

# Deep Monocular Depth Estimation Leveraging a Large-scale Outdoor Stereo Dataset

Jaehoon Cho, *Student Member, IEEE*, Dongbo Min<sup>†</sup>, *Senior Member, IEEE*,  
Youngjung Kim, *Member, IEEE*, and Kwanghoon Sohn, *Senior Member, IEEE*

**Abstract**—Current self-supervised methods for monocular depth estimation are largely based on deeply nested convolutional networks that leverage stereo image pairs or monocular sequences during training phase. However, they often exhibit inaccurate results around occluded regions and depth boundaries. In this paper, we present a simple yet effective approach for monocular depth estimation using stereo image pairs. The study aims to propose a student-teacher strategy in which a shallow student network is trained with the auxiliary information obtained from a deeper and more accurate teacher network. Specifically, we first train the *stereo* teacher network by fully utilizing the binocular perception of 3-D geometry, and then use the depth predictions of the teacher network to train the student network for monocular depth inference. This enables us to exploit all available depth data from massive unlabeled stereo pairs. We propose a strategy that involves the use of a data ensemble to merge the multiple depth predictions of the teacher network to improve the training samples by collecting non-trivial knowledge beyond a single prediction. To refine the inaccurate depth estimation that is used when training the student network, we further propose stereo confidence guided regression loss that handles the unreliable pseudo depth values in occlusion, texture-less region, and repetitive pattern. To complement the existing dataset comprising outdoor driving scenes, we built a novel large-scale dataset consisting of one million outdoor stereo images taken using hand-held stereo cameras. Finally, we demonstrate that the monocular depth estimation network provides feature representations that are suitable for high-level vision tasks. The experimental results for various outdoor scenarios demonstrate the effectiveness and flexibility of our approach, which outperforms state-of-the-art approaches.

**Index Terms**—Monocular depth estimation, convolutional neural network, student-teacher strategy, outdoor stereo dataset, stereo confidence maps.

## I. INTRODUCTION

**O**BTAINING 3-D depth of a scene is essential to alleviate a number of challenges in robotics and computer vision tasks including 3-D reconstruction [1], autonomous driving [2]–[4], intrinsic image decomposition [5], and scene understanding [6]. The human visual system (HVS) can understand the 3-D structure by perceiving the depth value of the scene using binocular fusion. Such a mechanism has been widely adopted in computational stereo approaches that establish correspondence maps across two (or more) images

taken of the same scene [7]. This approach has achieved outstanding performance in recent studies [8]–[13]. Furthermore, numerous monocular depth estimation approaches have been developed based on monocular cues, for example, inclusion of the object contour [14], and segmentation [15]. However, most methods rely heavily on handcrafted rules based on one or a few monocular cues, and thus they often fail to capture plausible depth from a single image and are effective only at very restricted environments.

Recently, owing to advances in deep neural networks, depth prediction from a single image has advanced considerably with the aid of convolutional neural networks (CNNs) [17]–[23]. However, supervised learning approaches for monocular depth estimation necessarily have several limitations outdoors. Monocular depth estimation needs (semi-)dense depth maps as pixel-level supervision for training the deep network, and constructing such large-scale training data with depth maps is extremely challenging. An active depth sensor, LiDAR, is commonly used to acquire depth maps, however, it is usually of low resolution and very sparse, for example, less than 6% in the KITTI dataset [24]. Owing to its sparsity, it cannot cover all salient objects in a scene. Additionally, the sensing device is very expensive and is often adversely affected by several internal degradations, such as imperfect sensor calibration and photometric distortions. Thus, most of the existing public datasets provide only a small number of depth maps for rather limited scenes. For example, they mostly consist of driving scenes obtained from the depth sensor mounted on a vehicle [24], [25].

To address the lack of ground truth depth maps, recent works introduce a self-supervised paradigm [21], [23], [27] that uses stereo image pairs during the training phase. These methods impose the left-right photometric consistency via view synthesis. Since it simply measures on RGB color or hand-crafted feature (e.g., SSIM) spaces, they inherently suffer from texture-less region, occlusion, and repeated pattern, resulting in blurry depth boundaries and fattening artifacts. The method [22] uses both the supervised loss using ground truth depth maps and the self-supervised image reconstruction loss. However, it still faces a performance degradation due to the supervised loss when the trained model feeds an image from a novel domain. Owing to the self-supervised loss, it yields blurry depth results similar to [21], [23], [27], as shown in Fig. 1. Another proposed method [20] first synthesizes the right view using deep networks and then performs stereo matching to produce a final depth map inside the monocular depth network. However, this two-step method substantially

J. Cho, and K. Sohn are with the School of Electrical and Electronic Engineering, Yonsei University, Seoul 120-749, South Korea (e-mail: {je-hoon.khsohn}@yonsei.ac.kr).

D. Min is with the Department of Computer Science and Engineering, Ewha Womans University, Seoul 03-760, South Korea (e-mail: dbmin@ewha.ac.kr).

Y. Kim is with the Agency for Defense Development, Daejeon 34186, South Korea (e-mail: read12300@add.re.kr).

<sup>†</sup> Corresponding author

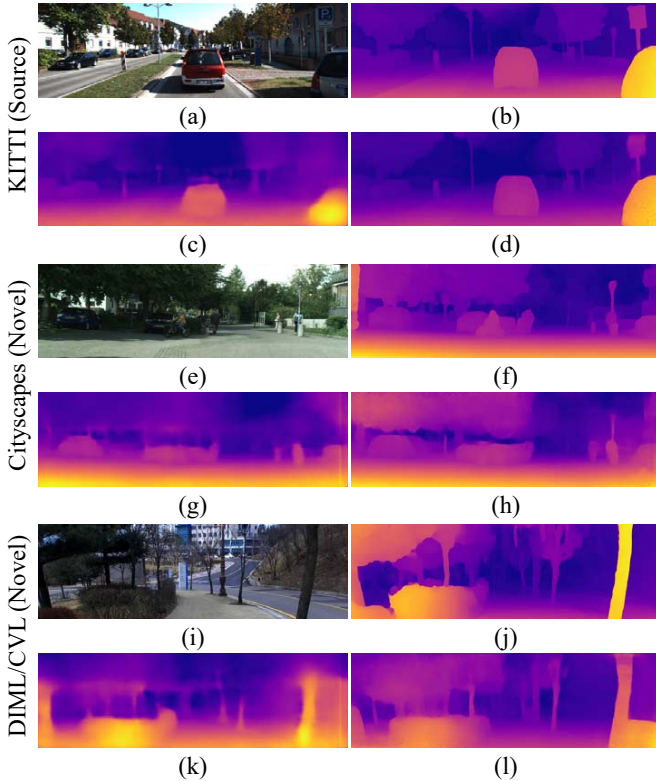


Figure 1. Sample images collected from various datasets and estimated depth maps: (a), (e), (i) input images, (b), (f), (j) depth maps predicted by deep stereo matching network [32], (c), (g), (k) depth maps obtained using state-of-the-art method [22], and (d), (h), (l) depth maps obtained using our method.

increases the computational complexity.

In this paper, we present a novel approach based on a student-teacher strategy. Our approach obviates the need to use massive ground truth depth maps. We propose the use of a massive number of stereo image pairs, which are relatively easy to obtain, and a small amount of training data with ground truth depth maps provided in existing public datasets. First, the teacher network for stereo matching is trained using a small amount of training data with ground truth depth maps. Subsequently, multiple estimated depth maps are generated on various scales, and then these maps are fused. This ensemble approach improves the accuracy of the pseudo-ground-truth depth maps by collecting non-trivial knowledge beyond a single prediction. Given the fused depth maps, the associated stereo confidence maps are generated via a confidence measure network. The stereo confidence maps encourage the use of only the reliable depth values in the fused depth maps. With these maps, we propose a stereo confidence-guided regression loss that employs a mask function. To use diverse stereo image pairs when training the monocular depth networks, we built a new large-scale dataset named the DIML/CVL dataset, by capturing stereo image pairs of various scenes, including parks, brooks, and apartments (Fig. 7). The experimental results demonstrate that our method outperforms the state-of-the-art method, and the DIML/CVL database is complementary to existing outdoor driving scenes [24], [25]. We additionally demonstrate that our pre-trained model for monocular depth prediction can be used as a powerful proxy task for scene understanding tasks.

Recent works [28], [29] are conceptually similar to the proposed method in that they also exploit stereo matching as a rich knowledge for monocular depth estimation. In contrast, our method is more effective than previous works by handling unreliable areas such as texture-less regions and occlusions in stereo matching using confidence guided loss.

Our main contributions are as follows:

- We propose a novel framework for monocular depth estimation based on a student-teacher strategy.
- We introduce a data ensemble and stereo confidence-guided regression loss to improve the usage of the pseudo-ground-truth.
- We introduce a new large-scale RGB-D dataset, named the DIML/CVL dataset, which is complementary to the existing datasets of outdoor driving images. The dataset is publicly available at <https://dimlrgb-d.github.io/>.
- We demonstrate that the feature representation of our monocular depth estimation provides rich knowledge for scene understanding tasks.

## II. RELATED WORK

In this section, we briefly review and discuss three lines of work that are the most relevant to our study.

### A. Stereo matching

The objective of stereo matching is to find a set of corresponding points between two (or more) images. The correspondence map is converted into a depth map using stereo calibration parameters. Early studies based on CNN attempted to measure the similarity between patches of two images. Han *et al.* [8] proposed a Siamese network that extracts features from patches followed by a similarity measure. Zbontar and LeCun [9] computed the matching cost using CNNs and applied it to classical stereo matching pipelines consisting of cost aggregation, depth optimization, and post-processing. Luo *et al.* [10] proposed computing the matching cost by learning a probability distribution over all depth values and then computing the inner product between two feature maps. Note that these approaches focused on computing the matching cost using CNNs and the remaining procedures for the stereo matching still rely on conventional handcrafted approaches.

Recent approaches have attempted to predict a depth map in an end-to-end fashion, achieving a substantial performance gain. Mayer *et al.* [11] proposed a new method, named DispNet, which uses a series of convolutional layers for cost aggregation and then employs regression to obtain the depth map. Pang *et al.* [32] introduced a two-stage network, named cascade residual learning (CRL), which is an extension of DispNet [11]. The first and second stages are used to calculate the depth maps and their multi-scale residuals, and then the outputs of both stages are combined to form a final depth map. Kendall *et al.* [12] introduced a new end-to-end approach that performs cost aggregation using a series of 3-D convolutions. Chang *et al.* [13] incorporated a contextual information through 3-D convolutions using stacked multiple hourglass networks over cost volume. The above-mentioned stereo matching approaches can be utilized as the teacher network in our framework.

## B. Monocular depth estimation

The performance of the monocular depth estimation has been advanced dramatically through the supervised learning approach that uses depth maps acquired from active sensors as ground truth for input images. Eigen *et al.* [17] designed a multi-scale deep network that predicts a coarse depth map and then progressively refines the depth map. Liu *et al.* [19] casted the monocular depth estimation into a continuous conditional random field (CRF) learning problem that jointly learns the unary and pairwise potentials of the CRF in a unified deep CNN framework. Fu *et al.* [30] introduced a discrete paradigm that uses an ordinal regression loss to encourage the ordinal competition among depth values in an end-to-end manner.

