

Unsupervised Learning-based Depth Estimation aided Visual SLAM Approach

Mingyang Geng¹, Suning Shang¹, Bo Ding^{1*}, Huaimin Wang¹, Pengfei Zhang¹, and Lei Zhang²

National Key Laboratory of Parallel and Distributed Processing,
College of Computer, National University of Defense Technology, China¹
National Key Laboratory of Integrated Automation of Process Industry,
Northeastern University, China²

Abstract. Existing visual-based SLAM systems mainly utilize the three-dimensional environmental depth information from RGB-D cameras to complete the robotic synchronization localization and map construction task. However, the RGB-D camera maintains a limited range for working and is hard to accurately measure the depth information in a far distance. Besides, the RGB-D camera will easily be influenced by strong lighting and other external factors, which will lead to a poor accuracy on the acquired environmental depth information. Recently, deep learning technologies have achieved great success in the visual SLAM area, which can directly learn high-level features from the visual inputs and improve the estimation accuracy of the depth information. Therefore, deep learning technologies maintain the potential to extend the source of the depth information and improve the performance of the SLAM system. However, the existing deep learning-based methods are mainly supervised and require a large amount of ground-truth depth data, which is hard to acquire because of the realistic constraints. In this paper, we first present an unsupervised learning framework, which not only uses image reconstruction for supervising but also exploits the pose estimation method to enhance the supervised signal and add training constraints for the task of monocular depth and camera motion estimation. Furthermore, we successfully exploit our unsupervised learning framework to assist the traditional ORB-SLAM system when the initialization module of ORB-SLAM method could not match enough features. Qualitative and quantitative experiments have shown that our unsupervised learning framework performs the depth estimation task comparably to the supervised methods and outperforms the previous state-of-the-art approach by 13.5% on KITTI dataset. Besides, our unsupervised learning framework could significantly accelerate the initialization process of ORB-SLAM system and effectively improve the accuracy on environmental mapping in strong lighting and weak texture scenes.

Keywords: Robotic visual SLAM, monocular depth estimation, pose estimation, unsupervised learning

1 Introduction

Simultaneous Localization and Mapping (SLAM) has attracted increasing attention in the robotic areas. SLAM technologies have wide applications in area such as autonomous driving, localization and navigation. The goal of a SLAM system is to construct the map of an unknown environment incrementally based on the perception information, i.e., scene information acquired by a radar or depth sensor when the robot is performing a complex task and confronted with an unknown environment. In order to achieve a satisfying performance in the visual SLAM tasks, the quality of the perception on the environmental depth, i.e., the distance of the objects in the environment, will play an indispensable role. Therefore, how to extract valuable depth information from the visual inputs is an important problem in the visual SLAM systems.

Existing visual-based SLAM systems mainly utilize the three-dimensional environmental depth information from RGB-D cameras. However, the RGB-D camera maintains a limited range for working and is hard to accurately measure the depth information in a far distance. Besides, in some special scenes, i.e., strong lighting and weak texture environments, robotic visual SLAM always faces the problems of scale drift or scale error because of the inaccurate accuracy on the acquired depth information. The reason of obtaining the imprecise depth information is that most of the existing visual SLAM algorithms design sparse image features manually, while the manually designed features often contain certain assumptions about the environment, i.e., sufficient illumination, material determination, which will lead to a poor performance on environmental depth estimation when the environmental factors change.

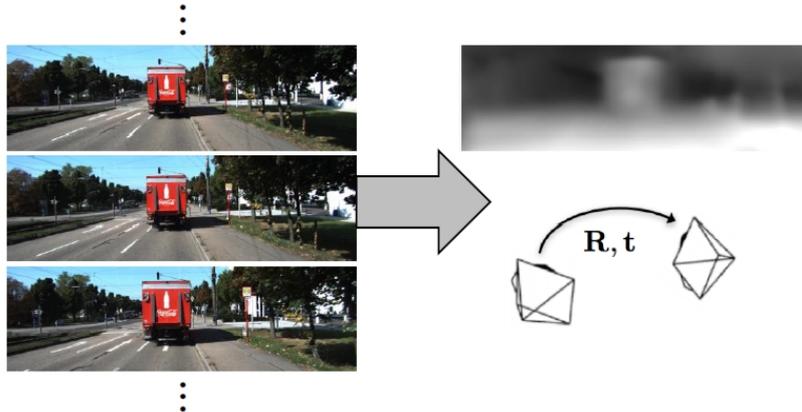


Fig. 1: Illustration of our learning framework. The input to our system consists solely of unlabeled video clips. Our learning framework estimates the depth in the first image and the camera motion.

Deep learning technologies have recently emerged as a powerful tool for improving the accuracy on monocular depth estimation [10, 15, 23, 34]. One of the advantages of deep learning technologies over common alternatives is that features are learned directly from data, and do not have to be chosen or designed by the algorithm developers for the specific problem on which they are applied. However, existing methods [10, 15, 23] are mainly supervised and need a large amount of ground-truth data, which is hard to acquire because of the expensive radar sensors and the limited working range. A promising branch in the depth estimation field is unsupervised learning [34], which exploits image reconstruction as supervision signal and significantly reduce the burden to collect high-quality depth training data in advance. However, the existing unsupervised method [34] does not fully exploit the heuristic knowledge during the image acquisition process, which can strengthen the supervised signal and further improve the accuracy on depth estimation. Therefore, how to enhance the supervised signal and utilize the unsupervised learning technologies to assist the traditional visual SLAM systems remains a great challenge.

In this paper, we first propose a novel unsupervised learning framework which exploit the pose estimation method to enhance the supervised signals and further promote the accuracy of extracting the depth information from monocular image sequences. In concrete, we utilize a large number of scene image sequences to train a model for camera motion prediction and scene structure prediction (shown in Fig. 1). In the pose estimation stage, we set up a continuous frame window and exploit the pose transformation relationships to construct the pose graph, which can partially eliminate the cumulative error. Furthermore, we successfully exploit our unsupervised learning framework to assist the traditional ORB-SLAM system, a widely used visual SLAM system, when the initialization module of ORB-SLAM method could not match enough features. Extensive experiments have shown that our method can significantly accelerate the initialization process of ORB-SLAM system and effectively improve the accuracy on environmental mapping in strong lighting and weak texture scenes.

The rest of this paper is organized as follows. Section 2 introduces the background and the highly related work. Section 3 describes the methodology of our work as well as the architecture designed for training and prediction. Section 4 describes the details of our unsupervised learning-based depth estimation-aided visual SLAM system. The validation and evaluation of our work based on different public datasets are described in Section 5. Section 6 presents the experimental results implemented in the real-world settings. We conclude and provide our future direction in Section 7.

2 Related Work

Our method covers three research areas, including depth estimation optimization, monocular depth estimation and motion estimation from images.

2.1 Depth Estimation optimization

The existing mechanism of depth estimation optimization can be mainly divided into three categories based on how to utilize the deep learning technologies.

The first kind of methods directly substitute the depth information acquisition module with deep learning technologies [12]. This method can effectively suit for the environments which are hard for traditional visual SLAM methods to deal with, i.e., strong lighting and weak texture environments. However, this method will require a more complex system with a stronger computing and processing ability. In addition, deep learning technologies are easily over-fit to the dataset and will lead to a poor performance in the unfamiliar environments. Last but not least, there does not exist a method which fully substitute the depth information acquisition module with deep learning technologies in realistic applications. Therefore, this kind of method needs further validation.

The second kind of methods exploit the environmental depth information from the deep learning technologies and the traditional SLAM system simultaneously [22]. This method can optimize the environmental depth information used by the visual SLAM system and indirectly improve the accuracy on mapping and localization. However, this method needs to implement the two depth information acquisition methods simultaneously and require a stronger ability of computing and processing. In addition, the complex evaluation process to select the optimal depth information needs to be accurately implemented, which is hard to satisfy the quality of service.

