

Virtual Normal: Enforcing Geometric Constraints for Accurate and Robust Depth Prediction

Wei Yin, Yifan Liu, Chunhua Shen

Abstract—Monocular depth prediction plays a crucial role in understanding 3D scene geometry. Although recent methods have achieved impressive progress in the evaluation metrics such as the pixel-wise relative error, most methods neglect the geometric constraints in the 3D space. In this work, we show the importance of the high-order 3D geometric constraints for depth prediction. By designing a loss term that enforces a simple geometric constraint, namely, *virtual normal* directions determined by randomly sampled three points in the reconstructed 3D space, we significantly improve the accuracy and robustness of monocular depth estimation.

Importantly, the virtual normal loss can not only improve the performance of learning metric depth, but also disentangle the scale information and enrich the model with better shape information. Therefore, when not having access to absolute metric depth training data, we can use virtual normal to learn a robust affine-invariant depth generated on diverse scenes. Our experiments demonstrate state-of-the-art results of learning metric depth on NYU Depth-V2 and KITTI. From the high-quality predicted depth, we are now able to recover good 3D structures of the scene such as the point cloud and surface normal directly, eliminating the necessity of relying on additional models as was previously done. To demonstrate the excellent generalization capability of learning affine-invariant depth on diverse data with the virtual normal loss, we construct a large-scale and diverse dataset for training affine-invariant depth, termed Diverse Scene Depth dataset (DiverseDepth), and test on five datasets with the zero-shot test setting. Code is available at: <https://git.io/Depth>

Index Terms—Monocular depth estimation, 3D from single images, surface normal, virtual normal

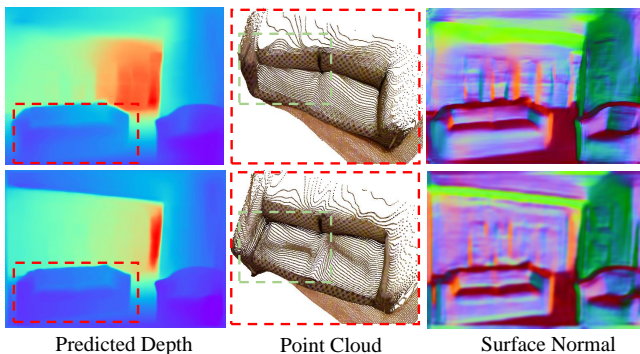


Fig. 1 – Example results of Hu *et al.* [5] (first row) and our method (second row). By enforcing the geometric constraints of virtual normals, our reconstructed 3D point cloud can represent better shape of sofa (see the part in the green dash box) and the recovered surface normal shows much fewer errors (in green) even though the absolute relative error (rel) of our predicted depth is only slightly better than Hu *et al.* (0.108 vs. 0.115).

1 INTRODUCTION

Monocular depth estimation aims to predict distances between scene objects and the camera from a still monocular image. Depth provides crucial information for understanding 3D scenes. As it is not trivial to enforce geometric constraints to recover depth from monocular still images, it is a challenging problem and various data-driven approaches were proposed to exploit comprehensive cues [1], [2], [3], [4].

Monocular depth prediction is an ill-posed problem because multiple 3D scenes can be projected to a same 2D

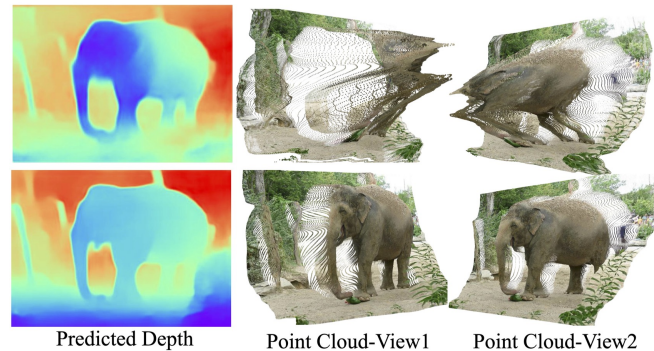


Fig. 2 – Qualitative comparison of depth and reconstructed 3D point cloud between our method and that of the recent learning relative depth method [6]. The first row is the predicted depth and reconstructed 3D point cloud from the depth of Xian *et al.* [6], while the second row is ours. The relative depth model fails to recover the 3D geometric shape of the scene (see the distorted elephant and ground area). Ours does much better. Note that this test image is sampled from the DIW dataset, which does not overlap with our training data.

image. Deep convolutional neural networks (CNN) based methods [7], [8], [9], [10] have achieved impressive performance on public benchmark datasets. The majority of works in the literature employ pixel-wise metric supervision to produce metric depth maps typically on some specific scenes, such as indoor environments, but in general do not work well on diverse scenes. The second line of works aim to address the issue of generalization to multiple scene data by learning with relative depth, such that large-scale

Work was done when all the authors were with The University of Adelaide, Australia. C. Shen is with Monash University, Australia. Corresponding author: C. Shen (email: chunhua@me.com).

datasets of diverse scenes may be collected. A typical example is the depth-in-the-wild (DIW) dataset [11]. Such methods often only explore the pair-wise ordinal relations for learning, and only the relative depth can be predicted. A clear drawback is that these models fail to recover the high-quality geometric 3D shapes, as only ordinal relations are used in learning.

These learning metric depth methods often formulate the optimization problem as either point-wise regression or classification. That is, with the *i.i.d.* assumption, the overall loss is summed over all pixels. To improve the performance, endeavors have been made to employ other loss terms besides the pixel-wise term. For example, a continuous conditional random field (CRF) [12] is used for depth prediction, which takes pair-wise information into account. Other geometric relations [13] are also exploited, incorporating the depth-to-surface-normal mutual transformation inside the optimization pipeline [13]. Note that, for the above methods, the geometric relations are ‘local’ in the sense that they are extracted from a small neighborhood in either 2D or 3D. Surface normal is ‘local’ by nature, as it is defined by the local tangent plane. As the ground-truth depth maps of most datasets are captured by consumer-level sensors, such as the Kinect, depth values can fluctuate considerably. Such noisy measurement would adversely affect the precision and subsequently the effectiveness of those local constraints, inevitably. Moreover, local constraints calculated over a small neighborhood have not fully exploited the structure information of the scene geometry that may be possibly used to boost the performance.

To address these limitations, here we propose a more stable geometric constraint from a global perspective to take long-range relations into account for predicting depth, termed *virtual normal*. A few previous methods already made use of 3D geometric information in depth estimation, almost all of which focus on using surface normal. *We instead reconstruct the 3D point cloud from the estimated depth map explicitly.* In other words, we generate the 3D scene by lifting each RGB pixel in the 2D image to its corresponding 3D coordinate with the estimated depth map. This 3D point cloud serves as an intermediate representation. With the reconstructed point cloud, we can exploit many kinds of 3D geometry information, not limited to the surface normal. Here we consider the long-range dependency in the 3D space by randomly sampling three non-colinear points with large distances to form a *virtual plane*, of which the normal vector is the proposed *virtual normal* (VN). The direction divergence between ground-truth and predicted VN can serve as a high-order 3D geometry loss. Owing to the long-range sampling of points, the adverse impact caused by noises in depth measurement is much alleviated compared to the computation of the surface normal, making VN significantly more accurate. Moreover, with random sampling we can obtain numerous such constraints, encoding the global 3D geometric information. Second, *by converting estimated depth maps from images to 3D point cloud representations, it opens many possibilities of exploiting algorithms of dealing with 3D point cloud for processing 2D images and 2.5D depth data.* Here we show one instance of such possibilities.

By combining the high-order geometric supervision and the pixel-wise depth supervision, our network can predict

not only an accurate depth map but also the high-quality 3D point cloud, subsequently other geometry information such as the surface normal. It is worth noting that we do not use an extra model or introduce network branches for estimating the surface normal. Instead, we directly compute surface normal from the reconstructed point cloud. The second row of Fig. 1 demonstrates an example of our results. By contrast, although the previously state-of-the-art method [5] predicts the depth with low errors, the reconstructed point cloud is far away from the original shape (see, *e.g.*, left part of ‘sofa’). The surface normal also contains many errors. *We are probably the first to achieve high-quality monocular depth and surface normal prediction with a single network, without training the separate networks or decoders for two tasks.*

As the model supervised with the virtual normal loss can recover high-quality 3D shape information, we propose to use it on diverse data to explore the 3D structure of scenes and disentangle the scale information. This can produce a robust model by leveraging the 3D structure information on diverse data. By contrast, previous methods mainly exploit the uniformity of pair-wise ordinal relations [6], [11], [14] on diverse scenes, *i.e.*, learning relative depth. Such predicted depth is not sufficient to encode rich geometric information. For example, in the first row of Fig. 2, we observe that learning ordinal relations fails to recover the shape of the flat ground and elephant in the image. By contrast, directly minimizing the pixel-wise divergence [1], [2], [7] may better recover scene geometry. However, such learning objective cannot be fulfilled on diverse scenes without metric depth annotation available, which are of different scales. In contrast, we propose to alleviate this difficulty of depth prediction by explicitly removing the depth scales during training on diverse scenes. We propose to apply the virtual normal loss to learn affine-invariant depth.

A large-scale diverse dataset is important for improving the model’s generalization capability. Existing RGB-D datasets can be summarized into two categories: 1) RGB-depth pairs captured by a depth sensor of high precision, typically accommodating only few scenes as it can be very costly to acquire a very large dataset of diverse scenes. For example, the KITTI dataset [15] is captured with LIDAR on road scenes only, while the NYU dataset [16] only contains several indoor rooms. 2) Images with much more diverse scenes that are available online and can be annotated with coarse depth with reasonable effort. The large-scale DIW dataset is manually annotated with only one pair of ordinal depth relations for each image [11]. In contrast, to construct our large and diverse dataset, we harvest stereoscopic videos and images with diverse contents and use stereo matching methods to obtain depth maps. The dataset contains both rigid and non-rigid foregrounds, such as humans, animals, and cars. Ours is considerably more diverse than metric depth datasets, while it contains more scene structure information than existing relative depth datasets because depth in our dataset is metric depth up to an affine transformation. Furthermore, to enrich more indoor and street scenes, we have sampled some images from Taskonomy [17] and DIML [18].

For learning metric depth, experimental results on NYUD-v2 [16] and KITTI [15] show that the virtual normal loss can boost performance significantly and achieve state-

of-the-art accuracy. Besides, from the reconstructed point cloud, we directly calculate the surface normal, with the accuracy being *on par* with that of specific CNN based surface normal estimation methods. Secondly, for learning affine-invariant depth, by training on our proposed diverse data with the virtual normal loss, our method outperforms previous methods on five datasets by a large margin with the zero-shot test setting, demonstrating the excellent generalization capacity of the learned model on a wide range of scenes.

In summary, our main contributions are as follows.

- We demonstrate the effectiveness of enforcing a high-order geometric constraint in the 3D space for the depth prediction task. Such global geometry information is instantiated with a simple yet effective concept termed *virtual normal* (VN). By enforcing a loss defined on VNs, we demonstrate the importance of 3D geometry information in depth estimation, and design a simple loss to exploit it.
- Our method can reconstruct high-quality 3D scene point clouds, from which other 3D geometry features may be calculated, such as the surface normal. In essence, we show that for depth estimation, one should not consider the information represented by depth only. Instead, converting depth into 3D point clouds and exploiting 3D geometry is likely to improve many tasks, including depth estimation.
- We propose to apply the virtual normal loss to enforce the model to learn affine-invariant depth on diverse scene data, which ensures both good generalization and high-quality geometric shapes of scenes. Experiments on five zero-shot datasets demonstrate that our method outperforms previous methods noticeably.
- To facilitate the learning on diverse scenes, we construct a new large scale and high-diversity RGB-D dataset, termed DiverseDepth.

2 RELATED WORK

Monocular depth estimation. Monocular depth estimation is important for many robotic and vision applications. According to the supervision, monocular depth estimation can be categorized into supervised methods [2], [6], [19] and unsupervised/self-supervised learning methods [20], [21], [22]. Saxena *et al.* [23] is among the first ones proposing to predict depth from a single image. They construct a Markov Random Field (MRF) model that incorporates multi-scale local and global image features. Later, a few methods [24], [25] based on the probabilistic model are proposed. When the powerful deep convolutional neural network emerges and benefits various computer vision tasks, many CNN-based methods are also proposed. Eigen *et al.* [7], [26] propose the first multi-scale network for dense prediction, including monocular depth prediction, surface normal estimation, and semantic estimation. Liu *et al.* [27] proposes to combine the CNN and CRF for depth estimation. Besides the study on the network architecture, many endeavours [1], [2], [4], [10], [13], [28] have been done on leveraging supervisions to improve the performance. Some works [1], [2], [10], [28] model the depth prediction as a classification problem. Qi *et al.* [13] propose to jointly predict the surface normal

and depth, which can refine the depth map based on the constraints from the surface normal.

