

SC-DepthV3: Robust Self-supervised Monocular Depth Estimation for Dynamic Scenes

Libo Sun*, Jia-Wang Bian*, Huangying Zhan, Wei Yin, Ian Reid, Chunhua Shen

Abstract—Self-supervised monocular depth estimation has shown impressive results in static scenes. It relies on the multi-view consistency assumption for training networks, however, that is violated in dynamic object regions and occlusions. Consequently, existing methods show poor accuracy in dynamic scenes, and the estimated depth map is blurred at object boundaries because they are usually occluded in other training views. In this paper, we propose SC-DepthV3 for addressing the challenges. Specifically, we introduce an external pretrained monocular depth estimation model for generating single-image depth prior, namely *pseudo-depth*, based on which we propose novel losses to boost self-supervised training. As a result, our model can predict sharp and accurate depth maps, even when training from monocular videos of highly dynamic scenes. We demonstrate the significantly superior performance of our method over previous methods on six challenging datasets, and we provide detailed ablation studies for the proposed terms. Source code and data have been released at https://github.com/JiawangBian/sc_depth_pl

Index Terms—Monocular Depth Estimation, Unsupervised Learning, Self-supervised Learning, Knowledge Distillation

1 INTRODUCTION

MONOCULAR depth estimation [2] has attracted great attention in computer vision. It provides valuable cues for various downstream tasks, such as semantic image segmentation [3], salient object detection [4], 3D reconstruction [5], novel view synthesis [6], and visual odometry [1], [7]. Early work [2], [8] solves the monocular depth estimation problem by using supervised learning. However, these methods rely on ground-truth depth labels that are not always available in real-world scenes. To address this limitation, self-supervised monocular depth estimation methods were proposed and showed that a depth network could be trained from stereo image pairs [9] or monocular videos with ego-motions [10] without the need for ground-truth depth labels. We focus on self-supervised learning of monocular depth from videos since only a single camera is required to collect training data in this setup, which has great potential for advancing real-world applications.

Self-supervised methods typically rely on the multi-view consistency assumption for training networks, *e.g.*, the photometric loss [10] and geometry consistency loss [1] that were used in previous methods. This assumption provides effective constraints for learning scene geometry, while it is violated at regions with occlusion (*e.g.*, object boundaries) and moving objects. Therefore, existing methods often show only excellent results in (almost) static scenes such as KITTI [11] and NYUv2 [12] datasets. When training on more challenging dynamic datasets that have an amount of fast-moving objects, previous state-of-the-art methods [1], [13], [14] show poor accuracy. Moreover, the estimated depth map is blurred at object boundaries because they are usually occluded in other training views. We illustrate several examples of

qualitative monocular depth results in Fig. 1.

To address the issues caused by moving objects and occlusions, existing approaches usually detect these bad regions and then exclude them from training. The methods can be categorized into four classes according to how they detect dynamic regions, involving in the prediction-based [10], semantic-based [15]–[18], flow-based [19]–[22], and geometry-based [1]. These methods can reduce corruption from noisy losses during training and generally improves overall accuracy, however, it leads to poor results on dynamic regions at inference time because these regions are not sufficiently regularized in training. There are also more sophisticated approaches [23], [24] that model the velocity of each moving object in multiple views, but they rely on solving a challenging problem in themselves.

We propose SC-DepthV3 in this paper, which addresses the above-mentioned issues by leveraging external single-image constraints. Specifically, we leverage an off-the-shelf monocular depth estimation model [25] to generate the single-image depth prior, which we term *pseudo-depth*. Based on it, we propose effective losses to constrain the depth estimation network in self-supervised learning. Here, we use LeReS [25] for generating pseudo-depth, which is trained in large-scale datasets with supervised learning and enables zero-shot generalization in previously unseen data. The excellent qualitative results have been demonstrated in [25], while we find that pseudo-depth may show low quantitative accuracy. Fig. 3 gives an example, where we visualize the error map of pseudo-depth by comparing it with the ground truth. This phenomenon makes supervised zero-shot methods unsuitable for accuracy-sensitive tasks such as visual SLAM and 3D Reconstruction. Furthermore, as pseudo-depth is not quantitatively accurate, it is non-trivial to use it for boosting self-supervised learning. In this paper, our technical contribution is designing effective losses that use imperfect pseudo-depth. It is also worth mentioning that although we use external depth estimation networks, they are only trained once and can be used as off-the-shelf tools in new scenes. Therefore, in practice, our method does not add extra cost to purely self-supervised methods.

- First two authors contributed equally. J.-W. Bian is the corresponding author. He is with the University of Oxford, United Kingdom;
- L. Sun and I. Reid are with The University of Adelaide, Australia;
- H. Zhan is with the OPPO US Research Center, United States.
- W. Yin is with the DJI Technology, China;
- C. Shen is with Zhejiang University, China.
- Part of this work was done when J.-W. Bian, H. Zhan, W. Yin, and C. Shen were with the University of Adelaide;

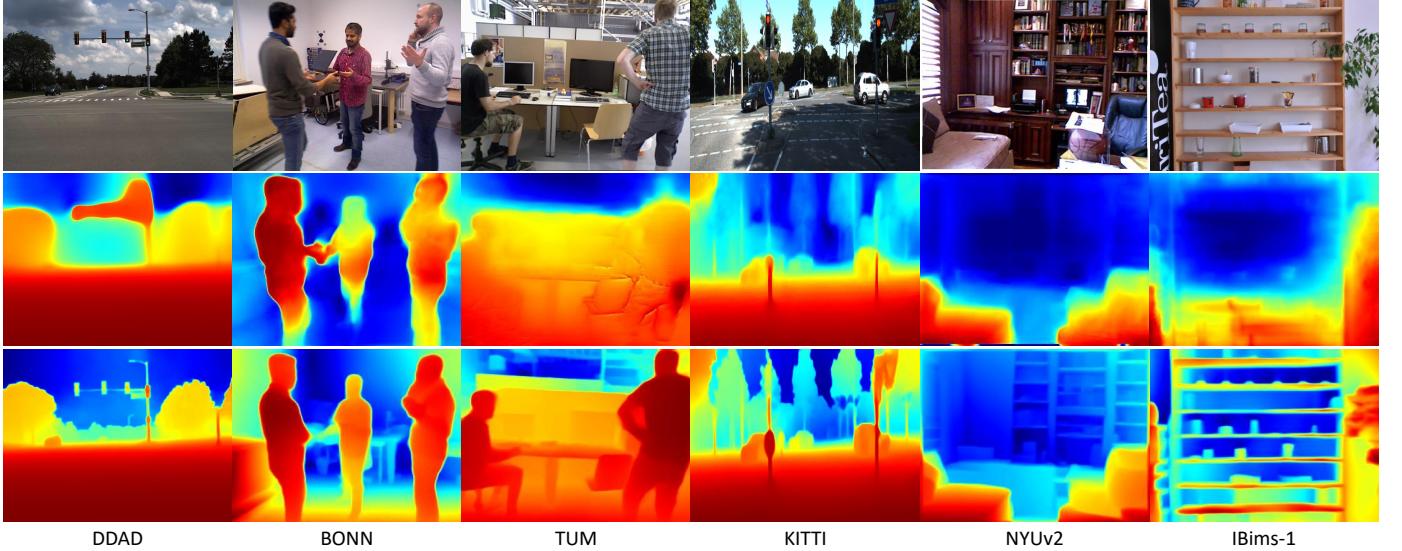


Fig. 1. Qualitative monocular depth estimation results on six datasets. We compare our method (third row) with SC-Depth [1] (second row), which is one of the previous state-of-the-art self-supervised methods. Compared with it, our method enables more robust learning in dynamic scenes (left three columns) and generates sharper depth maps, particularly at object boundary areas.

The key to solving the dynamic region issue is the proposed **Dynamic Region Refinement** (DRR) module. The method is inspired by an observation, *i.e.*, we find that pseudo-depth maintains excellent depth ordinal (the further/nearer relations) between any two objects or pixels. To capitalize on these findings, we propose to extract the “ground-truth” depth ordinal information between dynamic and static regions (from pseudo-depth) and use it to regularize the self-supervised depth estimation in dynamic regions. Specifically, we sample point pairs between two regions and apply depth ranking loss [26]. This is effective because the static backgrounds have already been well-supervised by multi-view losses, and the dynamic regions could be uniquely localized by sampling sufficient point pairs between dynamic and static regions. Our method is also based on the fact that the depth ordinal in pseudo-depth is sufficiently accurate [25]. Furthermore, to segment dynamic regions from static backgrounds, we use the self-discovered mask that was proposed in SC-Depth [1] and generated by computing forward-backward depth inconsistency in self-supervised training, so the external segmentation networks are not required. Fig. 4 illustrates the proposed DRR module.

Moreover, we observe that pseudo-depth shows smooth local structures and clean object boundaries. This motivates us to propose a **Local Structure Refinement** (LSR) module to improve the self-supervised depth estimation *w.r.t.* depth details. The proposed module contains two parts. On the one hand, we extract the surface normal from both pseudo-depth and network-predicted depth, and we constrain them to be consistent by applying a normal matching loss. This improves the overall depth significantly. On the other hand, we constrain depth estimation at object boundary areas by applying our proposed relative normal angle loss. More specifically, we sample point pairs around image edges and enforce their relative normal angles to be consistent between pseudo-depth and self-supervised depth. As a result, our method improves qualitative depth estimation results significantly, particularly at object boundaries. Fig. 1 shows several examples of the qualitative depth estimation results.

Our contributions are as follows:

- We propose SC-DepthV3 for robust self-supervised learning of monocular depth in highly dynamic scenes, which allows for predicting accurate and sharp depth maps.
- We propose Dynamic Region Refinement (DRR) and Local Structure Refinement (LSR) modules, which are based on pseudo-depth to boost self-supervised learning.
- We conduct comprehensive experiments and ablation studies on six challenging datasets. The results demonstrate the efficacy of our proposed methods.

