

# Deep Depth from Focal Stack with Defocus Model for Camera-Setting Invariance

Yuki Fujimura<sup>1</sup>Masaaki Iiyama<sup>2</sup>Takuya Funatomi<sup>1</sup>Yasuhiro Mukaigawa<sup>1</sup><sup>1</sup>Nara Institute of Science and Technology, Japan<sup>2</sup>Shiga University, Japan

fujimura.yuki@is.naist.jp iiyama@iiyama-lab.org

{funatomi,mukaigawa}@is.naist.jp

## Abstract

We propose a learning-based depth from focus/defocus (DFF), which takes a focal stack as input for estimating scene depth. Defocus blur is a useful cue for depth estimation. However, the size of the blur depends on not only scene depth but also camera settings such as focus distance, focal length, and f-number. Current learning-based methods without any defocus models cannot estimate a correct depth map if camera settings are different at training and test times. Our method takes a plane sweep volume as input for the constraint between scene depth, defocus images, and camera settings, and this intermediate representation enables depth estimation with different camera settings at training and test times. This camera-setting invariance can enhance the applicability of learning-based DFF methods. The experimental results also indicate that our method is robust against a synthetic-to-real domain gap, and exhibits state-of-the-art performance.

## 1. Introduction

In computer vision, depth estimation from two-dimensional (2D) images is an important task and used for many applications such as VR, AR, or autonomous driving. Defocus blur is a useful cue for such depth estimation because the size of the blur depends on scene depth. Depth from focus/defocus (DFF) takes defocus images as input for depth estimation. Typical inputs for DFF are stacked images, *i.e.*, *focal stack*, each of which is captured with a different focus distance.

DFF methods are roughly divided into two categories, model-based and learning-based. Model-based methods use a thin-lens model for modeling defocus blurs [13, 29, 30] or define focus measures [21, 28] to estimate scene depth. One of the drawbacks of such methods is difficulty in estimating scene depth with texture-less surfaces. Learning-based methods have been proposed to tackle the above drawback [9, 17, 33]. For example, Hazirbas *et al.* [9] proposed a con-

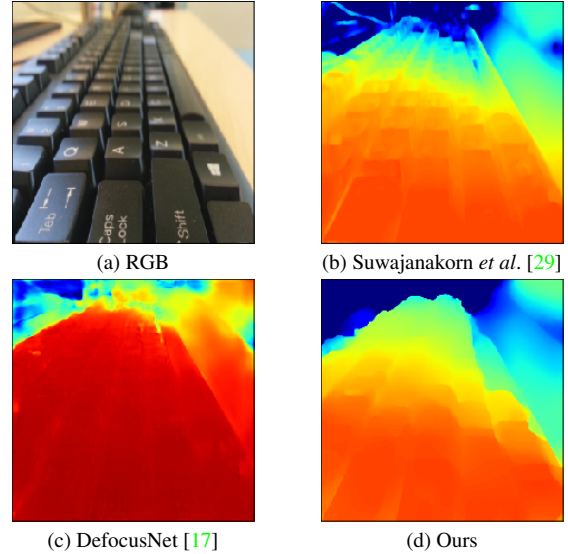


Figure 1. (a) One of input images in focal stack, (b) output depth of [29], (c) output depth of DefocusNet [17], and (d) our result. Our model and DefocusNet were trained on a dataset with camera settings that differed from those of the test data. Our method with camera-setting invariance can estimate correct depth map.

volutional neural network (CNN) taking a focal stack as input without any explicit defocus models. This is an end-to-end method that allows efficient depth estimation. It also enables the depth estimation of texture-less surfaces with learned semantic cues.

General learning-based methods often have limited generalization due to a domain gap between training and test data. Learning-based DFF methods suffer from the difference of capture settings of a camera at training and test times. The amount of a defocus blur depends on not only scene depth but also camera settings such as focus distance, focal length, and f-number. Different depths and camera settings can generate defocus images with the same appearance; thus this difference cannot be compensated with often

used domain adaptation method such as neural style transfer [15, 37]. If camera settings are different at training and test times, the estimated depth has some ambiguity, which is similar to the scale-ambiguity in monocular depth estimation [10]. Current learning-based DFF methods [9, 17, 33] do not take into account the latent defocus model, thus the estimated depth is not correct if the camera settings at test time differ from those at training time, as shown in Fig. 1(c). On the other hand, this problem does not matter for model-based methods with explicit defocus models under given camera settings.

We propose learning-based DFF with a lens defocus model. Our method is inspired by recent learning-based multi-view stereo (MVS) [32], where a cost volume is constructed on the basis of a plane sweep volume [4]. The proposed method also constructs a cost volume, which is passed through a CNN to estimate scene depth. Each defocus image in a focal stack is deblurred at each sampled depth in the plane sweep volume, then the consistency is evaluated between deblurred images. We found that scene depth is effectively learned from the cost volume in DFF. Our method has several advantages over the other learning-based methods directly taking a focal stack as input without an explicit defocus model [9, 17, 33]. First, output depth satisfies the defocus model because the cost volume imposes an explicit constraint among the scene depth, defocus images, and camera settings. Second, the camera settings, such as focus distances and f-number are absorbed into the cost volume as intermediate representation. This enables depth estimation with different camera settings at training and test times, as shown in Fig. 1(d).

The primary contributions of this paper are summarized as follows:

- To the best of our knowledge, this is the first study to combine a learning framework and model-based DFF through a plane sweep volume.
- Our method with camera-setting invariance can be applied to datasets with different camera settings at training and test times, which improves the applicability of learning-based DFF methods.
- Similar to the previous learning-based method [17], our method is also robust against a synthetic-to-real domain gap and achieves state-of-the-art performance.

## 2. Related work

**Depth from focus/defocus** Depth from focus/defocus (DFF) estimates scene depth from focus or defocus cues in captured images and is a major task in computer vision. In general, depth from focus takes many images captured with different focus distances and determines scene depth from an image with the best focus. On the other hand, depth from

defocus aims to estimate scene depth from a small number of images, which do not necessarily need to include focused images [35]. Recently, depth estimation from a focal stack implicitly uses both focus and defocus cues; thus, we use unified terminology, *depth from focus/defocus*.

Traditional DFF methods propose focus measures to evaluate the amount of a defocus blur [19, 21, 28, 38]. If we have a focal stack as input, we can simply refer to the image with noticeable edges and its focus distance. Other methods formulate the amount of defocus blur with a lens defocus model and solve an optimization problem to obtain a depth map together with an all-in-focus image [13, 29]. We refer to these methods as model-based methods. One of the drawbacks of such methods is difficulty in estimating scene depth with texture-less surfaces. Learning-based methods have been proposed to tackle these issues [9, 17, 33]. These methods enable depth estimation at texture-less surfaces and the depth estimation is achieved efficiently in an end-to-end manner. Other learning-based methods leveraged defocus cues as additional information [2, 3] or supervision [8, 26] for monocular depth estimation.

However, current learning-based DFF methods, which directly take a focal stack as input, do not take into account the latent defocus model [9, 17, 33]. For example, Hazirbas *et al.* [9] proposed a CNN that directly takes a focal stack as input. Maximov *et al.* [17] and Wang *et al.* [33] simply used focus distances as intermediate inputs of neural networks. These methods require the same camera settings at training and test times to obtain a correct depth map due to the lack of explicit defocus models. This characteristic reduces the applicability of learning-based DFF methods. On the other hand, our method is a combination of model-based and learning-based methods through a cost volume, which is computed with a lens defocus model, allowing depth estimation with camera-setting invariance.