Apart from these supervised learning methods, some work formulate unsupervised learning or self-supervised learning approaches [20], [21], [29], [30], [31] to address the lack of massive ground-truth depth training data. Zhou *et al.* [29] are among the first ones to demonstrate an approach to jointly predict the depth and the ego-motion from the monocular video. They use an image alignment loss, which is obtained by warping the source image to the neighboring frames with the predicted depth and ego-motion, to supervise the network. To improve the scale consistency between consecutive frames, Bian *et al.* [21] propose the geometry consistency loss. Ranjan *et al.* [32] propose to solve multiple low-level vision problems simultaneously, including depth, camera motion, optical flow, and moving objects segmentation, because such fundamental problems are coupled together through geometric constraints. Furthermore, several works [33], [34], [35] propose to leverage the geometric relations between consecutive frames.

Existing methods for metric depth estimation are difficult to generalize to diverse test scenes, mainly due to the lack of sufficiently large training datasets. To improve generalization, methods that learn relative depth [6], [11], [19], [36], [37], [38], [39] are proposed, as relative depth is much easier to obtain than metric depth. Chen *et al.* [11] construct the first large-scale and highly diverse dataset for learning the relative depth. As they use the ordinal relations and there is a large-scale dataset for training, their method can produce a model with good generalization. To construct better quality training data, Xian *et al.* [6], [40] collect stereo images and use stereo matching methods to obtain the inverse depth. Ranftl *et al.* [38] propose the scale-shift invariant loss to leverage the training on multi-source data, which can achieve promising generalization on diverse scenes.

RGB-D datasets. Datasets [3], [15], [16], [24], [41] are significant for the advancement of data-driven depth prediction methods. According to the quality of the ground truth depth, these datasets can be summarized into two categories. Depth sensors are used to directly collect high-quality RGB-D pairs, which can construct accurate metric depth dataset. Make3D [24] is the first outdoor RGB-D dataset constructed for monocular depth prediction study. KITTI [15] and NYU [16] are captured by LIDAR on outdoor streets and Kinect in indoor rooms. Larger-scale RGB-D datasets are also constructed, such as ScanNet [3], Taskonomy [17], DIML [18], DIODE [41]. These datasets usually only contain very limited scenes.

To improve the generalization of depth estimation methods on diverse scenes, several large-scale and diverse datasets are constructed, but the depth is not of high quality. Chen *et al.* [11] construct the largest RGB-D dataset, where the ground-truth depth maps are manually annotated with only one pair of ordinal relations. Similarly, Youtube3D [14] is also constructed to learn the relative depth but with more pairs of ordinal relations. MegaDepth [36] employs structure from motion to construct the depth supervision on the still and rigid scenes. To include more non-rigid and diverse scenes, Xian *et al.* [6] and Wang *et al.* [42] employ optical flow methods to construct datasets of relative depth. Chen *et*

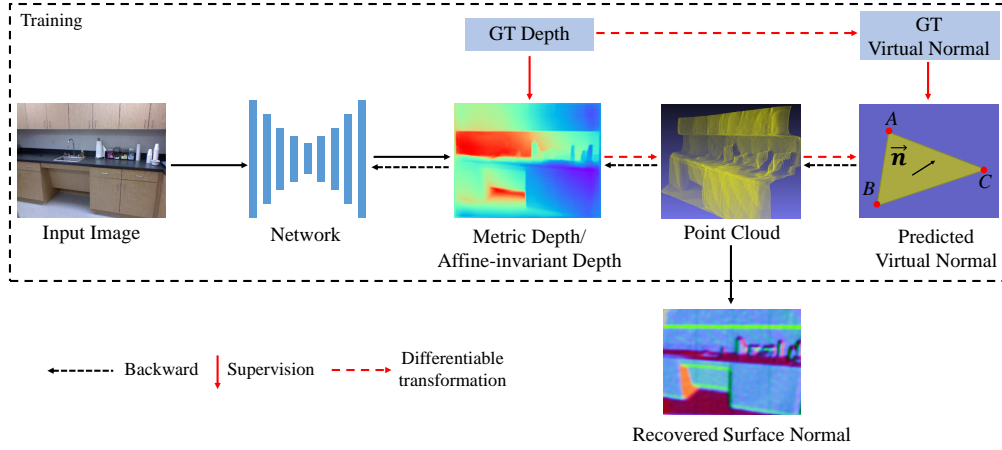


Fig. 3 – Overview of our framework. Our monocular depth prediction method inputs an image to an encoder-decoder network and outputs the depth. During training, we reconstruct the 3D point cloud from the depth and construct a virtual normal loss on the 3D point cloud to supervise the network. During the inference, the surface normal can be directly recovered from the point cloud.

al. [19] propose the diverse OASIS dataset, which includes both depth ordinal annotations and camera intrinsic parameters.

Curriculum learning. For many applications, introducing concepts in ascending difficulty to the learner is a common practice. Several works have demonstrated that curriculum learning [43], [44], [45] can boost the performance of deep learning methods. Weinshall *et al.* [43] combine the transfer learning and curriculum learning methods to construct a better curriculum, which can improve both the speed of convergence and the final accuracy. Hacohen and Weinshall [44] propose a bootstrapping method to train the network by self-tutoring.

3 OUR METHOD

The overall pipeline is illustrated in Fig. 3. We take an RGB image I_{in} as the input of an encoder-decoder network and predict the affine-invariant depth map D_{pred} . From the D_{pred} , the 3D scene point cloud P_{pred} can be reconstructed. The ground truth point cloud P_{gt} is reconstructed from D_{gt} . In order to improve the generation of the method, we firstly construct a large-scale and highly diverse dataset, DiverseDepth. Furthermore, we enforce a geometric loss, virtual normal loss, on the 3D point cloud and a scale and shift invariant loss on the depth to lead the model to learn the affine-invariant depth.

3.1 Virtual Normal

In order to employ the affine-invariant supervision, we propose a geometric loss. We consider a more stable geometric constraint from a global perspective to take long-range relations into account for predicting affine-invariant depth, termed virtual normal. With the predicted depth, the 3D point cloud can be reconstructed based on the pinhole camera model. For each pixel $p_i(u_i, v_i)$, the 3D location $P_i(x_i, y_i, z_i)$ in the world coordinate can be obtained by the perspective projection. We set the camera coordinate as the

world coordinate. Then the 3D coordinate P_i is denoted as follows:

$$\begin{cases} z_i = d_i \\ x_i = \frac{d_i(u_i - u_0)}{f_x} \\ y_i = \frac{d_i(v_i - v_0)}{f_y} \end{cases} \quad (1)$$

where d_i is the depth. f_x and f_y are the focal length along the x and y coordinate axis, respectively. u_0 and v_0 are the 2D coordinate of the optical center.

We randomly sample N groups of points from the depth map, with three points in each group. The corresponding 3D points are $\mathcal{S} = \{(P_A, P_B, P_C)_i | i = 0 \dots N\}$. We take two restrictions to make three points in a group non-colinear and long-range, i.e., \mathcal{R}_1 and \mathcal{R}_2 . $\angle(\cdot)$ denotes the angle between two vectors.

$$\mathcal{R}_1 = \{\alpha \geq \angle(\overrightarrow{P_A P_B}, \overrightarrow{P_A P_C}) \geq \beta, \alpha \geq \angle(\overrightarrow{P_B P_C}, \overrightarrow{P_B P_A}) \geq \beta | (P_A, P_B, P_C) \in \mathcal{S}\} \quad (2)$$

$$\mathcal{R}_2 = \{\|\overrightarrow{P_k P_m}\| > \theta | k, m \in [A, B, C], (P_A, P_B, P_C) \in \mathcal{S}\} \quad (3)$$

where α, β, θ are hyper-parameters.

Therefore, three 3D points in each group can establish a virtual plane in 3D space. We compute the normal vector of the plane to encode geometric relations, which can be written as

$$\mathcal{N} = \{\mathbf{n}_i = \frac{\overrightarrow{P_{Ai} P_{Bi}} \times \overrightarrow{P_{Ai} P_{Ci}}}{\|\overrightarrow{P_{Ai} P_{Bi}} \times \overrightarrow{P_{Ai} P_{Ci}}\|}, (P_A, P_B, P_C)_i \in \mathcal{S}, i = 0 \dots N\} \quad (4)$$

where \mathbf{n}_i is the normal vector of the virtual plane i .

The normal is an important geometric quantity, which is theoretically scale-and-shift invariant.

Robustness to depth noise. Surface normal is also a widely-used geometric feature. However, our proposed virtual normal is more robust to noise than it. In Fig. 4, we sample three 3D points with large distance. P_A and P_B are assumed to be located on the XY plane, P_C is on the Z axis. When P_C varies to P_C' , the direction of the virtual normal changes from \mathbf{n} to \mathbf{n}' . P_C'' is the intersection point between plane

$P_A P_B P_C'$ and Z axis. Because of restrictions \mathcal{R}_1 and \mathcal{R}_2 , the difference between \mathbf{n} and \mathbf{n}' is usually very small, which is simple to show:

$$\angle(\mathbf{n}, \mathbf{n}') = \angle(\overrightarrow{OP_C}, \overrightarrow{OP_C'}) = \arctan \frac{\|\overrightarrow{P_C P_C'}\|}{\|\overrightarrow{OP_C}\|} \approx 0, \quad (5)$$

$$\|\overrightarrow{P_C P_C'}\| \ll \|\overrightarrow{OP_C}\|$$

In contrast, the surface normal is a typical ‘local’ feature, which is widely used for many point-cloud based applications such as registration [46] and object detection [47], [48]. It appears to be a promising 3D cue for improving depth prediction. One can apply the angular difference between ground-truth and calculated surface normal to be a geometric constraint. One major issue of this approach is, when computing surface normal from either a depth map or 3D point cloud, it is sensitive to noise. Moreover, surface normal only considers short-range local information. We follow [49] to calculate the surface normal. It assumes that local 3D points locate in the same plane, of which the normal vector is the surface normal. In practice, ground-truth depth maps are usually captured by a consumer-level sensor with limited precision, so depth maps are contaminated by noise. The reconstructed point clouds in the local region can vary considerably due to noises as well as the size of local patch for sampling (Fig. 5a). We experiment on the NYUD-V2 dataset to test the robustness of the surface normal computation. Five different sampling sizes around the target pixel are employed to sample points, which are used to calculate its surface normal. The sample area is $a = (2i + 1) \cdot (2i + 1), i = 1, \dots, 5$. The Mean Difference Error (Mean) [26] between calculated surface normals is evaluated. From Fig. 5b, we can learn that the surface normal varies significantly with different sampling sizes. For example, the Mean between 3×3 and 11×11 is 22° . Such unstable surface normal negatively affects its effectiveness for learning. Likewise, other 3D geometric constraints demonstrating the ‘local’ relative relations also encounter this problem.

Furthermore, we conduct a simple experiment to verify the robustness of our proposed virtual normal against data noise. We create a unit sphere and then add Gaussian noise to simulate the ideal noise-free data and the real noisy data (see Fig. 6a). We then sample 100K groups of points

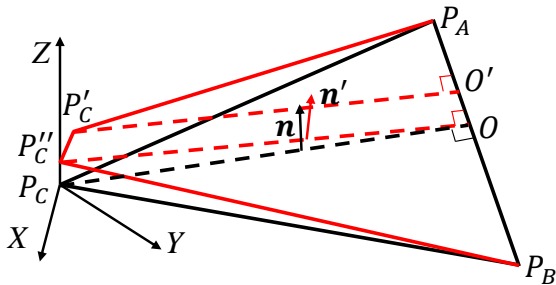


Fig. 4 – Robustness of VN to depth noise. Because of noise, point P_C may vary to P_C' . However, as there is a long distance constraint for virtual normal, the direction of virtual normal will not vary significantly.

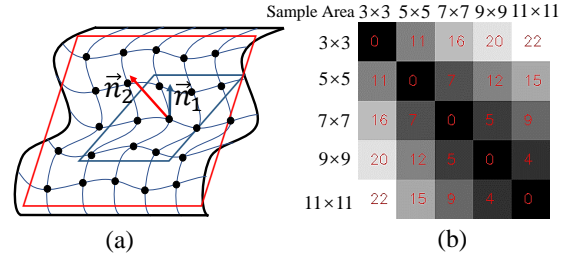


Fig. 5 – Illustration of fitting point clouds to obtain the local surface normal. The directions of the surface normals is fitted with different sampling sizes on a real point cloud (a). Because of noise, the surface normals vary significantly. (b) compares the angular difference between surface normals computed with different sample sizes in Mean Difference Error. The error can vary significantly.

from the noisy surface and the ideal one to compute the virtual normal respectively, while 100K points are sampled to compute the surface normal as well. For the Gaussian noise, we use different deviations to simulate different noise levels by varying deviation $\sigma = [0.0002, \dots, 0.01]$, and the mean being $\mu = 0$. The experimental results are shown in Fig. 6b. We can learn that the local feature, surface normal, is much more sensitive to data noise than our proposed virtual normal. Other local constraints are also sensitive to data noise.

Virtual normal loss. We can sample many triplets and compute corresponding VNs. With the sampled VNs, we compute the divergence as the Virtual Normal Loss (VNL):

$$\ell_{\text{VN}} = \frac{1}{N} \sum_{i=0}^N \|\mathbf{n}_i^{\text{pred}} - \mathbf{n}_i^{\text{gt}}\|_1, \quad (6)$$

where N is the number of valid sampling groups satisfying $\mathcal{R}_1, \mathcal{R}_2$. In experiments, we have employed online hard example mining to remove easy samples. For each batch, we remove 15% samples with the lowest virtual normal loss values.