The third kind of methods complement the drawbacks of the deep learning technologies and the depth information acquisition module of traditional SLAM algorithm [32]. In concrete, deep learning technologies is only applied when the traditional SLAM algorithms can not obtain high-accuracy depth information. This method can effectively improve the accuracy of the depth estimation while maintaining the ability to guarantee the QoS requirement. There are very limited works which exploit the deep learning technologies to acquire the environmental depth estimation. A deep learning-aided LSD-SLAM algorithm is proposed in [11], which achieves a better result than the traditional LSD-SLAM algorithm and a stronger adaptability to the strong lighting and weak texture environments. We choose ORB-SLAM as our baseline method, which has wide applications in visual SLAM area but performs an unsatisfied performance in strong lighting and weak texture environments. To the best of our knowledge, this is the first work which combines the deep learning technologies and the ORB-SLAM system.

2.2 Monocular Depth Estimation

Monocular depth estimation is a basic low-level challenge problem which has been studied for decades. Early works on depth estimation using RGB images usually relied on hand-crafted features and probabilistic graphical models. [16] introduced photo pop-up, a fully automatic method for creating a basic 3D model from a single photograph. In [18], the authors design Depth Transfer, a

non-parametric approach where the depth of an input image is reconstructed by transferring the depth of multiple similar images and then applying some warping and optimizing procedures. Delage et al. in [7] proposed a dynamic Bayesian framework for recovering 3D information from indoor scenes. A discriminatively-trained multi-scale Markov Random Field (MRF) was introduced in [29], in order to optimally fuse local and global features. Depth estimation is considered as an inference problem in a discrete-continuous CRF in [24].

More recent approaches for depth estimation are based on convolutional neural network(CNN). As a pioneer work, Eigen et al. proposed a multi-scale approach for depth prediction in [10]. It considers two deep networks, one performing a coarse global prediction based on the entire image, and the other refining predictions locally. This approach was extended in [9] to handle multiple tasks (e.g. semantic segmentation, surface normal estimation). In [23], authors combine a deep CNN and a continuous conditional random field, and attain visually sharper transitions and local details. In [21], a deep residual network is developed, based on the ResNet and achieved higher accuracy than [23]. Unlike our approach, these methods require explicit depth for training. Unsupervised learning setups have also been explored for disparity image prediction. For instance, Godard et al. formulate disparity estimation as an image reconstruction problem in [15], where neural networks are trained to warp left images for matching the right one. Though these methods show similarity with ours, which are unsupervised without requiring ground-truth depth data for training, they assume camera poses known in advance, which is treated a large simplification. Our work is inspired by that of [34], which proposes to use view synthesis as the supervisory signal. However, the further advantage of our approach which demonstrated in the following evaluations is that, the idea of continuous frame window used in traditional SLAM approach is applied to enhance the supervisory signal and capture more constraints which can guide the training process for more accuracy results.

2.3 Motion Estimation from Images

The motion estimation has a long history in computer vision. The underlying 3D geometry is a consolidated field. They consist of a long pipeline of methods, start from descriptor matching for finding a sparse set of correspondences between images [26], to estimating the essential matrix to determine the camera motion. Bundle adjustment [33] is used in the pipeline of method to refine the final structure and camera position. The bundle adjustment minimizes the reprojection error of the three-dimensional point in the two-dimensional image sequence by Levenberg-Marquardt (LM) nonlinear algorithm to get the optimal motion model [25]. The accuracy of the bundle adjustment method is related to the number of frames of the image. The more the number of image frames, the more accurate the camera motion parameters can obtain.

Recent works [8] propose learning frame-to-frame motion fields with deep neural networks supervised with ground-truth motion obtained from simulation or synthetic movies. This enables efficient motion estimation that learns to deal

with lack of texture using training examples rather than relying only on smoothness constraints of the motion field, as previous optimization methods [31]. Our approach draws on the respective advantages of the geometry-based motion estimation in SLAM and the learning-based motion estimation. Multiview pose network is used to estimate pose transformation matrix between adjacent frames. We set up a continuous frame window to construct the pose graph and use the pose transform relationship to calculate more pose transform matrices which perfect the pose graph.

3 Pose Estimation-based Monocular Depth Estimation Method

The accuracy of existing monocular depth estimation methods is hard to satisfy the requirement realistic applications. Therefore, it is meaningful to improve the accuracy of monocular depth estimation. In this section, we introduce our pose estimation-based monocular depth estimation method from the following three aspects: the basic framework, learning and geometry-based pose estimation and image recovery-based training method.

3.1 Framework

Given a single image frame I , the goal of our method is to provide two functions f_1 and f_2 which can predict the per-pixel scene depth $d' = f_1 I$ and the camera pose $p' = f_2 I$. We design two deep neural networks (depth estimation neural network and pose estimation neural network) to learn these two functions. Most existing methods treat the learning task as a supervised learning problem, where the color input images, the corresponding target depth and pose values are provided. However, it is not practical to acquire such large amount of ground-truth depth and pose data in various scenes because of the expensive lidar sensor and the limited working range. Besides, existing methods always neglect the traditional pose estimation methods and do not take the prior knowledge from traditional algorithms into account.

We propose an unsupervised learning method, which exploits the pose estimation approach in traditional SLAM algorithms to augment the supervised signals by image reconstruction. In concrete, based on a short image sequences I_i, I_{i+1}, I_{i+2} captured by a moving camera, we can reconstruct the image I'_i by the predicted the depth image D_i and the predicted pose estimation matrix. The difference between the image I_i and the reconstructed image I'_i can be used as the supervised signals to train the depth estimation and the pose estimation neural networks. The framework of our unsupervised method is illustrated in Fig. 2. The monocular video sequences are used for training the single-view depth estimation and the multi-view pose estimation networks. The output of the single-view depth estimation network is the depth map of the input image. For the pose estimation part, a continuous frame window are used as the input and the output of the multi-view pose estimation network is the pose transform

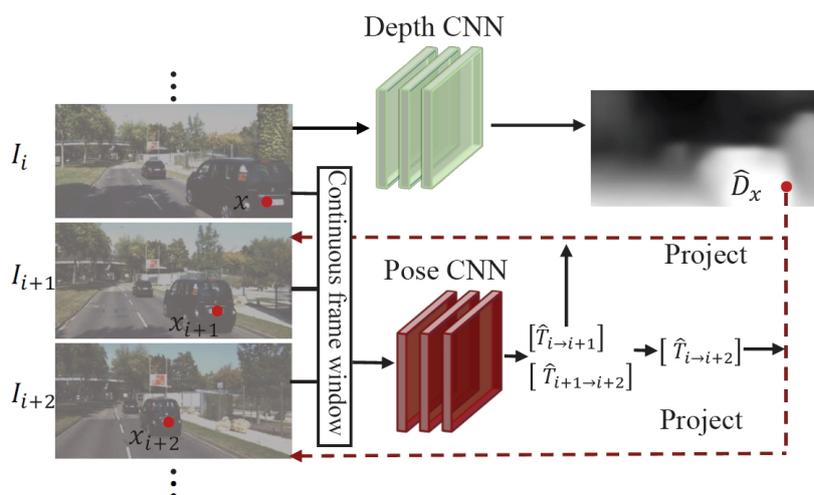


Fig. 2: The overview of the training pipeline based on image reconstruction. The depth network takes only the first image I_i as input and outputs a per-pixel depth map \hat{D}_x . The continuous frame window takes images (e.g., I_i , I_{i+1} , I_{i+2}) as input, through the pose network, outputs the relative camera pose matrices of adjacent frames ($\hat{T}_{i \rightarrow i+1}$, $\hat{T}_{i+1 \rightarrow i+2}$) and we can use the camera pose matrices to calculate more camera poses ($\hat{T}_{i+1 \rightarrow i+2}$). The outputs of both models have then used to inverse warp images to reconstruct the target image, and the photometric reconstruction loss is used for training the CNNs. By utilizing image reconstruction as supervision, we are able to train the entire framework in an unsupervised manner from videos.