**Learning from cost volume** Learning from a cost volume is efficient in many applications. A cost volume is constructed by sampling solution space and evaluating costs at each sampled point. Examples of learning-based methods with a cost volume are optical flow [11, 27] and disparity [12, 18] estimation. Learning-based MVS methods [5, 16, 32, 36] are also major examples, where a cost volume is constructed on the basis of a plane sweep volume [4]. Our method also constructs a plane sweep volume and evaluates consistency between defocus images in an input focal stack. We found that learning from a cost volume is also efficient for learning-based DFF.

## 3. Deep depth from focal stack

Our method combines a learning framework and model-based DFF through a cost volume for depth estimation with camera-setting invariance. We first give an overview of the

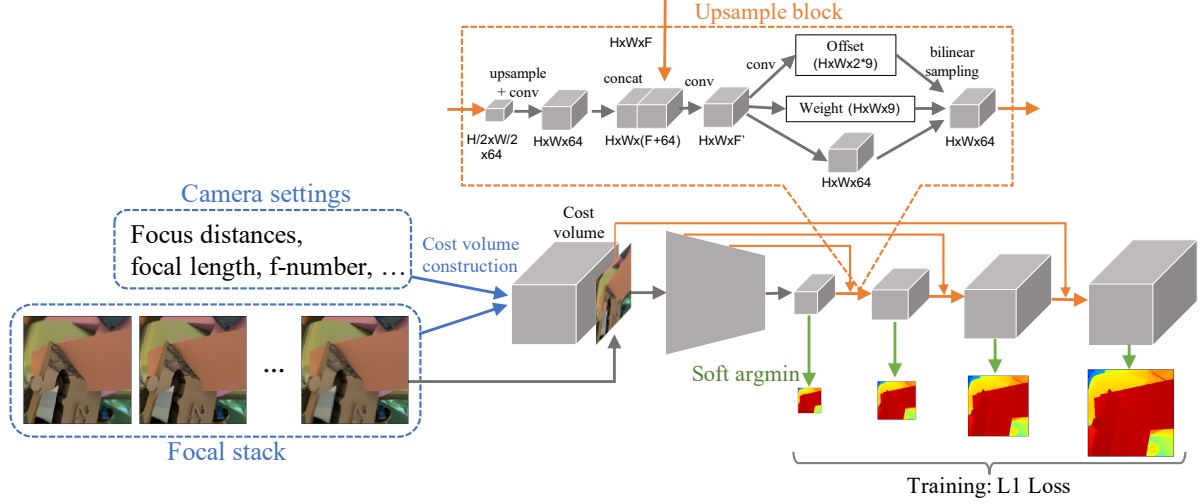


Figure 2. Overview of our method. Our method takes focal stack and camera settings as input then constructs cost volume as intermediate representation, which absorbs differences in camera settings. CNN takes this cost volume together with additional image as input then estimates refined cost volume in coarse-to-fine manner. Depth maps are computed by applying soft argmin operator at each resolution. Each upsample block has adaptive cost aggregation module.

proposed method then describe the lens defocus model and ambiguity of estimated depth in DFF, followed by details of cost volume construction. This cost volume as intermediate representation enables depth estimation with different camera settings at training and test times. The network architecture and loss function are also discussed at the end of this section.

### 3.1. Overview

Figure 2 shows an overview of the proposed method. Our method is inspired by recent learning-based MVS [32], where a cost volume is constructed on the basis of a plane sweep volume [4]. Our cost volume is constructed from an input focal stack by evaluating deblurred images at each depth hypothesis. This intermediate representation absorbs the difference in camera settings. The computed cost volume and an additional defocus image are passed through a CNN with an encoder-decoder architecture. At the decoder part, the cost volume is gradually upsampled for coarse-to-fine estimation. Output depth maps are obtained by applying a differentiable soft argmin operator [12] to intermediate refined cost volumes. Each upsample block includes a cost aggregation module for learning local structures adaptively.

### 3.2. Lens defocus model

Our cost volume construction is based on a lens defocus model, with which the size of a defocus blur is formulated as a circle of confusion (CoC) [38], as shown in Fig. 3. Let  $d$  and  $d_f$  be the scene depth and focus distance of a camera, respectively. CoC can be computed as

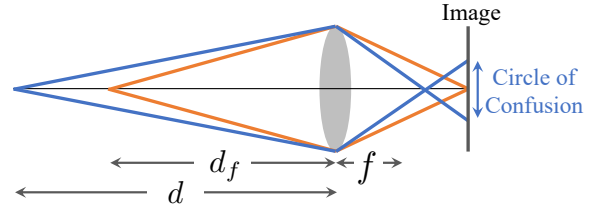


Figure 3. Circle of confusion (CoC) that corresponds to size of defocus blur

$$c = b \frac{\|d - d_f\|}{d} \frac{f^2}{N(d_f - f)}, \quad (1)$$

where  $f$  is the focal length of the lens and  $N$  is the f-number.  $b$  [px/m] converts the unit of the CoC from [m] to [px]. When  $d$  is equal to  $d_f$ , the light rays from the scene point converge on the image plane; otherwise, defocus blur results as the size of the diameter of the CoC. The blurred image can be computed as a convolution of an all-in-focus image with the point spread function (PSF), the kernel size of which corresponds to the size of the CoC.

The CoC can be computed from the scene depth  $d$  and the camera settings  $f$ ,  $d_f$ ,  $N$ , and  $b$  in Eq. (1). Note that these parameters can easily be extracted from EXIF properties [17] or calibrated beforehand [30], and the state-of-the-art methods assume these parameters are known [17, 33]; thus, this paper also follows the same assumption. Our method realizes depth estimation with camera-setting invariance about these parameters, and this improves the

applicability of learning-based DFF methods because our method with camera-setting invariance can be applied to datasets with different camera settings at training and test times.

Now, we discuss two ambiguities in DFF due to the camera settings. The first one is scale-ambiguity. From Eq. (1), the following relationship holds:

$$\begin{aligned} c &= b \frac{\|d - d_f\|}{d} \frac{f^2}{N(d_f - f)} \\ &= b^{-*} \frac{\|d^* - d_f^*\|}{d^*} \frac{f^{*2}}{N(d_f^* - f^*)}, \end{aligned} \quad (2)$$

where  $(\cdot)^* = (\cdot)\sigma$ ,  $(\cdot)^{-*} = (\cdot)/\sigma$ ,  $\forall \sigma \in \mathbb{R}$ . This means scaled camera settings and depth give the same CoC as that of the original ones.

The other ambiguity is affine-ambiguity. From Eq. (1), we can obtain

$$\begin{aligned} c &= b \frac{f^2}{N(d_f - f)} \left\| 1 - \frac{d_f}{d} \right\| \\ &= A(f, d_f, N) + \frac{B(f, d_f, N)}{d}, \end{aligned} \quad (3)$$

where  $A(f, d_f, N)$  and  $B(f, d_f, N)$  are constants. Thus, different camera settings and inverse depths can give the same CoC as follows:

$$\begin{aligned} c &= A(f, d_f, N) + \frac{B(f, d_f, N)}{d} \\ &= A(f', d'_f, N') + \frac{B(f', d'_f, N')}{d'}. \end{aligned} \quad (4)$$

This means the estimated inverse depth has affine-ambiguity (Similar discussion can be found in the previous study [7]). In the experiments, we evaluate the proposed method with respect to the scale-ambiguity in the depth space and the affine-ambiguity in the inverse depth space.