3.2 Learning Metric Depth

Pixel-wise depth supervision. To predict high-quality metric depth map, we combine a pixel-wise loss and virtual normal loss together to supervise the network output. Following Li *et al.* [10], we quantize the real-valued depth in the log space uniformly and formulate the depth prediction as a classification problem instead of regression, and employ the cross-entropy loss. In particular, we follow [50] to use the weighted cross-entropy loss (WCEL), with the weight being the information gain. See [50] for details. The overall loss is:

$$\ell = \ell_{\text{WCE}} + \lambda \cdot \ell_{\text{VN}}, \quad (7)$$

where λ is a trade-off parameter, which is set to 5 in all experiments to make the two terms roughly of the same scale. ℓ_{WCE} is the weighted cross-entropy loss.

Note that the above overall loss function is differentiable. The gradient of the ℓ_{VN} loss can be easily computed as Eq. (4) and Eq. (6) are both differentiable.

3.3 Learning Affine-invariant Depth

In order to train a model with good generalization, we construct a large-scale and diverse dataset and enforce the

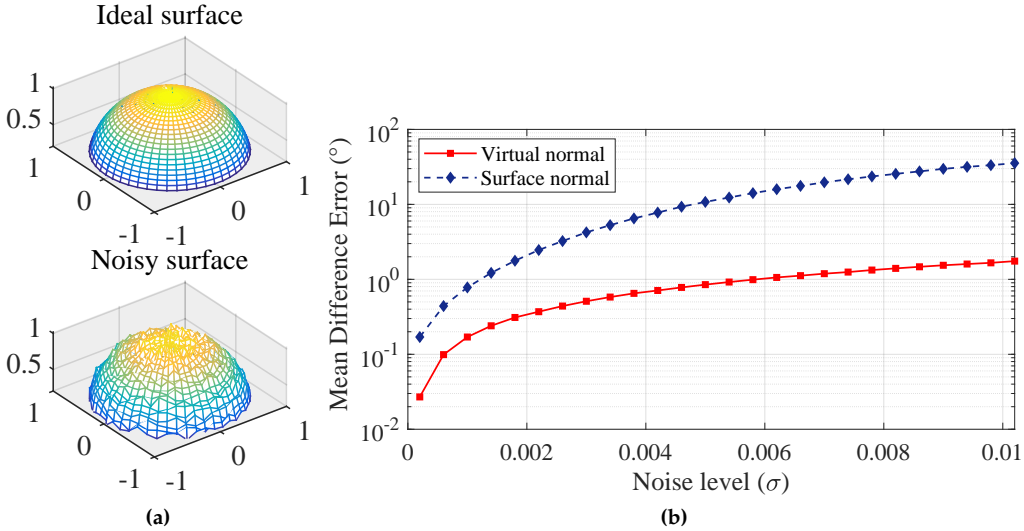


Fig. 6 – Robustness of virtual normal and surface normal against data noise. (a) The ideal surface and noisy surface. (b) The Mean Difference Error (Mean) is applied to evaluate the robustness of virtual normal and surface normal against different noise level. Our proposed virtual normal is more robust.

model to learn affine-invariant depth with virtual normal loss.

Diverse data for training. Table 1 compares the released popular RGB-D datasets. RGB-D sensors can capture high-precision depth data, but it is costly to build a large dataset. By contrast, crawling large-scale online images can boost scene diversity. Previous datasets only have sparse ordinal depth annotations, such as DIW [11] and Youtube3D [14]. Although RedWeb [6] and MegaDepth improve the ground-truth depth quality, RedWeb only has 3600 images and MegaDepth only contains static scenes.

Therefore, to feature diversity, quality, we collect large-size diverse web-stereo data and sample some data from Taskonomy [17] and DIML [18] to construct our large-scale training data, termed DiverseDepth. Most of existing data only contains limited indoor and outdoor scenes. In comparison, we harvest large-scale web-stereo images and videos, which cover diverse foreground objects and people. Following [6], we compute the depth maps from uncalibrated stereo images. This data part is termed *Part-fore*. Besides, we sample some images from Taskonomy [17] and DIML [18] to constitute the indoor and outdoor background part, termed *Part-in* and *Part-out*.

Predicting affine-invariant depth. The geometric model of the monocular depth estimation system is shown in Fig. 7. The ground-truth object in the scene is A^* , and the real camera system is O - XYZ (the black one in Fig. 7). When learning the metric depth, the model $\mathcal{G}(\mathbf{I}, \theta)$ may predict the object at location A . \mathbf{I} is the input image. The learning objective of such methods is to minimize the divergence between A and A^* , i.e., $\min_{\theta} |\mathcal{G}(\mathbf{I}, \theta) - d^*|$, where d^* is the ground-truth depth and θ is the network parameters. As previous learning metric depth methods mainly train and test the model on the same benchmark, where the camera system and the scale remain almost the same, the model can implicitly learn the camera system and produce accurate depth on the testing data [52]. The typical loss functions for learning metric depth are MSE loss and L_1 loss. However,

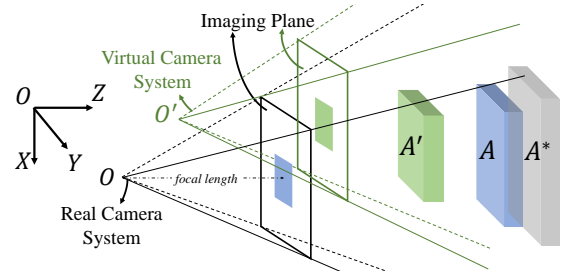


Fig. 7 – The geometric model of an imaging system. A^* is the ground-truth location for an object. A is the predicted location by learning metric depth method, while A' is the predicted location by our learning affine-invariant depth method.

when training and testing on diverse dataset, where the camera system and scale vary, it is theoretically not possible for the model to accommodate multiple camera parameters. The tractable approach is to feed camera parameters of different camera systems to the network as part of the input in order to predict metric depth. This requires the access to camera parameters, which are often not available when harvesting online image data. Our experiments show failure cases of learning metric depth on the diverse dataset (see Table 6, Table 8, and Fig. A1 of the supplementary document). Therefore, learning metric depth method cannot produce a robust model to work on diverse scenes.

Learning the relative depth reduces the difficulty of depth prediction from predicting the accurate metric depth to the ordinal relations. With enough diverse training data, such method can predict relative depth on diverse scenes, but it loses geometric information of the scene, such as the geometric shape. For example, the reconstructed 3D point cloud from the relative depth in Fig. 1 and Fig. 2 cannot represent the shape of the sofa and elephant respectively.

In this paper, we propose to learn the affine-invariant depth from the diverse dataset. On the diverse dataset, we define a virtual camera system, O' - $X'Y'Z'$ (the green one in Fig. 7), which has the same viewpoint as the real one but

has the different optical center location and the focal length. Therefore, there is an affine transformation, *i.e.*, translation T and scaling s , between the real camera system $O\text{-}XYZ$ and the virtual one $O'\text{-}X'Y'Z'$. For the predicted depth under the virtual camera system, it has to take an affine transformation to recover the metric depth under the real camera system, *i.e.*, $P_A = s \cdot (P_{A'} + T)$, where $P = (x, y, d)^T$. The learning objective is defined as follows.

$$L = \min_{\theta} |\mathcal{K}(\mathcal{G}(\mathbf{I}, \theta)) - d^*| \quad (8)$$

where $\mathcal{K}(\cdot)$ is the affine transformation to recover the scaling and translation.

Through explicitly defining a virtual camera system and disentangling the affine transformation between the diverse real camera system and the virtual one, we simplify the objective of monocular depth prediction. The predicted depth will be invariant to various scales and translations. Therefore, it will be easier to generalize to diverse scenes by learning affine-invariant depth than metric depth. Besides, such learning objective can maintain more geometric information than that of learning relative depth.

Pixel-wise depth supervision. Besides, we also apply a pixel-wise loss function to supervise the network. We use the scale-and-shift-invariant loss (SSIL) proposed in [38] to supervise the network. Note that we use this loss to supervise the depth output by our model instead of the inverse depth. The loss function is written as:

$$\begin{cases} \ell_{\text{SSIL}} = \frac{1}{2N} \sum_{i=1}^N \left(\vec{\mathbf{d}}_i^T \mathbf{h} - d_i^* \right)^2 \\ \mathbf{h} = \left(\sum_{i=1}^N \vec{\mathbf{d}}_i \vec{\mathbf{d}}_i^T \right)^{-1} \left(\sum_{i=1}^N \vec{\mathbf{d}}_i d_i^* \right) \\ \vec{\mathbf{d}}_i = (d_i, 1)^T \end{cases} \quad (9)$$

Thus, the overall loss function is as follows:

$$\ell = \ell_{\text{VN}}(d, d^*) + \lambda \cdot \ell_{\text{SSIL}}(d, d^*). \quad (10)$$

Multi-curriculum learning. Training the model on the large-scale and diverse scenes dataset effectively poses challenges. Most existing methods uniformly sample a sequence of mini-batches $\{\mathbb{B}_0, \dots, \mathbb{B}_M\}$ from the whole dataset for training. However, as our DiverseDepth has a wide range of

TABLE 1 – Comparison with previous RGB-D datasets. Our dataset features both diverse scenes and high-quality ground-truth depth.

Dataset	Diversity	Dense	Accuracy	Images
Captured by RGB-D sensor				
NYU [16]	Low	✓	High	407K
KITTI [15]	Low	✓	High	93K
SUN-RGBD [51]	Low	✓	High	10K
ScanNet [3]	Low	✓	High	2.5M
Make3D [24]	Low	✓	High	534
Taskonomy [17]	Low	✓	High	4.5M
DIML [18]	Low	✓	High	2M
DIODE [41]	Low	✓	High	26K
Crawled online				
DIW [11]	High		Low	496K
Youtube3D [14]	High		Low	794K
RedWeb [6]	Medium	✓	Medium	3.6K
WSVD [42]	Medium	✓	low	1.5M
MegaDepth [36]	Medium	✓	Medium	130K
Ours	High	✓	Medium	320K

Algorithm 1: Multi-curriculum learning algorithm

Input : scoring function \mathcal{F} , pacing function \mathcal{H} , dataset \mathbb{X}
Output: mini-batches sequence $\{\mathbb{B}_i | i = 0 \dots M\}$.

- 1 train the model \mathcal{G}_j on the data part \mathbb{D}_j as the teacher
- 2 sort each data part \mathbb{D}_j with ascending difficulty according to \mathcal{F} , the ranked data is \mathbb{C}_j
- 3 **for** $k = 0$ to K **do**
- 4 **for** $i = 0$ to M **do**
- 5 **for** $j = 0$ to P **do**
- 6 subset size $s_{kj} = \mathcal{H}(k, j)$
- 7 subset $\mathbb{S}_{kj} = \mathbb{C}_j[0, \dots, s_{kj}]$
- 8 uniformly sample batch \mathbb{B}_{ij} from \mathbb{S}_{kj}
- 9 **end**
- 10 concatenate P batches sampled from different data parts together $\mathbb{B}_i = \{\mathbb{B}_{ij}\}_{j=0}^P$
- 11 append \mathbb{B}_i to the mini-batches sequence
- 12 **end**
- 13 **end**

scenes, experiments illustrate that such training paradigm cannot effectively optimize the network. We propose a multi-curriculum learning method to solve this problem. We sort the training data by the increasing difficulty and sample a series of mini-batches that exhibit an increasing level of difficulty. Therefore, there are two problems that should be solved: 1) how to construct the curriculum; 2) how to yield a sequence of easy-to-hard mini-batches for the network. Pseudo-code for multi-curriculum algorithm is shown in Algorithm 1.

Three parts of DiverseDepth, *i.e.*, *part-fore*, *part-in* and *part-out*, are termed as $\mathbb{X} = \{\mathbb{D}_j\}_{j=0}^P$. Let $\mathbb{D}_j = \{(x_{ij}, y_{ij}) | i = 0, \dots, N\}$ represents the N data points of the part j , where x_{ij} denotes a single data, y_{ij} is the corresponding label. We train three models, \mathcal{G}_j , separately on 3 parts as teachers. The absolute relative error (Abs-Rel) is chosen as the *scoring function* $\mathcal{F}(\cdot)$ to evaluate the difficulty of each training sample. If $\mathcal{F}(\mathcal{G}_j(x_{ij}, y_{ij})) > \mathcal{F}(\mathcal{G}_j(x_{(i+1)j}, y_{(i+1)j}))$, then we define the data (x_{ij}, y_{ij}) is more difficult to learn. Finally, we sort 3 parts according to the ascending Abs-Rel error and the ranked datasets are $\mathbb{C}_j = \{(x_{ij}, y_{ij}) | i = 0, \dots, N\}$.

