# Exploring the Mutual Influence between Self-Supervised Single-Frame and Multi-Frame Depth Estimation

Jie Xiang[1,2], Yun Wang[1], Lifeng An[1], Haiyang Liu[1] and Jian Liu[1]

arXiv:2304.12685v2 [cs.CV] 28 Aug 2023

*Abstract*—Although both self-supervised single-frame and multi-frame depth estimation methods only require unlabeled monocular videos for training, the information they leverage varies because single-frame methods mainly rely on appearance-based features while multi-frame methods focus on geometric cues. Considering the complementary information of single-frame and multi-frame methods, some works attempt to leverage single-frame depth to improve multi-frame depth. However, these methods can neither exploit the difference between single-frame depth and multi-frame depth to improve multi-frame depth nor leverage multi-frame depth to optimize single-frame depth models. To fully utilize the mutual influence between single-frame and multi-frame methods, we propose a novel self-supervised training framework. Specifically, we first introduce a pixel-wise adaptive depth sampling module guided by single-frame depth to train the multi-frame model. Then, we leverage the minimum reprojection based distillation loss to transfer the knowledge from the multi-frame depth network to the single-frame network to improve single-frame depth. Finally, we regard the improved single-frame depth as a prior to further boost the performance of multi-frame depth estimation. Experimental results on the KITTI and Cityscapes datasets show that our method outperforms existing approaches in the self-supervised monocular setting.

*Index Terms*—Deep Learning for Visual Perception; Visual Learning; Deep Learning Methods.

## I. Introduction

DEPTH estimation is an essential and challenging problem in 3D vision, which can be applied in a wide range of applications such as autonomous driving [1] and augmented reality [2]. Although active depth sensors e.g. Lidar and binocular cameras-based methods [3] exist, monocular depth estimation (MDE) methods that only use a single RGB camera to estimate depth still attract much attention due to their flexibility and low cost.

In the past years, many deep learning-based MDE methods [4]–[8] have emerged. Among these methods, self-supervised

methods [6], [8] that use unlabeled monocular video sequences as training data to eliminate the dependence on ground-truth depth made exciting progress. Early self-supervised methods [6]–[8] mainly focus on single-frame depth estimation, which refers to inferring the corresponding depth map given a single image. Whilst flexible, single-frame approaches ignore that more than one frame may be available at test time in many practical applications. Therefore, a few recent works [8]–[11] take multiple frames as input for depth estimation. Different from single-frame methods, these multi-frame methods mainly utilize the geometric matching features between multiple frames. Considering that dense feature matching is easily affected by occlusions, moving objects, and textureless regions, these methods try to utilize the single-frame information to improve multi-frame depth estimation unilaterally and show promising results.

However, these "one-way" methods do not take full advantage of the "mutual" influence between single-frame depth estimation and multi-frame depth estimation due to ignoring the two issues: (1) the potential benefits of the difference between single-frame depth and multi-frame depth, and (2) the effect of multi-frame depth on single-frame depth estimation. Since current multi-frame depth methods [10], [12] usually predict more accurate depth than single-frame methods [7], [13], the output of multi-frame depth models can be regarded as pseudo labels. Then we can model the uncertainty of single-frame depth using the difference between single-frame depth and multi-frame depth. Based on the above idea, we propose the Pixel-wise Adaptive Depth Sampling (PADS) module to determine the depth candidates used for multi-frame depth estimation, in which the single-frame depth is used to determine the geometric center of sampling range following [10], [11], and the difference between single-frame depth and multi-frame depth is used to determine the width of the sampling range. In this way, we improve the efficiency of depth sampling and form effective cost volumes for more accurate multi-frame depth estimation. Regarding the second issue, we adopt the multi-frame depth model as the teacher to train another single-frame depth network via distillation learning. To alleviate the impact that the teacher model generates inaccurate labels, we combine the photometric loss [7] and the distillation loss to form a minimum reprojection based distillation loss, ignoring the pseudo labels with large reprojection errors. Thus, the single-frame depth network also produces better results.

To allow these two ideas to work in a compatible manner, we further propose a novel self-supervised training framework.

Specifically, an uncertainty map is iteratively updated using the PADS module when training the teacher model and then fixed at test time. After distillation learning, the learned uncertainty map and the improved single-frame depth are further regarded as the input of the PADS module to determine the sampling range for cost volume generation so that the multi-frame network also benefits from the improved single-frame network.

In summary, the contributions of the paper are as follows:

- A novel self-supervised distillation learning framework for MDE that fully utilizes the mutual benefits between self-supervised single-frame and multi-frame depth estimation.
- A pixel-wise adaptive depth sampling module to use the single-frame depth and the difference between single-frame depth and multi-frame depth as priors for multi-frame depth estimation.
- A distillation loss based on the minimum reprojection error to filter out the multi-frame depth predictions that may have large errors.
- A new state of the art on the KITTI and Cityscapes datasets in self-supervised monocular depth estimation. The code and models will be available at https://github.com/xjixzz/MISM.

## II. RELATED WORKS

### A. Single-Frame Depth Estimation

Single-frame depth estimation refers to inferring the corresponding pixel-wise depth from a single image, which is an ill-posed problem because there are an infinite number of possible 3D scenes that can correspond to the same image. Early single-frame depth estimation studies [4], [14] focused on supervised methods, which suffer from collecting ground truth depth. To avoid the heavy work of collecting labels, Garg et al. [5] proposed the first self-supervised single-frame depth estimation model supervised by view synthesis loss from rectified stereo image pairs. Zhou et al. [6] extended the self-supervised stereo training into a more general form, i.e., self-supervised monocular training, which jointly estimates depth and poses to form view synthesis loss only using unlabeled monocular videos. Monodepth2 [7] further improved the training loss to alleviate the problems caused by occlusions and stationary pixels and provided a strong baseline for the following works including our method.

Following [7], more powerful or efficient network architectures [12], [15] and more effective data augmentation strategies [13] have been proposed to improve single-frame depth. In addition, other tasks such as flow estimation [16], bird's-eye-view scene layout estimation [17], and semantic segmentation [18], [19] were introduced to provide extra information for single-frame depth estimation. Furthermore, some works [20], [21] attempted to use knowledge distillation to improve the results of depth estimation, and we will review these works in Section II-C.

