

Normal Assisted Stereo Depth Estimation

Uday Kusupati^{1*} Shuo Cheng² Rui Chen^{3*} Hao Su²

¹The University of Texas at Austin ²University of California San Diego

³Tsinghua University

uday@cs.utexas.edu, scheng@eng.ucsd.edu, chenr17@mails.tsinghua.edu.cn, haosu@eng.ucsd.edu

Abstract

Accurate stereo depth estimation plays a critical role in various 3D tasks in both indoor and outdoor environments. Recently, learning-based multi-view stereo methods have demonstrated competitive performance with limited number of views. However, in challenging scenarios, especially when building cross-view correspondences is hard, these methods still cannot produce satisfying results. In this paper, we study how to leverage a normal estimation model and the predicted normal maps to improve the depth quality. We couple the learning of a multi-view normal estimation module and a multi-view depth estimation module. In addition, we propose a novel consistency loss to train an independent consistency module that refines the depths from depth/normal pairs. We find that the joint learning can improve both the prediction of normal and depth, and the accuracy & smoothness can be further improved by enforcing the consistency. Experiments on MVS, SUN3D, RGBD and Scenes11 demonstrate the effectiveness of our method and state-of-the-art performance.

1. Introduction

Multi-view stereo (MVS) is one of the most fundamental problems in computer vision and has been studied over decades. Recently, learning-based MVS methods have witnessed significant improvement against their traditional counterparts [45, 23, 47, 7]. In general, these methods formulate the task as an optimization problem, where the target is to minimize the overall summation of pixel-wise depth discrepancy. However, the lack of geometric constraints leads to bumpy depth prediction especially in areas with low texture or that are textureless as shown in Fig. 1. Compared with depth that is a property of the global geometry, surface normal represents a more local geometric property and can be inferred more easily from visual appearance. For instance, it is much easier for

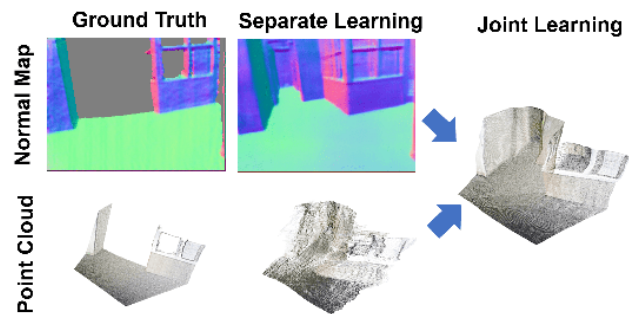


Figure 1. Illustration of results of separate learning and joint learning of depth and normal. While the normal prediction is smooth and accurate, existing state-of-the-art stereo depth prediction result is noisy. Our method improves the prediction quality significantly by joint learning of depth and normal and enforcing consistency. **Color format. Best viewed on screen.**

humans to estimate whether a wall is flat or not than the absolute depth. Fig. 1 shows an example where learning-based MVS methods perform poorly on depth estimation but significantly better on normal prediction.

Attempts have been made to incorporate the normal based geometric constraints into the optimization to improve the monocular depth prediction [49, 59]. One simple form of enforcing a consistency constraint between depth and normal is to enforce orthogonality between the predicted normal and the tangent directions computed from the predicted depths at every point. However, for usage as regularizing loss function during training, we find that a naive consistency method in the world coordinate space is a very soft constraint as there are many sub-optimal depth solutions that are consistent with a given normal. Optimizing depths to be consistent with the normal as a post processing[59] ensures local consistency; however, not only this is an expensive step during inference time, but also the post-processed result may lose grounding from the input images. Therefore, we strive to propose a new formulation of depth-normal consistency that can improve

*Work done while visiting University of California San Diego

the training process. Our consistency is defined in the pixel coordinate space and we show that our formulation is better than the simple consistency along with better performance than previous methods to make the geometry consistent in Section 4.4. This constraint is independent of the multi-view formulation and can be used to enforce consistency on any pair of depth and normal even in the single view setting. To this end, our contributions are mainly in the following aspects:

First, we propose a novel cost-volume-based multi-view surface normal prediction network (NNet). By constructing a 3D cost volume by plane sweeping and accumulating the multi-view image information to different planes through projection, our NNet can learn to infer the normal accurately using the image information at the correct depth. The construction of a cost volume with image features from multiple views contains the information of available features in addition to enforcing additional constraints on the correspondences and thus the depths of each point. We show that the cost volume is a better structural representation that facilitates better learning on the image features for estimating the underlying surface normal. While in single image setting, the network tends to overfit the texture and color and demonstrates worse generalizability, we show that our method of normal estimation generalizes better due to learning on a better abstraction than single view images.

Further, we demonstrate that learning a normal estimation model on the cost volume jointly with the depth estimation pipeline facilitates both tasks. Both traditional and learning-based stereo methods suffer from the noisy nature of the cost volume. The problem is significant in textureless surfaces when the image feature based matching doesn't offer enough cues. We show that enforce the network to predict accurate normal maps from the cost volume results in regularizing the cost volume representation, and thereby assists in producing better depth estimation. Experiments on MVS, SUN3D, RGBD, Scenes11, and Scene Flow datasets demonstrate that our method achieves state-of-the-art performance.

2. Related Work

In this section we review the literature relevant to our work concerned with stereo depth estimation, normal estimation and multi-task learning for multi-view geometry.

Classical Stereo Matching. A stereo algorithm typically consists of the following steps: matching cost calculation, matching cost aggregation and disparity calculation. As the pixel representation plays a critical role in the process, previous literature have exploited a variety of representations, from the simplest RGB colors to hand-craft feature descriptors [50, 44, 39, 31, 2]. Together with post-processing techniques like Markov random fields [38] and

semi-global matching [20], these methods can work well on relative simple scenarios.

Learning-based Stereo. To deal with more complex real world scenes, recently researchers leverage CNNs to extract pixel-wise features and match correspondences [25, 51, 27, 6, 30, 29, 19]. The learned representation shows more robustness to low-texture regions and various lightings [22, 45, 47, 7, 32]. Rather than directly estimating depth from image pairs as done in many previous deep learning methods, some approaches also tried to incorporate semantic cues and context information in the cost aggregation process [43, 8, 24] and achieved positive results. While other geometry information such as normal and boundary [56, 16, 26] are widely utilized in traditional methods for further improving the reconstruction accuracy, it is non-trivial to explicitly enforce the geometry constraints in learning-based approaches [15]. To the best of our knowledge, this is the first work that tries to solve depth and normal estimation in multi-view scenario in a joint learning fashion.

