

Distortion-aware Monocular Depth Estimation for Omnidirectional Images

Hong-Xiang Chen[†], Kunhong Li[†], Zhiheng Fu[§], Mengyi Liu[‡], Zonghao Chen[‡], Yulan Guo^{†*}
School of Electronics and Communication Engineering, Sun Yat-sen University
{chenhx97, likh25}@mail2.sysu.edu.cn, guoyulan@sysu.edu.cn

University of Western Australia[§]
22907304@student.uwa.edu.au

Alibaba Group[‡]
{suqing.lmy, czh190502}@alibaba-inc.com@alibaba-inc.com

Abstract

A main challenge for tasks on panorama lies in the distortion of objects among images. In this work, we propose a Distortion-Aware Monocular Omnidirectional (DAMO) dense depth estimation network to address this challenge on indoor panoramas with two steps. First, we introduce a distortion-aware module to extract calibrated semantic features from omnidirectional images. Specifically, we exploit deformable convolution to adjust its sampling grids to geometric variations of distorted objects on panoramas and then utilize a strip pooling module to sample against horizontal distortion introduced by inverse gnomonic projection. Second, we further introduce a plug-and-play spherical-aware weight matrix for our objective function to handle the uneven distribution of areas projected from a sphere. Experiments on the 360D dataset show that the proposed method can effectively extract semantic features from distorted panoramas and alleviate the supervision bias caused by distortion. It achieves state-of-the-art performance on the 360D dataset with high efficiency.

1. Introduction

3D scene perception and understanding is a fundamental technique for many applications such as robotics and intelligent vehicles. Among various 3D vision tasks, depth estimation is highly important since it forms the basis for many downstream tasks such as obstacle avoidance and object fetching. Due to the high costs of 3D sensors (e.g., LiDARs), inferring 3D information from 2D RGB images captured by cheap consumer-level cameras is significantly important.

Dense depth estimation is a challenging pixel-level task in 3D scene understanding. Different from stereo matching and Structure from Motion (SfM) methods, monocular depth estimation is an ill-posed problem for its information

loss caused by the projection from a 3D space to a 2D image plane. That is, one pixel in a 2D image may correspond to multiple points in a 3D space. Thanks to the power of deep learning, progressive advances have been achieved using prior geometric constraints extracted either explicitly or implicitly from annotated data [21, 9, 33]. However, most of these methods focus on predicting depth from general perspective images.

Distortion is a major challenge for tasks working on panoramas, such as classification, saliency detection, semantic segmentation and depth estimation [7, 22, 34]. Directly applying conventional CNNs on panoramas (e.g., represented in equirectangular formats) is hard to achieve promising performance. Since a panorama is usually produced by stitching several perspective images captured by a perspective camera located at the same place, equirectangular projection can be considered as a transformation from a non-Euclidean space to an Euclidean space and thus introducing distortion to panoramas. As a consequence, the projection of objects has irregular shapes and the distortion becomes extremely significant for pixels close to the poles or image plane. Therefore, standard convolution is unsuitable for panorama processing.

In this work, we propose a Distortion-Aware Monocular Omnidirectional (DAMO) network by combining strip pooling and deformable convolution to generate accurate depth maps from panoramas with distortion. We first present a distortion-aware feature extraction block to handle distortion introduced by equirectangular projection. Specifically, we utilize deformable convolution to learn offsets for the sampling grids, resulting in feature maps that are much denser than regular grids. We then exploit strip pooling to capture anisotropy context information from irregular regions (i.e., the distorted projection of objects) and preserve the integral distortion information for convolution sampling. In addition, to mitigate supervision bias caused by uneven sampling in different areas, we also propose an easy-to-use spherical-aware weight matrix for the objective

[†]indicates the corresponding author.

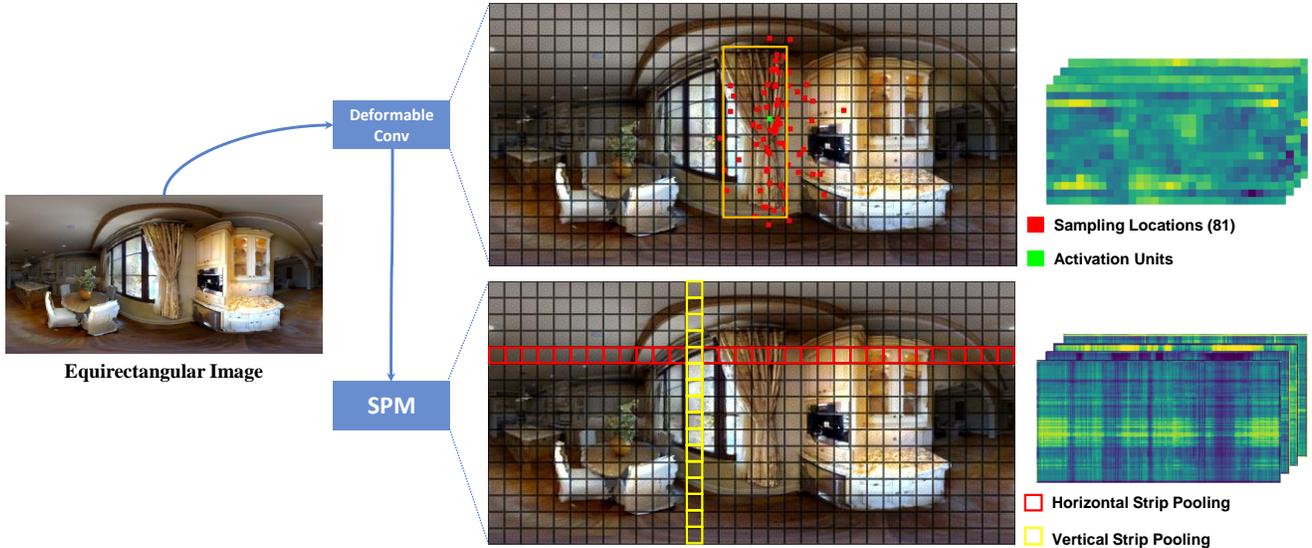


Figure 1. The paradigm of our distortion-aware feature extraction block. To illustrate the learned distortion, we use the area around the curtain in the image as an example. 81 locations (see red dots) are sampled by our DAMO network, it is clear that these sampling locations for an activation unit are mainly around the curtain (see green dot). Then, we utilize SPM to help deformable convolution focus mainly on informative regions and thus reduce the impact of distortion in panorama.

function. Experiments on the 360D dataset demonstrate that our DAMO network achieves the state-of-the-art performance with high efficiency.

