

# Unsupervised Monocular Depth Estimation in Highly Complex Environments

Chaoqiang Zhao, Yang Tang, *Senior Member, IEEE*, Qiyu Sun

**Abstract**—With the development of computational intelligence algorithms, unsupervised monocular depth and pose estimation framework, which is driven by warped photometric consistency, has shown great performance in the day-time scenario. While in some challenging environments, like night and rainy night, the essential photometric consistency hypothesis is untenable because of the complex lighting and reflection, so that the above unsupervised framework cannot be directly applied to these complex scenarios. In this paper, we investigate the problem of unsupervised monocular depth estimation in highly complex scenarios and address this challenging problem by adopting an image transfer-based domain adaptation framework. We adapt the depth model trained on day-time scenarios to be applicable to night-time scenarios, and constraints on both feature space and output space promote the framework to learn the key features for depth decoding. Meanwhile, we further tackle the effects of unstable image transfer quality on domain adaptation, and an image adaptation approach is proposed to evaluate the quality of transferred images and re-weight the corresponding losses, so as to improve the performance of the adapted depth model. Extensive experiments show the effectiveness of the proposed unsupervised framework in estimating the dense depth map from highly complex images.

**Index Terms**—Unsupervised estimation, domain adaptation, monocular depth estimation, night, rainy night.

## I. INTRODUCTION

Depth is one of the most important information for autonomous systems in perceiving their surroundings and their own states [1], [2]. Therefore, the accurate estimation of the depth information from monocular images has become a hot topic in recent years and been used to improve other perception tasks [3]. Structure from motion and stereo matching are two main ways to recover the depth information based on the geometric relationship between images, and these methods are widely used in traditional SLAM methods to map the environments [4], [5]. With the development of computational

©20XX IEEE. Personal use of this material is permitted. Permission from IEEE must be obtained for all other uses, in any current or future media, including reprinting/republishing this material for advertising or promotional purposes, creating new collective works, for resale or redistribution to servers or lists, or reuse of any copyrighted component of this work in other works.

This work was supported in part by National Natural Science Foundation of China (Basic Science Center Program: 61988101), in part by National Natural Science Fund for Distinguished Young Scholars (61725301), in part by the Programme of Introducing Talents of Discipline to Universities (the 111 Project) under Grant B17017, in part by the Program of Shanghai Academic Research Leader (20XD1401300), in part by Innovation Research Funding of China National Petroleum Corporation (2021D002-0902), and in part by Shanghai AI Lab. (*Corresponding author: Yang Tang.*)

C. Zhao, Y. Tang, and Q. Sun are with the Key Laboratory of Smart Manufacturing in Energy Chemical Process, Ministry of Education, East China University of Science and Technology, Shanghai, 200237, China (e-mail: zhaocqilc@gmail.com, yangtang@ecust.edu.cn, qysun291@163.com).

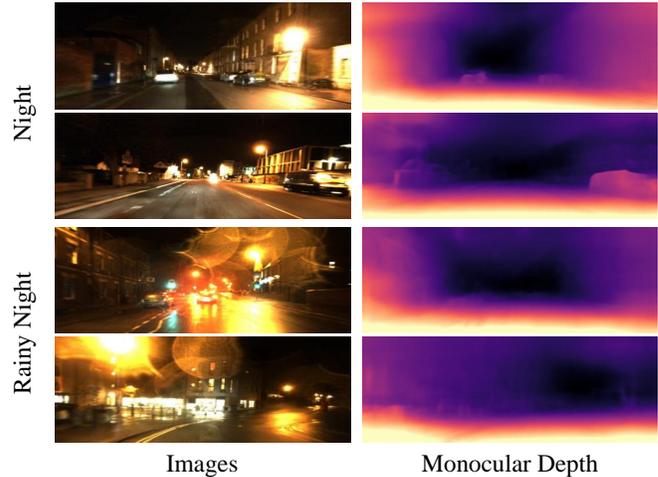


Fig. 1. The monocular depth predictions of the proposed method at night and rainy night. Note that the encoders of the above three depth models are adapted from the same basic encoder, and their decoders are the same.

intelligence algorithms [6]–[8], deep learning algorithms have shown great performance in many tasks [9], [10], and using deep neural networks to estimate the pixel-level dense depth from only a single image is becoming possible and has attracted much attention [11]. Recently, different kinds of deep learning-based monocular depth estimation frameworks have been proposed, including supervised methods, semi-supervised and unsupervised methods [12]–[14]. Because of the costly ground truth, geometric constraints are gradually replacing ground truth for the training of depth networks, and the unsupervised framework has become a promising direction for monocular depth estimation [15].

In the unsupervised framework of monocular depth estimation [16], [17], the geometric constraints between adjacent images are considered to supervise the network training. Therefore, only monocular image sequences and camera parameters are needed during the training process. This unsupervised framework is mainly composed of two deep neural networks, including a depth network to regress the dense depth from single images and a pose network to estimate the pose between two frames. Based on the estimated depth map and pose, the geometric relationship between images is built on the projection function. The mainly supervised signal is calculated from the photometric error of corresponding pixels between adjacent images by using view reconstruction [16], [18].

However, since the unsupervised signal is built on projection consistency, the above unsupervised framework [16],

[18] suffers from two big limitations, *the static scenario hypothesis* and *the photometric consistency hypothesis*. For *the static scenario hypothesis*, since the pixels on moving objects do not satisfy the projection function of camera ego-motion, which leads to incorrect calculation of the loss during training, thereby affecting the accuracy of the depth network [19]. Incorporating semantic information into the unsupervised framework is an effective way to recognize moving objects and eliminate their influence, and relevant results have emerged [19], [20]. For *the photometric consistency hypothesis*, since the training of the unsupervised framework relies heavily on the photometric error, the photometric consistency of the same pixel on different images is crucial to the overall framework [21]. Since all objects are illuminated by the same light source (sun) during the day, the photometric consistency assumption is basically valid, just like that in traditional direct visual odometer methods [22], [23]. Therefore, almost all of the current unsupervised monocular depth estimation methods are trained and tested on day-time images. When the environment changes into other highly complex conditions, like *night* and especially *rainy night*, static objects with non-Lambertian or high reflectance surfaces in the night/rainy night with dynamic lighting conditions violate the photometric consistency between frames. Besides, objects with low luminance in the night lack reliable cues for providing accurate correspondences. The problem of estimating monocular depth in such varying environments is challenging but practical and important. Meanwhile, the perception of changing and complex environments is crucial for autonomous systems [24], [25], like robots and autonomous driving cars, and this problem has only received some initial attention [21], [26].

Because of the limitations of the unsupervised framework in complex scenarios, we tackle this challenging problem by using an image transfer based domain adaptation framework. Instead of training the unsupervised framework on the images from complex environments, we adapt the model trained by day-time images to other complex environments, thereby circumventing the photometric inconsistency in complex environments and achieving satisfactory depth estimation in an unsupervised way in the highly complex environments. We only adapt the encoder of the depth network by following [26]. An additional encoder is designed to encode the images of complex environments, and after adaptation, the encoders for complex environments share the same feature space with the day-time encoder. Besides, the adapted encoders for different scenarios share the same decoder for monocular depth estimation in different complex scenarios, which is meaningful for practical applications. Therefore, this method not only reduces computational complexity but also facilitates practical applications: switching different encoders for adapting various environments. Different from [26] using adversarial domain adaptation in *feature space* for night-time depth estimation, we propose to adopt the image transfer-based domain adaptation framework and constrain the training from both *feature space* and *output space*, which is more stable and accurate [27], [28]. Besides, for image transfer based domain adaptation framework, the generated paired images are used to get pseudo labels from the known models and supervise the adaptation

process. Poor generated images result in wrong pseudo labels and affect the training process, which is not considered in previous works. Therefore, in this paper, we consider the errors introduced by the unstable image transfer, and an image quality adaptation approach is proposed to evaluate the quality of the transferred images and reduce their effects. The proposed image transfer-based domain feature adaptation (ITDFA) framework not only can be used for night-time depth estimation, but also shows outstanding performance on more challenging rainy night-time images, as shown in Fig. 1.