### 3.3. Cost volume

The proposed method computes a cost volume from the focal stack for the input of a CNN to impose a constraint between the defocus images and scene depth. This has several advantages over current learning-based methods that directly takes a focal stack as input [9, 17, 33]. First, output depth satisfies the lens defocus model because the cost volume imposes an explicit constraint between the defocus images and scene depth. Second, the camera settings are absorbed into the cost volume. This enables inference with camera settings that differ from those at training, and even in this case, the output depth satisfies the lens defocus model without any ambiguities.

Figure 4 shows a diagram of our cost volume construction. We first sample the 3D space in the camera coordinate

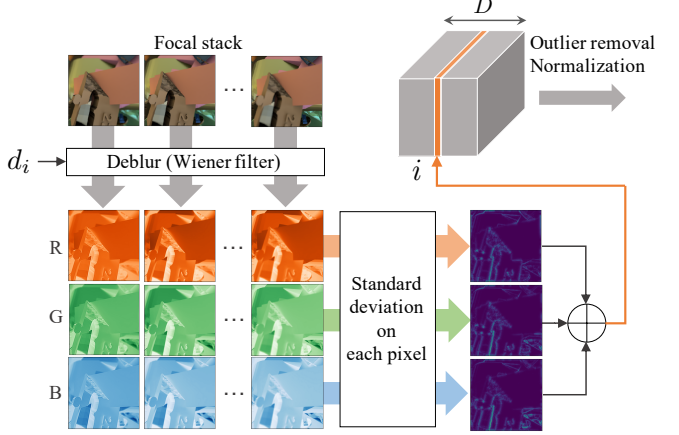


Figure 4. Cost volume construction. We first sweep fronto-parallel plane in camera coordinate system. At each swept plane, input image in focal stack is deblurred with Wiener–Hunt deconvolution [20] on each color channel. Standard deviation is applied for computing cost, which is followed by outlier removal and normalization.

system by sweeping a fronto-parallel plane. To evaluate each depth hypothesis, we deblur each image in the input focal stack. Let the cost volume be  $C : \{1, \dots, W\} \times \{1, \dots, H\} \times \{1, \dots, D\} \rightarrow \mathbb{R}$ , and the focal stack be  $\{I_{d_i}\}_{i=1}^F$ , where  $I_{d_i}$  is a captured image with focus distance  $d_i$ . Each element of the cost volume  $C$  is computed as follows:

$$C(u, v, d) = \sum_{ch \in \{r, g, b\}} \rho\left(\tilde{I}_{d_1}^{ch}(u, v), \dots, \tilde{I}_{d_F}^{ch}(u, v)\right), \quad (5)$$

$$\tilde{I}_{d_i}^{ch} = k(d, d_i) *^{-1} I_{d_i}^{ch}, \quad (6)$$

where  $k(d, d_i)$  is a blur kernel, the size of which is defined by Eq. (1) with the scene depth  $d$  and focus distance  $d_i$ . We used a disk-shaped PSF [23, 34], while any types of PSFs can be used at training and test time. The operator  $*^{-1}$  indicates a deblurring process applied to each color channel of the input image. We used Wiener–Hunt deconvolution [20] as this process. The function  $\rho$  evaluates the consistency between deblurred images. We adopt a standard deviation for  $\rho$ , which allows an arbitrary number of inputs. Note that a similar cost volume computation was proposed in model-based methods [13, 29]. However, these methods require an all-in-focus image, which leads to iterative optimization for the scene depth and all-in-focus image; thus these methods cannot be directly incorporated into sequential learning frameworks.

The process mentioned above is the essential part of our cost volume construction. However, differing from a learning-based MVS method [32], which is based on differentiable image warping, our cost volume construction re-



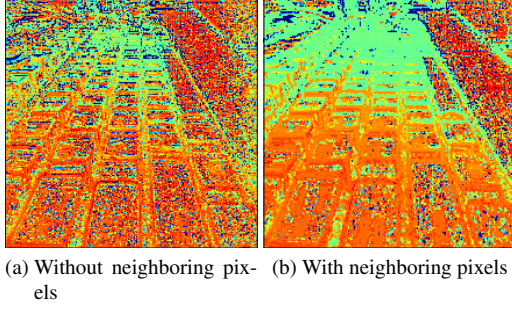


Figure 5. Estimated depth from the index of the minimum cost (a) without and (b) with the neighboring pixels.

quires careful design because the difference between images due to focus distances is smaller than that due to camera positions. Thus, for robustness and learning stability, the standard deviation in Eq. (5) is computed considering neighboring pixels as follows:

$$\rho\left(\tilde{I}_{d_1}^{ch}(u, v), \dots, \tilde{I}_{d_F}^{ch}(u, v)\right) = \sqrt{\frac{1}{F} \sum_{i=1}^F \sum_{(u', v') \in \mathcal{N}(u, v)} \gamma_{u', v'} (\tilde{I}_{d_i}^{ch}(u', v') - \mu(u', v'))^2}, \quad (7)$$

$$\mu(u, v) = \frac{1}{F} \sum_{i=1}^F \sum_{(u', v') \in \mathcal{N}(u, v)} \gamma_{u', v'} \tilde{I}_{d_i}^{ch}(u', v'), \quad (8)$$

where  $\mathcal{N}(u, v)$  is a set of neighboring pixels centered at  $(u, v)$  and  $\gamma_{u', v'}$  is a 2D spatial Gaussian weight. Figure 5 shows an example of the estimated depth only from the index of the minimum cost. The neighboring information can reduce noise, especially for the real captured data.

We also remove outliers by applying a nonlinear function  $f(\cdot)$  that bounds the cost by 1 after computing Eq. (5). We use a tanh-like function as follows:

$$f(x) = \frac{e^{ax} - e^{-ax}}{e^{ax} + e^{-ax}}, \quad (9)$$

$$a = \frac{1}{2C_{max}} \log \frac{1 + f_1}{1 - f_1}, \quad (10)$$

where  $C_{max}$  is the upper bound of the cost.  $f(x)$  is converged to  $f_1$  as  $x$  approaches  $C_{max}$ . We set  $C_{max} = 0.3$  and  $f_1 = 0.999$ . Finally, the cost  $f(C(u, v, d))$  at each pixel is normalized in  $[0, 1]$ . As shown in Fig. 6, this post-processing produces a sharp peak at the ground-truth depth. However, this normalization includes the possibility that such sharp peaks also appear at texture-less pixels where defocus cues are not effective, thus have negative effects on training. Nevertheless, we found that our network automatically learns effective regions and dramatically im-

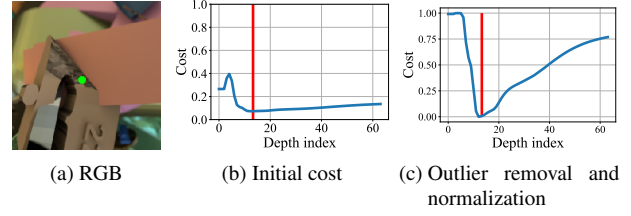


Figure 6. Cost plots (b) without and (c) with outlier removal and normalization at green dot in (a). Red lines indicate positions of ground-truth depth indices.

proves the accuracy of the estimated depth. We describe the ablation study on this in Section 4.4.

### 3.4. Architecture and loss function

As shown in Fig. 2, the cost volume and an additional defocus image, which helps the network to learn semantic cues [32], are concatenated and passed through the network. The input image is selected from the focal stack and we found that the selection of the input image does not affect the performance of the proposed method. During the training of our model, we selected the image with the farthest focus distance.