The *pacing function* $\mathcal{H}(\cdot)$ determines a sequence of subsets of the dataset so that the likelihood of the easier data would decrease in this sequence, *i.e.*, $\{\mathbb{S}_{0j}, \dots, \mathbb{S}_{Kj}\} \subseteq \mathbb{C}_j$, where \mathbb{S}_{kj} represents the first $\mathcal{H}(k, j)$ elements of \mathbb{C}_j . From each subset \mathbb{S}_{kj} , a sequence of mini-batches $\{\mathbb{B}_{0j}, \dots, \mathbb{B}_{Mj} | j = 0, 1, 2\}$ are uniformly sampled. Here we utilize the stair-case function as the *pacing function*, which is determined by the starting sampling percentage p_j , the current *step* k , and the fixed *step length* I_o (the number of iterations in each step). In each *step* k , there are I_o iterations and the $\mathcal{H}(k, j)$ remains constant, thus the step $k = \lfloor \frac{iter}{I_o} \rfloor$, where *iter* is the iteration index. $\mathcal{H}(k, j)$ is defined as follows:

$$\mathcal{H}(k, j) = \min(p_j \cdot k, 1) \cdot N_j \quad (11)$$

where N_j is the size of part \mathbb{D}_j .

4 IMPLEMENTATION DETAILS

4.1 Network Configures

Learning metric depth. We use ResNeXt-101 as the backbone. When comparing the performance with other state-of-the-art methods, while using ResNeXt-50 for ablation

studies. Following the network used in [2], the weights of ResNeXt in the encoding layers for depth estimation are initialized with models pretrained on the ImageNet dataset. A polynomial decaying strategy with a base learning rate of 0.01 and power of 0.9 is applied for SGD. The weight decay and the momentum are set to 0.0005 and 0.9 respectively. The batch size is 4 in our experiments.

Learning affine-invariant depth. We employ ResNeXt-50 backbone for evaluating the generalization of learning the affine-invariant depth. SGD’s initial learning rate is 0.0005 for all layers. The learning rate is decayed every 5K iterations with the ratio 0.9. The batch size is set to 12. Note that we evenly sample images from three data parts of DiverseDepth to constitute a batch. During the training, images are flipped horizontally, resized with the ratio from 0.5 to 1.5, and cropped with the size of 385×385 . In testing, we resize, pad, and crop the image to keep a similar aspect ratio.

4.2 Datasets

NYUD-V2. The NYUD-V2 dataset consists of 464 different indoor scenes, which are further divided into 249 scenes for training and 215 for testing. We randomly sample 29K images from the training set to form NYUD-Large. Apart from the whole dataset, there are officially annotated 1449 images (NYUD-Small), in which 795 images are split for training and others are for testing. In the ablation study, we use the NYUD-Small data.

KITTI. The KITTI dataset contains over 93K outdoor images and depth maps with the resolution around 1240×374 . All images are captured on driving cars by stereo cameras and a Lidar. We test on images from 29 scenes split by Eigen *et al.* [7].

DIW. The DIW [11] dataset is used for evaluating the generalization of learning affine-invariant depth on diverse data. We test on the testing set, which contains 74441 images.

ETH3D. To evaluate the affine-invariant depth, we sample 105 images from several scenes of ETH3D dataset for testing.

ScanNet. We sample 2692 images with 155 scenes from the ScanNet validation set to evaluate the affine-invariant depth.

DiverseDepth. Apart from sampling data from Taskonomy [17] and DIML [18], we follow [6] to collect large-scale and diverse web-stereo images and videos to construct the ‘Part-fore’ data part. The steps to construct the data is as follows: 1) Crawling online stereoscopic images and videos. We use three websites for data collection: Flickr, 3DStreaming and YouTube. Through comparing the similarity of left/right parts and manually inspection, we remove outliers. 2) Retrieving disparities from stereo materials, then reversing and scaling them to obtain relative depths. Following [6], we utilize the optical flow [53] method to match the paired pixels in stereo samples and take the horizontal matching as the disparity. 3) Filtering depth maps. As many outliers and noises contained in depths, we take 3 metrics to mask out such noises. Firstly, pixels with vertical disparities larger than 5 are removed. Secondly, pixels with the left-right disparity difference greater than 2 are removed.

Furthermore, images with valid pixels less than 30% are discarded. After these filtering processes, we collect more than 90K RGB-D pairs for the *Part-fore* in total. To enrich the diverse background environments, we sample 100K images from Taskonomy [17] and 100K images from DIML [18]. The training data for learning affine-invariant depths has around 300K RGB-D pairs.

4.3 Metrics

To evaluate the performance of learning metric depth on KITTI and NYUD-V2, we follow previous methods [9] take mean absolute relative error (Abs-Rel), mean \log_{10} error (\log_{10}), root mean squared error (RMS), root mean squared log error (RMS (log)), the accuracy under threshold ($\delta_i < 1.25^i, i = 1, 2, 3$), and weighted human disagreement rate (WHDR) [6] metrics.

To evaluate surface normal performance, we follow [13] take the mean error (Mean), median error (Median), and the percent of pixels whose normal degree error is lower than $11.2^\circ, 22.5^\circ$, and 30° .

To evaluate the generalizability of learning affine-invariant depth on diverse scenes, we conduct testing on zero-shot data in training, including ScanNet, DIW, KITTI, NYUD-V2, and ETH3D. Following [38], we explicitly scale and translate the depth to recover the metric depth before evaluating the affine-invariant depth. The scaling and translation factors are obtained by the least-squares method.

5 EXPERIMENTS

In this section, we conduct two groups of experiments. Firstly, we carry out some experiments to demonstrate the effectiveness of virtual normal loss for metric depth estimation, and analyze its property. Furthermore, we conduct several experiments to illustrate the effectiveness of learning affine-invariant depth on the proposed DiverseDepth dataset with the supervision of virtual normal loss.

5.1 Virtual Normal Loss for Learning Metric Depth

5.1.1 Comparison with State-of-the-art Methods

In this section, we compare our methods with state-of-the-art methods. We combine the virtual normal loss with a pixel-wise loss, weighted cross-entropy loss to supervise the network. Note, all the experiments in this section conduct the training on NYUD-V2 or KITTI dataset.

Comparison with state-of-the-art methods on NYU. In this experiment, we compare with a few state-of-the-art methods on the NYUD-V2 dataset. Table 2 demonstrates that our proposed method outperforms other state-of-the-art methods across all evaluation metrics significantly. Compare to DORN, we have improved the accuracy from 0.2% to 18% over all evaluation metrics.

In addition to the quantitative comparison, we demonstrate some visual results between our method and the state-of-the-art DORN in Fig. A4 of the supplementary document.¹ Clearly, the predicted depth by the proposed method is much more accurate. The plane of ours is much smoother and has fewer errors (see the wall regions colored

1. Also available at <https://arxiv.org/abs/2103.04216>

TABLE 2 – Results on NYUD-V2. Our method outperforms other state-of-the-art methods over all evaluation metrics.

Method	Abs-Rel	log10	RMS	δ_1	δ_2	δ_3
	Lower is better			Higher is better		
Saxena <i>et al.</i> [54]	0.349	-	1.214	0.447	0.745	0.897
Karsch <i>et al.</i> [55]	0.349	0.131	1.21	-	-	-
Liu <i>et al.</i> [56]	0.335	0.127	1.06	-	-	-
Ladicky <i>et al.</i> [57]	-	-	-	0.542	0.829	0.941
Li <i>et al.</i> [58]	0.232	0.094	0.821	0.621	0.886	0.968
Roy <i>et al.</i> [59]	0.187	0.078	0.744	-	-	-
Liu <i>et al.</i> [12]	0.213	0.087	0.759	0.650	0.906	0.974
Wang <i>et al.</i> [60]	0.220	0.094	0.745	0.605	0.890	0.970
Eigen <i>et al.</i> [26]	0.158	-	0.641	0.769	0.950	0.988
Chakrabarti [61]	0.149	-	0.620	0.806	0.958	0.987
Li <i>et al.</i> [62]	0.143	0.063	0.635	0.788	0.958	0.991
Laina <i>et al.</i> [9]	0.127	0.055	0.573	0.811	0.953	0.988
DORN [1]	0.115	0.051	0.509	0.828	0.965	0.992
DenseDepth [63]	0.123	0.053	0.465	0.846	0.974	0.994
DSN [64]	0.132	0.056	0.429	0.834	0.959	0.987
Chen [65]	0.111	0.048	0.514	0.878	0.977	0.994
Huynh <i>et al.</i> [66]	0.108	-	0.412	0.882	0.980	0.996
Ours (ResNet101)	0.112	0.051	0.465	0.859	0.970	0.993
Ours (ResNeXt101)	0.108	0.048	0.416	0.875	0.976	0.994

with red in the 1st, 2nd, and 3rd row). Furthermore, the last row in Fig. A4 manifests that our predicted depth is more accurate in the complicated scene. We have fewer errors in shelf and desk regions.

Comparison with state-of-the-art methods on KITTI. In order to demonstrate that our proposed virtual normal can also generalize to outdoor scenes, we test our method on the KITTI dataset and compare with previous state-of-the-art methods. Results in Table 3 show that our method has outperformed all other methods on all evaluation metrics except root-mean-square (RMS) error. The RMS error is only slightly behind that of DORN.

TABLE 3 – Depth prediction results on the KITTI dataset. Our method outperforms other methods over all evaluation metrics except RMS.

Method	δ_1	δ_2	δ_3	Abs-Rel	RMS	RMS (log)
	Higher is better			Lower is better		
Make3D [54]	0.601	0.820	0.926	0.280	8.734	0.361
Eigen <i>et al.</i> [7]	0.692	0.899	0.967	0.190	7.156	0.270
Liu <i>et al.</i> [12]	0.647	0.882	0.961	0.114	4.935	0.206
Semi. [67]	0.862	0.960	0.986	0.113	4.621	0.189
Guo <i>et al.</i> [8]	0.902	0.969	0.986	0.090	3.258	0.168
DORN [1]	0.932	0.984	0.994	0.072	2.727	0.120
DenseDepth [63]	0.886	0.965	0.986	0.093	4.170	0.171
DSN [64]	0.934	0.986	0.996	0.075	3.253	0.119
Ours	0.938	0.990	0.998	0.072	3.258	0.117

5.1.2 Ablation Studies for Virtual Normal Loss

Effectiveness of geometrical loss. In this study, in order to prove the effectiveness of the proposed VNL we compare it with two types of pixel-wise depth map supervision, a pair-wise geometric supervision, and a high-order geometric supervision: 1) the L_1 loss (L_1); 2) the surface normal loss (SNL); 3) the pair-wise geometric loss (PL). We reconstruct the point cloud from the depth map and further recover the surface normal from the point cloud. The angular discrepancy between the ground truth and recovered surface normal is defined as the surface normal loss, which is a high-order geometric supervision in 3D space. The pair-wise loss is the direction difference of two vectors in 3D, which are established by randomly sampling paired points in ground-truth and predicted point cloud. The loss function of PL is

TABLE 4 – The effectiveness of VNL. With virtual normal supervision, the performance can improve significantly on NYUD-V2 dataset.

Metrics	Abs-Rel	log10	RMS	δ_1	δ_2	δ_3
Pixel-wise Depth Supervision						
WCEL	0.1427	0.060	0.511	0.8117	0.9611	0.9895
WCEL+L1	0.1429	0.061	0.626	0.8098	0.9539	0.9858
Pixel-wise Depth Supervision + Geometric Supervision						
WCEL+PL [†]	0.1380	0.059	0.504	0.8212	0.9643	0.9913
WCEL+PL+VNL	0.1341	0.056	0.485	0.8336	0.9671	0.9913
WCEL+SNL [‡]	0.1406	0.059	0.599	0.8209	0.9602	0.9886
WCEL+VNL [‡] (Ours)	0.1337	0.056	0.480	0.8323	0.9669	0.9920

[†] ‘Local’ geometric supervision in 3D.

[‡] ‘Global’ geometric supervision in 3D.

as follows,

$$\ell_{PL} = \frac{1}{N} \sum_{i=0}^N \left[1 - \frac{\overrightarrow{P_{A_i}^* P_{B_i}^*} \cdot \overrightarrow{P_{A_i} P_{B_i}}}{\left\| \overrightarrow{P_{A_i}^* P_{B_i}^*} \right\| \cdot \left\| \overrightarrow{P_{A_i} P_{B_i}} \right\|} \right] \quad (12)$$

where $(P_A^*, P_B^*)_i$ and $(P_A, P_B)_i$ are paired points sampled from the ground truth and the predicted point cloud, respectively. N is the total number of pairs.

We also employ the long-range restriction \mathcal{R}_2 for the paired points. Therefore, similar to VNL, PL can also be seen as a global geometric supervision in 3D space. The experimental results are reported in Table 4. Weighted cross-entropy loss (WCEL) is the baseline for all following experiments.

Firstly, we analyze the effect of pixel-wise depth supervision on prediction performance. Most of previous methods have demonstrated the effectiveness of using the pixel-wise loss to supervise depth prediction. However, when we combine two pixel-wise supervision (WCEL+ L_1) on the depth map, the performance cannot improve anymore. Although they are different mathematically, their constraints are similar. Thus, using two pixel-wise loss terms does not help.

Secondly, we analyze the effectiveness of the supplementary 3D geometric constraint (PL, SNL, VNL). Compared with the baseline (WCEL), the three supplementary 3D geometric constraints can promote the network performance with varying degrees. However, our proposed VNL combining with WCEL has the best performance, which has improved the baseline performance by up to 8%.