In summary, in this paper, we analyze and tackle the unsupervised monocular depth estimation problem in three typical and challenging scenarios (night, rainy night), including proposing the ITDFA framework, constructing novel training/testing sets on different scenarios, digging into the continuous adaptation ability of the ITDFA, and exploring the influence of the image transfer model on ITDFA. Our main contributions are as follows:

- This paper analyzes the major reason for the limited performance and application of the current unsupervised monocular depth estimation framework, and we tackle the problem of unsupervised monocular depth estimation in highly complex environments by using domain adaptation.
- Image transfer-based unsupervised domain adaptation is applied to estimate monocular depth from challenging scenarios, like night and rainy night. To reduce the effects of the unstable image transfer quality, we propose an image quality adaptation approach to evaluate the quality of the transferred images and re-weight the corresponding losses.
- Extensive experiments and results on the RobotCar dataset [29] show the effectiveness of our proposed method in highly complex environments.

## II. RELATED WORK

In this section, we introduce the popular unsupervised monocular depth estimation framework [16], which is trained on monocular sequences. Firstly, many recent research results for improving this unsupervised framework are briefly reviewed, from the perspectives of occlusions, static scenario hypothesis and photometric consistency hypothesis. Then, we review the framework combined with domain adaptation, in which depth models are trained on synthetic datasets and then adapted to real-world scenarios through domain adaptation.

**Unsupervised framework.** To circumvent the need for costly ground truth, Zhou *et al.* [16] propose to use geometric constraints between frames instead of ground truth to train a depth network. Their framework contains a depth network for monocular depth estimation and a pose network for inter-frame pose estimation. Then, based on the projection function established by the estimated pose and depth, the view reconstruction is designed to warp and construct the target frame from its adjacent frame. The photometric error between the warped and real target images is used to supervise the training process, so that the depth and pose networks are trained in an unsupervised manner. To improve the accuracy of depth estimation, several

novel loss functions and network frameworks are proposed, which are well reviewed in [11].

**Photometric inconsistency.** Photometric inconsistency is one of the major reasons for the limited performance and application of the unsupervised framework. For some highly complex environments, like *night-time environments* and more challenging *rainy night-time environments*, due to the complex lighting conditions, e.g., street lamps, car lights, and especially the reflection of light from the road caused by rain, the essential photometric consistency hypothesis is untenable, so the unsupervised framework shows unsatisfied accuracy and robustness. The unsupervised monocular depth estimation in such complex scenarios is a largely under explored domain, and only a few methods for night-time monocular depth estimation have been proposed most recently [21], [26]. To overcome the photometric consistency of images, Spencer et al. [21] propose to use dense feature representation of images for unsupervised training. Since the corresponding features between different images are consistent and unaffected by light, this unsupervised framework can well adapt to night-time scenarios. Nevertheless, since the whole framework is unsupervised, their framework still needs the help of photometric error during training. Different from [21], Vankadari et al. [26] regard this challenge as a domain adaptation problem. The depth and pose networks are trained on day-time scenarios by following [18] at first. Then, an additional encoder is designed to encode the night-time images, and an adversarial domain feature adaptation method is used to adapt the features encoded by the day-time encoder and night-time encoder. Since the output of the encoder are multi-scale high-dimensional feature maps, they design multiple discriminators to constrain each scale feature map. However, adjusting the adversarial framework consisting of multiple discriminators and a generator is extremely complex, which influences its stability for applying to other scenarios [30], [31]. Moreover, although adversarial learning helps to reduce the distance between the distributions of day and night feature spaces, the key features for depth decoding are not valued because the decoder is not involved in their domain adaptation process.

**Domain adaptation.** Due to the domain shift, like differences in the background, lighting, weather, and so on of images between different datasets/domains, the performance of the trained model may degrade significantly when it was applied to other datasets [32]–[34]. Therefore, a domain adaptation framework is proposed to transfer the model from one domain to another for the same task [35], and most recently, adaptation between multiple domains has received a lot of attention [36], [37].

In monocular depth estimation, domain adaptation algorithms are mainly applied to adapt the model trained on synthetic datasets to real-world datasets [38]–[40]. Compared with the ground truth obtained by different sensors in the real-world, the ground truth obtained from virtual environments is cheaper and easier. The depth model is trained on synthetic and real-world datasets and supervised by the ground truth of synthetic datasets. Since supervised training can get more cues than unsupervised training, this method achieves better accuracy on monocular depth estimation than unsupervised

methods, and it provides a new way to circumvent the need for costly ground truth at the same time. Most recently, the LAB-based images transfer approach is proposed to transfer images between domains for domain adaptation [41]. While for the night-time scenario, the above frameworks cannot work because it is difficult to generate synthetic night-times images that can capture all the vagaries of real-world night conditions [26], let alone rainy night-time scenarios and even more complex scenarios. Therefore, to tackle the unsupervised monocular depth estimation in highly complex environments, we use unsupervised domain adaptation to transfer the model trained on day-time images to work for night-time and rainy night-time images.

To increase the applicability and reduce the computational complexity, our model only adapts the encoders during training by following [26], but our model ITDFA does not need to consider the stability of adversarial learning. Meanwhile, Vankadari et al. [26] only consider the adversarial constraint on feature space, while this paper constrains the training of the encoder from both feature space and output space. The constraints on output space help the encoder to focus on learning the key features of depth decoding. Similar to the concurrent work [42], the CycleGAN-based image transfer method is introduced to generate paired images between day and night because it can well mimic nighttime lighting conditions, and the transferred images are used to generate pseudo labels and then supervise the training process. However, as shown in Fig. 3 (a), poor transferred images will generate wrong pseudo labels and supervised signals, which will affect the performance of the network. We want to tackle the above problems by evaluating the transferred images and thus reducing their effect on training. However, since the transfer models are trained in an unsupervised manner, and there are no real paired images between domains, it is difficult to demonstrate the quality of transferred images. After lots of experiments, we find a new method to reflect the quality of the transferred images, which is described in the next section in detail. Hence, our proposed image transferred domain adaptation model achieves a better performance than [26], [42] in the night-time scenario. Moreover, they [26], [42] only address the unsupervised monocular depth estimation on night-time images, while our framework can also do well in more challenging rainy night-time scenario.

**Discussion:** *Why not adopt a framework that directly combines image style transfer with monocular depth estimation:* the images from complex environments are first transferred to normal day-time style and then use the day-time depth model to estimate the depth. After testing, we find this is a possible way to solve the problem, as shown in lines 1-2 of Fig. 3 (a). Nevertheless, poor real-time performance will limit the application of this framework because of the two-step process. Besides, since the accuracy of depth estimation relies heavily on the quality of transferred images, this approach has great instability in practical applications, as shown in lines 3-4 of Fig. 3 (a). The proposed ITDFA framework can get a new model for the new scene and predict the depth in an end-to-end manner. Moreover, since only the encoder is trained, our method can get multiple encoders for multiple scenes, and