The cost volume and input image are passed through the encoder, the architecture of which is the same as for MVDepthNet [32]. The outputs of the decoder are refined cost volumes  $C_{out}^s$  at different resolutions  $s \in \{1/8, 1/4, 1/2, 1\}$ .

At each upsample block, we implement an adaptive cost aggregation module inspired by Wang *et al.* [31] to aggregate neighboring information, and this enables depth estimation with clear boundaries by aggregating focus cues on edge pixels. The cost aggregation module is given as

$$\tilde{C}_{out}^s(u, v, d_i) = \sum_{(u_j, v_j) \in \mathcal{N}(u, v)} w_j C_{out}^s(u_j + \Delta u_j, v_j + \Delta v_j, d_i), \quad (11)$$

where the weight  $w_j$  and offset  $(\Delta u_j, \Delta v_j)$  are learnable parameters to aggregate neighboring information. As shown in Fig. 2, our upsample block first upsamples the input cost volume by the scale factor of 2. The feature map from the encoder is then concatenated to this upsampled cost volume. From this volume, offsets and weights for adaptive cost aggregation are learned together with a refined cost volume. The final cost volume is obtained by aggregating the neighboring costs following Eq. (11). Figure 7 shows an example of the learned offsets and output depth with the cost aggregation module, which yields clear boundaries in the estimated depth.

The refined cost volume at each resolution is obtained through softmax layers. Thus, the output depth at each resolution can be computed by applying a differentiable soft

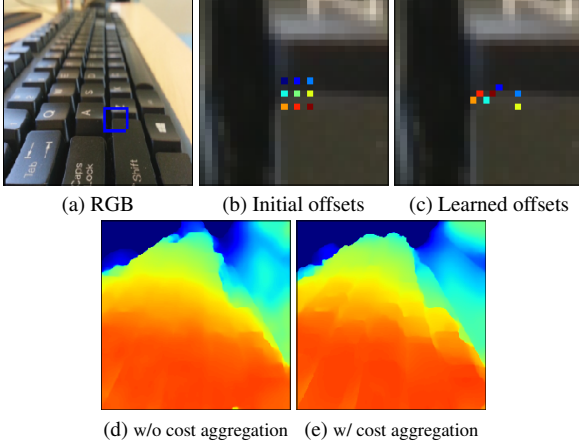


Figure 7. Example of learned offsets for cost aggregation in blue boxed region in (a). (b) At beginning of training, cost is aggregated from nearby grid points. (c) After training, cost is adaptively aggregated by considering local structures. (e) This yields clear boundaries in estimated depth.

argmin operator [12] as follows:

$$d_s(u, v) = \sum_i \tilde{C}_{out}^s(u, v, d_i) d_i. \quad (12)$$

**Training loss** The training loss is defined as the sum of L1 loss between the estimated depth maps  $d_s$  and ground-truth depth maps  $d_s^*$  at different resolutions as follows:

$$\mathcal{L} = \frac{1}{4} \sum_s \frac{1}{H_s W_s} \|d_s - d_s^*\|_1. \quad (13)$$

## 4. Experiments

We evaluated the proposed method for its camera-setting invariance and comparison it with the state-of-the-art learning-based DFF. Our method can be applied to datasets with camera settings that differ from those of a training dataset.

### 4.1. Implementation

Our network was implemented in PyTorch. The training was done on a NVIDIA RTX 3090 GPU with 24-GB memory. The size of a minibatch was 8 for the training of our model. We trained our network from scratch, and the optimizer was Adam [14] with a learning rate of  $1.0 \times 10^{-4}$ .

During the cost volume construction, we uniformly sampled the depth between 0.1 and 3, and set the number of samples to  $D = 64$ .

### 4.2. Dataset

This section describes the datasets for training and evaluation. We used three datasets with the meta data of full camera-settings.

**DefocusNet dataset [17]** This dataset consists of synthetic images, which were generated with physics-based rendering shaders on Blender. The released subset of this dataset has 400 and 100 samples for training and evaluation, respectively. The focal stack of each sample has five images with  $256 \times 256$  resolution. *Note that all models were trained only on this synthetic dataset unless otherwise noted.*

**NYU Depth V2 [24] synthetically blurred by [3]** Carvalho *et al.* [3] generated this dataset by adding synthetic blurs to the NYU Depth V2 dataset [24] that consists of pairs of RGB and depth images. The defocus model was based on Eq. (1) and takes into account object occlusions. The official training and test splits of the NYU Depth V2 dataset are 795 and 654 samples. We extracted  $256 \times 256$  patches from the original  $640 \times 480$  images and finally obtained 9540 and 7848 samples for training and evaluation. As with [17], we rescaled the depth range from  $[0, 10]$  to  $[0, 3]$ . Table 1 lists the camera settings of the DefocusNet dataset [17] and this NYU Depth V2 dataset [3].

**Mobile Depth [29]** This dataset consists of real focal stacks captured with a mobile phone camera. The images in each focal stack were aligned and the authors estimated the camera parameters and depth (*i.e.*, there are no actual ground-truth depth maps.). This dataset contains several scenes; thus, we used this dataset only for evaluation.

### 4.3. Data augmentation

In the DefocusNet dataset, defocus cues are effective only a short distance from a camera [17]. Therefore, we found that our cost volume learned on this dataset is effective only on small depth indices. To enhance the scalability of our cost volume, we scaled the depth maps in the DefocusNet dataset by a scale factor of  $\sigma \in \{1.0, 1.5, 2.0, \dots, 9.0\}$  when we trained our model on this dataset. We should also scale the camera parameters together with the depth map, *i.e.*, if each data sample consists of  $\{\{I_{d_1}, \dots, I_{d_F}\}, \{d_1, \dots, d_F\}, f, N, d^*, b\}$ , the scaled sample is  $\{\{I_{d_1}, \dots, I_{d_F}\}, \{\sigma d_1, \dots, \sigma d_F\}, \sigma f, N, \sigma d^*, b/\sigma\}$ . Note that in both samples, the depth and camera parameters give the same amount of defocus blurs; thus the original focal stack can be used in the scaled sample. This data augmentation is essential for applying our method to other datasets.

Table 1. Camera settings of datasets

Dataset	Size of focal stack	Focus distances [m]	Focal length [m]	f-number	[m/px]
DefocusNet [17]	5	{0.1, 0.15, 0.3, 0.7, 1.5}	$2.9 \times 10^{-3}$	1	$1.2 \times 10^{-5}$
NYU Depth V2 [3]	3	{2, 4, 8}	$15 \times 10^{-3}$	2.8	$5.6 \times 10^{-6}$

Table 2. Ablation study for cost volume construction on DefocusNet dataset [17]. Error metric is RMSE and errors were computed on datasets with different scales of data augmentation.

OR	Norm.	$\sigma = 1.0$	3.0	5.0	7.0	9.0
✓		0.261	0.415	0.450	0.463	0.475
	✓	<b>0.232</b>	<b>0.351</b>	0.377	0.396	0.422
✓	✓	0.239	0.363	<b>0.373</b>	<b>0.380</b>	<b>0.403</b>

Table 3. Experimental results on different focus distances at train and test time on DefocusNet dataset [17]. Both methods were trained on focal stacks with focus distances {0.1, 0.3, 1.5} then tested with focus distances {0.15, 0.7}.