2 RELATED WORK

Self-supervised Monocular Depth Estimation. Garg et al. [9] proposed to train monocular depth estimation models on stereo image pairs by using the photometric loss. Zhou et al. [10] proposed to train the depth estimation model on videos by jointly training a pose estimation model. Following them, many advanced techniques [1], [13], [14], [17], [19], [22], [27]–[30] were proposed to boost the performance. However, multi-view ambiguities make the self-supervised method hard to handle dynamic objects and object boundaries. Previous methods either excluded these regions from training [1], [15], [16], [19] or modeled the object motions [23], [24], but both solutions have their drawbacks. More specifically, simply excluding dynamic regions would result in poor accuracy on these regions at the inference time, and modeling each object’s motion is ill-posed and may not be robust in dynamic scenes. Compared with them, our method leverages pretrained single-image prior for resolving multi-view ambiguities, leading to a SOTA self-supervised depth estimation method. More recent methods include [31]–[33].

SC-Depth Series Methods. This paper is the third version of the SC-Depth series methods. In the SC-Depth [1], we addressed the scale inconsistency issue, so our method enables scale-consistent depth estimation over the video, which is beneficial to video-based tasks such as Visual SLAM. In the SC-DepthV2 [30], we analyzed the rotation issue in videos that are captured by handheld

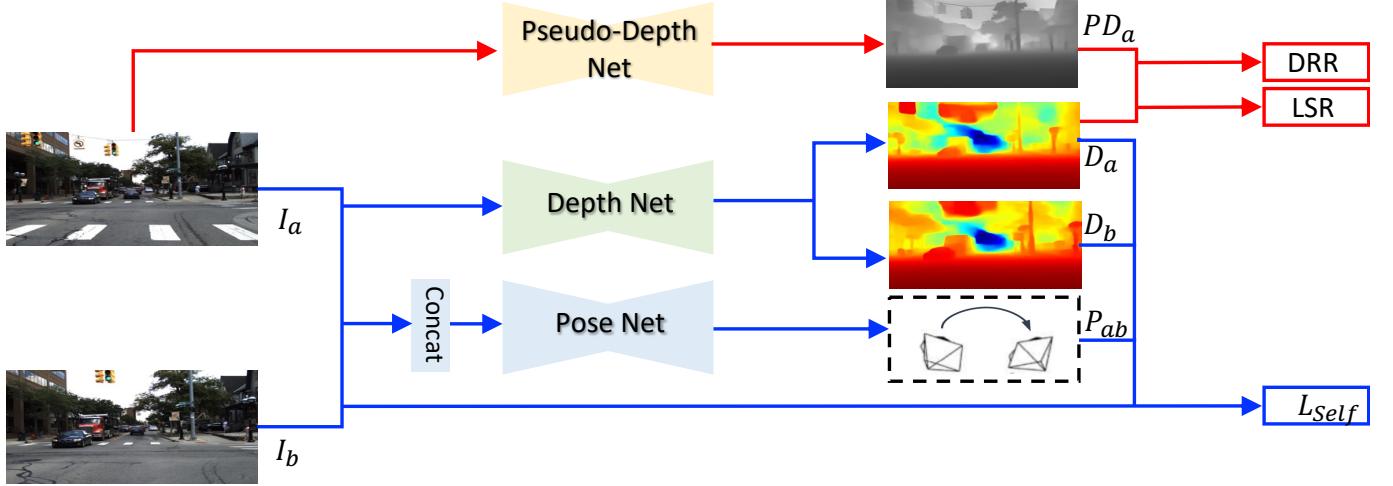


Fig. 2. Method overview. Firstly, given a training sample (i.e., I_a and I_b two images), we follow SC-Depth [1] to compute self-supervised losses L_{Self} (Eqn. 6), which is described in Sec. 3.1. Secondly, we generate pseudo-depth PD_a using a pretrained depth estimation network, which is discussed in Sec. 3.2. Finally, we propose DRR and LSR modules to constrain the network prediction (D_a) by using PD_a , which are presented in Sec. 3.3 and Sec. 3.4, respectively.

cameras, and we proposed an auto-rectify network to handle the large rotation. The V1 and V2 have shown great accuracy in both indoor and outdoor scenes. However, their predicted depth maps are blurred at object boundaries, and they suffer in highly dynamic scenes. In this paper, we propose SC-DepthV3 address the issue of dynamic objects and blurred object boundaries.

Zero-shot Monocular Depth Estimation. Many existing methods leverage large-scale datasets and supervised training [25], [34]–[40] to train monocular depth estimation models towards zero-shot generalization on unseen data. For example, [34]–[39] collect stereo images/videos from the internet and use geometric reconstruction tools [41], [42] to generate dense ground-truth depth labels. [40] export perfect ground-truth depths from the synthetic 3D movies [43]. Recently, LeReS [25] and DPT [44] achieve the state-of-the-art performance. However, note that their predicted depths are scale-shift-invariant, due to the high diversity of different scenes, which show low quantitative accuracy in out-of-distribution data and cannot be used for 3D reconstruction. Nevertheless, we find that their predicted depths carry good attributes that could be leveraged for boosting self-supervised learning of monocular depth estimation. Compared with these methods, our method enables consistent and accurate depth estimation for video-based tasks such as Visual SLAM, which has been demonstrated in SC-Depth [1], thanks to the scale-consistency constraints.

Knowledge Transfer. Our method is also related to knowledge transfer approaches, because the proposed method can be regarded as transferring the knowledge of pretrained monocular depth estimation models [25] to our self-supervised trained models. However, we argue that our method is very different from previous knowledge transfer or distillation methods. On the one hand, knowledge is often transferred by finetuning pretrained models in new datasets [45], which is not our case and cannot solve the challenges in our problem. The main issue in our problem is the imperfect self-supervised loss, so even if we finetune pretrained models (i.e., it provides a good initialization), the model would become worse and worse with training due to the deficient self-

supervised loss functions. On the other hand, knowledge transfer could also be achieved by conducting semi-supervised learning on mixed datasets. Specifically, we can train models on both previous large-scale datasets with ground-truth labels and new datasets without annotations, and then we apply supervised loss in the former and self-supervised loss in the latter. However, this involves new challenges of mix-data training, long training time, and the maintenance of large-scale previous data. In contrast, our method is more elegant than semi-supervised training, since the teacher model is trained only once on large-scale datasets and can be used as an off-the-shelf tool to generate pseudo-depth in new scenes. Moreover, our student model shows significantly higher accuracy than the teacher model, which is rare in the field of knowledge distillation.

3 METHOD

Fig. 2 illustrates an overview of the proposed method. First, our method is based on SC-Depth [1] for basic self-supervised training, which we describe in detail in Sec. 3.1. Second, we discuss the single-image depth prior in Sec. 3.2 that is generated by using the off-the-shelf monocular depth estimation methods and used in our method for generating auxiliary supervision signals. Finally, we describe the Dynamic Region Refinement (DRR) in Sec. 3.3 and Local Structure Refinement (LSR) modules in Sec. 3.4, respectively, which are the proposed terms to boost self-supervised training.

3.1 Self-supervised Depth Learning (SC-Depth)

In the self-supervised learning framework, a monocular depth estimation network (DepthNet) and a relative 6-DoF camera pose estimation network (PoseNet) are jointly trained on a large number of monocular videos. First, given a consecutive image pair (I_a , I_b) randomly sampled from a training video, we predict their depths (D_a , D_b) by forwarding the DepthNet and estimate their relative 6-DoF camera pose P_{ab} by forwarding the PoseNet. Then, we generate the warping flow between two images using the predicted depth and pose, followed by synthesizing the I'_a using the flow and I_b via bi-linear interpolation. Finally, we penalize the

color inconsistencies between I_a and I'_a , and we also constrain the geometry consistency between D_a and D_b , which back-propagates the gradients to the networks. The objective function is described below.

First, we use the geometry consistency loss L_G [1] to encourage the predicted depths (D_a, D_b) to be consistent with each other in 3D space. Formally,

$$L_G = \frac{1}{|\mathcal{V}|} \sum_{p \in \mathcal{V}} D_{\text{diff}}(p), \quad (1)$$

where V stands for valid points that are projected inside the image. D_{diff} stands for the pixel-wise depth inconsistency between D_a and D_b , which is detailed explained in [1]. With it, we can obtain the self-discovered mask:

$$M_s = 1 - D_{\text{diff}}, \quad (2)$$

which assigns lower weights to dynamics and occlusions than static regions, since the former is geometrically inconsistent across multiple views. We use this mask in our proposed DRR module (Sec. 3.3) to localize dynamic regions.

Second, we use the weighted photometric loss L_P^M to constrain the warping flow between I_a and I_b that is generated by the D_a and P_{ab} . Formally,

$$L_P^M = \frac{1}{|\mathcal{V}|} \sum_{p \in \mathcal{V}} (M_s(p) \cdot L_P(p)), \quad (3)$$

$$L_P = \frac{1}{|\mathcal{V}|} \sum_{p \in \mathcal{V}} (\lambda \|I_a(p) - I'_a(p)\|_1 + (1 - \lambda) \frac{1 - \text{SSIM}_{aa'}(p)}{2}), \quad (4)$$

where I'_a is synthesized from I_b using the warping flow, and SSIM [46] is a widely-used metric to measure image similarity. We set λ to 0.15 as in [1].

Third, we use the edge-aware smoothness loss to regularize the predicted depth map. Formally,

$$L_S = \sum_p (e^{-\nabla I_a(p)} \cdot \nabla D_a(p))^2, \quad (5)$$

where ∇ is the first derivative along spatial directions, which guides smoothness by image edges.

Overall, our objective function is formulated as follows:

$$L_{\text{Self}} = \alpha L_P^M + \beta L_G + \gamma L_S. \quad (6)$$

We set $\alpha = 1$, $\beta = 0.5$, and $\gamma = 0.1$ as in [1]. Note that we will replace L_S with the proposed normal loss L_N in Sec. 3.4. Moreover, we also use the auto-masking and per-pixel minimum reprojection loss that are proposed in [13] to filter stationary and non-best points during training.