Our contributions can be summarized as follows:

- We propose a DAMO network to handle distortion in panoramas using both deformable convolution and strip pooling module. Experiments on the 360D dataset show that DAMO is superior to the state-of-the-art.
- We introduce a plug-and-play spherical-aware weight for our objective function to make the network focus on informative areas. This weight helps our network to achieve fast convergence and improved performance.

2. Related Work

We will briefly describe several existing methods related to our work in this section.

2.1. Depth Estimation

Depth estimation has been a hot topic for a long time. Early studies [29, 26, 23] in this area focused on developing algorithms to generate point correspondences in stereo images. Different from these methods, Delage et al. [6] developed a Bayesian framework to perform 3D indoor reconstruction from one single perspective image based on a strong floor-wall assumption. Saxena et al. [27, 28] used Markov Random Fields (MRFs) to incorporate multiscale

and global image features to predict depth from a single RGB image.

Eigen et al. [10] proposed the first deep learning based network. They used a multiscale convolutional architecture to predict results in a coarse-to-fine manner. Eigen et al. [9] then adopted a multi-task training scheme to further improve the performance of their model. Laina et al. [21] proposed a regularization concerning loss and a uniform up-projection module for monocular depth estimation, which have been frequently used in subsequent methods. Fu et al. [12] considered the depth estimation task as an ordinal regression problem by applying a spacing-increasing discretization strategy and a well-designed ordinal regression loss. Yin et al. [33] improved the supervision capability of the objective function by randomly selecting a number of ternary points and producing a virtual plane for each ternary. With its geometric supervision, the virtual normal loss improves the convergence of the depth estimation model. However, all these methods focus on perspective images and may easily be trapped into suboptimal results while being directly applied to panoramas.

2.2. Representations on Panorama

Equirectangular image is one of the most widely used representations for panoramas and distortion has been a major challenge for years. Su et al. [31] used an adaptive kernel to handle the distortion near the pole. Following this idea, Zioulis et al. [36] designed a set of rectangular filterbanks to deal with the horizontal distortion introduced by

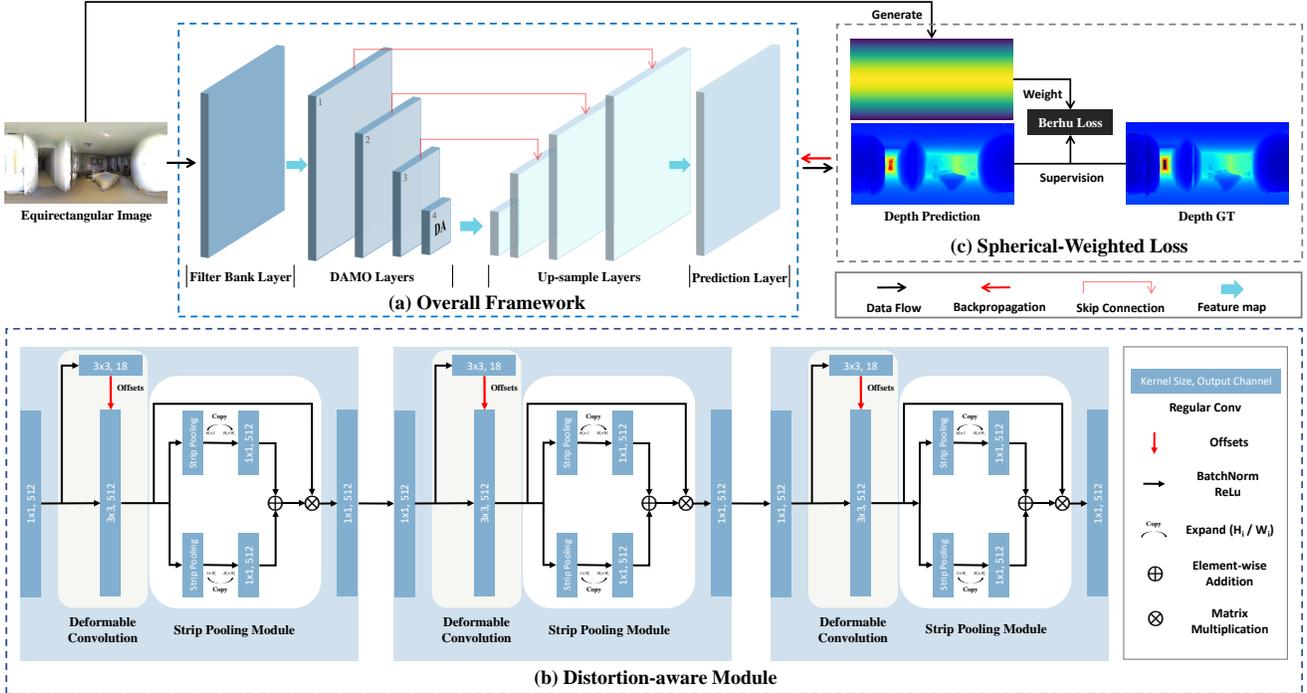


Figure 2. An overview of our DAMO network. We build a filter-bank layer utilizing a group of parallel rectangular convolutions to extract horizontal distortion-aware features following [36, 32]. Then, these features are fed into an encoder-decoder based architecture [24] with a skip connection at each resolution. Besides, we also add SPM in the last building block for each of first three DAMO layers to help deformable convolution focus on contextual regions.

equirectangular projection by increasing the receptive field of conventional convolution kernels on horizon. Although the representation capability of CNNs on panoramas has been improved by these methods, the gap between omnidirectional and perspective images still exists.