Thirdly, we analyze the difference of three geometric constraints. As SNL can only exploit geometric relations of homogeneous local regions, its performance is the lowest among the three constraints over all evaluation metrics. Compared with SNL, since PL constrains the global geometric relations, its performance is clearly better. However, the performance of WCEL+PL is not as good as our proposed WCEL+VNL. When we further add our VNL on top of WCEL+PL, the precision can further be slightly improved and is comparable to WCEL+VNL. Therefore, although PL is a global geometric constraint in 3D, the pair-wise constraint cannot encode as strong geometry information as our proposed VNL.

At last, in order to further demonstrate the effectiveness of VNL, we analyze the results of network trained with and without VNL supervision on the KITTI dataset. The visual comparison is shown in Fig. A2. One can see that VNL can improve the performance of the network in ambiguous

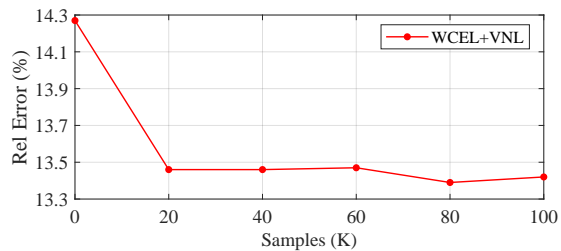


Fig. 8 – Illustration of the impact of the samples size. The more samples will promote the performance. The experiment is conducted on the NYUD-V2 dataset.

regions and distant regions. For example, the distant wall in the first row, the sign (2nd row), the distant pedestrian (3rd row), and the traffic light in the last row of the figure can demonstrate the effectiveness of the proposed VNL.

In conclusion, the geometric constraints in the 3D space can significantly boost the network performance. Moreover, the global and high-order constraints can enforce stronger supervision than the ‘local’ and pair-wise ones in 3D space. **Impact of the amount of samples.** In our proposed virtual normal loss, the amount of samples to construct the virtual normals is a hyper-parameter. Here, the impact of the size of samples for VNL is discussed. We sample six different sizes of point groups, 0K, 20K, 40K, 60K, and 80K and 100K, to establish VNL. ‘0K’ means that the model is trained without VNL supervision. The rel error is reported for evaluation. Fig. 8 demonstrates that ‘rel’ slumps by 5.6% with 20K point groups to establish VNL. However, it only drops slightly when the samples for VNL increase from 20K to 100K. Therefore, the performance saturates with more samples, when samples reach a certain number in that the diversity of samples is enough to construct the global geometric constraint.

Recovered surface normal from the depth. To demonstrate that virtual normal loss can lead the model to learn better shape, we directly recover the surface normal from the predicted depth map. The ground truth is obtained as described in [26]. The quantitative comparison is reported in Table 5. We first compare our geometrically calculated results with CNN-based optimization methods. Although we do not optimize a sub-model to achieve the surface normal, our results can outperform most of such methods and even are the best on 30° metric.

Furthermore, we compare the surface normals directly computed from the reconstructed point cloud with that

TABLE 5 – Evaluation of the surface normal on NYUD-V2. The surface normal can be directly recovered from point cloud. The performance is on par with previous learning-based methods.

Method	Mean	Median	11.2°	22.5°	30°
	Lower is better		Higher is better		
Predicted Surface Normal from the Network					
3DP [68]	33.0	28.3	18.8	40.7	52.4
Ladicky <i>et al.</i> [69]	35.5	25.5	24.0	45.6	55.9
Fouhey <i>et al.</i> [70]	35.2	17.9	40.5	54.1	58.9
Wang <i>et al.</i> [71]	28.8	17.9	35.2	57.1	65.5
Eigen <i>et al.</i> [26]	23.7	15.5	39.2	62.0	71.1
Calculated Surface Normal from the Point cloud					
GT-GeoNet [†] [13]	36.8	32.1	15.0	34.5	46.7
DORN [‡] [1]	36.6	31.1	15.7	36.5	49.4
Ours	24.6	17.9	34.1	60.7	71.7

[†] Cited from the original paper.

[‡] Using authors’ released models.

of DORN [1] and GeoNet [13]. Note that we run the released code and model of DORN to obtain depth maps and then calculate surface normals from the depth, while the evaluation of GeoNet is cited from the original paper. In Table 5, we can see that, with high-order geometric supervision, our method outperforms DORN and GeoNet by a large margin, and even is close to the method of Eigen *et al.* which is trained to output surface normals. It suggests that our method can well learn the shape from images.

Apart from the quantitative comparison, some results are shown in Fig. A5 for visualization, demonstrating that our directly calculated surface normals are not only accurate in planes (the 1st row), but also are of higher quality in regions with sophisticated curved surface (others).

5.2 Virtual Normal Loss for Learning Affine-invariant Depth

5.2.1 Comparison with State-of-the-art Methods

Quantitative comparison on standard benchmarks. The quantitative comparison is reported in Table 6. Apart from Chen *et al.* [11] and Xian *et al.* [6], whose performance is retrieved by re-implementing the ranking loss and training with our model, the performances of other methods are obtained by running their released codes and models. For all methods, we scale and translate the depth before evaluation. Those results whose models have been trained on the testing scene are marked with an underline.

Firstly, from Table 6, we can see that previous state-of-the-art methods, which enforce the model to learn accurate metric depth, cannot generalize to other scenes. For example, the well-trained models of Yin *et al.* [2] and Alhashim and Wonka [63] cannot perform well on other zero-shot scenes.

Secondly, although learning the relative depth methods can predict high-quality ordinal relations on the diverse DIW dataset, *i.e.*, one point being closer or further than another one, the discrepancy between the relative depth and the ground-truth metric depth is very large, see Abs-Rel on other datasets. Such high Abs-Rel results in these methods not being able to recover high-quality 3D shape of scenes, see Fig. 2.

By contrast, through enforcing the model to learn the affine-invariant depth and constructing a high-quality diverse dataset for training, our method can predict high-quality depths on various zero-shot scenes. Our method can outperform previous methods by up to 70%. Noticeably, on NYU, our performance is even on par with existing state-of-the-art methods which have trained on NYU (ours 11.7% *vs.* Alhashim and Wonda’s 12.3%).

Qualitative comparison on zero-shot datasets. Fig. A1 illustrates the qualitative comparison on five zero-shot datasets. The transparent white masks denote the method has trained the model on the corresponding dataset. We can see that the learning metric depth methods, Yin *et al.* [2] and Alhashim and Wonka [63], cannot work well on unseen scenes, while learning relative depth methods, see Ranftl *et al.*, cannot recover high-quality depth map, especially for distant regions (see the marked regions on KITTI, NYU, and ScanNet) and regions with high texture difference (see marked head and colorful wall on DIW). On the DIW dataset, our method can

TABLE 6 – Comparison with state-of-the-art methods on five zero-shot datasets. Our method outperforms previous learning the relative depth or metric depth methods significantly.

Method	Training dataset	Backbone	Testing on zero-shot datasets				
			DIW WHDR	NYU	KITTI	ETH3D Abs-Rel	ScanNet
Learning Metric Depth + Single-scene Dataset							
Yin <i>et al.</i> [2]	NYU	ResNeXt-101	27.0	10.8	35.1	29.6	13.66
Alhashim & Wonka [63]	NYU	DenseNet-169	26.8	12.3	33.4	34.5	12.5
Yin <i>et al.</i> [2]	KITTI	ResNeXt-101	30.8	26.7	7.2	31.8	23.5
Alhashim & Wonka [63]	KITTI	DenseNet-169	30.9	23.5	9.3	32.1	20.5
Learning Relative Depth + Diverse-scene Dataset							
Li & Snavely [36]	MegaDepth	ResNet-50	24.6	19.1	19.3	29.0	18.3
Chen <i>et al.</i> [11]	DIW	ResNeXt-50	11.5	16.7	25.6	25.7	16.0
OASIS [19]	OASIS	ResNet-50	18.9	21.9	31.7	29.2	19.8
MegaDepth [36]	MegaDepth	ResNet-50	31.3	19.4	20.1	26.0	19.0
Xian <i>et al.</i> [6]	RedWeb	ResNeXt-50	21.0	26.6	44.4	39.0	18.2
Learning Affine-invariant Depth + Diverse-scene Dataset							
MiDaS [38]	MIX 5	ResNeXt-101	14.7	11.1	23.6	18.4	11.1
Ours	DiverseDepth-W/o Part-fore	ResNeXt-50	21.5	11.2	13.1	21.0	10.7
Ours	DiverseDepth	ResNeXt-50	14.3	11.7	12.6	22.5	10.4

'_' means that the model was trained on the corresponding dataset. MiDaS results are computed by us with the released V2.0 model.

predict more accurate depth on diverse DIW scenes, such as the forest and sign. Besides, on popular benchmarks, such as ScanNet, KITTI, and NYU, our method can also produce more accurate depth maps.

Furthermore, We test some images captured by a mobile phone. The predicted affine-invariant depth results are shown in Fig. 9. We can see that the depth maps are of high quality.

5.2.2 Ablation Study for Learning Affine-invariant Depth

In this section, we carry out several experiments to analyze the effectiveness of the proposed multi-curriculum learning method, the effectiveness of different loss functions on the

diverse data, the comparison of the reconstructed 3D point cloud among different methods, and the linear relations between the predicted affine-invariant depth and ground truth.

Effectiveness of multi-curriculum learning. To demonstrate the effectiveness of multi-curriculum learning method, we take three settings for the comparison: (1) sampling a sequence of mini-batches uniformly for training, termed Baseline; (2) using the reverse *scoring function*, *i.e.*, $\mathcal{F}' = -\mathcal{F}$, thus the training samples are sorted in the descending order on difficulty and the harder examples are sampled more than easier ones, termed MCL-R; (3) using the proposed multi-curriculum learning method for

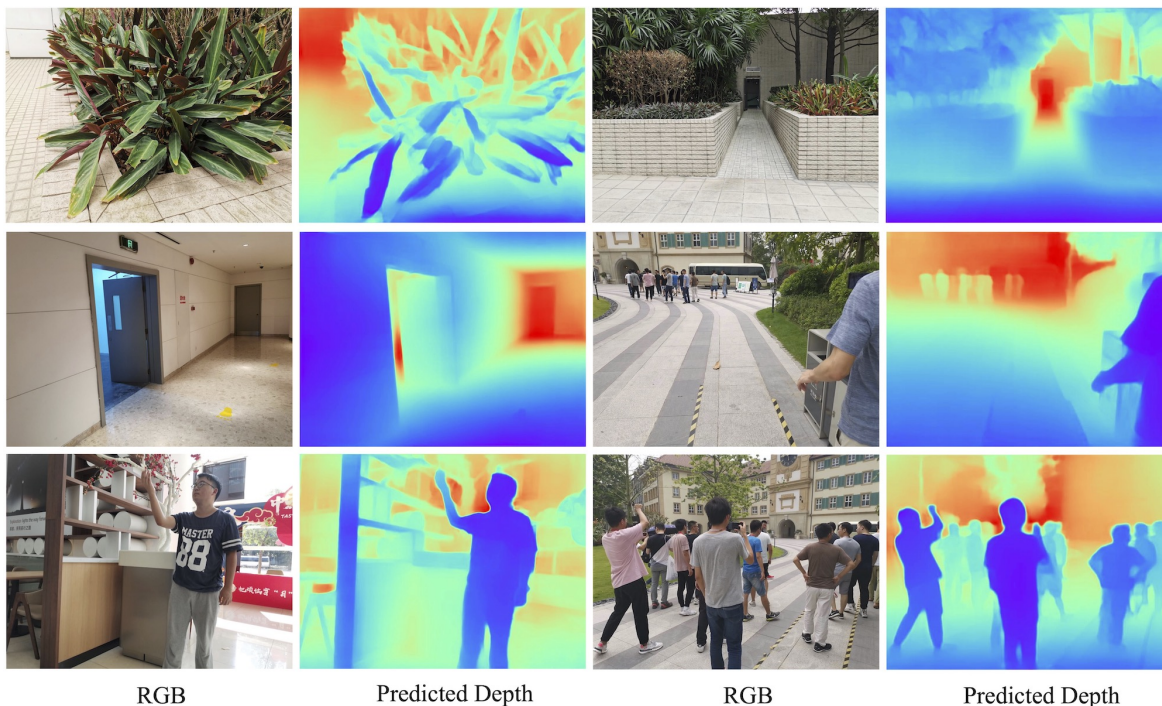


Fig. 9 – Testing on images captured by a phone.

TABLE 7 – Comparison of different training methods on five zero-shot datasets and our DiverseDepth dataset. The proposed multi-curriculum learning method outperforms the baseline noticeably, while MCL-R can also promote the performance.

Method	DIW [†]	NYU [†]	KITTI [†]	ETH3D [†]	ScanNet [†]	DiverseDepth	
	WHDR			Abs-Rel		Abs-Rel	WHDR
Baseline	14.5	11.7	17.9	26.1	11.2	26.0	16.4
MCL-R	15.0	11.8	15.8	24.7	11.0	24.4	15.9
MCL	14.3	11.7	12.6	22.5	10.4	20.6	15.0

[†] Testing on zero-shot datasets.

training, termed MCL. We make comparisons on 5 zero-shot datasets and our proposed DiverseDepth dataset. In Table 7, it is clear that MCL outperforms the baseline by a large margin over all testing datasets. Although MCL-R can also promote the performance, it cannot equal MCL. Furthermore, we demonstrate the validation error along the training in Fig. 10. It is clear that the validation error of MCL is always lower than the baseline and MCL-R over the whole training process. Therefore, the MCL method with an easy-to-hard curriculum can effectively train the model on diverse datasets.