3.2 Single-Image Depth Prior

Our idea is to leverage the pretrained monocular depth estimation network for generating single-image depth prior, which is then used to boost self-supervised learning. Here we use LeReS [25] to generate *pseudo-depth*, which is trained on large-scale datasets with ground-truth depth labels. Thanks to the supervised training on large-scale data, it shows excellent zero-shot generalization performance on unseen scenes. Note that LeReS was not trained on datasets that we use in this paper for evaluation. An example of LeReS outputs is shown in Fig. 3, where it shows plausible visual results on the DDAD dataset but poor quantitative accuracy.

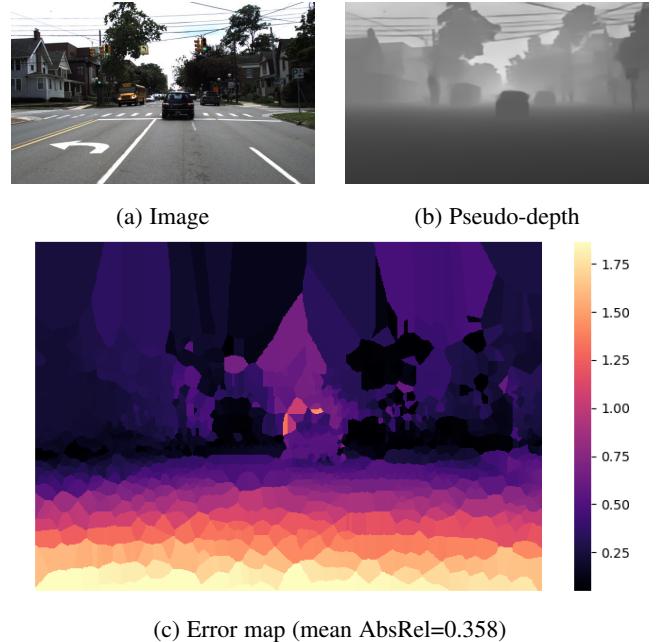


Fig. 3. Visualization of pseudo-depth (LeReS [25]) on the DDAD dataset. For the error map, we show the AbsRel error, and we use the nearest interpolation for pixels where the ground-truth depth labels (sparse LiDAR points) are unavailable. It shows that LeReS [25] can generalize to previously unseen data with plausible visual results (b), however, high quantitative accuracy is not guaranteed. Here “AbsRel=0.358” is averaged over all testing images. This indicates that our idea of leveraging pseudo-depth for boosting self-supervised training is motivated, and it is also non-trivial to use it.

These phenomena echo our motivation, *i.e.*, pseudo-depth is not accurate enough but has potential that can be leveraged for boosting self-supervised learning. More specifically, we find that the good visual results are contributed to several aspects, including (a) correct depth ordinal (nearer/further relation) between objects; (b) excellent smoothness in predicted depth; (c) sharp depth prediction at object boundaries.

Therefore, based on the above observation, we propose two modules to extract effective supervision signals from pseudo-depth. First, we propose a *Dynamic Region Refinement* (DRR) module that regularizes self-supervised training with a depth ranking constraint, particularly boosting depth prediction on moving objects. Second, we propose a *Local Structure Refinement* (LSR) module that constrains the smoothness and object boundaries of the predicted depth. The proposed two modules are presented in Sec. 3.3 and Sec. 3.4, respectively.

3.3 Dynamic Region Refinement

The key to our proposed dynamic region refinement (DRR) module is constraining depth estimation on dynamic regions by enforcing the nearer/further relation *w.r.t.* the predicted depths on static regions. Specifically, this is based on two assumptions including (i) The accurate depth ranking relations between any two pixels can be extracted from pseudo-depth; (ii) depth prediction on static regions is sufficiently accurate thanks to the self-supervised losses. The assumptions are valid, as demonstrated in prior work [13], [25], so our idea is generally effective. In the proposed DRR module, we first sample point pairs between dynamic and static regions, where the segmentation of images is obtained in a self-supervised manner (Eqn. 2). Then we compute depth ranking loss

on sampled point pairs to regularize the predicted depth map. The proposed sampling method and loss function are presented in the following paragraphs, and Fig. 4 illustrates a training example.

Dynamic-focused Sampling. To sample point pairs between static and dynamic regions, we need the segmentation of training images. This could be achieved with the use of pretrained semantic segmentation networks, *e.g.*, we can assume objects of certain classes such as vehicles as moving objects and others as static backgrounds. However, it involves extra data preprocessing and pretrained networks. Moreover, the dynamic of an object does not necessarily rely on the semantic classes, *e.g.*, a chair can be dynamic while a person is moving it around. Instead, we derive dynamic/static segmentation from the self-discovered mask (Eqn. 2) that is computed based on the geometric consistency. It is a soft weight mask and assigns smaller values for depth-inconsistent regions (dynamics or occlusions) than others (static regions). To obtain binary segmentation, we propose to rank weights and pick the lowest 20% as potential dynamic regions, rather than doing hard thresholding. Here we assume that the ratio of moving objects pixels is around 20% or less, which is true in most real-world scenes. Then for each point in the dynamic regions, we pair it with a point that is randomly sampled from static regions. Moreover, other than constructing dynamic-static pairs as discussed above, we also sample point pairs randomly from the whole image, which serves as an additional global regularization.

Confident Depth Ranking Loss. We compute the depth ranking loss on the sampled point pairs in training. The original loss function was proposed in [26]. Formally, for a pair of points with predicted depth values $[p_0, p_1]$, the loss is

$$\phi'(p_0, p_1) = \begin{cases} \log(1 + \exp(-\ell(p_0 - p_1))), & \ell \neq 0 \\ (p_0 - p_1)^2, & \ell = 0 \end{cases} \quad (7)$$

where ℓ is the ground truth ordinal label, which can be induced by a ground truth depth map:

$$\ell = \begin{cases} +1, & p_0^*/p_1^* \geq 1 + \tau, \\ -1, & p_0^*/p_1^* \leq \frac{1}{1+\tau}, \\ 0, & \text{otherwise.} \end{cases} \quad (8)$$

Here τ is a threshold, which is 0.03 in previous work [47], and p^* denotes pseudo-depth.

We empirically find that Eqn. 7 is sub-optimal in our method since pseudo-depth is not as accurate as the ground-truth depth. Therefore, we have to take the confidence of pseudo-depth ordinals into consideration. Specifically, we observe that the ordinal is often sufficiently reliable when two points have sufficiently different depth values, *i.e.*, when $p_0^*/p_1^* \gg 1$ or $p_0^*/p_1^* \ll 1$, and otherwise, it may be unreliable when two depth values are very close, *i.e.*, when $p_0^*/p_1^* \approx 1$.

Based on the above observation, we propose to (a) increase τ from 0.03 to 0.15 for higher tolerance; and (b) ignore point pairs that have $\ell = 0$. Formally, we reformulate Eqn. 7 as

$$\phi(p_0, p_1) = \log(1 + \exp(-\ell(p_0 - p_1))). \quad (9)$$

Therefore, our Confident Depth Ranking Loss is defined as:

$$L_{CDR} = \frac{1}{|\Omega|} \sum_{p \in \Omega} \phi(p), \quad (10)$$

where Ω stands for the sampled point pairs that have $\ell \neq 0$.

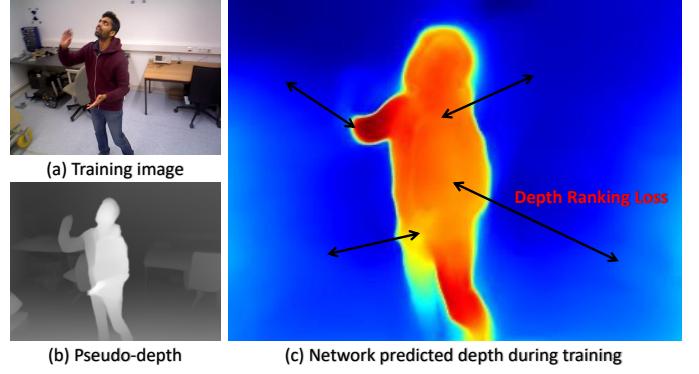


Fig. 4. Dynamic region refinement. We sample point pairs between dynamic and static regions, and then we apply depth ranking loss to constrain the network-predicted depth (c) during training. The “ground-truth” depth ordinal is extracted from pseudo-depth (b). To segment dynamic regions from static backgrounds, we use the self-discovered mask (Eqn. 2), so there is no extra computational cost.

3.4 Local Structure Refinement

As mentioned above in Sec. 3.2, pseudo-depth provides excellent depth smoothness. In this section, we propose to leverage such attributes to regularize the self-supervised depth. Our idea is to (i) constrain the surface normals that are derived from predicted depths and pseudo-depths to be matched; and (ii) constrain two depth maps to be consistent w.r.t. relative normal angles of sampled point pairs around edges. Here, the first step is focused on refining the overall depth structures, and the second step is focused on improving object boundary regions. The details are provided below.

Normal Matching Loss. The edge-aware depth smoothness loss (Eqn. 5) is often used in self-supervised monocular depth estimation, and it is also used in our baseline. Here, we propose to replace it with the normal matching loss:

$$L_N = \frac{1}{N} \sum_{i=1}^N \|n_i - n_i^*\|_1, \quad (11)$$

where n_i is the surface normal derived from the predicted depth, and n_i^* is the normal derived from pseudo-depth. N stands for the total number of pixels in the image. The pixel-wise loss function provides strong supervision for overall depth structures.

Edge-aware Relative Normal Loss. Not only do overall structure refinement, but also we focus on object boundary areas. Specifically, we sample point pairs around image edges and constrain the relative normal angles of sampled point pairs to be consistent with pseudo-depth. Here, we use edge-guided sampling that was proposed in [47] to construct point pairs $\langle A, B \rangle$, and we define the Edge-aware Relative Normal Loss as:

$$L_{ERN} = \frac{1}{N} \sum_{i=1}^N \|n_{Ai} \cdot n_{Bi} - n_{Ai}^* \cdot n_{Bi}^*\|_1, \quad (12)$$

where n_A denotes the normal of a sampled point from the predicted depth, and $*$ denotes pseudo-depth. Combining the edge-guided sampling and relative normal loss, we can effectively constrain the depth estimation on object boundary regions.