To address the lack of massive ground truth depth maps, several approaches have attempted to learn the monocular depth inference networks using image pairs only in a self-supervised manner. Using a left-right consistency constraint, Godard *et al.* [21] proposed an improved architecture for training the monocular depth estimation using stereo image pairs. Zhou *et al.* [33] designed a model to jointly estimate depth and camera pose in a self-supervised manner by leveraging temporally adjacent frames acquired by a single moving camera. Similarly, Godard *et al.* [23] exploited consecutive frames to estimate both the depth and the camera pose. One drawback of these approaches is that reconstruction loss based on image matching is an ill-posed problem on its own. To address this problem, Kuznetsov *et al.* [22] exploited both supervised and self-supervised losses with ground truth depth maps and stereo image pairs. This approach still suffers from the performance degradation in novel environments, as presented in Fig. 1. Luo *et al.* [20] formulated monocular depth estimation with two sub-networks: a view synthesis network and a stereo matching network. They first synthesize stereo pairs from an input image, and then apply the stereo matching network to produce the depth map. Zhao *et al.* [34] designed a geometry aware symmetric domain adaptation framework using the ground truth depth maps in synthetic data and epipolar geometry in the real data.

Recent studies attempted to use pseudo ground truth depth maps obtained by exploiting stereo matching as a rich knowledge. Guo *et al.* [28] trained the stereo matching network using a synthetic stereo dataset. The above-mentioned methods did not consider performance degradation incurred by unreliable pseudo depth values in occlusion, texture-less region, and repetitive pattern that can not be addressed using stereo matching methods. In contrast, our method is more effective than recent studies by handling unreliable values using confidence guided loss.

## C. Student-teacher strategy

Models that are deep and wide are known to tend to be more accurate than shallow models and they provide a large capacity [36], [37]. The student-teacher strategy concentrates on learning a much smaller model (student) from a large deep network (usually referred to as the teacher). With this method, the shallow network can be as accurate as the deep teacher network [38]. Ye *et al.* [39] introduced a joint scene

parsing and depth estimation network with a shallow model using heterogeneous-task deep teacher networks. Pilzer *et al.* [40] designed a refinement process via cycle consistency and distillation strategies for monocular depth estimation. In our method, the pseudo ground truth depth maps computed from the existing stereo matching network, which acts as a deep teacher network, are used to train the student network for monocular depth inference.

## D. Feature learning via pretext task

Several approaches have attempted to leverage a pretext task as an alternative form of supervision in some applications where it is difficult to construct massive ground truth data. Noroozi and Favaro [41] proposed solving jigsaw puzzles for representing object parts and their spatial arrangements, and then applied it to object classification and detection tasks. Pathak *et al.* [42] proposed an unsupervised image inpainting approach that generates the contents of an arbitrary image region conditioned on its surroundings. The learned encoder features are applied to object classification/detection and semantic segmentation tasks. Larsson *et al.* [43] investigated image colorization as a proxy task in the replacement of ImageNet [44].

These studies have been successfully transferred to various high-level tasks. In our study, we demonstrate that the network pre-trained for monocular depth prediction is a powerful proxy task for learning feature representations in scene understanding tasks.

# III. PROPOSED METHOD

## A. Motivation and overview

Owing to the lack of scene diversity, deep networks for monocular depth estimation often undergo performance degradation in novel environments. For instance, feeding a single (monocular) image from a source domain, in which the deep networks are trained, yields satisfactory results (Fig. 1 (c)). However, when we test an image from a novel domain, the output depth map produces inaccurate results around occluded regions and depth boundaries. The monocular depth estimation network trained with the KITTI dataset [24] does not work well on the Cityscapes dataset [25] and our new dataset. In contrast, the stereo matching network using [32] produces fine-grained depth maps on both the source and novel domains, as shown in Fig. 1 (b), (f), and (j), even though it is trained with the KITTI dataset only.

Stereo matching aims to find similar patches from a number of candidates extracted from two images. Thus, it is sufficient to train the network with similar patches (positive samples) and dissimilar patches (negative samples) [9]. The stereo matching network attempts to learn the local patch matching from the cost-volume (*i.e.*, correlation layer, which explicitly encodes the geometric relationship [11], [32]). Some methods propose training the stereo matching network using two images at once to additionally leverage a global context on stereo matching [13]. However, the underlying principle is to locally explore the patch-level similarity for two-view matching. Contrarily, monocular depth estimation, which infers a

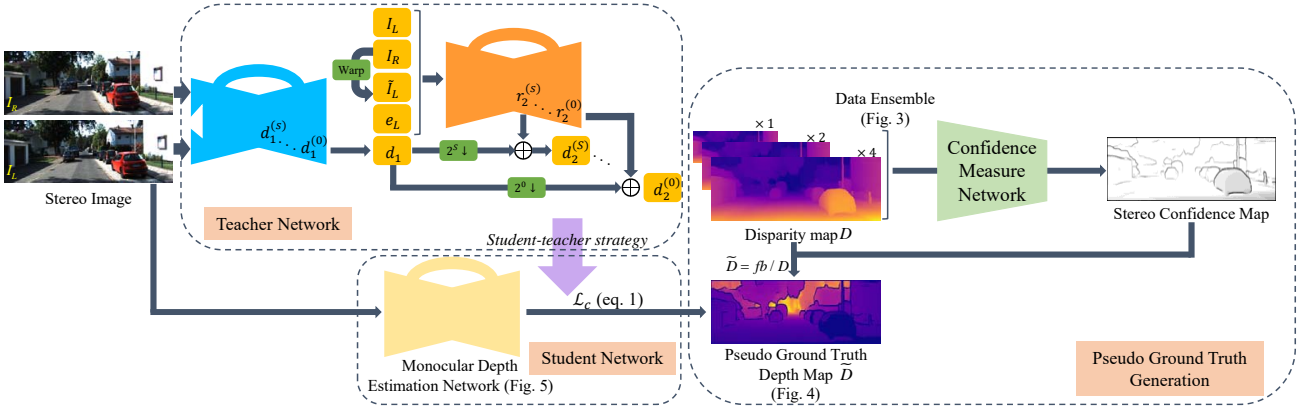


Figure 2. Proposed framework for monocular depth estimation. We first generate depth maps from stereo image pairs using a deep teacher network [32]. Here, we adopt the data ensemble [46] to fuse depth maps estimated on multiple scales. A stereo confidence map [45] is generated to identify inaccurate estimated stereo depth values. Subsequently, the monocular student network is trained with the pseudo-ground-truth depth maps via stereo confidence guided regression loss.

depth value from a single image by making use of monocular cues, is highly ill-posed and more challenging than stereo matching. Thus, the global context is crucial for predicting the overall 3-D structure of scenes. In this regard, the monocular depth estimation network is trained using the image and depth map, rather than a pair of patches extracted from them, to consider the global context. This increases the sensitivity of the monocular depth estimation network to the domain difference problem compared with the stereo matching network. Thus, a great variety of scenes is needed to train the monocular depth estimation network, whereas the stereo matching network is relatively free from such constraints [28], [29].

Our overall framework is illustrated in Fig. 2. We propose a simple yet effective approach for monocular depth estimation by leveraging the student-teacher strategy. The shallow student network learns from the more informative deep teacher network. Our method involves the following steps. Given a number of stereo images, we generate depth maps using the deep stereo matching network trained with ground truth data. When generating depth maps, we fuse depth maps that were estimated on multiple scales to provide non-trivial knowledge from multiple predictions. A stereo confidence map is then estimated as auxiliary data to avoid inaccurate stereo depth values being utilized when training the monocular depth estimation network. The pseudo-ground-truth depth maps are used to supervise the monocular student network via stereo confidence guided regression loss. The experimental results show that our framework is easily extendable to various domains by leveraging only stereo image pairs. Furthermore, the monocular depth estimation induces feature representation that improves scene-understanding tasks such as semantic segmentation and road detection.

Differently from other frameworks following a similar work [28], [29] that uses existing stereo matching to generate pseudo-ground-truth depth maps and then use these maps for training monocular depth network, we can handle unreliable areas such as texture-less regions and occlusions in stereo matching using confidence guided loss yielding a notable accuracy improvement compared to the existing solutions.

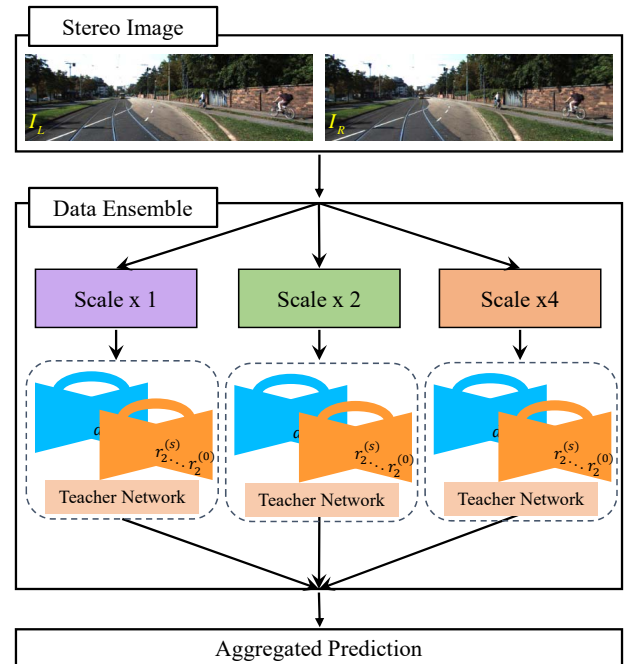


Figure 3. Data ensemble approach. We generate depth maps of stereo images at different scales via the deep teacher network, and then ensemble them on the smallest scale.

### B. Pseudo-ground-truth generation

The stereo matching network takes stereo images ( $I_l, I_r$ ) as input and outputs a depth map  $D$  aligned with the left image  $I_l$ . We adopted cascade residual learning (CRL) [32] as the stereo teacher network. To further improve the depth map of the teacher network, we adopt an ensemble prediction method that merges output depth maps on various scales, as shown in Fig. 3. It has been shown that the generated data can be improved by applying the same model to multiple transformations (e.g., scale, rotation, and flipping) of the input and then aggregating the results [51], [52]. We generate depth maps on three different scales and average them on the smallest scale. The data ensemble can provide auxiliary information from multiple predictions beyond a single prediction.

The pseudo-ground-truth depth maps inevitably contain

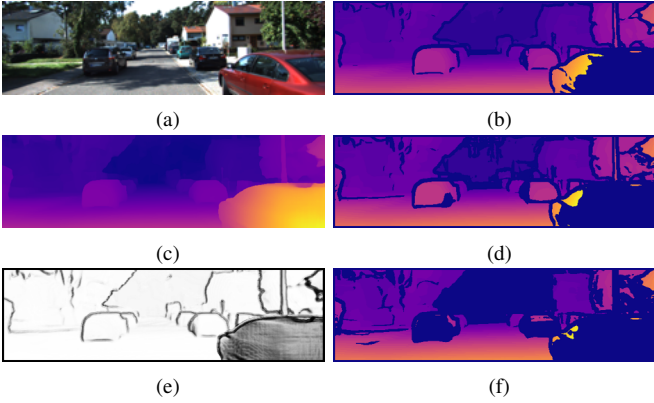


Figure 4. Pseudo ground truth data samples. (a) input image, (c) estimated depth map with data ensemble, and (e) predicted stereo confidence map. (b), (d), (f): pseudo ground truth depth maps thresholded by the stereo confidence map (e) with  $\tau = 0.3, 0.55, \text{ and } 0.75$ , respectively. Black pixels indicate unreliable pixels detected by the stereo confidence map.