### B. Multi-Frame Depth Estimation

In contrast to single-frame methods, multi-frame depth estimation methods can exploit consecutive multiple frames at test time. Among them, some works e.g. [2] iteratively finetune the pretrained single-frame network at test time for global temporal consistency, which suffers from the running speed. A second group of works e.g. [22] introduce recurrent networks to exploit temporal information for online depth estimation but are limited by implicitly geometric reasoning.

To explicitly reason about the geometry, deep learning-based multi-view stereo (MVS) methods [23]–[25] adopt plane-sweep stereo architectures to build 3D cost volumes from the features of multiple 2D images via differentiable warping. However, these methods focus on simple static scenes and assume that the camera pose is known in advance, which limits their applicability in more complex scenarios.

To infer the depth from multiple images with unknown poses, Watson et al. [8] proposed a self-supervised model with the improved multi-view plane-sweep stereo architecture. Manydepth [8] leverages the estimated single-frame depth to update the minimum and maximum depth values of the whole scene to alleviate the scale ambiguity, and provides supervision for the multi-frame depth network in the region where matching costs do not work. Based on [8], [9] disentangles object motions to overcome the mismatching problem using dynamic category segmentation masks and single-frame depth, and [26] uses an attention-based matching mechanism to improve multi-frame matching for cost volume generation. More recently, [10] leverages the single-frame depth and the magnitude of the estimated velocity as prior information to determine the search space for multi-frame depth and further uses an additional uncertainty-based network to fuse the single-frame depth and multi-frame depth. Similar to [10], [11] also utilizes the single-frame depth as prior depth but fuses the single-frame and multi-frame information in a multi-scale manner. However, none of these works exploit the mutual influence between single-frame and multi-frame depth estimation in a comprehensive way as our method.

### C. Knowledge Distillation

Knowledge distillation was originally proposed to compress a large model into a lightweight model without a large performance drop via the teacher-student architecture, which has been applied to many vision tasks [27], including MDE [45]. As for self-supervised depth estimation, research on applying distillation learning to self-supervised MVS methods [24], [25] or self-supervised stereo training methods [28]–[30] has been going on for several years. However, until recently, a few works [20], [21], [31] started to exploit knowledge distillation to improve single-frame depth estimation in the self-supervised monocular setting. Among these methods, [31] learns two task-dependent uncertainty maps to weight the pseudo label loss and self-supervised photometric loss respectively for more accurate single-frame depth, and [20] selects the optimal prediction from multiple predictions of the multi-stream ensemble network to help train the student network. Nevertheless, both of these two works ignore temporal information. Closest to our model in spirit is the work of Petrovai et al. [21], which leverages a self-distillation training strategy to distill the high-resolution pseudo labels with the 3D consistency filtering

Fig. 1. Overview of our pipeline. First, we train the teacher model consisting of Single-frame Depth Network (S-DepthNet), PoseNet, and Multi-frame Depth Network (M-DepthNet), in which the Pixel-wise Adaptive Depth Sampling (PADS) module generates the hypothesized depth $D_{hypoth}$ for Group-wise Correlation (GwC) based cost volume generation from single-frame depth $D_t^s$ and pixel-wise sampling width adjuster $\delta$, and simultaneously updates $\delta$ with $D_t^m$ and $D_t^s$. Then, the trained teacher model generates pseudo labels to guide the student model with the distillation loss. Last, the trained student model replaces the S-DepthNet and PoseNet of the teacher model to help M-DepthNet to produce better results during inference.

strategy. Instead of leveraging temporal information to distill pseudo labels via post-processing as [21], we directly use the multi-frame depth estimation network to generate pseudo labels for training the single-frame depth network and the trained single-frame depth network is further used to boost the performance of multi-frame depth estimation.

## III. METHOD

In this section, we first introduce the overall pipeline. Then, we introduce the paradigms for self-supervised single-frame and multi-frame depth estimation. After that, we describe the PADS module, which provides effective hypothesized depth for cost volume generation of multi-frame depth estimation. Finally, we introduce the distillation learning to train the student model with the supervision from photometric loss and pseudo labels generated by the teacher model.

### A. Method Overview

To exploit the mutual influence between single-frame and multi-frame depth estimation, we propose a two-stage training pipeline as shown in Fig. 1. In the first stage, we train the Multi-frame Depth Network (M-DepthNet) in a self-supervised manner. Following MOVEDepth [10], we additionally introduce Single-frame Depth Network (S-DepthNet) and Pose Network (PoseNet) when training M-DepthNet. The output of S-DepthNet is used as the input of the PADS module, which generates the pixel-wise depth range for M-DepthNet. And the output of PoseNet is used to construct cost volumes via warping. Thus, the three networks together form the teacher model and can be jointly optimized. In the second stage, the student model consisting of S-DepthNet and PoseNet is trained by combing the supervision of the teacher model and self-supervision, in which pseudo labels only work when they produce small photometric errors. During inference, the output of the student model is used to guide the cost volume generation for M-DepthNet, which helps improve the accuracy of M-DepthNet.

### B. Self-Supervised Single-Frame Depth Estimation

The objective of self-supervised single-frame depth estimation is to minimize the photometric error between the target

image $I_t$ and the synthesized image derived from the predicted target depth map $D_t^s$. As shown in the upper left part of Fig. 1, self-supervised monocular training jointly optimizes S-DepthNet and PoseNet. S-DepthNet takes $I_t$ as input and outputs the corresponding depth map $D_t^s$, while PoseNet takes $I_t$ and the context image $I_c$ as input and estimates the ego-motion $T_{t \to c}$. Assuming that the camera intrinsic matrix $K$ is known, then we can project the 2D pixel coordinates $(u, v)$ of $I_t$ to $(u_c, v_c)$ in $I_c$ as follows:

$$(u_c, v_c) = K T_{t \to c} D_t^s(u, v) K^{-1}(u, v, 1)^T, \quad (1)$$

where the conversion between homogeneous and inhomogeneous coordinates is omitted for notational simplicity. As in [7], we use bilinear interpolation denoted as $\langle \cdot \rangle$ to sample the context pixel $(u_c, v_c)$ to obtain the synthesized target image:

$$I_{c \to t}(u, v) = I_c \langle (u_c, v_c) \rangle. \quad (2)$$

Following [32], we compute the weighted sum of L1 loss and structural similarity (SSIM) to form the photometric loss:

$$PE(I_a, I_b) = \alpha ||I_a - I_b||_1 + (1 - \alpha) \frac{1 - SSIM(I_a, I_b)}{2}, \quad (3)$$

where $\alpha = 0.15$. To address the occlusion problem, we adopt the per-pixel minimum reprojection loss as in [7], i.e.