Cube map is another commonly used representation for panoramas. This representation faces the challenge of inconsistency between different faces. Cheng et al. [3] proposed cube padding to reduce the information loss along edges between faces. Wang et al. [32] further extended [3] to spherical padding and propose a two-branch encoder-decoder based network to predict depth maps for panoramas. However, their model was hard to train due to its multiple training settings and high time cost. Inspired by rectangular convolution in [36], we exploit strip pooling [16] to preserve more context details for convolution.

2.3. Dynamic Mechanism

Existing deep learning based dynamic mechanisms can be divided into two categories: weight based methods [19, 17, 25] and offset based methods [18, 5].

Weight based methods focus on adaptively generating weights for either feature map selection or channel-wise selection. For instance, Jia et al. [19] used a flexible filter

generation network to produce a set of filter operators that dynamically conditioned on an individual input, resulting in improved performance on the video and stereo prediction tasks. Besides, attention is investigated to generate weights. Hu et al. [17] proposed a light-weight gating mechanism to explicitly model channel-wise dependencies and further improve the model representation ability using global information.

Offset based methods aim at providing offsets for filters to aggregate more geometric information. Jeon et al. [18] proposed an Active Convolution Unit (ACU) to learn its own shape adaptively during training. However, the shapes of filters have to be fixed after training. Moreover, Dai et al. [5] introduced deformable convolution to learn offsets for each spatial location dynamically, resulting in higher generalization capability than ACU.

3. Proposed Method

The overall pipeline of our DAMO network for monocular dense depth estimation on panorama is shown in Fig. 2. In this section, we will first introduce our Distortion-Aware (DA) module for calibrated semantic feature extraction, including strip pooling and deformable convolution. Then, we will introduce our objective function with a spherical-aware

weight matrix.

3.1. Distortion-aware Module

The DA module is shown in Fig. 2(b). We first utilize deformable convolution to extract learnable distorted information and calibrate semantic features. After that, the learned distortion knowledge is fed into the Strip Pooling Module (SPM) [16]. The feature maps are further activated by multiple times and thus the distortion can be sufficiently learned by our DA module.

3.1.1 Deformable Convolution

To learn distortion knowledge and model the spatial transformation of convolution kernels on panorama, we adopt deformable convolution [5] in this work. Given a set of sampling locations on a regular grid R , the input feature map \mathbf{t} , the output feature map F and the weights for kernel \mathbf{w} , a conv layer is applied in conventional convolution neural networks. Here, we define the grid for 3×3 convolution kernel with a dilation of 1 as $R = \{(\pm 1, \pm 1), (\pm 1, 0), (0, \pm 1), (0, 0)\}$.

According to [5], deformable convolution can be formulated as:

$$F(\mathbf{p}_0) = \sum_{\mathbf{p}_n \in R} \mathbf{w}(\mathbf{p}_n) \cdot \mathbf{t}(\mathbf{p}_0 + \mathbf{p}_n + \Delta \mathbf{p}_n) \quad (1)$$

where \mathbf{p}_0 represents a location on output feature map F , \mathbf{p}_n is one position on the regular sampling grid R and $\Delta \mathbf{p}_n$ represents the offset corresponding to the position of \mathbf{p}_n .

The i -th input feature map (with a size of $C_i \times B_i \times H_i \times W_i$) is first fed into a convolutional layer to generate a group of 2D offsets for each sampling location on grid R . The offsets for both vertical and horizontal translation on R have a size of $18 \times B_i \times H_i \times W_i$. Here, we use a 3×3 sampling grid R as an example. Then, the sampling grid for deformable convolution is generated through Eq. 1 using the 2D offsets $\Delta \mathbf{p}_n$. Since most learnable offsets are fractional, we use bilinear interpolation to generate integer offsets [5, 35].

3.1.2 Strip Pooling

To preserve more distortion information and help the network to focus on informative regions, we adopt strip pooling [16] in our DA module. Different from standard spatial max pooling that adopts two-dimensional sampling grids, sampling along vertical and horizontal orientations are performed separately in strip pooling. That is, contextual information in feature maps are selected either in a row or a column according to its input (as shown in Fig. 2). Similar to the origin definition in [16], strip pooling is defined as:

$$\begin{cases} y_{c,i}^h = \max_{0 \leq j < W} x_{c,i,j} \\ y_{c,j}^v = \max_{0 \leq i < H} x_{c,i,j} \end{cases} \quad (2)$$

where $x_{c,i,j}$ is a location of the sampling grid on input (i.e., $\mathbf{x} \in \mathbb{R}^{C \times H \times W}$), H and W also represent the kernel size of strip pooling along horizontal and vertical orientations, respectively. $y_{c,i}^h \in \mathbf{y}^h$ and $y_{c,j}^v \in \mathbf{y}^v$ represent the i -th and j -th grids on the feature map of channel c for horizontal and vertical strip pooling, respectively.

Both vertical and horizontal strip pooling layers are used in SPM. Once these two expanded output feature maps (with the same size as their inputs) are added, an element-wise multiplication $G(\cdot, \cdot)$ is used to activate the contextual geometric areas of the input. SPM is sensitive to anisotropy context and produces denser distortion feature maps than the regular max pooling module. Given the combination of \mathbf{y}^h and \mathbf{y}^v , the output \mathbf{z} of SPM can be formulated as:

$$\mathbf{z} = G(\mathbf{x}, \sigma(f(\mathbf{y}))) \quad (3)$$

where σ is the sigmoid function, f is the last 1×1 convolution layer, the fusion of horizontal and vertical strip pooling information $y_{c,i,j} \in \mathbf{y}$ is defined as $y_{c,i,j} = y_{c,i}^h + y_{c,j}^v$.