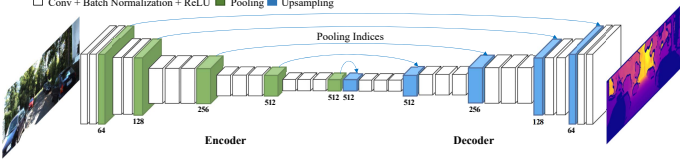


Figure 5. Monocular depth estimation network. We designed a variant of the U-Net architecture [55] as a baseline network.

erroneous estimates. To prevent such inaccurate depth values from being used when training the monocular depth estimation network, it is necessary to identify them in the pseudo-ground-truth depth maps. To this end, we additionally estimate the stereo confidence map [45]. The confidence measure network extracts a square patch from the depth map and forwards it to a CNN to infer a normalized confidence value  $C(p) \in [0, 1]$ . We denote the confidence threshold as the hyper-parameter  $\tau$ . The depth value of each pixel is set to be reliable when  $C(p) \geq \tau$ , and vice versa. By adjusting  $\tau$ , we can control the sparseness and reliability of the depth map. As  $\tau$  increases, unreliable areas such as texture-less regions and occlusions are removed effectively, however, the depth map becomes sparse. Sample images of depth maps and a stereo confidence map are shown in Fig. 4. Because the confidence measure network uses a square patch extracted from a depth map without using padding or stride, it assigns zero to the boundary of the confidence map. Note that a rich line of research that successfully improves the performance when imperfect training data are used has been published [28], [29], [53], [54]. They showed that the use of massive pseudo ground truth data achieved outstanding results. Although the use of massive training data does not deliver perfect results, the performance can be significantly boosted by identifying erroneous data well. Our method is conceptually similar to these approaches.

### C. Network Architecture

We designed monocular depth estimation networks based on a variant of the U-Net architecture [55]. As shown in Fig 5, the encoder network consists of the first 13 convolutional layers in the VGG [56] network, similar to [57]. We discarded the fully connected layers in favor of maintaining

spatial information. The convolutional layers consist of  $3 \times 3$  convolutions with batch normalization [59] and a rectified linear unit (ReLU). Max-pooling with a  $2 \times 2$  window and a stride of 2 is performed at the encoder. The indices of the max locations are computed and stored during pooling. Each convolutional layer at the encoder has a corresponding convolutional layer at the decoder; thus, the decoder network has 13 layers. The decoder upsamples the feature maps through unpooling with the memorized max-pooling indices [57] from the corresponding encoder feature map. The sparse feature map is densified by convolving it with a trainable decoder filter bank. Using such pooling indices boosts the performance and enables more efficient training. The final decoder output is fed to a regression loss.

### D. Training loss

Given a monocular input image and pseudo-ground-truth depth map  $\tilde{D}(p)$ , we propose to use the *stereo confidence guided regression loss*  $\mathcal{L}_c$ :

$$\mathcal{L}_c = \frac{1}{\sum_p M_p} \sum_p M_p \cdot \left| \hat{D}(p) - \tilde{D}(p) \right|_1, \quad (1)$$

$$M_p = \begin{cases} 1, & \text{if } C(p) \geq \tau \\ 0, & \text{if } C(p) < \tau \end{cases}, \quad (2)$$

where  $\hat{D}(p)$  denotes the depth map predicted by the monocular depth estimation network. In Section 4.4, we validate the effect of the stereo confidence measure by adjusting the hyper-parameter  $\tau$ . The *stereo confidence guided regression loss* in Eq. 1 not only ensures the use of reliable depth values when training the monocular depth estimation network but also propagates reliable depth predictions from highly confident pixels into low-confidence pixels. This loss could handle the challenging elements of stereo matching such as occlusion, texture-less areas, and depth edges and effectively identifies reliable depth values. In Fig. 6, although there is no valid depth value around depth boundaries in terms of stereo confidence map, our results equipped with the stereo confidence guided loss show an excellent boundary preserving capability.

### E. Large-scale outdoor stereo dataset

Our new outdoor stereo dataset, the DIML/CVL RGB-D dataset, is complementary to existing RGB-D datasets such as the KITTI and Cityscapes datasets. To ensure the diversity of training data, we attempted to obtain non-driving scenes (*e.g.*, parks, buildings, apartments, trails, and streets) using hand-held stereo cameras, unlike the existing dataset, which consists mostly of driving scenes (*e.g.*, road and traffic scenes). Fig. 7 shows sample RGB-D pairs of our dataset. The total number of distinct scenes is 1,053. Each scene contains a different number of images, ranging from 13,971 to 55,577. The images in this dataset were taken from fall 2015 to summer 2017 in four different cities (Seoul, Daejeon, Cheonan, and Sejong) in South Korea.

Two types of stereo cameras, a ZED stereo camera <sup>1</sup> and a custom-built stereo camera, were used with different

<sup>1</sup><https://www.stereolabs.com/>

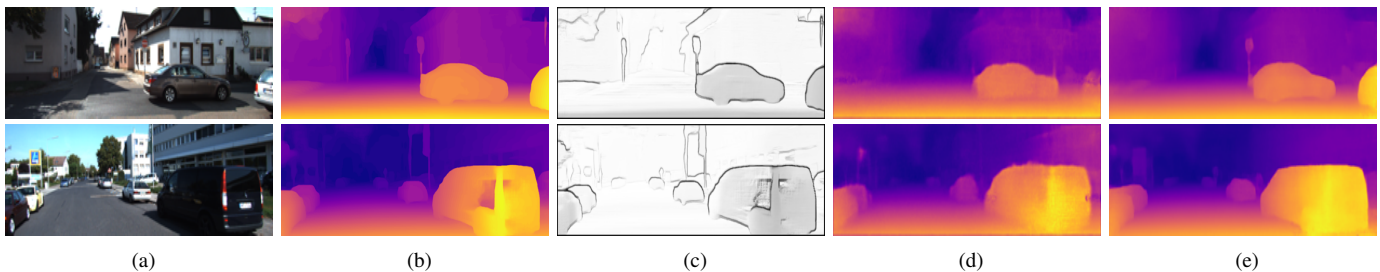


Figure 6. Visual results on the Eigen split [17] of the KITTI dataset. (a) monocular left image, (b) pseudo-ground-truth depth maps, (c) estimated stereo confidence maps, (d) results of depth maps trained without confidence guided loss, and (e) results of depth maps trained with confidence guided loss.

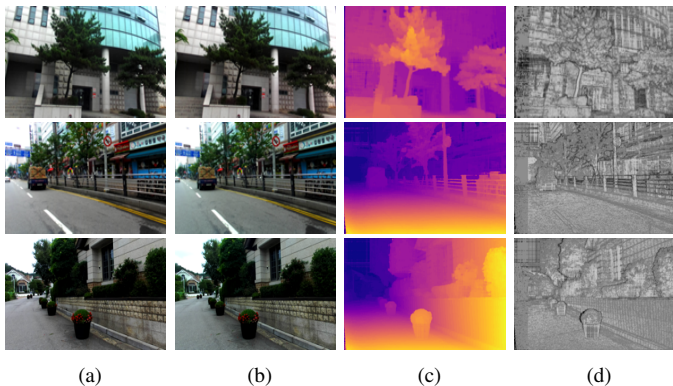


Figure 7. Samples of RGB-D pairs with stereo confidence map from DIML/CVL dataset [26]. (a) left image, (b) right image, (c) depth map using [9], and (d) stereo confidence map using [49].

camera configurations for the baseline and focal length. The commercial ZED camera has a small baseline (12 cm), thus its sensing range is rather limited (up to 20 m). We designed a custom-built stereo system with mvBlueFox3 sensors<sup>2</sup> with a baseline of 40 cm to increase the maximum sensing range to 80 m. The stereo image was captured with a resolution of either  $1920 \times 1080$  or  $1280 \times 720$ . Additional details of our dataset can be found in our technical report<sup>3</sup>. All scenes were captured steadily with a tripod and slider in a hand-held fashion.

The DIML/CVL RGB-D dataset comprises 1 million stereo images, depth maps computed using a stereo matching algorithm [9], and a stereo confidence map [49], [50]. The depth maps in the DIML/CVL RGB-D dataset were generated by MC-CNN [9], which was a state-of-the-art stereo matching algorithm at the time of stereo image acquisition (from fall 2015 to summer 2017). Note that any type of stereo matching network can be adopted to obtain depth maps. The experiment we conducted using a more advanced stereo matching network is described in Section IV-C3.

Our dataset differs from existing datasets in the following respects:

- 1) It is comprised of 1 million RGB-D data for outdoor scenes.
- 2) Unlike existing outdoor datasets for driving scenes, ours was taken using hand-held stereo cameras for non-driving scenes.
- 3) Stereo confidence maps are provided together to quantify the accuracy of depth maps.

<sup>2</sup><https://www.matrix-vision.com/USB3-vision-camera-mvbluefox3.html>

<sup>3</sup>[https://dimlrgbd.github.io/downloads/technical\\_report.pdf](https://dimlrgbd.github.io/downloads/technical_report.pdf)

## F. Transfer of feature representation

In addition, we studied the effectiveness of our pre-trained monocular depth estimation network by transferring its feature representations as a pretext task for training other similar tasks such as road detection and semantic segmentation. The experimental results demonstrate that our pre-trained model is comparable to the ImageNet pre-trained model, which often serves as the pretext task for various vision applications [22], [57], [58].

We fine-tuned our pre-trained model using the KITTI road benchmark [24] for road detection and Cityscapes [25] for semantic segmentation, respectively. Both datasets include a small amount of manually annotated training data. We transfer both encoder and decoder weights of the pre-trained model into road detection and semantic segmentation. The softmax loss is used to fine-tune the network.

## G. Implementation details

We adopted the stereo matching network (CRL) [32] using the pre-trained model provided by the author and implemented the confidence measure network [45]. We implemented the monocular depth estimation network using the VLFeat MatConvNet library [60].

1) *Confidence measure network*: The network was trained using 50 image pairs consisting of the ground truth depth maps and stereo image pairs provided in the KITTI 2012 dataset. Following recent work on stereo confidence estimation [45], we trained the stereo confidence network with a relatively small number of image pairs. This is possible because the confidence estimation can be seen as a local inference procedure. It determines whether an estimated depth value at a reference pixel is reliable when a small patch centered at the reference pixel is given. Specifically, the confidence measure network is trained with a set of small patches (*i.e.*,  $9 \times 9$ ) that are paired with the pseudo-ground-truth depth map and ground truth confidence value [45]. The ground truth stereo confidence map is obtained by comparing the absolute difference between the predicted depth map and ground truth depth map (KITTI LiDAR points). Following the literature [24], we set the confidence value to 1 when the absolute difference is smaller than 3 pixels, and 0 otherwise.

During the training phase, we use a binary cross-entropy loss after applying a sigmoid function to the output of the network. We carried out 100 training epochs with an initial learning rate of 0.001, decreased by a factor of 10 every 10 epochs, and a momentum of 0.9.