Method	Train	Test	RMSE
DefocusNet [17]	{0.1, 0.3, 1.5}	{0.15, 0.7}	0.299
Ours	{0.1, 0.3, 1.5}	{0.15, 0.7}	<b>0.242</b>

#### 4.4. Ablation study

Table 2 lists the results from the ablation study on the cost volume construction. We separately computed the RMSE on the DefocusNet dataset with a different scale factor of the data augmentation. The experimental results demonstrate that normalization (Norm.) dramatically improved the accuracy of depth estimation. Outlier removal (OR) also improved the accuracy, especially at a large depth scale, where the depth estimation will be more difficult than at a small depth scale, as mentioned in Section 4.3.

#### 4.5. Evaluation on different camera settings

We then evaluated the performance of depth estimation with different camera settings at training and test times. Table 3 lists the experimental results on the DefocusNet dataset. DefocusNet [17], which is a state-of-the-art learning-based DFF method, was compared with our method. We first decomposed each focal stack into two subsets, one with focus distances {0.1, 0.3, 1.5} and the other with {0.15, 0.7}. Both methods were trained only on the subset with focus distances {0.1, 0.3, 1.5} and evaluated on the other subset with different focus distances. Our method outperformed DefocusNet, demonstrating the camera-setting invariance of our method.

We also evaluated the proposed method on the NYU Depth V2 dataset, which has different scene statistics and different camera settings from the DefocusNet dataset, as shown in Table 1. Table 4 and Fig. 8 show the experimen-

tal results when comparing the proposed method other with state-of-the-art learning-based methods, *i.e.*, DDFF [9], AiFDepthNet [33], and DefocusNet [17]. For AiFDepthNet, we used the authors’ trained model, and the other methods were re-trained on the DefocusNet dataset. The parameters of DDFF were initialized by VGG16 [25] as in the original paper [9]. For error metrics, we used MAE, RMSE, absolute relative L1 error (Abs Rel), scale-invariant error (sc-inv) [6], and affine- (scale- and shift-) invariant error in the inverse depth space denoted by ssitrim [22].

As shown in the upper part of Table 1, our method outperformed the other methods trained on the DefocusNet dataset by large margins on most evaluation metrics, and is comparable to DefocusNet on the affine-invariant error metric in the inverse depth space (ssitrim). This is because the camera settings of the DefocusNet and NYU Depth V2 datasets are different. The other methods cannot handle this difference, and the estimated depths have ambiguity.

We also computed the errors on the depths rescaled by the median of the ratios between the output and the ground-truth depths followed by [17] to compensate the scale-ambiguity. The compensation has been done also on our results for fair comparison. The errors are presented in the middle part of the table. Our method also outperformed the other methods in this comparison. In addition, our method without scaling (Ours) still outperformed the rescaled previous methods (\*) in most evaluation metrics. Figure 8 shows examples of the estimated depths. In this figure, the affine-ambiguity of the other methods are compensated by estimating the scales and biases in a least-squares manner (+). Note that the output depths of our method were not rescaled, *i.e.*, our method can estimate depths without any ambiguities. In the bottom part of the table, we show the experimental results trained on the NYU Depth V2 dataset. Although DefocusNet performed better than our method, the accuracy of both methods improved dramatically as shown in Figs. 8(g) and (h), and DefocusNet is heavily affected by the difference of the camera settings in training and test datasets.

Figure 9 shows the experimental results on Mobile Depth with real focal stacks. We set the size of an input focal stack to 3 except for AiFDepthNet [33], which used from about 10 to 30 images for the size of a focal stack, and the model was trained on the synthetically blurred FlyingThings3D dataset [18]. The figure shows the qualitative comparison with the state-of-the-art learning-based methods, the out-

Table 4. Experimental results on blurred NYU Depth V2 dataset [3]. We computed errors on output depth and its rescaled version\* because scales of output depths of current learning-based methods largely differ due to camera setting difference at training and test times. Scale-invariant errors in the depth space (sc-inv [6]) and affine-invariant errors in inverse depth space (ssitrim [22]) were also computed for fair comparison.

Method	Train dataset	MAE	RMSE	Abs Rel	sc-inv	ssitrim
DDFF [9]	DefocusNet	0.719	0.773	0.793	0.199	0.318
AiFDepthNet [33]	DefocusNet	0.425	0.491	0.412	0.319	0.509
DefocusNet [17]	DefocusNet	0.599	0.621	0.706	0.213	<b>0.209</b>
Ours	DefocusNet	<b>0.139</b>	<b>0.186</b>	<b>0.181</b>	<b>0.157</b>	<b>0.209</b>
*DDFF [9]	DefocusNet	0.138	0.313	0.165	0.199	0.318
*AiFDepthNet [33]	DefocusNet	0.239	0.312	0.276	0.319	0.509
*DefocusNet [17]	DefocusNet	0.184	0.322	0.188	0.213	<b>0.209</b>
*Ours	DefocusNet	<b>0.097</b>	<b>0.141</b>	<b>0.126</b>	<b>0.157</b>	<b>0.209</b>
DefocusNet [17]	NYU Depth V2	0.016	0.029	0.018	0.030	0.033
Ours	NYU Depth V2	0.032	0.054	0.034	0.050	0.062

\*Rescaled by median of ratios between output and ground-truth depths.

Table 5. Runtime comparison

Cost volume construction	Depth estimation	DefocusNet [17]
4.278 sec.	0.0252 sec.	0.0285 sec.

put depths of which were rescaled by the median of the ratios between them and the outputs of Suwajanakorn *et al.* [29] (\*) followed by [17]. Note that the output depths of our method were not rescaled. The output depths of our method are qualitatively plausible and satisfy the defocus model under different camera settings. Figure 10 shows the quantitative errors between our method and Suwajanakorn *et al.* [29] under different sizes of input focal stacks, demonstrating that a few images are enough to obtain effective results with our method.

Finally, we show an example of applying the proposed method to real focal stacks captured with our camera, Nikon D5300 with f-number of 1.8. The focal stacks were captured with “Focus Stacking Simple” in digiCamControl [1]. All parameters required for the cost volume computation were extracted from EXIF properties, and the focal stack size was 3. Figure 11 shows the qualitative evaluation results. The values of the estimated depth maps are in meters. These results indicate the applicability of our method to real focal stacks.

#### 4.6. Computation time

Table 5 shows the runtime comparison. We measured the processing time for each test sample in the DefocusNet dataset [17]. The cost volume construction was done on AMD EPYC 7232P@3.1 GHz with 128GB RAM. The number of the depth samples in the cost volume is 64 and the image resolution is  $256 \times 256$ . Although the cost vol-

ume construction takes a few seconds, the costs at different depth slices in our cost volume can be computed in parallel to reduce the computation time.

#### 4.7. Limitations

We finally discuss the limitations of the proposed methods due to the explicit lens defocus model.

**Dynamic scenes and focus breathing** Similar to AiFDepthNet [33], our cost volume computation allows only static scenes. Focus breathing also affects our method. However, as mentioned in [33], simple preprocessed alignment can solve this problem (In the experiments with real data (Fig. 9), we used aligned focal stacks).

**Trade-off between defocus and semantic cues** We finally discuss the trade-off between model- and learning-based approaches. Table 6 and Fig. 12 show the results on the DefocusNet dataset. The other learning-based methods outperformed our method. This is because the defocus cues in the DefocusNet dataset are effective only at a short distance from a camera, as mentioned in Section 4.3. The other learning-based methods handle this limitation through semantic cues. Although our method also learns semantic cues, our method with the explicit lens defocus model is more affected by this limitation. For future work, a network architecture should be designed to effectively learn defocus and semantic cues simultaneously.