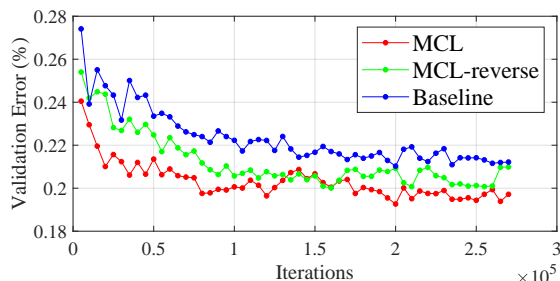


Fig. 10 – Validation error during the training process. The validation error of the proposed multi-curriculum learning method is always lower than that of the MCL-R and baseline.

Impact of different losses. In this section, we analyze the effectiveness of various loss functions for depth estimation on diverse datasets, including virtual normal loss (VNL), scale-shift-invariant loss (SSIL) [38], Silog [7], Ranking, and MSE. We enforce each loss for the network individually in experiments. We sample 10K images from each part of DiverseDepth separately for faster training then test the performances on 5 zero-shot datasets. All the experiments take a multi-curriculum learning method. In Table 8, the VNL and SSIL outperform others over five zero-shot datasets significantly, which demonstrates the effectiveness of learning the affine-invariant depth on diverse datasets. By contrast, as the MSE loss enforces the network to learn the accurate metric depth, it fails to generalize to unseen scenes, thus cannot perform well on zero-shot datasets. Although Ranking can make the model predict good relative depth on diverse DIW, Abs-Rel errors are very high on other datasets because it cannot enrich the model with any geometric information. By contrast, as Silog considers the varying scale in the dataset, it performs slightly better than Ranking and MSE.

Comparison of the recovered 3D shape. In order to further demonstrate learning affine-invariant depth can maintain the geometric information, we reconstruct the 3D point cloud from the predicted depth of a random ScanNet image. We compare our methods with MiDaS *et al.* [38] and Yin-NYU [2]. We take four viewpoints for visual comparison,

TABLE 8 – The effectiveness of different losses for zero-shot evaluation on five datasets. The model is supervised with each loss individually for each experiment. VNL and SSIL outperform others noticeably. By contrast, the model supervised by MSE fails to generalize to diverse scenes, while Ranking can only enforce the model to learn the relative depth. Although Silog considers the varying scale in the dataset, its performance cannot equal VNL and SSIL.

Loss	Testing on zero-shot datasets				
	DIW	NYU	KITTI	ETH3D	ScanNet
VNL+SSIL (Ours)	14.3	11.7	12.6	22.5	10.4
VNL	15.2	12.2	21.0	28.9	11.5
SSIL	17.5	16.5	16.3	26.8	15.6
Silog	19.6	20.8	30.8	29.4	17.6
Ranking	24.3	23.4	47.9	39.5	18.1
MSE	35.3	33.2	36.0	30.2	21.6

i.e., front, up, left, and right viewpoints.

In Fig. A6, it is clear that our reconstructed point cloud can clearly represent the shape of the sofa and the wall from four views, while the sofa shapes of the other two methods are distorted noticeably and the wall is not flat.

Furthermore, we randomly select several images from DIW and reconstruct the 3D point cloud from the predicted affine-invariant depth. We can see that our method can recover high quality 3D shape from a single image. See Fig. 2 and the supplementary document.

Illustration of the affine transformation relation. To illustrate the affine transformation between the predicted affine-invariant depth and the ground-truth metric depth, we randomly select two images from KITTI and NYU respectively, and uniformly sample around 15K points from each image. The predicted depth has been scaled and translated for visualization. In Fig. 11, the red line is the ideal linear relation, while the blue points are the sampled points. We can see the ground-truth depth and the predicted depth have a roughly linear relation. Note that as the precision of the sensor declines with the increase of depth, as expected.

6 CONCLUSION

We have proposed methods to solve the generalization issue of monocular depth estimation, at the same time maintaining as much geometric information as possible. Firstly, we construct a large-scale and highly diverse RGB-D dataset. Compared with previous diverse datasets, which only have sparse depth ordinal annotations, our dataset is annotated with dense and high-quality depth. Besides, we have proposed methods to learn the affine-invariant depth on our DiverseDepth dataset, which can ensure both good generalization and high-quality geometric shape reconstruction from the depth. In order to enable learning affine-invariant depth, we propose the high-order geometric loss, namely, virtual normal loss, which is more robust to noise and enables learning high-quality shapes from a single image. Furthermore, we propose a multi-curriculum learning method to train the model effectively on this diverse dataset. Experiments on NYU and KITTI have demonstrated the effectiveness of virtual normal loss for monocular depth estimation. Besides, experimental results on 8 unseen datasets have shown the usefulness of our dataset for learning affine-invariant depth on diverse scenes.

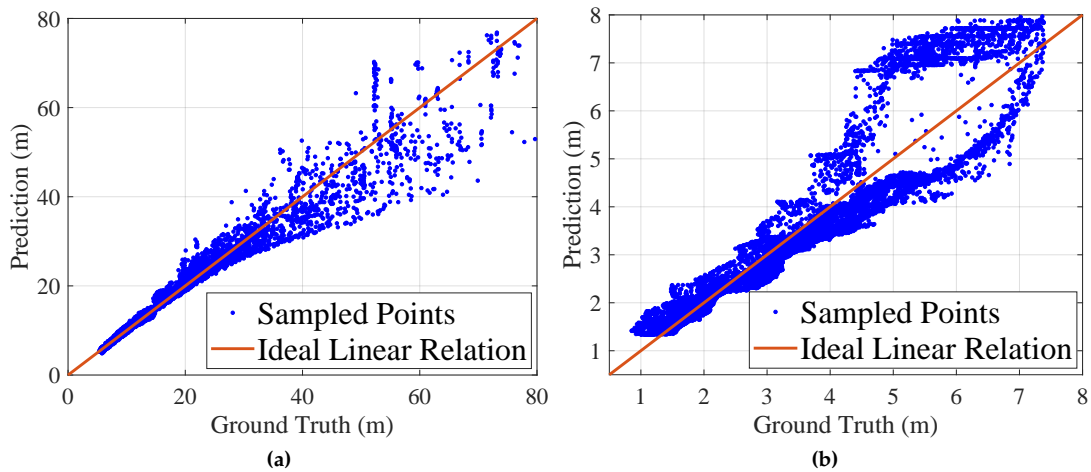


Fig. 11 – Testing the linear relation between the ground-truth and predicted depth. (a) Testing on KITTI. (b) Testing on NYU. Predicted depth has been scaled and translated for visualization. Blue points are the sampled points, while the red line is the ideal linear relation. There is roughly linear relation between the ground-truth and predicted depth.