Table I

QUANTITATIVE EVALUATION OF MONOCULAR DEPTH ESTIMATION ON THE EIGEN SPLIT [17] OF KITTI [24] DATASET. FOR DATASET, K = KITTI, CS = CITYSCAPES, FT = FLYINGTHINGS, AND DC = DIML/CVL.

Method	Training data	Approach	Dataset	RMSE(lin)	RMSE(log)	Abs rel	Sqr rel	$\delta < 1.25$	$\delta < 1.25^2$	$\delta < 1.25^3$
				Lower is better				Higher is better		
cap 80m										
Eigen <i>et al.</i> [17]	Left + LiDAR	<i>Sup.</i>	K	7.156	0.270	0.215	1.515	0.692	0.899	0.967
Godard <i>et al.</i> [21]	Stereo	<i>Self-sup.</i>	K	5.927	0.247	0.148	1.344	0.803	0.922	0.964
Godard <i>et al.</i> [21]	Stereo	<i>Self-sup.</i>	K + CS	5.311	0.219	0.124	1.076	0.847	0.942	0.973
Kuznetsov <i>et al.</i> [22]	Left + LiDAR	<i>Sup.</i>	K	4.815	0.194	0.122	0.763	0.845	0.957	<b>0.987</b>
Kuznetsov <i>et al.</i> [22]	Stereo + LiDAR	<i>Semi-sup.</i>	K	4.621	0.189	0.113	0.741	0.862	0.960	0.986
Luo <i>et al.</i> [20]	(Synthetic) Stereo + GT	<i>Semi-sup.</i>	K + FT	4.681	0.200	0.102	0.700	0.872	0.954	0.978
Monodepth2 [23]	Stereo	<i>Self-sup.</i>	K	4.750	0.196	0.109	0.873	0.874	0.957	0.979
Guo <i>et al.</i> [28]	Left + Pseudo GT	<i>Semi-sup.</i>	K	4.634	0.189	0.105	0.811	0.874	0.959	0.982
Tosi <i>et al.</i> [29]	Left + Pseudo GT	<i>Semi-sup.</i>	K	4.714	0.199	0.111	0.867	0.864	0.954	0.979
Tosi <i>et al.</i> [29]	Left + Pseudo GT	<i>Semi-sup.</i>	K + CS	4.351	0.184	0.096	0.673	0.890	0.961	0.981
Our Method	Left + Pseudo GT	<i>Semi-sup.</i>	K	4.599	0.183	0.099	0.748	0.880	0.959	0.983
Our Method	Left + Pseudo GT	<i>Semi-sup.</i>	K + DC	4.333	0.181	0.098	0.644	0.881	0.963	0.984
Our Method	Left + Pseudo GT	<i>Semi-sup.</i>	K + CS	4.286	0.177	0.097	0.641	0.882	0.963	0.984
Our Method	Left + Pseudo GT	<i>Semi-sup.</i>	K + CS + DC	<b>4.129</b>	<b>0.175</b>	<b>0.095</b>	<b>0.613</b>	<b>0.884</b>	<b>0.964</b>	0.986
Upper Bound			K	3.475	0.158	0.055	0.418	0.941	0.969	0.982
cap 50m										
Garg <i>et al.</i> [27]	Stereo	<i>Self-sup.</i>	K	5.104	0.273	0.169	1.080	0.740	0.904	0.962
Godard <i>et al.</i> [21]	Stereo	<i>Self-sup.</i>	K	4.471	0.232	0.140	0.976	0.818	0.931	0.969
Godard <i>et al.</i> [21]	Stereo	<i>Self-sup.</i>	K + CS	3.972	0.206	0.117	0.762	0.860	0.948	0.976
Kuznetsov <i>et al.</i> [22]	Stereo + LiDAR	<i>Semi-sup.</i>	K	3.518	0.179	0.108	0.595	0.875	0.964	<b>0.988</b>
Luo <i>et al.</i> [20]	(Synthetic) Stereo + GT	<i>Semi-sup.</i>	K + FT	3.503	0.187	0.097	0.539	0.885	0.960	0.981
Monodepth2 [23]	Stereo	<i>Self-sup.</i>	K	3.868	0.127	0.092	0.537	0.892	0.962	0.986
Our Method	Left + Pseudo GT	<i>Semi-sup.</i>	K + CS + DC	<b>3.162</b>	<b>0.162</b>	<b>0.091</b>	<b>0.505</b>	<b>0.901</b>	<b>0.969</b>	0.986

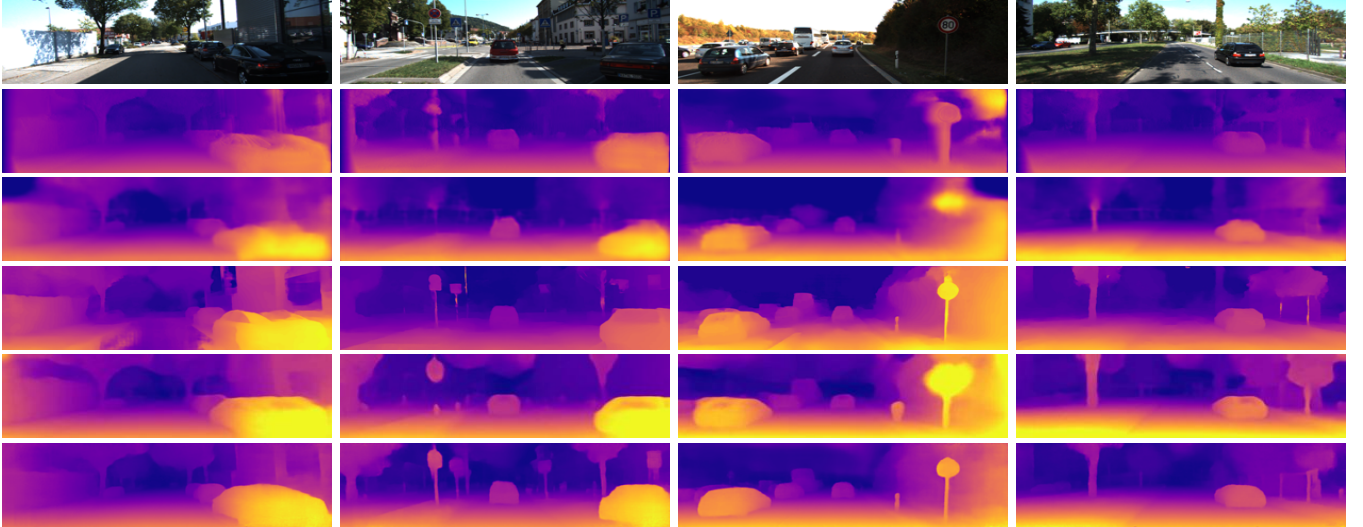


Figure 8. Qualitative results on the Eigen split [17] of the KITTI dataset [24] (from top to bottom): input image, Godard *et al.* [21], Kuznetsov *et al.* [22], Luo *et al.* [20], monodepth2 [23], and the proposed method.

2) *Monocular depth network*: To train the monocular depth estimation network, we collected stereo images from the KITTI, Cityscapes, and DIML/ CVL datasets, respectively. The pseudo-ground-truth training data were generated using the stereo matching network and stereo confidence map. The monocular student network was trained for 50 epochs with a batch size of 4. The Adam solver [61] was adopted for efficient stochastic optimization with a fixed learning rate of 0.001 and a momentum of 0.9.

#### IV. EXPERIMENTAL RESULTS

We conducted to validate the effectiveness of our semi-supervised monocular depth estimation via quantitative and

qualitative comparisons with state-of-the-art methods in outdoor scenes. For quantitative comparisons, we employ several metrics that have been used previously [17], [21]–[23]:

- Threshold: % s.t.  $\max\left(\frac{d_i}{u_i}, \frac{u_i}{d_i}\right) = \delta < thr$
- Abs rel:  $\frac{1}{N} \sum_i |d_i - u_i| / d_i$
- Sqr rel:  $\frac{1}{N} \sum_i \|d_i - u_i\|^2 / d_i$
- RMSE(lin):  $\sqrt{\frac{1}{N} \sum_i \|d_i - u_i\|^2}$
- RMSE(log):  $\sqrt{\frac{1}{N} \sum_i \|\log d_i - \log u_i\|^2}$

where  $u_i$  denotes the predicted depth at pixel  $i$ , and  $N$  is the total number of pixels. The monocular depth network was trained by selecting outdoor scenes similar to recently reported state-of-the-art methods [20]–[23]. Indoor datasets were not used in our experiments, however, training with indoor data

can be accomplished in the same manner. In our method, the stereo matching (teacher) networks are trained with the small amount of labeled training data. Then, the massive unlabeled training data, i.e., stereo image pairs, are used for training the monocular depth estimation (student) network. Therefore, the proposed method can be categorized as a semi-supervised learning approach.

### A. Dataset

We generated pseudo-ground-truth depth maps using stereo images provided in the KITTI [24], Cityscapes [25], and DIML/CVL datasets. All images were resized to  $620 \times 188$  for training and testing.

1) *KITTI (K)*: This dataset consists of outdoor driving scenes with sparse depth maps captured by the Velodyne LiDAR. The depth map is very sparse (density of less than 6%), and depth values are available only at the bottom parts of a color image. The dataset contains 42,382 rectified stereo pairs from 61 scenes, with a typical image sized  $1242 \times 375$  pixels. Following the Eigen split [17], we split the stereo image pairs into 22,600 images for training, 888 images for validation, and 697 images for testing.

2) *Cityscapes (CS)*: This dataset was originally constructed for semantic segmentation and provides manually annotated segmentation maps for 19 semantic classes, consisting of 2,975 images for training, 500 images for validation, and 1,525 images for testing. Additionally, they provide 22,973 stereo image pairs with spatial dimensions of  $2048 \times 1024$ . We split the stereo image pairs into 21,283 for training, 500 for validation, and 3,215 for testing. Following the literature [21], we cropped the stereo images by discarding the bottom 20% (the car hood) and then resized them.

3) *DIML/CVL (DC)*: Our outdoor dataset consists of 1 million stereo image pairs, depth maps, and stereo confidence maps. The original spatial resolution of this dataset is  $1920 \times 1080$  or  $1280 \times 720$ . Following previous practices [24], [25], of 1 million image pairs, we selected 23,500 image pairs. We split them into 22,000 images for training, 800 images for validation, and 700 images for testing. Input images were cropped and resized by discarding the bottom parts containing mostly ground.