We apply our SPM to each conv k layer of the DAMO layer (i.e., conv2, conv3, conv4 and conv5, as defined in ResNet [15]). Specifically, SPMs are used in the last building block of the first three DAMO layers and are stacked for all building blocks of the DA module (see Fig. 2).

3.2. Spherical-aware Weighted Loss

Since an equitangular image represents a panorama in 2D space, distortion is extremely large in the areas around the pole of a spherical space. Specifically, for two areas with the same coverage in a spherical surface, the area near the pole is much larger than the area near the equator in an equitangular image due to uneven projection. Loss functions defined in perspective images [36, 8, 4] are hard to produce optimal results due to the overfitting (weighting) in sparse areas near the pole and the underfitting in dense areas near the equator.

3.2.1 Weight Matrix

To achieve balanced supervision in different areas, we introduce a spherical-aware weighting strategy on objective function. Specifically, the Cartesian coordinates of a pixel $p_E = D(x, y)$ in the equirectangular image can be converted to spherical coordinates $p_S = \Pi(\theta, \phi)$ in a spherical surface, where longitude $\theta \in [0, 2\pi]$ and latitude $\phi \in [0, \pi]$. That is, $\phi_{(x,y)} = \frac{\pi x}{H}$, $\theta_{(x,y)} = \frac{2\pi y}{W}$.

Considering a sphere $\Pi(\theta, \phi)$ with a unit radius, we generate a weight matrix for objective functions according to the sphere angle of a pixel along the latitude. Take the north

hemisphere as an example, the weight is defined as the ratio of the area from the north pole to the current latitude to the total area of the sphere surface. The weights in the south hemisphere can be calculated following a similar way.

$$W_{(x,y)} = \int_0^{\phi(x,y)} \sin\phi d\phi \quad (4)$$

Here, W is the weight matrix for the objective function in each location. $\phi_{(x,y)}$ denotes the angle of the position (x, y) along the vertical axis.

3.2.2 Loss Function

Following [33, 11], we adopt the reverse Huber loss (also called the Berhu loss) [21] in this work.

$$\mathcal{L}_d = \begin{cases} |d_{pre} - d_{gt}|, & \text{if } |d_{pre} - d_{gt}| \leq \tau \\ \frac{(d_{pre} - d_{gt})^2 + \tau^2}{2\tau}, & \text{if } |d_{pre} - d_{gt}| > \tau \end{cases} \quad (5)$$

where d_{pre} and d_{gt} are the predicted and groundtruth depth values, respectively. The Berhu loss can achieve a good balance between L1 and L2 norms. Specifically, pixels with high gradient residuals will be assigned with large weights using the L2 term. Meanwhile, the L1 term pays more attention to regions with small gradients. In our experiments, we set the threshold τ as 20% of the maximum error between prediction and groundtruth.

Finally, our loss function is defined as:

$$\mathcal{L} = \mathcal{L}_d \otimes W \quad (6)$$

Compared to the original Berhu loss \mathcal{L}_d , our weighted Berhu loss \mathcal{L} can help the network to focus more on informative regions (i.e., regions near the equator) and mitigate the effects introduced by severe distortion (which is significant in areas near the pole) in panoramas.

4. Experiments

We conduct extensive experiments on a widely used omnidirectional dataset to evaluate the performance of our DAMO network. We will first describe the dataset and the evaluation toolbox, and then present the implementation details of our experiments. We further compare our method to the state-of-the-art.

4.1. Experimental Settings

4.1.1 Dataset

We adopt a large-scale indoor omnidirectional RGBD dataset (i.e., the 360D Dataset [36]) to conduct experiments. This dataset contains two real-world datasets (i.e., Stanford2D3D [1] and Matterport3D [2]), and two synthetic datasets (i.e., SunCG [30] and SceneNet [14]). Following

the original split in [36], the training and test sets are listed as follows:

- Training set: a cross-domain set of both real-world and synthetic images from the Stanford2D3D, Matterport3D and SunCG datasets were first obtained. Then, scenes with very close or far regions were removed, resulting in a training dataset with 34,679 RGBD image pairs.
- Test set: 1,289 omnidirectional image pairs were collected from the Stanford2D3D, Matterport3D and SunCG dataset for test. The remaining image pairs from the SceneNet dataset were used for validation.

4.1.2 Implementation Details

Our network was implemented in Pytorch with a single Nvidia RTX Titan GPU. Each RGB image has a resolution of 512×256 while invalid depth values are removed by a mask. The batch size was set to 8 and the initial learning rate was set to 1×10^{-4} . We used the Adam optimizer [20] with its default settings and a poly learning rate policy [13] for training. All models were trained for 20 epochs on the 360D dataset for fair comparison.

4.1.3 Evaluation Metrics

We adopted the same metrics as previous works [36, 4] for fair comparison, including Absolute average Relative Error (Abs.REL), Root Mean Squared Error (RMSE), Root Mean Squared Error in logarithmic space (RMSElog) and accuracy with a threshold δ_t , where $t \in \{1.25, 1.25^2, 1.25^3\}$. Note that, we used the same evaluation strategy as [36]. That is, depth maps were estimated by dividing a median scalar \bar{s} to achieve direct comparison among multiple datasets with different range scales, where $\bar{s} = \text{median}(D_{GT})/\text{median}(D_{Pred})$.

4.2. Comparison to the State-of-the-art

We compare our DAMO to two existing methods. Note that, only the Caffe model is provided by [36] while the source codes are provided by [32]. We used their released model or source codes for comparison in this work.