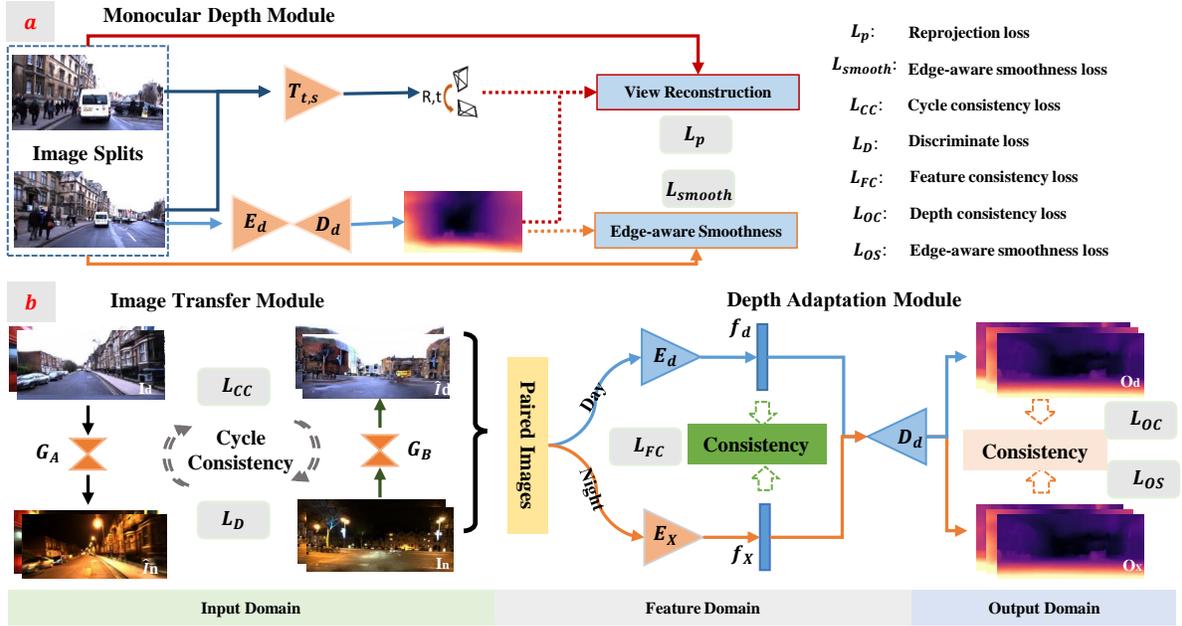


Fig. 2. The framework of ITDFA for unsupervised monocular depth estimation in highly complex environments. In (a), the unsupervised monocular depth estimation framework follows monodepth2 [18]. In (b), our ITDFA framework contains two modules, image transfer module and depth adaptation module.  $I_d$  and  $I_X$  refer to the image from day-time scenario and one of the complex scenarios, and  $O_d$  and  $O_X$  stand for their corresponding depth maps. The encoder  $E_d$  and decoder  $D_d$  trained by day-time images do not update their weights during training.

these encoders share the same decoder, which is very practical.

### III. METHODS

In this section, we will introduce the overall ITDFA framework, loss functions for training as well as the image quality adaptation strategy proposed in this paper.

#### A. Unsupervised depth estimation part

We use the famous unsupervised framework, monodepth2 [18], to acquire the trained depth network in an unsupervised manner, which is shown in Fig. 2-a. Monodepth2 [18] has been widely used as the basic unsupervised framework in this field because of its high practicability and accuracy [20], [21], [43]. For the day time scenario, the unsupervised monocular depth estimation is formulated as the minimization of the per-pixel minimum reprojection error:

$$L_p = \min \Psi(I_t, I_{s \rightarrow t}), \quad (1)$$

$$\Psi(I_a, I_b) = \frac{\alpha}{2} (1 - SSIM(I_a, I_b)) + (1 - \alpha) \|I_a - I_b\|_1, \quad (2)$$

and

$$I_{s \rightarrow t} = I_s \langle \text{proj}(O_t, T_{t \rightarrow s}, K) \rangle, \quad (3)$$

where  $\|\cdot\|_1$  refers to the L1 distance in pixel space, and  $\text{proj}$  stands for the 2D coordinate projection based on the predicted dense depth map  $O_t$  of the target image  $I_t$  and the related pose  $T_{t \rightarrow s}$  between target  $I_t$  and source images  $I_s$ . Meanwhile, the edge-aware smoothness loss is also used to improve the depth map  $O_t$ :

$$L_{smooth} = |\partial_x o_t^*| e^{\partial_x I_t} + |\partial_y o_t^*| e^{\partial_y I_t}, \quad (4)$$

where  $d_t^* = o_t / \hat{o}_t$  represents the mean-normalized inverse depth.

The monodepth2 framework is supervised by combining the per-pixel smoothness loss and masked photometric loss on the day-time scenario. The depth model consisting of an encoder  $E_d$  and a decoder  $D_d$  is used in the proposed ITDFA. The depth model learns a mapping from the day-time images  $I_d$  to the pixel-level depth maps  $O_d$ :

$$O_d = D_d(E_d(I_d)). \quad (5)$$

#### B. Image transfer part

Since the LAB-based image transfer approach [41] cannot well simulate the complex and heterogeneous lighting conditions at night, we utilize the CycleGAN-based framework [44] to transfer the images between scenarios. During training, the full objective is:

$$G_{d2X}, G_{X2d} = \arg \min_G \max_D L(G_{d2X}, G_{X2d}, D_d, D_X), \quad (6)$$

$$L(G_{d2X}, G_{X2d}, D_d, D_X) = L_{CC} + L_D, \quad (7)$$

where  $G_{d2X}$  tries to generate images  $G_{d2X}(I_d)$  that look similar to images from domain  $X$ , while  $D_X$  aims to distinguish between translated samples  $G_{d2X}(I_d)$  and real samples  $X$ .  $L_{CC}$  and  $L_D$  refer to the cycle consistency loss and adversarial loss:

$$L_{CC}(G_{d2X}, G_{X2d}) = \mathbb{E}_{I_d \sim p_{data}(I_d)} [\|G_{X2d}(G_{d2X}(I_d)) - I_d\|_1] + \mathbb{E}_{I_X \sim p_{data}(I_X)} [\|G_{d2X}(G_{X2d}(I_X)) - I_X\|_1], \quad (8)$$

and

$$L_D = L_D(G_{d2X}, D_X, I_d, I_X) + L_D(G_{X2d}, D_d, I_X, I_d), \quad (9)$$

and all of the above constraints used for training image transfer models ( $G_{d2X}$  and  $G_{X2d}$ ) are following Zhu *et al.* [44].

In this paper, for different complex environments, different image transfer models ( $G_{d2X}$  and  $G_{X2d}$ ) based on CycleGAN [44] are needed to transfer the images between day-time style  $d$  and different complex environmental styles  $X$ :

$$\begin{cases} \hat{I}_{d2X} = G_{d2X}(I_d), & \text{day to X} \\ \hat{I}_{X2d} = G_{X2d}(I_X), & \text{X to day} \\ \hat{I}_{d2X2d} = G_{X2d}(G_{d2X}(I_d)), & \text{cycle transfer} \\ \hat{I}_{X2d2X} = G_{d2X}(G_{X2d}(I_X)), & \text{cycle transfer} \end{cases}, \quad (10)$$

where  $X$  refers to night  $n$  or rainy night  $r$ . In addition, to verify the continuous transfer ability of the model obtained through domain adaptation, we also train an additional image transfer model between night-time style and rainy night-time style ( $G_{nr}$  and  $G_{rnr}$ ).