The proposed L_{ERN} is similar to the pair-wise normal loss that is proposed in [25], while the latter samples point pairs from

edges, planes, and whole images. In contrast, we sample points solely from edges because sampling from other regions requires high-quality ground-truth depth. In our case, pseudo-depth is not accurate enough to maintain high-quality global structures, hence we only constrain the local structure. We analyze this effect with a detailed ablation study in Sec. 4.3.

3.5 Training

Losses. Based on the proposed two refinement modules, we rewrite the overall loss function (Eqn. 6) of our baseline and obtain the new objective function as:

$$L = \alpha L_P^M + \beta L_G + \gamma L_N + \delta L_{CDR} + \epsilon L_{ERN}, \quad (13)$$

where we set $\alpha = 1$, $\beta = 0.5$, and $\gamma = \delta = \epsilon = 0.1$ in training based on empirical tuning.

Networks. Our depth and pose networks are the same as previous work [1], [13], where we use ResNet-18 [48] backbone for both depth and pose estimation networks. The depth network is a U-Net structure [49] with a DispNet [10] as the decoder. The activations are sigmoids at the output layer and ELU nonlinearities [50] elsewhere. We convert the sigmoid output x to depth with $D = 1/(ax + b)$, where a and b are chosen to constrain D between 0.1 and 100 units. The pose network accepts two RGB frames as input and outputs the 6D relative pose. We modify the first layer of ResNet-18 to have six channels for accepting two-frame inputs, and features are decoded to 6-DoF parameters via four convolutional layers.

Training Details. We implement the proposed method using the PyTorch library [51]. Following [10], [21], [52], we use a snippet of three sequential video frames as a training sample. The images are augmented with random scaling, cropping, and horizontal flips during training. We use the Adam [53] optimizer and set the learning rate to be 10^{-4} . We initialize the encoder by using the pre-trained model on ImageNet [54]. We train our networks in $100k$ iterations on each dataset.

4 EXPERIMENT

4.1 Datasets and Evaluation Metrics

The proposed method focuses on boosting self-supervised monocular depth estimation in challenging dynamic scenes, so we mainly evaluate our methods on three dynamic datasets, including DDAD driving dataset [14], BONN dynamic dataset [55], and TUM dataset [56] (dynamic object split). Note that these datasets contain fast-moving objects, which are much more challenging than the widely-used KITTI [11] and NYUv2 [12] datasets. We assume that the latter two datasets are almost static in this paper, and we also report results on them. All the mentioned self-supervised methods are trained on each dataset individually for a fair comparison. Moreover, following previous methods, we analyze the depth results at object boundaries and plane regions in the IBims-1 dataset [57]. In the following paragraphs, the details of each dataset are described.

DDAD. The dataset contains 200 driving videos that are captured in urban scenes. The LiDAR scanned point clouds are provided, which we use to generate sparse ground-truth depths for evaluation. In this dataset, almost vehicles are moving on

the road, and there are fewer stopping cars than KITTI, making it more challenging to train self-supervised models. We use the standard training/testing split, which has 150 training scenes (12650 images) and 50 validation scenes (3950 images). We use the validation scenes for evaluation. Depth ranges are capped to at most 200 meters, and images are resized to the resolution of 640×384 for training depth and pose networks.

BONN. The dataset contains 26 dynamic indoor videos that have fast-moving people or other objects. The Kinect captured depth maps are provided as the ground truth for evaluation. We manually find 4 challenging video sequences with fast-moving people (1785 images) for testing, and we use the remaining videos for training. Depth ranges are capped at 10 meters, and images are resized to the resolution of 320×256 for training networks.

TUM. The dataset provides a collection of indoor videos with Kinect-captured depth maps as the ground truth. We choose only videos that belong to the *Dynamic Objects* category, making sure that the model is trained in dynamic scenes. There are in total 11 sequences, and we use the last two sequences that contain moving people (1375 images) for testing. The remaining 9 dynamic videos are used for training, and images are resized to the resolution of 320×256 for feeding to networks.

KITTI. The dataset provides driving videos in urban scenes, and it is the most widely-used dataset in self-supervised monocular depth estimation problems. Following previous work [1], [10], [13], [14], we use the Eigen's split that has 697 images for testing, and we use the remaining video sequences for training. Depth ranges are capped at 80 meters, and images are resized to the resolution of 832×256 for training networks. Note that KITTI contains a large number of stopping cars that help self-supervised methods learn depth estimation on cars, so the results on this dataset cannot reflect our main contributions, *i.e.*, robust learning of monocular depth from dynamic scenes.

NYUv2. The dataset provides a large collection of indoor videos, and it is widely-used in the computer vision community. There are 654 testing images of static scenes for depth evaluation, and we use the remaining videos that do not contain testing images for training neural networks. Images are resized to the resolution of 320×256 before feeding to the network. Note that this dataset contains almost-static scenes.

IBims-1. The dataset provides 100 accurate and dense ground truths for analyzing depth details, including object boundaries and planes. Images are collected in different kinds of indoor environments, and it does not provide a training set. For a fair comparison with previous work, we use the model trained on the NYUv2 dataset for all methods.

Depth Evaluation Metrics. We use standard depth evaluation metrics, including mean absolute relative error (AbsRel), root mean squared error (RMS), root mean squared log error (RM-Slog), and the accuracy under threshold ($\delta_i < 1.25^i$, $i = 1, 2, 3$). The detailed definition of these depth metrics can be found in [2]. Besides, following previous work [1], [10], we multiply the predicted depth maps by a scalar that matches the median with that of the ground truth for evaluation, *i.e.*, $s = \text{median}(D_{gt})/\text{median}(D_{pred})$, since self-supervised methods cannot recover the metric scale. For the evaluation on iBims,

TABLE 1

Self-supervised monocular depth estimation results on the DDAD driving dataset [14]. We segment vehicles and pedestrians as dynamic objects and consider the remaining regions as static backgrounds. This dataset is more challenging than KITTI due to more complex scenes, fewer stopping cars, and longer depth ranges (200m vs 80m). Note that DynamicDepth [58] uses two frames for depth estimation.

| Methods | Full Image | | | | | | | Dynamic | | Static | |
|-------------------|--------------|--------------|---------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|
| | AbsRel | SqRel | RMS | RMSlog | δ_1 | δ_2 | δ_3 | AbsRel | δ_1 | AbsRel | δ_1 |
| Monodepth2 [13] | 0.239 | 12.547 | 18.392 | 0.316 | 0.752 | 0.899 | 0.949 | 0.747 | 0.432 | 0.188 | 0.771 |
| PackNet [14] | 0.182 | 7.945 | 15.021 | 0.259 | 0.828 | 0.925 | 0.961 | 0.564 | 0.520 | 0.137 | 0.843 |
| SGDepth [18] | 0.200 | 7.944 | 17.149 | 0.289 | 0.769 | 0.911 | 0.957 | 0.619 | 0.446 | 0.170 | 0.786 |
| DynamicDepth [58] | 0.156 | 3.305 | 15.612 | 0.258 | 0.785 | 0.914 | 0.962 | 0.258 | 0.612 | 0.149 | 0.792 |
| SC-Depth [1] | 0.169 | 3.877 | 16.290 | 0.280 | 0.773 | 0.905 | 0.951 | 0.345 | 0.546 | 0.155 | 0.783 |
| Ours w/o DRR | 0.153 | 3.124 | 15.237 | 0.252 | 0.799 | 0.920 | 0.963 | 0.259 | 0.612 | 0.146 | 0.806 |
| Ours w/o LSR | 0.149 | 3.094 | 16.198 | 0.262 | 0.794 | 0.913 | 0.956 | 0.210 | 0.666 | 0.146 | 0.799 |
| Ours | 0.142 | 3.031 | 15.868 | 0.248 | 0.813 | 0.922 | 0.963 | 0.199 | 0.697 | 0.140 | 0.813 |

TABLE 2

Self-supervised monocular depth estimation results on the BONN dynamic dataset [57]. This dataset is super-challenging because all training and testing videos contain fast-moving objects, which occupy a large proportion of pixels.

| Methods | Full Image | | | | | Dynamic | | Static | |
|-----------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|
| | AbsRel | RMS | δ_1 | δ_2 | δ_3 | AbsRel | δ_1 | AbsRel | δ_1 |
| Monodepth2 [13] | 0.565 | 2.337 | 0.352 | 0.591 | 0.728 | 0.474 | 0.172 | 0.594 | 0.383 |
| SC-Depth [1] | 0.272 | 0.733 | 0.623 | 0.858 | 0.948 | 0.704 | 0.166 | 0.180 | 0.714 |
| SC-DepthV2 [30] | 0.211 | 0.619 | 0.714 | 0.873 | 0.936 | 0.488 | 0.247 | 0.152 | 0.803 |
| Ours w/o DRR | 0.138 | 0.396 | 0.885 | 0.951 | 0.974 | 0.248 | 0.690 | 0.106 | 0.939 |
| Ours w/o LSR | 0.130 | 0.382 | 0.874 | 0.951 | 0.977 | 0.274 | 0.613 | 0.097 | 0.937 |
| Ours | 0.126 | 0.379 | 0.889 | 0.961 | 0.980 | 0.220 | 0.720 | 0.102 | 0.931 |

TABLE 3

Self-supervised monocular depth estimation results on the TUM dataset [56]. We use the videos under the category of "Dynamic Objects" for training and testing, in which moving objects occupy a large proportion of pixels in each image.

| Methods | Full Image | | | | | Dynamic | | Static | |
|-----------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|
| | AbsRel | RMS | δ_1 | δ_2 | δ_3 | AbsRel | δ_1 | AbsRel | δ_1 |
| Monodepth2 [13] | 0.312 | 1.408 | 0.474 | 0.793 | 0.905 | 0.431 | 0.348 | 0.262 | 0.526 |
| SC-Depth [1] | 0.257 | 0.283 | 0.616 | 0.814 | 0.909 | 0.512 | 0.274 | 0.176 | 0.715 |
| SC-DepthV2 [30] | 0.223 | 0.282 | 0.643 | 0.862 | 0.932 | 0.283 | 0.494 | 0.206 | 0.686 |
| Ours w/o DRR | 0.185 | 1.163 | 0.744 | 0.889 | 0.970 | 0.272 | 0.593 | 0.161 | 0.775 |
| Ours w/o LSR | 0.195 | 1.498 | 0.715 | 0.864 | 0.899 | 0.264 | 0.575 | 0.174 | 0.759 |
| Ours | 0.163 | 0.265 | 0.797 | 0.882 | 0.937 | 0.165 | 0.796 | 0.171 | 0.780 |

the depth boundary errors (DBE) and planarity errors (PE) are used to evaluate the accuracy of depth boundaries and planarity respectively. The detailed definitions of DBE and PE are in [57].

