880.
- [18] J. Valentin *et al.*, “Exploiting uncertainty in regression forests for accurate camera relocalization,” in *CVPR*, 2015, pp. 4400–4408.
- [19] A. Guzman-Rivera *et al.*, “Multi-output learning for camera relocalization,” in *CVPR*, 2014, pp. 1114–1121.
- [20] L. Meng, F. Tung, J. J. Little, J. Valentin, and C. W. de Silva, “Exploiting points and lines in regression forests for rgb-d camera relocalization,” in *IROS*. IEEE, 2018, pp. 6827–6834.
- [21] E. Brachmann *et al.*, “Uncertainty-driven 6d pose estimation of objects and scenes from a single rgb image,” in *CVPR*, 2016, pp. 3364–3372.
- [22] L. Meng *et al.*, “Backtracking regression forests for accurate camera relocalization,” in *IROS*. IEEE, 2017, pp. 6886–6893.
- [23] E. Brachmann and C. Rother, “Learning less is more-6d camera localization via 3d surface regression,” in *CVPR*, 2018, pp. 4654–4662.
- [24] E. Brachmann *et al.*, “Visual camera re-localization from rgb and rgb-d images using dsac,” *IEEE TPAMI*, 2021.
- [25] R. Clark *et al.*, “Vidloc: A deep spatio-temporal model for 6-dof video-clip relocalization,” in *CVPR*, 2017, pp. 6856–6864.
- [26] N. Radwan *et al.*, “Vlocnet++: Deep multitask learning for semantic visual localization and odometry,” *IEEE RA Letters*, vol. 3, no. 4, pp. 4407–4414, 2018.
- [27] F. Xue *et al.*, “Local supports global: Deep camera relocalization with sequence enhancement,” in *ICCV*, 2019, pp. 2841–2850.
- [28] L. Zhou *et al.*, “Kfnet: Learning temporal camera relocalization using kalman filtering,” in *CVPR*, 2020, pp. 4919–4928.
- [29] D. Eigen *et al.*, “Depth map prediction from a single image using a multi-scale deep network,” in *NeurIPS*, 2014, pp. 2366–2374.
- [30] N. Mayer *et al.*, “A large dataset to train convolutional networks for disparity, optical flow, and scene flow estimation,” in *CVPR*, 2016, pp. 4040–4048.
- [31] A. Dosovitskiy *et al.*, “Flownet: Learning optical flow with convolutional networks,” in *ICCV*, 2015, pp. 2758–2766.
- [32] B. UmmeHofer *et al.*, “Demon: Depth and motion network for learning monocular stereo,” in *CVPR*, 2017, pp. 5038–5047.
- [33] R. Garg *et al.*, “Unsupervised cnn for single view depth estimation: Geometry to the rescue,” in *ECCV*. Springer, 2016, pp. 740–756.
- [34] C. Godard *et al.*, “Unsupervised monocular depth estimation with left-right consistency,” in *CVPR*, 2017, pp. 270–279.
- [35] M. Jaderberg *et al.*, “Spatial transformer networks,” *NeurIPS*, vol. 28, pp. 2017–2025, 2015.
- [36] V. Casser *et al.*, “Depth prediction without the sensors: Leveraging structure for unsupervised learning from monocular videos,” in *AAAI*, 2019, pp. 8001–8008.
- [37] M. Klodt *et al.*, “Supervising the new with the old: learning sfm from sfm,” in *ECCV*, 2018, pp. 698–713.
- [38] R. Mahjourian *et al.*, “Unsupervised learning of depth and ego-motion from monocular video using 3d geometric constraints,” in *CVPR*, 2018, pp. 5667–5675.
- [39] C. Wang *et al.*, “Learning depth from monocular videos using direct methods,” in *CVPR*, 2018, pp. 2022–2030.
- [40] L. Zhou *et al.*, “Unsupervised learning of monocular depth estimation with bundle adjustment, super-resolution and clip loss,” *arXiv preprint arXiv:1812.03368*, 2018.
- [41] Y. Zou *et al.*, “Df-net: Unsupervised joint learning of depth and flow using cross-task consistency,” in *ECCV*, 2018, pp. 36–53.
- [42] N. Hirose *et al.*, “Plg-in: Pluggable geometric consistency loss with wasserstein distance in monocular depth estimation,” in *ICRA*. IEEE, 2021, pp. 12 868–12 874.
- [43] —, “Variational monocular depth estimation for reliability prediction,” in *3DV*. IEEE, 2021, pp. 637–647.
- [44] K. He *et al.*, “Deep residual learning for image recognition,” in *CVPR*, 2016, pp. 770–778.
- [45] A. Mallya *et al.*, “Packnet: Adding multiple tasks to a single network by iterative pruning,” in *CVPR*, 2018, pp. 7765–7773.
- [46] N. Yang *et al.*, “Deep virtual stereo odometry: Leveraging deep depth prediction for monocular direct sparse odometry,” in *ECCV*, 2018, pp. 817–833.
- [47] Z. Yin *et al.*, “Geonet: Unsupervised learning of dense depth, optical flow and camera pose,” in *CVPR*, 2018, pp. 1983–1992.
- [48] Z. Liang *et al.*, “Deep unsupervised learning based visual odometry with multi-scale matching and latent feature constraint,” in *IROS*, 2021.
- [49] O. Ronneberger *et al.*, “U-net: Convolutional networks for biomedical image segmentation,” in *MICCAI*. Springer, 2015, pp. 234–241.
- [50] O. Russakovsky *et al.*, “Imagenet large scale visual recognition challenge,” *IJCV*, vol. 115, no. 3, pp. 211–252, 2015.
- [51] S. Izadi *et al.*, “Kinectfusion: real-time 3d reconstruction and interaction using a moving depth camera,” in *UIST*, 2011, pp. 559–568.
- [52] S. Ioffe *et al.*, “Batch normalization: Accelerating deep network training by reducing internal covariate shift,” in *ICML*. PMLR, 2015, pp. 448–456.
- [53] D. P. Kingma *et al.*, “Adam: A method for stochastic optimization,” *arXiv preprint arXiv:1412.6980*, 2014.
- [54] A. Paszke *et al.*, “Automatic differentiation in pytorch,” in *CVPR Workshop*, 2017.
- [55] J. L. Schonberger *et al.*, “Structure-from-motion revisited,” in *CVPR*, 2016, pp. 4104–4113.
- [56] E. Brachmann *et al.*, “On the limits of pseudo ground truth in visual camera re-localisation,” in *ICCV*, 2021, pp. 6218–6228.
- [57] D. Gálvez-López *et al.*, “Bags of binary words for fast place recognition in image sequences,” *IEEE Trans. on Robotics*, vol. 28, no. 5, pp. 1188–1197, 2012.
- [58] E. Rublee *et al.*, “Orb: An efficient alternative to sift or surf,” in *ICCV*. Ieee, 2011, pp. 2564–2571.
- [59] M. Grupp, “evo: Python package for the evaluation of odometry and slam.” <https://github.com/MichaelGrupp/evo>, 2017.