### 5. Conclusion

We proposed learning-based DFF with a lens defocus model. We combined a learning framework and defocus model with the construction of a cost volume. This method



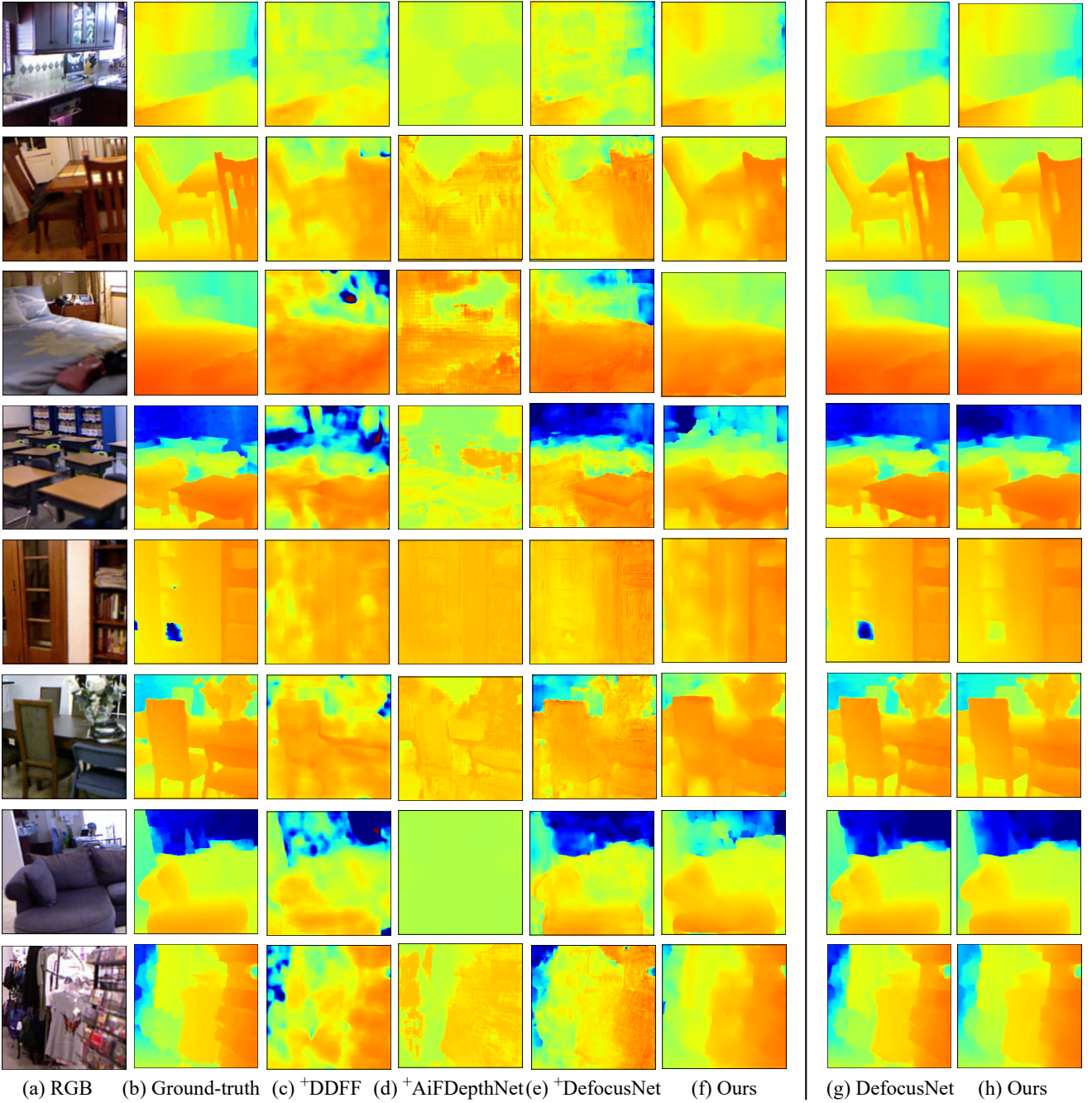


Figure 8. Qualitative comparison on NYU Depth V2 [3]. In (a)-(f), all models were trained on DefocusNet dataset [17]. In (g) and (h), both methods were trained on NYU Depth V2 dataset. Superscript <sup>+</sup> means that affine-ambiguity is compensated by estimating scales and biases in least-squares manner between output and ground-truth.

can absorb the difference in camera settings through the cost volume, which allows the method to estimate the scene depth from a focal stack with different camera settings at training and test times. The experimental results indicate that our model trained only on a synthetic dataset can be

applied to other datasets including real focal stacks with different camera settings. This camera-setting invariance will enhance the applicability of learning-based DFF methods.

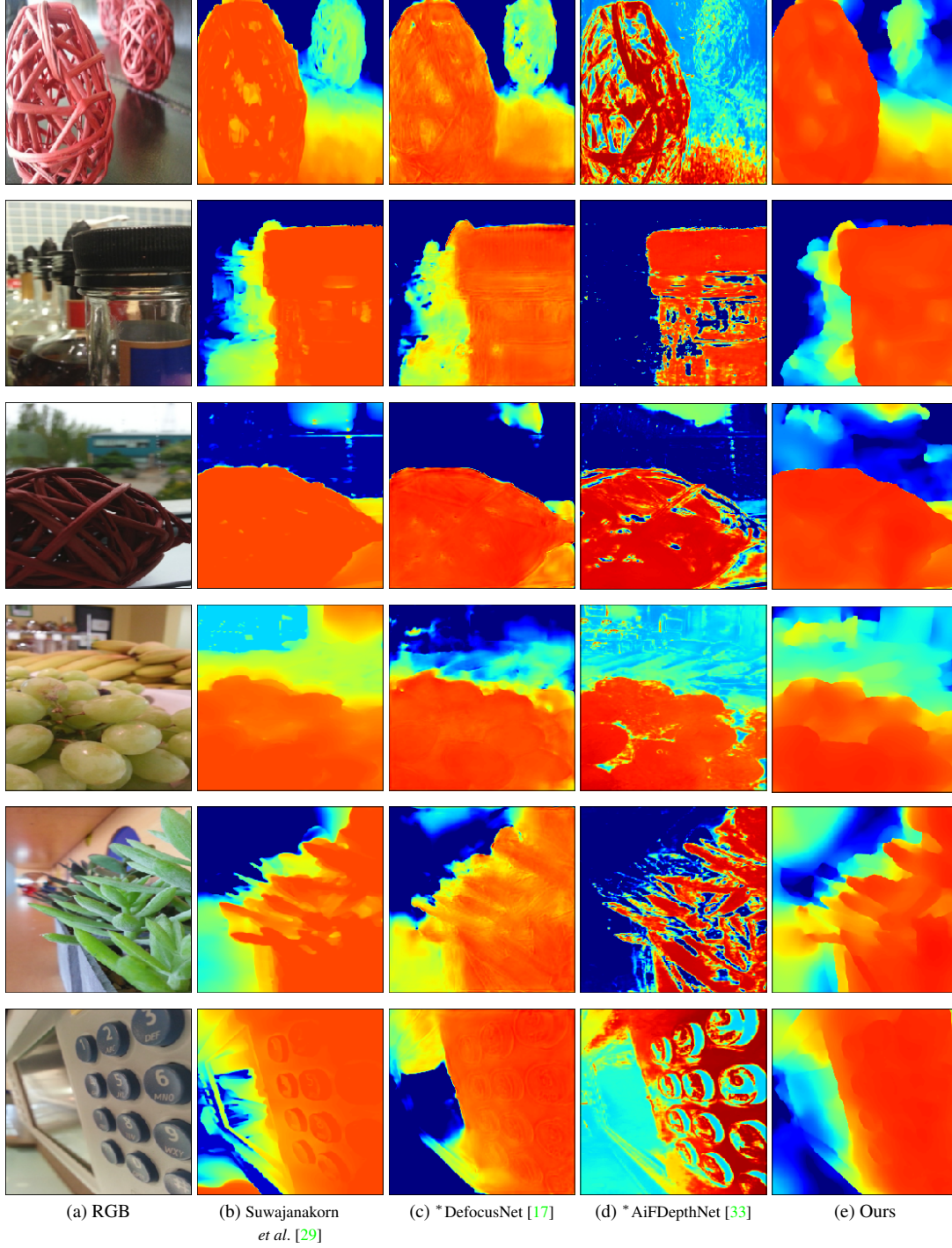


Figure 9. Experimental results on Mobile Depth [29]. Superscript \* means that depth is rescaled by median of ratios between output and Suwajanakorn et al. [29].