$$L_p = \min_{c \in \{t-1, t+1\}} PE(I_t, I_{c \to t}). \quad (4)$$

Similar to [7], we also apply the auto-masking strategy to generate the mask $\mu$ for removing the stationary pixels from $L_p$. Following [32], we also use the edge-aware smooth loss:

$$L_{sm} = |\partial_u d_t^*| e^{-|\partial_u I_t|} + |\partial_v d_t^*| e^{-|\partial_v I_t|}, \quad (5)$$

where $d_t^*$ is the mean-normalized inverse depth. Like [7], we also compute the multi-scale photometric loss when multi-scale depth predictions are available. Thus, the final self-supervised loss for single-frame depth $D_t^s$ is formulated as:

$$L_{self}(D_t^s) = \frac{1}{S} \sum_{i=0}^{S-1} \mu L_p + \lambda_{sm} L_{sm}, \quad (6)$$

where $S$ represents the number of multi-scale depth maps and $\lambda_{sm}$ is set to $10^{-3}$ as in [7].

## C. Self-Supervised Multi-Frame Depth Estimation

The diagram of the Multi-frame Depth Network (M-DepthNet) is shown at the bottom of Fig. 1. As in [10], M-DepthNet takes as input two $H \times W \times 3$ images, $I_t$ and $I_c$, and uses a shared Feature Network (Feat-Net) to extract the $h \times w \times C$ features $F_t$ and $F_c$ respectively, where $h = H/4$ and $w = W/4$. Then, similar to (1) and (2), the context feature $F_c$ is warped into the target view to obtain the feature volume $FV_{c \to t}$ according to the estimated relative pose $T_{t \to c}$ and the hypothesized discrete depth candidates $D_{hypoth} \in \mathbb{R}^{N \times h \times w}$, where $N$ is the number of depth candidates for each pixel. Next, group-wise correlation [3] is applied to construct the cost volume $CV_t \in \mathbb{R}^{N \times G \times h \times w}$, where $G$ is the number of groups that $C$-Channel feature volume $FV_{c \to t}$ is divided into. Subsequently, a 3D UNet [33] is used to regularize the cost volume to obtain the probability volume $P_t \in \mathbb{R}^{N \times h \times w}$, and local-max operation [34] is performed to generate the low-resolution depth map $D_t^l \in \mathbb{R}^{h \times w}$ as

$$D_t^l(u,v) = \sum_{i=x-r}^{x+r} D_{hypoth}(i,u,v) \frac{P_t(i,u,v)}{\sum_{j=x-r}^{x+r} P_t(j,u,v)}, \quad (7)$$

where $x$ is the index of the maximum value of the 1D vector $P_t(:,u,v)$ and $r$ is the radius of the local window. Finally, the convex upsampling layer [35] is used to interpolate the $D_t^l$ to output the final multi-frame depth map $D_t^m \in \mathbb{R}^{H \times W}$. We can calculate the self-supervised loss $L_{self}(D_t^m)$ like (6). Thus, the loss function to jointly optimize all networks of the teacher model is formulated as

$$L_{teacher} = L_{self}(D_t^s) + L_{self}(D_t^m). \quad (8)$$

## D. Pixel-wise Adaptive Depth Sampling

As described in the previous subsection, generating cost volume requires sampling the depth candidates $D_{hypoth}$. [8] iteratively updates the depth range $[d_{min}, d_{max}]$ for the whole scene according to the estimated single-view depth during training and fixed the two parameters $d_{min}$ and $d_{max}$ during inference, which is computationally expensive since the learned depth range needs to cover the depths of all viewpoints in the scene. To narrow the depth range for different views, some approaches [10], [11] take the estimated single-view depth as the geometric center for depth sampling, and additionally use one or more predefined hyperparameters to determine the width of the depth range. Furthermore, [10] leverages the magnitude of the velocity estimated by PoseNet to adjust the width of the image-wise depth range but suffers from the challenging task to estimate the absolute scale of the predicted velocity before training. More importantly, all these methods ignore the distribution difference for the width of the sampling range in the pixel space and fail to provide pixel-wise adaptive depth range to build efficient cost volumes.

To better exploit the spatial distribution of scene depth, we propose the PADS module, which adopts a learnable uncertainty map $\delta \in \mathbb{R}^{h \times w}$ to indicate the pixel-wise relative width of the sampling range. All elements in $\delta$ are initialized to one. Following [10], [11], we also use the single-frame depth $D_t^s$ as prior depth to determine the geometric center of the search space for per-pixel depth candidates. Considering the difference in the resolution of the predicted depth map and cost volume, we first downsample $D_t^s$ and $D_t^m$ to obtain $D_t^{s,l}$ and $D_t^{m,l}$ respectively. Let $D_{min}$ and $D_{max}$ denote the minimum and maximum depth map with a resolution of $h \times w$, respectively. Then we can specify the sampling range:

$$\begin{cases} D_{min} = D_t^{s,l}/(1+\delta), \\ D_{max} = D_t^{s,l}(1+\delta). \end{cases} \quad (9)$$

When training M-DepthNet, we adopt the exponential moving average strategy to update $\delta$ according to the difference between $D_t^{s,l}$ and $D_t^{m,l}$:

$$\delta \leftarrow 0.99\delta + 0.01\delta', \quad (10)$$

$$\delta' = \beta(\max(D_t^{s,l}/D_t^{m,l}, D_t^{m,l}/D_t^{s,l}) - 1). \quad (11)$$

Here, $\beta$ is a hyperparameter greater than 1 to avoid the estimated multi-frame depth falling on the boundary or even out of the sampling range. In our setting, $\beta$ is set to 1.2. The learned $\delta$ is visualized as Fig. 3(a), which reflects the estimated uncertainty distribution of single-frame depth for the target scene. Similar to [8], we save $\delta$ as part of the model weights after training and keep $\delta$ fixed during inference.

As in [10], according to $D_{min}$ and $D_{max}$ determined by (9), we then uniformly sample in the inverse depth space to obtain $D_{hypoth}$, i.e.

$$D_{hypoth}(i) = 1/(\frac{i}{N-1}(\frac{1}{D_{min}} - \frac{1}{D_{max}}) + \frac{1}{D_{max}}), \quad (12)$$

where $i = 0, 1, ..., N-1$. Thus, the depth candidates $D_{hypoth}$ are used to generate the pixel-wise adaptive cost volume for multi-frame depth estimation as described in the previous subsection. Compared to the previous sampling strategies [8], [10], the PADS module is capable of adjusting the sampling range at a finer granularity, which helps improve the accuracy of multi-frame depth estimation.