REFERENCES

- [1] H. Fu, M. Gong, C. Wang, K. Batmanghelich, and D. Tao, “Deep ordinal regression network for monocular depth estimation,” in *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, pp. 2002–2011, 2018.
- [2] W. Yin, Y. Liu, C. Shen, and Y. Yan, “Enforcing geometric constraints of virtual normal for depth prediction,” in *Proc. IEEE Int. Conf. Comp. Vis.*, 2019.
- [3] A. Dai, A. X. Chang, M. Savva, M. Halber, T. Funkhouser, and M. Nießner, “ScanNet: Richly-annotated 3d reconstructions of indoor scenes,” in *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, pp. 5828–5839, 2017.
- [4] Y. Cao, Z. Wu, and C. Shen, “Estimating depth from monocular images as classification using deep fully convolutional residual networks,” *IEEE Trans. Circuits Syst. Video Technol.*, 2017.
- [5] J. Hu, M. Ozay, Y. Zhang, and T. Okatani, “Revisiting single image depth estimation: Toward higher resolution maps with accurate object boundaries,” in *Proc. Winter Conf. Appl. Comp. Vis.*, 2019.
- [6] K. Xian, C. Shen, Z. Cao, H. Lu, Y. Xiao, R. Li, and Z. Luo, “Monocular relative depth perception with web stereo data supervision,” in *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, pp. 311–320, 2018.
- [7] 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.
- [8] X. Guo, H. Li, S. Yi, J. Ren, and X. Wang, “Learning monocular depth by distilling cross-domain stereo networks,” in *Proc. Eur. Conf. Comp. Vis.*, pp. 484–500, 2018.
- [9] I. Laina, C. Rupprecht, V. Belagiannis, F. Tombari, and N. Navab, “Deeper depth prediction with fully convolutional residual networks,” in *Proc. Int. Conf. 3D Vis.*, pp. 239–248, IEEE, 2016.
- [10] R. Li, K. Xian, C. Shen, Z. Cao, H. Lu, and L. Hang, “Deep attention-based classification network for robust depth prediction,” vol. abs/1807.03959, 2019.
- [11] W. Chen, Z. Fu, D. Yang, and J. Deng, “Single-image depth perception in the wild,” in *Proc. Advances in Neural Inf. Process. Syst.*, pp. 730–738, 2016.
- [12] F. Liu, C. Shen, G. Lin, and I. D. Reid, “Learning depth from single monocular images using deep convolutional neural fields,” *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 38, no. 10, pp. 2024–2039, 2016.
- [13] X. Qi, R. Liao, Z. Liu, R. Urtasun, and J. Jia, “Geonet: Geometric neural network for joint depth and surface normal estimation,” in *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, pp. 283–291, 2018.
- [14] W. Chen, S. Qian, and J. Deng, “Learning single-image depth from videos using quality assessment networks,” in *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, pp. 5604–5613, 2019.
- [15] A. Geiger, P. Lenz, C. Stiller, and R. Urtasun, “Vision meets robotics: The KITTI dataset,” *Int. J. Robot. Research*, vol. 32, no. 11, pp. 1231–1237, 2013.
- [16] N. Silberman, D. Hoiem, P. Kohli, and R. Fergus, “Indoor segmentation and support inference from rgb-d images,” in *Proc. Eur. Conf. Comp. Vis.*, pp. 746–760, Springer, 2012.
- [17] A. R. Zamir, A. Sax, W. Shen, L. J. Guibas, J. Malik, and S. Savarese, “Taskonomy: Disentangling task transfer learning,” in *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, pp. 3712–3722, 2018.
- [18] J. Cho, D. Min, Y. Kim, and K. Sohn, “A large RGB-D dataset for semi-supervised monocular depth estimation,” *arXiv: Comp. Res. Repository*, 2019.
- [19] W. Chen, S. Qian, D. Fan, N. Kojima, M. Hamilton, and J. Deng, “OASIS: A large-scale dataset for single image 3d in the wild,” in *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, 2020.
- [20] 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. Comp. Vis.*, pp. 740–756, Springer, 2016.
- [21] J.-W. Bian, Z. Li, N. Wang, H. Zhan, C. Shen, M.-M. Cheng, and I. Reid, “Unsupervised scale-consistent depth and ego-motion learning from monocular video,” in *Proc. Advances in Neural Inf. Process. Syst.*, 2019.
- [22] H. Li, A. Gordon, H. Zhao, V. Casser, and A. Angelova, “Unsupervised monocular depth learning in dynamic scenes,” *arXiv: Comp. Res. Repository*, 2020.
- [23] A. Saxena, S. H. Chung, and A. Y. Ng, “Learning depth from single monocular images,” in *Proc. Advances in Neural Inf. Process. Syst.*, pp. 1161–1168, 2006.
- [24] A. Saxena, M. Sun, and A. Ng, “Make3d: Learning 3d scene structure from a single still image,” *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 31, no. 5, pp. 824–840, 2008.
- [25] B. Liu, S. Gould, and D. Koller, “Single image depth estimation from predicted semantic labels,” in *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, 2010.
- [26] D. Eigen and R. Fergus, “Predicting depth, surface normals and semantic labels with a common multi-scale convolutional architecture,” in *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, pp. 2650–2658, 2015.
- [27] F. Liu, C. Shen, and G. Lin, “Deep convolutional neural fields for depth estimation from a single image,” in *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, pp. 5162–5170, 2015.
- [28] R. Diaz and A. Marathe, “Soft labels for ordinal regression,” in *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, pp. 4738–4747, 2019.
- [29] T. Zhou, M. Brown, N. Snavely, and D. Lowe, “Unsupervised learning of depth and ego-motion from video,” in *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, pp. 1851–1858, 2017.
- [30] C. Godard, O. M. Aodha, M. Firman, and G. Brostow, “Digging into self-supervised monocular depth estimation,” in *Proc. IEEE Int. Conf. Comp. Vis.*, pp. 3828–3838, 2019.
- [31] C. Shu, K. Yu, Z. Duan, and K. Yang, “Feature-metric loss for self-supervised learning of depth and egomotion,” in *Proc. Eur. Conf. Comp. Vis.*, pp. 572–588, 2020.
- [32] 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. Comp. Vis. Patt. Recogn.*, pp. 12240–12249, 2019.
- [33] X. Cheng, P. Wang, and R. Yang, "Learning depth with convolutional spatial propagation network," *IEEE Trans. Pattern Anal. Mach. Intell.*, 2019.
- [34] J. Jiao, Y. Cao, Y. Song, and R. Lau, "Look deeper into depth: Monocular depth estimation with semantic booster and attention-driven loss," in *Proc. Eur. Conf. Comp. Vis.*, pp. 53–69, 2018.
- [35] J.-Y. Zhu, T. Park, P. Isola, and A. Efros, "Unpaired image-to-image translation using cycle-consistent adversarial networks," in *Proc. IEEE Int. Conf. Comp. Vis.*, pp. 2223–2232, 2017.
- [36] Z. Li and N. Snavely, "Megadepth: Learning single-view depth prediction from internet photos," in *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, pp. 2041–2050, 2018.
- [37] W. Yin, J. Zhang, O. Wang, S. Niklaus, L. Mai, S. Chen, and C. Shen, "Learning to recover 3d scene shape from a single image," in *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, 2021.
- [38] R. Ranftl, K. Lasinger, D. Hafner, K. Schindler, and V. Koltun, "Towards robust monocular depth estimation: Mixing datasets for zero-shot cross-dataset transfer," *IEEE Trans. Pattern Anal. Mach. Intell.*, 2020.
- [39] W. Yin, X. Wang, C. Shen, Y. Liu, Z. Tian, S. Xu, C. Sun, and D. Renyin, "Diversedepth: Affine-invariant depth prediction using diverse data," *arXiv: Comp. Res. Repository*, vol. abs/2002.00569, 2020.
- [40] K. Xian, J. Zhang, O. Wang, L. Mai, Z. Lin, and Z. Cao, "Structure-guided ranking loss for single image depth prediction," in *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, pp. 611–620, 2020.
- [41] I. Vasiljevic, N. Kolkin, S. Zhang, R. Luo, H. Wang, F. Z. Dai, A. F. Daniele, M. Mostajabi, S. Basart, M. R. Walter, and G. Shakhnarovich, "DIODE: A Dense Indoor and Outdoor DDepth Dataset," vol. abs/1908.00463, 2019.
- [42] C. Wang, S. Lucey, F. Perazzi, and O. Wang, "Web stereo video supervision for depth prediction from dynamic scenes," in *Proc. Int. Conf. 3D Vis.*, pp. 348–357, 2019.
- [43] D. Weinshall, G. Cohen, and D. Amir, "Curriculum learning by transfer learning: Theory and experiments with deep networks," in *Proc. Int. Conf. Mach. Learn.*, pp. 5235–5243, 2018.
- [44] G. Hacohen and D. Weinshall, "On the power of curriculum learning in training deep networks," in *Proc. Int. Conf. Mach. Learn.*, pp. 2535–2544, 2019.
- [45] Y. Bengio, J. Louradour, R. Collobert, and J. Weston, "Curriculum learning," in *Proc. Int. Conf. Mach. Learn.*, pp. 41–48, ACM, 2009.
- [46] R. B. Rusu, N. Blodow, Z. C. Marton, and M. Beetz, "Aligning point cloud views using persistent feature histograms," in *Proc. IEEE Int. Conf. Intell. Robots & Syst.*, pp. 3384–3391, IEEE, 2008.
- [47] S. Hinterstoisser, S. Holzer, C. Cagniard, S. Ilic, K. Konolige, N. Navab, and V. Lepetit, "Multimodal templates for real-time detection of texture-less objects in heavily cluttered scenes," in *Proc. IEEE Int. Conf. Comp. Vis.*, pp. 858–865, 2011.
- [48] S. Gupta, R. Girshick, P. Arbeláez, and J. Malik, "Learning rich features from rgb-d images for object detection and segmentation," in *Proc. Eur. Conf. Comp. Vis.*, pp. 345–360, 2014.
- [49] K. Klasing, D. Althoff, D. Wollherr, and M. Buss, "Comparison of surface normal estimation methods for range sensing applications," in *Proc. IEEE Conf. Robot. & Automation*, pp. 3206–3211, 2009.
- [50] Y. Cao, Z. Wu, and C. Shen, "Estimating depth from monocular images as classification using deep fully convolutional residual networks," *IEEE Trans. Circuits Syst. Video Technol.*, 2017.
- [51] S. Song, S. P. Lichtenberg, and J. Xiao, "Sun RGB-D: A RGB-D scene understanding benchmark suite," in *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, pp. 567–576, 2015.
- [52] T. v. Dijk and G. d. Croon, "How do neural networks see depth in single images?," in *Proc. IEEE Int. Conf. Comp. Vis.*, pp. 2183–2191, 2019.
- [53] E. Ilg, N. Mayer, T. Saikia, M. Keuper, A. Dosovitskiy, and T. Brox, "FlowNet 2.0: Evolution of optical flow estimation with deep networks," in *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, pp. 2462–2470, 2017.
- [54] A. Saxena, M. Sun, and A. Y. Ng, "Make3d: Learning 3d scene structure from a single still image," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 31, no. 5, pp. 824–840, 2009.
- [55] K. Karsch, C. Liu, and S. B. Kang, "Depth transfer: Depth extraction from video using non-parametric sampling," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 36, no. 11, pp. 2144–2158, 2014.
- [56] M. Liu, M. Salzmann, and X. He, "Discrete-continuous depth estimation from a single image," in *IEEE Trans. Pattern Anal. Mach. Intell.*, pp. 716–723, 2014.
- [57] L. Ladicky, J. Shi, and M. Pollefeys, "Pulling things out of perspective," in *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, pp. 89–96, 2014.
- [58] B. Li, C. Shen, Y. Dai, A. van den Hengel, and M. He, "Depth and surface normal estimation from monocular images using regression on deep features and hierarchical CRFs," in *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, pp. 1119–1127, 2015.
- [59] A. Roy and S. Todorovic, "Monocular depth estimation using neural regression forest," in *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, pp. 5506–5514, 2016.
- [60] P. Wang, X. Shen, Z. Lin, S. Cohen, B. Price, and A. L. Yuille, "Towards unified depth and semantic prediction from a single image," in *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, pp. 2800–2809, 2015.
- [61] A. Chakrabarti, J. Shao, and G. Shakhnarovich, "Depth from a single image by harmonizing overcomplete local network predictions," in *Proc. Advances in Neural Inf. Process. Syst.*, pp. 2658–2666, 2016.
- [62] J. Li, R. Klein, and A. Yao, "A two-streamed network for estimating fine-scaled depth maps from single rgb images," in *Proc. IEEE Int. Conf. Comp. Vis.*, pp. 22–29, 2017.
- [63] I. Alhashim and P. Wonka, "High quality monocular depth estimation via transfer learning," *arXiv: Comp. Res. Repository*, 2018.
- [64] R. de Queiroz Mendes, E. G. Ribeiro, N. dos Santos Rosa, and V. Grassi Jr, "On deep learning techniques to boost monocular depth estimation for autonomous navigation," *Robotics and Autonomous Systems*, vol. 136, p. 103701, 2021.
- [65] X. Chen, X. Chen, and Z.-J. Zha, "Structure-aware residual pyramid network for monocular depth estimation," 2019.
- [66] L. Huynh, P. Nguyen-Ha, J. Matas, E. Rahtu, and J. Heikkilä, "Guiding monocular depth estimation using depth-attention volume," in *Proc. Eur. Conf. Comp. Vis.*, pp. 581–597, Springer, 2020.
- [67] Y. Kuznetsov, J. Stückler, and B. Leibe, "Semi-supervised deep learning for monocular depth map prediction," in *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, pp. 2215–2223, IEEE, 2017.
- [68] D. F. Fouhey, A. Gupta, and M. Hebert, "Data-driven 3d primitives for single image understanding," in *Proc. IEEE Int. Conf. Comp. Vis.*, pp. 3392–3399, 2013.
- [69] L. Ladicky, B. Zeisl, and M. Pollefeys, "Discriminatively trained dense surface normal estimation," in *Proc. Eur. Conf. Comp. Vis.*, pp. 468–484, Springer, 2014.
- [70] D. F. Fouhey, A. Gupta, and M. Hebert, "Unfolding an indoor origami world," in *Proc. Eur. Conf. Comp. Vis.*, pp. 687–702, Springer, 2014.
- [71] X. Wang, D. Fouhey, and A. Gupta, "Designing deep networks for surface normal estimation," in *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, pp. 539–547, 2015.

Authors' photographs and biographies not available at the time of publication.

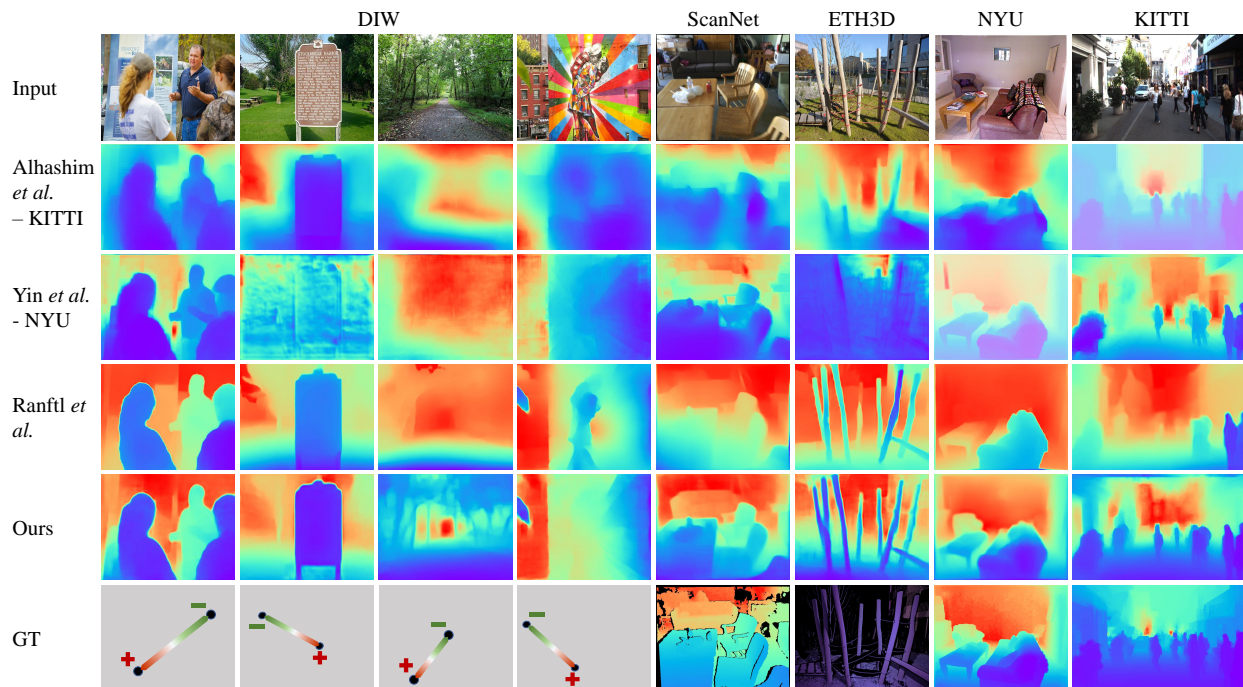


Fig. A1 – Qualitative comparison with state-of-the-art methods on zero-shot datasets. The transparent masks on images denote the method has been trained on the corresponding testing data. The black rectangles highlight the comparison regions. Learning metric depth on NYU and KITTI cannot generalize to diverse scenes, while learning relative depth can generalize to diverse scenes but the details are not good (see black box). Our method not only predicts more accurate depth on diverse DIW, but also recovers better details on indoor and outdoor scenes. Note that ground truth of DIW only annotates the ordinal relation between two points.

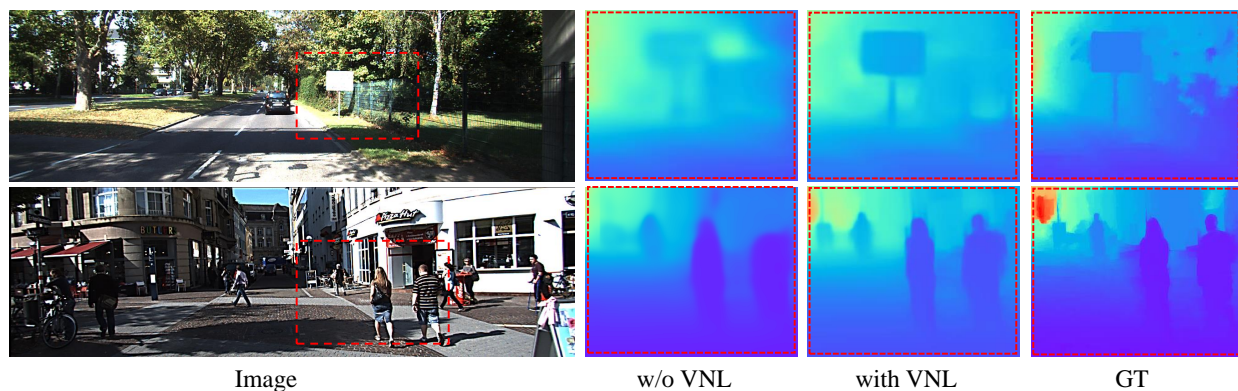


Fig. A2 – Depth maps in the red dashed boxes with sign, pedestrian and traffic lights are zoomed in. We can see that with the help of virtual normal, predicted depth maps in these ambiguous regions are considerably more accurate.