Surface Normal Estimation. Surface normal is an important geometry information for 3D scene understanding. Recently, several data-driven methods have achieved promising results [12, 1, 14, 41, 3, 52]. While these methods learn the image level features and textures to address normal prediction in a single image setting, we propose a multi-view method that generalizes better and reduces the learning complexity of the task.

Joint Learning of Normal and Depth. With deep learning, numerous methods have been proposed for joint learning of normal and depth [34, 18, 54, 11, 48, 57, 12]. Even though these methods achieved some progress, all these methods focus on single image scenario, while there are still few works exploring joint estimation of normal and depth in multi-view setting. Gallup *et al.* [17] estimated candidate plane directions for warping during plane sweep stereo and further enforce an integrability penalty between the predicted depths and the candidate plane for better performance on slanted surfaces, however, the geometry constraints are only applied in a post processing or optimization step (e.g., energy model or graph cut). The lack of end-to-end learning mechanism make their methods easier to stuck in sub-optimal solutions. In this work, our experiments demonstrate that with careful design, the joint learning of normal and depth is favorable for both sides, as the geometry information is easier to be captured. Benefited from the learned powerful representations, our approach achieves competitive results compared with previous methods.

3. Approach

We propose an end-to-end pipeline for multi-view depth and normal estimation as shown in Fig. 2. The entire

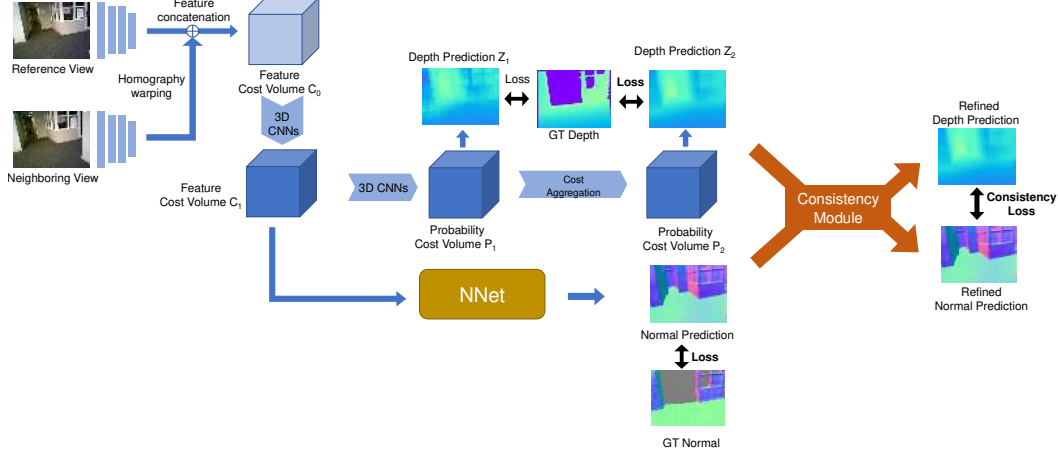


Figure 2. **Illustration of the pipeline of our method.** We first extract deep image features from viewed images and build a feature cost volume by using feature wrapping. The depth and normal are jointly learned in a supervision fashion. Further we use our proposed consistency module to refine the depth and apply a consistency loss.

pipeline can be viewed as two modules. The first module consists of joint estimation of depth and normal maps from the cost volume built from multi-view image features. The subsequent module refines the predicted depth by enforcing consistency between the predicted depth and normal maps using the proposed consistency loss. In the first module, joint prediction of normal from the cost volume implicitly improves the learned model for depth estimation. The second module is explicitly trained to refine the estimates by enforcing consistency.

3.1. Learning based Plane Sweep Stereo

First, we describe our depth prediction module. In terms of the type of target, current learning-based stereo methods can be categorized into : single object reconstruction [46, 7] and scene reconstruction [23, 22]. Compared with single object reconstruction, scene reconstruction, where multiple objects are included, requires larger receptive field for the network to better infer the context information. Because our work also aims at scene reconstruction, we take DPSNet [23], a state-of-the-art scene reconstruction method, as our depth prediction module.

The inputs to the network are a reference image I_1 and a neighboring view image I_2 of the same scene along with the intrinsic camera parameters and the extrinsic transformation between the two views. We first extract deep image features using a spatial pyramid pooling module. Then a cost volume is built by plane sweeping and 3D CNNs are applied on it. Multiple cost volumes can be built and averaged when multiple views are present. Further context-aware cost aggregation [23] is used to regularize the noisy cost volume. The final depth is regressed using soft argmin [28] from the final cost volume.

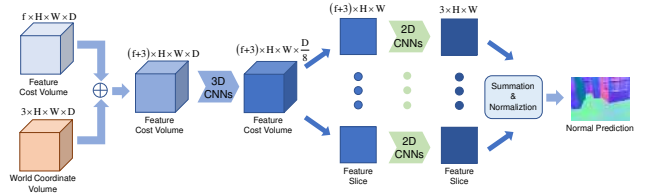


Figure 3. Network architecture of cost volume based surface normal estimation, NNet

3.2. Cost Volume based Surface Normal Estimation

In this section, we describe the network architecture of cost volume based surface normal estimation (Fig. 3). The cost volume contains all the spatial information in the scene as well as image features in it. The probability volume models a depth distribution across candidate planes for each pixel. In the limiting case of infinite candidate planes, an accurately estimated probability volume turns out to be the implicit function representation of the underlying surface *i.e.* takes value 1 where a point on the surface exists and 0 everywhere else. This motivates us to use the cost volume C_1 which also contains the image-level features to estimate the surface normal map \vec{n} of the underlying scene.

Given the cost volume C_1 we concatenate the world coordinates of every voxel to its feature. Then, we use three layers of 2-strided convolution along the depth dimension to reduce the size of this input to $((f+3) \times H \times W \times D/8)$ and call this C_n . Consider a fronto-parallel slice S_i of size $((f+3) \times H \times W)$ in C_n . We pass each slice through a normal-estimation network (NNet). NNet contains 7 layers of 2D convolutions of 3×3 with increasing receptive field as the layers go deep using dilated convolutions (1, 2, 4, 6, 8, 1, 1). We add the output of all slices and normalize the

sum to obtain the estimate of the normal map.

$$\vec{n} = \frac{\sum_{i=1}^{D/8} \mathbf{NNet}(S_i)}{\left\| \sum_{i=1}^{D/8} \mathbf{NNet}(S_i) \right\|_2} \quad (1)$$