## References

- [1] digicamcontrol. <http://digicamcontrol.com/>. 8
- [2] Saeed Anwar, Zeeshan Hayder, and Fatih Porikli. Depth estimation and blur removal from a single out-of-focus image. In *BMVC*, 2017. 2
- [3] Marcela Carvalho, Bertrand Le Saux, Pauline Trounev-Peloux, Andres Almansa, and Frederic Champagnat. Deep depth from defocus: how can defocus blur improve 3d estimation using dense neural networks? In *ECCVW*,

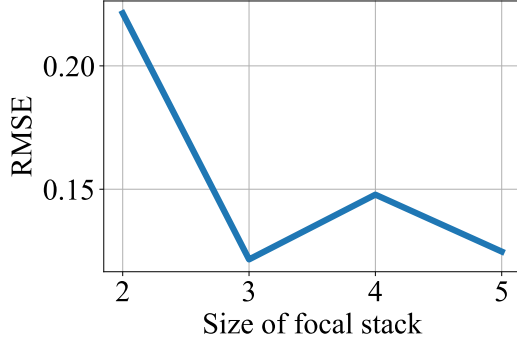


Figure 10. Ablation study of focal stack size on Mobile Depth [29]. Horizontal and vertical axes represent RMSE and size of input focal stack.

Table 6. Experimental results on DefocusNet dataset

Method	RMSE
AiFDepthNet [33]	0.156
DefocusNet [17]	0.177
Ours	0.239

2018. [https://github.com/marcelampc/d3net\\_depth\\_estimation](https://github.com/marcelampc/d3net_depth_estimation) (GPLv3 license). 2, 6, 7, 8, 9

- [4] R.T. Collins. A space-sweep approach to true multi-image matching. In *CVPR*, pages 358–363, 1996. 2, 3
- [5] Arda Duzceker, Silvano Galliani, Christoph Vogel, Pablo Speciale, Mihai Dusmanu, and Marc Pollefeys. Deep-videomvs: Multi-view stereo on video with recurrent spatio-temporal fusion. In *CVPR*, pages 15324–15333, 2021. 2
- [6] David Eigen, Christian Puhrsch, and Rob Fergus. Depth map prediction from a single image using a multi-scale deep network. In *NeurIPS*, volume 2, pages 2366–2374, 2014. 7, 8
- [7] Rahul Garg, Neal Wadhwa, Sameer Ansari, and Jonathan T. Barron. Learning single camera depth estimation using dual-pixels. In *ICCV*, pages 7628–7637, 2019. 4
- [8] Shir Gur and Lior Wolf. Single image depth estimation trained via depth from defocus cues. In *CVPR*, pages 7683–7692, 2019. 2
- [9] Caner Hazirbas, Sebastian Georg Soyer, Maximilian Christian Staab, Laura Leal-Taixé, and Daniel Cremers. Deep depth from focus. In *ACCV*, 2018. <https://github.com/soyers/ddff-pytorch> (GNU General Public License v3.0). 1, 2, 4, 7, 8
- [10] Junjie Hu, Mete Ozay, Yan Zhang, and Takayuki Okatani. Revisiting single image depth estimation: Toward higher resolution maps with accurate object boundaries. In *IEEE Winter Conference on Applications of Computer Vision (WACV)*, pages 1043–1051, 2019. 2
- [11] Eddy Ilg, Nikolaus Mayer, Tonmoy Saikia, Margret Keuper, Alexey Dosovitskiy, and Thomas Brox. FlowNet 2.0: Evolution of optical flow estimation with deep networks. In *CVPR*, pages 2462–2470, 2017. 2
- [12] 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 *ICCV*, 2017. 2, 3, 6
- [13] Hyeonwoo Kim, Christian Richardt, and Christian Theobalt. Video depth-from-defocus. In *International Conference on 3D Vision (3DV)*, pages 370–379, 2016. 1, 2, 4
- [14] Diederik P. Kingma and Jimmy Lei Ba. Adam: A method for stochastic optimization. In *ICLR*, 2015. 6
- [15] Yanghao Li, Naiyan Wang, Jiaying Liu, and Xiaodi Hou. Demystifying neural style transfer. In *IJCAI*, pages 2230–2236, 2017. 2
- [16] Xiaoxiao Long, Lingjie Liu, Wei Li, Christian Theobalt, and Wenping Wang. Multi-view depth estimation using epipolar spatio-temporal networks. In *CVPR*, pages 8258–8267, 2021. 2
- [17] Maxim Maximov, Kevin Galim, and Laura Leal-Taixé. Focus on defocus: bridging the synthetic to real domain gap for depth estimation. In *CVPR*, pages 1071–1080, 2020. <https://github.com/dvl-tum/defocus-net> (MIT License). 1, 2, 3, 4, 6, 7, 8, 9, 10, 11
- [18] Nikolaus Mayer, Eddy Ilg, Philip Hausser, Philipp Fischer, Daniel Cremers, Alexey Dosovitskiy, and Thomas Brox. A large dataset to train convolutional networks for disparity, optical flow, and scene flow estimation. In *CVPR*, pages 4040–4048, 2016. 2, 7
- [19] Michael Moeller, Martin Benning, Carola Schönlieb, and Daniel Cremers. Variational depth from focus reconstruction. *IEEE TPAMI*, 24(12):5369–5378, 2015. 2
- [20] François Orieux, Jean-François Giovannelli, and Thomas Rodet. Bayesian estimation of regularization and point spread function parameters for wiener–hunt deconvolution. *Journal of the Optical Society of America A*, 27(7):1593–1607, 2010. 4
- [21] Said Pertuz, Domenec Puig, and Miguel Angel Garcia. Analysis of focus measure operators for shape-from-focus. *PR*, 46:1415–1432, 2013. 1, 2
- [22] René Ranftl, Katrin Lasinger, David Hafner, Konrad Schindler, and Vladlen Koltun. Towards robust monocular depth estimation: Mixing datasets for zero-shot cross-dataset transfer. *IEEE TPAMI*, 2020. 7, 8
- [23] Jianping Shi, Xin Tao, Li Xu, and Jiaya Jia. Break ames room illusion: depth from general single images. *ACM TOG*, 34(6):1–11, 2015. 4
- [24] Nathan Silberman, Derek Hoiem, Pushmeet Kohli, and Rob Fergus. Indoor segmentation and support inference from rgb-d images. In *ECCV*, pages 746–760, 2012. 6
- [25] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. In *ICLR*, 2015. 7
- [26] Pratul P. Srinivasan, Rahul Garg, Neal Wadhwa, Ren Ng, and Jonathan T. Barron. Aperture supervision for monocular depth estimation. In *CVPR*, pages 6393–6401, 2018. 2
- [27] Deqing Sun, Xiaodong Yang, Ming-Yu Liu, and Jan Kautz. Pwc-net: Cnns for optical flow using pyramid, warping, and cost volume. In *CVPR*, pages 8934–8943, 2018. 2