4.2.1 Quantitative Comparison

As shown in Table 1, our method outperforms the baseline method [36] by a large margin. Specifically, the Abs_Rel and RMSE values of DAMO are much better than OmniDepth-RectNet* by **36.5%** and **22.9%**, respectively. Although the same backbone (i.e., ResNet-50 [15]) is used in BiFuse-Equi [32], BiFuse-Fusion [32], and

Table 1. Comparison to the state-of-the-art 360° monocular depth estimation methods on the 360D Dataset. The method with * represents a model provided by the author (with its Caffe based weights being converted to Pytorch based weights), and the method with † denotes our reproduction. Note that, the results and metrics reported in [32] are different from [36], to follow the baseline method of the 360D Dataset, we convert its $RMSE(log)$ results from base-10 logarithm to natural logarithm. Besides, the results in [36] are updated at the authors’ github repository¹.

Method	RMSE	Abs_REL	RMSE(log)	δ_1	δ_2	δ_3
	Lower the better			Higher the better		
OmniDepth-UResNet [36]	0.3084	0.0946	0.1315	0.9133	0.9861	0.9962
OmniDepth-RectNet [36]	0.2432	0.0687	0.0999	0.9583	0.9936	0.9980
BiFuse-Equi [32]	0.2667	-	0.1006	0.9667	0.9920	0.9966
BiFuse-Cube [32]	0.2739	-	0.1029	0.9688	0.9908	0.9956
BiFuse-Fusion [32]	0.2440	-	0.0985	0.9699	0.9927	0.9969
OmniDepth-RectNet*	0.2297	0.0641	0.0993	0.9663	0.9951	0.9984
BiFuse-Equi*	0.2415	0.0573	0.1000	0.9681	0.9928	0.9972
DAMO	0.1769	0.0406	0.0733	0.9865	0.9966	0.9987

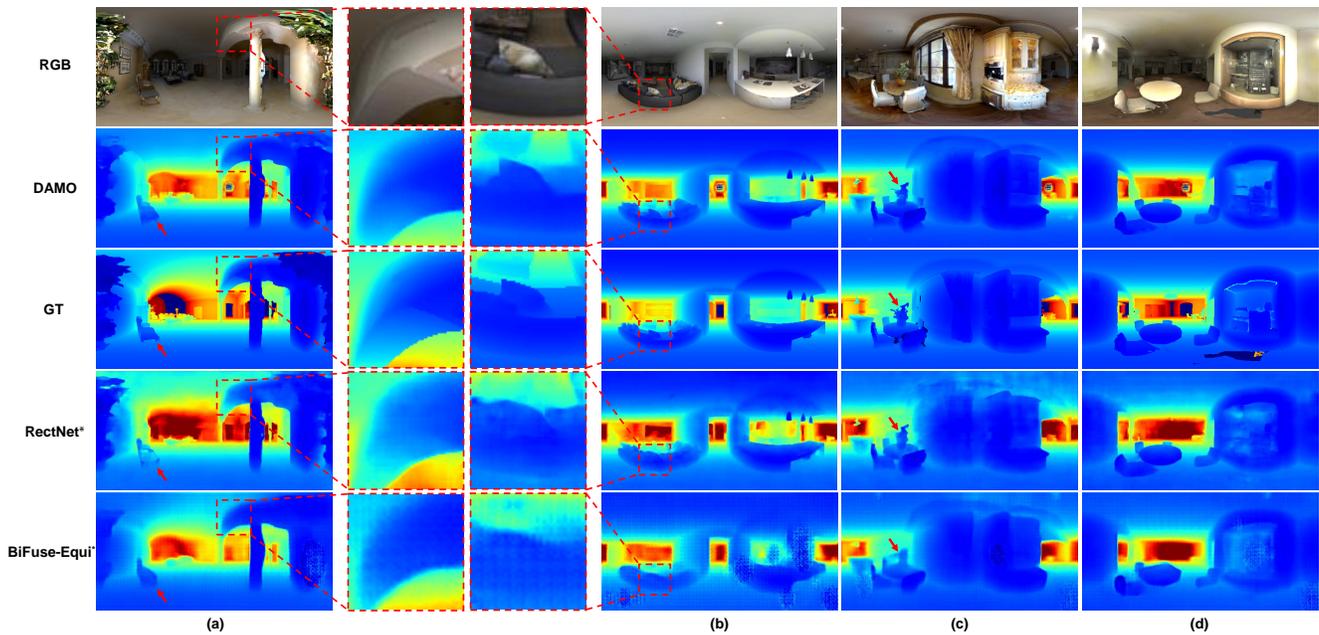


Figure 3. Qualitative Comparison on the 360D Dataset. Here, we colorize these depth maps to better distinguish the effectiveness of different methods. Points with dark color (blue) are closer than those with light color (red). The first row shows equirectangular RGB images, while invalid areas are masked with black in groundtruth (in the third row). It is obvious that our DAMO generates more accurate depth maps than existing methods, especially for those distorted areas.

our DAMO model, our model achieves significant performance improvement over the BiFuse-Equi and BiFuse-Fusion methods in almost all metrics. Note that, the BiFuse-Fusion method adopts an additional cube map representation branch and thus has more parameters to tune. Besides, the cube map representation has obvious discontinuity between neighboring faces. Although different padding schemes have been proposed to mitigate the edge effect [3, 32], the computational efficiency of this representation is still low. In contrast, our DAMO network has only one equirectangular branch to predict depth on panorama

and outperforms BiFuse-Fusion by nearly **26.7%** in RMSE.

The superiority of our DAMO network can be attributed to two reasons. **First**, DA module can adjust irregular sampling grids among distorted projection of objects in panoramas automatically and extracts rich geometric information in challenging areas (e.g., pole areas). **Second**, our weighted Berhu loss can help our model focusing on informative areas (especially for these areas near the equator) and alleviate the supervision bias caused by distortion. Consequently, the representation capability of our network

¹<https://github.com/VCL3D/360Vision/tree/master/SingleImageDepthMetrics>.