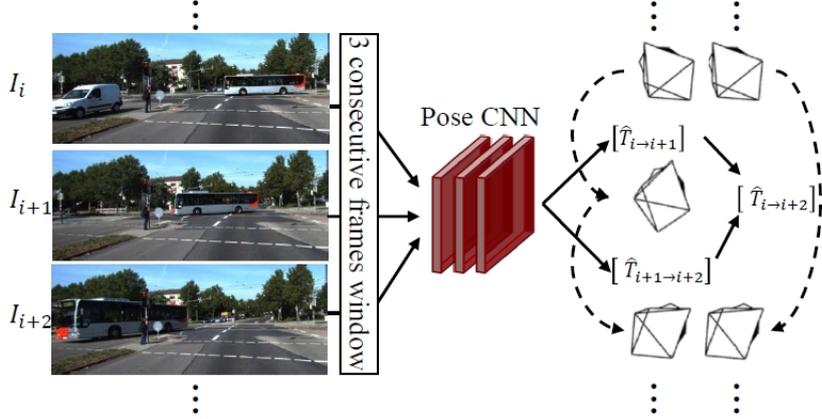


Fig. 3: Illustration of the camera pose estimation pipeline. For example, we maintain a continuous frame window whose length is 3 frames for estimating the camera pose transform matrix. For each sequence images (e.g., I_i, I_{i+1}, I_{i+2}), we will get the adjacent camera pose transform matrix ($\hat{T}_{i \rightarrow i+1}, \hat{T}_{i+1 \rightarrow i+2}$). We can use the camera pose transformation to get more camera poses $\hat{T}_{i \rightarrow i+2}$.

matrices between all adjacent frames in the continuous frame window. We then optimize the pose graph in the continuous frame window by calculating more nonadjacent pose transform matrices using a pose transform relationship. Then, we can reconstruct the image by the depth map and the pose transform matrices and train the two neural networks by calculating the difference between the input and the reconstructed images.

3.2 Estimation based on Learning and Geometry Pose Graph

Given a frame I_i , the single-view depth estimation network can directly predict the corresponding depth map \hat{D}_i . For the pose estimation part, our method is based both on the unsupervised learning technologies and the traditional geometry pose graph. A continuous frame window is set up before the multi-view pose estimation network in order to make the network learn the pose relationship between the continuous images $\langle I_1, I_2, \dots, I_n \rangle$ simultaneously. The length of the continuous frame window stands for the number of input images in a training episode. In other words, the continuous frame window sequentially reads n images from the training set and then sends the training data to the multi-view pose estimation network for further processing. During the training process, the multi-view pose estimation network will sequentially predict the transformation matrix between two adjacent frames in the continuous frame window (shown in Fig. 3). For a better illustration, denote the input image sequences in the continuous frame window as $\langle I_1, I_2, \dots, I_n \rangle$, the output of the multi-view pose estimation network is $\hat{T}_{i \rightarrow i+1}$. Then, we can build a preliminary pose graph us-

ing the pose transformation matrices between the adjacent frames. However, the preliminary pose graph lacks the pose transform matrices between non-adjacent frames, so we calculate the non-adjacent pose transformation relationships by using the following function:

$$\hat{T}_{i \rightarrow i+1} \times \hat{T}_{i+1 \rightarrow i+2} = \hat{T}_{i \rightarrow i+2} \quad (1)$$

Similarly, we can increase the length of the frame window to get more non-adjacent pose transformation relationships (e.g., $\hat{T}_{i \rightarrow i+5}$) and improve the pose graph. The acquired pose transformation relationships maintain the following two advantages. First, the cumulative error is partially eliminated. In the continuous frame window, the error of pose estimation between adjacent frames will accumulate gradually. But if we use the calculated pose matrix to reconstruct the image, we can sequentially adjust the parameters of the pose network and partially eliminate the cumulative error. The second advantage is that this mechanism can avoid the estimation errors between the frames which are far apart in the frame sequence. Experiments [34] have shown that learning-based methods could not predict a satisfying relationship between the two frames which maintain a far distance in the frame sequence. Our method solve this problem by calculating the far apart pose relationships based on the pose estimation of the adjacent frames.

3.3 Geometry-based Image Reconstruction

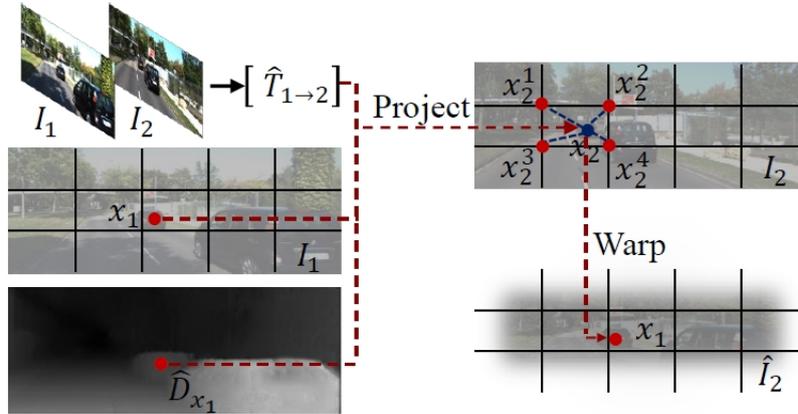


Fig. 4: Illustration of image reconstruction based on camera pose matrix. For each point (e.g., x_1) in the first image, we project it onto the other image base on the predicted depth and camera pose and then use bilinear interpolation to obtain the value of the warped image (\hat{I}_2) at location (x_1).

Image reconstruction through the means of warps and camera projection is an important application of geometric scene understanding. The goal of image reconstruction is to reconstruct a new viewpoints image from other viewpoints through warps and camera projection. In our learning framework, we reconstruct the target image I_t by sampling pixels from the other images I_r based on the predicted depth map \hat{D}_t and the predicted 4×4 camera pose transformation matrix $\hat{T}_{t \rightarrow r}$.

Our camera model is the pinhole model. Denote K as the camera intrinsic matrix, I_1 and I_2 , as the first and the second image in a training episodes. The transformation matrices of the two images to the world coordinates are represented as $T_{1 \rightarrow w}$ and $T_{2 \rightarrow w}$, and the homogeneous coordinates of a pixel in the first image is represented as x_1 . We can acquire the projected coordinates of x_1 onto the second image x_2 by $x_2 \sim K T_{2 \rightarrow w} T_{2 \rightarrow w}^{-1} K^{-1} \hat{D}_1(x_1) x_1$. Notice that the camera transformation matrix between I_1 and I_2 is equal to $T_{2 \rightarrow w} T_{1 \rightarrow w}^{-1}$, i.e., $\hat{T}_{1 \rightarrow 2} = T_{2 \rightarrow w} T_{1 \rightarrow w}^{-1}$. We substitute $T_{1 \rightarrow w} T_{2 \rightarrow w}^{-1}$ with $\hat{T}_{1 \rightarrow 2}$ so the formula becomes $x_2 \sim K \hat{T}_{1 \rightarrow 2} K^{-1} \hat{D}_1(x_1) x_1$. Furthermore, when we get a short image sequences $\langle I_1, I_2, \dots, I_n \rangle$ at the training time, denote I_t as the target view image, and I_r as the other images. The pixel project procedure can be formulated as:

$$x_r \sim K \hat{T}_{t \rightarrow r} K^{-1} \hat{D}_t(x_t) x_t. \quad (2)$$

Based on Eq.2, we can project the pixels on the target image I_t onto other images I_r . After that, our image reconstruction model uses the image sampler from the spatial transformer network (STN) to sample the projected image I_r . The STN uses bilinear sampling where the output pixel is the weighted sum of the four pixel neighbors $x_r^{(1)}, x_r^{(2)}, x_r^{(3)}, x_r^{(4)}$ of x_r , i.e., $\hat{I}_r(x_t) = I_r(x_r) = \sum_{i=1}^4 w^{(i)} I_r(x_r^i)$, where $w^{(i)}$ is linearly proportional to the spatial proximity between x_r and x_r^i , and $\sum_{i=1}^4 w^{(i)} = 1$ (shown in Fig. 4). Contrast with the alternative approaches [13], the bilinear sampler is locally fully differentiable and integrates seamlessly into our fully convolutional network, which means that we do not require any simplification or approximation of our cost function.