We explain the intuition behind this choice as follows. Each slice contains information corresponding to the patch match similarity of each pixel in all the views conditioned on the hallucinated depths in the receptive field of the current slice. In addition, due to the strided 3D convolutions, the slice features accumulate information about features of a group of neighboring planes. The positional information of each pixel in each plane is explicitly encoded into the feature when we concatenated the world coordinates. So $\mathbf{NNet}(S_i)$ is an estimate of the normal at each pixel conditional to the depths in the receptive field of the current slice. For a particular pixel, slices close to the ground truth depth predict good normal estimates, where as slices far from the ground truth predict zero estimates. One way to see this is, if the normal estimate from each slice for a pixel is \vec{n}_i , the magnitude of \vec{n}_i can be seen as the correspondence probability at that slice for that pixel. The direction of \vec{n}_i can be seen as the vector aligning with the strong correspondences in the local patch around the pixel in that slice.¹

We train the first module with ground truth depth (\mathbf{Z}_{gt}) supervision on \mathbf{Z}_1 & \mathbf{Z}_2 along with the ground truth normal (\vec{n}_{gt}) supervision on (\vec{n}). The loss function (\mathcal{L}) is defined as follows.

$$\begin{aligned} \mathcal{L}_z &= |Z_2 - Z_{gt}|_{\mathbf{H}} + \lambda_z |Z_1 - Z_{gt}|_{\mathbf{H}} \\ \mathcal{L}_n &= |\vec{n} - \vec{n}_{gt}|_{\mathbf{H}} \\ \mathcal{L} &= \mathcal{L}_z + \lambda_n \mathcal{L}_n \end{aligned} \quad (2)$$

where $|\cdot|_{\mathbf{H}}$ denotes the Huber norm².

3.3. Depth Normal Consistency

In addition to estimating depth and normal jointly from the cost volume, we use a novel consistency loss to enforce consistency between the estimated depth and normal maps. We utilize the camera model to estimate the spatial gradient of the depth map in the pixel coordinate space using the depth map and normal map. We compute two estimates for $(\frac{\partial Z}{\partial u}, \frac{\partial Z}{\partial v})$ and enforce them to be consistent.

A pinhole model of the camera is adopted as shown in Figure 4).

$$\begin{bmatrix} u \\ v \\ 1 \end{bmatrix} = \frac{1}{Z} \begin{bmatrix} f_x & 0 & u_c \\ 0 & f_y & v_c \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} X \\ Y \\ Z \end{bmatrix} \quad (3)$$

where (u, v) is the corresponding image pixel coordinate

¹Refer to the suppl. material for Visualisation of the NNet slices

²Also referred to as Smooth L1Loss.

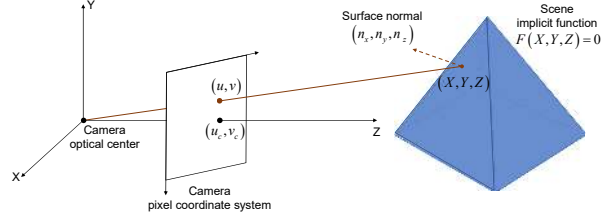


Figure 4. Camera Model

of 3D point (X, Y, Z) , (u_c, v_c) is the pixel coordinate of camera optical center, and f_x, f_y are the focal lengths for X-axis and Y-axis, respectively.

From the camera model, we can yield:

$$\begin{aligned} X &= \frac{Z(u - u_c)}{f_x} \implies \frac{\partial X}{\partial u} = \frac{u - u_c}{f_x} \frac{\partial Z}{\partial u} + \frac{Z}{f_x} \\ Y &= \frac{Z(v - v_c)}{f_y} \implies \frac{\partial Y}{\partial v} = \frac{v - v_c}{f_y} \frac{\partial Z}{\partial v} \end{aligned} \quad (4)$$

Estimate 1:

The spatial gradient of the depth map can first be computed from the depth map by using a Sobel filter:

$$\left(\frac{\partial Z}{\partial u}, \frac{\partial Z}{\partial v} \right)_1 = \left(\frac{\Delta Z}{\Delta u}, \frac{\Delta Z}{\Delta v} \right) \quad (5)$$

Estimate 2:

We assume the underlying scene to be of a smooth surface which can be expressed as an implicit function $F(X, Y, Z) = 0$. The normal map \vec{n} is an estimate of the gradient of this surface.

$$\begin{aligned} \vec{n} &= (n_x, n_y, n_z) = \left(\frac{\partial F}{\partial X}, \frac{\partial F}{\partial Y}, \frac{\partial F}{\partial Z} \right) \\ \implies \frac{\partial Z}{\partial X} &= \frac{-n_x}{n_z}, \frac{\partial Z}{\partial Y} = \frac{-n_y}{n_z} \end{aligned} \quad (6)$$

Therefore, we can derive the second estimate of the depth spatial gradient by:

$$\begin{aligned} \left(\frac{\partial Z}{\partial u} \right)_2 &= \frac{\partial Z}{\partial X} \frac{\partial X}{\partial u} + \frac{\partial Z}{\partial Y} \frac{\partial Y}{\partial u} \\ &= \frac{\left(\frac{-n_x}{n_z} \right)}{1 + \left[\frac{n_x(u - u_c)}{n_z f_x} \right] + \left[\frac{n_y(v - v_c)}{n_z f_y} \right]} \end{aligned} \quad (7)$$

Similarly,

$$\left(\frac{\partial Z}{\partial v}\right)_2 = \frac{\left(\frac{-n_y Z}{n_z f_y}\right)}{1 + \left[\frac{n_x(u-u_c)}{n_z f_x}\right] + \left[\frac{n_y(v-v_c)}{n_z f_y}\right]} \quad (8)$$

The consistency loss \mathcal{L}_c is given as the Huber norm of the deviation between the two estimates $\left(\frac{\partial Z}{\partial u}, \frac{\partial Z}{\partial v}\right)_1$ and $\left(\frac{\partial Z}{\partial u}, \frac{\partial Z}{\partial v}\right)_2$:

$$\mathcal{L}_c = \left| \left(\frac{\partial Z}{\partial u}, \frac{\partial Z}{\partial v}\right)_1 - \left(\frac{\partial Z}{\partial u}, \frac{\partial Z}{\partial v}\right)_2 \right|_{\mathbf{H}} \quad (9)$$

The second estimate of the depth spatial gradient depends only the absolute depth of the pixel in question and not the depths of the neighboring pixels. We obtain the local surface information from the normal map, which we deem more accurate and easier to estimate. Our consistency formulation not only enforces constraints between the relative depths in a pixel’s neighborhood but also the absolute depth. Previous approaches like [34], [58] enforce consistency between depth and normal maps by constraining the local surface tangent obtained from the depth map to be orthogonal to the estimated normal. These approaches typically enforce constraints on the spatial depth gradient in the world coordinate space, where as we enforce them in the pixel coordinate space. We convert the previous approach into a depth gradient consistency formulation and provide a detailed comparison along with experiments on SUN3D dataset in the appendix.