Table 2. Ablation study with strip pooling and deformable convolution. The second and third parts show the performance achieved by networks with SPM and deformable convolution (where ‘D’ represents deformable convolution), respectively. \mathcal{L} indicates the spherical-weighted loss. All models are trained in the same settings and the best results in each part are shown in boldface.

Method	Param.	RMSE	Abs_REL	RMSE(log)
		Lower the better		
Base Model		0.2129	0.0598	0.0959
+ \mathcal{L}	61.28M	0.2068	0.0562	0.0904
+ SPM		0.1879	0.0433	0.0769
+ \mathcal{L}	+ 5.83M	0.1803	0.0419	0.0749
+ D (conv5)		0.1906	0.0471	0.0789
+ \mathcal{L}	+ 0.24M	0.1882	0.0455	0.0794

on panoramas is improved.

4.2.2 Qualitative Comparison

We present several predicted depth maps from the 360D dataset in Fig. 3. It can be observed from the zoom-in regions that DAMO can predict more accurate and clearer depth maps than RectNet [36] and BiFuse-Equi [32] on panoramas. It is also shown that fine details such as chair and vase (as denoted by a small arrow in Fig. 3) can be successfully estimated. This further demonstrates the representation capability of our DAMO network on various types of objects in equirectangular images.

4.3. Ablation Study

In this section, we conduct experiments to illustrate the effectiveness of each component of our DAMO network.

4.3.1 Effectiveness of SPM

We will first analyze the benefit of exploiting SPM in our DA module. As illustrated in Table 2, it can be observed that the improvement of applying SPM is significant. While the number of parameters of the model is increased by **5.83M** (nearly 9.5% as compared to the base model), the performance is improved by about **27.5%** and **11.7%** in Abs_REL and RMSE, respectively. We argue that the trade-off between the complexity of network and advancement of its performance is well balanced. As shown in Fig.4, the base model produces many artifacts, especially on the south and north poles. In contrast, the base model with SPM can extract rich contextual information and predicts fine depth map on panoramas.

4.3.2 Deformable vs. Regular Convolution

We compare deformable convolution to its regular counterpart in this section. As shown in Table 2, the performance of our baseline is improved significantly with de-

formable convolution. Specifically, Abs_REL is improved from 0.0598 to 0.0471 (nearly by **21.2%**) and RMSE is improved from 0.2129 to 0.1906 (nearly by **10.4%**) by integrating deformable convolution into our DAMO network.

As shown in Figs. 4(c)(e), we can observe that the model with deformable convolution can generate much cleaner depth maps on walls, ceilings and floors. Note that, distortion exists mainly in these areas since ceilings and floors are always located in the north and south pole in images of the 360D dataset. Besides, long walls usually exist from left to right in indoor scenes. It is clear that deformable convolution can learn reasonable offsets to model the transformation of sampling grids and mitigate the harmful distortion effects for CNNs introduced in panoramas. Moreover, deformable convolution can produce sharper object boundaries and more accurate depth results in some difficult regions (e.g., see Fig. 4(e)).

4.3.3 Deformable Convolution with SPM

Here, we investigate how these two components of our DA module work together in synergy. As shown in Tables 1 and 2, although deformable convolution or SPM improves our base model for more than **10.0%** in Abs_REL and RMSE, the improvement of their combination is still significant. Specifically, our DAMO network outperforms the base model by nearly **31.9%** in Abs_REL and **16.9%** in RMSE. We believe the main reason for our improvement on panorama is that the DA module learns the distortion in this domain. The synergism of SPM and deformable convolution is obvious. That is, SPM activates informative regions and helps deformable convolution to focus on challenging areas by learning the transformation of sampling grids among panoramas.

4.3.4 Spherical-aware Weight vs. Regular Weight

We conduct ablation experiments to demonstrate the effectiveness of our spherical-aware weight for objective functions. Table 2 show the results of each component with/without the proposed weight. The largest improvement is achieved on our base model, its Abs_REL is improved from 0.0598 to 0.0562 (by nearly **6.0%**) and RMSE is improved from 0.2129 to 0.2068 (by nearly **2.8%**). When SPM or deformable convolution is used in the model, the spherical-aware weight can still introduce further performance improvement.

4.3.5 Generalization Analysis

To further demonstrate the generalization capability of our DAMO network, we tested our method without finetune on the validation split of the 360D dataset (i.e., SceneNet [14]). It is clear that DAMO outperforms other methods. As

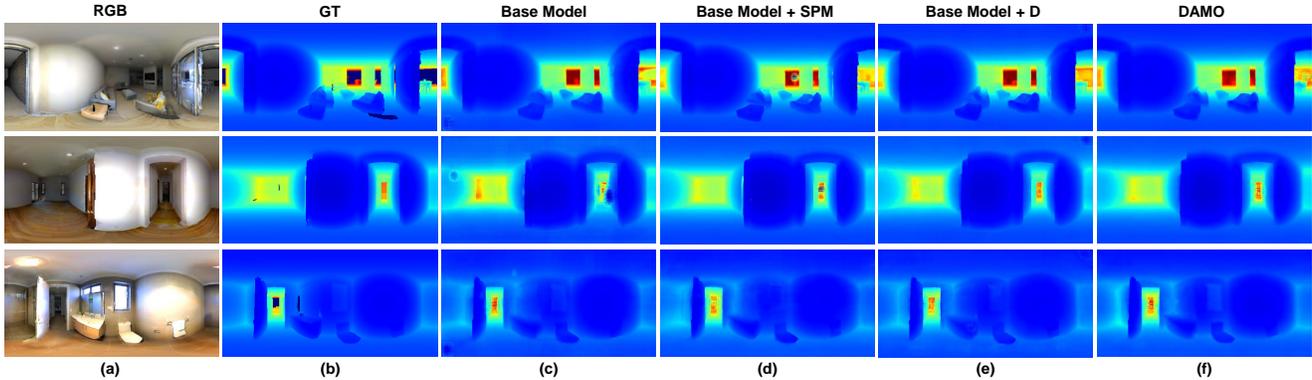


Figure 4. Ablation study with strip pooling and deformable convolution on the 360D Dataset. Models with strip pooling can predict clearer depth regions than the base model, while models with deformable convolution generate much sharper boundaries than those with regular convolution. Integrating both strip pooling and deformable convolution in our DA module can further improve the performance.