Finally, we use the predicted depth map \hat{D}_t and the predicted 4×4 camera pose transformation matrix $\hat{T}_{t \rightarrow r}$ of the previous step to reconstruct the target image through projection and the differentiable bilinear sampling mechanism.

3.4 Image Reconstruction as Supervision

Image reconstruction has been used to learn end-to-end unsupervised optical flow [17], disparity flow in a stereo rig [15] and video prediction [28]. These methods reconstruct the images by transforming the input based on depth maps or flow fields. Our work considers dense structure estimation and uses monocular videos to obtain the necessary self-supervision, instead of static images. The depth information could also be predicted from a single image supervised by photometric error [13]. However, the methods above do not infer camera pose transform or object motion and require stereo pairs with known baseline in the training process. Our work estimates the camera motion between frames, which

the activation blocks of the encoder, which enables the ability to obtain more representative features. The disparity are predicted at four different scales (from disp4 to disp1). The function of the pose network is to predict the relative poses between the target image and other input images, which are accurately described by the 6-dimensional camera pose transform matrix (3-dimensional euler angles and 3-dimensional translation).

In our learning framework, the gradients are mainly derived from the pixel intensity difference between the four-pixel neighbors of x_r and x_t . Therefore, the training process will be inhibited when the correct x_r (projected using the ground-truth depth and pose) is located in a low-texture region or far from the current estimation. We solve this problem by using multi-scale and smoothness loss which allows gradients to be derived from larger spatial regions directly.

Denote the loss at each output scale as C_s , so the total loss can be represented as $C = \sum_{s=1}^4 C_s$. Our loss module calculates C_s as a combination of the two main terms:

$$C_s = C_{vr} + \lambda C_{smooth}, \quad (4)$$

which encourages the reconstructed image to appear similar to the corresponding training input, indexes the minimized norm of the second-order gradients for the predicted depth maps. λ denotes the weighting for the depth smoothness loss.

4 Visual SLAM System with the Assistance of Unsupervised Learning-based Depth Estimation

Existing visual SLAM system always acquires the depth information from the depth sensor. The depth sensor can directly obtain the environmental depth information within a certain distance. However, the depth sensors suffer from the limited working range and are sensitive to the interference, which will decrease the acquired accuracy. In Section 3, we propose an image depth estimation method based on camera pose transformation relationship, which can directly obtain the depth information from the image through unsupervised learning. In this section, we introduce our method which extends the source of the three-dimensional environmental depth information and improves the accuracy of the depth information on the basis of ORB-SLAM algorithm.

4.1 Traditional ORB-SLAM Algorithm

The ORB-SLAM algorithm is mainly composed of the following three parts: the tracking module, the local map building module and the closed-loop detection [3] module. The tracking module aims to locate the camera through each frame and determines whether the frame should join the key-frame set. The tracking task first extracts the features of the frame, then initializes the camera pose and map and tracks the local map, finally determines the key frame. The local map

building module aims to process the new key-frames and rebuild the map through local BA [2]. The closed-loop detection module mainly judges the newly added key-frame and determines whether the scene has been encountered before.

In the initialization phase of ORB-SLAM, the monocular camera first reads the RGB image, and the depth sensor acquires the depth data. Then, feature point matching process is performed on the two consecutive frames in the time series to determine the motion of the camera. The image and depth acquisition process is implemented until the number of matched feature points on the two consecutive frames reaches a specified threshold. Then, the ORB-SLAM algorithm utilizes the existing geometric relationship between the matched feature points to calculate the current pose of the camera, and further creates an initial map or updates the local map.

4.2 Depth Estimation Optimizing Mechanism

By analyzing the ORB-SLAM algorithm, we decide to optimize the depth information acquisition process in the initialization phase. In concrete, our unsupervised learning-based depth estimation mechanism is applied only when the ORB-SLAM system does not match enough feature points. In this way, the accuracy of the ORB-SLAM will be effectively optimized while maintaining a reasonable computing ability and satisfying the QoS requirement.

The initialization phase of ORB-SLAM algorithm is optimized as follows. When there are not enough matched feature points between two consecutive frames, the RGB image will be transmitted to our monocular depth estimation network instead of making the system read the next frame. Then, our depth estimation network will re-estimate the environmental depth information of the input image and further implement the feature point matching process. The complete depth information optimization mechanism on the basis of ORB-SLAM algorithm is illustrated below:

4.3 Implement of Depth Estimation Assisted Visual SLAM System based on Unsupervised Learning

We now introduce the details to realize the optimized ORB-SLAM system using our unsupervised learning-based depth estimation method. For the training process of our monocular depth estimation networks, the training dataset includes samples from various scenes (i.e., indoor scenes, outdoor scenes) as well as samples of weak texture and strong light. After the training process, we decoupled the trained network because the optimization mechanism only uses the depth estimation network. In addition, we establish the transmission channel of the ORB-SLAM initialization module and the trained depth estimation network. When the ORB-SLAM algorithm does not match enough feature points, the RGB image is transmitted to the depth estimation network, which keeps the running state to guarantee a immediate output.

Algorithm 1 The initialization process of the optimized ORB-SLAM algorithm using our monocular depth estimation network.

Require: The threshold M , which stands for the number of the matched feature points required by the ORB-SLAM algorithm.

- 1: The monocular camera reads the RGB image and the depth sensor acquires the depth information.
 - 2: Feature point matching process is implemented on the two consecutive frames in time series based on the RGB image and the depth information.
 - 3: **if** the number of matched feature points is less than the threshold M **then**
 - 4: Transmit the current image to the monocular depth estimation network and get the new depth image.
 - 5: Implement the feature point matching process on the two consecutive frames in time series based on the RGB image and the new depth image.
 - 6: **if** the matched feature points still less than M **then**
 - 7: Read the next RGB image and get the corresponding depth information.
 - 8: **else**
 - 9: Get the pose of the camera by utilizing the geometric relationship between the matched feature points.
 - 10: **end if**
 - 11: **else**
 - 12: Get the pose of the camera by utilizing the geometric relationship between the matched feature points.
 - 13: **end if**
 - 14: Create the initial map or update the local map.
-

5 Evaluation on the Testing Set

In this section, we evaluate the performance of our approach and make comparison with the existing methods on single-view depth and ego-motion estimation. We choose KITTI dataset as the test benchmark. To evaluate the cross-dataset generalization ability of our approach and demonstrate the superiority on the strong lighting and weak texture environments, we also use the Make3D dataset for a better illustration.

5.1 Training Details

We implement our system in TensorFlow [1]. In all experiments, the value of λ is set to 0.5. During the training process, we use the Adam optimizer [19] with β_1 of 0.9, learning rate of 0.0002 and mini-batch size of 4. All the experiments are performed with image sequences captured with a monocular camera and the images are resized to 128×416 during training. In the test phase, the depth and pose networks can be applied independently for images of arbitrary size.