## E. Distillation Learning

Considering the performance gap between single-frame and multi-frame depth networks, we further transfer the knowledge from M-DepthNet to S-DepthNet. As shown in Fig. 1, the student model is composed of S-DepthNet and PoseNet, and the outputs of S-DepthNet and PoseNet are $\widetilde{D}_t^s$ and $\widetilde{T}_{t \to c}$ respectively. The teacher model generates pseudo labels $D_t^m$ for supervising the student model.

Since the teacher network might produce results with large errors for some pixels, it is necessary to filter out the pixels with large errors. Inspired by [7], we introduce the minimum reprojection error to construct distillation loss for filtering out the multi-frame depth values that generate larger errors than the student single-frame depth. Given $\widetilde{T}_{t \to c}$, we can synthesize the images $\widetilde{I}_{c \to t}^s$ and $\widetilde{I}_{c \to t}^m$ according to $\widetilde{D}_t^s$ and $D_t^m$ respectively. Then, we can compare their photometric errors and generate the mask:

$$M = \left[ PE(I_t, \widetilde{I}_{c \to t}^m) < PE(I_t, \widetilde{I}_{c \to t}^s) \right], \quad (13)$$

TABLE I
QUANTITATIVE RESULTS ON THE EIGEN SPLIT OF KITTI DATASET WITH THE RAW AND IMPROVED GROUND TRUTH

| | Method | Train | Test Frames | #Params. | MACs | Time | The lower the better | | | | The higher the better | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | | | | | Abs Rel | Sq Rel | RMSE | RMSE log | $\delta_1$ | $\delta_2$ | $\delta_3$ |
| Raw GT | Monodepth2 [7] | M | 1 | 14.3M | 8.0G | 1.4ms | 0.115 | 0.903 | 4.863 | 0.193 | 0.877 | 0.959 | 0.981 |
| | PackNet-SfM [12] | M | 1 | 128.3M | 205.2G | 27.4ms | 0.111 | 0.785 | 4.601 | 0.189 | 0.878 | 0.960 | 0.982 |
| | VADepth [15] | M | 1 | 18.8M | 9.7G | 3.0ms | 0.104 | 0.774 | 4.552 | 0.181 | 0.892 | 0.965 | 0.983 |
| | Ma et al. [19] | M+Sem | 1 | 30.3M | - | - | 0.104 | 0.690 | 4.473 | 0.179 | 0.886 | 0.965 | 0.984 |
| | SD-SSMDE (ResNet50) [21] | M | 1 | - | 18.6G | - | 0.100 | 0.661 | 4.264 | 0.172 | 0.896 | 0.967 | 0.985 |
| | SUB-Depth [31] | M | 1 | - | - | - | 0.099 | 0.695 | 4.326 | 0.175 | 0.900 | 0.966 | 0.984 |
| | RA-Depth [13] | M | 1 | 10.0M | 10.8G | 3.4ms | 0.096 | 0.632 | 4.216 | 0.171 | 0.903 | 0.968 | 0.985 |
| | ManyDepth [8] | M | 2 (-1, 0) | 26.9M | 15.1G | 5.2ms | 0.098 | 0.770 | 4.459 | 0.176 | 0.900 | 0.965 | 0.983 |
| | DynamicDepth [9] | M+Sem | 2 (-1, 0) | - | - | - | 0.096 | 0.720 | 4.458 | 0.175 | 0.897 | 0.964 | <u>0.984</u> |
| | Long et al. [11] | M | 2 (-1, 0) | - | 15.6G | - | 0.097 | 0.731 | 4.392 | 0.176 | 0.901 | 0.965 | 0.983 |
| | MOVEDepth (ResNet18) [10]† | M | 2 (-1, 0) | 28.2M | 20.2G | 5.0ms | 0.094 | 0.704 | 4.389 | 0.175 | 0.902 | 0.965 | 0.983 |
| | DepthFormer [26] | M | 2 (-1, 0) | 28.7M | 174.7G | - | 0.090 | <u>0.661</u> | <u>4.149</u> | 0.175 | <u>0.905</u> | <u>0.967</u> | <u>0.984</u> |
| | MOVEDepth (PackNet) [10] | M | 2 (-1, 0) | 142.2M | 217.3G | 28.4ms | <u>0.089</u> | 0.663 | 4.216 | <u>0.169</u> | 0.904 | 0.966 | <u>0.984</u> |
| | Ours (ResNet18) | M | 2 (-1, 0) | 28.2M | 20.2G | 4.1ms | 0.092 | 0.683 | 4.331 | 0.172 | <u>0.905</u> | 0.966 | <u>0.984</u> |
| | **Ours** | M | 2 (-1, 0) | 23.9M | 22.9G | 6.1ms | **0.086** | **0.613** | **4.096** | **0.165** | **0.915** | **0.969** | **0.985** |
| Improved GT | Eigen et al. [4] | D | 1 | - | - | - | 0.190 | 1.515 | 7.156 | 0.270 | 0.692 | 0.899 | 0.967 |
| | DORN [36] | D | 1 | - | - | - | 0.072 | 0.307 | 2.727 | 0.120 | 0.932 | 0.984 | <u>0.994</u> |
| | Adabins [37] | D | 1 | - | - | - | <u>0.058</u> | 0.190 | 2.360 | 0.088 | 0.964 | 0.995 | **0.999** |
| | NeW CRFs [38] | D | 1 | - | - | - | **0.052** | **0.155** | **2.129** | **0.079** | **0.974** | **0.997** | **0.999** |
| | Monodepth2 [7] | M | 1 | 14.3M | 8.0G | - | 0.090 | 0.545 | 3.942 | 0.137 | 0.914 | 0.983 | 0.995 |
| | PackNet-SfM [12] | M | 1 | 128.3M | 205.2G | - | 0.078 | 0.420 | 3.485 | 0.121 | 0.931 | 0.986 | 0.996 |
| | RA-Depth [13]† | M | 1 | 10.0M | 10.8G | - | 0.074 | 0.362 | 3.345 | 0.113 | 0.940 | 0.990 | 0.997 |
| | Patil et al. [22] | M | N | - | 16.9G | - | 0.087 | 0.495 | 3.775 | 0.133 | 0.917 | 0.983 | 0.995 |
| | ManyDepth [8] | M | 2 (-1, 0) | 26.9M | 15.1G | - | 0.070 | 0.399 | 3.455 | 0.113 | 0.941 | <u>0.989</u> | <u>0.997</u> |
| | Long et al. [11] | M | 2 (-1, 0) | - | 15.6G | - | 0.068 | <u>0.366</u> | <u>3.338</u> | 0.110 | <u>0.946</u> | 0.989 | 0.997 |
| | MOVEDepth (ResNet18) [10]† | M | 2 (-1, 0) | 28.2M | 20.2G | - | 0.065 | 0.377 | 3.449 | 0.112 | 0.942 | 0.988 | 0.996 |
| | Ours (ResNet18) | M | 2 (-1, 0) | 28.2M | 20.2G | - | <u>0.064</u> | 0.369 | 3.390 | <u>0.108</u> | <u>0.946</u> | 0.988 | 0.996 |
| | **Ours** | M | 2 (-1, 0) | 23.9M | 22.9G | - | **0.058** | **0.302** | **3.070** | **0.098** | **0.955** | **0.992** | **0.998** |