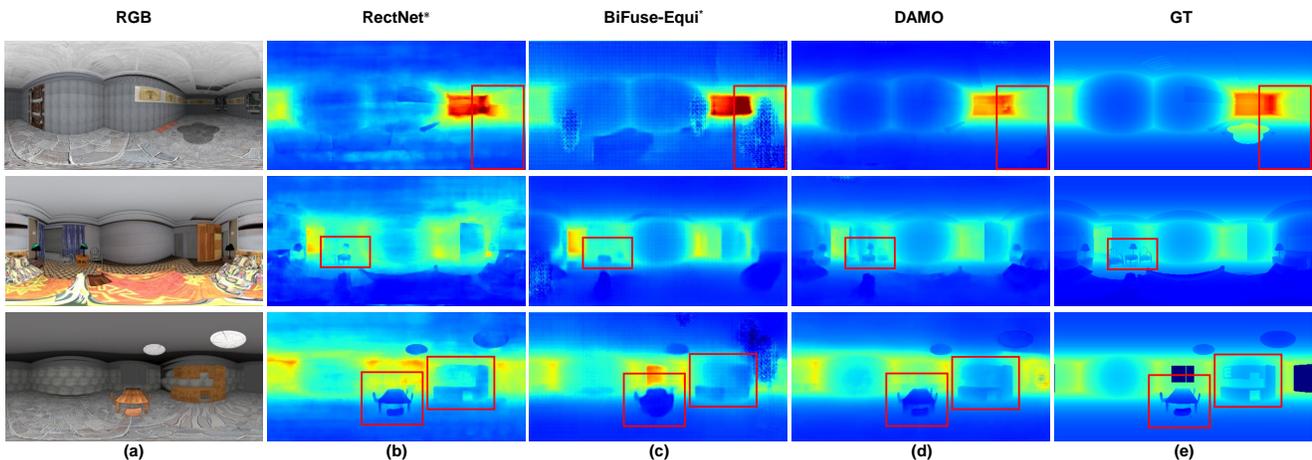


Figure 5. Generalization analysis for the proposed methods. We evaluate existing methods in SceneNet to test the generalization ability of the model.

shown in Fig. 5, DAMO predicts more fine details and its results are closer to groundtruth than other methods. Compared to the test set of 360D, the overall performance drop of our DAMO network on the validation split is lower than other methods. For instance, our RMSE is increased by 0.1088. In contrast, the RMSEs of BiFuse-Equi and OmniDepth-RectNet are increased by 0.1437 and 0.1471, respectively. That is because, RectNet and BiFuse-Equi do not consider the distortion of panorama and are therefore overfitted on the training set of 360D. More quantitative details are reported in our supplementary material.

The superiority of DAMO can be attributed to our DA module, which consists of SPM and deformable convolution. Specifically, reasonable offsets are learned with deformable convolution and its sampling grids can target distorted projection of objects at different locations of a panorama. Therefore, our DAMO obtains transformation capability of sampling grids against distortion. Besides, the

element-wise multiplication in SPM generates the connection among different channels in each input. Consequently, the transferability of DAMO is also improved. The DA module improves the overall representation capability of our network, especially along occlusion boundaries of objects (see the frame in Fig. 5).

5. Conclusion

We presented a method for omnidirectional dense depth estimation. We introduce deformable convolution to handle distortion of panorama and use strip pooling to improve the generalization ability of our DAMO network. To alleviate the supervision bias caused by distortion, we further introduce a spherical-weighted objective function to aggregate abundant information near the equator of the sphere. Experiments show that the proposed method outperforms the state-of-the-art monocular omnidirectional depth estimation

methods.