5.2 Single-view Depth Estimation

We present results for the KITTI dataset [14] using two different test splits, to enable comparison to existing works. In its raw form, the dataset contains 42382 images from 61 scenes. The length of the continuous frame window is set to 3.

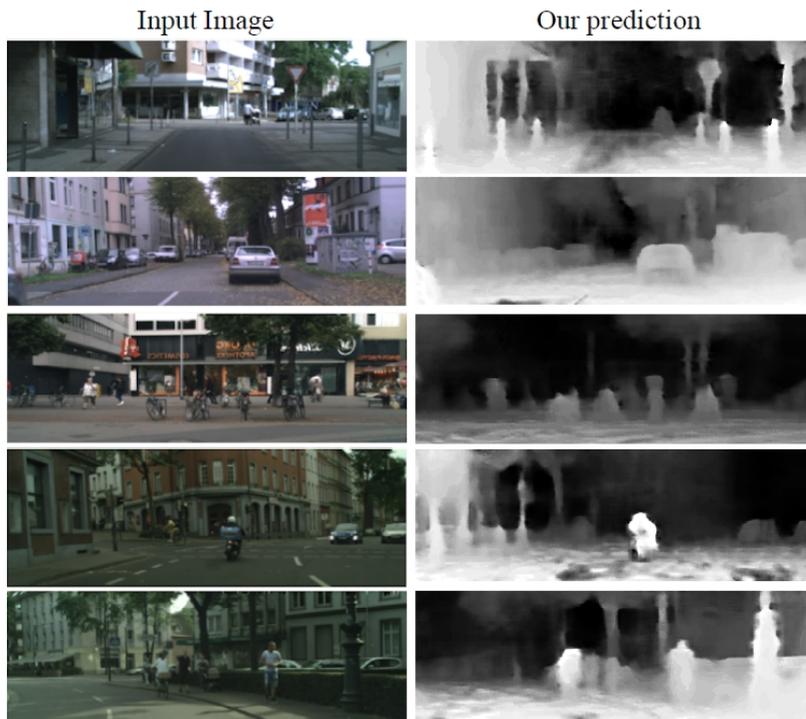


Fig. 6: The sample prediction on Cityscapes dataset using our approach trained on Cityscapes only.

To the best of our knowledge, among the methods which learn single-view depth estimation from monocular videos using unsupervised mechanism, state-of-the-art performance is achieved in Zhou [34]. We also make comparison with methods using supervised mechanism (depth ground-truth with depth supervision or calibrated stereo images with pose supervision) for training. Our method uses a scale factor to define the predicted depth so in the test phase, we multiply the predicted depth maps with a scalar which matches the median with the ground-truth data. Fig. 6 illustrates the predictions of our approach training on Cityscapes dataset [6]. We also make comparison with Godard [15] by taking the same training strategy which first pre-train the system on the Cityscapes dataset and then fine-tune on the KITTI dataset.

Table 1: Single-view depth results on the KITTI dataset and Cityscapes dataset.

Method	Dataset	Supervision		Error metric				Accuracy metric		
		Depth	Pose	Abs Rel	Sq Rel	RMSE	RMSE	$\log \delta < 1.25$	$\delta < 1.25^2$	$\delta < 1.25^3$
Train set mean	K	✓		0.403	5.530	8.709	0.403	0.593	0.776	0.878
Eigen <i>et al.</i> Coarse	K	✓		0.214	1.605	6.563	0.292	0.673	0.884	0.957
Eigen <i>et al.</i> Fine	K	✓		0.203	1.548	6.307	0.282	0.702	0.890	0.958
Liu <i>et al.</i>	K	✓		0.202	1.614	6.523	0.275	0.678	0.895	0.965
Godard <i>et al.</i>	K		✓	0.148	1.344	5.927	0.247	0.803	0.922	0.964
Godard <i>et al.</i>	CS+K		✓	0.124	1.076	5.311	0.219	0.847	0.942	0.973
Zhou <i>et al.</i>	K			0.208	1.768	6.856	0.283	0.678	0.885	0.957
Zhou <i>et al.</i>	CS+K			0.198	1.836	6.565	0.275	0.718	0.901	0.960
Ours	K			0.180	1.510	6.349	0.256	0.741	0.906	0.966
Ours	CS			0.236	2.476	7.249	0.307	0.645	0.861	0.946
Ours	CS+K			0.170	1.429	6.082	0.245	0.786	0.927	0.969

KITTI We follow the experimental settings proposed by [10] with the test set of 697 images covering 29 scenes. Table 1 shows the performance comparison between our method and the baseline methods. Here, *K* stands for the KITTI dataset and *CS* stands for the *Cityscapes* dataset. Compared with the methods using depth supervision [10, 23], our method performs better. However, our unsupervised method performs a little worse than the methods using pose supervision mechanism [13, 15]. [15] uses calibrated stereo images with left-right cycle consistency loss for training. In future work, we will apply the similar cycle consistency loss to our framework.

Compared with previous state-of-the-art method [34] using unsupervised mechanism, our method decreases the depth estimation error of nearly 13.5%. This validates that our learning framework can effectively take advantage of the knowledge gained from the camera pose in traditional SLAM algorithms. We further compare the depth images obtained by our method with the baseline methods. From Fig. 7, we can see that our results have no explicit difference with those of the supervised approaches. Furthermore, our method can even better represent the depth of the boundary information in some special scenes. Fig. 8 visualizes the testing results of our method using the strategy of [15] (first pre-train on the Cityscapes dataset and then fine-tune the model on the KITTI dataset).

In order to show the relationship between the length of the continuous frame window and the performance of our system, we set the length of the window to 5 and repeat the experiments above. As shown in Table 2, the performance of 5-frame version is better than that of the 3-frame version on all the metrics, which is consistent with our suppose. The benefit is attributed to the abundant transformation matrices between frames, which can further augment the supervised signals.

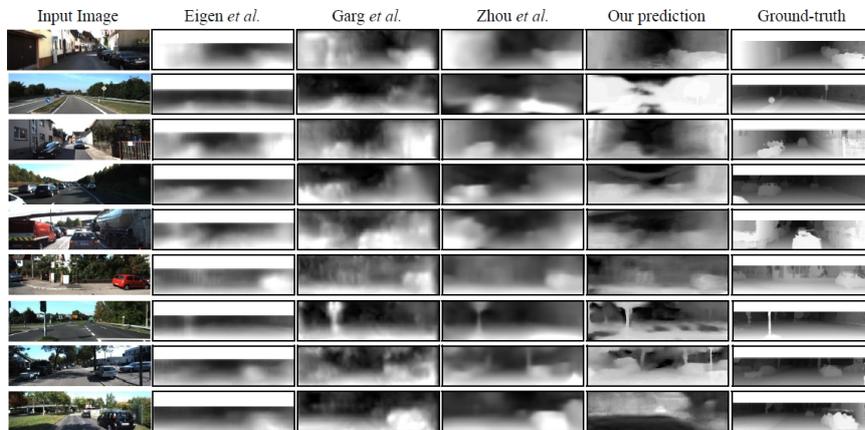


Fig. 7: Comparison of single-view depth estimation between Eigen *et al.* [10] (with ground-truth depth supervision), Garg *et al.* [34] (with ground-truth pose supervision), Zhou *et al.* [34] (unsupervised), and ours (unsupervised). The ground-truth depth map is interpolated from sparse measurements for visualization purpose.