APPENDIX A ADDITIONAL RESULTS

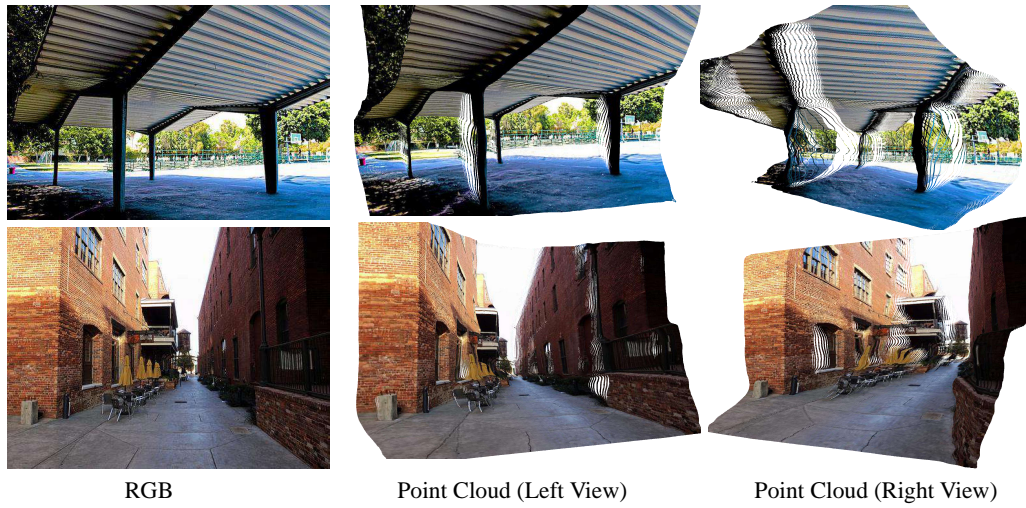


Fig. A3 – Qualitative comparison of the reconstructed 3D point cloud from the predicted affine-invariant depth. The images are randomly selected from the DIW dataset.

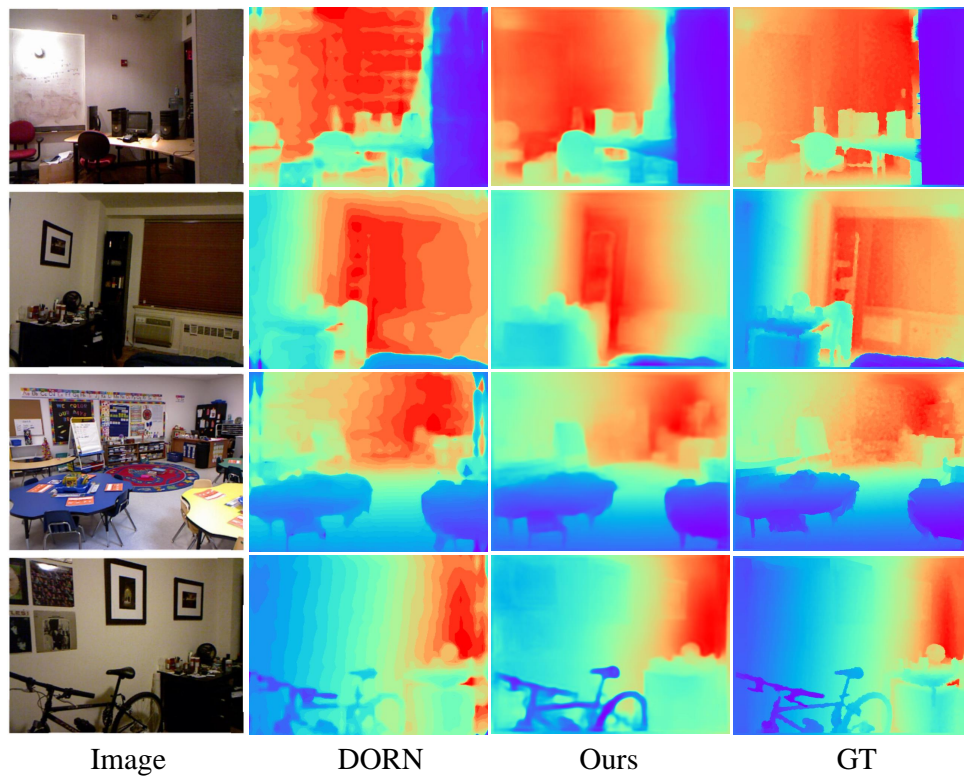


Fig. A4 – Examples of predicted depth maps by our method and the DORN method on NYUD-V2. Color indicates the depth (red is far, purple is close). Our predicted depth maps have fewer errors in planes (e.g., walls) and have high-quality details in complicated scenes (e.g., the desk and shelf in the last row).

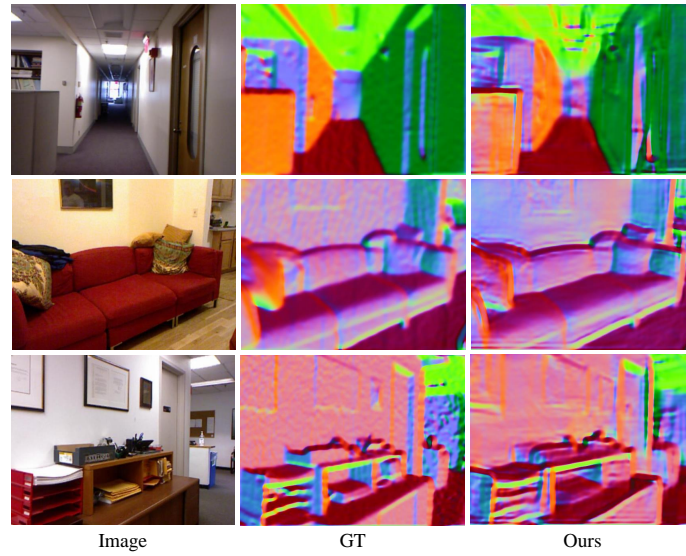


Fig. A5 – Recovered surface normal from the 3D point cloud. According to the visual effect, the surface normal is in high-quality in planes (1st row) and the complicated curved surface (2nd and last row).

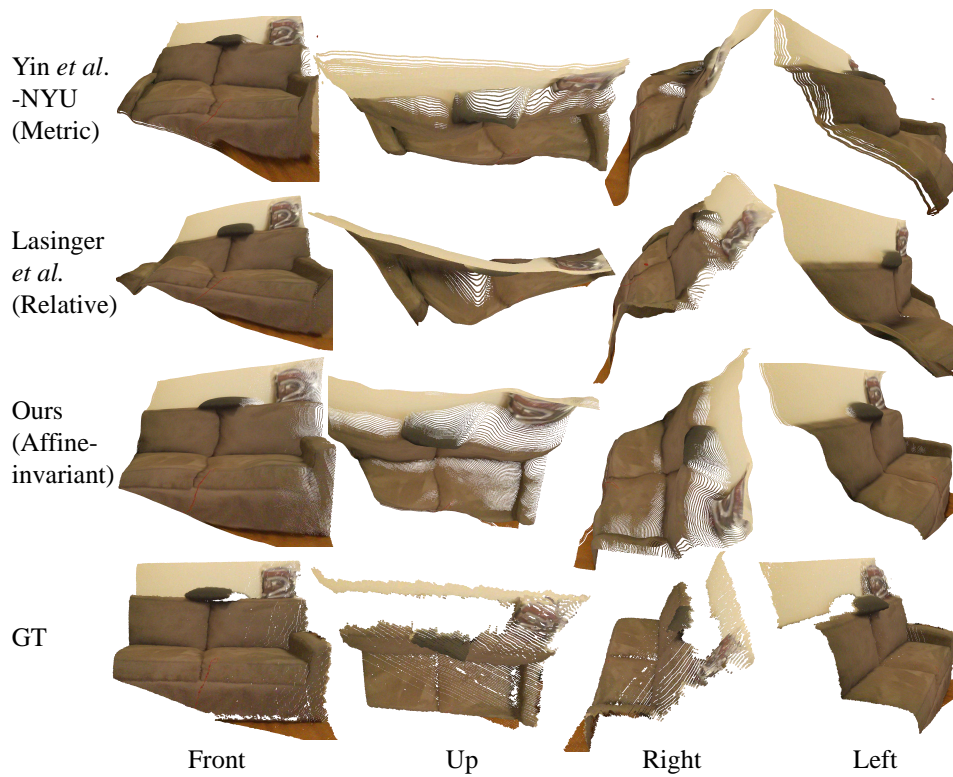


Fig. A6 – Qualitative comparison of the reconstructed 3D point cloud from the predicted depth of a ScanNet image. Our method can clearly recover the shapes of the sofa and wall, while the shape of other methods distort noticeably.

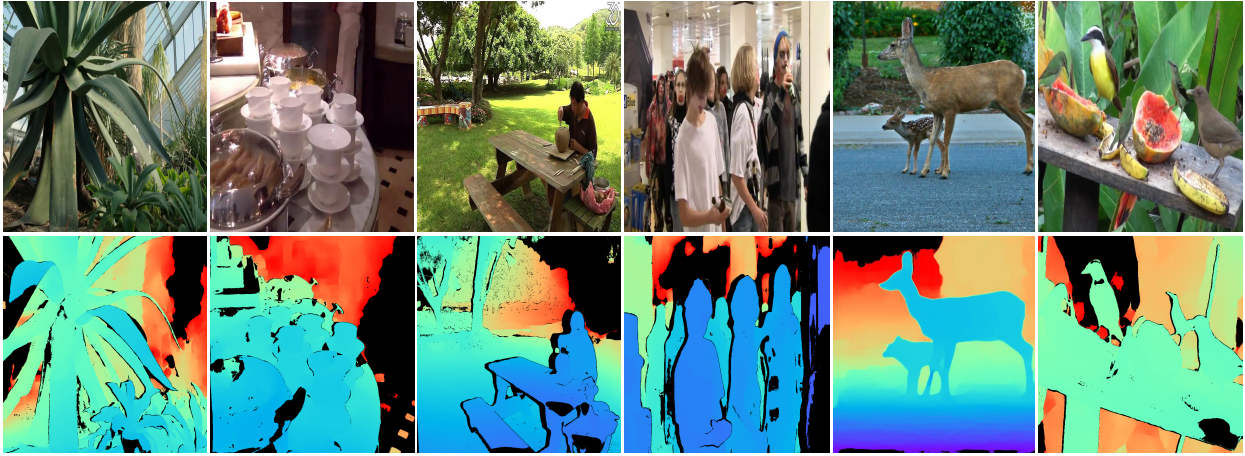


Fig. A7 – Dataset examples. Some examples of our constructed diversified depth dataset.

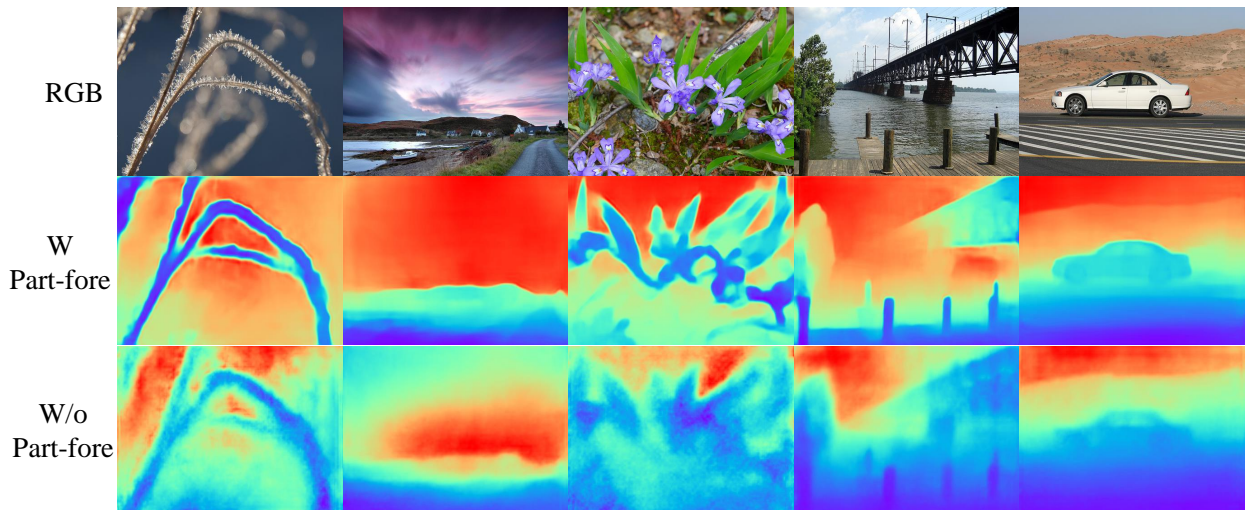


Fig. A8 – Quantitative comparison of depth results on in-the-wild scenes. The model is trained with or without ‘Part-fore’ data. We observe that ‘Part-fore’ data can improve the model’s generalization to diverse scenes.