### B. Comparison with state-of-the-art methods

We compared existing monocular depth estimation approaches including supervised [17], [22], self-supervised [21], [23], and semi-supervised methods [20], [22] with the proposed approach. Of the 1 million stereo image pairs, we selected 22,000 stereo images from various scenes as training data, similar to Eigen split [17]. The results in Table I indicate that our method consistently outperforms recent approaches, except for  $\delta < 1.25^3$ . Following the literature [20]–[23], the depth value was truncated at 80m or 50m. The proposed method was trained using various combinations of KITTI, Cityscapes, and our dataset. Eigen *et al.* [17] was trained using ground-truth depth maps augmented from the KITTI 2015 dataset. Godard *et al.* [21] and monodepth2 [23] proposed a self-supervised approach in which stereo images of the KITTI

and/or Cityscapes dataset were used. Although this method requires no ground truth depth maps during training, it is difficult to handle occlusions effectively and does not obtain a sharp depth boundary owing to the limitation of the image reconstruction loss. Kuznetsov *et al.* [22] employed both supervised regression loss and unsupervised reconstruction loss [21] to achieve a performance gain over existing supervised and self-supervised approaches [17], [21]. However, their method still requires ground truth depth maps as supervision. Luo *et al.* [20] proposed the use of a view synthesis network and stereo matching network in a unified framework. However, their model was mainly trained with synthetic FlyingThings3D data [11], and thus incurs domain gaps between synthetic and realistic images and cannot generalize well on real data. Such two-step method substantially increases the computational complexity (0.53 sec for an inference).

Guo *et al.* [28] trained the stereo matching network (teacher) using a synthetic stereo dataset. The pseudo labels for monocular depth estimation may thus incur undesired artifacts due to a large domain gap between synthetic and real datasets when training the monocular depth estimation network. To alleviate this domain gap, they performed the unsupervised fine-tuning. However, they do not fully consider the inherent errors such as occlusion, texture-less region, and repetitive patterns that occur in the stereo matching. Therefore, the performance of the generated pseudo label is limited. The method in Tosi *et al.* [29] is somewhat similar to the proposed method in that pseudo labels are refined by removing matching outliers through a left-right consistency check. However, the simple consistency check often fails to exclude the matching outliers of the pseudo depth labels which are constructed by a handcrafted stereo matching algorithm. In contrast, our method attempts to leverage such stereo knowledge more effectively with the proposed confidence guided loss and data ensemble strategy. The confidence guided loss can resolve the performance degeneration incurred by unreliable pseudo depth values in occlusion, texture-less region, and repetitive pattern that cannot be addressed using stereo matching methods only. The data ensemble method also helps improve the quality of the pseudo labels by combining multiple depth predictions beyond a single scale prediction. Furthermore, Tosi *et al.* [29] leverages a correlation layer to generate 3D cost volume from the input left image and synthesized right image during an inference, which demands a high computational complexity (0.48 sec for an inference). In contrast, our model can retain a relatively low complexity (0.21 sec) as the stereo knowledge is imposed only on the label (depth) space by training the student network with pseudo depth maps already computed from the teacher network.

Furthermore, qualitative results are provided in Fig. 8. Our monocular depth network achieves more accurate and edge-preserved depth maps than state-of-the-art methods [20]–[23].

### C. Discussion

1) *Impact of scene diversity of DIML/CVL dataset*: We conducted an ablation study to demonstrate that the DIML/CVL dataset is able to effectively train the algorithm in our semi-supervised approach. Fig. 9 shows the qualitative results



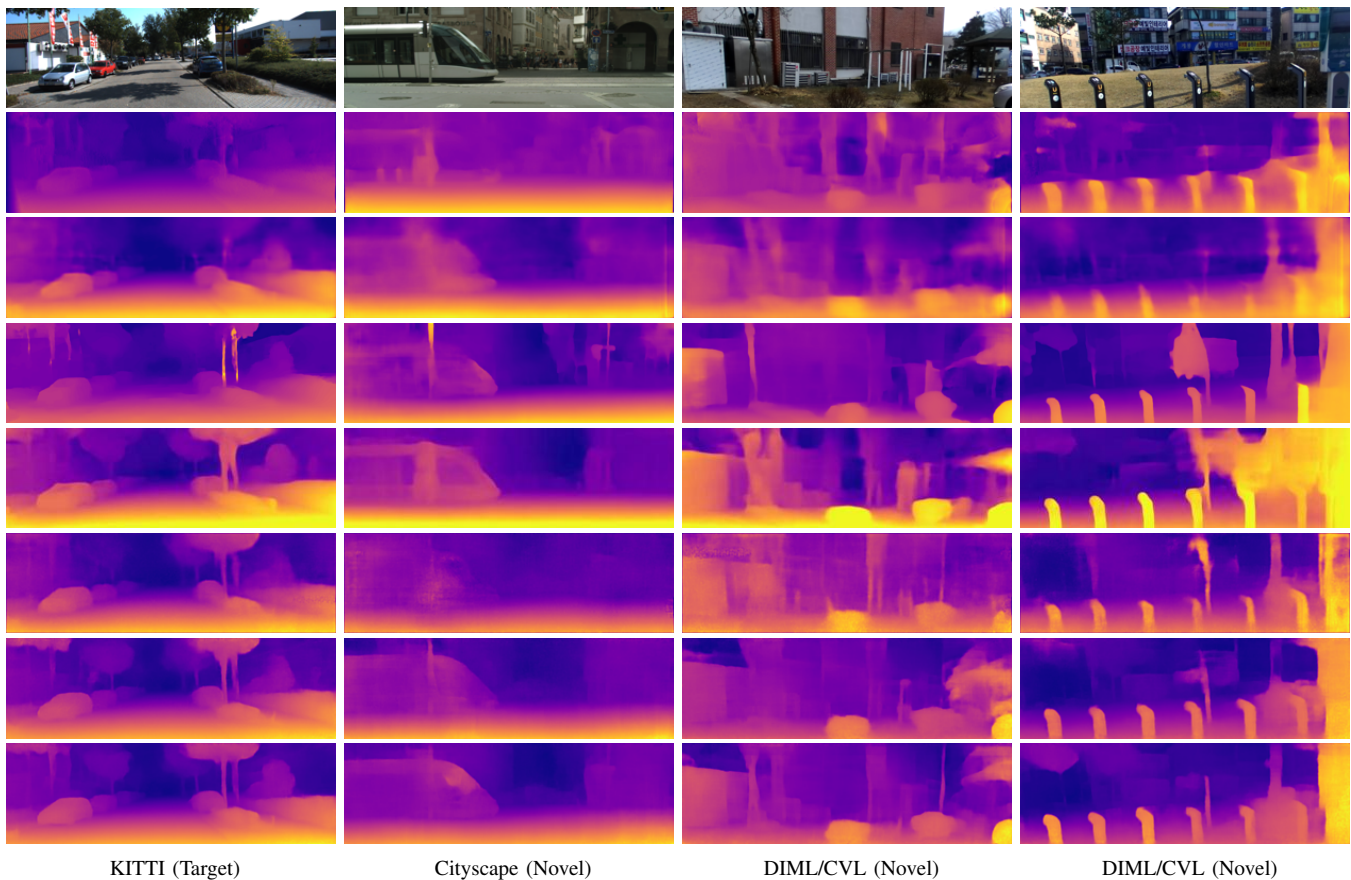


Figure 9. Qualitative comparison for the impact of scene diversity (from top to bottom): input image, Godard *et al.* [21], Kuznetsov *et al.* [22], Luo *et al.* [20], monodepth2 [23], and the proposed method trained with K, K + DC, KITTI + CS + DC, respectively.

Table II  
QUANTITATIVE EVALUATION OF MONOCULAR DEPTH ESTIMATION  
TRAINED WITH THE AMOUNT OF DATA DIFFERENTLY. THE NUMBER IN  
THE BRACKET INDICATES THE AMOUNT OF TRAINING DATA.

Method	Training data	RMSE(lin)	RMSE(log)	Abs rel
Our Method	CS (20k)	6.328	0.284	0.206
Our Method	DC (20k)	7.944	0.339	0.361
Our Method	DC (100k)	5.647	0.236	0.161
Our Method	DC (200k)	5.250	0.217	0.143
Our Method	K + DC (200k)	4.194	0.176	0.096

obtained by the proposed method for two different combinations (6<sup>th</sup> to 8<sup>th</sup> rows) of pseudo-ground-truth training data. Similar results were obtained in both the source domain and a novel domain when using the DIML/CVL dataset for training. This indicates that the proposed method addresses the generalization issue well when used in a novel domain. Specifically, using the DIML/CVL dataset for training leads to a performance gain in the novel domain, for example, when comparing the results of the proposed method trained with KITTI (6<sup>th</sup> row) and KITTI+DIML/CVL (7<sup>th</sup> row). The proposed method consistently generates smooth depth maps with sharp edges and recovers the overall scene layout well. In addition, we included the results (2<sup>nd</sup> to 5<sup>th</sup> rows) of state-of-the-art approaches [20]–[23] for a qualitative comparison. When inferring with the Cityscapes and DIML/CVL datasets (novel domain), they exhibited significant decreases in the performance. It was possible to train two of the methods [21], [23] with the Cityscapes and DIML/CVL dataset providing stereo

image pairs, as these methods do not require ground truth depth maps for training. In contrast, the other two methods we included in our evaluation [20] and [22] could not be trained with these datasets. The self-supervised approaches [21], [23] can be learned using training data from both the source and novel domains, but their self-supervised loss incurs blurry depth boundaries and inaccurate estimates in the occlusion, as revealed in the 2<sup>nd</sup> row and 5<sup>th</sup> row in Fig. 9.

Table II reports the depth accuracy using the Eigen split in the KITTI dataset. We conducted an experiment to verify the effectiveness of the DIML/CVL dataset according to the amount of training data. Because ground truth depth maps are provided in the KITTI dataset, the objective evaluation was conducted using this dataset. First, we trained our monocular depth network with the Cityscape and DIML/CVL datasets. With the same amount of data, the model trained with the Cityscape dataset is more effective than that trained with the DIML/CVL dataset. This is because the Cityscape dataset mostly contains driving scenes that are relatively similar to those in KITTI, whereas the DIML/CVL dataset includes both driving and non-driving scenes. Interestingly, when we increase the amount of DIML/CVL training data, the depth accuracy gradually improves even though the majority of images in the DIML/CVL dataset consist of non-driving scenes. This demonstrates the complementarity of the DIML/CVL data. Although we trained the monocular depth network using up to 200,000 images for training, the use of a larger amount of

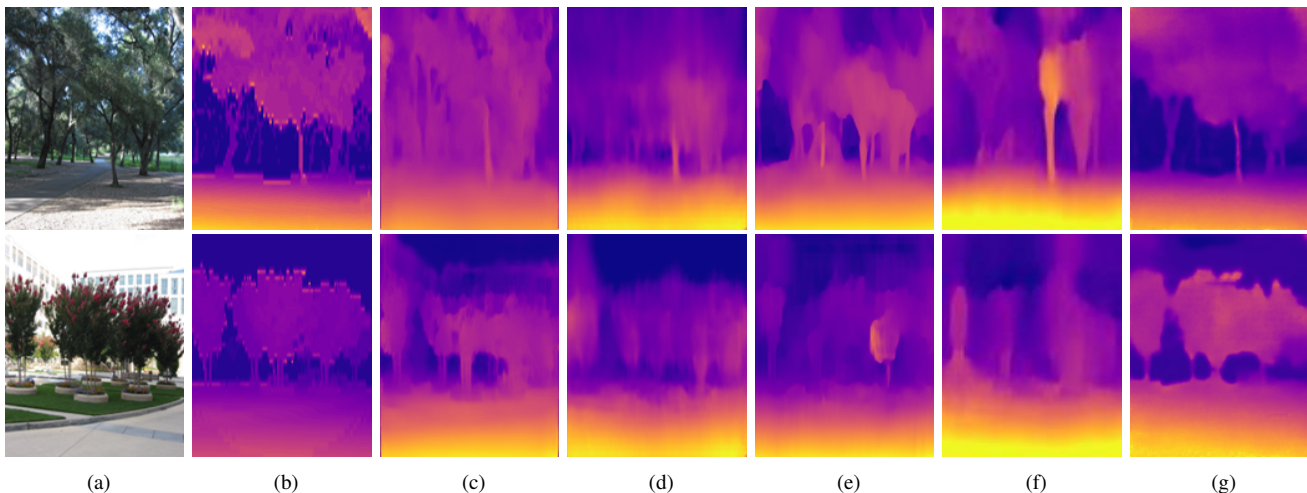


Figure 10. Qualitative results on the Make3D dataset [63]: (a) input image, (b) ground truth depth maps, (c) Godard *et al.* [21], (d) Kuznietsov *et al.* [22], (e) Luo *et al.* [20], (f) monodepth2 [23], and (g) the proposed method. None of the models were trained on Make3D.