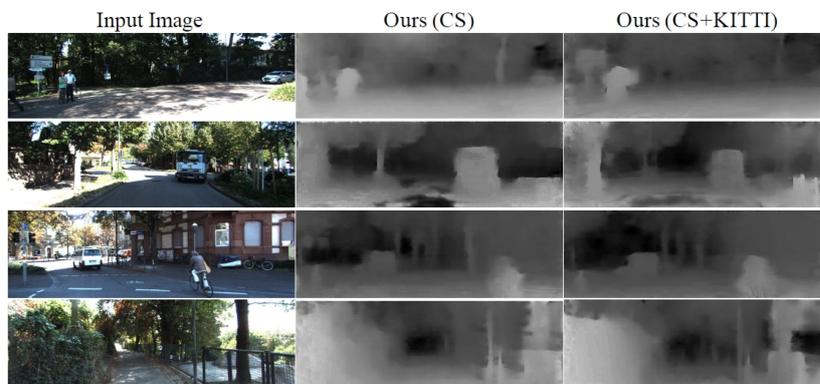


Fig. 8: Comparison of single-view depth predictions on the KITTI dataset by our initial Cityscapes model and the final model (pretrained on Cityscapes and then fine-tuned on KITTI). The Cityscapes model sometimes ignores structural mistakes (e.g. roadside billboards and lamp posts) likely due to the domain gap between the two datasets.

Table 2: Results of different continuous frame window length versions of our system.

Method	The length of the continuous frame window	Error metric			
		Abs Rel	Sq Rel	RMSE	RMSE log
Zhou <i>et al.</i>	3	0.208	1.768	6.856	0.283
Ours	3	0.180	1.510	6.349	0.256
Ours	5	0.176	1.455	5.940	0.248

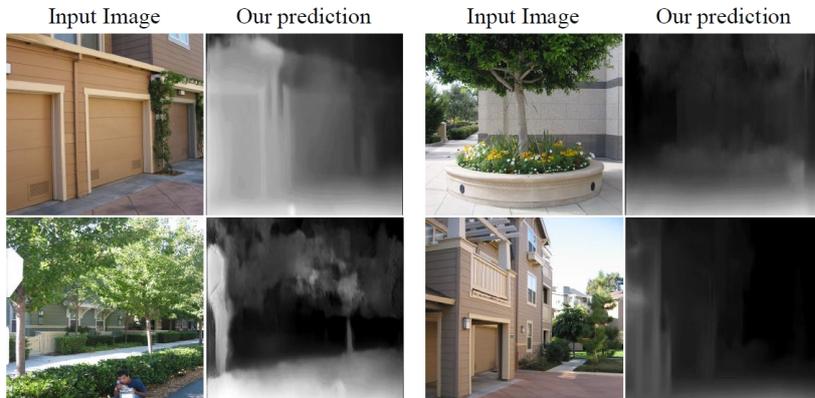


Fig. 9: Our sample predictions on the Make3D dataset. Note that our model is trained on KITTI+Cityscapes only, and directly tested on Make3D.

Make3D In order to evaluate the generalization ability and the adaptability of our proposed method on the strong lighting and weak texture environments, we choose Make3D [30] dataset for further experiments. In concrete, we pre-train our network on Cityscapes dataset and fine-tune on KITTI dataset. Then, we evaluate our model on Make3D dataset, which contains abundant strong lighting and weak texture samples and maintains an explicit difference between the other two datasets. As shown in Table 3, the results on Make3D dataset are similar to those of KITTI dataset. In concrete, compared with the methods using depth supervision mechanism, our method achieves a better performance than [10,23], but a little worse than [15], which indicates that our method maintains a satisfying generalization ability. Compared with the unsupervised method [34], our method still achieves a better performance. We further visualize the sample predictions of our method in Fig. 9 for a better illustration.

Table 3: Results on the Make3D dataset.

Method	Supervision		Error metric			
	Depth	Pose	Abs Rel	Sq Rel	RMSE	RMSE log
Train set mean	✓		0.876	13.98	12.27	0.307
Karsch <i>et al.</i>	✓		0.428	5.079	8.389	0.149
Liu <i>et al.</i>	✓		0.475	6.562	10.05	0.165
Laina <i>et al.</i>	✓		0.204	1.840	5.683	0.084
Godard <i>et al.</i>		✓	0.544	10.94	11.76	0.193
Zhou <i>et al.</i>			0.383	5.321	10.47	0.478
Ours			0.343	4.739	8.201	0.455

5.3 Pose Estimation

We choose ORB-SLAM [27] as the baseline method to illustrate the effectiveness of our proposed pose estimation network. We follow the experimental settings in [34] and use the official KITTI odometry split method to guarantee a fair comparison. The odometry benchmark is composed of 11 driving sequences with ground-truth odometry. We choose the first 9 driving sequences (00-08) for training and the last 2 driving sequences (09-10) for testing. The ground-truth odometry is used to evaluate our ego-motion estimation performance and the length of the frame window is set to 5. We compare our ego-motion estimation with two variants of monocular ORB-SLAM algorithm. The first one is ORB-SLAM (full) which uses all frames of the driving sequence to recover odometry. The second one is ORB-SLAM (short) which is lack of the loop closure and re-localization modules and maintains the same input setting as our system (5-frame snippets). Besides, the unsupervised method [34] is also selected as the baseline.

Notably, due to the reason that different methods have different scales, we optimize the scaling factor for the predictions made by each method to make all

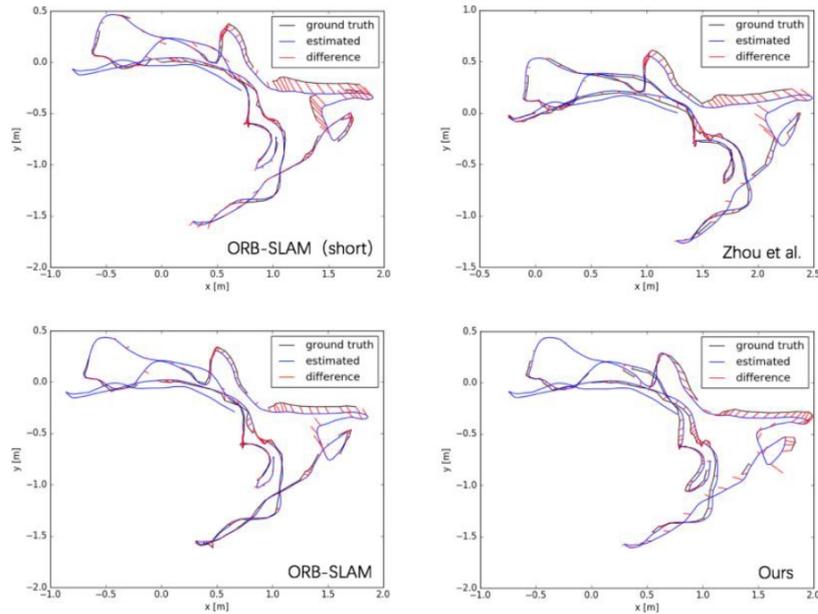


Fig. 10: Pose estimation trajectories comparison.

the scaling factors consistent with the ground-truth. The Absolute Trajectory Error (ATE) of the ground-truth and the estimated trajectory is chosen for evaluation. All the methods are computed on 5-frame snippets except for ORB-SLAM (full). For the ORB-SLAM (full) method, we break down the trajectory of the full sequence into 5-frame snippets by adjusting the reference coordinate frame to the central frame of each snippet.

Table 4: Absolute Trajectory Error (ATE) on the KITTI odometry (lower is better).

Method	Seq.09	Seq.10
ORB-SLAM(full)	0.014±0.008	0.012±0.011
ORB-SLAM(short)	0.064 ± 0.141	0.064 ± 0.130
Mean Odom	0.032 ± 0.026	0.028 ± 0.023
ORB-SLAM(short)	0.064 ± 0.141	0.064 ± 0.130
Zhou <i>et al.</i>	0.021 ± 0.017	0.021 ± 0.017
Ours	0.017±0.008	0.015±0.017

As shown in Table 4, our approach performs comparably with the ORB-SLAM (full) method, which utilizes the whole image sequences for loop closure and re-localization to improve the pose estimation accuracy. The ATE value

of our approach is about a quarter of that acquired by ORB-SLAM (short). In future, it would be interesting to use our learned ego-motion instead of the local estimation modules in monocular SLAM systems. Meanwhile, our pose estimation outperforms the previous state-of-the-art unsupervised method [34], which is conceptually similar to ours. Fig. 10 illustrates the pose estimation trajectories comparison between our method and the baseline methods.