References

- [1] Iro Armeni, Sasha Sax, Amir R. Zamir, and Silvio Savarese. Joint 2D-3D-semantic data for indoor scene understanding. *arXiv preprint arXiv:1702.01105*, 2017.
- [2] Angel Chang, Angela Dai, Thomas Funkhouser, Maciej Halber, Matthias Niessner, Manolis Savva, Shuran Song, Andy Zeng, and Yinda Zhang. Matterport3D: Learning from RGB-D data in indoor environments. In *Proceedings of the International Conference on 3D Vision (3DV)*, 2017.
- [3] Hsien-Tzu Cheng, Chun-Hung Chao, Jin-Dong Dong, Hao-Kai Wen, Tyng-Luh Liu, and Min Sun. Cube padding for weakly-supervised saliency prediction in 360° videos. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1420–1429, 2018.
- [4] Xinjing Cheng, Peng Wang, Yanqi Zhou, Chenye Guan, and Ruigang Yang. Omnidirectional depth extension networks. In *Proceedings of the International Conference on Robotics and Automation (ICRA)*, 2020.
- [5] Jifeng Dai, Haozhi Qi, Yuwen Xiong, Yi Li, Guodong Zhang, Han Hu, and Yichen Wei. Deformable convolutional networks. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, pages 764–773, 2017.
- [6] Erick Delage, Honglak Lee, and Andrew Y Ng. A dynamic bayesian network model for autonomous 3d reconstruction from a single indoor image. In *IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR)*, volume 2, pages 2418–2428, 2006.
- [7] Marc Eder and Jan-Michael Frahm. Convolutions on spherical images. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*, pages 1–5, 2019.
- [8] Marc Eder, Mykhailo Shvets, John Lim, and Jan-Michael Frahm. Tangent images for mitigating spherical distortion. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2020.
- [9] David Eigen and Rob Fergus. Predicting depth, surface normals and semantic labels with a common multi-scale convolutional architecture. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2650–2658, 2015.
- [10] David Eigen, Christian Puhrsch, and Rob Fergus. Depth map prediction from a single image using a multi-scale deep network. In *Advances in Neural Information Processing Systems (NIPS)*, pages 2366–2374, 2014.
- [11] Zhicheng Fang, Xiaoran Chen, and Luc Van Gool. Towards good practice for cnn based monocular depth estimation. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision (WACV)*, 2020.
- [12] Huan Fu, Mingming Gong, Chaohui Wang, Kayhan Batmanghelich, and Dacheng Tao. Deep ordinal regression network for monocular depth estimation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2002–2011, 2018.
- [13] Jun Fu, Jing Liu, Haijie Tian, Yong Li, Yongjun Bao, Zhiwei Fang, and Hanqing Lu. Dual attention network for scene segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 3146–3154, 2019.
- [14] Ankur Handa, Viorica Patraucean, Simon Stent, and Roberto Cipolla. Scenenet: An annotated model generator for indoor scene understanding. In *Proceedings of the International Conference on Robotics and Automation (ICRA)*, pages 5737–5743, 2016.
- [15] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 770–778, 2016.
- [16] Qibin Hou, Li Zhang, Ming-Ming Cheng, and Jiashi Feng. Strip Pooling: Rethinking spatial pooling for scene parsing. In *CVPR*, 2020.
- [17] Jie Hu, Li Shen, and Gang Sun. Squeeze-and-excitation networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 7132–7141, 2018.
- [18] Yunho Jeon and Junmo Kim. Active convolution: Learning the shape of convolution for image classification. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 4201–4209, 2017.
- [19] Xu Jia, Bert De Brabandere, Tinne Tuytelaars, and Luc V Gool. Dynamic filter networks. In *Advances in Neural Information Processing Systems (NIPS)*, pages 667–675, 2016.
- [20] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *Proceedings of the 2nd International Conference on Learning Representations (ICLR)*, 2014.
- [21] Iro Laina, Christian Rupprecht, Vasileios Belagiannis, Federico Tombari, and Nassir Navab. Deeper depth prediction with fully convolutional residual networks. In *Proceedings of the International Conference on 3D Vision (3DV)*, pages 239–248. IEEE, 2016.
- [22] Yeonkun Lee, Jaeseok Jeong, Jongseob Yun, Wonjune Cho, and Kukjin Yoon. Spherphd: Applying cnns on a spherical polyhedron representation of 360° images. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 9181–9189, 2019.
- [23] Zhengfa Liang, Yulan Guo, Yiliu Feng, Wei Chen, Linbo Qiao, Li Zhou, Jianfeng Zhang, and Hengzhu Liu. Stereo matching using multi-level cost volume and multi-scale feature constancy. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, pages 1–1, 2019.
- [24] Fangchang Ma and Sertac Karaman. Sparse-to-dense: Depth prediction from sparse depth samples and a single image. In *Proceedings of the International Conference on Robotics and Automation (ICRA)*, 2018.
- [25] Ningning Ma, Xiangyu Zhang, Jiawei Huang, and Jian Sun. Weightnet: Revisiting the design space of weight networks. In *Proceedings of the European Conference on Computer Vision (ECCV)*, 2020.
- [26] A N Rajagopalan, Subhasis Chaudhuri, and Uma Mudenagudi. Depth estimation and image restoration using de-

- focused stereo pairs. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 26(11):1521–1525, 2004.
- [27] Ashutosh Saxena, Sung H Chung, and Andrew Y Ng. Learning depth from single monocular images. In *Advances in Neural Information Processing Systems (NIPS)*, pages 1161–1168, 2005.
- [28] Ashutosh Saxena, Min Sun, and Andrew Y Ng. Make3d: Learning 3d scene structure from a single still image. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 31(5):824–840, 2009.
- [29] Daniel Scharstein, Richard Szeliski, and Ramin Zabih. A taxonomy and evaluation of dense two-frame stereo correspondence algorithms. *International Journal of Computer Vision*, 47(1):7–42, 2001.
- [30] Shuran Song, Fisher Yu, Andy Zeng, Angel X Chang, Manolis Savva, and Thomas Funkhouser. Semantic scene completion from a single depth image. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*, pages 190–198, 2017.
- [31] Yuchuan Su and Kristen Grauman. Flat2sphere: Learning spherical convolution for fast features from 360° imagery. In *Advances in Neural Information Processing Systems (NIPS)*, pages 529–539, 2017.
- [32] Fu-En Wang, Yu-Hsuan Yeh, Min Sun, Wei-Chen Chiu, and Yi-Hsuan Tsai. Bifuse: Monocular 360 depth estimation via bi-projection fusion. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2020.
- [33] Wei Yin, Yifan Liu, Chunhua Shen, and Youliang Yan. Enforcing geometric constraints of virtual normal for depth prediction. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 5684–5693, 2019.
- [34] Ziheng Zhang, Yanyu Xu, Jingyi Yu, and Shenghua Gao. Saliency detection in 360° videos. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 504–520, 2018.
- [35] Xizhou Zhu, Han Hu, Stephen Lin, and Jifeng Dai. Deformable convnets v2: More deformable, better results. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 9308–9316, 2019.
- [36] Nikolaos Zioulis, Antonis Karakottas, Dimitrios Zarpalas, and Petros Daras. Omnidepth: Dense depth estimation for indoors spherical panoramas. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 453–471, 2018.