### C. ITDFA part

As shown in Fig. 2-c, an encoder  $E_X$  is designed to encode the features of images from highly complex scenarios to the same feature space as the features of day-time images encoded by the day-time encoder  $E_d$ . In ITDFA, the pre-trained day-time encoder  $E_d$  is used to encode the day-time images and obtain their corresponding feature maps  $f$ :

$$\begin{cases} f_d = E_d(I_d) \\ f_{X2d} = E_d(\hat{I}_{X2d}) \end{cases}, \quad (11)$$

where  $I_d$  refers to the real day-time images, and  $\hat{I}_{X2d}$  stands for the fake day-time images generated by CycleGAN model  $G_{X2d}$  from the highly complex scenario  $X$ . The encoder  $E_X$  for complex environments has the same network framework as the day-time encoder  $E_d$ , and it is used to encode the real and fake images of highly complex scenarios and obtain their feature maps  $f$ :

$$\begin{cases} f_X = E_X(I_X) \\ f_{d2X} = E_X(\hat{I}_{d2X}) \end{cases}, \quad (12)$$

where  $I_X$  refers to the real images from highly complex scenarios, and  $\hat{I}_{d2X}$  stands for the fake images transferred from day-time scenario  $d$ . The pre-trained decoder  $D_d$  is used to decode the features from  $E_d$  and  $E_X$  and obtain their corresponding depth maps  $O$ :

$$\begin{cases} O_d = D_d(f_d) \\ O_{X2d} = D_d(f_{X2d}) \\ O_X = D_d(f_X) \\ O_{d2X} = D_d(f_{d2X}) \end{cases}. \quad (13)$$

During training, the weights of  $E_d$  and  $D_d$  are fixed, and only  $E_X$  is updated.

During testing, the depth map can be estimated from the images of highly complex scenarios in one-step:

$$O_X = D_d(E_X(I_X)). \quad (14)$$

**Training losses:** The ITDFA framework is an unsupervised framework, and neither ground truth nor real paired images are used to train the depth model, CycleGAN model and domain adaptation model. As shown in Fig. 2, different constraints

are designed to supervise the training process, including the feature consistency loss  $L_{FC}$  on feature space, and the depth consistency loss  $L_{OC}$  as well as smoothness loss  $L_{OS}$  on output (depth) space. Therefore, the overall loss function for training the encoder  $E_X$  is formulated as:

$$L_{DA} = L_{FC} + \beta L_{OC} + \gamma L_{OS}. \quad (15)$$

**Feature consistency loss:** Based on the pre-trained CycleGAN model, we can get the paired images from day-time scenario and highly complex scenario, like  $I_d$  with  $I_{d2X}$  and  $I_X$  with  $I_{X2d}$ . Therefore, to promote the consistency of different encoders in feature space, we direct minimize the error of the feature maps, which are encoded by  $E_d$  and  $E_X$  from these image pairs:

$$L_{FC_{L1}} = L1(f_d, f_{d2X}) + L1(f_X, f_{X2d}). \quad (16)$$

Moreover, inspired by style transfer methods and related works [42], [45], [46], to enhance the consistency of correlations between features, the Gram Matrices  $\mathcal{G}$  between features are calculated to further improve the feature consistency:

$$\begin{aligned} L_{FC} &= L_{FC_{L1}} + \alpha L_{FC_{Gram}} \\ &= L1(f_d, f_{d2X}) + L1(f_X, f_{X2d}) \\ &\quad + L1(\mathcal{G}(f_d), \mathcal{G}(f_{d2X})) + L1(\mathcal{G}(f_X), \mathcal{G}(f_{X2d})) \end{aligned}. \quad (17)$$

$\alpha$ ,  $\beta$  and  $\gamma$  are the weights of each loss function for training.

**Depth consistency loss:** Although the above loss can help to promote the consistency of feature space between the two encoders, the ultimate goal is the depth map rather than the feature map, and the contribution of different features to the depth decoding is different. To further constrain the key features of feature maps for depth decoding, we design a depth consistency loss in output space:

$$L_{OC} = L1(O_d, O_{d2X}) + L1(O_X, O_{X2d}). \quad (18)$$

**Smoothness loss:** Moreover, to promote the smoothness of the generated depth map, we propose to utilize the edge-aware smoothness during training, which is widely used in previous unsupervised depth framework [18], [47], [48]:

$$L_{OS} = |\partial_x O_{d2X}^*| e^{\partial_x I_d} + |\partial_y O_{d2X}^*| e^{\partial_y I_d}, \quad (19)$$

where  $O_{d2X}^* = O_{d2X} / \bar{O}_{d2X}$  represents the mean-normalized inverse depth. Note that this loss is established between the real day-time image  $I_d$  and the depth map  $O_{d2X}$  of its corresponding transferred images  $\hat{I}_{d2X}$ .

### D. Image quality adaptation part

Inspired by the cycle consistency of CycleGAN during the training, we try to use the consistency between the raw image and cycle transferred image to report the performance. Nevertheless, this method cannot work well because of the overfitting, as shown in columns 1 and 3 of Fig. 3 (b). After a series of tests, we found that using the models saved from different epoch can demonstrate the quality of the transferred images. As shown in column 4 of Fig. 3 (b), for the good generated images, our method can generate good cycle images; while for the poor generated images, our method can

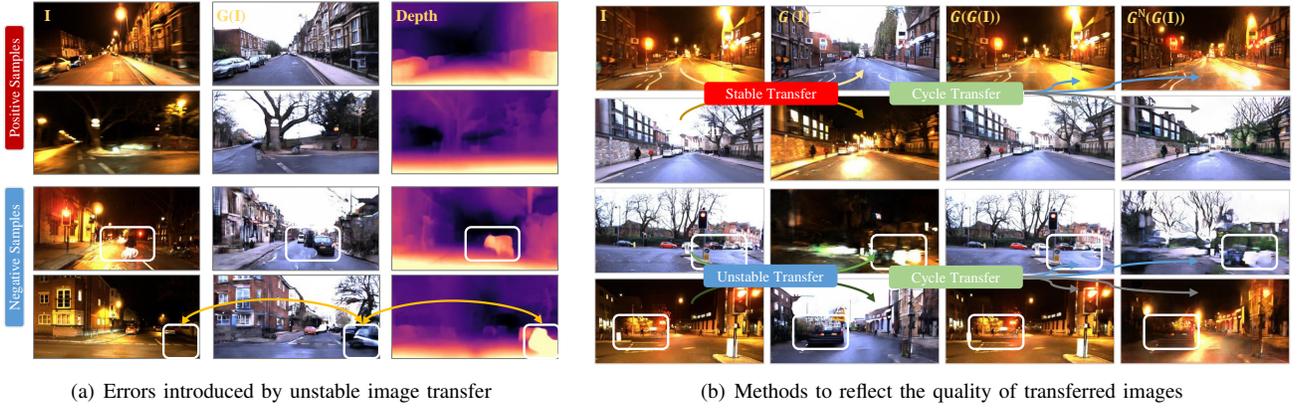


Fig. 3. Samples of applying CycleGAN [44] for monocular depth estimation in complex environments. As shown in (a), since the CycleGAN-based image transfer method is unstable, the wrong image conversion will introduce the incorrect depth estimation into training. To reduce the effects of unstable transfer on domain adaptation, the quality of transferred image (2nd column in (b)) should be accurately evaluated, and our proposed method reflects the quality effectively (4th column in (b)).

well reflect the poor regions of the images. Therefore, we compute the SSIM loss between the raw image ( $I$ ) and cycle transferred image  $\hat{I}_{cycle} = G^N(G(I))$  to quantify the quality of the transferred image  $\hat{I} = G(I)$ .  $G^N$  refers to a fixed transfer model saved at initial epoch  $N$ , in this paper, we set ‘ $N$ ’ = 30 in the training process.

During training, unpaired day and night time images are sent to the transfer network and depth network during training. Since the quality of transferred images varies from model to model and image to image, we evaluate the image quality and re-weight their domain adaptation loss  $L_{DA}$  based on our image adaptation method:

$$L_{total} = (1 - \eta)L_{DA}^d + \eta L_{DA}^x, \quad (20)$$

where  $L_{DA}^d$  and  $L_{DA}^x$  refer to the losses calculated from day-time input image flow and night/rainy night image flow.  $\eta$  stands for:

$$\eta = \frac{(1 - SSIM(I^d, \hat{I}_{cycle}^d))}{((1 - SSIM(I^d, \hat{I}_{cycle}^d)) + (1 - SSIM(I^x, \hat{I}_{cycle}^x))}. \quad (21)$$

However, the above weight  $\eta$  can only adjust the training process to pay more attention to good transferred image samples in each input pair, and it cannot completely eliminate the effects of poor transferred images. In the training framework, the pseudo depth labels are the depth maps predicted by pre-trained daytime models from real and transferred day-time images. As shown in the negative samples in Fig. 3 (a), if we directly use  $L1$  to constrain the consistency of depth maps (Eq. 18), the night-time depth model will be supervised by wrong labels.