In our pipeline, we implement this loss in an independent module. We use a UNet [35] with the raw depth and normal estimates as inputs to predict the refined depth and normal estimates. We train the entire pipeline in a modular fashion. We initially train the first module with loss \mathcal{L} and then add the second module to train with the consistency loss \mathcal{L}_c in conjunction with the previous losses. Moreover, our loss function can also be used for any depth refinement/completion method including single-view estimation given an estimate of the normal map.

4. Experiments

4.1. Datasets

We use SUN3D [42], RGBD [37] and Scenes11 [4] datasets for training our end-to-end pipeline from scratch. The train set contain 166,285 image pairs from 50420 scenes (SUN3D: 29294, RGBD: 3373, Scenes11: 17753). Scenes11 is a synthetic dataset whereas SUN3D and RGBD consist of real word indoor environments. We test on the same split as previous methods and report the common quantitative measures of depth quality: absolute relative error (Abs Rel), absolute relative inverse error (Abs R-Inv), absolute difference error (Abs diff), square relative error (Sq

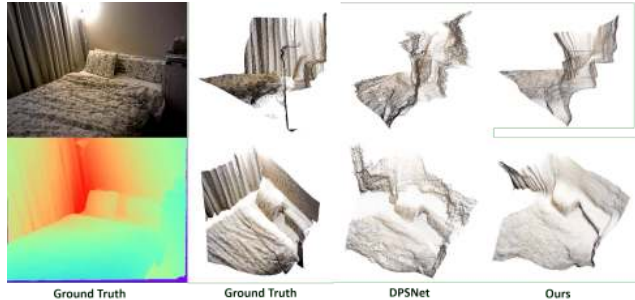


Figure 5. **Visualizing the depths in 3D for SUN3D.** Two views for the point cloud from depth prediction.

Rel), root mean square error and its log scale (RMSE and RMSE log) and inlier ratios ($\delta < 1.25^i$ where $i \in \{1, 2, 3\}$).

We also evaluate our method on a different class of datasets used by other state-of-the-art methods. We train and test on the Scene Flow datasets [33] which consist of 35454 training and 4370 test stereo pairs in 960×540 resolution with both synthetic and natural scenes. The metrics we use on this dataset are the popularly used average End Point Error (EPE) and the 1-pixel threshold error rate.

Further, we evaluate our task on ScanNet [10]. The dataset consists of 94212 image pairs from 1201 scenes. We use the same test split as in [53]. We follow [45] for neighboring view selection, and generate ground-truth normal map following [13]. We use ScanNet to evaluate the performance of the surface normal estimation task too. We use the mean angle error (mean) and median angle error (median) per pixel. In addition, we also report the fraction of pixels with absolute angle difference with ground truth less than t where $t \in \{11.25^\circ, 22.5^\circ, 30^\circ\}$. For all the tables, we represent if a lower value of a metric is better with (\downarrow) and if an upper value of a metric is better with (\uparrow).

For all the experiments, to be consistent with other works, we use only two views to train and evaluate. Please refer the supplementary material for View Selection and ground truth normal generation

4.2. Comparison with state-of-the-art

For comparisons on the DeMoN datasets (SUN3D, RGBD, Scenes11 and MVS), we choose state-of-the-art approaches of a diverse kind. We also evaluate on another dataset MVS [36] containing outdoor scenes of buildings which is not used for training to evaluate generalizability. The complete comparison on all the metrics is presented in Table 1, and some qualitative results are shown in Fig. 5. Our method outperforms existing methods in terms of all the metrics. Also, our method generates more accurate and smooth point cloud with fine details even in textureless regions (eg. bed, wall).

Dataset	Method	Abs Rel(\downarrow)	Abs diff(\downarrow)	Sq Rel(\downarrow)	RMSE (\downarrow)	RMSE log(\downarrow)	$\delta < 1.25$ (\uparrow)	$\delta < 1.25^2$ (\uparrow)	$\delta < 1.25^3$ (\uparrow)
MVS (Outdoor)	COLMAP [36]	0.3841	0.8430	1.257	1.4795	0.5001	0.4819	0.6633	0.8401
	DeMoN [40]	0.3105	1.3291	19.970	2.6065	0.2469	0.6411	0.9017	0.9667
	DeepMVS [21]	0.2305	0.6628	0.6151	1.1488	0.3019	0.6737	0.8867	0.9414
	DPSNet-U [23]	0.0813	0.2006	0.0971	0.4419	0.1595	0.8853	0.9454	0.9735
	Ours	0.0679	0.1677	0.0555	0.3752	0.1419	0.9054	0.9644	0.9879
SUN3D (Indoor)	COLMAP [36]	0.6232	1.3267	3.2359	2.3162	0.6612	0.3266	0.5541	0.7180
	DeMoN [40]	0.2137	2.1477	1.1202	2.4212	0.2060	0.7332	0.9219	0.9626
	DeepMVS [21]	0.2816	0.6040	0.4350	0.9436	0.3633	0.5622	0.7388	0.8951
	DPSNet-U [23]	0.1469	0.3355	0.1165	0.4489	0.1956	0.7812	0.9260	0.9728
	Ours	0.1271	0.2879	0.0852	0.3775	0.1703	0.8295	0.9437	0.9776
RGBD (Indoor)	COLMAP [36]	0.5389	0.9398	1.7608	1.5051	0.7151	0.2749	0.5001	0.7241
	DeMoN [40]	0.1569	1.3525	0.5238	1.7798	0.2018	0.8011	0.9056	0.9621
	DeepMVS [21]	0.2938	0.6207	0.4297	0.8684	0.3506	0.5493	0.8052	0.9217
	DPSNet-U [23]	0.1508	0.5312	0.2514	0.6952	0.2421	0.8041	0.8948	0.9268
	Ours	0.1314	0.4737	0.2126	0.6190	0.2091	0.8565	0.9289	0.9450
Scenes11 (Synthetic)	COLMAP [36]	0.6249	2.2409	3.7148	3.6575	0.8680	0.3897	0.5674	0.6716
	DeMoN [40]	0.5560	1.9877	3.4020	2.6034	0.3909	0.4963	0.7258	0.8263
	DeepMVS [21]	0.2100	0.5967	0.3727	0.8909	0.2699	0.6881	0.8940	0.9687
	DPSNet-U [23]	0.0500	0.1515	0.1108	0.4661	0.1164	0.9614	0.9824	0.9880
	Ours	0.0380	0.1130	0.0666	0.3710	0.0946	0.9754	0.9900	0.9947

Table 1. Comparative evaluation of our model on SUN3D, RGBD, Scenes11 and MVS datasets. For all the metrics except the inlier ratios, lower the better. We use the performance of COLMAP, DeMoN, and DeepMVS reported in [23].