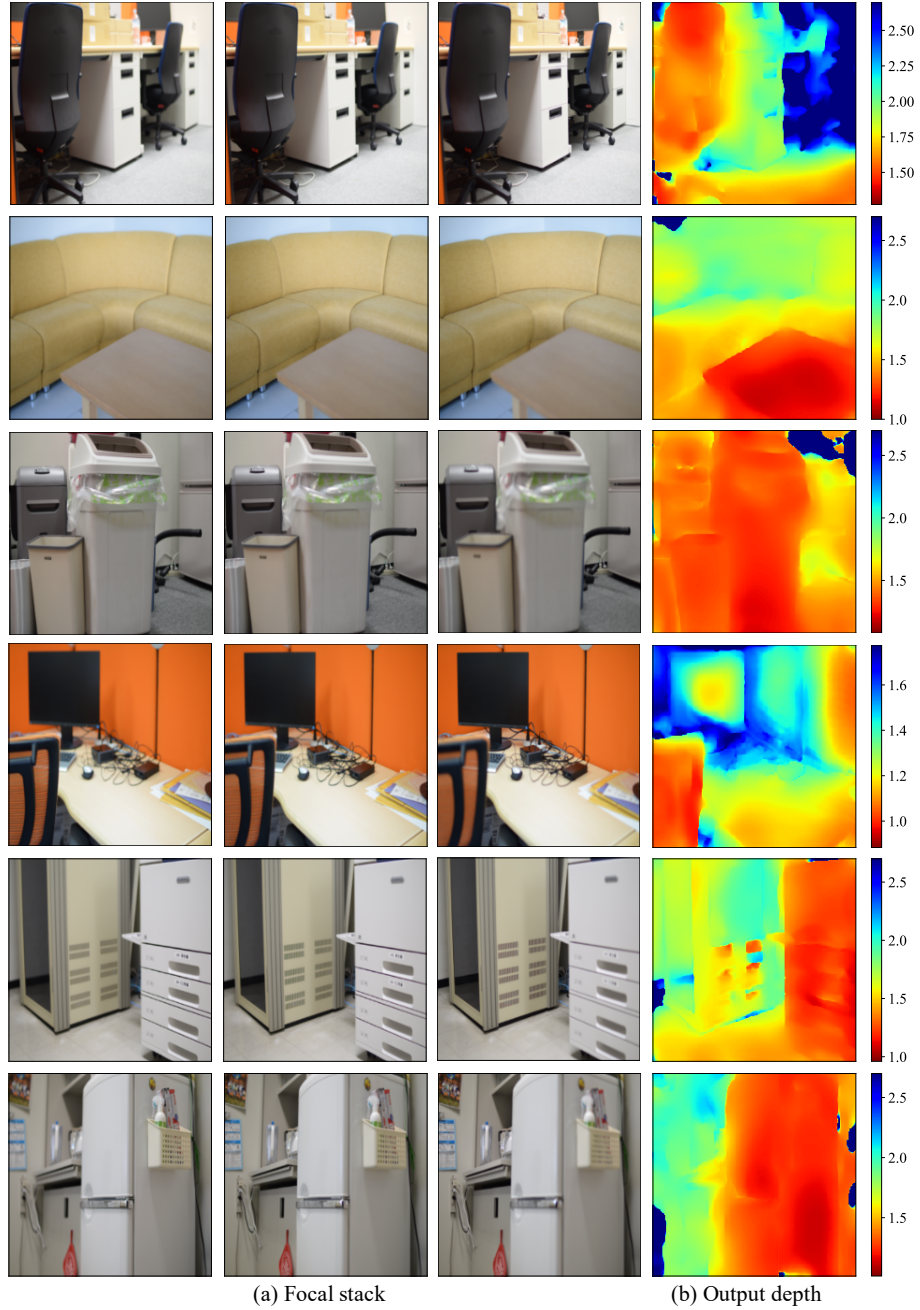


Figure 11. Experimental results on focal stacks captured with our camera. (a) Images in focal stack and (b) output depth of our method.

- [28] Jaeheung Surh, Hae-Gon Jeon, Yunwon Park, Sunghoon Im, Hyowon Ha, and In So Kweon. Noise robust depth from focus using a ring difference filter. In *CVPR*, pages 6328–6337, 2017. 1, 2
- [29] Supasorn Suwajanakorn, Carlos Hernandez, and Steven M. Seitz. Depth from focus with your mobile phone. In *CVPR*, pages 3497–3506, 2015. <https://www.supasorn.com/dffdownload.html>. 1, 2, 4, 6, 8, 10, 11
- [30] Huixuan Tang, Scott Cohen, Brian Price, Stephen Schiller, and Kiriakos N. Kutulakos. Depth from defocus in the wild. In *CVPR*, pages 2740–2748, 2017. 1, 3
- [31] Fangjinhua Wang, Silvano Galliani, Christoph Vogel, Pablo Speciale, and Marc Pollefeys. Patchmatchnet: Learned multi-view patchmatch stereo. In *CVPR*, pages 14194–14203, 2021. 5
- [32] Kaixuan Wang and Shaojie Shen. Mvdepthnet: Real-time multiview depth estimation neural network. In *International Conference on 3D Vision (3DV)*, pages 248–257, 2018. 2, 3,



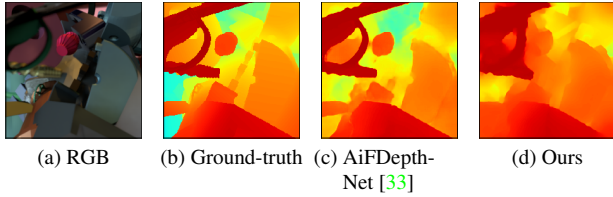


Figure 12. Limitations of our method. (a) One of input images in focal stack, (b) Ground-truth depth, (c) output depth of AiFDepthNet [33], and (d) that of our method. Defocus cues in the DefocusNet dataset are effective only at a short distance from the camera, and our method with explicit defocus model is more affected by this limitation.

4, 5

- [33] Ning-Hsu Wang, Ren Wang, Yu-Lun Liu, Yu-Hao Huang, Yu-Lin Chang, Chia-Ping Chen, and Kevin Jou. Bridging unsupervised and supervised depth from focus via all-in-focus supervision. In *ICCV*, 2021. <https://github.com/albert100121/AiFDepthNet>. 1, 2, 3, 4, 7, 8, 10, 11, 13
- [34] Masahiro Watanabe and Shree K. Nayar. Rational filters for passive depth from defocus (article) author. *IJCV*, 27(3):203–225, 1998. 4
- [35] Y. Xiong and S.A. Shafer. Depth from focusing and defocusing. In *CVPR*, pages 68–73, 1993. 2
- [36] Yao Yao, Zixin Luo, Shiwei Li, Tian Fang, and Long Quan. Mvsnet: Depth inference for unstructured multi-view stereo. In *ECCV*, pages 767–783, 2018. 2
- [37] Jun-Yan Zhu, Taesung Park, Phillip Isola, and Alexei A. Efros. Unpaired image-to-image translation using cycle-consistent adversarial networks. In *ICCV*, pages 2223–2232, 2017. 2
- [38] Shaojie Zhuo and Terence Sim. Defocus map estimation from a single image. *PR*, 44(9):1852–1858, 2011. 2, 3