Therefore, to filter the errors introduced by instable image transfer, inspired by the minimization loss used in monodepth2 [18], we propose a minimization loss to solve the above errors:

$$L_{OC} = \text{Min}(\langle O_d, O_{d2X} \rangle, \langle O_{d2X}, O_{d2X2d} \rangle) + \text{Min}(\langle O_x, O_{x2d} \rangle, \langle O_{x2d}, O_{x2d2X} \rangle), \quad (22)$$

where  $O_{d2X2d}$  and  $O_{x2d2X}$  stand for the depth maps predicted by the cycle transferred images  $\hat{I}_{d2X2d}^n = G_{x2d}(G_{d2X}^n(I_d))$  and  $\hat{I}_{x2d2X}^n = G_{d2X}^n(G_{x2d}(I_x))$ . As shown in the negative samples

in Fig. 3 (b), the proposed cycle transferred method can reflect the quality of transferred images, which means that the depth of transferred images is consistent with the depth of raw image or the depth of cycle transferred images. With the help of minimization loss,  $L_{OC}$  helps the network learn from the more accurate pseudo labels.

## IV. EXPERIMENTS

### A. Datasets

Since this paper focuses on the unsupervised monocular depth estimation in multiple highly complex environments, we choose the publicly available Oxford RobotCar dataset [29] as our training and testing sets. RobotCar dataset [29] is one of the most famous outdoor datasets, and it contains the image sequences collected in all weather conditions, including rain, night, direct sunlight and snow. The image sequences captured by the left camera of Bumblebee XB3 are used for the experiments of this paper. The images are manipulated to RGB style from the raw recordings with the resolution of 1280x960, and we crop the car-hood of the images and resize them to 512x256. For the day-time and night-time scenarios, we use the sequences from 2014-12-09-13-21-02 and 2014-12-16-18-44-24, which are the same as [26] for a fair comparison. For the rainy night-time scenario that have not received attention in recent research, we choose the sequences from 2014-12-17-18-18-43.

### B. Training and testing sets setup

**Training sets:** For the day-time depth model, the 5 splits of the day-time sequence are used to train the unsupervised framework [18], and the basic pre-trained depth model, monodepth2 (day), is obtained for ITDFA. To improve the performance of this depth model, 15,000 images are uniformly selected from 5 splits for training, and the training set does not include the images taken while parking. For the image transfer model, 5000 images of each scenario are random selected to obtain the image transfer models between different scenarios by using CycleGAN [44]. The selection of these

TABLE I

Comparison with the unsupervised depth estimation methods for night-time scenarios. “M” means that the supervisory signals mainly come from monocular sequences.

Method	Supervision	Depth-range (m)	Error Metrics (Lower is better)				Accuracy (Higher is better)		
			Abs Rel	Sq Rel	RMSE	RMSE log	$\delta < 1.25^1$	$\delta < 1.25^2$	$\delta < 1.25^3$
Monodepth2 (d) [18]	M	40	0.4240	3.8665	8.3071	0.4562	0.317	0.595	0.828
Monodepth2 (n) [18]	M	40	0.2484	5.0838	8.3473	0.3215	0.746	0.882	0.931
Vankadari <i>et al.</i> [26]	M	40	0.2005	2.5750	7.172	0.278	0.735	0.883	0.942
Liu <i>et al.</i> [42]	M	40	0.233	2.344	6.859	0.270	0.631	0.908	0.962
ITDFA <sub>d2n</sub> (JT)	M	40	0.2274	1.4327	5.5849	0.2855	0.573	0.889	0.953
ITDFA <sub>d2n</sub> (ST)	M	40	<b>0.1469</b>	<b>0.9963</b>	<b>4.6851</b>	<b>0.2065</b>	<b>0.778</b>	<b>0.928</b>	<b>0.973</b>
Monodepth2 (d) [18]	M	60	0.5009	6.2895	11.5469	0.5100	0.295	0.538	0.756
Monodepth2 (n) [18]	M	60	0.2899	6.8150	11.9479	0.3562	0.665	0.846	0.914
Vankadari <i>et al.</i> [26]	M	60	0.2327	3.783	10.089	0.319	0.668	0.844	0.924
Liu <i>et al.</i> [42]	M	60	0.231	2.674	8.800	0.286	0.620	<b>0.892</b>	0.956
ITDFA <sub>d2n</sub> (JT)	M	60	0.2789	2.4936	8.5216	0.3350	0.453	0.821	0.937
ITDFA <sub>d2n</sub> (ST)	M	60	<b>0.1869</b>	<b>1.7752</b>	<b>7.370</b>	<b>0.252</b>	<b>0.692</b>	0.889	<b>0.961</b>

TABLE II

Comparison with the unsupervised depth estimation methods for rainy night-time scenarios.

Method	Supervision	Depth-range (meter)	Error Metrics (Lower is better)				Accuracy (Higher is better)		
			Abs Rel	Sq Rel	RMSE	RMSE log	$\delta < 1.25^1$	$\delta < 1.25^2$	$\delta < 1.25^3$
Monodepth2 (d) [18]	M	40	0.4297	4.7547	8.780	0.463	0.369	0.640	0.828
Monodepth2 (r) [18]	M	40	0.3902	50.8281	14.945	0.386	0.616	0.837	0.915
ITDFA <sub>d2r</sub>	M	40	<b>0.1642</b>	<b>0.9688</b>	<b>4.5737</b>	<b>0.2252</b>	<b>0.733</b>	<b>0.932</b>	<b>0.983</b>
ITDFA <sub>d2n</sub>	M	40	0.1678	1.0283	4.6862	0.232	0.733	0.922	0.976
ITDFA <sub>d2n2r</sub>	M	40	<b>0.1495</b>	<b>0.9116</b>	<b>4.3658</b>	<b>0.2117</b>	<b>0.780</b>	<b>0.940</b>	<b>0.983</b>
Monodepth2 (d) [18]	M	60	0.4838	6.7168	11.357	0.509	0.343	0.596	0.781
Monodepth2 (r) [18]	M	60	0.4211	48.4135	17.129	0.423	0.541	0.794	0.897
ITDFA <sub>d2r</sub>	M	60	<b>0.1990</b>	<b>1.7338</b>	<b>6.9703</b>	<b>0.272</b>	<b>0.654</b>	<b>0.882</b>	<b>0.961</b>
ITDFA <sub>d2n</sub>	M	60	0.2060	1.8000	7.0483	0.280	0.644	0.878	0.958
ITDFA <sub>d2n2r</sub>	M	60	<b>0.1788</b>	<b>1.5625</b>	<b>6.5729</b>	<b>0.2530</b>	<b>0.716</b>	<b>0.902</b>	<b>0.967</b>

TABLE III

Quantitative results for ablation study on RobotCar dataset [29] using the night-time images. Depth range is 40m.