All self-supervised methods are tested with the resolution of $192 \times 640$. "†" means evaluation on the pretrained models from github. The best scores for each subsection are in **bold** and the second are underlined. In the "Train" column, we list the training data for each method with D — ground truth Depth, M — unlabeled Monocular videos, Sem — Semantic labels. In the "Test Frames" column, "N" refers to taking a long sequence of frames as input to predict the target depth map. In the "Time" column, we list the inference time to generate one depth map by averaging the inference time of all 697 test images with a batch size of 16.

where $[\cdot]$ is the Iverson bracket. Following [4], we adopt the scale-invariant error between $D_t^m$ and $\widetilde{D}_t^s$ as the pseudo-label based regression loss:

$$L_{si} = \sqrt{\frac{1}{n}\sum_{u,v}(d(u,v))^2 - \frac{\gamma}{n^2}\Big(\sum_{u,v}d(u,v)\Big)^2}, \quad (14)$$

where $d(u,v) = (log(\widetilde{D}_t^s(u,v)) - log(D_t^m(u,v)))M(u,v)$, $n$ represents the number of elements with a value of 1 in $M$, and $\gamma = 1.0$. To provide supervision for the pixels where pseudo labels do not work, we also introduce the self-supervised loss $L_{self}(\widetilde{D}_t^s)$ as in (6) to construct the final distillation loss:

$$L_{distill} = L_{self}(\widetilde{D}_t^s) + \lambda_{si}L_{si}, \quad (15)$$

where $\lambda_{si} = 0.1$. In this way, the trained student model produces more accurate single-view depth than the S-DepthNet of the teacher model, thanks to the distilled geometric matching knowledge from the multi-frame network. Ultimately, we use the distilled student model to guide the M-DepthNet to obtain better depth estimation at test time.

## IV. EXPERIMENTAL RESULTS

### A. Datasets

We conduct experiments on the KITTI [1] and Cityscapes [39] datasets to verify the effectiveness of our method with the metrics proposed in [4]. KITTI is one of the most widely-used datasets for depth estimation, which covers various outdoor scenes. We follow [6] to adopt the Eigen split [4] and remove the static frames, which results in 39810/4424/697 training, validation, and test images. We also evaluate our model on the improved depth maps from [40], containing 652 test images. As for Cityscapes [39], it is a large dataset that is comprised of video sequences captured in streets from 50 cities. Following [6], we train on the 69731 monocular triplets and evaluate on the 1525 test images. As in [6], the maximum depth of evaluation on both datasets is restricted to 80m.

### B. Implementation Details

In our experiments, we adopt the HRNet [41] based architecture in [13] as S-DepthNet unless otherwise stated, where the number of output scales $S$ is set to 1. The PoseNet is a modified ResNet18 [42] as in [7]. Both the backbones of S-DepthNet and PoseNet are initialized with weights pretrained on ImageNet [43]. Following [10], we adopt a four-stage FPN [44] as the Feat-Net, where the number of feature channels $C = 32$. As for cost volume generation, both the number of depth candidates $N$ and the number of groups $G$ are set to 16. Similar to [8], we set the input resolution as $192 \times 640$ for KITTI and $128 \times 416$ for Cityscapes. Following [7], we use random color-jitter and flip for data augmentations, and further apply the random image mask strategy when training M-DepthNet as in [10]. All experiments are performed on a single Nvidia RTX 3090 GPU. Our models are implemented in Pytorch and the batch size is set as 12 and 16 for training and

Fig. 2. Qualitative results on the Eigen test split of KITTI dataset. Rows 2, 4, 6 provide the equivalently colormapped error maps for the metric Abs. Rel. relative to the improved depth [40], from small (blue) to large (red) errors. The GT depth is interpolated for better visualization. Compared to other methods [8], [10], our model not only preserves better details for various objects but also predicts depth maps with small errors. White and black boxes highlight the difference for the predicted depth and error maps, respectively. Best viewed in color and zoom in.

TABLE II
QUANTITATIVE RESULTS ON THE CITYSCAPES DATASET

| Method | AbsRel↓ | SqRel↓ | RMSE↓ | RMSE log↓ |
|---|---|---|---|---|
| Monodepth2 [7] | 0.129 | 1.569 | 6.876 | 0.187 |
| InstaDM [18]* | 0.111 | 1.158 | 6.437 | 0.182 |
| ManyDepth [8] | 0.114 | 1.193 | 6.223 | 0.170 |
| Long et al. [11] | 0.113 | 1.093 | 6.119 | 0.170 |
| DynamicDepth [9]* | <u>0.103</u> | <u>1.000</u> | <u>5.867</u> | <u>0.157</u> |
| **Ours** | **0.102** | **0.948** | **5.788** | **0.154** |

All methods are tested with the resolution of $128 \times 416$, except InstaDM [18] with a resolution of $256 \times 832$. "*" means that the method requires semantic labels for training.