6 Evaluation in the Realistic Settings

In this section, we first introduce the details on our platform for implementing the unsupervised learning-based depth estimation aided ORB-SLAM system using cloud robotic infrastructure. Then, we present the testing results of our system in various scenes.

6.1 Experimental Platform

Our system is composed of two parts: the robot and the server from the perspective of hardware.

The robot is deployed with multiple sensors, i.e., camera, ultrasonic radar and sound sensor. In our system, the robot interacts with the real-world settings by collecting the images from the RGB-D sensor and transmitting them to the server. The server mainly deals with the data saving and data processing tasks. In concrete, for our system, the depth estimation network and the pose estimation networks are both deployed in the server and the server will also process the computing task and the simultaneous localization and mapping task. Fig. 11 illustrates the sparse point cloud image acquired by our system and the green lines stand for the pose trajectory predicted by our system.



Fig. 11: The sparse point cloud image acquired by our system on a indoor desk scene.

The depth estimation network and the pose estimation network deployed in the server are both well trained on KITTI, Cityscapes and Make3D datasets. We choose TUM RGBD dataset to evaluate the performance and the server uses NVIDIA GeForce 1080P GPU for the processing tasks.

6.2 Experimental Results in Strong Lighting and Weak Texture Environments

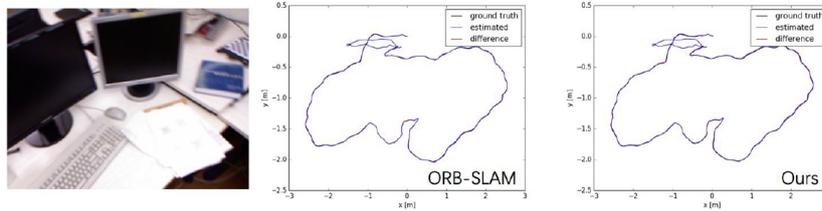


Fig. 12: Illustration of pose estimation in a normal scene (desk).

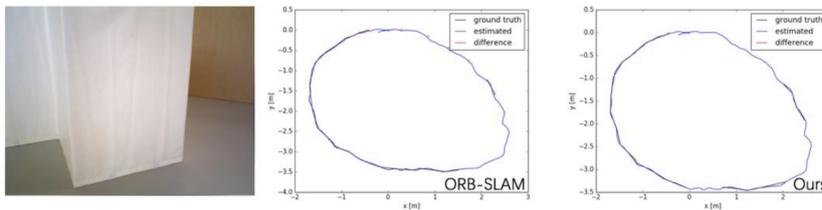


Fig. 13: Illustration of pose estimation in a weak texture scene Fr3/nst.

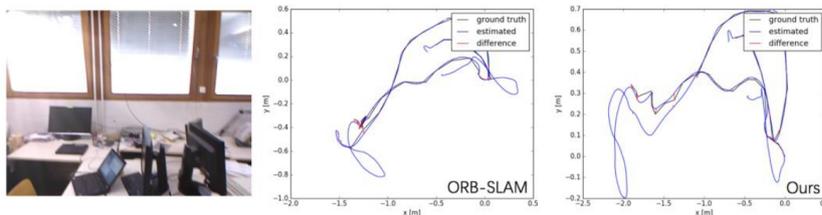


Fig. 14: Illustration of pose estimation in a strong lighting scene Fr3/stf.

The standard for evaluating the performance of our system and the ORB-SLAM system is ATE, which is widely used for testing in the SLAM area. We compare the performances of the two methods in various scenes of TUM RGBD dataset. The results are reported in Table 5, the first 5 scenes are in the normal environments and the last 3 scenes are in the strong lighting or weak texture environments. We can see that in normal scenes, our system performs comparably to the traditional ORB-SLAM system. In the strong lighting and weak

texture environments, our system achieves a better performance than that of the traditional ORB-SLAM system. For a better illustration, we report the pose estimation trajectories in normal, weak texture and strong lighting environments respectively (shown in Fig. 12, Fig. 13, and Fig. 14).

Table 5: ATE Comparison of our method and traditional ORB-SLAM system in different scenes.

Scene	ORB-SLAM	Ours
<i>Fr1/desk</i>	0.018490	0.017181
<i>Fr1/desk2</i>	0.021034	0.021564
<i>Fr2/desk</i>	0.011296	0.010978
<i>Fr1/room</i>	0.061536	0.062283
<i>Fr2/office</i>	0.010999	0.010901
<i>Fr2/stf</i>	0.013178	0.012105
<i>Fr2/stn</i>	0.012706	0.012294
<i>Fr2/nst</i>	0.023507	0.022203

6.3 Testing on the Speed for Initialization

Due to the reason that our method optimizes the initialization process of the traditional ORB-SLAM system, we design a series of experiments to test the initialization speed of the two systems by recording the number of images used for initialization.

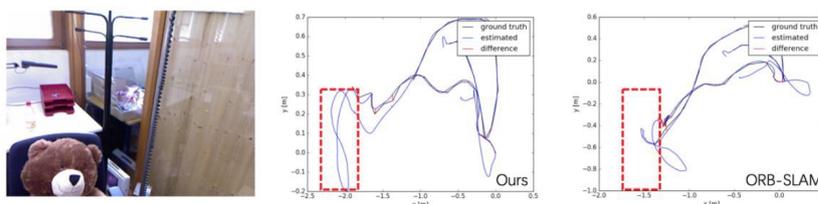


Fig. 15: Initialization speed comparison in a strong lighting scene.

Fig. 15 and Fig. 16 report the testing results in a strong lighting and weak texture environments of the TUM RGBD dataset respectively. The lines in the red square stand for the estimated trajectories in the initialization process. We can see that the estimated trajectory of our method is longer than that of the traditional ORB-SLAM system in both the two situations, which indicates that our method can complete the initialization process and start to build the map in a

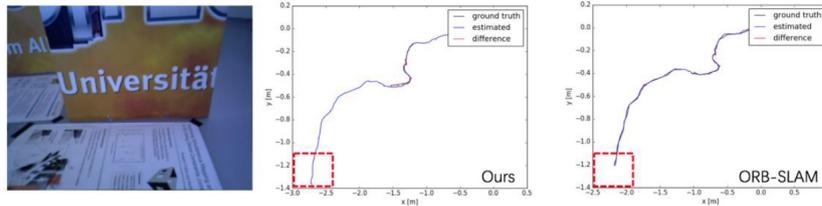


Fig. 16: Initialization speed comparison in a weak texture scene.

faster speed. Table 6 reports the comparison on the number of images used in the initialization process. For the first 5 normal scenes, our system uses comparable number of images to the traditional ORB-SLAM system. However, towards the last 3 scenes in the strong lighting and weak texture environments, the number of images used by our system is explicitly less than that of the traditional ORB-SLAM system, which indicates the effectiveness of our approach.

Table 6: Comparison on the number of images used in the initialization process.

Scene	ORB-SLAM	Ours
<i>Fr1/desk</i>	4	2
<i>Fr1/desk2</i>	5	2
<i>Fr2/desk</i>	4	4
<i>Fr1/room</i>	7	2
<i>Fr2/office</i>	6	5
<i>Fr2/stf</i>	38	10
<i>Fr2/stn</i>	145	84

7 Conclusion

We present an unsupervised learning framework for single-view depth and ego-motion estimation. The proposed method exploits the pose estimation method to enhance the supervised signal and add training constraints for the task of monocular depth and camera motion estimation. The system is trained on unlabeled videos and performs comparably to approaches that require ground-truth depth or pose for training. Furthermore, our method outperforms the previous state-of-the-art unsupervised learning method by 13.5% on KITTI dataset. Finally, we successfully exploit our unsupervised learning framework to assist the traditional ORB-SLAM system when the initialization module of ORB-SLAM method could not match enough features. Experiments have shown that our method can significantly accelerate the initialization process of traditional ORB-

SLAM system and effectively improve the accuracy on environmental mapping in strong lighting and weak texture scenes.