Table III

QUANTITATIVE EVALUATION OF MONOCULAR DEPTH ESTIMATION ON THE MAKE3D TEST SET FOR VARIOUS APPROACHES.

Method	Dataset	RMSE(lin)	RMSE(log)	Absrel
Godard <i>et al.</i> [21]	K	11.762	0.193	0.544
Godard <i>et al.</i> [21]	K + CS	10.369	0.188	0.536
Kuznietsov <i>et al.</i> [22]	K	8.237	0.191	0.421
Luo <i>et al.</i> [20]	K + FT	8.184	0.185	0.428
monodepth2 [23]	K	8.238	0.201	0.374
Our Method	K	7.992	0.181	0.423
Our Method	DC	5.512	0.118	0.297
Our Method	K + DC	5.347	0.113	0.278
Our Method	K + CS + DC	5.136	0.107	0.265

training data would also be possible.

2) *Generalization to other datasets*: Following the existing literatures [21], [23], [34], we used two datasets to validate the generalization ability of our method to other datasets, Make3D [63] and camVid.

Similar to previous work [21], we cropped the top and bottom of the input images in Make3D dataset for testing. As for the ground truth depth maps, we reshaped and cropped the top and bottom to match our predictions. Owing to the limitation of the depth-sensing device that was used for capturing depth data, we removed depth values larger than 70 m. Table III lists the numerical results, including comparisons to several methods trained on non-Make3D training data. Our approach achieves the best performance among all non-Make3D-trained models. It should be noted that our results are more accurate than or comparable to those of the state-of-the-art methods even when trained with the KITTI dataset only. This indicates the effectiveness of using pseudo-ground-truth depth maps and a stereo confidence map for training the monocular depth network. Using DIML/CVL leads to a meaningful performance gain, showing the importance of our training data. Fig. 10 visualizes the results of the depth predictions. Our results preserve the overall structure and boundary of objects more accurately than other depth maps, and the resulting images we obtain with our method are the most similar to the ground truth depth maps.

Because the CamVid driving dataset [64] only provides

Table IV

QUANTITATIVE EVALUATION OF THE IMPACT OF THE STEREO MATCHING NETWORK.

Method	RMSE(lin)	Absrel	Time(s)
MC-CNN [9]	4.236	0.074	67
CRL [32]	3.578	0.055	0.47
GA-NET [65]	3.457	0.050	1.5

ground truth maps for semantic segmentation, we validate the effectiveness of our method by carrying out qualitative comparisons with state-of-the-art methods. In Fig. 11, despite the domain gap between the training and test images in terms of location, image characteristics, and camera parameters, our model still produces visually plausible depths in a novel domain.

3) *Impact of the stereo matching network*: Any stereo matching methods can be utilized as the teacher network in our framework. For more solid analysis, we conducted the impact of the stereo matching network using various stereo matching algorithms. Table IV shows the quantitative performance of depth accuracy and runtime of the stereo matching network. The accuracy of pseudo ground truth depth results from MC-CNN and CRL shows a large performance gap. This performance gap also brings a large gap in monocular depth accuracy (MC-CNN (RMSE (lin)): 5.294/ CRL (RMSE (lin)): 4.995). Although the pseudo depth maps of GA-Net has a higher accuracy than those of CRL, the performance gain in terms of the monocular depth accuracy was relatively marginal (GA-Net (RMSE (lin)): 4.946). This is because the confidence map excludes unreliable pseudo depth values effectively. Additionally, the fast inference time of CRL enables us to construct the pseudo RGB+D dataset efficiently. Thus, we adopted CRL as a teacher network in our framework.

#### D. Ablation Study

1) *Impact of label quality*: We first investigated the performance gain of the data ensemble and the tradeoff between accuracy and density that determines the performance of the pseudo-ground-truth depth maps. The KITTI dataset using

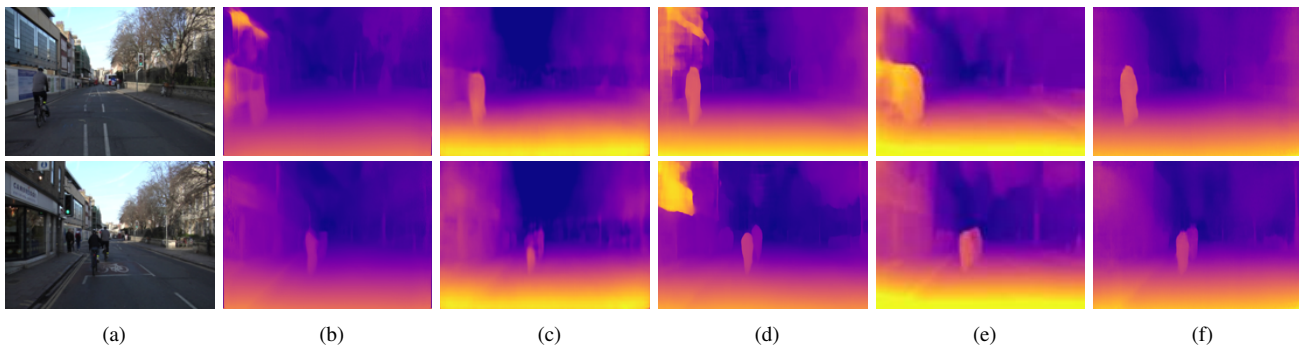


Figure 11. Qualitative results on the CamVid dataset [64]: (a) input image, (b) Godard *et al.* [21], (c) Kuznetsov *et al.* [22], (d) Luo *et al.* [20], (e) monodepth2 [23], and (f) the proposed method. None of the models were trained on CamVid.

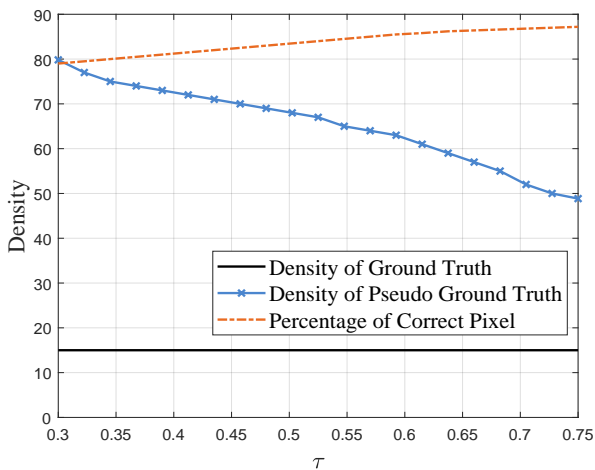


Figure 12. Quality analysis of pseudo-ground-truth depth maps controlled by the confidence threshold  $\tau$ .

the Eigen split [17] was used for an objective evaluation. The density and accuracy of the pseudo-ground-truth depth maps are controlled by the confidence threshold. However, the confidence map of the pseudo-ground-truth depth maps is not perfect. Nevertheless, excluding the depth outliers with the confidence map greatly improves the performance of the monocular depth estimation. This can be explained by two factors: 1) the quality of the pseudo-ground-truth depth maps in Fig. 12 and 2) the accuracy of the monocular depth estimation in Fig. 13.

In Fig. 12, we present the tradeoff between the accuracy and density of the pseudo-ground-truth depth maps. The black horizontal line shows the density of the KITTI ground truth depth maps. When the estimated confidence value  $C(p) \in [0, 1]$  at pixel  $p$ , we denote a set of pixels estimated to be correct ( $C(p) \geq \tau$ ) in the pseudo-ground-truth depth map as  $P$ . Note that the ground truth depth map is sparse, and thus, this is computed only at a set of valid depth pixels  $G$  in the ground truth depth map. Subsequently, the density of the pseudo-ground-truth depth map is computed using  $|P|/|G|$ . The accuracy is computed using  $|A|/|P|$ , where  $A$  indicates a set of correct pixels among  $P$ . The larger the value of  $\tau$ , the higher the accuracy of the pseudo-ground-truth depth maps becomes. However, this reduces the density of the depth maps. For instance, when  $\tau = 0.75$ , only half of the depth pixels are chosen as reliable.

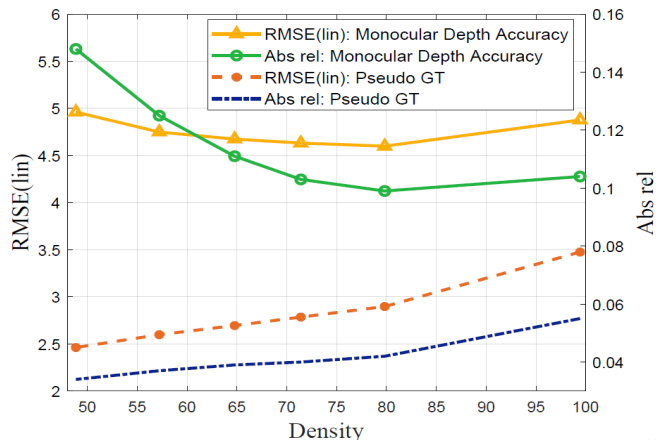


Figure 13. Analysis of the trade-off between the accuracy and density of pseudo-ground-truth depth maps.

Fig. 13 shows the accuracy measurements of the monocular depth network according to the confidence threshold  $\tau$ . Fig. 13 shows the RMSE (lin) and Abs rel of pseudo ground truth depth maps and monocular depth maps. We use a stereo confidence map to identify inaccurate depth values (Fig. 8) and avoid these values when training the monocular depth estimation network. Here, we investigate the inference accuracy of the monocular depth estimation according to the density of the pseudo-ground-truth depth maps. Training with more accurate pseudo-ground-truth depth maps (yet with a lower density) does not necessarily increase the accuracy of the monocular depth estimation. We achieved the best monocular depth accuracy when the density of the pseudo-ground-truth depth maps was approximately 80% with the confidence threshold  $\tau = 0.3$ . The monocular depth accuracy deteriorates when the density becomes 100%, that is, confidence is not used in (1) with the setting  $\tau = 0$ . Refer to Fig. 12 for the relationship between the density and the stereo confidence threshold  $\tau$ . As shown in Fig. 12, a pseudo-ground-truth depth map with a lower density tends to be more accurate. However, using more accurate pseudo-ground-truth depth maps does not necessarily increase the accuracy of a monocular depth network. This is because semi-dense pseudo-ground-truth depth maps often have no valid depth values around object boundaries or thin objects, in which case the monocular depth networks trained with these semi-dense depth maps may fail to recover reliable depth values around these regions.