TABLE III
GENERALIZATION PERFORMANCE ON THE CITYSCAPES DATASET

| Method | AbsRel↓ | SqRel↓ | RMSE↓ | RMSE log↓ |
|---|---|---|---|---|
| ManyDepth [8] | 0.170 | 1.789 | 8.357 | 0.236 |
| MOVEDepth [10]† | 0.164 | 1.780 | 8.678 | 0.238 |
| **Ours** | **0.150** | **1.492** | **7.810** | **0.216** |

"†" means evaluation on the pretrained models from github.

test respectively. Both the teacher model and student model are trained with Adam optimizer for $E$ epochs. $E$ is set to 20 for KITTI and 5 for Cityscapes. The initial learning rate is set to $2 \times 10^{-4}$ for the teacher model and $1 \times 10^{-4}$ for the student model, dropping by a factor of 10 after $Q$ epochs. $Q$ is 15 for KITTI and 1 for Cityscapes. For KITTI, it takes approximately 13 and 10 hours to train the teacher and student models, respectively. For Cityscapes, it takes about 5 and 4 hours respectively.

*C. Depth Evaluation*

To evaluate the performance of our method, we first conduct a comparison of its performance relative to the state-of-the-art self-supervised MDE methods on KITTI with the raw [1] and improved ground truth [40]. As shown in Table I, our method establishes a new state of the art in self-supervised MDE, with competitive model complexity and inference speed. Compared with the best-performing single-frame method [13] and distillation learning based methods [21], [31], our model improves the performance by more than $10.4\%$ (on Abs. Rel. with the raw GT). Note that our method only adopts the same S-DepthNet as RA-Depth, but does not adopt the data augmentation strategy and the cross-scale depth consistency loss proposed in [13]. The multi-frame depth estimation methods that perform closest to our method

are DepthFormer [26] and MOVEDepth [10]. Although neither using the computationally expensive transformer architecture [26] nor adopting the velocity-guided depth sampling strategy and additional depth fusing network [10], our model still outperforms these methods in all metrics. For a fair comparison with [10], we also list the results using the same ResNet18-based architecture (number of output scales $S = 4$) for our method and [10], where our method still performs better. Furthermore, the bottom half of Table I shows that our method even compares favorably to some single-frame supervised methods [36], [37], and narrows the performance gap between self-supervised monocular training methods and the best-performing supervised approach [38].

In addition, we also present the qualitative results in Fig. 2, where our method better preserves the shape of objects and outputs depth maps with smaller errors.

Moreover, we also compare the results with the current state-of-the-art methods [8], [9], [11] on Cityscapes. As shown in Table II, our method performs best again, even though DynamicDepth [9] leverages semantic labels.

*D. Generalization Performance*

To evaluate the generalization capability across datasets, the model trained on the KITTI dataset is used to test on the Cityscapes dataset without finetuning. Table III compares the generalization performance with the current self-supervised multi-frame depth estimation methods [8], [10], from which we can see that our method achieves better results. These data demonstrate that digging into the complementary information

TABLE IV
ABLATION STUDY ON KITTI EIGEN SPLIT FOR MULTI-FRAME DEPTH

| PADS | Distill | Min. Reproj. | The lower the better | | | | GPU (GB) | |
|---|---|---|---|---|---|---|---|---|
| | | | AbsRel | SqRel | RMSE | R log | train* | test |
| Baseline (HRNet18) | | | 0.090 | 0.704 | 4.293 | 0.171 | 15.1 | 4.2 |
| ✓ | | | 0.088 | 0.673 | 4.257 | 0.169 | 15.1 | 4.2 |
| | ✓ | | 0.088 | 0.640 | 4.196 | 0.168 | 15.1/10.4 | 4.2 |
| | ✓ | ✓ | 0.088 | 0.655 | 4.183 | 0.168 | 15.1/10.5 | 4.2 |
| ✓ | ✓ | | 0.087 | 0.637 | 4.133 | 0.166 | 15.1/10.4 | 4.2 |
| ✓ | ✓ | ✓ | **0.086** | **0.613** | **4.096** | **0.165** | 15.1/10.5 | 4.2 |
| Baseline (ResNet18) | | | 0.096 | 0.760 | 4.499 | 0.178 | 14.7 | 4.0 |
| ✓ | | | 0.094 | 0.748 | 4.457 | 0.176 | 14.7 | 4.0 |
| | ✓ | | 0.094 | 0.714 | 4.379 | 0.175 | 14.7/10.2 | 4.0 |
| | ✓ | ✓ | 0.093 | 0.684 | 4.335 | 0.173 | 14.7/10.3 | 4.0 |
| ✓ | ✓ | | **0.092** | 0.691 | 4.353 | **0.172** | 14.7/10.2 | 4.0 |
| ✓ | ✓ | ✓ | **0.092** | **0.683** | **4.331** | **0.172** | 14.7/10.3 | 4.0 |

"*": separate GPU memory required for the two stages of training.

TABLE V
ABLATION STUDY ON KITTI EIGEN SPLIT FOR SINGLE-FRAME DEPTH

| Distill | Min. Reproj. | The lower the better | | | | $\delta_1 \uparrow$ |
|---|---|---|---|---|---|---|
| | | AbsRel | SqRel | RMSE | R log | |
| Baseline (HRNet18) | | 0.102 | 0.741 | 4.470 | 0.179 | 0.896 |
| ✓ | | 0.100 | 0.702 | 4.328 | 0.175 | **0.900** |
| ✓ | ✓ | **0.099** | **0.676** | **4.287** | **0.174** | **0.900** |
| Baseline (ResNet18) | | 0.115 | 0.885 | 4.799 | 0.192 | 0.876 |
| ✓ | | 0.111 | 0.799 | 4.634 | **0.185** | 0.880 |
| ✓ | ✓ | **0.110** | **0.786** | **4.603** | **0.185** | **0.884** |
| SUB-Depth [31] | | **0.110** | 0.821 | 4.648 | **0.185** | **0.884** |



(a) The learned uncertainty map $\delta$

(b) Target image

Baseline / PADS / Distill w/o Min. Reproj. / Distill with Min. Reproj. / PADS + Distill w/o Min. Reproj. / PADS + Distill with Min. Reproj.

Baseline / Distill w/o Min. Reproj. / Distill with Min. Reproj.

(c) Error map of S-Depth

(d) error map of M-Depth

Fig. 3. Visualization results of ablation study. The learned $\delta$ is visualized in (a), where high uncertainty corresponding to large sampling range is white, otherwise black. The error maps of single-frame depth and multi-frame depth for different settings are visualized as (c) and (d) respectively.

between single-frame and multi-frame may contribute to improving the generalization capability across datasets.