Method	EPE(\downarrow)	1-pixel error rate(\downarrow)
GCNet	1.80	15.6
PSMNet	1.09	12.1
DPSNet	0.80	8.4
GANet-15	0.84	9.9
GANet-deep	0.78	8.7
GANet-NNet	0.77	8.0
Ours	0.69	7.0

Table 2. Comparative evaluation of our model on Scene Flow datasets. For all the metrics, lower the better.

We compare our performance against similar cost-volume based approaches GCNet [28], PSMNet [5] and GANet [55] which have different choices of cost aggregation. Since we use the same testing protocol, we use the performance of GCNet, PSMNet and GANet-15 as reported in [55]. We obtain the performance of GANet-deep which uses a deeper network with more 3D convolutions from the authors’ website. Further, we append our NNet branch to the existing GANet architecture by passing the cost volume of GANet through our NNet and train this branch simultaneously with the full GANet architecture. We call this GANet-NNet. Finally, we also train DPSNet on scene flow datasets to confirm that the better performance is due to normal supervision rather than

better cost aggregation or a better architecture.

Dataset	Method	Abs Rel(\downarrow)	Abs diff(\downarrow)	Sq Rel(\downarrow)	RMSE (\downarrow)
ScanNet	DPSNet	0.1258	0.2145	0.0663	0.3145
	Ours	0.1150	0.2068	0.0577	0.3009
	Ours-\mathcal{L}_c	0.1070	0.1946	0.0508	0.2807
SUN3D	DPSNet	0.1470	0.3234	0.1071	0.4269
	Ours	0.1332	0.3038	0.0910	0.3994
	Ours-\mathcal{L}_c	0.1247	0.2848	0.0791	0.3671

Table 3. Comparative evaluation of our consistency loss.

We also evaluate the performance of our consistency loss on SUN3D and ScanNet datasets. We train DPSNet only on SUN3D as well as ScanNet independently along with our method and present the results and present them in Table 3. We observe that our pipeline achieves significantly better performance on all the metrics on the MVS, SUN3D, RGBD, Scenes11 & SceneFlow datasets. We find that joint multi-view normal and depth estimation helps improve performance on indoor, outdoor, real and synthetic datasets. We further show that our consistency module significantly improves the performance on top of our existing pipeline. We further evaluate the performance on planar and textureless surfaces and visualise the changes in the cost volume due to the addition of NNet.

4.3. Surface Normal Estimation

Table 4 compares our cost volume based surface normal estimation with existing RGB-based, depth-based and RGB-D methods. We perform significantly better than the depth completion based method and perform similar to the RGB based method. The RGB-D based method performs the best, because of using the ground truth depth data.

Method	Mean (↓)	Median (↓)	11.25° (↑)	22.5° (↑)	30° (↑)
RGB-D [53]	14.6	7.5	65.6	81.2	86.2
DC [59]	30.6	20.7	39.2	55.3	60.0
RGB [59]	31.1	17.2	37.7	58.3	67.1
Ours	23.1	18.0	31.1	61.8	73.6

Table 4. Comparison of normal estimation on ScanNet with single view normal estimation. Note that the RGB-D and depth completion (DC) based methods use ground truth depth. The performances of DC & RGB-D are from [53] and RGB from [59].

We evaluate the surface normal estimation performance in the wild by comparing our method against RGB based methods [59]. We use models trained on ScanNet and test them on images from the SUN3D dataset. We present the results in Table 5 and visualize a few cases in Fig. 6. We notice that our method generalizes much better in the wild when compared to the single-view RGB based methods. NNet estimates normals accurately not only in regions of low texture but also in regions with high variance in texture (the bed’s surface). We attribute this performance to using two views instead of one which reduces the learning complexity of the task and thus generalizes better.

We also observe that irrespective of the dataset, the normal estimation loss as well as the validation accuracies saturate within 3 epochs, showing that the task of normal estimation from cost volume is much easier than depth estimation.

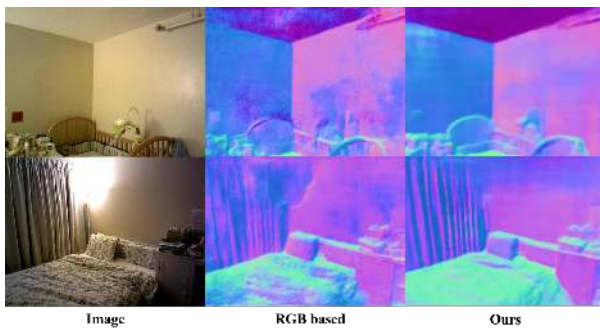


Figure 6. **Surface Normal Estimation.** Test on SUN3D after training on ScanNet. The RGB-based method is from [59]

Method	Mean (↓)	Median (↓)	11.25° (↑)	22.5° (↑)	30° (↑)
RGB - SUN3D	31.6	25.7	17.9	45.6	57.6
Ours - SUN3D	22.9	17.0	34.5	63.2	73.6
RGB - MVS	33.3	27.8	11.8	42.4	55.1
Ours - MVS	27.7	22.4	23.1	52.0	63.9

Table 5. Generalization performance. Both the models were trained on ScanNet (indoor) and tested on SUN3D (indoor) and MVS (outdoor) datasets

Method	Abs Rel(↓)	Abs diff(↓)	Sq Rel(↓)	RMSE (↓)
Raw	0.1332	0.3038	0.0910	0.3994
UNet	0.1307	0.2863	0.0878	0.3720
UNet- \mathcal{L}_t	0.1288	0.2980	0.0842	0.3820
UNet- \mathcal{L}_c	0.1247	0.2848	0.0791	0.3671

Table 6. Ablation Study of Consistency Loss on SUN3D

4.4. Consistency Loss

We perform a few experiments to analyse the performance gains due to our novel consistency loss. We freeze the stereo pipeline and train the UNet that takes the raw estimates of depth and normal maps and refines the depth estimate. We train three configurations of it: (1) Pure network based refinement with just ground truth supervision, (2) Naive consistency loss, \mathcal{L}_t as regularizer³, (3) Our consistency loss \mathcal{L}_c as regularizer. We analyse the performance of the configurations on the SUN3D dataset which consists of indoor environments with a lot of scope for planar, textureless surfaces. Results in Table 6 suggest experimentally that using our consistency loss as regularizer improves the depth estimation accuracy superior to the other approaches.