$L_{FC}$		$L_O$		Image adaptation	Error Metrics (Lower is better)				Accuracy (Higher is better)		
$L_{FC1}$	$L_{FC_{Gram}}$	$Loc$	$Los$		Abs Rel	Sq Rel	RMSE	RMSE log	$\delta < 1.25^1$	$\delta < 1.25^2$	$\delta < 1.25^3$
✓				✓	0.1675	1.3213	5.2446	0.2342	0.749	0.908	0.962
✓	✓			✓	0.1561	1.1678	4.9162	0.2161	0.767	0.921	0.967
		✓		✓	0.1637	1.1603	5.2037	0.2264	0.727	0.918	0.970
		✓	✓	✓	0.1664	1.0535	4.8925	0.2222	0.739	0.927	0.972
✓	✓	✓	✓	✓	<b>0.1469</b>	<b>0.9963</b>	<b>4.6851</b>	<b>0.2065</b>	<b>0.778</b>	<b>0.928</b>	<b>0.973</b>
✓	✓	✓	✓	✓	0.1545	1.0954	4.8664	0.2140	0.764	0.925	0.970
$L_{FC}$	$L_O$	Image adaptation		“N=”	Abs Rel	Sq Rel	RMSE	RMSE log	$\delta < 1.25^1$	$\delta < 1.25^2$	$\delta < 1.25^3$
✓	✓		✓	latest	0.1501	1.0664	4.8121	0.2105	0.775	0.926	0.971
✓	✓		✓	50	0.1510	1.0585	4.7674	0.2118	0.775	0.926	0.970
✓	✓		✓	30	<b>0.1466</b>	<b>0.9824</b>	<b>4.6649</b>	<b>0.2062</b>	<b>0.780</b>	<b>0.929</b>	<b>0.972</b>
✓	✓		✓	20	0.1582	1.0953	4.8872	0.2174	0.750	0.922	0.971

images follows the rules before. To address the problem of unsupervised monocular depth estimation in night-time and rainy night-time scenarios, three image transfer models should be pre-trained for ITDFA: between day and night ( $G_{d2n}$  and  $G_{n2d}$ ), and between day and rainy night ( $G_{d2r}$  and  $G_{r2d}$ ). In addition, to verify the continuous transfer ability of the model obtained through domain adaptation, we also train an additional image transfer model between night and rainy night ( $G_{n2r}$  and  $G_{r2n}$ ). During training the ITDFA, the training sets used for domain adaptation are the same as that of the image transfer model.

**Testing sets:** During testing, for the night scenario, the testing set is the same as [26] for a fair comparison, which contains 500 night images<sup>1</sup>. While for the more complex

rainy night scenario, 300 rainy night-time images are randomly selected from the remaining splits of each sequence to test the model obtained by ITDFA. The evaluation metrics used in this paper follow previous monocular methods [11], [18], [26], and we evaluate the depth models from the aspect of error and accuracy with different depth range (40m and 60m).

### C. Experimental setup

The experiments are implemented by using Pytorch framework on an NVIDIA RTX 2080 Ti GPU. The network frameworks of depth models (including encoder and decoder) and image transfer model are the same with previous work monodepth2 [18] and CycleGAN [44]. To pre-train the depth model and image transfer model, we use the original setting proposed in monodepth2 [18] and CycleGAN [44] on the new

<sup>1</sup>[https://github.com/zxcqf/RobotCar\\_DepthGT\\_Generate](https://github.com/zxcqf/RobotCar_DepthGT_Generate)

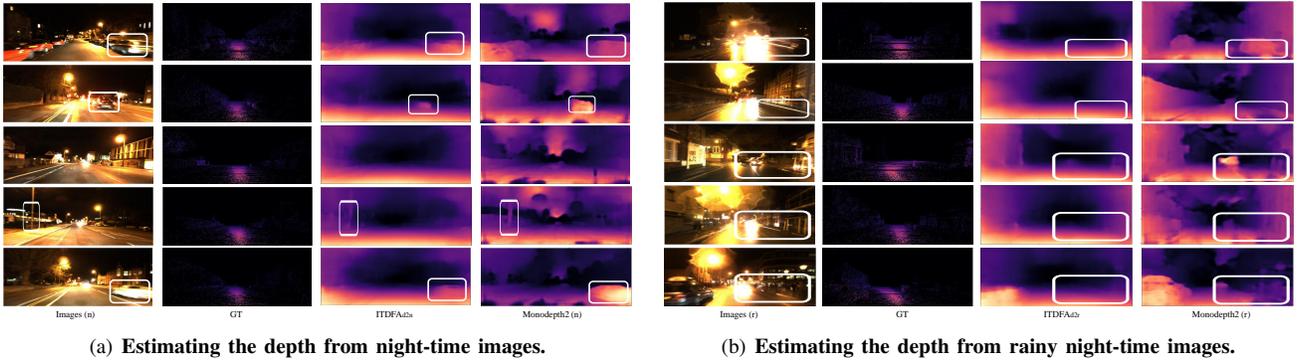


Fig. 4. In (a), “Monodepth2(n)” refers to the depth model trained on night-time images by the unsupervised monocular framework, monodepth2. For night-time scenario, it suffers from complex lighting conditions, like the street lights, car lights, and the reflection of light from the road, which bring big challenges to the unsupervised training framework. In (b), “Monodepth2(r)” refers to the depth model trained on rainy night-time images by the unsupervised monocular framework, monodepth2. For the rainy night scenario, it suffers from not only the complex lighting conditions at night but also the complex reflections on the road and camera caused by rain.

datasets, and the resolution of images in training process is resized to 512x256. For the ITDFA, the framework is trained by using Adam optimizer [49] with a batch size of 1. The learning rate is set as 0.0002 during training. The weights of the three loss components are set to  $\alpha = 15.0$ ,  $\beta = 0.01$  and  $\gamma = 0.001$ . The experimental setup of ITDFA is the same when applied for different complex scenarios, including the night and rainy night. Moreover, the encoders trained for different complex scenarios share the same day-time decoder. Therefore, in practice, this framework only needs to switch the corresponding encoder to cope with the change of environments, which is practical and meaningful for applying in changing environments and all-day depth estimation.

As shown in Fig. 2 (b), the image transfer module and the depth adaptation module can adopt two modes: joint training (JT) and separate training (ST). For the ST mode, the image transfer module is trained in advance, and then the latest saved image transfer model is used to transfer the image between different scenario styles in depth adaptation. To get better cycle transfer models, we firstly test the image transfer models by transferring many sample images, and then we select those models without over-fitting in cycle transfer. According to the experiments, we use the transfer model saved at 30 epoch ( $N=30$ ) to reflect the quality of the transferred images. For the JT mode, the image transfer module and the depth adaptation module are jointly trained. To get a more stable training process, we detach the gradient propagation between the two modules. Moreover, in the JT mode, we directly use the error between the reconstructed image and raw image to reflect the quality of the transferred image. ITDFA(ST) is trained for 50 epoches because the image transfer model is pretrained, while ITDFA(JT) is trained for 200 epoches.

#### D. Results

1) **Annotation:** Results of related experiments are shown in Table I and Table II. To the best of authors’ knowledge, since the proposed work may be the first attempt at solving the monocular depth estimation problem in rainy night-time scenario, and no priors are available in the literature, we

compare the results of different models based on the well known unsupervised framework, monodepth2 [18].

In the tables, “Monodepth2 (X)” refers to the depth model trained on the images from ‘X’ by the monodepth2 framework [18], and X refers to day (d), night (n), rainy night (r). ITDFA<sub>d2X</sub> represents the depth model trained by the proposed ITDFA for the complex scenario ‘X’, and these models (ITDFA<sub>d2n</sub>, ITDFA<sub>d2r</sub>) are adapted from the same day-time model, “Monodepth2 (d)”. Moreover, ITDFA<sub>d2n2r</sub> stands for the rainy night-time model adapted from the night-time model, ITDFA<sub>d2n</sub>, which means that the model is obtained through two domain adaptations by ITDFA. Note that all the models obtained by the proposed ITDFA framework share the same decoder with the “Monodepth2 (d)”. The qualitative results are shown in Fig. 4.