Evaluation on Static/Dynamic Regions. We use MSeg [75] to generate the semantic segmentation mask of testing images. The model is trained on a composite dataset, so it is able to generate segmentation results for both indoor and outdoor driving scenes. In driving datasets (*i.e.*, KITTI and DDAD), all vehicle and pedestrian segments are regarded as dynamic objects, and other regions are regarded as static backgrounds. In indoor datasets (*i.e.*, TUM and BONN), we consider all human segments as dynamic

regions. Note that we align the global scale to the ground-truth depth first, and then we evaluate depth accuracy on static regions, dynamic regions, and full images, individually.

4.2 Evaluation Results

Results on Dynamic Datasets. We use three dynamic datasets mentioned above to evaluate the proposed method, and the quantitative depth estimation results are reported in Tab. 1, 2, and 3, respectively. We show the qualitative comparison results in Fig. 5, and demo videos for depth estimation are in the supplementary. A more detailed analysis is conducted below.

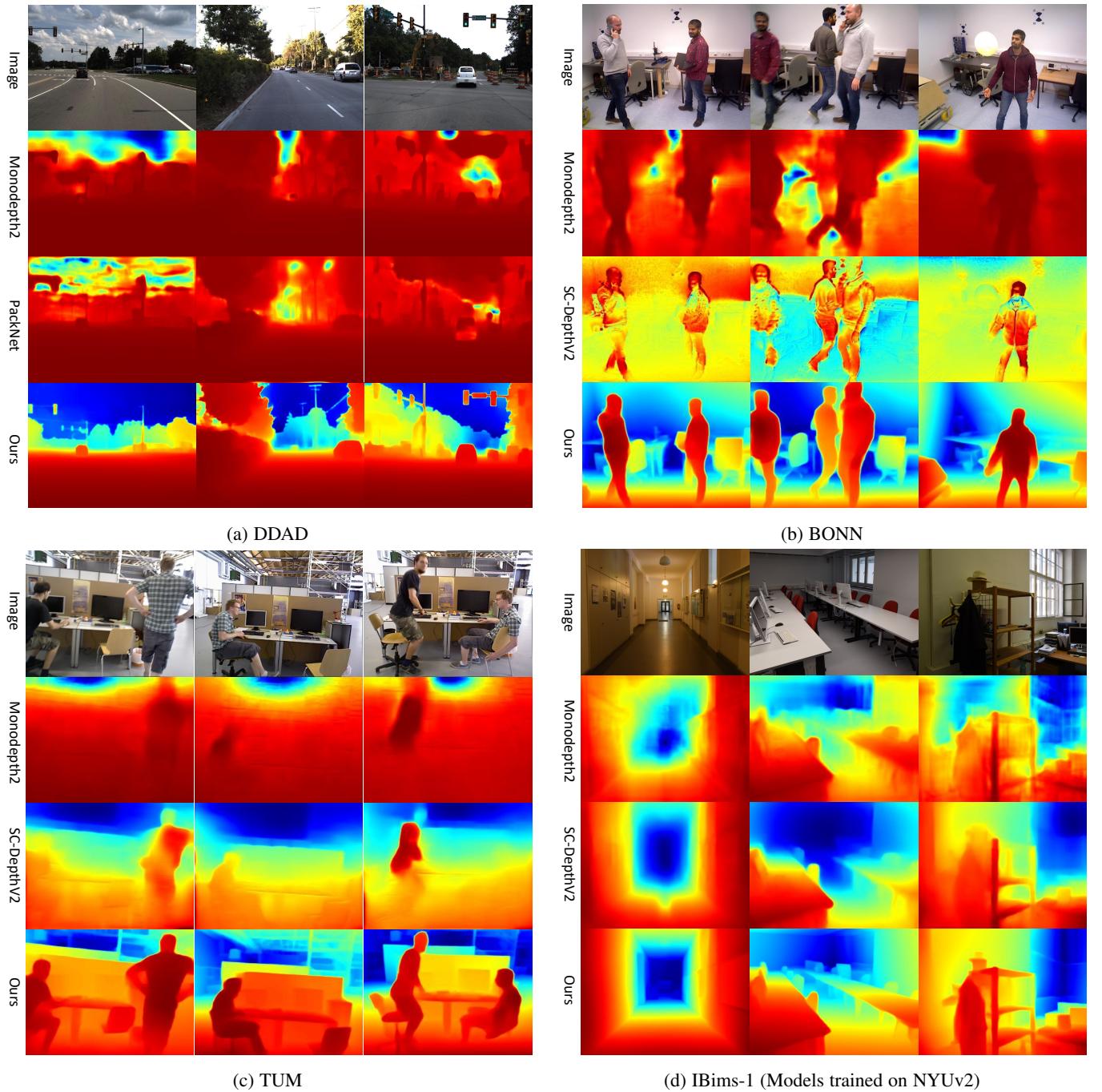


Fig. 5. Qualitative depth estimation results. Existing methods show poor results in dynamic scenes (a-c) because they are hard to handle fast-moving objects during training. Even though they show good accuracy in (d), where models are trained in static scenes, the depth is blurred at object boundaries. By contrast, our method predicts sharp and accurate depth robustly.

Tab. 1 shows the results on DDAD dataset, where we compare our method with previous state-of-the-art methods, including Monodepth2 [13], PackNet [14], and SC-Depth [1]. The results show that our method outperforms previous methods by a large margin, and particularly on dynamic regions. Note that our method outperforms PackNet [14], although the latter uses a significantly larger network backbone than ours. This demonstrates our main contribution in this paper, *i.e.*, robust learning of monocular depth in dynamic scenes. Besides, we also report the result without our proposed DRR and LSR modules. Here our baseline method is a modified version of SC-Depth, and it incorporates the advantages

of Monodepth2. The results show that the performance of these models is significantly lower than that of our full model, which demonstrates the efficacy of our proposed losses.

Tab. 2 and Tab. 3 show the depth estimation results on BONN and TUM datasets, respectively. These indoor datasets are more challenging than driving datasets such as DDAD since the ratio of dynamic regions to the full image of the former is significantly larger than that of the latter. Consequently, previous methods such as Monodepth2 [13] and SC-Depth [1] show poor accuracy in BONN and TUM datasets. Compared with these approaches, our method presents significantly better results. This is contributed

TABLE 4

Self-supervised monocular depth estimation results on KITTI [14]. Note that the KITTI dataset has many stopping vehicles that help learn depth on cars, which is not the case of *learning dynamic object depth from dynamic video* that we addressed in this paper. Besides, note that PackNet uses a large backbone, while other methods including ours use the ResNet-18 encoder.

| Methods | Full Image | | | | | | | Dynamic | | Static | |
|-----------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|
| | AbsRel | SqRel | RMS | RMSlog | δ_1 | δ_2 | δ_3 | AbsRel | δ_1 | AbsRel | δ_1 |
| Monodepth2 [13] | 0.114 | 0.848 | 4.986 | 0.198 | 0.869 | 0.956 | 0.980 | 0.187 | 0.731 | 0.104 | 0.884 |
| PackNet [14] | 0.109 | 0.839 | 4.696 | 0.188 | 0.884 | 0.961 | 0.981 | 0.208 | 0.737 | 0.099 | 0.901 |
| SGDepth [18] | 0.111 | 0.857 | 4.739 | 0.189 | 0.884 | 0.962 | 0.982 | 0.209 | 0.728 | 0.101 | 0.899 |
| SC-Depth [1] | 0.118 | 0.870 | 4.997 | 0.196 | 0.860 | 0.956 | 0.981 | 0.242 | 0.698 | 0.108 | 0.878 |
| Ours | 0.118 | 0.756 | 4.709 | 0.188 | 0.864 | 0.960 | 0.984 | 0.205 | 0.703 | 0.108 | 0.881 |

TABLE 5

Monocular depth estimation results on the NYUv2 [12] dataset. Our method outperforms a majority of supervised methods (first row) and all the self-supervised methods (second row).

| Methods | Error ↓ | | Accuracy ↑ | | |
|-------------------------|--------------|--------------|--------------|--------------|--------------|
| | AbsRel | RMS | δ_1 | δ_2 | δ_3 |
| Make3D [59] | 0.349 | 1.214 | 0.447 | 0.745 | 0.897 |
| DepthTransfer [60] | 0.349 | 1.210 | - | - | - |
| Liu et al. [61] | 0.335 | 1.060 | - | - | - |
| Ladicky et al. [62] | - | - | 0.542 | 0.829 | 0.941 |
| Li et al. [63] | 0.232 | 0.821 | 0.621 | 0.886 | 0.968 |
| Roy et al. [64] | 0.187 | 0.744 | - | - | - |
| Wang et al. [65] | 0.220 | 0.745 | 0.605 | 0.890 | 0.970 |
| Eigen et al. [66] | 0.158 | 0.641 | 0.769 | 0.950 | 0.988 |
| Chakrabarti et al. [67] | 0.149 | 0.620 | 0.806 | 0.958 | 0.987 |
| Li et al. [68] | 0.143 | 0.635 | 0.788 | 0.958 | 0.991 |
| DORN [69] | 0.115 | 0.509 | 0.828 | 0.965 | 0.992 |
| VNL [70] | 0.108 | 0.416 | 0.875 | 0.976 | 0.994 |
| Zhou et al. [71] | 0.208 | 0.712 | 0.674 | 0.900 | 0.968 |
| Zhao et al. [72] | 0.189 | 0.686 | 0.701 | 0.912 | 0.978 |
| Monodepth2 [13] | 0.169 | 0.614 | 0.745 | 0.946 | 0.987 |
| SC-Depth [1] | 0.159 | 0.608 | 0.772 | 0.939 | 0.982 |
| P2Net [73] | 0.150 | 0.561 | 0.796 | 0.948 | 0.986 |
| SC-DepthV2 [30] | 0.138 | 0.532 | 0.820 | 0.956 | 0.989 |
| MonoIndoor [74] | 0.134 | 0.526 | 0.823 | 0.958 | 0.989 |
| Ours | 0.123 | 0.486 | 0.848 | 0.963 | 0.991 |

to our proposed losses, which enables our method to learn depth estimation robustly from dynamic videos.