*E. Ablation Study*

To understand how much each component of our method contributes to the overall performance of multi-frame depth estimation, we conduct the ablation experiments on KITTI with the raw GT depth. In our experiments, the baseline model is the trained teacher model without distillation learning, in which single-frame depth is set as the geometric center of the sampling range and all elements of $\delta$ are fixed as 0.3 following [10]. The top half of Table IV shows that using the PADS module leads to better scores in all metrics, which proves the effectiveness of the PADS module. Adopting the PADS module does not change the computation complexity but only introduces 7.68K extra parameters of $\delta$, which is negligible compared to the 23.9M parameters of the baseline model. Thus, using the PADS module requires similar GPU memory as the baseline. Adopting distillation learning can achieve greater performance gains than only introducing the PADS module, which reveals the effect of the proposed two-stage training pipeline. Combining the distillation learning loss with the masking strategy based on minimum reprojection error further brings an improvement in the overall performance at the cost of a little more GPU memory (10.4 v.s. 10.5) required for training the student model, which reflects the necessity to mask out the false labels generated by the teacher model. Taking all three components together leads to the most accurate depth predictions and disabling any component may result in performance degradation. These results suggest that all these components are compatible with each other and lead to more utilization of the mutual influence between single-frame and multi-frame depth estimation. In addition, to verify the compatibility of our method with different single-frame depth network architectures, we adopt the ResNet-based architecture in [7] for the ablation study. From the bottom half of Table IV, we can observe consistent results and draw the same conclusion.

Moreover, we perform an ablation study on the single-frame depth to further analyze the effectiveness of the minimum reprojection based distillation learning. As listed in Table V, the minimum reprojection based distillation learning brings a considerable performance gain for S-DepthNet. When applying the same backbone (ResNet18), our method performs better than the self-distillation method [31], which demonstrates that distillation learning from multi-frame depth estimation is more effective than self-distillation. Considering the performance gap between the distilled S-DepthNet and the M-DepthNet (even the baseline), we further use the trained student model to guide the M-DepthNet to output the final depth map. The consistent visualization results of the HRNet18-based models corresponding to the settings of Table IV and Table V are shown in Fig. 3. Taking all results of the ablation study together, we find that both single-frame depth estimation and multi-frame depth estimation do help improve each other.

## V. CONCLUSION

In this paper, we presented a distillation learning pipeline for self-supervised MDE so that single-frame and multi-frame depth networks can benefit from each other. Thanks to the proposed PADS module and minimum reprojection based distillation loss, our model achieves state-of-the-art performance and generalizes better than the previous methods.

However, the performance of our method on Cityscapes is still limited, which may suffer from more moving objects compared to KITTI which captures more stationary dynamic objects. In addition, the generalization performance of our method is also unsatisfactory. Note that dynamic scenes or new scenes and cameras not only directly affect M-DepthNet, but

also indirectly affect M-DepthNet by affecting PoseNet. Thus, further combining S-DepthNet and M-DepthNet to adapt to highly dynamic scenes or new scenes and cameras might be worthwhile, especially taking PoseNet into account.

## REFERENCES

[1] A. Geiger, P. Lenz, and R. Urtasun, "Are we ready for autonomous driving? the kitti vision benchmark suite," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.* IEEE, 2012, pp. 3354–3361.

[2] X. Luo, J.-B. Huang, R. Szeliski, K. Matzen, and J. Kopf, "Consistent video depth estimation," *ACM Trans. Graph.*, vol. 39, no. 4, p. 71, 2020.

[3] X. Guo, K. Yang, W. Yang, X. Wang, and H. Li, "Group-wise correlation stereo network," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2019, pp. 3273–3282.

[4] D. Eigen, C. Puhrsch, and R. Fergus, "Depth map prediction from a single image using a multi-scale deep network," *Proc. Adv. Neural Inf. Process. Syst.*, pp. 2366–2374, 2014.

[5] R. Garg, V. K. Bg, G. Carneiro, and I. Reid, "Unsupervised cnn for single view depth estimation: Geometry to the rescue," in *Proc. Eur. Conf. Comput. Vis.* Springer, 2016, pp. 740–756.

[6] T. Zhou, M. Brown, N. Snavely, and D. G. Lowe, "Unsupervised learning of depth and ego-motion from video," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2017, pp. 1851–1858.

[7] C. Godard, O. Mac Aodha, M. Firman, and G. J. Brostow, "Digging into self-supervised monocular depth estimation," in *Proc. IEEE Int. Conf. Comput. Vis.*, 2019, pp. 3828–3838.

[8] J. Watson, O. Mac Aodha, V. Prisacariu, G. Brostow, and M. Firman, "The temporal opportunist: Self-supervised multi-frame monocular depth," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, June 2021, pp. 1164–1174.

[9] 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 *Proc. Eur. Conf. Comput. Vis.* Springer, 2022, pp. 228–244.

[10] X. Wang, Z. Zhu, G. Huang, X. Chi, Y. Ye, Z. Chen, and X. Wang, "Crafting monocular cues and velocity guidance for self-supervised multi-frame depth learning," in *Proc. AAAI Conf. Artif. Intell.*, vol. 37, no. 3, 2023, pp. 2689–2697.

[11] Y. Long, H. Yu, and B. Liu, "Two-stream based multi-stage hybrid decoder for self-supervised multi-frame monocular depth," *IEEE Robot. Autom. Lett.*, vol. 7, no. 4, pp. 12 291–12 298, 2022.

[12] V. Guizilini, R. Ambrus, S. Pillai, A. Raventos, and A. Gaidon, "3d packing for self-supervised monocular depth estimation," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2020, pp. 2485–2494.

[13] M. He, L. Hui, Y. Bian, J. Ren, J. Xie, and J. Yang, "Ra-depth: Resolution adaptive self-supervised monocular depth estimation," in *Proc. Eur. Conf. Comput. Vis.* Springer, 2022, pp. 565–581.

[14] A. Saxena, S. H. Chung, and A. Y. Ng, "Learning depth from single monocular images," in *Proc. Adv. Neural Inf. Process. Syst.*, 2005, pp. 1161–1168.

[15] J. Xiang, Y. Wang, L. An, H. Liu, Z. Wang, and J. Liu, "Visual attention-based self-supervised absolute depth estimation using geometric priors in autonomous driving," *IEEE Robot. Autom. Lett.*, vol. 7, no. 4, pp. 11 998–12 005, 2022.