4.5. Visualizing the Cost Volume

Regularization Existing stereo methods, both traditional and learning-based ones perform explicit cost aggregation on the cost volume. This is because each pixel has good correspondences only with a few pixels in its neighborhood in the other view. But the cost volume contains many more candidates, to be specific, the number of slices per pixel. Further, in textureless regions, there are no distinctive features and thus all candidates have similar features. This induces a lot of noise in the cost volume also is responsible for to false-positives. We show that normal supervision during training regularizes the cost volume. Fig. 7 visualises the probability volume slices and compares it against those of DPSNet. We consider the un-aggregated probability volume P_1 in both the cases. We visualise

³Refer to the supplementary material for elaborate description

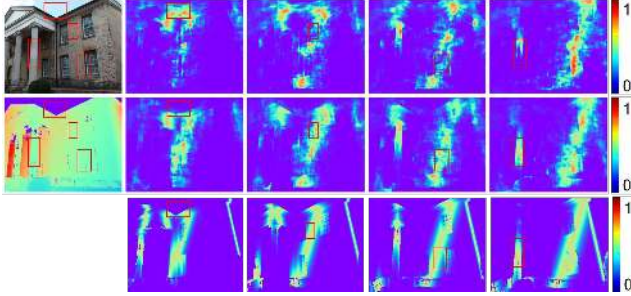


Figure 7. **Cost slice visualization:** The first column contains the reference image and the ground truth depth map. The first row contains the cost volume slices from DPSNet. The second row contains the same from our network. The third row contains the estimates of ground truth cost slices. This can be seen as a $\mathcal{N}(0, 0.01)$ distribution around the ground truth depth corresponding to each slice.

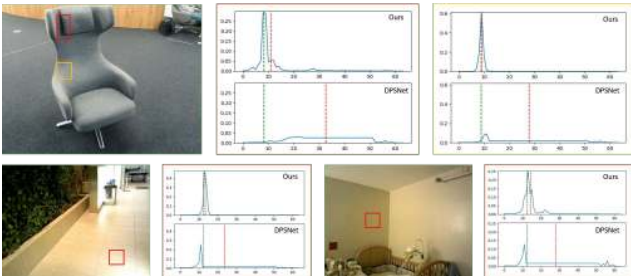


Figure 8. **Post-softmax probability distributions on disparity** Green lines illustrate the ground truth disparities while the red lines illustrate the predicted disparities.

the slices at disparities 14, 15, 16 & 17 (corresponding to depths 2.28, 2.13, 2.0, 1.88) which encompass the wall of the building. The slices of dpsnet are very noisy and do not actually produce good outputs in textureless regions like the walls & sky and reflective regions like the windows.

Planar and Textureless Regions We also visualise the softmax distribution at a few regions in Fig. 8. Challenging regions that are planar or textureless or both are chosen. (a) The chair image consists of very less distinctive textures and the local patches on the chair look the same as those on the floor. But given two views, estimation of normal in the regions with curvature is much easier than estimating depth. This fact allows our model to perform better in the these regions (red & yellow boxes). Cost volume based methods that estimate a probability for each of the candidate depths struggle in textureless regions and usually have long tails in the output distribution. Normal supervision provides additional local constraints on the cost volume and suppresses the tails. This further justifies our understanding (from Section 3.2) that the correspondence probability is related to the slice’s contribution to the normal map.

We quantify these observations by evaluating the

performance of depths D_1 obtained from P_1 against the performance of DPSNet without cost aggregation in Table 7. It shows that normal supervision helps to regularize the cost volume by constraining the cost volume better both qualitatively in challenging cases and quantitatively across the entire test data.

Method	Abs Rel(\downarrow)	Abs diff(\downarrow)	Sq Rel(\downarrow)	RMSE (\downarrow)
DPSNet	0.1274	0.3388	0.1957	0.6230
Ours	0.1114	0.3276	0.1466	0.5631

Table 7. Test performance without cost aggregation on the DeMoN datasets.

Further, we quantify the performance of our methods on planar and textureless surfaces by evaluating on semantic classes on ScanNet test images. Specifically we use the eigen13 classes [9] and report the depth estimation metrics of our methods against DPSNet on the top-2 occurring classes in Table 8. The performance on the remaining classes can be found in the supplementary material. We show that our methods perform well on all semantic categories and quantitatively show the improvement on planar and textureless surfaces which are usually found on walls and floors.

Label	Method	Abs Rel(\downarrow)	Abs diff(\downarrow)	Sq Rel(\downarrow)	RMSE (\downarrow)
Wall	DPSNet	0.1340	0.2968	0.0871	0.3599
	Ours	0.1255	0.2835	0.0799	0.3436
	Ours- \mathcal{L}_c	0.1173	0.2690	0.0721	0.3215
Floor	DPSNet	0.1116	0.2472	0.0777	0.2973
	Ours	0.1092	0.2242	0.0509	0.2642
	Ours- \mathcal{L}_c	0.1037	0.2061	0.0474	0.2561

Table 8. Semantic class specific evaluation on ScanNet

5. Conclusion

In this paper, we proposed to leverage multi-view normal estimation and apply geometry constraints between surface normal and depth at training time to improve stereo depth estimation. We jointly learn to predict the depth and the normal based on the multi-view cost volume. Moreover, we proposed to refine the depths from depth, normal pairs with an independent consistency module which is trained independently using a novel consistency loss. Experimental results showed that joint learning can improve both the prediction of normal and depth, and the accuracy & smoothness can be further improved by enforcing the consistency. We achieved state-of-the-art performance on MVS, SUN3D, RGBD, Scenes11 and Scene Flow datasets.