Results on static Datasets. Although our main contribution in this paper is boosting self-supervised monocular depth in dynamic scenes, we show that our method is also working well in almost-static scenes. The results are reported in the widely-used KITTI driving dataset and NYUv2 indoor dataset. Sampled qualitative results are illustrated in Fig. 6.

Tab. 4 shows the depth estimation results on KITTI, where our method is comparable but does not outperform the previous state-of-the-art methods. The reasons are two folds. First, the dataset

TABLE 6
Evaluation of depth boundaries (DBE) and planes (PE) on iBims-1 [57]. All models are trained on NYUv2.

| Method | iBims-1 | | | | |
|-----------------|--------------------------------------|---------------------------------------|--------------------------------------|--------------------------------------|--------------|
| | $\varepsilon_{DBE}^{acc} \downarrow$ | $\varepsilon_{DBE}^{comp} \downarrow$ | $\varepsilon_{PE}^{plan} \downarrow$ | $\varepsilon_{PE}^{orie} \downarrow$ | AbsRel↓ |
| Monodepth2 [13] | 4.269 | 89.771 | 10.943 | 29.327 | 0.202 |
| SC-DepthV2 [30] | 4.206 | 69.846 | 7.049 | 23.109 | 0.172 |
| Ours w/o LSR | 3.138 | 65.692 | 3.684 | 14.696 | 0.152 |
| Ours | 3.001 | 48.047 | 2.701 | 13.372 | 0.146 |

TABLE 7
Ablation studies of the proposed DRR on DDAD dataset. RS denotes random sampling used in [26], and RL denotes ranking loss used in [26]. The decreased performance demonstrates the effectiveness of our proposed methods.

| Methods | Full | | Dynamic | | Static | |
|--------------|--------------|--------------|--------------|--------------|--------------|--------------|
| | AbsRel | δ_1 | AbsRel | δ_1 | AbsRel | δ_1 |
| Baseline | 0.179 | 0.753 | 0.355 | 0.536 | 0.163 | 0.761 |
| B+DRR (Ours) | 0.149 | 0.794 | 0.210 | 0.666 | 0.146 | 0.799 |
| DRR w/ RS | 0.154 | 0.785 | 0.219 | 0.654 | 0.151 | 0.790 |
| DRR w/ RL | 0.159 | 0.767 | 0.214 | 0.659 | 0.159 | 0.765 |

contains a large number of stopping cars that help self-supervised methods learn depth on vehicles, so our method is hard to further improve the performance when previous methods have obtained good results on dynamic regions. Second, PackNet [14] uses a large network backbone, while other methods, including ours, use ResNet-18, which is much smaller than the former. Overall, we argue that the existing methods have reached a bottleneck in the KITTI dataset, and due to the low impact of dynamic objects on self-supervised learning here, our method is hard to further improve the performance. Moreover, we show qualitative results in Fig. 6 (a), which shows that our method generates sharper depth maps than other methods.

Tab. 5 shows the depth results on NYUv2, where we compare our method with previous state-of-the-art methods such as SC-DepthV2 [30] and MonoIndoor [74]. The results show that our method outperforms previous approaches significantly. This is mainly contributed to the single-image depth prior, which we use to constrain the normal smoothness and sharp object boundaries

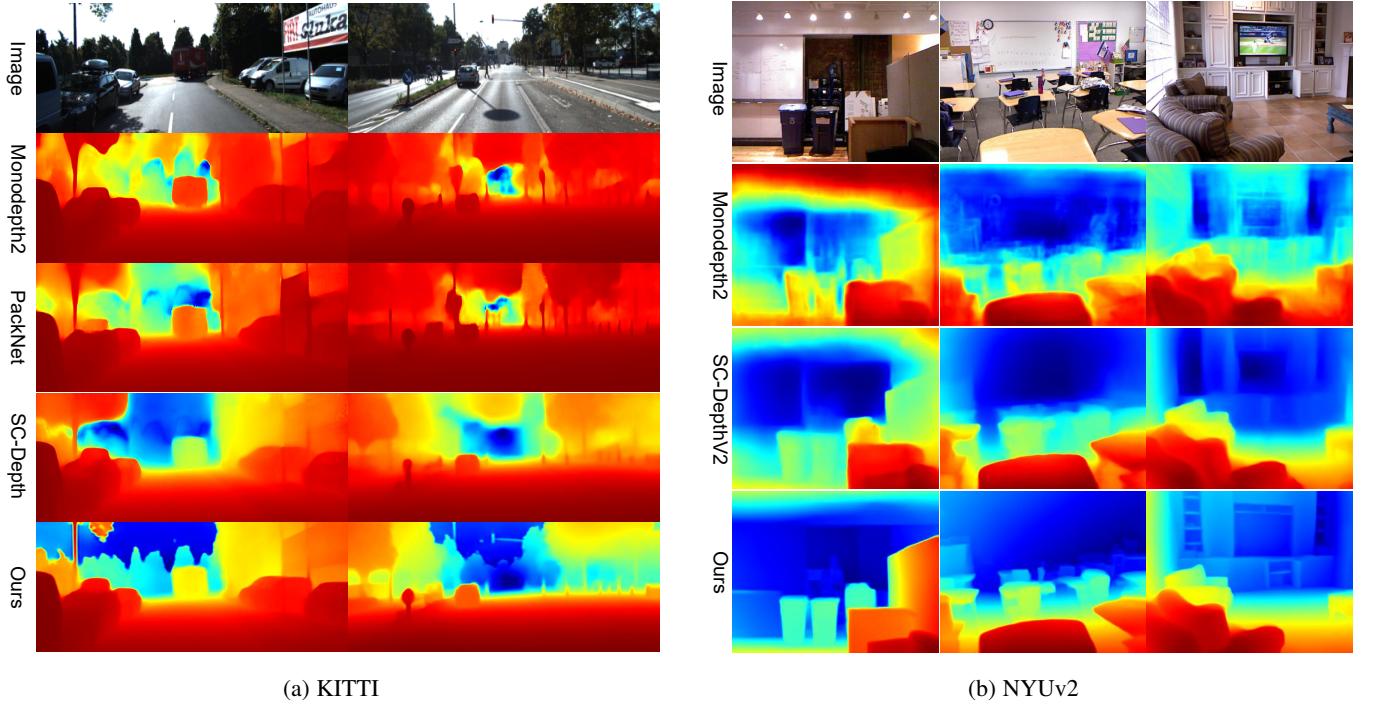


Fig. 6. Qualitative monocular depth estimation results on static datasets. Our method allows for generating sharper depth maps than previous methods—See object boundaries.

TABLE 8
ad LSR an

Ablation studies of the proposed LSR on DDAD dataset. EDS denotes edge-aware depth smoothness [1], and RS denotes additional random sampling beside edge-based sampling. The decreased performance demonstrates the importance of our proposed methods.

| Methods | Full | | Dynamic | | Static | |
|------------|--------------|--------------|--------------|--------------|--------------|--------------|
| | AbsRel | δ_1 | AbsRel | δ_1 | AbsRel | δ_1 |
| Baseline | 0.179 | 0.753 | 0.355 | 0.536 | 0.163 | 0.761 |
| LSR (Ours) | 0.142 | 0.813 | 0.199 | 0.697 | 0.140 | 0.813 |
| LSR w/ EDS | 0.148 | 0.793 | 0.200 | 0.694 | 0.145 | 0.796 |
| LSR w/ RS | 0.146 | 0.802 | 0.200 | 0.688 | 0.143 | 0.806 |

TABLE 9

Evaluation results on DDAD dataset. We compare different methods for generating pseudo-depth. “+Self” means training models with our proposed self-supervised method.

| Methods | Full Image | | Dynamic | | Static | |
|---------------------------|--------------|--------------|--------------|--------------|--------------|--------------|
| | AbsRel | δ_1 | AbsRel | δ_1 | AbsRel | δ_1 |
| DPT [44] | 0.224 | 0.632 | 0.296 | 0.492 | 0.220 | 0.636 |
| DPT+Self | 0.151 | 0.788 | 0.218 | 0.662 | 0.147 | 0.791 |
| LeReS(Res50) [25] | 0.385 | 0.411 | 0.354 | 0.380 | 0.390 | 0.402 |
| LeReS(Res50)+Self | 0.147 | 0.797 | 0.188 | 0.726 | 0.145 | 0.798 |
| LeReS(Res101) [25] | 0.358 | 0.434 | 0.341 | 0.386 | 0.363 | 0.424 |
| LeReS(Res101)+Self (Ours) | 0.142 | 0.813 | 0.199 | 0.697 | 0.140 | 0.813 |

4.3 Ablation Studies

We have shown results with and without our proposed DRR and LSR in Tab. 1, 2 and 3. The results demonstrate the efficacy of the proposed modules. In this section, we make a more detailed analysis of the proposed methods, and we also discuss the performance by using different methods to generate pseudo-depth.

Dynamic Region Refinement. The proposed DRR module consists of dynamic-focused sampling and confident depth ranking loss. We make ablation studies by comparing our method with random sampling (RS) and original ranking loss (RL) that are used in [26]. Tab. 7 shows the evaluation results, which show that the performance is significantly degraded when replacing our proposed terms with the existing methods. This demonstrates the efficacy of our proposed methods.