[16] A. Ranjan, V. Jampani, L. Balles, 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 *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2019, pp. 12 240–12 249.

[17] H. Zhao, J. Zhang, S. Zhang, and D. Tao, "Jperceiver: Joint perception network for depth, pose and layout estimation in driving scenes," in *Proc. Eur. Conf. Comput. Vis.* Springer, 2022, pp. 708–726.

[18] S. Lee, S. Im, S. Lin, and I. S. Kweon, "Learning monocular depth in dynamic scenes via instance-aware projection consistency," in *Proc. AAAI Conf. Artif. Intell.*, vol. 35, no. 3, 2021, pp. 1863–1872.

[19] J. Ma, X. Lei, N. Liu, X. Zhao, and S. Pu, "Towards comprehensive representation enhancement in semantics-guided self-supervised monocular depth estimation," in *Proc. Eur. Conf. Comput. Vis.* Springer, 2022, pp. 304–321.

[20] W. Ren, L. Wang, Y. Piao, M. Zhang, H. Lu, and T. Liu, "Adaptive co-teaching for unsupervised monocular depth estimation," in *Proc. Eur. Conf. Comput. Vis.* Springer, 2022, pp. 89–105.

[21] A. Petrovai and S. Nedevschi, "Exploiting pseudo labels in a self-supervised learning framework for improved monocular depth estimation," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, June 2022, pp. 1578–1588.

[22] V. Patil, W. Van Gansbeke, D. Dai, and L. Van Gool, "Don't forget the past: Recurrent depth estimation from monocular video," *IEEE Robot. Autom. Lett.*, vol. 5, no. 4, pp. 6813–6820, 2020.

[23] Y. Yao, Z. Luo, S. Li, T. Fang, and L. Quan, "Mvsnet: Depth inference for unstructured multi-view stereo," in *Proc. Eur. Conf. Comput. Vis.*, 2018, pp. 767–783.

[24] J. Yang, J. M. Alvarez, and M. Liu, "Self-supervised learning of depth inference for multi-view stereo," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2021, pp. 7526–7534.

[25] Y. Ding, Q. Zhu, X. Liu, W. Yuan, H. Zhang, and C. Zhang, "Kd-mvs: Knowledge distillation based self-supervised learning for multi-view stereo," in *Proc. Eur. Conf. Comput. Vis.* Springer, 2022, pp. 630–646.

[26] V. Guizilini, R. Ambruș, D. Chen, S. Zakharov, and A. Gaidon, "Multi-frame self-supervised depth with transformers," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, June 2022, pp. 160–170.

[27] J. Gou, B. Yu, S. J. Maybank, and D. Tao, "Knowledge distillation: A survey," *Int. J. Comput. Vis.*, vol. 129, pp. 1789–1819, 2021.

[28] A. Pilzer, S. Lathuiliere, N. Sebe, and E. Ricci, "Refine and distill: Exploiting cycle-inconsistency and knowledge distillation for unsupervised monocular depth estimation," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2019, pp. 9768–9777.

[29] X. Ye, X. Fan, M. Zhang, R. Xu, and W. Zhong, "Unsupervised monocular depth estimation via recursive stereo distillation," *IEEE Trans. Image Process.*, vol. 30, pp. 4492–4504, 2021.

[30] Z. Zhou and Q. Dong, "Self-distilled feature aggregation for self-supervised monocular depth estimation," in *Proc. Eur. Conf. Comput. Vis.* Springer, 2022, pp. 709–726.

[31] H. Zhou, S. Taylor, D. Greenwood, and M. Mackiewicz, "Self-distillation and uncertainty boosting self-supervised monocular depth estimation," in *Proc. Brit. Mach. Vis. Conf.* BMVA Press, 2022, p. 7.

[32] C. Godard, O. Mac Aodha, and G. J. Brostow, "Unsupervised monocular depth estimation with left-right consistency," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2017, pp. 270–279.

[33] O. Ronneberger, P. Fischer, and T. Brox, "U-net: Convolutional networks for biomedical image segmentation," in *Proc. Int. Conf. Med. Image Comput. Comput.-Assisted Intervention.* Springer, 2015, pp. 234–241.

[34] F. Wang, S. Galliani, C. Vogel, and M. Pollefeys, "Itermvs: Iterative probability estimation for efficient multi-view stereo," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, June 2022, pp. 8606–8615.

[35] Z. Teed and J. Deng, "Raft: Recurrent all-pairs field transforms for optical flow," in *Proc. Eur. Conf. Comput. Vis.* Springer, 2020, pp. 402–419.

[36] H. Fu, M. Gong, C. Wang, K. Batmanghelich, and D. Tao, "Deep ordinal regression network for monocular depth estimation," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2018, pp. 2002–2011.

[37] S. F. Bhat, I. Alhashim, and P. Wonka, "Adabins: Depth estimation using adaptive bins," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, June 2021, pp. 4009–4018.

[38] W. Yuan, X. Gu, Z. Dai, S. Zhu, and P. Tan, "Neural window fully-connected crfs for monocular depth estimation," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, June 2022, pp. 3916–3925.

[39] M. Cordts, M. Omran, S. Ramos, T. Rehfeld, M. Enzweiler, R. Benenson, U. Franke, S. Roth, and B. Schiele, "The cityscapes dataset for semantic urban scene understanding," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2016, pp. 3213–3223.

[40] J. Uhrig, N. Schneider, L. Schneider, U. Franke, T. Brox, and A. Geiger, "Sparsity invariant cnns," in *Proc. Int. Conf. 3D Vis.* IEEE, 2017, pp. 11–20.

[41] J. Wang, K. Sun, T. Cheng, B. Jiang, C. Deng, Y. Zhao, D. Liu, Y. Mu, M. Tan, X. Wang *et al.*, "Deep high-resolution representation learning for visual recognition," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 43, no. 10, pp. 3349–3364, 2020.

[42] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2016, pp. 770–778.

[43] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, "Imagenet: A large-scale hierarchical image database," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.* Ieee, 2009, pp. 248–255.

[44] T.-Y. Lin, P. Dollár, R. Girshick, K. He, B. Hariharan, and S. Belongie, "Feature pyramid networks for object detection," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2017, pp. 2117–2125.

[45] J. Baek, G. Kim, and S. Kim, "Semi-supervised learning with mutual distillation for monocular depth estimation," in *IEEE Int. Conf. Robot. Autom.*, 2022, pp. 4562–4569.