2) **Night-time depth estimation:** To verify the effectiveness of the proposed ITDFA framework in night-time scenario, we compare our model with the state-of-the-art method [26], [42], which only focuses on unsupervised night-time depth estimation. As shown in Table I, for estimating the depth from night-time images, the depth model trained by ITDFA gets a much better performance than the current state-of-the-art method [26], [42] with lower error and higher accuracy. Compared with the joint training mode (ITDFA<sub>d2n</sub>(JT)), the separate training mode (ITDFA<sub>d2n</sub>(ST)) gets more accurate adaptation results, because the pretrained image transfer model can provide more stable transferred images. Besides, the training time of the separate training mode for domain adaptation (1 day) is shorter than that of the joint training mode (1 week). Therefore, for the experiments of rainy night-time depth estimation, we adopt the separate training mode to train the depth models.

3) **Rainy night-time depth estimation:** As shown in Table II, because of the domain drift, the monodepth2(d) does not have a good performance when testing on the rainy night-time images. Compared with the day-time and night-time scenarios, the rainy night-time scenario suffers from not only photometric inconsistency but also the reflection of road and camera caused by rain. Therefore, the monodepth2(r) cannot get an accurate depth estimation because of the limitation of the unsupervised

framework [18] in such complex scenario. The rainy night-time depth model trained by ITDFA (ITDFA<sub>d2r</sub>) shows a much better performance than the monodepth2 models, which proves the effectiveness of the framework proposed in this paper.

**Double jump domain adaptation:** Because the domain gap between night and rainy night is smaller than that between day and rainy night, ITDFA<sub>d2n</sub> shows more accurate depth estimation than the monodepth2(d) on rainy night-time images. To study the influence of different domain gaps on adaptation, we additionally design a two-step adaptation method, in which the day-time depth model is firstly adapted to night-time scenario, and then the night-time model is adapted to rainy night scenario, shown as ITDFA<sub>d2n2r</sub>. As shown in Table II, even if after the second round adaptation, ITDFA<sub>d2n2r</sub> achieves an outstanding accuracy than others in rainy night-time scenario, which means that multi-step adaptation will be helpful for the domain adaptation between large domain gap. Besides, it also proves that the proposed ITDFA framework has the ability to effectively learn the key features for depth decoding, and these key features can be well adapted to new scenarios. The qualitative results are shown in Fig. 4.

4) **Ablation study:** To analyze the effects of each component in the overall loss function  $L_{total}$ , Eq. (15), we design a series of ablations to analyze our approach, and quantitative results are shown in Table III. Experiments show that the constraint of feature space is more effective than that of output space in promoting the consistency of feature maps, because in our framework, the night-time network shares the same decoder with the pretrained day-time models. Besides, the model trained by the constraints from both feature space and output space outperforms the others, which means that the constraints of output space help encoder to focus on learning the key features of depth decoding during training. The smoothness loss helps to improve the accuracy of depth estimation in complex environments. Moreover, the introduction of the proposed image adaptation method effectively improves the accuracy and reduces the error of monocular depth estimation.

## E. Discussion

The ITDFA is an unsupervised framework, and the depth models for highly complex scenarios are trained in a completely unsupervised manner. Neither the image style transfer model nor the depth model use paired images or ground truth labels during training process. Note that all the monocular depth models trained by ITDFA for different environments share the same decoder during testing, which has practical significance. For example, in autonomous driving, facing different weather conditions [50], the vehicle can independently switch to the corresponding encoder to obtain better environmental perception. Although this paper focuses on the effects of unstable transfer on image transfer-based domain adaptation, the proposed method cannot completely eliminate the effects, and the choices of transfer and cycle transfer models are important for the performance of the overall adaptation framework.

## V. CONCLUSION

In this paper, we tackle the problem of unsupervised monocular depth estimation in highly complex environments, which is important and practical for autonomous systems. We survey the related research on solving the limitations of current unsupervised monocular depth estimation framework, and analyze the reason why the unsupervised framework cannot do well in certain highly complex environments. A novel domain adaptation framework, called ITDFA, is proposed in this paper to address the above problem. The proposed ITDFA framework is totally unsupervised and does not use any ground truth labels in the training process. Our method considers the shortcomings of image transfer-based domain adaptation approach and achieves more accurate depth estimation in night-time scenario than the state-of-the-art [26], [42]. Moreover, the performance of ITDFA in the more challenging rainy-time scenario proves the practicability and effectiveness of ITDFA. Therefore, ITDFA is able to provide a way to address the complex environmental change problems faced by monocular depth estimation during practical application. Furthermore, there are still shortcomings that need to be addressed, like enhancing the depth perception of some small objects in complex environments, which is also a promising direction for future work.

## REFERENCES

- [1] W. Zhou, S. Pan, J. Lei, and L. Yu, "TMFNet: Three-input multilevel fusion network for detecting salient objects in RGB-D images," *IEEE Transactions on Emerging Topics in Computational Intelligence*, 2021, accepted.
- [2] C. Cadena, L. Carlone, H. Carrillo, Y. Latif, D. Scaramuzza, J. Neira, I. Reid, and J. J. Leonard, "Past, present, and future of simultaneous localization and mapping: Toward the robust-perception age," *IEEE Transactions on Robotics*, vol. 32, no. 6, pp. 1309–1332, 2016.
- [3] T. Sharma and N. K. Verma, "Estimating depth and global atmospheric light for image dehazing using type-2 fuzzy approach," *IEEE Transactions on Emerging Topics in Computational Intelligence*, 2020, accepted.
- [4] R. Hartley and A. Zisserman, *Multiple View Geometry in Computer Vision*. Cambridge University Press, 2003.
- [5] R. Mur-Artal and J. D. Tardós, "ORB-SLAM2: An open-source SLAM system for monocular, stereo, and RGB-D cameras," *IEEE Transactions on Robotics*, vol. 33, no. 5, pp. 1255–1262, 2017.
- [6] A. Gupta, Y.-S. Ong, and L. Feng, "Insights on transfer optimization: Because experience is the best teacher," *IEEE Transactions on Emerging Topics in Computational Intelligence*, vol. 2, no. 1, pp. 51–64, 2017.
- [7] X. Li, M. Chen, F. Nie, and Q. Wang, "A multiview-based parameter free framework for group detection," in *Thirty-First AAAI Conference on Artificial Intelligence*, 2017.
- [8] Y. Zhang, P. Tiño, A. Leonardis, and K. Tang, "A survey on neural network interpretability," *IEEE Transactions on Emerging Topics in Computational Intelligence*, 2021, accepted.
- [9] X. Li, M. Chen, F. Nie, and Q. Wang, "Locality adaptive discriminant analysis," in *IJCAI*, 2017, pp. 2201–2207.
- [10] S. P. Sahoo, S. Ari, K. Mahapatra, and S. P. Mohanty, "HAR-Depth: A novel framework for human action recognition using sequential learning and depth estimated history images," *IEEE Transactions on Emerging Topics in Computational Intelligence*, 2020, accepted.
- [11] C. Zhao, Q. Sun, C. Zhang, Y. Tang, and F. Qian, "Monocular depth estimation based on deep learning: An overview," *Science China Technological Sciences*, vol. 63, no. 9, pp. 1612–1627, 2020.
- [12] D. Eigen, C. Puhrsch, and R. Fergus, "Depth map prediction from a single image using a multi-scale deep network," in *Advances in Neural Information Processing Systems*, 2014, pp. 2366–2374.
- [13] C. Zhang, Y. Tang, C. Zhao, Q. Sun, Z. Ye, and J. Kurths, "Multitask gans for semantic segmentation and depth completion with cycle consistency," *IEEE Transactions on Neural Networks and Learning Systems*, 2021, accepted.