Local Structure Refinement. The proposed LSR module consists of normal matching loss and edge-guided relative normal ranking loss. The ablation study results are summarized in Tab. 8. We re-

and Fig. 5 (d), and the quantitative evaluation results on object boundaries are summarized in Tab. 6,

Depth Quality at Object Boundaries. Tab. 6 shows the detailed analysis of depth results on the IBims-1 dataset, where we compare our method with Monodepth2 [13] and SC-DepthV2 [30]. All models are trained on NYUv2 for a fair comparison. The AbsRel metric shows the overall accuracy of depth estimation results on full images, and other metrics reflect the detailed depth quality at object boundaries and plane regions. The results show that our method significantly outperforms previous methods. We also remove the proposed LSR from our full model for ablation study purposes, and the results in Tab. 6 show that the performance is clearly degraded. This demonstrates the efficacy of our proposed LSR module. The qualitative depth estimation results on the IBims-1 dataset are illustrated in Fig. 5 (d).

place the normal matching loss with edge-aware depth smoothness loss (EDS), and we also add random sampling (RS) to the edge-guided sampling. These variants degenerate the depth accuracy, which demonstrates that our proposed methods are better than existing solutions.

Pseudo-depth. We use LeReS [25] (ResNet-101) in this paper for generating pseudo-depth, while it is also possible to use other monocular depth estimation networks. Tab. 9 shows the ablation study results on DDAD dataset, where we also include DPT [44] and ResNet-50 version of LeReS. The results show that the pseudo-depths that are generated by all three variants are not accurate in the DDAD dataset. However, when applying our proposed method that uses pseudo-depth for training self-supervised models, high-accuracy depth estimation results can be obtained. This demonstrates that our proposed method is not limited to one specific method for generating pseudo-depth. The results also show that our method with LeReS (ResNet-101) [25] outperforms other variants, including DPT. We hypothesize that the reason is that our method incorporates the normal information in training which is shared with LeReS but not with DPT.

Discussion. We use pseudo-depth to boost self-supervised monocular depth estimation, which somewhat degrades our claim of self-supervised learning. However, in practice, the monocular depth estimation models such as [25], [44] are only trained once in large-scale datasets and can be used as off-the-shelf tools in new unseen scenes, so our method has almost no extra cost compared to pure self-supervised depth estimation methods [1], [13].

5 CONCLUSION

We propose SC-DepthV3 for robust self-supervised learning of monocular depth from challenging dynamic videos. The key to our method is that we use pseudo-depth, which is generated by a pretrained monocular depth estimation network, for addressing the challenges in self-supervised monocular depth estimation framework. More specifically, we address the issues of dynamic objects and blurred object boundaries. As a result, our proposed method can predict sharp and accurate depth maps, even when the model is trained from highly dynamic videos. We comprehensively evaluate our method on six challenging datasets, including both dynamic and static scenes. The results show that our method significantly outperforms previous alternatives, and the ablation study results demonstrate that the proposed modules are effective.

ACKNOWLEDGMENTS

This work was in part supported by National Key R&D Program of China (No. 2022ZD0118700). This work was in part supported by the Australian Centre of Excellence for Robotic Vision CE140100016, and the ARC Laureate Fellowship FL130100102 to I. Reid. We thank anonymous reviewers for their valuable suggestions.

REFERENCES

- [1] J.-W. Bian, H. Zhan, N. Wang, Z. Li, L. Zhang, C. Shen, M.-M. Cheng, and I. Reid, “Unsupervised scale-consistent depth learning from video,” *Int. J. Comput. Vis.*, 2021.
- [2] D. Eigen, C. Puhrsch, and R. Fergus, “Depth map prediction from a single image using a multi-scale deep network,” in *Adv. Neural Inform. Process. Syst.*, 2014.
- [3] S.-J. Park, K.-S. Hong, and S. Lee, “Rdfnet: Rgb-d multi-level residual feature fusion for indoor semantic segmentation,” in *Int. Conf. Comput. Vis.*, 2017, pp. 4980–4989.
- [4] T. Zhou, D.-P. Fan, M.-M. Cheng, J. Shen, and L. Shao, “Rgb-d salient object detection: A survey,” *Computational Visual Media*, pp. 1–33, 2021.
- [5] R. A. Newcombe, S. Izadi, O. Hilliges, D. Molnyneaux, D. Kim, A. J. Davison, P. Kohi, J. Shotton, S. Hodges, and A. Fitzgibbon, “Kinectfusion: Real-time dense surface mapping and tracking,” in *IEEE international symposium on mixed and augmented reality*. IEEE, 2011, pp. 127–136.
- [6] J.-W. Bian, H. Zhan, and I. Reid, “Nvss: High-quality novel view selfie synthesis,” in *International Conference on 3D Vision (3DV)*, 2021.
- [7] H. Zhan, C. S. Weerasekera, J. Bian, and I. Reid, “Visual odometry revisited: What should be learnt?” *arXiv preprint arXiv:1909.09803*, 2019.
- [8] F. Liu, C. Shen, G. Lin, and I. Reid, “Learning depth from single monocular images using deep convolutional neural fields,” *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 38, no. 10, 2016.
- [9] R. Garg, V. K. BG, G. Carneiro, and I. Reid, “Unsupervised cnn for single view depth estimation: Geometry to the rescue,” in *Eur. Conf. Comput. Vis.* Springer, 2016.
- [10] T. Zhou, M. Brown, N. Snavely, and D. G. Lowe, “Unsupervised learning of depth and ego-motion from video,” in *IEEE Conf. Comput. Vis. Pattern Recog.*, 2017.
- [11] A. Geiger, P. Lenz, C. Stiller, and R. Urtasun, “Vision meets Robotics: The kitti dataset,” *International Journal of Robotics Research (IJRR)*, 2013.
- [12] N. Silberman, D. Hoiem, P. Kohli, and R. Fergus, “Indoor segmentation and support inference from rgbd images,” in *Eur. Conf. Comput. Vis.*, 2012.
- [13] C. Godard, O. Mac Aodha, M. Firman, and G. J. Brostow, “Digging into self-supervised monocular depth prediction,” in *Int. Conf. Comput. Vis.*, 2019.
- [14] V. Guizilini, R. Ambrus, S. Pillai, A. Raventos, and A. Gaidon, “3d packing for self-supervised monocular depth estimation,” in *IEEE Conf. Comput. Vis. Pattern Recog.*, 2020.
- [15] V. Casser, S. Pirk, R. Mahjourian, and A. Angelova, “Depth prediction without the sensors: Leveraging structure for unsupervised learning from monocular videos,” in *AAAI*, 2019.
- [16] A. Gordon, H. Li, R. Jonschkowski, and A. Angelova, “Depth from videos in the wild: Unsupervised monocular depth learning from unknown cameras,” in *Int. Conf. Comput. Vis.*, 2019.
- [17] V. Guizilini, R. Hou, J. Li, R. Ambrus, and A. Gaidon, “Semantically-guided representation learning for self-supervised monocular depth,” in *Int. Conf. Learn. Represent.*, 2020.
- [18] M. Klingner, J.-A. Termöhlen, J. Mikolajczyk, and T. Fingscheidt, “Self-supervised monocular depth estimation: Solving the dynamic object problem by semantic guidance,” in *Eur. Conf. Comput. Vis.* Springer, 2020, pp. 582–600.
- [19] Z. Yin and J. Shi, “GeoNet: Unsupervised learning of dense depth, optical flow and camera pose,” in *IEEE Conf. Comput. Vis. Pattern Recog.*, 2018.
- [20] Y. Zou, Z. Luo, and J.-B. Huang, “DF-Net: Unsupervised joint learning of depth and flow using cross-task consistency,” in *Eur. Conf. Comput. Vis.*, 2018.
- [21] A. Ranjan, V. Jampani, K. Kim, D. Sun, J. Wulff, and M. J. Black, “Competitive Collaboration: Joint unsupervised learning of depth, camera motion, optical flow and motion segmentation,” in *IEEE Conf. Comput. Vis. Pattern Recog.*, 2019.
- [22] Y. Chen, C. Schmid, and C. Sminchisescu, “Self-supervised learning with geometric constraints in monocular video: Connecting flow, depth, and camera,” in *Int. Conf. Comput. Vis.*, 2019, pp. 7063–7072.
- [23] S. Lee, S. Im, S. Lin, and I. S. Kweon, “Learning monocular depth in dynamic scenes via instance-aware projection consistency,” in *Proceedings of the AAAI Conference on Artificial Intelligence (AAAI)*, 2021.
- [24] H. Li, A. Gordon, H. Zhao, V. Casser, and A. Angelova, “Unsupervised monocular depth learning in dynamic scenes,” in *Conference on Robot Learning*, 2020.
- [25] W. Yin, J. Zhang, O. Wang, S. Niklaus, L. Mai, S. Chen, and C. Shen, “Learning to recover 3d scene shape from a single image,” in *IEEE Conf. Comput. Vis. Pattern Recog.*, 2021, pp. 204–213.
- [26] W. Chen, Z. Fu, D. Yang, and J. Deng, “Single-image depth perception in the wild,” *Adv. Neural Inform. Process. Syst.*, 2016.
- [27] C. Godard, O. Mac Aodha, and G. J. Brostow, “Unsupervised monocular depth estimation with left-right consistency,” in *IEEE Conf. Comput. Vis. Pattern Recog.*, 2017.
- [28] H. Zhan, R. Garg, C. Saroj Weerasekera, K. Li, H. Agarwal, and I. Reid, “Unsupervised learning of monocular depth estimation and visual