The monocular depth accuracy at 80% is higher than that

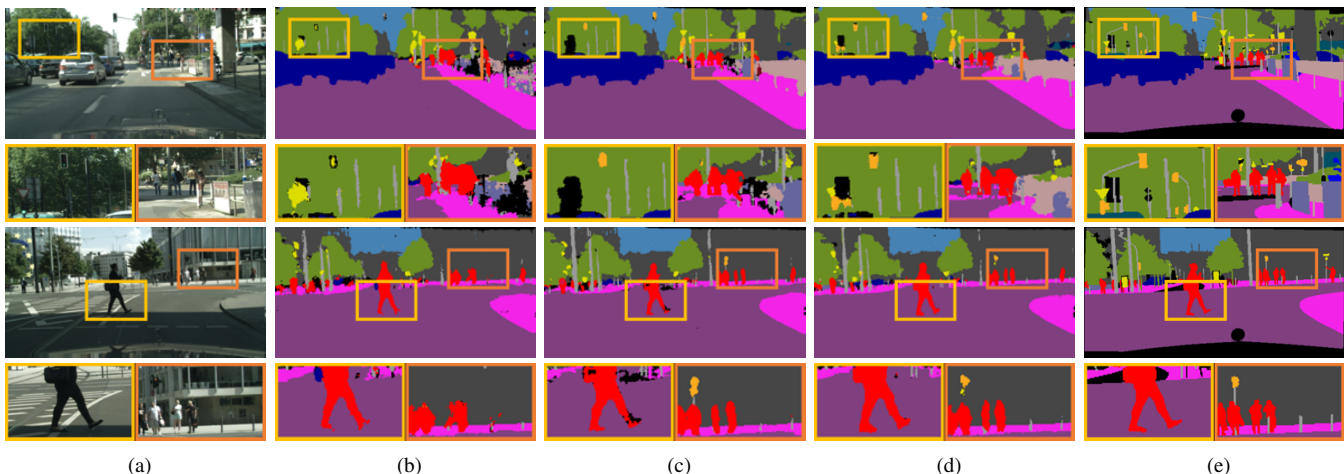


Figure 14. Semantic segmentation results on the Cityscapes dataset: (a) input images, (b)  $\sim$  (d) results of fine-tuning with different initialization methods. (b) scratch, (c) ImageNet pre-trained model [56], (d) our pre-trained model, and (e) ground truth image. The results indicate the benefit of using our pre-trained model.

at 100%. However, when the density of the pseudo-ground-truth depth maps is approximately 72%, the monocular depth accuracy deteriorates even though the pseudo-ground-truth depth maps are more accurate. This indicates that a tradeoff exists between the density and accuracy of the pseudo-ground-truth depth maps. In our experiment, the monocular depth network achieved the best accuracy when the density was approximately 80% ( $\tau = 0.3$ ). We generated pseudo-ground-truth training data using  $\tau = 0.3$  for all datasets.

2) *Impact of data ensemble*: Furthermore, we studied the effectiveness of the data ensemble when generating pseudo-ground-truth depth maps. The monocular depth accuracy was measured with the RMSE (lin) and the Abs rel. Without the data ensemble when setting  $\tau = 0$ , the RMSE(lin) and absolute relative error (Abs rel) are 4.995 and 0.109, respectively. When we used three scales for the data ensemble, the RMSE(lin) and the absolute relative error (Abs rel) were 4.877 and 0.104, respectively. The data ensemble achieves a meaningful gain in both metrics. In addition, we used the data ensemble for four scales but observed no marginal improvement, attaining an RMSE (lin) and absolute relative error (Abs rel) of 4.879 and 0.104, respectively. Therefore, we used three scales in the data ensemble for training.

### E. Transfer to high-level tasks

To investigate the applicability of our model trained for monocular depth prediction, we transferred the network parameters to scene understanding tasks such as semantic segmentation and road detection. This is in line with that the depth information can play an important role of guiding the semantic segmentation task as reported in the literatures of multi-task learning [67].

1) *Semantic segmentation*: We adopted the Cityscapes [25] dataset for training and evaluation. The methods were validated with the mean intersection over-union (IoU), which computes the mean value over all the classes including the background. Experiments were conducted with half-resolution images for fast computation. Following the approach in the literature [44], [52], we augmented the training data with

Table V  
QUANTITATIVE COMPARISON OF THE PRE-TRAINED MODEL USING THE SAME NETWORK ARCHITECTURE FOR THE CITYSCAPES BENCHMARK.

Semantic Segmentation		
Initialization	Pretext	mean IoU
Scratch	-	52.27
ImageNet pre-trained model [56]	Classification	66.27
K	Depth	62.82
K + DC	Depth	64.54
K + CS	Depth	65.02
K + CS + DC	Depth	65.47

images subjected to random scaling, random cropping, and horizontal flipping. The network was trained for 300 epochs with a batch size of 4. We used the Adam solver [61] for training with a weight decay of 0.0005 and an initial learning rate of 0.0001, which decreases by a factor of 10 every 10 epochs.

Table V and Fig. 14 report the quantitative and qualitative evaluation results obtained with three methods: a model learned from scratch, a model pre-trained with ImageNet [56], and ours. All results were obtained with the same encoder-decoder architecture that was used for the monocular depth estimation. Our pre-trained model significantly outperformed the model learned from scratch, with performance comparable to that of the model that was pre-trained with ImageNet, which is a massive manually labeled dataset. In addition, the results in Table V indicate that, the more accurate the monocular depth network, the higher the IoU in the semantic segmentation. Note that the experiments show that our monocular depth network based on the simple encoder-decoder architecture is a powerful proxy task for semantic segmentation.

2) *Road detection*: We investigated the effectiveness of our pre-trained model for road detection against the KITTI road benchmark [24], which provides 289 training images with annotated ground truth data and 290 test images. The benchmark is divided into three categories: single-lane roads with markings (UM), single-lanes road without markings (UU), and multi-lane roads with markings (UMM). Following the literature [68], the training data were augmented by subjecting them to various transformations. More specifically,

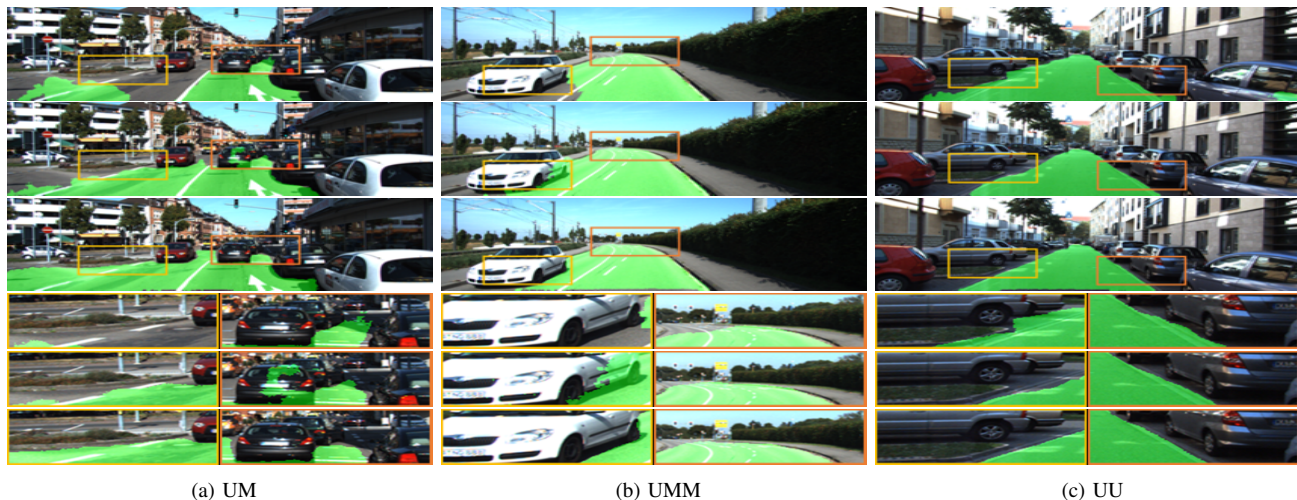


Figure 15. Road detection results on the KITTI dataset for different scene categories: (from top to bottom) model learned from scratch, model pre-trained on ImageNet [56], and our pre-trained model. Corresponding enlarged parts of the boxes are shown together.

Table VI  
QUANTITATIVE COMPARISON OF THE PRE-TRAINED MODEL USING THE SAME NETWORK ARCHITECTURE FOR THE KITTI ROAD BENCHMARK.

Road Detection			
Initialization	Pretext	Fmax	AP
Scratch	-	93.82	90.87
ImageNet pre-trained model [56]	Classification	94.28	92.25
K	Depth	94.41	92.04
K + DC	Depth	94.92	92.28
K + CS	Depth	95.12	93.09
K + CS + DC	Depth	95.65	94.46

we randomly scaled each image by a factor between 0.7 and 1.4, and performed color augmentation by adding values of -0.1 and 0.1 to the hue channel of the HSV space. For efficient stochastic optimization, we used the Adam optimizer [61], a fixed learning rate of 0.0001, and a weight decay of 0.0005. The road detection network was trained for 40 epochs with a batch size of 4.

For the quantitative comparison, we measured both the maximum F1-measurement (Fmax) and average precision (AP). Table VI indicates that our pre-trained model consistently outperforms both the model learned from scratch and the model pre-trained using ImageNet [56]. Similar to semantic segmentation, the accuracy of road detection is correlated with that of the monocular depth estimation network used as the pretext. Fig. 15 shows the excellent ability of our method to distinguish between roads and sidewalks.

## V. CONCLUSION

In this study, we propose a novel and effective approach to achieve monocular depth estimation. We adopt the student-teacher strategy, in which a shallow student network is trained by leveraging a deep and accurate teacher network. From massive stereo image pairs consisting of diverse outdoor scenes provided in the DIML/CVL dataset, we generated pseudo ground truth depth maps using the deep stereo matching network that serves as a teacher network. To improve the depth map of the teacher network, we applied the data ensemble and generated a stereo confidence map. The monocular depth network, which served as the student network,

was trained with the pseudo-ground-truth depth maps and stereo confidence-guided regression loss. We demonstrate that the proposed method is capable of leveraging stereo pairs on various domains, achieving state-of-the-art performance. Additionally, we show that training our model for monocular depth estimation provides semantically meaningful feature representations for high-level vision tasks. We expect that the proposed method serves as a key component in addressing the domain difference issue in various vision tasks.