- [14] C. Zhao, G. G. Yen, Q. Sun, C. Zhang, and Y. Tang, "Masked GAN for unsupervised depth and pose prediction with scale consistency," *IEEE Transactions on Neural Networks and Learning Systems*, vol. 32, no. 12, pp. 5392–5403, 2020.
- [15] M. Poggi, F. Aleotti, F. Tosi, and S. Mattoccia, "On the uncertainty of self-supervised monocular depth estimation," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020, pp. 3227–3237.
- [16] T. Zhou, M. Brown, N. Snavely, and D. G. Lowe, "Unsupervised learning of depth and ego-motion from video," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017, pp. 1851–1858.
- [17] Q. Sun, Y. Tang, C. Zhang, C. Zhao, F. Qian, and J. Kurths, "Unsupervised estimation of monocular depth and vo in dynamic environments via hybrid masks," *IEEE Transactions on Neural Networks and Learning Systems*, 2021, accepted.
- [18] C. Godard, O. Mac Aodha, M. Firman, and G. J. Brostow, "Digging into self-supervised monocular depth estimation," in *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, 2019, pp. 3828–3838.
- [19] S. Lee, S. Im, S. Lin, and I. S. Kweon, "Learning monocular depth in dynamic scenes via instance-aware projection consistency," *arXiv preprint arXiv:2102.02629*, 2021.
- [20] M. Klingner, J.-A. Termöhlen, J. Mikolajczyk, and T. Fingscheidt, "Self-supervised monocular depth estimation: Solving the dynamic object problem by semantic guidance," in *Proceedings of the European Conference on Computer Vision (ECCV)*. Springer, 2020, pp. 582–600.
- [21] J. Spencer, R. Bowden, and S. Hadfield, "DeFeat-Net: General monocular depth via simultaneous unsupervised representation learning," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020, pp. 14 402–14 413.
- [22] J. Engel, V. Koltun, and D. Cremers, "Direct sparse odometry," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 40, no. 3, pp. 611–625, 2017.
- [23] N. Yang, R. Wang, X. Gao, and D. Cremers, "Challenges in monocular visual odometry: Photometric calibration, motion bias, and rolling shutter effect," *IEEE Robotics and Automation Letters*, vol. 3, no. 4, pp. 2878–2885, 2018.
- [24] Y. Tang, C. Zhao, J. Wang, C. Zhang, Q. Sun, W. Zheng, W. Du, F. Qian, and J. Kurths, "An overview of perception and decision-making in autonomous systems in the era of learning," *arXiv preprint arXiv:2001.02319*, 2020.
- [25] C. Zhang, J. Wang, G. G. Yen, C. Zhao, Q. Sun, Y. Tang, F. Qian, and J. Kurths, "When autonomous systems meet accuracy and transferability through ai: A survey," *Patterns (Cell Press)*, vol. 1, no. 4, art. No. 100050, 2020.
- [26] M. Vankadari, S. Garg, A. Majumder, S. Kumar, and A. Behera, "Unsupervised monocular depth estimation for night-time images using adversarial domain feature adaptation," in *Proceedings of the European Conference on Computer Vision (ECCV)*. Springer, 2020, pp. 443–459.
- [27] Y.-C. Chen, Y.-Y. Lin, M.-H. Yang, and J.-B. Huang, "Crdoco: Pixel-level domain transfer with cross-domain consistency," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019, pp. 1791–1800.
- [28] T. Kim, M. Jeong, S. Kim, S. Choi, and C. Kim, "Diversify and match: A domain adaptive representation learning paradigm for object detection," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019, pp. 12 456–12 465.
- [29] W. Maddern, G. Pascoe, C. Linegar, and P. Newman, "1 Year, 1000km: The Oxford RobotCar Dataset," *The International Journal of Robotics Research (IJRR)*, vol. 36, no. 1, pp. 3–15, 2017.
- [30] T. Salimans, I. Goodfellow, W. Zaremba, V. Cheung, A. Radford, and X. Chen, "Improved techniques for training gans," *arXiv preprint arXiv:1606.03498*, 2016.
- [31] K. Wang, C. Gou, Y. Duan, Y. Lin, X. Zheng, and F.-Y. Wang, "Generative adversarial networks: introduction and outlook," *IEEE/CAA Journal of Automatica Sinica*, vol. 4, no. 4, pp. 588–598, 2017.
- [32] T. Tommasi, N. Patricia, B. Caputo, and T. Tuytelaars, "A deeper look at dataset bias," in *Domain adaptation in computer vision applications*. Springer, 2017, pp. 37–55.
- [33] X. Ye, Z. Li, B. Sun, Z. Wang, R. Xu, H. Li, and X. Fan, "Deep joint depth estimation and color correction from monocular underwater images based on unsupervised adaptation networks," *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 30, no. 11, pp. 3995–4008, 2019.
- [34] A. Gurram, A. F. Tuna, F. Shen, O. Urfalioglu, and A. M. López, "Monocular depth estimation through virtual-world supervision and real-world sfm self-supervision," *IEEE Transactions on Intelligent Transportation Systems*, pp. 1–14, 2021.
- [35] G. Wilson and D. J. Cook, "A survey of unsupervised deep domain adaptation," *ACM Transactions on Intelligent Systems and Technology (TIST)*, vol. 11, no. 5, pp. 1–46, 2020.
- [36] S. Roy, E. Krivosheev, Z. Zhong, N. Sebe, and E. Ricci, "Curriculum graph co-teaching for multi-target domain adaptation," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, pp. 5351–5360.
- [37] S. M. Ahmed, D. S. Raychaudhuri, S. Paul, S. Oymak, and A. K. Roy-Chowdhury, "Unsupervised multi-source domain adaptation without access to source data," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, pp. 10 103–10 112.
- [38] J. N. Kundu, P. K. Uppala, A. Pahuja, and R. V. Babu, "Adadepth: Unsupervised content congruent adaptation for depth estimation," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018, pp. 2656–2665.
- [39] X. Gu, Y. Guo, F. Deligianni, and G.-Z. Yang, "Coupled real-synthetic domain adaptation for real-world deep depth enhancement," *IEEE Transactions on Image Processing*, vol. 29, pp. 6343–6356, 2020.
- [40] S. Zhao, H. Fu, M. Gong, and D. Tao, "Geometry-aware symmetric domain adaptation for monocular depth estimation," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019, pp. 9788–9798.
- [41] J. He, X. Jia, S. Chen, and J. Liu, "Multi-source domain adaptation with collaborative learning for semantic segmentation," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, pp. 11 008–11 017.
- [42] L. Liu, X. Song, M. Wang, Y. Liu, and L. Zhang, "Self-supervised monocular depth estimation for all day images using domain separation," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 12 737–12 746.
- [43] A. Johnston and G. Carneiro, "Self-supervised monocular trained depth estimation using self-attention and discrete disparity volume," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020, pp. 4756–4765.
- [44] J.-Y. Zhu, T. Park, P. Isola, and A. A. Efros, "Unpaired image-to-image translation using cycle-consistent adversarial networks," in *2017 IEEE International Conference on Computer Vision (ICCV)*, 2017, pp. 2242–2251.
- [45] S. Lee, S. Cho, and S. Im, "DRANet: Disentangling representation and adaptation networks for unsupervised cross-domain adaptation," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, pp. 15 252–15 261.
- [46] J. Cheng, A. Jaiswal, Y. Wu, P. Natarajan, and P. Natarajan, "Style-aware normalized loss for improving arbitrary style transfer," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, pp. 134–143.
- [47] C. Godard, O. Mac Aodha, and G. J. Brostow, "Unsupervised monocular depth estimation with left-right consistency," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017, pp. 270–279.
- [48] C. Wang, J. M. Buenaposada, Z. Rui, and S. Lucey, "Learning depth from monocular videos using direct methods," in *Proceedings of the 2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, 2018, pp. 2022–2030.
- [49] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," *arXiv preprint arXiv:1412.6980*, 2014.
- [50] W.-T. Chu, X.-Y. Zheng, and D.-S. Ding, "Camera as weather sensor: Estimating weather information from single images," *Journal of Visual Communication and Image Representation*, vol. 46, pp. 233–249, 2017.