References

- [1] Aayush Bansal, Bryan Russell, and Abhinav Gupta. Marr revisited: 2d-3d alignment via surface normal prediction. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 5965–5974, 2016. 2
- [2] Herbert Bay, Tinne Tuytelaars, and Luc Van Gool. Surf: Speeded up robust features. In *European conference on computer vision*, pages 404–417. Springer, 2006. 2
- [3] Alexandre Boulch and Renaud Marlet. Deep learning for robust normal estimation in unstructured point clouds. In *Computer Graphics Forum*, volume 35, pages 281–290. Wiley Online Library, 2016. 2
- [4] Angel X. Chang, Thomas A. Funkhouser, Leonidas J. Guibas, Pat Hanrahan, Qi-Xing Huang, Zimo Li, Silvio Savarese, Manolis Savva, Shuran Song, Hao Su, Jianxiong Xiao, Li Yi, and Fisher Yu. Shapenet: An information-rich 3d model repository. *CoRR*, abs/1512.03012, 2015. 5
- [5] Jia-Ren Chang and Yong-Sheng Chen. Pyramid stereo matching network. *CoRR*, abs/1803.08669, 2018. 6
- [6] Jia-Ren Chang and Yong-Sheng Chen. Pyramid stereo matching network. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2018. 2
- [7] Rui Chen, Songfang Han, Jing Xu, and Hao Su. Point-based multi-view stereo network. In *The IEEE International Conference on Computer Vision (ICCV)*, 2019. 1, 2, 3
- [8] Ian Cherabier, Johannes L Schonberger, Martin R Oswald, Marc Pollefeys, and Andreas Geiger. Learning priors for semantic 3d reconstruction. In *Proceedings of the European conference on computer vision (ECCV)*, pages 314–330, 2018. 2
- [9] Camille Couprie, Clément Farabet, Laurent Najman, and Yann LeCun. Indoor semantic segmentation using depth information. In *1st International Conference on Learning Representations, ICLR 2013, Scottsdale, Arizona, USA, May 2-4, 2013, Conference Track Proceedings*, 2013. 8
- [10] Angela Dai, Angel X. Chang, Manolis Savva, Maciej Halber, Thomas A. Funkhouser, and Matthias Nießner. Scannet: Richly-annotated 3d reconstructions of indoor scenes. In *2017 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2017, Honolulu, HI, USA, July 21-26, 2017*, pages 2432–2443, 2017. 5
- [11] Thanuja Dharmasiri, Andrew Spek, and Tom Drummond. Joint prediction of depths, normals and surface curvature from rgb images using cnns. In *2017 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 1505–1512. IEEE, 2017. 2
- [12] David Eigen and Rob Fergus. Predicting depth, surface normals and semantic labels with a common multi-scale convolutional architecture. In *The IEEE International Conference on Computer Vision (ICCV)*, December 2015. 2
- [13] David F Fouhey, Abhinav Gupta, and Martial Hebert. Data-driven 3d primitives for single image understanding. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 3392–3399, 2013. 5
- [14] David Ford Fouhey, Abhinav Gupta, and Martial Hebert. Unfolding an indoor origami world. In *European Conference on Computer Vision*, pages 687–702. Springer, 2014. 2
- [15] Yasutaka Furukawa, Carlos Hernández, et al. Multi-view stereo: A tutorial. *Foundations and Trends® in Computer Graphics and Vision*, 9(1-2):1–148, 2015. 2
- [16] Silvano Galliani, Katrin Lasinger, and Konrad Schindler. Gipuma: Massively parallel multi-view stereo reconstruction. *Publikationen der Deutschen Gesellschaft für Photogrammetrie, Fernerkundung und Geoinformation e. V.*, 25:361–369, 2016. 2
- [17] David Gallup, Jan-Michael Frahm, Philippos Mordohai, Qingxiong Yang, and Marc Pollefeys. Real-time plane-sweeping stereo with multiple sweeping directions. In *2007 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR 2007), 18-23 June 2007, Minneapolis, Minnesota, USA, 2007*. 2
- [18] Christian Hane, Lubor Ladicky, and Marc Pollefeys. Direction matters: Depth estimation with a surface normal classifier. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 381–389, 2015. 2
- [19] Wilfried Hartmann, Silvano Galliani, Michal Havlena, Luc Van Gool, and Konrad Schindler. Learned multi-patch similarity. *2017 IEEE International Conference on Computer Vision (ICCV)*, pages 1595–1603, 2017. 2
- [20] Heiko Hirschmüller. Stereo processing by semiglobal matching and mutual information. *IEEE Transactions on pattern analysis and machine intelligence*, 30(2):328–341, 2007. 2
- [21] Po-Han Huang, Kevin Matzen, Johannes Kopf, Narendra Ahuja, and Jia-Bin Huang. Deepmvs: Learning multi-view stereopsis. In *2018 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2018, Salt Lake City, UT, USA, June 18-22, 2018*, pages 2821–2830, 2018. 6
- [22] Po-Han Huang, Kevin Matzen, Johannes Kopf, Narendra Ahuja, and Jia-Bin Huang. Deepmvs: Learning multi-view stereopsis. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018. 2, 3
- [23] Sunghoon Im, Hae-Gon Jeon, Stephen Lin, and In So Kweon. Dpsnet: End-to-end deep plane sweep stereo. In *7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019*, 2019. 1, 3, 6
- [24] Sunghoon Im, Hae-Gon Jeon, Stephen Lin, and In So Kweon. Dpsnet: end-to-end deep plane sweep stereo. *arXiv preprint arXiv:1905.00538*, 2019. 2
- [25] Mengqi Ji, Juergen Gall, Haitian Zheng, Yebin Liu, and Lu Fang. Surfacer: An end-to-end 3d neural network for multiview stereopsis. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2307–2315, 2017. 2
- [26] Michael Kazhdan, Matthew Bolitho, and Hugues Hoppe. Poisson surface reconstruction. In *Proceedings of the fourth Eurographics symposium on Geometry processing*, volume 7, 2006. 2
- [27] Alex Kendall, Hayk Martirosyan, Saumitro Dasgupta, Peter Henry, Ryan Kennedy, Abraham Bachrach, and Adam Bry. End-to-end learning of geometry and context for deep stereo regression. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 66–75, 2017. 2