## REFERENCES

- [1] B. Yang, S. Rosa, A. Markham, N. Trigoni, and H. Wen, "Dense 3D object reconstruction from a single depth view," in *IEEE Trans. Pattern Anal. Mach. Intell.*, 2018.
- [2] T. Cao, Z.-Y. Xiang, and J.-L. Liu, "Perception in disparity: An efficient navigation framework for autonomous vehicles with stereo cameras," in *IEEE Intell. Transp. Syst.*, 2015.
- [3] Y. Choi et al., "KAIST multi-spectral day/night data set for autonomous and assisted driving," in *IEEE Intell. Transp. Syst.*, 2018.
- [4] J. Cho, Y. Kim, H. Jung, C. Oh, J. Youn, and K. Sohn, "Multi-task self-supervised visual representation learning for monocular road segmentation," in *IEEE Int. Conf. on Multi. and Expo.*, 2018.
- [5] S. Kim, K. Park, K. Sohn, and S. Lin, "Unified depth prediction and intrinsic image decomposition from a single image via joint convolutional neural fields," in *Eur. Conf. on Comput. Vis.*, pp. 143–159, 2016.
- [6] S. Gupta, P. Arbelaez, R. Girshick, and J. Malik, "Indoor scene understanding with rgb-d images: Bottom-up segmentation, object detection and semantic segmentation," in *IEEE Int. J. Comput. Vis.*, vol. 112, no. 2, pp. 133–149, 2015.
- [7] D. Scharstein and R. Szeliski, "A taxonomy and evaluation of dense two-frame stereo correspondence algorithms," in *IEEE Int. J. Comput. Vis.*, vol. 47, pp. 7–42, 2002.
- [8] X. Han, T. Leung, Y. Jia, and R. Sukthankar, "Matchnet: Unifying feature and metric learning for patch-based matching," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, pp. 3279–3286, 2015.
- [9] J. Zbontar, Y. LeCun, "Computing the stereo matching cost with a convolutional neural network," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, pp. 1592–1599, 2015.
- [10] W. Luo, A.G. Schwing, and R. Urtasun, "Efficient deep learning for stereo matching," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, pp. 5695–5703, 2016.
- [11] N. Mayer, E. Ilg, P. Hausser, P. Fischer, D. Cremers, A. Dosovitskiy, and T. Brox, "A large dataset to train convolutional networks for disparity, optical flow, and scene flow estimation," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, pp. 4040–4048, 2016.
- [12] A. Kendall, H. Martirosyan, S. Dasgupta, and P. Henry, "End-to-end learning of geometry and context for deep stereo regression," in *Proc. Int. Conf. Comput. Vis.*, pp. 65–75, 2017.

- [13] JR. Chang, and YS. Chen, "Pyramid stereo matching network," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, pp. 5410–5418, 2018.
- [14] DC. Lee, M. Hebert, and T. Kanade, "Geometric reasoning for single image structure recovery," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2009.
- [15] D. Hoiem, AA. Efros, and M. Hebert, "Geometric context from a single image," in *Int. Conf. on Comput. Vis.*, pp. 654–661, 2005.
- [16] Y. Kim, H. Jung, D. Min, and K. Sohn, "Deep Monocular Depth Estimation via Integration of Global and Local Predictions," *IEEE Trans. Image Process.*, vol. 27, no. 8, pp. 4131–4144, 2018.
- [17] D. Eigen, C. Puhrsch, and R. Fergus, "Depth map prediction from a single image using a multi-scale deep network," in *Proc. Advances in Neural Inf. Process. Syst.*, pp. 2366–2374, 2014.
- [18] N. Silberman, P. Kohli, D. Hoiem, and R. Fergus, "Indoor segmentation and support inference from RGBD images," in *Proc. Eur. Conf. Comput. Vis.*, pp. 746–760, 2012.
- [19] F. Liu, C. Shen, and G. Lin, "Deep convolutional neural fields for depth estimation from a single image," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, pp. 5162–5170, 2015.
- [20] Y. Luo, J. Ren, M. Lin, J. Pang, and W. Sun, "Single view stereo matching," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, pp. 155–163, 2018.
- [21] 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.*, 2016.
- [22] Y. Kuznetsov, J. Tsai, J. Stuckler, and B. Leibe, "Semi-supervised deep learning for monocular depth map prediction," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, pp. 6647–6655, 2017.
- [23] C. Godard, O.M. Aodha, M. Firman, and G.J. Brostow, "Digging into self-supervised monocular depth estimation," in *Proc. Int. Conf. Comput. Vis.*, 2019.
- [24] 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.*, 2012.
- [25] M. Cordts, M. Omran, S. Ramos, and T. Rehfeld, "The Cityscapes Dataset for Semantic Urban Scene Understanding," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, pp. 3213–3223, 2016.
- [26] [http://diml.yonsei.ac.kr/DIML\\_rgb\\_d\\_dataset/paper/technical\\_report.pdf](http://diml.yonsei.ac.kr/DIML_rgb_d_dataset/paper/technical_report.pdf).
- [27] R. Garg, VK. BG, G. Carneiro, and I. Reid, "Unsupervised cnn for single view depth estimation: Geometry to the rescue," in *Proc. Eur. Conf. Comput. Vis.*, pp. 740–756, 2016.
- [28] X. Guo, H. Li, S. Yi, J. Ren, and X. Wang, "Learning monocular depth by distilling cross-domain stereo networks," in *Proc. Eur. Conf. Comput. Vis.*, 2018.
- [29] F. Tosi, F. Aleotti, M. Poggi, and S. Mattoccia, "Learning monocular depth estimation infusing traditional stereo knowledge," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2019.
- [30] 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.
- [31] J. Xie, R. Girshick, and A. Farhadi, "Deep3d: Fully automatic 2d-to-3d video conversion with deep convolutional neural networks," in *Proc. Eur. Conf. Comput. Vis.*, pp. 842–857, 2016.
- [32] J. Pang, W. Sun, JSJ. Ren, C. Yang, and Q. Yan, "Cascade Residual Learning: A Two-Stage Convolutional Neural Network for Stereo Matching," in *Proc. Int. Conf. Comput. Vis. Work.*, vol. 7, no. 8, pp. 887–895, 2017.
- [33] T. Zhou, M. Brown, N. Snavely, and DG. Lowe, "Unsupervised learning of depth and ego-motion from video," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2017.
- [34] S. Zhao, H. Fu, M. Gong, and D. Tao, "Geometry-aware symmetric domain adaptation for monocular depth estimation," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2019.
- [35] K. Lasinger, R. Ranftl, K. Schindler, and V. Koltun, "Towards robust monocular depth estimation: Mixing datasets for zero-shot cross-dataset transfer," in *arXiv*, 2019.
- [36] J. Ba, and R. Caruana, "Do deep nets really need to be deep," in *Proc. Advances in Neural Inf. Process. Syst.*, pp. 2654–2662, 2014.
- [37] G. Hinton, O. Vinyals, and J. Dean, "Distilling the knowledge in a neural network," in *arXiv*, 2015.
- [38] Z. Xu, YC. Hsu, and J. Huang, "Training Student Networks for Acceleration with Conditional Adversarial Networks," in *Proc. Brit. Mach. Vis. Conf.*, 2018.
- [39] J. Ye, Y. Ji, X. Wang, K. Ou, D. Tao, and M. Song, "Student Becoming the Master: Knowledge Amalgamation for Joint Scene Parsing, Depth Estimation, and More," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2019.
- [40] 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.
- [41] M. Noroozi, and P. Favaro, "Unsupervised learning of visual representations by solving jigsaw puzzles," in *Eur. Conf. on Comput. Vis.*, pp. 69–84, 2016.
- [42] D. Pathak, P. Krahenbuhl, and J. Donahue, "Context encoders: Feature learning by inpainting," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, pp. 2536–2544, 2016.
- [43] G. Larsson, M. Maire, and G. Shakhnarovich, "Colorization as a proxy task for visual understanding," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2017.
- [44] O. Russakovsky et al., "ImageNet Large scale visual recognition challenge," in *Int. Journ. Comput. Vis.*, vol. 115, no. 3, pp. 211–252, 2015.
- [45] M. Poggi, and S. Mattoccia, "Learning from scratch a confidence measure," in *Proc. Brit. Mach. Vis. Conf.*, vol. 10, pp. 1–13, 2016.
- [46] I. Radosavovic, P. Dollár, R. Girshick, G. Gkioxari, and K. He, "Data distillation: Towards omni-supervised learning," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2018.
- [47] Stereolabs. Zed camera. URL <https://www.stereolabs.com/>.
- [48] <https://www.matrix-vision.com/USB3-vision-camera-mvbluefox3.html>.
- [49] M. Park and K. Yoon, "Leveraging stereo matching with learning-based confidence measures," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2015.
- [50] S. Kim, D. Min, B. Ham, S. Kim, and K. Sohn, "Deep stereo confidence prediction for depth estimation," in *Proc. IEEE Conf. Image. Process.*, 2017.
- [51] K. He, X. Zhnag, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, pp. 770–778, 2016.
- [52] J. Long, E. Shelhamer, and T. Darrell, "Fully convolutional networks for semantic segmentation," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, pp. 3431–3440, 2015.
- [53] Y. Ding, L. Wang, D. Fan, and B. Gong, "A semi-supervised two-stage approach to learning from noisy labels," in *Proc. IEEE Winter Conf. Appl. Comput. Vis.*, 2018.
- [54] J. Li, Y. Wong, Q. Zhao, and M.S. Kankanhalli, "Learning to learn from noisy labeled data," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2019.
- [55] O. Ronneberger, P. Fischer, and T. Brox, "U-net: Convolutional networks for biomedical image segmentation," in *Int. Conf. on Medical Image Computing and computer-assisted intervention.*, 2015.
- [56] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition?," in *arXiv*, 2014.
- [57] V. Badrinarayanan, A. Kendall, and R. Cipolla, "Segnet: A deep convolutional encoder-decoder architecture for image segmentation," in *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 39, no. 12, pp. 2481–2495, 2017.
- [58] M. Teichmann, M. Weber, M. Zoellner, and R. Cipolla, "MultiNet: Real-time Joint Semantic Reasoning for Autonomous Driving," in *CORR.*, 2016.
- [59] S. Ioffe and C. Szegedy, "Batch normalization: Accelerating deep network training by reducing internal covariate shift," in *Proc. Int. Conf. Mach. Learn.*, pp. 448–456, 2016.
- [60] A. Vedaldi, and K. Lnc, "Matconvnet: Convolutional neural networks for matlab," in *Proc. ACM Int. Conf. Multimedia*, pp. 689–692, 2015.
- [61] D. Kingma and J. Ba, "Adam: A method for stochastic optimization," in *arXiv*, 2014.
- [62] Velodyne. Accessed: Feb. 15, 2017. [Online]. Available: <http://velodynelidar.com/>
- [63] A. Saxena, M. Sun, and A. Y. Ng, "Make3D: Learning 3D scene structure from a single still image," in *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 31, no. 5, pp. 824–840, May 2009.
- [64] G. J. Brostow, J. Fauqueur, and R. Cipolla, "Semantic object classes in video: A high-definition ground truth database," in *Pattern Recognit. Letters*, 2009.
- [65] F. Zhang, V. Prisacariu, R. Yang, and P.H.S. Torr, "Ga-net: Guided aggregation net for end-to-end stereo matching," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2019.
- [66] Cityscapes Benchmark Suite. URL <https://www.cityscapes-dataset.com/benchmarks/#scene-labeling-task>.
- [67] H. Jiang, G. Larsson, M. Maire Greg Shakhnarovich and E. Learned-Miller, "Self-supervised relative depth learning for urban scene understanding," in *Proc. Eur. Conf. Comput. Vis.*, 2018.
- [68] GL. Oliveira, W. Burgard, and T. Brox, "Efficient Deep Methods for Monocular Road Segmentation," in *IEEE/RSJ Int. Conf. on Intell. Robot. and Syst.*, 2016.