- odometry with deep feature reconstruction," in *IEEE Conf. Comput. Vis. Pattern Recog.*, 2018.
- [29] R. Mahjourian, M. Wicke, and A. Angelova, "Unsupervised learning of depth and ego-motion from monocular video using 3d geometric constraints," in *IEEE Conf. Comput. Vis. Pattern Recog.*, 2018.
- [30] J.-W. Bian, H. Zhan, N. Wang, T.-J. Chin, C. Shen, and I. Reid, "Auto-rectify network for unsupervised indoor depth estimation," *IEEE Trans. Pattern Anal. Mach. Intell.*, 2021.
- [31] R. Wang, Z. Yu, and S. Gao, "Planedepth: Self-supervised depth estimation via orthogonal planes," in *IEEE Conf. Comput. Vis. Pattern Recog.*, June 2023, pp. 21425–21434.
- [32] H. Si, B. Zhao, D. Wang, Y. Gao, M. Chen, Z. Wang, and X. Li, "Fully self-supervised depth estimation from defocus clue," in *IEEE Conf. Comput. Vis. Pattern Recog.*, 2023, pp. 9140–9149.
- [33] A. Bangunharcana, A. M. Aly, and K. Kim, "Dualrefine: Self-supervised depth and pose estimation through iterative epipolar sampling and refinement toward equilibrium," in *IEEE Conf. Comput. Vis. Pattern Recog.* IEEE, 2023.
- [34] K. Xian, C. Shen, Z. Cao, H. Lu, Y. Xiao, R. Li, and Z. Luo, "Monocular relative depth perception with web stereo data supervision," in *IEEE Conf. Comput. Vis. Pattern Recog.*, 2018, pp. 311–320.
- [35] Z. Li and N. Snavely, "Megadepth: Learning single-view depth prediction from internet photos," in *IEEE Conf. Comput. Vis. Pattern Recog.*, 2018, pp. 2041–2050.
- [36] Z. Li, T. Dekel, F. Cole, R. Tucker, N. Snavely, C. Liu, and W. T. Freeman, "Learning the depths of moving people by watching frozen people," in *IEEE Conf. Comput. Vis. Pattern Recog.*, 2019, pp. 4521–4530.
- [37] C. Wang, S. Lucey, F. Perazzi, and O. Wang, "Web stereo video supervision for depth prediction from dynamic scenes," in *International Conference on 3D Vision (3DV)*. IEEE, 2019, pp. 348–357.
- [38] W. Chen, S. Qian, and J. Deng, "Learning single-image depth from videos using quality assessment networks," in *IEEE Conf. Comput. Vis. Pattern Recog.*, 2019, pp. 5604–5613.
- [39] W. Yin, J. Zhang, O. Wang, S. Niklaus, L. Mai, S. Chen, and C. Shen, "Learning to recover 3d scene shape from a single image," in *IEEE Conf. Comput. Vis. Pattern Recog.*, 2020.
- [40] R. Ranftl, K. Lasinger, D. Hafner, K. Schindler, and V. Koltun, "Towards robust monocular depth estimation: Mixing datasets for zero-shot cross-dataset transfer," *IEEE Trans. Pattern Anal. Mach. Intell.*, 2020.
- [41] H. Hirschmuller, "Accurate and efficient stereo processing by semi-global matching and mutual information," in *IEEE Conf. Comput. Vis. Pattern Recog.*, vol. 2, 2005, pp. 807–814.
- [42] J. L. Schönberger, E. Zheng, M. Pollefeys, and J.-M. Frahm, "Pixelwise view selection for unstructured multi-view stereo," in *Eur. Conf. Comput. Vis.*, 2016.
- [43] D. J. Butler, J. Wulff, G. B. Stanley, and M. J. Black, "A naturalistic open source movie for optical flow evaluation," in *Eur. Conf. Comput. Vis.*, Oct. 2012, pp. 611–625.
- [44] R. Ranftl, A. Bochkovskiy, and V. Koltun, "Vision transformers for dense prediction," *ArXiv preprint*, 2021.
- [45] X. Luo, J.-B. Huang, R. Szeliski, K. Matzen, and J. Kopf, "Consistent video depth estimation," *ACM Transactions on Graphics (Proceedings of ACM SIGGRAPH)*, 2020.
- [46] Z. Wang, A. C. Bovik, H. R. Sheikh, E. P. Simoncelli *et al.*, "Image Quality Assessment: from error visibility to structural similarity," *IEEE Trans. Image Process.*, vol. 13, no. 4, 2004.
- [47] K. Xian, J. Zhang, O. Wang, L. Mai, Z. Lin, and Z. Cao, "Structure-guided ranking loss for single image depth prediction," in *IEEE Conf. Comput. Vis. Pattern Recog.*, 2020.
- [48] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *IEEE Conf. Comput. Vis. Pattern Recog.*, 2016, pp. 770–778.
- [49] O. Ronneberger, P. Fischer, and T. Brox, "U-net: Convolutional networks for biomedical image segmentation," in *International Conference on Medical Image Computing and Computer-Assisted Intervention (MICCAI)*. Springer, 2015.
- [50] D.-A. Clevert, T. Unterthiner, and S. Hochreiter, "Fast and accurate deep network learning by exponential linear units (elus)," *arXiv preprint arXiv:1511.07289*, 2015.
- [51] A. Paszke, S. Gross, S. Chintala, G. Chanan, E. Yang, Z. DeVito, Z. Lin, A. Desmaison, L. Antiga, and A. Lerer, "Automatic differentiation in pytorch," in *NIPS-W*, 2017.
- [52] C. Wang, J. Miguel Buenaposada, R. Zhu, and S. Lucey, "Learning depth from monocular videos using direct methods," in *IEEE Conf. Comput. Vis. Pattern Recog.*, June 2018.
- [53] D. P. Kingma and J. Ba, "ADAM: A method for stochastic optimization," *arXiv preprint arXiv:1412.6980*, 2014.
- [54] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, "ImageNet: A Large-Scale Hierarchical Image Database," in *IEEE Conf. Comput. Vis. Pattern Recog.*, 2009.
- [55] E. Palazzolo, J. Behley, P. Lottes, P. Giguère, and C. Stachniss, "ReFusion: 3D Reconstruction in Dynamic Environments for RGB-D Cameras Exploiting Residuals," in *Int. Conf. Intelligent Robots and Systems*, 2019.
- [56] J. Sturm, N. Engelhard, F. Endres, W. Burgard, and D. Cremers, "A benchmark for the evaluation of rgbd slam systems," in *Int. Conf. Intelligent Robots and Systems*, 2012.
- [57] T. Koch, L. Liebel, M. Körner, and F. Fraundorfer, "Comparison of monocular depth estimation methods using geometrically relevant metrics on the ibims-1 dataset," *Computer Vision and Image Understanding*, 2020.
- [58] Z. Feng, L. Yang, L. Jing, H. Wang, Y. Tian, and B. Li, "Disentangling object motion and occlusion for unsupervised multi-frame monocular depth," in *Eur. Conf. Comput. Vis.* Springer, 2022, pp. 228–244.
- [59] A. Saxena, S. H. Chung, and A. Y. Ng, "Learning depth from single monocular images," in *Adv. Neural Inform. Process. Syst.*, 2006.
- [60] K. Karsch, C. Liu, and S. B. Kang, "Depth transfer: Depth extraction from video using non-parametric sampling," *IEEE Trans. Pattern Anal. Mach. Intell.*, 2014.
- [61] M. Liu, M. Salzmann, and X. He, "Discrete-continuous depth estimation from a single image," in *IEEE Conf. Comput. Vis. Pattern Recog.*, 2014.
- [62] L. Ladicky, J. Shi, and M. Pollefeys, "Pulling things out of perspective," in *IEEE Conf. Comput. Vis. Pattern Recog.*, 2014.
- [63] B. Li, C. Shen, Y. Dai, A. Van Den Hengel, and M. He, "Depth and surface normal estimation from monocular images using regression on deep features and hierarchical crfs," in *IEEE Conf. Comput. Vis. Pattern Recog.*, 2015.
- [64] A. Roy and S. Todorovic, "Monocular depth estimation using neural regression forest," in *IEEE Conf. Comput. Vis. Pattern Recog.*, 2016.
- [65] P. Wang, X. Shen, Z. Lin, S. Cohen, B. Price, and A. L. Yuille, "Towards unified depth and semantic prediction from a single image," in *IEEE Conf. Comput. Vis. Pattern Recog.*, 2015.
- [66] D. Eigen and R. Fergus, "Predicting depth, surface normals and semantic labels with a common multi-scale convolutional architecture," in *Int. Conf. Comput. Vis.*, 2015.
- [67] A. Chakrabarti, J. Shao, and G. Shakhnarovich, "Depth from a single image by harmonizing overcomplete local network predictions," in *Adv. Neural Inform. Process. Syst.*, 2016.
- [68] J. Li, R. Klein, and A. Yao, "A two-streamed network for estimating fine-scaled depth maps from single rgb images," in *Int. Conf. Comput. Vis.*, 2017.
- [69] H. Fu, M. Gong, C. Wang, K. Batmanghelich, and D. Tao, "Deep ordinal regression network for monocular depth estimation," in *IEEE Conf. Comput. Vis. Pattern Recog.*, 2018, pp. 2002–2011.
- [70] W. Yin, Y. Liu, C. Shen, and Y. Yan, "Enforcing geometric constraints of virtual normal for depth prediction," in *Int. Conf. Comput. Vis.*, 2019.
- [71] J. Zhou, Y. Wang, K. Qin, and W. Zeng, "Moving indoor: Unsupervised video depth learning in challenging environments," in *Int. Conf. Comput. Vis.*, 2019.
- [72] W. Zhao, S. Liu, Y. Shu, and Y.-J. Liu, "Towards better generalization: Joint depth-pose learning without posenet," in *IEEE Conf. Comput. Vis. Pattern Recog.*, 2020.
- [73] Z. Yu, L. Jin, and S. Gao, "P²net: Patch-match and plane-regularization for unsupervised indoor depth estimation," in *Eur. Conf. Comput. Vis.*, 2020.
- [74] P. Ji, R. Li, B. Bhanu, and Y. Xu, "Monoindoor: Towards good practice of self-supervised monocular depth estimation for indoor environments," in *IEEE Conf. Comput. Vis. Pattern Recog.*, 2021.
- [75] J. Lambert, Z. Liu, O. Sener, J. Hays, and V. Koltun, "Mseg: A composite dataset for multi-domain semantic segmentation," in *IEEE Conf. Comput. Vis. Pattern Recog.*, 2020, pp. 2879–2888.