- [28] Alex Kendall, Hayk Martirosyan, Saumitro Dasgupta, Peter Henry, Ryan Kennedy, Abraham Bachrach, and Adam Bry. End-to-end learning of geometry and context for deep stereo regression. *CoRR*, abs/1703.04309, 2017. 3, 6
- [29] Sameh Khamis, Sean Fanello, Christoph Rhemann, Adarsh Kowdle, Julien Valentin, and Shahram Izadi. Stereonet: Guided hierarchical refinement for real-time edge-aware depth prediction. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 573–590, 2018. 2
- [30] Zhengfa Liang, Yiliu Feng, Yulan Guo, Hengzhu Liu, Wei Chen, Linbo Qiao, Li Zhou, and Jianfeng Zhang. Learning for disparity estimation through feature constancy. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2018. 2
- [31] David G Lowe et al. Object recognition from local scale-invariant features. 2
- [32] Keyang Luo, Tao Guan, Lili Ju, Haipeng Huang, and Yawei Luo. P-mvsnet: Learning patch-wise matching confidence aggregation for multi-view stereo. In *The IEEE International Conference on Computer Vision (ICCV)*, October 2019. 2
- [33] N. Mayer, E. Ilg, P. Häusser, 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 *IEEE International Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016. arXiv:1512.02134. 5
- [34] Xiaojuan Qi, Renjie Liao, Zhengzhe Liu, Raquel Urtasun, and Jiaya Jia. Geonet: Geometric neural network for joint depth and surface normal estimation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 283–291, 2018. 2, 5
- [35] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. *CoRR*, abs/1505.04597, 2015. 5
- [36] Johannes L. Schönberger and Jan-Michael Frahm. Structure-from-motion revisited. In *2016 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2016, Las Vegas, NV, USA, June 27-30, 2016*, pages 4104–4113, 2016. 5, 6
- [37] Jürgen Sturm, Nikolas Engelhard, Felix Endres, Wolfram Burgard, and Daniel Cremers. A benchmark for the evaluation of RGB-D SLAM systems. In *2012 IEEE/RSJ International Conference on Intelligent Robots and Systems, IROS 2012, Vilamoura, Algarve, Portugal, October 7-12, 2012*, pages 573–580, 2012. 5
- [38] Richard Szeliski, Ramin Zabih, Daniel Scharstein, Olga Veksler, Vladimir Kolmogorov, Aseem Agarwala, Marshall Tappen, and Carsten Rother. A comparative study of energy minimization methods for markov random fields. In *European conference on computer vision*, pages 16–29. Springer, 2006. 2
- [39] Engin Tola, Vincent Lepetit, and Pascal Fua. Daisy: An efficient dense descriptor applied to wide-baseline stereo. *IEEE transactions on pattern analysis and machine intelligence*, 32(5):815–830, 2009. 2
- [40] Benjamin Ummenhofer, Huizhong Zhou, Jonas Uhrig, Nikolaus Mayer, Eddy Ilg, Alexey Dosovitskiy, and Thomas Brox. Demon: Depth and motion network for learning monocular stereo. In *2017 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2017, Honolulu, HI, USA, July 21-26, 2017*, pages 5622–5631, 2017. 6
- [41] Xiaolong Wang, David Fouhey, and Abhinav Gupta. Designing deep networks for surface normal estimation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 539–547, 2015. 2
- [42] Jianxiong Xiao, Andrew Owens, and Antonio Torralba. SUN3D: A database of big spaces reconstructed using sfm and object labels. In *IEEE International Conference on Computer Vision, ICCV 2013, Sydney, Australia, December 1-8, 2013*, pages 1625–1632, 2013. 5
- [43] Guorun Yang, Hengshuang Zhao, Jianping Shi, Zhidong Deng, and Jiaya Jia. Segstereo: Exploiting semantic information for disparity estimation. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 636–651, 2018. 2
- [44] Qingxiong Yang, Liang Wang, Ruigang Yang, Henrik Stewénius, and David Nistér. Stereo matching with color-weighted correlation, hierarchical belief propagation, and occlusion handling. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 31(3):492–504, 2008. 2
- [45] Yao Yao, Zixin Luo, Shiwei Li, Tian Fang, and Long Quan. Mvsnet: Depth inference for unstructured multi-view stereo. *European Conference on Computer Vision (ECCV)*, 2018. 1, 2, 5
- [46] Yao Yao, Zixin Luo, Shiwei Li, Tian Fang, and Long Quan. Mvsnet: Depth inference for unstructured multi-view stereo. In *Computer Vision - ECCV 2018 - 15th European Conference, Munich, Germany, September 8-14, 2018, Proceedings, Part VIII*, pages 785–801, 2018. 3
- [47] Yao Yao, Zixin Luo, Shiwei Li, Tianwei Shen, Tian Fang, and Long Quan. Recurrent mvsnet for high-resolution multi-view stereo depth inference. *Computer Vision and Pattern Recognition (CVPR)*, 2019. 1, 2
- [48] Wei Yin, Yifan Liu, Chunhua Shen, and Youliang Yan. Enforcing geometric constraints of virtual normal for depth prediction. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 5684–5693, 2019. 2
- [49] Zhichao Yin and Jianping Shi. Geonet: Unsupervised learning of dense depth, optical flow and camera pose. In *CVPR*, 2018. 1
- [50] Jae-Chern Yoo and Tae Hee Han. Fast normalized cross-correlation. *Circuits, systems and signal processing*, 28(6):819, 2009. 2
- [51] Jure Zbontar, Yann LeCun, et al. Stereo matching by training a convolutional neural network to compare image patches. *Journal of Machine Learning Research*, 17(1-32):2, 2016. 2
- [52] Jin Zeng, Yanfeng Tong, Yunmu Huang, Qiong Yan, Wenxiu Sun, Jing Chen, and Yongtian Wang. Deep surface normal estimation with hierarchical rgb-d fusion. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 6153–6162, 2019. 2
- [53] Jin Zeng, Yanfeng Tong, Yunmu Huang, Qiong Yan, Wenxiu Sun, Jing Chen, and Yongtian Wang. Deep surface normal estimation with hierarchical RGB-D fusion. In *IEEE Conference on Computer Vision and Pattern Recognition*,

- CVPR 2019, Long Beach, CA, USA, June 16-20, 2019*, pages 6153–6162, 2019. [5](#), [7](#)
- [54] Huangying Zhan, Chamara Saroj Weerasekera, Ravi Garg, and Ian Reid. Self-supervised learning for single view depth and surface normal estimation. *arXiv preprint arXiv:1903.00112*, 2019. [2](#)
- [55] Feihu Zhang, Victor Adrian Prisacariu, Ruigang Yang, and Philip H. S. Torr. Ga-net: Guided aggregation net for end-to-end stereo matching. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2019, Long Beach, CA, USA, June 16-20, 2019*, pages 185–194, 2019. [6](#)
- [56] Shuangli Zhang, Weijian Xie, Guofeng Zhang, Hujun Bao, and Michael Kaess. Robust stereo matching with surface normal prediction. In *2017 IEEE International Conference on Robotics and Automation (ICRA)*, pages 2540–2547. IEEE, 2017. [2](#)
- [57] Yinda Zhang and Thomas Funkhouser. Deep depth completion of a single rgb-d image. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 175–185, 2018. [2](#)
- [58] Yinda Zhang and Thomas A. Funkhouser. Deep depth completion of a single RGB-D image. In *2018 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2018, Salt Lake City, UT, USA, June 18-22, 2018*, pages 175–185, 2018. [5](#)
- [59] Yinda Zhang, Shuran Song, Ersin Yumer, Manolis Savva, Joon-Young Lee, Hailin Jin, and Thomas A. Funkhouser. Physically-based rendering for indoor scene understanding using convolutional neural networks. In *2017 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2017, Honolulu, HI, USA, July 21-26, 2017*, pages 5057–5065, 2017. [1](#), [7](#)