Acknowledgments. This work was supported by the National Natural Science Foundation of China (grant numbers 61751208, 61502510, and 61773390), the Outstanding Natural Science Foundation of Hunan Province (grant number 2017JJ1001).

References

1. Abadi, M., Barham, P., Chen, J., Chen, Z., Davis, A., Dean, J., Devin, M., Ghemawat, S., Irving, G., Isard, M., et al.: Tensorflow: a system for large-scale machine learning. In: OSDI. vol. 16, pp. 265–283 (2016)
2. Agarwal, S., Snavely, N., Seitz, S.M., Szeliski, R.: Bundle adjustment in the large. In: European conference on computer vision. pp. 29–42. Springer (2010)
3. Canutescu, A.A., Dunbrack Jr, R.L.: Cyclic coordinate descent: A robotics algorithm for protein loop closure. *Protein science* 12(5), 963–972 (2003)
4. Cireşan, D., Giusti, A., Gambardella, L.M., Schmidhuber, J.: Deep neural networks segment neuronal membranes in electron microscopy images. In: Advances in neural information processing systems. pp. 2843–2851 (2012)
5. Cireşan, D., Meier, U., Schmidhuber, J.: Multi-column deep neural networks for image classification. arXiv preprint arXiv:1202.2745 (2012)
6. Cordts, M., Omran, M., Ramos, S., Rehfeld, T., Enzweiler, M., Benenson, R., Franke, U., Roth, S., Schiele, B.: The cityscapes dataset for semantic urban scene understanding. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 3213–3223 (2016)
7. Delage, E., Lee, H., Ng, A.Y.: A dynamic bayesian network model for autonomous 3d reconstruction from a single indoor image. In: Computer Vision and Pattern Recognition, 2006 IEEE Computer Society Conference on. vol. 2, pp. 2418–2428. IEEE (2006)
8. Dosovitskiy, A., Fischer, P., Ilg, E., Hausser, P., Hazirbas, C., Golkov, V., Van Der Smagt, P., Cremers, D., Brox, T.: FlowNet: Learning optical flow with convolutional networks. In: Proceedings of the IEEE International Conference on Computer Vision. pp. 2758–2766 (2015)
9. Eigen, D., Fergus, R.: Predicting depth, surface normals and semantic labels with a common multi-scale convolutional architecture. In: Proceedings of the IEEE International Conference on Computer Vision. pp. 2650–2658 (2015)
10. Eigen, D., Puhrsch, C., Fergus, R.: Depth map prediction from a single image using a multi-scale deep network. In: Advances in neural information processing systems. pp. 2366–2374 (2014)
11. Engel, J., Schöps, T., Cremers, D.: Lsd-slam: Large-scale direct monocular slam. In: European Conference on Computer Vision. pp. 834–849. Springer (2014)
12. Gao, X., Zhang, T.: Unsupervised learning to detect loops using deep neural networks for visual slam system. *Autonomous robots* 41(1), 1–18 (2017)
13. Garg, R., BG, V.K., Carneiro, G., Reid, I.: Unsupervised cnn for single view depth estimation: Geometry to the rescue. In: European Conference on Computer Vision. pp. 740–756. Springer (2016)

14. Geiger, A., Lenz, P., Urtasun, R.: Are we ready for autonomous driving? the kitti vision benchmark suite. In: *Computer Vision and Pattern Recognition (CVPR), 2012 IEEE Conference on*. pp. 3354–3361. IEEE (2012)
15. Godard, C., Mac Aodha, O., Brostow, G.J.: Unsupervised monocular depth estimation with left-right consistency. In: *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. pp. 6602–6611. IEEE (2017)
16. Hoiem, D., Efros, A.A., Hebert, M.: Automatic photo pop-up. In: *ACM transactions on graphics (TOG)*. vol. 24, pp. 577–584. ACM (2005)
17. Jason, J.Y., Harley, A.W., Derpanis, K.G.: Back to basics: Unsupervised learning of optical flow via brightness constancy and motion smoothness. In: *European Conference on Computer Vision*. pp. 3–10. Springer (2016)
18. Karsch, K., Liu, C., Kang, S.B.: Depth transfer: Depth extraction from videos using nonparametric sampling. In: *Dense Image Correspondences for Computer Vision*, pp. 173–205. Springer (2016)
19. Kingma, D.P., Ba, J.: Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980* (2014)
20. Krizhevsky, A., Sutskever, I., Hinton, G.E.: Imagenet classification with deep convolutional neural networks. In: *Advances in neural information processing systems*. pp. 1097–1105 (2012)
21. Laina, I., Ruppel, C., Belagiannis, V., Tombari, F., Navab, N.: Deeper depth prediction with fully convolutional residual networks. In: *3D Vision (3DV), 2016 Fourth International Conference on*. pp. 239–248. IEEE (2016)
22. Li, R., Wang, S., Long, Z., Gu, D.: Undeepvo: Monocular visual odometry through unsupervised deep learning. In: *2018 IEEE International Conference on Robotics and Automation (ICRA)*. pp. 7286–7291. IEEE (2018)
23. Liu, F., Shen, C., Lin, G., Reid, I.D.: Learning depth from single monocular images using deep convolutional neural fields. *IEEE Trans. Pattern Anal. Mach. Intell.* 38(10), 2024–2039 (2016)
24. Liu, M., Salzmann, M., He, X.: Discrete-continuous depth estimation from a single image. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. pp. 716–723 (2014)
25. Lourakis, M.I.: A brief description of the levenberg-marquardt algorithm implemented by levmar. *Foundation of Research and Technology* 4(1), 1–6 (2005)
26. Lowe, D.G.: Distinctive image features from scale-invariant keypoints. *International journal of computer vision* 60(2), 91–110 (2004)
27. Mur-Artal, R., Montiel, J.M.M., Tardos, J.D.: Orb-slam: a versatile and accurate monocular slam system. *IEEE Transactions on Robotics* 31(5), 1147–1163 (2015)
28. Patraucean, V., Handa, A., Cipolla, R.: Spatio-temporal video autoencoder with differentiable memory. *arXiv preprint arXiv:1511.06309* (2015)
29. Saxena, A., Chung, S.H., Ng, A.Y.: 3-d depth reconstruction from a single still image. *International journal of computer vision* 76(1), 53–69 (2008)
30. Saxena, A., Sun, M., Ng, A.Y.: Make3d: Learning 3d scene structure from a single still image. *IEEE transactions on pattern analysis and machine intelligence* 31(5), 824–840 (2009)
31. Sun, D., Roth, S., Black, M.J.: Secrets of optical flow estimation and their principles. In: *Computer Vision and Pattern Recognition (CVPR), 2010 IEEE Conference on*. pp. 2432–2439. IEEE (2010)
32. Tateno, K., Tombari, F., Laina, I., Navab, N.: Cnn-slam: Real-time dense monocular slam with learned depth prediction. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. vol. 2 (2017)

33. Triggs, B., McLauchlan, P.F., Hartley, R.I., Fitzgibbon, A.W.: Bundle adjustment a modern synthesis. In: International workshop on vision algorithms. pp. 298–372. Springer (1999)
34. Zhou, T., Brown, M., Snavely, N., Lowe, D.G.: Unsupervised learning of depth and ego-motion from video. In: CVPR. vol. 2, p. 7 (2017)